

LAST WEEK

Two result (No proof)

$$1) \text{ LLN } S_N = \frac{1}{N} \sum_i x_i \xrightarrow{P} E[m]$$

$$2) \text{ CLT } P\left[\left|S_N - E[m]\right| > \epsilon\right] \rightarrow 0$$

$$\frac{S_N - N\mu}{\sqrt{\Delta}} \xrightarrow{D} \mathcal{N}(0, 1)$$

Comparing Average : how far S_N is from $\mu = E[x]$?

$$\text{In General} : P[|S_N - \mu| \geq k\sqrt{\Delta}] \leq \frac{1}{k^2} \quad (\text{Chebyshev})$$

Can we much better; Ex X Gaussian ~ N(0, 1)

$$\begin{aligned} P\left[\frac{1}{N} \sum_i x_i \geq N(\mu + \epsilon)\right] &= P\left[\sum_i x_i \geq N(\mu + \epsilon)\right] = P\left[e^{t \sum_i x_i} \geq e^{Nt(\mu + \epsilon)}\right] \\ &\leq \frac{E[e^{t \sum_i x_i}]}{e^{Nt(\mu + \epsilon)}} = \left(E[e^{tx}]\right)^n = \left(e^{t\mu + \frac{t^2}{2\Delta}}\right)^n = e^{n\left[\frac{t^2}{2\Delta} - t\epsilon\right]} \end{aligned}$$

$$\frac{t^2}{\Delta} = \epsilon$$

$$P\left[\frac{1}{N} \sum_i x_i \geq N(\mu + \epsilon)\right] \leq e^{-\frac{n\epsilon^2}{2\Delta}}$$

$$\text{sim } P\left[\frac{1}{N} \sum_i x_i \leq N(\mu - \epsilon)\right] \leq e^{-\frac{n\epsilon^2}{2\Delta}}$$

$$P\left[|\frac{1}{N} \sum_i x_i - \mu| \geq \epsilon\right] \leq 2e^{-\frac{n\epsilon^2}{2\Delta}}$$

$$P\left[|\frac{1}{N} \sum_i x_i - \mu| \leq \epsilon\right] \geq 1 - 2e^{-\frac{n\epsilon^2}{2\Delta}}$$

Accuracy is cheap

$$\epsilon = \sqrt{\frac{2\Delta}{n} \log \frac{2}{\delta}}$$

with prob at least $1 - \delta$, $|p - \hat{p}| \leq \sqrt{\frac{2\Delta}{m} \log \frac{2}{\delta}}$

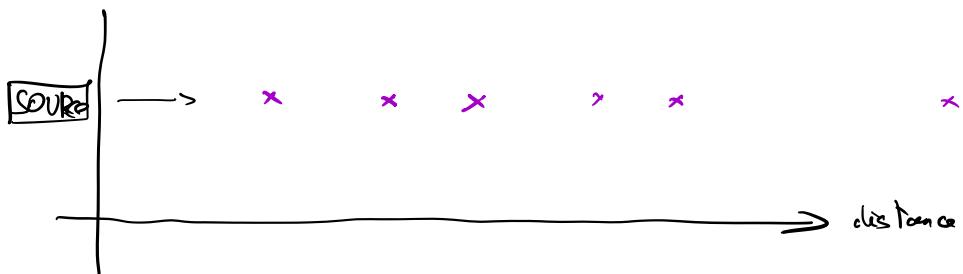
Precision
is Expressive

PA C Learning

~

Lesson II All of Statistics

Learn α single parameter



$$\hat{\mu}_{EM}(\{x_i\}) = \frac{1}{m} \sum_{i=1}^m x_i$$

How do you estimate the
"mean-path" typical distance $\hat{\mu}^*$

$$P(x; \hat{\mu}^*) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\hat{\mu}^*)^2}{2\sigma^2}}$$

