

MLE part I:

Example: Biased coin.

Assume we have a biased coin whose probability of head is unknown.

$$\begin{cases} p(X=H) = \theta \\ p(X=T) = 1-\theta. \end{cases}$$

we observe the following (iid): HHHHT.

Problem: estimate θ .

Intuitively, $\hat{\theta} = 4/5$.

supposed $\theta = 0.9$ vs $\theta = 0.1$:

$$\begin{aligned} \text{If } \theta = 0.9, p(\text{HHHHHT} \mid \theta = 0.9) &= \theta^4(1-\theta)^1 \\ &= (0.9)^4(0.1)^1 \end{aligned}$$

$$\text{If } \theta = 0.1, p(\text{HHHHHT} \mid \theta = 0.1) = (0.1)^4(0.9)^1$$

Likelihood Function (L):

$$L(\theta) = p(\text{HHHHHT} \mid \theta) = \theta^4(1-\theta)^1$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}}(L(\theta)).$$

$$= \underset{\theta}{\operatorname{argmax}}(\theta^4(1-\theta)^1).$$

Maximizing likelihood function & maximizing the log-likelihood function,

$$\begin{aligned}\ell(\theta) &= \log L(\theta) = \log (\theta^4(1-\theta)^1) \\ &= 4 \log \theta + 1 \log(1-\theta)\end{aligned}$$

$$\hat{\theta} = \arg \max_{\theta} \underbrace{(4 \log \theta + \log(1-\theta))}_{\ell(\theta)}.$$

To find the θ that maximizes ℓ , need to set $\nabla \ell = 0$:

$$\nabla \ell(\theta) = \frac{4}{\theta} + \frac{-1}{1-\theta} = 0.$$

$$\Rightarrow \frac{4}{\theta} = \frac{1}{1-\theta} \Rightarrow \frac{1-\theta}{\theta} = \frac{1}{4} \Rightarrow \frac{1}{\theta} - 1 = \frac{1}{4}$$

$$\Rightarrow \frac{1}{\theta} = \frac{5}{4} \Rightarrow \boxed{\hat{\theta} = \frac{4}{5}} \text{ as expected intuitively.}$$

Parameter Estimation by MLE : Algorithm

Problem : Given set of observations $\{x_i\}_{i=1}^n$, drawn iid from unknown distribution P_* , from a parametric family of distributions: $\{p(\cdot | \theta) : \theta \in \Theta\}$. Estimate θ .

Algorithm:

1) Form the likelihood function L :

$$L(\theta) = p(\{x_i\}_{i=1}^n | \theta) = \prod_{i=1}^n p(x_i | \theta)$$

2) Calculate $\ell(\theta) = \log L(\theta)$:

$$\begin{aligned}\ell(\theta) &= \log(L(\theta)) = \log\left(\prod_{i=1}^n p(x_i | \theta)\right) \\ &= \sum_{i=1}^n \log p(x_i | \theta)\end{aligned}$$

3) Calculate the maximum likelihood estimation:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \ell(\theta) = \underset{\theta}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \log p(x_i | \theta) \right\}$$

4) Solve the optimization problem:

- (1) closed form, by setting $\nabla \ell(\theta) = 0$
- (2) Numerical method,

Example : Biased Coin, Revisited :

- Given $\{x_i\}_{i=1}^n$, where $x_i \in \{0, 1\}$ iid drawn from Bernoulli distribution;

$$\begin{cases} p(X=1) = \frac{\exp(\omega)}{1 + \exp(\omega)} = \theta \\ p(X=0) = \frac{1}{1 + \exp(\omega)} = 1 - \theta \end{cases}$$

- Problem : Estimate ω .

