

# Probabilistic Model Setup and the MLE

Liutong Zhou

November 1, 2016

## Content:

1. Sample Mean and Sample Variance in Matrix Form
2. MLE of Multivariate Gaussian
3. Expectation and Variance of Gaussian Vectors (notes P8)
4. Probabilistic Model Setup and MLE of Parameters
5. Why MLE?

## 1 Sample Mean and Sample Variance in Matrix Form

Given a sample of random variables  $\{X_n\} \stackrel{iid}{\sim} \text{dist}$ , denoted by vector  $\vec{X}_n \in \mathbb{R}^n$

$$\text{Sample Mean (scalar)} : \quad \bar{X}_n = \frac{1}{n} \vec{X}_n^T \mathbf{1}_n \quad (1.1)$$

$$\text{Mean Squared} : \quad \bar{X}_n^2 = \frac{1}{n^2} \vec{X}_n^T \mathbb{J}_{n \times n} \vec{X}_n \quad (1.2)$$

$$\text{Sample Variance} : \quad \Sigma_{i=1}^n (x_i - \bar{X}_n)^2 = \vec{X}_n^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) \vec{X}_n \quad (1.3)$$

Given a sample of random vectors  $\{\vec{X}_{p(n)}\} \stackrel{iid}{\sim} \text{dist}$ , denoted by matrix  $X_{n \times p}$  with each row  $\vec{X}_{p(i)}^T$ .

$$\text{Sample Mean (vector)} : \quad \vec{\bar{X}}_p = \frac{1}{n} X_{n \times p}^T \mathbf{1}_n \quad (1.4)$$

$$\text{Sample Variance (matrix)} : \quad \bar{\Sigma}_{p \times p} = X_{n \times p}^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) X_{n \times p} \quad (1.5)$$

*Proof. 1.2*

$$\bar{X}_n^2 = \|\frac{1}{n} \vec{X}_n^T \mathbf{1}_n\|^2 = \frac{1}{n^2} \vec{X}_n^T \mathbf{1}_n \mathbf{1}_n^T \vec{X}_n = \frac{1}{n^2} \vec{X}_n^T \mathbb{J} \vec{X}_n \quad \square$$

*Proof. 1.3*

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{X}_n)^2 &= \|\vec{X}_n - \bar{X}_n \mathbf{1}_n\|^2 = \vec{X}_n^T \vec{X}_n + n \bar{X}_n^2 - 2 \bar{X}_n \vec{X}_n^T \mathbf{1}_n = \vec{X}_n^T \vec{X}_n - n \bar{X}_n^2 \\ &= \vec{X}_n^T \vec{X}_n - n \frac{1}{n^2} \vec{X}_n^T \mathbb{J} \vec{X}_n = \vec{X}_n^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) \vec{X}_n \end{aligned}$$

$\square$

*Proof.* 1.5

$$\begin{aligned}
\bar{\Sigma}_{p \times p} &= \frac{1}{n} \sum_{i=1}^n (\overrightarrow{X_{p(i)}} - \overrightarrow{X_p})(\overrightarrow{X_{p(i)}} - \overrightarrow{X_p})^T = \frac{1}{n} (X_{n \times p} - \overrightarrow{\mathbb{1}_n} \overrightarrow{X_p}^T)^T (X_{n \times p} - \overrightarrow{\mathbb{1}_n} \overrightarrow{X_p}^T) \\
&= \frac{1}{n} (X_{n \times p} - \frac{1}{n} \overrightarrow{\mathbb{1}_n} \overrightarrow{\mathbb{1}_n}^T X_{n \times p})^T (X_{n \times p} - \frac{1}{n} \overrightarrow{\mathbb{1}_n} \overrightarrow{\mathbb{1}_n}^T X_{n \times p}) \\
&= \frac{1}{n} \left[ (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) X_{n \times p} \right]^T \left[ (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) X_{n \times p} \right] = \frac{1}{n} X_{n \times p}^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J})^2 X_{n \times p} \\
&= X_{n \times p}^T (\mathbf{I}_n - \frac{2}{n} \mathbb{J} + \frac{1}{n^2} n \mathbb{J}) X_{n \times p} = X_{n \times p}^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) X_{n \times p}
\end{aligned}$$

