



UNIVERSIDAD  
DE GRANADA

Escuela Técnica Superior de Ingenierías Informática y de  
Telecomunicación

MÁSTER UNIVERSITARIO EN CIENCIA DE DATOS E  
INGENIERÍA DE COMPUTADORES

TRABAJO DE FIN DE MÁSTER

# Procesos Gaussianos para problemas MIL. Aplicación a la detección de hemorragias intracraneales.

Presentado por:

Francisco Miguel Castro Macías

Tutores:

Rafael Molina Soriano

*Departamento de Ciencias de la Computación e Inteligencia Artificial*

Pablo Morales Álvarez

*Departamento de Estadística e Investigación Operativa*

Curso académico 2021-2022



Procesos Gaussianos para problemas MIL.  
Aplicación a la detección de hemorragias  
intracraneales.

Francisco Miguel Castro Macías

Francisco Miguel Castro Macías *Procesos Gaussianos para problemas MIL. Aplicación a la detección de hemorragias intracraneales..*

Trabajo de fin de Máster. Curso académico 2021-2022.

**Responsable de  
tutorización**

Rafael Molina Soriano  
*Departamento de Ciencias de la Computación  
e Inteligencia Artificial*

Pablo Morales Álvarez  
*Departamento de Estadística e Investigación  
Operativa*

Máster Universitario en  
Ciencia de Datos e  
Ingeniería de  
Computadores

Escuela Técnica Superior de  
Ingenierías Informática y  
de Telecomunicación

Universidad de Granada

DECLARACIÓN DE ORIGINALIDAD

D. Francisco Miguel Castro Macías

Declaro explícitamente que el trabajo presentado como Trabajo de Fin de Máster (TFM), correspondiente al curso académico 2021-2022, es original, entendida esta, en el sentido de que no ha utilizado para la elaboración del trabajo fuentes sin citarlas debidamente.

En Granada a 4 de septiembre de 2022

Fdo: Francisco Miguel Castro Macías



*A mis padres, por creer siempre en mí.*





# Índice general

<b>Resumen</b>	<b>XI</b>
<b>English summary</b>	<b>XIII</b>
<b>1. Introducción y objetivos</b>	<b>1</b>
1.1. Introducción	1
1.2. Objetivos	2
<b>2. Aprendizaje a partir de Múltiples Instancias</b>	<b>5</b>
2.1. Motivación y orígenes	5
2.2. Definición formal	5
2.2.1. Algunos tipos de hipótesis	6
2.3. Aproximaciones al problema	7
2.4. Algunas aplicaciones	8
<b>3. Procesos Gaussianos</b>	<b>11</b>
3.1. De distribuciones gaussianas a procesos gaussianos	11
3.2. Procesos gaussianos para regresión	12
3.3. Procesos gaussianos para clasificación	15
3.4. Ajuste de hiperparámetros	16
3.5. Sparse GP	16
3.6. Aproximación variacional	18
3.6.1. Inferencia variacional	19
<b>4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias</b>	<b>23</b>
4.1. Notación y ambiente	24
4.2. VGPMIL: modelo basado en la desigualdad de Jaakkola	24
4.2.1. Relación entre las instancias y las bolsas	25
4.2.2. Distribuciones predictivas	26
4.2.3. Inferencia	26
4.3. PG-VGPMIL: modelo basado en variables Pólya-Gamma	30
4.3.1. Inferencia	32
4.3.2. Equivalencia con VGPMIL	36
4.4. G-VGPMIL: modelo basado en variables Gamma	36
4.4.1. Inferencia	37
4.4.2. Sobre la función $\theta$ y el parámetro $k$	39
<b>5. Aplicaciones</b>	<b>41</b>
5.1. Sobre los conjuntos de datos y la metodología empleada	41
5.2. Un ejemplo ilustrativo: MNIST	42
5.3. Los conjuntos musk	44
5.4. Detección de hemorragias intracraneales	45
5.4.1. Conjuntos utilizados	46

## Índice general

5.4.2. Preprocesamiento . . . . .	48
5.4.3. Metodología y arquitecturas . . . . .	49
5.4.4. Resultados y discusión . . . . .	50
<b>6. Conclusiones y trabajo futuro</b>	<b>57</b>
6.1. Conclusiones . . . . .	57
6.2. Trabajo futuro . . . . .	58
<b>A. Tablas</b>	<b>59</b>
<b>Bibliografía</b>	<b>67</b>

## Resumen

Un hematoma o hemorragia intracraneal (*IntraCranial Hemorrhage*, ICH) es una acumulación de sangre dentro del cráneo que provoca daños graves en el cerebro o incluso la muerte. Es esencial que su diagnóstico sea rápido y efectivo, ya que casi la mitad de las muertes se producen en las primeras 24 horas. Su detección se realiza a través de las tomografías computarizadas (*Computed Tomography*, CT) o TAC, que obtienen imágenes ordenadas del cráneo llamadas cortes, de forma que todas juntas componen un escáner global. Estudios recientes muestran que tras varias horas analizando estas pruebas, la fatiga lleva a los radiólogos a equivocarse al valorarlos.

Para crear métodos de triaje masivos y fiables, en la actualidad se están desarrollando Sistemas de Diagnóstico Asistido por Ordenador (*Computer Aided Diagnosis Systems*, CADs) mediante técnicas de Aprendizaje Automático, lo cual requiere grandes bases de datos etiquetadas. Anotar cada uno de los cortes de un TAC es una tarea muy costosa, por lo que es impensable solicitar el etiquetado de un gran volumen de ellos. Sin embargo, las etiquetas a nivel del escáner completo se almacenan de forma automática en el historial médico, por lo que obtenerlas no requiere ningún esfuerzo. ¿Es posible desarrollar modelos de diagnóstico usando sólo las etiquetas globales (a nivel del escáner completo) en lugar de las etiquetas de cada corte? Este es el reto al que nos enfrentamos en este trabajo.

Nuestro objetivo es tratar la detección de hemorragias intracraneales como un problema de Aprendizaje a partir de Múltiples Instancias (*Multiple Instance Learning*, MIL). Los ejemplos de entrenamiento, llamados instancias, se encuentran particionados en bolsas de instancias, y sólo conocemos la etiqueta global de cada bolsa. La etiqueta de cada instancia es desconocida y tiene asociada una gran incertidumbre. En nuestro caso, cada corte de un escáner es una instancia, y el escáner completo es la bolsa. Un escáner se clasifica como positivo si al menos uno de sus cortes es positivo (contiene evidencias de la lesión), y negativo en caso contrario.

Los Procesos Gaussianos (*Gaussian Processes*, GPs) son la herramienta que vamos a emplear para tal fin. Son modelos bayesianos que permiten aproximar funciones complejas y cuya formulación probabilística permite modelar la incertidumbre presente en los problemas de tipo MIL. El modelo *Variational Gaussian Process Multiple Instance Learning* (VGPMIL) es una de las aproximaciones más interesantes de la literatura para este problema. Emplea la *desigualdad de Jaakkola* para calcular distribuciones variacionales que estiman la probabilidad de que una bolsa o una instancia sean positivos, permitiendo obtener expresiones analíticas, pero puede conducir a soluciones subóptimas.

En este trabajo se realizan tres contribuciones dentro de esta línea de investigación.

1. La primera de ellas recibe el nombre de *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* (PG-VGPMIL) y propone, por primera vez en la literatura, el uso de variables aleatorias Pólya-Gamma en el contexto de MIL. En este trabajo se prueba su equivalencia, en el sentido *Mean Field*, con VGPMIL.
2. La segunda recibe el nombre de *Gamma Variational Gaussian Process Multiple Instance Learning* (G-VGPMIL) y propone, por primera vez en la literatura, sustituir las variables Pólya-Gamma por variables Gamma, siendo validado experimentalmente en tres problemas distintos.

## *Resumen*

3. La tercera contribución, en este caso de tipo práctico, consiste en el uso de G-VGPMIL junto con redes neuronales y mecanismos de atención para resolver el problema de la detección de hemorragias intracraneales. El modelo G-VGPMIL obtiene los mejores resultados de la literatura usando sólo etiquetas a nivel de escáner completo, lo cual es muy prometedor y establece el camino a seguir para desplegar esta solución en centros médicos.

**Palabras clave:** Hemorragia intracraneal, Tomografía Computarizada, Aprendizaje a partir de Múltiples Instancias, Procesos Gaussianos, Inferencia Variacional, desigualdad de Jaakkola, Pólya-Gamma, Gamma, Aprendizaje Automático, Aprendizaje Profundo, Redes Neuronales, Atención.

## English summary

An IntraCranial Hemorrhage (ICH) is a collection of blood inside the skull that results in severe brain damage or even death. Early and effective diagnosis is essential, as almost half of the deaths occur within the first 24 hours. Its detection is done through computed tomography (CT) scans, which obtain several ordered images of the skull called slices, so that all together they make up an overall CT scan. Recents studies show that after several hours analyzing these tests, fatigue leads radiologists to make mistakes when evaluating them.

To create massive and reliable triage methods, Computer Aided Diagnosis Systems (CADs) are currently being developed using Machine Learning techniques, which requires large labeled databases. Annotating every single CT slice is a very expensive task, so it is not possible to request the labeling of a large volume of CT slices. However, the label of the entire scan is automatically stored in the medical record, so obtaining them is effortless. Is it possible to develop diagnostic models using only the scan-level labels instead of the slice-level labels? This is the challenge we face in this work.

Our goal is to treat ICH detection as a Multiple Instance Learning (MIL) problem. The training examples, called instances, are partitioned into bags of instances, and we only know the global label of each bag. The label of each instance is unknown and has a large uncertainty associated with it. In our case, each slice of a scan is an instance, and the entire scan is the bag. A scan is classified as positive if at least one of its slices is positive (contains evidence of the lesion), and negative otherwise.

Gaussian Processes (GPs) are the tool we are going to use for this purpose. They are Bayesian models that allow us to approximate complex functions and whose probabilistic formulation allows us to model the uncertainty present in MIL type problems. The Variational Gaussian Process Multiple Instance Learning (VGPMIL) model is one of the most interesting approaches in the literature to solve this problem. It uses Jaakkola's inequality to compute variational distributions that estimate the probability that a bag or instance will be positive, which allows analytical expressions to be obtained, but may lead to suboptimal solutions.

This paper makes three contributions within this line of research.

1. The first one is called *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* (PG-VGPMIL) and proposes, for the first time in the literature, the use of Pólya-Gamma random variables in the context of MIL. In this paper we prove its equivalence, in the sense of *Mean Field*, with VGPMIL.
2. The second one is called *Gamma Variational Gaussian Variational Process Multiple Instance Learning* (G-VGPMIL) and proposes, for the first time in the literature, to replace the Pólya-Gamma variables by Gamma variables, being validated experimentally in three different problems.
3. The third one consists of the use of G-VGPMIL together with neural networks and attention mechanisms to solve the problem of ICH detection. G-VGPMIL model obtains the best results in the literature using only labels at scan-level, which is very promising and sets the way forward for deploying this solution in medical centers.

*English summary*

**Keywords:** Intracranial hemorrhage, Computed Tomography, Multiple Instance Learning, Gaussian Processes, Variational Inference, Jaakkola Inequality, Pólya-Gamma, Gamma, Machine Learning, Deep Learning, Neural Networks, Attention.

# 1. Introducción y objetivos

El objetivo de este capítulo es ofrecer una visión general del propósito del trabajo y aclarar los objetivos inicialmente planteados y el grado en que estos se han cumplido. Primero ofrecemos una introducción sobre los problemas que tratamos y después planteamos los objetivos, conectando cada uno con la parte de la memoria que lo cubre.

## 1.1. Introducción

Un hematoma o hemorragia intracraneal (*IntraCranial Hemorrhage*, ICH) es una acumulación de sangre dentro del cráneo que normalmente se debe a la rotura de un vaso sanguíneo [35]. La sangre que entra al cráneo provoca que aumente la presión cerebral, lo que conduce a daños graves en el cerebro o incluso a la muerte. Alrededor del 40 % de los pacientes fallecen durante el primer mes, la mayoría en los dos primeros días, y sólo un 20 % son independientes al cabo de 6 meses [16]. De hecho, casi el 30 % de las muertes se producen en las primeras 24 horas [40], por lo que el diagnóstico rápido y eficaz de esta lesión es esencial.

El diagnóstico se suele realizar a través de las tomografías computarizadas (*Computed Tomography*, CT) o TAC [31]. Esta técnica obtiene varios escáneres del cráneo llamados cortes, de forma que todos juntos componen un escáner global. Sin embargo, diversos estudios muestran que tras varias horas analizando estas pruebas, la fatiga lleva a los radiólogos a equivocarse en el diagnóstico [41]. Para conseguir mejores métodos de triaje se están desarrollando Sistemas de Diagnóstico Asistido por Ordenador (*Computer Aided Diagnosis Systems*, CADs) basados en técnicas de Aprendizaje Automático (*Machine Learning*, ML) [8].

Muchas de las soluciones que se están proponiendo se basan en modelos de Aprendizaje Profundo (*Deep Learning*, DL) [18] que son entrenados de forma directa, usando las etiquetas de todos los cortes. Sin embargo, para que los modelos sean fiables y obtengan buenos resultados, necesitamos grandes base de datos de cortes de TAC etiquetados. Anotar cada uno de los cortes de un TAC es una tarea muy costosa, por lo que, debido al escaso personal de radiología y al poco tiempo del que disponen, es impensable solicitar el etiquetado de un gran volumen de datos. Sin embargo, las etiquetas a nivel del escáner completo se almacenan de forma automática en el historial médico al diagnosticar a los pacientes, por lo que obtenerlas no requiere ningún esfuerzo adicional. La solución pasa por desarrollar modelos que usen sólo las etiquetas globales (a nivel del escáner completo) en lugar de las etiquetas de cada corte.

En este trabajo tratamos el problema de la detección de hemorragias intracraneales como un problema de Aprendizaje a partir de Múltiples Instancias (*Multiple Instance Learning*, MIL) [50]. Este es una variación del aprendizaje supervisado clásico que alivia el esfuerzo de anotación requerido. Los ejemplos de entrenamiento, llamados instancias, se encuentran particionados en bolsas de instancias, y sólo conocemos la etiqueta global de cada bolsa. La etiqueta de cada instancia es desconocida y tiene asociada una gran incertidumbre. En el problema al que nos enfrentamos cada corte de un escáner es una instancia, y el escáner completo es la bolsa. Un escáner se clasifica como positivo si al menos uno de sus cortes es positivo, y negativo en caso contrario.

## 1. Introducción y objetivos

Para construir una solución, nos serviremos de los Procesos Gaussianos (*Gaussian Processes*, GPs) [34]. Son modelos bayesianos que permiten aproximar funciones complejas y cuya formulación probabilística permite modelar la incertidumbre presente en los problemas de tipo MIL. El modelo *Variational Gaussian Process Multiple Instance Learning* (VGPMIL) [19] es una de las aproximaciones más fructíferas de la literatura. Utiliza el modelo de observación logístico y destaca por su forma de tratar la relación entre las etiquetas de las instancias y las bolsas. Aproxima la distribución a posteriori usando inferencia variacional, para lo que emplea la *desigualdad de Jaakkola* [22]. Esto le permite obtener expresiones analíticas para las actualizaciones de los parámetros, aunque puede conducir a soluciones subóptimas.

En este trabajo estamos interesados en el uso de técnicas de *data augmentation*<sup>1</sup> para realizar inferencia variacional en estos escenarios. Motivados por el buen comportamiento de las variables aleatorias Pólya-Gamma en problemas de clasificación clásicos [51], nos preguntamos qué ocurre al introducir estas variables en un modelo de GPs para MIL. La respuesta a esta pregunta es la primera contribución del trabajo, que recibe el nombre de *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* (PG-VGPMIL). Uno de los resultados más importantes de este trabajo es que es equivalente, en sentido *Mean Field*, a VGPMIL. Por otra parte, la sustitución de variables Pólya-Gamma por variables Gamma conduce al nuevo modelo *Gamma Variational Gaussian Process Multiple Instance Learning* (G-VGPMIL), la segunda contribución del trabajo.

Este último modelo es validado experimentalmente en tres problemas distintos antes de proponer una solución definitiva al problema de la detección de hemorragias intracraneales. Siguiendo la metodología de [53], proponemos una arquitectura que combina redes neuronales con mecanismos de atención para extraer características, junto con G-VGPMIL para realizar predicciones a nivel de corte y de escáner completo. Los resultados obtenidos se comparan con las aproximaciones existentes en la literatura.

## 1.2. Objetivos

Cuando se propuso el trabajo, inicialmente se plantearon los siguientes objetivos:

- A. Estudiar los modelos basados en GPs y comprender los problemas derivados de utilizar datos masivos y qué soluciones existen.
- B. Estudiar los problemas tipo MIL y cómo los GPs se pueden usar en ellos.
- C. Desarrollar un modelo basado en GPs para detectar hemorragias intracraneales a partir de tomografías computarizadas.

Todos han sido completamente satisfechos, como justificamos a continuación.

- El **Capítulo 2** cubre la primera parte del objetivo B. En él se recoge una definición general del concepto de MIL [50], así como un estudio de las hipótesis existentes para modelar la relación entre las instancias y las bolsas. También repasamos los esfuerzos realizados en la literatura para encontrar una solución. Por último, ofrecemos una visión general de las tareas o aplicaciones en las que estos problemas están presentes.

---

<sup>1</sup>En el contexto de este trabajo, *data augmentation* se refiere a una técnica consistente en añadir variables aleatorias para aumentar la probabilidad conjunta del modelo. Marginalizando estas variables recuperamos la distribución original. En nuestro caso, condicionar sobre estas variables nos permite calcular analíticamente la esperanza en el modelo de observación logístico.



- El estudio de GPs se aborda con total rigurosidad en el **Capítulo 3**. En él motivamos los GPs y presentamos su formulación en problemas de regresión y clasificación [34]. El cálculo de la distribución predictiva conlleva un coste computacional muy elevado, de donde se derivan todos los problemas de usar grandes volúmenes de datos con estos modelos. La aproximación más exitosa para solventar este asunto recibe el nombre de *Sparse GPs* y se estudia en profundidad en la **Sección 3.5**. Con esto, el objetivo A queda totalmente cubierto.
- La segunda parte del objetivo B se completa en el **Capítulo 4**. En él usamos a los GPs como herramienta para resolver los problemas de MIL a los que nos enfrentamos. Este es un capítulo muy importante ya que en él se presentan de forma teórica dos de nuestras contribuciones.
  1. Empezamos estudiando el trabajo VGPMIL [19], que es uno de los modelos que componen el estado del arte a la hora de usar GPs para resolver problemas de MIL. Como ya hemos aclarado, VGPMIL emplea la desigualdad de Jaakkola [22], que le permite llegar a expresiones analíticas de las mismas.
  2. Motivados por su buen comportamiento en problemas de clasificación clásicos [51], proponemos por primera vez el uso de variables Pólya-Gamma para poder realizar inferencia. Como resultado obtenemos PG-VGPMIL, del cual probamos su equivalencia, en el sentido *Mean Field*, con VGPMIL.
  3. La definición de las variables Pólya-Gamma nos conduce a considerar variables Gamma, obteniendo el nuevo modelo G-VGPMIL, que usa variables Gamma en lugar de Pólya-Gamma.
- El **Capítulo 5** es el más práctico. En él se valida experimentalmente el nuevo modelo G-VGPMIL, comparándolo con VGPMIL en tres problemas distintos. El objetivo C se corresponde con la **Sección 5.4**, en la que presentamos una arquitectura inspirada en [53] para la detección de hemorragias intracraneales. Además, se hace un estudio exhaustivo de los resultados obtenidos.



## 2. Aprendizaje a partir de Múltiples Instancias

En este capítulo se ofrece una introducción general al concepto de *Aprendizaje a partir de Múltiples Instancias (Multiple Instance Learning)*, el cual juega un papel clave en este trabajo. Empezaremos motivando este concepto y ofreciendo la intuición detrás del mismo. Definiremos rigurosamente qué es un problema de este tipo, revisaremos las aproximaciones más populares para resolverlo y señalaremos algunas de las aplicaciones en las que se emplea.

### 2.1. Motivación y orígenes

Uno de los mayores inconvenientes del aprendizaje supervisado es que no siempre podemos disponer de suficientes ejemplos etiquetados con los que entrenar un modelo. En muchos ámbitos, como por ejemplo la medicina, el proceso de etiquetado puede ser costoso y lento, por lo que es necesario explorar alternativas que alivien esta necesidad de anotación masiva.

El *aprendizaje a partir de múltiples instancias (Multiple Instance Learning, MIL)* es una variación del aprendizaje supervisado clásico que surge para dar respuesta a este problema. En lugar de etiquetar cada una de las instancias por separado, se etiquetan por grupos o *bolsas*, de forma que se ofrece una sola etiqueta para cada bolsa. A partir de esta información, se pretende aprender a etiquetar correctamente nuevas instancias y bolsas. De esta forma, el esfuerzo de anotación se reduce notablemente.

El término MIL fue usado por primera vez en [10] mientras se investigaba un problema relacionado con el desarrollo de nuevos fármacos. Se pretendía predecir si una molécula de un fármaco podría o no vincularse a una proteína relacionada con una enfermedad. Se disponían de ejemplos de moléculas que sí consiguen vincularse y de moléculas que no lo consiguen. El factor más importante para que esto ocurra es la forma (como ocurre con una llave y una cerradura). Así pues, cada bolsa representa a una molécula, y las instancias dentro de ella son todas las formas que puede adoptar. Una etiqueta positiva en la bolsa quiere decir que existe una forma con la que la molécula se adhiere a la proteína, mientras que una etiqueta negativa quiere decir que no hay ninguna forma con la que ocurra esto.

El problema al que nos enfrentamos en este trabajo encaja de forma natural en este paradigma. Cada escáner o tomografía completa es una bolsa de instancias. Una instancia se corresponde con un corte. El experto (en este caso un especialista médico) etiqueta el escáner completo, lo cual consiste en determinar si hay presencia o no de hemorragia. De forma natural, si un paciente sufre una hemorragia, al menos uno de los cortes así lo evidenciará. Obsérvese que el etiquetado del escáner completo (bolsa) es mucho más fácil que el de los cortes y que, además, la información sobre la presencia o no de hemorragia puede encontrarse en el historial clínico del paciente.

### 2.2. Definición formal

Presentamos a continuación la definición de MIL proporcionada en [50]. Denotamos por  $\mathcal{X}$  al espacio de instancias y por  $\mathcal{T}$  al espacio de etiquetas de las bolsas, respectivamente. Una *bolsa*

## 2. Aprendizaje a partir de Múltiples Instancias

es un elemento de  $\mathcal{P}(\mathcal{X})$ <sup>1</sup>. Un *concepto* es una función ideal

$$v: \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{T},$$

que asigna una etiqueta a cada bolsa. El objetivo de un problema MIL es aproximar esta función, en caso de que exista. Para ello, dispondremos de un conjunto de entrenamiento finito compuesto por bolsas etiquetadas, es decir, pares de la forma  $(B, T) \in \mathcal{P}(\mathcal{X}) \times \mathcal{T}$ . Las bolsas tienen intersección vacía y se supone que  $T = v(B)$ .

En la mayoría de problemas de MIL las instancias también tienen una etiqueta, aunque esta es desconocida. Es el caso de la detección de hemorragias intracrániales a partir de escáneres de TAC, donde cada corte puede presentar o no evidencias de la enfermedad, pero sólo conocemos la etiqueta del escáner completo. En estos casos existe un espacio  $\mathcal{Y}$  con las etiquetas de las instancias y cada instancia  $x$  tiene una etiqueta  $y_x \in \mathcal{Y}$ . Diremos que una bolsa  $B$  contiene un *concepto*  $c \in \mathcal{Y}$  si existe  $x \in B$  tal que  $y_x = c$ .

### 2.2.1. Algunos tipos de hipótesis

La forma en que se relacionan las bolsas y las instancias individuales se denomina *supuesto*, *suposición* o *hipótesis* y es fundamental en el problema. Vamos a enumerar algunas de las más populares.

**Hipótesis normal.** La hipótesis más común, y con la que nosotros trabajaremos en el problema al que nos enfrentamos, ocurre en problemas de clasificación binarios, en los que  $\mathcal{Y} = \mathcal{T} = \{0, 1\}$ . Consiste en tomar  $T = \max_{x \in B} \{y_x\}$ . Es decir, una bolsa se etiqueta como positiva si y sólo si al menos una de las instancias tiene etiqueta positiva. Recordemos que la etiqueta de las instancias es desconocida.

**Jerarquía de Weidman et al.** En [50] se propusieron tres paradigmas de complejidad incremental para modelar problemas de clasificación en los que  $\mathcal{T} = \{0, 1\}$  y el espacio  $\mathcal{Y}$  es discreto.

- **Hipótesis basada en la presencia de instancias.** Existe un subconjunto de conceptos requeridos  $\mathcal{C} \subset \mathcal{Y}$ . Una bolsa se etiqueta como positiva si y sólo si contiene al menos uno de estos conceptos. Esto es,

$$T = \begin{cases} 1 & \text{si } |\{y_x : x \in B\} \cap \mathcal{C}| \geq 1, \\ 0 & \text{en otro caso} \end{cases}$$

Tomando  $\mathcal{C} = \{1\}$  estamos en la hipótesis normal.

- **Hipótesis basada en un umbral.** Como en el anterior, se define un subconjunto de conceptos requeridos  $\mathcal{C} \subset \mathcal{Y}$  y un número entero positivo  $n$ . Una bolsa  $B$  se etiqueta como positiva si y sólo si contiene al menos  $n$  de estos conceptos. Es decir,

$$T = \begin{cases} 1 & \text{si } |\{y_x : x \in B\} \cap \mathcal{C}| \geq n, \\ 0 & \text{en otro caso} \end{cases}$$

Tomando  $n = 1$  y  $\mathcal{C} = \{1\}$  estamos en el caso anterior.

---

<sup>1</sup> $\mathcal{P}(\mathcal{X})$  denota las partes de  $\mathcal{X}$ , es decir, el conjunto de todos los subconjuntos de  $\mathcal{X}$ .

- **Hipótesis basada en dos umbrales.** Nuevamente existe un subconjunto de conceptos requeridos  $\mathcal{C} \subset \mathcal{Y}$ . Se definen dos enteros positivos  $n_1$  y  $n_2$ , con  $n_1 < n_2$ . Una bolsa  $B$  se etiqueta como positiva si y sólo si contiene al menos  $n_1$  y como máximo  $n_2$  de estos conceptos. Esto es,

$$T = \begin{cases} 1 & \text{si } n_1 \leq |\{y_x : x \in B\} \cap \mathcal{C}| \leq n_2, \\ 0 & \text{en otro caso} \end{cases}$$

Tomando  $n_2$  suficientemente grande se tiene el caso anterior.

**Hipótesis colectiva.** En lugar de ver a las bolsas como subconjuntos fijos, se modela una bolsa  $B$  como una distribución de probabilidad con densidad  $p(x | B)$ . También se modelan las etiquetas de las instancias  $y$  como una distribución con densidad  $p(y | x)$ . El objetivo es modelar la distribución  $p(y | B) = \int_{\mathcal{X}} p(y | x)p(x | B)dx$ .

## 2.3. Aproximaciones al problema

A partir de ahora vamos a centrarnos en problemas de clasificación binarios donde se cumple la hipótesis normal. En ellos disponemos de dos tipos de etiquetas (positiva o negativa), tanto a nivel de bolsa como de instancia. La presencia o no de una instancia positiva determina la etiqueta de la bolsa. A continuación, usando [23], revisamos algunas de las soluciones que se han propuesto para solventar los problemas MIL de este tipo.

Las primeras aproximaciones se centraron en diseñar algoritmos a medida para el dominio del problema que trataban. Dos ejemplos de ello son las soluciones que se presentan en [10] y en [26]. Posteriormente se comenzaron a diseñar técnicas que pretendían predecir las etiquetas de las bolsas empleando núcleos o *kernels* que medían la similitud entre las bolsas. También se comenzó a tratar como una aproximación para modelar el problema de la falta de etiquetas o la presencia de ruido en ellas. Surgieron así soluciones que modifican la formulación original de los modelos SVM (*Support Vector Machine*) para adaptarlos a MIL, como por ejemplo mi-SVM [1]. Sin embargo, requieren técnicas de optimización costosas que impiden usarlos con conjuntos grandes de datos.

Una de las tendencias más fructíferas consiste en encontrar el ejemplo "más positivo" dentro de una bolsa positiva, considerando que este es el responsable de la etiqueta de la bolsa. A este ejemplo se le suele denominar "testigo". MI-SVM [1] adapta la formulación de SVM con el objetivo de maximizar la distancia de tal instancia al hiperplano separador. En MICA [25] se emplea programación lineal y se construye un "testigo" como una combinación convexa de los ejemplos de una bolsa. En [58] se generaliza el trabajo de [26] incluyendo esta idea.

En los últimos años, muchos de los trabajos se centran en adaptar modelos de Aprendizaje Profundo (*Deep Learning, DL*). Este tipo de modelos son muy populares en todos los problemas de aprendizaje debido a su alta flexibilidad y su capacidad para aproximar funciones complejas. Por ello, los esfuerzos se centran en encontrar nuevas arquitecturas y mecanismos que permitan entrenar redes neuronales complejas usando bolsas de instancias [54, 48]. Dentro de esta línea, la aproximación que constituye el estado del arte emplea mecanismos de atención [45] para ponderar cada una de las instancias [20].

Junto a estos mecanismos de atención, los modelos probabilísticos constituyen el estado del arte para MIL. Los Procesos Gaussianos (*Gaussian Processes, GPs*) [34] son modelos bayesianos

## 2. Aprendizaje a partir de Múltiples Instancias

que aproximan funciones muy complejas a la vez que otorgan una estimación de la incertidumbre presente en las predicciones. Su formulación se puede adaptar a los problemas MIL para diseñar potentes modelos [23, 19, 47] que posteriormente se combinan con mecanismos de atención, y que constituyen el estado del arte en los problemas MIL.

### 2.4. Algunas aplicaciones

El enfoque MIL es muy potente ya que reduce drásticamente la necesidad de datos completamente etiquetados. Por ello se ha intentado aplicar en diversos campos, algunos de los cuales comentamos a continuación. Puede consultarse más información en [7].

**Biología y química.** A raíz de la ya comentada tarea de clasificación de moléculas han surgido numerosos problemas del ámbito de la biología en los que se puede aplicar este enfoque. Las entidades biológicas complejas, como compuestos, moléculas o genes, son modeladas como bolsas cuyas partes pueden inducir un efecto u otro. Es mucho más difícil observar cada una de las partes de estas entidades que observarlas como un conjunto y ofrecer una sólo etiqueta. Algunos ejemplos concretos son la predicción de funciones de genes [13], la predicción de la afinidad de péptidos [56] (ejemplo de un problema de regresión) y el descubrimiento de lugares de enlace en la expresión de nuevos genes [3].

**Visión por computador.** En este campo el enfoque MIL permite caracterizar patrones complejos a partir de otros más simples que los componen. Además, permite entrenar modelos complejos reduciendo el esfuerzo de anotación, lo que lo está convirtiendo en un paradigma popular dentro del diagnóstico asistido por ordenador.

- La *recuperación de imágenes basada en el contenido* (*Content Based Image Retrieval, CBIR*) [24] pretende categorizar las imágenes basándose en los objetos y conceptos que contiene. La imagen completa se considera como una bolsa que dentro contiene diversos objetos, cuya localización exacta no es importante [58]. Así, las imágenes se particionan mediante alguna técnica (parches, *keypoints*, segmentación, ...) para conformar las instancias [49]. En la [Figura 2.1](#) se muestra una imagen de un oso que se divide en varias regiones. De estas regiones sólo la que contiene al animal tiene etiqueta positiva, mientras que las demás tienen etiqueta negativa.
- En la *localización de objetos* [43] y *segmentación* [15] típicamente se han requerido imágenes fuertemente anotadas. Es decir, en el caso de la localización, era necesario dar la posición de una caja que englobase al objeto, mientras que en segmentación se necesitaban etiquetas a nivel de píxel. Estas tareas requieren un tiempo y esfuerzo considerables, por lo que la investigación se centra actualmente usar imágenes débilmente etiquetadas [57]. La aparición de MIL ha conseguido aliviar estos requerimientos y la comunidad investigadora pretende desarrollar soluciones para estos problemas usando sólo etiquetas a nivel de la imagen completa [2, 46].
- En tareas de *diagnóstico asistido por computador* (*Computer Aided Diagnosis, CAD*) [11] el paradigma MIL está adquiriendo una grandísima popularidad. Hay que tener en cuenta que el etiquetado de imágenes médicas requiere un grandísimo esfuerzo por parte de los médicos, por lo que usar soluciones basadas en imágenes fuertemente etiquetadas es impracticable. El enfoque MIL es muy apropiado en este tipo de situaciones ya que en

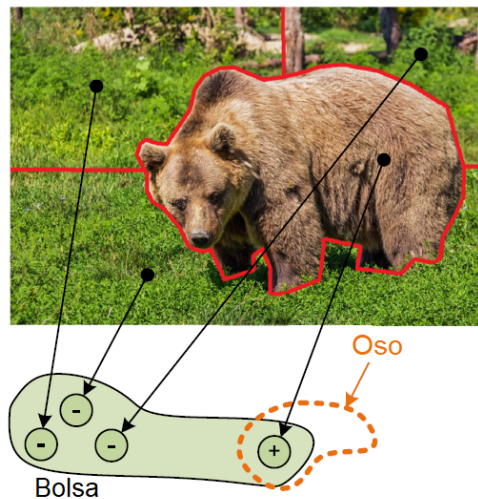


Figura 2.1.: La imagen se segmenta para conformar las instancias. Cada una de estas instancias tiene una etiqueta oculta. La etiqueta de la bolsa (imagen completa) indica la presencia de un oso. Imagen adaptada de [7].

un paciente existen regiones anormales y sanas, mientras que en una persona sana sólo existen regiones sanas. Algunas aplicaciones concretas son la clasificación de imágenes histológicas [36] o la detección de demencia a partir de resonancias magnéticas del cerebro [44].

Como vemos, el paradigma MIL no sólo sirve para formular nuevos problemas de aprendizaje, sino también para reformular los problemas ya existentes de una forma alternativa y natural. La definición de un problema MIL sólo exige el etiquetado de grupos de instancias, lo cual reduce el esfuerzo de anotación y permite utilizar conjuntos de datos débilmente anotados.

En este capítulo hemos presentado qué son los problemas MIL y la motivación detrás de ellos. Como ya hemos señalado, el objetivo de este trabajo es resolver un problema concreto de este tipo. Los siguientes capítulos preparan el camino para ello. En particular, el capítulo (Capítulo 3) presenta los Procesos Gaussianos, que será la herramienta que emplearemos para hacer frente al problema. Tras haberlos comprendido, podremos adaptarlos a los problemas MIL en el Capítulo 4, en donde señalaremos nuestras contribuciones al estado del arte.





### 3. Procesos Gaussianos

Un Proceso Gaussiano (*Gaussian Process, GP*) puede entenderse como una distribución de probabilidad normal sobre el espacio infinito dimensional de todas las funciones [34]. Numerosos problemas dentro del aprendizaje automático y estadístico consisten en aproximar funciones desconocidas. Para ello, es necesario realizar alguna suposición *a priori* sobre la naturaleza de la función. Estas suposiciones se codifican en un GP *a priori* que, tras un procedimiento de optimización, dará lugar a un proceso *a posteriori* que se ajusta a los datos [9].

