# A/B Testing

**Md Emdadul Sadik**

**Md Enamul Huq Sarker**

**Summer 2020**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Overview
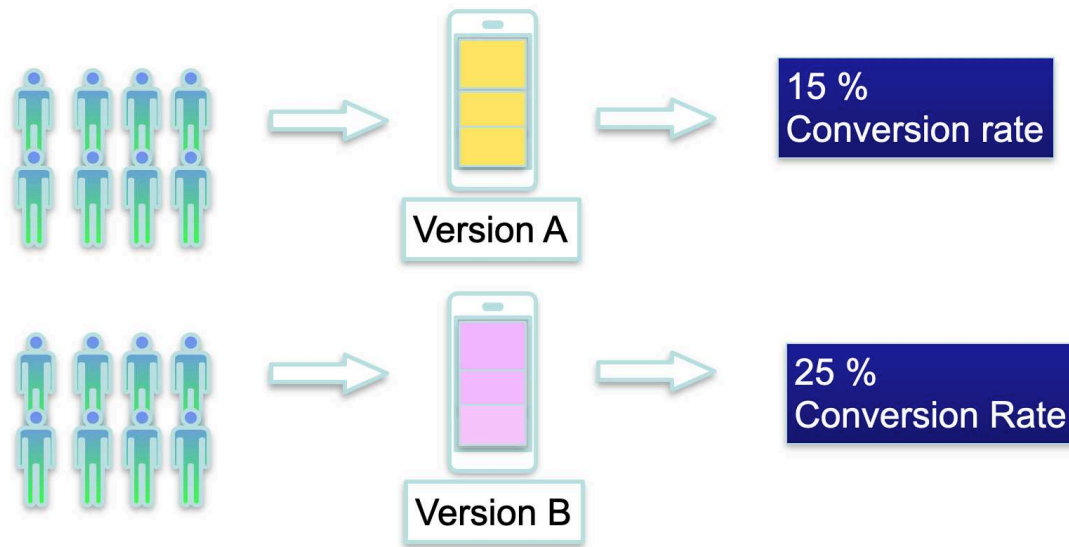
1. A/B testing

   - What is it?

   - Why is that used?

   - When (or not) to use A/B test?

