

Bayesian computation

Dr. Yoann Altmann

*B39AX – Fall 2023
Heriot-Watt University*

Plan

- Likelihood and Maximum Likelihood
- Linear regression
- Bayesian modelling
- Bayesian estimation
 - MAP estimation
 - MMSE estimation

Bayesian “vs” frequentist

- Common problems
 - Decision between hypotheses given some observations y
 - Choosing the “best” model able to explain the observations y
 - Estimation of some parameters x given some observations y
- Alternative notations
 - x for observations and θ for parameters

Bayesian “vs” frequentist

- Different interpretations
 - Frequentist: the unknown parameters are deterministic
 - Bayesian: the unknown parameters are RVs
 - Pro/cons for both approaches
- Here, we will adopt a Bayesian approach and use **prior distributions**
- **But first, let’s adopt a frequentist approach**

PMF/PDF

- Until now, we made the difference between
 - pmf $\mathbb{p}(\boldsymbol{x})$ (discrete variables)
 - pdf $f_X(\boldsymbol{x})$ (continuous variables)
- But a vector \boldsymbol{x} can also contain both discrete and continuous variables
- So the notations above can quickly become impractical

Bayes rule(s)

- y discrete, x discrete

$$\mathbb{P}_{X|Y}(x|y) = \frac{\mathbb{P}_{Y|X}(y|x)\mathbb{P}_X(x)}{\sum_{x'} \mathbb{P}_{Y|X}(y|x')\mathbb{P}_X(x')} = \frac{\mathbb{P}_{Y|X}(y|x)\mathbb{P}_X(x)}{\mathbb{P}_Y(y)}$$

- y discrete, x continuous

- $f_{X|Y}(x|y) = \frac{\mathbb{P}_{Y|X}(y|x)f_X(x)}{\int_{x'} \mathbb{P}_{Y|X}(y|x')f_X(x')dx'} = \frac{\mathbb{P}_{Y|X}(y|x)f_X(x)}{\mathbb{P}_Y(y)}$

- y continuous, x discrete

- $\mathbb{P}_{X|Y}(x|y) = \frac{\mathbb{P}_{Y|X}(y|x)\mathbb{P}_X(x)}{\sum_{x'} f_{Y|X}(y|x')\mathbb{P}_X(x')} = \frac{f_{Y|X}(y|x)\mathbb{P}_X(x)}{f_Y(y)}$

- y continuous, x continuous

- $f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x)f_X(x)}{\int_{x'} f_{Y|X}(y|x')f_X(x')dx'} = \frac{f_{Y|X}(y|x)f_X(x)}{f_Y(y)}$

Notations

- In the remainder of this chapter, we will use
 - $f_X(x)$ for the pmf/pdf of the RV x except in specific examples
 - Marginalisation will be denoted by integrals (\int) which can be replaced by sums for discrete RVs

Notations

- In the remainder of this chapter, unless stated otherwise, we will use
 - y : observations
 - x : unknown parameters of interest
 - θ : additional (hyper)-parameters
- Other classical notations (not used here)
 - x : observations
 - θ : unknown parameters of interest

Likelihood

- Goal: Estimation of some parameters x given some observations y
- **Model for noisy measurements**
- $y|x \sim f_{Y|X}(y|x)$
- Examples:

$$y = x + n$$

$$y = Ax + n$$

$$y = g(x) + n$$

Likelihood

- AWGN noise model
 - $\mathbf{y} = \mathbf{x} + \mathbf{n}, \mathbf{y} \in \mathbb{R}^N$
 - $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}) \rightarrow \mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{y}; \mathbf{x}, \sigma^2 \mathbf{I})$
- Likelihood function

$$f_{Y|X}(\mathbf{y}|\mathbf{x}) = \left(\frac{1}{2\pi\sigma^2}\right)^{N/2} e^{-\frac{\|\mathbf{y}-\mathbf{x}\|_2^2}{2\sigma^2}}$$