El lector puede preguntarse, de forma natural, por qué hemos decidido usar GPs para solucionar los problemas MIL. Un aspecto clave en el escenario MIL es la gran incertidumbre presente en las etiquetas de las instancias. Por ello, los métodos bayesianos son una buena elección ya que permiten tenerla en cuenta. Como muestra este capítulo y el siguiente, los GPs poseen una formulación probabilística en la que la hipótesis de los problemas MIL se puede incluir directamente y que permite cuantificar esa incertidumbre.

En este capítulo presentamos formalmente los GPs y aquellos elementos teóricos necesarios para modelizar el problema al que nos enfrentamos. Comenzaremos motivando y definiendo los GPs para después analizar su uso en problemas de regresión y clasificación. Posteriormente presentaremos una aproximación conocida como *Procesos Gaussianos Dispersos (Sparse GPs, SGPs)* y que posibilita el uso de GPs con grandes conjuntos de datos. Este capítulo está construido en torno al contenido de [34, 9, 4].

#### 3.1. De distribuciones gaussianas a procesos gaussianos

Consideremos una colección de variables aleatorias reales  $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n]$  tales que su distribución conjunta es una normal multivariante, también llamada distribución gaussiana multivariante. Una importante propiedad es que la marginalización de  $\mathbf{X}$  a cualquier subconjunto de variables es también una distribución gaussiana. Nos hacemos la siguiente pregunta: ¿y si  $\mathbf{X}$  fuese la marginalización de una colección (posiblemente infinita) de variables? Surge así la siguiente definición:

**Definición 3.1** (Proceso Gaussiano, GP). Un Proceso Gaussiano (GP) es una colección de variables aleatorias tales que cualquier subconjunto finito suyo sigue una distribución gaussiana. Si  $f$  es un GP en  $\mathcal{X}$ , escribiremos  $f \sim \mathcal{GP}(m, k)$ , donde

$$\begin{aligned} m: \mathcal{X} &\rightarrow \mathbb{R}, & m(x) &= \mathbb{E}[f(x)], \\ k: \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R}, & k(\mathbf{x}_1, \mathbf{x}_2) &= \mathbb{E}[(f(\mathbf{x}_1) - m(\mathbf{x}_1))(f(\mathbf{x}_2) - m(\mathbf{x}_2))]. \end{aligned}$$

Las aplicaciones  $m$  y  $k$  reciben el nombre de *media* y *covarianza* o *núcleo*, respectivamente.

*Observación 3.1.* De la definición deducimos que un GP induce una distribución de probabilidad sobre el espacio de todas las funciones. Esto es, si  $f: \mathbb{R}^D \rightarrow \mathbb{R}$  es tal que  $f \sim \mathcal{GP}(m, k)$  y  $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_N^\top] \in \mathbb{R}^{D \times N}$  entonces  $\mathbf{f} = f(\mathbf{X}) = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$  sigue una distribución gaussiana de media  $m(\mathbf{X}) = [m(\mathbf{x}_1), \dots, m(\mathbf{x}_N)]^\top$  y matriz de covarianza

### 3. Procesos Gaussianos

$(k(\mathbf{X}, \mathbf{X}))_{ij} = k(\mathbf{X}, \mathbf{X})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . Es común escribir  $\mathbf{f} \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X}))$ , omitiendo así la dependencia del conjunto de entrenamiento.

Como veremos a lo largo de este trabajo, no perdemos generalidad al elegir  $m(\mathbf{x}) = 0$  para cada  $\mathbf{x} \in \mathcal{X}$ . El núcleo suele elegirse de entre una familia paramétrica de los mismos. Los parámetros del núcleo reciben el nombre de *hiperparámetros*. Un ejemplo ampliamente usado es el núcleo de funciones de base radial, también llamado núcleo gaussiano.

**Definición 3.2** (Núcleo de funciones de base radial). Sea  $l \in \mathbb{R} \setminus \{0\}$ . El núcleo de funciones de base radial (*Radial Basis Function kernel, RBF*), se define como

$$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2l^2}\right), \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^D$$

El hiperparámetro  $l$  se denomina *lengthscale*.

## 3.2. Procesos gaussianos para regresión

Supongamos que disponemos de un conjunto de entrenamiento con  $N$  observaciones,  $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \{1, \dots, N\}\} \subset \mathbb{R}^D \times \mathbb{R}$ , donde  $\mathbf{x}$  denota el vector de las variables observadas de un ejemplo e  $y$  denota la variable objetivo. Podemos recoger los ejemplos en una matriz  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  y en un vector  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$ . Para modelizar la relación entre  $\mathbf{x}$  e  $y$  suponemos

$$y(\mathbf{x}) = f(\mathbf{x}) + \epsilon, \quad (3.1)$$

donde  $f: \mathbb{R}^D \rightarrow \mathbb{R}$  es una función desconocida y  $\epsilon \sim \mathcal{N}(0, \sigma_r^2)$ . Suponemos la existencia de cierto ruido  $\epsilon$  que hace diferir los valores observados de los que produce  $f$ . Cuando tomamos  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  nos encontramos en el escenario de la regresión lineal.

En lugar de suponer una forma funcional para  $f$ , podemos inducir una distribución sobre el espacio de las funciones usando un proceso gaussiano. Tomamos así  $f \sim \mathcal{GP}(0, k)$ , por lo que  $\mathbf{f} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, k(\mathbf{X}, \mathbf{X}))$ . De (3.1) tenemos  $\mathbf{y} | \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma_r^2 \mathbf{I})$ . Es importante conocer la siguiente definición.

**Definición 3.3** (Verosimilitud marginal). En las condiciones anteriores, la verosimilitud marginal es

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y}, \mathbf{f} | \mathbf{X}) d\mathbf{f} = \int p(\mathbf{y} | \mathbf{f}, \mathbf{X}) p(\mathbf{f} | \mathbf{X}) d\mathbf{f}.$$

*Observación 3.2.* Se cumple que

$$\log p(\mathbf{f} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^\top k(\mathbf{X}, \mathbf{X})^{-1} \mathbf{y} - \frac{1}{2} \log |k(\mathbf{X}, \mathbf{X})| - \frac{N}{2} \log(2\pi),$$

de donde

$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^\top \left(k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I}\right)^{-1} \mathbf{y} - \frac{1}{2} \log |k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I}| - \frac{N}{2} \log(2\pi). \quad (3.2)$$

Podemos deducir entonces que  $\mathbf{y} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I})$ . Atendiendo a (3.1) esto prueba que  $y$  puede ser considerado otro GP con media 0 y covarianza dada, para cada par  $\mathbf{x}_1, \mathbf{x}_2$ , por  $k(\mathbf{x}_1, \mathbf{x}_2) + \sigma_r^2 \delta(\mathbf{x}_1, \mathbf{x}_2)$ , donde  $\delta(\mathbf{x}_1, \mathbf{x}_2) = 1$  si  $\mathbf{x}_1 = \mathbf{x}_2$  y  $\delta(\mathbf{x}_1, \mathbf{x}_2) = 0$  en otro caso.

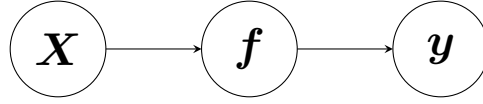


Figura 3.1.: Diagrama del modelo probabilístico para un problema de regresión.

El modelo que hemos construido se ilustra en la **Figura 3.1** y queda determinado como sigue,

$$p(\mathbf{f} | \mathbf{X}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, k(\mathbf{X}, \mathbf{X})),$$

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma_r^2 \mathbf{I}).$$

El objetivo en un problema de regresión es hacer predicciones para nuevos ejemplos basándonos en el conocimiento obtenido del conjunto de entrenamiento. Denotemos por  $\mathbf{x}_*$  a un nuevo ejemplo (ejemplo de *test*) y por  $y_* = y(\mathbf{x}_*)$  al valor de la función objetivo que deseamos predecir. Para poder hacer esto necesitamos calcular la distribución predictiva  $p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ . De la **Observación 3.2** deducimos

$$p\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \mid \mathbf{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix}\right)$$

$$p\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix} \mid \mathbf{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I} & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) + \sigma_r^2 \end{bmatrix}\right).$$

A partir de una distribución gaussiana podemos calcular fácilmente las distribuciones condicionadas a los valores de una o varias variables. Esto da lugar a los siguientes resultados:

**Proposición 3.1.** *La distribución predictiva de un GP sin considerar ruido, cumple*

$$p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \mathcal{N}(f_* | \mu_*, \sigma_*),$$

donde

$$\mu_* = k(\mathbf{x}_*, \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1} \mathbf{y},$$

$$\sigma_* = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}k(\mathbf{X}, \mathbf{x}_*).$$

**Proposición 3.2.** *La distribución predictiva de un GP considerando ruido, cumple*

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \mathcal{N}(y_* | \mu_*, \sigma_*),$$

donde

$$\mu_* = k(\mathbf{x}_*, \mathbf{X}) \left[ k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I} \right]^{-1} \mathbf{y},$$

$$\sigma_* = k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X}) \left[ k(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I} \right]^{-1} k(\mathbf{X}, \mathbf{x}_*).$$

La distribución predictiva depende del ejemplo  $\mathbf{x}_*$  y el cálculo de la misma conlleva la inversión de una matriz de dimensiones  $N \times N$ , lo cual tiene una complejidad  $O(N^3)$ . Esto es

### 3. Procesos Gaussianos

una gran limitación para los GPs ya que, en su formulación original, no pueden trabajar con conjuntos de datos grandes. Sabemos que cualquier modelo predictivo aumenta su capacidad de generalizar a nuevos ejemplos conforme mayor es el conjunto de entrenamiento. Para solucionar esto, se utilizan técnicas consistentes en evaluar la covarianza en un conjunto de  $M$  puntos (con  $M$  fijo) que resumen el conjunto de entrenamiento.

En la **Figura 3.2** se muestra un GP unidimensional. En la subfigura (a) se muestra la verdadera función que queremos aproximar. En la subfigura (b) se muestra la distribución a priori, antes de incorporar el conocimiento sobre los ejemplos observados. Una vez que se observa el valor de la función para algunos valores, la distribución a posteriori toma la forma de la subfigura (c). En la zona donde se han observado los ejemplos la varianza es menor, mientras que donde no hay ejemplos observados hay una gran incertidumbre. Si consideramos que hay cierto ruido a la hora de generar los ejemplos, se obtiene la distribución predictiva de la subfigura (d).

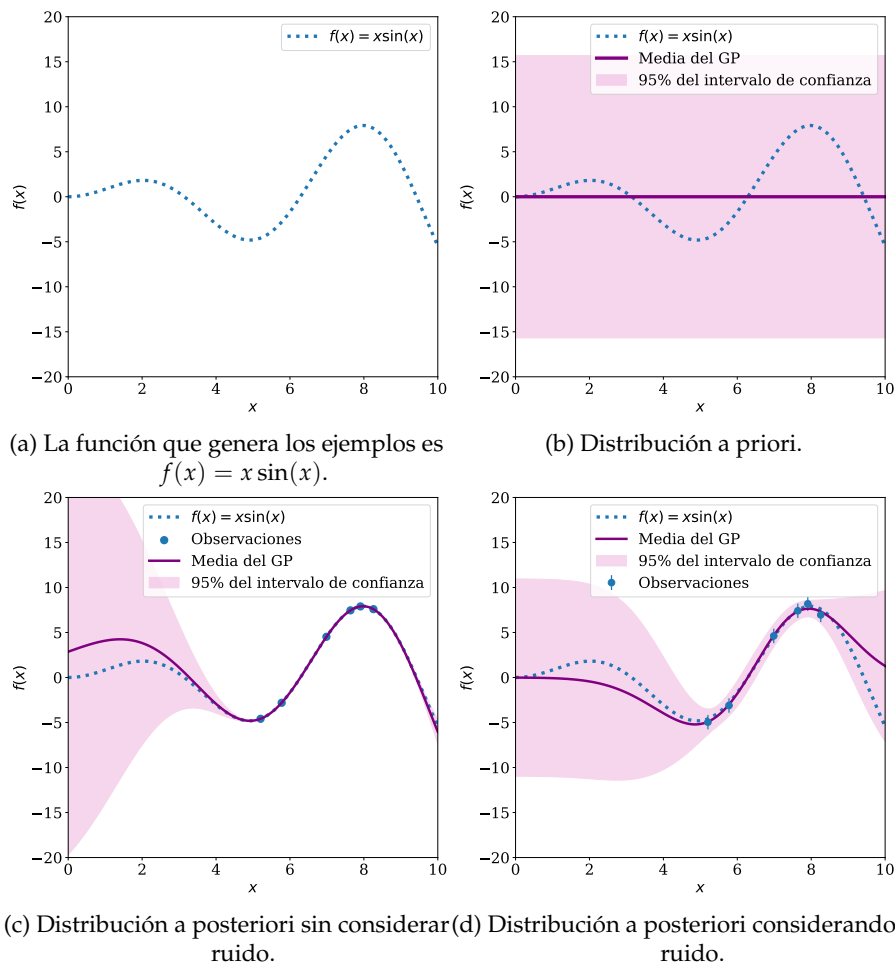


Figura 3.2.: Visualización de un GP para regresión unidimensional.

### 3.3. Procesos gaussianos para clasificación

En un problema de clasificación supervisado, las etiquetas de los ejemplos no son valores reales, sino que pertenecen a algún conjunto discreto. La modelización probabilística de este problema pretende modelar la probabilidad a posteriori de que un nuevo ejemplo pertenezca a una clase o a otra, una vez observado el conjunto de entrenamiento. Para poder aplicar GPs en este escenario debemos transformar la salida del GP mediante alguna función de activación no lineal que nos permita realizar esa interpretación probabilística.

Fijemos en primer lugar la notación. Nuevamente, sea  $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \{1, \dots, N\}\} \subset \mathbb{R}^D \times \mathcal{Y}$  el conjunto de entrenamiento, donde  $\mathcal{Y}$  es el espacio de las etiquetas. Consideremos el caso de un problema binario,  $\mathcal{Y} = \{0, 1\}$ . Modelamos la relación entre  $x$  e  $y$  suponiendo

$$p(y | f) = \sigma(f)^y (1 - \sigma(f))^{1-y}$$

donde  $\sigma: \mathbb{R} \rightarrow [0, 1]$  es una función sigmoideal y  $f \sim \mathcal{GP}(0, k)$ . La función  $f$  juega un papel peculiar ya que no la observamos directamente (sólo a través de las etiquetas) y tampoco estamos interesados en los valores que toma. Simplemente nos sirve para modelar el problema usando un GP. El modelo de GP para clasificación se puede representar gráficamente como en la [Figura 3.1](#). Es común suponer que los ejemplos de entrenamiento son independientes, por lo que la distribución conjunta queda completamente determinada por

$$p(\mathbf{f} | \mathbf{X}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, k(\mathbf{X}, \mathbf{X})),$$

$$p(\mathbf{y} | \mathbf{f}) = \prod_{n=1}^N \text{Bernouilli}(y_n | \sigma(f_n)).$$

Sin embargo, esta última suposición no tiene por qué darse en problemas MIL. Obsérvese que estamos suponiendo que las etiquetas de las instancias son independientes, lo cual no tiene por qué ocurrir en todos los problemas MIL. Un ejemplo de ello es el problema del cáncer de piel [36]. El modelo que hemos construido se ilustra en la [Figura 3.1](#).

Para clasificar un ejemplo  $\mathbf{x}_*$  necesitamos calcular la distribución predictiva  $p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ . Llamando  $f_* = f(\mathbf{x}_*)$ , se tiene

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int p(y_* | f_*) p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) df_*.$$

En un problema binario es suficiente calcular  $p(y_* = 1 | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$  ya que  $p(y_* = 0 | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = 1 - p(y_* = 1 | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ , luego basta concentrarse en

$$p(y_* = 1 | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int \sigma(f_*) p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) df_*.$$

Tanto esta expresión como la que define a  $p(f_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$  son intratables de forma analítica. En el caso de la regresión sí que podíamos obtener expresiones cerradas para la distribución predictiva ya que los términos involucrados se correspondían a distribuciones gaussianas. Para salvar esto, debemos aproximar dichas integrales. Algunas aproximaciones comunes son la *aproximación de Laplace* [52], la conocida como *propagación de la esperanza* (*Expectance Propagation, EP*) [27], o aquellas basadas en *inferencia variacional* [17], aunque se han propuesto muchas otras [37, 21]. Todos estos métodos tienen complejidad  $O(N^3)$ , siendo  $N$  el número de ejemplos en entrenamiento. Para poder utilizar conjuntos grandes, al igual que en el caso

### 3. Procesos Gaussianos

de los GPs para regresión, debemos integrar estas técnicas con otras que permitan reducir el coste computacional.

## 3.4. Ajuste de hiperparámetros

La distribución predictiva de un GP es la que nos permite predecir la variable objetivo para un ejemplo nuevo. Sin embargo, el comportamiento de esta distribución y la capacidad de un GP viene determinado por el núcleo que elijamos. Este núcleo, la mayoría de las ocasiones, dependerá de la elección de varios hiperparámetros, los cuales debemos ajustar cuidadosamente.

Para tal fin, recurrimos a la verosimilitud marginal, cuya expresión habíamos calculado para el problema de regresión en la ecuación (3.2) y volvemos a escribir a continuación haciendo explícita la dependencia de los hiperparámetros  $\theta$ ,

$$\log p(\mathbf{y} | \mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T \left( k_\theta(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I} \right)^{-1} \mathbf{y} - \frac{1}{2} \log |k_\theta(\mathbf{X}, \mathbf{X}) + \sigma_r^2 \mathbf{I}| - \frac{N}{2} \log(2\pi).$$

Si la expresión que define al núcleo,  $k_\theta$  es diferenciable, podremos aplicar métodos de gradiente para optimizar la verosimilitud marginal. Esto se traducirá en encontrar aquellos hiperparámetros que maximizan la verosimilitud de la muestra.

El caso de la clasificación es más complicado ya que no disponemos de una expresión cerrada para la verosimilitud marginal. En la práctica, esta cantidad se aproxima de alguna forma y se efectúa la optimización sobre la aproximación obtenida.

## 3.5. Sparse GP

La solución más popular para salvar el coste computacional de los GPs recibe el nombre de *Procesos Gaussianos Dispersos* (*Sparse GPs*, *SGPs*). Consiste en asumir la existencia de un conjunto de *inducing points*  $\{z_1, \dots, z_M\} \subset \mathbb{R}^D$  procedentes del mismo espacio que los ejemplos de entrenamiento. El número de puntos considerados,  $M$ , se suele fijar mucho más pequeño que  $N$ .

La idea consiste en considerar este conjunto,  $\mathbf{Z} = [z_1^\top, \dots, z_M^\top] \in \mathbb{R}^{D \times M}$  y los valores que producen a través del GP  $f$ ,  $\mathbf{u} = f(\mathbf{Z}) = [f(z_1), \dots, f(z_M)]^\top$ . A través de estas variables pretendemos sustituir la matriz  $k(\mathbf{X}, \mathbf{X})$  por otra de rango  $M$ , cuya inversión sea computacionalmente menos costosa.

Las relaciones que se supongan entre  $\mathbf{X}$ ,  $\mathbf{Z}$ ,  $\mathbf{u}$  y  $f$  dan lugar a diferentes aproximaciones. Siguiendo [9], suponemos que  $\mathbf{Z}$  y  $\mathbf{u}$  se relacionan de la misma forma que lo hacen  $\mathbf{X}$  y  $f$ . Esto nos permite escribir

$$p \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \right) = \mathcal{N} \left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}, \mathbf{0}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{Z}) \\ k(\mathbf{Z}, \mathbf{X}) & k(\mathbf{Z}, \mathbf{Z}) \end{bmatrix} \right),$$

lo cual da lugar a que las distribuciones involucradas tomen la siguiente forma,

$$\begin{aligned}
 p(\mathbf{u} \mid \mathbf{Z}) &= \mathcal{N}(\mathbf{u} \mid \mathbf{0}, k(\mathbf{Z}, \mathbf{Z})) \\
 p(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u}) &= \mathcal{N}(\mathbf{f} \mid \mathbf{a}, \tilde{\mathbf{K}}), \quad \text{donde} \\
 \mathbf{a} &= k(\mathbf{X}, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} \mathbf{u}, \\
 \tilde{\mathbf{K}} &= k(\mathbf{X}, \mathbf{X}) - k(\mathbf{X}, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} k(\mathbf{Z}, \mathbf{X}).
 \end{aligned} \tag{3.3}$$

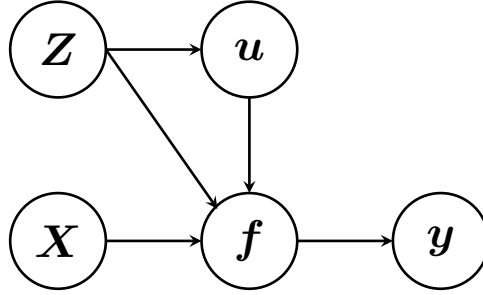


Figura 3.3.: Diagrama del modelo probabilístico SGP.

Este modelo se representa gráficamente en la **Figura 3.3**. Obsérvese que las variables de  $\mathbf{u}$  se introdujeron como variables auxiliares, luego las podemos marginalizar para volver al espacio de probabilidad de partida. Siguiendo [33], nos damos cuenta de que para aproximar  $k(\mathbf{X}, \mathbf{X})$  basta considerar una aproximación  $q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u})$  de la verdadera  $p(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u})$ . Las diferentes aproximaciones a los SGP se engloban dentro de esta interpretación sin más que cambiar las formas funcionales de  $q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u})$ . En general, se considera una distribución a posteriori aproximada de la forma

$$q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u}) = \mathcal{N}(\mathbf{f} \mid \mathbf{a}, \tilde{\mathbf{Q}}),$$

siendo  $\tilde{\mathbf{Q}}$  una matriz aún por elegir. Podemos entonces marginalizar las variables  $\mathbf{u}$  para obtener

$$q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}) = \int q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}, \mathbf{u}) p(\mathbf{u} \mid \mathbf{Z}) d\mathbf{u} = \mathcal{N}\left(\mathbf{f} \mid \mathbf{0}, \tilde{\mathbf{Q}} + \underbrace{k(\mathbf{X}, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} k(\mathbf{Z}, \mathbf{X})}_{\mathbf{Q}_X}\right),$$

donde, en la última igualdad, se ha usado la ecuación (3.3). La matriz  $\mathbf{Q}_X$  se conoce como la *aproximación de Nyström* de  $k(\mathbf{X}, \mathbf{X})$ . Para un problema de regresión podemos obtener la verosimilitud marginal aproximada como

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{Z}) = \int p(\mathbf{y} \mid \mathbf{f}, \mathbf{X}, \mathbf{Z}) q(\mathbf{f} \mid \mathbf{X}, \mathbf{Z}) = \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \tilde{\mathbf{Q}} + \mathbf{Q}_X + \sigma_r^2 \mathbf{I})$$

Aún no hemos especificado la matriz  $\tilde{\mathbf{Q}}$ , la cual tiene las mismas dimensiones que  $k(\mathbf{X}, \mathbf{X})$ . Para algunas elecciones de la misma podemos expresar su inversa de forma que involucre sólo la inversión de una matriz  $M \times M$ . Por ejemplo, el enfoque *deterministic training conditional* (DTC) [38] emplea  $\tilde{\mathbf{Q}} = \mathbf{0}$ . La aproximación *fully independent training conditional* (FITC) [39] usa  $\tilde{\mathbf{Q}} = \mathbf{Q}_X - \text{diag}(\mathbf{Q}_X - k(\mathbf{X}, \mathbf{X}))$ . Nosotros vamos a centrarnos en la propuesta de [42],

la cual introduce la inferencia variacional en el contexto de los GPs.

### 3.6. Aproximación variacional

La metodología que hemos descrito anteriormente se basa, esencialmente, en modificar la distribución a priori. De esta forma, los *inducing points* actúan como hiperparámetros del núcleo que hay que optimizar, lo cual puede llevar al modelo al sobreajuste. En lugar de esto, en [42] se propone aproximar la distribución a posteriori exacta empleando inferencia variacional. Así, los puntos  $z_n$  son parámetros que hay que elegir cuidadosamente.

En lo que sigue vamos a omitir la dependencia respecto de  $\mathbf{X}$  y  $\mathbf{Z}$ . Obsérvese que para calcular la distribución predictiva para un nuevo ejemplo  $\mathbf{x}_*$  podemos escribir

$$\begin{aligned} p(f_* | \mathbf{y}, \mathbf{x}_*) &= \int p(f_*, \mathbf{f}, \mathbf{u} | \mathbf{y}, \mathbf{x}_*) d(\mathbf{f}, \mathbf{u}) = \\ &= \int p(f_* | \mathbf{f}, \mathbf{u}, \mathbf{y}, \mathbf{x}_*) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u} | \mathbf{y}) d(\mathbf{f}, \mathbf{u}). \end{aligned}$$

Si suponemos que  $\mathbf{u}$  es un estadístico suficiente para  $\mathbf{f}$ , esto es,  $p(f_* | \mathbf{f}, \mathbf{u}) = p(f_* | \mathbf{u})$ , entonces

$$\begin{aligned} p(f_* | \mathbf{y}, \mathbf{x}_*) &= \int p(f_* | \mathbf{u}, \mathbf{y}, \mathbf{x}_*) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u} | \mathbf{y}) d(\mathbf{f}, \mathbf{u}) = \\ &= \int p(f_* | \mathbf{u}, \mathbf{y}, \mathbf{x}_*) p(\mathbf{u} | \mathbf{y}) d\mathbf{u}. \end{aligned}$$

En la práctica será difícil encontrar  $\mathbf{u}$  que cumpla esta propiedad, por lo que la anterior será sólo una aproximación a la verdadera distribución predictiva. En tal caso podemos aproximar la intratable  $p(\mathbf{u} | \mathbf{y})$  por  $q(\mathbf{u})$ . Si elegimos  $q(\mathbf{u}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$  y hacemos uso de la definición de GP (como hacíamos en la proposición **Proposición 3.2**, pero ahora considerando también los *inducing points*), entonces ocurrirá que  $p(f_* | \mathbf{y}, \mathbf{x}_*) \approx \mathcal{N}(\bar{f}_*, c_*)$  siendo

$$\begin{aligned} \bar{f}_* &= k(\mathbf{x}_*, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} \mathbf{m} \\ c_* &= k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} k(\mathbf{Z}, \mathbf{x}_*) + \\ &\quad + k(\mathbf{x}_*, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} \mathbf{S} k(\mathbf{Z}, \mathbf{Z})^{-1} k(\mathbf{Z}, \mathbf{x}_*). \end{aligned}$$

El cálculo de esta aproximación involucra la inversión de matrices de orden  $M \times M$ , de forma que la complejidad es  $O(NM^2)$ .

El problema ahora reside en cómo seleccionar los parámetros  $(\mathbf{m}, \mathbf{S})$  y la localización de los *inducing points*. Empleando inferencia variacional, minimizamos la *divergencia de Kullback-Leibler* entre la verdadera distribución  $p(\mathbf{f}, \mathbf{u} | \mathbf{y}) = p(\mathbf{f} | \mathbf{u}) p(\mathbf{u} | \mathbf{y})$  y la aproximación  $Q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} | \mathbf{u}) q(\mathbf{u})$ . Esto es equivalente a maximizar el ELBO, el cual actúa como una cota inferior al logaritmo de la verosimilitud marginal,

$$\begin{aligned} \log(p(\mathbf{y})) &\geq ELBO(Q) = \mathbb{E}_{Q(\mathbf{f}, \mathbf{u})} \left[ \frac{p(\mathbf{f}, \mathbf{u}, \mathbf{y})}{Q(\mathbf{f}, \mathbf{u})} \right] = \\ &= \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) \log \left( \frac{p(\mathbf{y} | \mathbf{f}) p(\mathbf{u})}{q(\mathbf{u})} \right) d(\mathbf{f}, \mathbf{u}). \end{aligned}$$



Como mostramos en la siguiente subsección, el ELBO es una aproximación inferior a la verosimilitud marginal, con la que podemos trabajar para ajustar los hiperparámetros. Cuanto mayor sea el ELBO, mejor aproximará a la verdadera verosimilitud marginal, por lo que su optimización tiene también un significado intuitivo. Tanto el ELBO como la solución óptima al problema de su optimización tienen una expresión cerrada, que se deduce de (3.4) y que puede consultarse en [42].

En esta sección hemos presentado la formulación de los SGPs para un problema de regresión. En un problema de clasificación es necesario, además, emplear algún tipo de técnica con la que poder tratar la función sigmoïdal. En el **Capítulo 4** vamos a estudiar dos de estas técnicas para adaptar los GPs al problema MIL, por lo que no consideramos necesario mostrarlas otra vez aquí. Para más información puede consultarse [4, 9].

### 3.6.1. Inferencia variacional

La inferencia variacional [5] consiste en la aplicación de métodos de optimización variacionales para resolver el problema de la inferencia bayesiana. Esta subsección recoge los resultados principales para poder aplicarla con toda generalidad. Esta técnica es esencial en nuestro trabajo, ya que la solución que proporciona constituye el estado del arte en los SGPs. Pueden consultarse algunos ejemplos en [4].

La divergencia de Kullback-Leibler jugará un papel muy importante. Empezamos definiéndola cuidadosamente.

**Definición 3.4** (Divergencia de Kullback-Leibler). Sean  $P$  y  $Q$  dos medidas de probabilidad en el mismo espacio de las que se conocen funciones de densidad o funciones masa de probabilidad,  $p$  y  $q$  respectivamente. Se define la *divergencia Kullback-Leibler* de  $p$  y  $q$  como

$$\text{KL}(p \parallel q) = \mathbb{E}_{p(x)} [\log p(x) - \log q(x)].$$

La desigualdad de la información hace referencia a uno de los resultados más importantes sobre la divergencia de Kullback-Leibler. Para probarla podemos recurrir a la desigualdad de Jensen [12]. La enunciamos a continuación.

**Teorema 3.1** (Desigualdad de la información). Sean  $P$  y  $Q$  dos medidas de probabilidad en el mismo espacio de las que se conocen funciones de densidad,  $p$  y  $q$  respectivamente. Entonces

1.  $\text{KL}(p \parallel q) \geq 0$
2.  $\text{KL}(p \parallel q) = 0$  si y sólo si  $p = q$  en casi todo punto.

Gracias a este resultado podemos identificar que la divergencia de Kullback-Leibler es una medida del parecido entre dos distribuciones de probabilidad. Veamos cómo podemos usar esta idea en el contexto de la inferencia bayesiana. Supongamos que tenemos un modelo bayesiano definido por unos parámetros a los que se les ha dado una distribución a priori. Supongamos un conjunto de variables observadas  $\mathbf{X}$  y otro conjunto de variables latentes (no observadas)  $\mathbf{Z}$ , dentro de las cuales entran los parámetros. El modelo especifica una probabilidad conjunta  $p(\mathbf{X}, \mathbf{Z})$ , y nuestro objetivo es encontrar una aproximación para la distribución a posteriori  $p(\mathbf{Z} \mid \mathbf{X})$  y para la verosimilitud  $p(\mathbf{X})$ .

Consideremos una distribución de  $q(\mathbf{Z})$ , en principio arbitraria, con la que queremos aproximar a  $p(\mathbf{Z} \mid \mathbf{X})$ . Para cualquier elección de la misma la verosimilitud se puede descomponer

### 3. Procesos Gaussianos

como sigue,

$$p(\mathbf{X}) = \underbrace{\int q(\mathbf{Z}) \log \left( \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right) d\mathbf{Z}}_{\text{ELBO}(q(\mathbf{Z}))} + \underbrace{\int q(\mathbf{Z}) \log \left( \frac{q(\mathbf{Z})}{p(\mathbf{Z} | \mathbf{X})} \right) d\mathbf{Z}}_{\text{KL}(q(\mathbf{Z}) || p(\mathbf{Z} | \mathbf{X}))}.$$

El primer sumando recibe el nombre de *Evidence Lower Bound* (ELBO) y jugará un papel muy importante a la hora de aproximar la verosimilitud. Su nombre se debe a la siguiente desigualdad, consecuencia de la desigualdad de la información,

$$p(\mathbf{X}) = \text{ELBO}(q(\mathbf{Z})) + \text{KL}(q(\mathbf{Z}) || p(\mathbf{Z} | \mathbf{X})) \geq \text{ELBO}(q(\mathbf{Z})).$$

Obsérvese que se da la igualdad si y sólo si las distribuciones  $q(\mathbf{Z})$  y  $p(\mathbf{Z} | \mathbf{X})$  coinciden. Minimizar la divergencia de Kullback-Leibler es equivalente a maximizar el ELBO, lo cual lo convierte en una aproximación a la verosimilitud, cuya calidad dependerá de cómo de bien  $q(\mathbf{Z})$  aproxime a  $p(\mathbf{Z} | \mathbf{X})$ . Así pues, para cumplir nuestro objetivo, debemos minimizar la divergencia de Kullback-Leibler o, equivalentemente, maximizar el ELBO.

Para llevar a cabo esta tarea no podemos trabajar con todos los tipos de distribuciones existentes, sino que debemos de restringirlos para poder encontrar una solución cerrada. Una forma de trabajar es parametrizar la distribución variacional, de forma que el ELBO sea una función de estos parámetros. En ese caso podemos emplear técnicas de optimización no lineal (como por ejemplo gradiente descendente en caso de que la expresión sea diferenciable) para maximizarlo.

Alternativamente, vamos a considerar la aproximación *Mean Field* [30]. Supongamos que dividimos las variables latentes en grupos disjuntos, de forma que  $\mathbf{Z} = \{\mathbf{Z}_i: i \in \{1, \dots, M\}\}$ . La teoría *Mean Field* consiste en suponer que estos grupos son independientes, de forma que la distribución variacional se escribe como

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i).$$

Obsérvese que no se están imponiendo restricciones sobre la forma funcional de las distribuciones  $q_i$ . En su lugar, se está asumiendo que la distribución variacional conjunta factoriza de una forma concreta.

Para maximizar el ELBO podemos pensar en optimizar su expresión respecto de cada distribución  $q_i$  mientras que se dejan el resto fijas. Escribamos la expresión del ELBO haciendo

uso de esta factorización,

$$\begin{aligned}
 \text{ELBO}(q(\mathbf{Z})) &= \int \prod_{i=1}^M q_i(\mathbf{Z}_i) \left( \log p(\mathbf{X}, \mathbf{Z}) - \sum_{i=1}^M \log q_i(\mathbf{Z}_i) \right) d\mathbf{Z} = \\
 &= \int q_j(\mathbf{Z}_j) \left( \underbrace{\int \log p(\mathbf{X}, \mathbf{Z}) \prod_{i=1, i \neq j}^M q_i(\mathbf{Z}_i) d\mathbf{Z}_i}_{\log \tilde{p}(\mathbf{X}, \mathbf{Z}_j)} \right) d\mathbf{Z}_j - \int q_j(\mathbf{Z}_j) \log q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{const} = \\
 &= \int q_j(\mathbf{Z}_j) \log \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j(\mathbf{Z}_j) \log q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{const} = \\
 &= -\text{KL}(q_j(\mathbf{Z}_j) \parallel \tilde{p}(\mathbf{X}, \mathbf{Z}_j)) + \text{const}.
 \end{aligned}$$

donde const hace referencia a una constante que no depende de  $\mathbf{Z}_j$  y hemos definido una nueva distribución mediante la relación

$$\log \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{q(\mathbf{Z})/q_j(\mathbf{Z}_j)} [\log p(\mathbf{X}, \mathbf{Z})] + \text{const}.$$

Si queremos maximizar el ELBO respecto de  $q_j(\mathbf{Z}_j)$  dejando el resto de las distribuciones fijas, a la vista de la igualdad obtenida basta minimizar la divergencia de Kullback-Leibler entre  $q_j(\mathbf{Z}_j)$  y  $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ , lo cual ocurre cuando ambas coinciden. Por tanto, hemos obtenido una expresión para la solución óptima,

$$\log q_j^*(\mathbf{Z}_j) = \mathbb{E}_{q(\mathbf{Z})/q_j(\mathbf{Z}_j)} [\log p(\mathbf{X}, \mathbf{Z})] + \text{const}. \quad (3.4)$$

Otra forma de expresar este resultado, sin involucrar a la constante aditiva, es

$$q_j^*(\mathbf{Z}_j) = \frac{\exp\left(\mathbb{E}_{q(\mathbf{Z})/q_j(\mathbf{Z}_j)} [\log p(\mathbf{X}, \mathbf{Z})]\right)}{\int \exp\left(\mathbb{E}_{q(\mathbf{Z})/q_j(\mathbf{Z}_j)} [\log p(\mathbf{X}, \mathbf{Z})] d\mathbf{Z}_j\right)}.$$

Las ecuaciones dadas por (3.4) no representan una solución explícita para la distribución variacional porque cada una depende del resto. Por ello, para obtener una solución en la práctica, debemos aplicarlas iterativamente hasta llegar a la solución óptima. Puede probarse que el ELBO es convexo respecto de cada uno de los factores, por lo que la convergencia a la solución óptima está garantizada [6].