□

## 2 MLE of Multivariate Gaussian

**Fact 2.1.**  $\nabla_X \ln |X| = \frac{1}{|X|} \nabla_X |X| = \frac{1}{|X|} |X| X^{-T} = X^{-T}$

**Fact 2.2.**  $\nabla_X \text{tr}(AX) = \nabla_X \text{tr}(XA) = A^T$

**Fact 2.3.** The quadratic form  $\vec{a}^T X \vec{a} = \text{tr}(X \vec{a} \vec{a}^T)$

Given data  $= \{\overrightarrow{X_{p(n)}}\} \stackrel{iid}{\sim} N_P(\vec{\mu}_p, \Sigma_{p \times p})$

$$\text{Likelihood} : f_{X_{n \times p}}(X_{n \times p}) = \prod_{i=1}^n (2\pi)^{-\frac{p}{2}} \cdot |\Sigma|^{-\frac{1}{2}} \cdot \exp \left\{ -\frac{1}{2} (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)^T \Sigma^{-1} (\overrightarrow{X_{p(i)}} - \vec{\mu}_p) \right\}$$

$$\text{LogLikelihood} : L_{X_{n \times p}}(X_{n \times p}) = -\frac{np}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)^T \Sigma_{p \times p}^{-1} (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)$$

$$0 = \nabla_{\vec{\mu}_p} L_{X_{n \times p}}(X_{n \times p}) = -\frac{1}{2} \sum_{i=1}^n 2 \Sigma_{p \times p}^{-1} (\overrightarrow{X_{p(i)}} - \vec{\mu}_p) (-\mathbf{I}_p) = \Sigma_{p \times p}^{-1} \left( \sum_{i=1}^n \overrightarrow{X_{p(i)}} - n \vec{\mu}_p \right) \Rightarrow$$

$$\widehat{\vec{\mu}}_{MLE} = \frac{1}{n} \sum_{i=1}^n \overrightarrow{X_{p(i)}} = \boxed{\frac{1}{n} X_{n \times p}^T \overrightarrow{\mathbb{1}_n}}$$

$$0 = \nabla_{\Sigma^{-1}} L_{X_{n \times p}}(X_{n \times p}) = \nabla_{\Sigma^{-1}} \left\{ \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n \text{tr} \left[ \Sigma^{-1} (\overrightarrow{X_{p(i)}} - \vec{\mu}_p) (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)^T \right] \right\}$$

$$= \frac{n}{2} \Sigma_{p \times p} - \frac{1}{2} \sum_{i=1}^n (\overrightarrow{X_{p(i)}} - \vec{\mu}_p) (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)^T \Rightarrow$$

$$\widehat{\Sigma}_{MLE} = \frac{1}{n} \sum_{i=1}^n (\overrightarrow{X_{p(i)}} - \vec{\mu}_p) (\overrightarrow{X_{p(i)}} - \vec{\mu}_p)^T = \boxed{X_{n \times p}^T (\mathbf{I}_n - \frac{1}{n} \mathbb{J}) X_{n \times p}}$$

## 3 to be added

## 4 Probabilistic Model Setup and the MLE of Parameters

### 4.1 Model Setup

$$\vec{Y}_n = X_{n \times p} \vec{\beta}_p + \vec{\epsilon}_n$$

**Assumptions:**

1.  $\vec{\epsilon}_n \sim N_n(\vec{0}_n, \sigma^2 \mathbf{I}_n)$ , where  $\sigma^2$  is known.
2. For  $i = 1, \dots, n$ ,  $\overrightarrow{X_{p(i)}}^T \stackrel{iid}{\sim}$  some known dist with pdf  $f_{\vec{X}_p}(\vec{\xi}_p)$

3.  $\vec{\epsilon}_n \perp X_{n \times p}$ , or equivalently  $\overrightarrow{X_{p(i)}} \perp \vec{\epsilon}_j \quad \forall i, j = 1, \dots, n$  (easily violated)

**MLE:** Given data =  $(\vec{y}_n, X_{n \times p})$ ,  $\vec{Y}_n | X_{n \times p} \sim N_n(X_{n \times p} \vec{\beta}_p, \sigma^2 \mathbf{I}_n)$

