


Lect 2: Gaussian Distribution and 3.



Gaussian Distribution: Extremely important.

Univariate Gaussian:

Toss a fair coin (head/tail $P=0.5$)

Define the random variable

$$X = \begin{cases} 1 & \text{if head} \\ 0 & \text{if tail.} \end{cases}$$

We will toss the coin N times
and count the # of heads, m .

calculate m/N .

The probability density function is

$$N(x; \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

↑ ↑
mean variance.

$$\theta = (\mu, \sigma^2)$$

In MLE:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \log P(\mathcal{D}|\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{\sum_{i=1}^N \log P(x^{(i)} | \theta)}$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(x^{(i)} - \mu)^2}{2\sigma^2}\right)$$

① Find μ_{MLE} first.

$$\underset{\mu}{\operatorname{argmax}} - \sum_{i=1}^N (x^{(i)} - \mu)^2$$

$$\frac{\partial}{\partial \mu} \sum_{i=1}^N (x^{(i)} - \mu)^2 = 0 \Rightarrow \mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

② Find σ_{MLE} .

$$\operatorname{argmax} \sum_{i=1}^N \left[-\log \sigma - \frac{1}{2\sigma^2} (x^{(i)} - \mu)^2 \right]$$

$$\frac{\partial}{\partial \sigma} \left[\sum_{i=1}^N \log \sigma + \frac{1}{2\sigma^2} (x^{(i)} - \mu)^2 \right] = 0$$

$$\Rightarrow \sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)^2$$

③ μ_{MLE} is unbiased

$$\begin{aligned} E_D[\mu_{MLE}] &= E_D\left[\frac{1}{N} \sum_{i=1}^N x^{(i)}\right] \\ &= \frac{1}{N} \sum_{i=1}^N E_D[x^{(i)}] = \mu. \end{aligned}$$

④ σ_{MLE} is biased. if we replace μ by μ_{MLE} .

$$\begin{aligned} E_D[\sigma_{MLE}^2] &= E_D\left[\frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu_{MLE})^2\right] \\ &= E_D\left[\frac{1}{N} \sum_{i=1}^N (x^{(i)2} - 2x^{(i)}\mu_{MLE} + \mu_{MLE}^2)\right] \\ &= E_D\left[\frac{1}{N} \sum_{i=1}^N x^{(i)2}\right] - E_D[\mu_{MLE}^2] \\ &= E_D\left[\frac{1}{N} \sum_{i=1}^N x^{(i)2} - \mu^2\right] - E_D[\mu_{MLE}^2 - \mu^2] \\ &= \sigma^2 - (E_D[\mu_{MLE}^2] - E_D^2[\mu_{MLE}]) \\ &= \sigma^2 - \text{var}[\mu_{MLE}] = \sigma^2 - \text{var}\left[\frac{1}{N} \sum_{i=1}^N x^{(i)}\right] \\ &= \sigma^2 - \frac{1}{N^2} \sum_{i=1}^N \text{var}[x^{(i)}] = \sigma^2 - \frac{1}{N} \sigma^2 \\ &= \frac{N-1}{N} \sigma^2 \end{aligned}$$

The unbiased estimator is

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (X^{(i)} - \mu_{MLE})^2$$

Multi variate Gaussian Distribution.

pdf is

$$P(\vec{x}; \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1} (\vec{x}-\vec{\mu})\right)$$

$$\vec{x}, \vec{\mu} \in \mathbb{R}^D \text{ and } \underline{\Sigma} \in \mathbb{R}^{D \times D}.$$

Σ : positive semi-definite matrix (PSD)

i.e. for any $\vec{v} \in \mathbb{R}^D$, $\vec{v}^T \Sigma \vec{v} \geq 0$

Here we consider Σ as PD for now.