## 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

En este punto del trabajo ya hemos estudiado los fundamentos de MIL y de los GPs. Ahora, vamos a usar los GPs como herramienta para resolver el problema MIL de clasificación binaria bajo la hipótesis normal. Recuérdese que la elección de los GPs estaba justificada por su formulación probabilística, que nos permite tener en cuenta la incertidumbre a la hora de modelar el problema.

Para poder usar los GPs en este problema, primero estudiamos el estado del arte. Es decir, analizamos cuidadosamente qué enfoques son los que presentan mejores resultados en clasificación con GPs y en MIL. Las ideas obtenidas de este estudio dan lugar a dos grandes contribuciones, que reciben el nombre de PG-VGPMIL y G-VGPMIL. La justificación teórica de ambas se recogen en este capítulo, por lo que es uno de los más importantes. Por ello creemos necesario explicar la estructura del mismo.

1. En la primera sección se estudia el modelo *Variational Gaussian Process Multiple Instance Learning* (VGPMIL) [19], que constituye el estado del arte en GPs para problemas MIL bajo el modelo de observación logístico. Para poder calcular las distribuciones variacionales, VGPMIL usa la *desigualdad de Jaakkola* que proporciona una cota cuadrática de la función logística. La desventaja de usar una desigualdad es que podemos perder información, por lo que conviene estudiar alternativas.
2. Motivados por el buen comportamiento de las variables aleatorias Pólya-Gamma en problemas de clasificación [51], nos preguntamos qué ocurre si las introducimos en la formulación de los GPs. Obtenemos así la primera contribución del trabajo, en la que, por primera vez en la literatura, se propone el uso de variables aleatorias Pólya-Gamma para realizar inferencia variacional en un nuevo modelo basado en GPs para MIL. Este modelo recibirá el nombre de *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* (PG-VGPMIL) y probaremos que es equivalente, en el *sentido Mean Field*, a VGPMIL.
3. PG-VGPMIL es un *modelo aumentado*, en el que se introducen ciertas variables que facilitan el cálculo de las distribuciones variacionales. La definición de las variables Pólya-Gamma motiva el uso de variables Gamma. Llegamos así a la segunda contribución, la cual propone por primera vez en la literatura el uso de variables aleatorias Gamma para realizar inferencia variacional en un nuevo modelo basado en GPs para MIL. Este modelo recibirá el nombre de *Gamma Variational Gaussian Process Multiple Instance Learning* (G-VGPMIL).

La primera sección de este capítulo fija la notación y el ambiente en el que nos vamos a mover. Hemos extraído aquellos elementos comunes a todos los modelos y los hemos apartado en esta sección introductoria.

## 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

### 4.1. Notación y ambiente

Antes de entrar en materia debemos fijar la notación y el ambiente en el que nos vamos a mover. Consideramos un problema MIL de clasificación binaria. Sea el conjunto de entrenamiento,

$$\mathcal{D} = \{(\mathbf{x}_n, y_n) : n \in \{1, \dots, N\}\} \subset \mathbb{R}^D \times \{0, 1\},$$

formado por las variables observadas  $\mathbf{x}_n$  y las etiquetas de las instancias que no han sido observadas,  $y_n$ . Este conjunto se recoge en una matriz  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  y en un vector  $\mathbf{y} = [y_1, \dots, y_N]^\top \in \{0, 1\}^N$ . El conjunto de entrenamiento se particiona en  $B$  bolsas  $\{\text{Bag}_1, \dots, \text{Bag}_B\}$  con intersección vacía y se observa una etiqueta  $T_b \in \{0, 1\}$  para cada bolsa. Usamos el operador  $\{\cdot\}_b$  para referirnos a los elementos de la bolsa  $b$ . Por ejemplo,  $\{y\}_b = \{y_i : i \in \text{Bag}_b\}$  y  $\{y\}_{b \setminus n} = \{y\}_b \setminus \{y_n\}$ . Asumimos que se cumple  $T_b = \max \{y\}_b$ . Dado una aplicación núcleo  $k: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  y dos conjuntos de datos  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{Z} \in \mathbb{R}^{M \times D}$ , llamamos  $\mathbf{K}_{\mathbf{X}\mathbf{Z}} \in \mathbb{R}^{N \times M}$  a la matriz dada por  $(\mathbf{K}_{\mathbf{X}\mathbf{Z}})_{ij} = k(\mathbf{x}_i, \mathbf{z}_j)$ . Recordemos que la clasificación binaria mediante GPs se efectúa tomando

$$\begin{aligned} p(\mathbf{f} | \mathbf{X}) &= \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K}_{\mathbf{X}\mathbf{X}}), \\ p(\mathbf{y} | \mathbf{f}) &= \prod_{n=1}^N p(y_n | f_n) = \prod_{n=1}^N \text{Bernouilli}(y_n | \sigma(f_n)), \end{aligned} \quad (4.5)$$

donde  $\sigma: \mathbb{R} \rightarrow ]0, 1[$  es la función logística. Así, el signo de  $f$  determina la clase de la predicción y el valor absoluto determina la confianza que se tiene en ella. Obsérvese, y esto es muy importante, que se está asumiendo independencia entre todas las instancias y, particularmente, independencia entre las instancias de una bolsa. Para poder trabajar con grandes cantidades de datos adoptamos la aproximación *fully independent training conditional (FITC, [39])*, a través de la cual tomamos un conjunto de *inducing points*  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_M]^\top \in \mathbb{R}^{D \times M}$  y su evaluación mediante  $f$ ,  $\mathbf{u} = f(\mathbf{Z}) = [f(\mathbf{z}_1), \dots, f(\mathbf{z}_M)]^\top \in \mathbb{R}^M$ . Suponemos que  $\mathbf{Z}$  y  $\mathbf{u}$  siguen la misma relación que  $\mathbf{X}$  y  $\mathbf{f}$ . Esto es,

$$\begin{aligned} p(\mathbf{u} | \mathbf{Z}) &= \mathcal{N}(\mathbf{u} | \mathbf{0}, k(\mathbf{Z}, \mathbf{Z})) & (4.6) \\ p(\mathbf{f} | \mathbf{X}, \mathbf{Z}, \mathbf{u}) &= \mathcal{N}(\mathbf{f} | \mathbf{a}, \widetilde{\mathbf{K}}), \quad \text{donde} & (4.7) \\ \mathbf{a} &= k(\mathbf{X}, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} \mathbf{u}, \\ \widetilde{\mathbf{K}} &= k(\mathbf{X}, \mathbf{X}) - k(\mathbf{X}, \mathbf{Z}) k(\mathbf{Z}, \mathbf{Z})^{-1} k(\mathbf{Z}, \mathbf{X}). \end{aligned}$$

En un problema MIL no conocemos las etiquetas de las instancias, por lo que debemos extender el modelo probabilístico para tener en cuenta la relación existente entre las instancias y las bolsas. Esto es lo que hacemos a continuación.

### 4.2. VGPMIL: modelo basado en la desigualdad de Jaakkola

El modelo *Variational Gaussian Process Multiple Instance Learning (VGPMIL)* fue presentado en [19] y muestra, por primera vez, que es posible predecir las etiquetas a nivel de instancia de un clasificador basado en GPs usando inferencia variacional. Para ello, aproxima la hipótesis normal de los problemas MIL mediante una versión ruidosa. Esto, junto con la desigualdad de Jaakkola, permite obtener una expresión analítica de las distribuciones variacionales, lo

que asegura que el modelo es escalable y las distribuciones convergen rápidamente al óptimo.

### 4.2.1. Relación entre las instancias y las bolsas

Recordemos que en el escenario MIL las etiquetas  $\mathbf{y}$  no son conocidas. En su lugar, hemos observado las etiquetas a nivel de bolsa,  $\mathbf{T} = [T_1, \dots, T_B]^\top$ . La clave para adaptar la formulación de un GP de clasificación a MIL es cómo modelamos la relación entre estas etiquetas y las instancias. Por ejemplo, en [23] se toma  $p(T_b | \{\mathbf{f}\}_b) = \text{Bernouilli}(T_b | \sigma(\max\{\mathbf{f}\}_b))$ , suponiendo que la instancia con mayor confianza es la responsable de la clase de la bolsa.

En su lugar, para modelizar la relación entre la etiqueta de la bolsa y las etiquetas de las instancias en VGPMIL se propone usar

$$p(T_b | \{\mathbf{y}\}_b) = \left(\frac{H}{H+1}\right)^{G_b} \left(1 - \frac{H}{H+1}\right)^{1-G_b} = \frac{H^{G_b}}{H+1}, \quad (4.8)$$

siendo  $H \in \mathbb{R}^+$  y  $G_b = T_b \max\{\mathbf{y}\}_b + (1 - T_b)(1 - \max\{\mathbf{y}\}_b)$ . Obsérvese que  $G_b$  es igual a uno si se cumple la hipótesis MIL, en cuyo caso se asigna una probabilidad cercana a 1. Si no se cumple,  $G_b$  es cero y se asigna una probabilidad cercana a cero. La constante  $H$  determina cómo de cercanas serán estas probabilidades a uno o a cero, respectivamente. Obsérvese que

$$\lim_{H \rightarrow +\infty} p(T_b | \{\mathbf{y}\}_b) = \begin{cases} 1 & \text{si } T_b = \max\{\mathbf{y}\}_b, \\ 0 & \text{en otro caso,} \end{cases}$$

lo cual convierte a  $p(T_b | \{\mathbf{y}\}_b)$  en una versión ruidosa de la hipótesis MIL. Esto justifica que la constante  $H$  deba fijarse en torno a valores altos. Por último, suponemos independencia entre las bolsas, lo que da lugar a

$$p(\mathbf{T} | \mathbf{y}) = \prod_{b=1}^B \frac{H^{G_b}}{H+1}.$$

El modelo completo se obtiene tomando el método de observación logístico para las etiquetas de las instancias (ecuación (4.5)) junto con la aproximación FITC (ecuaciones (4.6) y (4.7)) y la versión ruidosa de la hipótesis normal que acabamos de presentar (ecuación (4.8)). Se ilustra en la [Figura 4.1](#) y queda definido por las siguientes ecuaciones,

$$\begin{aligned} p(\mathbf{u} | \mathbf{Z}) &= \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{\mathbf{Z}\mathbf{Z}}) \\ p(\mathbf{f} | \mathbf{X}, \mathbf{Z}, \mathbf{u}) &= \mathcal{N}(\mathbf{f} | \mathbf{a}, \widetilde{\mathbf{K}}), \\ p(\mathbf{y} | \mathbf{f}) &= \prod_{n=1}^N \text{Bernouilli}(y_n | \sigma(f_n)), \\ p(\mathbf{T} | \mathbf{y}) &= \prod_{b=1}^B \frac{H^{G_b}}{H+1}. \end{aligned}$$

Para poder realizar predicciones con este modelo, debemos calcular las distribuciones predictivas de las etiquetas de las instancias y de las bolsas. Veamos cómo hacerlo a continuación.

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

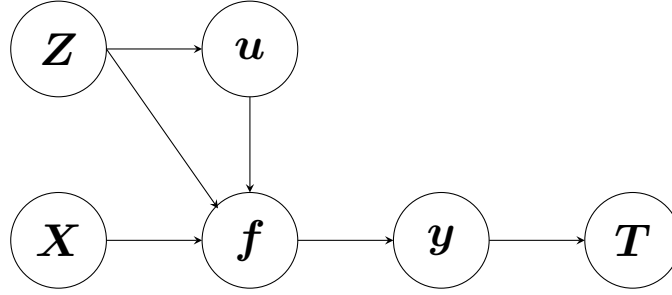


Figura 4.1.: Diagrama del modelo probabilístico VGPMIL.

#### 4.2.2. Distribuciones predictivas

Para obtener la distribución predictiva empleamos un procedimiento análogo al que se utiliza en los SGPs. Dado un nuevo ejemplo  $\mathbf{x}_*$ , se tiene

$$\begin{aligned}
 p(f_* | \mathbf{T}) &= \int p(f_* | \mathbf{f}, \mathbf{u}) p(\mathbf{f}, \mathbf{u} | \mathbf{T}) d\mathbf{f} d\mathbf{u} \approx \\
 &\approx \int p(f_* | \mathbf{f}, \mathbf{u}) p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) d\mathbf{f} d\mathbf{u} = \\
 &= \int p(f_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} = \mathcal{N}(f_* | \mu_*, \sigma_*^2)
 \end{aligned} \tag{4.9}$$

con  $\mu_* = \mathbf{K}_{\mathbf{x}_* \mathbf{Z}} \mathbf{K}_{\mathbf{Z} \mathbf{Z}}^{-1} \mathbf{m}$  y  $\sigma_*^2 = \mathbf{K}_{\mathbf{x}_* \mathbf{x}_*} + \mathbf{K}_{\mathbf{x}_* \mathbf{Z}} \mathbf{K}_{\mathbf{Z} \mathbf{Z}}^{-1} (\mathbf{S} \mathbf{K}_{\mathbf{Z} \mathbf{Z}}^{-1} - \mathbf{I}) \mathbf{K}_{\mathbf{Z} \mathbf{x}_*}$ . La distribución de las etiquetas se calcula mediante

$$p(y_* = 1 | \mathbf{T}) = \int \sigma(f_*) p(f_* | \mathbf{T}) df_*$$

la cual es una integral que no admite una expresión analítica pero que se puede calcular mediante métodos numéricos de forma muy rápida. Para predecir la etiqueta de una bolsa  $\{\mathbf{x}_*^1, \dots, \mathbf{x}_*^{N_*}\}$  con etiquetas desconocidas  $\mathbf{y}_* = \{y_*^1, \dots, y_*^{N_*}\}$ , calculamos la distribución de la etiqueta de la bolsa mediante

$$\begin{aligned}
 p(T_* = 1 | \mathbf{T}) &= \prod_{\mathbf{y}_* \in \{0,1\}^{N_*}} p(T_* = 1 | \mathbf{y}_*) p(\mathbf{y}_* | \mathbf{T}) = \\
 &= \prod_{\mathbf{y}_* \in \{0,1\}^{N_*}} \frac{H^{\max \mathbf{y}_*}}{H+1} \left( \prod_{i=1}^{N_*} p(y_*^i | \mathbf{T}) \right).
 \end{aligned}$$

Llegados a este punto ya hemos establecido el modelo probabilístico subyacente y la forma de realizar predicciones a partir de las distribuciones predictivas. Ahora pasamos a estudiar el proceso de inferencia de VGPMIL.

#### 4.2.3. Inferencia

El proceso de inferencia consiste en calcular la distribución a posteriori de las variables no observadas dadas las variables observadas. Esta distribución es  $p(\mathbf{u}, \mathbf{f}, \mathbf{y} | \mathbf{T}, \mathbf{X}, \mathbf{Z})$  y es



#### 4.2. VGPMIL: modelo basado en la desigualdad de Jaakkola

esencial para poder calcular las distribuciones predictivas. Como esta distribución no admite una expresión analítica, se decide aproximarla empleando inferencia variacional. El objetivo de esta subsección es obtener una expresión analítica para tal aproximación.

En lo que sigue omitiremos la dependencia respecto de  $\mathbf{X}$  y  $\mathbf{Z}$ . Para obtener la aproximación a posteriori  $p(\mathbf{u}, \mathbf{f}, \mathbf{y} \mid \mathbf{T})$  se emplea inferencia variacional y el enfoque Mean Field, considerando así la distribución variacional dada por

$$\begin{aligned} Q(\mathbf{u}, \mathbf{f}, \mathbf{y}) &= q(\mathbf{u})p(\mathbf{f} \mid \mathbf{u})q(\mathbf{y}), \quad \text{donde} \\ q(\mathbf{u}) &= \mathcal{N}(\mathbf{u} \mid \mathbf{m}, \mathbf{S}), \\ q(\mathbf{y}) &= \prod_{n=1}^N q(y_n) = \prod_{n=1}^N \pi_n^{y_n} (1 - \pi_n)^{1-y_n}. \end{aligned}$$

Buscamos minimizar la divergencia de Kullback-Leibler,

$$Q^*(\mathbf{u}, \mathbf{f}, \mathbf{y}) = q^*(\mathbf{u})p(\mathbf{f} \mid \mathbf{u})q^*(\mathbf{y}) = \arg \min_Q \text{KL}(Q(\mathbf{u}, \mathbf{f}, \mathbf{y}) \parallel p(\mathbf{u}, \mathbf{f}, \mathbf{y} \mid \mathbf{T})).$$

Como sabemos, la solución a este problema viene dada, para cada factor  $q^*$  de  $Q^*$ , por (3.4), que en nuestro caso se traduce en

$$\log q^* = \mathbb{E}_{Q^*/q^*} [\log p(\mathbf{T}, \mathbf{u}, \mathbf{f}, \mathbf{y})] + \text{const},$$

donde  $Q^*/q^*$  se refiere al cociente entre la distribución  $Q^*$  y el factor  $q^*$ . Es decir, de  $Q^*$  eliminamos el factor  $q^*$ . La clave de VGPMIL reside en la *desigualdad de Jaakkola*, que nos permite aproximar  $\mathbb{E}_{p(\mathbf{f} \mid \mathbf{u})q(\mathbf{y})} [p(\mathbf{y} \mid \mathbf{f})]$ .

**Lema 4.1** (Desigualdad de Jaakkola [22]). *Para todo  $x \in \mathbb{R}$  y  $\xi \in \mathbb{R}^+$  se cumple*

$$\sigma(x) \geq \sigma(\xi) \exp\left(\frac{x - \xi}{2} - \lambda(\xi)(x^2 - \xi^2)\right),$$

donde  $\lambda(\xi) = (\sigma(\xi) - 1/2) / 2$  y  $\sigma$  denota a la función logística.

Obsérvese que se da la igualdad si tomamos  $x = \xi$ , luego la bondad de esta aproximación vendrá dada por cómo se elija el valor de  $\xi$ . De ella deducimos

$$p(y_n \mid f_n) \geq \exp\left(-\frac{f_n + \xi_n}{2} - \lambda(\xi)(f_n^2 - \xi_n^2)\right) \exp(y_n f_n) \sigma(\xi_n).$$

Esta última desigualdad conduce a una cota inferior del ELBO, que es la que debemos maximizar para elegir los parámetros  $\xi_n$ . Esto es equivalente a minimizar la cota inferior de la siguiente expresión,

$$\mathbb{E} [\log p(y_n \mid f_n)] \geq \mathbb{E} [f_n] (\pi_n - 1/2) - \mathbb{E} [f_n^2] \lambda(\xi_n) - \frac{\xi_n}{2} + \lambda(\xi_n) \xi_n^2 - \log(1 + e^{-\xi_n}),$$

donde la esperanza se toma respecto a la distribución  $p(f_n \mid \mathbf{u})q(\mathbf{u})$ . Es fácil comprobar, y así lo hacemos en la prueba del resultado principal de esta sección, que el máximo se alcanza cuando  $\xi_n = \sqrt{\mathbb{E}[f_n^2]}$ . El resultado clave para poder realizar inferencia es el siguiente.

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

**Resultado 4.1.** Sea  $p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \mathbf{T})$  la distribución conjunta dada por el modelo VGPMIL. Entonces la solución óptima al problema

$$\min_{q(\mathbf{u})q(\mathbf{y})} \text{KL}(q(\mathbf{u})p(\mathbf{f} | \mathbf{u})q(\mathbf{y}) || p(\mathbf{u}, \mathbf{f}, \mathbf{y} | \mathbf{T}))$$

viene dada por las siguientes distribuciones:

$$\begin{aligned} q^*(\mathbf{u}) &\approx \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}), \\ q^*(y_n) &\approx \text{Bernouilli}(y_n | \pi_n) \end{aligned}$$

siendo

$$\begin{aligned} \pi_n &= \sigma \left[ \kappa_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) \right], \\ \mathbf{m} &= \mathbf{S} \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} \mathbf{K}_{\mathbf{Z}\mathbf{X}} \left( \pi - \frac{1}{2} \right), \\ \mathbf{S} &= \left( \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} \mathbf{K}_{\mathbf{Z}\mathbf{X}} \mathbf{\Lambda} \mathbf{K}_{\mathbf{X}\mathbf{Z}} \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} + \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} \right)^{-1}, \\ \xi_n &= \sqrt{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]} = \sqrt{\widetilde{\mathbf{K}}_{nn} + \kappa_n (\mathbf{m} \mathbf{m}^\top + \mathbf{S}) \kappa_n^\top}, \end{aligned}$$

donde  $\kappa_n = \mathbf{K}_{\mathbf{x}_n \mathbf{Z}} \mathbf{K}_{\mathbf{Z}\mathbf{Z}}^{-1} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n]$  y  $\mathbf{\Lambda} = \text{diag}(2\lambda(\xi_1), \dots, 2\lambda(\xi_N))$ .

Antes de proceder a la demostración, debemos enunciar y probar dos lemas que serán muy útiles en lo que sigue.

**Lema 4.2.** Consideremos una instancia  $n$  y sea  $b$  la bolsa donde se encuentra esta instancia. Entonces

$$\max \{ \mathbf{y} \}_b = y_n + \max \{ \mathbf{y} \}_{b \setminus n} - y_n \max \{ \mathbf{y} \}_{b \setminus n}.$$

*Demostración.* Repetimos el proceso de [19],

$$\begin{aligned} \max \{ \mathbf{y} \}_b &= 1 - \prod_{i \in \text{Bag}_b} (1 - y_i) = 1 - (1 - y_n) \prod_{i \in \text{Bag}_b \setminus \{n\}} (1 - y_i) \\ &= 1 - \prod_{i \in \text{Bag}_b \setminus \{n\}} (1 - y_i) + y_n - y_n + y_n \prod_{i \in \text{Bag}_b \setminus \{n\}} (1 - y_i) \\ &= \max \{ \mathbf{y} \}_{b \setminus n} + y_n - y_n \left( 1 - \prod_{i \in \text{Bag}_b \setminus \{n\}} (1 - y_i) \right) \\ &= y_n + \max \{ \mathbf{y} \}_{b \setminus n} - y_n \max \{ \mathbf{y} \}_{b \setminus n} \end{aligned}$$

□

**Lema 4.3.** Sean  $a \in \mathbb{R}^+$  y  $b \in \mathbb{R}$ . Consideramos la función  $f: \mathbb{R}^+ \rightarrow \mathbb{R}$  dada por

$$f(x) = b + \lambda(x) \left( -a + x^2 \right) - \frac{x}{2} - \log(1 + e^{-x}).$$

Entonces  $f$  alcanza un máximo absoluto cuando  $x = \sqrt{a}$ .

#### 4.2. VGPMIL: modelo basado en la desigualdad de Jaakkola

*Demostración.* La función  $f$  es derivable con

$$f'(x) = \lambda'(x) \left( -a + x^2 \right).$$

Como  $\lambda$  es estrictamente decreciente en  $\mathbb{R}^+$ , la función  $f'$  sólo se anula en  $x = \sqrt{a}$  y este es un máximo absoluto.  $\square$

*Demostración.* Apliquemos la desigualdad de Jaakkola,

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{f}) &= \prod_{n=1}^N p(y_n \mid f_n) \geq \\ &\geq \exp \left( \sum_{n=1}^N f_n (y_n - 1/2) - f_n^2 \lambda(\xi_n) \right) \underbrace{\exp \left( \sum_{n=1}^N -\frac{\xi_n}{2} + \lambda(\xi_n) \xi_n^2 \right) \prod_{n=1}^N \sigma(\xi_n)}_{C(\boldsymbol{\xi})} = \\ &= \exp \left[ (\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \boldsymbol{\Lambda} \mathbf{f} \right] C(\boldsymbol{\xi}). \end{aligned}$$

**Actualización de  $q(y_n)$ .** Denotamos  $q(y_{j \neq n}) = \prod_{j \neq n} \pi_j^{y_j} (1 - \pi_j)^{1-y_j}$ . Fijado  $n$ , sea  $b$  el índice de la bolsa en la que se encuentra la instancia  $n$ . Se tiene:

$$\begin{aligned} \log q^*(y_n) &= \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})} [\log p(\mathbf{y} \mid \mathbf{f})]}_{A_1} + \\ &\quad + \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})} [\log p(\mathbf{T} \mid \mathbf{y})]}_{A_2} + \text{const.} \end{aligned}$$

Analizamos cada sumando por separado,

$$\begin{aligned} A_1 &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})} \left[ (\mathbf{y} - 1/2)^\top \mathbf{f} \right] + \text{const} = \\ &= y_n \underbrace{\mathbf{K}_{x_i \mathbf{Z}} \mathbf{K}_{\mathbf{Z} \mathbf{Z}}^{-1} \mathbf{m}}_{\boldsymbol{\kappa}_n} + \text{const} \\ A_2 &= \log H \mathbb{E}_{q(y_{j \neq n})} [G_b] + \text{const} = \\ &= y_n \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) + \text{const.} \end{aligned}$$

Juntando todo nos queda,

$$\log q^*(y_n) \geq y_n \left[ \boldsymbol{\kappa}_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) \right] + \text{const},$$

de donde deducimos que  $q(y_n)$  debe aproximarse a una Bernouilli ( $y_n \mid \pi_n$ ) con

$$\pi_n = \sigma \left\{ \boldsymbol{\kappa}_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) \right\}$$

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

**Actualización de  $q(\mathbf{u})$ .**

$$\log q^*(\mathbf{u}) = \underbrace{\log p(\mathbf{u})}_{B_1} + \underbrace{E_{p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{f} | \mathbf{u})]}_{\text{const}} + \underbrace{E_{p(\mathbf{f}|\mathbf{u})q(\mathbf{y})} [\log p(\mathbf{y}|\mathbf{f})]}_{B_2} + \text{const}$$

Analizamos cada sumando por separado,

$$\begin{aligned} B_1 &= -\frac{1}{2}\mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{u} + \text{const}, \\ B_2 &\geq E_{q(\mathbf{y})p(\mathbf{f}|\mathbf{u})} \left[ (\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \Lambda \mathbf{f} \right] + \text{const} = \\ &= (\pi - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \left( \text{tr}(\Lambda \tilde{\mathbf{K}}) + \mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \Lambda \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} \right) + \text{const}. \end{aligned}$$

Juntando todo:

$$\log q^*(\mathbf{u}) \geq (\pi - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \mathbf{u}^\top \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \Lambda \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right) \mathbf{u} + \text{const},$$

de donde  $q(\mathbf{u})$  debe aproximarse a una  $\mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S})$  con

$$\begin{aligned} \mathbf{m} &= \mathbf{S} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} (\pi - 1/2) \\ \mathbf{S} &= \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \Lambda \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right)^{-1} \end{aligned}$$

**Actualización de  $\xi_n$ .** Para obtener la expresión de  $\xi_n$  debemos de maximizar la cota inferior del ELBO que produce la desigualdad de Jaakkola. Esto es equivalente a maximizar la cota inferior de la siguiente expresión,

$$\mathbb{E} [\log p(y_n | f_n)] \geq \mathbb{E} [f_n] (\pi_n - 1/2) - \mathbb{E} [f_n^2] \lambda(\xi_n) - \frac{\xi_n}{2} + \lambda(\xi_n) \xi_n^2 - \log(1 + e^{-\xi_n}),$$

en la que la esperanza se toma respecto a la distribución  $p(f_n | \mathbf{u})q(\mathbf{u})$ . Aplicando el **Lema 4.3** se tiene que el máximo se alcanza cuando  $\xi = \sqrt{\mathbb{E} [f_n^2]}$ .  $\square$

Obsérvese que no hemos escrito una igualdad al dar la forma de las distribuciones óptimas. Esto se debe al uso que se hace de la desigualdad de Jaakkola. Es natural pensar que al usar una desigualdad estamos perdiendo información, lo cual puede conducir a que las distribuciones variacionales no sean óptimas. La pregunta es, ¿cómo podemos realizar inferencia sin emplear una desigualdad para tratar la función logística? Este trabajo encontró la respuesta a tal pregunta en las variables Pólya-Gamma, que permiten expresar la función  $\sigma$  de forma cerrada. Así es como surge la primera contribución de este trabajo, que presentamos a continuación.

### 4.3. PG-VGPMIL: modelo basado en variables Pólya-Gamma

Como hemos comentado, el principal inconveniente de la desigualdad de Jaakkola es que puede conducir a distribuciones que no son óptimas. En [51] se usan variables aleatorias Pólya-Gamma para derivar un modelo de GPs para clasificación binaria. La técnica usada

### 4.3. PG-VGPMIL: modelo basado en variables Pólya-Gamma

recibe el nombre de *data augmentation* y consiste en modificar la distribución conjunta incluyendo nuevas variables. Al marginalizarlas, recuperamos el modelo original. En el caso de las variables Pólya-Gamma, estas facilitan el tratamiento de la función logística a la hora de llevar a cabo inferencia variacional ya que permiten expresarla como una integral sobre ellas.

La pregunta que surge es: ¿qué ocurre si empleamos este procedimiento en la formulación de un GP para MIL? La respuesta a esta pregunta es el modelo *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* PG-VGPMIL. Por primera vez en la literatura, empleamos variables aleatorias Pólya-Gamma para aumentar la probabilidad conjunta y construir un nuevo modelo. Estas variables permiten obtener las distribuciones variacionales sin emplear ninguna desigualdad.

**Definición 4.1** (Pólya-Gamma, [32]). Decimos que una variable aleatoria  $\omega$  sigue una distribución Pólya-Gamma con parámetros  $b \in \mathbb{R}^+$  y  $c \in \mathbb{R}$ , y lo denotamos por  $\omega \sim \text{PG}(b, c)$ , si

$$\omega \stackrel{D}{=} \frac{1}{2\pi^2} \sum_{k=1}^{\infty} \frac{g_k}{(k - \frac{1}{2})^2 + \frac{c^2}{4\pi^2}},$$

donde  $g_k \sim \text{Gamma}(b, 1)$  son variables aleatorias independientes que siguen una distribución Gamma, y  $\stackrel{D}{=}$  indica que las distribuciones de ambas variables aleatorias son iguales.

Obsérvese que hemos definido las nuevas variables exactamente como se hizo originalmente. Para este trabajo estamos interesados en las propiedades que se enuncian en el siguiente resultado. La prueba de las mismas puede consultarse en [32].

**Teorema 4.1** (Propiedades de la distribución Pólya-Gamma, [32]).

1. Sea  $\omega \sim \text{PG}(b, 0)$ . Entonces

$$\mathbb{E}_{\text{PG}(\omega|b,0)}[\exp(-\omega t)] = \frac{1}{\cosh^b(\sqrt{t/2})}.$$

2. Se cumple la siguiente igualdad,

$$\text{PG}(\omega | b, c) = \cosh^b\left(\frac{c}{2}\right) \exp\left(-\frac{c^2}{2}\omega\right) \text{PG}(\omega | b, 0)$$

3. Si denotamos a la función logística por  $\sigma$ , para cada  $x \in \mathbb{R}$  se cumple

$$\sigma(x) = \frac{1}{2} \mathbb{E}_{\text{PG}(\omega|1,0)} \left[ \exp\left(\frac{x}{2} - \frac{x^2}{2}\omega\right) \right]. \quad (4.10)$$

La igualdad (4.10) nos permitirá incorporar las variables aleatorias Pólya-Gamma en nuestro modelo de observación. Para ello, basta escribir  $p(y_n | f_n) = \sigma((2y_n - 1)f_n)$  y aplicar (4.10) considerando  $\omega = [\omega_1, \dots, \omega_N]$  con  $\omega_n \sim \text{PG}(1, 0)$ ,

$$p(\mathbf{y} | \mathbf{f}) = \prod_n \frac{1}{2} \mathbb{E}_{p(\omega_n)} \left[ \exp\left(\frac{(2y_n - 1)f_n}{2} - \frac{f_n^2}{2}\omega_n\right) \right].$$

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

Condicionando a  $\omega$  se obtiene

$$p(\mathbf{y} | \mathbf{f}, \omega) \propto \exp \left( \sum_n (y_n - 1/2) f_n - \frac{f_n^2}{2} \omega_n \right) = \exp \left( (\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \Omega \mathbf{f} \right),$$

siendo  $\Omega = \text{diag}(\omega_1, \dots, \omega_N)$ . Llegamos así al nuevo modelo PG-VGPMIL,

$$\begin{aligned} p(\mathbf{u} | \mathbf{Z}) &= \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{\mathbf{Z}\mathbf{Z}}) \\ p(\mathbf{f} | \mathbf{X}, \mathbf{Z}, \mathbf{u}) &= \mathcal{N}(\mathbf{f} | \mathbf{a}, \widetilde{\mathbf{K}}), \\ p(\mathbf{y} | \mathbf{f}, \omega) &\propto \exp \left( (\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \Omega \mathbf{f} \right), \\ p(\omega) &= \prod_{n=1}^N p(\omega_n), \quad \omega_n \sim \text{PG}(1, 0) \\ p(\mathbf{T} | \mathbf{y}) &= \prod_{b=1}^B \frac{H^{G_b}}{H + 1}. \end{aligned}$$

Una representación gráfica del mismo se encuentra en la [Figura 4.2](#). Obsérvese que el modelo de observación no ha cambiado ya que

$$p(y_n | f_n) = \int p(y_n | f_n, \omega_n) p(\omega_n) d\omega_n = \sigma((2y_n - 1) f_n).$$

Además, al marginalizar en la distribución conjunta sobre  $\omega$  recuperamos el modelo VGPMIL.

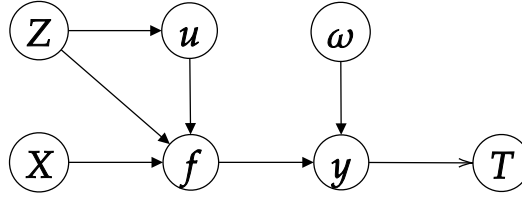


Figura 4.2.: Diagrama del modelo probabilístico PG-VGPMIL.

