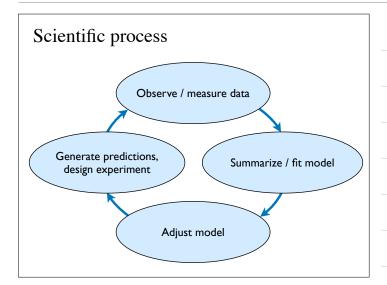
Mathematical Tools for Neural and Cognitive Science

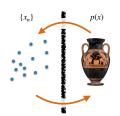
Fall semester, 2025

Section 5: Statistical Inference and Model Fitting



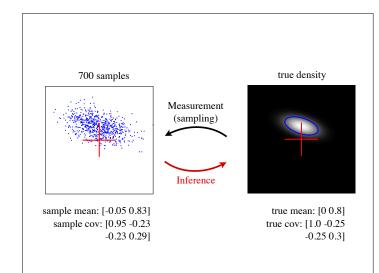
The sample average

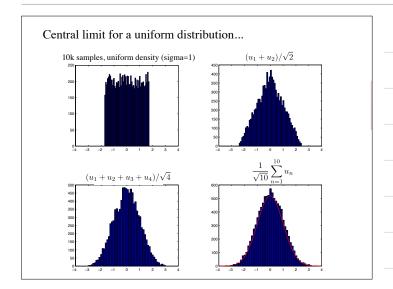
$$\bar{x} = \frac{1}{N} \sum_{n=1}^{N} x_n$$

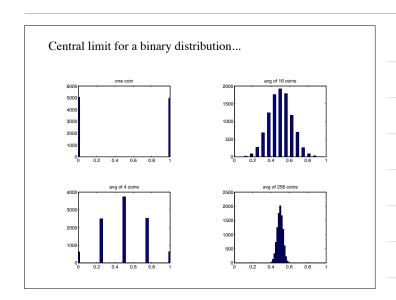


What happens as N increases?

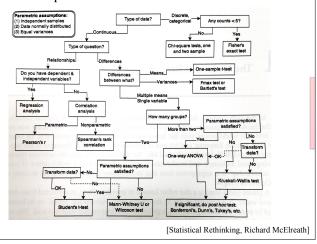
- Variance of \bar{x} is σ_x^2/N (the "standard error of the mean", or SEM), and so converges to zero <code>[on board]</code>
- "Unbiased": \bar{x} converges to the true mean, $\mu_x = \mathbb{E}(\bar{x})$ (formally, the "law of large numbers") [on board]
- The distribution $p(\bar{x})$ converges to a Gaussian (mean μ_x and variance σ_x^2/N): formally, the "Central Limit Theorem"







Classical "frequentist" statistical tests



Classical/frequentist approach - z

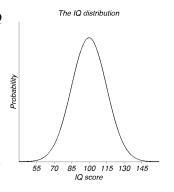
• In the general population, IQ is known to be distributed normally with:

$$\mu = 100, \ \sigma = 15$$

- We give a drug to 30 people and test their IQ
- Hypotheses:

 H_1 : NZT improves IQ

 H_0 ("null"): it does nothing



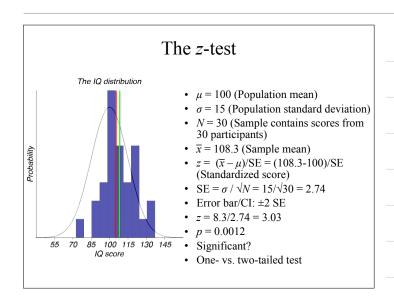
Test statistic

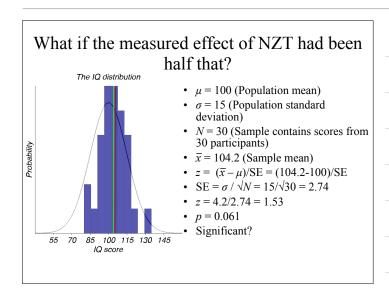
- We calculate how far the observed value of the sample average is away from its expected value.
- In units of standard error.
- In this case, the test statistic is

$$z = \frac{\overline{x} - \mu}{SE} = \frac{\overline{x} - \mu}{\sigma / \sqrt{N}}$$