“how likely is the observation \mathbf{y} given some value of \mathbf{x} ”

Maximum likelihood estimator

$$\begin{aligned}\hat{x}_{MLE} &= \operatorname{argmax}_x f_{Y|X}(\mathbf{y}|\mathbf{x}) \\ &= \operatorname{argmax}_x \log \left(f_{Y|X}(\mathbf{y}|\mathbf{x}) \right)\end{aligned}$$

- Example: $y_n = x_0 + n_n, \forall n$ and $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}_N)$ (σ^2 known, x_0 unknown)

- $f_{Y|X}(\mathbf{y}|x_0) = \left(\frac{1}{2\pi\sigma^2} \right)^{N/2} e^{-\frac{\sum_{n=1}^N (y_n - x_0)^2}{2\sigma^2}}$

- $\hat{x}_{MLE} = \frac{1}{N} \sum_{n=1}^N y_n$ (sample mean)

Maximum likelihood estimator

- Some properties

- Consistency

$\hat{\mathbf{x}}_{MLE} \rightarrow \mathbf{x}$ in probability when $N \rightarrow \infty$

$$\lim_{N \rightarrow \infty} \mathbb{P}(\|\hat{\mathbf{x}}_{MLE} - \mathbf{x}_0\| > \epsilon) = 0, \forall \epsilon > 0$$

- Asymptotic Normality

$\sqrt{N}(\hat{\mathbf{x}}_{MLE} - \mathbf{x}_0) \rightarrow \mathcal{N}(0, F^{-1})$ (in distribution)

F : Fisher information matrix (advanced concept)

Rq: $\hat{\mathbf{x}}_{MLE}$ is asymptotically unbiased

Computation of the MLE

$$\hat{x}_{MLE} = \operatorname{argmax}_x f_{Y|X}(\mathbf{y}|\mathbf{x})$$

- If $f_{Y|X}(\mathbf{y}|\cdot)$ is differentiable (and \mathbf{x} continuous)
 - Find analytically the zero(s) of the gradient $\nabla f_{Y|X}(\mathbf{y}|\cdot)$ (local extrema)
- If not possible, find a solution numerically (we will discuss this later)
- Sometimes it is easier to work with $\log(f_{Y|X}(\mathbf{y}|\mathbf{x}))$.

Examples (revisited)

- AWGN noise model
 - $y_n = x_0 + n_n, \forall n = 1, \dots, N$
 - $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}) \rightarrow \mathbf{y} | x_0 \sim \mathcal{N}(\mathbf{y}; x_0 \mathbf{1}, \sigma^2 \mathbf{I})$
 - If σ^2 is known but x_0 is unknown

$$\hat{x}_{MLE} = \frac{1}{N} \sum_{n=1}^N y_n, \hat{x}_{MLE} \sim \mathcal{N}\left(x_0, \frac{\sigma^2}{N}\right)$$

$$\mathbb{E}[\hat{x}_{MLE}] = x_0$$

unbiased

Examples (revisited)

- AWGN noise model
 - $y_n = x_0 + n_n, \forall n = 1, \dots, N$
 - $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}) \rightarrow \mathbf{y} | x_0 \sim \mathcal{N}(\mathbf{y}; x_0 \mathbf{1}, \sigma^2 \mathbf{I})$
 - If σ^2 is unknown but x_0 is known

$$\hat{\sigma}_{MLE}^2 = \frac{\sum_{n=1}^N (y_n - x_0)^2}{N}$$

$$\mathbb{E}[\hat{\sigma}_{MLE}^2] = \sigma^2$$

unbiased

Examples (revisited)

- AWGN noise model

- $y_n = x_0 + n_n, \forall n = 1, \dots, N$

- $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}) \rightarrow \mathbf{y} | x_0 \sim \mathcal{N}(\mathbf{y}; x_0 \mathbf{1}, \sigma^2 \mathbf{I})$

