



UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA

Aprendizaje Automático Probabilístico

Optimización

Alfons Juan

*DSIC*

Departamento de Sistemas  
Informáticos y Computación

# Índice

<b>8.0 Resumen</b>	<b>1</b>
<b>8.1 Introducción</b>	<b>8</b>
8.1.1 Optimización local vs global . . . . .	9
8.1.1.1 Condiciones de optimalidad local . . . . .	11
8.1.2 Optimización con o sin restricciones . . . . .	13
8.1.3 Optimización convexa vs no-convexa . . . . .	15
8.1.3.1 Conjuntos convexos . . . . .	15
8.1.3.2 Funciones convexas . . . . .	16
8.1.3.3 Caracterización de funciones convexas . . . . .	20
8.1.3.4 Funciones fuertemente convexas . . . . .	22
8.1.4 Optimización suave vs no-suave . . . . .	24
8.1.4.1 Subgradientes . . . . .	27
<b>8.2 Métodos de primer orden</b>	<b>29</b>
8.2.1 Dirección de descenso . . . . .	30
8.2.2 Tamaño de paso o factor de aprendizaje . . . . .	31
8.2.2.1 Tamaño de paso constante . . . . .	31
8.2.2.2 Búsqueda lineal . . . . .	35

8.2.3	Ratios de convergencia . . . . .	38
8.2.4	Momentum . . . . .	43
8.2.4.1	Momentum . . . . .	44
8.2.4.2	Momentum Nesterov . . . . .	46
<b>8.3</b>	<b>Métodos de segundo orden</b>	<b>49</b>
8.3.1	Método de Newton . . . . .	50
8.3.2	BFGS y otros métodos quasi-Newton . . . . .	55
8.3.3	Métodos en regiones de confianza . . . . .	57
<b>8.4</b>	<b>Descenso por gradiente estocástico</b>	<b>61</b>
8.4.1	Aplicación a problemas de sumas finitas . . . . .	62
8.4.2	Ejemplo: SGD para ajustar regresión lineal . . . . .	63
8.4.3	Elección del tamaño de paso . . . . .	65
8.4.4	Promediado iterativo . . . . .	72
8.4.6	SGD preconditionado . . . . .	73
8.4.6.1	ADAGRAD . . . . .	74
8.4.6.2	RMSPROP y ADADelta . . . . .	75
8.4.6.3	ADAM . . . . .	77
8.4.6.4	Problemas con los factores adaptativos . . . . .	79
8.4.6.5	Matrices de preconditionado no diagonales . . . . .	80

<b>8.5 Optimización con restricciones</b>	<b>81</b>
8.5.1 Multiplicadores de Lagrange . . . . .	83
8.5.1.1 Ejemplo 2D cuadrático con una restricción . . . . .	86
8.5.2 Las condiciones KKT . . . . .	87
8.5.3 Programación lineal . . . . .	90
8.5.3.1 El algoritmo simplex . . . . .	91
8.5.3.2 Aplicaciones . . . . .	91
8.5.4 Programación cuadrática . . . . .	92
8.5.4.1 Ejemplo: objetivo cuadrático 2d . . . . .	93
8.5.4.2 Aplicaciones . . . . .	96
<b>8.7 Optimización acotada</b>	<b>97</b>
8.7.1 El algoritmo general . . . . .	98
8.7.2 El algoritmo EM . . . . .	101
8.7.2.1 Cota inferior . . . . .	102
8.7.2.2 Paso E . . . . .	104
8.7.2.3 Paso M . . . . .	106
8.7.3 Ejemplo: EM para un GMM . . . . .	108
8.7.3.1 Paso E . . . . .	108
8.7.3.2 Paso M . . . . .	109

8.7.3.3	Ejemplo . . . . .	111
8.7.3.4	Estimación MAP . . . . .	112
8.7.3.5	Noconvexidad de la NLL . . . . .	117
<b>8.8</b>	<b>Optimización sin derivadas y caja-negra</b>	<b>119</b>

## 8.0. Resumen

- ▶ **Introducción:** optimización continua, no discreta!
  - ▷ **Dicotomía 1:** local o global?
    - Global para problemas convexos... ahora son no convexos
    - Local es lo que se hace... condiciones de optimalidad
  - ▷ **Dicotomía 2:** con o sin restricciones?
    - Mejor sin restricciones (p.e. con ayuda de la softmax)
    - Las de igualdad requieren multiplicadores de Lagrange; las de desigualdad, si son pocas, quizás las podemos ignorar
  - ▷ **Dicotomía 3:** convexa o no?
    - Si un problema es convexo, un óptimo local es global!
  - ▷ **Dicotomía 4:** suave o no?
    - Suave si objetivo y restricciones son continuamente diferenciables... constante de Lipschitz
    - No suave “por poco”... subgradiente

- ▶ **Métodos de primer orden:** basados en derivadas de primer orden del objetivo. . .  $\theta_{t+1} = \theta_t + \eta_t d_t$
- ▷ **Dirección de descenso:**  $d_t$  tal que  $\mathcal{L}(\theta + \eta d_t) < \mathcal{L}(\theta)$ 
  - ↳ **Descenso por gradiente:** negativo del gradiente,  $d_t = -g_t$
- ▷ **Tamaño de paso o factor de aprendizaje:**  $\{\eta_t\}$ ?
  - ↳ **Constante:**  $\eta_t = \eta$  . . . difícil de ajustar en la práctica
  - ↳ **Búsqueda lineal:** exacta o aproximada . . . **Armijo-Goldstein**
- ▷ **Ratios de convergencia:**  $\mu : |\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_*)| \leq \mu |\mathcal{L}(\theta_t) - \mathcal{L}(\theta_*)|$ 
  - ↳ **Objetivo cuadrático:**  $\mathbf{A} \succ 0$ ,  $\mu = \left( \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1} \right)^2$ ,  $\kappa(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}$
  - ↳ **No cuadrático:**  $\approx$  cuadrático cerca de local  $\rightarrow \kappa(\text{Hessiana})$
- ▷ **Momentum:** heurístico para acelerar la convergencia
  - ↳ **Estándar:**  $m_t = \beta m_{t-1} + g_{t-1}$  y  $\theta_t = \theta_{t-1} - \eta_t m_t$ ,  $\beta < 1$ 
    - **EWMA:**  $m_t = \sum_{\tau=0}^{t-1} \beta^\tau g_{t-\tau-1} \stackrel{g_{t-\tau-1}=g}{=} g \sum_{\tau=0}^{t-1} \beta^\tau \stackrel{t \rightarrow \infty}{=} \frac{g}{1-\beta}$
  - ↳ **Nesterov:** extrapola  $\theta_{t+1}$  para amortiguar oscilaciones