##### 4.3.1. Inferencia

A diferencia de VGPMIL, ahora la distribución a posteriori debe tener en cuenta las nuevas variables  $\omega$ . En lo que sigue omitiremos la dependencia respecto de  $\mathbf{X}$  y  $\mathbf{Z}$ . Para obtener la aproximación a posteriori  $p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \omega | \mathbf{T})$  vamos a emplear inferencia variacional. Empleamos el enfoque *Mean Field* considerando la distribución variacional

$$\begin{aligned} Q(\mathbf{u}, \mathbf{f}, \mathbf{y}, \omega) &= q(\mathbf{u}) p(\mathbf{f} | \mathbf{u}) q(\mathbf{y}) q(\omega), \quad \text{donde} \\ q(\mathbf{y}) &= \prod_{n=1}^N q(y_n), \\ q(\omega) &= \prod_{n=1}^N q(\omega_n). \end{aligned}$$

### 4.3. PG-VGPML: modelo basado en variables Pólya-Gamma

Buscamos minimizar la divergencia de Kullback-Leibler,

$$\begin{aligned} Q^*(\mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega}) &= q^*(\mathbf{u})p(\mathbf{f} | \mathbf{u})q^*(\mathbf{y})q^*(\boldsymbol{\omega}) = \\ &= \arg \min_Q \text{KL}(Q(\mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega}) || p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega} | \mathbf{T})). \end{aligned}$$

Como sabemos, la solución a este problema viene dada, para cada factor  $q^*$  de  $Q^*$ , por (3.4), que en nuestro caso se traduce en

$$\log q^* = \mathbb{E}_{Q^*/q^*} [\log p(\mathbf{T}, \mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega})] + \text{const}, \quad (4.11)$$

donde  $Q^*/q^*$  se refiere al cociente entre la distribución  $Q^*$  y el factor  $q^*$ . Es decir, de  $Q^*$  eliminamos el factor  $q^*$ . En la práctica, esto se lleva a cabo iterativamente, actualizando los parámetros de un factor manteniendo el resto fijos. El siguiente resultado recoge la forma de las distribuciones y es consecuencia de aplicar (4.11).

**Resultado 4.2.** Sea  $p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega}, \mathbf{T})$  la distribución conjunta dada por el modelo PG-VGPML. Entonces la solución óptima al problema

$$\min_{q(\mathbf{u})q(\mathbf{y})q(\boldsymbol{\omega})} \text{KL}(q(\mathbf{u})p(\mathbf{f} | \mathbf{u})q(\mathbf{y})q(\boldsymbol{\omega}) || p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega} | \mathbf{T}))$$

viene dada por las siguientes distribuciones:

$$\begin{aligned} q^*(\mathbf{u}) &= \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}), \\ q^*(y_n) &= \text{Bernouilli}(y_n | \pi_n), \\ q^*(\omega_n) &= \text{PG}(\omega_n | 1, c_n), \end{aligned}$$

siendo

$$\begin{aligned} \pi_n &= \sigma \left[ \kappa_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) \right], \\ \mathbf{m} &= \mathbf{S} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \left( \pi - \frac{1}{2} \right), \\ \mathbf{S} &= \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right)^{-1}, \\ c_n &= \sqrt{\widetilde{\mathbf{K}}_{nn} + \kappa_n (\mathbf{m} \mathbf{m}^\top + \mathbf{S}) \kappa_n^\top}, \end{aligned}$$

donde  $\kappa_n = \mathbf{K}_{\mathbf{x}_n \mathbf{Z}} \mathbf{K}_{ZZ}^{-1} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n]$ ,  $\boldsymbol{\Theta} = \text{diag}(\theta(c_1), \dots, \theta(c_N))$  y  $\theta(c) = \tanh(c/2)/(2c)$ .

*Demostración.* Calculamos la solución óptima para cada una de las distribuciones fijando el resto,

$$\log q^* = \mathbb{E}_{Q^*/q^*} [\log p(\mathbf{T}, \mathbf{u}, \mathbf{f}, \mathbf{y}, \boldsymbol{\omega})] + \text{const}.$$

**Actualización de  $q(y_n)$ .** Denotamos  $q(y_{j \neq n}) = \prod_{j \neq n} \pi_j^{y_j} (1 - \pi_j)^{1 - y_j}$ . Fijado  $n$ , sea  $b$  el índice de

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

la bolsa en la que se encuentra la instancia  $n$ . Se tiene:

$$\log q^*(y_n) = \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})q(\omega)} [\log p(\mathbf{y} | \mathbf{f}, \omega)]}_{A_1} + \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})q(\omega)} [\log p(\mathbf{T} | \mathbf{y})]}_{A_2} + \text{const.}$$

Analizamos cada sumando por separado,

$$\begin{aligned} A_1 &= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(y_{j \neq n})} \left[ (\mathbf{y} - 1/2)^\top \mathbf{f} \right] + \text{const} = \\ &= \mathbb{E}_{q(y_{j \neq n})} [\mathbf{y}]^\top \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\mathbf{f}] + \text{const} \\ &= y_n \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n] + \text{const} = \\ &= y_n \underbrace{\mathbf{K}_{x_i Z} \mathbf{K}_{ZZ}^{-1} \mathbf{m}}_{\kappa_n} + \text{const} \\ A_2 &= \log H \mathbb{E}_{q(y_{j \neq n})} [G_b] + \text{const} = \\ &= \log H(2T_b - 1) \mathbb{E}_{q(y_{j \neq n})} [\max \{\mathbf{y}\}_b] + \text{const} = \\ &= y_n \log H(2T_b - 1) \left( 1 - \mathbb{E} [\max \{\mathbf{y}\}_{b \setminus n}] \right) + \text{const}, \end{aligned}$$

donde hemos empleado el [Lema 4.2](#) para calcular la esperanza involucrada y hemos aproximado  $\max \mathbb{E} [x] \approx \mathbb{E} [\max x]$ . Juntando todo nos queda,

$$\log q^*(y_n) = y_n \left[ \kappa_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} [\max \{\mathbf{y}\}_{b \setminus n}] \right) \right] + \text{const},$$

de donde deducimos que  $q(y_n) = \text{Bernouilli}(y_n | \pi_n)$  con

$$\pi_n = \sigma \left\{ \kappa_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} [\max \{\mathbf{y}\}_{b \setminus n}] \right) \right\}$$

**Actualización de  $q(\mathbf{u})$ .**

$$\log q^*(\mathbf{u}) = \underbrace{\log p(\mathbf{u})}_{B_1} + \underbrace{E_{p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{f} | \mathbf{u})]}_{\text{const}} + \underbrace{\mathbb{E}_{p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega)} [\log p(\mathbf{y} | \mathbf{f}, \omega)]}_{B_2} + \text{const}$$

Analizamos cada sumando por separado,

$$B_1 = -\frac{1}{2} \mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{u} + \text{const}$$



### 4.3. PG-VGPMIL: modelo basado en variables Pólya-Gamma

$$\begin{aligned}
B_2 &= \mathbb{E}_{q(\mathbf{y})p(\mathbf{f}|\mathbf{u})q(\boldsymbol{\omega})} \left[ (\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \boldsymbol{\Omega} \mathbf{f} \right] + \text{const} = \\
&= \mathbb{E}_{p(\mathbf{f}|\mathbf{u})q(\boldsymbol{\omega})} \left[ (\boldsymbol{\pi} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \boldsymbol{\Omega} \mathbf{f} \right] + \text{const} = \\
&= \mathbb{E}_{q(\boldsymbol{\omega})} \left[ (\boldsymbol{\pi} - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \left( \text{tr}(\boldsymbol{\Omega} \widetilde{\mathbf{K}}) + \mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Omega} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} \right) \right] + \text{const} = \\
&= (\boldsymbol{\pi} - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \left( \text{tr}(\boldsymbol{\Theta} \widetilde{\mathbf{K}}) + \mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} \right) + \text{const} \\
&= (\boldsymbol{\pi} - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \left( \text{tr}(\boldsymbol{\Theta} \widetilde{\mathbf{K}}) + \mathbf{u}^\top \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} \right) + \text{const},
\end{aligned}$$

donde  $\boldsymbol{\Theta} = \text{diag}(\theta_1, \dots, \theta_N)$ , siendo  $\theta_n = \mathbb{E}_{q(\omega_n)} [\omega_n]$ . Juntando todo:

$$\log q^*(\mathbf{u}) = (\boldsymbol{\pi} - 1/2)^\top \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} \mathbf{u} - \frac{1}{2} \mathbf{u}^\top \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right) \mathbf{u} + \text{const}$$

Esto prueba que  $q(\mathbf{u}) = N(\mathbf{u} \mid \mathbf{m}, \mathbf{S})$  con

$$\begin{aligned}
\mathbf{m} &= \mathbf{S} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} (\boldsymbol{\pi} - 1/2) \\
\mathbf{S} &= \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \boldsymbol{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right)^{-1}
\end{aligned}$$

**Actualización de  $q(\omega_n)$ .** Fijamos  $n$ . Denotamos  $q(\omega_{j \neq n}) = \prod_{j \neq n} q(\omega_j)$ .

$$\log q^*(\boldsymbol{\omega}) = \underbrace{\mathbb{E}_{q(\omega_{j \neq n})} [\log p(\boldsymbol{\omega})]}_{C_1} + \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega_{j \neq n})} [\log p(\mathbf{y} \mid \mathbf{f}, \boldsymbol{\omega})]}_{C_2}$$

$$C_1 = \int \prod_{j \neq n} q(\omega_j) \left( \sum_i \log p(\omega_i) \right) = \log p(\omega_n) + \sum_{j \neq n} \mathbb{E}_{q(\omega_j)} [\log p(\omega_j)]$$

$$\begin{aligned}
C_2 &= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega_{j \neq n})} \left[ \sum_i \log p(y_i \mid \mathbf{f}_i, \omega_i) \right] + \text{const} = \\
&= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(y_n \mid \mathbf{f}_n, \omega_n)] + \text{const} = \\
&= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[ -\frac{1}{2} \omega_n f_n^2 \right] + \text{const} = \\
&= -\frac{1}{2} \omega_n \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2] + \text{const} \\
&= -\frac{1}{2} \omega_n \left[ \widetilde{\mathbf{K}}_{nn} - \boldsymbol{\kappa}_n \left( \mathbf{S} + \mathbf{m} \mathbf{m}^\top \right) \boldsymbol{\kappa}_n^\top \right] + \text{const}
\end{aligned}$$

Por tanto nos queda:

$$\log q^*(\omega_n) = -\frac{1}{2} \omega_n \left[ \widetilde{\mathbf{K}}_{nn} + \boldsymbol{\kappa}_n \left( \mathbf{S} + \mathbf{m} \mathbf{m}^\top \right) \boldsymbol{\kappa}_n^\top \right] + \log p(\omega_n) + \text{const},$$

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

de donde deducimos que  $q(\omega_n) = \text{PG}(\omega_n \mid 1, c_n)$  con  $c_n = \sqrt{\widetilde{\mathbf{K}}_{nm} + \kappa_n (\mathbf{S} + \mathbf{m}\mathbf{m}^\top) \kappa_n^\top}$ . Como consecuencia  $\theta_n = \tanh(c_n/2) / 2c_n$ .  $\square$

Si uno revisa la prueba del **Resultado 4.2** se da cuenta que las nuevas variables intervienen sólo a la hora de calcular la distribución variacional de  $q(\mathbf{u})$ . Estas variables se introducen para sustituir el cálculo de  $\mathbb{E}_{p(\mathbf{f}|\mathbf{u})q(\mathbf{y})} [p(\mathbf{y} \mid \mathbf{f})]$  por el de  $\mathbb{E}_{p(\mathbf{f}|\mathbf{u})q(\mathbf{y})} [p(\mathbf{y} \mid \mathbf{f}, \boldsymbol{\omega})]$ . Alternativamente, VGPMIL aproxima la función logística, que es la responsable de que  $\mathbb{E}_{p(\mathbf{f}|\mathbf{u})q(\mathbf{y})} [p(\mathbf{y} \mid \mathbf{f})]$  no admita una expresión cerrada.

#### 4.3.2. Equivalencia con VGPMIL

Si comparamos los resultados 4.2 y 4.1 nos damos cuenta de que ambos modelos usan actualizaciones muy parecidas y sólo se diferencian en las matrices  $\mathbf{\Lambda}$  y  $\mathbf{\Theta}$ . Pero estas matrices, en realidad, son las mismas. Basta tener en cuenta que  $2\sigma(x) - 1 = \tanh(x/2)$ , por lo que  $\mathbf{\Lambda} = \mathbf{\Theta}$ . Llegamos así al resultado clave de esta sección, que es la verdadera contribución de PG-VGPMIL.

**Resultado 4.3.** *VGPMIL y PG-VGPMIL son equivalentes, en el sentido de que las actualizaciones para obtener las distribuciones óptimas usando la aproximación Mean Field a la hora de emplear inferencia variacional son las mismas.*

Obsérvese que la equivalencia no se da en un sentido estricto. Lo que estamos afirmando es que ambos modelos son equivalentes si las distribuciones variacionales factorizan como indica la aproximación Mean Field. Pudiera ocurrir que la optimización de estas distribuciones usando otro método, como por ejemplo gradiente descendente, condujera a otras formas funcionales.

#### 4.4. G-VGPMIL: modelo basado en variables Gamma

Hasta ahora hemos obtenido dos modelos equivalentes por caminos distintos. PG-VGPMIL emplea *data augmentation*, modificando la probabilidad conjunta para introducir las nuevas variables Pólya-Gamma. Intuitivamente, podemos decir que se pivota en estas nuevas variables para poder obtener expresiones cerradas de las distribuciones variacionales. Por su parte, VGPMIL emplea una aproximación cuadrática a la función logística. Obsérvese que gracias a ser cuadrática podemos calcular la esperanza de la cota obtenida respecto de una distribución normal. Lo que hemos obtenido es que esta aproximación es equivalente a introducir variables Pólya-Gamma.

Una pregunta natural, y es la que nosotros nos hacemos, es: ¿qué ocurre si cambiamos la distribución de las variables  $\omega$ ? De ella surgen muchas otras cuestiones, como ¿qué distribución elegimos? o ¿podríamos interpretar que estamos aproximando la función logística de otra forma? Estas preguntas conducen a la segunda contribución del trabajo.

Motivados por la **Definición 4.1** de las variables Pólya-Gamma, en principio vamos a pedirle a las variables  $\omega_n$  que sigan una distribución Gamma. La justificación sigue a continuación. Obsérvese que una variable  $\omega \sim \text{PG}(1, 0)$  se define como una suma ponderada de variables Gamma(1, 1). Esto es,

$$\omega = \frac{1}{2\pi^2} \sum_{k=1}^{\infty} \frac{g_k}{(k - \frac{1}{2})^2},$$

#### 4.4. G-VGPMIL: modelo basado en variables Gamma

con  $g_k \sim \text{Gamma}(1, 1)$ . Si todas las variables fueran iguales a  $g \sim \text{Gamma}(1, 1)$ , ocurriría

$$\omega = \frac{g}{2\pi^2} \underbrace{\sum_{k=1}^{\infty} \frac{1}{(k - \frac{1}{2})^2}}_{\pi^2/2} = \frac{g}{4},$$

de donde  $\omega \sim \text{Gamma}(1, 4)$ . De forma general, supondremos  $\omega \sim \text{Gamma}(1, k)$  y estudiaremos el valor de  $k$ . El nuevo modelo, que llamamos G-VGPMIL, se especifica como sigue,

$$\begin{aligned} p(\mathbf{u} | \mathbf{Z}) &= \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{\mathbf{Z}\mathbf{Z}}) \\ p(\mathbf{f} | \mathbf{X}, \mathbf{Z}, \mathbf{u}) &= \mathcal{N}(\mathbf{f} | \mathbf{a}, \widetilde{\mathbf{K}}), \\ p(\mathbf{y} | \mathbf{f}, \omega) &\propto \exp\left((\mathbf{y} - 1/2)^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \boldsymbol{\Omega} \mathbf{f}\right), \\ p(\omega) &= \prod_{n=1}^N p(\omega_n), \quad \omega_n \sim \text{Gamma}(1, k) \\ p(\mathbf{T} | \mathbf{y}) &= \prod_{b=1}^B \frac{H^{G_b}}{H + 1}. \end{aligned}$$

Por primera vez en la literatura, usamos variables aleatorias Gamma para aumentar el modelo probabilístico y poder realizar inferencia con GPs en el escenario MIL. A continuación veremos que el proceso de inferencia, a pesar de ser análogo a los anteriores, conduce a un nuevo modelo con nuevas características.

##### 4.4.1. Inferencia

El proceso de inferencia es completamente análogo al de PG-VGPMIL. No vamos a repetir toda la formulación, sino que especificamos directamente la forma que toma cada una de las distribuciones variacionales.

**Resultado 4.4.** Sea  $p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \omega, \mathbf{T})$  la distribución conjunta dada por el modelo G-VGPMIL. Entonces la solución óptima al problema

$$\arg \min_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega)} \text{KL}(q(\mathbf{u})p(\mathbf{f} | \mathbf{u})q(\mathbf{y})q(\omega) || p(\mathbf{u}, \mathbf{f}, \mathbf{y}, \omega | \mathbf{T}))$$

viene dada por las siguientes distribuciones:

$$\begin{aligned} q^*(\mathbf{u}) &= \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}), \\ q^*(y_n) &= \text{Bernouilli}(y_n | \pi_n), \\ q^*(\omega_n) &= \text{Gamma}(\omega_n | 1, \beta_n), \end{aligned}$$

#### 4. Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias

siendo

$$\begin{aligned}\pi_n &= \sigma \left[ \kappa_n \mathbf{m} + \log H(2T_b - 1) \left( 1 - \mathbb{E} \left[ \max \{ \mathbf{y} \}_{b \setminus n} \right] \right) \right], \\ \mathbf{m} &= \mathbf{S} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \left( \pi - \frac{1}{2} \right), \\ \mathbf{S} &= \left( \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{ZX} \mathbf{\Theta} \mathbf{K}_{XZ} \mathbf{K}_{ZZ}^{-1} + \mathbf{K}_{ZZ}^{-1} \right)^{-1}, \\ \beta_n &= k + \frac{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]}{2} = k + \frac{\widetilde{\mathbf{K}}_{nn} + \kappa_n (\mathbf{m} \mathbf{m}^\top + \mathbf{S}) \kappa_n^\top}{2},\end{aligned}$$

donde  $\kappa_n = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n] = \mathbf{K}_{x_n Z} \mathbf{K}_{ZZ}^{-1}$ ,  $\mathbf{\Theta} = \text{diag}(\theta(\beta_1), \dots, \theta(\beta_N))$  y  $\theta(\beta) = 1/\beta$ .

*Demostración.* Las pruebas de las diferentes actualizaciones son muy similares a las de PG-VGPMIL. La de  $q(y_n)$  es la misma y no la vamos a repetir. En la actualización de  $q(\mathbf{u})$  debemos observar que la esperanza de  $\omega_n$  cambia. La actualización de  $q(\omega_n)$  es como sigue,

$$\log q^*(\omega_n) = \underbrace{\mathbb{E}_{q(\omega_j \neq n)} [\log p(\omega)]}_{A_1} + \underbrace{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega_j \neq n)} [\log p(\mathbf{y} | \mathbf{f}, \omega)]}_{A_2} + \text{const}$$

Analizamos cada término por separado,

$$\begin{aligned}A_1 &= \log p(\omega_n) + \sum_{j \neq n} \mathbb{E}_{q(\omega_j)} [\log p(\omega_j)] = \\ &= \log(k) - k\omega_n + \text{const} \\ A_2 &= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y})q(\omega_j \neq n)} \left[ \sum_i \log p(y_i | f_i, \omega_i) \right] + \text{const} = \\ &= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(y_n | f_n, \omega_n)] + \text{const} = \\ &= \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[ -\frac{1}{2} \omega_n f_n^2 \right] + \text{const} = \\ &= -\omega_n \frac{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]}{2} + \text{const}\end{aligned}$$

Por tanto nos queda:

$$\log q^*(\omega) = -\omega_n \left( k + \frac{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]}{2} \right) + \text{const}$$

de donde deducimos que  $q(\omega_n) = \text{Gamma}(\omega_n | 1, \beta_n)$  con

$$\beta_n = k + \frac{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]}{2} = k + \frac{\widetilde{\mathbf{K}}_{nn} + \kappa_n (\mathbf{S} + \mathbf{m} \mathbf{m}^\top) \kappa_n^\top}{2}.$$

Por tanto  $\theta_n = \mathbb{E}_{q(\omega_n)} [\omega_n] = 1/\beta_n$ . □

El modelo G-VGPMIL guarda cierto parecido con sus predecesores. Dedicamos la siguiente subsección a comprender las diferencias entre ellos.

#### 4.4.2. Sobre la función $\theta$ y el parámetro $k$

Al comparar las actualizaciones de PG-VGPMIL y G-VGPMIL es fácil darse cuenta de que sólo se producen cambios en la matriz diagonal  $\Theta$ . Esta matriz se forma a partir de la evaluación de una función  $\theta$ . En PG-VGPMIL esta función se evalúa sobre los parámetros  $c_n$ , mientras que en G-VGPMIL se evalúa sobre  $\beta_n$ . Sin embargo, en ambos casos esta es una función de  $\sqrt{\mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [f_n^2]}$ . Podemos escribir

$$\theta_{\text{PG}}(x) = \frac{\tanh(x/2)}{2x}, \quad (4.12)$$

$$\theta_{\text{G}k}(x) = \frac{1}{k + x^2/2}, \quad (4.13)$$

para las funciones de PG-VGPMIL y G-VGPMIL (con parámetro  $k$ ), respectivamente. Podemos analizar las diferencias entre ambos comprobando en qué se diferencian estas aplicaciones. En la Figura 4.3(a) se representa  $\theta_{\text{G}k}$  para distintos valores de  $k$ . Obsérvese que  $\theta_{\text{G}k}$  alcanza un máximo absoluto en cero con valor  $1/k$  por lo que al variar el valor de  $k$  controlamos el valor máximo que alcanza esta función.

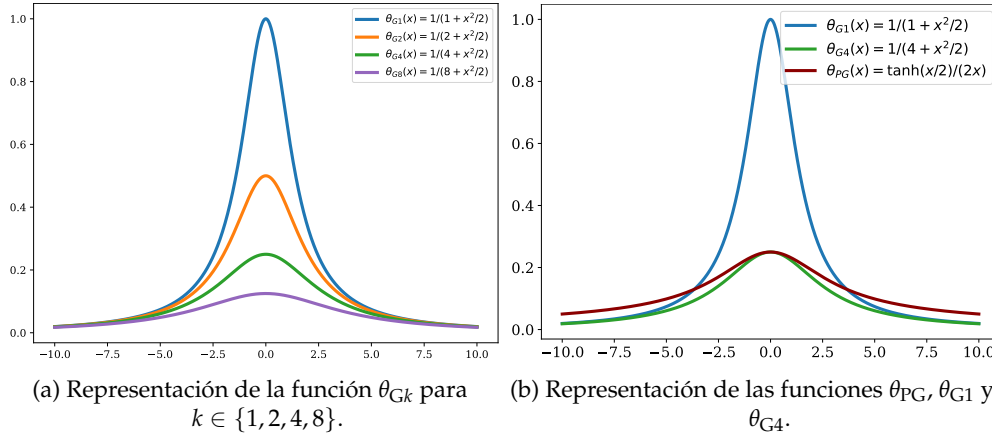


Figura 4.3.: Representación de las funciones  $\theta_{\text{PG}}$  y  $\theta_{\text{G}k}$ .

En la Figura 4.3(b) representamos  $\theta_{\text{PG}}$  junto a  $\theta_{\text{G}4}$ . No nos sorprende que al tomar  $k = 4$  ambas funciones alcancen el mismo valor máximo. Recuérdese lo que ocurría si en la definición de una variable Pólya-Gamma todas las variables Gamma eran iguales.

Debemos señalar una idea importante en relación al papel que juegan estas funciones. Si tomásemos  $\theta = 0$ , estaríamos imponiendo que la matriz de covarianza de  $q(\mathbf{u})$  fuese igual que la de la distribución a priori  $p(\mathbf{u})$ . De alguna forma, a través de la función  $\theta$  podemos cuantificar cuánto se diferenciarán las distribuciones a priori y a posteriori, o dicho de otra forma, la incertidumbre que tenemos sobre estas distribuciones. La distribución que le damos a la variable aleatoria  $\omega$  nos permite cuantificar esta incertidumbre.

Para acabar el capítulo, ofrecemos una visión general de lo que hemos conseguido hasta ahora. El estudio del modelo VGPMIL revela que el uso de la desigualdad de Jaakkola puede conducir a soluciones subóptimas para las distribuciones variacionales. Para intentar solucionar esto, hemos usado por primera vez las variables aleatorias Pólya-Gamma en un nuevo modelo basado en GPs para MIL. Esto nos ha llevado a la primera contribución de este trabajo,

#### 4. *Procesos Gaussianos para Aprendizaje a partir de Múltiples Instancias*

el modelo PG-VGPMIL. El resultado clave es la prueba de la equivalencia en el sentido Mean Field con VGPMIL. Sustituir las variables Pólya-Gamma por variables Gamma nos ha conducido a la segunda contribución del trabajo, que recibe el nombre de G-VGPMIL, y propone por primera vez el uso de variables Gamma en un modelo basado en GPs para MIL.

En el siguiente capítulo vamos a validar experimentalmente el comportamiento de G-VGPMIL, comparándolo con el ya existente VGPMIL. Para ello, nos enfrentaremos a tres problemas de clasificación binaria de MIL, siendo el último de ellos la detección de hemorragias intracraneales a partir de escáneres TAC.

## 5. Aplicaciones

En el capítulo anterior hemos estudiado tres modelos basados en GPs para problemas MIL de clasificación: VGPMIL, PG-VGPMIL y G-VGPMIL. VGPMIL ya estaba propuesto en la literatura, siendo el estado del arte a la hora de emplear el modelo de observación logístico. PG-VGPMIL y G-VGPMIL son las dos contribuciones que presenta este trabajo. Ambas están basadas en la inclusión de nuevas variables en la probabilidad conjunta. PG-VGPMIL es el primer modelo para problemas MIL de clasificación propuesto en la literatura basado en variables Pólya-Gamma. Hemos probado que es equivalente, en el sentido Mean Field, a VGPMIL. Por su parte, G-VGPMIL usa por primera vez variables Gamma y no es equivalente a los anteriores.

El objetivo de este capítulo es doble. Por una parte, evaluamos G-VGPMIL experimentalmente. VGPMIL ya está validado por la comunidad investigadora, por lo que aquí se usará para compararlo con G-VGPMIL. Como VGPMIL y PG-VGPMIL son equivalentes, sólo nos referiremos al primero de ellos. Por otra parte, resolvemos el problema de detección de hemorragias intracraneales que motiva este trabajo y constituye la tercera contribución del mismo.

### 5.1. Sobre los conjuntos de datos y la metodología empleada

La primera parte de este capítulo presenta un experimento desarrollado sobre un conjunto de datos sintético. Partiendo del clásico conjunto MNIST, diseñamos un conjunto para MIL repartiendo las imágenes en bolsas. Evaluaremos el funcionamiento de VGPMIL y G-VGPMIL de forma cuantitativa y cualitativa (mostrando qué predicciones son capaces de obtener).

En la segunda parte usamos dos conjuntos de MIL clásicos como son musk1 y musk2 para comparar ambos algoritmos. Con estos conjuntos tratamos el problema de clasificación de moléculas que explicábamos en el primer capítulo.

La tercera parte presenta el caso real de detección de hemorragias intracraneales. Usaremos dos conjuntos de imágenes de TAC (RSNA y CQ500) para desarrollar un modelo capaz de predecir la presencia de hemorragia en una persona a partir de bolsas de imágenes. Este experimento muestra que el nuevo modelo propuesto, G-VGPMIL, encaja de forma natural en las arquitecturas de Aprendizaje Profundo y son esenciales para mejorar los resultados que estas obtienen.

La información de cada uno de los conjuntos utilizados se encuentra en la [Tabla 5.1](#). Puede observarse que son de muy diversa naturaleza. Como es normal, este tipo de problemas son profundamente desbalanceados a nivel de instancia ya que contamos con muchas más instancias negativas que positivas. A nivel de bolsa se encuentra mucho más equilibrado, aunque en escenarios reales suelen predominar las bolsas negativas.

En todos los casos entrenamos los modelos VGPMIL y G-VGPMIL con el núcleo RBF. El número de iteraciones no es fijo, ya que el proceso de entrenamiento actualiza los parámetros de las distribuciones hasta que la función de pérdida (ELBO) desciende menos de  $10^{-4}$  durante 10 iteraciones. Para el número de inducing points consideramos los valores  $\{50, 100, 200\}$ , mientras que para el parámetro *lengthscale* del núcleo consideramos  $\{1, (1 + \sqrt{D})/2, \sqrt{D}, (\sqrt{D} + D)/2, D\}$ ,

## 5. Aplicaciones

	Núm. instancias	Instancias positivas	Instancias negativas	Núm. bolsas	Bolsas positivas	Bolsas negativas
MNIST	70000	6903	63097	7778	4741	3037
musk1	476	?	?	92	47	45
musk2	6598	?	?	102	39	63
RSNA	39750	5782	33968	1150	483	667
CQ500	193317	?	?	490	205	285

Tabla 5.1.: Distribución de las etiquetas a nivel de instancia y bolsa para cada conjunto de datos.

siendo  $D$  el número de rasgos o características. Finalmente, para el parámetro  $k$  de G-VGPMIL se consideran los valores  $\{0.5, 1.0, 2.0, 4.0\}$ . Los resultados que se presentan, a menos que se indique lo contrario, corresponden a la media y a la desviación estándar obtenida tras realizar validación cruzada de cinco capas.

Por último, aclaramos qué métricas vamos a considerar. Como ya hemos explicado, los problemas MIL de clasificación binaria tienen características especiales que debemos tener en cuenta a la hora de evaluar los clasificadores obtenidos. A nivel de instancia hay un gran desequilibrio, por lo que la métrica F1 tendrá una gran importancia. Junto a esta evaluaremos también la tasa de acierto o *accuracy*, el *log loss* y el área bajo la curva ROC (AUC). A nivel de bolsa obtendremos las mismas métricas, pero ahora F1 no será tan determinante.

### 5.2. Un ejemplo ilustrativo: MNIST

MNIST<sup>1</sup> es un conjunto clásico de Visión por Computador consistente en 70.000 imágenes de dígitos (del 0 al 9) manuscritos y etiquetados. Para obtener un conjunto apropiado para MIL repartimos aleatoriamente las imágenes en bolsas de 9 dígitos. Si una bolsa contiene un cero tendrá etiqueta positiva, y negativa en caso contrario. En la [Figura 5.1](#) se muestran dos bolsas de ejemplo. Obtenemos así un conjunto con 7778 bolsas, de las cuales 4741 son positivas y 3037 son negativas. El objetivo de este experimento es ilustrar el correcto funcionamiento de VGPMIL y el nuevo modelo que proponemos, G-VGPMIL.

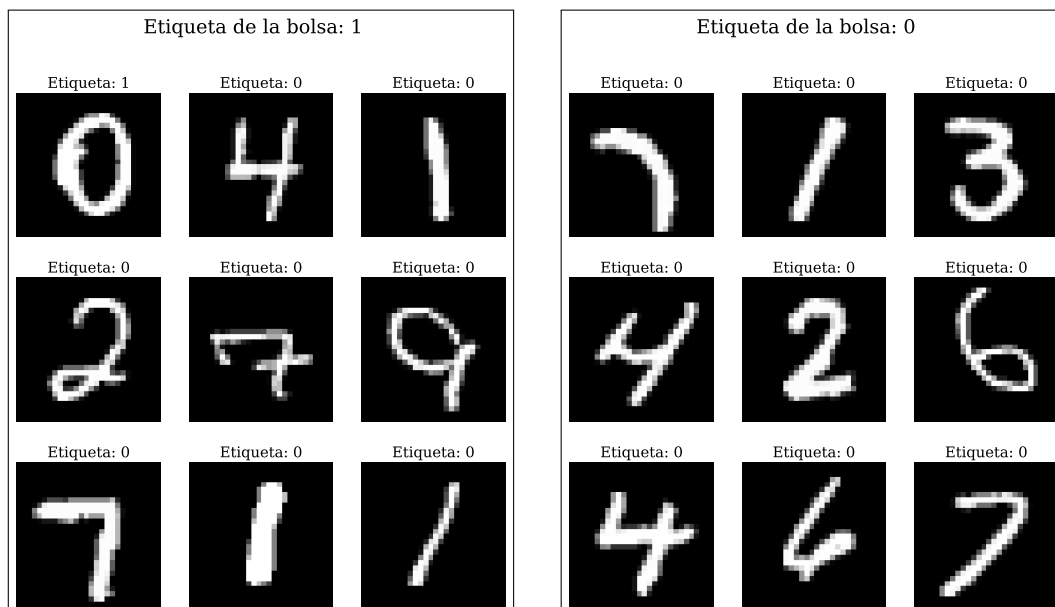
Dado que los ejemplos de MNIST son de alta dimensionalidad ( $28 \times 28 = 784$ ) evaluamos el rendimiento de VGPMIL y G-VGPMIL sin aplicar ningún preprocesamiento ([Tabla A.2](#)) y aplicando PCA (*Principal Components Analysis*, [55]) para quedarnos con las 30 componentes principales ([Tabla A.1](#)). En las tablas mencionadas podemos comprobar que G-VGPMIL es un método muy competitivo con VGPMIL y que en algunas ocasiones mejora su rendimiento.

Los resultados obtenidos los hemos resumido en las [tablas 5.2a](#) y [5.2b](#) para facilitar su discusión. Como esperábamos, G-VGPMIL presenta los mejores resultados cuando  $k = 4$  (recuérdese que se obtenía  $k = 4$  al sustituir las variables Gamma por una sola variable Gamma en la definición de las variables Pólya-Gamma). Si no aplicamos PCA, G-VGPMIL obtiene mejores resultados en cuanto a la métrica AUC se refiere. Si aplicamos PCA, los papeles se intercambian y VGPMIL se comporta mejor. Por otra parte, en ambos modelos el rendimiento aumenta conforme crece el número de inducing points. Esto tiene sentido ya que cuantos más inducing points introduzcamos más información podemos extraer del conjunto de entrenamiento. En el caso de MNIST esto es especialmente beneficioso debido al elevado número de ejemplos de entrenamiento.

En la [Figura 5.2](#) hemos representado las métricas del mejor clasificador de cada tipo. El

<sup>1</sup><http://yann.lecun.com/exdb/mnist/>





(a) Ejemplo de una bolsa con etiqueta positiva.

(b) Ejemplo de una bolsa con etiqueta negativa.

Figura 5.1.: Dos de las bolsas generadas a partir del conjunto MNIST.

mejor clasificador se elige teniendo en cuenta la métrica AUC a nivel de bolsa e instancia. Observamos que VGPMIL y G-VGPMIL obtienen resultados muy similares. Si no aplicamos PCA, G-VGPMIL se comporta mejor en las métricas AUC y Log Loss. Si aplicamos PCA ambos se encuentran muy igualados, aunque G-VGPMIL sigue presentando valores menores en Log Loss.

La [Figura 5.3](#) muestra cómo los métodos estudiados permiten realizar predicciones a nivel de instancia proporcionando también el nivel de incertidumbre asociado a la predicción. Este es el tipo de comportamiento que podemos esperar de este tipo de clasificadores probabilísticos. Para cada bolsa que tengamos, podemos predecir la probabilidad de que cada una de sus instancias pertenezcan a la clase positiva o a la clase negativa. Junto a esta probabilidad obtendremos también la incertidumbre asociada a la predicción (4.9), la cual puede interpretarse como la seguridad del clasificador en su decisión. Obsérvese que en la subfigura 5.3a se asocia una incertidumbre mucho mayor al cero que al resto de instancias. Esto puede deberse a la forma extraña que presenta el dígito. Aquí se ve reflejada la utilidad de este tipo de clasificadores. Como las etiquetas a nivel de instancia suelen ser desconocidas, identificar las bolsas positivas de esta forma nos lleva a conocer qué instancias son positivas.

Los resultados obtenidos muestran que G-VGPMIL es un clasificador para MIL mucho más que aceptable y que es capaz de competir con el estado del arte. Hemos observado que presenta un comportamiento distinto a VGPMIL que aún debemos explorar en conjuntos más complejos. Además hemos ilustrado cómo interpretar las predicciones de estos modelos.

