

REFERENCES

WAINWRIGHT (2019) SECTIONS 1.1, 1.2, 1.3

CANDES, LECTURE NOTES: STATS 800C - THEORY OF STATISTICS

HIGH-DIMENSIONAL STATISTICS

CLASSICAL THEORY

- CONCERNS THE BEHAVIOUR WHEN SAMPLE SIZE $n \rightarrow \infty$
- Y_1, \dots, Y_n IID Y WITH $\mu = \mathbb{E}(Y)$ AND $\Sigma = \text{Var}(Y)$ FINITE.
- LAW OF LARGE NUMBERS
SAMPLE MEAN $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{P} \mu$
- CLT $\sqrt{n}(\hat{\mu}_n - \mu) \rightarrow N(0, \Sigma)$
- CONSISTENCY OF MLE
-

SUPPOSE DATA $n = 1000$ SAMPLE SIZE
 $m = 500$ DIMENSION

QUESTION: IS CLASSICAL THEORY (REQUIRES $n \rightarrow \infty$ n FIXED) PROVIDING
USEFUL INFORMATION?

HD DATA $m \gg n$
CLASSICAL "LARGE n , FIXED m " FAILS IN HD
STATISTICAL METHODS BREAK DOWN IN HD.

CLASSIFICATION PROBLEM

- $Y \in \{A, B\}$ CLASSES

- DETERMINE WHETHER

$$x = (x_1, \dots, x_m)' \in \mathbb{R}^m$$

FROM

$$\begin{cases} A: f_A(x) = P(X | Y = A) \\ B: f_B(x) = P(X | Y = B) \end{cases}$$

DENSITY FUNCTION
OF X FROM CLASS A

$$b: f_B(x) = P(X|Y=B)$$

Prior Prob.

$$\pi_A = P(Y=A)$$

$$\pi_B = P(Y=B)$$

Posterior Prob

$$\frac{P(Y=B|X=x)}{P(Y=A|X=x)} = \frac{\pi_B f_B(x)}{\pi_A f_A(x)}$$

LR

ALLOCATE TO THE CLASS WITH HIGHER POSTERIOR PROB.

$$X_A \equiv X|Y=A \sim N(\mu_A, I_m)$$

$$X_B \equiv X|Y=B \sim N(\mu_B, I_m)$$

$$\pi_A = \pi_B = 1/2$$

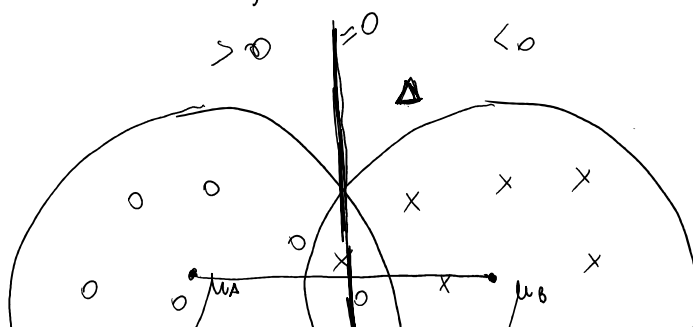
$$\log \frac{f_B(x)}{f_A(x)} = \Psi(x) = \alpha_0 + \alpha^T x$$

