

# Monografía sobre Aprendizaje Automático

Nicolás Caro

Departamento de Ingeniería Matemática, Universidad de Chile, Santiago.

`ncaro@dim.uchile.cl`

September 4, 2018



# Contents

Preface	vii
Chapter 1. Modelos gráficos probabilísticos	1
Introducción	1
1.1. Modelos gráficos dirigidos	1
1.1.1. Naive Bayes	3
1.1.2. Regresión polinomial	3
1.1.3. Modelos gráficos dirigidos gaussianos	4
1.2. Independencia condicional en modelos gráficos dirigidos	5
1.2.1. d-separación	6
1.2.2. Markov blankets	7
Chapter 2. Métodos de inferencia aproximada	9
Introducción	9
2.1. Inferencia Markov Chain Monte Carlo	9
2.2. Cadenas de Markov	9
2.3. Algoritmo Metropolis Hastings (MH)	10
2.3.1. MH caminata aleatoria (RWMH)	10
2.3.2. Muestreo de independencia	11
2.3.3. Metropolis simétrico	11
2.4. Muestreo de Gibbs	11
2.5. Convergencia en métodos MCMC	12
2.5.1. Tamaño efectivo de muestra	12
2.5.2. Varianza entre cadenas	13
2.6. Monte Carlo Hamiltoniano	13
2.7. Inferencia variacional	15
2.8. Descripción del método	15
2.8.1. GMM Bayesiano I	16
2.9. Cota inferior para la evidencia	17
2.10. Familia variacional mean-field	18
2.10.1. GMM Bayesiano II	18
2.11. Inferencia variacional mean-field de coordenadas ascendentes	18
2.11.1. Formulación	19
2.11.2. Actualización de coordenadas	19
2.11.3. GMM Bayesiano III	19
2.11.4. Densidad variacional para las asignaciones de cluster	20
2.11.5. Densidad variacional para las medias	21
2.11.6. CAVI para GMM	21
2.12. Inferencia variacional con familias exponenciales	22
2.12.1. Condicionales completas en familias exponenciales	22

2.12.2.	Conjugación condicional y modelos bayesianos	23
2.13.	Inferencia variacional en modelos condicionalmente conjugados	24
2.14.	Inferencia variacional estocástica	25
2.14.1.	Gradiente natural de ELBO	25
2.14.2.	Optimización estocástica de ELBO.	26
Chapter 3.	Aprendizaje con Kernels	27
3.1.	Introducción	27
3.2.	Terminología y propiedades	27
3.3.	Espacios de Hilbert con kernel reproductor - RKHS	29
3.4.	Kernels	31
3.4.1.	Kernel RBF	31
3.4.2.	Kernels lineales	32
3.4.3.	Matern Kernels	32
3.5.	Kernels derivados de modelos probabilísticos generativos	32
3.5.1.	Kernels producto de probabilidad	32
3.5.2.	Kernels Fisher	33
3.6.	Kernels en modelos lineales generalizados	33
3.6.1.	Máquinas de kernel	33
3.6.2.	LIVMs, RVMs y Maquinas de vectores sparse	33
3.6.3.	Maquinas de vectores de soporte (SVM)	34
3.6.4.	LS-SVM	37
3.6.5.	Clasificación kernel knn	38
3.6.6.	Clustering de K-medoides Kernelizado	38
3.6.7.	Kernel PCA	39
3.7.	Kernels para construir modelos generativos	40
3.7.1.	Kernel de suavizado	40
3.7.2.	kernel density estimation (KDE)	41
3.7.3.	De KDE a KNN	42
3.7.4.	Regresión por kernel	42
3.7.5.	Regresión localmente ponderada	43
Chapter 4.	Discusión	45
Appendix A.	Anexos	47
A.0.1.	Tópicos de Análisis	47
A.0.1.1.	Normas	47
A.0.1.2.	Cauchy Schwarz	47
A.0.2.	Tópicos de Teoría de la Medida	48
A.0.3.	Funciones Simples	49
A.0.4.	Definición de la integral	49
A.0.5.	Complejidad de un espacio de medida	49
A.0.6.	Probabilidad	49
A.0.7.	Tópicos de Optimización	51
A.0.7.1.	Problema de Optimización con Restricciones, Lagrangiano	51
A.0.7.2.	Condiciones KKT	51
A.0.7.3.	Problema de Programación Cuadrática	52
Appendix.	Bibliography	53

## Preface

El presente trabajo reúne resultados referentes al aprendizaje automático, y busca servir como material de apoyo en el estudio de esta disciplina. En esta versión se estudian de manera autocontenida tópicos sobre modelos gráficos probabilísticos, métodos de inferencia aproximada y aprendizaje con kernels.

Nicolás Caro



## Modelos gráficos probabilísticos

### Introducción

El eje central de este capítulo se basa en la búsqueda de una representación compacta, para distribuciones de probabilidad conjunta de la forma  $p(\mathbf{x}|\boldsymbol{\theta})$ . Esto, con la intención de realizar inferencia sobre variables y aprendizaje de parámetros de manera eficiente.

#### 1.1. Modelos gráficos dirigidos

Toda distribución de probabilidad conjunta  $p(\mathbf{x}) = p(x_1, x_2, \dots, x_v)$  se puede representar de la forma:

$$(1.1) \quad p(\mathbf{x}) = p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \dots p(x_v | x_1, x_2, \dots, x_{v-1})$$

El problema con esta expresión es la dificultad computacional subyacente al cálculo de distribuciones condicionales de la forma  $p(x_t | x_1, \dots, x_{t-1})$  cuando el número de variables incidentes  $t$  aumenta.

No obstante, la representación (1.1) reduce su complejidad en presencia de **independencia condicional**.

En efecto, si se asume  $x_{t+1} \perp x_1, \dots, x_{t-1} | x_t$ . Es decir, las observaciones futuras  $x_{t+1}$  son independientes del pasado  $x_1, \dots, x_{t-1}$ , dado el estado presente  $x_t$ . La probabilidad conjunta se reduce entonces a:

$$(1.2) \quad p(\mathbf{x}) = p(x_1) \prod_{t=2}^v p(x_t | x_1, \dots, x_{t-1}) = p(x_1) \prod_{t=2}^v p(x_t | x_{t-1})$$

De lo cual se obtiene una expresión más simple.

Modelar la independencia condicional entre las variables permite entonces reducir la complejidad de representación para la distribución conjunta. En particular, la elección tomada en (1.2) se conoce como **propiedad de Markov** de primer orden. En un contexto general, las relaciones de independencia condicional entre variables aleatorias de dimensión arbitraria, se modelan utilizando *diagramas de independencia* o **modelos gráficos**. Estos se valen de un grafo  $G = (\mathcal{V}, \mathcal{E})$ <sup>1</sup> para representar mediante nodos  $v = 1, \dots, \mathcal{V}$  las variables aleatorias del modelo, mientras que la presencia o ausencia de aristas entre estos nodos, permite modelar las relaciones de dependencia condicional subyacentes.

Una *red bayesiana* o **modelo gráfico dirigido** es un modelo gráfico probabilístico, cuyo grafo subyacente es un **grafo dirigido acíclico** (DAG por sus siglas en inglés). Todo DAG posee un *ordenamiento topológico*, es decir, los nodos de cualquier DAG pueden ser numerados de manera tal, que todo nodo padre posea una numeración inferior a sus nodos hijos. Esta característica permite enriquecer

---

<sup>1</sup>Conjunto consistente de  $\mathcal{V} = \{1 \dots, V\}$  vértices (o nodos) y  $\mathcal{E} = \{(s, t) : s, t \in \mathcal{V}\}$  aristas.

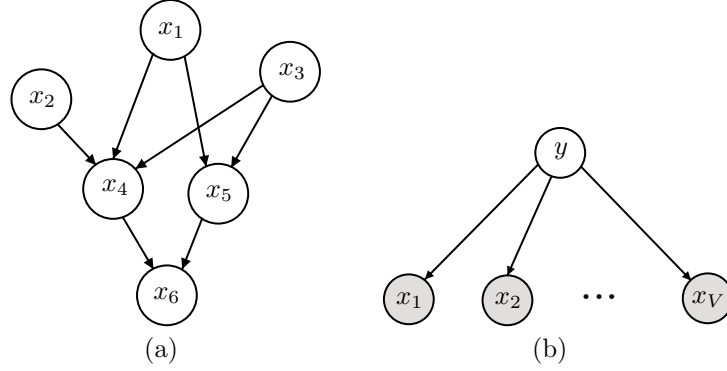


FIGURE 1. (a) Ejemplo de modelo gráfico dirigido. (b) Relaciones de dependencia condicional en el clasificador naive Bayes como un modelo gráfico dirigido, las variables aleatorias observadas se denotan por nodos grises.

la formulación de la propiedad de Markov (1.2), usando la estructura grafica como componente adicional. De esta forma, se puede formular la **propiedad ordenada de Markov** en modelos gráficos dirigidos:

$$(1.3) \quad x_s \perp \mathbf{x}_{pred(s) \setminus pa(s)} \mid \mathbf{x}_{pa(s)}$$

Es decir, un nodo  $x_s$  es independiente de aquellos predecesores, menores en orden topológico, a sus padres  $\mathbf{x}_{pred(s) \setminus pa(s)}$ , dados sus nodos padres  $\mathbf{x}_{pa(s)}$ . De manera equivalente, un nodo  $x_s$  solo depende de sus padres inmediatos  $\mathbf{x}_{pa(s)}$  y no de todos sus predecesores.

De esta forma, la probabilidad conjunta de un modelo gráfico dirigido, que cumple la propiedad ordenada de Markov, se puede descomponer de la forma:

$$(1.4) \quad p(\mathbf{x}) = \prod_{t=1}^V p(x_t \mid \mathbf{x}_{pa(t)})$$

**EJEMPLO 1.1** (Modelo grafico asociado a  $p(\mathbf{x})$ ). *Si se estudia un modelo probabilístico, donde la probabilidad conjunta de las variables estudiadas  $p(\mathbf{x})$  esta dada por:*

$$(1.5) \quad p(\mathbf{x}) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3)p(x_5|x_1, x_3)p(x_6|x_4, x_5)$$

*Entonces, un grafo dirigido asociado a tal factorización es el de la figura 1(a). Para construir dicho grafo, se consideran las relaciones de independencia condicional en la factorización (1.5), para luego establecer aristas  $s \rightarrow t$  si la probabilidad condicional del nodo  $x_s$  depende de  $x_t$ . En este caso, no hay aristas incidentes hacia  $x_1$ ,  $x_2$  ni  $x_3$ . Por otra parte, se deben crear aristas desde  $x_1, x_2$  y  $x_3$  hacia  $x_4$ , desde  $x_1$  y  $x_3$  hacia  $x_5$  y desde  $x_4, x_5$  hacia  $x_6$ .*

*En general, es posible reconstruir la probabilidad conjunta subyacente a un modelo gráfico probabilístico conociendo el grafo y haciendo el proceso inverso al descrito anteriormente.*



Con el fin de explorar las posibilidades de este tipo de modelos e introducir conceptos referentes a la notación de estos, se pasan a estudiar los siguientes ejemplos:

**1.1.1. Naive Bayes.** Dado un problema de clasificación de vectores  $\mathbf{x} = (x_1, \dots, x_V)$  en  $C$  clases. Es posible modelar las variables de decisión  $x_t$  como condicionalmente independientes dada la categoría de clasificación:

$$(1.6) \quad x_i \perp x_j \mid y = c, \quad i \neq j$$

Si se usa este enfoque, se obtiene que la densidad condicional de clases toma la forma:

$$(1.7) \quad p(\mathbf{x} \mid y = c) = \prod_{t=1}^V p(x_t \mid y = c)$$

Al parametrizar las distribuciones de densidad condicional, es posible obtener un modelo de clasificación conocido como **clasificador naive Bayes**. La estructura de las relaciones de independencia inducidas por (1.6) se pueden expresar según (1.7) y el modelo gráfico dirigido de la figura 1(b).

**1.1.2. Regresión polinomial.** las variables aleatorias son el vector de coeficientes polinomiales  $\mathbf{w}$  y los datos observados  $\mathbf{y} = (y_1, \dots, y_N)^T$ . Adicionalmente, se parametriza el ruido del modelo a través de  $\sigma_\varepsilon^2$  y la varianza de la distribución a priori <sup>2</sup> de  $\mathbf{w}$  por  $\sigma_w^2$ . Finalmente, los datos de entrada se denotan por  $\mathbf{x} = (x_1, \dots, x_N)^T$ .

La probabilidad conjunta de este modelo es el producto de la probabilidad a priori  $p(\mathbf{w})$  con las distribuciones condicionales  $p(y_i \mid \mathbf{w})$  para  $i = 1, \dots, N$ :

$$(1.8) \quad p(\mathbf{y}, \mathbf{w}) = p(\mathbf{w}) \prod_{i=1}^N p(y_i \mid \mathbf{w})$$

El grafo de tal factorización es similar al del clasificador naive Bayes 1(b). Para representarlo de manera compacta, se usa la notación de de placas o *plates*, en la figura 2(a) se muestra el grafo de (1.8) usando esta convención. Aquí  $N$  es la cantidad de nodos del modelo, de los cuales se muestra el representante  $y_i$ .

Si por otra parte, si se quiere estudiar la interacción de los parámetros en el modelo, es posible explicitarlos en la probabilidad conjunta para luego agregarlos al grafo:

$$(1.9) \quad p(\mathbf{y}, \mathbf{w} \mid \mathbf{x}, \sigma_\varepsilon^2, \sigma_w^2) = p(\mathbf{w} \mid \sigma_w^2) \prod_{i=1}^N p(y_i \mid \mathbf{w}, x_i, \sigma_\varepsilon^2)$$

La figura 2(b) muestra el grafo correspondiente a (1.9). Por convención, las variables deterministas se incluyen en el grafo como círculos pequeños, mientras que las variables aleatorias observadas se muestran como nodos grises, los nodos incoloros representan variables latentes o no observadas, finalmente las aristas, al igual que en los ejemplos anteriores, representan la dependencia condicional en la factorización de la probabilidad conjunta.

---

<sup>2</sup>Considerándose una distribución a priori, gaussiana y esférica sobre  $\mathbf{w}$ .

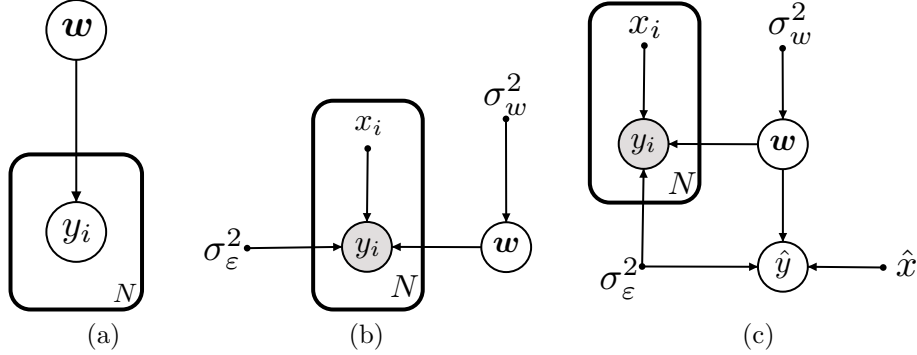


FIGURE 2. Modelo grafico dirigido para regresión polinomial usando notación de placas (o *plates*). En (a) se muestra el grafo correspondiente a (1.8). En (b) se añaden los parámetros deterministas y las variables aleatorias observadas. En (c) se añaden datos de entrada y predicciones.

Para realizar predicciones en datos nuevos  $\hat{x}$ , se desea encontrar la distribución de probabilidad para  $\hat{y}$  condicionada a la información que ya se posee. Esta corresponde a:

$$(1.10) \quad p(\hat{y}, \mathbf{y}, \mathbf{w} | \hat{x}, \mathbf{x}, \sigma_w^2, \sigma_\varepsilon^2) = \left[ \prod_{i=1}^N p(y_i | x_i, \mathbf{w}, \sigma_\varepsilon^2) \right] p(\mathbf{w} | \sigma_w^2) p(\hat{y} | \hat{x}, \mathbf{w}, \sigma_\varepsilon^2)$$

Finalmente, se deduce la distribución predictiva para  $\hat{y}$  :

$$(1.11) \quad p(\hat{y} | \hat{x}, \mathbf{y}, \mathbf{x}, \sigma_w^2, \sigma_\varepsilon^2) \propto \int p(\hat{y}, \mathbf{y}, \mathbf{w} | \hat{x}, \mathbf{x}, \sigma_w^2, \sigma_\varepsilon^2) d\mathbf{w}$$

El modelo gráfico dirigido que encapsula estas últimas ecuaciones se aprecia en 2(c).

**1.1.3. Modelos gráficos dirigidos gaussianos.** Sea  $\mathcal{M}$  un modelo grafico dirigido, en el cual todas las variables son reales y sus distribuciones de probabilidad condicional son lineal-gaussianas:

$$(1.12) \quad p(x_t | \mathbf{x}_{pa(t)}) = \mathcal{N}(x_t | \mu_t + \mathbf{w}_t^T \mathbf{x}_{pa(t)}, \sigma_t^2)$$

La estructura de  $\mathcal{M}$  permite modelar la probabilidad conjunta de las variables del modelo en la forma:

$$(1.13) \quad p(\mathbf{x} | \mathcal{M}) = \prod_{t=1}^V p(x_t | \mathbf{x}_{pa(t)}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Lo cual se conoce como **red bayesiana gaussiana**. Para este tipo de modelos, es posible inferir  $\boldsymbol{\mu}$  y  $\boldsymbol{\Sigma}$ . En efecto, según (1.13):

$$(1.14) \quad \log p(\mathbf{x} | \mathcal{M}) = - \sum_{t=1}^V \frac{1}{2\sigma_t^2} \left( x_t - \sum_{s \in pa(t)} w_{ts} x_s - \mu_t \right)^2 + K$$

Donde  $K$  representa una constante independiente de  $\mathbf{x}$ . Al ser la log-probabilidad conjunta, cuadrática en las componentes de  $\mathbf{x}$ , se obtiene que efectivamente la probabilidad conjunta es normal multivariada para  $\mathbf{x}$  en (1.13). Para estimar la media, se observa en primera instancia:

$$(1.15) \quad x_t = \sum_{s \in pa(t)} w_{ts} \mathbb{E}[x_s] + \mu_t + \sigma_t \varepsilon_t$$

Donde  $\varepsilon_t \sim \mathcal{N}(0, 1)$  y  $\mathbb{E}[\varepsilon_t, \varepsilon_s] = 0$ , para  $s \neq t$ . De esto se deduce:

$$(1.16) \quad \mathbb{E}[x_t] = \sum_{s \in pa(t)} w_{ts} \mathbb{E}[x_s] + \mu_t$$

Es posible entonces, encontrar las componentes de  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] = (\mathbb{E}[x_1], \dots, \mathbb{E}[x_V])^T$  utilizando la estructura gráfica dirigida de  $\mathcal{M}$  (y por tanto su ordenamiento topológico). Para ello, se comienza calculando  $\mathbb{E}[x_1]$  para luego continuar de manera recursiva según la numeración de los nodos.

Similarmente, es posible calcular el elemento  $\boldsymbol{\Sigma}_{st}$  de la matriz de covarianza, observando:

$$(1.17) \quad \begin{aligned} \text{cov}(x_s, x_t) &= \mathbb{E}[(x_s - \mathbb{E}[x_s])(x_t - \mathbb{E}[x_t])] \\ &= \left[ (x_s - \mathbb{E}[x_s]) \left\{ \sum_{k \in pa(x_t)} w_{tk} (x_k - \mathbb{E}[x_k]) + \sigma_t \varepsilon_t \right\} \right] \\ &= \sum_{k \in pa(x_t)} w_{tk} \text{cov}[x_s, x_t] + \sigma_t^2 \mathbf{I}_{st} \end{aligned}$$

De donde al igual que en (1.16), se calculan los elementos de  $\boldsymbol{\Sigma}$  recursivamente.

