# Maximum Likelihood Estimation: Theoretical Properties

**Qiang Liu**
UT Austin

- MLE estimator is random variable (as a function of the random data):

$$\hat{\theta} = \hat{\theta}(x_1, \ldots, x_n), \qquad \{x_i\} \overset{iid}{\sim} p(\cdot \mid \theta_*).$$

- Evaluation metrics: Bias, variance, mean square error (MSE).
- Unbiased estimators vs. consistent estimators.

$$\text{Bias}(\hat{\theta}) = \underset{\theta^*}{E}\left[\hat{\theta}(x_1 \cdots x_n)\right] - \theta^*$$

If $\text{Bias}(\hat{\theta}) = 0$, we call $\hat{\theta}$ "unbiased".

$\text{MSE}(\hat{\theta}) \to 0$, as $n \to \infty$, $\hat{\theta}$ "consistent".

$\text{Bias}(\hat{\theta}) \to 0$, as $n \to \infty$, "Asymptotic Unbiased"

2

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

Proof: 
$$\boxed{\text{MSE}(\hat{\theta})} = E\left[(\hat{\theta} - \theta_*)^2\right]$$
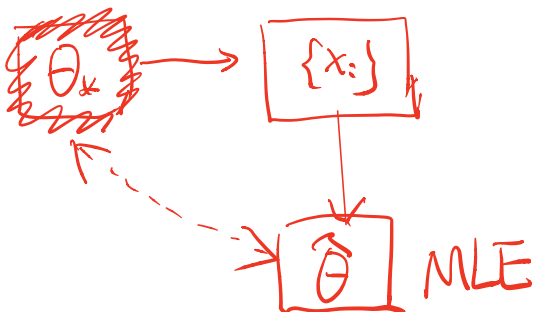$$= E\left[(\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta_*)^2\right]$$
$$= E\left[(\hat{\theta} - E(\hat{\theta}))^2 + (E(\hat{\theta}) - \theta_*)^2 + 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta_*)\right]$$
$$\underset{=0}{}$$
$$= \text{Var}(\hat{\theta}) + (\text{bias}(\hat{\theta}))^2$$

$$E\left[(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta_*)\right] = E\left[\hat{\theta} - E(\hat{\theta})\right](E(\hat{\theta}) - \theta^*)$$
$$= (E(\hat{\theta}) - E(\hat{\theta}))(E(\hat{\theta}) - \theta^*)$$
$$= 0$$

$\theta_* \longrightarrow \boxed{\{x_i\}}$

$\boxed{\hat{\theta}}$ MLE

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

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

"Unbiased"

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

$$\boxed{\text{var}(\hat{\mu})} = E\left[(\hat{\mu} - \mu)^2\right] \qquad\qquad = \frac{1}{n}\sum_{i=1}^n \mu - \mu = \boxed{0}$$

$$= E\left[\left(\frac{1}{n}\sum_{i=1}^n x_i - \mu\right)^2\right] = E\left[\frac{1}{n^2}\left(\sum_{i=1}^n (x_i - \mu)^2 + \sum_{i \neq j}(x_i - \mu)(x_j - \mu)\right)\right]$$

$$= \frac{1}{n^2}\sum_{i=1}^n E\left((x_i - \mu)^2\right) + \underbrace{\sum_{i \neq j} E\left[(x_i - \mu)(x_j - \mu)\right]}_{= E[(x_i - \mu)]\,E[(x_j - \mu)]}$$

$$= \boxed{\frac{1}{n}\sigma^2} + \qquad = E[(x_i - \mu)]\,E[(x_j - \mu)] \quad = 0$$

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

If $n \to \infty$, $\text{MSE}(\hat{\mu}) \to 0 \Rightarrow$ Consistent.

$$\boxed{E[\hat{\sigma}^2] \neq \sigma^2} \quad, \qquad \hat{S}^2 = \boxed{\frac{1}{n-1}}\sum_{i=1}^n (x_i - \hat{\mu})^2$$