• Compare to a distribution, in this case z or N(0,1)

Does NZT improve IQ scores or not?						
	Reality					
		Yes	No			
Decision	Yes	Correct	Type I error α-error "False alarm"			
	No	Type II error β -error "Miss"	Correct			





Significance levels

- Are denoted by the Greek letter α .
- In principle, we can pick anything that we consider unlikely.
- In practice, the consensus is that a level of 0.05 or 1 in 20 is considered as unlikely enough to reject H₀ and accept the alternative.
- A level of 0.01 or 1 in 100 is considered "highly significant" or "really unlikely".

Common misconceptions



Is "Statistically significant" a synonym for:

- Substantial
- Important
- Big
- · Real

Does statistical significance gives the

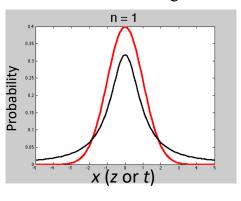
- probability that the null hypothesis is true
- probability that the null hypothesis is false
- probability that the alternative hypothesis is true
- probability that the alternative hypothesis is false

Meaning of *p*-value. Meaning of CI.

Student's *t*-test

- σ not assumed known
- Use $s^{2} = \frac{\sum_{i=1}^{N} (x_{i} \overline{x})^{2}}{N}$
- Why *N*-1? *s* is unbiased (unlike ML version), i.e., $\mathbb{E}(s^2) = \sigma^2$
- Test statistic is $t = \frac{\overline{x} \mu_0}{s / \sqrt{N}}$
- Compare to t distribution for CIs and NHST
- "Degrees of freedom" reduced by 1 to N-1

The t distribution approaches the normal distribution for large N



The z-test for binomial data

- Is the coin fair?
- Lean on central limit theorem
- Sample is *n* heads out of *m* tosses
- Sample mean: $\hat{p} = n / m$
- H_0 : p = 0.5
- Binomial variability (one toss): $\sigma = \sqrt{pq}$, where q = 1 p
- Test statistic: $z = \frac{\hat{p} p_0}{\sqrt{p_0 q_0 / m}}$
- Compare to z (standard normal)
- For CI, use $\pm z_{\alpha/2} \sqrt{\hat{p}\hat{q}/m}$

Other frequentist univariate tests

- χ^2 goodness of fit
- χ^2 test of independence
- test a variance using χ^2
- F to compare variances (as a ratio)
- Nonparametric tests (e.g., sign, rank-order, etc.)

Estimation of model parameters (outline)

- How do I estimate parameters from data?
- How "good" are my estimated parameters?
- How well does my model explain data to which it was fit? Other data (prediction/generalization)?
- How do I compare two models?

Estimation

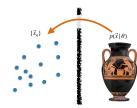
- An "estimator" is a function of the data, intended to provide an approximation of the "true" value of a parameter
- One can evaluate estimator quality in terms of squared error, MSE = bias^2 + variance
- Traditional statistics often aims for an unbiased estimator, with minimal variance ("MVUE")
- More nuanced view: trade off bias and variance, through model selection, "regularization", or Bayesian "priors" ...

The maximum likelihood estimator (MLE)

Sample average is appropriate when one has direct measurements of the thing being estimated. But one may want to estimate something that is *indirectly* related to the measurements...

