# Lecture 16:

$\hat{\hat{\theta}}_{1.} = \frac{n_{1.}}{n}$ proportion of people with black

HB

hair.

$$\Rightarrow \hat{\hat{\theta}}_{i.} = \frac{n_{i.}}{n} \quad , \quad \hat{\theta}_{.j} = \frac{n_{.j}}{n}$$

$$\Rightarrow \hat{\phi} = \sum_{i=1}^{r} \sum_{j=1}^{c} \frac{(\theta_{i.j} - n \hat{\theta}_{i.} \hat{\theta}_{.j})}{\underbrace{n \hat{\theta}_{i.} \hat{\theta}_{.j}}_{E_{ij}}} \xrightarrow{d} X_{(r-1)}^{\,} {}_{(c-1)}$$

a fact that could be
proved in 36 if time
permitted

" Chi-squared test of independence"

$\overset{\uparrow}{\phi} = 41.20 \quad \alpha = 5\% . \quad \bar{F}_{X_9}^2 = (16.99) = 95\%$

$\Rightarrow$ Reject Ho

$RET = [0, 16.99]$

we conclude that hair color and

eye color are dependent.

This class has focused mainly on three goals of inference: estimation, testing and confidence sets. We will continue to study these three goals but... we will also study some tangential "meta" concept that are classic.

Here is one such classic "meta concept". usually you are

(a) given a dataset $x_1, \ldots x_n$ then you

(b) assume a DGP, then you

(c) define one or (many) inferential target parameters, theta, the

(d) compute $\hat{\theta} = \omega_1(x_1 \ldots x_n)$ and then you

(e) make a CI/run a test at size $\alpha$.

(f) examine (b). How do you just "assume a DGP"? Sometimes you really do know the DGP e.g a coin flipped repeatedly is iid Bernoulli $(\theta)$, a die rolled repeatedly is iid uniform discrete. But what about " daily wind speeds at JFK airport" or " rat survival times" (like on midterm) or " daily percentage returns at the S&P 500"? The DGPs for the last three are very complicated and they're unknown.

What if we wanted to "guess" the DGP's? This is actually a really big part of what

staticians do. This is called "model fitting". DGP = model. Model you kinda make up and hopefully they are useful for whatever you are doing. Why don't we proceed as follows:

let's guess M candidate models / DGPs $m = 1, 2, ..., M$ and then

(1) pick the "best" model out of my M guesses and maybe

(2) provide a weighting score to each of the M guesses (low scores indicate bad guesses and high scores indicate good guesses). Goal (1) is famous and called "Model Selection". In 342, we do a

little of this atheoretically. Here we'll
do it more theoretically.

Model Selection is more fundamental than
you realize. It's actually the entire
problem of all of science. For example,
let's say you have some astronomical
data on movement of different planets,
stars, etc. You want to fit a model (guess
a DGP) for the force on celestial bodies
with masses $m_1$ and $m_2$ at a distance $r_0$
from eachother (i.e "gravity"). Consider
the following models,

Model 1: $F = h \frac{m_1 m_2}{r^2}$   Newton's law

Model 2: $F = a_1 \frac{m_1 m_2}{r^2} + a_2 \frac{m_1 m_2}{r^-}$   Newton's extension

Model 3: $F = a_1 \frac{m_1 m_2}{r^2} e^{-a_2 r}$   Laplace extension

Which model is best? We know all these are wrong because Einstein came and disprove them ~~and~~ with general relativity.

---

Lets talk about Model Selection techniques. Our data $x_1 \ldots x_n$ comes from an unknown DGP. Here are M candidate models:

Model 1: $f_1(x_1 \ldots x_n; \theta_{11}, \theta_{1K}) = \mathcal{L}(\theta_{11}, \ldots, \theta_{1K_1} \mid x_1 \ldots x_n)$

Model 2: $f_2(x_1 \ldots x_n; \theta_{21} \ldots \theta_{2K_2}) = \mathcal{L}(\theta_{21} \ldots, \theta_{2K_2} \mid x_1 \ldots x_n)$

$\vdots$

Model M: $f_m(x_1 \ldots x_n; \theta_{m1} \ldots \theta_{mK_m}) = \mathcal{L}(\theta_{m1} \ldots \theta_{mK_m} \mid x_1 \ldots x_n)$

$K_1$ is the # of parameters in model 1,

$K_2$ is the num of parameter in model $2$, $K_m$ is the num of parameter in to model $M$. Each $K_m$ could be different.

Why don't we just select the model that has the highest likelihood?

$$m_? := \underset{m}{\text{argmax}} \left\{ L_M\left(\theta_{M_1} \cdots \theta_{M_{K_m}}; x_1 \cdots x_n\right) \right\}$$

$$= \underset{m}{\text{argmax}} \left\{ l_m\left(\theta_{M_1} \cdots \theta_{M_{K_m}}; x_1 \cdots x_n\right) \right\}$$

The problem with this is we don't know the Values of $\theta$ for any of the models!

So let's do rechardize $K_1 + K_2 + \cdots K_m$ times! we'll estimate each of the parameters

using MLE'S and plug them all in
and then

$$m_d = \text{argmax} \left\{ l\left( \hat{\theta}^{MLE}_{M_1}, \ldots \hat{\theta}^{MLE}_{mk_m}; x_1..x_n \right) \right\}$$

You could do this. But... it will not
give you the best model. why?

$l\left( \hat{\theta}^{MLE}_{m_1}, \ldots \hat{\theta}^{MLE}_{mk_m}; x_1..x_n \right)$ is an estimator

for $l\left( \theta_{m_1}, \ldots \theta_{mk_m}; x_1..x_n \right)$

and it's biased.... with many
assumptions, you can prove that

$\text{Bias}\left[ l\left( \hat{\theta}^{MLE}_{m_1}, \ldots \hat{\theta}^{MLE}_{mk_m}; x_1..x_n \right) \right] = K_m > 0$

There is positive Bias (meaning
the log-likelihood would appear

higher on average) and this bias increases
you use more parameters in your
candidate models. This was figured out by
H. Akaike, a Japanese statistician
and he published it in 1974.
Once you have "the bias, you
just use it to correct your
estimate:

$$\ell(\theta_{m_1}, \dots \theta_{mk_m}) \approx \ell(\hat{\theta}_{m_1}^{MLE}, \dots \hat{\theta}_{mk_m}^{MLE}; X_1 \dots X_n)$$

$$- K_m.$$

Recall that log-likelihood is always
negative for discrete DUP's and almost
always negative for continuous DUP's
So let's flip it's sign and
multiply by 2:

Complexity penalty.

$$AIC_m = -2\ell(\hat{\theta}_{m_1}^{MLE}, \dots \hat{\theta}_{mk_m}^{MLE}, X_1 \dots X_n) + 2k_m$$

(Akaike's Information Criterion)

The "best" log likelihood is the largest i.e. closest to zero

so once negated, the "best" negative log likelihood is the smallest i.e. closest to zero.

$$m_\alpha = argmax ($$