$|\Sigma|$ is the determinant of Σ .

eigenvalue decomposition:

$$\begin{aligned} \Sigma &= U \Lambda U^T = [\vec{u}_1 | \vec{u}_2 | \dots | \vec{u}_D] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_D \end{bmatrix} [\vec{u}_1 | \vec{u}_2 | \dots | \vec{u}_D]^T \\ &= \sum_{i=1}^D \vec{u}_i \lambda_i \vec{u}_i^T \end{aligned}$$

$$\text{so } \Sigma^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^D \vec{u}_i \frac{1}{\lambda_i} \vec{u}_i^T$$

\uparrow

Orthogonal matrix i.e. $U U^T = I$.

$$\Delta = (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})$$

$$= \sum_{i=1}^D (\vec{x} - \vec{\mu})^T \vec{u}_i \frac{1}{\lambda_i} \vec{u}_i^T (\vec{x} - \vec{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

$\triangleq y_i$

y_i is the projected vector of $(\vec{x} - \vec{\mu})$ on the eigenvector \vec{u}_i

↑
equation of
ellipsoid

$D_M(\vec{x})$ = Mahalanobis distance.

$$= \sqrt{(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

MLE for θ :

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(x^{(i)} | \mu, \Sigma)$$

$$= -\frac{Nd}{2} \log(2\pi) + \frac{N}{2} \log(|\Sigma|)$$

$$- \frac{1}{2} \sum_{i=1}^N [(\vec{x}^{(i)} - \vec{\mu})^T \Sigma^{-1} (\vec{x}^{(i)} - \vec{\mu})] \triangleq \ell(\mu, \Sigma)$$

where

Σ^{-1} is the precision matrix

Fact:

① Trace is invariant under cyclic permutation of Matrix Product.

$$\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

$$\textcircled{2} \frac{\partial}{\partial A} \text{tr}(AB) = \frac{\partial}{\partial B} \text{tr}(BA) = B^T$$

$$\textcircled{3} \frac{\partial}{\partial A} \log |A| = A^{-T} \quad \textcircled{4} |A| = \frac{1}{|A^{-1}|}$$

$$\frac{\partial \ell(\mu, \Sigma)}{\partial \mu} = \sum_{i=1}^N \Sigma^{-1} (x^{(i)} - \mu) = 0$$

$$\Leftrightarrow \Sigma^{-1} \sum_{i=1}^N (x^{(i)} - \mu) = 0 \quad \text{assume } \Sigma \text{ is PD}$$

$$\Rightarrow \mu_{ME} = \frac{1}{N} \sum_{i=1}^N x^{(i)} = \text{empirical mean.}$$

$$\begin{aligned} \ell(\mu, \Sigma) = & -\frac{Nd}{2} \log(2\pi) + \frac{N}{2} \log(|\Sigma|) \\ & - \frac{1}{2} \sum_{i=1}^N \text{tr}((x^{(i)} - \mu)(x^{(i)} - \mu)^T \Sigma^{-1}) \end{aligned}$$

By trace trick.

$$= -\frac{Nd}{2} \log(2\pi) + \frac{N}{2} \log(|\Sigma|)$$

$$-\frac{1}{2} \text{tr}(S_{\bar{x}} \Sigma)$$

$$\text{where } S_{\bar{x}} = \sum_{i=1}^N (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

$$= \sum_{i=1}^N x^{(i)} (x^{(i)})^T - N \mu \mu^T$$

scatter matrix centered on μ

Rewrite scatter matrix into a compact form.

$$S_{\bar{x}} = \tilde{X}^T \tilde{X} = X^T C_N^T C_N X = X^T C_N X$$

$$\text{where } C_N = I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$$

is the centering matrix

$$\frac{\partial \ell}{\partial \Sigma} = -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} S_{\bar{x}} = 0$$

$$\text{Then } \Sigma_{MLE} = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

\Downarrow

$$\textcircled{1} E[\Sigma_{MLE}] = \frac{N-1}{N} \Sigma$$

② degree of freedom of θ is.

$$D + \frac{D(D+1)}{2} \sim O(D^2)$$

\downarrow
 μ

\downarrow
 Σ

\hookrightarrow overfitting -

possible way: assume Σ is a diagonal matrix.

Common theorems.

