



Computer Science Year 2

Algorithms & Data

Estimation, Regression, Classification

Prof Alin Achim



🔥 Last time ...

- If a MVU estimator does not exist, or can not be found, the parameters can be obtained from the likelihood function.
- A Maximum Likelihood (ML) parameter estimate is found by maximising the likelihood function $p(x;\theta)$ which is essentially the probability of the data given the parameters;
- ML estimators are asymptotically efficient, as the number of observations increase and the covariance of the estimates tends to CRLB
- A major advantage of the MLE is that we can find an estimate from the given data numerically since it requires only the maximum of a known function. The Newton-Raphson iterative techniques or the Expectation-Maximisation (EM) algorithm can be used for iterative estimation of the parameters.

🔥 Asymptotic properties of the MLE

- As $N \rightarrow \infty$ we have that $\hat{\theta} \rightarrow \theta$ (*consistent estimator*)

- Moreover:

$$\begin{aligned} E(\hat{\theta}) &\rightarrow \theta \\ \text{var}(\hat{\theta}) &\rightarrow CRLB \end{aligned}$$

- The MLE is *asymptotically unbiased* and *asymptotically efficient*.

🔥 Example: Exponential distribution

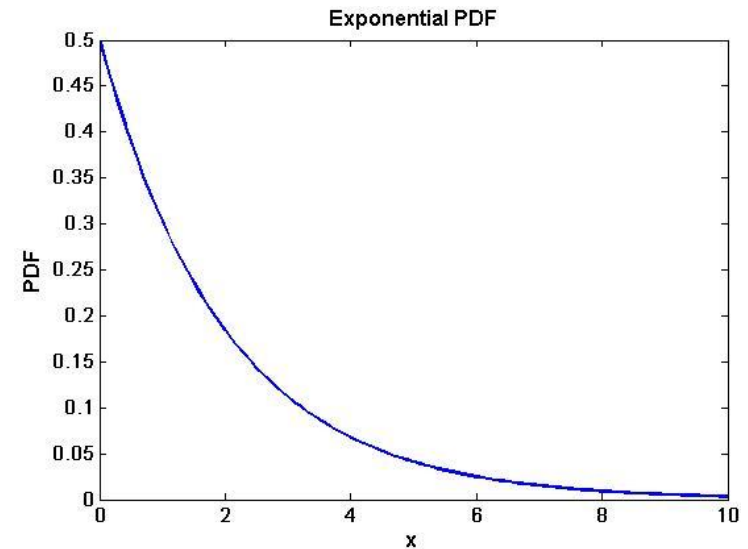
➤ Assume x_1, x_2, \dots, x_N is a random sample from an exponential distribution:-

$$f(x|\theta) = \begin{cases} \theta e^{-\theta x}, & x \geq 0 \\ 0 & x < 0 \end{cases}$$

➤ The likelihood function is

$$L(\theta) = \prod_{i=1}^N f(x_i|\theta) = \theta^N e^{-\theta \sum_{i=1}^N x_i}$$

$$l(\theta) = N \log \theta - \theta \sum_{i=1}^N x_i$$



➤ Taking the derivative and setting equal to zero:-

$$\hat{\theta} = \frac{N}{\sum_{i=1}^N x_i} = \frac{1}{\bar{x}}$$

🔥 Joint MLE for several parameters

- Often in practice, a statistical model has more than one unknown parameter
- If there are k parameters, then we have a vector parameter $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)^T$ and the PDF is written as $f(x|\boldsymbol{\theta})$.
- $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)^T$ are the values in the parameter space that jointly maximise the likelihood function
- If $l(\boldsymbol{\theta})$ is differentiable then the vector estimate satisfy k joint differential equations

$$\frac{\partial}{\partial \theta_j} l(\theta_1, \theta_2, \dots, \theta_k) = 0 \quad \text{for } j = 1, 2, \dots, k$$

