# Lecture 10

*) M.V.G. MLE / Bernoulli & M.V.

*) MAP

MAP estimate.

$$\underset{\Theta}{\max} \left[ \sum_{i=1}^{N} \ln P(x \mid \Theta) \right] + \ln P(\Theta)$$

$$\underset{\Theta}{\max} \left[ \sum_{i=1}^{N} \left\{ x_i \ln \Theta + (1-x_i) \ln(1-\Theta) \right\} + (\alpha-1) \ln \Theta \right.$$
$$\left. + (\beta-1) \ln(1-\Theta) \right.$$

$$P(\Theta) \prod_{i=1}^{N} P(x \mid \Theta)$$

$$P(\Theta) \underbrace{P(D \mid \Theta)}_{\prod_{i=1}^{N}}$$

$$\Theta_{MAP} = \frac{\sum_i x_i + \alpha - 1}{N + \alpha + \beta - 2}$$

$$\Theta_{MLE} = \frac{\sum_i x_i}{N}$$

$N$ iid's. $\quad x_1, x_2 \cdots x_4 \qquad N = 4$

$$AB \subseteq \quad C \ (AB)$$

$$x_1 \overset{H}{=} 1, \quad x_2 \overset{T}{=} 0, \quad x_3 \overset{T}{=} 0, \quad x_4 \overset{T}{=} 0$$

$$\Theta_{MLE} = \frac{1}{4} \qquad\qquad \alpha = 5, \ \beta = 5$$

$$\Theta_{MP} = \frac{1 + \alpha - 1}{4 + \alpha + \beta - 2} \quad = \quad \frac{5}{12}$$

# Demo

- Once you have ML/MAP estimates, you could plug them in likelihood function for a classification setting.
- Take an eg of two category case with 3-d and 4 samples.

$$P(\omega|x) = P(x|\theta)\,P(\omega)$$

$\omega_1$

MLE/MAP

$\mathcal{L}_{ML}$ GA Poet.

|       | $d_1$ | $d_2$ | $d_3$ |
|-------|-------|-------|-------|
| $x_1$ | 1     | 0     | 1     |
| $x_2$ | 0     | 0     | 0     |
| $x_3$ | 0     | 0     | 1     |
| $x_4$ | 1     | 0     | 1     |

$$P(x|\theta) = \theta_1^{x_1}(1-\theta_1)^{1-x_1}\cdot \theta_2(1-\theta_1)^{1-x_2}$$

$$\theta_3^{x_3}(1-\theta_3)^{1-x_3}$$

$$\theta_{1ML}=\frac{1}{2},\quad \theta_{2ML}=0,\quad \theta_{3ML}=3/4$$

$$\omega_2, \quad \begin{array}{ccc} d_1 & d_2 & d_3 \\ x_1 & 1 & 0 & 0 \end{array} \qquad \Theta_1 = \frac{1}{4}, \quad \Theta_2 = \frac{3}{4}, \quad \Theta_3 = \frac{1}{4}$$

$$x_1$$
$$x_2$$
$$x_3$$
$$x_4$$

$$P(x \mid \Theta) = (\Theta_1)^{x_1} (1 - \Theta_1)^{1 - x_1} \cdot \Theta_2 \cdot \Theta_3$$

$$P(\omega_1) = P(\omega_2) = \frac{1}{2}$$

$$x_0 \rightsquigarrow \text{test sample} \rightarrow \{1, 1, 1\} \quad \begin{array}{l} \nearrow \omega_1 ? \\ \searrow \omega_2 ? \end{array}$$

$$P(\omega_1 \mid x) \propto P(x \mid \Theta) \, P(\omega_1) = P(x \mid \omega_1) \, P(\omega_1)$$
$$P(\omega_2 \mid x) \propto P(x \mid \Theta) \, P(\omega_2) = P(x \mid \omega_2) \, P(\omega_2)$$

$$P(\omega_1 \mid x) = P(x \mid \omega_1) \, P(\omega_1)$$

- Consider data matrix X?
- What will be the mean?
- What will be mean of X – mu?
- $Y = a^T X$, mean of Y? var of Y?

