

WHERE TO GO NEXT

John Sandall

Data Science Consultant

@john_sandall

INTRODUCTION

BUSINESS FUNDAMENTALS

BUSINESS FUNDAMENTALS



- ▶ AARRR
- ▶ Unit Economics
- ▶ Retention/Churn
- ▶ Product/Pricing Optimisation
- ▶ Operations Research
- ▶ Demand Prediction
- ▶ Customer Segmentation
- ▶ ...

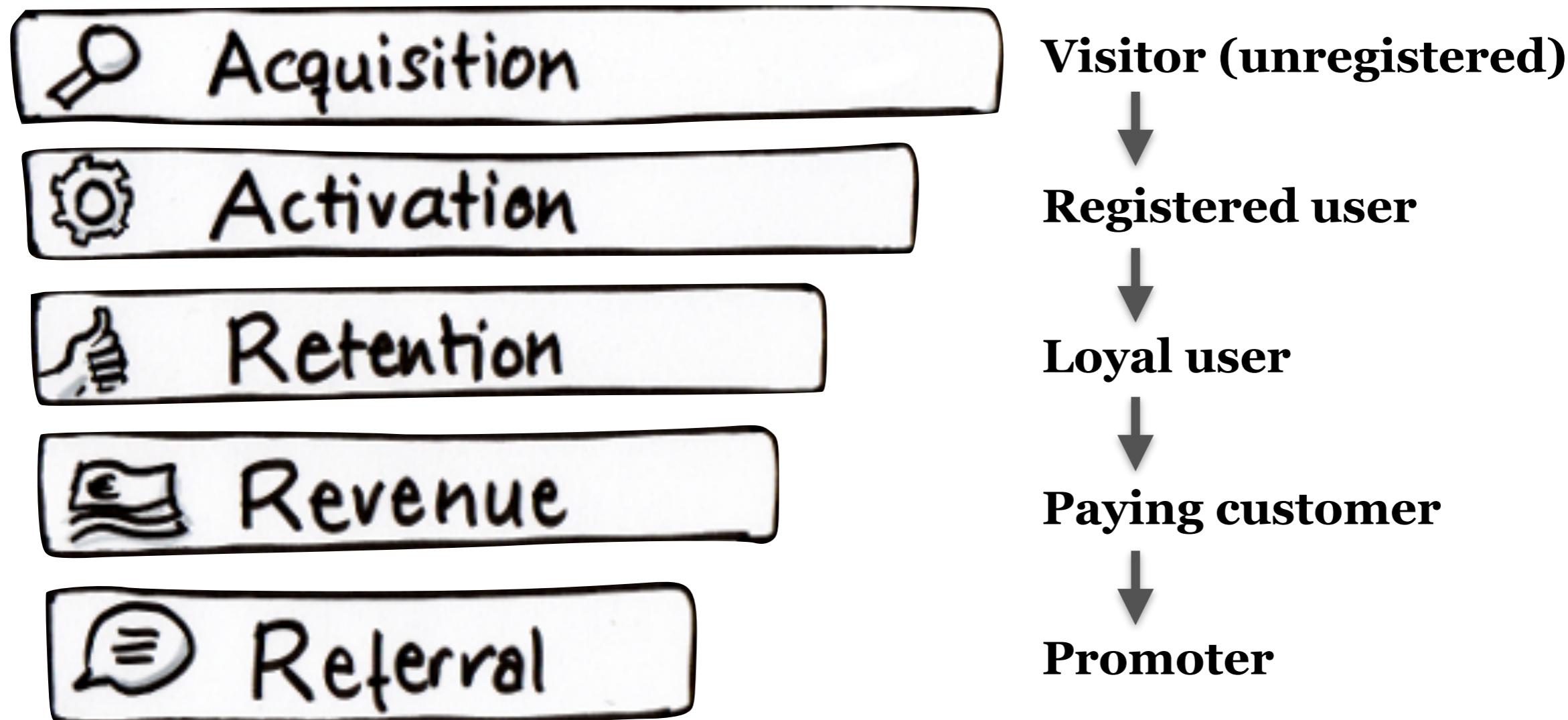
GETTING YOUR METRICS RIGHT

- **Goal:** Use data science to optimise social media strategy.
- **Questions:** What do you mean by "optimise"? What is "success" here?
- In order to optimise anything, we need a metric to optimise.
- My goal here is to provide insight into how a data scientist might approach some common problems in the world of business, marketing and social media.

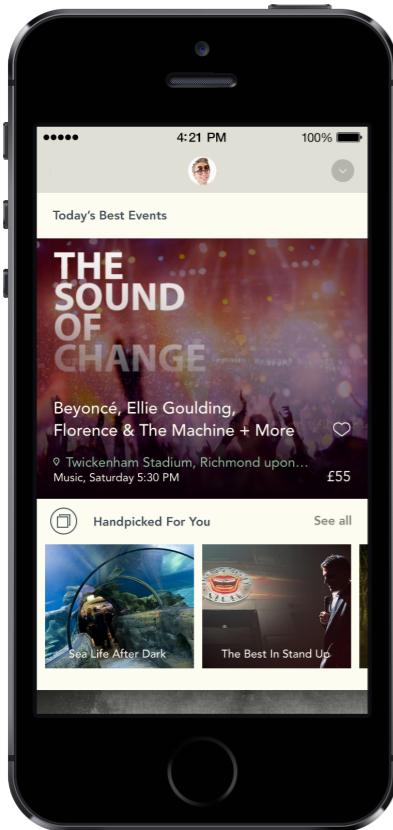
AARRR ("PIRATE METRICS")



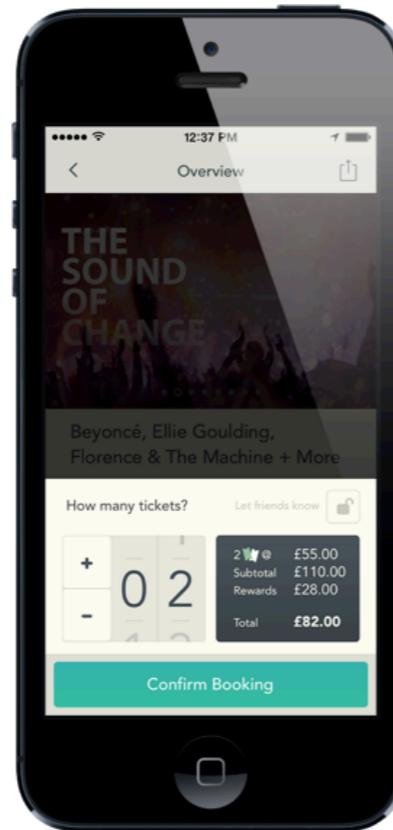
AARRR ("PIRATE METRICS")



