

2016-2017

Aprendizaje Automático

3. Técnicas de optimización



Francisco Casacuberta Nolla
(fcm@dsic.upv.es)

Enrique Vidal Ruiz
(evidal@dsic.upv.es)

Departament de Sistemes Informàtics i Computació (DSIC)
Universitat Politècnica de València (UPV)

Septiembre, 2016

Aprendizaje Automático. 2016-2017

Técnicas de optimización: 3.1

Index

- 1 *Introducción* ▷ 1
- 2 Optimización analítica: gradiente ▷ 5
- 3 Optimización con restricciones: multiplicadores de Lagrange y teorema Kuhn-Tucker ▷ 11
- 4 Técnicas de descenso por gradiente ▷ 20
- 5 Esperanza-Maximización (EM) ▷ 31
- 6 Notación ▷ 45

Clasificación, regresión y optimización

- Los modelos están parametrizados por un vector de parámetros Θ ; es decir, $\mathcal{F} = \{f_{\Theta} : \mathcal{X} \rightarrow \mathcal{Y}, \Theta \in \mathbb{R}^D\}$
- Clasificación: $f_{\Theta} : \mathcal{X} \rightarrow \{1, \dots, C\}$. En muchos problemas $\mathcal{X} \equiv \mathbb{R}^d$
- Regresión: $f_{\Theta} : \mathcal{X} \rightarrow \mathcal{Y}$. Típicamente $\mathcal{X} \equiv \mathbb{R}^d$ y $\mathcal{Y} \equiv \mathbb{R}$.
- Dos etapas:
 - **Aprendizaje:** Dado $S \subset \mathcal{X} \times \mathcal{Y}$, estimar $\hat{\Theta} \in \mathbb{R}^D$
 Técnicas de optimización para aprendizaje:
 - * Optimización analítica.
 - * Optimización con restricciones: Multiplicadores de Lagrange.
 - * Descenso (ascenso) por gradiente.
 - * Optimización probabilística: Algoritmo EM.
 - **Búsqueda o inferencia:** Dados Θ y $x \in \mathcal{X}$, estimar $\hat{y} = f_{\Theta}(x)$.
 Técnicas de optimización para búsqueda:
 - * Exhaustiva: por ejemplo, clasificación, si $C \ll$.
 - * Programación dinámica: algoritmo de Viterbi con modelos ocultos de Markov.
 - * Inteligente: ramificación y poda, A^* , etc.

Optimización y aprendizaje automático

- Datos:
 - N muestras de aprendizaje:
 $S = \{(x_1, y_1), \dots, (x_N, y_N)\}, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, 1 \leq n \leq N,$
 - un clasificador o regresor, $f_{\Theta} : \mathcal{X} \rightarrow \mathcal{Y}$, parametrizado por $\Theta \in \mathbb{R}^D$,
 - un criterio aprendizaje definido por una función objetivo, $q_S : \mathbb{R}^D \rightarrow \mathbb{R}$
- estimar $\hat{\Theta}$ mediante optimización de q_S ; es decir:

$$\hat{\Theta} \equiv \Theta^* = \arg \min_{\Theta} q_S(\Theta)$$

o bien:

$$\hat{\Theta} \equiv \Theta^* = \arg \max_{\Theta} q_S(\Theta)$$

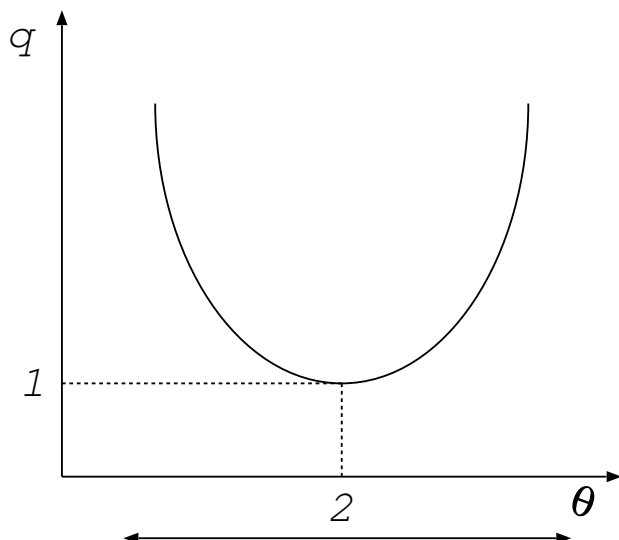
Técnicas generales de AA basadas en optimización

- f_{Θ} es una función cualquiera que depende de un vector de parámetros Θ : $q_S(\Theta)$ se basa en *funciones de error* y típicamente su optimización utiliza técnicas de *descenso/ascenso por gradiente* (caso particular de “hill-climbing”).
- f_{Θ} es una función cualquiera dependiente de Θ , pero hay ciertas restricciones en los valores posibles de los parámetros Θ : *optimización con restricciones* de $q_S(\Theta)$ mediante la técnica de los *multiplicadores de Lagrange*.
- f_{Θ} se basa en distribuciones (o densidades) de probabilidad: estimación de *máxima verosimilitud*. Frecuentemente hay restricciones en los parámetros a estimar y se requiere el uso de *multiplicadores de Lagrange*.
- f_{Θ} se basa en distribuciones (o densidades) de probabilidad, pero hay variables aleatorias “latentes” u “ocultas”: La estimación de *máxima verosimilitud* generalmente requiere una técnica de optimización llamada *esperanza-maximización* (EM).

Index

- 1 Introducción ▷ 1
- 2 Optimización analítica: gradiente ▷ 5
- 3 Optimización con restricciones: multiplicadores de Lagrange y teorema Kuhn-Tucker ▷ 11
- 4 Técnicas de descenso por gradiente ▷ 20
- 5 Esperanza-Maximización (EM) ▷ 31
- 6 Notación ▷ 45

Optimización analítica: ejemplo



- Dado $q(\theta) = 1 + (\theta - 2)^2$
- Calcular $\theta^* = \arg \min_{\theta \in \mathbb{R}} q(\theta)$
- Procedimiento: $\frac{dq(\theta)}{d\theta} = 2(\theta - 2) = 0$
- Solución: $\theta^* = 2$

Optimización analítica: gradiente

- Dada una función *convexa* $q: \mathbb{R}^D \rightarrow \mathbb{R}$, calcular $\arg \min_{\Theta \in \mathbb{R}^D} q(\Theta)$
- Procedimiento:
 1. Calcular el gradiente de q : $\nabla q(\Theta) \stackrel{\text{def}}{=} \left(\frac{\partial q(\Theta)}{\partial \Theta_1}, \dots, \frac{\partial q(\Theta)}{\partial \Theta_D} \right)^t$
 2. Resolver $\nabla q(\Theta) = 0$;
es decir, resolver el sistema de ecuaciones $\frac{\partial q(\Theta)}{\partial \Theta_i} = 0, 1 \leq i \leq D$.
Sean $\Theta_1^*, \dots, \Theta_D^*$ las soluciones obtenidas.
- Solución: $\Theta^* = (\Theta_1^*, \dots, \Theta_D^*)^t$
- Si q es convexa, $\nabla q(\Theta^*) = 0$ es una condición *necesaria y suficiente* para que Θ^* sea (la única) solución.

