

$$m_K := \arg \min \{AIC_m\}$$

$$AIC_1, AIC_2, \dots, AIC_M. \text{ Let } AIC_K = \min \{AIC_m\}$$

$$\text{let } w_m = \frac{e^{-(AIC_m - AIC_K)}}{\sum_{j=1}^M e^{-(AIC_j - AIC_K)}}$$

Akaike weights (sum to 1)

$$e^{-(AIC_m - AIC_K)/2} = e^{-(\ell_m - \ell_K) - (2K_m - 2K_K)/2}$$

$$= e^{(\ell_K - \ell_m) + (K_K - K_m)} = \frac{e^{\ell_K}}{e^{\ell_m}} \frac{e^{K_K}}{e^{K_m}} = \frac{\ell_K}{\ell_m} \frac{e^{K_K}}{e^{K_m}} \quad \text{kind of like a ratio of probabilities}$$

If the "true model/DGP" is one of the candidate models, then w_m is the probability that m is the true model.

Beyond the scope of this class, people use ~~the~~ Akaike weights to create a model which is ~~the~~ an average over the candidates:

$$f = \sum_{m=1}^M w_m f(x_1, \dots, x_n; \hat{\theta}_{m1}^{MLE}, \dots, \hat{\theta}_{mK_m}^{MLE}) \quad \text{mixed model}$$

It turns out that in low n situations, the bias is very incorrect, so there's a correction term that is employed to fix the bias and make the AIC more accurate, it's called "AIC-corrected" or

$$AICC: AICC_m = -2 \left(\ell_m \hat{\theta}_{m1}^{MLE}, \dots, \hat{\theta}_{mK_m}^{MLE}; x_1, \dots, x_n \right) + 2K_m \left(\frac{n}{n - K_m - 1} \right)$$

inflates the penalty.

Midterm II ↑
Final ↓

The concept of "practical significance" (or "clinical significance" if you happen to be in a medical/health context). Let's say you're testing

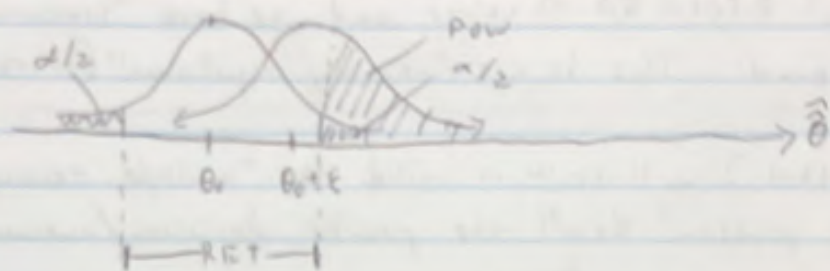
$$H_a: \theta \neq \theta_0 \text{ vs. } H_0: \theta = \theta_0$$

but the true value of the parameter is $\theta_0 + \epsilon$ where ϵ is a small number. So H_0 is technically false. Now you won't be able to reject H_0 unless your n is very high because power to find small effects is low. But ... given enough n , you always reject for any ϵ and any α .

Proof: = assume $\hat{\theta}$ is asymptotically normal and ϵ is positive (for HW you'll prove it for negative). This means:

$$\hat{\theta} | H_0 \sim N(\theta_0, \overbrace{SE[\hat{\theta}](\theta_0)^2}^{\sigma/\sqrt{n}}) = N(\theta_0, (\frac{\sigma}{\sqrt{n}})^2)$$

$$\hat{\theta} \sim N(\theta_0 + \epsilon, \overbrace{SE[\hat{\theta}](\theta_0 + \epsilon)^2}^{\sigma/\sqrt{n}}) = N(\theta_0 + \epsilon, (\frac{\sigma}{\sqrt{n}})^2)$$