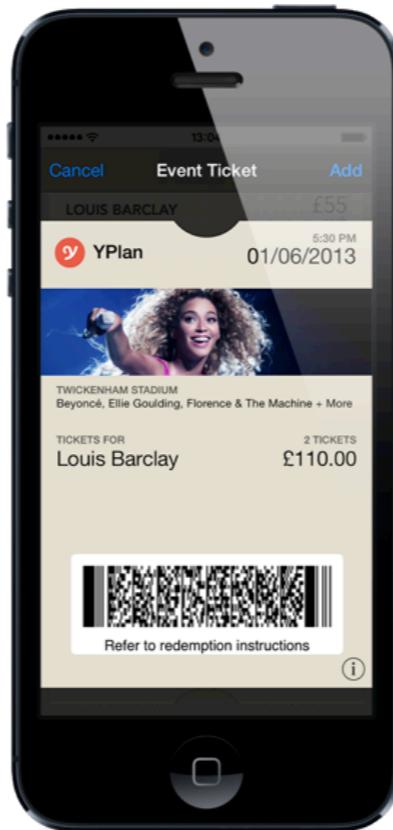
AARRR FOR YPLAN



Spontaneous
inspiration



2-tap booking
process



Paperless
ticketing



Discovery of new experiences



Curated events

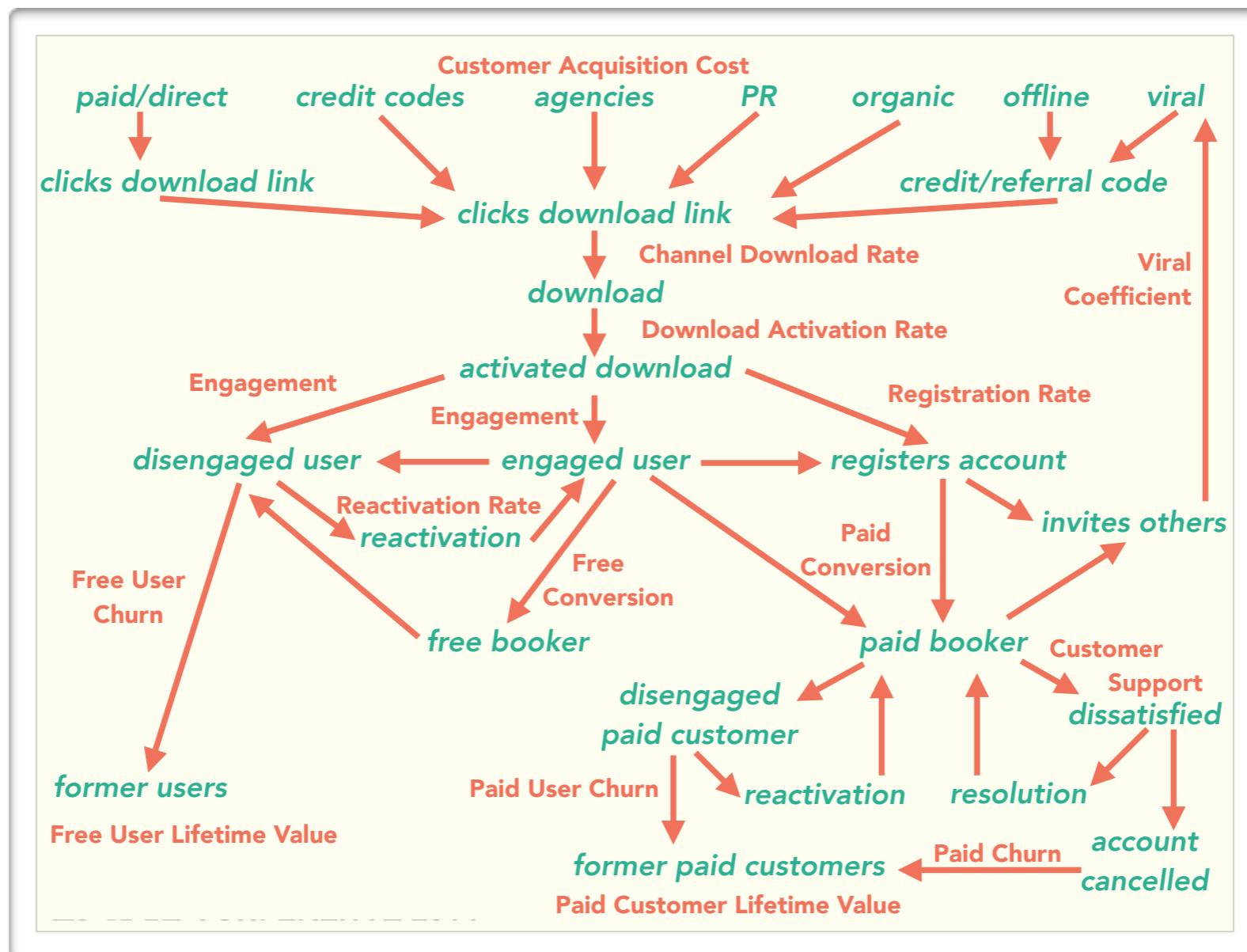


Social media integration



100% mobile

AARRR FOR YPLAN

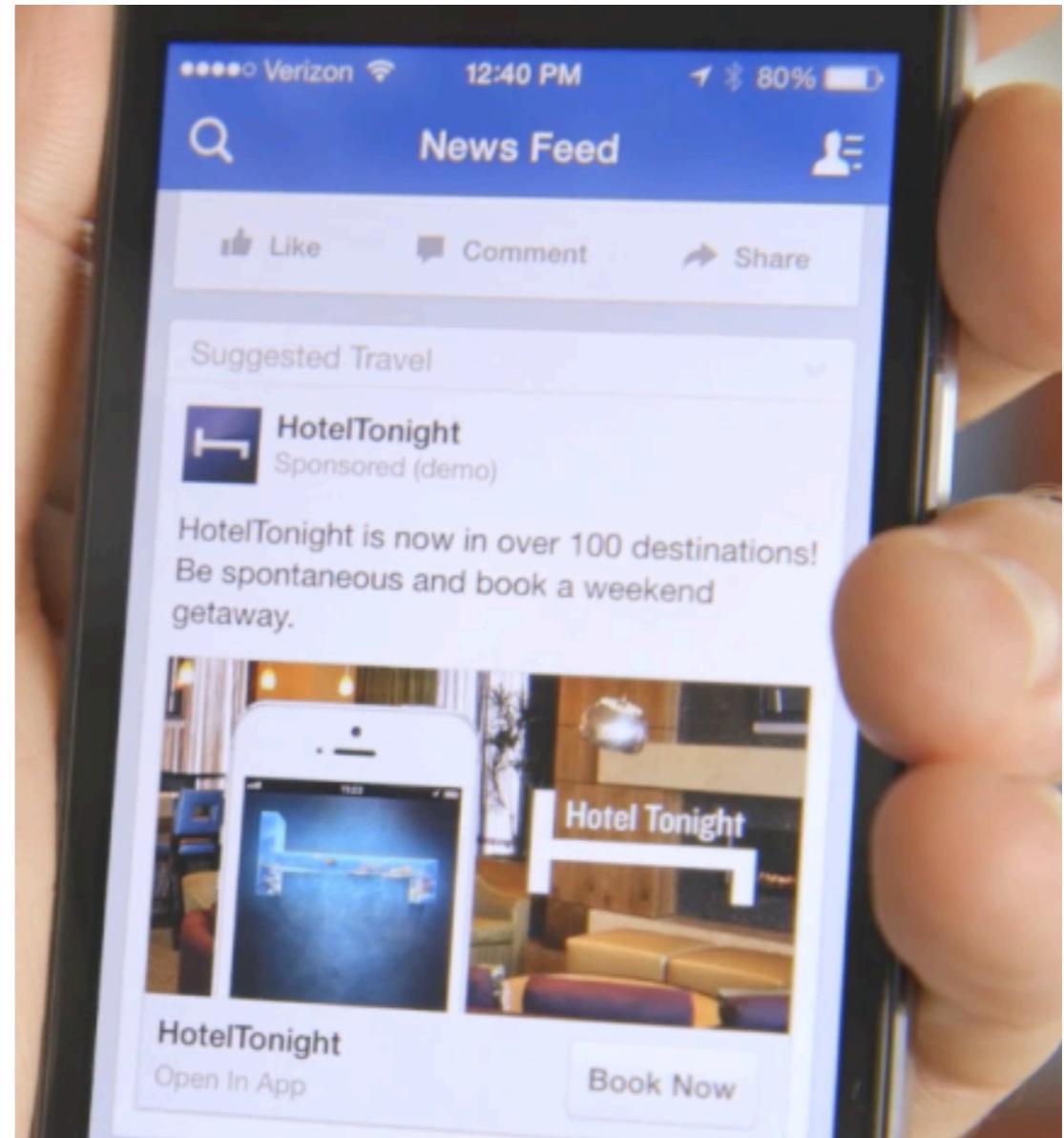


LESSON: BE VERY CLEAR ABOUT WHAT IT IS YOU'RE TRYING TO IMPACT



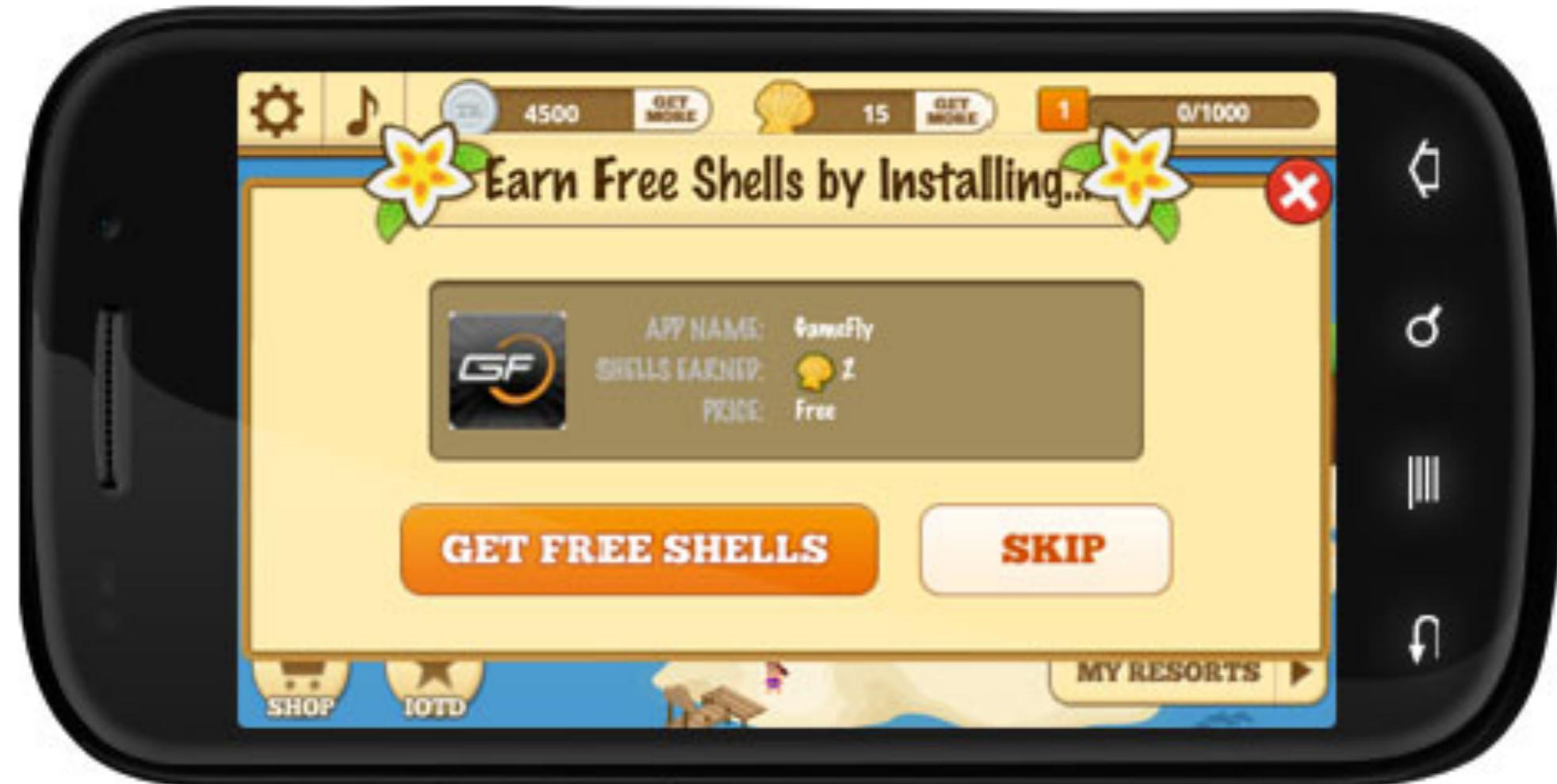
EXAMPLE: ADVERTISING ON GOOGLE/FACEBOOK/AD NETWORKS

- ▶ **Goals:** Awareness? Acquisition?
Both?
- ▶ **Track:** Each activity against these goals. Want to prove causality.
- ▶ **Proxy metrics:** Ad impressions are useful...if they follow through!



EXAMPLE: ADVERTISING ON GOOGLE/FACEBOOK/AD NETWORKS

- Proxy metrics going wrong: Incentivised downloads



EXAMPLE: ADVERTISING ON GOOGLE/FACEBOOK/AD NETWORKS

- **Proxy metrics going wrong:** Optimising subject lines for open rate only

What would you like to test?

Subject lines

From names

Delivery date/times

How should we split the campaign?

We'll run your test on a segment of the list. When the winner is determined, we'll send it to the remaining portion of the list.

The diagram illustrates the segmentation of a list. It features a horizontal bar divided into three sections. The leftmost section, colored teal, is labeled 'A' in red text. The middle section, also colored teal, is labeled 'B' in black text. To the right of these two sections is a grey rectangular area labeled 'Remainder segment'. Below the bar, the text 'Test segment: 40%' is positioned under section A, and 'Send the winner to: 60%' is positioned under section B.

Test segment: 40%

Send the winner to: 60%

EXAMPLE: SOCIAL MEDIA ACTIVITY

- **Goals:** Difficult to isolate!
- **Guiding question:** "If we stopped doing this, what would happen?"
- **Some example goals:**
 - Engaging new users (social media as an acquisition tool)
 - Engaging existing users (social media as a retention tool)
 - Persuading acquired users to activate into paying customers
 - Building communities shifts people from "like" to "love" to "promoters/defenders"

EXAMPLE: SOCIAL MEDIA ACTIVITY

- ▶ **Impacts:** Almost every part of AARRR funnel!
- ▶ **Example:** Engaging new users (social media as an acquisition tool)
- ▶ **Measure:**
 - ▶ # of new users from social media sources
 - ▶ % of new users from social media sources
 - ▶ segment by social media platform
 - ▶ segment by post
 - ▶ look at trends over time
 - ▶ look at best/worst posts from last month...what can we learn?

EXAMPLE: SOCIAL MEDIA ACTIVITY

- ▶ **Impacts:** Almost every part of AARRR funnel!
- ▶ **Example:** Measuring impact of social media activities on referrals
- ▶ **Problem:** Effectively word of mouth...referred users look like organics.
- ▶ **Measure:**
 - ▶ Total/proportion of new users from organic sources
 - ▶ Experiment: increase/decrease activities, look for correlations
 - ▶ Regular surveys of new users: "how did you hear about us?"
 - ▶ Directly link customer accounts with social media profiles
 - ▶ e.g. Insightly: "*We'll detect virtually every social media profile related to a contact's email address*"
 - ▶ Accurately measuring brand equity is notoriously difficult!

CASE STUDY

DRIVING DOWN CAC USING PREDICTIVE MODELS

WHAT IS CAC?



Acquisition

Cost Per Install (CPI)



Activation

Signup Conversion Rate %



Retention

Booker Conversion Rate %



Revenue



Referral

WHAT IS CAC?



Cost Per Install (CPI)

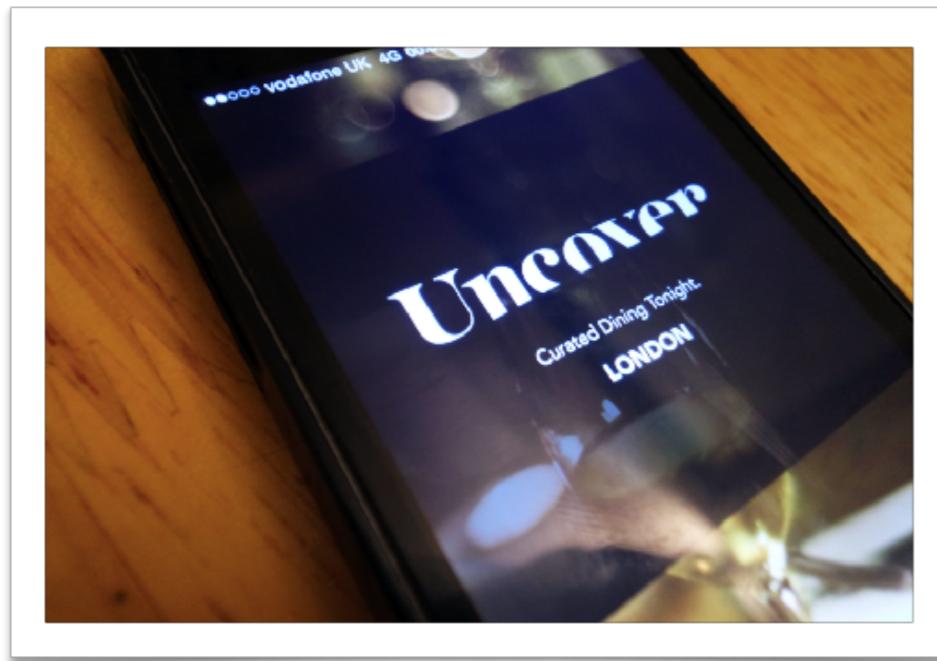
Signup Conversion
Rate %

Booker Conversion
Rate %

Customer
Acquisition
Cost (CAC)

$$\begin{aligned} \text{CAC} &= \text{CPI} \times \text{Signup Conversion Rate} \times \text{Signup-to-Booker Conversion Rate} \\ &= \text{Cost per Acquired Customer} \end{aligned}$$

FIRST CHALLENGE: BRING THE DATA TOGETHER



EXAMPLE: PREDICTING CUSTOMER CONVERSION

- ▶ **Metric:** CAC here means "cost to acquire an install that made a paid booking within 30 days of install"
- ▶ **Problem:** Have to wait 30 days before you can evaluate if a campaign is working or not
- ▶ **Goals:**
 - ▶ Determine leading indicators that are predictive for conversion to paid customer.
 - ▶ Build a classifier model using these indicators to predict conversion rates.
 - ▶ Is it possible to predict within 24h of install which new installs will convert?
 - ▶ Putting it all together...we can accurately estimate CAC within 24h of a new campaign going live.
 - ▶ This means much faster turnarounds, massively speeds up learning cycle of what works.

EXAMPLE: PREDICTING CUSTOMER CONVERSION

► How:

- What data is there?
- Throw this into some machine learning classifier models, e.g. decision tree, random forest
- Make predictions! Look at them 30 days later. How good were they? Learn, iterate, improve.

user_id	source	device	location	restaurants viewed	skipped tutorial	...	booked (within 30 days)
1	facebook	android	london	5	Y	...	
2	twitter	ios	new york	10	N	...	
3	organic	ios	bristol	15	N	...	
...	

