

## MATH 341/641 LEC 15

So far, this class focused on the 3 goals of inference: estimation, testing, confidence set creation.

We began by <sup>assuming</sup> a DGP that explained where the data was realized from and what the parameters were and went from there.

How do you assume a DGP? Sometimes you tentatively know e.g. coin flips are Binomial. Most often you don't. Wind speeds at JFK airport DGP? Survival times of lab rats? Unknown! We need to guess. Why not guess a few different DGPs and select the best one (Model Selection).

⇒ A DGP<sub>1</sub><sup>posited</sup> is a "model" which is a useful approximation to reality. We then "fit" the model by estimating its params. Once we have that, we know properties of the model and can make predictions (302 make one).

Consider  $M$  candidate models  $m=1, 2, \dots, M$ . We would like to pick the best one ("model selection").

Additionally, we would like to assign a probability measure score to each.

Model Selection is a fundamental problem in science!

$$m=1: F = G \frac{m_1 m_2}{r^2}$$

(Newton's Law)

$$m=2: F = G_1 \frac{m_1 m_2}{r^2} + G_2 \frac{m_1 m_2}{r^3}$$

(Newton's Extension)

$$m=3: F = G_1 \frac{m_1 m_2}{r^2} e^{-G_2 r}$$

(Laplace's Extension)

Which one is "best"? We know mod #1 is a good approx but wrong as Einstein's relativity supersedes it.

$$\text{Ob } \#1: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_1(x_i; \theta_{1,1}, \dots, \theta_{1,k_1}) = L_1(x_i; \theta_{1,1}, \dots, \theta_{1,k_1})$$

$$\vdots$$
$$\text{Ob } \#m: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_m(x_i; \theta_{m,1}, \dots, \theta_{m,k_m}) = L_m(x_i; \theta_{m,1}, \dots, \theta_{m,k_m})$$

The # of params could be different in each model  $k_1, k_2, \dots, k_m$

Why not just take the model with the highest likelihood?

$$\text{i.e. } m_{\hat{}} = \underset{m \in \{1, \dots, m\}}{\text{argmax}} \left\{ L_m(\theta_{m,1}, \dots, \theta_{m,k_m}; X_1, \dots, X_n) \right\}$$

$$= \underset{m \in \{1, \dots, m\}}{\text{argmax}} \left\{ L_m(\theta_{m,1}, \dots, \theta_{m,k_m}; X_1, \dots, X_n) \right\}$$

Can we do this? No! we don't know the values of the  $\theta_m$ 's for any  $m$ !

Why not use our best guess of the  $\theta$ 's i.e.  $\hat{\theta}_{MLE}$ ?

$$\Rightarrow m_x = \argmax_{m \in \{1, \dots, M\}} \{ \ell(\hat{\theta}_{m,1}^{MLE}, \dots, \hat{\theta}_{m,K_m}^{MLE}; X_1, \dots, X_n) \}$$

This works only if  $K_1 = K_2 = \dots = K_M$  i.e. the # params is the same in all models. If not, you're giving unfair advantage to models that have a higher # of params.

Eventually as  $K \rightarrow \infty$ , the model will fit perfectly ("overfitting", a core concept in the 392 class).

"With 4 parameters, I can fit an elephant; with 5 parameters,

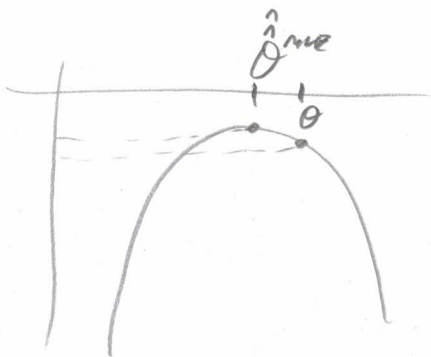
I can make it wiggle its trunk" - John von Neumann,

famous mathematician, computer scientist

So we need a way to "penalize" the model by  $K_m$ , the # of parameters in the model. The penalty is to

subtract the bias from  $\ell(\hat{\theta}_{m,1}^{MLE}, \dots, \hat{\theta}_{m,K_m}^{MLE}; X_1, \dots, X_n)$ .

Why is there bias? If  $K=1$



$$\hat{\theta}_{MLE} \sim N(\theta, \sqrt{\frac{1}{nI(\theta)}}^2)$$

$$\Rightarrow \hat{\theta}_{MLE} \neq \theta$$

$$\text{and by definition } \ell(\hat{\theta}; X_1, \dots, X_n) > \ell(\theta; X_1, \dots, X_n)$$

$\Rightarrow$  Over any dataset, this will always

$$\text{be the case } \Rightarrow E[\ell(\hat{\theta}_{MLE}; X)] > E[\ell(\theta; X)]$$

It can be shown that this bias is  $K_m$ . This

$$l(\hat{\theta}_{m_1}, \dots, \hat{\theta}_{m_{K_m}}; X_n, Y_n) \approx l(\hat{\theta}_{m_1}^{true}, \dots, \hat{\theta}_{m_{K_m}}^{true}; X_n, Y_n) - K_m \quad \text{Very approximate!}$$