## 5. Aplicaciones

Num. Inducing	Lengthscale	k	Modelo	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	406.0	0.5	G-VGPMIL	0.9221 ± 0.0022	0.3471 ± 0.0269	0.6504 ± 0.0001	<b>0.976 ± 0.0017</b>	0.5624 ± 0.015	0.4397 ± 0.0303	0.6855 ± 0.0005	<b>0.9394 ± 0.007</b>
	28.0	-	VGPMIL	<b>0.958 ± 0.0025</b>	<b>0.7959 ± 0.0132</b>	<b>0.2749 ± 0.0101</b>	0.9679 ± 0.002	<b>0.8557 ± 0.0103</b>	<b>0.8842 ± 0.0084</b>	<b>0.426 ± 0.0161</b>	0.9338 ± 0.007
100	406.0	1.0	G-VGPMIL	0.9417 ± 0.0015	0.5813 ± 0.0211	0.6112 ± 0.0004	<b>0.9837 ± 0.0015</b>	0.6999 ± 0.0146	0.6739 ± 0.0211	0.6698 ± 0.001	<b>0.9584 ± 0.0036</b>
	28.0	-	VGPMIL	<b>0.9684 ± 0.0013</b>	<b>0.8416 ± 0.0086</b>	<b>0.233 ± 0.0071</b>	0.9782 ± 0.0023	<b>0.8826 ± 0.0044</b>	<b>0.9043 ± 0.0032</b>	<b>0.3759 ± 0.0132</b>	0.9498 ± 0.0042
200	28.0	4.0	G-VGPMIL	<b>0.9757 ± 0.0014</b>	<b>0.8694 ± 0.0084</b>	<b>0.091 ± 0.0046</b>	0.9831 ± 0.0025	0.9009 ± 0.0066	0.9151 ± 0.0069	<b>0.2617 ± 0.0102</b>	<b>0.9687 ± 0.0042</b>
	28.0	-	VGPMIL	0.9744 ± 0.0024	0.8691 ± 0.0104	0.1984 ± 0.004	0.983 ± 0.0034	<b>0.9038 ± 0.0079</b>	<b>0.9205 ± 0.0067</b>	0.3275 ± 0.0059	0.961 ± 0.0049

(a) Sin aplicar PCA

Num. Inducing	Lengthscale	k	Modelo	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	17.738	0.5	G-VGPMIL	0.9449 ± 0.0015	0.6124 ± 0.0124	0.6484 ± 0.0002	0.9923 ± 0.002	0.72 ± 0.0124	0.7017 ± 0.0174	0.6788 ± 0.0004	0.9835 ± 0.0026
	5.477	-	VGPMIL	<b>0.9855 ± 0.0013</b>	<b>0.9257 ± 0.0077</b>	<b>0.1414 ± 0.0041</b>	<b>0.993 ± 0.0017</b>	<b>0.9445 ± 0.0072</b>	<b>0.954 ± 0.0058</b>	<b>0.2353 ± 0.0076</b>	<b>0.9846 ± 0.002</b>
100	3.238	4.0	G-VGPMIL	0.9875 ± 0.0009	0.9333 ± 0.0051	0.0509 ± 0.0011	0.9961 ± 0.0006	0.9469 ± 0.0063	0.9547 ± 0.0056	0.1598 ± 0.0073	0.9902 ± 0.0011
	3.238	-	VGPMIL	<b>0.9909 ± 0.0009</b>	<b>0.9527 ± 0.0059</b>	0.1102 ± 0.0031	<b>0.9968 ± 0.0007</b>	<b>0.9645 ± 0.0041</b>	<b>0.9705 ± 0.0034</b>	0.1833 ± 0.0041	<b>0.9921 ± 0.0011</b>
200	3.238	4.0	G-VGPMIL	0.9897 ± 0.0005	0.9453 ± 0.0027	<b>0.0429 ± 0.0008</b>	0.9974 ± 0.0003	0.9576 ± 0.0047	0.964 ± 0.0041	<b>0.1379 ± 0.0068</b>	0.9934 ± 0.0012
	3.238	-	VGPMIL	<b>0.9924 ± 0.0006</b>	<b>0.9607 ± 0.0034</b>	0.0957 ± 0.0014	<b>0.9978 ± 0.0004</b>	<b>0.9691 ± 0.0032</b>	<b>0.9744 ± 0.0027</b>	0.1597 ± 0.0025	<b>0.9944 ± 0.0008</b>

(b) Aplicando PCA.

Tabla 5.2.: Resumen de los resultados en MNIST. Para cada valor del número de inducing points considerado, se recogen las métricas del mejor clasificador de cada tipo. Se señala en negrita el mejor valor de cada métrica dentro de cada grupo, y subrayado el mejor de cada columna.

### 5.3. Los conjuntos musk

En esta sección vamos a trabajar en un escenario de MIL frecuentemente usado para evaluar modelos de MIL. Los conjuntos musk1<sup>2</sup> y musk2<sup>3</sup> están formados por datos tabulares. Cada fila (instancia) representa una de las posibles "formas" que puede adoptar una molécula. Las filas se agrupan en bolsas y cada bolsa conforma una molécula. El objetivo es predecir si esa molécula dará lugar a un almizcle o musco<sup>4</sup>. Esto queda determinado por las formas que puede adoptar esa molécula, por lo que habrá formas que den lugar a un musco (instancias positivas) y otras que no (instancias negativas). En estos conjuntos no se disponen de las etiquetas a nivel de instancia, por lo que sólo podemos evaluar los modelos a nivel de bolsa. Este ya no es un escenario sintético de MIL, que además tiene una gran utilidad. La identificación de bolsas positivas con modelos probabilísticos nos permite conocer qué tipo de formas pueden dar lugar a un musco.

Todos los resultados que vamos a discutir en esta sección se encuentra en las tablas A.3 y A.4. Nuevamente, nos centramos en los resultados del mejor clasificador (de acuerdo a las métricas AUC a nivel de instancia y de bolsa) de cada tipo para cada valor del número de *inducing points*:

- En musk1 (Tabla 5.3) ambos modelos son capaces de alcanzar el mismo rendimiento. Obsérvese de nuevo que predomina el valor de  $k = 4$  para G-VGPMIL y que las métricas crecen conforme aumenta el número de *inducing points* considerado. Algo que llama la atención es que al fijar 200 *inducing points* el mejor valor de *lengthscale* es 1.0. Por otra parte, también debemos notar que aunque al usar 200 *inducing points* los clasificadores son mejores, los valores más bajos de Log Loss se consiguen al usar 100 *inducing points*.
- Si ahora nos fijamos en musk2 (Tabla 5.4) observamos un comportamiento totalmente distinto al que habíamos detectado hasta ahora. Ahora G-VGPMIL presenta mejores resultados cuando  $k = 2$ . Además, las métricas sólo mejoran al aumentar el número de

<sup>2</sup>[https://archive.ics.uci.edu/ml/datasets/Musk+\(Version+1\)](https://archive.ics.uci.edu/ml/datasets/Musk+(Version+1))

<sup>3</sup>[https://archive.ics.uci.edu/ml/datasets/Musk+\(Version+2\)](https://archive.ics.uci.edu/ml/datasets/Musk+(Version+2))

<sup>4</sup><https://es.wikipedia.org/wiki/Almizcle>

## 5.4. Detección de hemorragias intracraneales

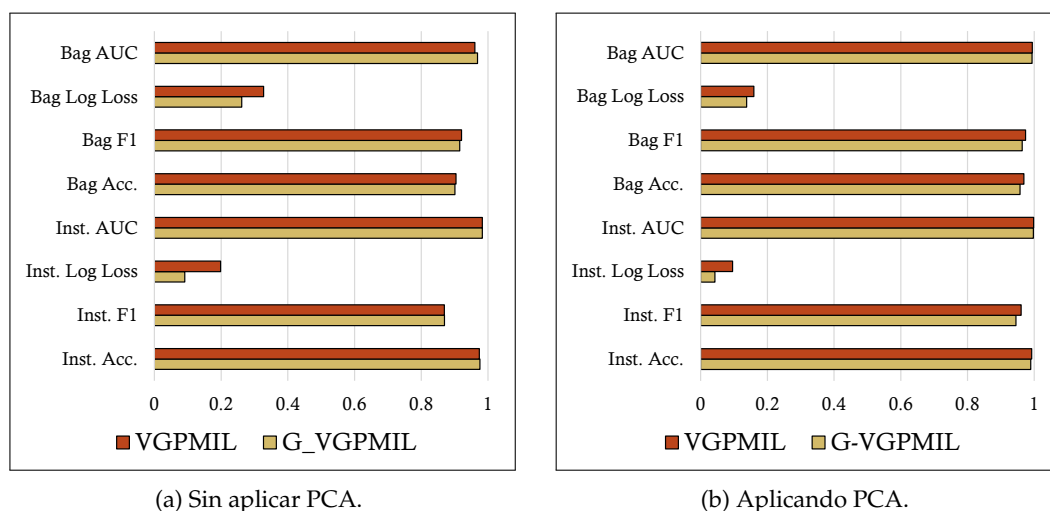


Figura 5.2.: Gráficos de barras de los experimentos de MNIST. Estos gráficos recogen las métricas del mejor clasificador de cada tipo.

*inducing points* en el caso de VGPMIL. Puede verse que G-VGPMIL se muestra muy por encima de VGPMIL en este conjunto de datos.

Num. Inducing	Lengthscale	k	Modelo	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	12.884	0.5	C-VGPMIL	<b>0.5968 ± 0.1239</b>	0.2775 ± 0.14	<b>0.6671 ± 0.0195</b>	0.7975 ± 0.0499	0.6427 ± 0.0957	0.4449 ± 0.2241	0.6643 ± 0.0101	0.859 ± 0.0631
	12.884	-	VGPMIL	0.5706 ± 0.1236	0.2118 ± 0.1242	0.7125 ± 0.1322	0.7148 ± 0.0953	0.6216 ± 0.1119	0.4162 ± 0.2313	<b>0.6568 ± 0.0656</b>	0.7985 ± 0.0493
100	12.884	0.5	G-VGPMIL	<b>0.5867 ± 0.1287</b>	<b>0.2643 ± 0.1293</b>	<b>0.6674 ± 0.0185</b>	<b>0.809 ± 0.0541</b>	<b>0.6316 ± 0.1077</b>	<b>0.4412 ± 0.2282</b>	0.665 ± 0.0089	<b>0.8615 ± 0.0644</b>
	12.884	-	VGPMIL	0.5683 ± 0.1238	0.2054 ± 0.1204	0.7086 ± 0.128	0.7191 ± 0.0941	0.6216 ± 0.1119	0.4162 ± 0.2313	<b>0.6522 ± 0.0629</b>	0.8081 ± 0.0457
200	12.884	0.5	G-VGPMIL	0.5917 ± 0.1236	0.2674 ± 0.1181	<b>0.6676 ± 0.018</b>	<b>0.8128 ± 0.0567</b>	0.6427 ± 0.0957	0.4449 ± 0.2241	<b>0.6651 ± 0.0079</b>	<b>0.8711 ± 0.0608</b>
	1.0	-	VGPMIL	<b>0.7432 ± 0.0786</b>	<b>0.6769 ± 0.081</b>	0.6931 ± 0.0	0.7889 ± 0.0369	<b>0.8053 ± 0.0606</b>	<b>0.7896 ± 0.0845</b>	0.6931 ± 0.0001	0.8567 ± 0.0292

Tabla 5.3.: Resumen de los resultados en musk1. Para cada valor del número de inducing points considerado, se recogen las métricas del mejor clasificador de cada tipo. Se señala en negrita el mejor valor de cada métrica dentro de cada grupo, y subrayado el mejor de cada columna.

Si analizamos las tablas A.3 y A.4, donde se presentan todos los resultados, nos damos cuenta de que G-VGPMIL es muy superior a VGPMIL. Obsérvese que para cada valor del número de *inducing points* casi todas las métricas son mejores al usar alguno de los clasificadores G-VGPMIL. Al fijar este valor y el de *lengthscale* estamos comparando cómo se comportan los dos clasificadores partiendo de la misma información. De alguna forma, estamos evaluando cuánto conocimiento es capaz de extraer cada uno cuando se encuentran exactamente en las mismas condiciones. La Figura 5.4 representa la métrica AUC obtenida por ambos clasificadores en cada uno de los conjuntos. Obsérvese que la longitud de las barras amarillas (G-VGPMIL) es en todos los casos mayor o igual que las barras rojas (VGPMIL). Esto prueba la superioridad de G-VGPMIL en este conjunto de datos.

## 5.4. Detección de hemorragias intracraneales

En esta sección vamos a presentar una solución para el problema de detección de hemorragias intracraneales basada en los modelos que hemos desarrollado. En la introducción a este trabajo

## 5. Aplicaciones

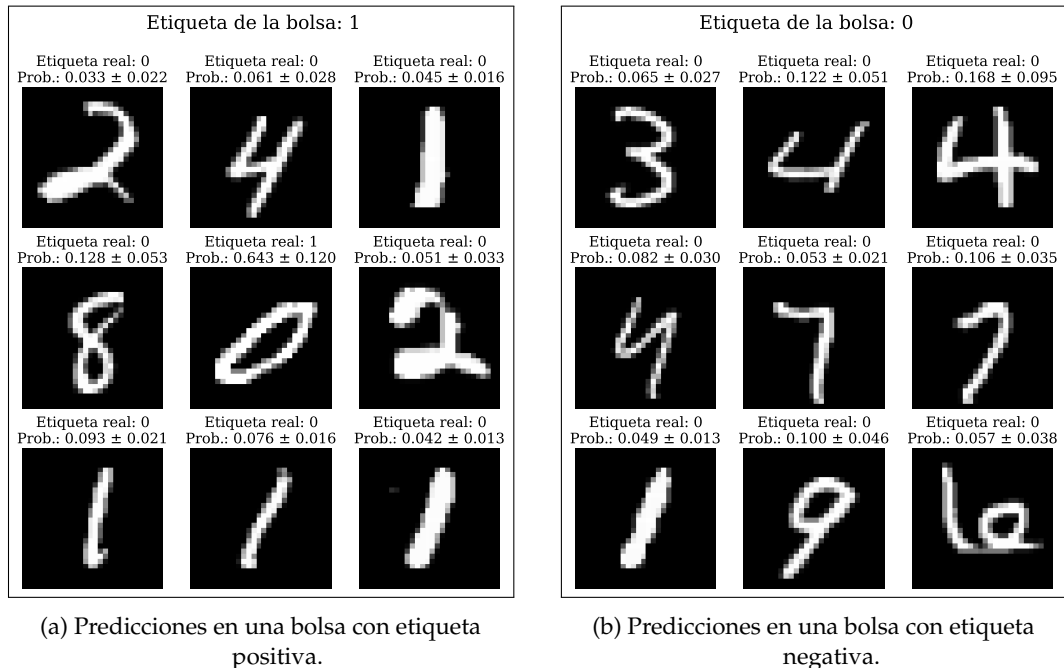


Figura 5.3.: Predicciones de dos de las bolsas generadas a partir del conjunto MNIST. Todas las instancias se clasifican correctamente. Como consecuencia, la etiqueta que se predice de cada bolsa también es correcta.

ya se explicó la importancia de obtener un diagnóstico temprano de este tipo de lesiones. Además, también se justificó la idoneidad de los modelos probabilísticos para poder obtener predicciones interpretables y fiables, como las que se muestran en la Figura 5.3 para el caso de MNIST.

La detección de este tipo de lesiones se enmarca de forma natural dentro del paradigma de MIL. Cada escáner es una bolsa, y cada uno de los cortes que lo componen es una instancia. Si un paciente padece la lesión, su escáner se etiqueta como positivo. Pero, como se muestra en la Figura 5.5, no todos los cortes de un escáner positivo presentarán una etiqueta positiva, sino sólo unos pocos de ellos. Esto se debe a que el hematoma no se produce en todo el cráneo, sino sólo en una parte localizada de él. Obsérvese que en la subfigura 5.5a se aprecia una "mancha", correspondiente a la hemorragia.

### 5.4.1. Conjuntos utilizados

Para poder entrenar y evaluar los modelos propuestos vamos a utilizar dos conjuntos de datos. El primero de ellos fue publicado por la Sociedad Norteamericana de Radiología (*Radiological Society of North America, RSNA*) [14] y se puede descargar desde Kaggle<sup>5</sup>. Está formado por 39750 cortes obtenidos de 1150 personas. De estas, se apartan 150 para test y 1000 para entrenamiento. El conjunto de entrenamiento contiene 589 escáneres normales (etiqueta negativa) y 411 escáneres con lesión (etiqueta positiva). El conjunto de test contiene 78 negativos y 72 positivos. El número de cortes en cada escáner varía entre desde 24 y 57, y cada uno tiene

<sup>5</sup><https://www.kaggle.com/c/rsna-intracranial-hemorrhage-detection>

#### 5.4. Detección de hemorragias intracraneales

Num. Inducing	Lengthscale	k	Modelo	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	83.0	0.5	G-VGPMIL	0.8144 ± 0.0909	0.0183 ± 0.0366	0.6256 ± 0.0257	<b>0.7883 ± 0.1846</b>	0.6276 ± 0.0196	0.0444 ± 0.0889	0.6649 ± 0.0052	<b>0.8281 ± 0.0871</b>
	12.884	-	VGPMIL	<b>0.8262 ± 0.1102</b>	<b>0.2166 ± 0.2576</b>	<b>0.4003 ± 0.1561</b>	0.7425 ± 0.1759	<b>0.6662 ± 0.0971</b>	<b>0.338 ± 0.2068</b>	<b>0.5475 ± 0.1046</b>	0.7987 ± 0.107
100	12.884	2.0	G-VGPMIL	<b>0.8251 ± 0.1105</b>	<b>0.2973 ± 0.2244</b>	0.4366 ± 0.1274	<b>0.7721 ± 0.1503</b>	<b>0.7167 ± 0.0598</b>	<b>0.5221 ± 0.1362</b>	<b>0.5472 ± 0.0397</b>	<b>0.8271 ± 0.0915</b>
	12.884	-	VGPMIL	0.8224 ± 0.104	0.212 ± 0.232	<b>0.4057 ± 0.1511</b>	0.7374 ± 0.1713	0.6667 ± 0.0724	0.3453 ± 0.1301	0.5484 ± 0.1048	0.8157 ± 0.1014
200	12.884	1.0	G-VGPMIL	<b>0.8343 ± 0.0894</b>	<b>0.2606 ± 0.2131</b>	0.5627 ± 0.0606	<b>0.7888 ± 0.1704</b>	<b>0.7257 ± 0.0577</b>	<b>0.4889 ± 0.1438</b>	0.6204 ± 0.0181	<b>0.8389 ± 0.1011</b>
	12.884	-	VGPMIL	0.8226 ± 0.1101	0.2278 ± 0.2362	<b>0.3964 ± 0.1476</b>	0.7497 ± 0.1668	0.6767 ± 0.0575	0.3756 ± 0.1034	<b>0.5421 ± 0.1112</b>	0.8181 ± 0.1064

Tabla 5.4.: Resumen de los resultados en musk2. Para cada valor del número de inducing points considerado, se recogen las métricas del mejor clasificador de cada tipo. Se señala en negrita el mejor valor de cada métrica dentro de cada grupo, y subrayado el mejor de cada columna.

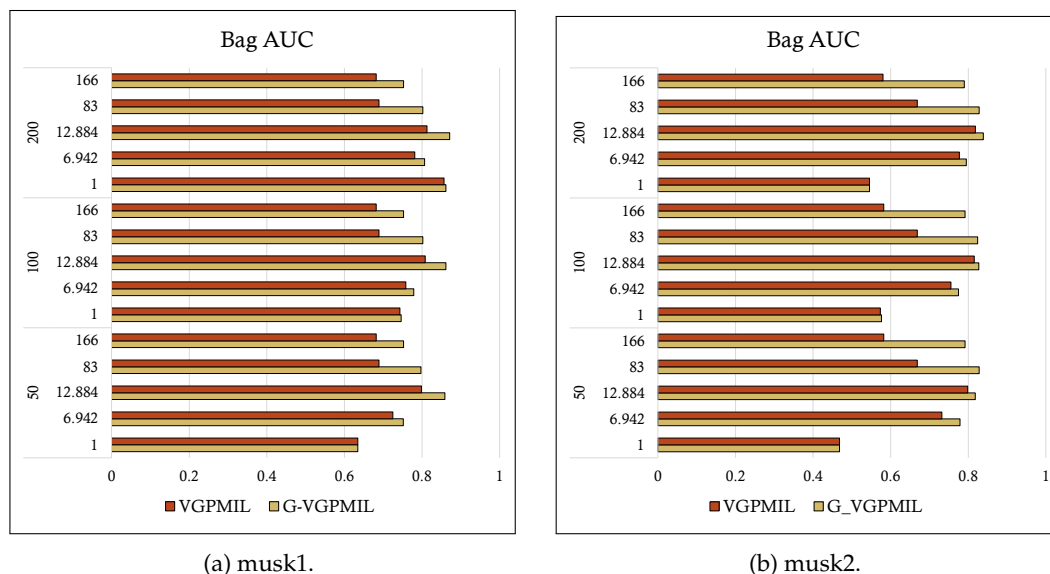


Figura 5.4.: Gráfico de barras de la métrica AUC en cada una de las configuraciones de los experimentos con los conjuntos musk1 y musk2. El eje vertical refleja el número de *inducing points* (50, 100 ó 200) seguido del parámetro *lengthscale* (1, 6.942, 12.884, 83 y 166). Obsérvese cómo la barra amarilla tiene en todos los casos una longitud mayor o igual a la roja.

dimensiones  $512 \times 512$ . En entrenamiento tenemos 29520 cortes negativos y 4976 cortes positivos. En test contamos con 4448 negativos y 806 positivos. Obsérvese que en este conjunto disponemos de las etiquetas tanto a nivel de instancia como a nivel de bolsa.

El otro conjunto que usaremos se denomina CQ500 y podemos acceder a él desde V7<sup>6</sup>. Está formado por escáneres de pacientes de varios centros médicos de India [8]. El número de cortes en cada escáner varía desde 16 hasta 128 cortes. En este conjunto sólo disponemos de etiquetas a nivel del escáner completo y lo usaremos para evaluación. En este tipo de trabajos es común evaluar los modelos desarrollados en este conjunto para poder compararlos con los de la comunidad.

<sup>6</sup><https://www.v7labs.com/open-datasets/cq500>

## 5. Aplicaciones

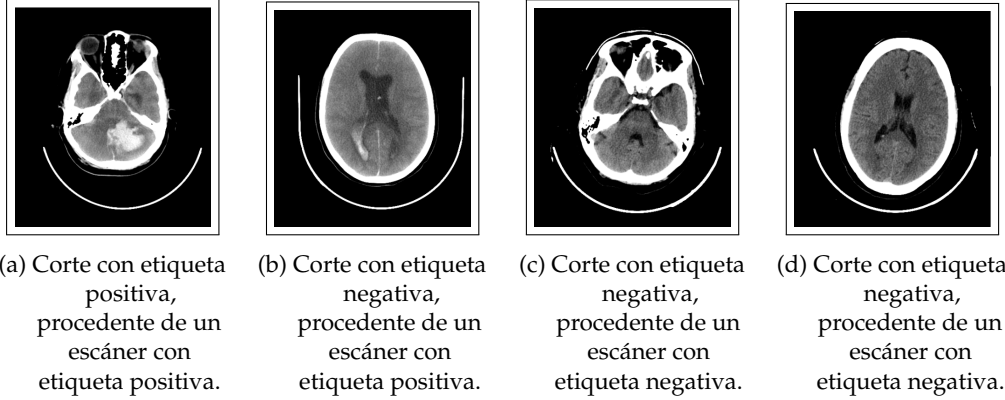


Figura 5.5.: Ejemplos de cortes del conjunto CQ500. Las subfiguras (a) y (b) muestran dos cortes de un escáner correspondiente a un paciente que presenta la lesión. Las subfiguras (c) y (d) muestran cortes de un escáner correspondiente a un paciente que no presenta la lesión.

### 5.4.2. Preprocesamiento

Para analizar los cortes de TAC, los radiólogos se fijan en diferentes intervalos de radiodensidad dentro de cada corte. Para imitar este comportamiento, a partir de cada imagen original  $I$  se extraen tres nuevas imágenes,  $I_1, I_2, I_3$  de las mismas dimensiones. Cada una de estas imágenes  $I_n$  está definida mediante un par  $(c_n, w_n)$  donde  $c_n$  es el centro del intervalo y  $w_n$  es la anchura del intervalo,

$$I_n(i, j) = \begin{cases} c_n - w_n/2 & \text{si } I(i, j) < c_n - w_n/2, \\ I(i, j) & \text{si } c_n - w_n/2 \leq I(i, j) \leq c_n + w_n/2, \\ c_n + w_n/2 & \text{si } c_n + w_n/2 < I(i, j). \end{cases}$$

Los centros y anchuras que definen a cada una de las nuevas imágenes son  $(40, 80)$ ,  $(80, 200)$  y  $(40, 380)$ , todos ellos en unidades Hounsfield. Cada imagen captura un tipo de tejido diferente, de forma que la primera representa el tejido del cerebro, la segunda el subdural, y la tercera el hueso. Este proceso transforma cada matriz  $I \in \mathbb{R}^{H \times W}$  (en escala de grises) en una matriz multidimensional o tensor  $(I_1, I_2, I_3) \in \mathbb{R}^{H \times W \times 3}$ . En la [Figura 5.6](#) se muestra un ejemplo de este preprocesamiento.

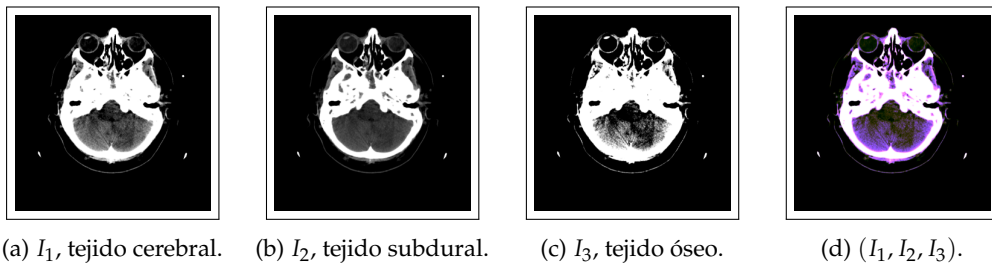


Figura 5.6.: Ejemplo de un corte de TAC preprocesado.

### 5.4.3. Metodología y arquitecturas

Para aumentar la capacidad predictiva de nuestros clasificadores, primero extraeremos características de cada uno de los cortes mediante una red convolucional. Con estas características, posteriormente se entrenarán los modelos G-VGPMIL y VGPMIL. De esta forma, siguiendo el modelo propuesto en [53], el entrenamiento se organiza en dos fases.

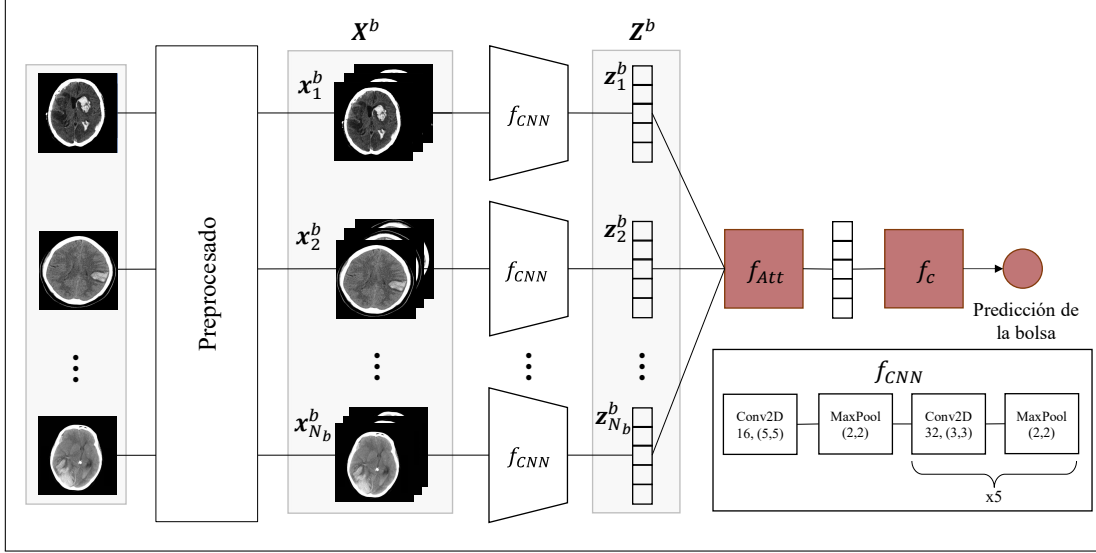


Figura 5.7.: Esquema de AttCNN. Imagen adaptada de [53].

**Fase 1 (AttCNN).** Usamos la arquitectura propuesta en [53], la cual está basada en capas convolucionales y una capa de atención. Un esquema de la misma puede observarse en la Figura 5.7. Este modelo recibe el nombre de AttCNN y responde a la composición de tres funciones, de forma que para cada bolsa  $\mathbf{X}^b$  se aproxima  $p(T_b = 1 | \mathbf{X}^b) \approx \text{AttCNN}(\mathbf{X}^b) = (f_c \circ f_{Att} \circ f_{CNN})(\mathbf{X}^b)$ . A continuación describimos cada una de estas funciones.

- En primer lugar se aplica una red convolucional, que sirve como un primer extractor de características. Esta red se corresponde con la aplicación de  $f_{CNN}: \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^D$  a cada instancia de la bolsa y la posterior concatenación de sus resultados. Esto es,

$$f_{CNN}(\mathbf{X}^b) = [f_{CNN}(x_1^b), \dots, f_{CNN}(x_{N_b}^b)],$$

donde  $N_b$  es el número de ejemplos en la bolsa  $b$ . Obsérvese que esta aplicación transforma cada bolsa  $\mathbf{X}^b$  en otra  $\mathbf{Z}^b = f_{CNN}(\mathbf{X}^b)$  que será la entrada a VGPMIL y G-VGPMIL. El valor de  $D$  determina el número de características que usarán estos clasificadores. Se mostrarán resultados con  $D \in \{8, 32, 128\}$ .

- A continuación se aplica una capa de atención. Esta capa se corresponde con la aplicación  $f_{Att}: \mathbb{R}^{D \times N_b} \rightarrow \mathbb{R}^D$  y asigna un peso  $\alpha_i \in \mathbb{R}$  a cada vector de entrada, de forma que

$$f_{Att}(z_1, \dots, z_{N_b}) = \sum_{i=1}^{N_b} \alpha_i z_i.$$

## 5. Aplicaciones

Los pesos se definen como

$$\alpha_i = \frac{\exp(\mathbf{w}^\top \tanh(Vz_i))}{\sum_{j=1}^{N_b} \exp(\mathbf{w}^\top \tanh(Vz_j))},$$

donde  $\mathbf{w} \in \mathbb{R}^{50}$  y  $V \in \mathbb{R}^{50 \times D}$  son parámetros entrenables. Obsérvese que estos pesos siempre suman uno.

- Por último se aplica un clasificador  $f_c: \mathbb{R}^D \rightarrow [0, 1]$  cuya salida se puede interpretar como la probabilidad de que ese ejemplo pertenezca a la clase 1. Este clasificador está formado por una capa totalmente conectada con activación sigmoidal.

En la primera fase AttCNN se entrena para minimizar la entropía cruzada entre las verdaderas etiquetas de cada bolsa y las que ella misma predice. Tras este proceso  $f_{CNN}$  será capaz de extraer características apropiadas para el problema de MIL. Las características extraídas por AttCNN se encuentran disponibles en Github<sup>7</sup>.

**Fase 2 (AttCNN+(G-)VGPMIL).** En la segunda fase AttCNN se modifica para poder usar los clasificadores basados en GPs. La nueva arquitectura se ilustra en la Figura 5.8. Como ya habíamos adelantado, cada bolsa  $\mathbf{X}^b$  se transforma en  $\mathbf{Z}^b = f_{CNN}(\mathbf{X}^b)$ , y estas últimas se usan para aprender las distribuciones variacionales de los GPs. Se predice la etiqueta de cada instancia  $y$ , a partir de ellas, se obtiene una predicción de la etiqueta de la bolsa. Hay que recalcar que en esta segunda fase,  $f_{CNN}$  se encuentra "congelada", en el sentido de que sus parámetros no se modifican.

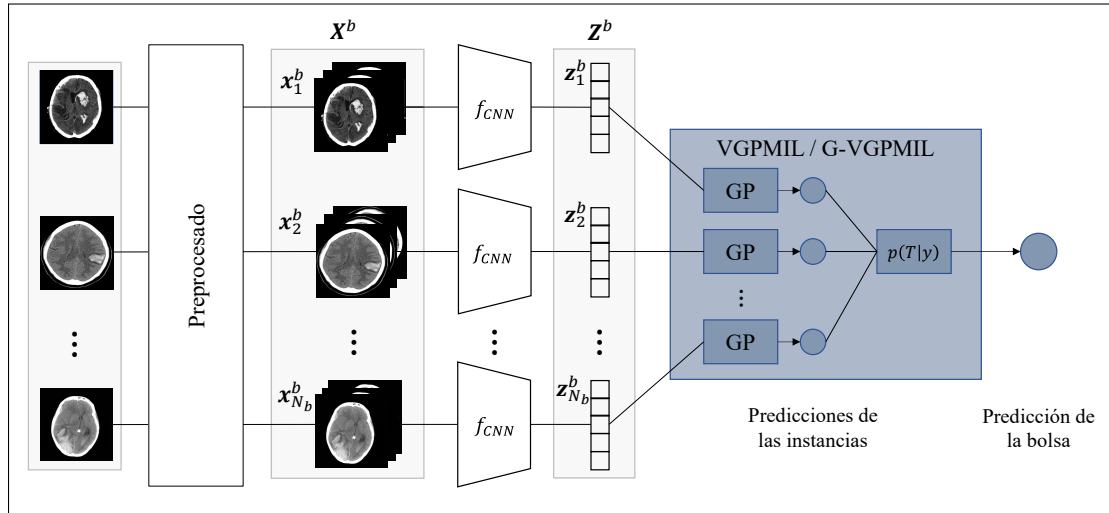


Figura 5.8.: Esquema de AttCNN+(G-)VGPMIL. Imagen adaptada de [53].

### 5.4.4. Resultados y discusión

En este problema contamos con tres clasificadores. El primero es AttCNN, y los otros dos son el resultado de combinar las características extraídas por AttCNN con VGPMIL y G-VGPMIL.

<sup>7</sup>[https://github.com/arneschmidt/cnn\\_plus\\_vgpmil](https://github.com/arneschmidt/cnn_plus_vgpmil)



Para cada valor de  $D$  se han repetido todos los experimentos, cuyos resultados se encuentran en las tablas A.5, A.6 y A.7.

#### 5.4.4.1. AttCNN vs VGPMIL vs G-VGPMIL

Comencemos discutiendo los resultados sobre el conjunto de RSNA, que se resumen en la [Tabla 5.5](#) y en la [Figura 5.9](#) (subfiguras 5.9a y 5.9b). Lo primero que observamos es que VGPMIL y G-VGPMIL suponen un salto de calidad respecto a AttCNN. Obsérvese que G-VGPMIL ofrece resultados muy competitivos con los de VGPMIL. En las métricas Accuracy y F1 nuestro modelo se encuentra más de una décima por encima. Sólo a nivel de instancia AttCN y G-VGPMIL se encuentran más igualados. AttCNN usa los pesos de atención para aproximar la probabilidad a este nivel, lo cual produce predicciones de etiqueta acertadas al fijar el umbral en 0.5 (métricas Accuracy y F1). Sin embargo, estos pesos no funcionan mejor que G-VGPMIL a la hora de estimar las probabilidades (métricas Log Loss y AUC).