EXAMPLE: PREDICTING CUSTOMER CONVERSION

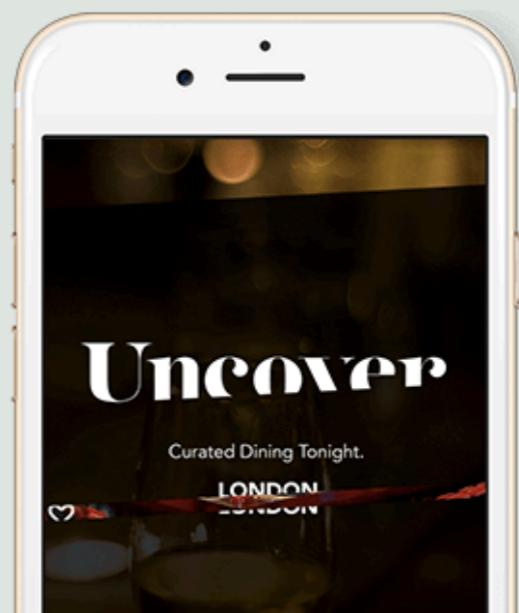
► Impact...

A appsumer®

Product Services Resources About Request a demo

PROVEN RESULTS

Restaurant booking app reduces paid acquisition cost by 74% and increased acquisition volume by 4x



EXAMPLE: PREDICTING CUSTOMER CONVERSION

► Impact...

- ▶ 74% reduction in CAC within 3 months
- ▶ Precise understanding of behavioural indicators that result in paying customers => product improvements
- ▶ High-probability customers who didn't convert are ideal "low hanging fruit" for personalised targeted offers
- ▶ Data warehouse infrastructure can be reused for other data science projects (e.g. churn prediction, recommendation systems, product optimisation)

INTRODUCTION

CUSTOMER SEGMENTATION

RFM MODELS

Traditional approach to customer segmentation by:

- ▶ **Recency:** How recently did we see this customer?
- ▶ **Frequency:** How often do we see this customer?
- ▶ **Monetisation:** How much does this customer spend per transaction?

ACTIVE

AT RISK

CHURNED

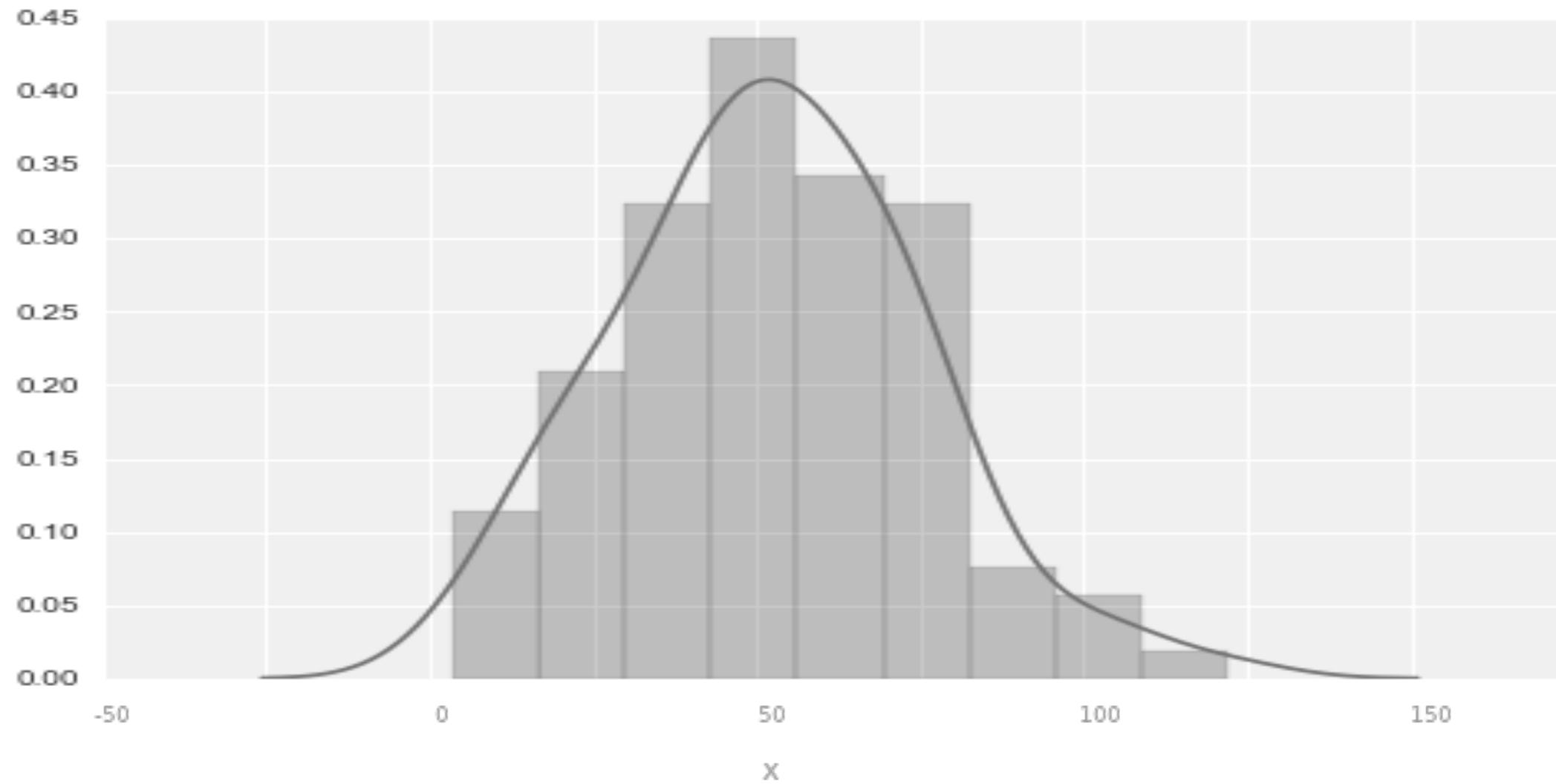
		M			
R	F	1	2	3	4
1	1				
	2				
	3				
	4				
2	1				
	2				
	3				
	4				
3	1				
	2				
	3				
	4				

MORE ADVANCED STILL: LIFECYCLE MODELS

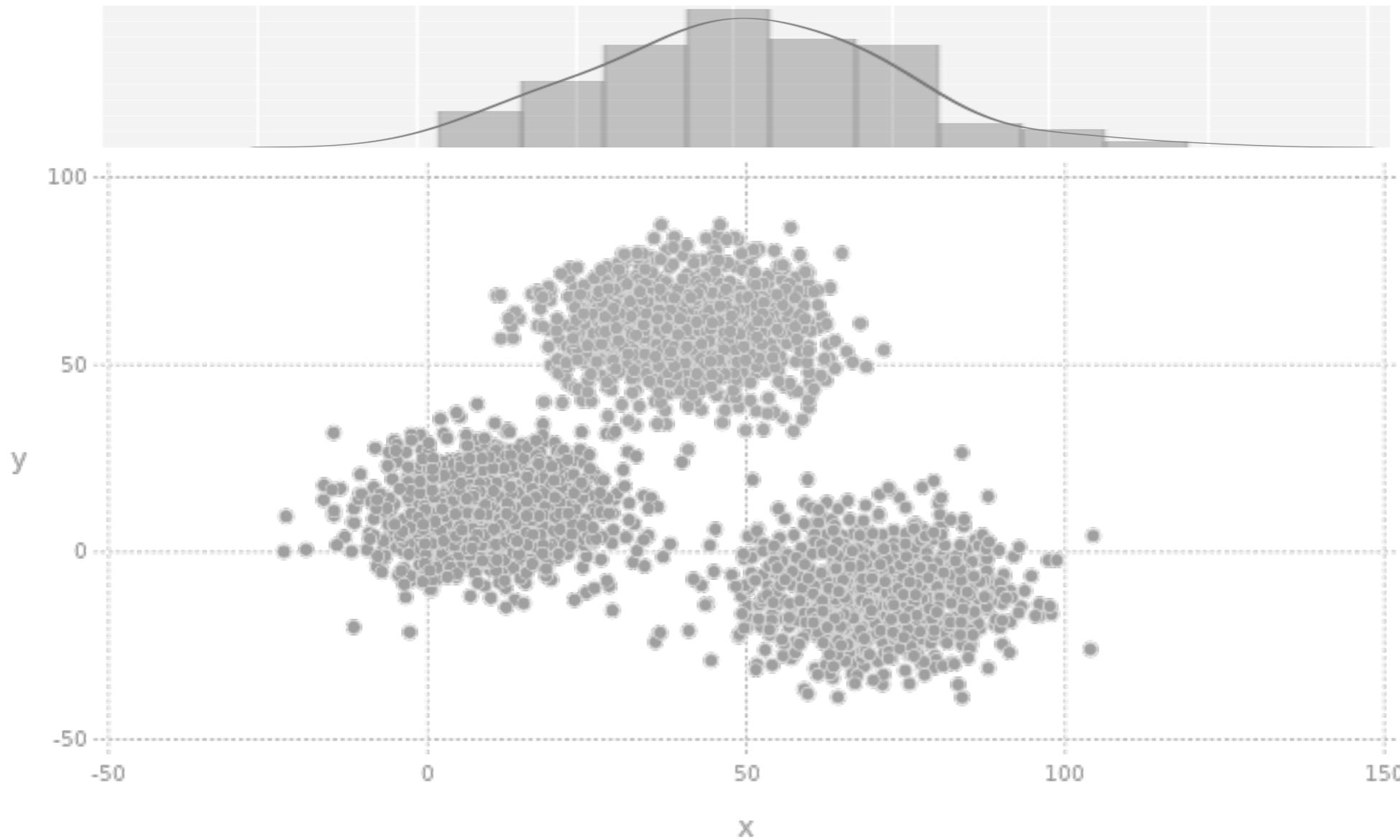
Dimensions to segment on:

- ▶ **Conversion Funnel:**
 - ▶ Visitor → Registered user → Paying customer → Repeat customer
- ▶ **Time since last engagement:**
 - ▶ 7 days ("active") → 30 days ("churn risk") → 90 days ("churned")
- ▶ **Monetisation:**
 - ▶ High vs Medium vs Low value customers
- ▶ And more...product, SKU, behavioural groupings

