

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

generative model for classification: (Fisher) linear discriminant analysis

FLD (instead of LDA)

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

$$X \in \mathbb{R}^d \quad \text{class condition}$$

$$\text{generative approach: } p(x, y; \theta) = \overbrace{p(x|y; \theta)}^{\text{shared across classes}} p(y; \theta)$$

vs.

$$\text{conditional approach: } p(y|x; \theta)$$

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

$$\theta = (\underbrace{\mu_0, \mu_1}_{\text{mean for class } 0}, \underbrace{\Sigma}_{\text{shared covariance}}, \pi)_{\text{prob=1}}$$

as before (see exponential family argument)

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

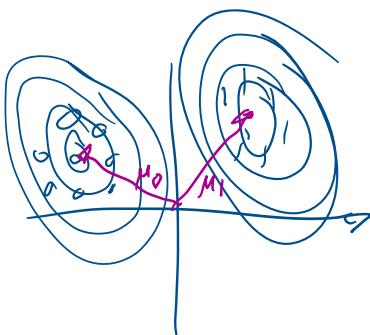
[note: if use $\Sigma \neq \Sigma_1$, get "quadratic discriminant analysis" (QDA)]

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

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

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



sidelong: MLE for multivariate Gaussian

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

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

- dim

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

- T

-

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

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

Σ is symmetric

$$\Sigma \geq 0$$

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

$$\Sigma^T = \Sigma$$

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

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

$$\text{tr}(\Sigma^{-1}(x-\mu)(x-\mu)^T) = \langle \Sigma^{-1}, (x-\mu)(x-\mu)^T \rangle$$

$$\langle A, B \rangle \triangleq \sum_{i,j} A_{ij} B_{ij} = \text{tr}(A^T B)$$

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

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

$$\triangleq \Sigma(\mu)$$

Vector derivative review:

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

$$h \in o(\|\Delta\|) \quad \text{"little oh"} \quad \Leftrightarrow \lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} = 0 \quad \text{e.g. } \|\Delta\|^2 \in o(\|\Delta\|)$$

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

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

"derivative"

$$\lim_{\|\Delta\| \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\|\Delta\|} = \lim_{\|\Delta\| \rightarrow 0} \left(\underbrace{df_{x_0}(\Delta)}_{\text{linear}} + \underbrace{o(\|\Delta\|)}_{\text{small}} \right)$$

$$df_{x_0}(\frac{\Delta}{\|\Delta\|})$$

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

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

$$Df(x_0) = (df_{x_0})^T$$

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 a way to get df_{x_0} for "anything" (matrix, tensor, 10-dim fct., etc.)
 n-dimensional

- 1) This gives a way to get df_{x_0} for "anything" (matrix, tensor, 10-dim fct., etc.)
 2) be careful with dimensions

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

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

Chain rule:

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

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

$\stackrel{\uparrow}{g(f(x_0))} \cdot \stackrel{\downarrow}{\text{matrix product of Jacobians}}$

e.g. $f(\mu) = x - \mu$

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

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

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

$$dg_{v_0} = v^T (A + A^T)$$

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

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

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

$$\nabla_{\mu} \left(\frac{-1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_i x_i$$

15h29

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

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

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

$$\begin{aligned} \log \det(A + \Delta) &\sim \log \det(A) \quad A > 0 \Rightarrow \text{invertible; has unique square root } A^{1/2} \\ &= \log \det(A^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) A^{1/2}) - \log \det(A) \\ &= \log \underbrace{|A|^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) |A|^{1/2}}_{\text{values of } B} - \log \det(A) \\ &= \log |I + A^{-1/2} \Delta A^{-1/2}| \quad \text{use } \det(B) = \prod_i \lambda_i(B) \quad Bv = \lambda v \\ &= \sum \log \lambda_i(I + A^{-1/2} \Delta A^{-1/2}) \quad \text{use } \lambda(I + B) = 1 + \lambda(B) \quad (I + B)v = (1 + \lambda)v \end{aligned}$$