Finalmente, se puede extender el modelo inducido por (1.12) a uno donde los nodos del modelo gráfico representen variables aleatorias gaussianas multivariantes. Para esto, se reescribe la distribución de probabilidad condicional para el nodo  $x_t$  en la forma:

$$(1.18) \quad p(\mathbf{x}_t \mid pa(\mathbf{x}_t)) = \mathcal{N} \left( \mathbf{x}_t \mid \sum_{s \in pa(\mathbf{x}_t)} \mathbf{W}_{ts} \mathbf{x}_s + \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t \right)$$

Donde  $\mathbf{W}_{ts}$  es una matriz de pesos entre los vectores de cada nodo.

## 1.2. Independencia condicional en modelos gráficos dirigidos

Como se mencionó anteriormente, los modelos gráficos encapsulan las relaciones de independencia condicional entre las variables aleatorias del fenómeno que se modela. En esta sección se estudian las propiedades de los modelos gráficos dirigidos en cuanto a estas relaciones.

En un grafo  $G$ , se escribe  $x_i \perp_G x_j \mid x_k$  si el nodo  $x_i$  es independiente de  $x_j$  dado  $x_k$ . Se denota por  $I(G)$  al conjunto de todas las relaciones de independencia condicional codificadas en el grafo  $G$ .

**Definición 1.2** (Diagrama de independencia). Sea  $p(\cdot)$  una distribución de probabilidad. Se dice que un grafo  $G$  es un diagrama de independencia o *I-map* para  $p$  si y solo si  $I(G) \subseteq I(p)$ . Donde  $I(p)$  es el conjunto de todas las relaciones de independencia condicional ciertas para las variables de  $p$ .

De la definición anterior, se deduce que un grafo  $G$  es un diagrama de independencia para la distribución de probabilidad  $p$ , si este no contiene más relaciones de independencia condicional que las permitidas por  $p$ . De esta forma, toda distribución de probabilidad  $p(\mathbf{x})$ , donde  $\mathbf{x} = (x_1, \dots, x_V)^T$ , posee al menos un diagrama de independencia. En efecto, si se considera un grafo  $G$  con nodos  $\mathcal{V} = \{x_1, \dots, x_V\}$  completamente conectados, entonces  $G$  es un diagrama de independencia para  $p$  pues no presenta aristas faltantes y por tanto se condiciona en todas las variables.

De la discusión anterior, tiene sentido hablar de un *diagrama de independencia minimal*  $G$  para  $p$ , es decir, un grafo  $G$ , tal que si  $G'$  es otro diagrama de independencia para  $p$  que cumple  $G' \subseteq G$ , entonces  $G' = G$ .

Finalmente, tal representación, permite extraer de su estructura gráfica, relaciones no triviales de independencia condicional, entre las variables de importancia. En el caso de un modelo gráfico dirigido, la noción de *separación dirigida* o *d-separación* facilita dicha tarea.

**1.2.1. d-separación.** Se dice que un camino no dirigido  $P$  está *separado de manera dirigida* o *d-separado* por un conjunto de nodos  $E$ , si y solo si, se cumple alguna de las siguiente condiciones:

- (1)  $P$  contiene una cadena,  $s \rightarrow e \rightarrow t$ , donde  $e \in E$ .
- (2)  $P$  contiene una estructura  $s \leftarrow e \rightarrow t$ , donde  $e \in E$ .
- (3)  $P$  contiene una estructura  $s \rightarrow e \leftarrow t$ , donde  $e \notin E$  o  $e$  **no** es descendiente de algún elemento de  $E$ .

Se dice que un conjunto de nodos  $A$  está d-separado de un conjunto de nodos  $B$ , dado un conjunto de nodos  $E$ , si y solo si, todo camino no dirigido desde cada nodo de  $A$  a cada nodo de  $B$  está d-separado por  $E$ .

En un grafo acíclico dirigido  $G$ , se aprecia la siguiente propiedad:

Para comprender las propiedades anteriores, se analizan los siguientes ejemplos:

- Sea  $x \rightarrow y \rightarrow z$  una cadena, tal grafo codifica la siguiente probabilidad conjunta:

$$p(x, y, z) = p(x)p(y|x)p(z|y)$$

Usando la propiedad (1), se puede deducir que  $x \perp z|y$ . Esto se comprueba pues:

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x)p(y|x)p(z|y)}{p(y)} = \frac{p(x, y)p(z|y)}{p(y)} = p(x|y)p(z|y)$$

- Sea la estructura  $x \leftarrow y \rightarrow z$ , según la propiedad (2),  $x \perp z|y$ . En efecto,

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(y)p(x|y)p(z|y)}{p(y)} = p(x|y)p(z|y)$$

- Sea finalmente la estructura  $x \rightarrow y \leftarrow z$ , en este caso  $x \not\perp z|y$ :

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x)p(y|x,z)p(z)}{p(y)} \neq p(x|y)p(z|y)$$

Es por tal motivo que en la propiedad (3), se requiere en este tipo de estructuras, que no existan nodos, ni descendientes de nodos de la familia condicionante  $E$ .

En el último caso, se puede comprobar que los nodos  $x$  e  $y$  son marginalmente independientes entre sí, es decir,  $p(x, z) = p(x)p(z)$ . Sin embargo, al condicionar ambos nodos por  $y$ , se vuelven dependientes, este efecto se denomina **paradoja de Berkson**. Finalmente, un modelo gráfico probabilístico que verifica la equivalencia (??) se dice que cumple la **propiedad global de Markov**.

**1.2.2. Markov blankets.** De la d-separación es posible concluir:

$$(1.19) \quad x_t \perp \mathbf{x}_{nd(t) \setminus pa(t)} | \mathbf{x}_{pa(t)}$$

Donde  $nd(t)$  son los **no-descendientes** del nodo  $x_t$ <sup>3</sup>. De esta forma, es posible concluir que en 1(a)  $x_4 \perp x_5 | x_1, x_2, x_3$ , pues en efecto,  $nd(4) \setminus pa(4) = x_5$  y  $pa(4) = x_1, x_2, x_3$ . La ecuación (1.19) se conoce como **propiedad dirigida local de Markov**.

En especial, dado que  $pred(t) \subseteq nd(t)$  se deriva la **propiedad ordenada de Markov**, ya presentada en (1.3). Sorprendentemente, estas tres propiedades son equivalentes.

Por otra parte, para cada nodo  $x_t$  es posible extraer el conjunto de nodos que lo separan del resto del grafo, es decir, se puede para cada nodo  $x_t$ , obtener el conjunto de todas las variables aleatorias que lo vuelven condicionalmente independiente a los demás nodos del modelo. El conjunto antes descrito se denomina **Markov blanket** y se denota por  $mb(t)$  este conjunto de nodos corresponde a:

$$(1.20) \quad mb(t) := ch(t) \cup pa(t) \cup copa(t)$$

Donde  $ch(t)$  son los nodos hijos de  $x_t$ , de manera análoga  $pa(t)$  son nodos padres y  $copa(t)$  sus copadres<sup>4</sup>. En la figura 1(a) se tiene por ejemplo  $mb(5) = \{x_6, x_1, x_3, x_4\}$ . Según la propiedad global de Markov, la presencia de los nodos copadres no parece ser necesaria en primera instancia (la dependencia condicional debería recaer únicamente en los nodos padres), sin embargo, al definir  $\mathbf{x}_{-t}$  como el conjunto de nodos distintos a  $x_t$ , es posible observar que la probabilidad conjunta adquiere la forma  $p(\mathbf{x}) = p(x_t, \mathbf{x}_{-t})$ . De donde, al marginalizar sobre el nodo  $x_t$ , se obtiene que  $p(\mathbf{x}_{-t})$  contiene sólo a aquellos nodos del modelo en los que la variable  $x_t$  no aparece como argumento, ni como condicionante (dada la factorización de la probabilidad conjunta codificada en el grafo). Lo anterior implica que en  $p(x_t | \mathbf{x}_{-t}) = p(\mathbf{x}) / p(\mathbf{x}_{-t})$  solo se podrán encontrar probabilidades condicionales donde  $x_t$  sea el argumento, lo que expresa con  $p(x_t | \mathbf{x}_{pa(t)})$ , o donde sea variable condicionante, es decir, sea padre o copadre de algún otro nodo. Se deduce:

$$(1.21) \quad p(x_t | \mathbf{x}_{-t}) \propto p(x_t | \mathbf{x}_{pa(t)}) \prod_{s \in ch(t)} p(x_s | \mathbf{x}_{pa(s)})$$

La expresión (1.21) se conoce como **condicional completa** del nodo  $x_t$ .

<sup>3</sup> $nd(t) = \mathcal{V} \setminus \{t \cup desc(t)\}$ , donde  $desc(t)$  son los descendientes del nodo  $x_t$ , es decir, aquellos nodos que provienen de un camino dirigido con origen en  $x_t$ .

<sup>4</sup> nodos que comparten hijos con  $x_t$



## CHAPTER 2

# Métodos de inferencia aproximada

### Introducción

En este capítulo se estudian métodos de inferencia aproximada, en particular se abordan algoritmos basados en Markov Chain Monte Carlo e inferencia variacional. La idea detrás de la inferencia aproximada se basa en la obtención de muestras de una distribución de la forma  $\mathbf{x}^s \sim p(\mathbf{x}|\mathcal{D})$  para calcular por ejemplo la distribución posterior predictiva de cierto modelo  $p(y|\mathcal{D})$  o cierta marginal posterior  $p(x_1|\mathcal{D})$ . Cuando no se conocen las formas cerradas de estos objetos, es necesario por tanto aproximarlos de manera tal que se tengan garantías sobre las cantidades aproximadas que se obtienen.

En términos generales, se considera una densidad conjunta de variables latentes  $\mathbf{z} = z_1, \dots, z_m$  y observaciones  $\mathbf{x} = x_1, \dots, x_n$  dada por,

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}).$$

Desde la perspectiva bayesiana, las variables latentes provienen de una densidad prior  $p(\mathbf{z})$ . Esta se relaciona con las observaciones a través de  $p(\mathbf{x}|\mathbf{z})$ . En este contexto, el problema de inferencia equivale a condicionar los datos y calcular la posterior  $p(\mathbf{z}|\mathbf{x})$ . Los métodos explorados a continuación, abordan tal problema a través de la aproximación de la probabilidad posterior, ya sea de manera numérica o a través de distribuciones base.

### 2.1. Inferencia Markov Chain Monte Carlo

En esta sección se estudian métodos de Monte Carlo basados en cadenas de Markov. Este tipo de método permite obtener representaciones numéricas (muestras) de distribuciones incluso de alta dimensionalidad.

### 2.2. Cadenas de Markov

Una cadena de Markov es una sucesión de variables aleatorias  $x_1, \dots, x_n$  que verifican  $x_{t+1} \perp x_1, \dots, x_{t-1} | x_t$ . Es decir, la distribución condicional de  $x_{t+1}$ , que puede ser interpretada como un estado futuro, depende únicamente del estado presente  $x_t$ . Si la probabilidad de transición es independiente de  $t$ , la cadena se denomina Markov homogénea. Tal tipo de cadenas se define especificando las probabilidades de transición de un estado a otro. Para el caso continuo, esto se traduce en:

$$(2.1) \quad K(x, y) = p(x_{t+1} = y | x_t = x)$$

En el caso discreto, se obtiene una matriz de transición  $K_{x,y}$ . Dado cierto espacio de estado, una cadena de Markov homogénea en el tiempo posee una distribución

estacionaria  $\pi$  si:

$$(2.2) \quad \pi(y) = \int \pi(x)K(x, y)dx$$

Una cadena de MARKOV se dice irreducible, si puede transitar desde cualquier estado  $x$  de un espacio de estado discreto a cualquier otro estado  $y$  en un numero finito de pasos, es decir, si existe  $t$  tal que  $K_{xy}^n > 0$ . Si una cadena con distribución estacionaria es irreducible, la distribución estacionaria es unica y la cadena se dice recurrente positiva. Una cadena recurrente positiva, aperiodica con distribución estacionaria  $\pi$ , cumple para toda distribución inicial  $\lambda$  sobre sus estados:

$$(2.3) \quad \lim_{n \rightarrow \infty} \|\lambda K^n - \pi\| = 0$$

donde  $K$  es el operado (o matriz) de transición. Para cadenas de Markov irreducibles con una única distribución estacionaria  $\pi$ , según la ley de los grandes números se debe cumplir que el valor esperado de una función  $g(x)$  sobre  $\pi$  verifica:

$$(2.4) \quad \mathbb{E}_\pi[g(x)] = \int g(x)\pi(x)dx = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n g(x_i)$$

Esta propiedad permite estimar calcular aproximaciones de cantidades específicas de interés a partir de una cadena de Markov. Los algoritmos que utilizan tal resultado se conoce como métodos MCMC.

Una cadena con distribución estacionaria  $\pi$  se dice reversible si el par  $(x_t, x_{t+1})$  tiene la misma distribución conjunta que  $(x_{t+1}, x_t)$ . En terminos del operador de transición  $K$ , esto se traduce como

$$(2.5) \quad \pi(x_t) K(x_t, x_{t+1}) = \pi(x_{t+1}) K(x_{t+1}, x_t)$$

Una cadena de Markov no necesariamente debe ser reversible para poseer una distribución estacionaria. Sin embargo, la reversibilidad de una cadena, garantiza la existencia de una distribución estacionaria. Esta es la razon por la cual se busca que en la mayoría de los algoritmos MCMC se cumpla esta propiedad (*balance detallado*).

### 2.3. Algoritmo Metropolis Hastings (MH)

El algoritmo canónico MCMC se conoce como Metropolis-Hastings. Su formulación pasa por obtener muestras de una distribución  $p(x)$  en un espacio de estados  $E$  donde  $x \in E$ . Se construye un kernel de transición  $k(x, y)$  para pasar del estado  $x$  al  $y$  mediante dos procesos: en primer lugar, se especifica una distribución  $q(y|x)$  propuesta (proposál). Para luego en segunda instancia obtener muestras de  $q(y|x)$  con un coeficiente de aceptación  $\alpha(x, y) = \min \left[ 1, \frac{p(y)q(x|y)}{p(x)q(y|x)} \right]$ . De esta forma, el kernel de transición queda definido como  $K(x, y) = q(y|x)\alpha(x, y)$ . El algoritmo se describe a continuación:

El kernel de transición correspondiente al algoritmo MH es reversible y por tanto cumple con la propiedad de balance detallado,  $p(x)K(x, y) = p(y)K(y, x)$ .

Existen múltiples formas de construir distribuciones proposál  $q$ , según la elección de esta, se construyen distintas versiones del algoritmo MH.

**2.3.1. MH caminata aleatoria (RWMH).** Seleccionando  $q(y|x) = q(y-x)$ , ocurre que la dirección y distancia de un nuevo punto con respecto al punto actual, es independiente del punto actual. Las distribuciones que cumplen esta condición de manera inmediata son  $\mathcal{N}(x, \sigma^2)$  y  $U(x - \sigma, x + \sigma)$



**Algorithm 1:** Algoritmo Metropolis Hastings

---

**Input:** Punto inicial  $x_1$ ,  $p(x)$ , kernel de transición  $q(y|x)$   
**Output:** Vector de  $N$  puntos  $x_1, \dots, x_N$

```

1 for  $t = 1, \dots, N-1$  do
2   | Obtener una muestra  $y$  de  $q(y|x_t)$ 
3   | Obtener una muestra de una v.a uniforme  $U$ 
4   | if  $U < \frac{p(y)q(x_t|y)}{p(x_t)q(y|x_t)}$  then  $x_{t+1} = y$ 
5   | else  $x_{t+1} = x_t$ 
6 end

```

---

**2.3.2. Muestreo de independencia.** Seleccionando  $q(y|x) = q(y)$ , es decir, el nuevo estado es independiente del estado actual. el coeficiente de aceptación en este caso corresponde a  $\min \left\{ 1, \frac{p(y)q(x)}{p(x)q(y)} \right\}$ . En general se busca que  $q(x)$  sea similar a  $p(x)$  pero con colas más pesadas.

**2.3.3. Metropolis simétrico.** Seleccionando  $q(y|x) = q(x|y)$ , esta opción simplifica la probabilidad de aceptación a  $\min\{1, f(y)/f(x)\}$ .

**2.4. Muestreo de Gibbs**

El muestreo de Gibbs, supone una distribución  $p(x)$  de la cual se desea obtener muestras, donde  $x \in \mathbb{R}^d$ . En este método, el kernel de transición  $K(x, y)$  es separado en múltiples pasos. En cada paso, se actualiza el valor de una coordenada basada en la densidad condicional con respecto a la demás coordenadas. Esto se observa en el siguiente algoritmo:

**Algorithm 2:** Algoritmo de Muestreo de Gibbs

---

**Input:** Punto inicial  $x_1$ , y distribución  $p(x)$   
**Output:** Vector de  $N$  puntos  $x_1, \dots, x_N$

```

1 for  $t = 1, \dots, N-1$  do
2   | Obtener una muestra  $x_1^{t+1} \sim p(x_1^{t+1} | \mathbf{x}_{-1}^t)$ 
3   | Obtener una muestra  $x_2^{t+1} \sim p(x_2^{t+1} | \mathbf{x}_{-2}^t)$ 
4   |  $\vdots$ 
5   | Obtener una muestra  $x_d^{t+1} \sim p(x_d^{t+1} | \mathbf{x}_{-d}^t)$ 
6 end

```

---

El kernel de transformación se puede escribir como:

$$(2.6) \quad K_{1 \rightarrow d}(x_{t+1}|x_t) = \prod_{i=1}^d p(x_{i+1}^i | x_{i+1}^1, \dots, x_{i+1}^{i-1}, x_t^{i+1}, \dots, x_{t+1}^d)$$

De esta relación, se puede observar que

$$(2.7) \quad p(x_t) K_{1 \rightarrow d}(x_{t+1}|x_t) = p(x_{t+1}) K_{d \rightarrow 1}(x_t|x_{t+1})$$

De donde al integrar en ambos lados

$$(2.8) \quad \int p(x_t) K_{1 \rightarrow d}(x_{t+1}|x_t) dx = p(y)$$

Luego,  $p$  es la distribución estacionaria de la cadena de Markov formada por el kernel de transición  $K(x_{t+1}|x_t)$ . El método de muestreo de Gibbs con el kernel propuesto no es reversible. Otra característica de este método, es que puede verse como un caso especial del método MH con coeficiente de aceptación  $\min\left(1, \frac{p(y)q(x|y)}{p(x)q(y|x)}\right) = 1$ , de esto, al observar  $\mathbf{y}_{-i} = \mathbf{x}_{-i}$  se tiene:

$$(2.9) \quad \begin{aligned} p(\mathbf{y})q(\mathbf{x}|\mathbf{y}) &= p(y_i|\mathbf{y}_{-i})p(\mathbf{y}_{-i})p(x_i|\mathbf{y}_{-i}) = p(y_i|\mathbf{x}_{-i})p(\mathbf{x}_{-i})f(x_i|\mathbf{x}_{-i}) \\ &= p(x_i|\mathbf{x}_{-i})p(\mathbf{x}_{-i})p(y_i|\mathbf{x}_{-i}) = p(\mathbf{x})q(\mathbf{y}|\mathbf{x}) \end{aligned}$$

Finalmente, cabe destacar que la falencia del Muestreo Gibbs recae en que en cada paso, se deben obtener muestras de la distribución conjunta de las demás variables, este paso puede ser costoso desde el punto de vista computacional. Una solución a este problema consiste en aplicar un muestreo de esta probabilidad condicional aplicando MH, esto da lugar al algoritmo MWG (Metropolis within Gibbs).

## 2.5. Convergencia en métodos MCMC

Si bien la formulación de los algoritmos MCMC estudiados hasta el momento garantizan convergencia a la distribución que se desea muestrear, no se especifica a priori el número de iteraciones necesarias para alcanzar esta convergencia. Para abordar este problema se estudian esquemas de diagnóstico capaces de detectar fallas en la convergencia en este tipo de algoritmos.

**2.5.1. Tamaño efectivo de muestra.** Sea  $p(\mathbf{x})$  la densidad de la variable aleatoria  $\mathbf{x}$  con desviación estándar  $\sigma_x$ , si se obtienen  $n$  muestras independientes, el error estándar Monte Carlo se aproxima a  $\sigma/\sqrt{n}$ . De esta forma, para medir la media de una cantidad con un error cercano al 3%, con respecto a la incertidumbre  $\sigma_x$ , es necesario obtener alrededor de  $n = 1000$  muestras.

