

today : • Fisher LDA
• math. tricks & MLE for Gaussian

generative model for classification: (Fisher) linear discriminant analysis
FLD (instead of LDA)

for classification $y \in \{0, 1\}$

$X \in \mathbb{R}^d$ class conditional

generative approach $p(x, y; \theta) = p(x|y; \theta)p(y; \theta)$
vs.

conditional approach $p(y|x; \theta)$

skewed across classes

For Fisher model : we assume $p(x|y; \theta) = N(x|\mu_y, \Sigma)$

$$\theta = (\mu_0, \mu_1, \Sigma, \pi)$$

mean
of class 0 skewed
 $p(y=1)$

as before (see exponential family argument)

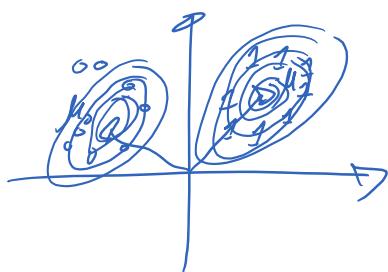
can show that $p(y|x; \theta) = \mathcal{O}(w^T x)$ where w is a ff. of $(\mu_0, \mu_1, \Sigma, \pi)$

[note : if use use $\Sigma_0 \neq \Sigma_1$, get "quadratic discriminant analysis"]

i.e. $\mathcal{O}(w^T Q(x))$ where $Q(x)$ is a quadratic ff. of x [see hwk 2]

gen. approach : do joint MLE to estimate

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_i \log p(x_i, y_i; \theta)$$



[vs. $\underset{w \in W}{\operatorname{argmax}} \sum_i \log p(y_i|x_i; w)$
for logistic regression]

sideneote : MLE for multivariate Gaussian

$$x_i \sim N(\mu, \Sigma)$$

$$\mu \in \mathbb{R}^d$$

$$\Sigma \in \mathbb{R}^{d \times d}$$

$$\Sigma \stackrel{\Delta}{=} E[(x-\mu)(x-\mu)^T]$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

$$\begin{array}{l} \mu \in \mathbb{R}^m \\ \Sigma \in \mathbb{R}^{d \times d} \\ \Sigma \text{ is symmetric} \\ \Sigma \geq 0 \end{array} \quad \begin{array}{l} \Sigma^T = \Sigma \\ \Sigma = \Sigma^T \\ \Sigma \geq 0 \\ \Sigma \geq 0 \end{array}$$

$$v^T \Sigma v = \mathbb{E}[v^T (x-\mu)(x-\mu)^T v]$$

$$(x-\mu)^T v)^2 \geq 0$$

$$p(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} \underbrace{(x-\mu)^T \Sigma^{-1} (x-\mu)}_{\text{tr}((x-\mu)^T \Sigma^{-1} (x-\mu))}\right) \quad \text{tr}(AB) = \text{tr}(BA)$$

$$\Sigma = (\mu, \Sigma)$$

$$\begin{aligned} & \text{tr}((x-\mu)^T \Sigma^{-1} (x-\mu)) \\ &= \text{tr}(\Sigma^{-1} (x-\mu)(x-\mu)^T) = \langle \Sigma^{-1}, (x-\mu)(x-\mu)^T \rangle \\ & \langle A, B \rangle \triangleq \sum_{i,j} A_{ij} B_{ij} = \text{tr}(A^T B) \end{aligned}$$

$$\text{log-likelihood: } \sum_{i=1}^n \log p(x_i; \theta) = \text{const.} - \frac{n \log |\Sigma|}{2} - \frac{1}{2} \sum_{i=1}^n \langle \Sigma^{-1}, (x_i - \mu)(x_i - \mu)^T \rangle + \frac{n}{2} \log |\Sigma|$$

$$|\Sigma^{-1}| = \frac{1}{|\Sigma|}$$

vector derivative review:

suppose $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$\mathcal{O}(\|\Delta\|)$ "little oh"
means that if $\Delta \rightarrow 0$, $f(x_0 + \Delta) \approx f(x_0) + \mathcal{O}(\|\Delta\|)$

f is differentiable at x_0 iff \exists a linear operator $df_{x_0}: \mathbb{R}^m \rightarrow \mathbb{R}^n$

s.t. $\forall \Delta \in \mathbb{R}^m$ $f(x_0 + \Delta) - f(x_0) = df_{x_0}(\Delta) + \mathcal{O}(\|\Delta\|)$