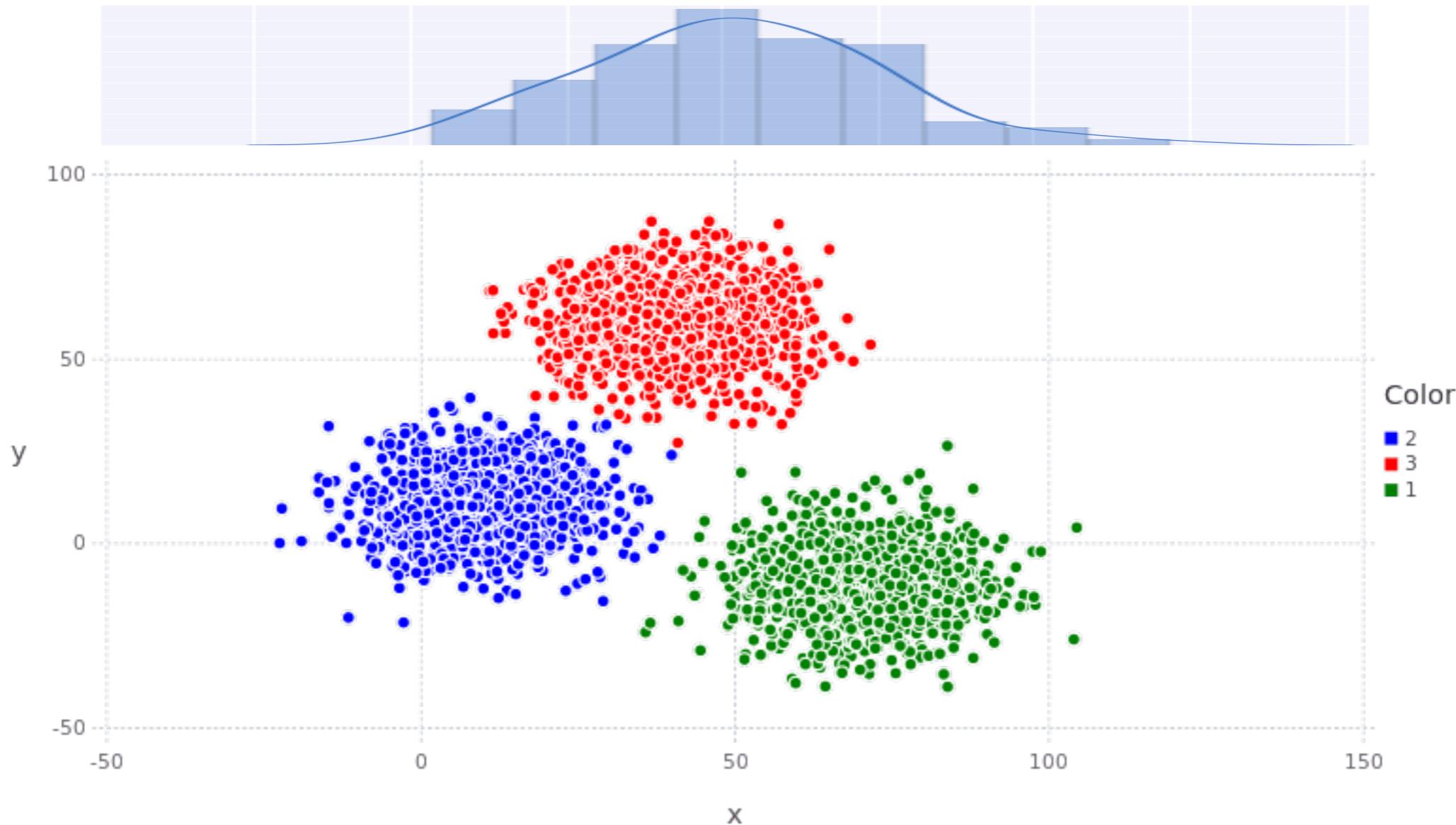
UNSUPERVISED TECHNIQUE: CLUSTERING



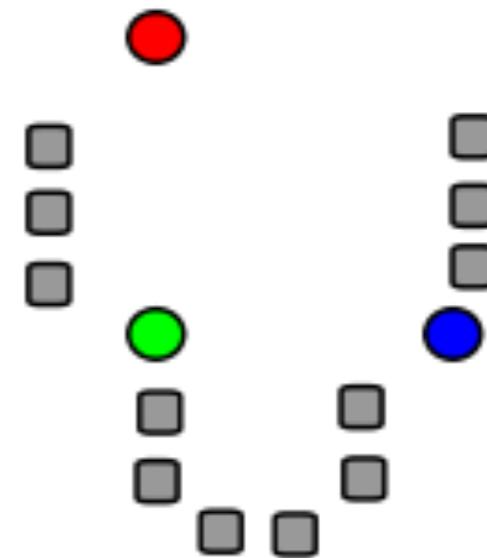
UNSUPERVISED TECHNIQUE: CLUSTERING



UNSUPERVISED TECHNIQUE: CLUSTERING



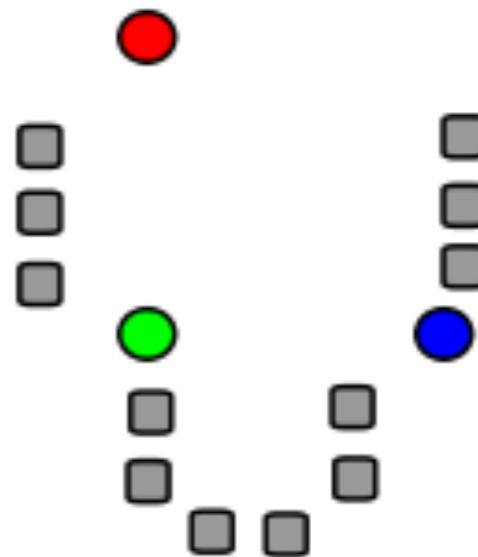
RECAP: K-MEANS CLUSTERING



RECAP: K-MEANS CLUSTERING

Steps

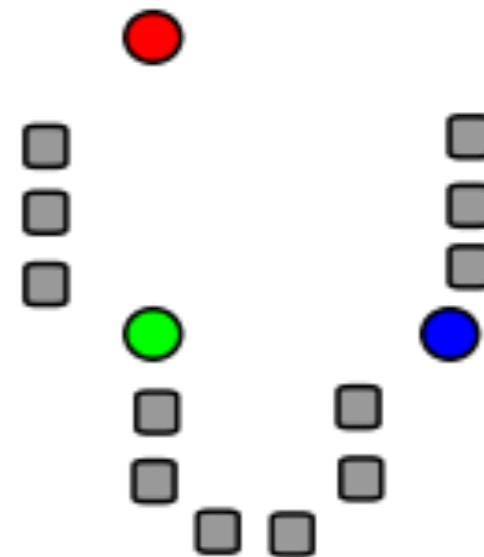
1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

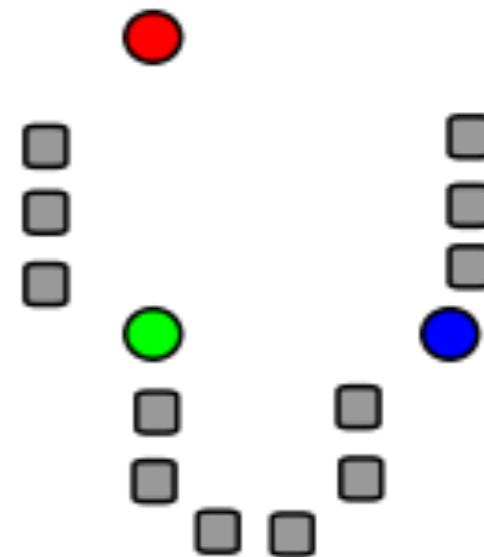
1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

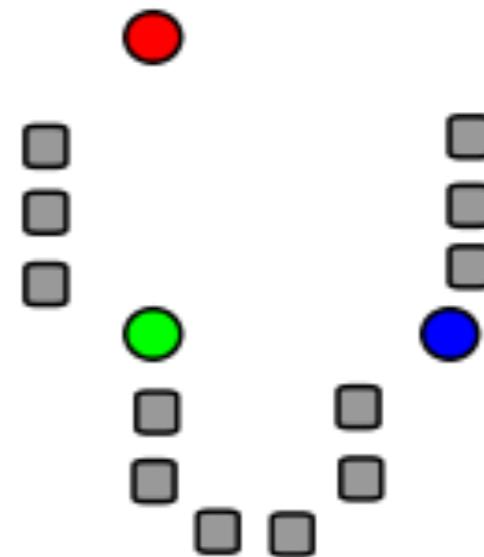
1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

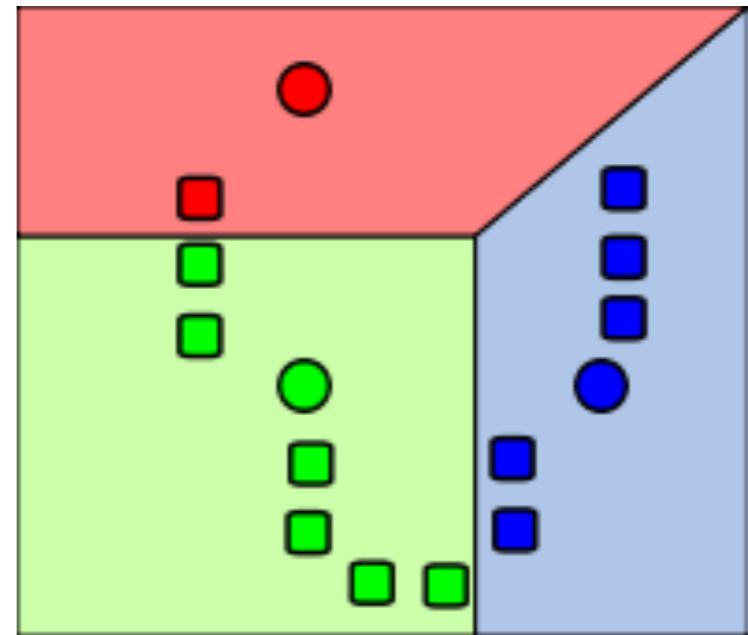
1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

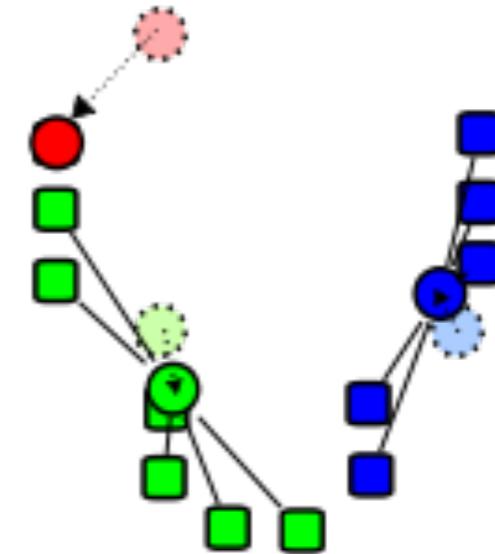
1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

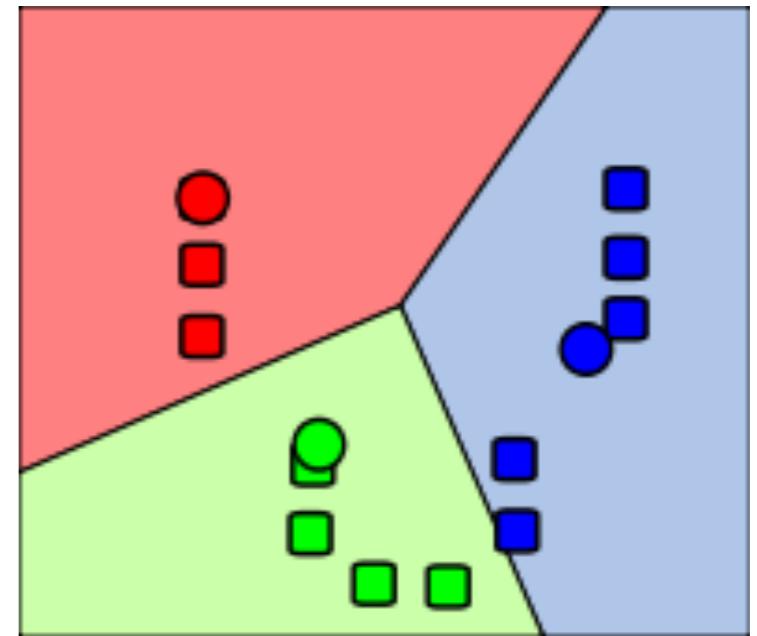
1. Choose k. Here, let's try to cluster into k=3 groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING

Steps

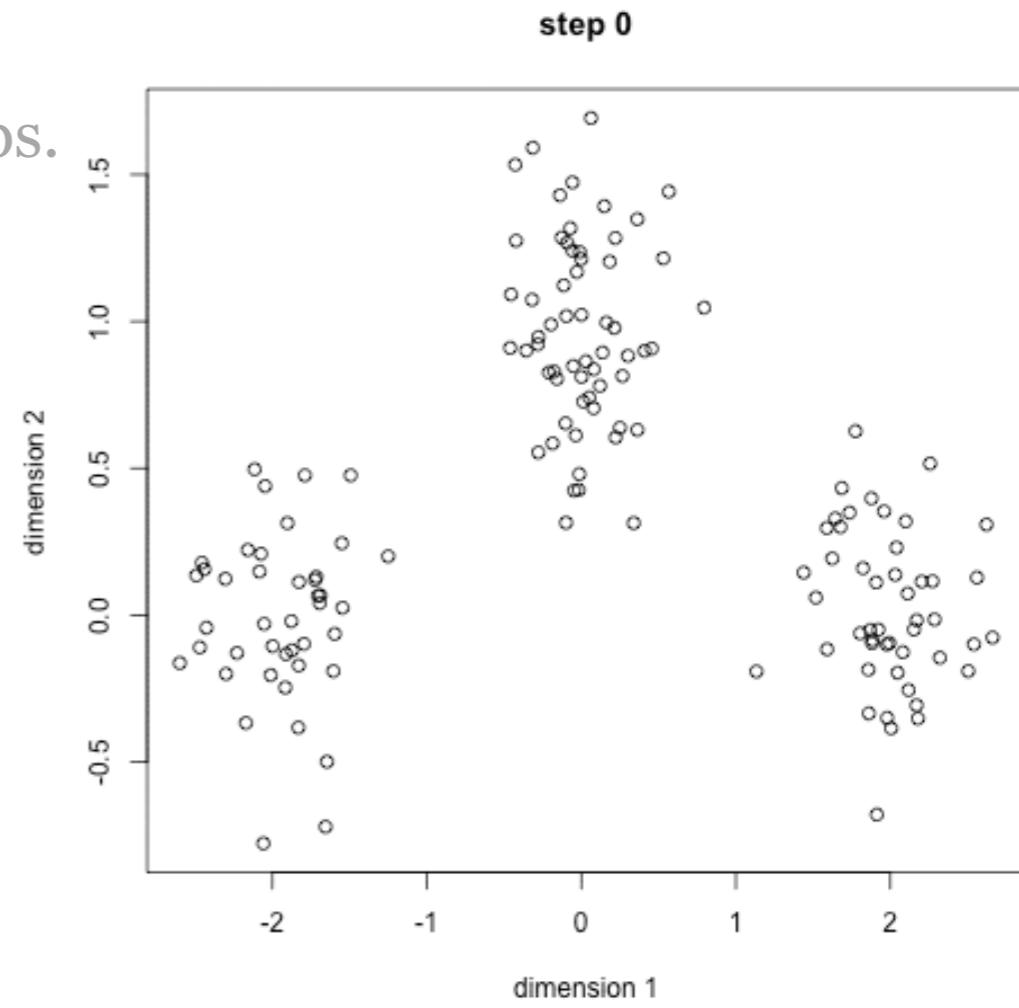
1. Choose k. Here, let's try to cluster into k=3 groups
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