- ▶ **Métodos de segundo orden:** añaden la Hessiana de  $\mathcal{L}$  o aprox.
  - ▷ **Método de Newton:**  $\theta_{t+1} = \theta_t - \eta_t \mathbf{H}_t^{-1} \mathbf{g}_t$ 
    - ↳ Primero halla  $\mathbf{d}_t$  tal que  $\mathbf{H}_t \mathbf{d}_t = -\mathbf{g}_t$
    - ↳ Luego  $\theta_{t+1} = \theta_t + \eta_t \mathbf{d}_t$  con  $\eta_t$  hallado por búsqueda lineal
  - ▷ **Métodos quasi-Newton:** aproximan  $\mathbf{H}_t$  con  $\mathbf{B}_t$ , obtenida iterativamente a partir de los gradientes hallados en cada paso
    - ↳ **BFGS (Broyden–Fletcher–Goldfarb–Shanno):** aplica actualizaciones sucesivas de rango dos ... **Wolfe;** alt  $\mathbf{C}_t \approx \mathbf{H}^{-1}$
    - ↳ **Limited memory BFGS (L-BFGS):** aproxima  $\mathbf{H}_t^{-1} \mathbf{g}_t$  con las  $M$  actualizaciones más recientes
  - ▷ **Métodos en regiones de confianza:** no fijan  $\mathbf{d}_t$  y luego  $\eta_t$ , sino al revés; aproximan  $\mathcal{L}$  alrededor de  $\theta_t$  y buscan una dirección óptima... **regularización de Tikhonov**



- **Descenso por gradiente estocástico:** descenso por gradiente aplicado a **optimización estocástica**,  $\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{q(z)}[\mathcal{L}(\boldsymbol{\theta}, z)]$

$$\boldsymbol{\theta}_{t+1} \stackrel{z_t \sim q}{=} \boldsymbol{\theta}_t - \eta_t \nabla \mathcal{L}(\boldsymbol{\theta}_t, z_t) \stackrel{q(z) \text{ indep } \boldsymbol{\theta}}{=} \boldsymbol{\theta}_t - \eta_t \mathbf{g}_t$$

- ▷ **Aplicación a problemas de sumas finitas:** en minimización del riesgo empírico con  $N$  muestras, observamos un **minibatch** de  $B \ll N$  muestras en la iteración  $t$
- ▷ **Ejemplo: SGD para ajustar regresión lineal:** se conoce como **mínimos cuadrados**, **regla delta** o **Widrow-Hoff**
- ▷ **Elección del tamaño de paso: Robbins-Monro**
- ▷ **Promediado iterativo:** EWMA para reducir la varianza
- ▷ **SGD preconditionado:**  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{M}_t^{-1} \mathbf{g}_t$  con **precondicionador** diagonal  $\mathbf{M}_t$  **ADAGRAD**, **RMSPROP**, **ADADelta** o **ADAM**

- ▶ **Optimización con restricciones:** de igualdad y desigualdad
  - ▷ **Multiplicadores de Lagrange:** solo restricciones de igualdad
  - ▷ **Las condiciones KKT:** para el caso general ... **estacionariedad, factibilidad primal y dual, y holgura complementaria**
  - ▷ **Programación lineal:** objetivo lineal con restricciones lineales
  - ▷ **Programación cuadrática:** objetivo cuadrático con restricciones lineales

- **Optimización acotada:** basada en una cota inferior del objetivo
- ▷ **Algoritmo majorize-minorize (MM):** basado en una **función sustituta**, cota inferior  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \leq \text{LL}(\boldsymbol{\theta})$  que toca el objetivo en  $\boldsymbol{\theta}^t$ ,  $Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = \text{LL}(\boldsymbol{\theta}^t)$ :  $\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$
- ▷ **Algoritmo expectation maximization (EM):** algoritmo MM para calcular el MLE o MAP de modelos con datos perdidos
 
$$\text{LL}(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{y}_n \mid \boldsymbol{\theta}) = \sum_{n=1}^N \log \left[ \sum_{\mathbf{z}_n} p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta}) \right]$$
  - ⇒ **Evidence lower bound (ELBO):**  $\mathbb{L}(\boldsymbol{\theta}, q_{1:N}) \leq \text{LL}(\boldsymbol{\theta})$  (Jensen)
  - ⇒ **Paso E:** cálculo de  $q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta}) \rightarrow \mathbb{L}(\boldsymbol{\theta}, q_n^*) = \log p(\mathbf{y}_n \mid \boldsymbol{\theta})$
  - ⇒ **Paso M:** maximización de la log-verosimilitud completa esperada,  $\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} \sum_n \mathbb{E}_{q_n^t(\mathbf{z}_n)} [\log p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})]$
  - ⇒ **Aplicación a mixturas de Gaussianas**
  - ⇒ **Estimación MAP:** para mixturas de Gaussianas (robusta)
  - ⇒ **Noconvexidad de la NLL:** label switching problem

- ▶ **Optimización sin derivadas y caja-negra:** optimización mediante búsqueda en rejilla para selección de modelos
  - ▷ Para exploración de hiperparámetros y cuesta “horrores”

## 8.2. Métodos de primer orden

- ▶ Los *métodos de primer orden* son métodos iterativos basados en derivadas de primer orden del objetivo

- ▶ Dado un punto de inicio  $\theta_0$ , la iteración  $t$  consiste en hacer:

$$\theta_{t+1} = \theta_t + \eta_t d_t \quad (22)$$

- ▶  $\eta_t$  es el *tamaño del paso (step size)* o *factor de aprendizaje (learning rate)*
- ▶  $d_t$  es la *dirección de descenso*, como el negativo del *gradiente*,  
 $g_t = \nabla_{\theta} \mathcal{L}(\theta)|_{\theta_t}$
- ▶ Se termina al alcanzar un punto estacionario, de gradiente nulo

## 8.2.1. Dirección de descenso

- $d$  es una ***dirección de descenso*** si existe un  $\eta_{\max} > 0$  tal que

$$\mathcal{L}(\boldsymbol{\theta} + \eta \mathbf{d}) < \mathcal{L}(\boldsymbol{\theta}) \quad \text{para todo } 0 < \eta < \eta_{\max} \quad (23)$$

- La dirección de máximo ascenso en  $f$  es la del gradiente actual:

$$\mathbf{g}_t \triangleq \nabla \mathcal{L}(\boldsymbol{\theta})|_{\boldsymbol{\theta}_t} = \mathcal{L}(\boldsymbol{\theta}_t) = \mathbf{g}(\boldsymbol{\theta}_t) \quad (24)$$

- $d$  es dirección de descenso si el ángulo  $\theta$  entre  $d$  y  $-\mathbf{g}_t$  es menor de 90 grados y satisface:

$$\mathbf{d}^t \mathbf{g}_t = \|\mathbf{d}\| \|\mathbf{g}_t\| \cos(\theta) < 0 \quad (25)$$