🔥 Example: Normal distribution

- Assume $X \sim N(\mu, \sigma^2)$ i.e.

$$p(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[\frac{-(x - \mu)^2}{2\sigma^2} \right]$$
$$-\infty < \mu < \infty, \quad \sigma > 0$$

- The log-likelihood is

$$l(\mu, \sigma^2) = -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \sigma^2 - \frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}$$

- Solving the likelihood equations yields:

$$\frac{\partial l}{\partial \mu} = -\frac{(-2) \sum_{i=1}^N (x_i - \mu)}{2\sigma^2} = 0 \Rightarrow \hat{\mu} = \bar{x}$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{\sigma^4} \frac{\sum_{i=1}^N (x_i - \mu)^2}{2} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\sum_i (x_i - \bar{x})^2}{N}$$

Objectives

- Least Square Estimation (LSE)
- Method of Moments (MoM)



🔥 Least Squares - The intuition

- Consider (again) the multiple observations:

$$x[n] = A + w[n], \text{ where } n = 0, 1, \dots, N-1 \text{ and } w[n] \sim N(0, \sigma^2)$$

- Or more generally:

$$x[n] = f(\theta_1, \theta_2, \dots, \theta_M) + w[n], \text{ i.e. } f \text{ is a function of } M \text{ parameters}$$

- The least square estimator (LSE) minimizes:

$$J = \sum_{n=0}^{N-1} (x[n] - f(\theta_1, \theta_2, \dots, \theta_M))^2$$

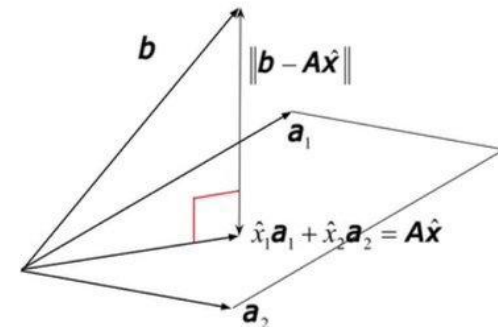
🔥 Least Squares - Properties

- LSE is widely used when estimating parameters for linear models
 - No assumptions about the data are made
 - If $w[n] \sim N(0, \sigma^2)$, LSE coincides with the MLE!
 - Geometric interpretation: the LS estimate is an orthogonal projection of the data vector onto the space defined by the independent variable.
-
- Inverse problems formulation:
 - $\mathbf{b} = \mathbf{A}\mathbf{x} + \mathbf{n}, \mathbf{n} \sim N(0, \sigma^2)$
 - $\hat{\mathbf{x}} = \min_{\mathbf{x}} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|^2$

Geometric interpretation

• $\mathbf{A}\hat{\mathbf{x}}$ is the orthogonal projection of \mathbf{b} onto $\text{range}(\mathbf{A})$

$$\Leftrightarrow \mathbf{A}^T(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = \mathbf{0} \Leftrightarrow \mathbf{A}^T \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}$$



🔥 Least Squares - A much simpler example

- Suppose we have a random sample $\{x_n\}, n = 0, \dots, N - 1$, drawn from a population with mean μ_x and standard deviation σ_x .

- We can express x_n using a linear model:

$$x_n = \mu_x + \varepsilon_n, E[\varepsilon_n] = 0 \text{ and } E[\varepsilon_n^2] = \sigma_x^2$$

- Estimate μ_x via LSE, that is minimise:

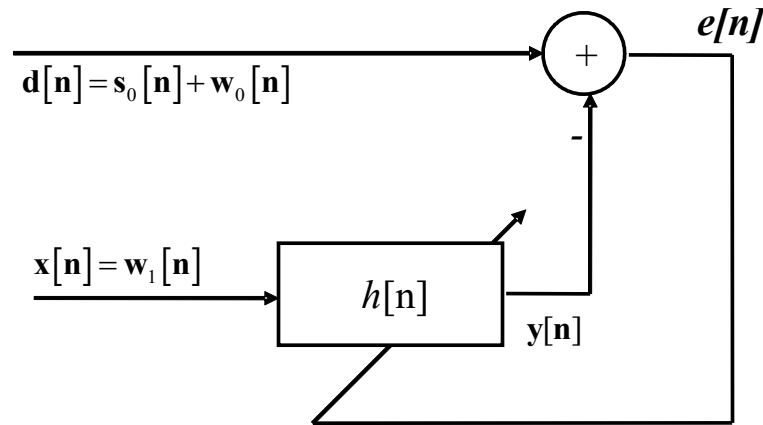
$$J(\mu_x) = \sum_{n=0}^{N-1} (x_n - \hat{\mu}_x)^2$$

$$\frac{\partial J(\mu_x)}{\partial \mu_x} = -2 \sum_{n=0}^{N-1} (x_n - \hat{\mu}_x) = 0$$

$$\Rightarrow \sum_{n=0}^{N-1} x_n = N \hat{\mu}_x \Rightarrow \hat{\mu}_x = \frac{1}{N} \sum_{n=0}^{N-1} x_n$$

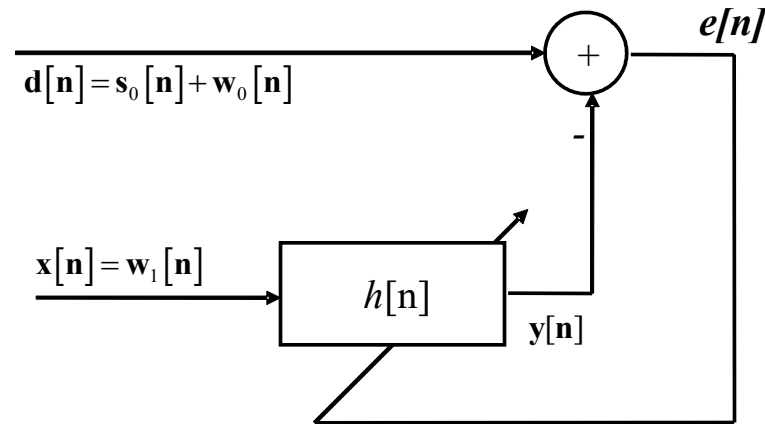
🔥 Least Squares - A signal processing example

- One of the main applications of LSE is in Adaptive Noise Cancellation (ANC).



- A signal s_0 is received with some uncorrelated noise w_0 . A second signal contains noise w_1 uncorrelated with s_0 but correlated with w_0 in some unknown way. The problem is to enhance the signal s_0 in the primary input.
- The principle of ANC is to estimate the noise in the primary input by filtering the reference noise signal with a linear filter. *How is the filter designed??*

🔥 Least Squares - A signal processing example



➤ Choose the weights at time n to minimize:

$$J[n] = \sum_{k=0}^n e^2[k] = \sum_{k=0}^n (d[k] - y[k])^2 = \sum_{k=0}^n \left(d[k] - \sum_{l=0}^{p-1} h[l]x[k-l] \right)^2$$

🔥 Method of Moments (MoM) - The intuition

- The method of moments seeks to equate the moments as implied by the underlying model of the population distribution $p(x)$ (mean, variance, skewness, kurtosis, etc) with the actual moments in the sample.

$$\text{population moment} = \text{sample moment}$$

- The population probability density function (pdf) is a function of the parameter we want to estimate, θ .

$$p(x) = f(x|\theta), \text{ e. g. } f(x|\theta) \propto e^{-\frac{x^2}{2\theta^2}}$$



