

# Modelos estadísticos con métodos variacionales

Doble Grado en Ingeniería Informática y Matemáticas

---

Luis Antonio Ortega Andrés

15 de septiembre de 2020

Trabajo Fin de Grado

*E.T.S. de Ingenierías Informática y de Telecomunicación  
Facultad de Ciencias*



**UNIVERSIDAD  
DE GRANADA**

## Inferencia variacional

- Algoritmo EM

- Algoritmo CAVI

- Familia exponencial

## Redes Bayesianas

- Definición

- Algoritmo paso de mensajes

## Modelos

- Mixtura de Gaussianas

- Asignación latente de Dirichlet

- Reducción de dimensionalidad

## Casos prácticos

- InferPy

- BayesPy

- Scikit-Learn

Los métodos variacionales permiten la resolución de problemas de inferencia con variables ocultas, tanto con enfoque Bayesiano como basado en verosimilitud. Han aumentado notablemente los modelos estadísticos en los que se puede hacer inferencia.

Las variables ocultas pueden ser parámetros bajo un modelo Bayesiano.

Modelos concretos como la mixtura de Gaussianas, asignación latente de Dirichlet y análisis de componentes principales pueden ser estudiados mediante la utilización de lenguajes de programación especializados.

## Divergencia Kullback-Leibler

Sean  $P$  y  $Q$  dos distribuciones de probabilidad sobre el mismo espacio probabilístico, su *divergencia de Kullback-Leibler*  $KL(Q \mid P)$  mide la “diferencia” de  $Q$  a  $P$

$$KL(Q \mid P) = \mathbb{E}_Q \left[ \log Q(x) - \log P(x) \right].$$

La divergencia de Kullback-Leibler es siempre no negativa.

## El problema

- Un conjunto de variables observadas i.i.d  $\mathbf{X} = (X_1, \dots, X_N)$ .
- Variables ocultas globales  $\theta$  y variables ocultas locales  $\mathbf{Z} = (Z_1, \dots, Z_N)$ .
- La distribución conjunta factoriza como

$$P(\mathbf{x}, \mathbf{z}, \theta) = P(\theta) \prod_{n=1}^N P(x_n, z_n \mid \theta).$$

Las variables ocultas pueden ser tratadas como variables o como parámetros.

**Objetivo:** calcular

$$P(\theta, \mathbf{z} \mid \mathbf{x}) = \frac{P(\theta, \mathbf{z}, \mathbf{x})}{\int_{\theta, \mathbf{z}} P(\theta, \mathbf{z}, \mathbf{x})}.$$

# Inferencia variacional

---

Resuelve el problema de inferencia mediante uno de *optimización*:

Dada una familia de distribuciones  $\mathcal{Q}$  sobre el conjunto de variables ocultas  $\mathbf{Z}$ , encontrar

$$Q^{opt} = \arg \min_{Q \in \mathcal{Q}} KL\left(Q(\mathbf{z}) \mid P(\mathbf{z} \mid \mathbf{x})\right).$$

El problema se afronta mediante técnicas de aprendizaje automático como *gradiente descendente* o *descenso coordinado*.

## Cota inferior para la evidencia

Utilizando la positividad de la divergencia de Kullback-Leibler, se obtiene una cota inferior para la evidencia:

$$\begin{aligned} KL(Q(z) \mid P(z \mid \mathbf{x})) &\geq 0 \\ \Updownarrow \\ \log P(\mathbf{x}) &\geq \underbrace{-\mathbb{E}_{Q(z)}[\log Q(z)]}_{\text{Entropía}} + \underbrace{\mathbb{E}_{Q(z)}[\log P(\mathbf{x}, z)]}_{\text{Energía}} = ELBO(Q). \end{aligned}$$



# Algoritmo Esperanza-Maximización

No considera distribución sobre los parámetros.

$$ELBO(Q, \theta) = \underbrace{-\mathbb{E}_{Q(z)} \left[ \log Q(z) \right]}_{\text{Entropía}} + \underbrace{\mathbb{E}_{Q(z)} \left[ \log P(\mathbf{x}, z \mid \theta) \right]}_{\text{Energía}}.$$

