

# Statistical Inference

1. Make a good guess of the parameter
2. Quantify our uncertainty

"What can be learned from the data"

"Data is given, parameter is random"

Bayesian view

• prior with marginal likelihood to get posterior

$$P(\theta = \theta_0 | Y=y) = \frac{P(Y=y | \theta = \theta_0) \cdot P(\theta = \theta_0)}{P(Y=y)}$$

← marginal over all  $\theta$

Point Estimates: Mean, Median, Mode (MAP)

Interval Estimates: Uncertainty Credible Intervals

A  $1-\alpha$  credible  $[L, U]$

$$P(\theta \text{ in between } L, U \text{ given data})$$

random  
not random

"Says something about parameter itself"

Standard Bayes: known from specified prior

Empirical Bayes: Parameterize the prior and maximize  $f(Y; k) = \int f(Y|\theta) f(\theta|k) d\theta$

\* For normal: the smaller the  $b$  (variance), flatter the prior the closer MAP = MLE and posterior  $\sim N(\theta^*, se(MLE))$

OR as  $n \rightarrow \infty$ , posterior shifts towards likelihood/MLE  
\* flat is not uniformity, because can use credible intervals of MLE

Frequentist View: data comes from some true distribution

Maximum Likelihood

$$\ln L(\theta) = \sum \ln f(Y_i; \theta)$$

$$P(\text{data} | \theta)$$

$$MLE = \arg \max_{\theta} \ln L(\theta) = \hat{\theta}_n$$

Wald Test "Approximate Normal"

$$H_0: \theta^* = \theta_0, \hat{\theta}_n \sim N(\theta^*, \hat{se}(\hat{\theta}_n))$$

$$\hat{\theta}_n - \theta_0 \approx N(0, 1)$$

$$\hat{\theta}_n \pm cv \cdot \hat{se}(\hat{\theta}_n)$$

Interval Estimates

"Don't know null? distribution"

Bootstrap

- 1 Simulate  $\hat{\theta}_n$  drawn from sample  $\hat{F}_n$  rather than  $F$
- 2 Estimate  $\hat{se}(\hat{\theta}_n)$
- 3 Construct CI

- Wald  $\pm 1.96 \hat{se}(\hat{\theta}_n)$
- Pivotal  $\alpha/2, 1-\alpha/2$  quantiles of  $2(\hat{\theta}_n - \theta_n)$
- Quantile:  $\alpha/2, 1-\alpha/2$  of  $\hat{\theta}_n$

Uncertainty Inference

- Student t
- P-value: probability of given result or more significant given null is true
- CI: over 95% of intervals will contain true parameter over new draws of data

Estimation: learn some statistic

$$Bias(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta$$

$$Var(\hat{\theta}_n) = E[\hat{\theta}_n^2] - E[\hat{\theta}_n]^2$$

$$MSE(\hat{\theta}_n) = Bias^2(\hat{\theta}_n) + Var(\hat{\theta}_n)$$



## Predictive Modeling "Understand relationships between variables"

SSR

- ① Finding Functions  $f$  st.  $y_i \approx f(x_i)$ , want  $f(x) = E[Y|X=x]$  minimizing  $E[Y - f(x)]^2$
- ② Let be  $F$  a blackbox, focus on predictive inference avg value of  $Y$  among  $X=x$

The Linear Model:  $Y = B^T X + \epsilon$   $E[\epsilon] = 0$ , becomes parametric when  $\epsilon \sim N(0, \sigma^2)$

unbiased for  $X$  with linearly independent columns

$$\hat{\beta}_n = (X^T X)^{-1} X^T Y \text{ called Ordinary Least Squares closed form solution to minimize } MSE \|Y - X\hat{\beta}_n\|^2$$

$$\text{Residuals } \hat{\epsilon}_i = y_i - X_i \hat{\beta}_n, \quad \hat{\sigma}_n^2 = \frac{1}{n-p} \| \hat{\epsilon} \|^2, \quad \hat{se}(\hat{\beta}_{n,i}) = \hat{\sigma}_n \sqrt{(X^T X)^{-1}_{ii}}$$

Predictive Inference (hypothesis of  $B$ )

Generalized Linear Models  
for binary can use LR

$$g(z) = \frac{\exp(z)}{1 + \exp(z)}$$

But use Wald Test  
Instead of Student  
T-test  
for Inference

Can use MLE:  $\hat{\beta}_{n,i} - \beta_i^* \rightarrow N(0, 1)$ ,  $H_0: \beta_i^* = 0$ ?