- If (x_0, σ^2) is **unknown**

$$\hat{x}_{MLE} = \frac{1}{N} \sum_{n=1}^N y_n, \quad \hat{\sigma}_{MLE}^2 = \frac{\sum_{n=1}^N (y_n - \hat{x}_{MLE})^2}{N}$$

$$\mathbb{E}[\hat{x}_{MLE}] = x_0, \quad \mathbb{E}[\hat{\sigma}_{MLE}^2] = \frac{N-1}{N} \sigma^2$$

unbiased

biased but asymptotically unbiased

Verification using Matlab

- Demo1.m

Linear regression

- AWGN noise model (known variance)
 - $y_n = \mathbf{a}_n^T \mathbf{x} + n_n, \forall n = 1, \dots, N$
 - $\{\mathbf{a}_n\}_n \in \mathbb{R}^D$ with $D < N$
 - $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N]^T$ is a **known** $N \times D$ matrix (full-rank)
 - How to estimate \mathbf{x} via MLE?

$$f_{Y|X}(\mathbf{y}|\mathbf{x}) = \left(\frac{1}{2\pi\sigma^2} \right)^{N/2} e^{-\frac{\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}{2\sigma^2}}$$

Linear regression

- Least-square regression

$$\hat{\mathbf{x}}_{MLE} = \operatorname{argmax}_x \log \left(f_{Y|X}(\mathbf{y}|\mathbf{x}) \right)$$

$$\hat{\mathbf{x}}_{MLE} = \operatorname{argmin}_x \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

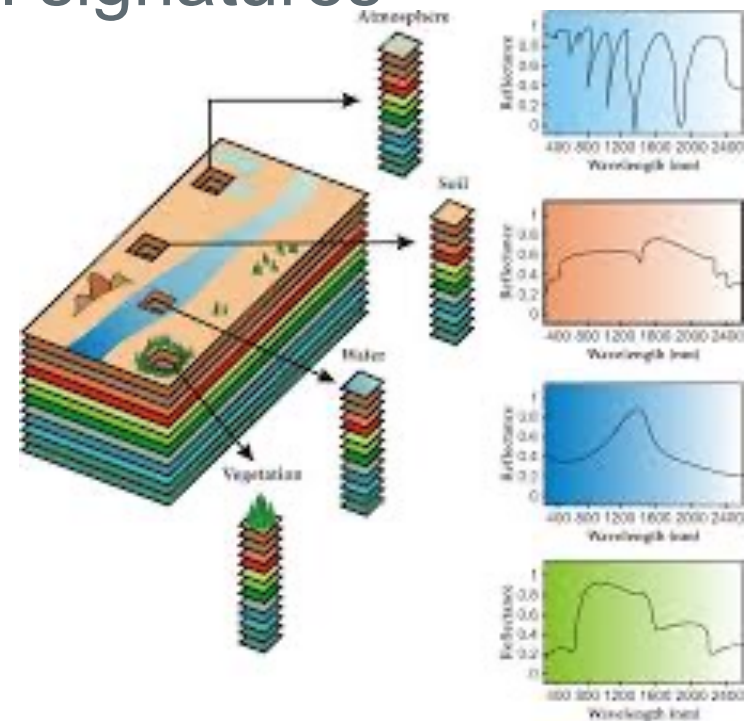
- Using the pseudo inverse

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T, D \times N \text{ matrix}$$