$$(Pdf) \sim \Pi$$

$$\text{Gaus.} \sim P(x \mid \mu, \sigma^2) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(x-\mu)^2_{ML}}$$

$$\text{Scalar.} \quad MLE \quad \prod_{i=1}^{N} P(x_i \mid \theta) \rightarrow \theta_{MLE}$$

$\sigma_{MLE}$  Gaussian  $\mu_{MLE}$

$$p(x; \mu; \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e\left\{ -\frac{1}{2\sigma^2}\left(x - \mu_{ML}\right)^2 \right\}$$

→ am I call. Pdf. of r.v. $x$

→ am I call. Pdfs of $N$ iids. they are
realization of $p(x; \mu, \sigma^2)$

$$\Theta_{MAP} = \frac{\sum_{i} x_i + \alpha - 1}{N + \alpha + \beta - 2}$$

$$= \frac{1 + \alpha - 1}{4 + \alpha + \beta - 2}$$

$$
\begin{array}{ccc}
 & d & \\
x_1 & 1 & x_1 = 1 \\
x_2 & 0 & x_2 = 0 \\
x_3 & 0 & x_3 = 1 \\
x_4 & 0 & x_4 = 1 \\
\omega_1 & & \omega_2
\end{array}
$$

$\alpha = \beta = 5$

$x_0 = 1$

$$\Theta_{MAP} = \frac{5}{12}$$

$$
\begin{array}{cc}
\omega_1 & \omega_2 \longrightarrow \Theta_{MAP} = \frac{1}{3} \\
Y & N
\end{array}
$$

$$P(\omega_1 | x) = \frac{P(x | \Theta) \, P(\omega_1)}{x}$$

$$= \Theta_{MAP}^{x} \, (1 - \Theta_{MAP})^{(1-x)} \, P(\omega_1)$$

Data matrix $X \in R^{d \times n}$

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ |      |
|-------|-------|-------|-------|-------|------|
| $d_1$ | 1     | 0     | 0     | 0     | $\frac{1}{4}$ |
| $d_2$ | 1     | 0     | 1     | 0     | $2/4$ |
| $d_3$ | 1     | 1     | 1     | 0     | ~~3/5~~ $3/4$ |

$\mu_X$

$X - \mu_X$

| $3/4$ | $-1/4$ | $-1/4$ | $-1/4$ | 0 |
|-------|--------|--------|--------|---|
| $1/2$ |        |        |        | 0 |
| $1/4$ |        |        |        | 0 |

$$X \in \mathbb{R}^{d \times n} \qquad X = \begin{bmatrix} x_{11} & x_{21} & x_{31} & x_{41} \\ x_{12} & x_{92} & x_{32} & x_{42} \end{bmatrix}_{2 \times 4}$$

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \end{matrix}$$

$$\mu_x = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \qquad Y = a^T X \qquad a \in \mathbb{R}^{2 \times 1}$$

$$\mu_y = \begin{bmatrix} a_1 & a_2 \end{bmatrix} \begin{bmatrix} x_{11} & \cdots & x_{41} \\ x_{12} & & x_{42} \end{bmatrix}$$

$$= a_1 \mu_1 + a_2 \mu_2$$

$$= a^T \mu_x$$

# Biased Estimate

$$a \in \mathbb{R}^{d \times 1}$$
$$X = \in \mathbb{R}^{d \times n}$$
$$Y \in \mathbb{R}^{1 \times n}$$

*) $\quad Y = a^T X$

$$var(Y) = \frac{1}{n-1} \sum_{i=1}^{n} \left( y_i - \mu_y \right)^2$$

$$= \frac{1}{n-1} \sum_{i=1}^{n} \left( a^T x_i - a^T \mu_x \right)^2$$

$$= \frac{1}{n-1} \sum_i \underbrace{a^T (x_i - \mu_x)} (x_i - \mu_x)^T a$$

$$= a^T \left[ \frac{1}{n-1} \sum_i (x_i - \mu_x)(x_i - \mu_x)^T \right] a = a^T S_x a$$

$$\text{ML estimate of mean is unbiased} \qquad E(x^2) = \text{var}(x) + (E(x))^2$$

*) is $E(\sigma^2_{ML}) = \sigma^2$ ?

