

Lecture 9

* MLE \rightarrow You may not have likelihood parameters.

$$P(x_i | \omega_i) \sim N(\mu_i, \sigma_i^2) \quad \mu_i, \sigma_i^2 \text{ are not known.}$$

N iids. x_i 's $\forall i = 1, \dots, N$

likelihood function $F(\theta) = \ln \prod_{i=1}^N P(x_i | \theta)$

$$\nabla_{\theta} F(\theta) = 0$$

$\hat{\theta}_{MLE}$

* $\hat{\mu}_{ML} = \frac{1}{N} \sum_i x_i \quad \hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{ML})^2$

* MLE for multivariate case,

→ Bernoulli distribution.

* MAP / BCA.

Multivariate case Gaussian.

Given: N i.i.d. $x_i \quad i=1, 2, \dots, N \quad x_i \in \mathbb{R}^d$

$$F(\theta) = \ln \prod_{i=1}^N \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{\left\{-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right\}}$$

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

$$= \sum_{i=1}^N \ln p(x_i | \mu, \Sigma)$$

$$= \sum_{i=1}^N -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

μ is unknown & Σ is known.

$$\mu_M = ?$$

$$\frac{\partial}{\partial \mu} F(\mu) = 0$$

$\frac{\partial}{\partial \mu}$

$$-\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu)$$

$$= -\frac{1}{2} \sum_{i=1}^N \mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i - \frac{1}{2} \mu^T \Sigma^{-1} \mathbf{x}_i + \frac{1}{2} \mathbf{x}_i^T \Sigma^{-1} \mu$$

$$= \underbrace{-\frac{1}{2} \sum_{i=1}^N \mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i}_{\text{cancel}} + \underbrace{\frac{1}{2} \mu^T \Sigma^{-1} \mu}_{\text{cancel}}$$

$$= \sum_{i=1}^N (\Sigma^{-1} \mathbf{x}_i)^T \mu - \frac{1}{2} \mu^T \Sigma^{-1} \mu$$

$$\begin{aligned}
 \mu^T \Sigma^{-1} \mathbf{x}_i &= (\mu^T \Sigma^{-1} \mathbf{x}_i)^T \\
 &= \mathbf{x}_i^T (\Sigma^{-1})^T \mu \\
 &= \mathbf{x}_i^T \Sigma^{-1} \mu \\
 &= (\Sigma^{-1} \mathbf{x}_i)^T \mu
 \end{aligned}$$

$$\begin{cases}
 \mu^T \Sigma^{-1} \mathbf{x}_i \\
 = \text{tr}(\mu^T \Sigma^{-1} \mathbf{x}_i)
 \end{cases}$$

$$\frac{\partial (b^T y)}{\partial y} = b \quad b \text{ & } y \text{ are vectors.} \quad \frac{\partial f(g)}{\partial g}$$

$$b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$b^T y = b_1 y_1 + b_2 y_2$$

$$\frac{\partial}{\partial y} b^T y = \begin{bmatrix} \frac{\partial b^T y}{\partial y_1} \\ \frac{\partial b^T y}{\partial y_2} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$\frac{\partial F(\theta)}{\partial \mu} = \sum_{j=1}^N \frac{g^{-1} x_j - \frac{1}{2} g^T \mu}{2} = 2 A M$$

~~$\frac{\partial}{\partial \mu} (\sum_i x_i - \mu) = 0$~~

$$M^T g^{-1} M \boxed{M^T A M}$$

$$M = \begin{bmatrix} M_1 \\ M_2 \end{bmatrix} \quad A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\frac{\partial M^T A M}{\partial \mu} = \frac{\partial (2M_1^2 + M_2^2)}{\partial \mu}$$

$$= \begin{bmatrix} 2 \cdot 2 M_1 \\ 2 \cdot M_2 \end{bmatrix}$$

$$\sum_i x_i = NM$$

$y \Rightarrow d$ -dim vector

$$\mu_{ML} = \frac{1}{N} \sum_i x_i$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}$$