Dos pasos iterativos:

- **Paso E:** Fijado el parámetro  $\theta$ ,

$$Q^{new}(z) = \arg \max_Q ELBO(Q, \theta) = P(z \mid \mathbf{x}, \theta).$$

- **Paso M:** Fijada la distribución  $Q$ ,

$$\theta^{new} = \arg \max_{\theta} ELBO(Q, \theta) = \arg \max_{\theta} \mathbb{E}_{Q(z)} \left[ \log P(\mathbf{x}, z \mid \theta) \right].$$

El algoritmo EM incrementa la verosimilitud.

$$\log P(\mathbf{x} \mid \boldsymbol{\theta}) \geq -\mathbb{E}_{Q(\mathbf{z})} \left[ \log Q(\mathbf{z}) \right] + \mathbb{E}_{Q(\mathbf{z})} \left[ \log P(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) \right].$$

⇓ Paso E

$$\log P(\mathbf{x} \mid \boldsymbol{\theta}) = -\mathbb{E}_{P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta})} \left[ \log P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}) \right] + \mathbb{E}_{P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta})} \left[ \log P(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}) \right].$$

⇓ Paso M

$$\log P(\mathbf{x} \mid \boldsymbol{\theta}^{new}) \geq -\mathbb{E}_{P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta})} \left[ \log P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}) \right] + \mathbb{E}_{P(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta})} \left[ \log P(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}^{new}) \right].$$

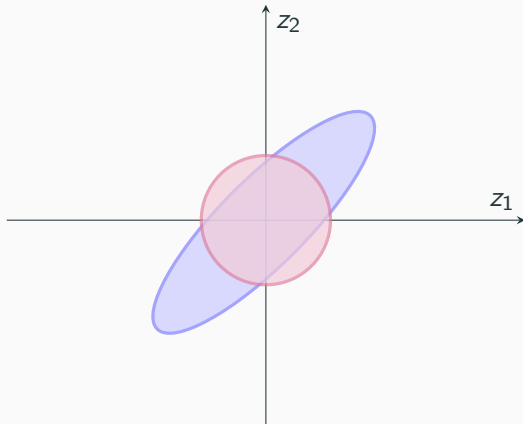
Converge a un máximo local.

# Algoritmo de ascenso coordinado en inferencia variacional

La familia de distribuciones de *campo medio* factoriza como producto de marginales:

$$Q(\mathbf{z}) = \prod_{m=1}^M Q(z_m).$$

- Pueden capturar cualquier distribución marginal.
- No pueden capturar correlación entre variables.



El algoritmo de ascenso coordinado se basa en:

- Los parámetros se consideran variables ocultas  $\theta \in \mathbf{Z}$ .
- La familia de distribuciones  $\mathcal{Q}$  del problema de optimización es la familia de campo medio.
- La distribución marginal de cada variable se actualiza cada vez.

$$Q^{new}(z_m) = \arg \max_{Q \in \mathcal{Q}} ELBO(Q)$$

$\Downarrow$

$$Q^{new}(z_m) \propto \exp \mathbb{E}_{Q_{\setminus m}} \left[ \log P(z_m \mid \mathbf{z}_{\setminus m}, \mathbf{x}) \right] \propto \exp \mathbb{E}_{Q_{\setminus m}} \left[ \log P(\mathbf{z}, \mathbf{x}) \right]$$

## Familia exponencial

Una variable aleatoria  $X$  sigue una distribución en la *familia exponencial* con parámetros  $\theta$  si y solo si existen  $h$ ,  $\mathbf{T}$ ,  $\eta$  y  $\psi$  tales que

$$P(x | \theta) = h(x) \exp \left( \eta(\theta)^T \mathbf{T}(x) - \psi(\theta) \right).$$

- $\mathbf{T}$  es el estadístico natural de  $X$ . *Teorema de Fisher-Neyman*.
- $\eta$  se denomina la función paramétrica de la distribución.
- $h$  se denomina medida base.
- $\psi$  asegura normalización logarítmica

$$\psi(\theta) = \log \int_x h(x) \exp \left( \eta(\theta)^T \mathbf{T}(x) \right).$$

## Modelos condicionalmente conjugados

Se consideran distribuciones conjugadas en la familia exponencial.