$$\hat{\mathbf{x}}_{MLE} = \mathbf{A}^+ \mathbf{y}$$

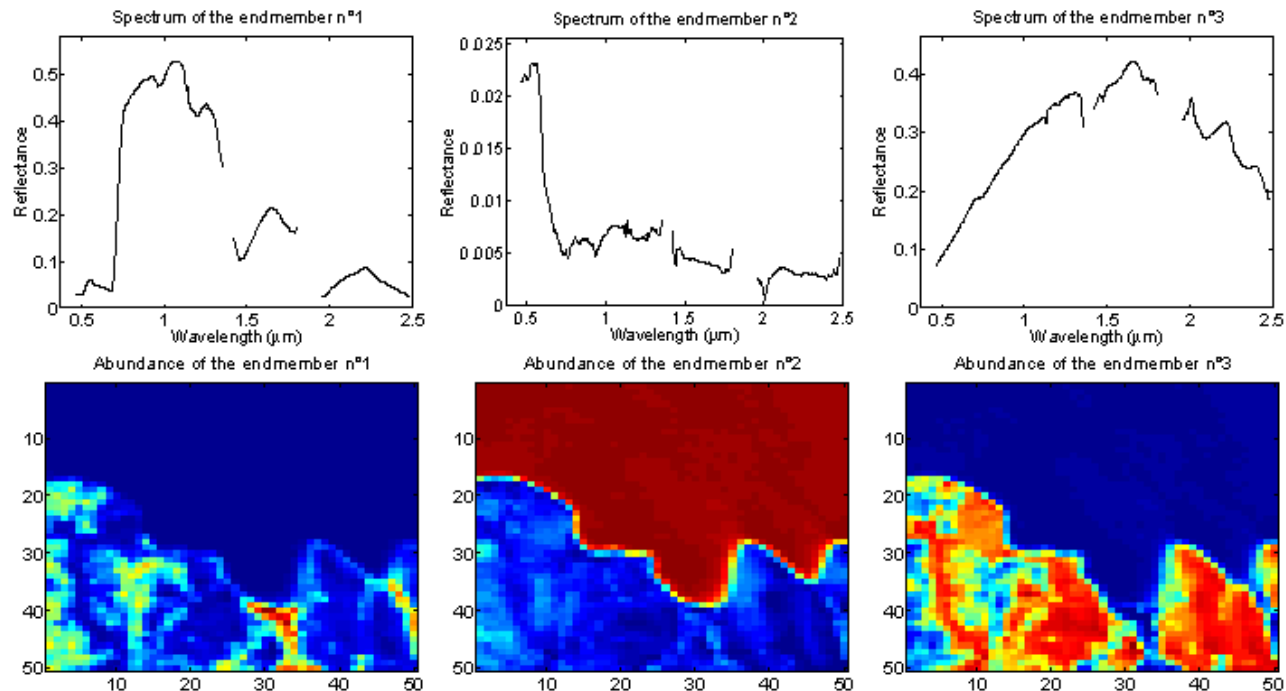
Example of linear regression

- Spectral unmixing for earth observation
 - $A = [\mathbf{a}_1, \dots, \mathbf{a}_N]^T$: spectral signatures
 - \mathbf{y} : observed spectrum
 - \mathbf{x} : material fractions or abundances



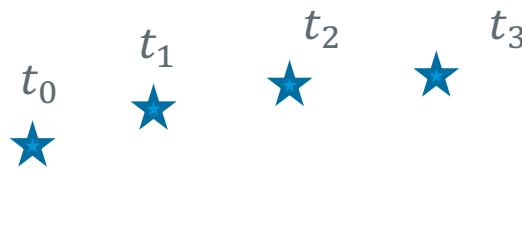
Example of linear regression

- Spectral unmixing for earth observation



Beyond MLE

- MLE: “Good” estimator if N is large
- But ... large variance if N “small”
- What if we have additional information about x_0 ?
- Can we improve the estimation of x ?
- Example: noisy trajectory



Should we
use/trust this
measurement?

Bayesian estimation

- Instead of treating x as a deterministic but unknown variable, it is seen as a RV
- Additional information expressed as a distribution, called **prior distribution** $f_X(x)$ (or $\mathbb{p}(x)$)
- $f_X(x)$: what we know about x *prior to observing y*
 - Range (e.g., positivity)
 - Mean
 - Smoothness, sparsity,...