* ~~$E(\mu_{ML}) \rightarrow \text{Biased?}$~~

$\mu_{ML} \rightarrow \text{Biased?}$

$$b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_d \end{bmatrix}$$

$$b^T y$$

$E(\mu_{ML}) = \mu \rightarrow \text{Unbiased.}$

$$E(\mu_{ML}) = \frac{1}{N} \sum_i E(x_i) = \mu$$

$$\mathcal{S}_{ML} = \frac{1}{N} \sum_{i=1}^N \underbrace{(x_i - \mu_{ML})(x_i - \mu_{ML})^T}_{d^2}$$

How many mult. do you need to obtain \mathcal{S}_{ML} ?

$$\mathcal{S}_{ML} = Nd^2$$

$$N = 10^3, d = 10^3$$

- Example of a specific case: unknown μ
 - $P(x_i | \mu) \sim N(\mu, \Sigma)$
 (Samples are drawn from a multivariate normal population)

$$\ln P(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$$

and $\nabla_{\theta\mu} \ln P(x_k | \mu) = \Sigma^{-1} (x_k - \mu)$

$\theta = \mu$ therefore:

- The ML estimate for μ must satisfy:

$$\sum_{k=1}^{k=n} \Sigma^{-1} (x_k - \hat{\mu}) = \mathbf{0}$$

- Multiplying by Σ and rearranging, we obtain:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

Just the arithmetic average of the samples of the training samples!

Conclusion:

If $P(\mathbf{x}_k | \omega_j)$ ($j = 1, 2, \dots, c$) is supposed to be Gaussian in a d -dimensional feature space; then we can estimate the vector $\theta = (\theta_1, \theta_2, \dots, \theta_c)^t$ and perform an optimal classification!

Bernoulli experiment

- Let's conduct N independent trials of the following Bernoulli experiment:
We will ask each student whether they will take SML or ML.
- Let p be the probability that an individual will vote SML. In this example, each observation x_i is a scalar. So it's better to represent it by x_i . For each i, the value of x_i is either SML (1) or ML (0).

$$x_i = \begin{cases} 1 & \text{or SML with } P = \theta \\ 0 & \text{or ML with } 1 - P = 1 - \theta \end{cases}$$

$$P(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$$

$$\hat{\theta}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$f(\theta) = \ln \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i}$$

$$= \sum_{i=1}^N x_i \ln \theta + (1-x_i) \ln(1-\theta)$$

$$= \sum_{i=1}^N \frac{x_i}{\theta} - \frac{(1-x_i)}{1-\theta} = 0 \Rightarrow \hat{\theta}_{MLE} =$$

* Multivariate Bernoulli $\theta_1, \theta_2, \theta_3$

* Student $X_1 = \{x_{1SM1}=1, x_{1GPA}=1, x_{1Prof}=0\}$

Student 2 $X_2 = \{x_{2SM1}=0, x_{2GPA}=0, x_{2Prof}=0\}$

... $X_N = \{x_{NSM1}=0, x_{NGPA}=1, x_{NProf}=1\}$

$$\theta_{ij} = \cancel{P(x(x_{ij}=1|\theta))}$$

$j' \rightarrow$ denoting the dim.

$$\theta_j = P(x_{j'}=1|\theta_j)$$

$$P(x|\mu, \sigma^2) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

μ, Σ

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

$\theta_1 \rightarrow$ dim. to SML

$\theta_2 \rightarrow$ dim to GFA

$\theta_3 \rightarrow$: Poetry

N samples.

$$P_x(x_i|\theta_j) = \theta_j^{x_i} (1-\theta_j)^{1-x_i}$$

$$\theta = \{\theta_1, \theta_2, \theta_3\}$$

$$P_x(x|\theta) = \prod_{j=1}^3 \theta_j^{x_j} (1-\theta_j)^{1-x_j}$$

$$X = \{x_{SML}, x_{GFA}, x_{Poet}\}$$

all dim. decrease ind.