   - Hypothesis testing & p-value

   - Type I & Type II error

2. Multivariate testing
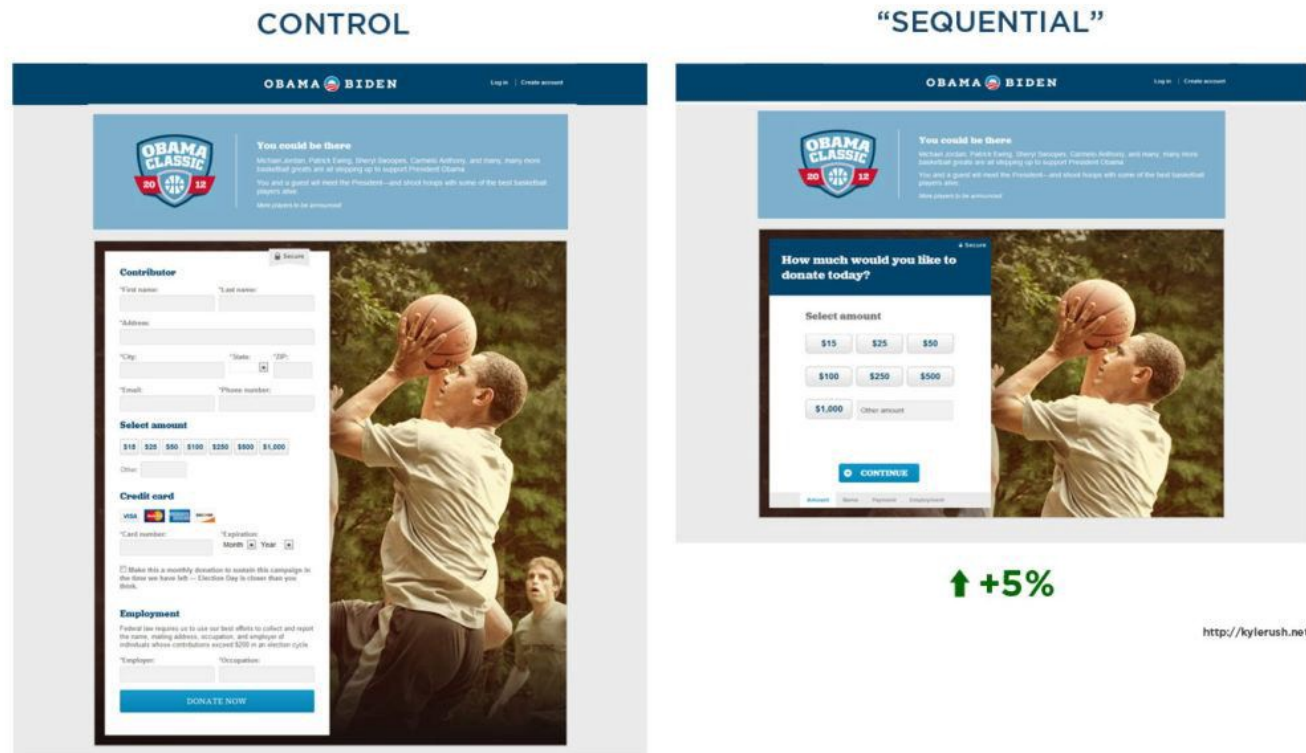
3. A/B Testing of ML Models

# What is A/B Testing?

- A user experience research methodology.

- Compares two versions of design alternatives (i.e two versions of a single variable)

# Obama campaign 2012

- A/B testing in Obama's 2012 presidential campaign
- 165 team digital team
- 500+ experiments
- Over 20 months
- $190 million extra



Image source

# Should I use A/B Test?

- All the big companies use A/B testing. But why?

  - Intuition can be often wrong! Reading user mind is complex.

  - Higher risk to roll out a features to all users.


- Think, if you should use A/B testing in below cases?

  - Changing colour or theme of a website

  - Changing company logo

  - Car sellers website

  - Movie preview

# When A/B test shouldn't be used?

- You shouldn't go for A/B test if
  - You don't have meaningful traffic
    - Statistically significant sample size is important.
  - You can't spend the mental bandwidth.
  - You don't have a solid hypothesis to start with.
    - Ex: Adding a 'Finish purchase' button will increase purchase by 20 percent.
  - Risk is too low to immediate action.
    - Implementation is preferable instead of wasting time on A/B testing

# Common terms

- What is a hypothesis?
  - Claim or idea to be tested
- Control group
  - Doesn't get special treatment.
- Experimental group
  - Gets special treatment.
- Null hypothesis ($H_0$)
  - Outcome from control and treatment are identical.
- Alternate Hypotheis ($H_a$)
  - Outcome from treatment is different.

Does adding fertiliser 'X' increases plant growth?

Control group

Experimental group

# Hypothesis Testing

- Average session time is 20 minutes

- Change website background colour from <span style="color:blue">Blue</span> to <span style="color:orange">Orange</span>

- How to do the hypothesis testing?

  1. Null hypothesis ($H_0$) : mean = 20 minutes <span style="color:orange">after the change</span>

  2. Alternate hypothesis ($H_a$) : mean > 20 minutes <span style="color:orange">after the change</span>

  3. Significance level (p-value threshold): $\alpha$ = .05

  4. Take sample, for example, n = 100, sample mean $\bar{X}$ = 25 minutes.

  5. p-value: $P(\bar{X} \geq 25$ minutes | $H_0$ is true)

     - If p-value < $\alpha$ then reject $H_0$, suggest $H_a$

     - If p-value >= $\alpha$ then don't reject $H_0$, (doesn't mean accept $H_0$)

# Hypothesis Testing (cont.)

- If p-value < α then reject $H_0$, suggest $H_a$

- If p-value >= α then don't reject $H_0$, (doesn't mean accept $H_0$)

- Example:

    - p_value is 0.03, reject $H_0$, suggest $H_a$

    - p_value is 0.05, Fail to reject reject $H_0$

- Why should you set significance value prior to the experiment?