Bayes rule revisited

- Likelihood: $f_{Y|X}(\mathbf{y}|\mathbf{x})$
- Prior distribution: $f_X(\mathbf{x})$
- Posterior distribution: $f_{X|Y}(\mathbf{x}|\mathbf{y}) = \frac{f_{Y|X}(\mathbf{y}|\mathbf{x})f_X(\mathbf{x})}{f_Y(\mathbf{y})}$

What we know about \mathbf{x} after having observed \mathbf{y}

- Evidence: $f_Y(\mathbf{y})$ (also marginal likelihood)

Bayesian estimation

$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x)f_X(x)}{f_Y(y)}$$

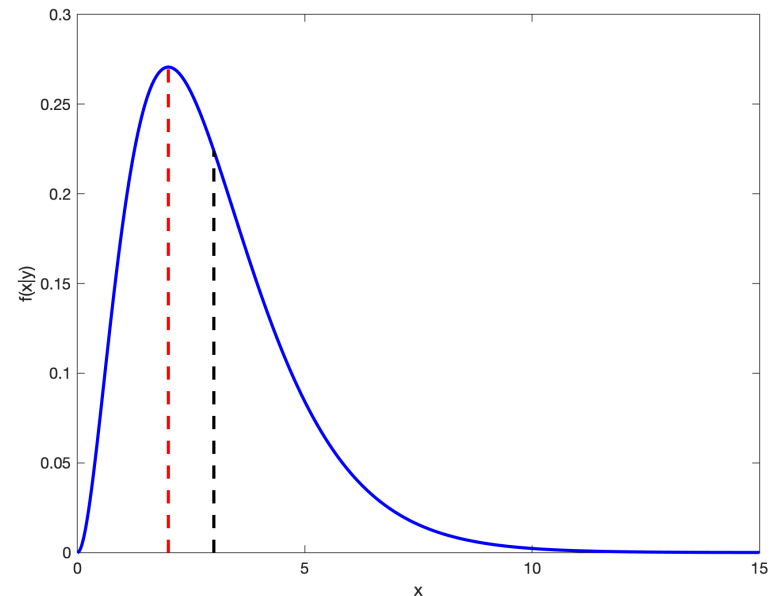
- What we know about x after having observed y
- Combines the data and our prior knowledge
- How can we use this information?

Bayesian estimation

- Example:

$f_{X|Y}(x|y)$: gamma distribution $G(x; 3, 1)$

Which point estimator or summary statistics should we use?



Maximum a posteriori (MAP)

$$\hat{\mathbf{x}}_{MLE} = \operatorname{argmax}_{\mathbf{x}} f_{Y|X}(\mathbf{y}|\mathbf{x})$$

$$\begin{aligned}\hat{\mathbf{x}}_{MAP} &= \operatorname{argmax}_{\mathbf{x}} f_{X|Y}(\mathbf{x}|\mathbf{y}) \\ &= \operatorname{argmax}_{\mathbf{x}} \log \left(f_{X|Y}(\mathbf{x}|\mathbf{y}) \right)\end{aligned}$$

$\hat{\mathbf{x}}_{MLE}$: most likely \mathbf{x} given the data \mathbf{y} only

$\hat{\mathbf{x}}_{MAP}$: most likely \mathbf{x} given the data \mathbf{y} and our prior knowledge

Maximum a posteriori (MAP)

- Important remark (for model selection):

If x is discrete and takes only a finite number of values, the MAP rule/estimator minimizes (over all decision rules) the probability of selecting an incorrect hypothesis.