mean and variance

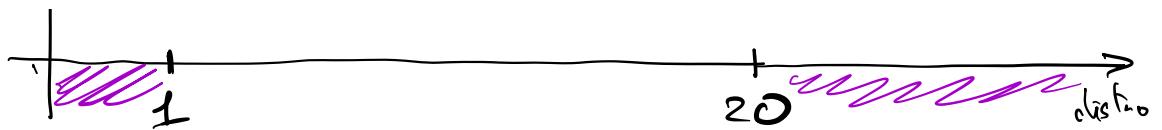
$$\mathbb{E}[\hat{\mu}_{EM}(\{x_i\})] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}x_i = \hat{\mu}^*$$

$$\text{VAR}[\hat{\mu}_{EM}] = \frac{1}{m} \text{VARIANCE}(x) = \frac{\sigma^2}{m}$$

$$\hat{\mu}_{nL} = \hat{\mu}^* \pm \frac{\sigma}{\sqrt{m}}$$

~





How to choose $\hat{\theta}(\xi \bar{x})$??

① Bayesian Approach

likelihood

Post

$$P(\Theta | \xi \bar{x}) = \frac{P(\xi \bar{x} | \Theta) P(\Theta)}{P(\xi \bar{x})}$$

Posterior Distribution

$$\tilde{P}(x | \Theta) = \begin{cases} \frac{e^{-x/\Theta}}{2\Theta} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$P(\xi \bar{x} | \Theta) = \prod_{i=1}^n \tilde{P}(x_i | \Theta)$$

Jeffreys Prior?

$$P(\Theta) = \underline{\text{const}}$$

"Maximum Entropy"
"minimum Assumption"

???

$$P(\Theta | \xi \bar{x}) \propto \text{likelihood} \Rightarrow \hat{\Theta}_{ML} = \text{ARGMAX } P(\xi \bar{x} | \Theta)$$

Maximizing likelihood

$$\hat{\Theta}_{ML} = \text{ARGMAX } \log P(\xi \bar{x} | \Theta)$$

② Frequentist

$$\text{Risk } [\hat{\theta}, \theta^*] = E(\hat{\theta}(\xi) - \theta^*)^2 \quad \text{QUADRATIC RISK}$$

MSE
Mean Squared Error

Study of Maximum Likelihood

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^n P(x_i | \theta)$$

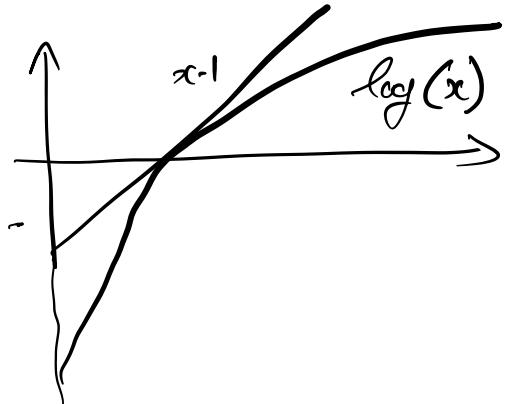
$$\begin{aligned} \ell(x, \theta) &= \log P(x | \theta) \\ \ell^n(\xi, \theta) &= \frac{1}{n} \sum_{i=1}^n \ell(x_i, \theta) \end{aligned} = \hat{\theta}(\xi) = \underset{\theta}{\operatorname{argmax}} \ell^n(\xi, \theta)$$

Consistent $n \rightarrow \infty$, then $\hat{\theta}_{ML} \rightarrow \theta^*$

$$\begin{aligned} \lim_{n \rightarrow \infty} \ell^n(\xi, \theta) - \ell^n(\xi, \theta^*) &= \frac{1}{n} \sum_{i=1}^n \log P(x_i | \theta) - \frac{1}{n} \sum_{i=1}^n \log P(x_i | \theta^*) \\ &= E \log P(x | \theta) - E \log P(x | \theta^*) \\ &= \int d\omega P(x | \theta^*) \log P(x | \theta) - \int d\omega P(x | \theta^*) \log P(x | \theta^*) \\ &= \int d\omega P(x | \theta^*) \log \frac{P(x | \theta)}{P(x | \theta^*)} = -D_{KL}(P(x | \theta) || P(x | \theta^*)) \leq 0 \end{aligned}$$