Wald

$$se(\hat{\beta}_{n,i})$$

Referred

Can use Student  $t$ , with  $n-p$  degrees freedom, "More conservative", larger intervals"

$$\text{If } \epsilon_i \text{ are normal, student } t \text{ is exact parametric test}$$

$$t_n = \frac{\hat{\beta}_{n,i} - \beta_{0,i}}{\hat{\sigma}_n / \sqrt{n}}$$

accounts for conditional var  $\sigma$

$$E[Y|X=x] = \hat{\beta}_n^T x \pm 1.96 \sigma \cdot \sqrt{x^T (X^T X)^{-1} x}$$

95% CI for  $Y$

## Dangers of Linear Models:

Collinearity: two copies of same variable, can't interpret Beta's

High Dimension: If data  $n <$  parameters  $p$  will overfit

Heteroscedastic Noise: Violates the assumption  $\epsilon_i \sim N(0, \sigma^2)$  and  $se(\hat{\beta}_{n,i})$  is wrong

We think significance but there is none because  $se$  is too small.

1/20 coefficients will not contain true

Multiple Testing: OLS confidence intervals are not joint CIs and need to use Bonferroni correction, can't consider statistical set simultaneously



## Improving the Linear Model (Baseline OLS)

### Cross Validation

More folds, lower bias but higher variance

Less folds: higher bias, lower variance

- Obtain unbiased estimates of MSE

Prediction error:  $\mathcal{L}^*(f) = E[(Y - f(X))^2]$  average squared loss MSE

$$R^2 = 1 - \frac{\mathcal{L}^*(f)}{\text{Var}(Y)}$$

percentage of  $Y$  variance reduced by  $X$  features

For homoskedastic noise

$$\text{Var}(Y|X) = \sigma^2 \leftarrow \text{is constant}$$

$$\text{so } R^2(f) = 1 - \frac{\sigma^2}{\text{Var}(Y)} \quad \text{minimum error we can aspire}$$

Try to improve  $R^2$  in-sample and out of sample simultaneously

Stepwise Regression on AIC  $\rightarrow$  AIC = maximum log likelihood of model - # of parameters   
  $\leftarrow$  expectation of error

### Methods

#### Regularization

Ridge + Lasso = Elastic Net

#### Shrinkage

Selection + shrinkage

- simple model w/ least variance

- Shrink coeff towards zero

- absolute sum of coeff penalty term  $\|B\|_1$

- makes terms  $O(\lambda)$  concave

$$P_n \in \arg \min_{\lambda} \frac{1}{2n} \|Y - XB\|^2 + \lambda \cdot \|B\|_2^2$$

Strength penalty term

of regularization = square of coefficient

$$\text{Prediction Error} = \text{Variance } \sigma^2 \text{ under homoskedasticity}$$

$$+ \text{bias}^2 \text{ estimating } f^*(x) + \text{Var}(\hat{f}(x))$$

uncontrolled error

irreducible

Bias-Variance Tradeoff:  $E[(Y - \hat{f}(X))^2 | X=X] \leftarrow$  expectation of error

$\rightarrow$  Bias is reduced when feature space is rich but variance is highly change on new data

$\rightarrow$  If Bias is high, model is too simple and variance is low on new data

$\rightarrow$  Error is minimized where model complexity balances error



Causal Inference: causality  $\neq$  association = prediction

Assignment at random? Yes "Average Treatment effect ATE =  $E[Y(1) - Y(0)]$

Assignment not at random

Parametric  
A/B testing

Nonparametric  
Permutation Test

Ignorability

Method 1

Method 2  
Instrumental Variables

$Y(0), Y(1) \perp T$   
independence

• Calculate ATE by two sample means

• Any observed difference causal effect or violation of non-interference consistency assumptions

• Use two sample t-test if equal variance

• Randomization ensures no alternate explanations for correlation except causation (idiosyncrasies of each unit)

• Can use Welch's for unequal var

Process of AB Test

1. Split n units randomly into two groups
2. Apply control, treatment
3. Measure outcomes
4. Use t/welch to compare difference

• Sharp null  
 $H_0: Y(0) = Y(1), TE = 0$

• Non parametric assumption when n is small and heavy tail is significant

• Big  $\hat{\Delta}_n$ ,  $H_0$  is false and small  $\hat{\Delta}_n$   $H_0 = 0$ , is plausible

• Simulate  $\Delta_n$  by random permutation of assignment to treatment

• Get a distribution over  $\hat{\Delta}_n$  and compare  $|\hat{\Delta}_n|$  to overall distribution

• Can perform for any statistic medians, ks etc

Process of Permutation

1. Calculate difference delta
2. Simulate  $\hat{\Delta}$  over many permutations / random shuffles
3. Calculate p value and show distribution

Controlling for all alternative explanations of non-causal difference

$Y(t) \perp T | X$   
(assignment at random within each  $X = x$  subpopulation)

$$E[Y(t) | X] =$$

$$E[Y | T=t, X]$$

and ATE =

$$E[E[Y | T=1, X]] - E[E[Y | T=0, X]]$$

using OLS

Process

1. Filter data
2. Run OLS
3. Treatment coeff is causal effect

• Add a variable to data to simulate randomness

• Must be ignorable, only affects outcome and affects treatment

Ignorable - does not share common effects with outcome  $Y$ . no confounding factors.

