

Date: Tuesday, February 5, 2019

# Bayesian A/B Testing

## \* Classic A/B Testing (aka Bucket testing)

- Setup:
- Two experimental treatments A & B
  - Collection of randomly sampled individuals
  - $n_A$  are exposed to treatment A
  - $n_B$  are exposed to treatment B

Original example: Drug efficacy.

(Inspired by clinical trials)

Treatment A is say a blood pressure medicine and treatment B is a placebo.

**QUESTION:** Is the blood pressure medicine working?

↳ Measurement: Blood Pressure of every individual  $x_j$  on day 31.

- Let  $\mu_A$  = True mean blood pressure of group A  
 $\mu_B$  = " " " B

Test the efficacy of the drug by investigating the hypothesis test:

$$H_0: \mu_A = \mu_B \quad (\text{It didn't work})$$

↳ Recast just a bit:  $\delta = \mu_B - \mu_A$

↳ If the drug is "useful", then we expect  $\delta > 0$ .

## \* Frequentist approach:

$$\text{Take } \bar{x}_A = \frac{1}{n_A} \sum_{j \in A} x_j$$

$$\bar{x}_B = \frac{1}{n_B} \sum_{j \in B} x_j$$

Run a t-test :  $T = \frac{\bar{x}_B - \bar{x}_A}{\sqrt{\frac{1}{n_A} \sum_{j \in A} (x_j - \bar{x}_A)^2 + \frac{1}{n_B} \sum_{j \in B} (x_j - \bar{x}_B)^2}}$

↳ Pooled standard deviation of the groups

In the frequentist perspective  $\mu_A$  and  $\mu_B$  are fixed unknown constants, so we estimate them and look at the sampling distributions of the estimates.

## \* Bayesian approach

$\mu_A$  and  $\mu_B$  are events that have a corresponding probability distribution (i.e.  $\mu_A$  and  $\mu_B$  are random!)

Prior beliefs on  $\mu_A$  and  $\mu_B$  are specified by

$$P(\mu_A); P(\mu_B)$$

Data:  $x_j, j \in A, B$  w/distribution  $f(x_j | \mu_A, \mu_B)$

Posterior:  $P(\mu_A | x_j, j \in A); P(\mu_B | x_j, j \in B)$

To answer our question of efficacy, we need to investigate the posterior distribution of

$$\delta = \mu_B - \mu_A$$

↳ key idea: If the distribution of  $\delta$  is stochastically greater than 0, drug A is useful.

↳ every possible value is always positive (best case scenario) Distribution is shifted to the right of 0.

- \* In the frequentist approach, we have a p-value to characterize significant difference.
- \* Bayesians though have probabilities on  $\delta$  (much stronger)!

In our example, we'd like to know  $P(\delta > 0 | X, \mu_A, \mu_B)$

NOTE:  $\delta$  is called the **average treatment effect (ATE)**

**Example:** User experience & click-thru rate  $\rightarrow$  conversion rate  
 actually buy a product

Version A vs Version B (think ~~Netflix~~ Amazon)

Version A: Shown to  $N_A$  people  
 $n_A$  people purchase a product  
 $P_A = \frac{n_A}{N_A} = \text{conversion rate}$

Version B: Shown to  $N_B$  people  
 $n_B$  people purchase a product  
 $P_B = \frac{n_B}{N_B} = \text{conversion rate}$

**Experimental setup:** A unique visitor comes to Amazon. With probability  $1/2$ , this visitor is shown version A and w/ prob  $1/2$  shown version B. Amazon keeps running this procedure until a desired number of visits have occurred.

QUESTION: Which version of the website led to a higher conversion rate?

→ To answer this, we look at posterior distributions for  $P_A$ ,  $P_B$ , and  $\delta = P_A - P_B$

Data: We observe conversion rates for each of our versions. Call them  $\hat{P}_A$  and  $\hat{P}_B$ .

We also know  $N_A = 1500$ ,  $N_B = 750$

$$\hat{P}_A = 0.05; \hat{P}_B = 0.04$$

Model:  $n_A \sim \text{Bin}(N_A=1500, P_A)$

$n_B \sim \text{Bin}(N_B=750, P_B)$

$P_A, P_B$ : "True" conversion rates.

prior:  $P_A \sim U(0,1)$

$P_B \sim U(0,1)$

No other info except we know they [lie] between 0, 1  
[1...]

Posteriors:

$P_A | \hat{P}_A, N_A$

$P_B | \hat{P}_B, N_B$

} Get posterior of  $\delta$  and answer question!

To answer our question, we calculate:

$$P(\delta \geq 0 | \hat{P}_A, \hat{P}_B, P_A, P_B, N_A, N_B) = 0.983$$

↳ weight of the histogram to the right of 0.

$$P(\delta < 0 | \dots) = 0.017$$

Good evidence that A is better for conversion than B.

# A little more about posterior inference

There are many ways to use the Posterior distribution once you have it:

- 1) Probabilities of the parameter of interest.

Ex: For A/B test,  $P(\delta > 0 | X)$

- 2) Credible intervals: the Bayesian analog to confidence intervals.

A  $(1-\alpha)100\%$ . credible interval  $[a, b]$  st  $P(\theta \in [a, b] | X) = 1-\alpha$

- 3) Maximum a posteriori (MAP) estimate:

The value of  $\theta$  that has the highest posterior probability (i.e. the mode of  $\theta | X$ ).

## SOME COMMENTS ABOUT GOODNESS OF FIT

- \* Recall that our overall aim is to develop a model for how the data we observe was generated!
- \* This is the same as any other statistical or ML model, so the same metrics can be used to evaluate our model/performance.
- \* As an example, we can split data into training vs test. Build Bayesian model on training, and simulate possible values for the test set. Finally, compare simulated values with observed values using your favorite metric.
- \* Similarly, to evaluate performance within the training sample, we can simulate in-training data and compare w/ what we observed using similarity metrics (ex: MSE, Kolmogorov, Smirnov, etc.)

## About the test:

Fill in the blank, multiple choice, Short answer

- Types of models to use — distribution  
how to use & parameters.
- Basic probability questions (result of calculation)
- Bayesian modeling
- No coding!