Proof

Gibbs Inequality



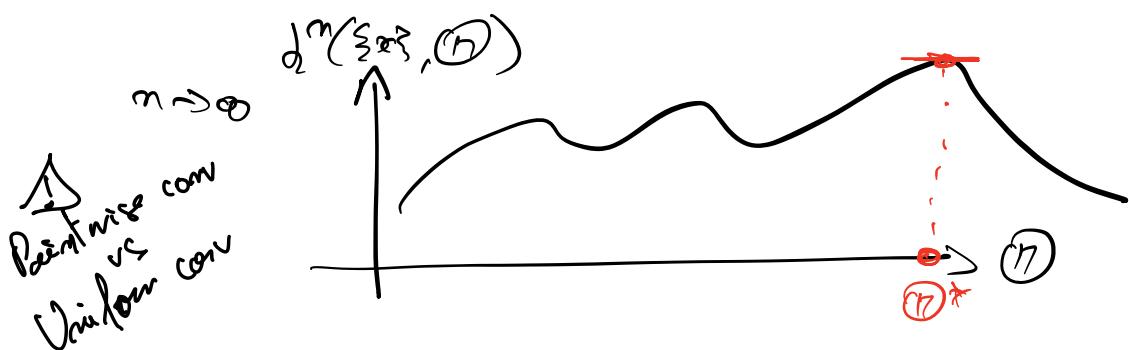
$$x-1 \geq \log x$$

$$P_i, q_i \quad i=1 \dots N$$

$$\pi = \frac{P_i}{q_i}$$

$$\sum_{i=1}^N q_i \left(\frac{P_i}{q_i} - 1 \right) \geq \sum_{i=1}^N q_i \log \frac{P_i}{q_i}$$

$$\Omega = \sum_{i=1}^N P_i - \sum_{i=1}^N q_i \geq \sum_{i=1}^N q_i \log \frac{P_i}{q_i}$$



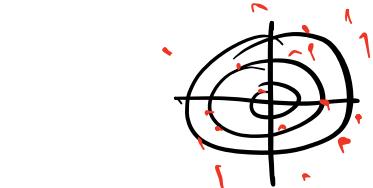
MAXIMUM LIKELIHOOD IS CONSISTENT

$$\lim_{n \rightarrow \infty} \hat{\theta}_L^n(\Sigma x_i^n) = \theta^*$$

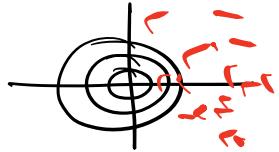
It limits of learning

BIGE-VARIATION

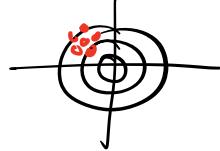
$$MSE = \mathbb{E} \left[(\hat{\lambda}(\sum x_i^2) - \lambda^*)^2 \right] = \underbrace{\mathbb{E}[\hat{\lambda}(\sum x_i^2)]}_{\text{Variance } (\hat{\lambda})} - \mathbb{E}[\hat{\lambda}(\sum x_i^2)]^2 + \underbrace{(\mathbb{E}[\hat{\lambda}(\sum x_i^2)] - \lambda^*)^2}_{\text{BIAS}}$$



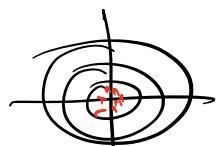
large variance
small bias



large variance
large bias



small variance
large bias



small variance, small bias

$$\text{bias } (\hat{\theta}, \theta) = \mathbb{E} (\hat{\theta} - \theta^*)$$

$$MSE(\hat{\theta}^*, \hat{\theta}) \geq b^2(\hat{\theta}, \theta^*) + \frac{1}{nI(\theta)} (\text{variance of } \hat{\theta})$$

Cramér-Rao

Assume n is large, and assume $\hat{\lambda}$ is constant when $n \rightarrow \infty$