OPTIMAL DECISION IS TO THRESHOLD THE LOG-LIK. RATIO

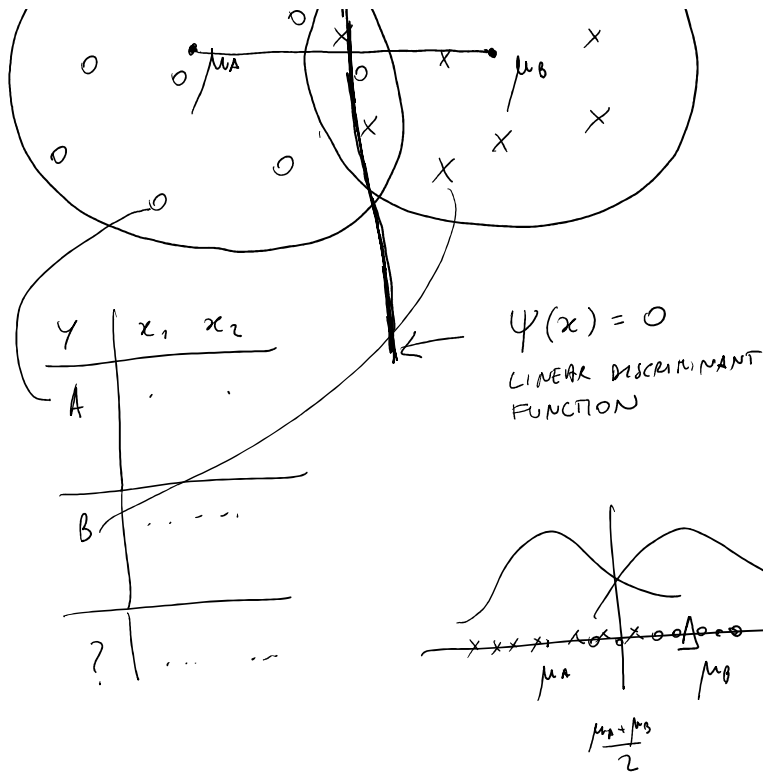
$$\Psi(x) = \langle \mu_A - \mu_B, (x - \frac{\mu_A + \mu_B}{2}) \rangle$$

WHERE $\langle x, z \rangle = x^T z = \sum_{j=1}^m x_j z_j$ EUCLIDEAN INNER PRODUCT.

$$\text{CLASSIFY} \begin{cases} A & \text{IF } \Psi(x) > 0 \\ B & \text{IF } \Psi(x) \leq 0 \end{cases}$$



x OBS FROM B
o OBS FROM A



ERROR PROBABILITY OF THE OPTIMAL RULE

$$\begin{aligned}
 \text{Err}(\Psi) &= \frac{1}{2} P(\Psi(X_A) < 0) + \frac{1}{2} P(\Psi(X_B) \geq 0) \\
 &= \Phi\left(-\frac{\gamma}{2}\right) \quad \Phi \text{ CDF of } N(0,1)
 \end{aligned}$$

WHERE

$$\gamma = \|\mu_A - \mu_B\|_2$$

$$\|\mu\|_2 = \sqrt{\mu^T \mu}$$

EUCL. NORM

FISHER'S LDA.

PLUG-IN PRINCIPLE BASED ON

$\begin{cases} m_A \text{ OBS FROM A} \\ m_B \text{ OBS FROM B} \end{cases}$

$$\hat{\Psi}(x) = \left\langle \hat{\mu}_A - \hat{\mu}_B, x - \frac{\hat{\mu}_A + \hat{\mu}_B}{2} \right\rangle$$

ERROR PROB. OF FISHER'S LDA

$$\text{Err}(\hat{\Psi}(x)) = \frac{1}{2} P(\hat{\Psi}(X_A) < 0) + \frac{1}{2} P(\hat{\Psi}(X_B) \geq 0)$$

IS A RANDOM
VARIABLE

CLASSICAL THEORY

$$\begin{aligned} (n_A, n_B) \rightarrow \infty \quad \Rightarrow \quad \begin{aligned} \hat{\mu}_A &\rightarrow \mu_A \\ \hat{\mu}_B &\rightarrow \mu_B \end{aligned} \quad \begin{aligned} E n(\hat{\psi}) &\rightarrow E n(\psi) \\ &= \Phi(-\delta/2) \end{aligned} \\ m \text{ FIXED} \end{aligned}$$