Use $p(x) = \frac{\exp(x\omega)}{1 + \exp(\omega)}$ $\xrightarrow{\text{If } x=1} \frac{\exp(\omega)}{1 + \exp(\omega)}$
 $\xrightarrow{\text{If } x=0} \frac{1}{1 + \exp(\omega)}$

$$\ell(\omega) = \log p(x_1, \dots, x_n | \omega) = \log \prod_{i=1}^n p(x_i | \omega)$$

$$= \sum_{i=1}^n \log p(x_i | \omega) = \sum_{i=1}^n \log \left(\frac{\exp(x_i \omega)}{1 + \exp(\omega)} \right)$$

$$= \sum_{i=1}^n [x_i \omega - \log(1 + \exp(\omega))]$$

$$= \left(\sum_{i=1}^n x_i \right) \omega - n \log(1 + \exp(\omega))$$

$$= n(\bar{x}\omega - \log(1 + \exp(\omega))), \text{ where } \bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

$$\nabla \ell(\omega) = n \left(\bar{x} - \frac{\exp(\omega)}{1 + \exp(\omega)} \right) = 0.$$

$$\text{Recall : } \theta = \frac{\exp(\omega)}{1 + \exp(\omega)}$$

$$\text{therefore } \hat{\theta} = \frac{\exp(\hat{\omega})}{1 + \exp(\hat{\omega})} = \bar{x} = \frac{\#(1)}{\#(1) + \#(0)}$$

$$\Rightarrow \bar{x} = \frac{\exp(\hat{\omega})}{1 + \exp(\hat{\omega})} \rightarrow \bar{x}(1 + \exp(\hat{\omega})) = \exp(\hat{\omega})$$

$$\rightarrow \bar{x} + \bar{x}\exp(\hat{\omega}) = \exp(\hat{\omega})$$

$$\rightarrow \bar{x} = \exp(\hat{\omega}) - \bar{x}\exp(\hat{\omega}) = \exp(\hat{\omega})(1 - \bar{x})$$

$$\rightarrow \frac{\bar{x}}{1 - \bar{x}} = \exp(\hat{\omega}) \rightarrow \boxed{\hat{\omega} = \log\left(\frac{\bar{x}}{1 - \bar{x}}\right)}$$

Example : Gaussian Distribution:

Given $\{x_i\}_{i=1}^n$ iid drawn from Gaussian distribution

$N(\mu, \sigma^2)$:

$$p(x | \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right), \quad \theta = \{\mu, \sigma\}.$$

$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \log p(x_i | \theta) \\ &= \sum_{i=1}^n \log \left[\frac{1}{(2\pi)^{1/2}\sigma} - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right] \\ &= \sum_{i=1}^n \left[\log \left(\frac{1}{(2\pi)^{1/2}} \right) - \log(\sigma) - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right] \\ &= \underbrace{-\frac{n}{2} \log(2\pi)}_{\text{constant}} - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

σ terms μ term

Finding $\hat{\mu}$ min instead of max bc of negative sign.

$$\hat{\mu} = \underset{\mu}{\operatorname{argmin}} \left[\sum_{i=1}^n (x_i - \mu)^2 \right]. \quad \text{Find the derivative wrt } \mu.$$

$$\nabla_{\mu} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) = \sum_{i=1}^n 2(x_i - \mu)(-1) = \sum_{i=1}^n 2(\mu - x_i) \stackrel{\text{set}}{=} 0 = \sum_{i=1}^n (\mu - x_i)$$

$$\rightarrow n\mu - \sum_{i=1}^n x_i = 0 \rightarrow \boxed{\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i}$$

= Sample /empirical mean is
the MLE of the mean.

finding $\hat{\sigma}$:

$$\hat{\sigma} = \underset{\sigma}{\operatorname{argmax}} \left[-n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \hat{\mu})^2 \right]$$

call this $f(\sigma)$.

$$\nabla_{\sigma} f(\sigma) = \frac{-n}{\sigma} - \frac{(-2)}{2\sigma^3} \sum_{i=1}^n (x_i - \hat{\mu})^2 \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{n}{\sigma}$$

$$\Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \hat{\mu})^2 = n$$

$$\Rightarrow \boxed{\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2}$$

empirical/sample variance is
the MLE of the variance!

Example: Uniform distribution.

