

3.4. Maximum Likelihood Estimation

x_1, \dots, x_n sample, ~~independent~~ observations from
pdf $f(x; \vartheta)$, ϑ parameter vector

$$L(x; \vartheta) = \prod_{i=1}^n f(x_i; \vartheta)$$

is called likelihood function.

$$\ell(x; \vartheta) = \log L(x; \vartheta) = \sum_{i=1}^n \log f(x_i; \vartheta)$$

is called log-likelihood fun.

Given x_1, \dots, x_n , consider L and ℓ as a fct. of ϑ .

Find ϑ which fits the data best, i.e.) solve

$$\hat{\vartheta} = \arg \max_{\vartheta} \ell(x; \vartheta)$$

$\hat{\vartheta}$ is called maximum likelihood estimator. (MLE)

Th. 3.6. $X \sim N_p(\mu, \Sigma)$, x_1, \dots, x_n iid samples from X .

The MLE of μ and Σ are

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = S_n$$

Proof. Density of $N_p(\mu, \Sigma)$

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

$$\ell(x_1, \dots, x_n; \mu, \Sigma) = \sum_{i=1}^n \left[\log \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right]$$

$$= n \underbrace{\log \frac{1}{(2\pi)^{p/2}}}_{\text{constant}} + \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Leave the constant, set $\Lambda = \Sigma^{-1}$

$$\ell^*(\mu, \Sigma) = \frac{n}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Lambda (x_i - \mu)$$

$$= \frac{n}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^n \ln \Lambda (x_i - \mu) (x_i - \mu)^T$$

Steiner rule:

$$= \frac{n}{2} \log |\Lambda| - \frac{1}{2} \ln \Lambda \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T$$

Steiner rule:

$$\sum_{i=1}^n (x_i - \mu) (x_i - \mu)^T$$

$$= \underbrace{\sum_{i=1}^n (x_i - \bar{x}) (x_i - \bar{x})^T}_{n S_n} + (\bar{x} - \mu) (\bar{x} - \mu)^T$$

$$\geq n S_n \quad (\text{equality if } \mu = \bar{x})$$

$$\leq \frac{n}{2} \log |\Lambda| - \frac{n}{2} \ln \Lambda S_n = \ell^*(\mu^*, \Sigma)$$

max

$$\max_{\Lambda} \ell^*(\mu^*, \Lambda)$$

needs

$$\frac{\partial}{\partial \Lambda} \log |\Lambda| = \Lambda^{-1}$$

$$\frac{\partial}{\partial \Lambda} \text{tr}(\Lambda A) = A^T$$

$$\frac{\partial}{\partial \Lambda} \ell^*(\mu^*, \Lambda) = \frac{n}{2} \Lambda^{-1} - \frac{n}{2} S_n = 0$$

$$\text{Solution } \Lambda^{-1} = S_n = \sum \mathbb{E}$$

4. Principal Dimensionality Reduction

Represent cluster in a low-dim. space in an "optimal" way.

4.1. Principal component analysis (PCA)

Lose as little information as possible.

Given cluster $x_1, \dots, x_n \in \mathbb{R}^p$

- Find a k -dim subspace s.t. the projections of x_1, \dots, x_n thereon represent the original cluster on its best.
- Preserve as much variance as possible in the projected points.

a) and b) are equivalent! \rightarrow later

x_1, \dots, x_n independently sampled from some distribution.

Sample mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Sample covariance matrix $S_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top$

$\left[\begin{array}{l} \bar{x} : \text{MLE, unbiased estimator } E\bar{x} \\ S_n : \text{near MLE, unbiased est. of } \text{Cov}(X) \end{array} \right]$

4.1.1. Find the best projection

Optimization problem:

$$\min_{a \in \mathbb{R}^p} \sum_{i=1}^n \|x_i - a - Q(x_i - a)\|_F^2$$

Q orth. proj.

on a k -dim. subspace

$$\min_{a, Q} \sum_{i=1}^n \|x_i - a - Q(x_i - a)\|^2$$

$$= \min_{a, Q} \sum_{i=1}^n \|(I - Q)(x_i - a)\|^2$$

$$= \min_{a, R} \sum_{i=1}^n \|R(x_i - a)\|^2, \quad R = I - Q \text{ orth. proj}$$

$$= \min_{a, R} \sum_{i=1}^n (x_i - a)^\top \underbrace{R^\top R}_{=R} (x_i - a)$$

$$\begin{aligned}
 &= \min_{\alpha, R} \sum_{i=1}^n \text{tr} (x_i - \alpha)^T R (x_i - \alpha) \\
 &= \min_{\alpha, R} \text{tr} R \sum_{i=1}^n (x_i - \alpha)(x_i - \alpha)^T \\
 &\geq \min_R \text{tr} \left(R \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \right) \\
 &= \min_R \text{tr} (R(S_n - I)) \\
 &= \min_Q \text{tr} (S_n(I - Q))
 \end{aligned}$$

It remains to solve

$$\begin{aligned}
 &\max_Q \text{tr}(S_n Q), \quad Q \text{ orth. proj.} \\
 Q &= \sum_{i=1}^k q_i q_i^T, \quad q_i \text{ orth.} \\
 Q &= \tilde{Q} \tilde{Q}^T, \quad \tilde{Q} = (q_1, \dots, q_k) \\
 &= \max_{\tilde{Q}^T \tilde{Q} = I_k} \text{tr}(\tilde{Q}^T S_n \tilde{Q}) = \sum_{i=1}^k \lambda_i(S_n) \\
 &\quad (\text{K. Fann, Th. 2.4})
 \end{aligned}$$

where $\lambda_1(S_n) \geq \dots \geq \lambda_p(S_n)$ are the eigenvalues of S_n .

The max attained if q_1, \dots, q_k are the orth. normal eigenvectors corresp. to ~~$\lambda_1(S_n), \dots, \lambda_k(S_n)$~~ .

4.1.2 Preserve most variance

Seek k-dim. projection preserving most variance.

$$\begin{aligned}
 & \max_Q \sum_{i=1}^n \|Qx_i - \frac{1}{n} \sum_{\ell=1}^n Qx_\ell\|_F^2, \quad Q = \tilde{Q}\tilde{Q}^\top \\
 & \quad \tilde{Q}^\top \tilde{Q} = I_k \\
 & \quad Q \text{ orth. proj.} \\
 & = \max_Q \sum_{i=1}^n \|Qx_i - Q\bar{x}\|^2 \\
 & = \max_Q \sum_{i=1}^n \|Q(x_i - \bar{x})\|^2 \\
 & = \max_Q \sum_{i=1}^n \text{tr}(x_i - \bar{x})^\top Q(x_i - \bar{x}) \\
 & = \max_Q \text{tr} Q \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top \\
 & = \max_Q (n-1) \text{tr} Q S_n
 \end{aligned}$$

with the same solution as above.

4.1.3 How to carry out PCA

Given $x_1, \dots, x_n \in \mathbb{R}^p$, Fix $k \leq p$

Compute $S_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top$

$S_n = V \Lambda V^\top$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$

$\lambda_1 \geq \dots \geq \lambda_p$, $V = (v_1, \dots, v_p) = \mathcal{O}(p)$

Spectral decomposition

① v_1, \dots, v_k are called principal eigenvectors to the principle eigenvalues $\lambda_1 \geq \dots \geq \lambda_k$.

Projected points $\hat{x}_i = \begin{pmatrix} v_1^T \\ \vdots \\ v_k^T \end{pmatrix} x_i$

$(k\text{-dim})$ $(k \times p)$ $(p\text{-dim})$