

Maximum Likelihood and the Multivariate Gaussian

Let $\vec{\mu}, \vec{x} \in \mathbb{R}^D$, $\Sigma \in \text{PSD}^{D \times D}$ then let

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

be the "multivariate normal density function" of RV \vec{x} given vector mean $\vec{\mu}$ and covariance matrix Σ .

An important thing to know about MVN distributed vectors is that

\vec{x} has the same distribution as

$$A \vec{z} + \vec{\mu} \quad \text{where}$$

$$\vec{z} = [z_1, \dots, z_D], \quad z_j \sim N(0, 1)$$

and A satisfies $AA^T = \Sigma$ (Cholesky)

Certainly we can note that

$$\mathbb{E}[A \vec{z} + \vec{\mu}] = A \mathbb{E}[\vec{z}] + \vec{\mu} = A \cdot \vec{0} + \vec{\mu} = \vec{\mu}$$

$$\begin{aligned} \text{and} \\ \text{Cov}[A \vec{z} + \vec{\mu}] &= \text{Cov}[A \vec{z}] = A \text{Cov}[\vec{z}] A^T \\ &= A A^T \\ &= \Sigma \end{aligned}$$

But higher-order moments could appear.

We can "derive" the MVN PDF starting from

$$P_{\vec{z}}(\vec{z}) = \prod_{i=1}^D (2\pi)^{-1/2} \exp\left\{-\frac{1}{2} z_i^2\right\} = (2\pi)^{-D/2} \exp\left\{-\frac{1}{2} \vec{z}^T \vec{z}\right\}$$

By the change of variable rule, let

$$\vec{x} = A\vec{z} + \vec{\mu} \Rightarrow A^{-1}(\vec{x} - \vec{\mu}) = \vec{z}$$

where $A A^T = \Sigma$

The multivariate change of variable rule

says

$$P_{\vec{x}}(\vec{x}) = P_{\vec{z}}(\vec{z}) \begin{vmatrix} dz_1/dx_1 & dz_1/dx_2 & \dots & dz_1/dx_D \\ dz_2/dx_1 & dz_2/dx_2 & \dots & dz_2/dx_D \\ \vdots & \vdots & \ddots & \vdots \\ dz_D/dx_1 & dz_D/dx_2 & \dots & dz_D/dx_D \end{vmatrix}$$

$$g(\vec{x}) = f(\vec{z}) \left| \det \left(\frac{d\vec{z}}{d\vec{x}} \right) \right|$$

$$\vec{x} = r(\vec{z})$$

$$r: \mathcal{Z} \rightarrow \mathcal{X}$$

$$\text{Noting } \frac{d}{d\vec{x}} A^{-1}(\vec{x} - \vec{\mu}) = \frac{d\vec{z}}{d\vec{x}}$$

$$A^{-1} = \frac{\partial \vec{z}}{\partial \vec{x}}$$

and

$$\begin{bmatrix} a_{11}^{-1} & a_{12}^{-1} \\ \vdots & \vdots \\ a_{D1}^{-1} & a_{D2}^{-1} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} z_1 \\ \vdots \\ z_D \end{bmatrix}$$

we have

$$P_{\vec{x}}(\vec{x}) = P_{\vec{z}}(\vec{z}) |A^{-1}| \quad \text{note } |A^{-1}| = |A|^{-1}$$

$$= P_{\vec{z}}(A^{-1}(\vec{x} - \vec{\mu})) |A|^{-1}$$

$$\begin{aligned}
 P_X(\vec{x}) &= P_Z(A^{-1}(\vec{x} - \mu)) |A|^{-1} \\
 &= (2\pi)^{-D/2} |A|^{-1} \exp \left\{ -\frac{1}{2} [A^{-1}(\vec{x} - \mu)]^T [A^{-1}(\vec{x} - \mu)] \right\} \\
 &= (2\pi)^{-D/2} |AA^T|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\vec{x} - \mu)^T (A^{-1})^T A^{-1} (\vec{x} - \mu) \right\} \\
 &= (2\pi)^{-D/2} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\vec{x} - \mu)^T (AA^T)^{-1} (\vec{x} - \mu) \right\} \\
 &= (2\pi)^{-D/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\vec{x} - \mu)^T \Sigma^{-1} (\vec{x} - \mu) \right\} \quad \checkmark
 \end{aligned}$$

