## Lecture 2 | Parameter Estimation

How do we find a prob. dist. for a r.v. $X$?

Three Steps:

1) Choose a parametric model (eg. Gaussian) call the parameters $\theta$.

2) Collect a set of observations (samples) from $X$.

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

we assume $x_i$'s are <u>independently</u> & <u>identically distributed</u> (iid)

3) <u>Maximum Likelihood principle</u>

"The optimal parameter $\theta^*$ is that which maximizes the probability (likelihood) of observing the training data $D$."

<u>ML estimate (MLE)</u>

$$\theta^* = \underset{\theta}{\text{argmax}} \; p(D|\theta)$$

↖ likelihood of the data $D$ w.r.t. parameter $\theta$. <u>likelihood function</u>

※ $D$ is known, so $p(D|\theta)$ is a function of $\theta$. It is not a probability distribution. It doesn't have the same shape as the pdf.

$$\theta^* = \underset{\theta}{\text{argmax}} \; \log p(D|\theta)$$

$\underbrace{\qquad\qquad}$ $\ell(\theta)$ = log-likelihood function (LL)

$$\theta^* = \underset{\theta}{\text{argmin}} \; -\log p(D|\theta)$$

$\underbrace{\qquad\qquad}$ negative log-likelihood (NLL); loss

---

$\log$ = natural logarithm ($\ln$), log base $e$.

<u>Data log likelihood</u>

$$\ell(\theta) = \log p(D|\theta)$$
$$= \log \prod_{i=1}^{N} p(x_i|\theta)$$
$$= \sum_{i=1}^{N} \log p(x_i|\theta)$$
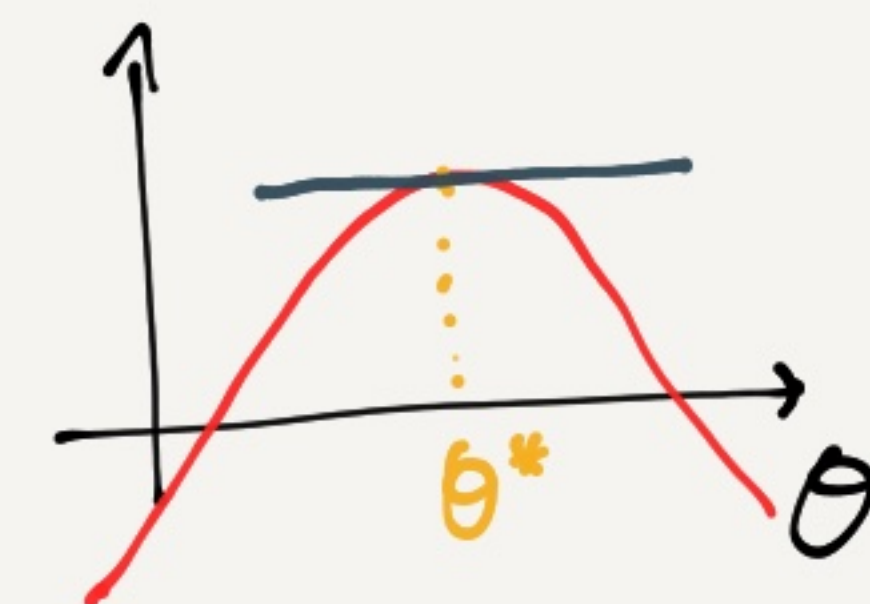
⎫ independence assumption

$\log(ab) = \log a + \log b$

<u>To get the ML solution</u>

If $\theta$ is a scalar, at local maximum:

1) $\dfrac{\partial}{\partial \theta} \log p(D|\theta) = 0$ at $\theta^*$

2) $\dfrac{\partial^2}{\partial \theta^2} \log p(D|\theta) < 0$ (at the max, it's concave)

3) check the boundary conditions of $\theta$ (i.e. it's a valid parameter)



If $\theta$ is a vector...

1) $\nabla_\theta \ell(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \ell(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} \ell(\theta) \end{bmatrix} = 0$

↖ gradient

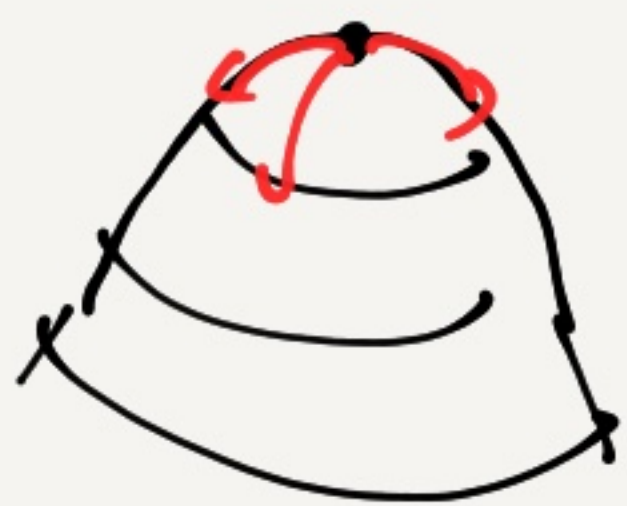2) $\nabla_\theta^2 \ell(\theta) < 0$ (negative definite)

