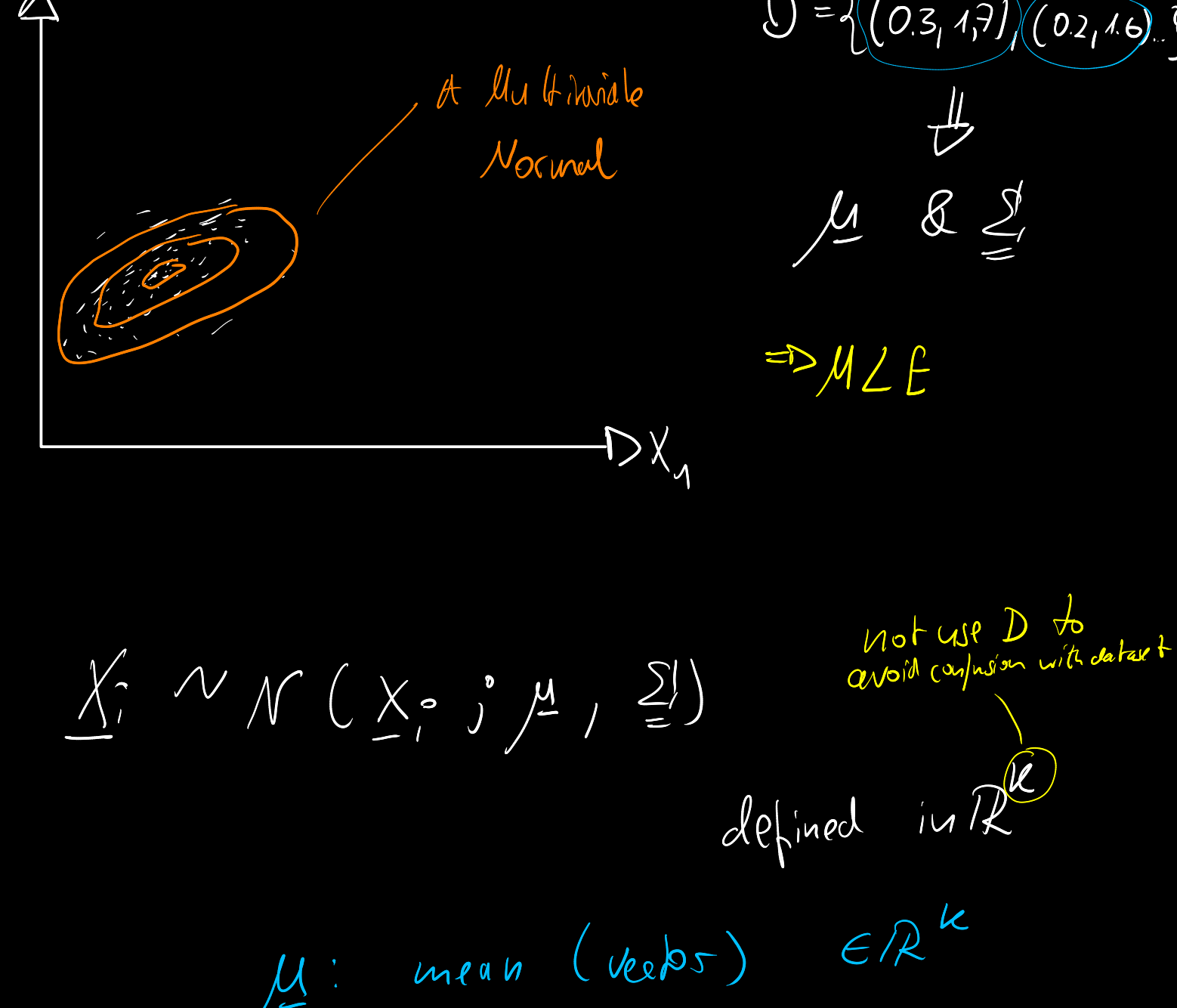


Maximum Likelihood Estimate for Multivariate Normal



$\underline{X}_i \sim \mathcal{N}(\underline{x}_0; \underline{\mu}, \underline{\Sigma})$

not use \underline{D} to avoid confusion with dataset

defined in \mathbb{R}^k

$\underline{\mu}$: mean (vector) $\in \mathbb{R}^k$

$\underline{\Sigma}$: covariance (matrix) $\in \mathbb{R}^{k \times k}$

$\underline{\Sigma} > 0 \rightarrow$ sym. positive definite

$(\underline{\Sigma} = \underline{\Sigma}^T)$

e.g. $k=2$

$$\underline{\Sigma} = \begin{bmatrix} 1.7 & 0.3 \\ 0.3 & 1.5 \end{bmatrix}$$

(Covariance)