Example: Bayesian model selection

- Likelihood: $y_n = x_0 + n_n, \forall n = 1, \dots, N$
 - $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I})$ with σ^2 known
 - $\mathbb{P}(x_0 = 1) = \pi, \mathbb{P}(x_0 = 2) = 1 - \pi$
- Given N observations y_1, \dots, y_N and our prior knowledge, we would like to decide between $x_0 = 1$ and $x_0 = 2$
- We decide $x_0 = 1$ if

$$\mathbb{P}(x_0 = 1|\mathbf{y}) > \mathbb{P}(x_0 = 2|\mathbf{y})$$

i.e.,

$$\mathbb{P}(x_0 = 1)f_{Y|X_0}(\mathbf{y}|x_0 = 1) > \mathbb{P}(x_0 = 2)f_{Y|X_0}(\mathbf{y}|x_0 = 2)$$

Computation of the MAP estimate

$$\hat{x}_{MAP} = \operatorname{argmax}_x f_{X|Y}(x|y)$$

- If $f_{X|Y}(\cdot | y)$ is differentiable (and x continuous)
 - Find analytically the zero(s) of the gradient $\nabla f_{X|Y}(\cdot | y)$ (local extrema)
- If not possible, find a solution numerically
 - Using optimization algorithms
 - Using simulation methods (later)
- Sometimes, it is easier to work with $\log(f_{X|Y}(x|y))$

Computation of the MAP estimate

$$\hat{\mathbf{x}}_{MAP} = \operatorname{argmax}_{\mathbf{x}} f_{X|Y}(\mathbf{x}|\mathbf{y}) = \operatorname{argmax}_{\mathbf{x}} f_{Y|X}(\mathbf{y}|\mathbf{x})f_X(\mathbf{x})$$

$$\hat{\mathbf{x}}_{MAP} = \operatorname{argmin}_{\mathbf{x}} -\log\left(f_{Y|X}(\mathbf{y}|\mathbf{x})\right) - \log(f_X(\mathbf{x}))$$

- Often easier to solve as nearly quadratic functions
- MAP estimation as the maximization of a **penalised likelihood**
- The term $-\log(f_X(\mathbf{x}))$ acts as a penalty or regularisation

Example

- AWGN noise model
 - $y_n = x_0 + n_n, \forall n = 1, \dots, N$
 - $\mathbf{n} \sim \mathcal{N}(\mathbf{n}; \mathbf{0}, \sigma^2 \mathbf{I}) \rightarrow \mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{y}; \mathbf{x}, \sigma^2 \mathbf{I})$
 - If σ^2 is known but x_0 is unknown
 - Prior distribution: $\mathcal{N}(x_0; m, s^2)$

$$\hat{x}_{MLE} = \frac{1}{N} \sum_{n=1}^N y_n, \hat{x}_{MAP} = ?$$

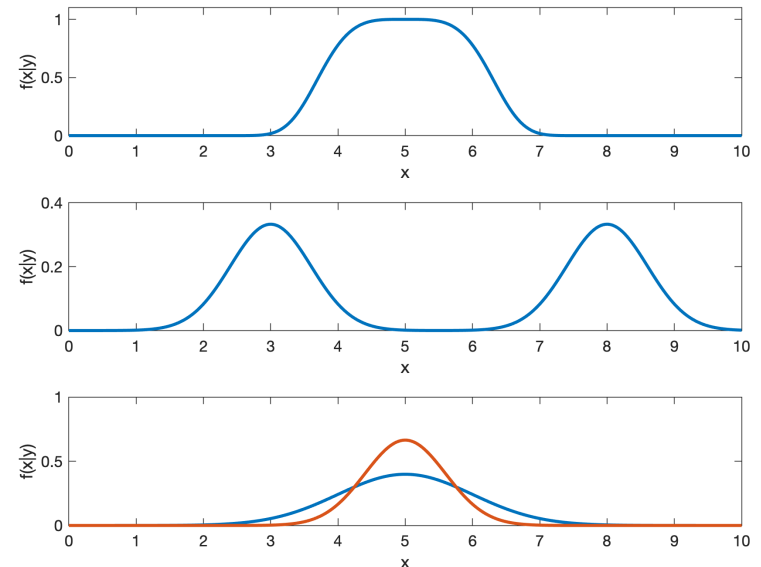
$$\hat{x}_{MAP} = \gamma^2 \left(\frac{\sum_{n=1}^N y_n}{\sigma^2} + \frac{m}{s^2} \right) \text{ with } \gamma^2 = \frac{\sigma^2 s^2}{\sigma^2 + N s^2}$$

Verification using Matlab

- Demo2.m

Limitations of the MAP estimator

- Useful estimator, obtained via optimization
- But...might be difficult to compute
 - Flat gradient / numerical errors
- Might not be unique
 - Multimodal/flat distribution
- Provides limited information
 - Other solutions almost as likely?