Debido a la correlación existente en los puntos de una cadena MCMC, el cálculo anterior no se verifica (las muestras no son independientes). No obstante, es posible medir que tan correlacionados están estos puntos entre sí, por medio del cálculo de la autocorrelación, esta se define por:

$$(2.10) \quad \rho_{xx}(t) = \frac{\mathbb{E}[(x_i - \bar{x})(x_{i+t} - \bar{x})]}{\mathbb{E}[(x_i - \bar{x})^2]}$$

la cual se aplica sobre dos puntos separados por una distancia fija (o lag)  $t$ . Es de esperar, que la autocorrelación disminuya con respecto a la distancia  $t$  con que se comparan los puntos, de hecho el valor de la autocorrelación decae de manera proporcional a  $\exp(-t/\tau_x)$  donde  $\tau_x$  se conoce como la *longitud de correlación*. En este contexto, se define la *correlación integrada* como  $\tau_{int,x} = \sum_t \rho_{xx}(t)$ . La variancia de la media  $\bar{x}$  de  $\mathbf{x}$ , para una muestra de tamaño  $n$  se puede escribir en función de la correlación integrada obteniéndose:

$$(2.11) \quad \text{Var}(\bar{x}) = (2 \tau_{int,x}) \frac{\mathbb{E}[(x_i - \bar{x})^2]}{n}$$

De esta forma, para muestras correlacionadas, la variancia pasa a ser 2  $\tau_{int,x}$  veces mayor con respecto a la variancia sobre muestras independientes. Usando

$\tau_{int,x}$  es posible medir el *numero efectivo* de muestras independientes en una cadena correlacionada calculado  $n/(2 \tau_{int,x})$  para luego usar esta cantidad para decidir si la cantidad de muestras es suficiente, en el caso de la media se busca como heurística un numero efectivo de muestras mayor a 1000.

**2.5.2. Varianza entre cadenas.** Sean  $M$  cadenas consistentes de  $2N$  iteraciones, de las cuales se utilizan las últimas  $N$  muestras. Dado un modelo probabilístico un parámetro de interés  $\theta$ , sea  $[\theta_{mt}]_{t=1}^N$  la cadena  $m$ -ésima para  $m = 1, \dots, M$ . En este contexto, sean  $\hat{\theta}_m$  y  $\hat{\sigma}_m^2$  la media muestral posterior y la variancia de la cadena  $m$ -ésima, sea además, la media posterior muestral total  $\hat{\theta} = (1/M) \sum_{m=1}^M \hat{\theta}_m$ . La variancia entre cadenas (B) e intra cadena (W) están dadas por

$$(2.12) \quad B = \frac{N}{M-1} \sum_{m=1}^M (\hat{\theta}_m - \hat{\theta})^2 \quad \text{y} \quad W = \frac{1}{M} \sum_{m=1}^M \hat{\sigma}_m^2$$

El estimador insesgado  $\hat{V}$  para la variancia posterior marginal de  $\theta$  se obtiene bajo ciertas condiciones de estacionalidad [16] por medio de:

$$(2.13) \quad \hat{V} = \frac{N-1}{N} W + \frac{M+1}{MN} B$$

El factor potencial de reducción de escala (PSRF) se define como el coeficiente entre  $\hat{V}$  y  $W$ . Si las  $M$  cadenas convergen a la distribución posterior, entonces el factor PSRF debe ser cercano a 1. Este factor estima la disminución potencial para la variabilidad entre cadenas  $B$  con respecto a la variabilidad intra cadenas  $W$ .

Utilizando una formulación similar [17], es posible obtener una corrección a este índice:

$$(2.14) \quad R_c = \sqrt{\frac{\hat{d} + 3}{\hat{d} + 1} \frac{\hat{V}}{W}}$$

donde  $\hat{d}$  representa los grados de libertad estimados de una distribución  $t$  de Student. Si  $R_c$  crece, a medida que las cadenas son iteradas una mayor cantidad de veces, se espera que tanto  $B$  disminuya o  $W$  crezca. Según se sugiere en [16] si  $R_c < 1.2$  para todos los parámetros de interés, se puede considerar que se ha alcanzado convergencia.

## 2.6. Monte Carlo Hamiltoniano

Una ventaja de MCMC sobre otro métodos de Monte Carlo, es su rendimiento en muestreo de probabilidades en altas dimensiones. Sin embargo, en presencia de dimensionalidad extremadamente alta, los algoritmo tradicionales MCMC comienzan a presentar problemas. La formulación HMC (Monte Carlo Hamiltoniano) introduce una variable auxiliar de *momentum*  $u$  para para cada variable  $x$ . La distribución log posterior  $\pi(x)$  se considera como el potencial de energía  $U(x) = -\ln \pi(x)$  donde el momentum define la energía cinética  $K(u)$ . De esta forma se define el *Hamiltoniano*  $H(x, u) = U(x) + K(u)$ , donde  $K(u) = u^2/2$ . La distribución a explorar corresponde a

$$(2.15) \quad p(x, u) = \exp[-H(x, u)] = \exp \left[ -\ln \pi(x) - \frac{1}{2}u^2 \right]$$

Siguiendo las ecuaciones de la dinámica hamiltoniana [18], se encuentra un nuevo punto/posición  $x'$  que finalmente es aceptado o rechazado en base al algoritmo MH. El uso de la dinámica hamiltoniana permite explorar área en posiciones más lejanas a la posición actual.

Cabe destacar que si bien HMC aborda el problema de muestreo en altas dimensiones. para su implementación se requiere el gradiente de la densidad buscada, además de dos parámetros extra que deben ser obtenidos por el usuario.

### 2.7. Inferencia variacional

En esta sección, se estudia un método de aproximación de distribuciones conocido como inferencia variacional (**VI** por sus siglas en inglés). En general, la aproximación por inferencia variacional tiende a ser más rápida y fácil de escalar a datos de gran tamaño en comparación al muestreo por **MCMC**. Por tal motivo, se ha hecho presente en análisis de documentos a gran escala y problemas de neurociencia computacional.

A diferencia de los métodos ya estudiados, la idea principal detrás de la inferencia variacional consiste en usar optimización sobre una familia  $\mathcal{Q}$  de densidades aproximadas sobre las variables latentes  $\mathbf{z} = z_1, \dots, z_m$ . Basándose en esto, se busca al miembro (distribución) de esa familia que minimice cierta noción de distancia o disimilitud  $D(\cdot, \cdot)$ <sup>1</sup> en el espacio de distribuciones, con respecto a la densidad posterior ideal,

$$(2.16) \quad q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} D(q(\mathbf{z}), p(\mathbf{z}|\mathbf{x})).$$

Finalmente, se utiliza la distribución aproximada  $q^*(\cdot)$  para aproximar la posterior buscada. De esta forma, el problema de inferencia pasa a ser un problema de optimización, donde la complejidad se ve regulada por la estructura  $\mathcal{Q}$ . Por tanto, se busca modelar  $\mathcal{Q}$  aumentando su expresividad, de manera tal, que se garantice una buena aproximación de  $p(\mathbf{z}|\mathbf{x})$ , pero a la vez, se busca disminuir su complejidad en cuanto al gasto computacional.

### 2.8. Descripción del método

La idea clave tras el proceso de aproximación por inferencia variacional, pasa por resolver su problema de optimización subyacente. Desde este punto de vista,  $\mathcal{Q}$  se puede escoger como una familia de funciones parametrizadas, por tanto, el problema de encontrar la función óptima  $q^*(\cdot)$  se transforma en encontrar los parámetros que la caracterizan. Para definir esta optimalidad, en la aproximación, se utiliza la divergencia de *Kullback Leibler* (KL), dada por:

$$(2.17) \quad \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \mathbb{E}[\log q(\mathbf{z})] - \mathbb{E}[\log p(\mathbf{z}|\mathbf{x})],$$

La densidad variacional obtenida de esta forma, sirve para representar la densidad condicional exacta.

En el contexto inicial, si  $\mathbf{x} = x_1, \dots, x_n$  es un conjunto de variables observadas y  $\mathbf{z} = z_1, \dots, z_m$  variables latentes, con densidad conjunta  $p(\mathbf{z}, \mathbf{x})$ . El problema de inferencia, consiste en calcular la densidad condicional de las variables latentes dadas las observaciones,  $p(\mathbf{z}|\mathbf{x})$ , esta se puede escribir como:

$$(2.18) \quad p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})}.$$

Donde el denominador contiene la densidad marginal de las observaciones, también llamada la *evidencia*. Al marginalizar sobre las variables latentes, se obtiene

---

<sup>1</sup>No se requiere estrictamente que  $D(\cdot, \cdot)$  sea una métrica.

$$(2.19) \quad p(\mathbf{x}) = \int p(\mathbf{z}, \mathbf{x}) d\mathbf{z}.$$

La dificultad de este proceso, recae en el calculo de esta integral, pues por lo general no se tiene su forma cerrada o se requiere de recursos computacionales prohibitivos para obtenerla. A modo de ejemplo se estudia el siguiente modelo:

**2.8.1. GMM Bayesiano I.** Se considera un **GMM** (modelo de mezcla de Gaussianas) de  $C$  componentes, univariado y de varianza unitaria. Cada una de sus componentes corresponden a distribuciones Gaussianas con medias  $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_C\}$ .

Los parámetros correspondientes a las medias se obtienen independientemente según su distribución prior  $p(\mu_k)$ , que se supondrá  $\mathcal{N}(0, \sigma^2)$ , donde la varianza prior  $\sigma^2$  corresponde a un hiperparámetro.

Para generar una observación  $x_i$  del modelo, es necesario generar una asignación de cluster, esta se denota por  $c_i$  e indica de qué cluster latente  $x_i$  proviene. Posteriormente, el valor de  $c_i$  se obtiene de una distribución categórica sobre  $\{1, \dots, C\}$ , siendo  $c_i$  un vector indicador de dimensión  $C$ . Finalmente, se asigna a  $x_i$  un valor obtenido de  $\mathcal{N}(c_i^\top \boldsymbol{\mu}, 1)$ . Esto se resume según el siguiente esquema:

$$(2.20) \quad \begin{aligned} \mu_k &\sim \mathcal{N}(0, \sigma^2), & k &= 1, \dots, C, \\ c_i &\sim \text{Cat}(1/k, \dots, 1/k), & i &= 1, \dots, n, \\ x_i | c_i, \boldsymbol{\mu} &\sim \mathcal{N}(c_i^\top \boldsymbol{\mu}, 1) & i &= 1, \dots, n. \end{aligned}$$

Para una muestra de tamaño  $n$ , la densidad conjunta de variables latentes y observadas es

$$(2.21) \quad p(\boldsymbol{\mu}, \mathbf{c}, \mathbf{x}) = p(\boldsymbol{\mu}) \prod_{i=1}^n p(c_i) p(x_i | c_i, \boldsymbol{\mu}).$$

En este modelo, las variables latentes corresponden a  $\mathbf{z} = \{\boldsymbol{\mu}, \mathbf{c}\}$ . De este modo, la evidencia se expresa por:

$$(2.22) \quad p(\mathbf{x}) = \int p(\boldsymbol{\mu}) \prod_{i=1}^n \sum_{c_i} p(c_i) p(x_i | c_i, \boldsymbol{\mu}) d\boldsymbol{\mu}.$$

El integrando en la ecuación (2.22) no se puede reducir a un producto de integrales unidimensionales sobre los  $\mu_k$ . Más aún, la complejidad del tiempo de evaluar numéricamente esta integral  $C$ -dimensional es  $\mathcal{O}(k^n)$ .

Si por otra parte, se distribuye el producto sobre la suma en (2.22), es posible escribir la evidencia como:

$$(2.23) \quad p(\mathbf{x}) = \sum_{\mathbf{c}} p(\mathbf{c}) \int p(\boldsymbol{\mu}) \prod_{i=1}^n p(x_i | c_i, \boldsymbol{\mu}) d\boldsymbol{\mu}.$$

Donde debido a la conjugación entre las distribuciones (prior gaussiana), cada integral individual pasa a ser computable. Sin embargo, se deben calcular  $k^n$  integrales

de este tipo: una para cada configuración de las asignaciones de cluster. El calculo de la evidencia permanece exponencial en  $C$  y por tanto intratable.

### 2.9. Cota inferior para la evidencia

Para abordar problemas como el anterior desde el punto de vista de inferencia variacional, se especifica una familia  $\mathcal{Q}$  de densidades sobre las variables latentes. Cada  $q(\mathbf{z}) \in \mathcal{Q}$  es una aproximación candidata densidad condicional exacta.

Para encontrar tal candidato, se minimiza la divergencia KL con respecto a la probabilidad condicional exacta, esto se traduce en resolver el siguiente problema:

$$(2.24) \quad q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})).$$

Esto pues, se requiere calcular:

$$(2.25) \quad \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \mathbb{E} [\log q(\mathbf{z})] - \mathbb{E} [\log p(\mathbf{z}|\mathbf{x})],$$

Que corresponde a

$$(2.26) \quad \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \mathbb{E} [\log q(\mathbf{z})] - \mathbb{E} [\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}).$$

Es decir, se requiere calcular la evidencia  $\log p(\mathbf{x})$ , que como se comprobó anteriormente, es por lo general difícil de calcular. Debido a que no se puede calcular la divergencia KL, se optimiza una función objetivo equivalente (hasta la suma una constante):

$$(2.27) \quad \text{ELBO}(q) = \mathbb{E} [\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E} [\log q(\mathbf{z})].$$

Esta función se denomina cota inferior para la evidencia (**ELBO** por sus siglas en ingles). De su definición, se observa que  $\text{ELBO}(q) = \log p(\mathbf{x}) - \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))$ . Cabe destacar, que  $p(\mathbf{x})$ , es una constante con respecto a  $q(\mathbf{z})$ .

Según lo anterior, maximizar la cota ELBO es equivalente a minimizar la divergencia KL. Para entender el problema de inferencia variacional, se estudian las propiedades de este funcional. Para ello, se reescribe de la siguiente forma:

$$\begin{aligned} \text{ELBO}(q) &= \mathbb{E} [\log p(\mathbf{z})] + \mathbb{E} [\log p(\mathbf{x}|\mathbf{z})] - \mathbb{E} [\log q(\mathbf{z})] \\ &= \mathbb{E} [\log p(\mathbf{x}|\mathbf{z}) - \text{KL}(q(\mathbf{z})||p(\mathbf{z}))]. \end{aligned}$$

Es decir, la cota ELBO corresponde a la resta entre log verosimilitud y la divergencia KL entre la distribución prior  $p(\mathbf{z})$  y  $q(\mathbf{z})$ .

El primer término de la expresión anterior la verosimilitud esperada por tanto, al maximizar la cota ELBO se están buscando densidades que distribuyan su masa en aquellos sectores donde las variables latentes explican mejor los datos observados. Por su parte, el segundo término es la divergencia negativa entre la densidad variacional y la prior, es decir, esta parte de la igualdad fomenta densidades cercanas a la prior. Se concluye que al maximizar la cota ELBO se busca ajuste a los datos, al mismo tiempo que se regulariza según la distribución prior.

Otra propiedad de la cota ELBO es que además acota inferiormente la log evidencia, es decir,  $\log p(\mathbf{x}) \geq \text{ELBO}(q)$  para todo  $q(\mathbf{z})$ . Para observar esto, se puede obtener de las ecuaciones (2.26) y (2.27) la siguiente expresión:

$$(2.28) \quad \log p(\mathbf{x}) = \text{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) + \text{ELBO}(q)$$

El resultado buscado se obtiene al observar que  $\text{KL}(\cdot) \geq 0$ .

### 2.10. Familia variacional mean-field

En cuanto a la familia  $\mathcal{Q}$  sobre la cual se busca aproximar, es necesario que su complejidad de sus complejidad permita resolver el problema de optimización subyacente. En esta contexto, se estudia la familia variacional *mean-field*, la cual se define de manera tal que las variables latentes son mutuamente independientes, cada una gobernada por un parametro distinto en la densidad variacional. Por tanto, si  $q \in \mathcal{Q}$  entonces  $q(\cdot)$  se podrá expresar como:

$$(2.29) \quad q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j).$$

De esta forma, latente  $z_j$  se rige por su propia componente variacional  $q_j(\cdot)$ .

En principio, no existe restricción para seleccionar el factor  $q_j(\cdot)$ , no obstante, su elección se relaciona directamente con el tipo de variable (continua/discreta) latente que se modela. Finalmente, los parámetros que caracterizan estos factores, definen al elemento de la familia  $\mathcal{Q}$ . Se continua con el ejemplo presentado anteriormente.

**2.10.1. GMM Bayesiano II.** Al modelo GMM bayesiano del ejemplo anterior, se agrega una familia variacional mean-field de la forma:

$$(2.30) \quad q(\boldsymbol{\mu}, \mathbf{c}) = \prod_{k=1}^C q(\mu_k; m_k, s_k^2) \prod_{i=1}^n q(c_i; \varphi_i)$$

En este caso, el factor  $q(\mu_k; m_k, s_k^2)$  se modela como una distribución gaussiana en la componente  $k$ -ésima con media  $m_k$  y varianza es  $s_k^2$ . El factor  $q(c_i; \varphi_i)$  corresponde a una distribución para la  $i$ -ésima asignación cluster y por tanto, sus probabilidades de asignación corresponden a un vector  $C$ -dimensional denotado por  $\varphi_i$ . De esta manera, se posee una familia paramétrica de distribuciones para los factores de  $q(\cdot)$  y se puede abordar la maximización de la cota **ELBO** mediante la ecuación (2.21), usando la familia mean-field de la ecuación (2.30). El problema de optimización variacional correspondiente maximiza la cota **ELBO** con respecto a los parámetros variacionales, es decir, los parámetros gaussianos para cada componente y los parámetros categóricos para cada asignación de cluster.

### 2.11. Inferencia variacional mean-field de coordenadas ascendentes

El método descrito anteriormente, permite por medio de la cota **ELBO** y una familia mean-field adecuada, establecer el problema de inferencia variacional y obtener su formulación como un problema de optimización. En esta sección, se describe el algoritmo de **inferencia variacional mean-field de coordenadas ascendentes** (CAVI por sus siglas en ingles), el cual permite resolver tal problema de optimización. El algoritmo **CAVI** optimiza iterativamente cada factor de la densidad variacional mean-field, mientras mantiene los demás fijos.



**2.11.1. Formulación.** Para la  $j$ -ésima variable latente  $z_j$ , se utiliza su *condicional completa*, es decir, su densidad condicional dadas las demás variables latentes y las observaciones  $p(z_j | \mathbf{z}_{-j}, \mathbf{x})$ . Manteniendo fijos los factores variacionales  $q_\ell(z_\ell)$ ,  $\ell \neq j$ , se obtiene que el factor óptimo  $q_j^*(z_j)$  verifica,

$$(2.31) \quad q_j^*(z_j) \propto \exp \{ \mathbb{E}_{-j} [\log p(\mathbf{z}_j | \mathbf{z}_{-j}, \mathbf{x})] \}.$$

El valor esperado en la ecuación (2.31) se calcula con respecto la densidad variacional de las variables  $\mathbf{z}_{-j}$  para todo  $j = 1, \dots, m$ , esta a su vez se expresa como  $\prod_{\ell \neq j} q_\ell(z_\ell)$ , de esta forma, se pueden incluir las variables  $\mathbf{z}_{-j}$  al calculo de la probabilidad condicional y excluir del calculo del valor esperado  $\mathbb{E}_j$ . Por tal motivo, se verifica la siguiente relación de proporcionalidad (sobre la conjunta):

$$(2.32) \quad q_j^*(z_j) \propto \exp \{ \mathbb{E}_{-j} [\log p(\mathbf{z}_j, \mathbf{z}_{-j}, \mathbf{x})] \}.$$

Debido a que todas las variables latentes son consideradas independientes (hipotesis mean-field) los valores esperados del lado derecho no involucran al  $j$ -ésimo factor variacional. Por tanto la ecuación anterior, permite actualizar los factores (coordenadas) de manera iterativa. A continuación se presenta este algoritmo:

---

**Algorithm 3:** CAVI

---

**Input:** Modelo  $p(\mathbf{x}, \mathbf{z})$  con datos  $\mathbf{x}$   
**Output:** Densidad variacional  $\prod_{j=1}^m q_j(z_j)$   
**1 Inicializar:** Factores variacionales  $q_j(z_j)$   
**2 while** ELBO *no converge* **do**  
**3     for**  $j = 1, \dots, m$  **do**  
**4         |** Calcular  $q_j^*(z_j) \propto \exp \{ \mathbb{E}_{-j} [\log p(\mathbf{z}_j | \mathbf{z}_{-j}, \mathbf{x})] \}$   
**5         end**  
**6 end**  
**7 return**  $q(\mathbf{z})$

---

**2.11.2. Actualización de coordenadas.**

El proposito de esta sección es obtener la actualización de coordenadas en la ecuación (2.32). Para ello se reescribe la cota ELBO en

$$(2.33) \quad \text{ELBO}(q_j) = \mathbb{E}_j [\mathbb{E}_{-j} [\log p(\mathbf{z}_j, \mathbf{z}_{-j}, \mathbf{x})]] - \mathbb{E}_j [\log q_j(z_j)] +$$

De esta forma, el primer término de ELBO se reescribe como una esperanza iterada. El segundo término, por su parte, se descompone, usando la hipótesis de independencia mean-field y conservando solo el término que depende de  $q_j(z_j)$ .