Variance

$$N(\underline{X}_i; \underline{\mu}, \underline{\Sigma}) = \frac{1}{\sqrt{\det(\underline{\Sigma}) \cdot (2\pi)^k}} \exp\left(-\frac{1}{2}(\underline{X}_i - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{X}_i - \underline{\mu})\right)$$

Likelihood:

$$\mathcal{L}(\underline{D}; \underline{\mu}, \underline{\Sigma}) \stackrel{\text{i.i.d.}}{=} \prod_{i=0}^{N-1} p(\underline{X}_i = \underline{x}^{(i)})$$
$$= \prod_{i=0}^{N-1} N(\underline{X}_i = \underline{x}^{(i)}; \underline{\mu}, \underline{\Sigma})$$

$$= \prod_{i=0}^{N-1} \frac{1}{\sqrt{\det(\underline{\Sigma}) \cdot (2\pi)^k}} \exp\left(-\frac{1}{2}(\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu})\right)$$

don't confuse

Log-likelihood

$$\ell(\underline{D}; \underline{\mu}, \underline{\Sigma}) = \log(\mathcal{L}(\underline{D}; \underline{\mu}, \underline{\Sigma}))$$
$$= -\frac{N}{2} \log \det(\underline{\Sigma}) - \frac{Nk}{2} \log(2\pi) - \frac{1}{2} \sum_{i=0}^{N-1} (\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu})$$

Maximum Likelihood Estimate

$$\underline{\mu}_1^*, \underline{\Sigma}_1^* = \underset{\substack{\underline{\mu} \in \mathbb{R}^k \\ \underline{\Sigma} \in \mathbb{R}^{k \times k} \\ \underline{\Sigma} > 0}}{\operatorname{argmax}} (\ell(\underline{D}; \underline{\mu}, \underline{\Sigma}))$$

$\underline{\Sigma} > 0$ — $\text{cov}(\delta \theta)$ naturally $*$

Take derivative & set to zero

* More detail on this (was not 100% in the video)

our MLE does not guarantee positive definiteness, although it guarantees symmetry.

In order to be on the safe side one can use sth called a "shrunk" covariance.

Essentially, this is: $\underline{\Sigma}_{\text{shrunk}} = (1-\alpha)\underline{\Sigma}_{\text{MLE}} + \alpha \frac{\operatorname{tr}(\underline{\Sigma}_{\text{MLE}})}{k} \cdot \underline{I}$

which generates more diagonal dominance

(1)

$$\frac{\partial \ell}{\partial \underline{\mu}} = -\frac{1}{2} \sum_{i=0}^{N-1} \left(\underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu}) + (\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1} \right) \stackrel{!}{=} \underline{0}$$

$$(\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1} = \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu})$$

because $\underline{\Sigma}^T = \underline{\Sigma}$

$\hookrightarrow \underline{\Sigma}^{-1} = (\underline{\Sigma}^T)^{-1}$

Small correction to the video (at ~ 14:40)

the shown equality is not correct

(left hand side is a row vector, right hand side is a column vector)

Problem is in the derivative, correct would be according to formula (A6) of the matrix cookbook

$$\frac{\partial (\underline{X}^T \underline{B} \underline{X})}{\partial \underline{X}} = \underline{B} \underline{X} + \underline{B}^T \underline{X}$$

in our case

$$\frac{\partial}{\partial \underline{\mu}} \left(-\frac{1}{2} \sum_{i=0}^{N-1} (\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu}) \right)$$
$$= -\frac{1}{2} \sum_{i=0}^{N-1} \left(\underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu}) + \underline{\Sigma}^{-T}(\underline{x}^{(i)} - \underline{\mu}) \right)$$

$\underline{\Sigma}^{-1} = \underline{\Sigma}^{-T}$ due to symmetry

$$= -\frac{1}{2} \sum_{i=0}^{N-1} \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu}) \stackrel{!}{=} \underline{0}$$

1 (-1)

$$= \sum_{i=0}^{N-1} \underline{\Sigma}^{-1} \underline{x}^{(i)} - \underline{\Sigma}^{-1} \underline{\mu} \stackrel{!}{=} \underline{0}$$

$\underline{\Sigma}^{-1}$

$$\hookrightarrow \sum_{i=0}^{N-1} \underline{x}^{(i)} - \underline{\mu} \stackrel{!}{=} \underline{0}$$
$$\sum_{i=0}^{N-1} \underline{x}^{(i)} = \underline{\mu}$$

$\underline{\mu}$

$$\Leftrightarrow \underline{\mu} = \frac{1}{N} \sum_{i=0}^{N-1} \underline{x}^{(i)}$$

(2)

$$\frac{\partial \ell}{\partial \underline{\Sigma}} = -\frac{N}{2} \log \det(\underline{\Sigma}) - \frac{1}{2} \sum_{i=0}^{N-1} (\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu})$$