    - Ethical reason

# How to calculate P-value

- P-Value means probability value which indicates how likely a result occurred by chance alone
- P-value is calculated as probability of the random chance that generated the data or (+) something else that is equal (probability) or (+) something rarer (less probability)

The p-value for getting **HHHHH**

A p-value is the probability that random chance generated the data, or something else that is equal or rarer.

Since there is nothing rarer, this is what we get:

$$\text{Pr(5 heads)} + \text{Pr(5 tails)}$$

$$\frac{1}{32} + \frac{1}{32} = 0.0625$$

Outcomes

HHHHH

THHHH
HTHHH
HHTHH
HHHTH
HHHHT

TTHHH   TTTHH
THTHH   TTHTH
THHTH   THTTH
THHHT   HTTTH
HTTHH   TTHHT
HTHTH   THTHT
HTHHT   HTTHT
HHTTH   THHTT
HHTHT   HTHTT
HHHTT   HHTTT

TTTTH
TTटHT
TTHTT
THTTT
HTTTT

TTTTT

# Type I and Type II error

| | Fail to Reject | Reject |
|---|---|---|
| $H_0$ is true | Correct conclusion | Type I error |
| $H_0$ is false | Type II error | Correct conclusion |

- How to reduce Type I error?
  - Lower the value α
  - Reducing value of α, increases type II error
- How to to reduce Type II error?
  - Increased sample size
  - Less variability
  - True parameter far from from $H_0$

# Multi variate & A/A testing

- **Multivariate testing :** Multiple variables are modified, also called full factorial testing.
  - Advantage: A lot of combinations can be tested
  - Limitation: Bigger sample size, complex, needs better understanding of interactions



- **A/A Testing:**
  - Identical version is compared against each other.
  - Used to validate the tool(s) being used.

# Factorial testing with PlanOut

- Factorial test is complex to realise and implement.
- Planout (https://facebook.github.io/planout/) a framework for online field experiment



**Figure 1: A factorial experiment in PlanOut. (a) PlanOut language script (b) a graphical interface for specifying simple PlanOut experiments (c) a JSON representation of serialized PlanOut code (d) an illustration of the proportion of `cookieids` allocated to each parameterization. Note that because we use `weightedChoice()` to assign `button_text`, more cookies are assigned to "Sign up" than ''Join Now''.**
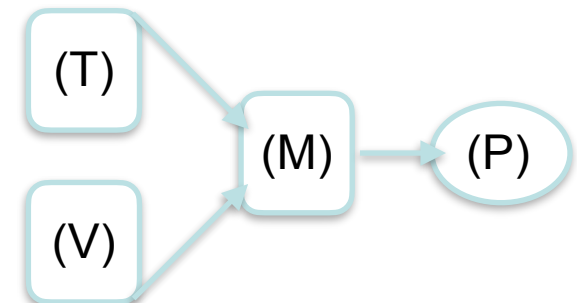
# Machine learning with A/B Testing

- Only relying on outcome from A/B testing sometimes doesn't lead to best decision.

- Applying machine learning, better insight on user behaviour
  - Possible to achieve alternate suggestion. I.e In order to achieve A, instead of adding a button 'X' focus on Y.



Image source

# A/B Testing of ML Models

- Model (M)
  - A model is artefact(s) created (trained) by AI creation algorithm(s). Example: MS [ONNX](ONNX) file.
- Model Predictions (Brings Output)
  - Predictions, (P) are the output of a model, (M) trained using AI algorithm(s).
- Model Deployment (Brings Outcome)
  - Means that model predictions are being consumed by an application that is directly affecting business operations.
- Predictive models are trained on historical data set (experiences), (T)
- Models are tested on holdout/validation data set (V). Presumably best performant model is deployed.
- Finding the best model post-deployment is the purpose.