Hasta una constante añadida, la función objetivo en la ecuación (2.33) es igual a la divergencia KL negativa entre  $q_j(z_j)$  y  $q_j^*(z_j)$  de la ecuación (2.32). Por lo tanto, se maximiza el ELBO con respecto a  $q_j$  cuando fijamos  $q_j(z_j) = q_j^*(z_j)$ .

**2.11.3. GMM Bayesiano III.** Se continua explorando el modelo GMM, para ello se consideran  $C$  componentes y  $n$  observaciones (datos)  $x_1, \dots, x_n$ . Las variables latentes son  $C$  parámetros  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_C)^T$  y  $n$  asignaciones de clase  $\mathbf{c} = c_1, \dots, c_n$ . Se supone nuevamente una varianza unitaria y una distribución prior uniforme sobre los componentes de la mezcla.

La densidad conjunta de las variables latentes y observadas se encuentra en la ecuación (2.21). La familia variacional está en la ecuación (2.30). En este caso hay

dos tipos de parámetros variacionales: los parámetros categóricos  $\varphi_i$  para aproximar la asignación de clúster posterior del  $i$ -ésimo punto de datos y los parámetros gaussianos  $m_k$  y  $s_k^2$  para aproximar la posterior de la  $k$ -ésima componente de mezcla.

Usando la familia mean-field en la definición de la cota ELBO, se obtiene una función de los parámetros variacionales  $\mathbf{m}$ ,  $\mathbf{s}^2$  y  $\boldsymbol{\varphi}$ :

$$\begin{aligned}
 \text{ELBO}(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi}) &= \sum_{k=1}^C \mathbb{E} [\log p(\mu_k); m_k, s_k^2] \\
 (2.34) \quad &+ \sum_{k=1}^C (\mathbb{E} [\log p(c_i); \varphi_i] + \mathbb{E} [\log p(x_i | c_i, \boldsymbol{\mu}); \varphi_i, \mathbf{m}, \mathbf{s}^2]) \\
 &- \sum_{k=1}^C \mathbb{E} [\log q(c_i; \varphi_i)] - \sum_{k=1}^C \mathbb{E} [\log q(\mu_k; m_k, s_k^2)].
 \end{aligned}$$

En cada término, se hace explícita la dependencia de los parámetros variacionales. Más aún, cada esperanza puede ser calculada en forma cerrada.

Debido a que el algoritmo CAVI actualiza cada parámetro variacional a su vez, es posible obtener una regla de actualización para cada factor de forma cerrada.

**2.11.4. Densidad variacional para las asignaciones de cluster.** Se derivamos la actualización variacional para la asignación de cluster  $c_i$ . Para ello se usa la ecuación (2.32):

$$(2.35) \quad q^*(c_i; \varphi_i) \propto \exp \{ \log p(c_i) + \mathbb{E} [\log p(x_i | c_i, \boldsymbol{\mu}); \mathbf{m}, \mathbf{s}^2] \}.$$

Los términos en el exponente son los componentes de la densidad conjunta que dependen de  $c_i$ , la esperanza en el segundo término se calcula sobre  $\boldsymbol{\mu}$ . El primer término de la ecuación (2.35) corresponde a la log prior de  $c_i$ , que al considerarla uniforme pasa a ser  $\log p(c_i) = -\log k$ .

El segundo término es la log esperanza de la  $c_i$ -ésima densidad gaussiana. Al ser  $c_i$  es un vector indicador, se obtiene:

$$p(x_i | c_i, \boldsymbol{\mu}) = \prod_{k=1}^C p(x_i | \mu_k)^{c_{ik}}.$$

Usando este resultado, y eliminando los términos constantes con respecto a  $c_i$ , se calcula la probabilidad de la log esperanza,

$$\begin{aligned}
 \mathbb{E} [\log p(x_i | c_i, \boldsymbol{\mu})] &= \sum_k c_{ik} \mathbb{E} [\log p(x_i | \mu_k); m_k, s_k^2] \\
 (2.36) \quad &= \sum_k c_{ik} \mathbb{E} [-(x_i - \mu_k)^2 / 2; m_k, s_k^2] + \text{const.} \\
 &= \sum_k c_{ik} (\mathbb{E} [\mu_k; m_k, s_k^2] x_i - \mathbb{E} [\mu_k^2; m_k, s_k^2] / 2) + \text{const}
 \end{aligned}$$

Este cálculo requiere de  $\mathbb{E}[\mu_k]$  y  $\mathbb{E}[\mu_k^2]$  para cada componente, los cuales son calculables a través de la distribución variacional gaussiana para la  $k$ -ésima componente.

Por lo tanto, la actualización variacional para la  $i$ -ésima asignación de cluster, depende únicamente de los parámetros variacionales y verifica:

$$(2.37) \quad \varphi_{ik} \propto \exp \left\{ \mathbb{E} [\mu_k; m_k, s_k^2] x_i - \mathbb{E} [\mu_k^2; m_k, s_k^2] / 2 \right\}$$

**2.11.5. Densidad variacional para las medias.** Se continua con el cálculo de la densidad variacional  $q(\mu_k; m_k, s_k^2)$  para la  $k$ -ésima componente. Nuevamente, se utiliza la ecuación (2.32) y se escribe la densidad conjunta hasta una constante de normalización:

$$(2.38) \quad q(\mu_k) \propto \exp \left\{ \log p(\mu_k) + \sum_{i=1}^n \mathbb{E} [\log p(x_i | c_i, \mu); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}_{-k}^2] \right\}.$$

Posteriormente, se calcula la log coordenada-optimal  $q(\mu_k)$ . En este caso, dado nuevamente que  $c_i$  es un vector indicador, se aprecia que  $\varphi_{ik} = \mathbb{E}[c_{ik}; \varphi_i]$ , luego

$$(2.39) \quad \begin{aligned} \log q(\mu_k) &= \log p(\mu_k) + \sum_i \mathbb{E} [\log p(x_i | c_i, \mu); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}_{-k}^2] + \text{const.} \\ &= \log p(\mu_k) + \sum_i \mathbb{E} [c_{ik} \log p(x_i | \mu_k); \varphi_i] + \text{const.} \\ &= -\mu_k^2 / 2\sigma^2 + \sum_i \mathbb{E} [c_{ik}; \varphi_i] \log p(x_i | \mu_k) + \text{const.} \\ &= -\mu_k^2 / 2\sigma^2 + \sum_i \varphi_{ik} (-(x_i - \mu_k)^2 / 2) + \text{const.} \\ &= -\mu_k^2 / 2\sigma^2 + \sum_i \varphi_{ik} x_i \mu_k - \varphi_{ik} \mu_k^2 / 2 + \text{const.} \\ &= \left( \sum_i \varphi_{ik} x_i \right) \mu_k - \left( 1/2\sigma^2 + \sum_i \varphi_{ik} / 2 \right) \mu_k^2 + \text{const.} \end{aligned}$$

El cálculo anterior muestra que la densidad variacional optima por coordenada de  $\mu_k$  corresponde a una familia exponencial con estadísticos suficientes  $\mu_k, \mu_k^2$  y parámetros  $\{\sum_{i=1}^n \varphi_{ik} x_i - 1/2\sigma^2 - \sum_{i=1}^n \varphi_{ik}/2\}$ , es decir, corresponde a una gaussiana. En términos de la media y la varianza variacionales, las reglas de actualización para  $q(\mu_k)$  son

$$(2.40) \quad m_k = \frac{\sum_i \varphi_{ik} x_i}{1/\sigma^2 + \sum_i \varphi_{ik}}, \quad s_k^2 = \frac{1}{1/\sigma^2 + \sum_i \varphi_{ik}}$$

Estas reglas se relacionan fuertemente con la densidad condicional completa de la para la  $k$ -ésima componente del modelo. La condicional completa es una gaussiana posterior dados los datos asignados a la  $k$ -ésima componente. La actualización variacional corresponde a una condicional completa ponderado, donde cada dato se pondera por su probabilidad variacional de ser asignado a la componente  $k$ .

**2.11.6. CAVI para GMM.** El algoritmo 4 muestra el algoritmo CAVI aplicado al GMM estudiado. En este, se combinan las actualizaciones variacionales en la ecuación (2.35) y la ecuación (2.40). El algoritmo requiere calcular el ELBO de la ecuación (2.34). Usando el ELBO para rastrear el progreso del algoritmo y evaluar cuándo ha convergido. Una vez se tiene una densidad variacional ajustada, esta se puede usar como reemplazo para la posterior. Por ejemplo, es posible obtener una

descomposición posterior de los datos. Asignando puntos a su cluster más probable  $\hat{c}_i = \arg \max_k \varphi_{ik}$  y se estiman las medias de cluster con sus medias variacionales  $m_k$ .

---

**Algorithm 4:** CAVI para GMM

---

**Input:** Datos  $x_1, \dots, x_n$ , numero de componentes  $C$ , varianza prior para las medias de las componentes  $\sigma^2$

**Output:** Densidades variacionales  $q(\mu_k; m_k, s_k^2)$  y  $q(c_i, \varphi_i)$

```

1 Inicializar:  $\mathbf{m} = m_1, \dots, m_k, \mathbf{s}^2 = s_1 \dots s_k, \boldsymbol{\varphi} = \varphi_1, \dots, \varphi_n$ 
2 while ELBO no converge do
3   for  $i = 1, \dots, n$  do
4     | Calcular  $\varphi_{ij} \propto \exp \{ \mathbb{E} [\mu_k; m_k, s_k^2] x_i - \mathbb{E} [\mu_k^2; m_k, s_k^2] / 2 \}$ 
5   end
6   for  $k = 1, \dots, C$  do
7     |  $m_k \leftarrow \frac{\sum_i \varphi_{ik} x_i}{1/\sigma^2 + \sum_i \varphi_{ik}}$ 
8     |
9     |  $s_k^2 \leftarrow \frac{1}{1/\sigma^2 + \sum_i \varphi_{ik}}$ 
10  end
11  Calcular ELBO( $\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi}$ )
12 end
13 return  $q(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi})$ 

```

---

También es posible usar la densidad variable ajustada para aproximar la densidad predictiva en nuevos datos:

$$(2.41) \quad p(x_{\text{nex}} | x_{1:n}) \approx \frac{1}{K} \sum_{k=1}^C p(x_{\text{new}} | m_k),$$

donde  $p(x_{\text{new}} | m_k)$  es gaussiana con media  $m_k$  y varianza unitaria.

### 2.12. Inferencia variacional con familias exponenciales

En el modelo GMM expuesto anteriormente, la distribución condicional completa pertenece a la clase de distribuciones conocida como *familia exponencial*, este tipo de distribuciones (normal, gamma, chi-cuadrado, beta, ...) han sido ampliamente estudiadas y poseen propiedades interesantes en cuanto a conjugación prior-posterior, entre otras. Los modelos cuyas distribuciones condicionales completas pertenecen a esta familia de distribuciones presentan grandes ventajas en cuando a al aproximación por medio de inferncia variacional. En esta sección se estudian tales propiedades.

**2.12.1. Condicionales completas en familias exponenciales.** Se considera el modelo  $p(\mathbf{z}, \mathbf{x})$  y se agrega la suposición de que cada condicional completa pertenece a la familia de distribuciones exponencial, es decir:

$$(2.42) \quad p(z_j | \mathbf{Z}_{-j}, \mathbf{x}) = h(z_j) \exp \{ \eta_j(z_j, \mathbf{x})^\top z_j - a(\eta_j(\mathbf{z}_{-j}, \mathbf{x})) \},$$

Donde  $z_j$  es su propio estadístico suficiente,  $h(\cdot)$  es una medida base, y  $a(\cdot)$  es el log normalizador. Debido a que esta es una densidad condicional, el parámetro  $\eta_j(\mathbf{z}_{-j}, \mathbf{x})$  es una función del conjunto condicionante.

En este contexto, la inferencia variacional por medio de distribuciones mean-field pasa por ajustar  $q(\mathbf{z}) = \prod_j q_j(z_j)$ . Al suponer pertenencia a la familia exponencial, se simplifica la actualización de coordenadas de la ecuación (2.31), en efecto:

$$\begin{aligned}
 q(z_j) &\propto \exp \{ \mathbb{E} [\log p(z_j | \mathbf{z}_{-j}, \mathbf{x})] \} \\
 (2.43) \quad &= \exp \{ \log h(z_j) + \mathbb{E} [\eta_j(\mathbf{z}_{-j}, \mathbf{x})]^\top z_j - \mathbb{E} [a(\eta_j(\mathbf{z}_{-j}, \mathbf{x}))] \} \\
 &\propto h(z_j) \exp \{ \mathbb{E} [\eta_j(\mathbf{z}_{-j}, \mathbf{x})]^\top z_j \}
 \end{aligned}$$

En esta actualización se observa la forma paramétrica de los factores variacionales óptimos. Más aún, cada uno está en la misma familia exponencial que su correspondiente condicional completo. Su parámetro tiene la misma dimensión y tiene la misma medida base  $h(\cdot)$  y log normalizador  $a(\cdot)$ .

Habiendo establecido sus formas paramétricas, se denota por  $\nu_j$  el parámetro variacional para el factor variacional  $j$ -ésimo. Cuando se actualiza cada factor, se establece su parámetro igual al parámetro esperado de la condicional completa:

$$(2.44) \quad \nu_j = \mathbb{E} [\eta_j(\mathbf{z}_{-j}, \mathbf{x})].$$

Esta transformación permite simplificar los calculos necesarios para implementar CAVI.

**2.12.2. Conjugación condicional y modelos bayesianos.** Un caso especial de modelos de familias exponenciales son los *modelos condicionalmente conjugados* con variables locales y globales. Este tipo de modelos son comunes en aprendizaje automático, donde las variables globales corresponden a “parámetros” y las variables locales a variables latentes.

En tal contexto, sea  $\beta$  un vector de *variables latentes globales* y  $\mathbf{z}$  un vector de *variables locales latentes*. La probabilidad conjunta en este caso es:

$$(2.45) \quad p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta).$$

En la perspectiva del modelo GMM ya explorado, las variables globales corresponden a las componentes de mezcla y la  $i$ -ésima variable local es la asignación de cluster para el dato  $x_i$ .

En el modelo de la ecuación (2.45) se supondrá, que los términos que la componen hacen que cada condicional completa pertenesca a la familia exponencial. Para ello, se supondrá que la densidad conjunta de cada par  $(x_i, z_i)$ , condicional a  $\beta$ , tiene una forma de familia exponencial,

$$(2.46) \quad p(z_i, x_i | \beta) = h(z_i, x_i) \exp \{ \beta^\top t(z_i, x_i) - a(\beta) \},$$

donde  $t(\cdot, \cdot)$  es el estadístico suficiente. Posteriormente, se escoge una prior sobre las variables globales de manera que sea su conjugada,

$$(2.47) \quad p(\beta) = h(\beta) \exp \{ \alpha^\top [\beta, -a(\beta)] - a(\alpha) \}.$$

Esta prior tiene el parámetro  $\alpha = [\alpha_1, \alpha_2]^\top$ , un vector de columna, y estadísticos suficientes que concatenan la variable global y su log normalizador en la densidad de las variables locales.

Con la prior conjugada, la condicional completa de las variables globales está en la misma familia. Su parámetro natural es

$$(2.48) \quad \hat{\alpha} = \left[ \alpha_1 + \sum_{i=1}^n t(z_i, x_i), \alpha_2 + n \right]^\top.$$

En cuanto al condicional completo de la variable local  $z_i$ , se observa que dado  $\beta$  y  $x_i$ , la variable local  $z_i$  es condicionalmente independiente de las otras variables locales  $z_{-i}$  y otros datos  $x_{-i}$ . Esto se deduce de la forma de la densidad conjunta en la ecuación 2.45). Así

$$(2.49) \quad p(z_i | x_i, \beta, \mathbf{z}_{-i}, \mathbf{x}_{-i}) = p(z_i | x_i, \beta).$$

Si se supone además que esta densidad está en una familia exponencial,

$$(2.50) \quad p(z_i | x_i, \beta) = h(z_i) \exp \{ \eta(\beta, x_i)^\top z_i - a(\eta(\beta, x_i)) \}.$$

### 2.13. Inferencia variacional en modelos condicionalmente conjugados

A continuación, se describe **CAVI** para esta clase general de modelos. Para ello, sea  $q(\beta | \lambda)$  la aproximación posterior variacional en  $\beta$ , se denomina además por  $\lambda$  al parámetro variacional global. El cual posee la misma densidad de familia exponencial que la prior. De forma similar, sea  $q(z_i | \varphi_i)$  la distribución variacional posterior en cada variable local, en este caso, cada  $z_i$  rije por un parámetro local variacional  $\varphi_i$ . Este posee la misma densidad de familia exponencial que el condicional local completo. El algoritmo **CAVI** itera entre la actualización de cada parámetro variacional local y la actualización del parámetro variacional global. La actualización variacional local es

$$(2.51) \quad \varphi_i = \mathbb{E}_\lambda [\eta(\beta, x_i)].$$

Esta es una aplicación de la Ecuación (2.44), donde se toma la esperanza del parámetro natural de la condicional completa en la ecuación (2.49). Para la actualización variacional global se aplica la misma técnica:

$$(2.52) \quad \lambda = \left[ \alpha_1 + \sum_{i=1}^n \mathbb{E}_{\varphi_i} [t(z_i, x_i)], \alpha_2 + n \right]^\top.$$

Aquí se calcula la esperanza del parámetro natural en la Ecuación (2.48). El algoritmo **CAVI** optimiza la cota **ELBO** al iterar entre las actualizaciones locales de cada parámetro local y las actualizaciones globales de los parámetros globales. Para evaluar la convergencia, se puede calcular la cota **ELBO** en cada iteración hasta una constante que no depende de los parámetros variacionales,

$$(2.53) \quad \text{ELBO} = \left( \alpha_1 + \sum_{i=1}^n \mathbb{E}_{\varphi_i} [t(z_i, x_i)] \right)^\top \mathbb{E}_\lambda [\beta] - (\alpha_2 + n) \mathbb{E}_\lambda [a(\beta)] - \mathbb{E}[\log q(\beta, \mathbf{z})].$$

Esta corresponde a la cota **ELBO** en la ecuación (2.27) aplicada en (2.45) con la correspondiente densidad variacional mean-field. El último término es

$$(2.54) \quad \mathbb{E}[\log q(\beta, \mathbf{z})] = \lambda^\top \mathbb{E}_\lambda [t(\beta)] - a(\lambda) + \sum_{i=1}^n \varphi_i^\top \mathbb{E}_{\varphi_i} [z_i] a(\varphi_i).$$

**CAVI** para GMM (Algoritmo 4) corresponde a una instancia de este método.

### 2.14. Inferencia variacional estocástica

En relación al escalamiento a conjuntos de datos masivos, la mayoría de los algoritmos de inferencia posterior, incluido CAVI sufren de problemas de escalamiento. La razón de esto, es que la estructura de ascenso de coordenadas del algoritmo requiere iterar a través de todo el conjunto de datos en cada iteración. Como alternativa a esto, surge la optimización basada en gradiente, esta perspectiva es la clave para ampliar la inferencia variacional a su variante estocástica (SVI). Esta se centra en la optimización de los parámetros variacionales globales  $\lambda$  en modelos condicionalmente conjugados.