Where we have used

$$\begin{aligned}
 |A|^{-1} &= |A|^{-\frac{1}{2}} |A|^{-\frac{1}{2}} = |A|^{-\frac{1}{2}} |A^T|^{-\frac{1}{2}} = (|A| |A^T|)^{-\frac{1}{2}} \\
 &= |AA^T|^{-\frac{1}{2}}
 \end{aligned}$$

$$\text{and } (A^{-1})^T A^{-1} = (A^T)^{-1} A^{-1} = (AA^T)^{-1}$$

with individual steps from

Sam Rowe's Matrix cheat sheets
or the "matrix cookbook"

$$\begin{aligned}
 \text{notably } |A| &= |A^T|, & (A^{-1})^T &= (A^T)^{-1} \\
 (AB)^{-1} &= B^{-1} A^{-1} & \text{and } |AB| &= |A| |B|
 \end{aligned}$$

So we can represent a multivariate Gaussian as an affine transformation of a vector of individual $N(0, I)$ distributed RV's

But, given a collection of N vectors $\vec{x}_1, \dots, \vec{x}_N$ which we assume to be iid $N(\mu, \Sigma)$ how do we get $\mu \in \Sigma$?

M.L. For MV Gaussians

$$\text{Let } X = \{ \vec{x}_1, \dots, \vec{x}_N \}$$

$$\mathcal{L}(X; \{\mu, \Sigma\}) = \prod_{n=1}^N \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x}_n - \mu)^T \Sigma^{-1} (\vec{x}_n - \mu)\right\}$$

like before we will estimate μ_{ML} & Σ_{ML} by taking derivatives of log likelihoods and setting the = equal to zero.

We will need the following, also from the Matrix cookbook.

with W symmetric

$$\frac{d}{ds} (x-s)^T W (x-s) = -2W(x-s)$$

$$\frac{\partial \ln |X|}{\partial X} = (X^{-1})^T = (X^T)^{-1}$$

$$\frac{\partial (a^T X b)}{\partial X} = b a^T$$

$$\frac{\partial a^T X^{-1} b}{\partial X} = -X^{-T} a b^T X^{-T}$$

First ML Est $\hat{\mu}$ for μ :

$$\frac{d}{d\vec{\mu}} \ln \mathcal{L}(X; \{\mu, \Sigma\})$$

$$= \sum_{n=1}^N \frac{\partial}{\partial \vec{\mu}} \left(-\frac{1}{2} (\vec{x}_n - \vec{\mu})^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}) \right) = 0$$

$$= \sum_{n=1}^N -\Sigma^{-1} (\vec{x}_n - \vec{\mu}) = 0$$

$$\Rightarrow \sum_{n=1}^N \vec{x}_n - N \vec{\mu} = 0$$

$$\Rightarrow \hat{\vec{\mu}} = \frac{1}{N} \sum_{n=1}^N \vec{x}_n$$

Note

$$\mathbb{E}[\hat{\vec{\mu}}] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\vec{x}_n] = \frac{1}{N} N \vec{\mu} \Rightarrow \text{Bias}(\hat{\vec{\mu}}) = 0$$

Now $\hat{\Sigma}_{ML} = \underset{\Sigma}{\operatorname{argmax}} \log \mathcal{L}(X; \hat{\mu}, \Sigma)$

$$\begin{aligned} & \frac{\partial}{\partial \Sigma} \log \mathcal{L}(X; \hat{\mu}, \Sigma) \\ &= -\frac{N}{2} \frac{\partial}{\partial \Sigma} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial \Sigma} (x_n - \hat{\mu})^T \Sigma^{-1} (x_n - \hat{\mu}) \\ &= -\frac{N}{2} (\Sigma^{-1})^T + \frac{1}{2} \sum_{n=1}^N \Sigma^{-T} (x_n - \hat{\mu}) (x_n - \hat{\mu})^T \Sigma^{-T} \end{aligned}$$

Setting this expression equal to zero and solving yields

Σ symmetric so $(\Sigma^{-1})^T = (\Sigma^{-1})$

$$N \Sigma^{-1} = \sum_{n=1}^N \Sigma^{-1} (x_n - \hat{\mu}) (x_n - \hat{\mu})^T \Sigma^{-1}$$

$$\Rightarrow \hat{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu}) (x_n - \hat{\mu})^T$$

note Λ 's

So now, given data, we can analytically compute the ML parameters of a multivariate Gaussian.

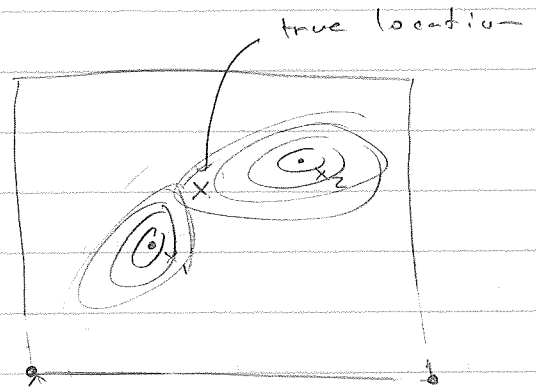
The maximum likelihood principle has other "uses".