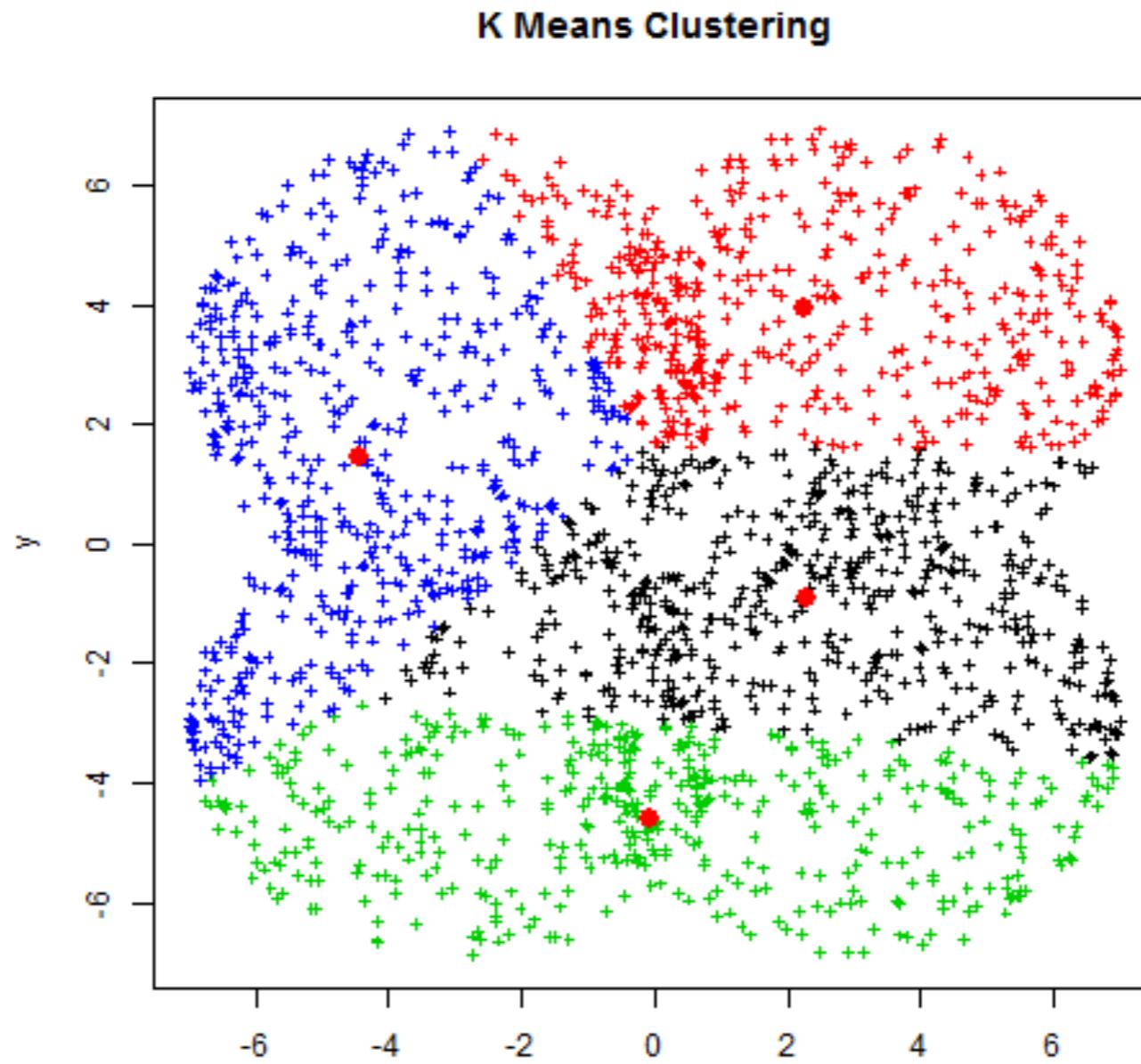
RECAP: K-MEANS CLUSTERING

Steps

1. Choose k. Here, let's try to cluster into $k=3$ groups.
2. Choose 3 initial "centroids"...random points.
3. For each grey point:
 - find distance to nearest centroid
 - assign the point to the nearest centroid's "team"
4. Recalculate centroid positions to become the centre of each cluster
5. Repeat steps 3 and 4 until the centroids don't move



RECAP: K-MEANS CLUSTERING



RECAP: DISTANCE METRICS

K-Means requires the concept of "distance" between two data points. Let's consider a few different metrics.

- ▶ **Euclidean distance...**

- ▶ ...between 1 and 2 is 1
- ▶ ...between point (1,1) and point (2,2) is $\sqrt{2}$
- ▶ ...this works for any "vector" of numbers

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

RECAP: DISTANCE METRICS

K-Means requires the concept of "distance" between two data points. Let's consider a few different metrics.

- ▶ **Euclidean distance...**

- ▶ ...between 1 and 2 is 1
- ▶ ...between point (1,1) and point (2,2) is $\sqrt{2}$
- ▶ ...this works for any "vector" of numbers

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

- ▶ What about the distance between two tweets?

DISTANCE METRICS

K-Means requires the concept of "distance" between two data points. Let's consider a few different metrics.

- **Jaccard distance...**

- ...between "rock music" and "classical music" is "one word shared out of three words total" or $1/3 = 0.33$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} :$$

- The same concept can determine distances between words, tweets, books...even the Jaccard distance between actors in two movies

CASE STUDY

BEHAVIOURAL TARGETING

EXAMPLE: BEHAVIOURAL TARGETING

- ▶ **Problem:** Can we find smarter customer segmentations than RFM/lifecycle models?
- ▶ **Goals:**
 - ▶ Build a dataset of indicators on a per customer basis
 - ▶ Demographic indicators (e.g. gender, age, location, device type)
 - ▶ Acquisition channel indicators (e.g. organic, paid, CAC)
 - ▶ Behavioural indicators (e.g. RFM metrics, genre/topics of interest)
 - ▶ Social indicators (e.g. number of friends using service)
 - ▶ Build a clustering model using this dataset that finds the best "distinct groups"
 - ▶ Interpret what each group represents
 - ▶ Can we use these groupings for targeted newsletter or social media campaigns, and other retention activities?

EXAMPLE: BEHAVIOURAL TARGETING

► **Results:**

- ▶ Six groupings emerged, upon investigation we had groups such as:
 - ▶ Users who open the app almost every day but never book anything
 - ▶ Users who only book free events
 - ▶ Users who book a "date night" style event once a month
 - ▶ Users who book a headline high value item once a year
 - ▶ And so on...
 - ▶ Subgroups within main groups, e.g.
 - ▶ Users who go to free lunchtime concerts
 - ▶ The clustering algorithm can automatically assign new users into one of these categories.
- **How to validate if this segmentation is better than RFM/lifecycle models?**

INTRODUCTION

DATA SCIENCE == SCIENCE

DATA + SCIENCE = DATA SCIENCE

- ▶ Two options for user segmentation for newsletter campaigns:
 - ▶ **Old method:** Most viewed category (e.g. sport, music)
 - ▶ **New method:** Algorithmically generated clusters
- ▶ **Question:** How to know which is better?
- ▶ **Answer:** Run an experiment!
- ▶ **Hypothesis:** "New method is significantly better"
- ▶ **Question:** What do we mean by "better"?
 - ▶ Email campaign opens? CTR? CTR + purchase? Long-term retention metrics?

DATA + SCIENCE = DATA SCIENCE

- ▶ Let's go with "the campaign that generated more revenue". Difficult to argue with that.
- ▶ Run the test:
 - ▶ Divide all customers into two groups at random
 - ▶ A gets old segmentation
 - ▶ B gets new segmentation
 - ▶ ...wait...
 - ▶ Crunch numbers once sample size is reached
- ▶ Welcome to the wonderful world of A/B testing!

A/B TESTING

heartwarming tale of a young boy's triumph over adversity. It's one of the most award-

Book Now

VS

He sweats profusely, never reaches a punchline, and often finds himself off topic.

Book Now

HOW NOT TO RUN AN A/B TEST

Statistical **significance** is for scientists!

Don't bother calculating a **sample size**

End the A/B test **early**

TOP TIP: EVANMILLER.ORG

Evan's Awesome A/B Tools ([home](#)):

[Sample Size Calculator](#) | [Chi-Squared Test](#) | [Sequential Sampling](#) | [Count Data](#) | [Survival Times](#) | [2 Sample T-Test](#)

Question: How many subjects are needed for an A/B test?

Baseline conversion rate: %  20% [\[link \]](#)

Minimum Detectable Effect: %  15% – 25%

The Minimum Detectable Effect is the smallest effect that will be detected $(1-\beta)\%$ of the time.

Absolute Relative

Conversion rates in the gray area will not be distinguishable from the baseline.

Sample size:

1,030

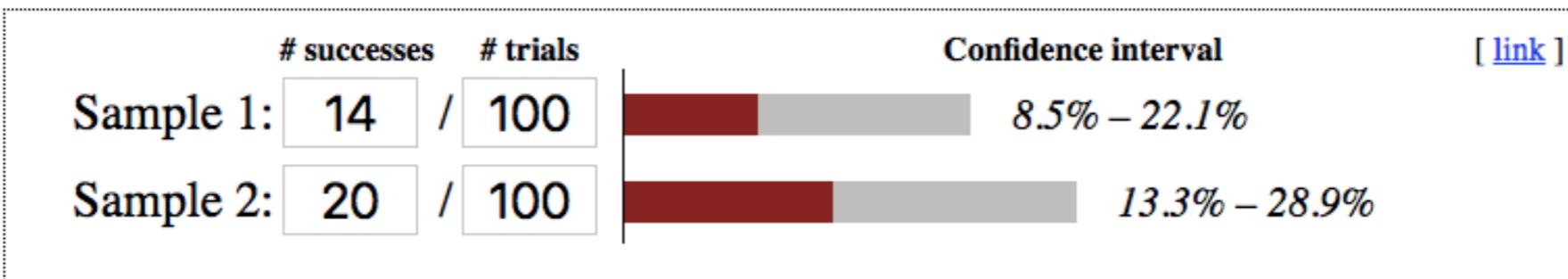
per variation

TOP TIP: EVANMILLER.ORG

Evan's Awesome A/B Tools ([home](#)):

[Sample Size Calculator](#) | [Chi-Squared Test](#) | [Sequential Sampling](#) | [Count Data](#) | [Survival Times](#) | [2 Sample T-Test](#)

Question: Does the rate of success differ across two groups?



Verdict:

No significant difference

($p = 0.26$)

Confidence level: 95%

TOP TIP: EVANMILLER.ORG

Evanmiller.org

How Not To Run An A/B Test

By [Evan Miller](#)

April 18, 2010

If you run A/B tests on your website and regularly check ongoing experiments for significant results, you might be falling prey to what statisticians call *repeated significance testing errors*. As a result, even though your dashboard says a result is statistically significant, there's a good chance that it's actually insignificant. This note explains why.

Background

When an A/B testing dashboard says there is a “95% chance of beating original” or “90% probability of statistical significance,” it’s asking the following question: Assuming there is no underlying difference

A/B TESTING = GROWTH

2048

different versions of the app

A/B TESTING = GROWTH



A/B TESTING = GROWTH

2048

different versions of the app

Fearless experimentation!

A/B TESTING = GROWTH

What Are You Into?



Tell Us What You Like



A/B TESTING = GROWTH

2048

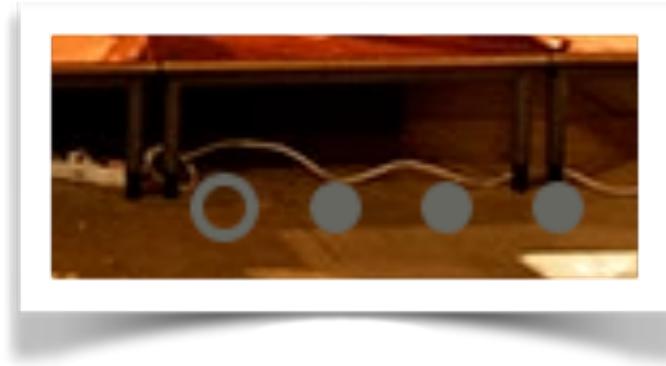
different versions of the app

Fearless experimentation!