El esquema de calculo se basa en mantener una estimación de los parámetros variables globales, para repetidamente sub-samplear un dato del conjunto completo y luego, usar los parámetros globales actuales para calcular los parámetros locales óptimos para el dato muestreado, finalmente, se ajustan los parámetros globales actuales de una manera apropiada.

**2.14.1. Gradiente natural de ELBO.** En la optimización basada en gradiente, el termino *gradiente natural* toma en cuenta la estructura geométrica de los parámetros de probabilidad. Específicamente, los gradientes naturales deforman el espacio de parámetros de una manera sensible, de modo que mover la misma distancia en diferentes direcciones equivale a un cambio igual en la divergencia KL simetrizada. El gradiente Euclidiano habitual no posee esta propiedad.

En familias exponenciales, se encuentra el gradiente natural con respecto al parámetro al premultiplicar el gradiente usual por la covarianza inversa del estadístico suficiente,  $a''(\lambda)^{-1}$ . Esta es la métrica Riemanniana inversa y la matriz de información Fisher inversa.

Los modelos conjugados condicionalmente poseen gradientes naturales simples para la cota ELBO. Para gradientes con respecto al parámetro global  $\lambda$  se conoce la formula del gradiente euclidiano de ELBO:

$$(2.55) \quad \nabla_{\lambda} \text{ELBO} = a''(\lambda)(\mathbb{E}_{\varphi}[\hat{\alpha}] - \lambda),$$

donde  $\mathbb{E}_{\varphi}[\hat{\alpha}]$  está en la ecuación (2.52). Premultiplicar por la información de Fisher inversa da el gradiente natural  $g(\lambda)$ ,

$$(2.56) \quad g(\lambda) = \mathbb{E}_{\varphi}[\hat{\alpha}] - \lambda.$$

La cual es la diferencia entre las actualizaciones de coordenadas  $E_{\varphi}[\hat{\alpha}]$  y los parámetros variacionales  $\lambda$  en los que se evalua el gradiente.

Es posible usar este gradiente natural en un algoritmo de optimización basado en gradiente. Donde en cada iteración, se actualizan los parámetros globales,

$$(2.57) \quad \lambda_t = \lambda_{t-1} + \epsilon_t g(\lambda_{t-1}),$$

donde  $\epsilon_t$  corresponde a un tamaño de paso. Al sustituir la ecuación (2.56) en el segundo término se obtiene una estructura especial,

$$(2.58) \quad \lambda_t = (1 - \epsilon_t)\lambda_{t-1} + \epsilon_t \mathbb{E}_{\varphi}[\hat{\alpha}].$$

Esto no requiere de cálculos adicionales a los de las actualizaciones de ascenso de coordenadas. En cada iteración, primero se calcula la actualización de coordenadas. Luego se ajusta la estimación actual para que sea una combinación ponderada de la actualización y el parámetro variacional actual.

**2.14.2. Optimización estocástica de ELBO.** Aunque es fácil de calcular, usar el gradiente natural tiene el mismo costo que la actualización de coordenadas en la ecuación (2.52), pues requiere sumar en todo el conjunto de datos y calcular los parámetros variables locales óptimos para cada uno de ellos.

La inferencia variacional estocástica resuelve este problema utilizando el gradiente natural en un algoritmo de optimización estocástico. Los algoritmos de optimización estocástica siguen gradientes ruidosos pero de bajo costo computacional para alcanzar el óptimo de una función objetivo.

El objetivo es construir un gradiente natural, ruidoso calculado a bajo costo computacional. Para ello, se expande el gradiente natural en la ecuación (2.56) usando (2.48):

$$(2.59) \quad g(\lambda) = \alpha + \left[ \sum_{i=1}^n \mathbb{E}_{\varphi_i^*}[t(z_i, x_i)], n \right]^\top - \lambda,$$

Se construye un gradiente natural ruidoso al muestrear un índice de los datos y luego reescalar el segundo término,

$$(2.60) \quad t \sim \text{Unif}(1, \dots, n)$$

$$(2.61) \quad \hat{g}(\lambda) = \alpha + n[\mathbb{E}_{\varphi_t^*}[t(z_t, x_t)], 1]^\top - \lambda.$$

El gradiente natural ruidoso  $\hat{g}(\lambda)$  es insesgado:  $E_t[\hat{g}(\lambda)] = g(\lambda)$ . Y es computacionalmente eficiente, solo involucra un único dato muestreado y un conjunto de parámetros locales optimizados.

Finalmente, se establece una sucesión de tamaño de paso. La cual debe cumplir las condiciones:

$$(2.62) \quad \sum_t \epsilon_t = \infty \quad ; \quad \sum_t \epsilon_t^2 < \infty$$

Muchas sucesiones cumplen estas condiciones, por ejemplo  $\epsilon_t = t^{-\kappa}$  para  $\kappa \in (0.5, 1]$ . **SVI** no requiere una nueva derivación más allá de lo que se necesita para **CAVI**. Cualquier implementación de **CAVI** puede escalarse inmediatamente a un algoritmo estocástico.



## CHAPTER 3

# Aprendizaje con Kernels

### 3.1. Introducción

Un kernel es una función simétrica y definida positiva  $k(\cdot, \cdot)$  que puede ser entendida como una medida de similitud entre los argumentos que opera. En el siguiente capítulo, se definen estos objetos matemáticos de manera formal, se investigan sus características y se derivan algunos métodos del aprendizaje de máquinas que toman ventaja sus propiedades.

### 3.2. Terminología y propiedades

El término **kernel** proviene del estudio de operadores integrales en el campo del análisis funcional. En tal contexto se les identifica como aquellas funciones  $k$  que determinan un operador  $T_k$  a través de:

$$(3.1) \quad (T_k f)(x) = \int_{\mathcal{X}} k(x, x') f(x') dx$$

En concordancia con la perspectiva que se desea abarcar, se denotará como kernel a toda función  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  en el espacio de características<sup>1</sup>. Dentro de tal clase de funciones, son de importancia a aquellas capaces de “generalizar” el concepto de *producto interno*, el **Teorema de Mercer** permite identificar tal subconjunto.

**Teorema 3.1** (Teorema de Mercer). Sea  $(\mathcal{X}, \mu)$  un espacio de medida finita y  $k \in L_{\infty}(\mathcal{X}^2)$  una función real y simétrica, tal que el operador integral:

$$(3.2) \quad \begin{aligned} T_k : L_2(\mathcal{X}) &\rightarrow L_2(\mathcal{X}) \\ (T_k f)(x) &:= \int_{\mathcal{X}} k(x, x') f(x') d\mu(x') \end{aligned}$$

es semidefinido positivo, es decir, para toda  $f \in L_2(\mathcal{X})$  se cumple:

$$(3.3) \quad \int_{\mathcal{X}^2} k(x, x') f(x) f(x') d\mu(x) d\mu(x') \geq 0$$

Si  $\psi_j \in L_2(\mathcal{X})$  son las funciones propias ortonormales de  $T_f$  asociadas a los valores propios  $\lambda_j > 0$ , ordenados de manera decreciente. Entonces,

$$(1) \quad (\lambda_j)_j \in l_1$$

---

<sup>1</sup>El condominio del kernel  $k$  no tiene por que restringirse a los reales, para ciertas aplicaciones puede ser conviene utilizar los complejos. En general las propiedades de interés se preservan en ambos cuerpos, por lo que por simplicidad se considera solo el caso real en esta monografía.

- (2)  $k(x, x') = \sum_{j=1}^{N_{\mathcal{H}}} \lambda_j \psi_j(x) \psi_j(x')$  se cumple para casi todos los elementos  $x, x' \in \mathcal{X}$ . Además  $N_{\mathcal{H}} \in \mathbb{N}$  o  $N_{\mathcal{H}} = \infty$ , en este último caso, la serie respectiva converge absoluta y uniformemente para casi todos los elementos  $x, x' \in \mathcal{X}$ .

La segunda implicación del teorema anterior, permite construir una aplicación  $\Phi : \mathcal{X} \rightarrow l_2^{N_{\mathcal{H}}}$  tal que  $x \mapsto (\sqrt{\lambda_j} \psi_j(x))_{j=1, \dots, N_{\mathcal{H}}}$ , es decir, a cada elemento de  $\mathcal{X}$  se le asocia una transformación por medio de las funciones propias asociadas a  $k(\cdot, \cdot)$  al espacio  $l_2^{N_{\mathcal{H}}}$ . Este último espacio está provisto de producto interno, por lo que es posible expresar  $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{l_2^{N_{\mathcal{H}}}}$  para casi todo  $x \in \mathcal{X}$ . Esto permite interpretar a  $\Phi$  como una aplicación a un espacio de *características*, más aún, en este espacio  $k(\cdot, \cdot)$  actúa como un producto interno, esto se concreta en el siguiente teorema.

**Teorema 3.2** (Aplicación kernel de Mercer). Si  $k$  es un kernel que cumple las condiciones del teorema (3.1), se puede construir una aplicación  $\Phi$  a un espacio donde  $k$  se comporta como un producto interno:

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

para casi todo  $x, x' \in \mathcal{X}$ . Más aún, para todo  $\varepsilon > 0$ , existe una aplicación  $\Phi_n$  a un espacio  $n$ -dimensional con producto interno tal que

$$|k(x, x') - \langle \Phi_n(x), \Phi_n(x') \rangle| < \varepsilon$$

para casi todo  $x, x' \in \mathcal{X}$ .

La convergencia uniforme en el teorema anterior implica que para cualquier precisión  $\varepsilon > 0$ , debe existir un  $n \in \mathbb{N}$  tal que  $k$  puede ser aproximado como un producto interno en  $\mathbb{R}^n$ . Lo anterior se observa al notar que para casi todo  $x, x' \in \mathcal{X}$ , se tiene  $|k(x, x') - \langle \Phi^n(x), \Phi^n(x') \rangle| < \varepsilon$ , donde  $\Phi^n(x) : x \mapsto (\sqrt{\lambda_1} \psi_1(x), \dots, \sqrt{\lambda_n} \psi_n(x))$ . En tal contexto, se puede interpretar al espacio de características como un espacio finito dimensional dentro de cierta precisión  $\varepsilon$ . Esta característica es una parte esencial en el aprendizaje con kernels y dará lugar a la técnica conocida como *truco del kernel*.

La condición de simetría para  $k(\cdot, \cdot)$ , necesaria en el teorema (3.1), permite dotar de esta propiedad al producto interno definido en (3.2). Por otra parte, es necesario caracterizar la positividad de  $k(\cdot, \cdot)$  de manera tal que se pueda contextualizar dentro del aprendizaje de máquinas, para ello se utiliza la noción de matriz de **Gram**.

**Definición 3.3** (Matriz de Gram). Dado un kernel  $k(\cdot, \cdot)$  y un conjunto de puntos  $\mathcal{X}_n = \{x_i \mid i = 1, \dots, n\}$ , se define la matriz de Gram como la matriz  $K$  que verifica  $K_{i,j} = k(x_i, x_j)$ .

Esta noción, permite definir la positividad de un kernel  $k(\cdot, \cdot)$  en función de las propiedades que sus matrices de Gram poseen, en este sentido, el concepto de **matriz semidefinida** juega un rol fundamental.

**Definición 3.4** (matriz semidefinida positiva). Una matriz real  $K$  de  $n \times n$  se dice semidefinida positiva (denotado SDP) si para todo vector  $v \in \mathbb{R}^n$  se cumple que  $Q(v) = v^T K v \geq 0$ . Si además se cumple que  $Q(v) = 0 \Leftrightarrow v = 0$ , entonces se dice que la matriz  $K$  es definida positiva (denotada DP).

Dentro de las propiedades de este tipo de matrices se encuentra la siguiente proposición:

**Proposición 1.** Una matriz simétrica es semidefinida positiva si y solo si todos sus valores propios son no negativos.

Finalmente, se podrá caracterizar la positividad de un kernel  $k(\cdot, \cdot)$  mediante la siguiente proposición:

**Proposición 2.** Un kernel  $k$  es semidefinido positivo si y sólo si toda matriz de Gram  $K$  generada con  $k$  es semidefinida positiva. En ese caso, se dice que  $k$  es una función de covarianza.

En particular, esto significa que cualquier función simétrica y semidefinida positiva, acepta una descomposición por medio del teorema de Mercer y por tanto induce una transformación  $\Phi$  a un espacio de características (posiblemente de dimensión infinita,  $N_{\mathcal{H}} = \infty$ ) donde actúa como un producto interno. Si por otra parte, se posee una transformación  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ , donde  $\mathcal{H}$  es un espacio con producto interno, es posible construir un kernel por medio de  $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$ . Este es por definición, simétrico (lo hereda del producto interno) y además para todo  $c_i \in \mathbb{R}$ ,  $x_i \in \mathcal{X}$ ,  $i = 1, \dots, m$ , cumple

$$\sum_{i,j} c_i c_j k(x_i, x_j) = \left\langle \sum_i c_i \Phi(x_i), \sum_j c_j \Phi(x_j) \right\rangle = \left\| \sum_i c_i \Phi(x_i) \right\|^2 \geq 0$$

Por tanto será semidefinido positivo y finalmente un kernel que cumple las condiciones del teorema (3.1). Esto permite modelar kernels por medio de transformaciones conocidas, a la vez que permite asumir la existencia de una transformación dado un kernel.

### 3.3. Espacios de Hilbert con kernel reproductor - RKHS

Hasta este punto, se puede entender un kernel como una función semidefinida positiva y simétrica compatible con una representación, ya sea implícita o explícita, de un espacio inicial  $\mathcal{X}$  en un espacio con producto interno (o *pre-Hilbertiano*). El propósito de esta sección corresponde a estudiar las características de este espacio con el fin de obtener herramientas para el uso de kernels en tareas de aproximación de funciones.

**Definición 3.5** (Espacio de Hilbert con Kernel Reprodutor). Sea  $\mathcal{X}$  conjunto no vacío y  $\mathcal{H}$  un espacio de Hilbert de funciones  $f : \mathcal{X} \rightarrow \mathbb{R}$ . Se denomina a  $\mathcal{H}$  espacio de Hilbert con kernel reproductor con producto punto  $\langle \cdot, \cdot \rangle$  (y por tanto la norma  $\|f\| := \sqrt{\langle f, f \rangle}$ ) si existe una función  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  con las siguientes propiedades.

- (1)  $k(\cdot, \cdot)$  tiene la propiedad de reproductor  
 $\langle f, k(x, \cdot) \rangle = f(x)$  para todo  $f \in \mathcal{H}$ , en particular  
se verifica:  $\langle k(x, \cdot), k(x', \cdot) \rangle = k(x, x')$
- (2)  $k(\cdot, \cdot)$  genera  $\mathcal{H}$ , es decir,  $\mathcal{H} = \overline{\text{gen} \{k(x, \cdot) \mid x \in \mathcal{X}\}}$  donde  $\overline{X}$  denota la completación del conjunto  $X$ .

Por definición (punto 2), un RKHS es un espacio formado por combinaciones lineales de funciones de la forma  $k_x(\cdot)$  para todo  $x \in \mathcal{X}$  y sus funciones límite (clausura). Esto último quiere decir, que toda función  $f \in \mathcal{H}$  puede ser aproximada

a través de sucesiones de combinaciones lineales de elementos de la forma  $k_x(\cdot)$ , que claramente están únicamente determinados por el kernel  $k(\cdot, \cdot)$ .

Por su parte, la propiedad de reproducción (punto 1), se puede comprender al observar que para todo  $f \in \mathcal{H}$ , espacio RKHS correspondiente al kernel  $k(\cdot, \cdot)$ , es posible escribir  $f$  como:

$$(3.4) \quad f(\cdot) = \sum_{i=1}^{\infty} \alpha_i k(x_i, \cdot)$$

Que por el teorema de Mercer verifica:

$$(3.5) \quad f(x) = \sum_{i=1}^{\infty} \alpha_i \sum_{j=1}^{N_{\mathcal{H}}} \lambda_j \psi_j(x_i) \psi_j(x)$$

Como todo funcional de la forma  $k_x(\cdot) = k(x, \cdot) \in \mathcal{H}$  y este espacio posee producto interno, es posible escribir  $\langle f(\cdot), k(x, \cdot) \rangle$  de la forma:

$$(3.6) \quad \langle f, k_x \rangle_{\mathcal{H}} = \sum_{i=1}^{\infty} \alpha_i \sum_{j=1}^{N_{\mathcal{H}}} \sum_{l=1}^{N_{\mathcal{H}}} \lambda_j \psi_j(x_i) \lambda_l \psi_l(x) \langle \psi_j, \psi_l \rangle_{L_2(\mathcal{X})}$$

$$(3.7) \quad = \sum_{i=1}^{\infty} \alpha_i \sum_{j=1}^{N_{\mathcal{H}}} \lambda_j(x_i) \psi_j(x_i) \psi_j(x)$$

$$(3.8) \quad = f(x)$$

Donde en la primera ecuación se utiliza la ortogonalidad  $\psi_i \perp \psi_j$  en todo  $i \neq j$ . Posteriormente, se redefine el producto interno en  $L_2(\mathcal{X})$  de manera que  $\langle \psi_j, \psi_l \rangle = \delta_{ij}/\lambda_j$ , finalmente se tiene que <sup>2</sup> :

$$(3.9) \quad \langle f, k(x, \cdot) \rangle = f(x)$$

para todo  $f \in \mathcal{H}$ . Es decir, en un espacio compuesto por funciones generadas por un kernel de Mercer se cumple la propiedad de reproducción, por lo tanto, si el espacio es completo (en caso contrario se considera la adherencia del conjunto) se posee un RKHS con respecto a tal kernel. Desde este punto vista, se puede también concluir que es posible construir un espacio RKHS con respecto a un espacio  $\mathcal{X}$  y un kernel de Mercer, definiendo el producto interno de las funciones generadas por el kernel de manera tal que se cumpla la propiedad de reproducción. Finalmente, se consideran las funciones límite de tal conjunto para construir el RKHS correspondiente.

Por definición para cada espacio RKHS existe una función kernel, por el argumento anterior, se puede intuir la existencia de cierta relación de equivalencia entre los kernels que describen a un espacio RKHS. La respuesta a tal inquietud es aún más fuerte, pues establece que a cada RKHS corresponde un único kernel [6].

**Teorema 3.6** (Moore-Aronszajn). Sea  $\mathcal{X}$  un conjunto, para cada función definida positiva  $k$  en  $\mathcal{X} \times \mathcal{X}$  existe un único RKHS asociado a esta. Más aún, cada RKHS posee un único kernel definido positivo asociado.

Debido a esta relación de unicidad, tiene sentido definir al espacio RKHS  $\mathcal{H}$  por medio de la base ortogonal conformada por las funciones propias de su kernel.

---

<sup>2</sup> $\delta_{ij} = 1$  si  $i = j$ , y  $\delta_{ij} = 0$  en caso contrario. Este objeto matemático se conoce como *delta de Kronecker*.

**Proposición 3.** Sea  $\mathcal{H}$  el espacio compuesto por las combinaciones lineales de las funciones propias ortogonales  $\{\psi_j\}_{j=1}^{N_{\mathcal{H}}}$  de un kernel  $k(\cdot, \cdot)$  con respecto a una medida  $\mu$ . Toda  $f \in \mathcal{H}$ , se expresa por tanto como  $f(x) = \sum_{i=1}^{N_{\mathcal{H}}} f_i \psi_i(x)$ , si se impone la restricción,  $\sum_{i=1}^N \frac{f_i^2}{\lambda_i} < \infty$ , es posible dotar a  $\mathcal{H}$  con el producto interno  $\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^N \frac{f_i g_i}{\lambda_i}$ , donde  $\lambda_j$  es el valor propio asociado a  $\psi_j$ . Si  $\mathcal{H}$  es completo, entonces es el RKHS asociado a  $k(\cdot, \cdot)$ .

De esta forma, se puede formular el siguiente teorema:

**Teorema 3.7** (Teorema representador). Dado un kernel  $k$  sobre  $\mathcal{H}$ , se considera el espacio de funciones:

$$\mathcal{H}_0 = \left\{ f(x) = \sum_{i=1}^{n_f} f_i k(x, x_i) \mid n_f \in \mathbb{N}, x_i \in \mathcal{X}, f_i \in \mathbb{R} \right\}$$

con el producto interno  $\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^{n_f} \sum_{j=1}^{n_g} f_i g_j k(x_i, x_j)$ . Luego  $\mathcal{H} = \overline{\mathcal{H}_0}$  es el RKHS asociado a  $k(\cdot, \cdot)$ .