- ***Descenso por gradiente (gradient descent)*** o ***más pronunciado (steepest descent)***: escogemos el negativo del gradiente

$$\mathbf{d}_t = -\mathbf{g}_t \quad (26)$$

## 8.2.2. Tamaño de paso o factor de aprendizaje

- *Learning rate schedule*: secuencia de tamaños de paso  $\{\eta_t\}$

### 8.2.2.1. Tamaño de paso constante

- La opción más simple consiste en usar un learning rate constante

$$\eta_t = \eta \quad (27)$$

- Si  $\eta$  es demasiado grande, el método puede no converger
- Si  $\eta$  es demasiado pequeño, el método convergerá muy lentamente
- *Ejemplo*:  $\mathcal{L}(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$ 
  - ▷ Con  $\eta = 0.1$  converge lentamente
  - ▷ Con  $\eta = 0.6$  oscila y no converge

## 8.2.2.2. Búsqueda lineal

- **Búsqueda lineal** consiste en hallar el tamaño de paso óptimo en la dirección escogida mediante optimización:

$$\eta_t = \arg \min_{\eta > 0} \phi_t(\eta) \quad (31)$$

$$= \arg \min_{\eta > 0} \mathcal{L}(\boldsymbol{\theta}_t + \eta \mathbf{d}_t) \quad (32)$$

- **Búsqueda lineal exacta** consiste en resolver analíticamente la optimización anterior, si se puede
  - ▷ En particular, si  $\mathcal{L}$  es convexa,  $\phi$  también lo es y sí se puede



## 8.2.3. Ratios de convergencia

- ▶ Queremos algoritmos que converjan rápidamente a un óptimo
- ▶ Descenso por gradiente converge con *ratio lineal* en problemas convexos con gradiente acotado por una constante de Lipschitz
- ▶ *Ratio de convergencia:*  $\mu \in (0, 1)$  tal que

$$|\mathcal{L}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}(\boldsymbol{\theta}_*)| \leq \mu |\mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}_*)| \quad (39)$$

- ▶ El ratio puede derivarse explícitamente en algunos problemas

## ► *Ratio de convergencia de un objetivo cuadrático:*

▷ Consideremos un objetivo cuadrático

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2}\boldsymbol{\theta}^t \mathbf{A} \boldsymbol{\theta} + \mathbf{b}^t \boldsymbol{\theta} + c \quad \text{con} \quad \mathbf{A} \succ 0 \quad (40)$$

▷ Aplicamos descenso por gradiente con búsqueda lineal exacta

▷ Se puede ver que el ratio de convergencia es

$$\mu = \left( \frac{\lambda_{\max}(\mathbf{A}) - \lambda_{\min}(\mathbf{A})}{\lambda_{\max}(\mathbf{A}) + \lambda_{\min}(\mathbf{A})} \right)^2 = \left( \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1} \right)^2 \quad (41)$$

donde el número de condición de  $\mathbf{A}$ ,

$$\kappa(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})} \quad (42)$$

mide la curvatura del objetivo (respecto a un bol simétrico)

## 8.7. Optimización acotada

- ▶ *Optimización acotada o MM (majorize-minimize)*: clase de algoritmos de optimización de gran interés en ML
- ▷ *Algoritmo expectation-maximization (EM)*: caso especial de algoritmo MM muy usado en ML

## 8.7.1. El algoritmo general

- **Objetivo:** maximizar  $LL(\boldsymbol{\theta})$
- **Función sustituta (surrogate):** construimos una función cota inferior del objetivo,  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$ , que lo iguala en un  $\boldsymbol{\theta}^t$  dado:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \leq LL(\boldsymbol{\theta}) \quad y \quad Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = LL(\boldsymbol{\theta}^t) \quad (137)$$

- Si se cumplen ambas condiciones, decimos que **minoriza**  $LL$

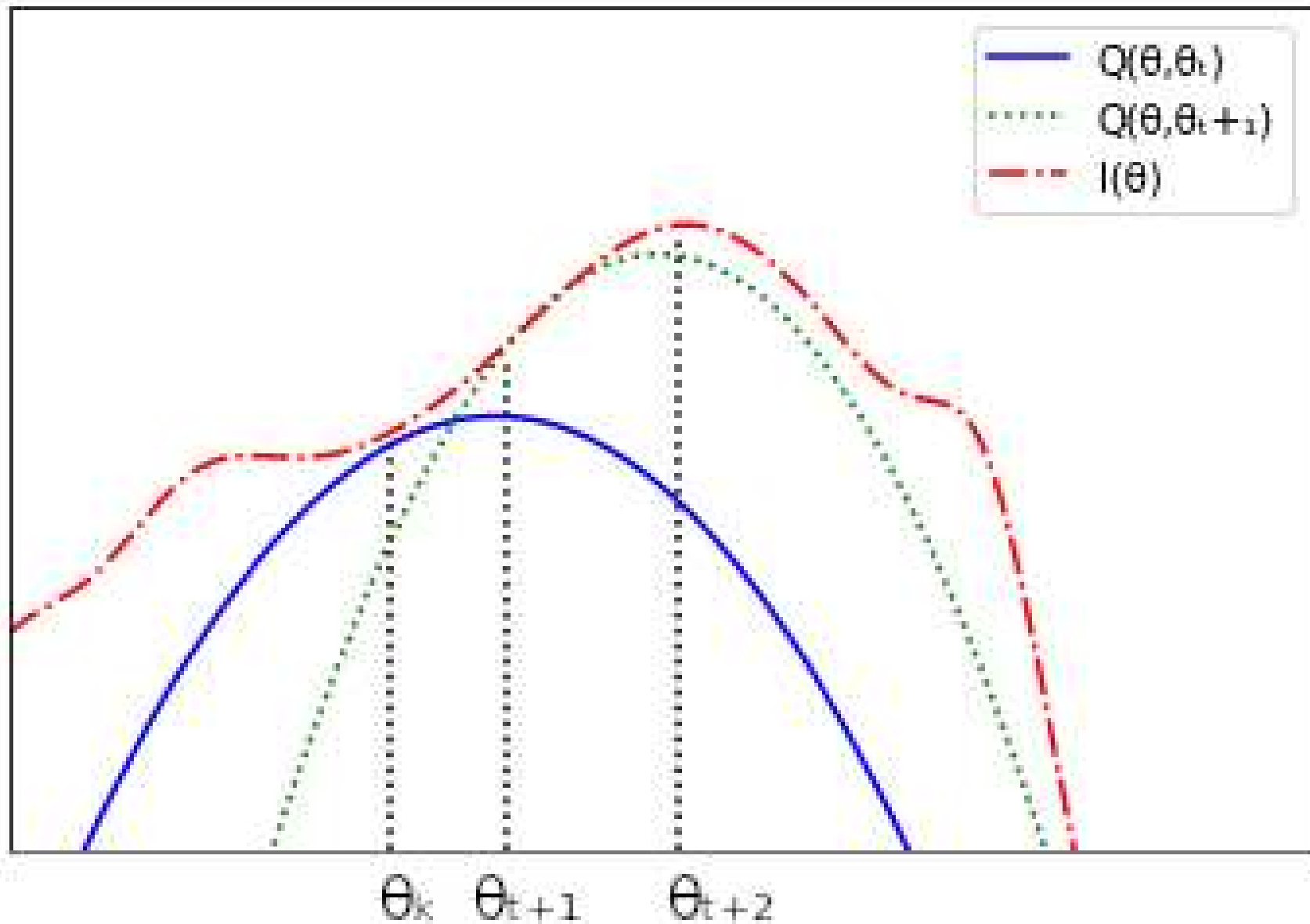
- **Algoritmo majorize-minimize (MM):** para  $t = 0, 1, \dots$