Relevance - instrument  $Z$  has effect on Treatment  $X$

Exclusion -  $Z$  affects outcome  $Y$  only through  $X$

$$\uparrow = E[Y | Z=1] - E[Y | Z=0]$$

If we get large effect  $\rightarrow P(T=1 | Z=1) - P(T=1 | Z=0)$

$\uparrow$  = Intention to treat causal effect

$\uparrow$  = Compliance w/ assigned treatment

$\uparrow$  = causal effect had everyone complied

Two Stage Least Squares

• Empirically verifies the relevance assumption

1. Predict expected value of treatment based on IV

2. Use new treatment on outcome to remove individual idiosyncrasies

Heterogeneous Effects LATE  $\neq$  ATE

IV only works for compliers

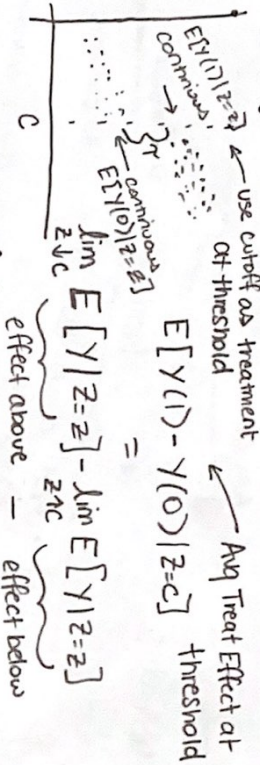
$$LATE = E[Y(1) - Y(0) | \text{compliers}]$$



"Closer to gold for causal effects"

## Causal Inference Method: Regression Discontinuity Designs

Analyze effects at natural discontinuity, as if random assignment



- the size of discontinuity is causal LATE
- Parameter in the subpopulation near threshold. treatment threshold ↑
- To estimate use OLS and  $E[Z \geq c]$  as variable read off by sharp threshold ↑
- Nonparametrically use KNN sample averages coefficient for nearest k-neighbors
- Can extrapolate near cutoff by adding covariates to help reduce errors in local extrapolation

What if treatment assignment not sharp at c?

Fuzzy RDD = RDD + IV ← nudge towards/against treatment where  $E[Z \geq c]$  at cutoff

Step 1: Run RDD to find effect of cutoff on treatment

$$\tau_1 = E_X \left( I[\text{class size} \geq 30] \sim I[\text{enrollment} > 40] \right) \text{ RDD } E[Z \geq c] \text{ cutoff}$$

Step 2: Run RDD to find effect on outcome  $E[Z \geq c]$  cutoff

$$\tau_2 = E_X \left[ \text{Mathscore} \mid \text{enrollment} = 40 \right] \text{ IN } \textcircled{a}$$

Step 3: Get effect of treatment on mathscore

$$\tau = \frac{\tau_2}{\tau_1} \text{ super local avg treatment effect of class size on enrollment}$$

## Causal Inference Online Decision Making

"Maximize reward of A/B test balancing exploration vs exploitation"

optimal arm →  $t^* = \text{argmax } E[Y(t)] \leftarrow \text{reward of arm}$

$$\text{Regret}(N) = \sum_{n=1}^N \mu(t^*) - \mu(\frac{n}{N})$$

Regret(N) should be sublinear in N,  $\text{Regret}(N) \rightarrow 0$

Multi-Armed Bandit Strategies

Goal:  $E \text{Regret}(N) \leq c' \log(N)$

Greedy: See all arms, pick highest reward for N epochs

$\epsilon$ -Greedy: See all arms, w.p.  $1-\epsilon$  be greedy,  $\epsilon$  pull random

\*  $\text{Regret}(N) \geq cN$ ,  $\text{Regret}(N) \geq \delta \epsilon N^{1/(1-\epsilon)}$

K- $\epsilon$ -Greedy: Experiment for K, then be greedy.

$\text{Regret}(N) \geq \delta K(m-1) + \delta p(N-K)$

Fix, N, K greedy: tune K to get regret N

$E[\text{Regret}(N)] \leq \Delta \frac{2m}{\epsilon} \log(Nm) + \Delta$

UCB: Create CI for each arm, pick arm that can potentially do best, update CI to not select in future.

Pull arm t w/ max  $\hat{\mu}_n(t) + \alpha \sqrt{2 \log(n)} \leq \hat{\mu}_n(t)$

For binary  $\hat{\mu}_n(t) \neq \alpha \sqrt{\frac{\log(n)}{2n(t)}} \leftarrow \text{epochs}$

$E(\text{Regret}) < c \log(N)$

Thompson Sampling: draw reward from posterior and pretend like true reward

$E(\text{Regret}) \leq c \log(N)$

\* posterior auto adjust best guess for reward

Contextual Bandit: observe context  $X_n$  before  $T_n$

$t^*(X) = \text{argmax } \mu(t, X) \quad t=1 \dots m$

$\mu(t, X) = E[Y(t) | X=X]$

reward

X needs to be iid,  $\mu(t, X)$  is linear