

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

note for hwk 2:

Newton's $w_{t+1} = w_t - \underbrace{H^{-1}(w_t)}_{\triangleq d_t} \nabla f(w_t)$

instead solve for $H_t d_t = \nabla f(w_t)$ for d_t

$\min_d \|H_t d - \nabla f(w_t)\|$

use `numpy.linalg.lstsq`

instead of checking $x=y$, $|x-y| < \delta$

$a-b \quad a \approx b$
↑
|digits of accuracy|

$a^2 + bx + c = 0$
if b is huge $b^2 \gg 4ac$
cancellation error

instead $\nabla \frac{-b + \sqrt{b^2 - 4ac}}{2a}$
 $\frac{-b + \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b - \sqrt{b^2 - 4ac}}{-b - \sqrt{b^2 - 4ac}}$
 $= \frac{b^2 - (b^2 - 4ac)}{2a(-b - \sqrt{b^2 - 4ac})} = \frac{4ac}{-2a(b + \sqrt{b^2 - 4ac})}$
 $= \frac{-2c}{b + \sqrt{b^2 - 4ac}}$

generative model for classification (Fisher) linear discriminant analysis

FLD (instead of LDA)

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

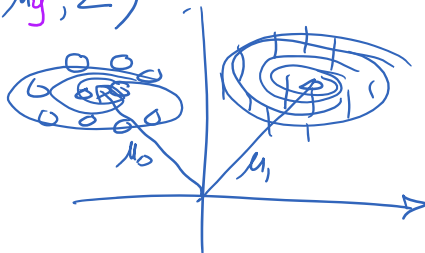
generative approach $p(x, y; \theta) = \overbrace{p(x|y; \theta)}^{\text{class conditional}} p(y; \theta)$

vs.

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

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

$\theta = (\underbrace{\mu_0}_{\text{mean of class 0}}, \underbrace{\mu_1}_{\text{shared } (y=1)}, \underbrace{\Sigma}_{\text{shared } (y=1)}, \underbrace{\pi_0}_{\text{shared } (y=1)})$



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

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

(QDA)

i.e. $\sigma(w^T \phi(x))$ where $\phi(x)$ is quadratic fct. of x
→ see hwk 2

⊗ gen. approach: do joint MLE to estimate

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \sum_i \log p(x_i, y_i; \theta) \quad [\text{vs. } \sum_i \log p(y_i | x_i) \text{ for log. reg.}]$$

side note: MLE for multivariate Gaussian

$x_i \text{ i.i.d. } N(\mu, \Sigma)$ $\mu \in \mathbb{R}^d$ $\Sigma \in \mathbb{R}^{d \times d}$
 Σ is symmetric $\Sigma > 0$
 $E[(x-\mu)(x-\mu)^T] = \Sigma$
 $\Sigma^T = \Sigma$

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

$$\text{tr}\left(\overbrace{(x-\mu)^T \Sigma^{-1} (x-\mu)}^{\text{tr}(AB) = \text{tr}(BA)}\right)$$

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

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

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

using linearity of $\langle \cdot, \cdot \rangle$

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

vector derivative review:

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

"little dh"

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

means that fct. $h(\|\Delta\|)$ s.t. $\lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} = 0$

df_{x_0} is linear

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

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

i th component of f

$$\text{standard representation (JCF)} \dots = \partial f_0$$

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

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

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

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

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

$= \begin{pmatrix} & \end{pmatrix} \begin{pmatrix} & \end{pmatrix}$
 matrix product of Jacobians

2:41pm -> back at 2:51pm!

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

$df_{x_0} = I$

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

$dg_{x_0} = x^T (A + A^T)$

$g \circ f(x) = (x - \mu)^T A (x - \mu)$ $d(g \circ f)_{x_0} = dg_{f(x_0)} df_{x_0}$
 $= (x - \mu)^T (A + A^T) \cdot I$

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

$\nabla_{\mu} \frac{1}{2} \sum_i 2 \Sigma^{-1} (x_i - \mu) \stackrel{\text{want}}{=} 0 \Rightarrow \mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$

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 a scalar, as a matrix

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

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

$\log \det(A + \Delta) - \log \det(A)$

$A > 0 \Rightarrow$ invertible & has unique square root $A^{1/2}$

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

$\log \det(A) = \log |A|$

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

e-values of A

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

use $\det(A) = \prod \lambda_i |A|$

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

use $\det(A) = \prod_i \lambda_i(A)$ (eigenvalues of A)
 $\log(1+x) = x + O(x^2)$ for $|x| < 1$
 other property $\lambda_i(A^{-1/2} \Delta A^{-1/2}) = O(\|\Delta\|)$
 $\text{tr}(A) = \sum_i \lambda_i(A)$

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

back to log-likelihood of Gaussian

$$\begin{aligned}
&+ \frac{1}{2} \log |\Sigma^{-1}| - \frac{n}{2} \langle \Sigma^{-1}, \tilde{\Sigma}(\mu) \rangle \quad (\text{concave function of } -\Lambda = \Sigma^{-1}) \\
&\text{take derivative w.r. to } \Sigma^{-1} = \Lambda \quad \frac{n}{2} (\Sigma^{-1})^{-1} - \frac{n}{2} \tilde{\Sigma}(\mu) \stackrel{\text{want}}{=} 0 \\
&\Rightarrow \boxed{\begin{aligned} \hat{\Sigma}_{MLE} &= \tilde{\Sigma}(\mu_{MLE}) \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})(x_i - \mu_{MLE})^T \end{aligned}} \\
&\text{(the empirical covariance matrix)}
\end{aligned}$$