*Likelihood* :  $f_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = f_{\vec{Y}_n}(\vec{Y}_n | X_{n \times p}, \vec{\theta}) \cdot f_{X_{n \times p}}(X_{n \times p} | \vec{\theta}) = (2\pi)^{-\frac{n}{2}} \cdot (\sigma^2)^{-\frac{n}{2}} \cdot$

$$\exp \left\{ -\frac{1}{2\sigma^2} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \right\} \prod_{i=1}^n f_{\vec{X}_{p(i)}}(\vec{X}_{p(i)})$$

*LogLikelihood* :  $L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) + \text{const.}$

$$\text{set } 0 = \nabla_{\vec{\beta}} L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = -\frac{1}{2\sigma^2} \cdot 2(-X_{n \times p}^T)(\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \Rightarrow X_{n \times p}^T \vec{Y}_n = X^T X \vec{\beta}_p \Rightarrow$$

$$\boxed{\widehat{\vec{\beta}}_{pMLE} = (X^T X)^{-1} X^T \vec{Y}_n}$$

$$\text{set } 0 = \nabla_{\sigma^2} L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \|\vec{Y}_n - X_{n \times p} \vec{\beta}_p\|^2 \Rightarrow \boxed{\widehat{\sigma}_{MLE}^2 = \frac{1}{n} \|\vec{Y}_n - X_{n \times p} \vec{\beta}_p\|^2 = \frac{1}{n} \|(\mathbf{I} - \mathbf{P}) \vec{Y}_n\|^2}$$

$$= \frac{1}{n} \vec{Y}_n^T (\mathbf{I} - \mathbf{P}) \vec{Y}_n \sim \frac{1}{n} \sigma^2 \chi_{n-p}^2$$

**Note:**

1. Even if in assumption 2 the rows of design matrix X are not independent,  $\hat{\vec{\beta}}_{MLE} = \hat{\vec{\beta}}_{LSE}$  = sample mean is still true.
2. Even if in assumption 1, the distribution is not Gaussian, it is eventually asymptotic Gaussian.  $\hat{\vec{\beta}}_{MLE} \xrightarrow{D} N_p(\beta_p, \frac{I_{\beta}^{-1}}{n})$
3. If assumption 3 is violated, the model setup is no longer valid.

## 4.2 Non-i.i.d errors

If assumption 1 is changed to  $\vec{\epsilon}_n \sim N_n(\vec{0}_n, \text{some known } \Sigma)$  (non i.i.d. errors), then

**MLE:** given data =  $(\vec{y}_n, X_{n \times p})$ ,  $\vec{Y}_n | X_{n \times p} \sim N_n(X_{n \times p} \vec{\beta}_p, \Sigma_{p \times p})$

*Likelihood* :  $f_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = f_{\vec{Y}_n}(\vec{Y}_n | X_{n \times p}, \vec{\theta}) \cdot f_{X_{n \times p}}(X_{n \times p} | \vec{\theta}) = (2\pi)^{-\frac{n}{2}} \cdot |\Sigma|^{-\frac{1}{2}} \cdot$

$$\exp \left\{ -\frac{1}{2} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T \Sigma_{p \times p}^{-1} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \right\} \prod_{i=1}^n f_{\vec{X}_{p(i)}}(\vec{X}_{p(i)})$$

*LogLikelihood* :  $L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = -\frac{n}{2} \ln(2\pi) + \frac{1}{2} \ln|\Sigma^{-1}| - \frac{1}{2} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T \Sigma^{-1} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) + \text{const.}$

$$\text{set } 0 = \nabla_{\vec{\beta}} L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = -\frac{1}{2} \cdot 2(-X_{n \times p}^T) \Sigma^{-1} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \Rightarrow X_{n \times p}^T \Sigma^{-1} \vec{Y}_n = X^T \Sigma^{-1} X \vec{\beta}_p \Rightarrow$$

$$\boxed{\widehat{\vec{\beta}}_{pMLE} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \vec{Y}_n}$$

$$\text{set } 0 = \nabla_{\Sigma^{-1}} L_{\vec{\theta}}(\vec{Y}_n, X_{n \times p} | \vec{\theta}) = \frac{1}{2} \Sigma^T - \frac{1}{2} \nabla_{\Sigma^{-1}} \text{tr} \left[ \Sigma^{-1} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \cdot (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T \right]$$

$$= \frac{1}{2} \Sigma^T - \frac{1}{2} (\vec{Y}_n - X_{n \times p} \vec{\beta}_p) \cdot (\vec{Y}_n - X_{n \times p} \vec{\beta}_p)^T \Rightarrow \boxed{\widehat{\Sigma}_{MLE} = (\vec{Y}_n - X_{n \times p} \widehat{\vec{\beta}}_{pMLE}) \cdot (\vec{Y}_n - X_{n \times p} \widehat{\vec{\beta}}_{pMLE})^T}$$