"derivative"

$$\lim_{\|\Delta\| \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\|\Delta\|} = \lim_{\|\Delta\| \rightarrow 0} \left(\underbrace{\frac{df_{x_0}(\Delta)}{\|\Delta\|}}_{df_{x_0}(\frac{\Delta}{\|\Delta\|})} + \underbrace{\mathcal{O}(\|\Delta\|)}_{\text{directional derivative}} \right)$$

→ directional derivative
of f at x_0 in direction d

"differential"

df_{x_0} is linear

means $df_{x_0}(\Delta_1 + b\Delta_2) = df_{x_0}(\Delta_1) + b df_{x_0}(\Delta_2)$

can represent as a $n \times m$ matrix
called the Jacobian matrix

standard representation $(df_{x_0})_{ij} = \frac{\partial f_j}{\partial x_i}$
then $df_{x_0}(\Delta) = df_{x_0} \cdot \Delta$

• if $n=1$

$$df_{x_0}(d) = \langle \nabla f(x_0), d \rangle$$

$$\nabla f(x_0) = (df_{x_0})^T$$

1) this gives a way to get df_{x_0} for
"anything" (matrix, tensor, scalar function, etc.)

2) be careful with dimensions

$f: \mathbb{R}^m \rightarrow \mathbb{R}$
 df_{x_0} is a row vector $(1 \times m)$
 $d \in \mathbb{R}^{m \times 1} = \mathbb{R}^{m \times 1 \times T}$

chain rule:

$$r \cdot m \times n \rightarrow m \times n$$

$$\underbrace{d(\cdot \circ \dots \circ \cdot)}_{dm \times \dots \times dm} \quad dm \times \dots \times dm$$

chain rule:

$$\begin{array}{l} f: \mathbb{R}^m \rightarrow \mathbb{R}^n \\ g: \mathbb{R}^n \rightarrow \mathbb{R}^q \end{array}$$

$$d(g \circ f)_{x_0} = dg_{f(x_0)} \circ df_{x_0}$$

\downarrow

$$g(f(x_0)) \quad = \quad (\quad)(\quad) \quad \text{matrix product}$$

\vec{c} is a row vector ($1 \times m$)

$$df_{x_0} = (\nabla f(x_0))^T$$

1

$$f(\mu) = x - \mu$$

$$df_{\mu_0} = -I$$

$$g(w) = w^T A w$$

$$d\hat{q}_{w_0} = w_0^T (A + A^T)$$

$$g_{\text{of}}(\mu) = (\mathbf{x} - \mu)^T A (\mathbf{x} - \mu)$$

$$d(g \circ f)_{\mu_0} = dg_{f(\mu_0)} \circ df_{\mu_0}$$

$$= (x - \mu_0)^T (A + A^T) (-I)$$

für Gaußkern: $\sim \frac{1}{\sqrt{\pi}} \mathcal{E}(\mathbf{x}_i - \mu)^T \mathcal{E}^{-1}(\mathbf{x}_i - \mu)$

$$\nabla_{\mu} :=$$

$$\frac{1}{2} \sum_i \mathbf{x}_i^\top (\mathbf{x}_i - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \\ + \frac{1}{2} \sum_i 2 \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) = 0 \stackrel{\text{want}}{\Rightarrow} \hat{\boldsymbol{\mu}}_{\text{MLE}} = \frac{1}{n} \sum_i \mathbf{x}_i$$