$$E\left( \frac{1}{N} \sum_i (x_i - \mu_{ML})^2 \right)$$

$$\frac{1}{N} E\left( \sum_i x_i^2 - 2 x_i \mu_{ML} + \mu^2_{ML} \right)$$

$$= \frac{1}{N} E\left( \sum_i x_i^2 - 2\mu_{ML} \sum_i x_i + \sum_i \mu^2_{ML} \right)$$

$$= \frac{1}{N} E\left( \sum_i x_i^2 - 2\mu^2_{ML} N + N \mu^2_{ML} \right)$$

$$= \ldots$$

$$\frac{1}{N} E\left( \sum_i x_i^2 \right)$$

$$- E(\mu^2_{ML})$$

$$= \frac{1}{N} \sum_i E(x_i^2) - \left\{ \text{var}(\mu_{ML}) + \mu^2 \right\}$$

$$= \frac{1}{N} \sum_i (\sigma^2 + \mu^2) - \cancel{\text{var}}$$

$$- \text{var}(\mu_{ML}) - \mu^2$$

$$= \sigma^2 + \mu^2 - \text{var}(\mu_{ML}) - \mu^2$$

$$= \sigma^2 - \text{var}(\mu_{ML})$$

$$= \sigma^2 - \text{Var}\left(\frac{1}{N}\sum_i x_i\right)$$

$$= \sigma^2 - \frac{1}{N^2}\text{Var}\left(\sum_i x_i\right) \swarrow$$

$$= \sigma^2 - \frac{1}{N^2}\sum_i \text{Var}(x_i)$$

$$= \sigma^2 - \frac{1}{N}\sigma^2$$

$$= \sigma^2\left(\frac{N-1}{N}\right) \neq \sigma^2$$

$\sigma^2_{MLE}$ is biased estimate.

$$\text{Var}(a+b+c) =$$
$$\text{Var}(a) + \text{Var}(b) + \text{Var}(c)$$
$$\text{if } a, b, c \text{ are ind.}$$

$$\sigma^2_{MLE} = \frac{1}{N-1}\sum_i (x_i - \mu_{ML})^2$$

$$E\left(\sigma^2_{MLE}\right) = \sigma^2$$

■ Bias

– ML estimate for $\sigma^2$ is biased

$$E\left[\frac{1}{n}\Sigma(x_i - \overline{x})^2\right] = \frac{n-1}{n} \cdot \sigma^2 \neq \sigma^2$$

– An elementary unbiased estimator for $\Sigma$ is:

$$C = \frac{1}{n-1}\sum_{k=1}^{k=n}(x_k - \hat{\mu})(x_k - \hat{\mu})^t$$

$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{Sample\ \text{covariance matrix}}$
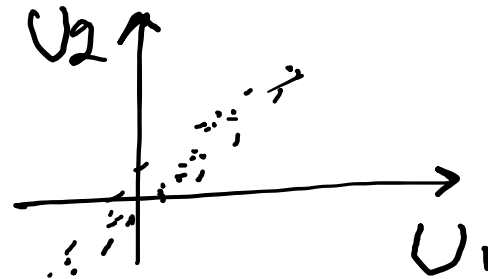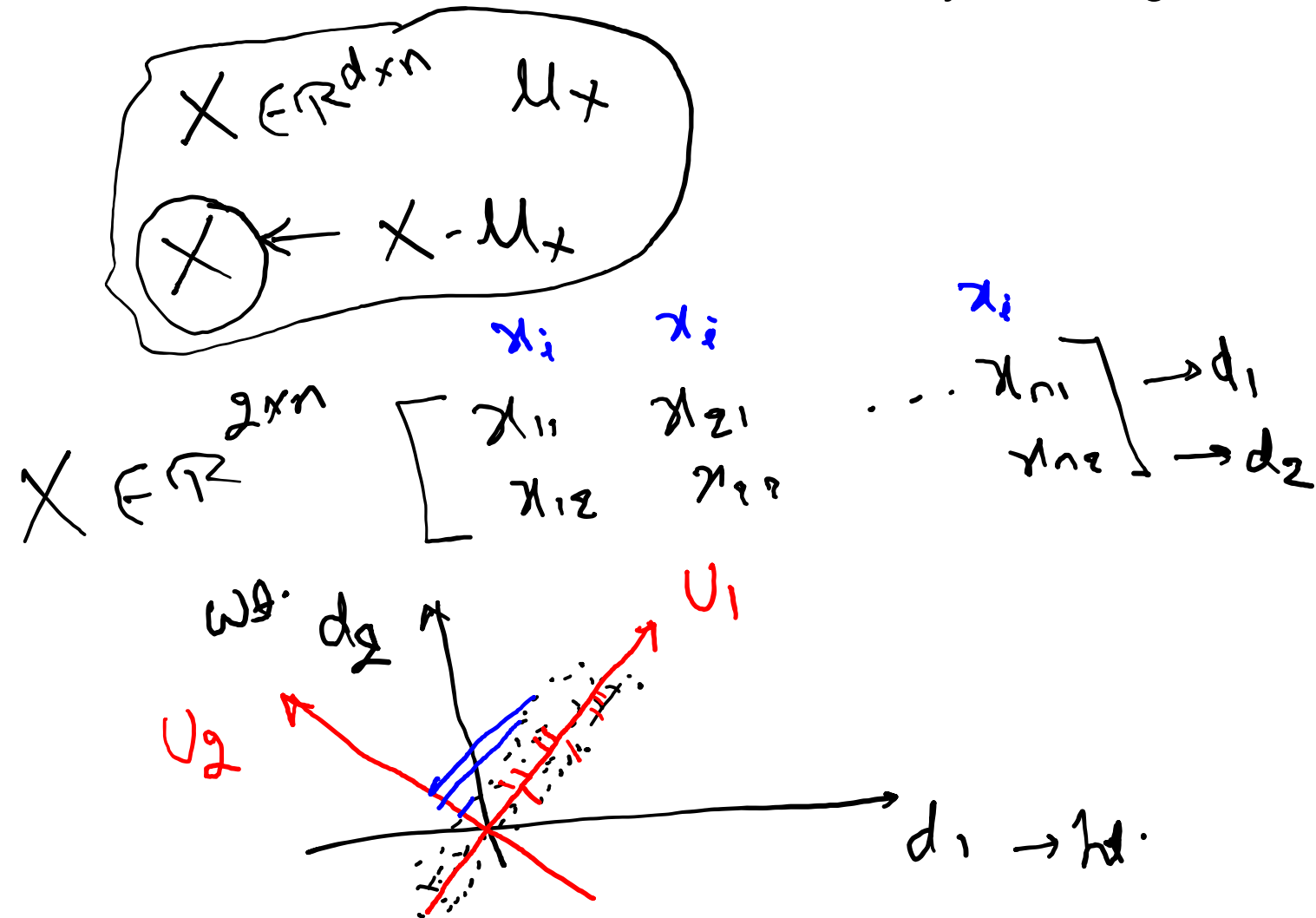
# PCA

→ Principal Component Analysis.

- Start with some data matrix, if you are using a discriminant based on MVG, what all things do you need?

- How many multiplications are needed?

- Can you reduce this complexity?

$$O(Nd^2)$$

- Goal: Project the data onto orthogonal space such that the projected points have max. variance

Consider data matrix X and its centralized version obtained by removing the mean.

$X \in \mathbb{R}^{d \times n} \quad \mu_x$

$X \leftarrow X - \mu_x$

$X \in \mathbb{R}^{2 \times n} \quad \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & & x_{n2} \end{bmatrix} \begin{array}{l} \rightarrow d_1 \\ \rightarrow d_2 \end{array}$

$x_i \quad x_i \quad x_i$

w.t. $d_2$ $U_1$

$U_2$

$d_1 \rightarrow h.t.$

$$Y_i = U_1^T x_i$$

$$y = U_1^T X$$

$$\mu_y = U_1^T \mu_X = 0$$

$$var(y) = U_1^T \delta_X U_1$$

$$\delta_X \rightarrow \frac{XX^T}{n-1}$$

$$\max_{U_1} \ U_1^T \delta_X U_1 \quad \text{s.t.} \ U_1^T U_1 = 1$$

$$\max_{U_1} \ U_1^T \delta_X U_1$$

$$2 \delta_X U_1 = 0$$

$$U_1 = 0$$