Method of Moments (MoM) - The intuition

population moment = sample moment

and

$$p(x) = f(x|\theta)$$

$$1. \quad E[x^n] = \int_{-\infty}^{\infty} x^n p(x) = \int_{-\infty}^{\infty} x^n f(x|\theta) = g(\theta)$$

$$2. \quad \theta = g^{-1}(E[x^n])$$

$$3. \quad \theta = g^{-1}(n^{th} \text{ sample moment })$$

MoM - Properties

- Main advantage: extremely easy to determine and implement
- The order of the moment(s) used depends on the parameter to be estimated, i.e. population mean/maximum/correlation is estimated using sample mean/maximum/correlation, etc
- To obtain estimates of several population parameters, several moments need to be employed
- As an example, consider again the case of the 2 parameter Gaussian distribution, in which case, we use the first and second order moments, i.e. the sample mean and sample variance:

$$E[x] = g_1(x|\theta_1) = \frac{1}{N} \sum_{n=1}^N x_n \text{ and } E[x^2] = g_2(x|\theta_2) = \frac{1}{N} \sum_{n=1}^N x_n^2$$

🔥 MoM - Another simple example

- Suppose the population distribution follows an uniform distribution with unknown parameter θ :

$$p(X) = \begin{cases} \frac{1}{\theta}, & \text{for } 0 \leq X \leq \theta \\ 0, & \text{otherwise} \end{cases}$$

- Given a random sample $\{x_n\}, n = 1, \dots, N$, estimate θ .