15h58

example 2 : derivative of $f(A) \triangleq \log \det(A)$ where assume A is symmetric $A > 0$

can represent the derivatives of a function from matrix to scalar, as a matrix

$$\begin{aligned} f(A + \Delta) - f(A) &= \text{tr} \left(f'(A)^T \Delta \right) + o(\|\Delta\|) \\ &= \underbrace{\langle f'(A), \Delta \rangle}_{+o(\|\Delta\|)} \end{aligned}$$

$$\begin{aligned}
 & \log \det(A + \Delta) - \log \det(A) \\
 &= \log \det \left(A^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) A^{1/2} \right) \sim \log \det(A) \\
 &= \log \underbrace{|A|^{1/2} |I + A^{-1/2} \Delta A^{-1/2}| |A|^{1/2}}_{\text{e-value of } B} - \log |A| \\
 &= \log |I + A^{1/2} \Delta A^{-1/2}| \quad \text{use } \det(B) = \prod_i \lambda_i(B) \\
 &\approx \sum_i \log \lambda_i(I + A^{-1/2} \Delta A^{-1/2}) \quad Bv = \lambda v \\
 &= \sum_i \log (1 + \lambda_i(A^{-1/2} \Delta A^{-1/2})) \\
 &= \sum_i \left[\lambda_i(A^{-1/2} \Delta A^{-1/2}) + \mathcal{O} \left(\underbrace{\lambda_i^2 (A^{-1/2} \Delta A^{-1/2})^2}_{\text{cat.}} \right) \right] \quad \log(1+x) = x + \mathcal{O}(x^2) \text{ for } |x| < 1 \\
 &\qquad\qquad\qquad \lambda \text{ is homogeneous fct.} \\
 &\qquad\qquad\qquad \text{i.e. } Bv = \lambda v \\
 &\qquad\qquad\qquad \left(\frac{B}{\lambda}\right)v = \left(\frac{\lambda}{\lambda}\right)v \\
 &= \mathfrak{H}(A^{-1/2} \Delta A^{-1/2}) + \mathcal{O}(|\Delta|) \\
 &\qquad\qquad\qquad \frac{|\mathfrak{H}(\Delta)|^2}{|\Delta|} = |\Delta| \xrightarrow{\Delta \rightarrow 0} 0 \quad \text{as } |\Delta| \rightarrow 0
 \end{aligned}$$

$$\begin{aligned}
 &= \text{tr} (A^{-1/2} \Delta A^{-1/2}) + o(\|\Delta\|) \\
 &= \underbrace{\text{tr}(A^{-1} \Delta)}_{\langle A^{-1}, \Delta \rangle} + o(\|\Delta\|) \\
 &\quad (\text{recall } A \text{ is symmetric}) \\
 &\Rightarrow \boxed{\frac{d}{dA} \log \det(A) = A^{-1}}
 \end{aligned}$$

see [Boyd's book](#) A.4.1 for the above proof

back to log-likelihood of Gaussian:

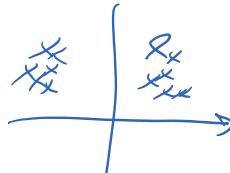
$$+ \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \langle \Sigma^{-1}, \tilde{\Sigma}(\mu) \rangle \quad (\text{concave fct. of } -L = \Sigma^{-1})$$

take derivative w.r.t.

$$\begin{aligned}
 \Sigma^{-1} = \Lambda &\quad \frac{n}{2} \underbrace{(\Sigma^{-1})^{-1}}_{\Sigma} - \frac{n}{2} \tilde{\Sigma}(\mu) \stackrel{\text{want}}{=} 0 \\
 &\Rightarrow \boxed{\begin{aligned} \hat{\Sigma}_{MLE} &= \tilde{\Sigma}(\mu_{MLE}) \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^\top \end{aligned}} \\
 &\quad (\text{the empirical covariance matrix})
 \end{aligned}$$

Unsupervised learning

here X without any label Y



consider the Gaussian mixture model (GMM)
(can be obtained from FLD)