$$\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \quad (138)$$

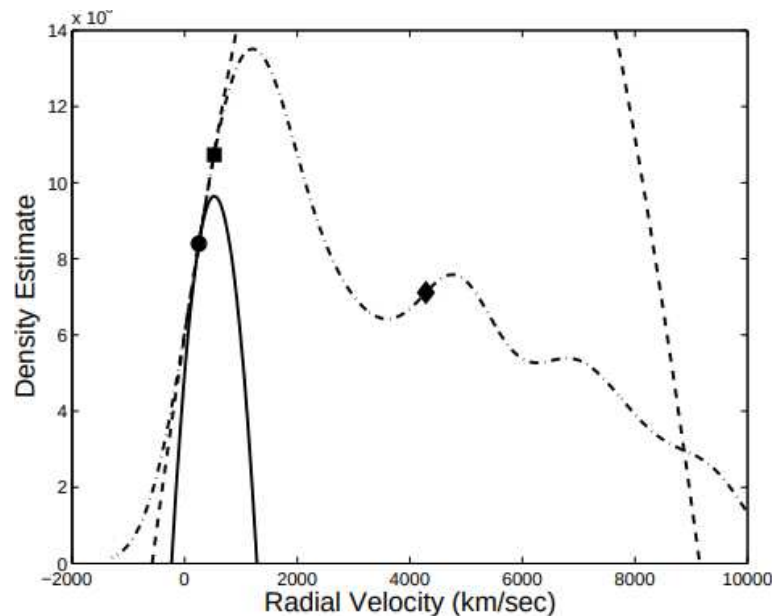
- Si  $\boldsymbol{\theta}^{t+1}$  se escoge tal que  $Q(\boldsymbol{\theta}^{t+1}, \boldsymbol{\theta}^t) \geq Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t)$ :

$$LL(\boldsymbol{\theta}^{t+1}) \geq Q(\boldsymbol{\theta}^{t+1}, \boldsymbol{\theta}^t) \geq Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = LL(\boldsymbol{\theta}^t) \quad (139)$$

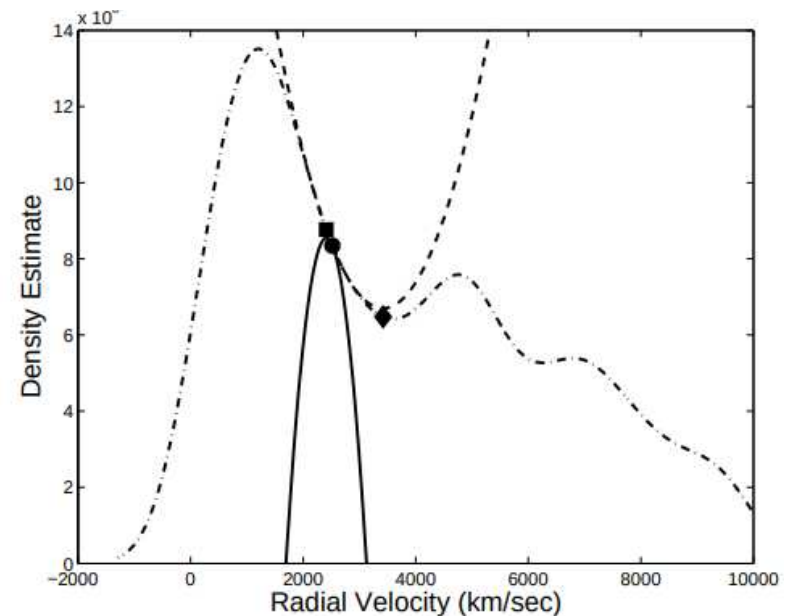
- *Ejemplo:*  $Q(\theta, \theta^t)$  toca  $LL(\theta)$  en  $\theta^t$ ; su maximización da lugar a  $\theta^{t+1}$  y  $Q(\theta, \theta^{t+1})$  toca  $LL(\theta)$  en  $\theta^{t+1}$ ; su maximización da lugar a  $\theta^{t+2}$ , etc.



- **Similitud con el método de Newton:** si  $Q$  es una cota inferior cuadrática, el MM se asemeja al método de Newton, pues ajusta y optimiza una aproximación cuadrática del objetivo repetidamente
- ▷ **Diferencia:** MM garantiza una mejora del objetivo en cada iteración, incluso si no es convexo, pero Newton no
- ▷ **Ejemplo:** a la izquierda Newton se “pasa de largo” buscando un máximo y a la derecha se va a un mínimo



(a) Overshooting.



(b) Seeking the wrong root.

## 8.7.2. El algoritmo EM

- ▶ **Algoritmo expectation maximization (EM):** algoritmo de optimización acotada para calcular el estimador MLE o MAP de modelos probabilísticos con **datos perdidos** o **variables ocultas**
- ▷ **Notación:** para cada dato  $n$ ,  $y_n$  denota su parte observada y  $z_n$  su parte perdida u oculta
- ▷ **Algoritmo EM básico:** repetir los siguientes dos pasos
  - ↳ **Paso E (expectation):** estimación de datos perdidos
  - ↳ **Paso M (maximization):** cálculo del MLE o MAP a partir de los datos completos
- ▷ **Convergencia:** veremos que el paso E calcula una función sustituta del objetivo, por lo que el EM converge a un máximo local

### 8.7.2.1. Cota inferior

► **Objetivo:** maximizar la log-verosimilitud de los datos observados

$$LL(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{y}_n \mid \boldsymbol{\theta}) \quad (140)$$

$$= \sum_{n=1}^N \log \left[ \sum_{\mathbf{z}_n} p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta}) \right] \quad (141)$$

► **Dificultad:** difícil de optimizar a causa del logaritmo delante del sumatorio



- **Evidence lower bound (ELBO):** dado un conjunto de distribuciones arbitrarias sobre cada  $\mathbf{z}_n$ ,  $q_n(\mathbf{z}_n)$ , la **desigualdad de Jensen** (sección 6.2.4) permite construir una función  $\mathbb{L}(\boldsymbol{\theta}, q_{1:N})$  cota inferior de la log-verosimilitud marginal o evidencia:

$$\text{LL}(\boldsymbol{\theta}) = \log p(\mathbf{y}_{1:N} \mid \boldsymbol{\theta}) \quad (142)$$

$$= \sum_{n=1}^N \log \left[ \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \frac{p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \right] \quad (143)$$

$$\geq \sum_{n=1}^N \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (144)$$

$$= \sum_n \underbrace{\mathbb{E}_{q_n}[\log p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})] + \mathbb{H}(q_n)}_{\mathbb{L}(\boldsymbol{\theta}, q_n \mid \mathbf{y}_n)} \quad (145)$$

$$= \sum_n \mathbb{L}(\boldsymbol{\theta}, q_n) \triangleq \mathbb{L}(\boldsymbol{\theta}, \{q_n\}) = \mathbb{L}(\boldsymbol{\theta}, q_{1:N}) \quad (146)$$