- Given $\{x_i\}_{i=1}^n$ iid drawn from uniform distribution $\text{Uniform}([0, \theta])$:

$$p(x|\theta) = \begin{cases} 1/\theta & \text{if } x \in [0, \theta] \\ 0 & \text{otherwise.} \end{cases}, \quad \theta > 0.$$

- What is the optimal θ ?

$$\rightarrow p(x|\theta) = \frac{1}{\theta} I(x \in [0, \theta])$$

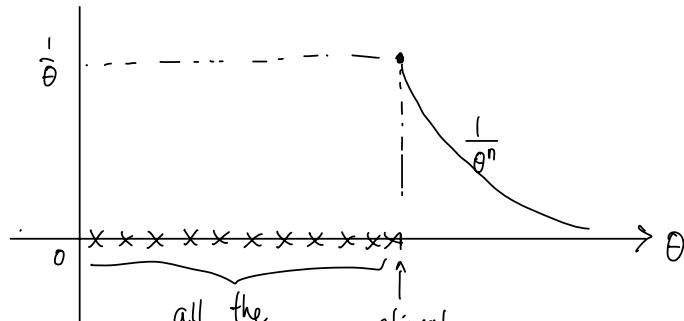
\uparrow indicator function: 1 if $[0, \theta]$, 0 otherwise,

$$L(\theta) = \prod_{i=1}^n p(x_i | \theta) = \prod_{i=1}^n \left(\frac{1}{\theta} \right) I(x_i \in [0, \theta])$$

$$= \frac{1}{\theta^n} \prod_{i=1}^n I(x_i \in [0, \theta])$$

$$= \begin{cases} \frac{1}{\theta^n} & \text{if } x_i \in [0, \theta] \quad \forall i \\ 0 & \text{otherwise (i.e., any } x_i \text{ is not } \in [0, \theta]). \end{cases}$$

$$\Rightarrow \hat{\theta} = \max \{x_1, \dots, x_n\}.$$



all the data points.

optimal point for θ ,

which covers the data points as tight as possible!

MLE for Regression:

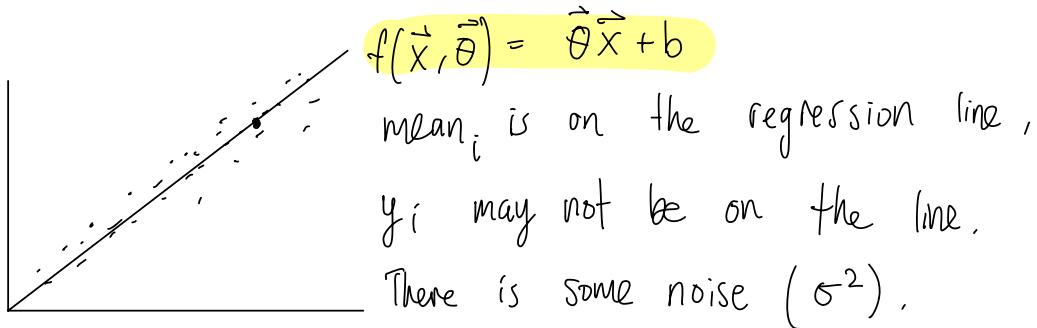
- Given $\{\vec{x}_i, y_i\}_{i=1}^n$, where \vec{x}_i , such that

$$p(y|\vec{x}_i; \vec{\theta}) = N(y | f(\vec{x}_i; \vec{\theta}), \sigma^2) \quad // \text{Normal distribution.}$$

assume y is following the normal distribution,

$$\Rightarrow p(y|\vec{x}_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - f(\vec{x}_i; \vec{\theta}))^2}{2\sigma^2}\right). \quad \text{Here } x \text{ and } \theta \text{ are vectors.}$$

- Estimate μ and σ^2 .



$$\begin{aligned} l(\theta, \sigma) &= \sum_{i=1}^n \log p(y_i | \vec{x}_i, \vec{\theta}) \\ &= \sum_{i=1}^n \log \underbrace{\left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - f(x_i; \theta))^2}{2\sigma^2}\right) \right]}_{p(y_i | \vec{x}_i, \vec{\theta})} \\ &= \underbrace{\left[-\frac{n}{2} \log(2\pi) - n \log \sigma - \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - f(\vec{x}_i; \theta))^2 \right]}_{\text{constant}} \end{aligned}$$

$$\hat{\theta} = \max_{\theta} l(\theta, \sigma) \Rightarrow \arg \min_{\theta} \sum_{i=1}^n (y_i - f(\vec{x}_i, \hat{\theta}))^2$$

least square estimator.