↖ Hessian

$$\nabla^2 \ell(\theta) = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} & \cdots & \frac{\partial^2}{\partial \theta_1 \partial \theta_p} \\ & \vdots & \\ \frac{\partial^2}{\partial \theta_p \partial \theta_1} & & \frac{\partial^2}{\partial \theta_p^2} \end{bmatrix}$$

$H < 0$ : $H$ is negative definite: $\Theta^T H \Theta < 0$, $\forall \Theta$. "mountain"

$H > 0$ : $H$ is positive definite:

$\Theta^T H \Theta > 0 \quad \forall \Theta$. "bowl"

---

## Example: Bernoulli

$\Theta = \pi$ , $0 \leq \pi \leq 1$ , $x = \{0, 1\}$

$\log a^b = b \log a$

### log-likelihood

$l(\Theta) = \sum_{i=1}^{N} \log p(x_i | \Theta) = \sum_{i=1}^{N} \log \left( \pi^{x_i} (1-\pi)^{1-x_i} \right)$

$= \sum_{i=1}^{N} \left[ x_i \log \pi + (1-x_i) \log(1-\pi) \right]$

$= \left( \sum_{i=1}^{N} x_i \right) \log \pi + \left[ \sum_{i=1}^{N} (1-x_i) \right] \log(1-\pi)$

$\underbrace{\qquad}_{\text{\# of 1s}} \qquad \underbrace{\qquad}_{\text{\# of 0s}}$

"sufficient statistic" — $l(\Theta)$ only depends on the data through this suff. statistic.

$m = \sum_{i=1}^{N} x_i$

$l(\Theta) = m \log \pi + (N-m) \log(1-\pi)$

---

Solve for $\Theta^*$ : compute the derivative & set to 0.

1) $\frac{\partial}{\partial \pi} l(\Theta) = \frac{m}{\pi} + (N-m) \frac{1}{1-\pi}(-1) = 0 \qquad \Big\downarrow \pi(1-\pi)$

$m(1-\pi) - (N-m)\pi = 0$

$m - m\pi - N\pi + m\pi = 0$

$\boxed{\hat{\pi} = \frac{m}{N} = \frac{1}{N}\sum_{i=1}^{N} x_i}$ ← "fraction of 1s observed" (sample mean)

2) $\frac{\partial^2}{\partial \pi^2} l(\Theta) = \frac{\partial}{\partial \pi}\left( \frac{\partial}{\partial \pi} l(\Theta) \right) = \frac{\partial}{\partial \pi}\left( \frac{m}{\pi} + (m-N)\frac{1}{1-\pi} \right)$ $\qquad \frac{\partial}{\partial x}\frac{1}{x} = \frac{-1}{x^2}$

$= -\frac{m}{\pi^2} - \frac{(m-N)}{(1-\pi)^2}(-1)$

$= -\underbrace{\frac{m}{\pi^2}}_{-} - \underbrace{\frac{(N-m)}{(1-\pi)^2}}_{-} < 0 \checkmark$

3) boundary condition:

$0 \leq m \leq N$

$\Rightarrow 0 \leq \pi \leq 1 \checkmark$

## Example: Gaussian

$\theta = \mu \qquad (\sigma^2 \text{ known})$

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

$$= \sum_i \log \left[ \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} \right]$$

$$= \sum_i -\frac{1}{2}\log 2\pi - \frac{1}{2}\log \sigma^2 - \frac{1}{2\sigma^2}(x_i - \mu)^2$$

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

what are the suff. stats.? $\left\{ \sum_{i=1}^{N} x_i, \sum_{i=1}^{N} x_i^2 \right\}$

## Solve for $\mu$

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

$$\sum_i x_i - \sum_i \mu = 0$$

$$N\mu = \sum_i x_i \implies \boxed{\hat{\mu} = \frac{1}{N}\sum_{i=1}^{N} x_i} \quad \text{"sample mean"}$$

$\theta = \sigma^2 \quad (\mu \text{ known})$

$$\frac{\partial}{\partial \sigma^2} l(\theta) = -\frac{N}{2}\frac{1}{\sigma^2} - \frac{1}{2}\left[\sum_i (x_i - \mu)^2\right]\frac{(-1)}{\sigma^4} = 0 \quad \Bigg\} \sigma^4 \cdot 2$$

$$= -N\sigma^2 + \sum_{i=1}^{N}(x_i - \mu)^2 = 0$$

$$\implies \boxed{\hat{\sigma}^2 = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2} \quad \text{"sample variance"}$$

## Multivariate Gaussian

$$\boxed{\begin{array}{l} \hat{\mu} = \frac{1}{N}\sum_{i=1}^{N} x_i \\[2mm] \hat{\Sigma} = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu)(x_i - \mu)^T \end{array}}$$

see the tutorial next week.