## 8.7.2.2. Paso E

► **Paso E:** cálculo de  $q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})$  y, así,  $\mathbb{L}(\boldsymbol{\theta}, q_n^*) = \log p(\mathbf{y}_n \mid \boldsymbol{\theta})$

$$\mathbb{L}(\boldsymbol{\theta}, q_n) = \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (147)$$

$$= \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta}) p(\mathbf{y}_n \mid \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} \quad (148)$$

$$= \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log \frac{p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})}{q_n(\mathbf{z}_n)} + \sum_{\mathbf{z}_n} q_n(\mathbf{z}_n) \log p(\mathbf{y}_n \mid \boldsymbol{\theta}) \quad (149)$$

$$= -\mathbb{KL}(q_n(\mathbf{z}_n) \parallel p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})) + \log p(\mathbf{y}_n \mid \boldsymbol{\theta}) \quad (150)$$

$$\stackrel{q_n=q_n^*}{=} \log p(\mathbf{y}_n \mid \boldsymbol{\theta}) \quad (151)$$

pues  $\mathbb{KL}(q_n(\mathbf{z}_n) \parallel p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})) = 0$  sii  $q_n \triangleq q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})$

► **Función sustituta:** dado que  $\mathbb{L}(\boldsymbol{\theta}, \{q_n^*\}) = \text{LL}(\boldsymbol{\theta})$ , la función

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \mathbb{L}(\boldsymbol{\theta}, \{q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta}^t)\}) \quad (152)$$

es ELBO por Jensen,

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \leq \text{LL}(\boldsymbol{\theta}) \quad (153)$$

y toca  $\text{LL}(\boldsymbol{\theta})$  en  $\boldsymbol{\theta}^t$  al tomar  $\{q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta}^t)\}$ ,

$$Q(\boldsymbol{\theta}^t, \boldsymbol{\theta}^t) = \text{LL}(\boldsymbol{\theta}^t) \quad (154)$$

► **Paso E aproximado:** si el cálculo de  $q_n^* = p(\mathbf{z}_n \mid \mathbf{y}_n, \boldsymbol{\theta})$  es muy costoso, podemos emplear una aproximación a la misma y la  $Q$ , aunque menos ajustada, sigue siendo ELBO

▷ **Aproximación directa:** comprobamos que la LL no decrece; cosa quizás sencilla si solo consideramos distribuciones delta

▷ **EM variacional:** EM generalizado en marco Bayesiano

### 8.7.2.3. Paso M

- **Expected complete data log likelihood:** el paso M maximiza  $\mathbb{L}(\boldsymbol{\theta}, \{q_n^t\})$  con respecto a  $\boldsymbol{\theta}$ , donde las  $\{q_n^t\}$  son las distribuciones halladas en el paso E de la iteración  $t$ ; ahora bien, como los términos de entropía  $\mathbb{H}(q_n)$  no dependen de  $\boldsymbol{\theta}$ , podemos ignorarlos,

$$\text{LL}^t(\boldsymbol{\theta}) = \sum_n \mathbb{E}_{q_n^t(\mathbf{z}_n)} [\log p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})] \quad (155)$$

- **Caso familia exponencial:** si la probabilidad conjunta pertenece a la familia exponencial, no necesitamos  $\{q_n^t\}$ ; bastan estadísticos suficientes esperados,  $\mathbb{E}[\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)]$ ,

$$\text{LL}^t(\boldsymbol{\theta}) = \sum_n \mathbb{E}[\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)^t \boldsymbol{\theta} - A(\boldsymbol{\theta})] \quad (156)$$

$$= \sum_n (\mathbb{E}[\mathcal{T}(\mathbf{y}_n, \mathbf{z}_n)]^t - A(\boldsymbol{\theta})) \quad (157)$$

► **Paso M:** maximización de la log-verosimilitud completa esperada

$$\boldsymbol{\theta}^{t+1} = \arg \max_{\boldsymbol{\theta}} \sum_n \mathbb{E}_{q_n^t(\mathbf{z}_n)} [\log p(\mathbf{y}_n, \mathbf{z}_n \mid \boldsymbol{\theta})] \quad (158)$$

▷ **Caso familia exponencial:** se resuelve en forma cerrada

## 8.7.3. Ejemplo: EM para un GMM

### 8.7.3.1. Paso E

- **Responsability:** el paso E calcula la **responsabilidad** del clúster  $k$  en la generación del dato  $n$ , según la estimación actual de los parámetros  $\theta^{(t)}$ ,

$$r_{nk}^{(t)} = p^*(z_n = k \mid \mathbf{y}_n, \theta^{(t)}) \quad (159)$$

$$= \frac{\pi_k^{(t)} p(\mathbf{y}_n \mid \theta_k^{(t)})}{\sum_{k'} \pi_{k'}^{(t)} p(\mathbf{y}_n \mid \theta_{k'}^{(t)})} \quad (160)$$

### 8.7.3.2. Paso M

► **Log-verosimilitud completa esperada:** versión ponderada de la LL para la Gaussiana multivariada; sea  $z_{nk} = \mathbb{1}(z_n = k)$ ,

$$LL^t(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_n \log p(z_n \mid \boldsymbol{\pi}) + \log p(\mathbf{y}_n \mid z_n, \boldsymbol{\theta}) \right] \quad (161)$$

$$= \mathbb{E} \left[ \sum_n \log \left( \prod_k \pi_k^{z_{nk}} \right) + \log \left( \prod_k \mathcal{N}(\mathbf{y}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_{nk}} \right) \right] \quad (162)$$

$$= \sum_n \sum_k \mathbb{E}[z_{nk}] \log \pi_k + \sum_n \sum_k \mathbb{E}[z_{nk}] \log \mathcal{N}(\mathbf{y}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (163)$$

$$\begin{aligned} &= \sum_n \sum_k r_{nk}^{(t)} \log(\pi_k) \\ &\quad - \frac{1}{2} \sum_n \sum_k r_{nk}^{(t)} [\log |\boldsymbol{\Sigma}_k| + (\mathbf{y}_n - \boldsymbol{\mu}_k)^t \boldsymbol{\Sigma}_k^{-1} (\mathbf{y}_n - \boldsymbol{\mu}_k)] + \text{const} \end{aligned} \quad (164)$$