Ejercicios:

a) ¿Qué ocurre si q no es convexa?

b) Encontrar el vector $\theta^* \in \mathbb{R}^2$ que minimiza la función $q(\theta) = (\theta_1 - 1)^2 + (\theta_2 - 2)^2$

Optimización analítica: otro ejemplo simple

- Estimar los parámetros $\Theta \equiv (\mu, \sigma)^1$ de una gaussiana univariada (en \mathbb{R}^1):

$$p(x | \Theta) \stackrel{\text{def}}{=} p(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

- Logaritmo de la verosimilitud de una muestra $S = \{x_1, \dots, x_N\}$:

$$\begin{aligned} q_S(\Theta) \equiv L_S(\mu, \sigma) &= \log \prod_{n=1}^N p(x_n | \mu, \sigma) = \sum_{n=1}^N \log p(x_n | \mu, \sigma) \\ &= N \log \frac{1}{\sigma\sqrt{2\pi}} - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \end{aligned}$$

- Para una estimación de máxima verosimilitud basta hacer $\nabla L_S(\Theta) = 0$. En nuestro caso unidimensional ([ejercicio](#)):

$$\frac{\partial L_S(\mu, \sigma)}{\partial \mu} = 0 \Rightarrow \hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n; \quad \frac{\partial L_S(\mu, \sigma)}{\partial \sigma} = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

¹ media y desviación típica

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Optimización analítica: otro ejemplo algo menos simple

Ejercicio propuesto:

Estimar los parámetros $\Theta \equiv (\mu_1, \mu_2)$ de una gaussiana bivariada (en \mathbb{R}^2), en la que la matriz de covarianza

$$\Sigma = \begin{bmatrix} \sigma_1 & \sigma_{12} \\ \sigma_{12} & \sigma_2 \end{bmatrix}$$

es conocida:

$$p(x_1, x_2 | \mu_1, \mu_2) = A \cdot \exp\left(-B \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} - \frac{2\sigma_{12}(x_1 - \mu_1)(x_2 - \mu_2)}{(\sigma_1\sigma_2)^2} \right]\right)$$

donde

$$A = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1 - (\sigma_{12}/\sigma_1\sigma_2)^2}}, \quad B = \frac{1}{2(1 - (\sigma_{12}/\sigma_1\sigma_2)^2)}$$

Otro ejercicio algo más complejo propuesto: Asumir que Σ tampoco es conocida; es decir estimar: $\Theta \equiv (\mu_1, \mu_2, \sigma_1, \sigma_2, \sigma_{12})$

Optimización analítica: otro ejemplo

- Estimar los parámetros de una gaussiana multivariada (en \mathbb{R}^D), con Σ dada:

$$p(\mathbf{x} \mid \Theta) = (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

donde $\Theta \equiv (\boldsymbol{\mu}, \Sigma)$. Si Σ está prefijada, entonces $\Theta \equiv \boldsymbol{\mu} \in \mathbb{R}^D$

- Logaritmo de la verosimilitud de una muestra $S = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$:

$$\begin{aligned} q_S(\Theta) \equiv L_S(\Theta) &= \log \prod_{n=1}^N p(\mathbf{x}_n \mid \Theta) = \sum_{n=1}^N \log p(\mathbf{x}_n \mid \Theta) \\ &= N \log \left((2\pi)^{-D/2} |\Sigma|^{-1/2} \right) - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

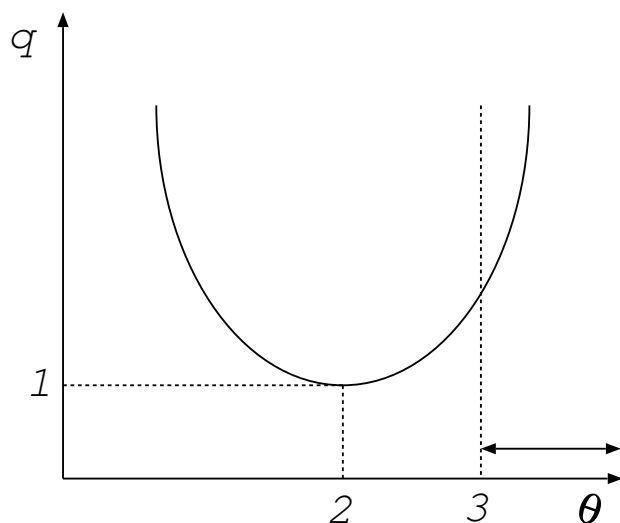
- Para una estimación de máxima verosimilitud basta hacer $\nabla L_S(\Theta) = \mathbf{0}$. Si Σ está prefijada, se obtiene (*[ejercicio propuesto](#)*):

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

Index

- 1 Introducción ▷ 1
- 2 Optimización analítica: gradiente ▷ 5
- 3 *Optimización con restricciones: multiplicadores de Lagrange y teorema Kuhn-Tucker* ▷ 11
- 4 Técnicas de descenso por gradiente ▷ 20
- 5 Esperanza-Maximización (EM) ▷ 31
- 6 Notación ▷ 45

Optimización con restricciones: ejemplo



- Dado $q(\theta) = 1 + (\theta - 2)^2$,
- Calcular: $\theta^* = \arg \min_{\theta \geq 3} q(\theta)$
- Solución: ??

Optimización con restricciones: multiplicadores de Lagrange

Consideremos un problema de optimización definido por:

$$\begin{array}{ll} \text{minimizar} & q(\Theta) \quad \Theta \in \mathbb{R}^D \\ \text{sujeto a} & v_i(\Theta) \geq 0 \quad 1 \leq i \leq k \\ & u_i(\Theta) = 0 \quad 1 \leq i \leq m \end{array}$$

donde q es una función convexa y v_i, u_i son funciones que expresan restricciones.

Equivalentemente, el problema consiste en calcular:

$$\begin{aligned} \Theta^* = \arg \min_{\substack{\Theta \in \mathbb{R}^D \\ v_i(\Theta) \geq 0, 1 \leq i \leq k \\ u_i(\Theta) = 0, 1 \leq i \leq m}} & q(\Theta) \end{aligned}$$

Para resolver este problema, se define la *función Lagrangiana*:

$$\Lambda(\Theta, \alpha, \beta) \stackrel{\text{def}}{=} q(\Theta) - \sum_{i=1}^k \alpha_i v_i(\Theta) + \sum_{i=1}^m \beta_i u_i(\Theta)$$

donde $\alpha_i \geq 0, 1 \leq i \leq k$ y $\beta_i, 1 \leq i \leq m$, son los *multiplicadores de Lagrange*.

La técnica de los multiplicadores de Lagrange

1. Definir multiplicadores de Lagrange y Lagrangiana:

$$\Lambda(\Theta, \alpha, \beta) \stackrel{\text{def}}{=} q(\Theta) - \sum_{i=1}^k \alpha_i v_i(\Theta) + \sum_{i=1}^m \beta_i u_i(\Theta)$$

2. Obtener el minimizador Θ^* de la Lagrangiana $\Lambda(\Theta, \alpha, \beta)$, en función de α, β (resolviendo $\nabla_{\Theta} \Lambda(\Theta, \alpha, \beta) = 0$):

$$\Theta^*(\alpha, \beta) = \arg \min_{\Theta} \Lambda(\Theta, \alpha, \beta)$$

3. Obtener la función dual de Lagrange (sustituir Θ por $\Theta^*(\alpha, \beta)$ en $\Lambda(\Theta, \alpha, \beta)$):

$$\Lambda_D(\alpha, \beta) \stackrel{\text{def}}{=} \Lambda(\Theta^*(\alpha, \beta), \alpha, \beta)$$