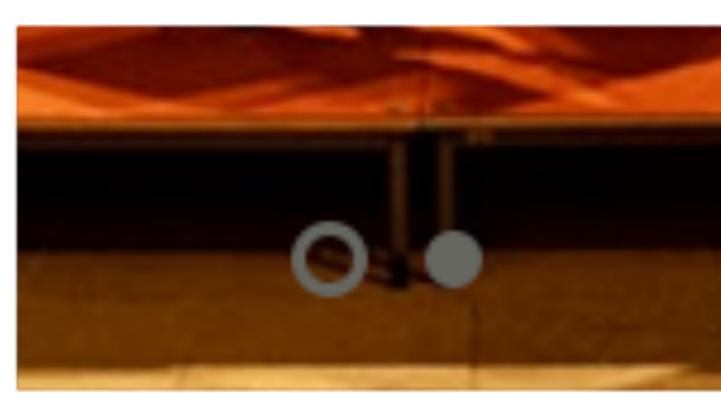
Not significant doesn't mean

boring

A/B TESTING = GROWTH



4 images



2 images



A WORD OF WARNING



JUNE 18, 2014



How Optimizely (Almost) Got Me Fired

<http://blog.sumall.com/journal/optimizely-got-me-fired.html>

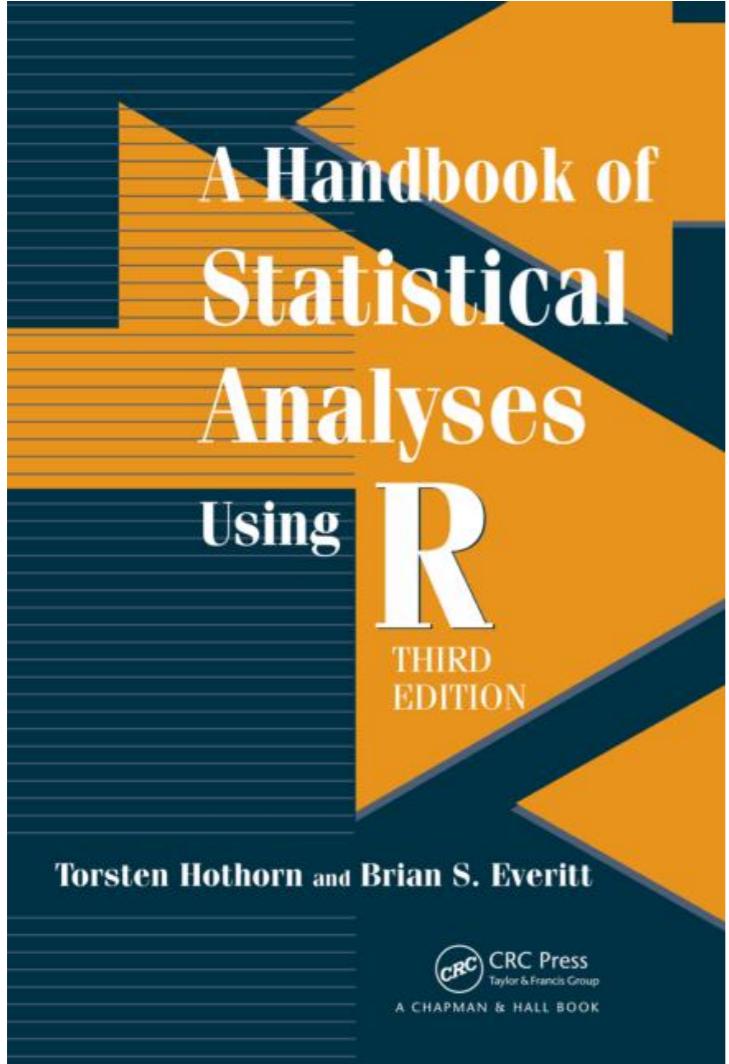
INTRODUCTION

INFERENCE TESTS

CHI-SQUARED

	Group A	Group B
Viewed signup button	100	100
Clicked signup button	20	25

HANDBOOK OF STATISTICAL ANALYSES USING R

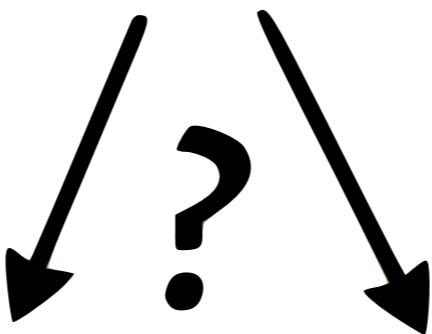


- ▶ **2nd edition is freely available online as a PDF**
- ▶ **#1 podcast recommendation:**
 - ▶ Data Skeptic

INTRODUCTION

P HACKING

HYPOTHESIS TESTING AND P-VALUES



BEWARE P-HACKING



BEWARE P-HACKING

WE FOUND NO
LINK BETWEEN
PURPLE JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
BROWN JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
PINK JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
BLUE JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
TEAL JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
GREY JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
TAN JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
YAN JELLY
JEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
SALMON JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
RED JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
TURQUOISE JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
MAGENTA JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
YELLOW JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
BEIGE JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
LILAC JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND NO
LINK BETWEEN
BLACK JELLY
BEANS AND ACNE
($P > 0.05$).



BEWARE P-HACKING

WE FOUND NO
LINK BETWEEN
CYAN JELLY
BEANS AND ACNE
($P > 0.05$).



WE FOUND A
LINK BETWEEN
GREEN JELLY
BEANS AND ACNE
($P < 0.05$).



= News =
GREEN JELLY BEANS LINKED TO ACNE!

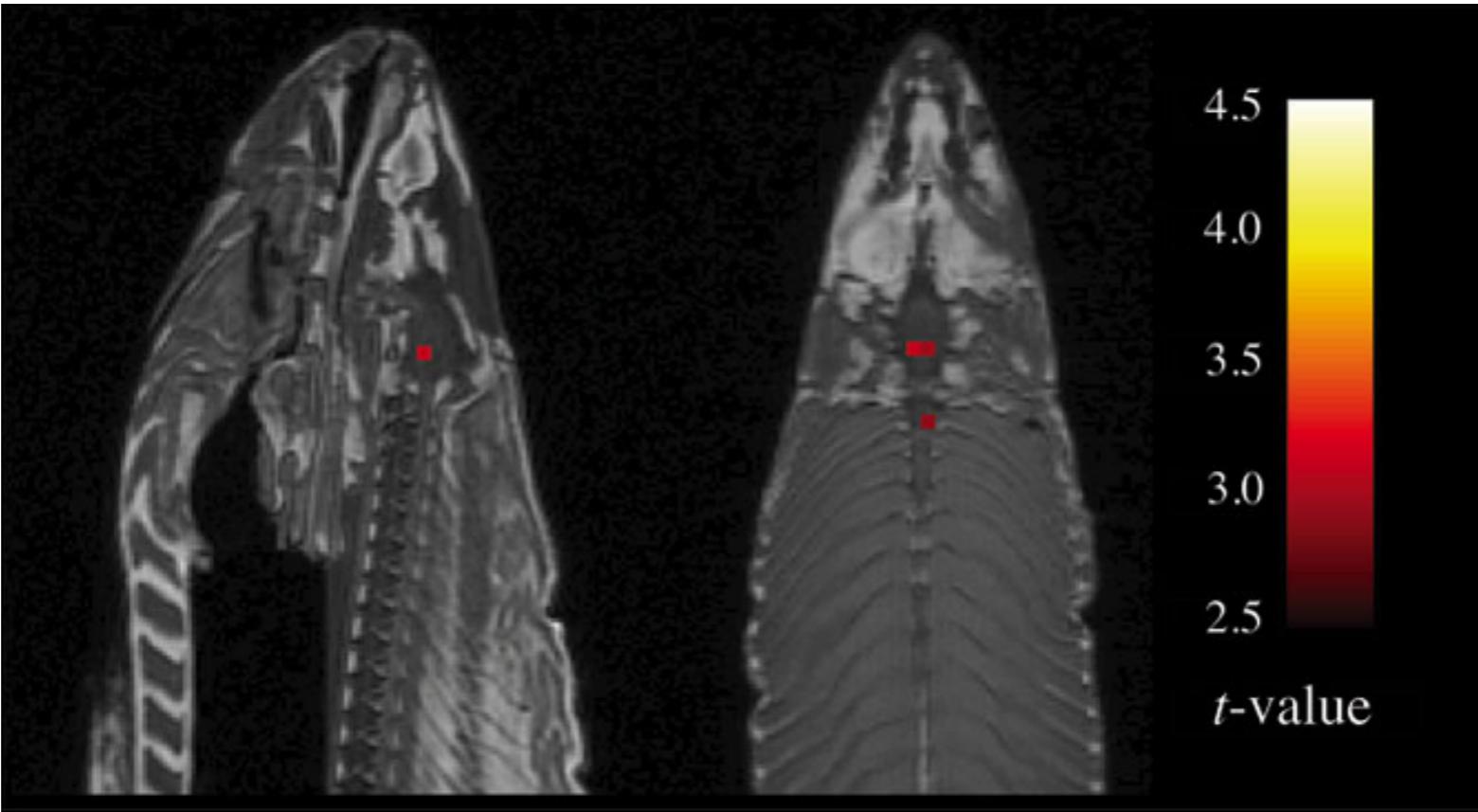
95% CONFIDENCE

ONLY 5% CHANCE OF COINCIDENCE!



SCIENTISTS

CASE STUDY: THE DEAD SALMON STUDY



Bennett et al. "Neural Correlates of Interspecies Perspective Taking in the Post-Mortem Atlantic Salmon: An Argument For Proper Multiple Comparisons Correction". Journal of Serendipitous and Unexpected Results, 2010. IgNobel Award in Neuroscience 2012.

CASE STUDY: THE DEAD SALMON STUDY



[Live Science](#) > [Strange News](#)

Dead Salmon 'Responds' to Pictures of People

By Robert Roy Britt | September 27, 2009 07:03am ET



A dead salmon has become a scientific celebrity after its brain supposedly lit up when shown pictures of humans during a brain scan.



CASE STUDY: THE DEAD SALMON STUDY



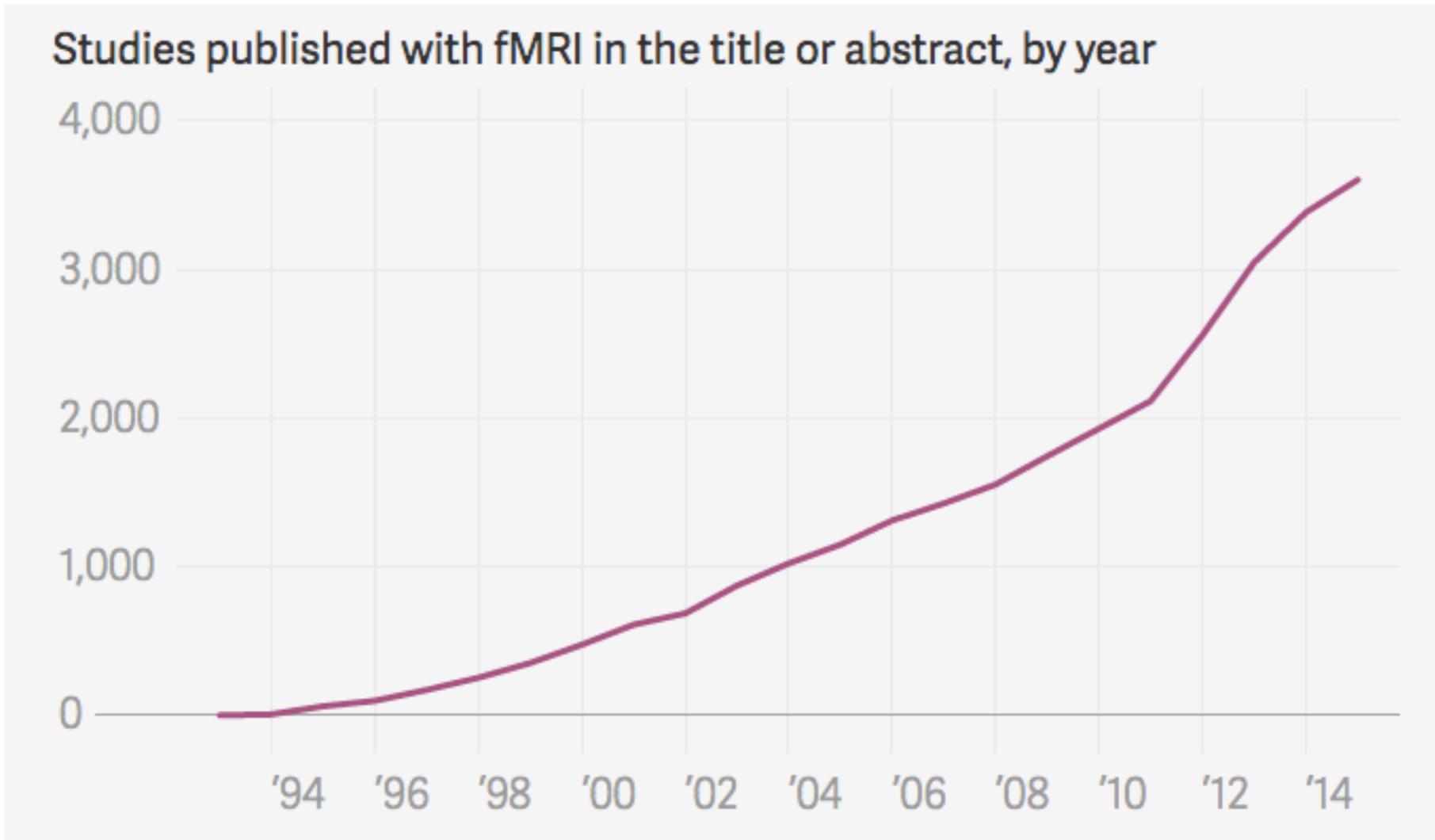
ALEXIS MADRIGAL SCIENCE 09.18.09 5:37 PM

WIRE.D

SUBSCRIBE

SCANNING DEAD SALMON IN FMRI MACHINE HIGHLIGHTS RISK OF RED HERRINGS

CASE STUDY: THE DEAD SALMON STUDY



P-HACKING: TRY IT OUT!