Lo anterior es directo de  $f(x) = k_{\bar{x}}(x) \in \mathcal{H}_0$ , donde  $n_f = 1, x_1 = \bar{x}$  y  $f_1 = 1$ , de donde se obtiene:

$$\langle f, k_{\bar{x}} \rangle_{\mathcal{H}_0} = \sum_{i=1}^{n_f} f_i k(\bar{x}, x_i) = f(\bar{x})$$

Dado que  $k(\cdot, \cdot)$  es único con respecto a su RKHS.

### 3.4. Kernels

En la presente sección se estudian distintos kernels. Se observan sus características y se derivan modelos basados en estas.

**3.4.1. Kernel RBF.** El *kernel exponencial cuadrático* (SE) o *gaussiano* para  $\mathbf{x} = (x_1, \dots, x_D)^T \in \mathbb{R}^D$  se define por:

$$(3.10) \quad k(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{1}{2} (\mathbf{x} - \mathbf{x}')^\top \Sigma^{-1} (\mathbf{x} - \mathbf{x}') \right)$$

donde si  $\Sigma$  es diagonal, se reduce a

$$(3.11) \quad k(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{1}{2} \sum_{j=1}^D \frac{1}{\sigma_j^2} (x_j - x'_j)^2 \right)$$

En este tipo de kernels, se puede interpretar  $\sigma_j$  como el **ancho de banda característico** de la componente (característica)  $j$ -ésima. En tal contexto, si  $\sigma_j = \infty$ , los valores de tal dimensión pasan a ser ignorados. Por lo tanto, se pueden identificar las características o componentes con menor influencia examinando la magnitud de su ancho de banda correspondiente. Este método se conoce como **Kernel ARD**.

Por otra parte, si  $\Sigma$  es esférica, se obtiene el kernel isotrópico

$$(3.12) \quad k(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right)$$

donde  $\sigma$  se conoce como **ancho de banda**. Las funciones del tipo (3.12), se conocen como **funciones de base radial**, y se caracterizan por depender únicamente de la diferencia entre los puntos que opera, es decir,  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ .

Por este motivo el kernel de la ecuación (3.12) se conoce como kernel RBF (Radial Basis Function).

**3.4.2. Kernels lineales.** Como se estudió en la sección pasada, la obtención de la aplicación de características  $\Phi(\cdot)$  asociada a un kernel requería de una descomposición espectral. Sin embargo, el proceso de obtener un kernel desde a partir de una aplicación de características requiere únicamente de:

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

El caso  $\phi = I_d$ , es decir  $\phi(\mathbf{x}) = \mathbf{x}$  para todo  $\mathbf{x}$ , se obtiene el **kernel lineal**, este corresponde al producto interno del espacio donde se opera. En el caso  $\mathcal{X} = \mathbb{R}^D$  este se define de manera natural por:

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$$

**3.4.3. Matern Kernels.** Los **Matern kernel**, corresponden funciones del tipo,

$$k(r) = \frac{2^{(1-v)}}{\Gamma(v)} \left( \frac{\sqrt{2v}r}{l} \right)^v K_v \left( \frac{\sqrt{2v}r}{l} \right)$$

Donde  $r = \|\mathbf{x} - \mathbf{x}'\|$ ,  $v > 0$ ,  $l > 0$  y  $K_v$  corresponde a una función de Bessel modificada. Cuando  $v \rightarrow \infty$  este kernel se aproxima al kernel SE. Por otra parte, si  $v = \frac{1}{2}$  el kernel se simplifica a

$$k(r) = \exp(-r/l)$$

### 3.5. Kernels derivados de modelos probabilísticos generativos

En presencia de un modelo generativo probabilístico de vectores representados por  $p(\mathbf{x}|\theta)$ . Es posible utilizar la estructura de tal modelo para construir funciones kernel.

**3.5.1. Kernels producto de probabilidad.** En el contexto anterior, se define un kernel por medio de la siguiente identidad

$$(3.13) \quad k(\mathbf{x}_i, \mathbf{x}_j) = \int p(\mathbf{x}|\mathbf{x}_i)^\rho p(\mathbf{x}|\mathbf{x}_j)^\rho d\mathbf{x}$$

donde  $\rho > 0$  y  $p(\mathbf{x}|\mathbf{x}_i)$  se aproxima por  $p(\mathbf{x}|\hat{\theta}(\mathbf{x}_i))$ , donde  $\hat{\theta}(\mathbf{x}_i)$  es un parámetro estimado obtenido usando el dato  $\mathbf{x}_i$ . Este kernel es llamado **kernel producto de probabilidad** [19].

El modelo ajustado será usado para ver cuan similar son los objetos que se operan. En particular, si se ajusta el modelo a  $\mathbf{x}_i$  se desea que este sea capaz de indicar cuando otro dato  $\mathbf{x}_j$  es similar. Si a modo de ejemplo, se supone  $p(\mathbf{x}|\theta) = \mathcal{N}(\mu, \sigma^2 \mathbf{I})$  donde  $\sigma^2$  es fijo. Con  $\rho = 1$  y  $\hat{\mu}(\mathbf{x}_i) = \mathbf{x}_i$  y  $\hat{\mu}(\mathbf{x}_j) = \mathbf{x}_j$ , el kernel

$$(3.14) \quad k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{(4\pi\sigma^2)^{D/2}} \exp \left( -\frac{1}{(4\pi\sigma^2)} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right)$$

el cual corresponde a kernel de base radial. Se puede calcular la ecuación 3.13 por una variedad de modelos generativos, incluyendo aquellos con variables latentes.

**3.5.2. Kernels Fisher.** Una manera más eficiente para usar modelos generativos en definición kernels es usar **Kernel Fisher** [8] definido por

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{g}(\mathbf{x})^\top \mathbf{F}^{-1} \mathbf{g}(\mathbf{x}')$$

donde  $\mathbf{g}$  es el gradiente de la log verosimilitud o **vector de puntaje**, evaluado en el estimador MLE  $\hat{\theta}$ :

$$\mathbf{g}(\mathbf{x}) \triangleq \nabla_{\theta} \log p(\mathbf{x}|\theta)|_{\hat{\theta}}$$

$\mathbf{F}$  corresponde a la matriz de información de Fisher, la cual es corresponde esencialmente a la matriz Hessiana:

$$\mathbf{F} = \nabla \nabla \log p(\mathbf{x}|\theta)|_{\hat{\theta}}$$

Debido a que el estimador  $\hat{\theta}$  se obtiene en función de todos los datos, se tiene que similitud de  $\mathbf{x}$  y  $\mathbf{x}'$  es por tanto calculada usando todos los datos.

Si  $\mathbf{g}(\mathbf{x})$  es la dirección en la cual los datos  $\mathbf{x}$  condicionan el valor de  $\hat{\theta}$  al momento de maximizar la verosimilitud. Se desea considerar como similares a los puntos  $\mathbf{x}$  y  $\mathbf{x}'$  si las direcciones para el gradiente que condicionan la verosimilitud en  $\theta$  son similares entre si.

### 3.6. Kernels en modelos lineales generalizados

Con las herramientas estudiadas, es posible definir una clase de modelos basados en kernel para tareas de clasificación y regresión.

**3.6.1. Máquinas de kernel.** Se define una **Máquina de kernel** como un modelo lineal generalizado (GLM) donde el vector de características tiene la forma:

$$(3.15) \quad \phi(\mathbf{x}) = [k(\mathbf{x}, \mu_1), \dots, k(\mathbf{x}, \mu_K)]$$

donde  $\mu_k \in \mathcal{X}$  son un conjunto de  $K$  **centroides**. Si  $k$  es un kernel RBF, este modelo se denota por una red RBF. El vector de características presente en la ecuación 3.15 se denota como **vector de características kernelizado**. En este enfoque, no se requiere que el kernel sea Mercer.

Se puede usar el vector de características kernelizado entrenar un modelo de regresión logística definiendo  $p(y|\mathbf{x}, \theta) = \text{Ber}(\mathbf{w}^\top \phi(\mathbf{x}))$ . Esta formulación permite definir bordes de decisión no lineales de manera sencilla a través de la elección del kernel.

De la misma forma, es posible usar el vector de características anterior para entrenar un modelo de regresión lineal al definir  $p(y|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}), \sigma^2)$ .

**3.6.2. LIVMs, RVMs y Maquinas de vectores sparse.** Al formular modelos basados en Maquinas de kernel, no se especifica la manera optima para la obtención de centroides  $\mu_k$ . Un enfoque consiste en encontrar clusters en los datos y luego asignarlos como prototipos. Sin embargo, no se puede garantizar que a través de esta búsqueda se encuentren representaciones convenientes para la predicción. A esto se suma la elección del número de clusters, que tampoco esta especificada. Otro enfoque es considerar hacer cada elemento  $\mathbf{x}_i$  como un prototipo, de esta forma se obtiene

$$\phi(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_N)]$$

De esta forma, se tienen tantos parámetros con datos. Si bien esto supone un problema desde el punto de vista computacional, es posible subsanarlo por medio del uso de distribuciones prior  $w$  que induzcan dispersión (sparsity) en el regresor final, de esta forma, no solo se resuelve el problema referente a la alta cantidad de parámetros, sino que también la dispersión inducida en el vector de características permite escoger aquellos puntos más relevantes, permitiendo seleccionar un subconjunto de los datos para predecir. Los modelos que hacen uso de este principio se denominan **Maquinas de vectores sparse**.

En tal sentido, la elección más natural al momento de inducir dispersión es el uso de regularización  $l_1$  [9]. A la maquina de vectores resultante del uso de este regularizador para inducir dispersión se le denominará **L1VM** (Maquina de vectores regularizada  $l_1$ ). Análogamente, el uso de un regularizador  $l_2$  da lugar a **L2VM** (máquina de vectores  $l_2$ -regularizada), en este último modelo no se induce dispersión (debido al tipo de regularización).

Es también posible implementar el método ARD para inducir dispersión, el método resultante se denomina **Máquina de vectores de relevancia** o **RVM** [10].

Otra manera maquina de vector sparse los lo modelos conocidos como **máquinas de vectores de soporte** (SVM)

**3.6.3. Maquinas de vectores de soporte (SVM).** Este método se desarrolló para resolver problemas de clasificación y aproximación estática de funciones [14] [15], tanto para problemas de regresión como para reconocimiento de patrones. La formulación de esta metodología, adaptada a la estimación no lineal de funciones, comienza considerando una regresión de la forma:

$$(3.16) \quad f(x) = w^T \varphi(x) + b$$

En este contexto, la aplicación  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^{N_{\mathcal{H}}}$  corresponde la transformación de características sobre los datos. El parametro  $w \in \mathbb{R}^{N_{\mathcal{H}}}$  corresponde a al vector de pesos asociado a la construcción del modelo y cuya dimensión depende directamente de  $\varphi$ , pudiendo ser infinito dimensional.

Se busca estimar minimizando el riesgo empírico  $\mathbf{R}_{emp}$  definido por:

$$(3.17) \quad \mathbf{R}_{emp} = \frac{1}{N} \sum_{i=1}^N |y_i - w^T \varphi(x_i) - b|_{\varepsilon}$$

Donde se emplea la función de perdida  $\varepsilon$ -insensible de Vapnik, definida por:

$$(3.18) \quad |y - f(x)|_{\varepsilon} = \begin{cases} 0, & \text{si } |y - f(x)| \leq \varepsilon \\ |y - f(x)| - \varepsilon, & \text{en caso contrario} \end{cases}$$

Una primera aproximación a un modelo de entrenamiento, viene dada por el siguiente problema de optimización:

$$(3.19) \quad \left[ \begin{array}{l} \boxed{P} : \min_{w,b} J_P(w) = \frac{1}{2} w^T w \\ \text{s.a } y_i - w^T \varphi(x_i) - b \leq \varepsilon, \quad i = 1, \dots, N \\ \quad \quad w^T \varphi(x_i) + b - y_i \leq \varepsilon, \quad i = 1, \dots, N \end{array} \right]$$



Dicha formulación corresponde al caso en que todos los datos de entrenamiento se encuentran dentro una banda de ancho  $\varepsilon$ . Como es de esperar, elegir una precisión muy alta puede dejar puntos de entrenamiento fuera de dicha banda, lo cual hace que el problema (3.19) se vuelva infactible. Para solucionar tal inconveniente y obtener un modelo más flexible, se agregan variables de holgura  $\xi_i$ ,  $\xi_i^*$  para  $i = 1, \dots, N$ , de esta manera se permiten ciertos puntos fuera de la precisión preestablecida, al agregar dichas variables, se obtiene el siguiente problema:

$$(3.20) \quad \left[ \begin{array}{l} \boxed{P} : \min_{w, b, \xi, \xi^*} J_P(w, \xi, \xi^*) = \frac{1}{2} w^T w + c \sum_{i=1}^N (\xi_i + \xi_i^*) \\ \text{s.a } y_i - w^T \varphi(x_i) - b \leq \varepsilon + \xi_i, \quad i = 1, \dots, N \\ w^T \varphi(x_i) + b - y_i \leq \varepsilon + \xi_i^*, \quad i = 1, \dots, N \\ \xi_i, \xi_i^* \geq 0, \quad k = 1, \dots, N \end{array} \right]$$

Donde la constante  $c > 0$  determina la proporción de *desviaciones* toleradas a partir de cierta precisión  $\varepsilon$ . El Lagrangiano de este problema es:

$$(3.21) \quad \begin{aligned} \mathcal{L}(w, b, \xi, \xi^*; \alpha, \alpha^*, \eta, \eta^*) = & \frac{1}{2} w^T w + c \sum_{i=1}^N (\xi_i + \xi_i^*) - \sum_{i=1}^N \alpha_i (\varepsilon + \xi_i - y_i + w^T \varphi(x_i) + b) \\ & - \sum_{i=1}^N \alpha_i^* (\varepsilon + \xi_i^* + y_i - w^T \varphi(x_i) + b) - \sum_{i=1}^N (\eta_i \xi_i + \eta_i^* \xi_i^*) \end{aligned}$$

Donde los factores  $\alpha_i, \alpha_i^*, \eta_i, \eta_i^* \geq 0$  son los multiplicadores de Lagrange correspondientes. El punto estacionario de (3.21) se caracteriza de la siguiente manera:

$$(3.22) \quad \max_{\alpha, \alpha^*, \eta, \eta^*} \min_{w, b, \xi, \xi^*} \mathcal{L}(w, b, \xi, \xi^*; \alpha, \alpha^*, \eta, \eta^*)$$

De las condiciones de optimalidad se observa:

$$(3.23) \quad \left[ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial w} = 0 \rightarrow w = \sum_{i=1}^N (\alpha_i - \alpha_i^*) \varphi(x_i) \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0 \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \rightarrow c - \alpha_i - \eta_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \xi_i^*} = 0 \rightarrow c - \alpha_i^* - \eta_i^* = 0 \end{array} \right]$$

De lo cual se obtiene el siguiente problema dual:

$$(3.24) \quad \left[ \begin{array}{l} \boxed{D} : \max_{\alpha, \alpha^*} J_D(\alpha, \alpha^*) = -\frac{1}{2} \sum_{i,j=1}^N (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \varphi(x_i)^T \varphi(x_j) \\ \quad - \varepsilon \sum_{i=1}^N (\alpha_i + \alpha_i^*) + \sum_{i=1}^N y_i (\alpha_i - \alpha_i^*) \\ \text{s.a} \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0 \\ \quad \alpha_i, \alpha_i^* \in [0, c], \quad i = 1, \dots, N \end{array} \right]$$

Finalmente la representación dual del modelo de regresión es:

$$(3.25) \quad f(x) = \sum_{i=1}^N (\alpha_i - \alpha_i^*) \varphi(x_i)^T \varphi(x) + b$$

Los datos de entrada  $x_i$  tales que la diferencia  $\alpha_i - \alpha_i^*$  es no nula, son los llamados *vectores de soporte*. La estructura de (3.25) permite hacer la regresión sólo sobre estos elementos, por lo que este método induce o dispersión en la estimación (3.25).

Esta formulación responde a un esquema de regresión no lineal donde de induce dispersión por medio de restricciones, la formulación presentada presenta la gran desventaja de ser incalculable cuando la dimensión del espacio de características  $N_{\mathcal{H}}$  es infinita. Usando los resultados obtenidos anteriormente, es posible reformular por medio de la herramienta conocida como el **truco del kernel** [12].

El truco del kernel consiste en reemplazar el calculo de productos internos por una función kernel haciendo uso del resultado (2). En el caso de la formulación SVM, es posible reemplazar  $\varphi(x_i)^T \varphi(x_j)$  por una función kernel de la forma  $k(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$ . La estructura de este kernel no se conoce a priori, esto pues la definición de  $\varphi(\cdot)$  fue netamente teórica. Sin embargo, se sabe que su producto interno se puede representar por un kernel. La idea detrás de este método consiste en seleccionar un kernel de manera tal que las funciones base de su RKHS subyacente permitan una aproximación lo suficientemente flexible y regular.

Al aplicar el truco del kernel a los modelos de entrenamiento y regresión obtenidos anteriormente para SVM, se obtiene la versión *kernelizada*:

$$(3.26) \quad \left[ \begin{array}{l} \boxed{D} : \max_{\alpha, \alpha^*} J_D(\alpha, \alpha^*) = -\frac{1}{2} \sum_{i,j=1}^N (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) k(x_i, x_j) \\ \quad - \varepsilon \sum_{i=1}^N (\alpha_i + \alpha_i^*) + \sum_{i=1}^N y_i (\alpha_i - \alpha_i^*) \\ \text{s.a} \sum_{i=1}^N (\alpha_i - \alpha_i^*) = 0 \\ \quad \alpha_i, \alpha_i^* \in [0, c], \quad i = 1, \dots, N \end{array} \right]$$

Con su respectiva formula de regresión dual:

$$(3.27) \quad f(x) = \sum_{i=1}^N (\alpha_i - \alpha_i^*) k(x_i, x) + b$$

**3.6.4. LS-SVM.** Una variante del método SVM es su versión basada en mínimos cuadrados, la cual se formula como una alternativa más eficiente[13].

Su formulación es similar a la expuesta en 3.20, considerando la formula habitual de regresión presentada en (3.16):

$$(3.28) \quad \left[ \begin{array}{l} \boxed{P} : \min_{w, b, e} J_P(w, e) = \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{i=1}^N e_i^2 \\ \text{s.a } y_i = w^T \varphi(x_i) + b + e_i, \quad i = 1, \dots, N \end{array} \right]$$

En esta nueva formulación, los términos de error  $e_i$  juegan un papel similar al de las variables de holgura  $\xi_i$ ,  $\xi_i^*$  en (3.20), con excepción de que la función de perdida que se aplica a estas variables es ahora cuadrática. Se comienza entonces por establecer el Lagrangiano:

$$(3.29) \quad \mathcal{L}(w, b, e; \alpha) = J_P(w, e) - \sum_{i=1}^N \alpha_i (w^T \varphi(x_i) + b - e_i - y_i)$$

Donde los coeficientes  $\alpha_i$  son los multiplicadores de Lagrange. Las condiciones de optimalidad obtenidas son las siguientes:

$$(3.30) \quad \left[ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial w} = 0 \rightarrow w = \sum_{i=1}^N \alpha_i \varphi(x_i) \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_{i=1}^N \alpha_i = 0 \\ \frac{\partial \mathcal{L}}{\partial e_i} = 0 \rightarrow \alpha_i = \gamma e_i \quad i = 1, \dots, N \\ \frac{\partial \mathcal{L}}{\partial \alpha_i} = 0 \rightarrow w^T \varphi(x_i) + b + e_i - y_i = 0 \quad i = 1, \dots, N \end{array} \right]$$

Al despejar  $w$ ,  $e_i$  y reemplazar en la ultima igualdad se obtiene:

$$(3.31) \quad \sum_{i=1}^N \alpha_i \varphi(x_i)^T \varphi(x_j) + b + \frac{\alpha_j}{\gamma} = y_j \quad j = 1, \dots, N$$

Usando el truco del kernel en la forma:

$$(3.32) \quad \mathbf{K}_{ij} = k(x_i, x_j) = \varphi(x_i)^T \varphi(x_j) \quad j = 1, \dots, N$$

Reescribiendo (3.31) y haciendo uso de la notación  $\mathbf{y} = [y_1, \dots, y_N]^T$ ,  $\mathbf{1}_N = [1, \dots, 1]^T$ ,  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]^T$ , la condición  $\sum_{i=1}^N \alpha_i = 0$  obtenida en (3.30) y el resultado (3.32), se elabora el siguiente problema:

$$(3.33) \quad \left[ \begin{array}{c} \boxed{D} : \text{Resolver para } \boldsymbol{\alpha}, b \\ \left[ \begin{array}{c|c} 0 & \mathbf{1}_N^T \\ \hline \mathbf{1}_N & \mathbf{K} + \gamma^{-1} \mathbf{I} \end{array} \right] \left[ \begin{array}{c} b \\ \boldsymbol{\alpha} \end{array} \right] = \left[ \begin{array}{c} 0 \\ \mathbf{y} \end{array} \right] \end{array} \right]$$

Haciendo nuevamente uso de las condiciones obtenidas en (3.30), se puede verificar que el estimador obtenido es bastante similar en estructura al expuesto en (3.27), siendo este:

$$(3.34) \quad f(x) = \sum_{i=1}^N \alpha_i k(x, x_i) + b$$

El proceso de entrenamiento a través de (3.33), se ve facilitado numéricamente, ya que se basa en la solución de un sistema de ecuaciones lineales, el cual presenta una solución única cuando es de rango completo y que además, en su forma dual, no depende de la dimensión del espacio de características. Estas ventajas, sin embargo, no son gratuitas, pues como es posible observar en la obtención de las condiciones de primer orden del Lagrangiano (3.30), la ecuación  $\alpha_i = \gamma e_i$  anula la característica de *dispersión* que posee la formulación tradicional, esto se debe a que en este nuevo escenario, se dificulta la existencia de  $\alpha_i$  exactamente iguales a 0, dando el carácter de *vector de soporte* a todos los elementos presentes en el conjunto de entrenamiento, cada uno con cierto factor  $\alpha_i$  de influencia en el modelo.

