

today :
• Fisher LDA

- Math. tricks & MLE for Gaussian

generative model for classification : Fisher's linear discriminant analysis

FLD (instead of LDA)

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

class conditional

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

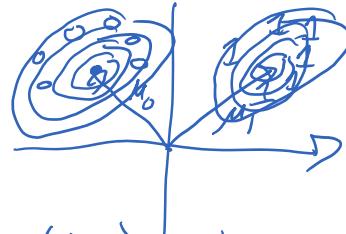
vs.

conditional approach $p(y|x; \theta)$ shared 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 shared $p(y=1)$



as before, could show that $p(y|x; \theta) = \sigma(w^T x)$ where w is a function of $(\mu_0, \mu_1, \Sigma, \pi)$

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

i.e. $\sigma(w^T \phi(x))$ where $\phi(x)$ is a quadratic function of x
 \rightarrow see hawk 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. $\sum_i \log p(y_i|x_i; \theta)$]

[or logistic regression]

Side note: MLE for multivariate Gaussian

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

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

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

Σ is symmetric
 $\Sigma > 0$

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

$$\Sigma^T = \Sigma = \Sigma$$

... \rightarrow $|$ $,$ \dots \rightarrow \rightarrow \rightarrow

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

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

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

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

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

vector derivative review:

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

"little Δ " means that for $\Delta \in \mathbb{R}^m$ s.t. $\lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} = 0$

f is differentiable at x_0 iff \exists a linear operator $df_{x_0}: \mathbb{R}^m \rightarrow \mathbb{R}^n$
 s.t. $\forall \Delta \in \mathbb{R}^m$ $f(x_0 + \Delta) - f(x_0) \approx df_{x_0}(\Delta) + O(\|\Delta\|)$
 "differential"

"derivative"

df_{x_0} is linear

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

$$\begin{aligned} \lim_{\|\Delta\| \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\|\Delta\|} &= \lim_{\|\Delta\| \rightarrow 0} \left(\underbrace{df_{x_0}(\Delta)}_{df_{x_0}(\Delta)} + \underbrace{o(\|\Delta\|)}_{\|\Delta\|} \right) \\ \hat{d} = \frac{\Delta}{\|\Delta\|} &\quad \text{can represent as a } n \times m \text{ matrix} \end{aligned}$$

$$\begin{aligned} \text{standard representation} \quad (df_{x_0})_{i,j} &\approx \frac{\partial f_i}{\partial x_j} \\ \text{then } df_{x_0}(\Delta) &= df_{x_0} \cdot \Delta \end{aligned}$$

• if $n=1$:

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

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

1) this gives a way to get df_{x_0} for "anything"
 (matrix, tensor, \otimes -dim \mathbb{R}^n)

2) be careful with dimensions

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

df_{x_0} is a row vector $(1 \times m)$

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

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

$g: \mathbb{R}^n \rightarrow \mathbb{R}^q$

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

$g(f(x)) = \begin{pmatrix} \quad & \end{pmatrix} \begin{pmatrix} \quad & \end{pmatrix}$
matrix product of Jacobians

e.g. $f(\mu) = x - \mu \quad df_{\mu_0} = -I$

$$g(w) = w^T A w \quad dg_w = w_0^T (A + A^T)$$

$$g \circ f(\mu) = (x - \mu)^T A (x - \mu) \quad d(g \circ f)_{\mu_0} = dg_{g(\mu_0)} \circ df_{\mu_0} \\ = (x - \mu)^T (A + A^T) (-I)$$

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

16h05 $\frac{\partial \mu}{\partial \mu} + \frac{1}{2} \sum_i 2 \Sigma^{-1} (x_i - \mu) \stackrel{\text{want}}{=} 0 \quad \Rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_i x_i$
 $(dg \circ f)_{\mu}^T \quad \text{by } x_i \in \mathbb{R}^q$

example 2: derivative of $f(A) \triangleq \log \det(A)$ where assume A is symmetric
can represent the derivative of a fn. from matrix to scalar, as a matrix

$$f(A + \Delta) - f(A) = \text{tr}(f'(A)^T \Delta) + o(\|\Delta\|) \\ = \underbrace{\langle f'(A), \Delta \rangle}_{df_A} + o(\|\Delta\|)$$

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

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

$$= \log |A|^{1/2} |\underbrace{I + A^{-1/2} \Delta A^{-1/2}}_{\text{e-values of } A} |^{1/2} - \log |A|$$

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

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

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

$$\log(1+x) = x + O(x^2) \quad \text{for } |x| < 1$$

$$\lambda(I + A) = I + \lambda(A)$$

$$(I + A)v = (I + \lambda(A))v$$

$$\begin{aligned}
&= \sum_i \log(1 + \lambda_i(A^{-1/2} \Delta A^{-1/2})) \\
&\approx \sum_i \lambda_i(A^{-1/2} \Delta A^{-1/2}) + O(\lambda_i(A^{-1/2} \Delta A^{-1/2})^2) \\
&\quad = O(\|A\|^2) \\
&\quad + O(\|A\|) \\
&\quad \leftarrow \text{tr}(A) = \sum_i \lambda_i(A) \\
&= \text{tr}(A^{-1/2} \Delta A^{-1/2}) + O(\|\Delta\|) \\
&= \underbrace{\text{tr}(A^{-1} \Delta)}_{\langle A^{-1}, \Delta \rangle} + O(\|\Delta\|) \\
&\quad \Rightarrow \boxed{\frac{\partial}{\partial A} \log \det(A) = A^{-1}}
\end{aligned}$$

$\log(1+x) = x + O(x^2) \text{ for } |x| < 1$
 $\lambda_i(A^{-1/2} \Delta A^{-1/2}) = O(\|A\|)$
 $Av = Av$
 $\frac{Av}{D} = \frac{A}{D}v$

(recall
 A is symmetric)

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

back to log-likelihood of Gaussians:

$$+\frac{n}{2} \log |\Sigma^{-1}| - \frac{n}{2} \langle \Sigma^{-1}, \tilde{\Sigma}(\mu) \rangle$$

(concave fct. of $\Sigma^{-1} = \tilde{\Sigma}$)

take the derivative w.r.t.

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

$$-\frac{n}{2} \tilde{\Sigma}(\mu) \stackrel{\text{want}}{=} 0$$

(the empirical covariance matrix)

$$\begin{aligned}
 \hat{\Sigma}_{MLE} &= \tilde{\Sigma}(\mu_{MLE}) \\
 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^T
 \end{aligned}$$