## 5 Why MLE

No matter what distribution  $\vec{\epsilon}_n$  has, under the assumption that  $data \stackrel{iid}{\sim}$  some dist. , MLE is asymptotically normal and efficient (unbiased and hits the Cramer-Rao bound, hence has the MV) when  $n \rightarrow \infty$  , thus the best estimator.

$\hat{\theta}_{MLE} \xrightarrow{D} N(\theta, \frac{\mathbf{I}^{-1}}{n})$  where  $\mathbf{I}$  is Fisher Information Matrix (see summary notes of Chapter 8 )

Note: How large is large? Empirically,  $n \geq 50p$  for the design matrix.

*Proof.* We need Theorem 1,2,3 in Notes of Chapter 6 Part II, Limit Theorems in Multivariate Case. Under the same assumptions in 4.1 except that assumption 1 is modified to  $\epsilon_i \stackrel{iid}{\sim}$  some unknown dist. with known mean 0 and variance  $\sigma^2$ , then

$$\widehat{\vec{\beta}}_{MLE} - \vec{\beta} = (X^T X)^{-1} X^T \vec{y} - \vec{\beta} = (X^T X)^{-1} X^T \vec{\epsilon}_n$$

$$\frac{1}{n} X_{n \times p}^T X_{n \times p} \xrightarrow{LLN} E(\vec{X}_p \vec{X}_p^T) = \tilde{\Sigma}_{p \times p} \quad \textcircled{1}$$

$$X_{n \times p}^T \vec{\epsilon}_n = \sum_{i=1}^n \vec{X}_{p(i)} \cdot \epsilon_i \xrightarrow{CLT} N_p \left( \sum_{i=1}^n E(\vec{X}_{p(i)} \cdot \epsilon_i), \sum_{i=1}^n Var(\vec{X}_{p(i)} \cdot \epsilon_i) \right) \quad \textcircled{2}$$

$$\text{but } E(\vec{X}_{p(i)} \cdot \epsilon_i) = E(\vec{X}_{p(i)}) \cdot E(\epsilon_i) = \vec{0}_p$$

$$Var(\vec{X}_{p(i)} \cdot \epsilon_i) = E \left[ (\vec{X}_{p(i)} \cdot \epsilon_i) (\vec{X}_{p(i)} \cdot \epsilon_i)^T \right] - E(\vec{X}_{p(i)} \cdot \epsilon_i) E(\vec{X}_{p(i)} \cdot \epsilon_i)^T = E(\vec{X}_{p(i)}) E(\epsilon_i \epsilon_i^T) E(\vec{X}_{p(i)})^T$$

$$= \sigma^2 \tilde{\Sigma}_{p \times p} \xrightarrow{\textcircled{2}} X_{n \times p}^T \vec{\epsilon}_n \xrightarrow{CLT} N_p(\vec{0}_p, n \sigma^2 \tilde{\Sigma}_{p \times p}) \quad \textcircled{3}$$

$$\textcircled{1}^{-1} \times \textcircled{3} \Rightarrow n(X^T X)^{-1} X^T \vec{\epsilon}_n \xrightarrow{\text{Slutsky's Theorem}} N_p(\vec{0}_p, n \sigma^2 \tilde{\Sigma}_{p \times p}^{-1}) \Rightarrow$$

$$\boxed{(X^T X)^{-1} X^T \vec{\epsilon}_n \xrightarrow{D} N_p(\vec{0}_p, \frac{\sigma^2}{n} \tilde{\Sigma}_{p \times p}^{-1})}, \text{ where } \tilde{\Sigma}_{p \times p} = E(\vec{X}_p \vec{X}_p^T)$$

□