**3.6.5. Clasificación kernel knn.** En un clasificador 1-NN se necesita calcular la distancia Euclidiana de un vector test a todos los puntos en entrenamiento, para luego encontrar al mas cercano, y asignar su etiqueta. Esto puede ser kernelizado al observar:

$$(3.35) \quad \|x_i x_{i'}\|_2^2 = \langle x_i, x_i \rangle + \langle x_{i'}, x_{i'} \rangle - 2 \langle x_i, x_{i'} \rangle$$

Esta reformulación permite la aplicación del clasificador knn una clase más amplia de problemas.

**3.6.6. Clustering de K-medoides Kernelizado.** El clustering K-medias al, igual que knn, usa la distancia euclidiana para medir similitud.

Para kernelizar este algoritmo, se reemplaza el algoritmo K-medias con el algoritmo **K-medoides**. El cual es similar a K-medias, pero en vez de representar cada centroide del cluster por la media de todos los vectores de datos asignados a este cluster, se hace que cada centroide sea uno de los vectores de datos por si mismo. En este caso se busca resolver:

$$m_k = \arg \min_{i: z_i = k} \sum_{i': z_{i'} = k} d(i, i')$$

donde

$$d(i, i') \triangleq \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2^2$$

Esto toma  $O(n_k^2)$  trabajo por cluster, mientras K-medias toma  $O(n_k D)$  para actualizar cada cluster. En el algoritmo (5) se aprecia la estructura del método.

Al calcular el medoide más cercano a la clase, el problema se transforma en **clasificación medoide mas cercano**[11] Este algoritmo puede ser kernelizado usando Ecuación (3.35) para reemplazar la distancia  $d(i, i')$ .

---

**Algorithm 5:** Algoritmo kernel K-medoids

---

```

1 Inicializar:  $m_1, \dots, m_K$  subconjunto aleatorio de tamaño  $K$  en
    $\{1, \dots, N\}$ ;
2 repeat
3    $z_i = \arg \min_k d(i, m_k)$  para  $i = 0, \dots, N$ 
4    $m_k \leftarrow \arg \min_{i: z_i=k} \sum_{i': z_{i'}=k} d(i, i')$ 
5 until convergencia

```

---

**3.6.7. Kernel PCA.** PCA requiere encontrar vectores propios de la matriz de covarianza muestral  $\mathbf{S} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top = (1/N) \mathbf{X}^\top \mathbf{X}$ . Sin embargo, se puede también calcular PCA, encontrando los vectores propios de la matriz de producto interno  $\mathbf{X} \mathbf{X}^\top$ , esto permitirá derivar el método **kernel PCA**[7].

Primero, sea  $\mathbf{U}$  la matriz ortogonal conteniendo los vectores propios de  $\mathbf{X} \mathbf{X}^\top$  con los correspondientes valores propios en  $\Lambda$ . Por definición, se tiene  $(\mathbf{X} \mathbf{X}^\top) \mathbf{U} = \mathbf{U} \Lambda$ . Premultiplicando por  $\mathbf{X}^\top$ :

$$(\mathbf{X}^\top \mathbf{X})(\mathbf{X}^\top \mathbf{U}) = (\mathbf{X}^\top \mathbf{U}) \Lambda$$

de donde se ve que los vectores propios de  $(\mathbf{X}^\top \mathbf{X})$  son  $\mathbf{V} = \mathbf{X}^\top \mathbf{U}$  con los valores propios dados por  $\Lambda$ . Sin embargo, tales vectores propios no se encuentran normalizados. Dado  $\|v_j\|^2 = \mathbf{u}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{u}_j = \lambda_j \mathbf{u}_j^\top \mathbf{u}_j = \lambda_j$ . por lo que es posible obtener un vector propio normalizado por medio de  $\mathbf{V}_{pca} = \mathbf{X}^\top \mathbf{U} \Lambda^{-1/2}$ .

Sea ahora  $\mathbf{K} = \mathbf{X} \mathbf{X}^\top$  la matriz de Gram del kernel  $k(\cdot, \cdot)$ . Sea  $\Phi$  la matriz diseño correspondiente a la transformación de características inducida por  $k(\cdot, \cdot)$ .  $\mathbf{S}_\phi = \frac{1}{N} \sum_i \phi_i \phi_i^\top$  la correspondiente matriz de covarianza en el espacio de características. Los vectores propios son dados por  $\mathbf{V}_{kpca} = \Phi^\top \mathbf{U} \Lambda^{-1/2}$  donde  $\mathbf{U}$  y  $\Lambda$  contienen los vectores propios y valores propios de  $\mathbf{K}$ . Claramente no es posible calcular  $\mathbf{V}_{kpca}$ , dado que  $\Phi_i$  es potencialmente de dimensión infinita. Sin embargo, se puede calcular la proyección de un vector  $\mathbf{x}$  sobre el espacio característica como sigue.

$$\phi_*^\top \mathbf{V}_{kpca} = \phi_*^\top \Phi \mathbf{U} \Lambda^{-\frac{1}{2}} = \mathbf{k}_*^\top \mathbf{U} \Lambda^{-\frac{1}{2}}$$

donde  $\mathbf{k}_* = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]$ .

Hasta este punto, se asume que los datos proyectados tienen media cero, lo cual en general no es el caso. Para sortear tal detalle, se define el vector de características centrado como  $\tilde{\phi}_i = \phi(\mathbf{x}_i) - \frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}_j)$ . La matriz Gram de vectores característicos centrados esta dada por:

$$\begin{aligned}\tilde{K}_{ij} &= \tilde{\phi}_i^T \tilde{\phi}_j \\ &= \phi_i^T \phi_j - \frac{1}{N} \sum_{k=1}^N \phi_i^T \phi_k - \frac{1}{N} \sum_{k=1}^N \phi_j^T \phi_k + \frac{1}{N^2} \sum_{k=1}^N \sum_{l=1}^M \phi_k^T \phi_l\end{aligned}$$

De donde

$$\tilde{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{N} \sum_{k=1}^N k(\mathbf{x}_i, \mathbf{x}_k) - \frac{1}{N} \sum_{k=1}^N k(\mathbf{x}_j, \mathbf{x}_k) + \frac{1}{N^2} \sum_{k=1}^N \sum_{l=1}^M k(\mathbf{x}_k, \mathbf{x}_l)$$

Lo cual puede ser expresado en notación matricial como:

$$\tilde{\mathbf{K}} = \mathbf{H}\mathbf{K}\mathbf{H}$$

donde  $\mathbf{H} \triangleq \mathbf{I} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$  es la **matriz centrada**. El proceso seguido anteriormente, se resume en el algoritmo 6.

---

**Algorithm 6:** kernel PCA

---

**Input:**  $\mathbf{K}$  de tamaño  $N \times N$ ,  $\mathbf{K}_*$  de tamaño  $N_* \times N$ , numero de dimensiones latentes  $L$

- 1  $\mathbf{O} = \mathbf{1}_N \mathbf{1}_N^T / N$
- 2  $\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{O}\mathbf{K} - \mathbf{K}\mathbf{O} + \mathbf{O}\mathbf{K}\mathbf{O}$
- 3  $[\mathbf{U}, \Lambda] = \text{eig} \tilde{\mathbf{K}}$
- 4 **for**  $i = 1, \dots, N$  **do**
- 5    $\mathbf{v}_i = \mathbf{u}_i / \sqrt{\lambda_i}$
- 6 **end**
- 7  $\mathbf{O}_* = \mathbf{1}_{N_*} \mathbf{1}_N^T / N$
- 8  $\tilde{\mathbf{K}}_* = \mathbf{K}_* - \mathbf{O}_* \mathbf{K}_* - \mathbf{K}_* \mathbf{O}_* + \mathbf{O}_* \mathbf{K}_* \mathbf{O}_*$
- 9  $Z = \tilde{\mathbf{K}} \mathbf{V}(:, 1:L);$

---

### 3.7. Kernels para construir modelos generativos

Otra case kernels, conocidos como kernels de suavizado, permiten caracterizar estimaciones no paramétricas de densidades. Lo cual puede ser utilizado para generar estimaciones de densidades de manera no supervisada, así como también para crear modelos generativos tanto para regresión como para clasificación.

**3.7.1. Kernel de suavizado.** Un **kernel suavizado** es una función de un argumento que satisface las siguientes propiedades:

$$\int k(x) dx = 1, \int x k(x) dx = 0, \int x^2 k(x) dx > 0$$

Un ejemplo sencillo es el **kernel gaussiano**:

$$k(x) \triangleq \frac{1}{(2\pi)^{\frac{1}{2}}} e^{-x^2/2}$$

En este caso, se puede controlar el ancho de banda del kernel introduciendo un parámetro  $h$ :

$$k_h(x) \triangleq \frac{1}{h} k\left(\frac{x}{h}\right)$$

Esto puede ser generalizado a un kernel RBF:

$$k_h(\mathbf{x}) = k_h(\|\mathbf{x}\|)$$

No obstante, en el caso del kernel gaussiano, el kernel de suavizado pasa a ser:

$$k_h(\mathbf{x}) = \frac{1}{h^D (2\pi)^{D/2}} \prod_{j=1}^D \exp\left(-\frac{1}{2h^2} x_j^2\right)$$

Si bien los kernels de la familia gaussiana poseen propiedades que permiten simplificar el gasto computacional, su soporte no es acotado, por lo que en modelos reales carece de expresividad. Una alternativa a esto es el **kernel Epanechnikov**, el cual posee soporte compacto y está definido por:

$$k(x) \triangleq \frac{3}{4} (1 - x^2) \mathbb{I}(|x| \leq 1)$$

Desafortunadamente, el kernel Epanechnikov no es diferenciable en los bordes del soporte. Una alternativa es el **kernel tri-cubo**, definido como sigue:

$$k(x) \triangleq \frac{70}{81} (1 - |x|^3)^3 \mathbb{I}(|x| \leq 1)$$

Este tiene soporte compacto y tiene dos derivadas continuas en los bordes del soporte. Por otra parte, se puede obtener un kernel a partir de la distribución uniforme:

$$k(x) \triangleq \mathbb{I}(|x| \leq 1)$$

Denominado **kernel boxcar**.

**3.7.2. kernel density estimation (KDE).** El modelo GMM estudiado anteriormente corresponde a un estimador de densidad paramétrico para datos en  $\mathbb{R}^D$ . Este requiere especificar el número de componentes  $K$  y ubicaciones  $\boldsymbol{\mu}_k$  de los cluster. Una alternativa para la estimación de  $\boldsymbol{\mu}_k$  es ubicar un cluster central por punto de datos, tal que  $\boldsymbol{\mu}_i = \mathbf{x}_i$ . En este caso el modelo se convierte:

$$p(\mathbf{x}|\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(\mathbf{x}|\mathbf{x}_i, \sigma^2 \mathbf{I})$$

Lo cual se puede generalizar escribiendo:

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i)$$

Modelo conocido como estimador **Parzen window density**, o **estimador de densidad de kernel (KDE)**, y corresponde a una formulación no paramétrica. La ventaja sobre modelos paramétricos es que no se requiere el modelo ajustado y no existe la necesidad de elegir  $K$ . La desventaja es que el modelo toma mucho almacenamiento de memoria y gran tiempo para evaluar.

**3.7.3. De KDE a KNN.** Se puede usar KDE para definir las clases de densidades condicionales en un clasificador generativo. Esto proporciona una derivación alternativa del clasificador KNN. En KDE con un kernel boxcar se fijara el ancho de banda y se contará cuantos puntos caen dentro del hyper cubo centrado sobre cada dato. Si se modela el ancho de banda  $h$  como un parámetro de modelo, se permite que el ancho de banda o volumen sea distinto para cada dato. Específicamente, se aumentará el volumen alrededor de  $\mathbf{x}$  hasta encontrar  $K$  datos, sin importar su clase. Sea el volumen resultante  $V$  con  $V(\mathbf{x})$  y  $N_c(\mathbf{x})$  ejemplos de clase  $c$  en tal volumen. Luego es posible estimar la densidad condicional la de clase  $c$  de como sigue:

$$p(\mathbf{x}|y = c, \mathcal{D}) = \frac{N_c(\mathbf{x})}{N_c V(\mathbf{x})}$$

Donde  $N_c$  es el número total de ejemplos en la clase  $c$  en aquel conjunto de datos. La clase a priori puede ser modelada como

$$p(y = c|\mathcal{D}) = \frac{N_c}{N}$$

Por lo tanto, la posterior de clase esta dada por:

$$p(y = c|\mathbf{x}, \mathcal{D}) = \frac{\frac{N_c(\mathbf{x})}{N_c V(\mathbf{x})} \frac{N_c}{N}}{\sum_{c'} \frac{N_{c'}(\mathbf{x})}{N_{c'} V(\mathbf{x})} \frac{N_{c'}}{N}} = \frac{N_c(\mathbf{x})}{\sum_{c'} N_{c'}(\mathbf{x})} = \frac{N_c(\mathbf{x})}{K}$$

Donde se usa el hecho que  $\sum_c N_c(\mathbf{x}) = K$ , dado que se elige un total de  $K$  puntos alrededor de cada punto.

**3.7.4. Regresión por kernel.** Es posible usar KDE para regresión. El objetivo es calcular la esperanza condicionada

$$f(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}] = \int y p(y|\mathbf{x}) dy = \frac{\int y p(\mathbf{x}, y) dy}{\int p(\mathbf{x}, y) dy}$$

Se puede usar KDE para aproximar la densidad conjunta  $p(\mathbf{x}, y)$  como sigue

$$p(\mathbf{x}, y) \approx \frac{1}{N} \sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i) k_h(y - y_i)$$

Por lo tanto

$$\begin{aligned} f(\mathbf{x}) &= \frac{\frac{1}{N} \sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i) \int y k_h(y - y_i) dy}{\frac{1}{N} \sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i) \int k_h(y - y_i) dy} \\ &= \frac{\sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i) y_i}{\sum_{i=1}^N k_h(\mathbf{x} - \mathbf{x}_i)} \end{aligned}$$

Lo cual se obtiene al observar que la integral  $\int k_h(y - y_i) dy = 1$ . Además del hecho que  $\int y k_h(y - y_i) dy = y_i$ . Esto sigue definiendo  $x = y - y_i$  y usando la propiedad de media igual a cero de kernel suavizado:

$$\int (x + y_i) k_h(x) dx = \int x k_h(x) dx + y_i \int k_h(x) dx = 0 + y_i = y_i$$



Se puede reescribir el resultado anterior de la siguiente manera:

$$f(\mathbf{x}) = \sum_{i=1}^N w_i(\mathbf{x}) y_i$$

$$w_i(\mathbf{x}) \triangleq \frac{k_h(\mathbf{x} - \mathbf{x}_i)}{\sum_{i'=1}^N k_h(\mathbf{x} - \mathbf{x}_{i'})}$$

La predicción es por tanto una suma ponderada de los resultados de los datos de entrenamiento, donde los pesos dependen de que tan similar es  $\mathbf{x}$  de los puntos de entrenamientos almacenados. Este método es llamado **regresión kernel**, **kernel suavizado** o el modelo **Nadaraya-Watson**.

**3.7.5. Regresión localmente ponderada.** Si se define  $k_h(\mathbf{x} - \mathbf{x}_i) = k(\mathbf{x}, \mathbf{x}_i)$  se puede reescribir la predicción hecha mediante regresión kernel como sigue

$$\hat{f}(\mathbf{x}_*) = \sum_{i=1}^N y_i \frac{k(\mathbf{x}_*, \mathbf{x}_i)}{\sum_{i'=1}^N k(\mathbf{x}_*, \mathbf{x}_{i'})}$$

Si en este caso  $k(\mathbf{x}, \mathbf{x}_i)$  fuera un kernel suavizado, se podría ignorar el termino normalización, en tal caso:

$$\hat{f}(\mathbf{x}_*) = \sum_{i=1}^N y_i k(\mathbf{x}_*, \mathbf{x}_i)$$

Este modelo consiste esencialmente en una función fija localmente constante. Este modelo se puede mejorar fijando un modelo de regresión lineal para cada punto  $\mathbf{X}_*$  resolviendo:

$$\min_{\beta(\mathbf{x}_*)} \sum_{i=1}^N k(\mathbf{x}_*, \mathbf{x}_i) \left[ y_i - \beta(\mathbf{x}_*)^T \phi(\mathbf{x}_i) \right]^2$$

donde  $\phi(\mathbf{x}) = [1, \mathbf{x}]$ . Esto es llamado **regresión ponderada localmente**. Un ejemplo de tal metodo es **LOESS** o **LOWES** que significa para “locally-weighted scatterplot smoothing”. Se puede calcular el parámetro  $\beta(\mathbf{x}_*)$  para cada caso de test resolviendo el siguiente problema de mínimos cuadrados ponderados:

$$\beta(\mathbf{x}_*) = (\Phi^T \mathbf{D}(\mathbf{x}_*) \Phi)^{-1} \Phi^T \mathbf{D}(\mathbf{x}_*) \mathbf{y}$$

donde  $\Phi$  es una matriz de diseño de  $N \times (D+1)$  y  $\mathbf{D} = \text{diag}(k(\mathbf{x}_*, \mathbf{x}_i))$ . La predicción correspondiente tiene la forma

$$\hat{f}(\mathbf{x}_*) = \phi(\mathbf{x}_*)^T \beta(\mathbf{x}_*) = (\Phi^T \mathbf{D}(\mathbf{x}_*) \Phi)^{-1} \Phi^T \mathbf{D}(\mathbf{x}_*) \mathbf{y} = \sum_{i=1}^N w_i(\mathbf{x}_*) y_i$$

El termino  $w_i(\mathbf{x}_*)$  combina el kernel local suavizado con el efecto de la regresión lineal, es llamado el **kernel equivalente**.



## CHAPTER 4

### Discusión

En este reporte se busco compilar ciertos tópicos de interés pertenecientes al aprendizaje automático. En específico se abordaron temas referentes a modelos gráficos probabilísticos, haciendo hincapié en modelos gráficos dirigidos. En este contexto, el paso natural es expandir dicha sección para tratar temas relacionados a modelos gráficos probabilísticos no dirigidos y grafos factor. En segunda instancia se abordaron métodos de inferencia aproximada basados tanto en MCMC como en inferencia variacional con enfoque en el algoritmo CAVI. En cuanto a la tercera y ultima sección, se revisaron las bases y propiedades referentes al los métodos basados en kernel.