Num. características	Modelo	RSNA inst. Acc.	RSNA inst. F1	RSNA inst. Log Loss	RSNA inst. AUC	RSNA bag Acc.	RSNA bag F1	RSNA bag Log Loss	RSNA bag AUC
8	AHCNN	0.9228 ± 0.0042	0.7732 ± 0.0071	0.3347 ± 0.006	0.8956 ± 0.003	0.7813 ± 0.0208	0.8108 ± 0.0155	7.5526 ± 0.7194	0.7888 ± 0.0202
	G-VGPMIL	0.9211 ± 0.0024	0.6662 ± 0.0142	0.221 ± 0.0052	0.9613 ± 0.0008	0.884 ± 0.009	0.8654 ± 0.0115	0.2817 ± 0.0135	0.9596 ± 0.0072
	VGPMIL	<b>0.9397 ± 0.0016</b>	<b>0.7778 ± 0.0077</b>	<b>0.1689 ± 0.0043</b>	<b>0.9649 ± 0.002</b>	<b>0.8973 ± 0.0161</b>	<b>0.8914 ± 0.0168</b>	<b>0.2583 ± 0.0204</b>	<b>0.9604 ± 0.007</b>
32	AHCNN	0.9022 ± 0.0172	0.7359 ± 0.0293	0.3561 ± 0.001	0.8931 ± 0.0073	0.704 ± 0.065	0.7619 ± 0.0388	10.2237 ± 2.2448	0.7144 ± 0.0623
	G-VGPMIL	0.9243 ± 0.0088	0.7274 ± 0.0519	0.401 ± 0.0171	0.9414 ± 0.0066	0.868 ± 0.0314	0.8668 ± 0.022	0.4468 ± 0.0653	0.9568 ± 0.0029
	VGPMIL	<b>0.9316 ± 0.0097</b>	<b>0.7576 ± 0.0342</b>	<b>0.1895 ± 0.0176</b>	<b>0.9501 ± 0.0058</b>	<b>0.8893 ± 0.0161</b>	<b>0.8862 ± 0.0158</b>	<b>0.2791 ± 0.0222</b>	<b>0.9583 ± 0.0048</b>
128	AHCNN	0.9129 ± 0.0117	0.7557 ± 0.0223	0.3501 ± 0.003	0.897 ± 0.0038	0.7307 ± 0.0443	0.7777 ± 0.0262	9.3026 ± 1.5316	0.7401 ± 0.0423
	G-VGPMIL	0.928 ± 0.0036	0.7328 ± 0.012	0.4011 ± 0.0183	0.9415 ± 0.0051	0.8973 ± 0.0229	0.8902 ± 0.0216	0.4757 ± 0.0527	0.9657 ± 0.0083
	VGPMIL	<b>0.9344 ± 0.0069</b>	<b>0.7647 ± 0.0243</b>	<b>0.1797 ± 0.0112</b>	<b>0.9534 ± 0.0049</b>	<b>0.908 ± 0.0186</b>	<b>0.9037 ± 0.0196</b>	<b>0.246 ± 0.0204</b>	<b>0.966 ± 0.0062</b>

Tabla 5.5.: Resultados sobre el conjunto de RSNA. Para cada valor del número de características,  $D$ , se recogen las métricas obtenidas por los mejores modelos. Se señala en negrita el mejor valor de cada métrica para cada valor de  $D$ , y subrayado el mejor de cada columna.

Debemos recalcar la superioridad de VGPMIL en el conjunto de RSNA, que obtiene grandes resultados tanto a nivel de bolsa como a nivel de instancia. Para todos los valores de  $D$  este clasificador es, indiscutiblemente, el que mejor funciona. Llama la atención que el mejor modelo a nivel de bolsa no lo sea a nivel de instancia, y viceversa. Aunque el rendimiento a ambos niveles está claramente relacionado, esto muestra que no debemos concentrarnos en uno de ellos.

Vamos a centrarnos ahora en el conjunto CQ500, cuyos resultados se encuentran en la [Tabla 5.6](#) y en la [Figura 5.9c](#). Como era de esperar, todas la métricas son inferiores a las de RSNA. Esto ocurre porque los modelos han sido entrenados en un conjunto diferente, y al evaluarlos en CQ500 el rendimiento disminuye debido a la nueva naturaleza de las imágenes. Aún así, los resultados para VGPMIL y G-VGPMIL siguen siendo satisfactorios. AttCNN vuelve a quedarse por detrás de VGPMIL y G-VGPMIL, pero ahora este último ocupa el primer puesto con diferencia para todos los valores de  $D$ . En métricas como la tasa de acierto, G-VGPMIL se sitúa hasta dos centésimas por encima. Esto pone de manifiesto que G-VGPMIL presenta una capacidad de generalización mayor que la de VGPMIL.

#### 5.4.4.2. Número de características

En la discusión anterior no hemos profundizado en el papel que juega el parámetro  $D$ . Recordemos que este número controla la dimensión del vector de características que se extrae de cada corte, por lo que esperamos que influya en la capacidad de decisión de cada clasificador.

## 5. Aplicaciones

Num. características	Modelo	CQ500 bag Acc.	CQ500 bag F1	CQ500 bag Log Loss	CQ500 bag AUC
8	AttCNN	0.6547 ± 0.0387	0.7004 ± 0.0204	0.6238 ± 0.0442	0.9061 ± 0.0088
	G-VGPMIL	<b>0.8645 ± 0.013</b>	<b>0.8216 ± 0.02</b>	<b>0.3457 ± 0.0223</b>	<b>0.9235 ± 0.0058</b>
	VGPMIL	0.8314 ± 0.0077	0.8096 ± 0.0093	0.3766 ± 0.011	0.922 ± 0.0061
32	AttCNN	0.6329 ± 0.0731	0.6924 ± 0.0536	0.6229 ± 0.0823	0.8717 ± 0.0255
	G-VGPMIL	0.8141 ± 0.0076	0.7409 ± 0.0265	0.4637 ± 0.0307	<b>0.8869 ± 0.03</b>
	VGPMIL	<b>0.8252 ± 0.0327</b>	<b>0.7893 ± 0.0439</b>	<b>0.4189 ± 0.05</b>	0.8812 ± 0.0311
128	AttCNN	0.6673 ± 0.0407	0.712 ± 0.0348	0.5956 ± 0.0736	0.8988 ± 0.0215
	G-VGPMIL	<b>0.8542 ± 0.0171</b>	<b>0.8199 ± 0.0332</b>	0.4148 ± 0.041	<b>0.9115 ± 0.0197</b>
	VGPMIL	0.8392 ± 0.0199	0.813 ± 0.0282	<b>0.3976 ± 0.035</b>	0.9096 ± 0.0207

Tabla 5.6.: Resultados sobre el conjunto de CQ500. Para cada valor del número de características,  $D$ , se recogen las métricas obtenidas por los mejores modelos. Se señala en negrita el mejor valor de cada métrica para cada valor de  $D$ , y subrayado el mejor de cada columna.

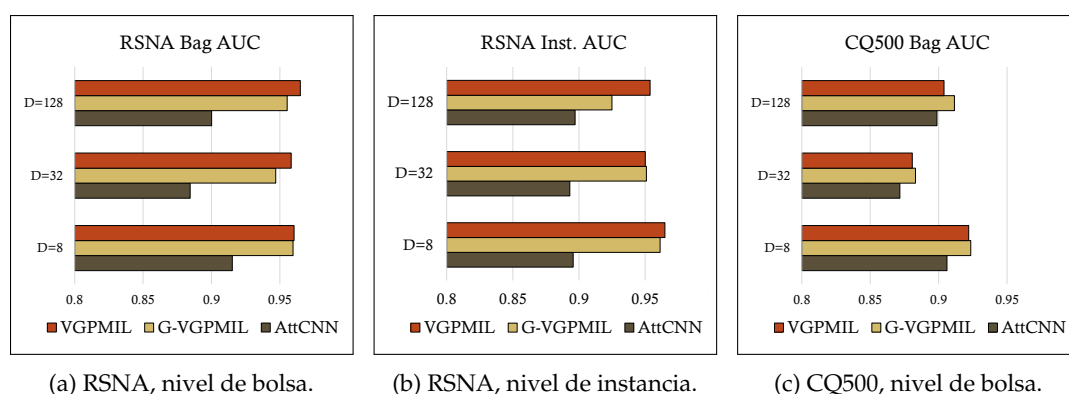


Figura 5.9.: Métrica AUC, tanto a nivel de bolsa como a nivel de instancia, en los conjuntos de RSNA y CQ500.

Para poder analizarlo hemos representado cómo varía la métrica AUC, tanto a nivel de instancia como de bolsa, en la Figura 5.10. El rendimiento empeora cuando fijamos a 32 el número de características, dibujándose una especie de valle en las gráficas. A nivel de bolsa es evidente que con 128 características se obtiene un mejor comportamiento, pero a nivel de instancia el mejor desempeño ocurre cuando este número se fija a 8. Sin embargo, hemos decidido evaluar estos valores de  $D$  para seguir el trabajo de [53].

### 5.4.4.3. El parámetro $k$ de G-VGPMIL

Uno de los aspectos más interesantes de G-VGPMIL es el parámetro  $k$ . Con el objetivo de analizar cómo influye este parámetro en el rendimiento del modelo, vamos a estudiar cómo influye en la convergencia de las distribuciones variacionales y su relación con los resultados obtenidos.

Para ello hemos representado en la Figura 5.12 cómo evoluciona la métrica Log Loss con el número de épocas. Una época es una actualización completa de todos los parámetros, y en ella se usa todo el conjunto de entrenamiento. Analizamos los resultados correspondientes a tomar  $D = 128$ ,  $\text{lengthscale} = \sqrt{D}$  y 200 inducing points. Usando los datos del Apéndice A, el lector

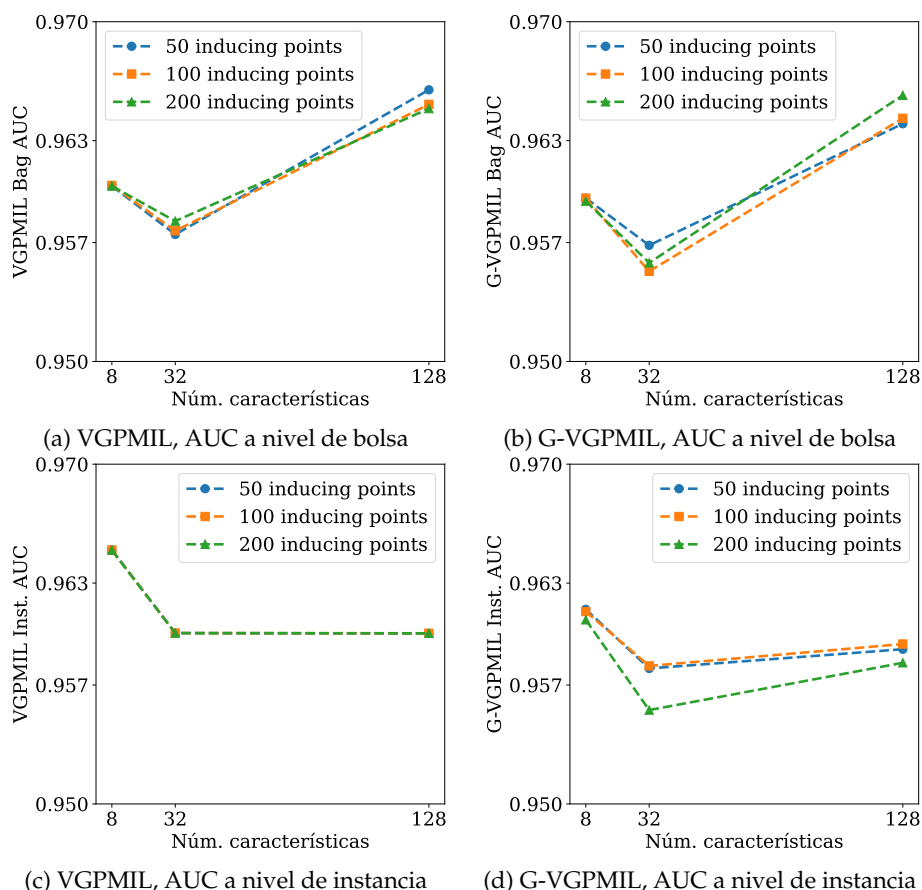


Figura 5.10.: Influencia del número de características en el rendimiento. Se muestra la evolución de la métrica AUC conforme aumenta el valor de  $D$ . Resultados obtenidos a partir del conjunto RSNA.

puede comprobar que para otras configuraciones de estos hiperparámetros el comportamiento es similar. Cuando  $k = 1$  la función de pérdida se estanca en un valor alto, pero conforme aumenta  $k$  el comportamiento mejora. En la subfigura 5.12a puede comprobarse que cuando  $k = 6$  alrededor de la sexta época ya se ha obtenido el mínimo, que se encuentra por debajo del que conseguirá VGPMIL. Sin embargo, las gráficas de test, subfiguras 5.12b y 5.12c, muestran que ocurre cierto sobreajuste y que este es más acusado cuanto mayor es  $k$ . Esto convierte a G-VGPMIL en un modelo cuyo proceso de entrenamiento es mucho más eficiente que VGPMIL si introducimos mecanismos para evitar el sobreaprendizaje. Un simple *early stopping* bastaría para solucionar este problema.

Para comprobar la relación de este hecho con la métrica AUC, hemos fijado el número de características en 128 y en la Figura 5.12 hemos recogido las gráficas de las métricas AUC y Log Loss frente al valor de  $k$  con el que se han obtenido. Los resultados en ambos conjuntos muestran que cuanto mayor es el parámetro  $k$  mejor rendimiento se obtiene. Las subfiguras 5.12a, 5.12b, 5.12d y 5.12e muestran que en el conjunto de RSNA la métrica AUC empeora a nivel de bolsa cuando  $k = 4$ . Sin embargo, Log Loss mejora considerablemente a ambos niveles. Las subfiguras 5.12c y 5.12f apoyan esta idea ya que la gráfica de Log Loss es estrictamente

## 5. Aplicaciones

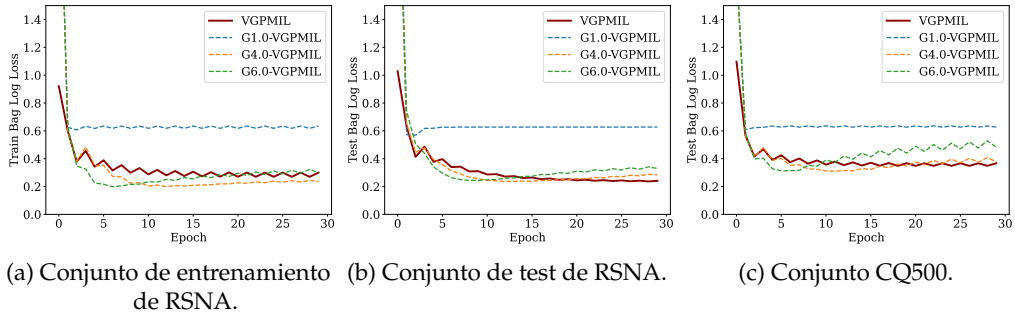


Figura 5.11.: Evolución de la métrica Log Loss a lo largo del proceso de entrenamiento. Resultados obtenidos tomando  $D = 128$  y 200 inducing points.

decreciente, mientras que la de AUC es estrictamente creciente.

Recordemos que modificando  $k$  podíamos alterar la forma de la función  $\theta$  (ecuaciones (4.12) y (4.13)) a partir de la cual se construye la matriz diagonal  $\Theta$ . Cuando  $k = 4$ , las funciones  $\theta$  y  $2\lambda$  presentaban el mismo máximo absoluto. El comportamiento que observamos ya se había dejado ver en las tablas de MNIST y musk, y justifica la sustitución de las variables Pólya-Gamma. Estos resultados convierten a G-VGPMIL en un potente modelo, que destaca por su flexibilidad frente a VGPMIL.

### 5.4.4.4. Comparación con el estado del arte

La [Tabla 5.7](#) compara los modelos estudiados en este trabajo con otras aproximaciones de la literatura. En todas ellas se entrenó un modelo sobre un conjunto de datos (en algunas también RSNA) para después evaluarlo sobre CQ500. Observamos que todas las propuestas emplean redes convolucionales, pero sólo una se limita a usar etiquetas a nivel de bolsa. En [?] emplean etiquetas a nivel de instancia sin imponer ninguna relación entre los objetos de una bolsa. En [29], sin embargo, usan el mecanismo de las capas LSTM para introducir relaciones entre las instancias contiguas. Ambos planteamiento obtienen una métrica AUC superior a nuestros clasificadores, pero no debemos olvidar que requieren un esfuerzo de anotación mucho mayor. Las propuestas basadas en MIL obtienen métricas muy competitivas (obsérvese que AttCNN + G-VGPMIL se encuentra a menos de 4 centesimas de [29]), usando sólo etiquetas a nivel de bolsa. Las aproximaciones presentadas en este trabajo, y en particular G-VGPMIL, son muy superiores al resto de trabajos que usan etiquetas a nivel de bolsa, como es [28].

Modelo	Tipo de etiquetas	CQ500 bag AUC
2D CNNs [8]	Instancia	0.94
2D CNN + LSTM [29]	Instancia	0.96
Segmentación con CNN [28]	Bolsa	0.83
AttCNN	Bolsa	0.9061
AttCNN + VGPMIL	Bolsa	0.922
AttCNN + G-VGPMIL	Bolsa	0.9235

Tabla 5.7.: Comparación con otras aproximaciones de la literatura.

#### 5.4. Detección de hemorragias intracraneales

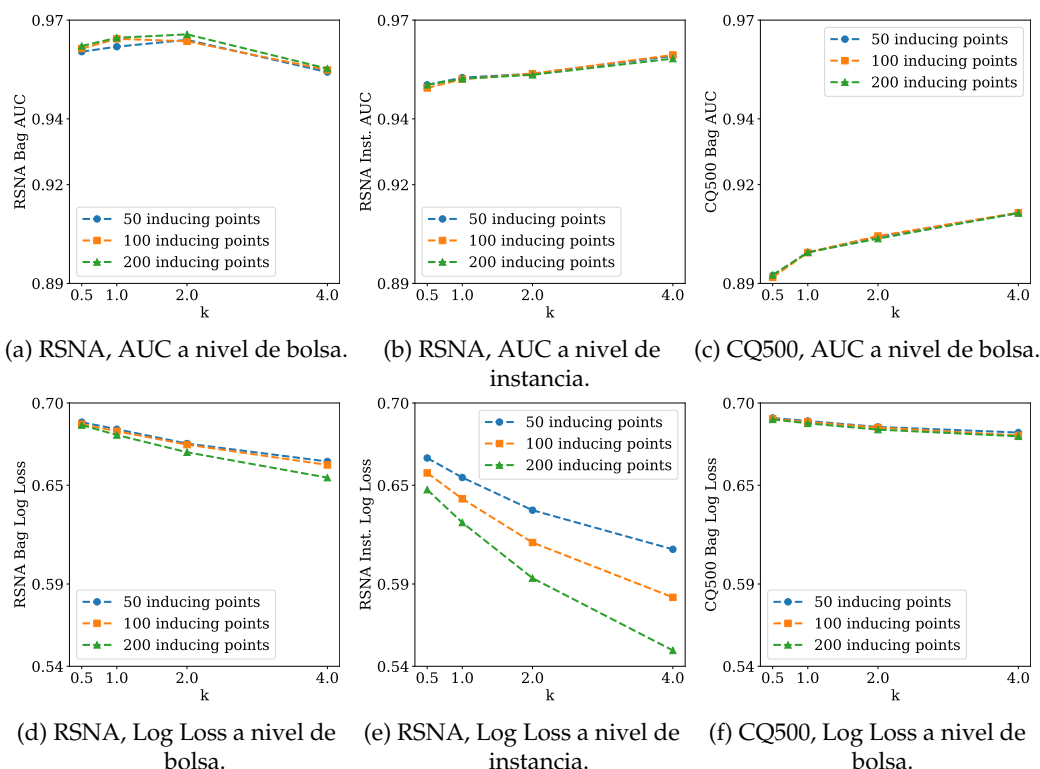


Figura 5.12.: Métricas AUC y Log Loss, tanto a nivel de bolsa como a nivel de instancia, en los conjuntos de RSNA y CQ500, frente al parámetro  $k$ . Resultados obtenidos tomando  $D = 128$ .

Para terminar, ofrecemos en la [Figura 5.13](#) un ejemplo de cómo funciona el sistema CAD desarrollado. Dado un escáner completo de un paciente, el modelo AttCNN + G-VGPMIL es capaz de estimar si el paciente sufre o no una hemorragia intracraneal. Esta estimación se da en forma de probabilidad, tanto a nivel de bolsa completa como a nivel de instancia. Además para cada probabilidad se incluye una cantidad que indica el intervalo de confianza de la predicción (obsérvese la ecuación (4.9)). Juntos cuantifican la incertidumbre que se tiene en la estimación ofrecida.

Ante la más que posible pregunta del radiólogo *¿por qué esa bolsa debe tener etiqueta positiva?* nuestro modelo sería capaz de señalar qué instancias le han llevado a tomar esa decisión (es decir, aquellas cuya etiqueta es positiva) y con qué probabilidad. Además, sería capaz de indicar qué seguridad se tiene en tal estimación, lo cual es esencial para advertir al experto sobre qué instancias o ejemplos debe revisar. Ninguno de los modelos basados exclusivamente en arquitecturas profundas, como los de la [Tabla 5.7](#), pueden satisfacer estos requerimientos como lo hacen los modelos probabilísticos.

## 5. Aplicaciones

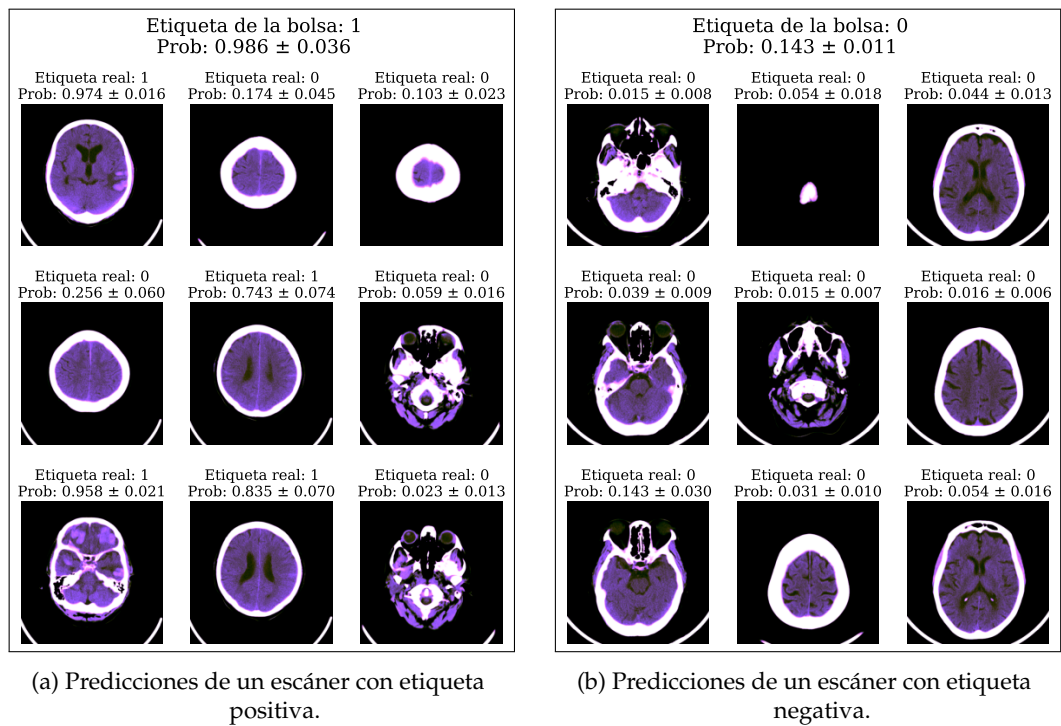


Figura 5.13.: Predicciones de dos escáneres completos del conjunto de RSNA.

## 6. Conclusiones y trabajo futuro

El objetivo de este capítulo es resumir las lecciones que hemos aprendido y el trabajo que estas puede motivar en el futuro. Para ello primero ofrecemos un resumen de los problemas que hemos tratado y las conclusiones a las que hemos llegado tras resolverlos. Posteriormente, señalamos dónde podemos centrar nuestros esfuerzos en el futuro.

### 6.1. Conclusiones

Comenzábamos este trabajo hablando sobre la importancia de una detección temprana y fiable de las hemorragias intracraneales. Justificábamos por qué era necesario elaborar métodos de triaje basados en Sistemas de Diagnóstico Asistido por Ordenador en los que los radiólogos pudieran apoyarse. La construcción de estos sistemas mediante las técnicas habituales de Aprendizaje Automático requiere disponer de bases de datos fuertemente anotadas. Esto es, los radiólogos deben etiquetar una gran cantidad de cortes de TAC, lo cual es una tarea inabordable debido al gran esfuerzo que requiere y al poco tiempo del que disponen.

Lo anterior motivaba el uso de técnicas de Aprendizaje a partir de Múltiples Instancias (*Multiple Instance Learning*, MIL), las cuales sólo requieren una etiqueta para cada escáner completo. Esta información es fácil de conseguir ya que se encuentra disponible en el historial clínico del paciente. Así es cómo surgía el primer objetivo de este trabajo, que era comprender los problemas de tipo MIL y revisar la literatura comprendiendo qué soluciones se han propuesto hasta el momento.

La incertidumbre presente en MIL nos sugería el uso de métodos probabilísticos capaces de cuantificarla. Esto justificó la elección de los Procesos Gaussianos (*Gaussian Processes*, GPs) como herramienta, cuya formulación probabilística permite adaptarlos para modelar la relación entre las instancias y las bolsas en los problemas de MIL. En este trabajo hemos estudiado los fundamentos de estos modelos bayesianos en problemas de aprendizaje supervisado clásicos, para después considerar las aproximaciones basadas en GPs para MIL existentes en la literatura. De todas ellas, nos hemos centrado en un modelo denominado *Variational Gaussian Process Multiple Instance Learning* (VGPMIL), cuyo proceso de inferencia se basa en una desigualdad que puede conducir a soluciones subóptimas.

En este proyecto nos proponíamos explorar técnicas de *data augmentation* para realizar este proceso de inferencia. Concretamente, inspirados por una propuesta previa que usa variables aleatorias Pólya-Gamma, nos preguntábamos qué ocurría al incluir este tipo de variables en el modelo de GPs para MIL. Así es cómo hemos llegado a las dos aportaciones teóricas de este trabajo:

- Hemos construido *Pólya-Gamma Variational Gaussian Process Multiple Instance Learning* (PG-VGPMIL) que propone usar variables Pólya-Gamma para aumentar la probabilidad conjunta y así expresar la función logística a través de ellas. Probamos que PG-VGPMIL es equivalente, en el sentido *Mean Field*, a VGPMIL. Esto significa que existe cierta relación entre la desigualdad de Jaakkola y este tipo de variables aleatorias.

## 6. Conclusiones y trabajo futuro

- La definición de las variables Pólya-Gamma nos sugiere sustituirlas por variables Gamma. Así hemos desarrollado *Gamma Variational Gaussian Process Multiple Instance Learning* (G-VGPMIL), que nos conduce a un nuevo modelo de observación. Hemos validado experimentalmente G-VGPMIL usando tres conjuntos diferentes. Los experimentos nos han llevado a concluir que es un modelo muy competitivo que en varias ocasiones se comporta mejor que VGPMIL.

La última parte del trabajo ofrece una solución al problema que lo motiva, la detección de hemorragias intracraneales. Esta solución, que se denomina AttCNN+G-VGPMIL se divide en dos fases. La primera se basa en el uso de redes neuronales y mecanismos de atención para extraer características. La segunda fase emplea G-VGPMIL sobre estas características para obtener predicciones a nivel de bolsa y de instancia. El modelo obtenido sólo necesita etiquetas a nivel de bolsa y ofrece los mejores resultados de todas las metodologías de la literatura de ese tipo, a la vez que se sitúa unas centésimas por debajo de modelos mucho más complejos que emplean etiquetas a nivel de instancia. Esto demuestra que usar metodologías de etiquetado débil es la vía a seguir para construir modelos que no dependan de conjuntos fuertemente anotados. Por otra parte, AttCNN+G-VGPMIL es capaz de informar sobre la incertidumbre de sus decisiones gracias al uso de GPs. Esto es una cualidad esencial dentro de la medicina.

### 6.2. Trabajo futuro

Los resultados de este Trabajo Fin de Máster se incluirán en un artículo que será enviado a una revista en el Journal Citation Report. Además, cada una de las aportaciones a la literatura que hemos realizado deja abierta una puerta a la investigación. Nos detenemos en cada una de ellas.

1. La equivalencia entre PG-VGPMIL y VGPMIL nos lleva a pensar en la relación existente entre la desigualdad de Jaakkola y las variables Pólya-Gamma. Pensamos que un hecho relacionado ya se apunta en las conclusiones de [51]. Aclarar esta cuestión puede suponer un gran avance a la hora de tratar el modelo de observación logístico.
2. ¿Qué ocurriría al cambiar las variables Gamma por otro tipo de variables aleatorias? La forma general de las distribuciones Pólya-Gamma y Gamma se definen como el producto por un factor exponencial de una versión simple de ellas mismas. Este hecho es el que nos permite calcular la distribución variacional de las nuevas variables. Explorar nuevas variables con esta propiedad se traduce en cambiar la función que define a la matriz diagonal  $\Theta$ .
3. El modelo AttCNN+G-VGPMIL obtiene grandes resultados, pero tiene un gran inconveniente: el proceso de entrenamiento no es *end-to-end*, sino que debe realizarse en dos fases. Por tanto, las características extraídas en la primera fase pueden no ser las mejores para G-VGPMIL. Una forma de hacer esto podría ser entrenar G-VGPMIL mediante gradiente descendente y usar alguna herramienta basada en diferenciación automática para guiar el proceso de entrenamiento por toda la red.