Las actualizaciones de las distribuciones variacionales  $Q^{new}(z_n)$  y  $Q^{new}(\theta)$  consisten en actualizar la *función paramétrica de la distribución*  $\eta$ .

Estas actualizaciones puede calcularse de forma eficiente, calculando esperanzas de los estadísticos suficientes.

# Redes Bayesianas

---

## Red Bayesiana

Una *red de creencia* o *red Bayesiana* es una pareja  $(G, P)$  formada por un grafo dirigido acíclico  $G$  y una distribución de probabilidad  $P$  tal que existe una correspondencia entre variables y nodos verificando:

$$P(x_1, \dots, x_N) = \prod_{n=1}^N P(x_n \mid pa(x_n)).$$



$$P(x_1, x_2, x_3) = P(x_1 \mid x_3)P(x_2 \mid x_3)P(x_3).$$



# Algoritmo de paso de mensajes

Proceso de paso de mensajes entre los nodos de la red.

- Familia variacional de campo medio.

$$Q(\mathbf{z}) = \prod_{m=1}^M Q(z_m).$$

- Modelo condicionalmente conjugado en la familia exponencial.
- Ascenso coordinado en redes Bayesianas.

La actualización de la distribución dada la red es:

$$\log Q^{new}(z) = \mathbb{E}_{Q_{\setminus z}} \left[ \sum_{n=1}^N \log P(x_n \mid pa(x_n)) \right] + \text{const.}$$

La aportación de  $Z$ :

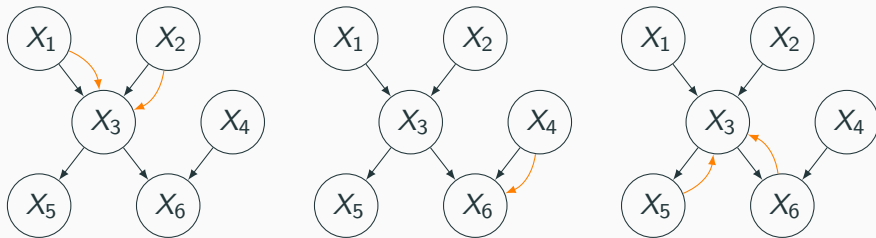
$$\log Q^{new}(z) = \mathbb{E}_{Q_{\setminus z}} \left[ \log P(z \mid pa(z)) \right] + \sum_{X \in ch(Z)} \mathbb{E}_{Q_{\setminus z}} \left[ \log P(x \mid z, cp(z, x)) \right] + \text{const.}$$

El mensaje de un nodo padre  $X$  a uno hijo  $Z$ :

$$\mathbf{m}_{X \rightarrow Z} = \mathbb{E}_Q \left[ \mathbf{T}_X(x) \right].$$

El mensaje de un nodo hijo  $X$  a un nodo padre  $Z$ :

$$\mathbf{m}_{X \rightarrow Z} = \bar{\eta}_{X,Z} \left( \mathbb{E}_Q \left[ \mathbf{T}_X(x) \right], \{ \mathbf{m}_{Y \rightarrow X} \}_{Y \in cp_{Z,X}} \right),$$