Para futuras versiones de esta monografía, se propone desarrollar teoría relacionada a modelos bayesianos en aprendizaje automático.

Trabajar en este proyecto me permitió conocer nuevas áreas y perspectivas en aprendizaje automático. Doy las gracias a Felipe Tobar por su iniciativa y por permitirme desarrollar mis habilidades en este campo.



## APPENDIX A

### Anexos

#### A.0.1. Tópicos de Análisis.

A.0.1.1. *Normas.* Una norma sobre  $\mathbb{R}^n$  es una función  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  que satisface las siguientes propiedades:

- (1)  $f(x) \geq 0$   $x \in \mathbb{R}^n$  ( $f(x) = 0$  si y solo si  $x = 0$ )
- (2)  $f(x + y) \leq f(x) + f(y)$   $x, y \in \mathbb{R}^n$
- (3)  $f(\alpha x) = |\alpha| f(x)$   $\alpha \in \mathbb{R}, x \in \mathbb{R}^n$

Una norma  $f(\cdot)$  cumpliendo lo anterior, se denota por  $f(x) = \|x\|$ . Ejemplos de normas son:

$$\begin{aligned} \|x\|_p &= (|x_1|^p + \dots + |x_n|^p)^{1/p}, \quad p \geq 1 \\ \|x\|_1 &= |x_1| + \dots + |x_n| \\ \|x\|_2 &= (x^T x)^{1/2} \\ \|x\|_\infty &= \max_{1 \leq i \leq n} |x_i| \end{aligned}$$

En la práctica, si no se especifica la norma  $\|x\|$  se asume por la norma-2  $\|x\|_2$ . Se puede además definir la noción de ángulo  $\theta$  entre los vectores  $x$  e  $y$  a través de la norma y el producto entre ellos:

$$\cos(\theta) = \frac{x^T y}{\|x\|_2 \|y\|_2}$$

A.0.1.2. *Cauchy Schwarz.* Sea  $E$  un espacio vectorial real y sea  $\langle \cdot, \cdot \rangle : E \rightarrow \mathbb{R}$  una función positiva, simétrica y bilinear, es decir:

- (1)  $\langle x, x \rangle \geq 0$
- (2)  $\langle x, y \rangle = \langle y, x \rangle$
- (3)  $\langle \alpha x + y, z \rangle = \alpha \langle x, z \rangle + \langle y, z \rangle$

Para todo  $x, y, z \in E$ ,  $\alpha \in \mathbb{R}$ . Se quiere demostrar que  $\langle \cdot, \cdot \rangle$  cumple la desigualdad de Cauchy-Schwarz, es decir:

$$|\langle x, y \rangle|^2 \leq \langle x, x \rangle \cdot \langle y, y \rangle \quad x, y \in E$$

PROOF. Sean  $x, y \in E$  y  $\alpha \in \mathbb{R}$ , debido a la positividad de  $\langle \cdot, \cdot \rangle$  se cumple:

$$0 \leq \langle \alpha x + y, \alpha x + y \rangle$$

Lo cual, aplicando la bilinearidad y simetría de  $\langle \cdot, \cdot \rangle$ , se puede reescribir en la forma:

$$0 \leq \alpha^2 \langle x, x \rangle + 2\alpha \langle x, y \rangle + \langle y, y \rangle$$

Si  $\langle x, x \rangle = 0$  o  $\langle y, y \rangle = 0$ , la desigualdad se cumple de manera trivial. Por otra parte, nuevamente debido a la positividad de  $\langle \cdot, \cdot \rangle$ , si  $\langle x, x \rangle > 0$ , tomar  $\alpha = -\frac{\langle x, y \rangle}{\langle x, x \rangle}$  reduce la expresión anterior a:

$$0 \leq -\frac{|\langle x, y \rangle|^2}{\langle x, x \rangle} + \langle y, y \rangle$$

De la arbitrariedad de  $x, y \in E$  se concluye la afirmación buscada.  $\square$

### A.0.2. Tópicos de Teoría de la Medida.

**Definición A.1** (medida). Sea  $X$  un conjunto y  $\mathcal{C} \subseteq \mathcal{P}(X)$  tal que  $\varphi \in \mathcal{C}$ . Una función:

$$\mu : \mathcal{C} \rightarrow \overline{\mathbb{R}}_+$$

Se dice que es una **medida** sobre  $\mathcal{C}$  si satisface:

- (1)  $\mu(\varphi) = 0$
- (2) Si  $(A_n)_{n \in \mathbb{N}} \subseteq \mathcal{C}$  es una sucesión de conjuntos disjuntos tales que  $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{C}$ , entonces:

$$\mu\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} \mu(A_n)$$

(esta propiedad se llama  $\sigma$ -aditividad).

**Definición A.2** ( $\sigma$ -álgebra).  $\mathcal{T} \subseteq \mathcal{P}(X)$  se dice que es una  $\sigma$ -álgebra si:

- (1)  $X \in \mathcal{T}$
- (2)  $A, B \in \mathcal{T} \Rightarrow A \setminus B \in \mathcal{T}$
- (3)  $(A_n)_{n \in \mathbb{N}} \subseteq \mathcal{T} \Rightarrow \bigcup_{n \in \mathbb{N}} (A_n) \in \mathcal{T}$

Si tenemos un conjunto  $X$ , una  $\sigma$ -álgebra  $\mathcal{T} \subseteq \mathcal{P}(X)$ , y una medida  $\mu$  sobre  $\mathcal{T}$ , lo identificaremos con la tupla:

$$(X, \mathcal{T}, \mu)$$

**Definición A.3.** (Medida  $\sigma$ -finita) Sea  $\mathcal{C} \subseteq \mathcal{P}(X)$  tal que  $\phi \in \mathcal{C}$ . Una medida  $\mu : \mathcal{C} \rightarrow \overline{\mathbb{R}}_+$  se dirá:

- (1)  $\sigma$ -finita si:

$$\exists (A_n)_{n \in \mathbb{N}} \subseteq \mathcal{C} \text{ tal que } X = \bigcup_{n \in \mathbb{N}} A_n \text{ y } \mu(A_n) < +\infty, \forall n \in \mathbb{N}$$

(2) *finita* si:

$$\mu(A_n) < +\infty, \forall A \in \mathcal{C}$$

### A.0.3. Funciones Simples.

**Definición A.4.** (Función Simple) Sea  $(X, \mathcal{T})$  un espacio medible. Una función:

$$f : X \rightarrow \mathbb{R}$$

se dirá simple si:

$$\exists (A_i)_{i=1}^n \subseteq \mathcal{T}, \exists (a_i)_{i=1}^n \subseteq \mathbb{R} \text{ tal que } f(x) = \sum_{i=1}^n a_i 1_{A_i}(x)$$

Denotaremos por  $\xi = \xi(X, \mathcal{T}, \mathbb{R})$  al conjunto de todas las funciones simples a valores reales. Análogamente denotaremos por  $\xi_+ = \xi(X, \mathcal{T}, \mathbb{R}_+)$  al conjunto de funciones simples a valores en  $\mathbb{R}_+$ .

Notemos que si  $f$  es una función simple, entonces existe  $(A_i)_{i=1}^n$  partición medible de  $X$  tal que:

$$f(x) = \sum_{i=1}^n a_i 1_{A_i}(x)$$

### A.0.4. Definición de la integral.

**Definición A.5.** (Integral) Sea  $(X, \mathcal{T}, \mu)$  un espacio de medida. Sea  $f \in \xi_+$ , digamos  $f(x) = \sum_{i=1}^n a_i 1_{A_i}(x)$ , con  $(A_i)_{i=1}^n$  partición medible de  $X$ . Se define la **integral** de  $f$  como:

$$\int f d\mu := \sum_{i=1}^n a_i \mu(A_i) \in \overline{\mathbb{R}}_+$$

### A.0.5. Completitud de un espacio de medida.

**Definición A.6.** (Conjunto Despreciable) Sea  $(X, \mathcal{T}, \mu)$  un espacio de medida. Un conjunto  $N \in \mathcal{P}(X)$  se dice que es un **despreciable** si:

$$\exists A \in \mathcal{T} \text{ tal que } N \subseteq A \text{ y } \mu(A) = 0$$

La  $\sigma$ -álgebra  $\mathcal{T}$  se dirá completa si contiene a todos los despreciables.

**Definición A.7.** Si  $(X, \mathcal{T}, \mu)$  un espacio de medida y  $p : X \rightarrow \{V, F\}$  una proposición (V=verdadero, F=falso). Se dirá que la proposición  $p$  **se satisface en casi todas partes** (lo cual se abrevia ctp ó  $\mu$ -ctp) si el siguiente conjunto:

$$\{x \in X / p(x) = F\}$$

es un despreciable para  $(X, \mathcal{T}, \mu)$ .

**A.0.6. Probabilidad.** Sea  $(\Omega, \mathcal{B}, \mathbb{P})$  espacio de probabilidades, donde :

- $\Omega$  : espacio muestral
- $\mathcal{B}$  sigma álgebra

- $\mathbb{P}$ : ley de probabilidades

A continuación se definen algunos conceptos de probabilidad de tal manera de explicar teoremas y proposiciones que se utilizarán a lo largo del texto.

**Definición A.8.** Probabilidad Condicional:  $\mathbb{P}(X|Y) = \frac{\mathbb{P}(X,Y)}{\mathbb{P}(Y)}$

**Definición A.9.** independencia:  $X \perp Y \iff \mathbb{P}(X, Y) = \mathbb{P}(X)\mathbb{P}(Y)$

**Definición A.10.** Independencia Condicional:  $X \perp Y|E \iff \mathbb{P}(X, Y|E) = \mathbb{P}(X|E)\mathbb{P}(Y|E)$

A partir de esto, se pueden formular los siguientes teoremas:

**Proposición 4.** Teorema de Bayes: Se define como  $\mathbb{P}(X|E) = \frac{\mathbb{P}(E|X)\mathbb{P}(X)}{\mathbb{P}(E)}$ , donde  $\mathbb{P}(X|E)$  es la distribución a posteriori de  $X$  dado  $E$ .  $\mathbb{P}(E|X)$  es la verosimilitud de los datos.  $\mathbb{P}(X)$  es la probabilidad a priori de  $X$  y  $\mathbb{P}(E)$  es la probabilidad de los datos.

**Proposición 5.** Regla de la Cadena: Sean  $X_1, \dots, X_N$  eventos medibles, entonces la probabilidad conjunta puede expresarse como:

$$\mathbb{P}(X_1, \dots, X_N) = \mathbb{P}(X_1) \prod_{i=2}^N \mathbb{P}(X_i|X_1, \dots, X_{i-1})$$

**Proposición 6.** Regla de la Marginación: Sea  $\{E_i\}_{i \in \mathbb{N}}$  una partición de  $\Omega$ , tales que  $\bigcup_{i \in \mathbb{N}} E_i = \Omega$  entonces:

$$\mathbb{P}(X) = \sum_{i \in \mathbb{N}} \mathbb{P}(X|E_i)\mathbb{P}(E_i)$$

**Proposición 7.** Dos eventos  $X$  e  $Y$  son independientes si y solo si la distribución a priori y a posteriori son iguales, es decir,  $X \perp Y \iff \mathbb{P}(X|Y) = \mathbb{P}(X)$

**Proposición 8.** Dos eventos  $X$  e  $Y$  son independientes condicionados al evento  $E$  si y solo si, la distribución a posteriori de  $X$  dado  $Y, E$  es igual a la distribución a posteriori de  $X$  dado solo  $E$ , es decir  $X \perp Y|E \iff \mathbb{P}(X|Y, E) = \mathbb{P}(X|E)$

**Definición A.11.** Una variable aleatoria  $x$  es una función de  $\Omega$  a valores en  $\mathbb{R}$ , es decir,  $x(\omega) \in \mathbb{R}$  con  $\omega \in \Omega$

**Definición A.12.** Dada una variable aleatoria  $x$ , se define su esperanza y varianza de la siguiente manera:

$$\begin{aligned} \mu_x &= \mathbb{E}[x] = \int_{\Omega} x(\omega) d\mathbb{P}(\omega) \\ \sigma_x^2 &= \mathbb{V}[x] = \int_{\Omega} (x(\omega) - \mu_x)^2 d\mathbb{P}(\omega) \end{aligned}$$

**Definición A.13.** Dada una variable aleatoria  $x$ , se dice que  $x$  sigue una distribución de probabilidades  $F(c)$  si se cumple que  $\mathbb{P}(\omega : x(\omega) \leq c) = F(c)$ . Se dice que  $x$  tiene función de densidad  $p(x)$  si se cumple que:

$$F(c) = \int_{-\infty}^c p(x) dx = \int_{\omega: x(\omega) \leq c} d\mathbb{P}(\omega)$$

Se conoce como el soporte de la distribución  $F$  a la imagen de la variable aleatoria  $im(x) = x(\Omega) \subseteq \mathbb{R}$ .

Si el soporte de  $F$  es numerable entonces se dirá que la variable aleatoria sigue una



distribucion discreta, en caso contrario, se dirá que es una distribución continua. Para el primer caso, se utilizan sumatorias:

$$F(c) = \sum_{i \leq c} p_i$$

$$p(x) = \sum_{i \leq c} p_i \delta_i(x)$$

Donde  $\delta_i$  corresponde a la función delta dirac.

**Proposición 9.** Dada una variable aleatoria  $x$ , su esperanza y varianza son:

$$\mu_x = \mathbb{E}[x] = \int_{\mathbb{R}} xp(x)dx$$

$$\sigma_x^2 = \mathbb{V}[x] = \int_{\mathbb{R}} (x - \mu_x)^2 p(x)dx$$

#### A.0.7. Tópicos de Optimización.

A.0.7.1. *Problema de Optimización con Restricciones, Lagrangiano.* Sea el problema de optimización con restricciones:

$$\begin{aligned} \min f(x) \quad & x \in \mathbb{R}^n \\ \text{s.a} \quad & c_i(x) = 0, \quad i \in \mathcal{C}_E \\ & c_i(x) \geq 0, \quad i \in \mathcal{C}_I \end{aligned}$$

donde  $f(x)$  es la función objetivo,  $c_i$  con  $i = 1, 2, \dots, p$  son las funciones de restricción,  $\mathcal{C}_E$  es el conjunto de índices de las restricciones de igualdad en el problema y  $\mathcal{C}_I$  es el conjunto de restricciones de desigualdad. Cualquier punto que satisface todas las restricciones es llamado un punto factible y el conjunto de todos aquellos puntos es referido como la región factible. Se definen las restricciones activas en un punto  $x'$  por el conjunto de índices  $\mathcal{A}(x') = \{i : c_i(x') = 0\}$  tales que cualquier restricción es activa en  $x'$  si  $x'$  esta en el borde o en la región factible. En este contexto se define el Lagrangiano por:

$$\mathcal{L}(x, \lambda) = f(x) - \sum_i \lambda_i c_i(x)$$

A.0.7.2. *Condiciones KKT.* En un problema de optimización con restricciones se tiene como condición necesaria de primer orden lo siguiente:

Si  $x^*$  es un mínimo local del problema de optimización con restricciones y se mantiene cierta regularidad alrededor de  $x^*$ , entonces existen multiplicadores de Lagrange  $\lambda^*$  tales que  $x^*, \lambda^*$  satisfacen el siguiente sistema:

$$\begin{aligned}
\nabla_x \mathcal{L}(x, \lambda) &= 0 \\
c_i(x) &= 0 \quad i \in \mathcal{C}_E \\
c_i(x) &\geq 0 \quad i \in \mathcal{C}_I \\
\lambda_i &\geq 0 \quad i \in \mathcal{C}_I \\
\lambda_i c_i(x) &= 0 \quad \forall i
\end{aligned}$$

El cual consiste en las llamadas condiciones de Karush-Kuhn-Tucker (KKT). Un punto  $x^*$  que satisface las condiciones se refiere como punto KKT. La última condición  $\lambda_i c_i(x) = 0$  es llamada *condición de complementariedad* y afirma que  $\lambda_i$  y  $c_i$  no pueden ser ambos distintos de cero, o equivalentemente que las restricciones inactivas tienen un multiplicador igual a cero. Si no existe  $i$  tal que  $\lambda_i^* = c_i^*(x) = 0$  entonces se mantiene una complementariedad estricta. En cambio si  $\lambda_i^* = 0$  se denomina fuertemente activa,  $c_i^* > 0$  débilmente activa  $\lambda_i^* = c_i^* = 0$  e inactiva si  $\lambda_i^* = 0, c_i^* > 0$ .

Las Condiciones de segundo orden pueden ser expresadas en términos de la matriz Hessiana con respecto a  $x$  del Lagrangiano  $\nabla_x^2 \mathcal{L}(x^*, \lambda^*) = \nabla^2 f(x^*) - \sum_i \lambda_i^* \nabla^2 c_i(x^*)$ .

A.0.7.3. *Problema de Programación Cuadrática.* Un problema de programación cuadrática es de la forma:

$$\begin{aligned}
\text{Minimizar } q(x) &= \frac{1}{2} x^T G x + g^T x \\
\text{sujeto a } a_i^T x &= b_i, \quad i \in \mathcal{C}_E \\
a_i^T x &\geq b_i, \quad i \in \mathcal{C}_I
\end{aligned}$$

Si la matriz Hessiana  $G$  es semidefinida positiva, la solución  $x^*$  es global. Si  $G$  es definida positiva la solución  $x^*$  es global y única. Cuando la Hessiana es indefinida puede existir otra solución local que la solución global.

## Bibliography

- [1] Murphy, K. P. (2012). Machine learning: a probabilistic perspective. MIT press.
- [2] Bishop, C. M. (2006). Pattern Recognition. Machine Learning.
- [3] Hastie, T., R. Tibshirani, and J. Friedman (2001). The Elements of Statistical Learning. Springer.
- [4] Schoelkopf, B. and A. Smola (2002). Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press.
- [5] Koller, D. and N. Friedman (2009). Probabilistic Graphical Models: Principles and Techniques. MIT Press.
- [6] Aronszajn, N. (1950). Theory of reproducing kernels. Transactions of the American mathematical society, 68(3), 337-404.
- [7] Schoelkopf, B., A. Smola, and K.-R. Mueller (1998). Nonlinear component analysis as a kernel eigenvalue problem. Neural Computation 10, 1299 – 1319.
- [8] Jaakkola, T. and D. Haussler (1998). Exploiting generative models in discriminative classifiers. In NIPS, pp.487–493
- [9] Krishnapuram, B., L. Carin, M. Figueiredo, and A. Hartemink (2005). Learning sparse bayesian classifiers: multi-class formulation, fast algorithms, and generalization bounds. IEEE Transaction on Pattern Analysis and Machine Intelligence.
- [10] Tipping, M. (2001). Sparse bayesian learning and the relevance vector machine. J. of Machine Learning Research 1, 211–244.
- [11] Hastie, T., R. Tibshirani, and J. Friedman (2009). The Elements of Statistical Learning. Springer. 2nd edition
- [12] Hofmann, T., Scholkopf, B., Smola, A. J. (2008). Kernel methods in machine learning. The annals of statistics, 1171-1220
- [13] Suykens, J. A., Van Gestel, T., De Brabanter, J. (2002). “Least squares support vector machines”. World Scientific.
- [14] Cortes, C., Vapnik, V. (1995). Support-vector networks. Machine learning, 20(3), 273-297.
- [15] Steinwart, I., Christmann, A. (2008). Support vector machines. Springer Science and Business Media.
- [16] Gelman, A., and D. B. Rubin. 1992. Inference from Iterative Simulation Using Multiple Sequences. Statistical Science 7: 457–511.
- [17] Brooks, S. P., and A. Gelman. 1997. General Methods for Monitoring Convergence of Iterative Simulations. Journal of Computational and Graphical Statistics 7: 434–455.
- [18] Duane, S., A. Kennedy, B. Pendleton, and D. Roweth (1987). Hybrid Monte Carlo. Physics Letters B 195(2), 216–222
- [19] Jebara, T., R. Kondor, and A. Howard (2004). Probability product kernels. J. of Machine Learning Research 5, 819–844
- [20] David M. Blei, Alp Kucukelbir and Jon D. McAuliffe (2017) Variational Inference: A Review for Statisticians, Journal of the American Statistical Association, 112:518, 859-877