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \leftarrow \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$$

$$\boldsymbol{x}_1 = \begin{bmatrix} x_{1d_1} & x_{1d_2} \end{bmatrix} \leftarrow \boldsymbol{\mu}, \boldsymbol{\Sigma},$$

$$\boldsymbol{x}_2 = \begin{bmatrix} x_{2d_1} & x_{2d_2} \end{bmatrix} \leftarrow \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2$$

$$\vdots$$

$$\boldsymbol{x}_N = \begin{bmatrix} x_{Nd_1} & x_{Nd_2} \end{bmatrix} \leftarrow (\boldsymbol{\mu}_{\text{avg}}, \boldsymbol{\Sigma}_N)$$

* N iids $x_{ij} \quad x_{ij} \in \mathbb{R}^d$

$$\ln \prod_{\delta=1}^N \rho(x_{\delta j} | \theta) = \ln \prod_{\delta=1}^N \prod_{j=1}^d \rho(x_{\delta j} | \theta_j)$$

$$= \ln \prod_{\delta=1}^N \prod_{j=1}^d \theta_j^{x_{\delta j}} (1-\theta_j)^{1-x_{\delta j}}$$

$$F(\theta) = \sum_{\delta=1}^N \sum_{j=1}^d x_{\delta j} \ln \theta_j + (1-x_{\delta j}) \ln (1-\theta_j)$$

$$\frac{\partial}{\partial \theta_j} F(\theta) = \sum_{\delta=1}^N \frac{x_{\delta j}}{\theta_j} - \frac{1-x_{\delta j}}{1-\theta_j} = 0 \Rightarrow \theta_j = \frac{1}{N} \sum_{\delta=1}^N x_{\delta j}$$

$$\Theta_1 = \frac{1}{N} \sum_{j=1}^N x_{j1} = \frac{1}{4}, \quad \Theta_2 = \frac{1}{4}, \quad \Theta_3 = \frac{1}{2}$$

gML, GFA, Poetry.

1 0 1

0 0 0

0 1 0

0 0 1

$$P(x|\theta) = \Theta_1^{x_1} (1-\Theta_1)^{1-x_1} \cdot \Theta_2^{x_2} (1-\Theta_2)^{1-x_2} \Theta_3^{x_3} (1-\Theta_3)^{1-x_3}$$

$$x_1, x_2, x_3 \in \{0, 1\}$$

$$\frac{\partial}{\partial \Theta_j} \sum_{j=1}^N \left\{ x_{j1} \ln \Theta_1 + x_{j2} \ln \Theta_2 + x_{j3} \ln \Theta_3 + x_{j4} \ln \Theta_4 \right\}$$

Example: Bernoulli distribution

If x_1, \dots, x_n are a set of Bernoulli observations,

with $f(1 | p) = p$ and $f(0 | p) = 1 - p$, i.e. $f(x_i | p) = p^{x_i} (1 - p)^{1-x_i}$

The likelihood function for observations x_1, \dots, x_n is

$$L(p | x_1, \dots, x_n) = \prod_{i=1}^n p^{x_i} (1 - p)^{1-x_i} = p^{\sum x_i} (1 - p)^{n - \sum x_i},$$

and the m.l.e \hat{p} is the value that maximizes this.

The log-likelihood is $\ln(L) = \ln(p) \sum x_i + (n - \sum x_i) \ln(1 - p)$

Set: $\frac{d \ln(L)}{dp} = \frac{\sum x_i}{p} - \frac{n - \sum x_i}{1 - p} = 0 \quad \text{solve for } p \Rightarrow \hat{p} = \frac{\sum x_i}{n}$

Drawback of MLE?

$$\hat{\theta}_{ML} = \frac{1}{N} \sum_{j=1}^N x_j = 1 \approx \frac{1}{2}$$

$$x = \{H, T\}$$

$$\hat{\theta} = \hat{\theta}_x(x=H|\theta)$$

$$x_1, x_2, x_3, \dots, x_{100},$$

~~x₁~~ | 1 1 1

Biased Estimate

/ Practice.

$\text{Var}(\mu_{ML}) = ?$ Summation of ind. rvs

Suppose MLE for mean for a dataset is taken as the first sample.
Is the estimator unbiased?

MAP

→ Max. a Posteriori

$$\downarrow f(\theta)$$

- $\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} P(D|\theta)$
- $\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} P(\theta|D) \equiv \operatorname{argmax}_{\theta} \underbrace{P(D|\theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}$

Likelihood prior.