Natural choice: assuming a probability model $p(\vec{x} \mid \theta)$ find the value of θ that maximizes this "likelihood" function

$$\begin{split} \hat{\theta}(\{\vec{x}_n\}) &= \arg\max_{\theta} \prod_{n} p(\vec{x}_n|\theta) \\ &= \arg\max_{\theta} \sum \log p(\vec{x}_n|\theta) \end{split}$$



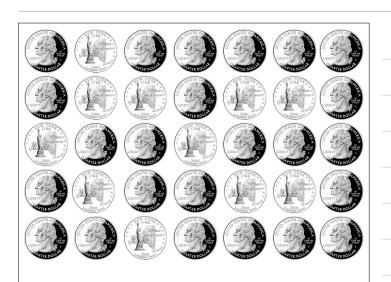
Example: Estimate the true probability of a flipped coin landing "heads" up, by observing some samples







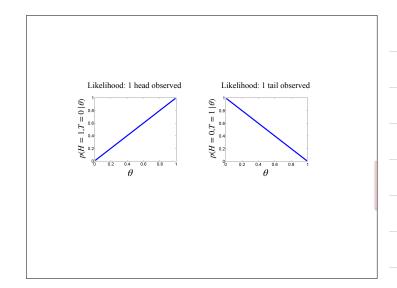
66%?

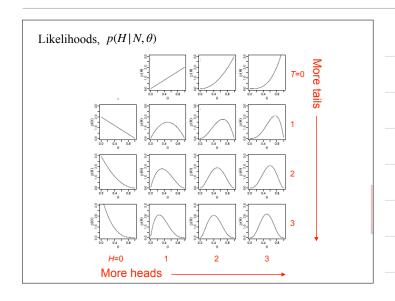


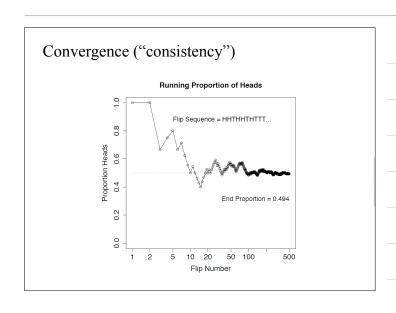
Example ML Estimators - discrete

Binomial: $p(H|N,\theta) = \binom{N}{H} \theta^{H} (1-\theta)^{N-H}$

(H = # heads observed, in N flips of a coin, with probability of heads θ)







Example ML Estimators - discrete

Binomial:
$$p(H|N,\theta) = \binom{N}{H} \theta^H (1-\theta)^{N-H}$$
 (H = # heads observed, in N flips of a coin, with probability of heads θ)

Poisson:
$$p(\{k_n\} | \theta) = \prod_{n=1}^{N} \frac{\theta^{k_n} e^{-\theta}}{k_n!}$$
 (k's are measured counts, with mean arrival rate of θ)

[on board]

Example ML Estimators - continuous

Uniform:
$$p(x|\theta) = \begin{cases} \frac{1}{\theta} & 0 \le x \le \theta \\ 0 & \text{otherwise} \end{cases}$$

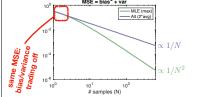
$$\hat{\theta} = \max_{n} \{x_n\}$$
 (Note: this is biased!)

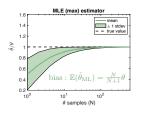
Two estimators for range of a uniform distribution

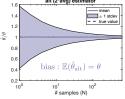
Given N samples $\{x_n\}$ from the uniform distribution over $[0,\theta]$, consider two estimators of θ :

$$\hat{\theta}_{\mathrm{ML}}(\{x_n\}) = \max_{n}(x_n)$$

$$\hat{\theta}_{\mathrm{alt}}(\{x_n\}) = \frac{2}{N} \sum_{n} x_n$$

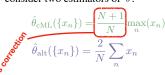


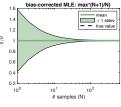


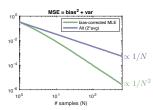


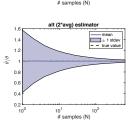
Two estimators for range of a uniform distribution

Given N samples $\{x_n\}$ from the uniform distribution over $[0, \theta]$ consider two estimators of θ :









Example ML Estimators - Continuous

Uniform: $p(x|\theta) = \begin{cases} \frac{1}{\theta} & 0 \le x \le \theta \\ 0 & \text{otherwise} \end{cases}$