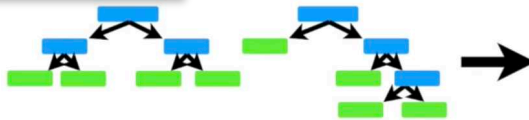(T) → (M) → (P)
(V) → (M)

# The Two Variants

Imagine, we have some clinical data that helps deciding whether a patient has heart disease or not.

| Chest Pain | Good Blood Circ. | Blocked Arteries | Weight | Heart Disease |
|---|---|---|---|---|
| No | No | No | 125 | No |
| Yes | Yes | Yes | 180 | Yes |
| Yes | Yes | No | 210 | No |
| ... | ... | ... | ... | ... |

# The Two Variants

We deploy Random Forest (**model A**) and K-Nearest (**model B**) and to find out.
TP looks good for **model A.**

Model A - RF

| | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 142 | 22 |
| Does Not Have Heart Disease | 29 | 110 |

| | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 107 | 53 |
| Does Not Have Heart Disease | 64 | 79 |

Model B - KNN

**K-Nearest Neighbors** was worse than the **Random Forest** at predicting patients *with* Heart Disease (**107** vs **142**)…

We deploy Random Forest (**model A**) and K-Nearest (**model B**) to find out. TN also looks good for **model A.**



…and worse at predicting patients
*without* Heart Disease (**79** vs **110**)…

We deploy Random Forest (**model A**) and K-Nearest (**model B**) to find out.
**Model A wins!**



| | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| **Has Heart Disease** | 142 | 22 |
| **Does Not Have Heart Disease** | 29 | 110 |

| | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 107 | 53 |
| Does Not Have Heart Disease | 64 | 79 |

…so if we had to choose between using the **Random Forest** and **K-Nearest Neighbors**, we would choose the **Random Forest**.

- **Hypothesis Test** ( between models A, B to find a winner)
- **model A** (control) is deployed and predicting sth. i.e **Null Hypothesis $H_0$**
- **model B** (test), challenging **model A**, predicts sth. even better i.e **Alternative $H_a$**

$$Sensitivity = TPR = \frac{TP}{TP + FN} = \frac{TP}{Actual\ Positives} \qquad Specificity = TNR = \frac{TN}{TN + FP} = \frac{TP}{Actual\ Negatives}$$

$$Confidence Level, CL = The\ probability\ of\ correctly\ retaining\ the\ H_0\ ;\ \ 95\%$$

$$Statistical\ significance\ \alpha = 1 - CL$$

$$Accuracy = \frac{Total\ Correct\ Predictions}{Total\ Data\ Set}$$

**Effect Size:** The difference between the two models' performance metrices.

# MQ: Sensitivity & Specificity



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Again we have a confusion matrix from that clinical data we saw. This time we apply LR (A) and RF (B) to measure models' performance w/ Sensitivity and Specificity.

|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 139 | 20 |
| Does Not Have Heart Disease | 32 | 112 |

Model A - Logistic Regression

Src: StatQuest

|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 142 | 22 |
| Does Not Have Heart Disease | 29 | 110 |

Model B - Random Forest

Src: StatQuest

# MQ: Sensitivity & Specificity

$$Sensitivity = TPR = \frac{TP}{TP \; + \; FN} \; = \; \frac{TP}{Actual \; Positives}$$

$$Sensitivity(LR) = \frac{139}{139 \; + \; 32} \; = \; 0.81$$

$$Sensitivity(RF) = \frac{142}{142 \; + \; 29} \; = \; 0.83$$



|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 139 | 20 |
| Does Not Have Heart Disease | 32 | 112 |

|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 142 | 22 |
| Does Not Have Heart Disease | 29 | 110 |

Logistic Regression

Sensitivity = 0.81
Specificity = 0.85

Random Forest

Sensitivity = 0.83
Specificity = 0.83

**Sensitivity** tells us that the **Random Forest** is slightly better at correctly identifying *positives*, which, in this case, are patients *with* heart disease.