\therefore This is equivalent to minimizing the MSE.

$$\hat{\sigma}^2 = \max_{\theta} l(\theta, \sigma)$$

$$= \max_{\theta} \left(-n \log \sigma - \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - f(\vec{x}_i, \hat{\theta}))^2 \right)$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - f(\vec{x}_i, \hat{\theta}))^2$$

MLE for logistic regression binary classification

- Given $\{\vec{x}_i \cdot y_i\}_{i=1}^n$, where $y_i \in \{0, 1\}$, s.t.

$$p(y | x, \theta) = \frac{\exp(y \cdot f(x, \theta))}{1 + \exp(y \cdot f(x, \theta))}$$

- Estimate θ .

$$\left\{ \begin{array}{l} p(y = 1 | x, \theta) = \frac{\exp(1 \cdot f(x, \theta))}{1 + \exp(y \cdot f(x, \theta))} \\ p(y = 0 | x, \theta) = \frac{1}{1 + \exp(y \cdot f(x, \theta))} \end{array} \right\} \quad \begin{array}{l} \text{probability} \\ \text{mass} \\ \text{functions} \end{array}$$

$$l(\theta) = \sum_{i=1}^n \log p(y_i | x_i, \theta).$$

$$= \sum_{i=1}^n \log \left[\frac{\exp(y_i \cdot f(x_i, \theta))}{1 + \exp(y_i \cdot f(x_i, \theta))} \right]$$

$$= \sum_{i=1}^n \left[y_i \cdot f(x_i, \theta) - \log(1 + \exp(f(x_i, \theta))) \right]$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \left[y_i \cdot f(x_i, \theta) - \log(1 + \exp(f(x_i, \theta))) \right] \right\}$$

use numeric methods (such as gradient descent).

MLE part II : Theoretical properties.

- MLE estimator is a random variable

$$\hat{\theta} = \hat{\theta}(x_1, \dots, x_n) \quad \{x_i\} \stackrel{iid}{\sim} p(\cdot | \theta_*)$$

- Evaluation metrics : Bias, variance, mean square error,
- Unbiased estimators vs. Consistent estimators.

$$\text{Bias}(\hat{\theta}) = (\text{Expected value of } \hat{\theta}) - (\text{Actual value of } \theta)$$

$$= E[\hat{\theta}(x_1, \dots, x_n)] - \underbrace{\theta_*}_{\text{true parameter}}$$

$$= \int \hat{\theta}(x_1, \dots, x_n) \underbrace{\prod_{i=1}^n p(x_i | \theta_*)}_{\text{pdf of } \hat{\theta}} dx - \theta_*$$

$$\text{Var}(\hat{\theta}) = E \left\{ \underbrace{[\hat{\theta}(x_1, \dots, x_n)]}_{\text{estimator}} - \underbrace{E(\hat{\theta}(x_1, \dots, x_n))}_{\text{expectation of the estimator}} \right\}^2$$

Recall :

$$\text{Var}(X) \stackrel{\text{def}}{=} E\{(X - EX)^2\}$$

expectation of the estimator : $E(\hat{\theta})$.

If unbiased : $E(\hat{\theta}) = \theta_*$.

$$\text{MSE}(\hat{\theta}) = E \left\{ \underbrace{(\hat{\theta}(x_1, \dots, x_n) - \theta_*)^2}_{\text{mean squared error.}} \right\}$$