 $\hat{\theta}_{\text{ML}} = \max_n \{x_n\}$ (Note: this is biased!) $\hat{\theta}_{\text{cML}} = \frac{N+1}{N} \hat{\theta}_{\text{ML}}$

Gaussian: $p(x|\mu,\sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

 $a = \frac{\sum_{n} x_n}{\sum_{n} x_n}$ (sample average, again)

 $\hat{\sigma}_{_{\rm ML}}^2 = \frac{\sum_n (x_n - \hat{\mu})^2}{N} \hspace{1cm} \text{(Note: this is biased!)} \\ \hat{\sigma}_{_{\rm cML}}^2 = \frac{N}{N-1} \hat{\sigma}_{_{\rm ML}}^2$

[on board]

Summarizing errors of ML estimators

Bias: the MLE is *asymptotically unbiased* and *Gaussian*, but can only rely on these if:

- the likelihood model is correct
- the likelihood can be maximized
- you have lots of data

Variance: (error bars)

- S.E.M. (relevant for sample averages only)
- second deriv of NLL (multi-D: "Hessian")
- simulation (resample from $p(x|\hat{\theta})$)
- bootstrapping (resample from *the data*, with replacement)

Bootstrapping

- "The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps" [Adventures of Baron von Munchausen, by Rudolph Erich Raspe]
- A (re)sampling method for computing estimator dispersion (eg., stdev or confidence intervals)
- Idea: instead of looking at distribution of estimates across repeated experiments, look across repeated resamplings (with replacement) from the existing data ("bootstrapped" data sets)

FOUND TO BE CUT BY TAKING ASPIRIN

LIFESAVING EFFECTS SEEN

Study Finds Benefit of Tablet Every Other Day Is Much Greater Than Expected

The summary statistics in the newspaper article are very simple: heart attacks subjects

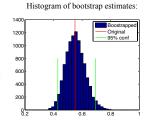
(fatal plus non-fatal) 104 189 aspirin group: placebo group:

 $\widehat{\theta} = \frac{104/11037}{189/11034} = .55.$

If this study can be believed, and its solid design makes it very believable, the aspirin-takers only have 55% as many heart attacks as placebo-takers. Of course we are not really interested in $\hat{\theta}$, the estimated ratio. What we would like to know is θ , the true ratio

[Efron & Tibshirani '98]

| HEART ATTACK RISK | [New York Times, 27 Jan 1987]



=> with 95% confidence,

 $0.43 < \theta < 0.7$

	strokes	subjects	
aspirin group:	119	11037	
placebo group:	98	11034	(1.3)

For strokes, the ratio of rates is

$$\widehat{\theta} = \frac{119/11037}{98/11034} = 1.21. \tag{1.4}$$

It now looks like taking aspirin is actually harmful. However the interval for the true stroke ratio θ turns out to be

$$.93 < \theta < 1.59$$
 (1.5)

with 95% confidence. This includes the neutral value θ = 1, at which aspirin would be no better or worse than placebo vis-à-vis strokes. In the language of statistical hypothesis testing, aspirin was found to be significantly beneficial for preventing heart attacks, but not significantly harmful for causing strokes.

[Efron & Tibshirani '98]

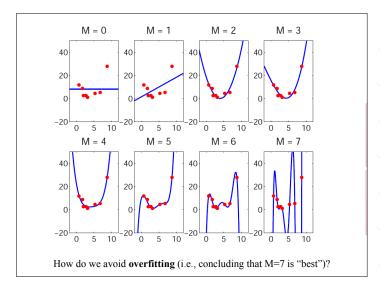
Permutation test

- Given {n1,n2} measurements under two different conditions, are they significantly different (i.e., can we reject null hypothesis?)
- Measure difference in means, m2-m1
- Construct permuted sets of {n1,n2} measurements, and compute difference in means for each of these
- Ask: How far in the tail is the true difference in means? One-sided p-value is proportion of permutation values > m2-m1