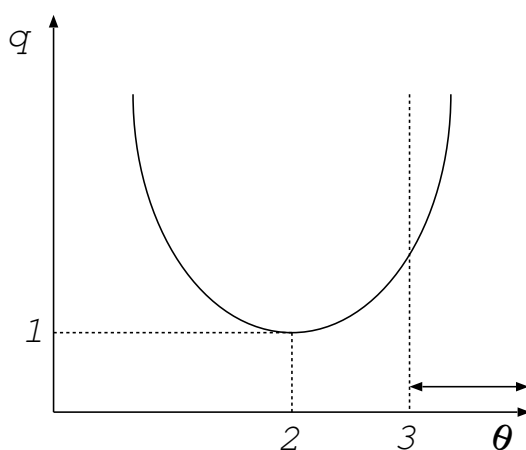
4. Optimizar la función dual de Lagrange (usualmente resolviendo $\nabla \Lambda_D(\alpha, \beta) = 0$):

$$(\alpha^*, \beta^*) = \arg \max_{\alpha, \beta: \alpha_i \geq 0} \Lambda_D(\alpha, \beta)$$

5. Solución final:

$$\Theta^* = \Theta^*(\alpha^*, \beta^*)$$

Multiplicadores de Lagrange: ejemplo



$$\text{minimizar } q(\theta) = 1 + (\theta - 2)^2 \quad \text{con } \theta \geq 3$$

$$\Lambda(\theta, \alpha) = 1 + (\theta - 2)^2 - \alpha (\theta - 3)$$

$$\frac{\partial \Lambda(\theta, \alpha)}{\partial \theta} = 2(\theta - 2) - \alpha = 0 \Rightarrow \theta^*(\alpha) = 2 + \frac{\alpha}{2}$$

$$\Lambda_D(\alpha) = \Lambda(\theta^*(\alpha), \alpha) = 1 + \frac{\alpha^2}{4} - \alpha \left(2 + \frac{\alpha}{2} - 3\right)$$

$$\frac{d \Lambda_D}{d \alpha} = \frac{\alpha}{2} - 2 - \frac{\alpha}{2} + 3 - \frac{\alpha}{2} = 1 - \frac{\alpha}{2} = 0 \Rightarrow$$

$$\alpha^* = 2 \geq 0 \rightarrow \theta^* = \theta^*(\alpha^*) = 3$$

Ejercicio:

a) ¿Y qué ocurre si la condición de desigualdad es $\theta \leq 3$?

b) minimizar $q(\theta) = 1 + (\theta - 2)^2$ con $q(\theta) + \theta = 4$

Multiplicadores de Lagrange: otro ejemplo

En una muestra S de una tarea de clasificación en *tres* clases se observan 4 datos de la clase $c = 1$, 2 datos de $c = 2$ y 1 dato de $c = 3$. Estimar por *máxima verosimilitud* las probabilidades a priori de las clases, $p_c, 1 \leq c \leq 3$.

- Modelo: $P(c = 1) = p_1, P(c = 2) = p_2, P(c = 3) = p_3,$
 $p_1 + p_2 + p_3 = 1, \Theta \equiv (p_1, p_2, p_3)^t$

- Verosimilitud y logaritmo de la verosimilitud:

$$P(S | \Theta) = \prod_{i=1}^4 p_1 \prod_{j=1}^2 p_2 \prod_{k=1}^1 p_3 = p_1^4 p_2^2 p_3$$

$$q_S(\Theta) = L_S(\Theta) = \log P(S | \Theta) = 4 \log p_1 + 2 \log p_2 + \log p_3$$

- Estimación de máxima verosimilitud:

$$\Theta^* = \arg \max_{\Theta} L_S(\Theta) = \arg \max_{\substack{p_1, p_2, p_3 \\ p_1 + p_2 + p_3 = 1}} (4 \log p_1 + 2 \log p_2 + \log p_3)$$

- Problema de optimización con restricciones al que aplicaremos la técnica de los multiplicadores de Lagrange

Ejemplo: aplicación de la técnica de multiplicadores de Lagrange

- Lagrangiana: $\Lambda(p_1, p_2, p_3, \beta) = 4 \log p_1 + 2 \log p_2 + \log p_3 + \beta (1 - p_1 - p_2 - p_3)$
- Soluciones óptimas en función del multiplicador de Lagrange:

$$\left. \begin{aligned} \frac{\partial \Lambda}{\partial p_1} &= \frac{4}{p_1} - \beta = 0 \\ \frac{\partial \Lambda}{\partial p_2} &= \frac{2}{p_2} - \beta = 0 \\ \frac{\partial \Lambda}{\partial p_3} &= \frac{1}{p_3} - \beta = 0 \end{aligned} \right\} \begin{aligned} p_1^*(\beta) &= \frac{4}{\beta} \\ p_2^*(\beta) &= \frac{2}{\beta} \\ p_3^*(\beta) &= \frac{1}{\beta} \end{aligned}$$

- Función dual de Lagrange:

$$\Lambda_D(\beta) = 4 \log \frac{4}{\beta} + 2 \log \frac{2}{\beta} + \log \frac{1}{\beta} + \beta \left(1 - \frac{4}{\beta} - \frac{2}{\beta} - \frac{1}{\beta}\right) = \beta - 7 \log \beta - 7 + 10 \log 2$$

- Valor óptimo del multiplicador de Lagrange: $\frac{d\Lambda_D}{d\beta} = 1 - \frac{7}{\beta} = 0 \Rightarrow \beta^* = 7$

- Solución final: $p_1^* = p_1^*(\beta^*) = \frac{4}{7} \quad p_2^* = p_2^*(\beta^*) = \frac{2}{7} \quad p_3^* = p_3^*(\beta^*) = \frac{1}{7}$

EJERCICIO: Demostrar que en cualquier problema de clasificación en C clases, la estimación de máxima verosimilitud de la probabilidad a priori de cada clase $c, 1 \leq c \leq C$, es $\hat{p}_c = n_c/N$, donde $N = \sum_c n_c$ es el número total de datos observados y n_c es el número de datos de la clase c .

Teorema de Kuhn-Tucker

Consideremos un problema de optimización, \mathcal{O} , definido por:

$$\begin{aligned} &\text{minimizar} && q(\Theta), \quad \Theta \in \mathbb{R}^D \\ &\text{sujecto a} && v_i(\Theta) \geq 0, \quad 1 \leq i \leq k \\ &&& u_i(\Theta) = 0, \quad 1 \leq i \leq m \end{aligned}$$

y la correspondiente función Lagrangiana:

$$\Lambda(\Theta, \alpha, \beta) = q(\Theta) - \sum_{i=1}^k \alpha_i v_i(\Theta) + \sum_{i=1}^m \beta_i u_i(\Theta)$$

Teorema de Kuhn-Tucker: si $\exists \Theta^*, \alpha^*, \beta^*$ tales que:

$$\nabla_{\Theta} \Lambda(\Theta, \alpha^*, \beta^*)|_{\Theta^*} = \mathbf{0};$$

$$\alpha_i^* \geq 0, \quad v_i(\Theta^*) \geq 0, \quad \alpha_i^* v_i(\Theta^*) = 0, \quad 1 \leq i \leq k,$$

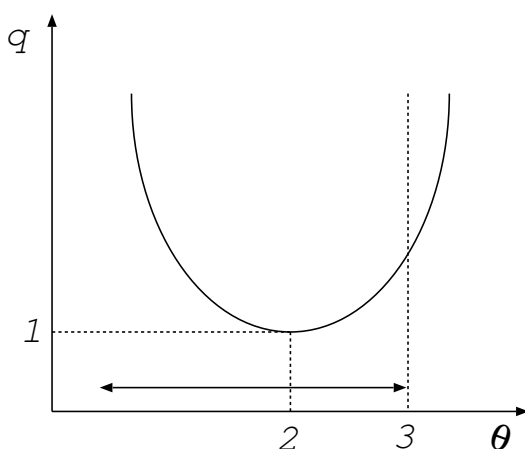
$$u_i(\Theta^*) = 0, \quad 1 \leq i \leq m$$

entonces $q(\Theta^*)$ es solución al problema \mathcal{O} .