Bernoulli Beta.

Posterior will be Beta distribution.

Likelihood \rightarrow Gaussian, prior Gaussian.
Post. Gaussian.

Conjugate prior - Beta distribution

ωικι

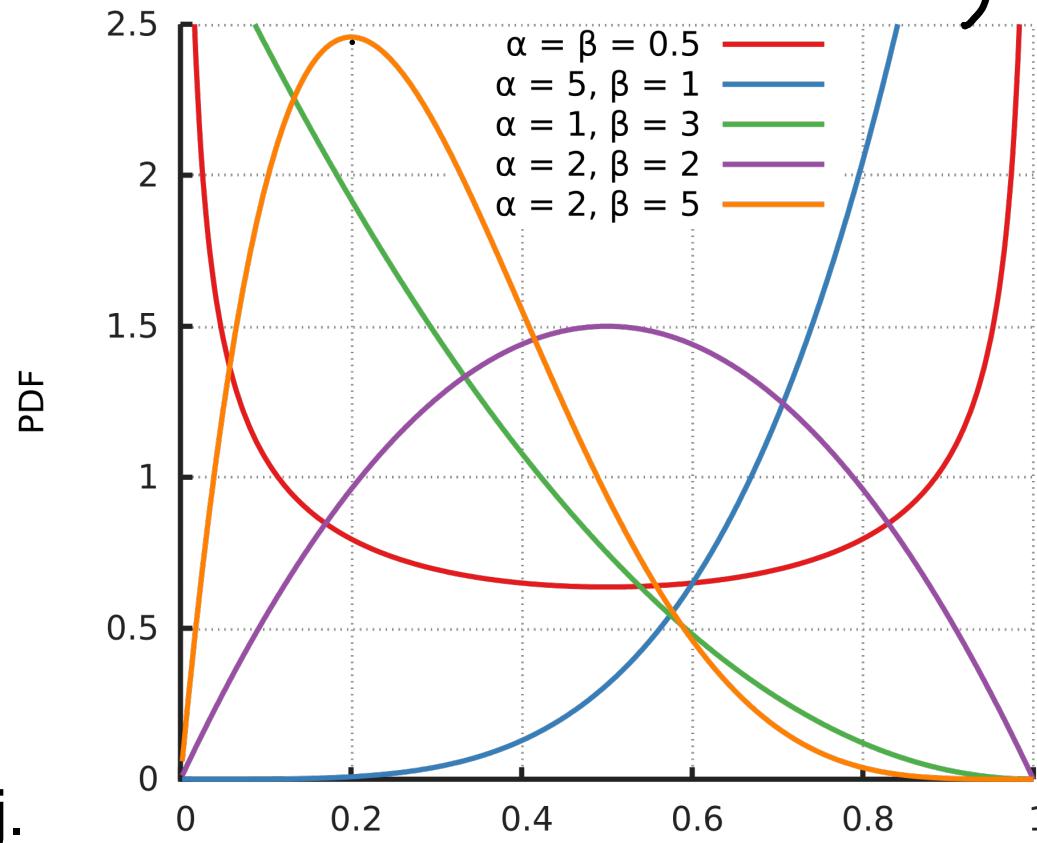
- Conjugate prior of likelihood gives the posterior distribution from the same family of prior

$$f(p; \alpha, \beta) = \frac{p^{\alpha-1}(1-p)^{\beta-1}}{B(\alpha, \beta)}$$

Beta distribution.

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

Posterior and prior are conj.
distributions



MAP estimate:

$$\max_{\theta} \left[\sum_{j=1}^N \ln p(x_j | \theta) \right] + \ln p(\theta)$$

$$\max_{\theta} \left[\sum_{j=1}^N \left[x_j \ln \theta + (1-x_j) \ln (1-\theta) \right] + (\alpha-1) \ln \theta + (\beta-1) \ln (1-\theta) \right]$$

$$P(\theta) \frac{\prod_{j=1}^N p(x_j | \theta)}{P(\theta) P(\theta | \alpha, \beta)}$$