# A. Tablas

Num. Inducing	Lengthscale	k	Model	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	1.0	0.5	G-VGPML	<b>0.9383 ± 0.0089</b>	<b>0.5541 ± 0.0089</b>	0.6874 ± 0.0089	<b>0.9038 ± 0.0089</b>	<b>0.6855 ± 0.0089</b>	<b>0.6546 ± 0.0089</b>	0.6931 ± 0.0089	<b>0.8025 ± 0.0089</b>
		1.0	G-VGPML	0.9382 ± 0.0089	0.5531 ± 0.0089	0.6851 ± 0.0089	0.9036 ± 0.0089	0.685 ± 0.0089	0.6538 ± 0.0089	0.6931 ± 0.0089	0.802 ± 0.0089
		2.0	G-VGPML	0.938 ± 0.0088	0.5507 ± 0.0088	0.6815 ± 0.0088	0.9032 ± 0.0088	0.6837 ± 0.0088	0.6519 ± 0.0088	0.6931 ± 0.0088	0.8009 ± 0.0088
		4.0	G-VGPML	0.9378 ± 0.0089	0.5481 ± 0.0089	<b>0.6764 ± 0.0089</b>	0.9026 ± 0.0089	0.6821 ± 0.0089	0.6492 ± 0.0089	<b>0.6931 ± 0.0089</b>	0.7992 ± 0.0089
	-	VGPML	0.9372 ± 0.0087	0.5422 ± 0.0087	0.6765 ± 0.0087	0.9019 ± 0.0087	0.677 ± 0.0087	0.6423 ± 0.0087	0.6931 ± 0.0087	0.7958 ± 0.0087	
	3.238	0.5	G-VGPML	0.9759 ± 0.0039	0.8623 ± 0.0039	0.6196 ± 0.0039	0.988 ± 0.0039	0.8919 ± 0.0039	0.9032 ± 0.0039	0.6476 ± 0.0039	0.9739 ± 0.0039
		1.0	G-VGPML	0.9824 ± 0.0011	0.906 ± 0.0011	0.5724 ± 0.0011	0.9906 ± 0.0011	0.9262 ± 0.0011	0.937 ± 0.0011	0.5994 ± 0.0011	0.9785 ± 0.0011
		2.0	G-VGPML	0.9827 ± 0.002	0.9129 ± 0.002	0.4482 ± 0.002	0.9915 ± 0.002	0.9347 ± 0.002	0.9465 ± 0.002	0.447 ± 0.002	0.9807 ± 0.002
		4.0	G-VGPML	0.9842 ± 0.0015	0.9144 ± 0.0015	<b>0.0657 ± 0.0015</b>	0.9921 ± 0.0015	0.9331 ± 0.0015	0.9426 ± 0.0015	<b>0.1901 ± 0.0015</b>	0.9807 ± 0.0015
	-	VGPML	<b>0.9862 ± 0.0008</b>	<b>0.9287 ± 0.0008</b>	0.1393 ± 0.0008	<b>0.9933 ± 0.0008</b>	<b>0.9466 ± 0.0008</b>	<b>0.9557 ± 0.0008</b>	0.2309 ± 0.0008	<b>0.9839 ± 0.0008</b>	
	5.477	0.5	G-VGPML	0.9638 ± 0.0055	0.7835 ± 0.0055	0.6392 ± 0.0055	0.9854 ± 0.0055	0.8388 ± 0.0055	0.8491 ± 0.0055	0.6678 ± 0.0055	0.9667 ± 0.0055
		1.0	G-VGPML	0.9764 ± 0.0059	0.8683 ± 0.0059	0.5911 ± 0.0059	0.9873 ± 0.0059	0.9006 ± 0.0059	0.9125 ± 0.0059	0.6313 ± 0.0059	0.9715 ± 0.0059
2.0		G-VGPML	0.981 ± 0.004	0.9013 ± 0.004	0.4646 ± 0.004	0.9886 ± 0.004	0.927 ± 0.004	0.9391 ± 0.004	0.4977 ± 0.004	0.9756 ± 0.004	
4.0		G-VGPML	0.9824 ± 0.0016	0.905 ± 0.0016	<b>0.0643 ± 0.0016</b>	0.9908 ± 0.0016	0.9285 ± 0.0016	0.9387 ± 0.0016	<b>0.2047 ± 0.0016</b>	0.9795 ± 0.0016	
-	VGPML	<b>0.9855 ± 0.0013</b>	<b>0.9257 ± 0.0013</b>	0.1414 ± 0.0013	<b>0.993 ± 0.0013</b>	<b>0.9445 ± 0.0013</b>	<b>0.954 ± 0.0013</b>	0.2353 ± 0.0013	<b>0.9846 ± 0.0013</b>		
17.738	0.5	G-VGPML	0.9449 ± 0.0015	0.6124 ± 0.0015	0.6484 ± 0.0015	<b>0.9923 ± 0.0015</b>	0.72 ± 0.0015	0.7017 ± 0.0015	0.6788 ± 0.0015	<b>0.9835 ± 0.0015</b>	
	1.0	G-VGPML	0.9684 ± 0.0013	0.8098 ± 0.0013	0.6068 ± 0.0013	0.9914 ± 0.0013	0.8538 ± 0.0013	0.864 ± 0.0013	0.6543 ± 0.0013	0.9816 ± 0.0013	
	2.0	G-VGPML	<b>0.9815 ± 0.0021</b>	<b>0.903 ± 0.0021</b>	0.522 ± 0.0021	0.9891 ± 0.0021	<b>0.9277 ± 0.0021</b>	0.9302 ± 0.0021	0.5685 ± 0.0021	0.9768 ± 0.0021	
	4.0	G-VGPML	0.9811 ± 0.0019	0.9006 ± 0.0019	<b>0.0781 ± 0.0019</b>	0.9893 ± 0.0019	0.9259 ± 0.0019	0.9374 ± 0.0019	<b>0.2177 ± 0.0019</b>	0.9766 ± 0.0019	
-	VGPML	0.9803 ± 0.0019	0.8994 ± 0.0019	0.2043 ± 0.0019	0.9884 ± 0.0019	0.9274 ± 0.0019	<b>0.9398 ± 0.0019</b>	0.3175 ± 0.0019	0.9756 ± 0.0019		
30.0	0.5	G-VGPML	0.9392 ± 0.0017	0.5546 ± 0.0017	0.649 ± 0.0017	<b>0.9905 ± 0.0017</b>	0.6835 ± 0.0017	0.6493 ± 0.0017	0.6805 ± 0.0017	<b>0.9795 ± 0.0017</b>	
	1.0	G-VGPML	0.9649 ± 0.0024	0.7837 ± 0.0024	0.6075 ± 0.0024	0.9892 ± 0.0024	0.8357 ± 0.0024	0.8445 ± 0.0024	0.6575 ± 0.0024	0.9765 ± 0.0024	
	2.0	G-VGPML	0.9791 ± 0.002	0.8898 ± 0.002	0.5194 ± 0.002	0.9869 ± 0.002	0.9185 ± 0.002	0.9309 ± 0.002	0.5743 ± 0.002	0.9713 ± 0.002	
	4.0	G-VGPML	0.979 ± 0.0029	0.8872 ± 0.0029	<b>0.1129 ± 0.0029</b>	0.9891 ± 0.0029	0.9171 ± 0.0029	0.9287 ± 0.0029	<b>0.2697 ± 0.0029</b>	0.976 ± 0.0029	
-	VGPML	<b>0.9798 ± 0.0006</b>	<b>0.895 ± 0.0006</b>	0.2367 ± 0.0006	0.9876 ± 0.0006	<b>0.9236 ± 0.0006</b>	<b>0.936 ± 0.0006</b>	0.3577 ± 0.0006	0.9734 ± 0.0006		
1.0	0.5	G-VGPML	0.9592 ± 0.0045	0.7378 ± 0.0045	0.6869 ± 0.0045	0.9408 ± 0.0045	0.8057 ± 0.0045	0.8094 ± 0.0045	0.6931 ± 0.0045	<b>0.8802 ± 0.0045</b>	
	1.0	G-VGPML	<b>0.9592 ± 0.0045</b>	<b>0.738 ± 0.0045</b>	0.6845 ± 0.0045	0.9408 ± 0.0045	<b>0.8059 ± 0.0045</b>	<b>0.8095 ± 0.0045</b>	0.6931 ± 0.0045	0.8801 ± 0.0045	
	2.0	G-VGPML	0.9591 ± 0.0046	0.7375 ± 0.0046	0.6806 ± 0.0046	<b>0.9409 ± 0.0046</b>	0.8059 ± 0.0046	0.8095 ± 0.0046	0.6931 ± 0.0046	0.8798 ± 0.0046	
	4.0	G-VGPML	0.959 ± 0.0046	0.7368 ± 0.0046	<b>0.6751 ± 0.0046</b>	0.9407 ± 0.0046	0.8051 ± 0.0046	0.8088 ± 0.0046	<b>0.6931 ± 0.0046</b>	0.8787 ± 0.0046	
-	VGPML	0.959 ± 0.0046	0.7368 ± 0.0046	0.6754 ± 0.0046	0.9407 ± 0.0046	0.8051 ± 0.0046	0.8088 ± 0.0046	0.6931 ± 0.0046	0.8788 ± 0.0046		
3.238	0.5	G-VGPML	0.9783 ± 0.0004	0.8771 ± 0.0004	0.6215 ± 0.0004	0.9927 ± 0.0004	0.9049 ± 0.0004	0.915 ± 0.0004	0.6494 ± 0.0004	0.9859 ± 0.0004	
	1.0	G-VGPML	0.987 ± 0.0009	0.9308 ± 0.0009	0.5739 ± 0.0009	0.9948 ± 0.0009	0.9456 ± 0.0009	0.9537 ± 0.0009	0.5955 ± 0.0009	0.989 ± 0.0009	
	2.0	G-VGPML	0.9881 ± 0.0011	0.9397 ± 0.0011	0.4391 ± 0.0011	0.9954 ± 0.0011	0.9581 ± 0.0011	0.9656 ± 0.0011	0.4193 ± 0.0011	0.9895 ± 0.0011	
	4.0	G-VGPML	0.9875 ± 0.0009	0.9333 ± 0.0009	<b>0.0509 ± 0.0009</b>	0.9961 ± 0.0009	0.9469 ± 0.0009	0.9547 ± 0.0009	<b>0.1598 ± 0.0009</b>	0.9902 ± 0.0009	
-	VGPML	<b>0.9909 ± 0.0009</b>	<b>0.9527 ± 0.0009</b>	0.1102 ± 0.0009	<b>0.9968 ± 0.0009</b>	<b>0.9645 ± 0.0009</b>	<b>0.9705 ± 0.0009</b>	0.1833 ± 0.0009	<b>0.9921 ± 0.0009</b>		
100	0.5	G-VGPML	0.9695 ± 0.0042	0.8181 ± 0.0042	0.638 ± 0.0042	0.9639 ± 0.0042	0.8592 ± 0.0042	0.8699 ± 0.0042	0.6641 ± 0.0042	0.9875 ± 0.0042	
	1.0	G-VGPML	0.9844 ± 0.0014	0.9156 ± 0.0014	0.5888 ± 0.0014	0.9947 ± 0.0014	0.9351 ± 0.0014	0.9442 ± 0.0014	0.6146 ± 0.0014	0.9887 ± 0.0014	
	2.0	G-VGPML	0.9882 ± 0.0009	0.9398 ± 0.0009	0.4544 ± 0.0009	0.9951 ± 0.0009	0.9558 ± 0.0009	0.9635 ± 0.0009	0.4555 ± 0.0009	0.9898 ± 0.0009	
	4.0	G-VGPML	0.9867 ± 0.0004	0.9282 ± 0.0004	<b>0.0476 ± 0.0004</b>	0.9958 ± 0.0004	0.9438 ± 0.0004	0.952 ± 0.0004	<b>0.1607 ± 0.0004</b>	0.9901 ± 0.0004	
-	VGPML	<b>0.9901 ± 0.0008</b>	<b>0.9493 ± 0.0008</b>	0.115 ± 0.0008	<b>0.9965 ± 0.0008</b>	<b>0.9635 ± 0.0008</b>	<b>0.9697 ± 0.0008</b>	0.189 ± 0.0008	<b>0.9917 ± 0.0008</b>		
17.738	0.5	G-VGPML	0.9535 ± 0.0026	0.6916 ± 0.0026	0.6479 ± 0.0026	<b>0.9954 ± 0.0026</b>	0.7704 ± 0.0026	0.7677 ± 0.0026	0.6763 ± 0.0026	<b>0.99 ± 0.0026</b>	
	1.0	G-VGPML	0.9735 ± 0.0016	0.8456 ± 0.0016	0.6058 ± 0.0016	0.9945 ± 0.0016	0.8784 ± 0.0016	0.8893 ± 0.0016	0.6498 ± 0.0016	0.9881 ± 0.0016	
	2.0	G-VGPML	<b>0.9834 ± 0.0013</b>	<b>0.9136 ± 0.0013</b>	0.5124 ± 0.0013	0.9913 ± 0.0013	<b>0.9346 ± 0.0013</b>	<b>0.9452 ± 0.0013</b>	0.5498 ± 0.0013	0.9811 ± 0.0013	
	4.0	G-VGPML	0.9803 ± 0.0004	0.8935 ± 0.0004	<b>0.0719 ± 0.0004</b>	0.9886 ± 0.0004	0.9191 ± 0.0004	0.9304 ± 0.0004	<b>0.2294 ± 0.0004</b>	0.9748 ± 0.0004	
-	VGPML	0.9817 ± 0.0007	0.9059 ± 0.0007	0.1955 ± 0.0007	0.9895 ± 0.0007	0.9317 ± 0.0007	0.9433 ± 0.0007	0.3073 ± 0.0007	0.9779 ± 0.0007		
30.0	0.5	G-VGPML	0.947 ± 0.0079	0.629 ± 0.0079	0.649 ± 0.0079	<b>0.9921 ± 0.0079</b>	0.7305 ± 0.0079	0.7137 ± 0.0079	0.6788 ± 0.0079	<b>0.9807 ± 0.0079</b>	
	1.0	G-VGPML	0.9677 ± 0.0045	0.8054 ± 0.0045	0.6071 ± 0.0045	0.9906 ± 0.0045	0.8511 ± 0.0045	0.8613 ± 0.0045	0.6513 ± 0.0045	0.9777 ± 0.0045	
	2.0	G-VGPML	0.98 ± 0.0015	0.895 ± 0.0015	0.5211 ± 0.0015	0.9888 ± 0.0015	0.9223 ± 0.0015	0.9346 ± 0.0015	0.5767 ± 0.0015	0.9747 ± 0.0015	
	4.0	G-VGPML	<b>0.9801 ± 0.002</b>	0.8943 ± 0.002	<b>0.112 ± 0.002</b>	0.9896 ± 0.002	0.9218 ± 0.002	0.9333 ± 0.002	<b>0.2619 ± 0.002</b>	0.9767 ± 0.002	
-	VGPML	0.9799 ± 0.0004	<b>0.8957 ± 0.0004</b>	0.2353 ± 0.0004	0.9877 ± 0.0004	<b>0.9245 ± 0.0004</b>	<b>0.9367 ± 0.0004</b>	0.3562 ± 0.0004	0.9738 ± 0.0004		
1.0	0.5	G-VGPML	0.9757 ± 0.0018	0.8604 ± 0.0018	0.6848 ± 0.0018	<b>0.9664 ± 0.0018</b>	<b>0.8916 ± 0.0018</b>	<b>0.9027 ± 0.0018</b>	0.6931 ± 0.0018	<b>0.9346 ± 0.0018</b>	
	1.0	G-VGPML	0.9757 ± 0.0018	0.8604 ± 0.0018	0.6816 ± 0.0018	0.9664 ± 0.0018	0.8914 ± 0.0018	0.9025 ± 0.0018	0.6931 ± 0.0018	0.9339 ± 0.0018	
	2.0	G-VGPML	<b>0.9757 ± 0.0017</b>	<b>0.8606 ± 0.0017</b>	0.6765 ± 0.0017	0.9663 ± 0.0017	0.8911 ± 0.0017	0.9023 ± 0.0017	0.6931 ± 0.0017	0.9333 ± 0.0017	
	4.0	G-VGPML	0.9756 ± 0.0017	0.86 ± 0.0017	<b>0.6691 ± 0.0017</b>	0.9662 ± 0.0017	0.8906 ± 0.0017	0.9019 ± 0.0017	<b>0.6931 ± 0.0017</b>	0.9329 ± 0.0017	
-	VGPML	0.9756 ± 0.0017	0.8601 ± 0.0017	0.6696 ± 0.0017	0.9662 ± 0.0017	0.8907 ± 0.0017	0.902 ± 0.0017	0.6931 ± 0.0017	0.933 ± 0.0017		
3.238	0.5	G-VGPML	0.981 ± 0.0017	0.8945 ± 0.0017	0.6216 ± 0.0017	0.9944 ± 0.0017	0.9148 ± 0.0017	0.9253 ± 0.0017	0.6438 ± 0.0017	0.9882 ± 0.0017	
	1.0	G-VGPML	0.9892 ± 0.0007	0.9428 ± 0.0007	0.5687 ± 0.0007	0.9961 ± 0.0007	0.9537 ± 0.0007	0.961 ± 0.0007	0.5807 ± 0.0007	0.9913 ± 0.0007	
	2.0	G-VGPML	0.9897 ± 0.0012	0.9478 ± 0.0012	0.4124 ± 0.0012	0.9963 ± 0.0012	0.9627 ± 0.0012	0.9695 ± 0.0012	0.374 ± 0.0012	0.9913 ± 0.0012	
	4.0	G-VGPML	0.9897 ± 0.0005	0.9453 ± 0.0005	<b>0.0429 ± 0.0005</b>	0.9974 ± 0.0005	0.9576 ± 0.0005	0.964 ± 0.0005	<b>0.1379 ± 0.0005</b>	0.9934 ± 0.0005	
-	VGPML	<b>0.9924 ± 0.0006</b>	<b>0.9607 ± 0.0006</b>	0.0957 ± 0.0006	<b>0.9978 ± 0.0006</b>	<b>0.9691 ± 0.0006</b>	<b>0.9744 ± 0.0006</b>	0.1597 ± 0.0006	<b>0.9944 ± 0.0006</b>		
200	0.5	G-VGPML	0.9747 ± 0.004	0.853 ± 0.004	0.6394 ± 0.004	0.9957 ± 0.004	0.8838 ± 0.004	0.8944 ± 0.004	0.6633 ± 0.004	0.9901 ± 0.004	
	1.0	G-VGPML	0.9876 ± 0.0014	0.9341 ± 0.0014	0.5898 ± 0.0014	0.9965 ± 0.0014	0.949 ± 0.0014	0.9565 ± 0.0014	0.6124 ± 0.0014	0.9916 ± 0.0014	
	2.0	G-VGPML	0.9907 ± 0.0009	0.9527 ± 0.0009	0.4481 ± 0.0009	0.9967 ± 0.0009	0.9653 ± 0.0009	0.9714 ± 0.0009	0.4349 ± 0.0009	0.9924 ± 0.0009	
	4.0	G-VGPML	0.9889 ± 0.0004	0.941 ± 0.0004							

Num. Inducing	Lengthscale	k	Modelo	Inst. Acc.	Inst. F1	Inst. Log Loss	Inst. AUC	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC	
50	1.0	0.5	G-VGPML	<b>0.9014 ± 0.0026</b>	<b>0.0 ± 0.0</b>	<b>0.6931 ± 0.0</b>	<b>0.5 ± 0.0</b>	<b>0.3905 ± 0.0003</b>	<b>0.0 ± 0.0</b>	<b>0.6931 ± 0.0</b>	<b>0.5 ± 0.0</b>	
		1.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		2.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		4.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
	14.5	0.5	G-VGPML	0.9173 ± 0.0151	0.2514 ± 0.2087	0.6365 ± 0.0008	0.83 ± 0.0312	0.5184 ± 0.1146	0.3138 ± 0.231	0.6915 ± 0.0025	0.6548 ± 0.082	
		1.0	G-VGPML	0.9438 ± 0.0116	0.5992 ± 0.1227	0.6069 ± 0.0016	0.8863 ± 0.0225	0.7155 ± 0.0809	0.6899 ± 0.1141	0.682 ± 0.0064	0.7903 ± 0.06	
		2.0	G-VGPML	<b>0.9492 ± 0.0046</b>	<b>0.7408 ± 0.0237</b>	0.5513 ± 0.0105	0.9361 ± 0.011	<b>0.821 ± 0.012</b>	<b>0.8518 ± 0.0156</b>	0.64 ± 0.0195	0.8964 ± 0.0247	
		4.0	G-VGPML	0.9238 ± 0.0144	0.6948 ± 0.0464	<b>0.3237 ± 0.076</b>	<b>0.9564 ± 0.0054</b>	0.7895 ± 0.0263	0.8474 ± 0.0159	<b>0.4712 ± 0.0574</b>	<b>0.9253 ± 0.0112</b>	
	28.0	0.5	G-VGPML	0.9314 ± 0.0087	0.4651 ± 0.0913	0.6399 ± 0.0005	0.9345 ± 0.0114	0.6291 ± 0.055	0.5589 ± 0.0894	0.6855 ± 0.0022	0.8343 ± 0.0255	
		1.0	G-VGPML	0.9532 ± 0.0072	0.6949 ± 0.0618	0.6026 ± 0.0025	0.9517 ± 0.0074	0.7744 ± 0.0409	0.7745 ± 0.0511	0.6675 ± 0.0057	0.8879 ± 0.0187	
		2.0	G-VGPML	<b>0.9621 ± 0.0019</b>	<b>0.8016 ± 0.0104</b>	0.5295 ± 0.0195	0.9652 ± 0.0032	<b>0.856 ± 0.0086</b>	0.8788 ± 0.008	0.6055 ± 0.0135	0.9267 ± 0.0084	
		4.0	G-VGPML	0.9569 ± 0.0077	0.793 ± 0.0261	<b>0.1663 ± 0.0384</b>	<b>0.9682 ± 0.0021</b>	0.8497 ± 0.0188	<b>0.8798 ± 0.0128</b>	<b>0.3486 ± 0.0263</b>	<b>0.9341 ± 0.0098</b>	
406.0	0.5	G-VGPML	0.9221 ± 0.0022	0.3471 ± 0.0269	0.6504 ± 0.0001	<b>0.976 ± 0.0017</b>	0.5624 ± 0.015	0.4397 ± 0.0303	0.6855 ± 0.0005	<b>0.9394 ± 0.007</b>		
	1.0	G-VGPML	0.9368 ± 0.0018	0.5295 ± 0.0216	0.6121 ± 0.0003	0.9753 ± 0.0013	0.6669 ± 0.0139	0.6244 ± 0.0218	0.6732 ± 0.0011	0.9381 ± 0.0066		
	2.0	G-VGPML	0.9545 ± 0.0022	0.7065 ± 0.0149	0.5479 ± 0.0004	0.9743 ± 0.0012	0.7832 ± 0.0116	0.7857 ± 0.0136	0.6385 ± 0.0024	0.936 ± 0.0055		
	4.0	G-VGPML	<b>0.9623 ± 0.0014</b>	<b>0.7834 ± 0.0076</b>	<b>0.4149 ± 0.0065</b>	0.9723 ± 0.0025	<b>0.8411 ± 0.0059</b>	<b>0.8569 ± 0.0066</b>	<b>0.5484 ± 0.0054</b>	<b>0.9317 ± 0.0031</b>		
784.0	0.5	G-VGPML	0.9091 ± 0.0021	0.145 ± 0.0171	0.6514 ± 0.0002	<b>0.9743 ± 0.0015</b>	0.458 ± 0.0089	0.1991 ± 0.0231	0.6884 ± 0.0005	<b>0.9356 ± 0.0044</b>		
	1.0	G-VGPML	0.9099 ± 0.0021	0.1588 ± 0.0287	0.6156 ± 0.0003	0.9729 ± 0.0025	0.4632 ± 0.0156	0.2175 ± 0.0392	0.6847 ± 0.0011	0.9323 ± 0.0028		
	2.0	G-VGPML	0.9104 ± 0.0027	0.1665 ± 0.0449	0.5591 ± 0.0005	0.9711 ± 0.0037	0.4695 ± 0.0242	0.2276 ± 0.0591	0.6799 ± 0.0026	0.9282 ± 0.0022		
	4.0	G-VGPML	<b>0.9134 ± 0.0042</b>	<b>0.2136 ± 0.0699</b>	<b>0.4921 ± 0.0013</b>	0.9692 ± 0.0046	<b>0.4934 ± 0.0373</b>	<b>0.2847 ± 0.0833</b>	<b>0.6733 ± 0.0056</b>	0.9236 ± 0.0029		
100	1.0	-	VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		14.5	-	VGPML	0.929 ± 0.007	0.7033 ± 0.0277	0.3866 ± 0.0307	0.9531 ± 0.0058	0.7989 ± 0.0152	0.8516 ± 0.0101	0.5217 ± 0.0357	0.9216 ± 0.0124
		28.0	-	VGPML	0.958 ± 0.0025	<b>0.7959 ± 0.0132</b>	<b>0.2749 ± 0.0101</b>	0.9679 ± 0.002	<b>0.8557 ± 0.0103</b>	<b>0.8842 ± 0.0084</b>	<b>0.426 ± 0.0161</b>	<b>0.9338 ± 0.007</b>
		406.0	-	VGPML	<b>0.9607 ± 0.0018</b>	0.7682 ± 0.0074	0.4063 ± 0.002	<b>0.9723 ± 0.0025</b>	0.829 ± 0.0033	0.8427 ± 0.0035	0.5526 ± 0.0053	0.9316 ± 0.0031
	784.0	-	VGPML	<b>0.9137 ± 0.0041</b>	<b>0.2183 ± 0.0687</b>	<b>0.4928 ± 0.0008</b>	<b>0.9691 ± 0.0046</b>	<b>0.4958 ± 0.0369</b>	<b>0.2904 ± 0.082</b>	<b>0.6729 ± 0.0054</b>	<b>0.9235 ± 0.0029</b>	
	1.0	0.5	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		1.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		2.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		4.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
	14.5	0.5	G-VGPML	0.948 ± 0.0053	0.645 ± 0.0523	0.6301 ± 0.001	0.9156 ± 0.0093	0.743 ± 0.0359	0.7332 ± 0.0485	0.6806 ± 0.0033	0.8255 ± 0.0206	
		1.0	G-VGPML	<b>0.9618 ± 0.0022</b>	<b>0.7783 ± 0.0169</b>	0.5958 ± 0.0035	0.9415 ± 0.0087	0.8376 ± 0.0145	0.8528 ± 0.0169	0.6559 ± 0.0092	0.8916 ± 0.0104	
		2.0	G-VGPML	0.9515 ± 0.0035	0.7747 ± 0.0167	<b>0.5125 ± 0.0169</b>	<b>0.9617 ± 0.0082</b>	<b>0.8406 ± 0.0121</b>	<b>0.8754 ± 0.009</b>	<b>0.5702 ± 0.0216</b>	<b>0.9327 ± 0.007</b>	
4.0		G-VGPML	0.9591 ± 0.0084	0.7988 ± 0.0261	<b>0.1784 ± 0.0218</b>	0.9559 ± 0.0127	0.853 ± 0.0195	0.8807 ± 0.0113	<b>0.3653 ± 0.0113</b>	0.9258 ± 0.0126		
28.0	0.5	G-VGPML	0.9363 ± 0.0046	0.5236 ± 0.0499	0.6375 ± 0.0004	0.9611 ± 0.0031	0.6644 ± 0.0286	0.6203 ± 0.0451	0.6817 ± 0.0019	0.8881 ± 0.0171		
	1.0	G-VGPML	0.9579 ± 0.0025	0.7324 ± 0.0184	0.5997 ± 0.0008	0.9704 ± 0.0018	0.799 ± 0.0182	0.8041 ± 0.0208	0.6606 ± 0.0036	0.9218 ± 0.0104		
	2.0	G-VGPML	<b>0.9691 ± 0.0019</b>	<b>0.8392 ± 0.0126</b>	0.5266 ± 0.0125	<b>0.9764 ± 0.0037</b>	<b>0.8789 ± 0.0094</b>	<b>0.8988 ± 0.0078</b>	0.5862 ± 0.0158	<b>0.9456 ± 0.0067</b>		
	4.0	G-VGPML	<b>0.969 ± 0.0041</b>	<b>0.8436 ± 0.0165</b>	<b>0.1214 ± 0.0129</b>	0.9786 ± 0.003	<b>0.8836 ± 0.0108</b>	<b>0.9045 ± 0.008</b>	<b>0.2943 ± 0.017</b>	0.9503 ± 0.0053		
406.0	0.5	G-VGPML	0.9255 ± 0.0019	0.3932 ± 0.0119	0.6499 ± 0.0002	<b>0.9863 ± 0.0011</b>	0.5872 ± 0.0073	0.488 ± 0.0134	0.6838 ± 0.0003	0.9554 ± 0.0041		
	1.0	G-VGPML	0.9417 ± 0.0015	0.5813 ± 0.0211	0.6112 ± 0.0004	0.9837 ± 0.0015	0.6999 ± 0.0146	0.6739 ± 0.0211	0.6698 ± 0.001	<b>0.9584 ± 0.0036</b>		
	2.0	G-VGPML	0.9588 ± 0.0015	0.7418 ± 0.0093	0.5416 ± 0.0065	0.9783 ± 0.003	0.8082 ± 0.0098	0.815 ± 0.0113	0.6292 ± 0.0061	0.9457 ± 0.0059		
	4.0	G-VGPML	<b>0.9583 ± 0.0047</b>	<b>0.7469 ± 0.037</b>	<b>0.3732 ± 0.0374</b>	0.9699 ± 0.0037	<b>0.8142 ± 0.0256</b>	<b>0.8246 ± 0.0316</b>	<b>0.5445 ± 0.0055</b>	0.9276 ± 0.009		
784.0	0.5	G-VGPML	0.9102 ± 0.0029	0.1633 ± 0.0314	0.6513 ± 0.0002	<b>0.978 ± 0.0023</b>	0.4675 ± 0.0155	0.2235 ± 0.04	0.6877 ± 0.0005	<b>0.9453 ± 0.0039</b>		
	1.0	G-VGPML	0.9122 ± 0.0029	0.1973 ± 0.0313	0.6154 ± 0.0004	0.9727 ± 0.0025	0.4842 ± 0.0157	0.2657 ± 0.0385	0.6814 ± 0.0012	0.9333 ± 0.0051		
	2.0	G-VGPML	0.9145 ± 0.0025	0.235 ± 0.0308	0.5586 ± 0.0006	0.9674 ± 0.0023	0.5035 ± 0.0166	0.3119 ± 0.0383	0.677 ± 0.0022	0.9217 ± 0.0059		
	4.0	G-VGPML	0.9205 ± 0.0021	0.324 ± 0.0361	0.4909 ± 0.0008	0.9631 ± 0.002	0.5501 ± 0.0202	0.4145 ± 0.0418	0.666 ± 0.0036	0.9122 ± 0.0064		
200	1.0	-	VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		14.5	-	VGPML	0.9527 ± 0.0026	0.782 ± 0.0129	0.2892 ± 0.0142	0.9641 ± 0.0082	0.8461 ± 0.0075	0.8802 ± 0.0061	0.4241 ± 0.0147	0.9377 ± 0.0067
		28.0	-	VGPML	<b>0.9684 ± 0.0013</b>	<b>0.8416 ± 0.0086</b>	<b>0.233 ± 0.0071</b>	<b>0.9782 ± 0.0023</b>	<b>0.8826 ± 0.0044</b>	<b>0.9043 ± 0.0032</b>	<b>0.3759 ± 0.0132</b>	<b>0.9498 ± 0.0042</b>
		406.0	-	VGPML	<b>0.9607 ± 0.0018</b>	<b>0.768 ± 0.0093</b>	<b>0.4047 ± 0.0046</b>	<b>0.9714 ± 0.003</b>	<b>0.8293 ± 0.0055</b>	<b>0.843 ± 0.006</b>	<b>0.5489 ± 0.0056</b>	<b>0.9308 ± 0.0043</b>
	784.0	-	VGPML	0.9206 ± 0.0022	0.3263 ± 0.0358	0.4913 ± 0.0008	0.9631 ± 0.002	0.5514 ± 0.02	0.4171 ± 0.0412	0.6658 ± 0.0035	0.9122 ± 0.0064	
	1.0	0.5	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		1.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		2.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
		4.0	G-VGPML	0.9014 ± 0.0026	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	0.3905 ± 0.0003	0.0 ± 0.0	0.6931 ± 0.0	0.5 ± 0.0	
	14.5	0.5	G-VGPML	0.9541 ± 0.0023	0.7032 ± 0.0142	0.6244 ± 0.0016	0.947 ± 0.006	0.7787 ± 0.0105	0.7812 ± 0.0122	0.6689 ± 0.005	0.8728 ± 0.0098	
		1.0	G-VGPML	<b>0.9648 ± 0.0025</b>	<b>0.7951 ± 0.0107</b>	<b>0.5862 ± 0.0044</b>	<b>0.9585 ± 0.0052</b>	<b>0.8464 ± 0.0028</b>	<b>0.8606 ± 0.0023</b>	<b>0.6348 ± 0.0105</b>	<b>0.9076 ± 0.0016</b>	
		2.0	G-VGPML	0.9641 ± 0.0037	0.8184 ± 0.0146	0.493 ± 0.0151	0.9663 ± 0.006	0.8692 ± 0.0086	0.8926 ± 0.0078	0.5242 ± 0.015	0.9316 ± 0.0074	
4.0		G-VGPML	<b>0.9703 ± 0.0029</b>	<b>0.8404 ± 0.0145</b>	<b>0.135 ± 0.0061</b>	0.9669 ± 0.0047	<b>0.8817 ± 0.0128</b>	<b>0.8988 ± 0.0118</b>	<b>0.3347 ± 0.0162</b>	0.9288 ± 0.0106		
28.0	0.5	G-VGPML	0.9461 ± 0.0026	0.6287 ± 0.0127	0.6364 ± 0.0006	0.9753 ± 0.0039	0.728 ± 0.0065	0.7149 ± 0.0083	0.6763 ± 0.001	0.9353 ± 0.0079		
	1.0	G-VGPML	0.9646 ± 0.0023	0.7863 ± 0.0109	0.597 ± 0.0007	<b>0.98 ± 0.0035</b>	0.8381 ± 0.0062	0.8487 ± 0.0067	0.6507 ± 0.0018	<b>0.9502 ± 0.0055</b>		
	2.0	G-VGPML	0.9738 ± 0.0024	0.8651 ± 0.0099	0.5185 ± 0.0063	0.9827 ± 0.0035	0.8993 ± 0.0075	<b>0.9164 ± 0.0066</b>	0.5613 ± 0.0088	0.9605 ± 0.0048		
	4.0	G-VGPML	<b>0.9757 ± 0.0014</b>	<b>0.8694 ± 0.0084</b>	<b>0.091 ± 0.0046</b>	0.9831 ± 0.0025	<b>0.9009 ± 0.0066</b>	0.9151 ± 0.0069	<b>0.2617 ± 0.0102</b>	<b>0.9687 ± 0.0042</b>		
406.0	0.5	G-VGPML	0.9334 ± 0.0068	0.4892 ± 0.0637	0.6494 ± 0.0004	<b>0.9893 ± 0.0016</b>	0.6443 ± 0.0358	0.5856 ± 0.0565	0.6815 ± 0.0012	0.9225 ± 0.002		
	1.0	G-VGPML	0.9496 ± 0.0059	0.6559 ± 0.0507	0.6096 ± 0.0015	0.9863 ± 0.0019	0.7454 ± 0.0321	0.7352 ± 0.0421	0.6644 ± 0.0037	0.9441 ± 0.0045		
	2.0	G-VGPML	<b>0.9615 ± 0.0016</b>	<b>0.7641 ± 0.0095</b>	0.5375 ±							

A. Tablas

Num. Inducing	Lengthscale	k	Model	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	1.0	0.5	G-VGPMIL	<b>0.6181 ± 0.0831</b>	<b>0.4415 ± 0.0831</b>	<b>0.6931 ± 0.0831</b>	<b>0.6343 ± 0.0831</b>
		1.0	G-VGPMIL	0.6181 ± 0.0831	0.4415 ± 0.0831	0.6931 ± 0.0831	0.6343 ± 0.0831
		2.0	G-VGPMIL	0.6181 ± 0.0831	0.4415 ± 0.0831	0.6931 ± 0.0831	0.6343 ± 0.0831
		4.0	G-VGPMIL	0.6181 ± 0.0831	0.4415 ± 0.0831	0.6931 ± 0.0831	0.6343 ± 0.0831
	-	VGPMIL	0.6181 ± 0.0831	0.4415 ± 0.0831	0.6931 ± 0.0831	0.6343 ± 0.0831	
	6.942	0.5	G-VGPMIL	0.6959 ± 0.1138	0.5537 ± 0.1138	0.6706 ± 0.1138	<b>0.7519 ± 0.1138</b>
		1.0	G-VGPMIL	0.6965 ± 0.1397	0.5678 ± 0.1397	0.6625 ± 0.1397	0.7486 ± 0.1397
		2.0	G-VGPMIL	<b>0.7076 ± 0.1434</b>	<b>0.5845 ± 0.1434</b>	<b>0.6554 ± 0.1434</b>	0.7375 ± 0.1434
		4.0	G-VGPMIL	0.6959 ± 0.1337	0.5633 ± 0.1337	0.6659 ± 0.1337	0.7158 ± 0.1337
	-	VGPMIL	0.6959 ± 0.1337	0.5633 ± 0.1337	0.6585 ± 0.1337	0.7249 ± 0.1337	
	12.884	0.5	G-VGPMIL	<b>0.6427 ± 0.0957</b>	<b>0.4449 ± 0.0957</b>	0.6643 ± 0.0957	<b>0.859 ± 0.0957</b>
		1.0	G-VGPMIL	0.6427 ± 0.0957	0.4449 ± 0.0957	0.6488 ± 0.0957	0.8415 ± 0.0957
		2.0	G-VGPMIL	0.6322 ± 0.1098	0.4418 ± 0.1098	<b>0.6359 ± 0.1098</b>	0.8133 ± 0.1098
		4.0	G-VGPMIL	0.6105 ± 0.1038	0.3971 ± 0.1038	0.669 ± 0.1038	0.7894 ± 0.1038
	-	VGPMIL	0.6216 ± 0.1119	0.4162 ± 0.1119	0.6568 ± 0.1119	0.7985 ± 0.1119	
	83.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.6885 ± 0.0129</b>	<b>0.7973 ± 0.0129</b>
		1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.695 ± 0.0129	0.7736 ± 0.0129
		2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7266 ± 0.0129	0.7254 ± 0.0129
		4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8649 ± 0.0129	0.6842 ± 0.0129
	-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8596 ± 0.0129	0.6889 ± 0.0129	
166.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.694 ± 0.0129</b>	<b>0.7523 ± 0.0129</b>	
	1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7028 ± 0.0129	0.7254 ± 0.0129	
	2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7379 ± 0.0129	0.699 ± 0.0129	
	4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8861 ± 0.0129	0.68 ± 0.0129	
-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8809 ± 0.0129	0.6822 ± 0.0129		
100	1.0	0.5	G-VGPMIL	<b>0.6953 ± 0.0743</b>	<b>0.5956 ± 0.0743</b>	0.6931 ± 0.0743	<b>0.7464 ± 0.0743</b>
		1.0	G-VGPMIL	0.6848 ± 0.0632	0.5874 ± 0.0632	0.6931 ± 0.0632	0.7431 ± 0.0632
		2.0	G-VGPMIL	0.6848 ± 0.0632	0.5874 ± 0.0632	0.6931 ± 0.0632	0.7431 ± 0.0632
		4.0	G-VGPMIL	0.6848 ± 0.0632	0.5874 ± 0.0632	<b>0.6931 ± 0.0632</b>	0.7431 ± 0.0632
	-	VGPMIL	0.6848 ± 0.0632	0.5874 ± 0.0632	0.6931 ± 0.0632	0.7431 ± 0.0632	
	6.942	0.5	G-VGPMIL	0.6538 ± 0.1069	0.483 ± 0.1069	0.6724 ± 0.1069	0.7716 ± 0.1069
		1.0	G-VGPMIL	0.6538 ± 0.1215	0.4983 ± 0.1215	0.662 ± 0.1215	<b>0.7788 ± 0.1215</b>
		2.0	G-VGPMIL	0.7175 ± 0.1316	0.6143 ± 0.1316	0.6488 ± 0.1316	0.7657 ± 0.1316
		4.0	G-VGPMIL	0.7287 ± 0.1366	0.6191 ± 0.1366	0.6427 ± 0.1366	0.7486 ± 0.1366
	-	VGPMIL	<b>0.7392 ± 0.1466</b>	<b>0.6307 ± 0.1466</b>	<b>0.6395 ± 0.1466</b>	0.7583 ± 0.1466	
	12.884	0.5	G-VGPMIL	0.6316 ± 0.1077	0.4412 ± 0.1077	0.665 ± 0.1077	<b>0.8615 ± 0.1077</b>
		1.0	G-VGPMIL	<b>0.6427 ± 0.0957</b>	<b>0.4449 ± 0.0957</b>	0.6487 ± 0.0957	0.8519 ± 0.0957
		2.0	G-VGPMIL	0.6322 ± 0.1098	0.4418 ± 0.1098	<b>0.6341 ± 0.1098</b>	0.8274 ± 0.1098
		4.0	G-VGPMIL	0.6216 ± 0.1119	0.4162 ± 0.1119	0.6623 ± 0.1119	0.7985 ± 0.1119
	-	VGPMIL	0.6216 ± 0.1119	0.4162 ± 0.1119	0.6522 ± 0.1119	0.8081 ± 0.1119	
	83.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.6885 ± 0.0129</b>	<b>0.802 ± 0.0129</b>
		1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.695 ± 0.0129	0.7736 ± 0.0129
		2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7266 ± 0.0129	0.7254 ± 0.0129
		4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8649 ± 0.0129	0.6842 ± 0.0129
	-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8596 ± 0.0129	0.6889 ± 0.0129	
166.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.694 ± 0.0129</b>	<b>0.7523 ± 0.0129</b>	
	1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7028 ± 0.0129	0.7254 ± 0.0129	
	2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7379 ± 0.0129	0.699 ± 0.0129	
	4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8861 ± 0.0129	0.68 ± 0.0129	
-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8809 ± 0.0129	0.6822 ± 0.0129		
200	1.0	0.5	G-VGPMIL	<b>0.8053 ± 0.0606</b>	<b>0.7896 ± 0.0606</b>	0.6931 ± 0.0606	<b>0.8616 ± 0.0606</b>
		1.0	G-VGPMIL	0.8053 ± 0.0606	0.7896 ± 0.0606	0.6931 ± 0.0606	0.8616 ± 0.0606
		2.0	G-VGPMIL	0.8053 ± 0.0606	0.7896 ± 0.0606	0.6931 ± 0.0606	0.8591 ± 0.0606
		4.0	G-VGPMIL	0.8053 ± 0.0606	0.7896 ± 0.0606	<b>0.6931 ± 0.0606</b>	0.8567 ± 0.0606
	-	VGPMIL	0.8053 ± 0.0606	0.7896 ± 0.0606	0.6931 ± 0.0606	0.8567 ± 0.0606	
	6.942	0.5	G-VGPMIL	0.7181 ± 0.0899	0.6298 ± 0.0899	0.6685 ± 0.0899	<b>0.8067 ± 0.0899</b>
		1.0	G-VGPMIL	0.707 ± 0.1184	0.6027 ± 0.1184	0.6588 ± 0.1184	0.797 ± 0.1184
		2.0	G-VGPMIL	0.7175 ± 0.1316	0.6143 ± 0.1316	0.6408 ± 0.1316	0.7916 ± 0.1316
		4.0	G-VGPMIL	0.7497 ± 0.1113	0.6741 ± 0.1113	0.6272 ± 0.1113	0.7788 ± 0.1113
	-	VGPMIL	<b>0.7608 ± 0.1162</b>	<b>0.6888 ± 0.1162</b>	<b>0.6254 ± 0.1162</b>	0.7815 ± 0.1162	
	12.884	0.5	G-VGPMIL	<b>0.6427 ± 0.0957</b>	0.4449 ± 0.0957	0.6651 ± 0.0957	<b>0.8711 ± 0.0957</b>
		1.0	G-VGPMIL	0.6322 ± 0.1098	0.4418 ± 0.1098	0.648 ± 0.1098	0.8615 ± 0.1098
		2.0	G-VGPMIL	0.6427 ± 0.0957	<b>0.47 ± 0.0957</b>	<b>0.6321 ± 0.0957</b>	0.8346 ± 0.0957
		4.0	G-VGPMIL	0.6322 ± 0.0992	0.4444 ± 0.0992	0.6584 ± 0.0992	0.8005 ± 0.0992
	-	VGPMIL	0.6322 ± 0.0992	0.4444 ± 0.0992	0.648 ± 0.0992	0.8128 ± 0.0992	
	83.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.6885 ± 0.0129</b>	<b>0.802 ± 0.0129</b>
		1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.695 ± 0.0129	0.7736 ± 0.0129
		2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7266 ± 0.0129	0.7254 ± 0.0129
		4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8649 ± 0.0129	0.6842 ± 0.0129
	-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8596 ± 0.0129	0.6889 ± 0.0129	
166.0	0.5	G-VGPMIL	<b>0.4895 ± 0.0129</b>	<b>0.0 ± 0.0129</b>	<b>0.694 ± 0.0129</b>	<b>0.7523 ± 0.0129</b>	
	1.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7028 ± 0.0129	0.7254 ± 0.0129	
	2.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.7379 ± 0.0129	0.699 ± 0.0129	
	4.0	G-VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8861 ± 0.0129	0.68 ± 0.0129	
-	VGPMIL	0.4895 ± 0.0129	0.0 ± 0.0129	0.8809 ± 0.0129	0.6822 ± 0.0129		

Tabla A.3.: Resultados obtenidos en musk1. Para cada configuración del número de inducing points, lengthscale y  $k$ , se señala en negrita el mejor resultado de cada métrica.