Hack Your Way To Scientific Glory

You're a social scientist with a hunch: **The U.S. economy is affected by whether Republicans or Democrats are in office.** Try to show that a connection exists, using real data going back to 1948. For your results to be publishable in an academic journal, you'll need to prove that they are "statistically significant" by achieving a low enough p-value.

1 CHOOSE A POLITICAL PARTY

Republicans

Democrats

2 DEFINE TERMS

Which politicians do you want to include?

Presidents

Governors

Senators

Representatives

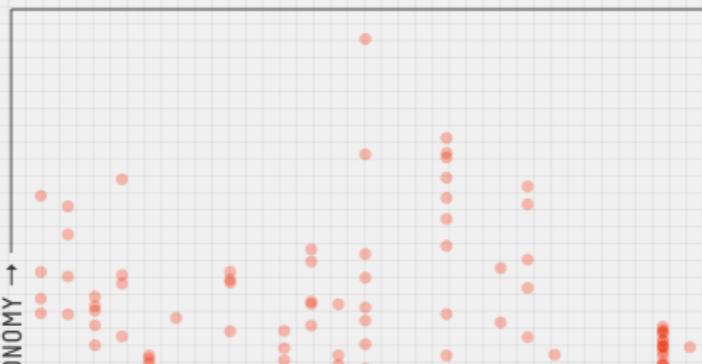
How do you want to measure economic performance?

Employment

...

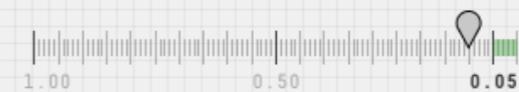
3 IS THERE A RELATIONSHIP?

Given how you've defined your terms, does the economy do better, worse or about the same when more Republicans are in office? Each dot below represents one month of data.



4 IS YOUR RESULT SIGNIFICANT?

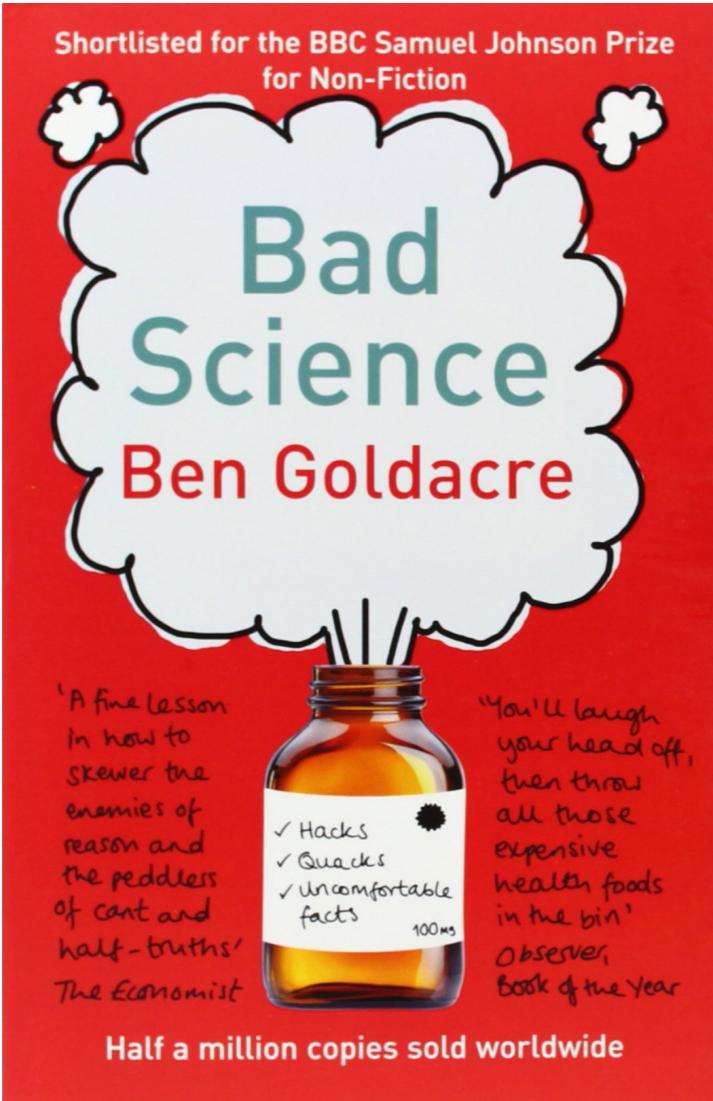
If there were no connection between the economy and politics, what is the probability that you'd get results at least as strong as yours? That probability is your p-value, and by convention, you need a **p-value of 0.05 or less** to get published.



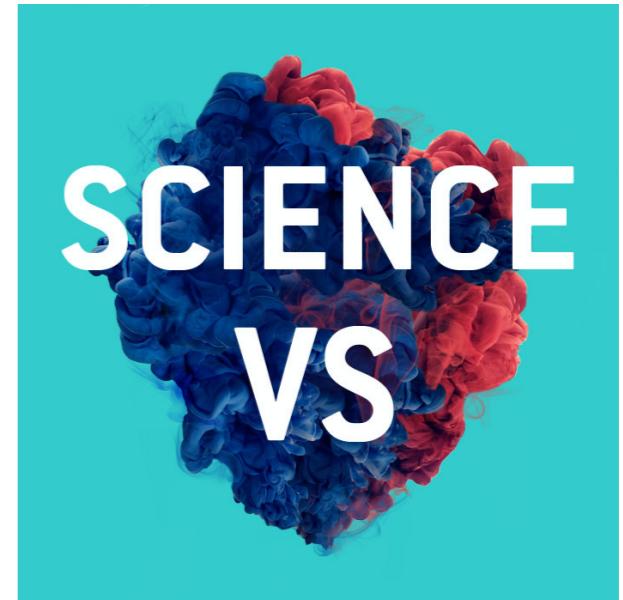
Result: Almost

Your **0.10 p-value** is close to the **0.05** threshold. Try tweaking your variables to see if you can push it

RECOMMENDED READING



Podcasts



FiveThirtyEight

CASE STUDY

BUILDING YPLAN'S RECOMMENDATION ENGINE

INTRODUCTION

CLOSING

DATA SCIENCE VS DATA ANALYTICS

"Data Analytics"

Historical reporting.

Metrics. KPIs. Segmentation.

Dashboards. BI tools. Pivot tables.

Necessary...keeps the engines running.

Tools: Excel, SQL, Tableau.

"Data Science"

Predictive forecasting.

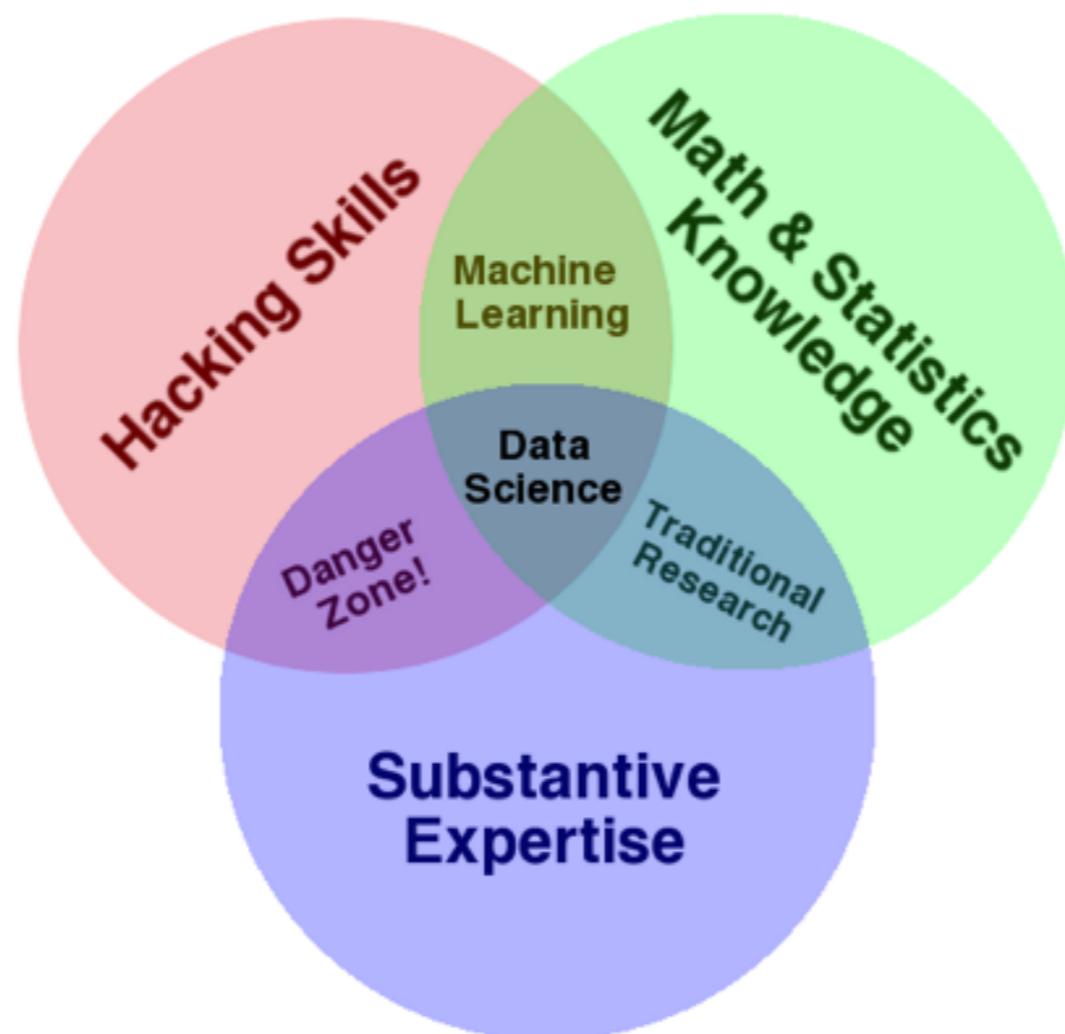
Statistics. Regression. Machine learning.

Coding. Flexibility. Automation.

Exciting...unexpected insights.

Tools: Python, R, scikit-learn.

DATA ANALYST...OR DATA SCIENTIST?



source: <http://www.dataists.com/2010/09/the-data-science-venn-diagram/>

WHAT IS DATA SCIENCE FOR ME?

"Data Analytics"

Historical reporting.

Metrics. KPIs. Segmentation.

Dashboards. BI tools. Pivot tables.

Necessary...keeps the engines running.

Tools: Excel, SQL, Tableau.

"Data Science"

Predictive forecasting.

Statistics. Regression. Machine learning.

Coding. Flexibility. Automation.

Exciting...unexpected insights.

Tools: Python, R, scikit-learn.

WHAT IS DATA SCIENCE FOR ME?

"Data Analytics"

Historical reporting.
Metrics. KPIs. Segmentation.
Dashboards. BI tools. Pivot tables.
Necessary...keeps the engines running.
Tools: Excel, SQL, Tableau.

"Data Science"

Predictive forecasting.
Statistics. Regression. Machine learning.
Coding. Flexibility. Automation.
Exciting...unexpected insights.
Tools: Python, R, scikit-learn.