$\alpha_i^* v_i(\Theta^*) = 0, 1 \leq i \leq k$: *condiciones complementarias de Karush-Kuhn-Tucker (KKT).*

Multiplicadores de Lagrange y KKT: ejemplo

En el ejemplo anterior, ¿qué ocurre si la condición de desigualdad es $\theta \leq 3$?



$$\text{minimizar } q(\theta) = 1 + (\theta - 2)^2 \quad \text{con } 3 - \theta \geq 0$$

$$\Lambda(\theta, \alpha) = 1 + (\theta - 2)^2 - \alpha (3 - \theta)$$

$$\frac{\partial \Lambda(\theta, \alpha)}{\partial \theta} = 2(\theta - 2) + \alpha = 0 \Rightarrow \theta^*(\alpha) = 2 - \frac{\alpha}{2}$$

$$\text{KKT: } \alpha^* v(\theta^*(\alpha^*)) = 0 \Rightarrow \alpha^* (3 - 2 + \frac{1}{2}\alpha^*) = 0$$

$$\Rightarrow \begin{cases} \alpha^* = 0 \\ \alpha^* = 2(2 - 3) < 0 \rightarrow \text{¡VIOLA } \alpha \geq 0! \end{cases}$$

$$\text{KKT} \Rightarrow \alpha^* = 0 \Rightarrow \theta^* = \theta^*(\alpha^*) = 2$$

Ejercicio: mediante el método de KKT, minimizar $q(\theta) = 1 + (\theta - 2)^2$ con $3 \leq \theta$.

Index

- 1 Introducción ▷ 1
- 2 Optimización analítica: gradiente ▷ 5
- 3 Optimización con restricciones: multiplicadores de Lagrange y teorema Kuhn-Tucker ▷ 11
- 4 *Técnicas de descenso por gradiente* ▷ 20
- 5 Esperanza-Maximización (EM) ▷ 31
- 6 Notación ▷ 45

Descenso por gradiente

Problema: **Minimización sin restricciones de una función objetivo** $q : \mathbb{R}^D \rightarrow \mathbb{R}$, cuando una solución analítica no es viable:

$$\Theta^* = \arg \min_{\Theta} q(\Theta)$$

- Una solución: construir una secuencia de puntos $\Theta(1), \dots, \Theta(k), \dots$, que converja a Θ^* .
- Cada valor $\Theta(k)$ se contruye a partir del anterior $\Theta(k-1)$ en la secuencia dependiendo de las derivadas de la función en el punto $\Theta(k)$.
- Recordatorio: Gradiente en q en el punto $\Theta(k)$: vector formado por las derivadas parciales de la función calculadas en $\Theta(k)$:

$$\nabla q|_{\Theta=\Theta(k)} \equiv \left(\left. \frac{\partial q}{\partial \Theta_1} \right|_{\Theta=\Theta(k)}, \dots, \left. \frac{\partial q}{\partial \Theta_D} \right|_{\Theta=\Theta(k)} \right)^t$$

Descenso por gradiente: algoritmo general

$$\Theta(1) = \text{arbitrario}$$

$$\Theta(k+1) = \Theta(k) - \rho_k \nabla q(\Theta) |_{\Theta=\Theta(k)}$$

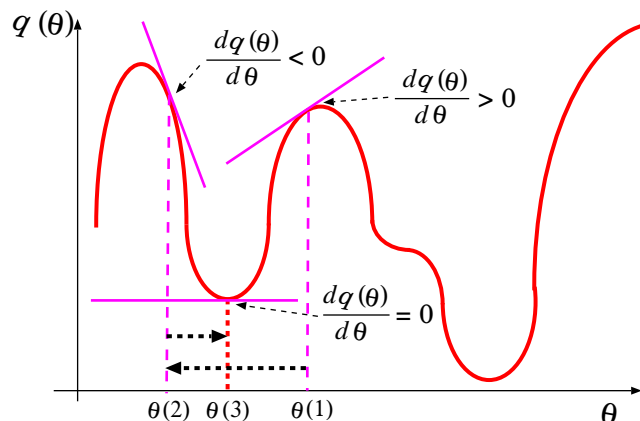
Donde $\rho_k \in \mathbb{R}^{>0}$ es un *factor de aprendizaje*

Ejemplo en \mathbb{R}^1 con $\Theta \stackrel{\text{def}}{=} \theta$:

$$\theta(2) = \theta(1) - \rho_1 \left. \frac{dq}{d\theta} \right|_{\theta(1)}$$

$$\theta(3) = \theta(2) - \rho_2 \left. \frac{dq}{d\theta} \right|_{\theta(2)}$$

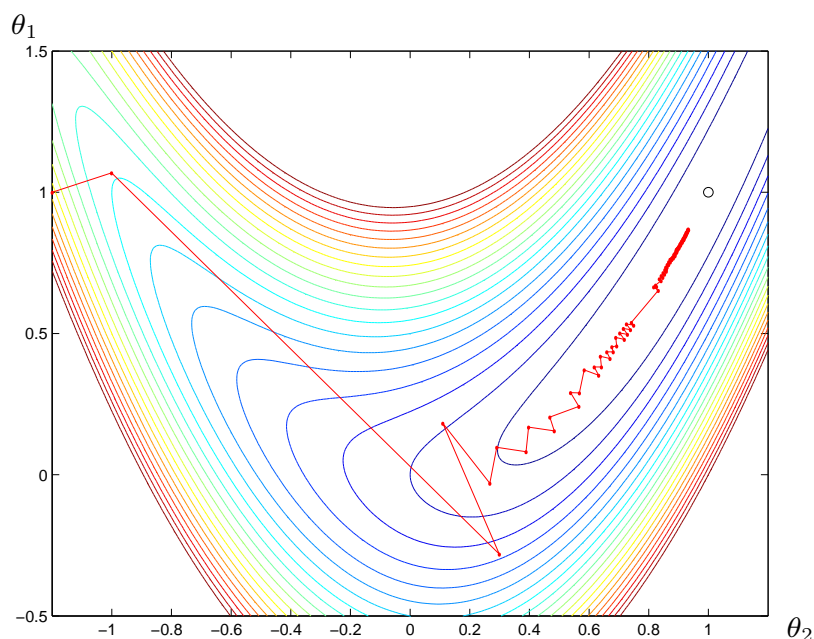
$$\theta(3) \equiv \theta^*, \quad \left. \frac{dq}{d\theta} \right|_{\theta(3)} = 0$$



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Descenso por gradiente: Ejemplo en \mathbb{R}^2



Curvas de nivel de la función de Rosenbrock¹ $q(\theta_1, \theta_2) = 10(\theta_1 - \theta_2^2)^2 + (\theta_1 - 1)^2$ y trayectoria seguida por el vector de parámetros $\theta = (\theta_1, \theta_2)^t$.

1. Figuras basadas en la presentación de R. Hauser http://people.maths.ox.ac.uk/hauser/hauser_lecture2.pdf.

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Convergencia y factor de aprendizaje

- **TEOREMA GENERAL DE CONVERGENCIA:**

Sea $H(q, \Theta)$ la matriz de segundas derivadas (*Hessiana*) de q evaluada en Θ :

$$H_{ij}(q, \Theta) \stackrel{\text{def}}{=} \frac{\partial^2 q(\Theta)}{\partial \Theta_i \partial \Theta_j}$$

Sean $\lambda_l(k)$ los valores propios de $H(q, \Theta(k))$ en el paso k -ésimo del algoritmo de *descenso por gradiente*.