A Note: the ML estimator, particularly for exponential families, has nice properties, particularly in the limit. It can be shown, for instance, that the sampling distribution of most ML estimators tends asymptotically to a normal distribution.

Extra reading: Fisher Information Matrix.

Example: ML For Sensor Fusion

Assume sensors are unbiased but produce "noisy" measurements of 2D location.



Assume 2 sensors that make independent measurements of an object's location in 2D space.

From the sensors we have

$$x_1, x_2 \in \mathbb{R}^2 \text{ the measurements}$$

From the sensor manufacturer we know the MVN observation noise variance for each sensor

$$x_i - x_{\text{true}} \sim \mathcal{N}(0, \Sigma_i) \Leftrightarrow x_{\text{true}} \sim \mathcal{N}(x_i, \Sigma_i)$$

$$\Leftrightarrow x_i \sim \mathcal{N}(x_{\text{true}}, \Sigma_i)$$

Let's say we believe both sensors equally, and that they are independent.

Then

$$P(x_1, x_2 | x_{\text{true}}) = P(x_1 | x_{\text{true}}) P(x_2 | x_{\text{true}})$$

$$= \mathcal{N}(x_1 | x_{\text{true}}, \Sigma_1) \mathcal{N}(x_2 | x_{\text{true}}, \Sigma_2)$$

The ML principle says that

$$\hat{x}_{\text{true}} = \underset{x_{\text{true}}}{\text{argmax}} \log \mathcal{N}(x_1 | x_{\text{true}}, \Sigma_1) + \log \mathcal{N}(x_2 | x_{\text{true}}, \Sigma_2)$$

Setting up as usual

$$\frac{d}{dx_{\text{true}}} \left((x_1 - x_{\text{true}})^T \Sigma_1^{-1} (x_1 - x_{\text{true}}) + (x_2 - x_{\text{true}})^T \Sigma_2^{-1} (x_2 - x_{\text{true}}) \right) = 0$$

$$-2 \Sigma_1^{-1} (x_1 - x_{\text{true}}) - 2 \Sigma_2^{-1} (x_2 - x_{\text{true}}) = 0$$

$$\Sigma_1^{-1} x_1 + \Sigma_2^{-1} x_2 = \left(\Sigma_1^{-1} + \Sigma_2^{-1} \right) x_{\text{true}}$$

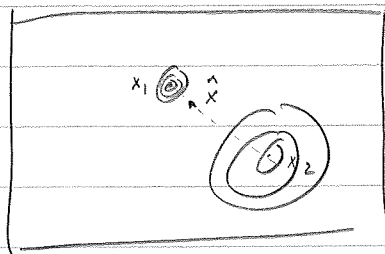
$$\hat{x}_{\text{true}} = \left(\Sigma_1^{-1} + \Sigma_2^{-1} \right)^{-1} \left(\Sigma_1^{-1} x_1 + \Sigma_2^{-1} x_2 \right)$$

In multiple dimensions this is a little difficult to see, but, let's assume that the derivation holds in 1-D, then

$$\hat{x}_{\text{true}} = \left(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} \right)^{-1} \left(\frac{1}{\sigma_1^2} x_1 + \frac{1}{\sigma_2^2} x_2 \right)$$

which is just a weighted average with higher weight on the smallest variance.

Graphically (back in 2d)



ML ≠ Linear Regression (is not what you think it is)

If you think linear regression is about drawing a line through points, think again! (OK, it's that too)

Linear regression as you know it-ish:

$$n \in \{1, \dots, N\}$$

$$y_n \in \mathbb{R}$$

$$x_n \in \mathbb{R}$$

$$w \in \mathbb{R}$$

$$b \in \mathbb{R}$$

$$Y = \{y_1, \dots, y_N\}$$

$$X = \{x_1, \dots, x_N\}$$

$$y_n | x_n \sim N(wx_n + b, \sigma^2)$$

$$\Rightarrow \mathcal{L}(X; w, b) = \prod_{n=1}^N N(y_n; wx_n + b, \sigma^2)$$

$$\underset{w}{\operatorname{argmax}} \log \prod \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} (y_n - wx_n + b)^2 / \sigma^2 \right\}$$

= take deri's wrt. w, b set = 0, = gross!