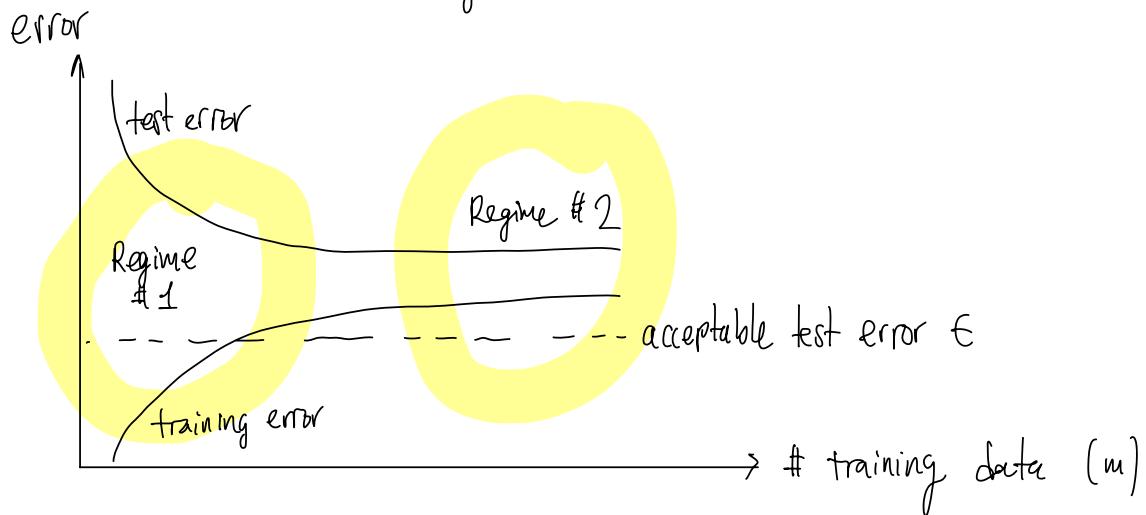
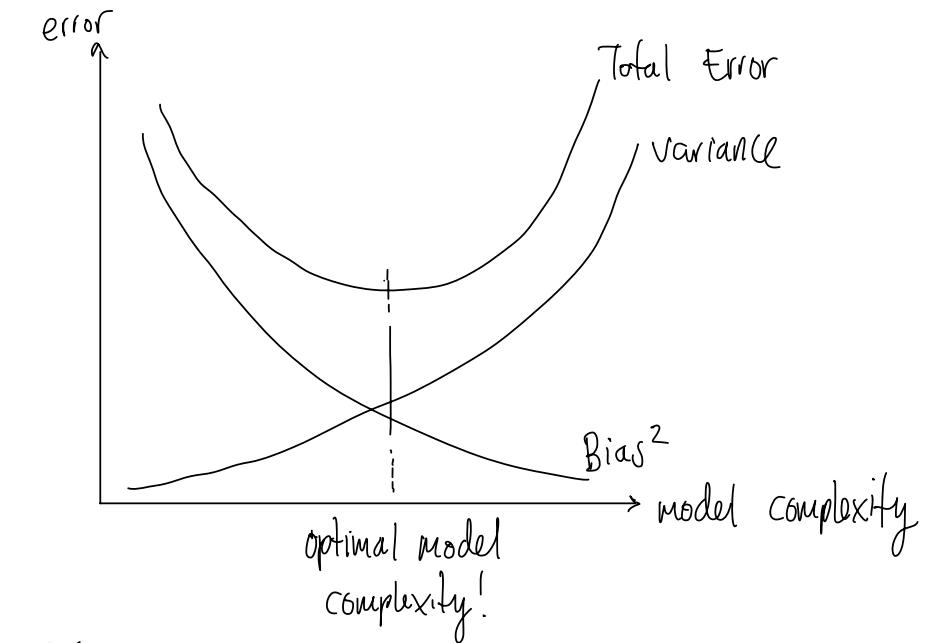
Bias-Variance decomposition :

$$\text{MSE}(\hat{\theta}) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$$

Proof : mean squared error.

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= E[(\hat{\theta} - \theta_*)^2] \quad // \text{definition of MSE} \\ &= E\left\{[(\hat{\theta} - \underbrace{E(\hat{\theta})}_{\text{subtract}}) + (\underbrace{E(\hat{\theta}) - \theta_*}_{\text{add it back}})]^2\right\} \end{aligned}$$

$$\begin{aligned} &= E\left[\underbrace{(\hat{\theta} - E(\hat{\theta}))^2}_{\text{variance}} + \underbrace{(E(\hat{\theta}) - \theta_*)^2}_{\text{bias}^2} + \right. \\ &\quad \left. 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta_*)\right] \\ &\quad \xrightarrow{\text{constant, move outside the expectation function,}} \\ &= E[(\hat{\theta} - E(\hat{\theta}))(\underbrace{E(\hat{\theta}) - \theta_*}_{\text{bias}})] = E[\hat{\theta} - E(\hat{\theta})](E(\hat{\theta}) - \theta_*) \\ &= [E(\hat{\theta}) - E(\hat{\theta})](E(\hat{\theta}) - \theta_*) \quad // E[E[X]] = E[X] (?) \\ &= 0. \end{aligned}$$



Regime #1 = High variance:

Symptoms: training error \ll test error, training error $< \epsilon$, test error $> \epsilon$.