Taxonomy of model-fitting errors

- Unexplainable variability (e.g., due to noisy measurements)
- Optimization failures (e.g., local minima)
- Overfitting (too many params, not enough data)
- Model failures (what you'd really like to know)

Optimization... Heuristics, exhaustive search, (pain & suffering) Iterative descent, possibly non-unique (local minima) Iterative descent, unique Closed-form, and unique



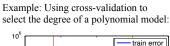
Model Comparison

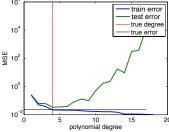
- If models are optimized according to some objective, it is natural to compare them based on the value of that objective...
 - for least squares regression, compare the residual squared error of two models (with different regressors).
 - for ML estimates, compute the likelihood (or log likelihood) ratio, and compare to 1 (or zero).
- **Problem**: evaluating the objective with the same data used to optimize the model leads to **over-fitting!**

Cross-validation

A resampling method for estimating predictive error of a model. Widely used to identify/avoid over-fitting, and to provide a fair comparison of models.

- (1) Randomly partition data into a "training" set, and a "test" set.
- (2) Fit model to training set. Measure error on test set.
- (3) Repeat (many times).
- (4) Choose model that minimizes the average "cross-validated" (test) error





Ridge regression

(a.k.a. *L*₂ regularization)

Ordinary least squares regression:

$$\arg\min_{\vec{\beta}} ||\vec{y} - X\vec{\beta}||^2$$

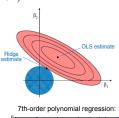
"Regularized" least squares regression:

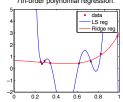
$$\arg\min_{\vec{\beta}}||\vec{y}-X\vec{\beta}||^2+\lambda||\vec{\beta}||^2$$

Equivalent formulation: MAP estimate, assuming Gaussian likelihood & prior!

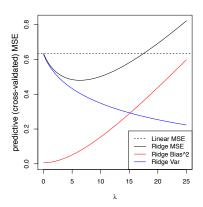
$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

Choose lambda by cross-validation:





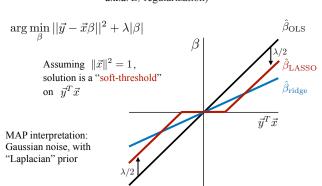
Ridge Regression trades off bias and variance:

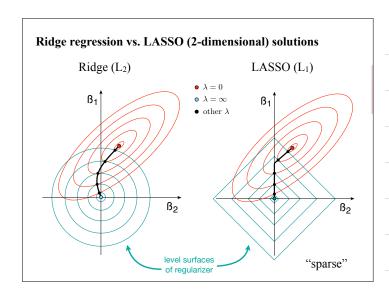


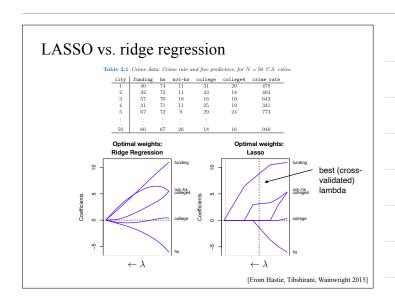
[from http://www.stat.cmu.edu/~ryantibs/datamining/]

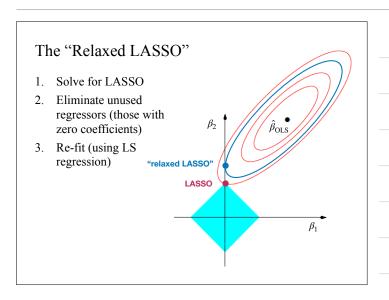
LASSO regularization

("least absolute shrinkage and selection operator", a.k.a. L_I regularization)









[Insert LNP neural fitting example here]

Bayesian Inference

"Posterior" "Likelihood" "Prior"
$$p(\theta \mid \text{data}) = \frac{p(\text{data} \mid \theta)p(\theta)}{p(\text{data})}$$
Normalization factor