► **Solución cerrada:** sea  $r_k^{(t)} \triangleq \sum_n r_{nk}(t)$

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{1}{r_k^{(t)}} \sum_n r_{nk}(t) \mathbf{y}_n \quad (165)$$

$$\Sigma_k^{(t+1)} = \frac{1}{r_k^{(t)}} \sum_n r_{nk}(t) (\mathbf{y}_n - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{y}_n - \boldsymbol{\mu}_k^{(t+1)})^t \quad (166)$$

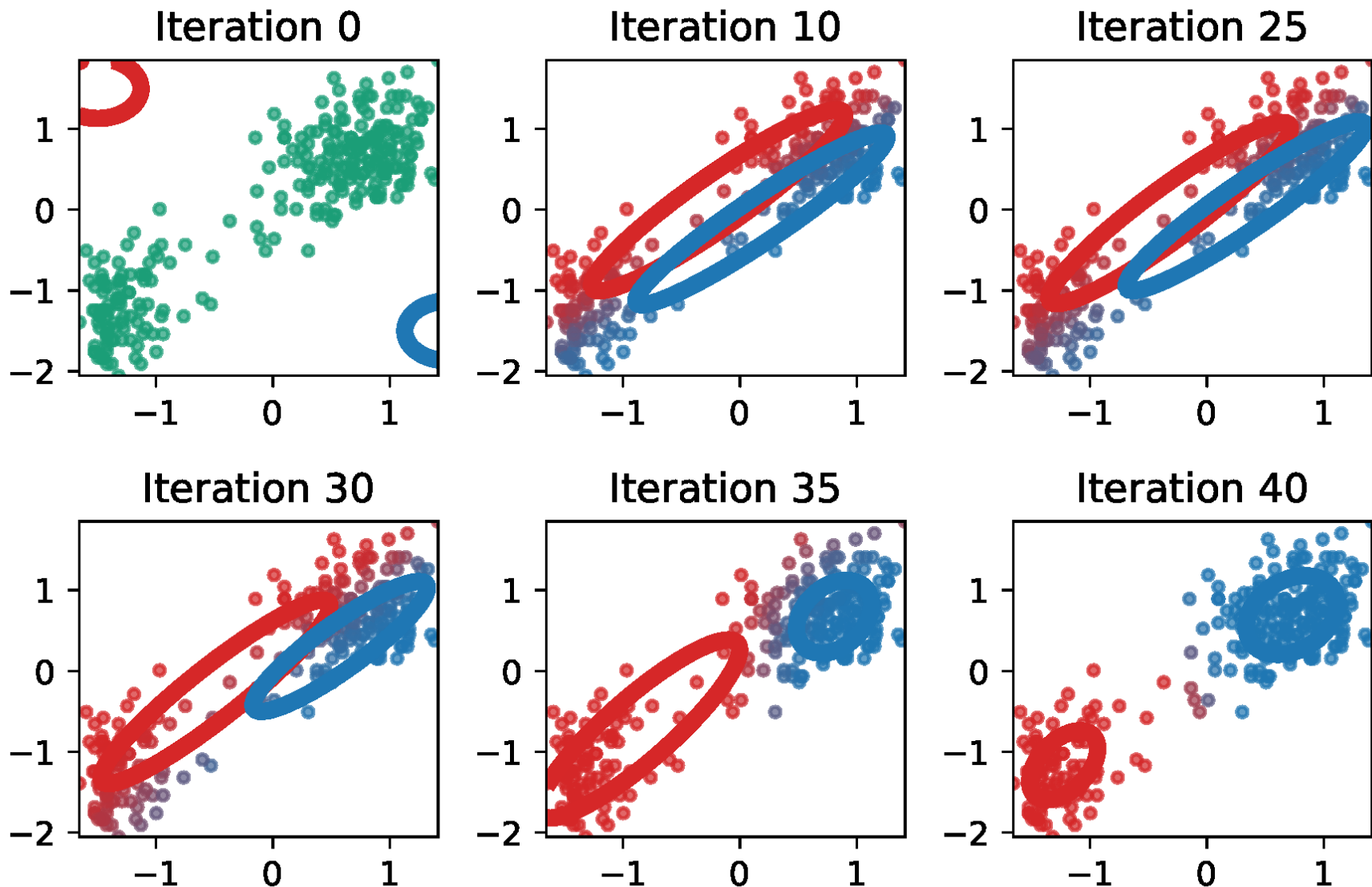
$$= \frac{1}{r_k^{(t)}} \left( \sum_n r_{nk}(t) \mathbf{y}_n \mathbf{y}_n^t \right) - \boldsymbol{\mu}_k^{(t+1)} (\boldsymbol{\mu}_k^{(t+1)})^t \quad (167)$$

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_n r_{nk}(t) = \frac{r_k^{(t)}}{N} \quad (168)$$



### 8.7.3.3. Ejemplo

- *Old Faithful*: GMM ajustado con el EM a datos 2d del géiser Old Faithful (minutos siguiente erupción vs duración; estandarizados)



### 8.7.3.4. Estimación MAP

- **Problema del colapso de la varianza:** si  $\Sigma_k = \sigma_k^2 \mathbf{I}$  y  $\mu_k$  se asigna a un único punto,  $\mathbf{y}_n$ , su verosimilitud diverge con  $\sigma_k \rightarrow 0$

$$\mathcal{N}(\mathbf{y}_n \mid \mu_k = \mathbf{y}_n, \sigma_k^2 \mathbf{I}) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^0 \quad (169)$$

- **Estimación MAP:** maximiza la log-verosimilitud completa esperada más un log-prior

$$\begin{aligned} \text{LL}^t(\boldsymbol{\theta}) = & \left[ \sum_n \sum_k r_{nk}^{(t)} \log \pi_k + \sum_n \sum_k r_{nk}^{(t)} \log p(\mathbf{y}_n \mid \boldsymbol{\theta}_k) \right] \\ & + \log p(\boldsymbol{\pi}) + \sum_k \log p(\boldsymbol{\theta}_k) \end{aligned} \quad (170)$$

## ► *Algoritmo EM:*

- ▷ *Paso E:* igual que para el MLE
- ▷ *Paso M para los coeficientes:* con prior *Dirichlet*,  $\pi \sim \text{Dir}(\alpha)$ , conjugada de la categórica,

$$\tilde{\pi}_k^{(t+1)} = \frac{r_k^{(t)} + \alpha_k - 1}{N + \sum_k \alpha_k - K} \quad (171)$$

⇒ Coincide con el MLE con un prior uniforme,  $\alpha_k = 1$

- ▷ **Paso M para las componentes:** con prior **Normal-Inverse-Wishart**, conjugada de la Gaussiana multivariada,

$$p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \text{NIW}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k \mid \widetilde{\boldsymbol{m}}, \widetilde{\kappa}, \widetilde{\nu}, \widetilde{\mathbf{S}}) \quad (172)$$

- ⇒ Con  $\widetilde{\kappa} = 0$ , las  $\boldsymbol{\mu}_k$  no se regularizan, por lo que el prior solo afecta a las  $\boldsymbol{\Sigma}_k$  y los estimadores MAP son:

$$\widetilde{\boldsymbol{\mu}}_k^{(t+1)} = \hat{\boldsymbol{\mu}}_k^{(t+1)} \quad (173)$$