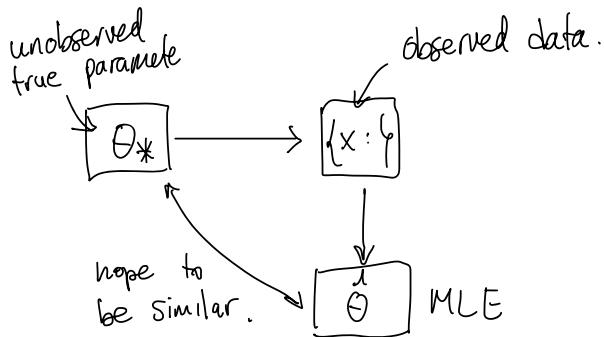
Remedy: increase m , reduce complexity of model, bagging,

Regime #2 = High bias.

Symptoms: training error $> \epsilon$.

Remedy: increase model complexity, add features, boosting.

Framework:



If $\text{Bias}(\hat{\theta}) = 0$, $\hat{\theta}$ is "unbiased", implies

If $\text{MSE}(\hat{\theta}) \rightarrow 0$ as $n \rightarrow \infty$: $\hat{\theta}$ is "consistent". implies

If $\text{Bias} \rightarrow 0$ as $n \rightarrow \infty$: $\hat{\theta}$ is "asymptotically unbiased."

Example: Gaussian Distribution,

for $\{x_i\}_{i=1}^n \sim N(\mu, \sigma^2)$, MLE is:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2.$$

$$\begin{aligned} \text{Bias}(\hat{\mu}) &= E[\hat{\mu}] - \mu = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] - \mu = \frac{1}{n} \sum_{i=1}^n E[x_i] - \mu \\ &= \frac{1}{n} \sum_{i=1}^n \mu - \mu = 0 \end{aligned}$$

"unbiased".

$$\begin{aligned}
 \hat{\sigma}^2 &= \text{Var}(\hat{\mu}) = E[(\hat{\mu} - \bar{\mu})^2] \stackrel{\text{calculated, } E(\hat{\mu}) = \mu}{=} E\left[\left(\overbrace{\frac{1}{n} \sum_{i=1}^n x_i}^{\hat{\mu}} - \mu\right)^2\right] \\
 &= E\left[\frac{1}{n^2} \left(\sum_{i=1}^n (x_i - \mu)^2 + \sum_{i \neq j} (x_i - \mu)(x_j - \mu) \right)\right] // \text{completing the square?} \\
 &= \underbrace{\frac{1}{n^2} \sum_{i=1}^n E[(x_i - \mu)^2]}_{\hat{\sigma}^2} + \underbrace{\sum_{i \neq j} E[(x_i - \mu)(x_j - \mu)]}_{\text{Cov}(x_i, x_j)}.
 \end{aligned}$$

Assume independence:

$$\Rightarrow E[(x_i - \mu)] E[(x_j - \mu)] = 0$$

$$= \frac{n \sigma^2}{n^2} = \boxed{\frac{\sigma^2}{n}}$$

$$MSE(\hat{\mu}) = \text{Bias}(\hat{\mu})^2 + \text{Var}(\hat{\mu}) = \frac{\sigma^2}{n}.$$

If $n \rightarrow \infty$, $MSE(\hat{\mu}) \rightarrow 0 \Rightarrow$ consistent.

$E[\hat{\sigma}^2] \neq \sigma^2 \Rightarrow \hat{\sigma}^2$ is a biased estimator for σ^2 .

$$\hat{\sigma}^2 = \underbrace{\frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2}_{\text{instead of } n.} \text{ is an unbiased estimator.}$$

$\text{bias}(\hat{\sigma}^2) \rightarrow 0$ as $n \rightarrow \infty \Rightarrow$ asymptotically Unbiased.

$\text{Var}(\hat{\sigma}^2) \rightarrow 0$ as $n \rightarrow \infty$

$\Rightarrow MSE(\hat{\sigma}^2) \rightarrow 0$ as $n \rightarrow \infty \Rightarrow$ consistent.

Why MLE?