Need large  $n$   
to be true

Therefore,  $\max_{\text{model } M} \{ \downarrow -K_m \}$

For historical reasons we use  $-2 \times$  and take the min:

$$AIC_m := -2 l(\hat{\theta}_{m_1}^{true}, \dots, \hat{\theta}_{m_{K_m}}^{true}; X_n, Y_n) + 2K_m$$

the model with lowest  $AIC_m$  is "selected",  $AIC_{m^*}$

"Akaike's Information Criterion" (Akaike, 1973)

Note: if all models have same # of params, the highest  $l(\hat{\theta}_{true})$  is selected still

$$\text{Foster } w_m := e^{-(AIC_m - AIC_{m^*})/2} = e^{-\frac{(-2l_n + 2K_m) - (-2l_{n^*} + 2K_{m^*})}{2}}$$

$$= e^{(l_n - l_{n^*}) + (K_m - K_{m^*})} = \frac{L_m}{L_{m^*}} \frac{e^{K_m}}{e^{K_{m^*}}}, \text{ i.e. prop of likelihood}$$

is called the "weight".

Thm: if one of the  $M$  models is correct, the weights are the prob's of the model being true.

For small  $n$ , the AIC is not accurate. There is a small  $n$  correction:

$$AICC_m := -2l_n + 2K_m \left( \frac{n}{n - K_m - 1} \right)$$

If all models have same # of params  $K$ , nothing changes

$\uparrow$   
inflates penalty for small  $n$ .



## Troubles... Paradoxes...

① Dummy pt. estimation, we sometimes get pt. estimates we know are silly.

Eg. ObP:  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(0)$

$$X_1=0, X_2=0, X_3=0 \quad \hat{\theta} = \bar{x} = \frac{0+0+0}{3} = 0.$$

This seems overly harsh... saying the Head on the coin is impossible.

Consider  $X_1, \dots, X_n \stackrel{iid}{\sim} N(0, 0.2)$

②  $\hat{CI}_{0,1-\alpha} = \left[ \bar{x} \pm t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right]$  which is exact.

Does this CI mean anything? You could say  $P(\theta \in \hat{CI}_{0,1-\alpha} | \bar{X}) = 1-\alpha!$

③ Consider  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(0)$

$$X_1=0, X_2=1, X_3=0$$

$$\hat{CI}_{0,95\%} = \left[ \frac{1}{3} \pm 1.96 \sqrt{\frac{\frac{1}{3} \frac{2}{3}}{3}} \right] = [-0.20, 0.67] \neq \textcircled{H} \text{ which is absurd}$$

This happens due to the approximate nature of the CI.

These approximations get us into trouble

④ Run a test. Get  $p_{val} = 2\% \Rightarrow$  Reject  $H_0$ .

$p_{val} = P(H_0 | \bar{X})$ ? Does it mean any prob! No, it's just the smallest  $\alpha$  that Rejects  $H_0$ .

⑤  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(0)$  but you know  $\textcircled{H} = [0.3, 0.73] \neq [0, 1]$

due to some physical limitation. No way to account for this!!

MANY MORE...

and more

All these problems could be solved, if you allow for  $P(\theta|\vec{x})$  to be a fully fledged legitimate probability density. Why the fuss? Let  $\theta_R$  be the real value of  $\theta$

Assume  $X_1, \dots, X_n$  are i.i.d.

proper prior

$$P(\theta|\vec{x}) = \frac{P(\vec{x}|\theta) P(\theta)}{P(\vec{x})} = \frac{P(\vec{x}|\theta) P(\theta)}{\sum_{\theta \in \Theta} P(\vec{x}|\theta) P(\theta)}$$

joint prob called likelihood prior distr marginal likelihood distr

$$P(\theta) = \text{Deg}(\theta) = \begin{cases} 1 & \text{if } \theta = \theta_R \\ 0 & \text{if } \theta \neq \theta_R \end{cases}$$

$$\Rightarrow P(\theta|\vec{x}) = \begin{cases} \frac{P(\vec{x}|\theta=\theta_R)}{P(\vec{x}|\theta=\theta_R)} = 1 & \text{if } \theta = \theta_R \\ 0 & \text{if } \theta \neq \theta_R \end{cases}$$

If  $X_1, X_2, \dots$  are i.i.d.

$$P(\theta|\vec{x}) = \frac{L(\vec{x}|\theta) P(\theta)}{L(\vec{x})} = \frac{L(\vec{x}|\theta) P(\theta)}{\sum_{\theta \in \Theta} L(\vec{x}|\theta) P(\theta)}$$

This is the source of the incoherence. If this is accepted as non-degenerate

$\Rightarrow P(\theta|\vec{x})$  becomes non-degenerate.

This is the big leap, the big pill to swallow! Can you call this a real rv! It represents uncertainty in  $\theta$  before data is seen - you can encode any ideas into this.

Frequentists scream: this is subjective! You can allow your own ideas to influence science? Absurd! We will address this concern later.

Can  $\theta$  still be considered fixed? Yes,  $P(\theta|\vec{x})$  just codes the uncertainty in it. Or... you can go beyond and consider  $\theta$  in a quantum state  $|\rightarrow| \leftarrow|$  that respects conditions.