# Estimators

The estimate is a number: $\mu_{ML} = \frac{1}{N} \sum_i x_i$

(value) (values)

The estimator is a r.v. over many possible datasets.

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} X_i$$

estimator (also a r.v.)   r.v. for each sample

$X_i \sim p(x_i | \theta)$ (the true distribution)

Since the estimator is a r.v., we can calculate __mean__ & __variance__. Hence we can quantify how __good__ the estimator is.

## Bias & Variance : $\hat{\theta} = f(X_1, \dots, X_N)$

1) Will it converge to the true value of $\theta$?

$$Bias(\hat{\theta}) = E_{X_1, \dots, X_N}[\hat{\theta} - \theta] = E(\hat{\theta}) - \theta$$

estimator   true value

Measures expressiveness: if the bias is non-zero, we can never get the true parameter (even $\infty$ samples)

2) How long will it take to converge? (How many samples do we need?)

$$var(\hat{\theta}) = E_{X_1 \dots X_N}\left[(\hat{\theta} - E\hat{\theta})^2\right]$$

measuring the uncertainty / variability.

---

# Example: Gaussian

Estimator: $\hat{\mu} = \frac{1}{N} \sum_i X_i$ , $X_i \sim N(\mu, \sigma^2)$

Mean $E_{X_1 \dots X_N}[\hat{\mu}] = E_{X_1 \dots X_N}\left[\frac{1}{N} \sum_{i=1}^{N} X_i\right] = \frac{1}{N} \sum_{i=1}^{N} E_{X_i}[X_i] = \frac{N\mu}{N} = \mu$

$$\Rightarrow \boxed{Bias(\hat{\mu}) = 0}$$

Var: $var(\hat{\mu}) = E_{X_1 \dots X_N}\left[(\hat{\mu} - \mu)^2\right] = E_{X_1 \dots X_N}\left[\left(\frac{1}{N}\sum_i X_i - \mu\right)^2\right]$

$\frac{1}{N}$ , $\frac{1}{N^2}$

$$= \frac{1}{N^2} E\left[\left(\sum_{i=1}^{N}(X_i - \mu)\right)^2\right]$$

$$= \frac{1}{N^2} E\left[\sum_{i=1}^{N}\sum_{j=1}^{N}(X_i - \mu)(X_j - \mu)\right]$$

$\begin{cases} i = j \Rightarrow \sigma^2 \\ i \neq j \Rightarrow 0 \end{cases}$

$$= \frac{1}{N^2}\left[\sum_{i=j} \sigma^2 + \sum_{i \neq j} 0\right] =$$

N    0

$$\boxed{var(\hat{\mu}) = \frac{1}{N}\sigma^2}$$

variance converges to zero as $N \Rightarrow \infty$

standard deviation bars

---

## for variance (PS 2-12)

$$E(\hat{\sigma}^2) = \frac{N-1}{N}\sigma^2 \Rightarrow Bias(\hat{\sigma}^2) = -\frac{1}{N}\sigma^2$$