$$\begin{aligned}
&= \log \lambda_i^*(I + A^{-1/2} \Delta A^{-1/2}) \quad \text{use } \lambda^*(I + B) = 1 + \lambda(B) \quad DV = \lambda V \\
&= \log(1 + \lambda_i^*(A^{-1/2} \Delta A^{-1/2})) \quad \log(1+x) = x + O(x^2) \text{ for } |x| < 1 \quad (I+B)V = (1+\lambda)V \\
&= \sum_i \lambda_i^*(A^{-1/2} \Delta A^{-1/2}) + O\left(\left[\sum_i \lambda_i^*(A^{-1/2} \Delta A^{-1/2})\right]^2\right) \quad \text{because } \lambda^* \text{ is homogenous fct.} \\
&\qquad \qquad \qquad \sup_{\|\Delta\|=1} \frac{\lambda_i^*(A^{-1/2} \Delta A^{-1/2})^2}{\|\Delta\|} \quad \text{i.e. } BV = \lambda V \\
&\qquad \qquad \qquad \text{const. with respect to } \|\Delta\| \quad \left(\frac{b}{a}\right)v = \left(\frac{a}{b}\right)v \\
&= \text{tr}(A^{-1/2} \Delta A^{-1/2}) + O(\|\Delta\|) \\
&= \underbrace{\text{tr}(A^{-1} \Delta)}_{\langle A^{-1}, \Delta \rangle} + o(\|\Delta\|) \quad \Rightarrow \boxed{\frac{\partial}{\partial A} \log \det(A) = A^{-1}} \\
&\quad (\text{recall } A \text{ is symmetric})
\end{aligned}$$

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

back to log-likelihood of Gaussian

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

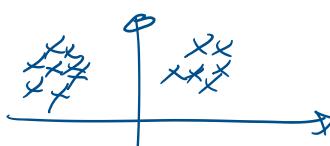
take derivative w.r.t.

$$\begin{aligned}
\Sigma^{-1} = \Lambda \quad & \frac{n}{2} \underbrace{\langle \Sigma^{-1} \rangle^{-1}}_{\tilde{\Sigma}} - \frac{n}{2} \tilde{\Sigma}(\mu) \stackrel{\text{want}}{=} 0 \\
& \Rightarrow \boxed{\begin{aligned} \hat{\Sigma}_{MLE} &= \tilde{\Sigma}(\mu_{MLE}) \\ &= \left[\frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^T \right] \end{aligned}}
\end{aligned}$$

(the empirical covariance matrix)

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]

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

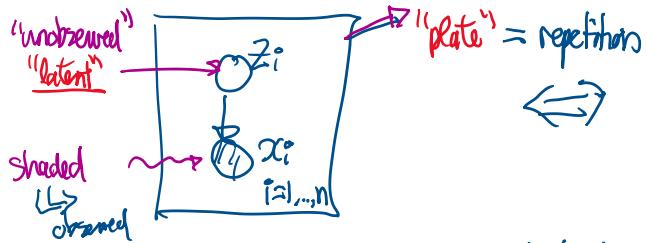
$$X|Y=j \sim N(\mu_j, \Sigma)$$

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

"GMM"

(more generally, can have Σ_j per class)

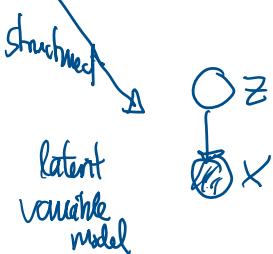
graphical model for this "latent variable model"



GMM model:

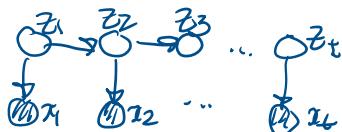
$$\begin{cases} z_i \sim \text{Mult}(\pi) \\ x_i | z_i \sim N(x_i | \mu_{z_i}, \Sigma_{z_i}) \end{cases}$$

two views on $p(x)$

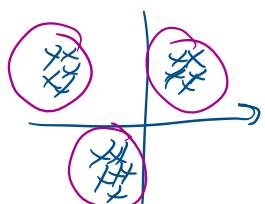


$$p(x) = \sum_z p_\theta(x|z) p_\theta(z)$$

(later in class, we will add time structure HMM)



k-means → to do clustering i.e. group data



we want to get a cluster assignment for every data point x_i

represent $z_{ij}=1$ to mean that x_i belongs to cluster j

$$j=1, \dots, k$$

of clusters (specified in advance for k-means)

applications: • vector quantization (compression)

• in computer vision : use k-mean to get

"bag of visual words" representation
of image patches

- many many others?