

- today:
- Fisher LDA
  - math tricks & MLE for Gaussian

generative model for classification : (Fisher) linear discriminant analysis

FLD (instead of LDA)

for classification  $\mathcal{Y} = \{0, 1\}$   
 $X \in \mathbb{R}^d$

generative approach  $p(x, y; \Theta) = \underbrace{p(x|y; \Theta)}_{\text{class conditional}} p(y; \Theta)$   
 vs.

conditional approach  $p(y|x; \Theta')$

⊕ for Fisher model: we assume  $p(x|y; \Theta) \sim N(x | \mu_y, \Sigma)$

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

$\mu_0$  mean for class 0  
 $\mu_1$  mean for class 1  
 $\Sigma$  shared covariance  
 $\pi$  shared across classes

$$\left( \begin{array}{c} \mu_0 \\ \mu_1 \\ \text{vec}(\Sigma) \\ \pi \end{array} \right)$$

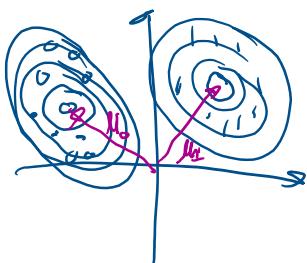
as before (see exponential family argument)

can show that  $p(y|x; \Theta) = \sigma(w^T x)$  where  $w$  is a fct. of  $(\mu_0, \mu_1, \Sigma, \pi)$

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

i.e.  $\sigma(w^T q(x))$  where  $q(x)$  is a quadratic fct. of  $x$  [see hwk. 2]

⊕ generative approach: do joint MLE to estimate  $\Theta$ :



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

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

sidenote: MLE for multivariate Gaussian

$$x_i \stackrel{iid}{\sim} N(\mu, \Sigma)$$

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

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

$\Sigma$  is symmetric

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

$$\Sigma^T = \mathbb{E}[ \quad ] = \Sigma$$

$$V^T \Sigma V = \mathbb{E}[\underbrace{y^T(x-\mu)}_{\sim \sim \sim} \underbrace{(x-\mu)^T V}_{\sim \sim \sim}]$$

$\Sigma \in \mathbb{R}^{d \times d}$   
 $\Sigma$  is symmetric  
 $\Sigma > 0$

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

$$\mathbb{E}[(x-\mu)^T v]^2 \geq 0 \Rightarrow \Sigma \geq 0$$

$$p(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\langle A, B \rangle = \sum_{i,j} A_{ij} B_{ij}$$

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

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

$$h \in o(\|\Delta\|) \quad \text{"little oh"} \quad \text{e.g. } \|\Delta\|^2 \text{ is } o(\|\Delta\|)$$

$$\Leftrightarrow \lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} = 0$$

Vector derivative review:

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

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

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

"differential"

"derivative"

$$\lim_{\delta \rightarrow 0} \frac{f(x_0 + \delta \vec{d}) - f(x_0)}{\delta} = \lim_{\delta \rightarrow 0} \frac{df_{x_0}(\delta \vec{d}) + o(\delta \|\vec{d}\|)}{\delta}$$

$$\underset{\delta \rightarrow 0}{\cancel{\frac{df_{x_0}(\delta \vec{d})}{\delta}}} + o(\|\vec{d}\|)$$

→ directional derivative at  $x_0$  in direction  $\vec{d}$

$$\text{if } n=1: \quad df_{x_0}(\vec{d}) = \langle \nabla f(x_0), \vec{d} \rangle$$

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

$df_{x_0}$  is linear

means that

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

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

standard representation

$$(df_{x_0})_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\text{then } df_{x_0}(\Delta) = df_{x_0} \cdot \Delta$$

1) this gives you a way to get  $df_{x_0}$  for "anything"  
(matrix, tensor,  $\infty$ -dim jet, etc...)

2) be careful with dimensions

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

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

chain rule

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

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

$g(f(x_0))$

$= (\quad)(\quad)$

matrix product of Jacobians

$$\text{eg } f(\mu) = x - \mu$$

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

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

$$dg_{v_0} = \sqrt{A + A^T}$$

$$g \circ f(\mu) = (x - \mu)^T A (x - \mu)$$

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

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

for Gaussian:  $\frac{1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$

$$\left( \frac{\partial g \circ f}{\partial \mu} \right)^T \left[ + \frac{1}{2} \sum_i 2 \Sigma^{-1} (x_i - \mu) \right] = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