Si $|1 - \lambda_l(k)\rho_k| < 1 \forall l$, entonces $\Theta(k)$ tiende a un mínimo local de $q(\Theta)$ cuando $k \rightarrow \infty$

- **INFLUENCIA DEL FACTOR DE APRENDIZAJE:**

- $\rho < 2/\lambda_{\max}$ garantiza la convergencia
- $\rho \gg \Rightarrow$ convergencia rápida y tendencia a oscilar
- $\rho \ll \Rightarrow$ convergencia lenta

Ejemplo: clasificador lineal en dos clases

- Clasificador en *dos* clases basado *funciones discriminantes lineales* (FDL):

$$f(\mathbf{x}) = \arg \max_{1 \leq c \leq 2} \phi_c(\mathbf{x}), \quad \mathbf{x} = (x_1, \dots, x_d)^t \in \mathbb{R}^d, \quad \phi_c: \mathbb{R}^d \rightarrow \mathbb{R}, \quad 1 \leq c \leq 2$$

Cada FDL ϕ_c está definida por un vector de pesos $\theta_c \in \mathbb{R}^D$ donde $D = d + 1$

En *notación homogénea* se añade una componente, $x_0 \equiv 1$, a \mathbf{x} , con lo que:

$$\phi_c(\mathbf{x}) = \sum_{j=1}^d \theta_{c_j} x_j + \theta_{c_0} = \sum_{j=0}^d \theta_{c_j} x_j = \theta_c^t \mathbf{x}, \quad 1 \leq c \leq 2$$

- Simplificación (si $C = 2$): etiquetar las clases $\{1, 2\}$ como $\{+1, -1\}$ y usar un único vector de pesos $\theta = \theta_1 - \theta_2$.

Clasificador, $f_\theta: \mathbb{R}^D \rightarrow \{-1, +1\}$:

$$f_\theta(\mathbf{x}) = \begin{cases} +1 & \text{si } \theta^t \mathbf{x} \geq 0 \\ -1 & \text{si } \theta^t \mathbf{x} < 0 \end{cases}$$

Aprendizaje de funciones discriminantes lineales

- Sea $S = \{(x_1, c_1), \dots, (x_N, c_N)\}$, $x_n \in \mathbb{R}^D$, $c_n \in \{+1, -1\}$ una muestra de entrenamiento. S es **linealmente separable** (LS) si $\exists \theta \in \mathbb{R}^D$ ($D = d + 1$) tal que:

$$\forall n, 1 \leq n \leq N, \quad \theta^t x_n \begin{cases} \geq 0 & \text{if } c_n = +1 \\ < 0 & \text{if } c_n = -1 \end{cases}; \quad \text{es decir, } c_n \theta^t x_n \geq 0$$

- Aprendizaje:** Dada S , encontrar un vector de pesos $\hat{\theta}$ que la separe; es decir, que satisfaga el sistema de N inecuaciones:

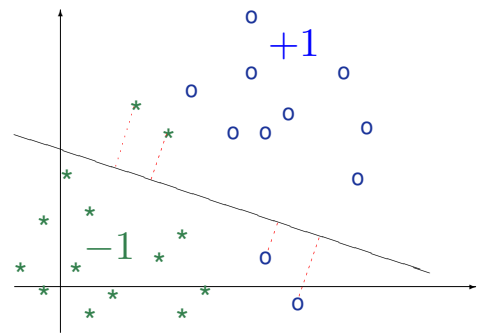
$$c_n \theta^t x_n \geq 0, \quad 1 \leq n \leq N$$

- Planteamiento equivalente:**

minimizar la función: $q_S : \mathbb{R}^D \rightarrow \mathbb{R}^{\geq 0}$:

$$q_S(\theta) = \sum_{\substack{(x,c) \in S \\ c \theta^t x < 0}} -c \theta^t x$$

proporcional a la suma de
segmentos punteados \rightarrow



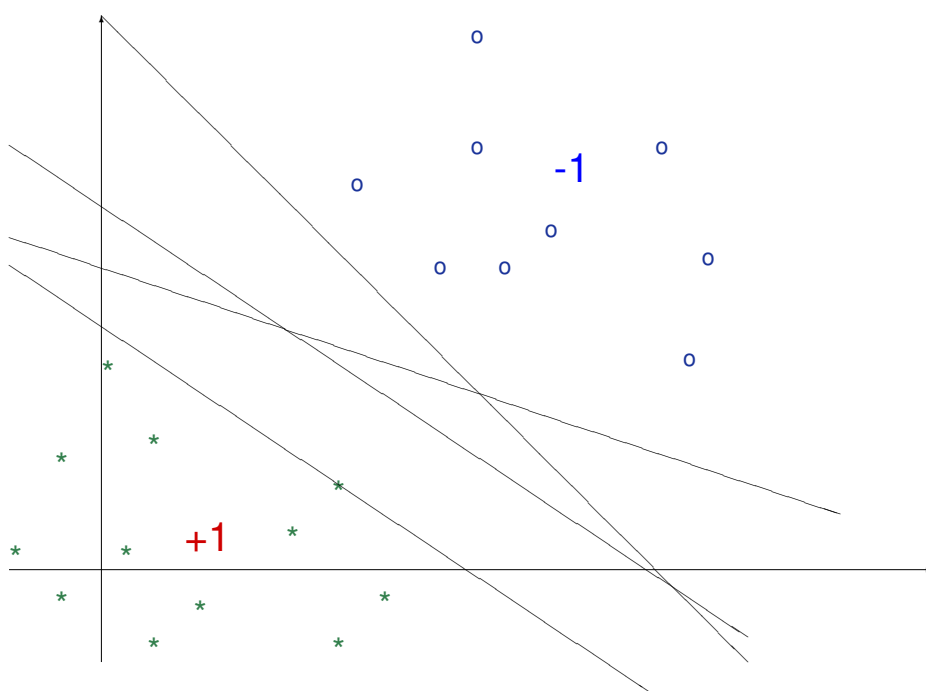
- Sea $\theta^* = \arg \min_{\theta} q_S(\theta)$. S es LS, $\Rightarrow q_S(\theta^*) = 0$, $\hat{\theta} = \theta^*$.