MLE is always consistent.

MLE is equivalent to minimizing Kullback-Leibler (KL) divergence (difference btwn two distributions, p and q):

$$KL(q \parallel p) \equiv E_q \left[\log q(x) - \log p(x) \right]$$

\neq

$$= \begin{cases} \sum_x q(x) \left(\log \left(\frac{q(x)}{p(x)} \right) \right) & \text{// discrete case.} \\ \int_x q(x) \left(\log \left(\frac{q(x)}{p(x)} \right) \right) dx & \text{// continuous case.} \end{cases}$$

not symmetric!

Claims

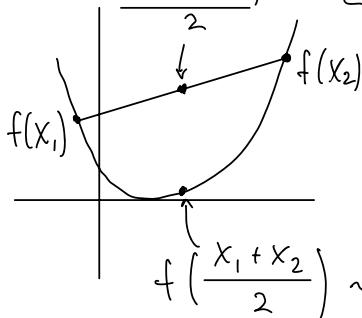
• $KL(q \parallel p) \geq 0$ for any q and p . (1)

• $KL(q \parallel p) = 0$ i.i.f. $q = p$. (2)

Proof of (1) : use Jensen Inequality :

If $f(x)$ is convex, then $E_q[f(x)] \geq f(E_q[x])$

$$\frac{f(x_1) + f(x_2)}{2} \sim E[f(x_{1,5})]$$

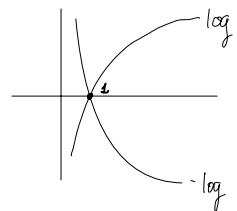


$$q(x) = \frac{f(x_1) + f(x_2)}{2}$$

$$f\left(\frac{x_1 + x_2}{2}\right) \sim f(E[x_{1,5}])$$

$$KL(q \parallel p) = E_q \left[\log \left(\frac{q}{p} \right) \right]$$

↑
concave function



$$\begin{aligned}
 &= E_q \left[-\underbrace{\log \left(\frac{p(x)}{q(x)} \right)}_{-\log = \text{convex.}} \right] \geq -\log \left(E_q \left[\frac{p(x)}{q(x)} \right] \right) \quad // \text{Jensen inequality} \\
 &= -\log \left(\sum_x q(x) \frac{p(x)}{q(x)} \right) \quad // \text{Definition of Expectation.} \\
 &= -\log \left(\sum_x p(x) \right) = -\log(1) = 0.
 \end{aligned}$$

This completes proof of (1).

Proof of (2) :

Suppose $KL(q \parallel p) = 0$. This means

$$0 = KL(q \parallel p) = -\log \left(E_q \left[\frac{p(x)}{q(x)} \right] \right) \Rightarrow \frac{p(x)}{q(x)} = \text{const.}$$

$$\Rightarrow p(x) = \text{const.} \cdot q(x) \Rightarrow \sum_x p(x) = \text{const.} \sum_x q(x)$$

$$\Rightarrow \underbrace{1}_{\substack{\text{sum of} \\ \text{probability} \\ \text{mass} = 1}} = \text{const.} \cdot \underbrace{1}_{\substack{\text{sum of} \\ \text{probability} \\ \text{mass} = 1}} \Rightarrow \text{const.} = 1 \Rightarrow p = q.$$

How to connect MLE with KL Divergence?

$$KL(q \parallel p) = E_q [\log q(x) - \log p(x)]$$

Assume q is the data distribution, // Given

p is the "model". // want to optimize.

$$\begin{aligned} \min_{\theta} [KL(q \parallel p_{\theta})] &\Leftrightarrow \min_{\theta} E_q [\log q(x) - \log p_{\theta}(x)] \\ &\Leftrightarrow \min_{\theta} -E_q [\log p_{\theta}(x)] \\ &\Leftrightarrow \max_{\theta} E_q [\log p_{\theta}(x)] \\ &= \underbrace{\max_{\theta} \left[\frac{1}{n} \sum_{i=1}^n \log p_{\theta}(x_i) \right]}_{\text{average log-likelihood.}} \quad \{x_i\} \sim q \end{aligned}$$

∴ Maximizing the log-likelihood is equivalent to
minimizing the KL divergence btwn the
data distribution q and the model p_{θ} .