# Modelos

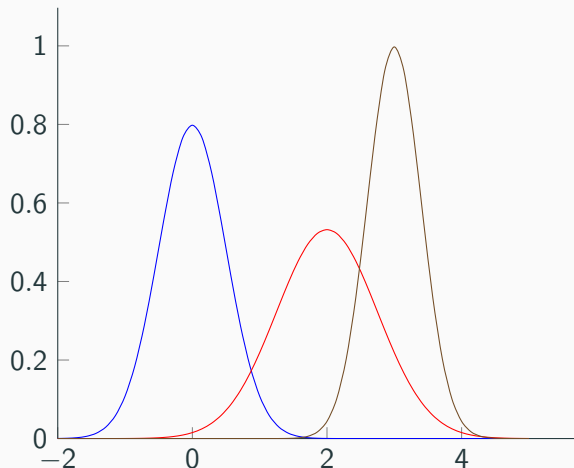
---

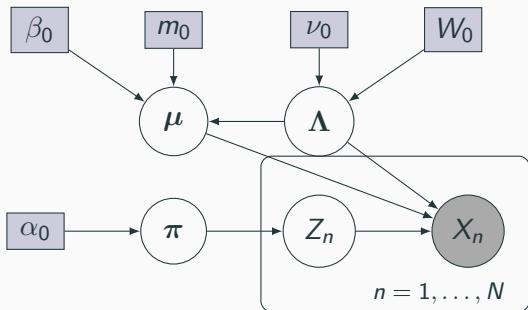
# Mixtura de Gaussianas

Mixtura de Gaussianas  $\mathcal{N}(0, 0.5)$ ,  
 $\mathcal{N}(2, 0.75)$  y  $\mathcal{N}(3, 0.4)$ .

Dados unos pesos 0.2, 0.3 y 0.5, la  
probabilidad de un punto sería

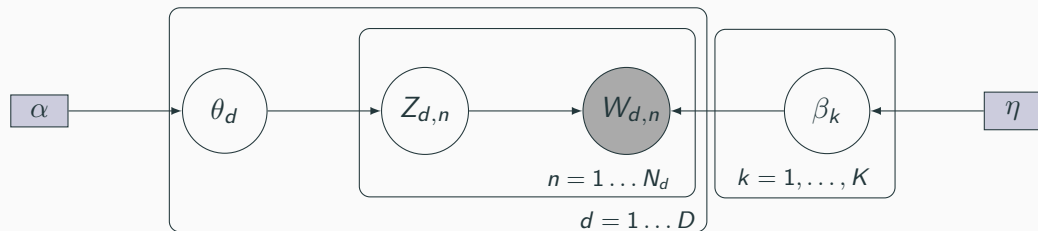
$$\begin{aligned} P(x) = & 0.2 \times \mathcal{N}(0, 0.5)(x) \\ & + 0.3 \times \mathcal{N}(2, 0.75)(x) \\ & + 0.5 \times \mathcal{N}(3, 0.4)(x). \end{aligned}$$





- **Pesos:**  $\pi \sim \text{Dirichlet}(\alpha_0)$ .
- **Medias y precisiones:**  
 $(\mu, \Lambda) \sim \text{Gaussian-Wishart}(\beta_0, m_0, \nu_0, W_0)$ .
- **Componentes:**  $Z_n \mid \pi \sim \text{Categorica}(\pi)$ .
- **Observaciones:**  
 $X_n \mid z_n, \mu, \Lambda \sim \mathcal{N}(\mu_{z_n}, \Lambda_{z_n})$ .

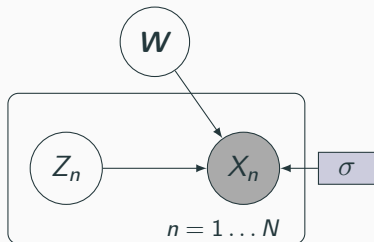
## Asignación latente de Dirichlet



- **Palabras por tema:**  $\beta_k \sim \text{Dirichlet}(\eta)$ .
- **Temas por documento:**  $\theta_d \sim \text{Dirichlet}(\alpha)$ .
- **Temas por documento y palabra:**  $Z_{d,n} \mid \theta_d \sim \text{Categorica}(\theta_d)$ .
- **Palabras:**  $W_{d,n} \mid z_{d,n}, \beta \sim \text{Categorica}(\beta_{z_{d,n}})$ .

# Reducción de dimensionalidad

Reducción de espacio  $D$  dimensional a uno  $K$  dimensional.



- **Transformación lineal:**  $W \sim \mathcal{N}_{D \times K}(0, I)$ .
- **Representación oculta:**  $Z_n \sim \mathcal{N}_K(0, I)$ .
- **Representación observada:**  
 $X_n \mid W, z_n \sim \mathcal{N}_d(Wz_n, \sigma I)$ .