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Aprendizaje de funciones discriminantes lineales

Ejemplo de muestras linealmente separables en \mathbb{R}^2 y posibles soluciones



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Algoritmo perceptrón

$$\nabla q_S(\theta) = \nabla \sum_{\substack{(\mathbf{x},c) \in S \\ c \theta^t \mathbf{x} < 0}} -c \theta^t \mathbf{x} = \sum_{\substack{(\mathbf{x},c) \in S \\ c \theta^t \mathbf{x} < 0}} -c \mathbf{x}$$

$$\begin{aligned} \theta(1) &= \text{arbitrario} \\ \theta(k+1) &= \theta(k) + \rho_k \sum_{\substack{(\mathbf{x},c) \in S \\ c \theta^t \mathbf{x} < 0}} c \mathbf{x} \end{aligned}$$

Algoritmo perceptrón muestra a muestra ("online"):

$$\begin{aligned} \theta(1) &= \text{arbitrario} \\ \theta(k+1) &= \begin{cases} \theta(k) & c(k) \theta^t \mathbf{x}(k) \geq 0 \\ \theta(k) + \rho_k c(k) \mathbf{x}(k) & c(k) \theta^t \mathbf{x}(k) < 0 \end{cases} \end{aligned}$$

TEOREMA DEL PERCEPTRÓN:

Si S es LS y ρ_k es positivo y decreciente o creciente sublinealmente con k , el algoritmo perceptrón converge a una solución en un número finito de iteraciones

Regresión lineal mediante descenso por gradiente

- Sea $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$ una función *lineal* ($D = d + 1$):

$$f_\theta(\mathbf{x}) \stackrel{\text{def}}{=} \theta^t \mathbf{x}, \quad \mathbf{x}, \theta \in \mathbb{R}^D$$

y S una muestra de entrenamiento:

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \quad \mathbf{x}_n \in \mathbb{R}^D, y_n \in \mathbb{R}$$

- Aprendizaje: Calcular $\hat{\theta}$ tal que;

$$\hat{\theta}^t \mathbf{x}_n \approx y_n, \quad 1 \leq n \leq N$$

- Aproximación por mínimos cuadrados:
minimizar la **función de Widrow-Hoff**:

$$q_S(\theta) = \frac{1}{2} \sum_{n=1}^N (\theta^t \mathbf{x}_n - y_n)^2$$

- Solución: descenso por gradiente.

Algoritmo de Widrow-Hoff (Adaline)

$$\nabla q_S(\theta) = \nabla \frac{1}{2} \sum_{n=1}^N (\theta^t x_n - y_n)^2 = \sum_{n=1}^N (\theta^t x_n - y_n) x_n$$

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) + \rho_k \sum_{n=1}^N (y_n - \theta(k)^t x_n) x_n$$

Algoritmo muestra a muestra:

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) + \rho_k \left(y(k) - \theta(k)^t x(k) \right) x(k)$$

TEOREMA:

*Si $\rho_k = \rho_1/k$, $\rho_1 > 0$, $\hat{\theta} = \lim_{k \rightarrow \infty} \theta(k)$ *satisface* $\nabla q_S(\theta)|_{\theta=\hat{\theta}} = 0$*

Index

- 1 Introducción ▷ 1
- 2 Optimización analítica: gradiente ▷ 5
- 3 Optimización con restricciones: multiplicadores de Lagrange y teorema Kuhn-Tucker ▷ 11
- 4 Técnicas de descenso por gradiente ▷ 20
- 5 *Esperanza-Maximización (EM)* ▷ 31
- 6 Notación ▷ 45

Aprendizaje de modelos probabilísticos con variables latentes

- Se suele usar el criterio de *máxima verosimilitud*; es decir, $q_S(\Theta) \equiv L_S(\Theta)$
- En ocasiones los datos observados *no* contienen suficiente información sobre cómo han sido generados por los modelos probabilísticos asumidos
- Por ejemplo, en los modelos ocultos de Markov los datos de entrenamiento son cadenas de símbolos, sin información sobre qué *secuencia de estados* ha producido cada cadena
- Otro ejemplo típico son los modelos definidos como combinación lineal (“mezcla” o “mixture”) de distribuciones de probabilidad. Los coeficientes de combinación son parámetros a aprender, pero los datos de entrenamiento no contienen información sobre la distribución con que se ha generado cada dato
- La información ausente en los datos de entrenamiento generalmente se denomina datos *perdidos*, o variables *latentes* u *ocultas*
- Las técnicas simples de optimización resultan insuficientes para la estimación de los parámetros de estos modelos

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Ejemplo: mezcla de gaussianas y modelo generador

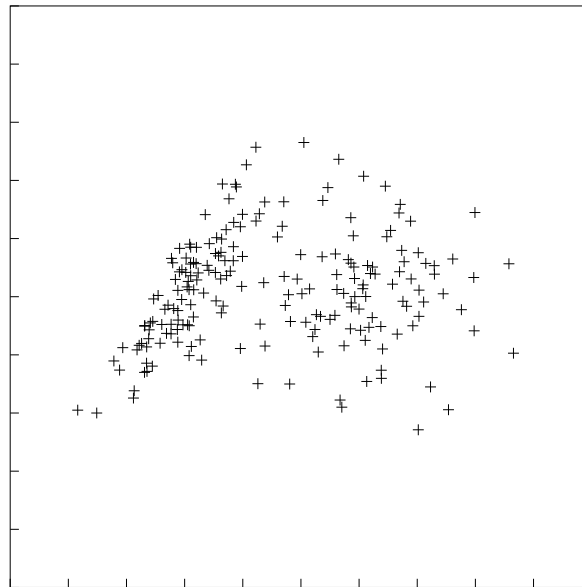
- En el caso de una única gaussiana las muestras se generan en un paso:
 1. Escoger x , de acuerdo con la distribución $p(x \mid \mu, \Sigma)$
- Si el modelo es una *mezcla* de K gaussianas, el proceso de generación se compone de dos etapas:
 1. De acuerdo con la distribución $P(k) = \pi_k$, escoger la componente k -ésima de la mezcla con la que se va a generar x
 2. Escoger x , según la distribución definida por la k -ésima gaussiana, $p(x \mid \mu_k, \Sigma_k)$
 - x es el dato *observable* y k una variable *oculta*. Los datos observables junto con los ocultos se denominan *datos completos*
 - Probabilidad con la que se genera x según este proceso:

$$\begin{aligned}
 p(x) &= \sum_{k=1}^K p(k, x) = \sum_{k=1}^K P(k) p(x \mid k) \equiv \sum_{k=1}^K \pi_k p(x \mid \mu_k, \Sigma_k) \\
 &\equiv p(x \mid \Theta); \quad \Theta \stackrel{\text{def}}{=} [\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K]
 \end{aligned}$$

Septiembre, 2016

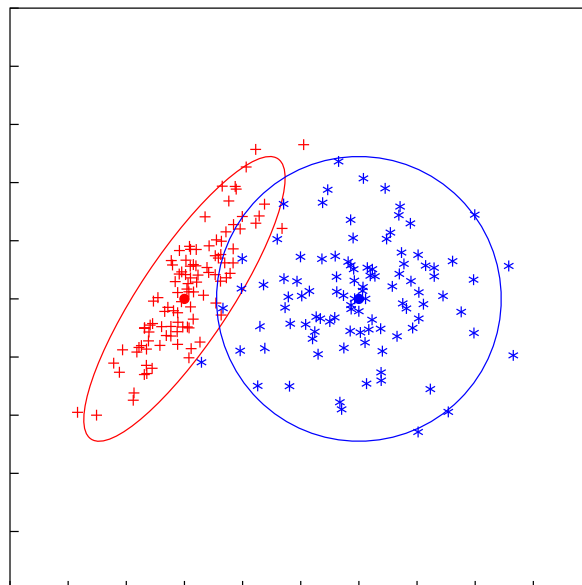
Departament de Sistemes Informàtics i Computació

Mezcla de gaussianas: ilustración



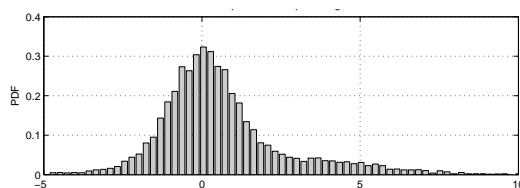
Datos *observables* generados por una mezcla de dos gaussianas.

Mezcla de gaussianas: ilustración

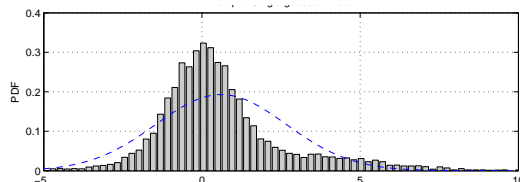


Datos de una mezcla de dos gaussianas con la variable oculta expuesta.
Las elipses muestran los parámetros de las gaussianas del modelo generador.

