

# Lecture 21, Nov 29, 2023

## Extracting Estimates from PDFs

### Maximum Likelihood (ML)

- This method is often used when  $\mathbf{x}$  is an unknown constant parameter without a known probabilistic description, i.e. we have no prior information about  $\mathbf{x}$ 
  - e.g. in Bayesian estimation, we had a prior (prediction) for  $\mathbf{x}$ , but here we are assuming no knowledge of that
- For a given observation  $\mathbf{y}$  and observation model  $f(\mathbf{y}|\mathbf{x})$ , the method seeks a value of  $\mathbf{x}$  that maximizes the likelihood of observing  $\mathbf{y}$ , i.e.  $\hat{\mathbf{x}}^{ML} = \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmax}} f(\mathbf{y}|\mathbf{x})$ 
  - $f(\mathbf{y}|\mathbf{x})$  as a function of  $\mathbf{x}$  is the *likelihood function*
  - $\mathbf{x}$  is a parameter of the observation model; e.g. the model can be a Gaussian, and  $\mathbf{x}$  may denote its mean or variance, etc
- Example: Consider two measurements of a scalar quantity  $x \in \mathbb{R}$ :  $y_1 = x + w_1, y_2 = x + w_2$  where  $w_1, w_2 \sim \mathcal{N}(0, 1)$ 
  - Note  $\mathcal{N}(\mu, \sigma)$  denotes a Gaussian with mean  $\mu$  and variance  $\sigma$
  - $f(w_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{w_i^2}{2}\right)$
  - We can consider  $w_1, w_2$  as additive noise parameters; this essentially makes  $y_i \sim \mathcal{N}(x, 1)$ 
    - \* Note formally we would use a change of variables:  $y_i = x + w_i \implies w_i = y_i - x$
    - \* Now we can just substitute  $w_i$  into the Gaussian equation since we have a linear relationship
  - $y_1, y_2$  are conditionally independent on  $x$ , so  $f(y_1, y_2|x) = f(y_1|x)f(y_2|x) = \frac{1}{2\pi} e^{-\frac{1}{2}((y_1-x)^2 + (y_2-x)^2)}$
  - This is now an unconstrained optimization problem; we can differentiate with respect to  $x$  and set this to 0
  - We get  $(y_1 - \hat{x}) + (y_2 - \hat{x}) = 0 \implies \hat{x} = \frac{y_1 + y_2}{2}$ , which is just the average
- Suppose we generalize the last example to a collection of measurements  $\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{w}$  where  $\mathbf{z}, \mathbf{w} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^n$  and  $m > n$ ; as above  $w_i \sim (0, 1)$  are independent
  - Let  $\mathbf{H} = \begin{bmatrix} \mathbf{h}_1^T \\ \vdots \\ \mathbf{h}_m^T \end{bmatrix}$  where  $\mathbf{h}_i^T = [h_{i1} \ \dots \ h_{in}]$
  - Then  $z_i = \mathbf{h}_i^T \mathbf{x} + w_i$
  - As before  $f(\mathbf{z}|\mathbf{x}) \propto \exp\left(-\frac{1}{2}((z_1 - \mathbf{h}_1^T \mathbf{x})^2 + \dots + (z_m - \mathbf{h}_m^T \mathbf{x})^2)\right)$
  - Differentiating with respect to each  $x_j$  we have  $(z_1 - \mathbf{h}_1^T \hat{\mathbf{x}})h_{1j} + \dots + (z_m - \mathbf{h}_m^T \hat{\mathbf{x}})h_{mj} = [h_{1j} \ \dots \ h_{mj}] (\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}) = 0$
  - With all the rows, we get  $\mathbf{H}^T (\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}) = \mathbf{0} \implies \hat{\mathbf{x}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{z}$ , which is the least squares solution
    - \* Note we can write  $\mathbf{w}(\mathbf{x}) = \mathbf{z} - \mathbf{H}\mathbf{x}$ , so  $\mathbf{w}$  is some error term; then  $\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathbf{w}^T \mathbf{w}$
    - \* If not all the errors have the same variance, then we have weighted least squares

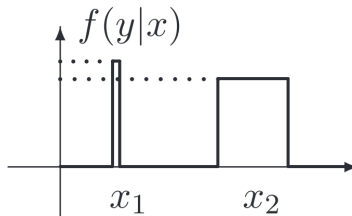


Figure 1: Undesirable case of maximum likelihood.