$$\begin{aligned} \text{Pow} &= P(\text{Reject } H_0) = P(\hat{\theta} > \theta_0 + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}) = P(\hat{\theta} - (\theta_0 + \epsilon) > \theta_0 + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} - (\theta_0 + \epsilon)) \\ &= P(z > \underbrace{-\epsilon}_{\frac{\sigma}{\sqrt{n}}} + \underbrace{z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}}_{\frac{z_{1-\frac{\alpha}{2}}}{\frac{\sigma}{\sqrt{n}}}}) = P(z > \underbrace{-\frac{\epsilon}{\sigma}}_{\frac{1}{\sqrt{n}}} + \underbrace{z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}}_{\frac{z_{1-\frac{\alpha}{2}}}{\frac{\sigma}{\sqrt{n}}}}) \end{aligned}$$

$$\lim_{n \rightarrow \infty} \text{Pow} = \lim_{n \rightarrow \infty} P(z > -\infty) = 1 \Rightarrow \text{Your estimate is "statistically significant" } P_{val} < \alpha$$

So one can argue that in the real world, θ is never exactly some value you propose. Take the case of a coin. You want to prove it's unfair. $H_0: \theta = 50\%$ exactly. But take a look at this coin... the real θ is probably 49.9999% so $\epsilon = 0.0001\%$. But... is a coin with $\text{Prob}(\text{heads}) = 49.9999\%$ actually "unfair"? NO. Practically this coin is fair (for all practical uses).

Thus if you flip enough times, you will get a "statistically significant" estimate that has "no practical significance."

No amount of math can tell you what "practical significance" means. You have to define it yourself based on your own context and your own objective.

As an example. Let's say you're testing a weight loss pill so you

randomly give n_T subjects the pill ($T = \text{treatment}$) and randomly give n_C subjects the placebo ($C = \text{control}$) and then you measure $\bar{X}_T - \bar{X}_C$ and run a test $H_0: \theta_T = \theta_C$ (no mean difference between pill group and control group.) Where θ_T is mean weight loss in the pill group and θ_C is mean weight loss in control group.

You get a p value of $0.001 < 5\% \Rightarrow$ reject and you have "statistical significance." But ... $\bar{X}_T - \bar{X}_C = 0.1$ pounds. This is not "clinically significant" (that's our feeling).

Next "meta concept" that I will cover is called the "multiple testing problem" or "multiple comparisons problem." Recall the possible decisions/outcomes from a hypothesis test:

		Decision	
		Retain H_0	Reject H_0
Truth	H_0	Justified RET	Type I error
	H_A	Type II error	Justified REJ

We "control" the probability of Type I errors by setting it to be at most ... $P(\text{Type I error}) \leq \alpha$

Let's say you're doing m hypothesis tests (many of them) each with α (controlled Type I error probability). This collection of tests is called a "family of tests." Among these tests, you reject r of them and retain f of them so that $r + f = m$. But unbeknownst to you, you could've made some Type I or Type II errors. Here's a contingency table with the number of each possibility:

		Decision	
		Retain H_0	Reject H_0
Truth	H_0	u	v
	H_A	t	s
		f	r
		observed/realizations	

v is the number of Type I errors
AKA "false rejections" AKA "false discoveries."

you don't know any of these values.

Furthermore, which quantities are random? And the randomness is due to the DGP. The ones with capital letters below:

Decision

Truth		Reject H_0	Reject H_0	
		V	S	
	H_0	V	S	m_0
	H_a	T	R	m_1

What if you want control over the number of Type I errors v i.e. you want to control the r.v. V ? Previously, in the context of $m=1$, you decided α , which controlled v/V to the level of your comfort. Parenthetically, control of Type II errors is done by maximizing the power for each test. So we won't talk about it.

Why do we care about controlling the Type I errors? Here's an example that should get you scared. Let's say you have m independent hypothesis tests each with size $\alpha = 5\%$. Also, let $m = m_0 = 30$ i.e. H_0 's are true. By chance alone,

$$R \sim \text{Bin}(m, \alpha)$$

I'm interested in the chance I make at least one Type I error. By simple 241, calculation,

$$\begin{aligned} P(R > 0) &= 1 - P(R = 0) = 1 - P(\text{all correct}) = 1 - (1 - \alpha)^m \\ &= 1 - (1 - 5\%)^{30} \approx 76\% \end{aligned}$$

So there is a ~~large~~ huge chance you make at least one "false discovery." Maybe this 76% probability is too high for you.