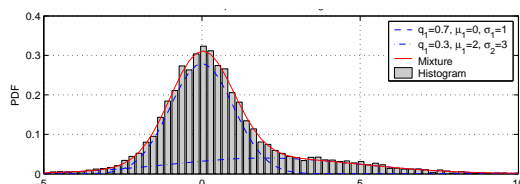
Mezcla de gaussianas: otro ejemplo en 1D



Histograma de una muestra de una variable unidimensional



Una única gaussiana estimada a partir de la muestra



Mezcla de dos gaussianas con la que se ha generado de la muestra

Parámetros de una mezcla de gaussianas

- Para simplificar, suponemos que la matriz de covarianzas, Σ , de todas las K gaussianas de la mezcla es la misma y es fija y conocida:

$$\begin{aligned}
 p(\mathbf{x} \mid \Theta) &= \sum_{k=1}^K \pi_k p(\mathbf{x} \mid \mu_k) \\
 &= (2\pi)^{-d/2} |\Sigma|^{-1/2} \sum_{k=1}^K \pi_k \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_k)^t \Sigma^{-1} (\mathbf{x} - \mu_k) \right\}
 \end{aligned}$$

$$\text{donde } \Theta \equiv [\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K], \quad \sum_{k=1}^K \pi_k = 1$$

Mezcla de K gaussianas en \mathbb{R}^d : variables latentes

Sea $\Theta = (\pi_k, \mu_k, 1 \leq k \leq K)$ los parámetros de la mezcla de gaussianas $p(\mathbf{x} | \Theta)$ y sea $S = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^d$, $1 \leq n \leq N$, una muestra de esta distribución.

- Para cada dato \mathbf{x}_n , hay K variables latentes $z_{kn} \in \{0, 1\}$, $1 \leq k \leq K$

- $z_{kn} = 1 \Leftrightarrow \mathbf{x}_n$ ha sido generado por la gaussiana k ; por tanto:

$$\sum_{k=1}^K z_{kn} = 1 \quad \forall n \quad \text{y} \quad \pi_k \text{ es la prob. "a priori" (independiente de } \mathbf{x}_n) \text{ de que } z_{kn} = 1$$

- **Valor esperado de z_{kn}** , dado \mathbf{x}_n y Θ :

$$\begin{aligned} \hat{z}_{kn} &\stackrel{\text{def}}{=} E(z_{kn}) = \sum_{z_{kn}=0}^1 z_{kn} P(z_{kn} | \mathbf{x}_n, \Theta) \\ &= P(z_{kn} = 1 | \mathbf{x}_n, \Theta) = \frac{P(z_{kn} = 1 | \Theta) p(\mathbf{x}_n | z_{kn} = 1, \Theta)}{p(\mathbf{x}_n | \Theta)} \\ &= \frac{\pi_k p(\mathbf{x}_n | \mu_k)}{p(\mathbf{x}_n | \Theta)} = \frac{\pi_k p(\mathbf{x}_n | \mu_k)}{\sum_{k'=1}^K \pi_{k'} p(\mathbf{x}_n | \mu_{k'})} \end{aligned}$$

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Estimación de parámetros de una mezcla de gaussianas

- El logaritmo de la verosimilitud de la muestra $S = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, dada una estimación de las variables latentes $\{\hat{z}_1, \dots, \hat{z}_N\}$ es:

$$L'_S(\Theta) = \log \prod_{n=1}^N p(\mathbf{x}_n | \hat{z}_n; \Theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \hat{z}_n; \mu_k)$$

- Estimación de $\Theta \equiv [\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K]$ por máxima verosimilitud:

$$\hat{\Theta} = \arg \max_{\substack{\Theta \\ \sum_{k=1}^K \pi_k = 1}} L'_S(\Theta)$$

Optimización analítica:

$$1 \leq k \leq N : \quad \nabla_{\mu_k} L'_S(\Theta) = 0 \rightarrow \mu_k = \frac{1}{\pi_k N} \sum_{n=1}^N \hat{z}_{kn} \mathbf{x}_n$$

Multiplicadores de Lagrange $[\Lambda(\Theta, \beta) = L'_S(\Theta) + \beta(1 - \sum_{k=1}^K \pi_k)]$:

$$1 \leq k \leq K : \quad \frac{\partial \Lambda}{\partial \pi_k} = 0 \rightarrow \pi_k = \sum_{n=1}^N \frac{\hat{z}_{kn}}{\beta} \rightarrow \left[\frac{\partial \Lambda_D}{\partial \beta} = 0 \right] \rightarrow \pi_k = \frac{1}{N} \sum_{n=1}^N \hat{z}_{kn}$$

Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Estimación de parámetros de una mezcla de gaussianas: Algoritmo EM

1. Inicializar $i = 0$ y $\Theta(i) \equiv (\pi_k(i), \mu_k(i), 1 \leq k \leq K)$ con valores adecuados
2. Iterar hasta convergencia:

- **Paso E** (“expectación” o cálculo de la esperanza de variables latentes):

$$\left. \begin{array}{l} 1 \leq n \leq N \\ 1 \leq k \leq K \end{array} \right\} : \hat{z}_{kn} = \frac{\pi_k(i) p(\mathbf{x}_n | \mu_k(i))}{\sum_{k'} \pi_{k'}(i) p(\mathbf{x}_n | \mu_{k'}(i))}$$

- **Paso M** (“maximización” de la log-verosimilitud, $L_S(\Theta)$):

$$1 \leq k \leq K : \left\{ \begin{array}{l} \pi_k(i+1) = \frac{1}{N} \sum_{n=1}^N \hat{z}_{kn} \\ \mu_k(i+1) = \frac{1}{\pi_k N} \sum_{n=1}^N \hat{z}_{kn} \mathbf{x}_n \end{array} \right.$$

- $i \leftarrow i + 1$

Algoritmo EM basado en función Q

- Se define una función de Θ a partir de valores dados (previos) de Θ' :

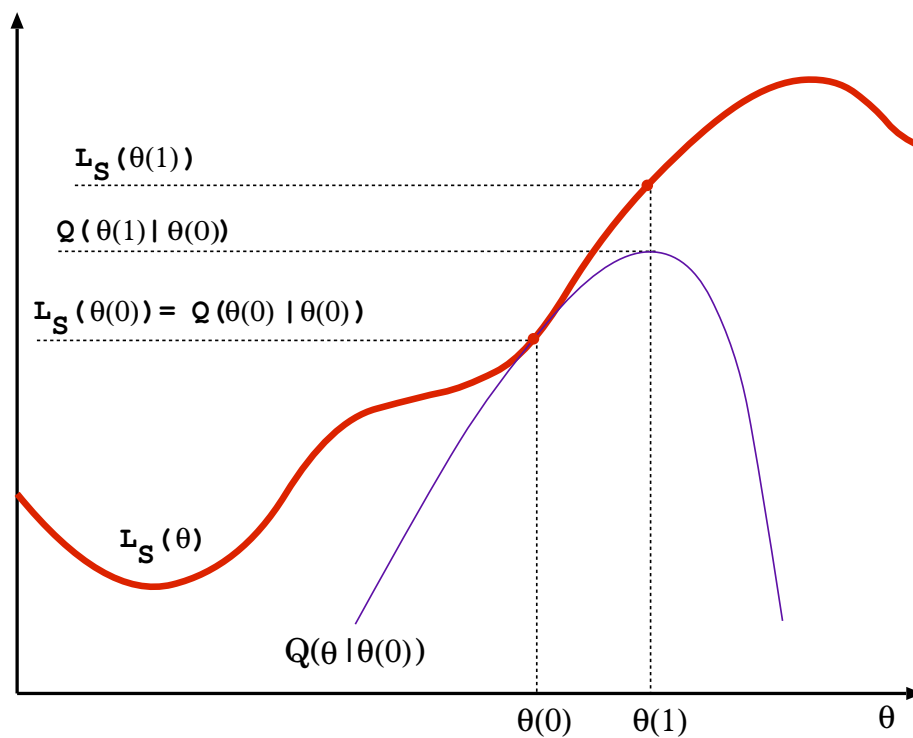
$$Q(\Theta, \Theta') \stackrel{\text{def}}{=} E_{z|x, \Theta'}(\log P(\mathbf{x}, z | \Theta)) = \sum_z \sum_{n=1}^N \log P(\mathbf{x}_n, z | \Theta) P(z | \mathbf{x}_n, \Theta')$$