Scalar

$\underline{\Sigma}^{-1}$ does not commute

\hookrightarrow commutes inside trace

$$\operatorname{tr}(\underline{A} \underline{B}) = \operatorname{tr}(\underline{B} \underline{A})$$

in general: $\underline{A} \underline{B} \neq \underline{B} \underline{A}$

$$= -\frac{N}{2} \log \det(\underline{\Sigma}) - \frac{1}{2} \sum_{i=0}^{N-1} \operatorname{tr}((\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}^{(i)} - \underline{\mu}))$$
$$= -\frac{N}{2} \log \det(\underline{\Sigma}) - \frac{1}{2} \sum_{i=0}^{N-1} \operatorname{tr}((\underline{x}^{(i)} - \underline{\mu})(\underline{x}^{(i)} - \underline{\mu})^T \underline{\Sigma}^{-1})$$

outer product

$$\underline{S} = (\underline{x}^{(0)} - \underline{\mu})(\underline{x}^{(0)} - \underline{\mu})^T$$
$$= \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix} \in \mathbb{R}^{k \times k}$$

Matrix cookbook

$$\frac{\partial \det(\underline{A})}{\partial \underline{A}} = \det(\underline{A}) \cdot \underline{A}^{-T}$$
$$\frac{\partial \operatorname{tr}(\underline{A} \underline{B}^{-1})}{\partial \underline{B}} = -(\underline{B}^{-1} \underline{A} \underline{B}^{-1})^T$$

$$\frac{\partial \ell}{\partial \underline{\Sigma}} = -\frac{N}{2} \frac{1}{\det(\underline{\Sigma})} \cdot \det(\underline{\Sigma}) \underline{\Sigma}^{-T} + \frac{1}{2} \sum_{i=0}^{N-1} (\underline{\Sigma}^{-1} \underline{S} \underline{\Sigma}^{-1})^T$$
$$\stackrel{!}{=} \underline{0}$$

$$= -\frac{N}{2} \underline{\Sigma}^{-1} + \frac{1}{2} \sum_{i=0}^{N-1} \underline{\Sigma}^{-1} \underline{S} \underline{\Sigma}^{-1} \stackrel{!}{=} \underline{0}$$

1 0

$$= -N \underline{\Sigma}^{-1} + \sum_{i=0}^{N-1} \underline{\Sigma}^{-1} \left((\underline{x}^{(i)} - \underline{\mu})(\underline{x}^{(i)} - \underline{\mu})^T \right) \underline{\Sigma}^{-1} \stackrel{!}{=} \underline{0}$$

$\underline{\Sigma}^{-1}$

$$= -N \underline{\Sigma}^{-1} + \sum_{i=0}^{N-1} (\underline{x}^{(i)} - \underline{\mu})(\underline{x}^{(i)} - \underline{\mu})^T \stackrel{!}{=} \underline{0}$$

$$\underline{\Sigma}_{\text{MLE}} = \frac{1}{N} \sum_{i=0}^{N-1} (\underline{x}^{(i)} - \underline{\mu})(\underline{x}^{(i)} - \underline{\mu})^T$$

MLE

$$\underline{\Sigma}_{\text{MLE}} = \frac{1}{N} \sum_{i=0}^{N-1} \begin{bmatrix} x_1^{(i)} - \mu_1 & \dots & x_k^{(i)} - \mu_k \\ \vdots & \ddots & \vdots \\ x_k^{(i)} - \mu_k & \dots & x_k^{(i)} - \mu_k \end{bmatrix}$$

$\in \mathbb{R}^{k \times k}$

data matrix

$$\underline{X} = \begin{bmatrix} \dots & x_1^{(0)} & \dots & x_k^{(0)} \\ \dots & x_1^{(1)} & \dots & x_k^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ \dots & x_1^{(N-1)} & \dots & x_k^{(N-1)} \end{bmatrix} \in \mathbb{R}^{N \times k}$$

k dim ($\hat{=}$ features)

N samples

Mean alongside rows

$$\underline{\mu} = [\dots \dots \dots] \in \mathbb{R}^k$$

centered data matrix

$$\underline{\tilde{X}} = \underline{X} - \begin{bmatrix} \dots & \underline{\mu} & \dots \\ \dots & \underline{\mu} & \dots \\ \vdots & \vdots & \ddots \\ \dots & \underline{\mu} & \dots \end{bmatrix}$$

$\in \mathbb{R}^{N \times k}$

$$\underline{\tilde{X}} = \underline{X} - \underline{1} \underline{\mu}^T$$

\downarrow

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\underline{\Sigma}_{\text{MLE}} = \frac{1}{N} \underline{\tilde{X}}^T \underline{\tilde{X}} \in \mathbb{R}^{k \times k}$$

$$= \frac{1}{N} \cdot k \cdot \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix} \cdot \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

$$= \frac{1}{N} \cdot k \cdot \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$