# MQ: Sensitivity & Specificity

$$Specificity = TNR = \frac{TN}{TN + FP} = \frac{TP}{Actual\ Negatives}$$

$$Specificity(LR) = \frac{112}{112 + 20} = 0.85$$

$$Specificity(RF) = \frac{110}{110 + 22} = 0.83$$



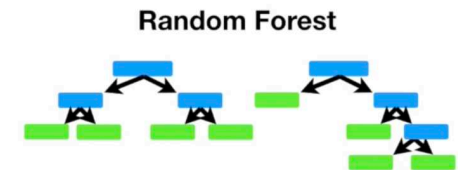|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 139 | 20 |
| Does Not Have Heart Disease | 32 | 112 |

|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 142 | 22 |
| Does Not Have Heart Disease | 29 | 110 |

Logistic Regression

Random Forest

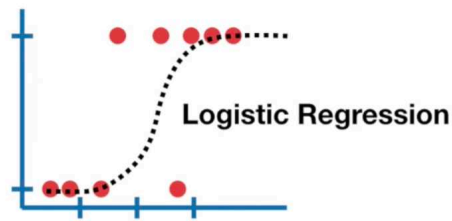Sensitivity = 0.81

Sensitivity = 0.83

Specificity = 0.85

Specificity = 0.83

**Specificity** tells us that **Logistic Regression** is slightly better at correctly identifying *negatives*, which, in this case, are patients *without* heart disease.

# MQ: Sensitivity & Specificity

|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 139 | 20 |
| Does Not Have Heart Disease | 32 | 112 |

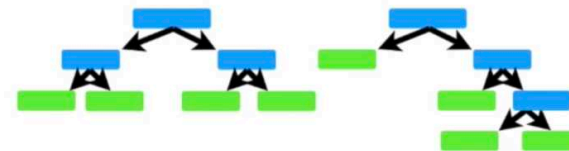|  | Has Heart Disease | Does Not Have Heart Disease |
|---|---|---|
| Has Heart Disease | 142 | 22 |
| Does Not Have Heart Disease | 29 | 110 |

**Random Forest**

Logistic Regression

**Sensitivity** = 0.81

**Specificity** = 0.85

We would choose the **Logistic Regression** model if correctly identifying patients **without** heart disease was more important than correctly identifying patients **with** heart disease.

**Sensitivity** = 0.83

**Specificity** = 0.83

Alternatively, we would choose the **Random Forest** model if correctly identifying patients **with** heart disease was more important than correctly identifying patients **without** heart disease.

# MQ: Accuracy

## Accuracy Comparison

$$Accuracy = \frac{Total\ Correct\ Predictions}{Total\ Data\ Set}$$

Model 1

| predictions | | | | |
|---|---|---|---|---|
| actual values | A | B | C | D |
| A | 10 | 0 | 0 | 0 |
| B | 0 | 5 | 3 | 2 |
| C | 0 | 1 | 8 | 1 |
| D | 0 | 1 | 0 | 9 |

(10 + 5 + 8 + 9) / 40 = **0.8**

Model 2

| predictions | | | | |
|---|---|---|---|---|
| actual values | A | B | C | D |
| A | 8 | 2 | 0 | 0 |
| B | 1 | 7 | 0 | 2 |
| C | 0 | 0 | 9 | 1 |
| D | 2 | 3 | 0 | 5 |