Denote $\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \in \mathbb{R}^D = \begin{bmatrix} \vec{x}_a \\ \vec{x}_b \end{bmatrix}$ $\vec{x}_a \in \mathbb{R}^m$
 $\vec{x}_b \in \mathbb{R}^n$
($m+n=D$)

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

$\vec{\mu} = (\vec{\mu}_a, \vec{\mu}_b)^T$ and $\vec{x} \sim N(\vec{\mu}, \Sigma)$

① If $\vec{x} \sim N(\vec{\mu}, \Sigma)$ and $\vec{y} = A\vec{x} + b$

then $\vec{y} \sim N(A\vec{\mu} + b, A\Sigma A^T)$

proof:

$$E[y] = E[Ax + b] = AE[x] + b = A\mu + b$$

$$\text{Var}[y] = \text{Var}[Ax + b] = \text{Var}[Ax] = A \cdot \text{Var}[x] \cdot A^T$$

② Find $P(x_a)$, $P(x_b)$, $P(x_a|x_b)$, $P(x_b|x_a)$

$$\vec{x}_a = \begin{bmatrix} I_{m \times m} & 0_{m \times n} \end{bmatrix} \begin{bmatrix} \vec{x}_a \\ \vec{x}_b \end{bmatrix}$$

\Downarrow

$$E[\vec{x}_a] = [I, 0] \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \mu_a$$

$$\begin{aligned} \text{Var}[\vec{x}_a] &= [I, 0] \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix} \\ &= \Sigma_{aa} \end{aligned}$$

$$\text{so } \vec{x}_a \sim N(\mu_a, \Sigma_{aa})$$

$$\text{similarly: } \vec{x}_b \sim N(\mu_b, \Sigma_{bb})$$

Define: ① $\vec{X}_{b|a} = \vec{X}_b - \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a$

② $\vec{\mu}_{b|a} = \vec{\mu}_b - \Sigma_{ba} \Sigma_{aa}^{-1} \vec{\mu}_a$

③ $\Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}$

Schur Complementary of Σ_{bb} .

$$X_{b|a} = \begin{bmatrix} -\Sigma_{ba} \Sigma_{aa}^{-1} & I_{n \times n} \end{bmatrix} \begin{bmatrix} x_a \\ x_b \end{bmatrix}$$

so $E[X_{b|a}]$

$$= \begin{bmatrix} -\Sigma_{ba} \Sigma_{aa}^{-1} & I_{n \times n} \end{bmatrix} \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix} = \vec{\mu}_{b|a}$$

$\text{var}[\vec{X}_{b|a}]$

$$= \begin{bmatrix} -\Sigma_{ba} \Sigma_{aa}^{-1} & I_{n \times n} \end{bmatrix} \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \begin{bmatrix} -\Sigma_{aa}^{-1} \Sigma_{ba}^T \\ I_{n \times n} \end{bmatrix}$$

$$= \Sigma_{b|a}$$

Then $\vec{X}_b = \vec{X}_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a$.

step 1: $\text{COV}(\vec{X}_{b|a}, \vec{X}_a)$

$$= \text{COV}(\vec{X}_b - \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a, \vec{X}_a)$$

$$= \Sigma_{ba} - \Sigma_{ba} \Sigma_{aa}^{-1} \cdot \Sigma_{aa} = 0$$

$\Rightarrow \vec{X}_{b|a}$ and \vec{X}_a are unrelated.

They are jointly normal so they are independent.

step 2:

$$E[\vec{X}_b | \vec{X}_a]$$

$$= E[\vec{X}_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a | \vec{X}_a]$$

$$= E[\vec{X}_{b|a}] + \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a = \vec{\mu}_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a$$

step 3:

$$= \vec{\mu}_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\vec{X}_a - \vec{\mu}_a)$$

$$\text{var}[\vec{X}_b | \vec{X}_a]$$

$$= \text{var}[\vec{X}_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} \vec{X}_a | \vec{X}_a]$$

$$= \text{var}[\vec{X}_{b|a}] = \Sigma_{b|a}$$