Num. Inducing	Lengthscale	k	Model	Bag Acc.	Bag F1	Bag Log Loss	Bag AUC
50	1.0	0.5	G-VGPML	<b>0.6657 ± 0.0808</b>	<b>0.2133 ± 0.0808</b>	0.6931 ± 0.0808	<b>0.4679 ± 0.0808</b>
		1.0	G-VGPML	0.6657 ± 0.0808	0.2133 ± 0.0808	0.6931 ± 0.0808	0.4679 ± 0.0808
		2.0	G-VGPML	0.6657 ± 0.0808	0.2133 ± 0.0808	0.6931 ± 0.0808	0.4679 ± 0.0808
		4.0	G-VGPML	0.6657 ± 0.0808	0.2133 ± 0.0808	<b>0.6931 ± 0.0808</b>	0.4679 ± 0.0808
	6.942	1.0	G-VGPML	0.7167 ± 0.0793	0.5508 ± 0.0793	0.661 ± 0.0793	0.7575 ± 0.0793
		2.0	G-VGPML	<b>0.7452 ± 0.1102</b>	<b>0.6239 ± 0.1102</b>	0.638 ± 0.1102	<b>0.7786 ± 0.1102</b>
		4.0	G-VGPML	0.6976 ± 0.0776	0.5254 ± 0.0776	0.6055 ± 0.0776	0.7238 ± 0.0776
		-	VGPML	0.6952 ± 0.1126	0.5078 ± 0.1126	<b>0.603 ± 0.1126</b>	0.732 ± 0.1126
	12.884	0.5	G-VGPML	0.6767 ± 0.0794	0.3143 ± 0.0794	0.6583 ± 0.0794	0.8068 ± 0.0794
		1.0	G-VGPML	0.6767 ± 0.0575	0.3356 ± 0.0575	0.6279 ± 0.0575	0.8166 ± 0.0575
		2.0	G-VGPML	<b>0.7557 ± 0.0488</b>	<b>0.5991 ± 0.0488</b>	<b>0.5434 ± 0.0488</b>	<b>0.8177 ± 0.0488</b>
		4.0	G-VGPML	0.6571 ± 0.0685	0.2035 ± 0.0685	0.6734 ± 0.0685	0.7775 ± 0.0685
83.0	0.5	G-VGPML	<b>0.6276 ± 0.0196</b>	<b>0.0444 ± 0.0196</b>	0.6649 ± 0.0196	<b>0.8281 ± 0.0196</b>	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6458 ± 0.0183</b>	0.8171 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6588 ± 0.0183	0.7661 ± 0.0183	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.6515 ± 0.0183	0.5398 ± 0.0183	
166.0	0.5	G-VGPML	<b>0.6176 ± 0.0183</b>	<b>0.0 ± 0.0183</b>	0.6688 ± 0.0183	<b>0.7917 ± 0.0183</b>	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6577 ± 0.0183</b>	0.7537 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6943 ± 0.0183	0.6809 ± 0.0183	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.769 ± 0.0183	0.4501 ± 0.0183	
100	1.0	0.5	G-VGPML	<b>0.7257 ± 0.0577</b>	<b>0.4625 ± 0.0577</b>	0.6931 ± 0.0577	<b>0.5766 ± 0.0577</b>
		1.0	G-VGPML	0.7157 ± 0.0473	0.4522 ± 0.0473	0.6931 ± 0.0473	0.5734 ± 0.0473
		2.0	G-VGPML	0.7157 ± 0.0473	0.4522 ± 0.0473	0.6931 ± 0.0473	0.5734 ± 0.0473
		4.0	G-VGPML	0.7157 ± 0.0473	0.4522 ± 0.0473	<b>0.6931 ± 0.0473</b>	0.5734 ± 0.0473
	6.942	0.5	G-VGPML	0.6981 ± 0.098	0.48 ± 0.098	0.6654 ± 0.098	0.7366 ± 0.098
		1.0	G-VGPML	0.7362 ± 0.0614	0.5788 ± 0.0614	0.6434 ± 0.0614	0.7579 ± 0.0614
		2.0	G-VGPML	0.7357 ± 0.0348	0.6012 ± 0.0348	0.6169 ± 0.0348	<b>0.7749 ± 0.0348</b>
		4.0	G-VGPML	0.7138 ± 0.147	0.4244 ± 0.147	0.6791 ± 0.147	0.7175 ± 0.147
	12.884	0.5	G-VGPML	<b>0.7538 ± 0.1338</b>	<b>0.6061 ± 0.1338</b>	<b>0.5952 ± 0.1338</b>	<b>0.7553 ± 0.1338</b>
		1.0	G-VGPML	0.6767 ± 0.08	0.3453 ± 0.08	0.6566 ± 0.08	0.7877 ± 0.08
		2.0	G-VGPML	0.6671 ± 0.0834	0.3408 ± 0.0834	0.6255 ± 0.0834	0.7887 ± 0.0834
		4.0	G-VGPML	<b>0.7167 ± 0.0598</b>	<b>0.5221 ± 0.0598</b>	<b>0.5472 ± 0.0598</b>	<b>0.8271 ± 0.0598</b>
83.0	0.5	G-VGPML	<b>0.6276 ± 0.0196</b>	<b>0.0444 ± 0.0196</b>	0.6649 ± 0.0196	<b>0.8241 ± 0.0196</b>	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6457 ± 0.0183</b>	0.8213 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6589 ± 0.0183	0.7661 ± 0.0183	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.6517 ± 0.0183	0.5356 ± 0.0183	
166.0	0.5	G-VGPML	<b>0.6176 ± 0.0183</b>	<b>0.0 ± 0.0183</b>	0.6688 ± 0.0183	<b>0.7916 ± 0.0183</b>	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6577 ± 0.0183</b>	0.7495 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6943 ± 0.0183	0.6809 ± 0.0183	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.7689 ± 0.0183	0.4501 ± 0.0183	
200	1.0	0.5	G-VGPML	0.7043 ± 0.0973	0.3941 ± 0.0973	0.6931 ± 0.0973	0.5369 ± 0.0973
		1.0	G-VGPML	<b>0.7138 ± 0.1129</b>	<b>0.4091 ± 0.1129</b>	0.6931 ± 0.1129	<b>0.5456 ± 0.1129</b>
		2.0	G-VGPML	0.7138 ± 0.1129	0.4091 ± 0.1129	0.6931 ± 0.1129	0.5456 ± 0.1129
		4.0	G-VGPML	0.7138 ± 0.1129	0.4091 ± 0.1129	<b>0.6931 ± 0.1129</b>	0.5456 ± 0.1129
	6.942	0.5	G-VGPML	0.7167 ± 0.0598	0.5012 ± 0.0598	0.6595 ± 0.0598	0.7801 ± 0.0598
		1.0	G-VGPML	0.7467 ± 0.0889	0.5659 ± 0.0889	0.6334 ± 0.0889	<b>0.795 ± 0.0889</b>
		2.0	G-VGPML	0.7367 ± 0.0671	0.5796 ± 0.0671	0.606 ± 0.0671	0.795 ± 0.0671
		4.0	G-VGPML	0.6871 ± 0.0607	0.3933 ± 0.0607	0.6502 ± 0.0607	0.7298 ± 0.0607
	12.884	0.5	G-VGPML	<b>0.7543 ± 0.1427</b>	<b>0.6006 ± 0.1427</b>	<b>0.5723 ± 0.1427</b>	<b>0.7772 ± 0.1427</b>
		1.0	G-VGPML	0.7162 ± 0.0694	0.4249 ± 0.0694	0.6532 ± 0.0694	0.8365 ± 0.0694
		2.0	G-VGPML	0.7257 ± 0.0577	0.4889 ± 0.0577	0.6204 ± 0.0577	<b>0.8389 ± 0.0577</b>
		4.0	G-VGPML	<b>0.7462 ± 0.0667</b>	<b>0.5664 ± 0.0667</b>	<b>0.5415 ± 0.0667</b>	0.8199 ± 0.0667
83.0	0.5	G-VGPML	0.6371 ± 0.0747	0.1944 ± 0.0747	0.6632 ± 0.0747	0.7999 ± 0.0747	
	1.0	G-VGPML	0.6767 ± 0.0575	0.3756 ± 0.0575	0.5421 ± 0.0575	0.8181 ± 0.0575	
	2.0	G-VGPML	<b>0.6276 ± 0.0196</b>	<b>0.0444 ± 0.0196</b>	0.6649 ± 0.0196	<b>0.828 ± 0.0196</b>	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6457 ± 0.0183</b>	0.8212 ± 0.0183	
166.0	0.5	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6589 ± 0.0183	0.7661 ± 0.0183	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.6514 ± 0.0183	0.5378 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.9816 ± 0.0183	0.6686 ± 0.0183	
	4.0	G-VGPML	<b>0.6176 ± 0.0183</b>	<b>0.0 ± 0.0183</b>	0.6688 ± 0.0183	<b>0.7895 ± 0.0183</b>	
200	0.5	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	<b>0.6577 ± 0.0183</b>	0.7516 ± 0.0183	
	1.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	0.6943 ± 0.0183	0.6809 ± 0.0183	
	2.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.7689 ± 0.0183	0.4543 ± 0.0183	
	4.0	G-VGPML	0.6176 ± 0.0183	0.0 ± 0.0183	1.0428 ± 0.0183	0.5804 ± 0.0183	

Tabla A.4.: Resultados obtenidos en musk2. Para cada configuración del número de inducing points, lengthscale y  $k$ , se señala en negrita el mejor resultado de cada métrica.

Num.Inducing	Lengthscale	k	Model	RSNA inst. Acc.	RSNA inst. F1	RSNA inst. Log Loss	RSNA inst. AUC	RSNA inst. Acc.	RSNA inst. F1	RSNA inst. Log Loss	RSNA inst. AUC	CQ500 bag. Acc.	CQ500 bag. F1	CQ500 bag. Log Loss	CQ500 bag. AUC		
50	0.5	G-VGPMIL	0.5	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			1.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			2.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			4.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
	2.82	G-VGPMIL	0.5	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			1.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			2.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			4.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
	100	0.5	G-VGPMIL	0.5	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032
				1.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032
				2.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032
				4.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032
2.82		G-VGPMIL	0.5	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			1.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			2.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	
			4.0	0.9389 ± 0.0032	0.7805 ± 0.0072	0.6205 ± 0.0032	0.9488 ± 0.0032	<b>0.802 ± 0.0032</b>	<b>0.801 ± 0.0032</b>	0.9488 ± 0.0032	0.6278 ± 0.0032	0.9487 ± 0.0032	0.8176 ± 0.0032	0.8015 ± 0.0032	0.6673 ± 0.0032	0.9064 ± 0.0032	

Tabla A.5.: Resultados obtenidos en los conjuntos RSNA y CQ500, fijando a 8 el número de características extraídas con AttCNN. Para cada configuración del número de inducing points, lengthscale y k, se señala en negrita el mejor resultado de cada métrica.



A. Tablas

Num. Inducing	Lengthscale	k	Model	RSNA Inst. Acc.	RSNA Inst. F1	RSNA Inst. Log. Loss	RSNA Inst. AUC	RSNA Inst. Acc.	RSNA Inst. F1	RSNA Inst. Log. Loss	RSNA Inst. AUC	CQ500 Inst. Acc.	CQ500 Inst. F1	CQ500 Inst. Log. Loss	CQ500 Inst. AUC
1.0	11.31	0.5	G-FCPML	0.9222±0.0129	0.666±0.0129	0.867±0.0129	0.905±0.0129	0.899±0.0129	0.881±0.0129	0.888±0.0129	0.884±0.0129	0.913±0.0129	0.799±0.0129	0.699±0.0129	0.863±0.0129
		1.0	G-FCPML	0.9222±0.0132	0.665±0.0132	0.867±0.0132	0.905±0.0132	0.899±0.0132	0.881±0.0132	0.888±0.0132	0.884±0.0132	0.913±0.0132	0.799±0.0132	0.699±0.0132	0.863±0.0132
		2.0	G-FCPML	0.9221±0.0134	0.668±0.0134	0.867±0.0134	0.905±0.0134	0.899±0.0134	0.881±0.0134	0.888±0.0134	0.884±0.0134	0.913±0.0134	0.799±0.0134	0.699±0.0134	0.863±0.0134
		4.0	G-FCPML	0.922±0.0136	0.668±0.0136	0.867±0.0136	0.905±0.0136	0.899±0.0136	0.881±0.0136	0.888±0.0136	0.884±0.0136	0.913±0.0136	0.799±0.0136	0.699±0.0136	0.863±0.0136
6.15	11.31	0.5	G-FCPML	0.9219±0.0135	0.67±0.0135	0.867±0.0135	0.905±0.0135	0.899±0.0135	0.881±0.0135	0.888±0.0135	0.884±0.0135	0.913±0.0135	0.799±0.0135	0.699±0.0135	0.863±0.0135
		1.0	G-FCPML	0.9219±0.0135	0.67±0.0135	0.867±0.0135	0.905±0.0135	0.899±0.0135	0.881±0.0135	0.888±0.0135	0.884±0.0135	0.913±0.0135	0.799±0.0135	0.699±0.0135	0.863±0.0135
		2.0	G-FCPML	0.9219±0.0135	0.67±0.0135	0.867±0.0135	0.905±0.0135	0.899±0.0135	0.881±0.0135	0.888±0.0135	0.884±0.0135	0.913±0.0135	0.799±0.0135	0.699±0.0135	0.863±0.0135
		4.0	G-FCPML	0.9219±0.0135	0.67±0.0135	0.867±0.0135	0.905±0.0135	0.899±0.0135	0.881±0.0135	0.888±0.0135	0.884±0.0135	0.913±0.0135	0.799±0.0135	0.699±0.0135	0.863±0.0135
50	11.31	0.5	G-FCPML	0.9235±0.0138	0.705±0.0138	0.867±0.0138	0.905±0.0138	0.899±0.0138	0.881±0.0138	0.888±0.0138	0.884±0.0138	0.913±0.0138	0.799±0.0138	0.699±0.0138	0.863±0.0138
		1.0	G-FCPML	0.9235±0.0138	0.705±0.0138	0.867±0.0138	0.905±0.0138	0.899±0.0138	0.881±0.0138	0.888±0.0138	0.884±0.0138	0.913±0.0138	0.799±0.0138	0.699±0.0138	0.863±0.0138
		2.0	G-FCPML	0.9235±0.0138	0.705±0.0138	0.867±0.0138	0.905±0.0138	0.899±0.0138	0.881±0.0138	0.888±0.0138	0.884±0.0138	0.913±0.0138	0.799±0.0138	0.699±0.0138	0.863±0.0138
		4.0	G-FCPML	0.9235±0.0138	0.705±0.0138	0.867±0.0138	0.905±0.0138	0.899±0.0138	0.881±0.0138	0.888±0.0138	0.884±0.0138	0.913±0.0138	0.799±0.0138	0.699±0.0138	0.863±0.0138
69.65	11.31	0.5	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		1.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		2.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		4.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
128.0	11.31	0.5	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		1.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		2.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		4.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
100	11.31	0.5	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		1.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		2.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141
		4.0	G-FCPML	0.9238±0.0141	0.694±0.0141	0.867±0.0141	0.905±0.0141	0.899±0.0141	0.881±0.0141	0.888±0.0141	0.884±0.0141	0.913±0.0141	0.799±0.0141	0.699±0.0141	0.863±0.0141

Tabla A.7.: Resultados obtenidos en los conjuntos RSNA y CQ500, fijando a 128 el número de características extraídas con AttCNN. Para cada configuración del número de inducing points, lengthscale y k, se señala en negrita el mejor resultado de cada métrica.



# Bibliografía

Las referencias se listan por orden alfabético. Aquellas referencias con más de un autor están ordenadas de acuerdo con el primer autor.

- [1] Andrews, Stuart, Tsochantaridis, Ioannis, & Hofmann, Thomas. 2002. Support vector machines for multiple-instance learning. *Advances in neural information processing systems*, **15**. [Citado en pág. 7]
- [2] Babenko, Boris, Dollár, Piotr, Tu, Zhuowen, & Belongie, Serge. 2008. Simultaneous learning and alignment: Multi-instance and multi-pose learning. *In: Workshop on Faces in 'Real-Life' Images: Detection, Alignment, and Recognition*. [Citado en pág. 8]
- [3] Bandyopadhyay, Sanghamitra, Ghosh, Dip, Mitra, Ramkrishna, & Zhao, Zhongming. 2015. MBS-TAR: multiple instance learning for predicting specific functional binding sites in microRNA targets. *Scientific reports*, **5**(1), 1–12. [Citado en pág. 8]
- [4] Bishop, Christopher M, & Nasrabadi, Nasser M. 2006. *Pattern recognition and machine learning*. Vol. 4. Springer. [Citado en págs. 11 and 19]
- [5] Blei, David M, Kucukelbir, Alp, & McAuliffe, Jon D. 2017. Variational inference: A review for statisticians. *Journal of the American statistical Association*, **112**(518), 859–877. [Citado en pág. 19]
- [6] Boyd, Stephen, Boyd, Stephen P, & Vandenberghe, Lieven. 2004. *Convex optimization*. Cambridge university press. [Citado en pág. 21]
- [7] Carbonneau, Marc-André, Cheplygina, Veronika, Granger, Eric, & Gagnon, Ghyslain. 2018. Multiple instance learning: A survey of problem characteristics and applications. *Pattern Recognition*, **77**, 329–353. [Citado en págs. 8 and 9]
- [8] Chilamkurthy, Sasank, Ghosh, Rohit, Tanamala, Swetha, Biviji, Mustafa, Campeau, Norbert G, Venugopal, Vasantha Kumar, Mahajan, Vidur, Rao, Pooja, & Warier, Prashant. 2018. Deep learning algorithms for detection of critical findings in head CT scans: a retrospective study. *The Lancet*, **392**(10162), 2388–2396. [Citado en págs. 1, 47, and 54]
- [9] Damianou, Andreas. 2015. *Deep Gaussian processes and variational propagation of uncertainty*. Ph.D. thesis, University of Sheffield. [Citado en págs. 11, 16, and 19]
- [10] Dietterich, Thomas G, Lathrop, Richard H, & Lozano-Pérez, Tomás. 1997. Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, **89**(1-2), 31–71. [Citado en págs. 5 and 7]
- [11] Doi, Kunio. 2007. Computer-aided diagnosis in medical imaging: historical review, current status and future potential. *Computerized medical imaging and graphics*, **31**(4-5), 198–211. [Citado en pág. 8]
- [12] Durrett, Rick. 2019. *Probability: theory and examples*. Vol. 49. Cambridge University Press. [Citado en pág. 19]
- [13] Eksi, Ridvan, Li, Hong-Dong, Menon, Rajasree, Wen, Yuchen, Omenn, Gilbert S., Kretzler, Matthias, & Guan, Yuanfang. 2013. Systematically Differentiating Functions for Alternatively Spliced Isoforms through Integrating RNA-seq Data. *PLOS Computational Biology*, **9**(11), 1–16. [Citado en pág. 8]

## Bibliografía

- [14] Flanders, Adam E, Prevedello, Luciano M, Shih, George, Halabi, Safwan S, Kalpathy-Cramer, Jayashree, Ball, Robyn, Mongan, John T, Stein, Anouk, Kitamura, Felipe C, Lungren, Matthew P, *et al.* . 2020. Construction of a machine learning dataset through collaboration: the RSNA 2019 brain CT hemorrhage challenge. *Radiology: Artificial Intelligence*, **2**(3), e190211. [Citado en pág. 46]
- [15] Garcia-Garcia, Alberto, Orts-Escolano, Sergio, Oprea, Sergiu, Villena-Martinez, Victor, & Garcia-Rodriguez, Jose. 2017. A review on deep learning techniques applied to semantic segmentation. *arXiv preprint arXiv:1704.06857*. [Citado en pág. 8]
- [16] Gebel, James M, & Broderick, Joseph P. 2000. Intracerebral hemorrhage. *Neurologic clinics*, **18**(2), 419–438. [Citado en pág. 1]
- [17] Gibbs, Mark N, & MacKay, David JC. 2000. Variational Gaussian process classifiers. *IEEE Transactions on Neural Networks*, **11**(6), 1458–1464. [Citado en pág. 15]
- [18] Goodfellow, Ian, Bengio, Yoshua, & Courville, Aaron. 2016. *Deep learning*. MIT press. [Citado en pág. 1]
- [19] Hausmann, Manuel, Hamprecht, Fred A, & Kandemir, Melih. 2017. Variational bayesian multiple instance learning with gaussian processes. *Pages 6570–6579 of: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. [Citado en págs. 2, 3, 8, 23, 24, and 28]
- [20] Ilse, Maximilian, Tomczak, Jakub, & Welling, Max. 2018. Attention-based deep multiple instance learning. *Pages 2127–2136 of: International conference on machine learning*. PMLR. [Citado en pág. 7]
- [21] Jaakkola, Tommi S, & Haussler, David. 1999. Probabilistic kernel regression models. *In: Seventh International Workshop on Artificial Intelligence and Statistics*. PMLR. [Citado en pág. 15]
- [22] Jaakkola, Tommi S, & Jordan, Michael I. 2000. Bayesian parameter estimation via variational methods. *Statistics and Computing*, **10**(1), 25–37. [Citado en págs. 2, 3, and 27]
- [23] Kim, Minyoung, & De la Torre, Fernando. 2010. Gaussian processes multiple instance learning. *In: ICML*. [Citado en págs. 7, 8, and 25]
- [24] Liu, Ying, Zhang, Dengsheng, Lu, Guojun, & Ma, Wei-Ying. 2007. A survey of content-based image retrieval with high-level semantics. *Pattern recognition*, **40**(1), 262–282. [Citado en pág. 8]
- [25] Mangasarian, Olvi L, & Wild, Edward W. 2008. Multiple instance classification via successive linear programming. *Journal of optimization theory and applications*, **137**(3), 555–568. [Citado en pág. 7]
- [26] Maron, Oded, & Lozano-Pérez, Tomás. 1997. A framework for multiple-instance learning. *Advances in neural information processing systems*, **10**. [Citado en pág. 7]
- [27] Minka, Thomas Peter. 2001. *A family of algorithms for approximate Bayesian inference*. Ph.D. thesis, Massachusetts Institute of Technology. [Citado en pág. 15]
- [28] Monteiro, Miguel, Newcombe, Virginia FJ, Mathieu, Francois, Adatia, Krishma, Kamnitsas, Konstantinos, Ferrante, Enzo, Das, Tilak, Whitehouse, Daniel, Rueckert, Daniel, Menon, David K, *et al.* . 2020. Multiclass semantic segmentation and quantification of traumatic brain injury lesions on head CT using deep learning: an algorithm development and multicentre validation study. *The Lancet Digital Health*, **2**(6), e314–e322. [Citado en pág. 54]
- [29] Nguyen, Nhan T, Tran, Dat Q, Nguyen, Nghia T, & Nguyen, Ha Q. 2020. A CNN-LSTM architecture for detection of intracranial hemorrhage on CT scans. *medRxiv*. [Citado en pág. 54]
- [30] Parisi, Giorgio, & Shankar, Ramamurti. 1988. Statistical field theory. [Citado en pág. 20]

- [31] Parizel, P, Makkat, S, Van Miert, E, Van Goethem, J, Van den Hauwe, L, & De Schepper, A. 2001. Intracranial hemorrhage: principles of CT and MRI interpretation. *European radiology*, **11**(9), 1770–1783. [Citado en pág. 1]
- [32] Polson, Nicholas G, Scott, James G, & Windle, Jesse. 2013. Bayesian inference for logistic models using Pólya–Gamma latent variables. *Journal of the American statistical Association*, **108**(504), 1339–1349. [Citado en pág. 31]
- [33] Quinonero–Candela, Joaquin, & Rasmussen, Carl Edward. 2005. A unifying view of sparse approximate Gaussian process regression. *The Journal of Machine Learning Research*, **6**, 1939–1959. [Citado en pág. 17]
- [34] Rasmussen, Carl Edward. 2003. Gaussian processes in machine learning. *Pages 63–71 of: Summer school on machine learning*. Springer. [Citado en págs. 2, 3, 7, and 11]
- [35] Rodríguez-Yáñez, M, Castellanos, M, Freijo, MM, Fernández, JC López, Martí-Fàbregas, J, Nombela, F, Simal, P, Castillo, J, Díez-Tejedor, E, Fuentes, B, *et al.* . 2013. Guías de actuación clínica en la hemorragia intracerebral. *Neurología*, **28**(4), 236–249. [Citado en pág. 1]
- [36] Schmidt, Arne, Silva-Rodríguez, Julio, Molina, Rafael, & Naranjo, Valery. 2022. Efficient Cancer Classification by Coupling Semi Supervised and Multiple Instance Learning. *IEEE Access*, **10**, 9763–9773. [Citado en págs. 9 and 15]
- [37] Seeger, Matthias. 1999. Bayesian model selection for support vector machines, Gaussian processes and other kernel classifiers. *Advances in neural information processing systems*, **12**. [Citado en pág. 15]
- [38] Seeger, Matthias W, Williams, Christopher KI, & Lawrence, Neil D. 2003. Fast forward selection to speed up sparse Gaussian process regression. *Pages 254–261 of: International Workshop on Artificial Intelligence and Statistics*. PMLR. [Citado en pág. 17]
- [39] Snelson, Edward, & Ghahramani, Zoubin. 2005. Sparse Gaussian processes using pseudo-inputs. *Advances in neural information processing systems*, **18**. [Citado en págs. 17 and 24]
- [40] Sobrino, Justin, & Shafi, Shahid. 2013. Timing and causes of death after injuries. *Pages 120–123 of: Baylor University Medical Center Proceedings*, vol. 26. Taylor & Francis. [Citado en pág. 1]
- [41] Strub, WM, Leach, JL, Tomsick, T, & Vagal, A. 2007. Overnight preliminary head CT interpretations provided by residents: locations of misidentified intracranial hemorrhage. *American journal of neuroradiology*, **28**(9), 1679–1682. [Citado en pág. 1]
- [42] Titsias, Michalis. 2009. Variational learning of inducing variables in sparse Gaussian processes. *Pages 567–574 of: Artificial intelligence and statistics*. PMLR. [Citado en págs. 17, 18, and 19]
- [43] Tompson, Jonathan, Goroshin, Ross, Jain, Arjun, LeCun, Yann, & Bregler, Christoph. 2015. Efficient object localization using convolutional networks. *Pages 648–656 of: Proceedings of the IEEE conference on computer vision and pattern recognition*. [Citado en pág. 8]
- [44] Tong, Tong, Wolz, Robin, Gao, Qinquan, Guerrero, Ricardo, Hajnal, Joseph V, Rueckert, Daniel, Initiative, Alzheimer’s Disease Neuroimaging, *et al.* . 2014. Multiple instance learning for classification of dementia in brain MRI. *Medical image analysis*, **18**(5), 808–818. [Citado en pág. 9]
- [45] Vaswani, Ashish, Shazeer, Noam, Parmar, Niki, Uszkoreit, Jakob, Jones, Llion, Gomez, Aidan N, Kaiser, Łukasz, & Polosukhin, Illia. 2017. Attention is all you need. *Advances in neural information processing systems*, **30**. [Citado en pág. 7]
- [46] Vezhnevets, Alexander, & Buhmann, Joachim M. 2010. Towards weakly supervised semantic segmentation by means of multiple instance and multitask learning. *Pages 3249–3256 of: 2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE. [Citado en pág. 8]

## Bibliografía

- [47] Wang, Fulton, & Pinar, Ali. 2021. The Multiple Instance Learning Gaussian Process Probit Model. *Pages 3034–3042 of: International Conference on Artificial Intelligence and Statistics*. PMLR. [Citado en pág. 8]
- [48] Wang, Xinggang, Yan, Yongluan, Tang, Peng, Bai, Xiang, & Liu, Wenyu. 2018. Revisiting multiple instance neural networks. *Pattern Recognition*, **74**, 15–24. [Citado en pág. 7]
- [49] Wei, Xiu-Shen, & Zhou, Zhi-Hua. 2016. An empirical study on image bag generators for multi-instance learning. *Machine learning*, **105**(2), 155–198. [Citado en pág. 8]
- [50] Weidmann, Nils, Frank, Eibe, & Pfahringer, Bernhard. 2003. A two-level learning method for generalized multi-instance problems. *Pages 468–479 of: European Conference on Machine Learning*. Springer. [Citado en págs. 1, 2, 5, and 6]
- [51] Wenzel, Florian, Galy-Fajou, Théo, Donner, Christan, Kloft, Marius, & Opper, Manfred. 2019. Efficient Gaussian process classification using Pólya-Gamma data augmentation. *Pages 5417–5424 of: Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33. [Citado en págs. 2, 3, 23, 30, and 58]
- [52] Williams, Christopher KI, & Barber, David. 1998. Bayesian classification with Gaussian processes. *IEEE Transactions on pattern analysis and machine intelligence*, **20**(12), 1342–1351. [Citado en pág. 15]
- [53] Wu, Yunan, Schmidt, Arne, Hernández-Sánchez, Enrique, Molina, Rafael, & Katsaggelos, Aggelos K. 2021. Combining attention-based multiple instance learning and gaussian processes for CT hemorrhage detection. *Pages 582–591 of: International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. [Citado en págs. 2, 3, 49, 50, and 52]
- [54] Yan, Yongluan, Wang, Xinggang, Guo, Xiaojie, Fang, Jiemin, Liu, Wenyu, & Huang, Junzhou. 2018. Deep multi-instance learning with dynamic pooling. *Pages 662–677 of: Asian Conference on Machine Learning*. PMLR. [Citado en pág. 7]
- [55] Yang, Jian, Zhang, David, Frangi, Alejandro F, & Yang, Jing-yu. 2004. Two-dimensional PCA: a new approach to appearance-based face representation and recognition. *IEEE transactions on pattern analysis and machine intelligence*, **26**(1), 131–137. [Citado en pág. 42]
- [56] Yasser, EL-Manzalawy, Dobbs, Drena, & Honavar, Vasant. 2010. Predicting MHC-II binding affinity using multiple instance regression. *IEEE/ACM Transactions on computational biology and bioinformatics*, **8**(4), 1067–1079. [Citado en pág. 8]
- [57] Zhang, Dingwen, Han, Junwei, Cheng, Gong, & Yang, Ming-Hsuan. 2021. Weakly supervised object localization and detection: A survey. *IEEE transactions on pattern analysis and machine intelligence*. [Citado en pág. 8]
- [58] Zhang, Qi, Goldman, Sally A, Yu, Wei, & Fritts, Jason E. 2002. Content-based image retrieval using multiple-instance learning. *Page 2 of: ICML*, vol. 1. Citeseer. [Citado en págs. 7 and 8]