(8 + 7 + 9 + 5) / 40 = **0.725**

Image Courtesy : Minsuk Heo

For balanced data accuracy could alone answer for the best model.
But the reality is not always ideal!
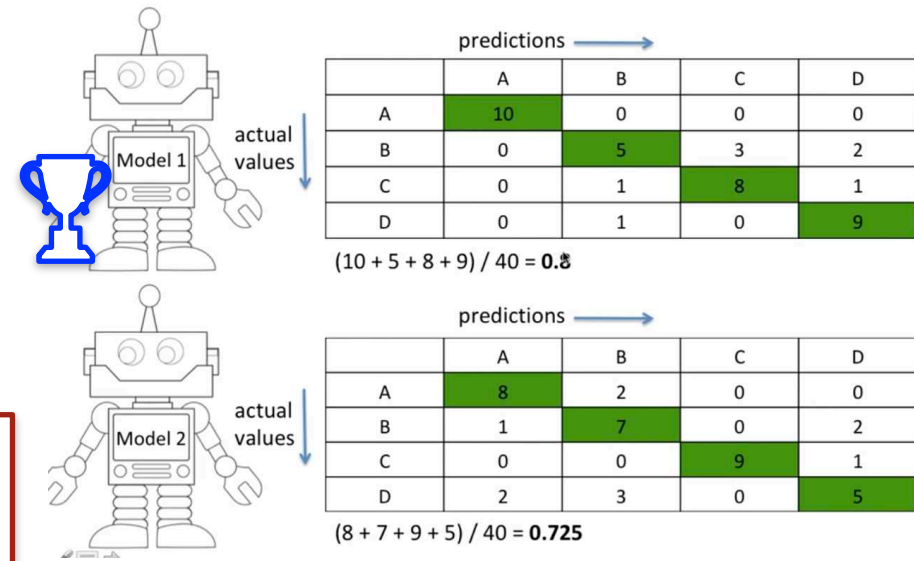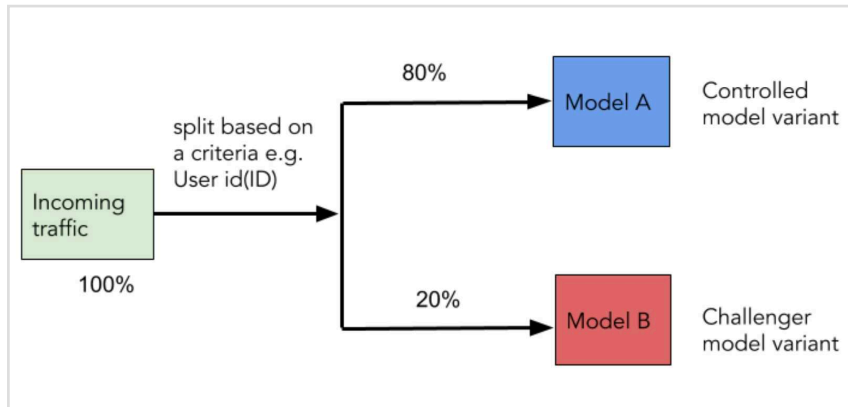
The picture shows two models deployed to classify multiple classes (A-D). By comparing the accuracies one could decide that Model 1 wins.

# Model Deployment - MD



Orcale White Paper on Model Testing

```python
# IDs of users in treatment group
TREAMENT_IDS = {}

@app.route("/predict")
def predict():
    features = request.get_json['features']
    if request.get_json['user_id'] in TREAMENT_IDS:
        return model_B.predict(features)
    else:
        return model_A.predict(features)
```
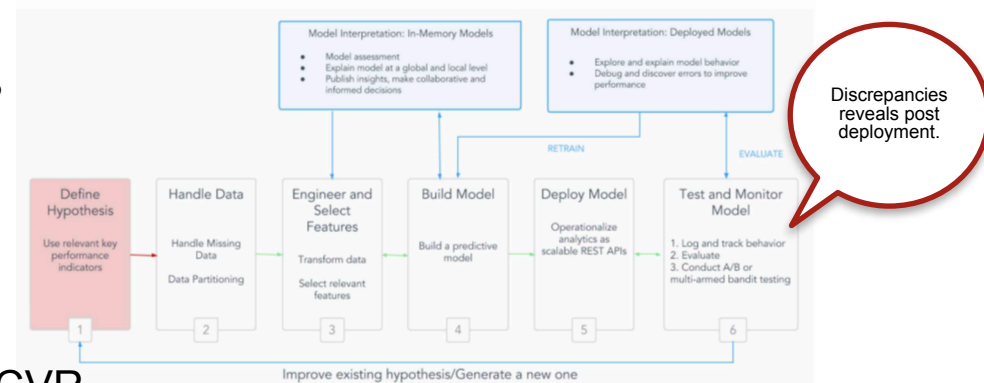
mlinproduction.com

The picture shows an A/B testing of two models. If we add more models C,D.. N in the same way the test would become a A/B/n or multivariate test.

A Trivial model deployment example using Python Flask http endpoint.

# MD - Post Deployment Discrepancies

- Predictors (features) changing
  - e.g. a CTR model sees a new acquisition channel.