To make it unbiased: $\hat{\sigma}^2 = \frac{N}{N-1}\hat{\sigma}^2 = \frac{N}{N-1}\frac{1}{N}\sum_i(X_i - \mu)^2 = \boxed{\frac{1}{N-1}\sum_i(X_i - \mu)^2}$
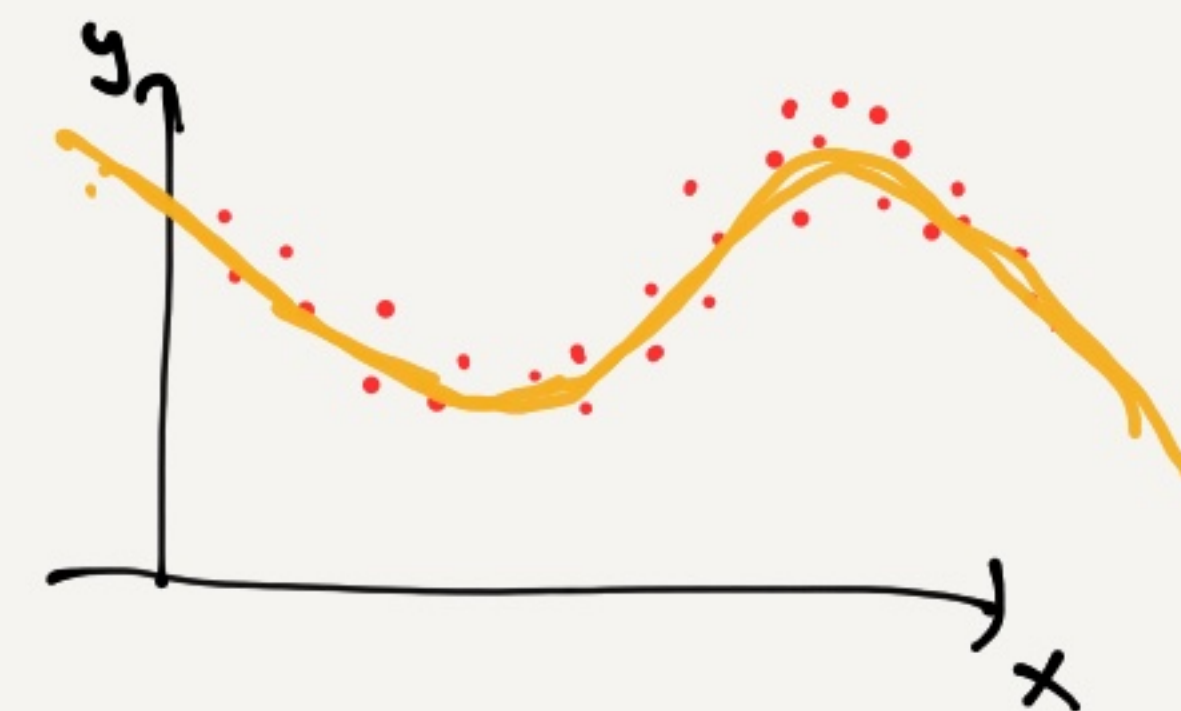
# Important Asymptotic Properties of MLE

1) **consistent** – asymptotically unbiased: as $N \to \infty$, the estimated value converges (in probability) to the true value.

2) **efficient** – achieves the Cramér-Rao Lower Bound (CRLB) as $N \to \infty$.
CRLB is a theoretical bound on the variance of $\underset{\wedge}{\text{any}}$ estimator for a particular $p(x|\theta)$.
  unbiased
(no unbiased estimator can get lower variance).

---

## MLE for regression (supervised learning)

$$D = \{(x_i, y_i)\}$$



$x \in \mathbb{R}$ input

polynomial function: $f(x, \theta) = \sum_{d=0}^{K} \theta_d x^d = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots \theta_K x^K$

$$= \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_K \end{bmatrix}^T \begin{bmatrix} 1 \\ x \\ \vdots \\ x^K \end{bmatrix} = \phi(x)^T \theta$$

$\underbrace{\quad}_{\theta} \quad \underbrace{\quad}_{\phi(x)}$

<span style="color:red">linear function of the parameters $\theta$.</span>

Observe a **noisy** output $y_i$ for a given $x_i$:

$$y_i = f(x_i, \theta) + \epsilon_i$$

<span style="color:red">$\epsilon_i \sim N(0, \sigma^2)$</span>   iid Gaussian noise.

**pdf of $y_i$**
$$p(y_i | x_i, \theta) = N(y_i | f(x_i, \theta), \sigma^2)$$

## Estimate $\theta$ using MLE

$$\hat{\theta} = \underset{\theta}{\arg\max} \sum_{i=1}^{N} \log p(y_i | x_i, \theta)$$

$\vdots$

<span style="color:red">Least-squares formulation</span>

$$= \underset{\theta}{\arg\min} \sum_i (y_i - f(x_i, \theta))^2$$

$\vdots$

$$= \underset{\theta}{\arg\min} \| y - \Phi^T \theta \|^2, \quad \Phi = \begin{bmatrix} \phi(x_1) \dots \phi(x_N) \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$\vdots$

$$\boxed{\hat{\theta} = (\Phi \Phi^T)^{-1} \Phi y}$$ <span style="color:red">←</span>

# Notes:

1) ML is more general than L.S.

2) assumptions are explicit

   i) Gaussian additive noise

   ii) iid samples (iid noise)

   iii) $\mu = 0$, $\sigma^2$ variance

3) ML can describe other LS formulations

   i) weighted LS    (PS 2-8)

   ii) regularized LS    (Lecture 3)

   iii) $L_p$ Norm    (PS 2-9)