$$1. \quad E[X] = \int_0^\theta x \frac{1}{\theta} dx = \frac{\theta}{2}$$

$$E[X^n] = g(\theta)$$

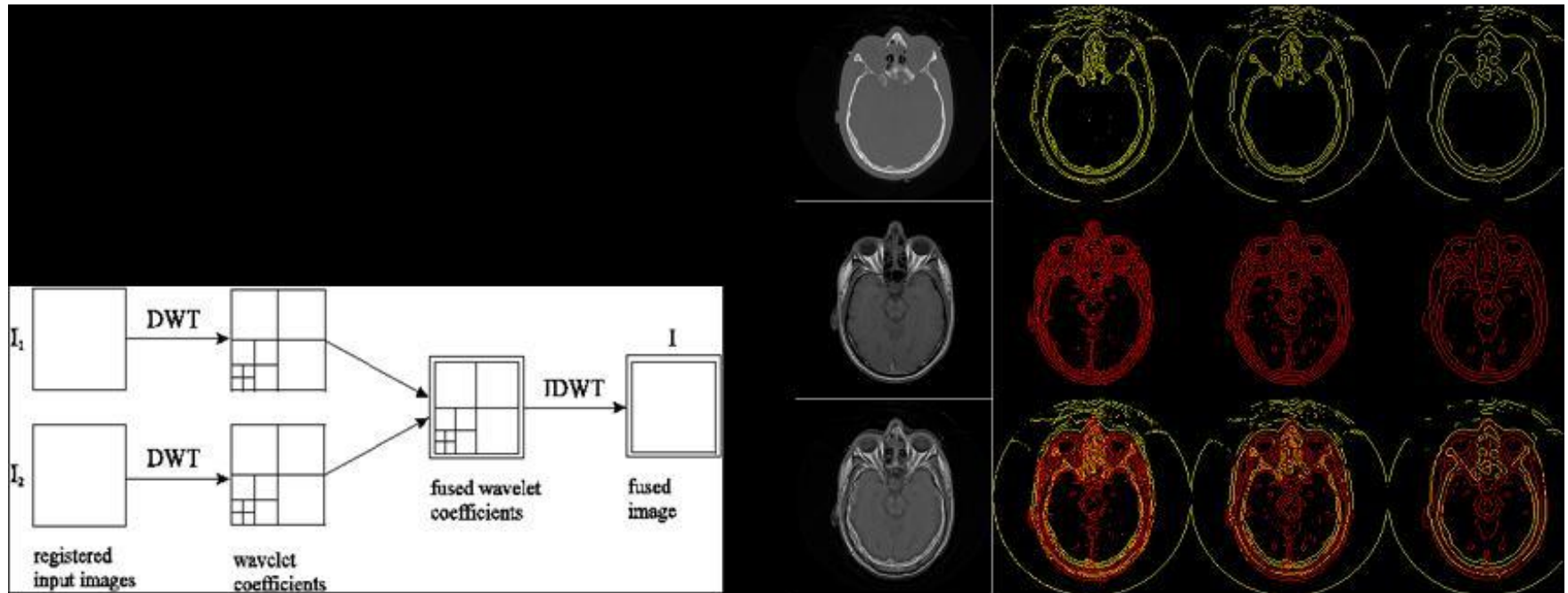
$$2. \quad \theta = 2E[X]$$

$$\theta = g^{-1}(E[X^n])$$

$$3. \quad \hat{\theta} = 2\bar{X}$$

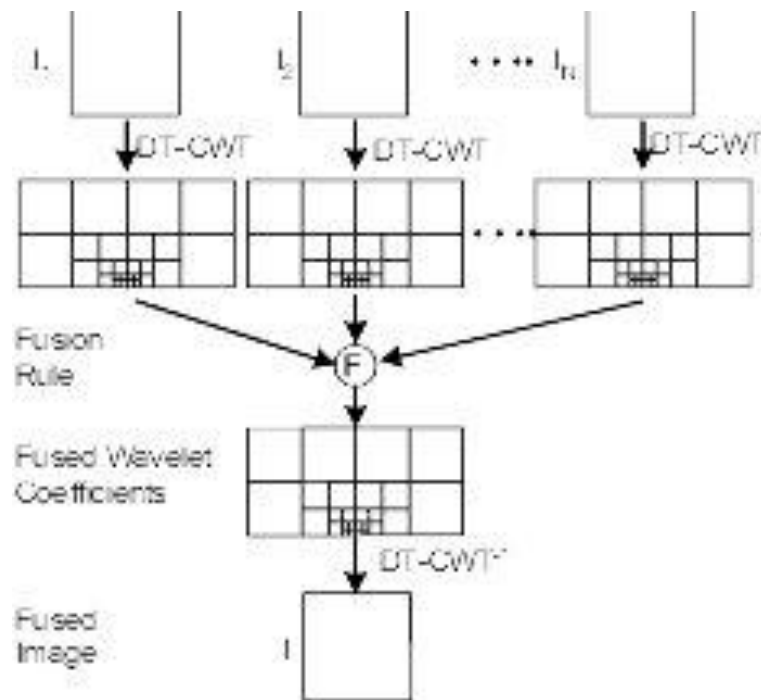
$$\hat{\theta} = g^{-1}(n^{\text{th}} \text{ sample moment})$$

🔥 Real world application - image fusion



🔥 Image fusion algorithms

- Pixel - level fusion
- Region - level fusion
- Compressive fusion



Pixel level approaches:

- Maximum selection scheme
- Average scheme
- Weighted average scheme

🔥 The FLOM weighted average (FLOM-WA) scheme

I. DTCWT of both images

II. For each subband pair

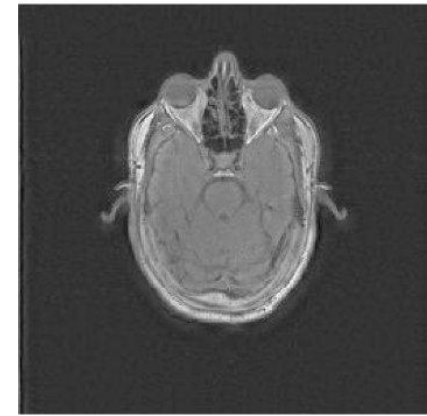
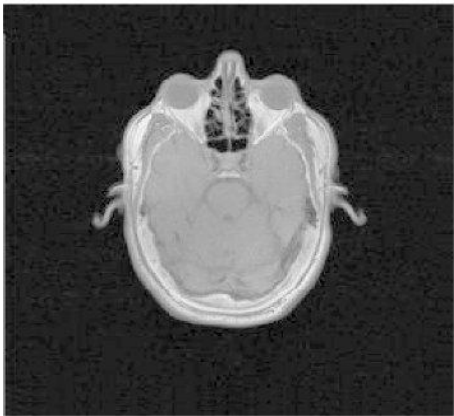
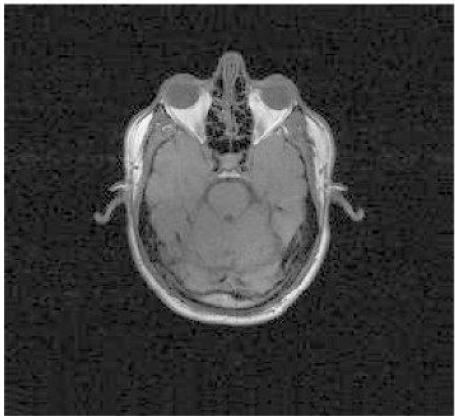
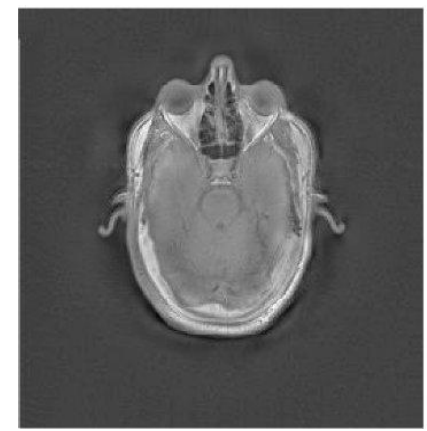
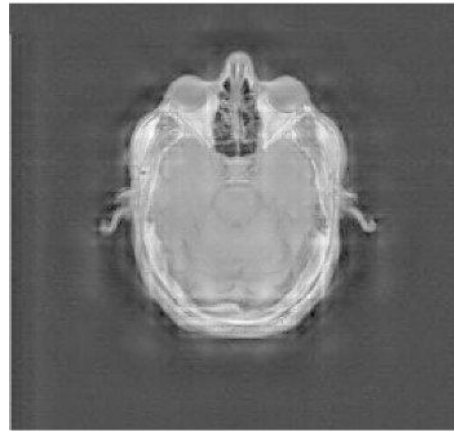
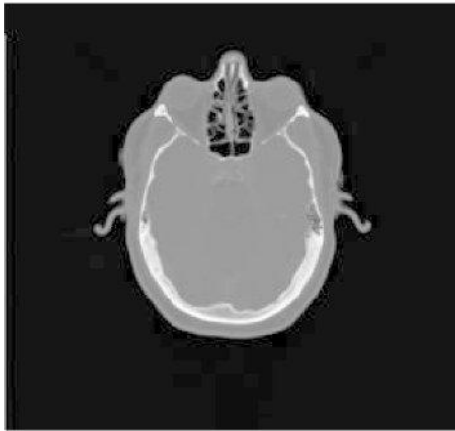
1. Estimate γ_{x_1} and γ_{x_2}
2. Compute $\text{Corr}_\alpha(X_1, X_2)$
3. Calculate the fused coefficient using

$$Y = w_1 X_1 + w_2 X_2$$

III. Average coefficients in lowpass residual

IV. IDTCWT

✦ Example WA-FLOM Fusion of MR & CT images



Summary of LSE & MoM

➤ LSE:

- ❖ minimises the sum of squares between the measurements and the model
- ❖ no assumption about the data — generally applicable estimators
- ❖ best in the context of linear models

➤ MoM:

- ❖ equate the sample with population moments
- ❖ the simplest, intuitive, works well in straightforward cases
- ❖ estimators not always with good properties (e.g. for small sample size)