Use matrix linear regression with

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \vec{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_N & 1 \end{bmatrix} \quad \vec{w} = \begin{bmatrix} w \\ b \end{bmatrix}$$

with likelihood

$$\vec{y} = \vec{X}\vec{w} \sim N(\vec{0}, I)$$

and

$$\underset{w}{\operatorname{argmax}} \log \mathcal{L}(Y; \vec{w})$$

$$= w \text{ s.t. } \frac{\partial}{\partial \vec{w}} \log \mathcal{L}(\vec{y} = \vec{X}\vec{w}; I) = \vec{0}$$

Plugging in to MUV PDF and taking logs we see

$$\begin{aligned}
 U &= \frac{d}{d\vec{\omega}} (\vec{y} - \vec{X}\vec{\omega})^T \mathbf{I} (\vec{y} - \vec{X}\vec{\omega}) \quad \leftarrow \text{least squares} \\
 &= \frac{d}{d\vec{\omega}} -2(\vec{X}\vec{\omega})^T \vec{y} + (\vec{X}\vec{\omega})^T (\vec{X}\vec{\omega}) \\
 &= \frac{d}{d\vec{\omega}} -2\vec{\omega}^T \vec{X}^T \vec{y} + \vec{\omega}^T \vec{X}^T \vec{X} \vec{\omega}
 \end{aligned}$$

$\sum_{i=1}^n (y_i - x_i \omega)^2$
 $= \sum_{i=1}^n (y_i - (\omega x_i + b))^2$

From matrix cookbook

$$\begin{aligned}
 &= -2\vec{X}^T \vec{y} + (\vec{X}^T \vec{X} + (\vec{X}^T \vec{X})^T) \vec{\omega} \\
 &= \vec{X}^T \vec{y} + \vec{X}^T \vec{X} \vec{\omega}
 \end{aligned}$$

$$\Rightarrow \vec{\omega}_{ML} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y} = \vec{\omega}_{\text{least squares}}$$

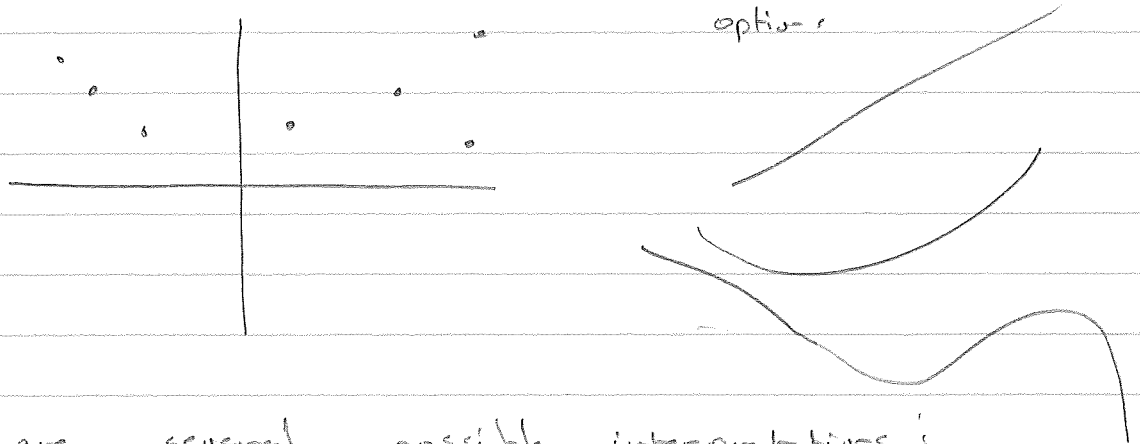
OK, so, we have that the linear regression estimator (linear in the parameters) is the same as the least squares estimator.

We can ask, when can this estimator be computed. Let's look at X

$$\begin{array}{c} \underbrace{\hspace{10em}}_p \\ \begin{bmatrix} x & x^2 & x^3 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \\ \begin{array}{l} N \\ \vdots \\ \vdots \end{array} \end{array}$$

$X^T X \in \mathbb{R}^{p \times p}$ and $(X^T X)^{-1}$ requires $X^T X$ nonsingular, i.e. $\text{rank}(X^T X) \geq p$.

You might ask, "When could the design matrix (X) be rank deficient?" Let's consider the following linear regression—



There are several possible interpretations of said data

Let's say I prefer $y = ax^3 + bx^2 + cx + d$ for some unknown $\vec{w} = [a \ b \ c \ d]^T$.

Note that nothing from the above need change except the definitions.

$$X = \begin{bmatrix} x_1^3 & x_1^2 & x_1 & 1 \\ x_2^3 & x_2^2 & x_2 & 1 \end{bmatrix} \quad \vec{w} = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

It is entirely feasible then that $p > n$ and $X^T X$ becomes rank deficient.