15h38

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

can represent the derivative of a ff. from matrix to scalar as a matrix

$$f(A + \Delta) - f(A) = \text{tr}(f'(A)^T \Delta) + o(\|\Delta\|)$$

$$= \langle f'(A), \Delta \rangle + o(\|\Delta\|)$$

$\frac{\partial f}{\partial A}$

$$\log \det(A + \Delta) - \log \det(A) \quad A > 0 \Rightarrow \text{invertible}; \text{ has unique square root } A^{1/2}$$

$$= \log \det(A^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) A^{1/2}) - \log \det(A)$$

$$= \log \frac{|A|^{1/2} |I + A^{-1/2} \Delta A^{-1/2}| |A|^{1/2}}{|A|} - \log |A| \quad \text{evalues of } B$$

$$= \log |I + A^{-1/2} \Delta A^{-1/2}| \quad \text{use } \det(B) = \prod_i \lambda_i(B)$$

$$= \sum_i \log \lambda_i(I + A^{-1/2} \Delta A^{-1/2}) \quad \text{use } \lambda_i(I + B) = 1 + \lambda_i(B) \quad B V = \Delta V \quad (I + B)V = (I + \Delta)V$$

$$= \sum_i \log(1 + \lambda_i(A^{-1/2} \Delta A^{-1/2})) \quad \text{use } \log(1 + x) = x + O(x^2) \quad \text{for } |x| < 1$$

$$= \sum_i \lambda_i(A^{-1/2} \Delta A^{-1/2}) + O\left(\frac{\sum_i \lambda_i(A^{-1/2} \Delta A^{-1/2})^2}{\|\Delta\|^2}\right) \quad \text{because } A \text{ is homogeneous ff.}$$

$\lim_{\Delta \rightarrow 0} \text{cat. with respect}$

$$i.e. \quad B V = \Delta V \quad (B^{-1}V = \Delta^{-1}V)$$

$$\begin{aligned}
 & \text{Let } h(\Delta) = \frac{\log \det(A + \Delta A^{-1})}{\|\Delta\|} \\
 & \sup_{\|\Delta\|=1} \text{cst. with respect to } \|\Delta\| \\
 & h(\Delta) = O(\|\Delta\|^2) \Rightarrow \frac{h(\Delta)}{\|\Delta\|} \leq C \frac{\|\Delta\|^2}{\|\Delta\|} = C \|\Delta\|
 \end{aligned}$$

i.e.  $BV = A V$   
 $(\frac{B}{b})V = (\frac{A}{b})V$

$$\begin{aligned}
 & = \text{tr}(A^{-1} \Delta A^{-1}) \rightarrow o(\|\Delta\|) \\
 & = \text{tr}(A^{-1} \Delta) \rightarrow o(\|\Delta\|) \\
 & \quad (\text{recall } A \text{ is symmetric}) \\
 & \Rightarrow \boxed{\frac{d}{dA} \log \det(A) = A^{-1}}
 \end{aligned}$$

see Boyd's book A.4.1 for the above proof

back to Log-Likelihood of Gaussian

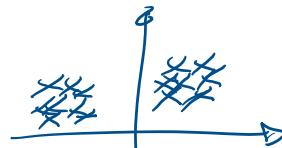
$$+ \frac{n}{2} \log |\Sigma^{-1}| - \frac{n}{2} \langle \Sigma^{-1}, \tilde{\Sigma}(\mu) \rangle \quad (\text{concave fct. of } \Sigma = \Sigma^{-1})$$

take derivative w.r.t.  $\Sigma^{-1} = \Sigma$

$$\begin{aligned}
 & \frac{n}{2} \underbrace{(\Sigma^{-1})^{-1}}_{\Sigma} - \frac{n}{2} \tilde{\Sigma}(\mu) \xrightarrow{\text{want}} \Sigma_{MLE} = \tilde{\Sigma}(MLE) \\
 & = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^T \\
 & \quad (\text{the empirical covariance matrix})
 \end{aligned}$$

unsupervised learning

here  $X$  without any label  $\Sigma$



Consider the Gaussian mixture model (GMM)  
 (can be obtained FLD)

$$Y \sim \text{mult}(\pi) \quad \pi \in \Delta_k$$

[extension of FLD to multiple classes]

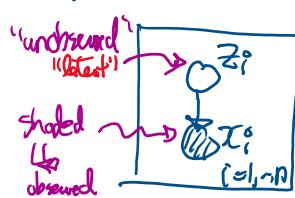
$$X | Y=j \sim N(\mu_j, \Sigma)$$

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

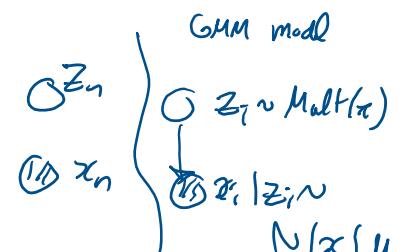
"GMM"

(more generally, can have  $\Sigma_j$  per class)

graphical model for this "latent variable model"



"latent" = repetition



observed  $\xrightarrow{\text{latent}}$

$x_1 \quad x_2$

$x_n$  }  $\mathcal{N}(x | \mu_{z_i}, \sigma_{z_i}^2)$