HIGH-DIM THEORY

$$(n_A, n_B, m) \rightarrow \infty$$

$$\frac{m}{n_A} \rightarrow \delta \geq 0$$

$$\frac{m}{n_B} \rightarrow \delta$$

$$\|\mu_A - \mu_B\|_2 \rightarrow \gamma > 0$$

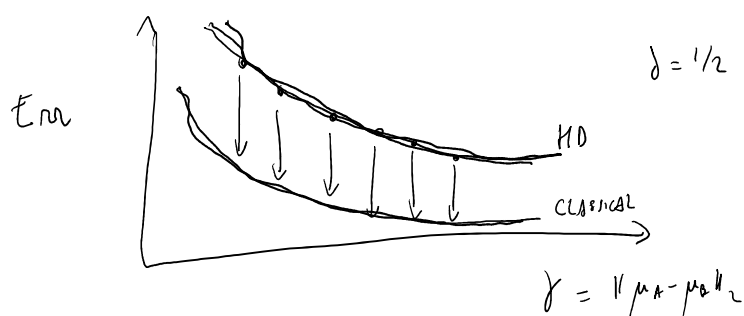
KOLMOGOROV (1960)

$$E n(\hat{\psi}) \xrightarrow{P} \Phi\left(-\frac{\gamma^2}{2\sqrt{\gamma^2 + 2\delta}}\right)$$

$$\xrightarrow{\text{WHEN } \delta=0} \Phi(-\gamma/2)$$

EXAMPLE

$$(m, n_A, n_B) = (400, 800, 800)$$



WHAT CAN HELP US IN HIGH-DIMENSIONS?

HOPE THAT THE "EFFECTIVE" DIMENSION OF THE PROBLEM IS LOW-DIMENSIONAL

SPARSITY

$\mu = (\underbrace{\mu_1, \dots, \mu_n}_{\neq 0}, \underbrace{\mu_{n+1}, \dots, \mu_m}_{=0})^T$ IS SPARSE WHEN ONLY n ENTRIES ARE $\neq 0$

ESTIMATOR
THRESHOLDED
MEAN

$$\tilde{\mu}_j = \hat{\mu}_j \mathbb{1}_{\{|\hat{\mu}_j| > \lambda\}}$$

$$= \begin{cases} \hat{\mu}_j & \text{IF } |\hat{\mu}_j| > \lambda \\ 0 & \text{o/w} \end{cases}$$

$$\lambda = \sqrt{\frac{2 \log m}{n}}$$

$$\tilde{\Psi}(x) = \langle \tilde{\mu}_A - \tilde{\mu}_B, x - \frac{\tilde{\mu}_A + \tilde{\mu}_B}{2} \rangle$$

$$\text{Err}(\tilde{\Psi}(x)) \xrightarrow{P} \text{Err}(\Psi) \quad \text{IF} \quad \log\left(\frac{m}{n}\right)/n \rightarrow 0$$

TESTING THE MEAN VECTOR

$$\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_m \end{pmatrix}, \Sigma\right)$$

\uparrow \uparrow
 y μ

PARAMETER OF INTEREST: $\mu = \mathbb{E}(y)$

NOISANCE PAR. : $\Sigma = \text{Var}(y)$

$\mu_j = 0$ NO EFFECT
 $\mu_j \neq 0$ EFFECT

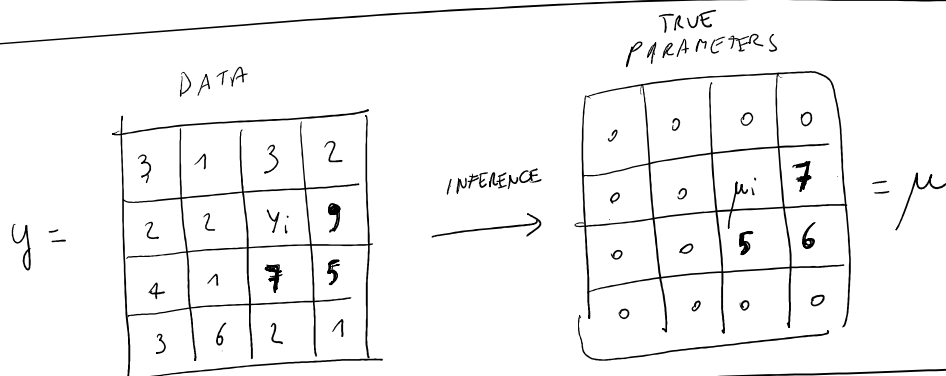
QUESTIONS

- 1) DETECT IF THERE IS AT LEAST ONE EFFECT : AT LEAST ONE $\mu_j \neq 0$?
- 2) COUNTING THE NUMBER OF EFFECTS : HOW MANY $\mu_j \neq 0$?
- 3) IDENTIFY EFFECTS : WHICH $\mu_j \neq 0$?