Example: Posterior for coin

infer whether a coin is fair by flipping it repeatedly here, x is the probability of heads (50% is fair) $y_{1\dots n}$ are the outcomes of flips

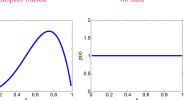


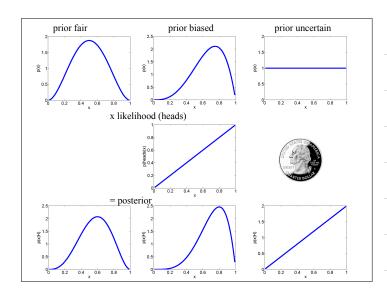
Consider three different priors: suspect fair

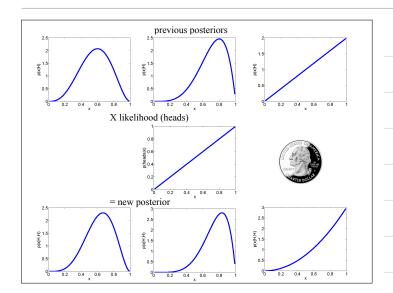
2

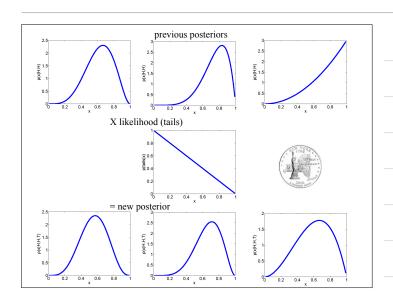


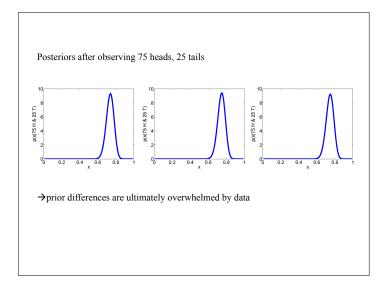
suspect biased







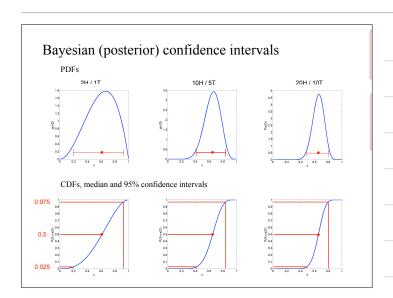




Bayesian estimators

Summarize **central tendency** of posterior:

- Posterior mode ("maximum aposteriori estimate" MAP)
- Posterior mean (minimizes squared error MMSE) Summarize dispersion with posterior variance
- Posterior median (minimizes abs error) Summarize dispersion with posterior quantiles



Bayesian inference: Gaussian case

For measurements with Gaussian noise, and assuming a Gaussian prior:

- posterior is Gaussian, allowing sequential updating
- precision is sum of measurement and prior precisions
- mean is precision-weighted average of prior mean and measurement
- explains "regression to the mean" as shrinkage toward the prior

Bayesian inference: Gaussian case

$$y = x + n$$
, $x \sim N(\mu_x, \sigma_x)$, $n \sim N(0, \sigma_n)$

$$p(x|y) \propto p(y|x)p(x)$$

$$\propto e^{-\frac{1}{2}\left[\frac{1}{\sigma_n^2}(x-y)^2\right]}e^{-\frac{1}{2}\left[\frac{1}{\sigma_x^2}(x-\mu_x)^2\right]}$$

$$= e^{-\frac{1}{2}\left[\left(\frac{1}{\sigma_n^2} + \frac{1}{\sigma_x^2}\right)x^2 - 2\left(\frac{y}{\sigma_n^2} + \frac{\mu_x}{\sigma_x^2}\right)x + \dots\right]}$$

This is Gaussian, with:

$$\frac{1}{\sigma^2} = \frac{1}{\sigma_n^2} + \frac{1}{\sigma_x^2}$$

$$\mu = \left(\frac{y}{\sigma_n^2} + \frac{\mu_x}{\sigma_x^2}\right) / \left(\frac{1}{\sigma_n^2} + \frac{1}{\sigma_x^2}\right)$$

The average of y and μ_x , weighted by inverse variances (a.k.a. "precisions")!

likelihood

Regression to the mean

"Depressed children treated with an energy drink improve significantly over a three-month period. I made up this newspaper headline, but the fact it reports is true: if you treated a group of depressed children for some time with an energy drink, they would show a clinically significant improvement....