- Performance Metrics may differ
  - e.g. Training set was measured against
  - Balanced data —> AUC, Accuracy
  - Imbalanced data —> F1-score
  - With which do we measure the winner?

- Experiments of models may hurt UX
  - shouldn't be the case in anyway.

- Deployed to measure a business KPI
  - e.g customer churn rate or to increase CVR.
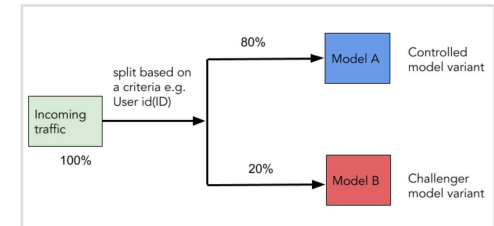  - But now it measures performance with AUC

Orcale White Paper on Model Testing

## Designing a Model A/B Test

At a high level, designing an A/B test for models involves the following steps

- Deciding on a performance metric. It could be the same as the one used during the model training phase (e.g., F1, AUC, RMSE, etc.)

- Deciding on test type based on your performance metric.

- Choosing a minimum effect size you want to detect.

- Determining the sample size N, based on your choice of selected minimum effect size, significance level, power, and computed/estimated sample variance.

- Running the test until N test units are collected.

**Effect Size:** The difference between the two models' performance metrices.

# MD - Mistake of Early Declaration

## It's a mistake, don't you pull the plug!

Declaring a model a resounding success before collecting N units of sample. The early significance could also be achieved by random chance!

| Test | Samples Collected | | | | |
|------|------|------|------|------|------|
| | 200 | 400 | 600 | 800 | 1000 |
| Test 1 | Not Sig | Sig - STOP | Not Sig | Not Sig | Not Sig |
| Test 2 | Sig - STOP | Not Sig | Not Sig | Not Sig | Not Sig |
| Test 3 | Not Sig | Not Sig | Not Sig | Not Sig | Not Sig |
| Test 4 | Not Sig | Not Sig | Not Sig | Not Sig | **Sig** |

- If we pick a significance level $\alpha=0.05$, we'd expect to see significant results (Sig) in one of 20 independent tests for a fixed and identical N (N=1000).
- If we stop as soon as significance is reached, we preferentially select spurious false positives.

# MD - Holy Grails of Model A/B Testing

- Perform an A/A test. At α=0.05, a significant result should be seen 5% of the time.
- Do not turn off the test as soon as you detect an effect. Stick to your pre-calculated sample size. Often, there is a novelty effect in first few days of model deployment and a higher risk of false positives.
- Use a two-tailed test instead of a one-tailed test.  (look both for $H_0$ and $H_a$)
- Control for multiple comparisons. ( use Bonferroni correction, stringent α to avoid Type I / FPs. )
- Beware of cross-pollination of users between experiments.  ( same user does not get both a/b)
- Make sure users are identically distributed. Any segmentation (traffic source, country, etc.) should be done before randomisation is needed.
- Run tests long enough to capture variability such as day of the week seasonality.
- If possible, run the test again. See if the results still hold. B
- Beware of Simpson's paradox. (Changing experiment during intervention settings skews result. Rollout new model instead.)
- Report confidence intervals; they are different for percent change or non-linear combinations of metrics.

# References

A / B Testing: The Most Powerful Way to Turn Clicks Into Customers. John Wiley & Sons. ISBN 978-1-118-65920-5.

Designing and Deploying Online Field Experiments
By Eytan Bakshy, Dean Eckles, Michael S. Bernstein

When A/B Testing Isn't Worth It

Oracle Whitepaper - Testing Predictive Models in Production
By Ruslana Dalinina Jean-René Gauthier, and Pramit Choudhary

A/B Testing Machine Learning Models (Deployment Series: Guide 08)
ML in production

Khan Academy - Unit: Significance tests (hypothesis testing)

Confusion Matrix _ StatQuest on Youtube
Sensitivity and Specificity - StatQuest on Youtube.
Statistical Significance in A/B Testing – a Complete Guide

# Licence

This work is licensed under a Creative Commons "AttributionShareAlike 4.0 International" license.