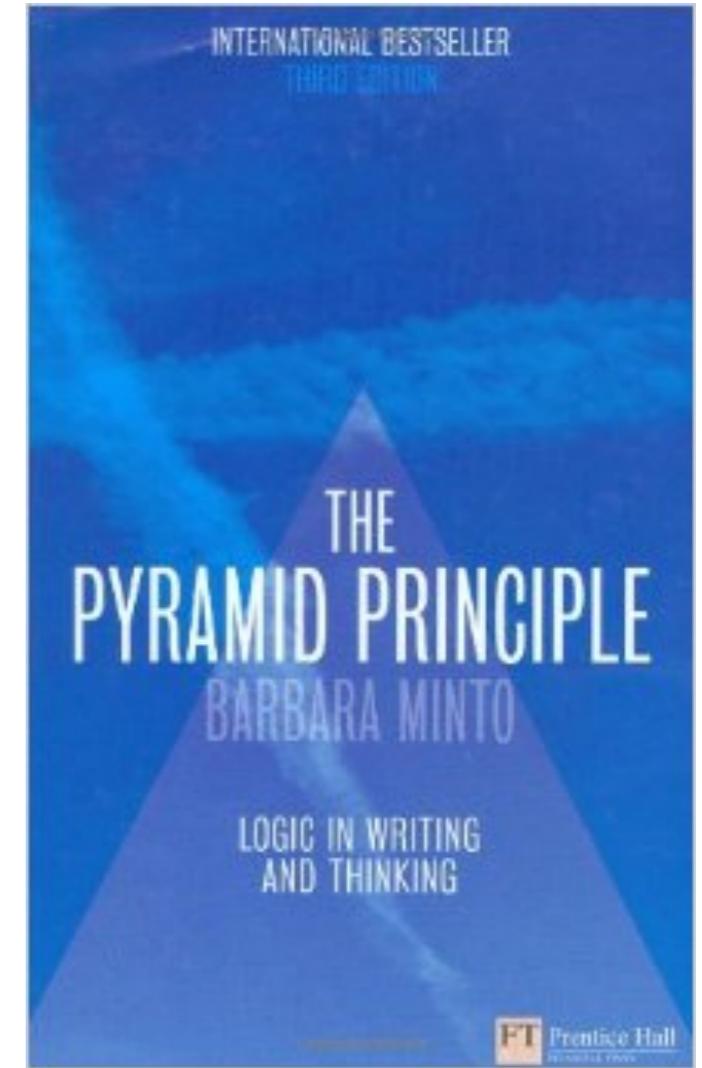
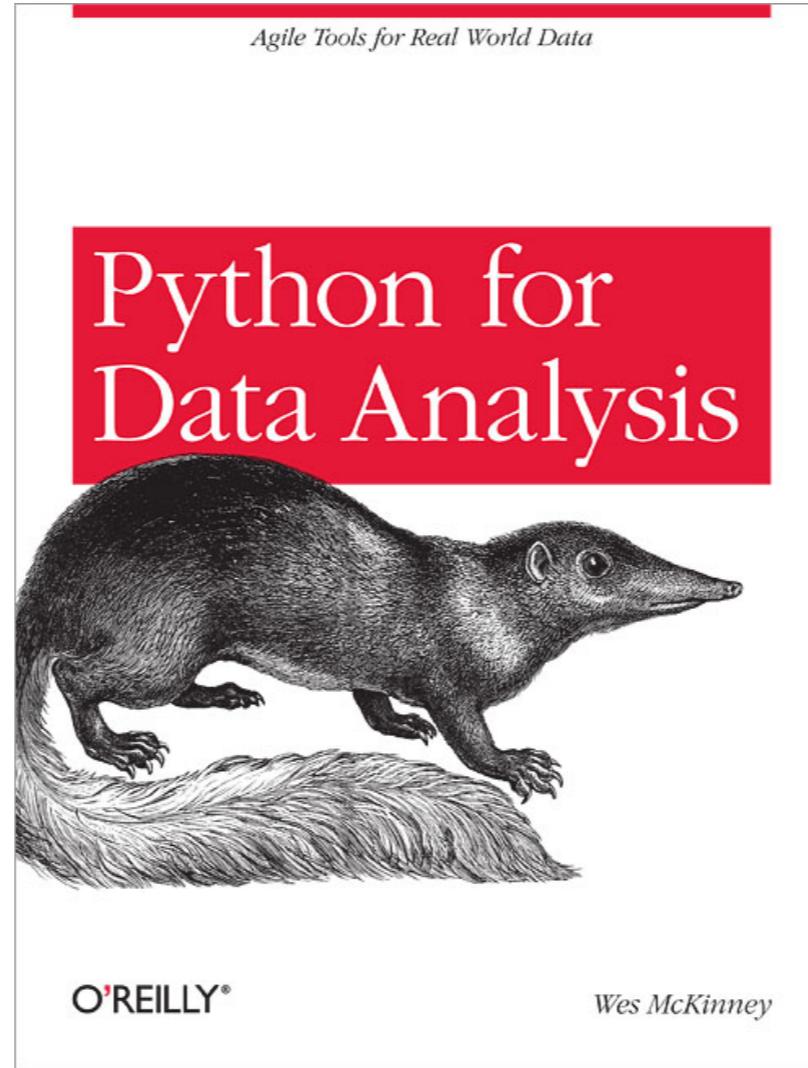
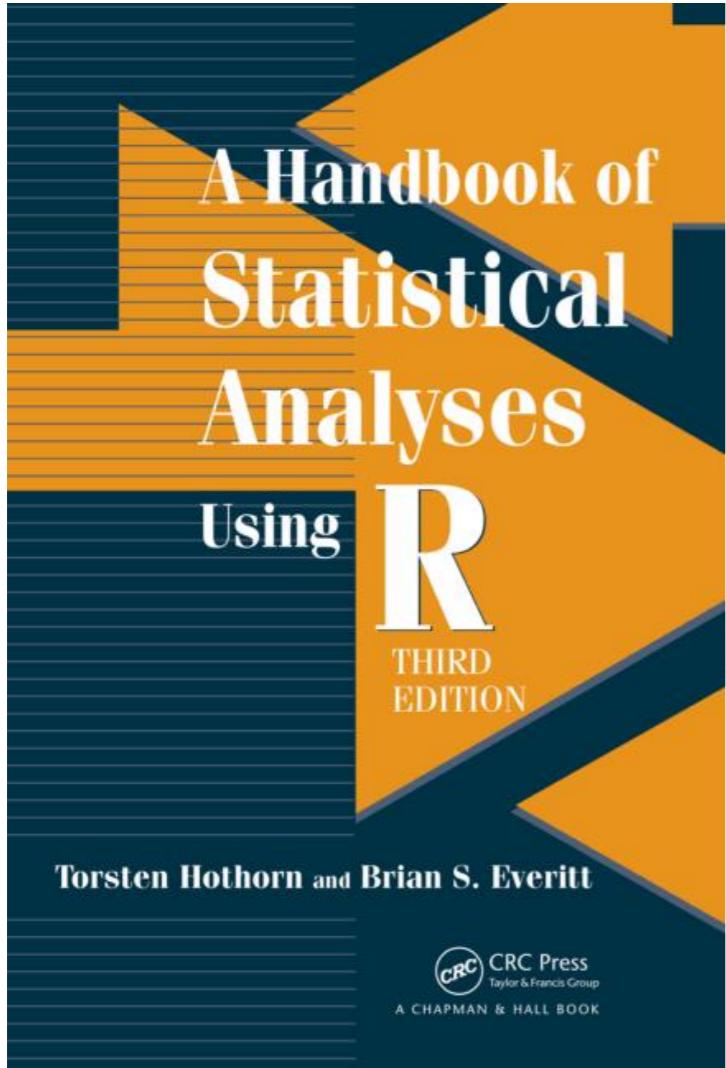
"Data Engineering"

Architecture. Devops. Cloud solutions.
Databases. Data warehouses. Big data.
Integrations (e.g. tracking, channel attribution).
BI tools. Automated reporting. Bespoke solutions.
Version control. Repo management. Code review.

"Strategic Analysis"

Business skills. Startup methodology. Working lean.
Measuring success. KPIs. Data-informed decisions.
Communication. Technical writing. Domain expertise.
Project management. Agile workflows. Problem solving.
Education. Hiring. Mentoring. Advisory.

MY TOP 3 BOOK RECOMMENDATIONS



PODCASTS

- ▶ **Data Skeptic** (Kyle Polich, I ❤️ the mini-explainer episodes!)
- ▶ **Partially Derivative** (light hearted)
- ▶ **Linear Digressions** (Udacity)
- ▶ **More or Less** (Tim Harford & BBC Radio 4)
- ▶ **O'Reilly Data Show** (Ben Lorica, technical with more focus on data engineering)
- ▶ **Planet Money** (NPR, economics/data/finance – A/B testing, multiple comparisons)
- ▶ **What's The Point** (FiveThirtyEight, how data is changing our lives)
- ▶ **Science Vs** (Gimlet Media, new this summer, controversial issues + rigour)

LONDON MEETUPS

1. **PyData London**
2. **LondonR**
3. **Data Science Meetup London**
4. **Big Data London**
5. **London Machine Learning Meetup**
6. **Quantified Self**
7. **Predictive Analytics London Meetup**
8. **Data Visualization Meetup**
9. **PyLadies London**
10. **Women in Data**
11. **Londata**
12. **Data Science Journal Club**

HACKATHONS AND DATADIVES

- ▶ **DataKind**
- ▶ **NHS Hack**
- ▶ **Kaggle**
- ▶ **UK Hackathons & James Meetup**
- ▶ **StartupWeekend**
- ▶ **Code for Good**

FOUR STEPS TO SUCCESS

1. Learn to code

Python. R. Professional software engineering practices.

2. Get statistical

Significance. Inference. Regression. Machine learning.

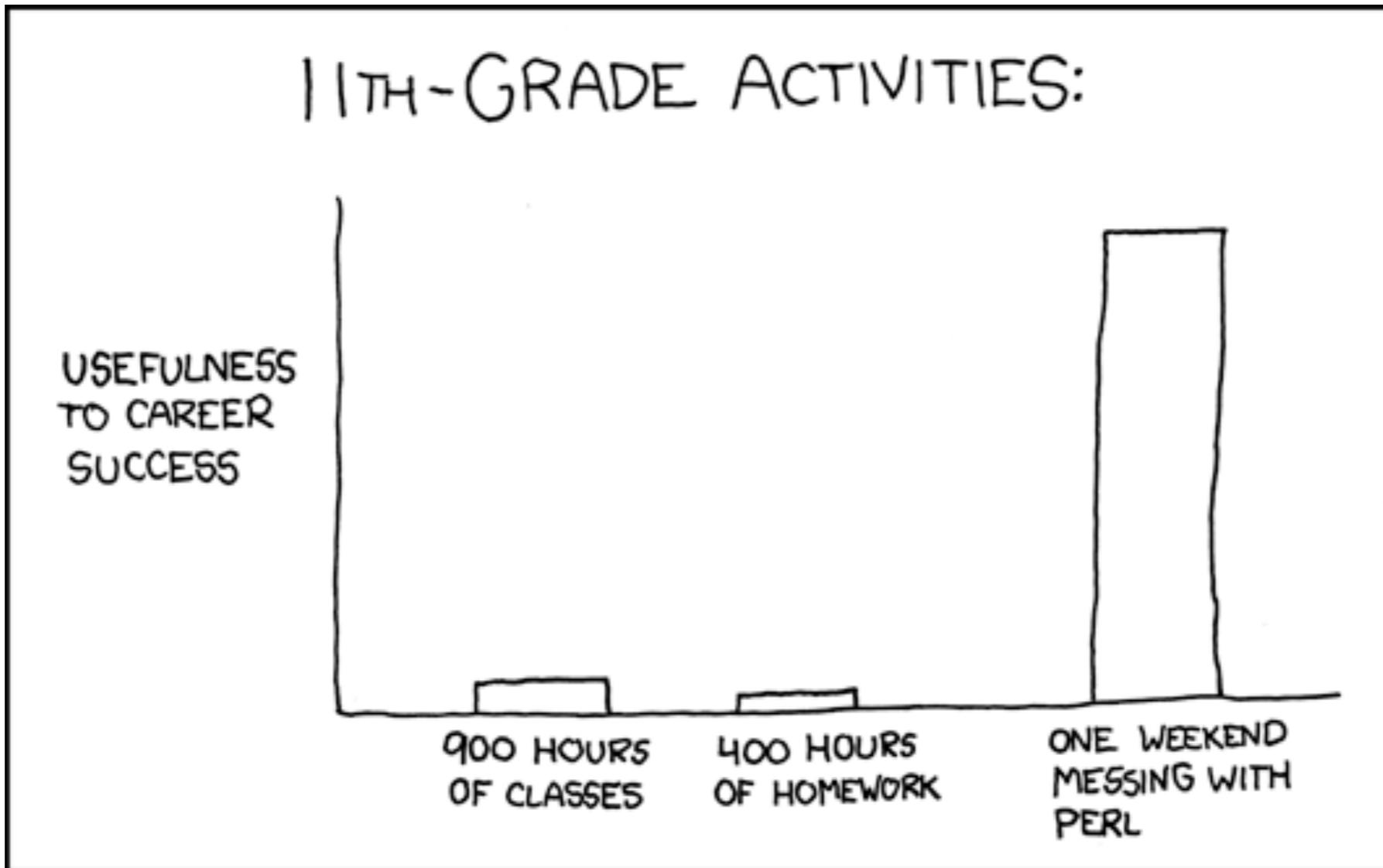
3. Learn lean

Business skills. Startup methodology. Communication.

4. Experience

Side projects. Github. Kaggle. Hackathons. Stand out.

"BECOME A TOP DATA SCIENTIST WITH THIS ONE WEIRD TIP..."



DATA SCIENCE IMMERSIVE

QUESTIONS?