GLOBAL NULL HYPOTHESIS

$$H_0 : \mu = 0 \Leftrightarrow \bigcap_{j=1}^m \{ \mu_j = 0 \}$$

$$H_1 : \mu \neq 0 \Leftrightarrow \bigcup_{j=1}^m \{ \mu_j \neq 0 \}$$



ASSUMPTIONS

- $\Sigma = I_m$

- $m = 1$

- $y \sim N(\mu, I_m)$

$H_0 : \mu = 0$

$H_1 : \bigcup_{j=1}^m \{ \mu_j > 0 \}$

$\mu \neq 0$

$$- T_j = y_j \sim N(\mu_j, 1) \quad \forall j = 1, \dots, m$$

$$(T_1, \dots, T_m)' \stackrel{H_0}{\sim} N(0, I_m)$$

$$- \text{MAX } T \quad T_{\max} = \max(T_1, \dots, T_m)$$

$$- \text{SUM } T \quad T_{\text{sum}} = \sum_{j=1}^m T_j$$

MAX T

$$P_0(T_{\max} > t_{1-\alpha}) = \alpha$$

$t_{1-\alpha}$ IS $(1-\alpha)$ QUANTILE OF THE DISTRIBUTION OF THE MAXIMUM OF m INDEPENDENT $N(0, 1)$

$$\int_{t_{1-\alpha}}^{\infty} m \phi(y) \bar{\Phi}(y)^{m-1} dy = \alpha$$

APPROXIMATION OF THE CRITICAL VALUE

$$P_0\left(T_{\max} \geq z_{1-\frac{\alpha}{m}}\right) = P_0\left(\bigcup_{j=1}^m \left\{T_j \geq z_{1-\frac{\alpha}{m}}\right\}\right)$$

$$\leq \sum_{j=1}^m P_0\left(T_j \geq z_{1-\frac{\alpha}{m}}\right)$$

Boole
INEQ.

$\left(1-\frac{\alpha}{m}\right)$ QUANTILE OF $N(0, 1)$

$$= m \cdot \frac{\alpha}{m} = \alpha$$

$$P_0\left(T_{\max} \geq z_{1-\frac{\alpha}{m}}\right) = 1 - \prod_{j=1}^m P_0\left(T_j < z_{1-\frac{\alpha}{m}}\right)$$

$$= 1 - \left(1 - \frac{\alpha}{m}\right)^m \xrightarrow{m \rightarrow \infty} 1 - e^{-\alpha}$$

$$\alpha = 0.05, \quad 1 - e^{-\alpha} = 0.0487$$

MAGNITUDE OF THE CRITICAL VALUE $z_{1-\frac{\alpha}{m}}$

FOR LARGE m

$$z_{1-\frac{\alpha}{m}} \approx \sqrt{2 \log m} - \frac{\log(2 \log m) + \log 2\pi}{2\sqrt{2 \log m}}$$

$$\underset{m \rightarrow \infty}{\approx} \sqrt{2 \log m} \quad \leftarrow \text{IT DOESN'T DEPEND ON } \alpha$$

NEEDLE IN A HAYSTACK PROBLEM

$$H_0: \mu_j = 0 \quad \forall j$$

$$H_1: \mu_1 = \underbrace{C_m}_{>0}, \mu_2 = \dots = \mu_m = 0$$

WHAT IS THE LIMITING POWER OF MAXT ?

$$\lim_{m \rightarrow \infty} P_1 \left(T_{\text{MAX}} > z_{1-\frac{\alpha}{m}} \right) = ?$$