$\hat{\lambda}$ should be unbiased

Bound for unbiased Estimators

$$MSE \geq \frac{1}{nI(\theta^*)}$$

Fisher Score

$$S(x, \theta) = \frac{\partial \ell(x, \theta)}{\partial \theta}$$

$$S^*(\bar{x}_i, \theta) = \frac{\partial \ln(\bar{x}_i, \theta)}{\partial \theta}$$

$$\mathbb{E}[S(x, \theta)] = 0$$

$$\mathbb{E}[S^2(x, \theta^*)] = -\mathbb{E}[\partial_\theta S(x, \theta)|_{\theta^*}] = I_{\text{F}}(\theta^*) \quad \text{Fisher Information}$$

$$\mathbb{E}[(S^m(\bar{x}_i, \theta^*))^2] = m I_{\text{F}}(\theta^*)$$

$$\int dx P(x|\theta^*) \frac{\partial P(x|\theta)}{\partial \theta}$$

$$= \partial_\theta \underbrace{\int dx P(x|\theta^*)}_{=1} = 0$$

$$\begin{aligned} \mathbb{E}[\partial_\theta S(x, \theta)|_{\theta^*}] &= \int dx P(x|\theta^*) \frac{(\partial_\theta P)P - (\partial_\theta P)^2}{P^2(x|\theta^*)} \\ &= \cancel{\int dx P} - \int dx P \left(\frac{\partial_\theta P}{P} \right)^2 = -\mathbb{E}[S^2] \end{aligned}$$

Proof C.R : Cauchy-Schwarz

$$\text{cov}[\hat{\theta}, S]^2 \leq \text{var}[\hat{\theta}] \text{var}[S] = \text{var}(\hat{\theta}) n I$$

$$\text{var}[\hat{\theta}] \geq \frac{(\mathbb{E}[\hat{\theta}S] - \mathbb{E}[\hat{\theta}]\mathbb{E}[S])^2}{n}$$

$$\text{If } E > 0, \quad \frac{\mathbb{E}[\hat{\theta}S]^2}{n} + \text{bias}^2$$

$$\text{But } \Rightarrow \partial_\theta \mathbb{E}[\hat{\theta}] \cdot 1 = \partial_\theta \int dx P(x|\theta) \hat{\theta}(x) \cdot 1 = \int dx P(x|\theta) \cdot 1$$

$$= \int dx P \frac{\partial}{\partial \theta} \hat{\theta} \cdot 1 = \mathbb{E}[\hat{\theta}] \cdot 1 \quad \blacksquare$$

$\Rightarrow \text{ML is efficient} \Rightarrow \text{It Achieves Crammer-Rao Bound}$
Asymptotically

Bound for Unbiased Estimators

$$\text{NCF} \geq \frac{1}{nI(\theta^*)}$$

$$O = \hat{\theta}_n (\sum x_i, \hat{\theta}_n) = \hat{\theta}_n (\sum x_i, \theta^*) + (\hat{\theta}_n - \theta^*) \hat{\theta}_n (\sum x_i, \theta^*) + O(n(\hat{\theta}_n - \theta^*)^2)$$

$$\sqrt{n}(\hat{\theta}_n - \theta^*) = -\frac{\partial \ln(\sum x_i, \theta^*)}{\partial \theta_n (\sum x_i, \theta^*)} + X \xrightarrow{\text{CLT}} \frac{\mathcal{N}(0, I)}{I(\theta^*)}$$

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{n \text{ large}} \mathcal{N}(0, \frac{1}{I})$$

$$\hat{\theta}_n \text{ distributed as } e^{-\frac{(\hat{\theta}_n - \theta^*)^2}{2 \frac{1}{I}}} / \sqrt{2\pi \frac{1}{I}}$$

$$\hat{\theta}_n \sim \mathcal{N}(\theta^*, \frac{1}{I})$$

- ② Maximum Likelihood is Asymptotically Gaussian
- ③ " " achieves " " " Crammer-Rao Bound

ML is efficient

Find Note \Rightarrow Work for Many Variables as Well

$$[I(\vec{\theta})]_{ij} = \mathbb{E} \left[\frac{\partial \log P}{\partial \theta_i} \frac{\partial \log P}{\partial \theta_j} \right] = \mathbb{E} \left[\frac{\partial}{\partial \theta_i \partial \theta_j} \log P \right]$$

Fisher Matrix

$$\text{Cov}[\vec{\theta}] \geq (I_{\vec{\theta}})^{-1} \quad \leftarrow \text{Cramers-Rao}$$