"It is also the case that depressed children who spend some time standing on their head or hug a cat for twenty minutes a day will also show improvement."

- D. Kahneman

Two noisy measurements of the same variable:

$$y_1 = x + n_1 x \sim N(0, \sigma_x)$$

$$y_2 = x + n_2$$
 $n_k \sim N(0, \sigma_n)$, independent

$$C_y = \begin{bmatrix} \sigma_x^2 + \sigma_n^2 & \sigma_x^2 \\ \sigma_x^2 & \sigma_x^2 + \sigma_n^2 \end{bmatrix}$$

LS Regression:

$$\hat{\beta} = \arg\min_{\beta} \mathbb{E} \left[\|y_2 - \beta y_1\|^2 \right]$$

$$= \frac{\mathbb{E}[y_1 y_2]}{\mathbb{E}[y_1^2]} \ = \ \frac{\sigma_x^2}{\sigma_x^2 + \sigma_n^2}$$

$$\mathbb{E}(y_2|y_1) = \hat{\beta} \ y_1$$

"regression to the mean"



y₁ Least-squares regression

TLS regression -3 (largest eigenvector) -3 0

Hierarchy of common statistical estimators

- Maximum likelihood (ML): $\hat{x}(\vec{d}) = \arg\max_{x} p(\vec{d}|x)$ (requires likelihood, $p(\vec{d}|x)$)
- Maximum a posteriori (MAP): $\hat{x}(\vec{d}) = \arg\max_{x} p(x|\vec{d})$ (requires prior, p(x))
- Bayes estimator (general): $\hat{x}(\vec{d}) = \arg\min_{\hat{x}} \mathbb{E}\left(L(x,\hat{x}) \mid \vec{d}\right)$ (requires loss, $L(x,\hat{x})$)
- Bayes least squares (BLS): $\hat{x}(\vec{d}) = \arg\min_{\hat{x}} \mathbb{E}\left((x-\hat{x})^2 \mid \vec{d}\right)$ (special case: squared loss) $= \mathbb{E}\left(x \mid \vec{d}\right)$

Bayesian Model Comparison

- Eg: Is the coin fair? Compared to what?
- Consider two models: $M_1: p = 0.5$ $M_2: p = 0.6$

$$p(M_k \mid D) = \frac{p(D \mid M_k)P(M_k)}{p(D)}$$

Compare their posterior ratio:

$$\frac{p(M_1 | D)}{p(M_2 | D)} = \frac{p(D | M_1)P(M_1)}{p(D | M_2)P(M_2)}$$

Comparing models' predictive performance

Option 1: Include a penalty for number of parameters:

For an ML estimate:
$$\hat{\theta} = \arg\min_{\theta} \left[-\ln p(\vec{d}|\theta) \right]$$

a. Akaike information criterion (AIC) [Akaike, 1974]

$$E_{\rm AIC}(\vec{d},\hat{\theta}) = 2 \ {\rm dim}(\hat{\theta}) - 2 \ln p(\vec{d}|\hat{\theta})$$

b. Bayesian information criterion (BIC) [Schwartz, 1978]

$$\begin{split} E_{\mathrm{BIC}}(\vec{d}, \hat{\theta}) &= \dim(\hat{\theta}) \; \ln \left[\dim(\vec{d}) \right] - 2 \ln p(\vec{d} | \hat{\theta}) \\ & \text{valid when } \dim(\vec{d}) \gg \dim(\hat{\theta}) \end{split}$$

Option 2: Cross-validation (evaluate generalization to held-out data)