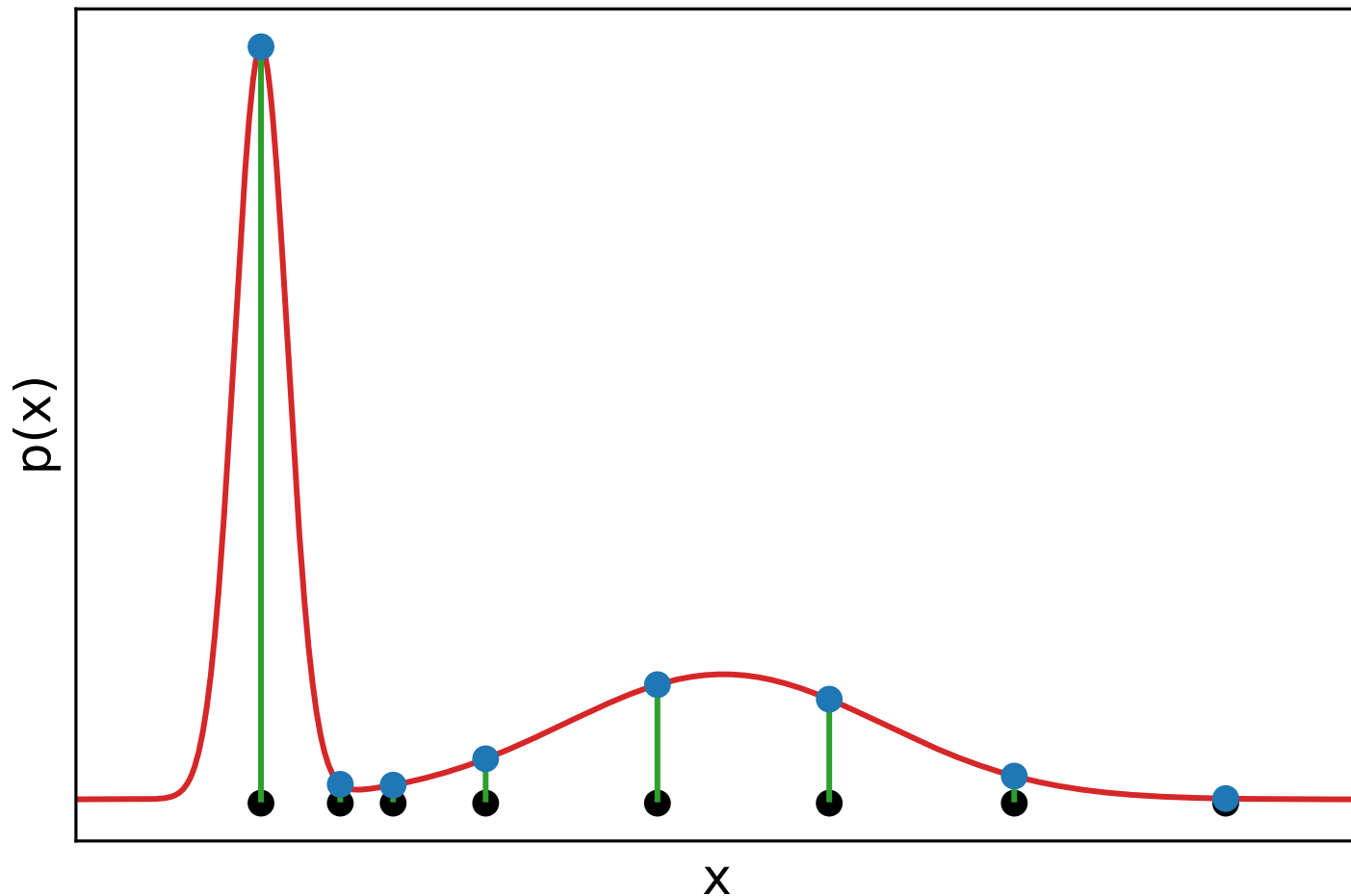
$$\widetilde{\boldsymbol{\Sigma}}_k^{(t+1)} = \frac{\widetilde{\mathbf{S}} + \hat{\boldsymbol{\Sigma}}_k^{(t+1)}}{\widetilde{\nu} + r_k^{(t)} + D + 2} \quad (174)$$

- ⇒ **Covarianza a priori:** si  $s_d = \frac{1}{N} \sum_n (x_{nd} - \bar{x}_d)^2$  es la varianza global en la dimensión  $d$ , una posibilidad consiste en usar

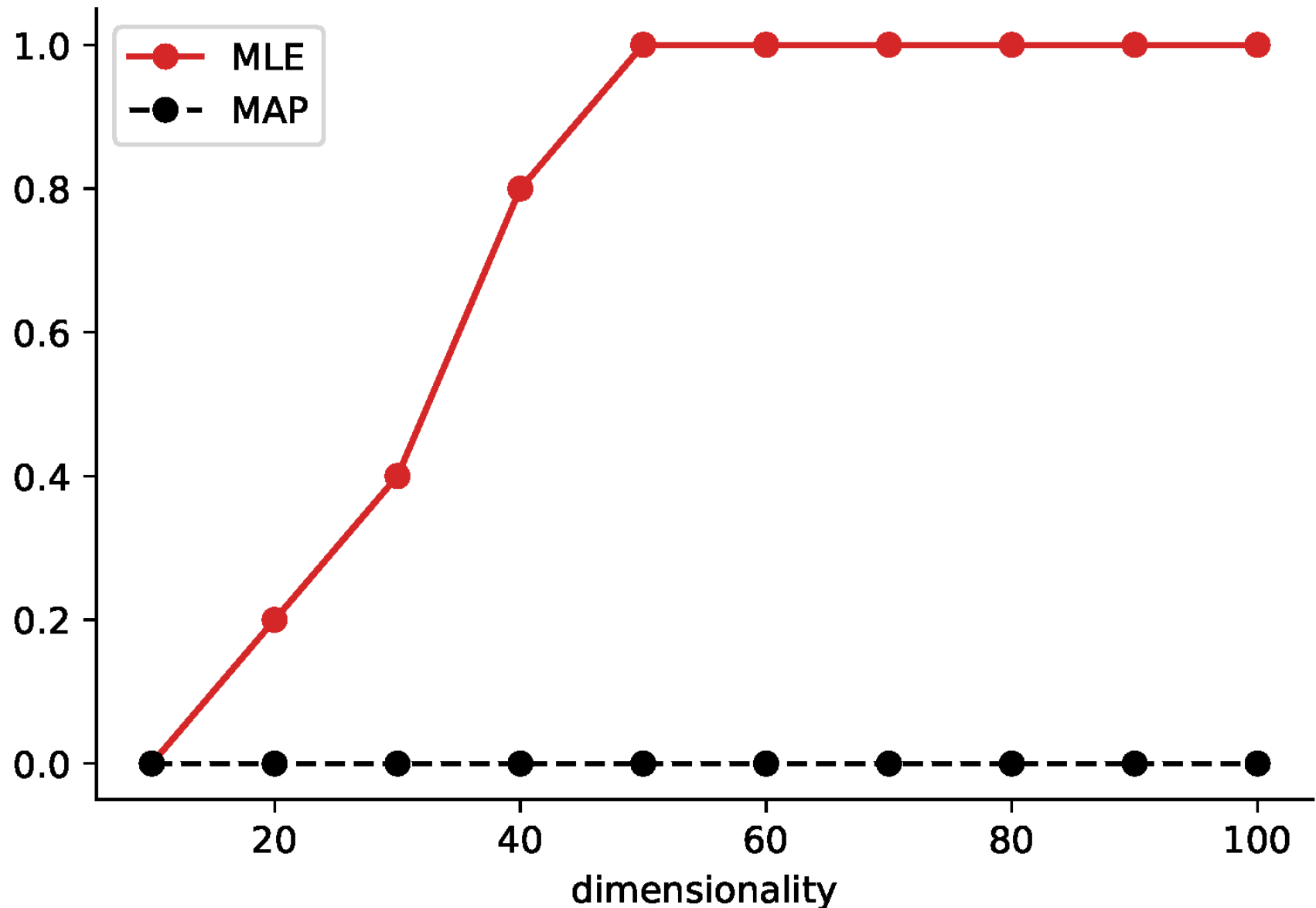
$$\widetilde{\mathbf{S}} = \frac{1}{K^{1/D}} \text{diag}(s_1^2, \dots, s_D^2) \quad (175)$$