"Biased"

$$E[\hat{S}^2] = \sigma^2$$

$$\text{Bias}(\hat{\sigma}^2) \to 0 \quad . \quad n \to \infty \ (\text{Asymp. Unbiased})$$
$$\text{var}(\hat{\sigma}^2) \to 0 \qquad n \to \infty$$
$$\text{MSE}(\hat{\sigma}^2) \to 0 \qquad n \to \infty \ (\text{Consistent})$$

- MLE is equivalent to minimizing Kullback-Leibler (KL) Divergence.

$$KL(q \| p) = \mathbb{E}_q[\log q(x) - \log p(x)]. = \begin{cases} \sum_x q(x)\left(\log\left(\frac{q(x)}{p(x)}\right)\right) \\ \int q(x) \log\left(\frac{q(x)}{p(x)}\right) dx \end{cases}$$

- $KL(q \| p) \geq 0$ for any $q$ and $p$. ☆
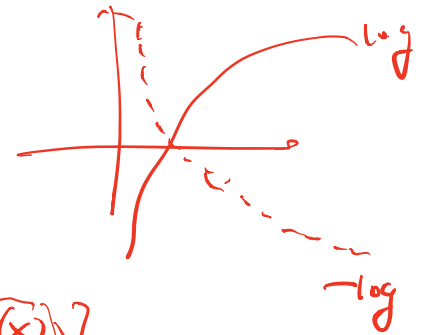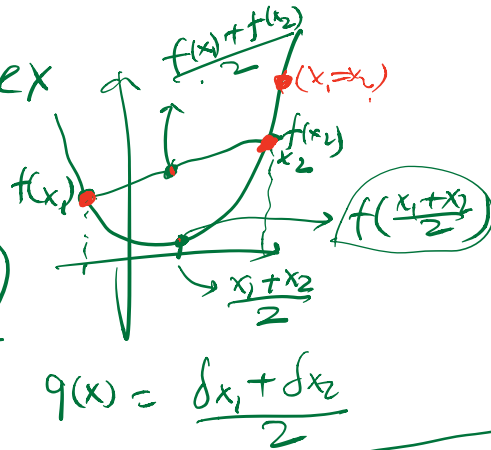- $KL(q \| p) = 0$ if and only if $q = p$.

$$KL(q \| P) \neq KL(P \| q)$$

**Jensen's Inequality:**

If $f(x)$ is convex

then

$$\mathbb{E}_q[f(x)] \geq f(\mathbb{E}_q(x))$$



$$f\left(\frac{x_1 + x_2}{2}\right)$$

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

$$KL(q \| P) = \mathbb{E}_q\left[\log\left(\frac{q}{p}\right)\right] = \mathbb{E}_q\left[-\log\left(\frac{P(x)}{q(x)}\right)\right]$$

$$\underset{\text{"}0}{\geq} -\log\left(\mathbb{E}_q\left[\frac{P(x)}{q(x)}\right]\right)$$

$$\Rightarrow \frac{P(x)}{q(x)} = Const.$$

$$= -\log\left(\sum_x q(x) \frac{P(x)}{q(x)}\right)$$

$$\Rightarrow P(x) = Const \cdot q(x)$$

$$= -\log\left(\sum_x P(x)\right) = -\log(1)$$

$$\Rightarrow \sum_x P(x) = Const \cdot \sum_x q(x) \qquad \boxed{= 0}$$

$$\Rightarrow 1 = Const \cdot 1 \qquad \Rightarrow Const = 1$$

$$\Rightarrow P = q$$

$$MLE \iff KL$$

$$KL(q \| P) = \mathbb{E}_q[\log q(x) - \log P(x)]$$

Assume $q$ is data distribution (Given)

$P_\theta$ is "model" $\theta \in \Theta$

$$\min_\theta\left(KL(q \| P_\theta)\right) \iff \min_\theta \mathbb{E}_q[\log q(x) - \log P_\theta(x)]$$

$$\Leftrightarrow \min_{\theta} -E_q\left[\log P_\theta(x)\right]$$

$$\Leftrightarrow \boxed{\max_{\theta} E_q\left[\log P_\theta(x)\right]}$$

$$\approx \max_{\theta} \frac{1}{n}\sum_{i=1}^{n} \log P_\theta(x_i) \qquad \{x_i\} \sim q$$

$$\underline{\text{Avg log-likelihood.}}$$