**Modelo no lineal:** Red neuronal totalmente conectada de 2 capas.

$$X_n \mid z_n \sim \mathcal{N}_D(f(z_n), \sigma I).$$



## Modelo paramétrico

Representación oculta

$$Z_n \sim \mathcal{N}_K(0, I).$$

Representación observada

$$X_n \mid z_n \sim \mathcal{N}_D(f(z_n), \sigma I).$$

## Modelo variacional

Representación observada:

$$X_n \sim \mathcal{N}_D(0, I)$$

Representación oculta

$$Z_n \mid x_n \sim \mathcal{N}_K(\mu, \sigma I) \quad (\mu, \sigma) = g(x_n).$$

## Caso práctico

---

InferPy, BayesPy y BayesianGaussianMixture de Scikit-Learn.

Promovieron la utilización de métodos variacionales.

Especificación de modelos con alta abstracción.

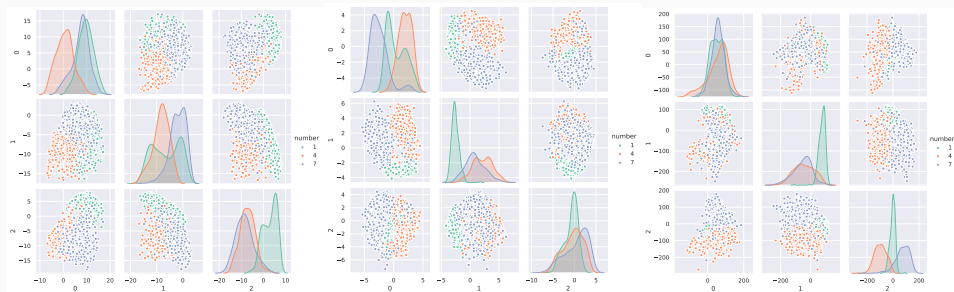
### PCA no lineal con InferPy:

```
nn = decoder(k, l, d)
with inf.datamodel():
    z = inf.Normal(loc=tf.zeros([k]), scale=1, name="z")
    output = nn(z)
    x = inf.Normal(loc=output, scale=1.0, name="x")
```

Integración con Keras.

No puede aprender variables categóricas.

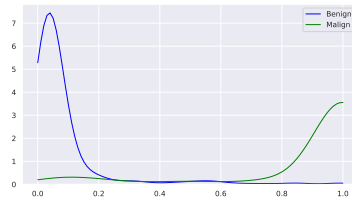
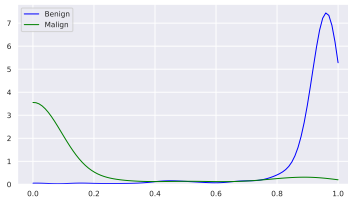
Algoritmo gradiente descendente.



Flexibilidad en diseño de modelos.

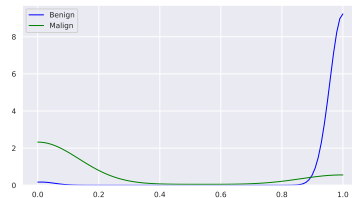
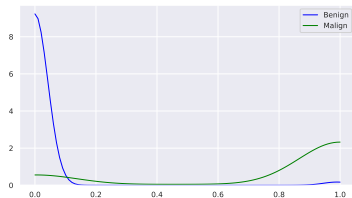
Altos requisitos de memoria.

Algoritmo de paso de mensajes variacional.



Modelo de mixtura pre-definido.

Algoritmo EM.



- La utilización de modelos gráficos como redes Bayesianas y modelos conjugados en la familia exponencial simplifican la inferencia con variables ocultas hasta su automatización.
- La integración de redes neuronales, permite aplicar inferencia a modelos fuera de los modelos conjugados.
- Cada software utilizado presenta ventajas e inconvenientes sobre los demás, haciendo que su elección dependa de la tarea que se desee realizar.

**Gracias por su atención**