IT DEPENDS ON

$$\lim_{m \rightarrow \infty} \frac{C_m}{\sqrt{2 \log m}}$$

> 1 POWER 1

< 1 NO POWER

Proof

②

$$C_m > (1+\epsilon)\sqrt{2 \log m}$$

$$\begin{aligned} \lim_{m \rightarrow \infty} P_1 \left(T_{\text{MAX}} > z_{1-\frac{\alpha}{m}} \right) &\geq P_1 \left(T_1 > z_{1-\frac{\alpha}{m}} \right) \\ &= P_1 \left(N(C_m, 1) > \sqrt{2 \log m} \right) \\ &= P_1 \left(N(0, 1) > \underbrace{\sqrt{2 \log m} - C_m}_{\rightarrow -\infty} \right) \\ &= 1 \end{aligned}$$

$$\textcircled{2} \quad c_m < (1-\epsilon) \sqrt{2 \log m}$$

$$\lim_{m \rightarrow \infty} \mathbb{P}_1 \left(T_{\max} > z_{1-\frac{\alpha}{m}} \right) \stackrel{\text{Boole}}{\leq} \mathbb{P}_1 \left(T_1 > z_{1-\frac{\alpha}{m}} \right) + \mathbb{P}_1 \left(\max_{j>1} T_j > z_{1-\frac{\alpha}{m}} \right)$$

$$= \underbrace{\mathbb{P}_1 \left(N(0,1) > \underbrace{\sqrt{2 \log m} - c_m}_{\rightarrow \infty} \right)}_0 + (1 - e^{-\alpha})$$

$$\leq 1 - e^{-\alpha} \approx \alpha$$

SUM T

$$T_{\text{sum}} = \sum_{j=1}^m T_j \sim N \left(\sum_{j=1}^m \mu_j, m \right)$$

$$\frac{T_{\text{sum}}}{\sqrt{m}} \stackrel{H_0}{\sim} N(0, 1), \quad \frac{T_{\text{sum}}}{\sqrt{m}} \stackrel{H_1}{\sim} N(g_m, 1)$$

with

$$g_m = \frac{\sum_{j=1}^m \mu_j}{\sqrt{m}}$$

$\xrightarrow{m \rightarrow \infty} 0$

NO POWER

$\rightarrow \infty$

POWER 1

$$H_0: \mu_j = 0 \quad \forall j$$

$$H_1: \mu_j = c_m > 0 \quad \forall j$$

$$g_m = \sqrt{m} c_m \xrightarrow{m \rightarrow \infty} 0 \quad \text{IF } c_m = \frac{1}{m}$$

$$\xrightarrow{m \rightarrow \infty} \infty \quad \text{IF } c_m = 1$$

Min P

$$- p_j = 1 - \Phi(T_j)$$

- $p_1, \dots, p_m \text{ i.i.d. } U(0,1) \text{ UNDER } H_0$

- $p_{\min} = \min(p_1, \dots, p_m)$

$p_{\min} \stackrel{H_0}{\sim} \text{Beta}(1, m)$

- THE TEST REJECTS IF $p_{\min} \leq \frac{1 - (1 - \alpha)^{\frac{1}{m}}}{1}$

$$P_0 \left(p_{\min} \leq 1 - (1 - \alpha)^{\frac{1}{m}} \right) = 1 - P_0 \left(\bigcap_{j=1}^m \left\{ p_j > 1 - (1 - \alpha)^{\frac{1}{m}} \right\} \right)$$

$$= 1 - \left[(1 - \alpha)^{\frac{1}{m}} \right]^m = \alpha$$

$$p_{\min} \leq \frac{\alpha}{m} < 1 - (1 - \alpha)^{\frac{1}{m}}$$

\uparrow APPROXIMATED CRITICAL VALUE \uparrow CRITICAL VALUE

FISHER'S COMBINATION METHOD

$$T_F = \sum_{j=1}^m 2 \log \left(\frac{1}{p_j} \right) \stackrel{H_0}{\sim} \chi_{2m}^2$$

SIMES' TEST

$$p_{(1)} \leq \dots \leq p_{(m)}$$

$$p_{(j)} \stackrel{H_0}{\sim} \text{Beta}(j, m - j + 1)$$

THE SIMES P-VALUE IS

$$p_s = \min_{j=1, \dots, m} \left\{ p_{(j)} \frac{m}{j} \right\} \stackrel{H_0}{\sim} U(0,1)$$

REJECTS H_0 IF $\exists j : p_{(j)} \leq \frac{\alpha j}{m}$