- Limitations of ML:
  - In general ML is more sensitive to outliers and modelling error
  - The maximum of the distribution may not always be what we want – we may lose robustness
    - \* In the example above, ML will give  $x_1$  if there are measurements on it, which is very sensitive to changes in the data or model – small variations in the model might cause  $x_1$  to have a likelihood of zero instead
    - \* Choosing  $x_2$  is more robust; since the distribution is wider, we're less sensitive to changes in the data or model
    - \* Outliers that happen to line up with a peak can give us an incorrect estimate
  - We might also have prior knowledge about  $x$  (i.e. its PDF), which ML cannot incorporate

### Maximum a Posteriori (MAP)

- If we have a PDF for  $\mathbf{x}$ , we can use MAP
- From Bayes's theorem:  $f(\mathbf{x}|\mathbf{y}) = \frac{\mathbf{f}(\mathbf{y}|\mathbf{x})\mathbf{f}(\mathbf{x})}{\mathbf{f}(\mathbf{y})}$
- With MAP, we have  $\hat{\mathbf{x}}^{MAP} = \underset{\mathbf{x}}{\operatorname{argmax}} f(\mathbf{y}|\mathbf{x})f(\mathbf{x})$ 
  - We want to maximize the choice of the parameter that makes both the observations and the parameter itself most likely
- If  $f(\mathbf{x})$  is constant, then  $\hat{\mathbf{x}}^{MAP} = \hat{\mathbf{x}}^{ML}$
- As with ML, we are still maximizing a function over  $\mathbf{x}$ , so the same sensitivity to outliers and modelling error still applies
- Example: consider the scalar observation  $y = x + w$ , where  $w \sim \mathcal{N}(0, 1)$ ,  $x \sim \mathcal{N}(\bar{x}, \sigma_x^2)$  and  $x, w$  independent
  - $f(x) \propto \exp\left(-\frac{1}{2} \frac{(x - \bar{x})^2}{\sigma_x^2}\right)$
  - $f(y|x) \propto \exp\left(-\frac{1}{2}(y - x)^2\right)$
  - $f(y|x)f(x) \propto \exp\left(-\frac{1}{2} \left(\frac{(x - \bar{x})^2}{\sigma_x^2} + (y - x)^2\right)\right)$
  - Differentiating with respect to  $x$  and setting to zero gives the following solution:
  - $\hat{x}^{MAP} = \frac{1}{1 + \sigma_x^2} \bar{x} + \frac{\sigma_x^2}{1 + \sigma_x^2} y$ 
    - \* Notice that this is a weighted sum between the mean of the prior distribution and the new measurement
  - Consider the extreme cases:
    - \*  $\sigma_x^2 = 0 \implies \hat{x}^{MAP} = \bar{x}$  (if we're certain about  $x$  before any measurements, we just get the max of the prior)
    - \*  $\sigma_x^2 \rightarrow \infty \implies \hat{x}^{MAP} = y$  (if we're uncertain about  $x$ , we just get the new measurement; note this is the same as maximum likelihood)
- This is most often used in state estimation

### Minimum Mean Squared Error (MMSE)

- The MMSE is the a posteriori estimate that minimizes the mean squared error
- $\hat{\mathbf{x}}^{MMSE} = \underset{\hat{\mathbf{x}}}{\operatorname{argmin}} E_{x|y} [(\mathbf{x} - \hat{\mathbf{x}})^T (\mathbf{x} - \hat{\mathbf{x}})]$ 
  - Expand and differentiate with respect to  $\hat{\mathbf{x}}$ :  $2\hat{\mathbf{x}} - 2E[x|y] = 0 \implies \hat{\mathbf{x}} = E[x|y]$
  - The MMSE estimate is the expected value of  $\mathbf{x}$  conditioned on  $\mathbf{y}$
- While MAP is the maximum of the posterior, MMSE is the mean of the posterior
- Note we did not constrain  $\hat{\mathbf{x}}$  in our minimization, but for some applications we might want to introduce constraints
  - e.g. for a discrete random variable with sample space  $\mathcal{X}$ , we need to constrain the minimization to  $\hat{\mathbf{x}} \in \mathcal{X}$

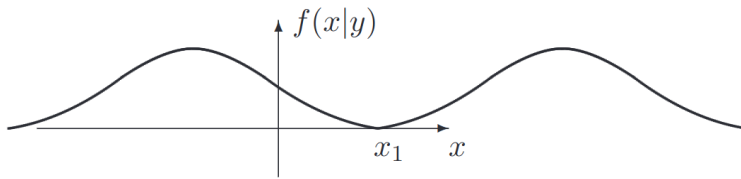


Figure 2: Undesirable case of MMSE.

- MMSE always takes the whole probability mass into consideration, whereas MAP and ML just pick the maximum probability – in some cases, this is desirable, while in other cases it is not
  - Consider the bimodal distribution of  $f(x|y)$  above; MMSE would give  $x_1$ , but the probability of having  $x$  actually being near  $x_1$  is zero
    - \* MAP would have picked one of the two peaks
  - On the other hand, the MMSE is typically more robust to modelling errors and outliers, since it is not as sensitive to sharp peaks