- ⇒ El hiperparámetro  $\widetilde{\nu}$  controla la fuerza del prior; una elección usual es el prior propio más débil:  $\widetilde{\nu} = D + 1$

- ▷ *Ejemplo:* mezcla de  $K = 2$  componentes a ajustar con  $N = 100$  datos sintéticos en  $D$  dimensiones ( $D = 1$  en la gráfica)
- La primera componente es un pico estrecho (con  $\sigma_1 \approx 0$ ) centrado en un único dato  $x_1$



- ▷ *Ejemplo (cont.):* fracción del número de veces (de 5 intentos) que el EM presenta problemas numéricos con MLE y MAP

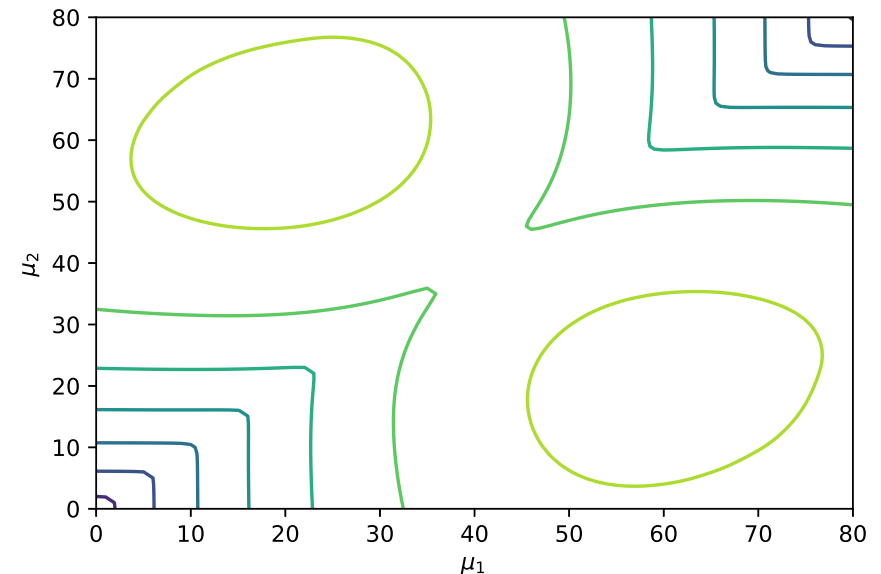
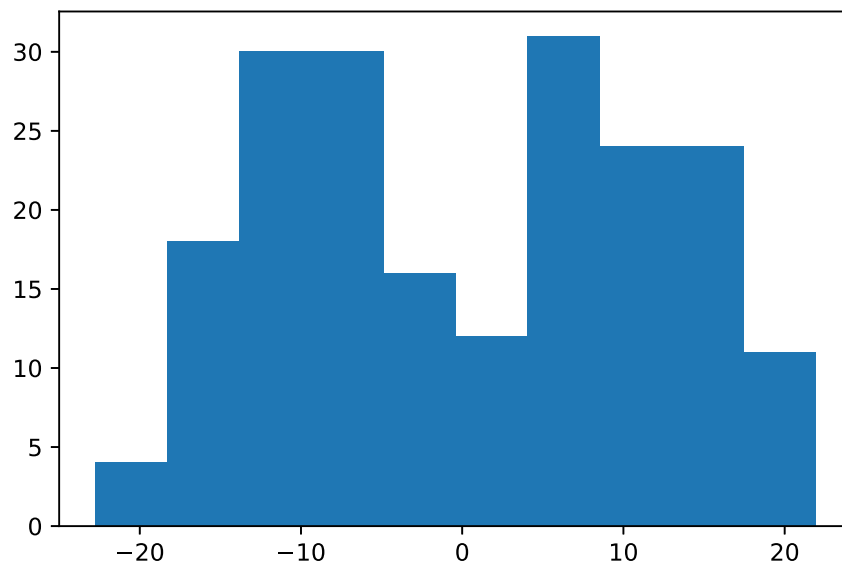


### 8.7.3.5. Noconvexidad de la NLL

- **Noconvexidad de la NLL:** la log-verosimilitud de una mixtura suele tener múltiples modas, esto es, más de un óptimo global

$$\text{LL}(\boldsymbol{\theta}) = \sum_{n=1}^N \log \sum_{z_n=1}^K p(\mathbf{y}_n, z_n \mid \boldsymbol{\theta}) \quad (176)$$

- **Ejemplo:** 200 puntos de una mixtura de 2 Gaussianas 1d con  $\pi_k = 0.5$ ,  $\sigma_k = 5$ ,  $\mu_1 = -10$  y  $\mu_2 = 10$ ; y verosimilitud  $p(\mathcal{D} \mid \mu_1, \mu_2)$



⇒ **Label switching problem:** 2 óptimos cambiando etiquetas

## ► *Complejidad del problema: label switching problem*

- ▷ Es difícil establecer el número de modas pues, aunque potencialmente hay  $K!$  etiquetados distintos, muchos picos pueden reducirse al mezclarse con otros cercanos
- ▷ *Número de modas exponencial:* en cualquier caso, puede haber un número de modas exponencial con  $K$ , por lo que el problema es NP-duro
- ▷ *Óptimo local:* únicamente podemos aspirar a encontrar un buen óptimo local, por lo general posible con una buena inicialización