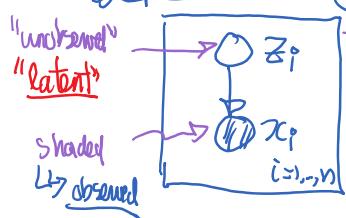
$Y \sim \text{mult}(\pi) \quad \pi \in \Delta_K$ [extension of FLD to multiple classes]

$X|Y=j \sim N(\mu_j, \Sigma)$

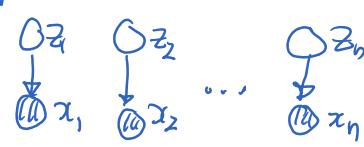
$$p(x) = \sum_j p(x, y) = \sum_j p(x|y)p(y) = \sum_{j=1}^K \pi_j N(x; \mu_j, \Sigma)$$

"GMM model" (more generally, can have $\leq j$ per class)

graphical model for this "latent variable model"



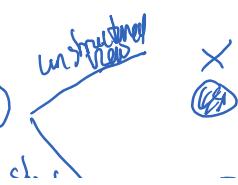
→ "plate" = repetition

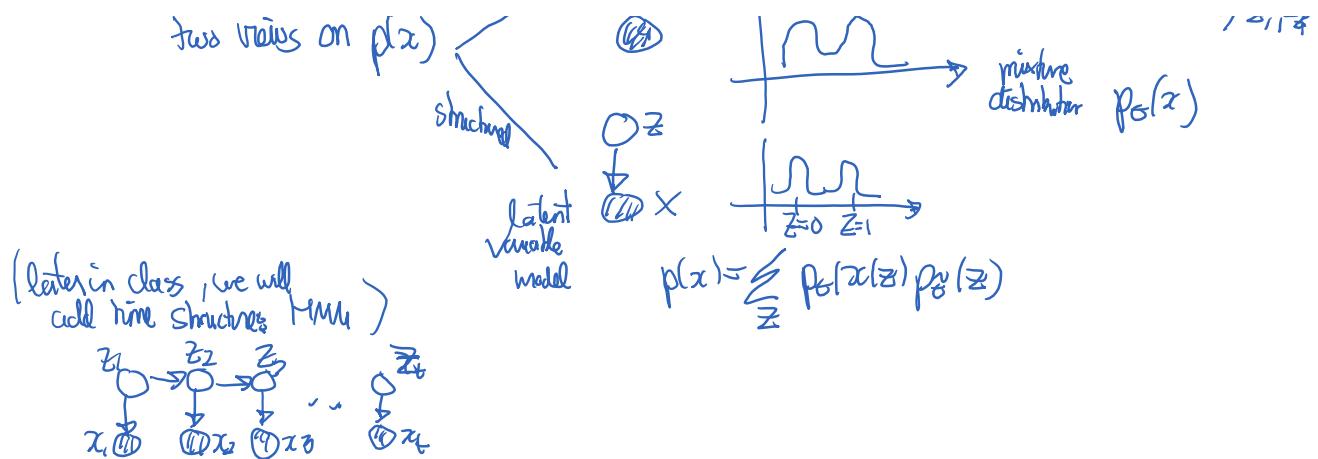


GMM model:

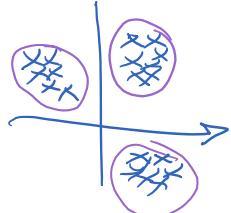
$$\begin{aligned}
 Z_i &\sim \text{Mult}(\pi) \\
 x_i | z_i &\sim N(x_i; \mu_{z_i}, \Sigma)
 \end{aligned}$$

two views on $p(x)$





K-means → to do clustering ie. Group data



we want to get a cluster assignment for every data pt. z_i :

represent $z_{i,j} = 1$ to mean \hat{z}_i belongs to cluster j

$j = 1, \dots, k$
k = # of clusters (specified in advance for K-means)

applications:

- vector quantization (compression)

- in computer vision: use K-means to get "bag of visual words" representation of image patches

- many many others!