Alternative estimator

- Posterior mean or minimum mean square error (MMSE) estimator

$$\hat{\mathbf{x}}_{MMSE} = \mathbb{E}_{f(\mathbf{x}|\mathbf{y})}[\mathbf{x}] = \int \mathbf{x} f(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

under weak conditions on $f(\mathbf{x}|\mathbf{y})$ (e.g., existence of mean and variance)

Properties of the MMSE estimator

By definition, it minimizes the mean square error

$$MSE = \mathbb{E}[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x})]$$

where the expectation is taken over \mathbf{x} and \mathbf{y} .

The MMSE estimator is **unbiased** and asymptotically Gaussian

Computation of the MMSE estimator

- Analytically,

$$\hat{x}_{MMSE} = \mathbb{E}_{f(x|y)}[x] = \int x f(x|y) dx.$$

- Often the integral is not tractable
 - One possibility is to approximate the expectations using simulation tools (Monte Carlo sampling)
 - Another possibility is to simplify the estimation by imposing additional constraints.

Principle of Monte Carlo sampling

To approximate

$$\mathbb{E}[g(\mathbf{x})] = \int g(\mathbf{x})f(\mathbf{x}|\mathbf{y})d\mathbf{x}$$

Monte Carlo sampling consists of generating a large number N_{MC} of random variables $\mathbf{x}_1, \dots, \mathbf{x}_{N_{MC}}$ from the distribution $f(\mathbf{x}|\mathbf{y})$. We then obtain

$$\frac{1}{N_{MC}} \sum_{n=1}^{N_{MC}} g(\mathbf{x}_n) \approx \mathbb{E}[g(\mathbf{x})]$$

Efficient methods to sample for complex distributions $f(\mathbf{x}|\mathbf{y})$ is still an open domain of research !

Principle of Monte Carlo sampling

- Estimation of mean and variance of arbitrary distribution
- Demo3.m

Linear MMSE estimator

- If we enforce $\hat{\mathbf{x}}_{MMSE} = \mathbf{W}\mathbf{y} + \mathbf{b}$ (linear function of \mathbf{y}), the problem

$$\min_{\hat{\mathbf{x}}} MSE, \quad \text{s.t. } \hat{\mathbf{x}}_{MMSE} = \mathbf{W}\mathbf{y} + \mathbf{b}$$

can become easier than the original problem, with a negligible performance degradation.

MAP vs MMSE estimation

- For discrete RVs, the MMSE estimator might not be meaningful

Ex1: $x \in \{0,1\}$, $\mathbb{E}[x] \in [0,1]$

can make sense.

Ex2: $x \in \{cat, dog\}$, $\mathbb{E}[x]$ does not make sense

In this case, the MAP estimator is more adapted

MAP vs MMSE estimation

- For continuous variables, MAP and MMSE estimation can be complementary
 - MAP: most probable solution
 - MMSE: minimizes the MSE
 - Flat, skewed or multimodal distributions
 - Sometimes, the posterior covariance can also be computed (e.g., using MC sampling)
 - If the mean and mode of the posterior distribution coincide, the two estimators are the same
 - Example: Gaussian distribution

Summary

- Likelihood and Maximum Likelihood
- Linear regression
- Bayesian modelling
- Bayesian estimation
 - MAP estimation
 - MMSE estimation
- Next chapter: Introduction to Information Theory