- Propiedades: 1) $Q(\Theta, \Theta')$ es convexa; 2) $Q(\Theta, \Theta') \leq L_S(\Theta) \forall \Theta, \Theta'$

Algoritmo EM basado en Q :

1. Inicializar $\Theta(0)$; $i = 0$
2. Repetir, hasta convergencia:
 - Paso E (esperanza): Obtener $Q(\Theta, \Theta(i))$ a partir de $\Theta(i)$
 - Paso M (maximización): Calcular $\Theta(i+1) = \arg \max_{\Theta} Q(\Theta, \Theta(i))$
 - $i = i + 1$

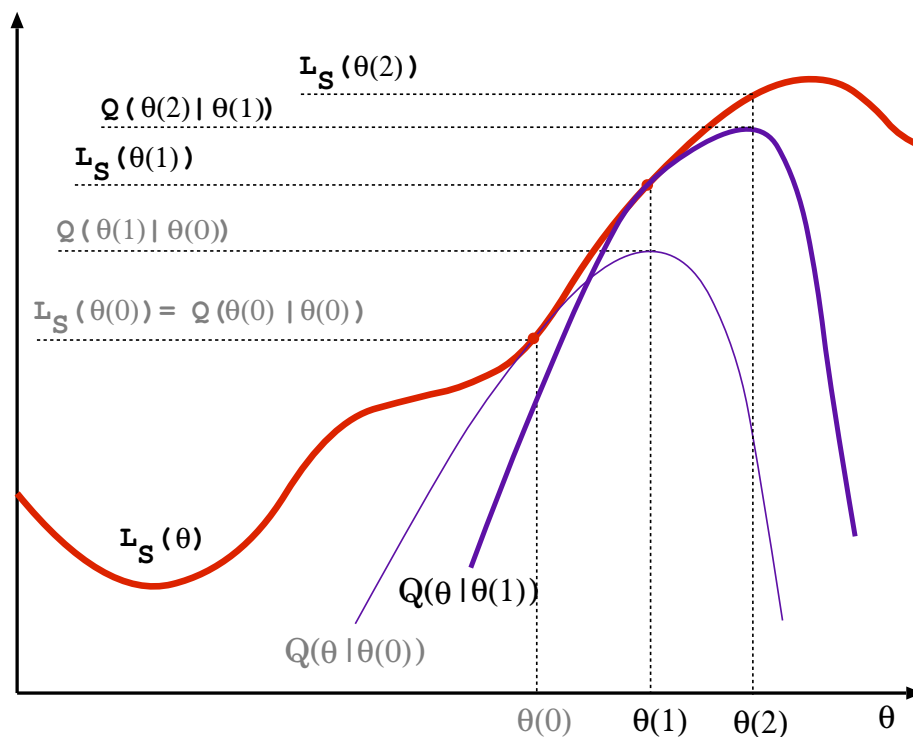
Propiedades y convergencia del EM



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

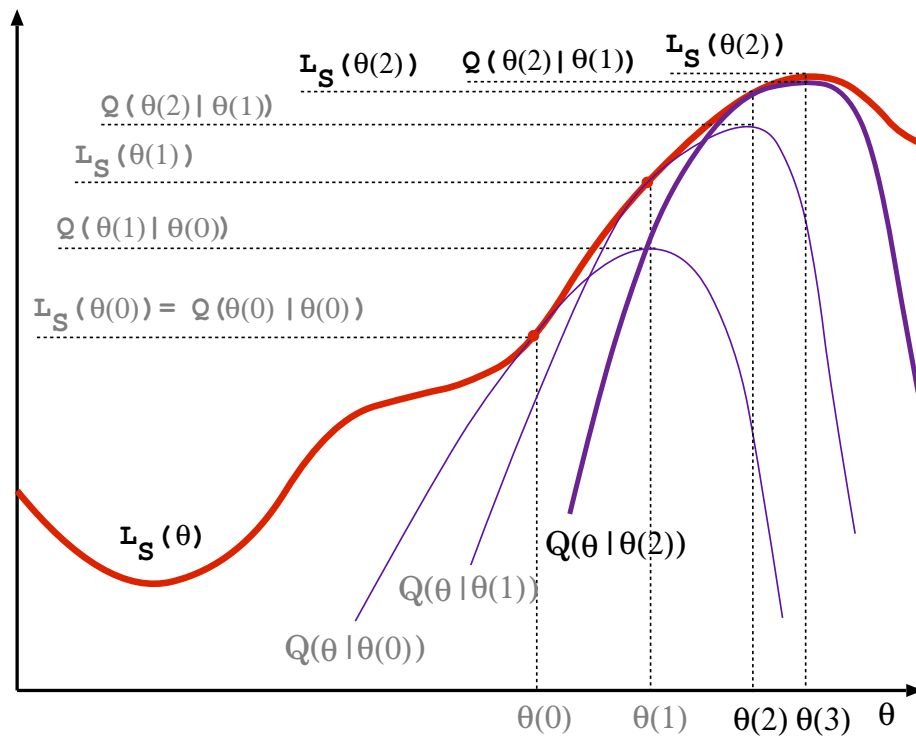
Propiedades y convergencia del EM



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Propiedades y convergencia del EM



Septiembre, 2016

Departament de Sistemes Informàtics i Computació

Notación

- $\Theta = (\Theta_1, \dots, \Theta_D)^t$: vector de parámetros. Como los vectores son matrices columna, para representar las componentes en fila se usa la t ("transpuesta").
- $q_S(\Theta)$: función objetivo a optimizar definida sobre un conjunto de entrenamiento S , cuyos de parámetros son Θ
- $\nabla q(\Theta) \stackrel{\text{def}}{=} \left(\frac{\partial q(\Theta)}{\partial \Theta_1}, \dots, \frac{\partial q(\Theta)}{\partial \Theta_D} \right)^t$: gradiente de la función q_s
(vector de derivadas parciales con respecto a cada componente de Θ)
- $\nabla q(\Theta) |_{\Theta=\Theta^{(k)}} \equiv \left(\frac{\partial q}{\partial \Theta_1} \Big|_{\Theta=\Theta^{(k)}}, \dots, \frac{\partial q}{\partial \Theta_D} \Big|_{\Theta=\Theta^{(k)}} \right)^t$: gradiente de la función q calculado en $\Theta^{(k)}$
- Σ : matriz de covarianzas de una distribución gaussiana

Septiembre, 2016

Departament de Sistemes Informàtics i Computació