

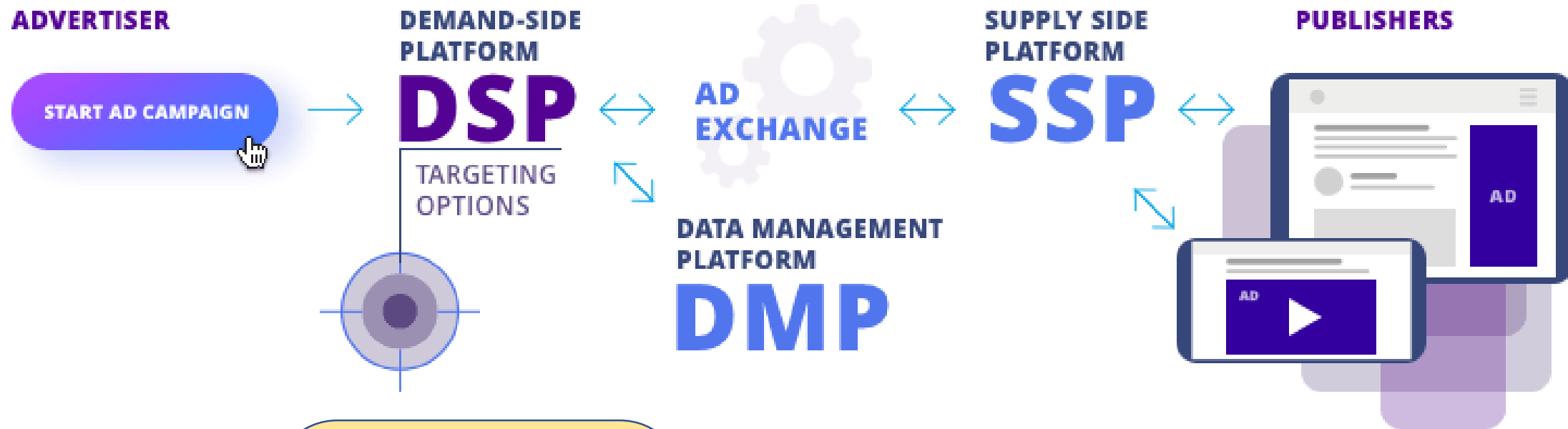
Doing online experiments the proper way: true effect or due to chance?

Uyen Nguyen

Data Scientist, Knorex

06/06/2019

Demand side platform (DSP)



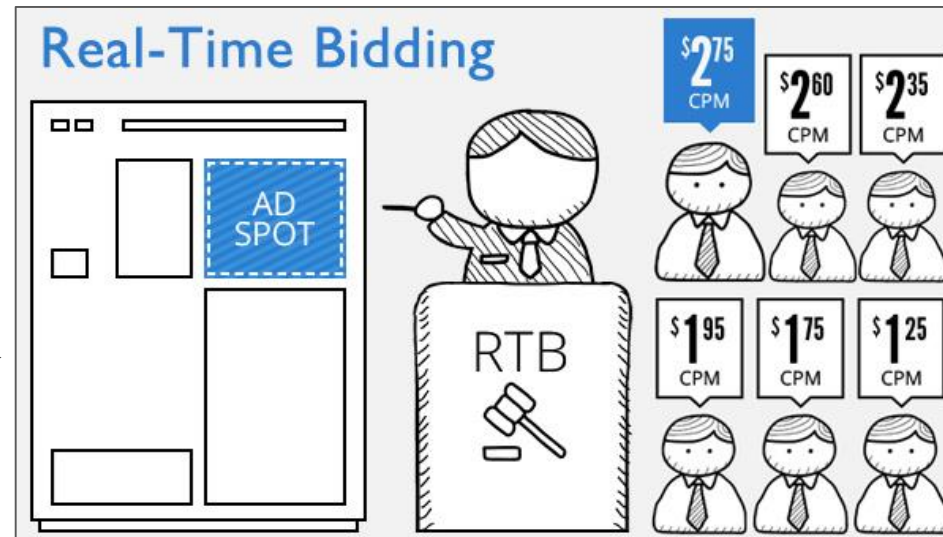
- Targeting
- Retargeting
- Audience profiling
- Brand safety
- Dynamic creative optimization
- Real time bidding

(Source: [Quora](#))

Bidding strategies to optimize for:

- cost per click (CPC)
- cost per acquisition (CPA)
- cost per installation (CPI)

Predicting chances of
click/conversion



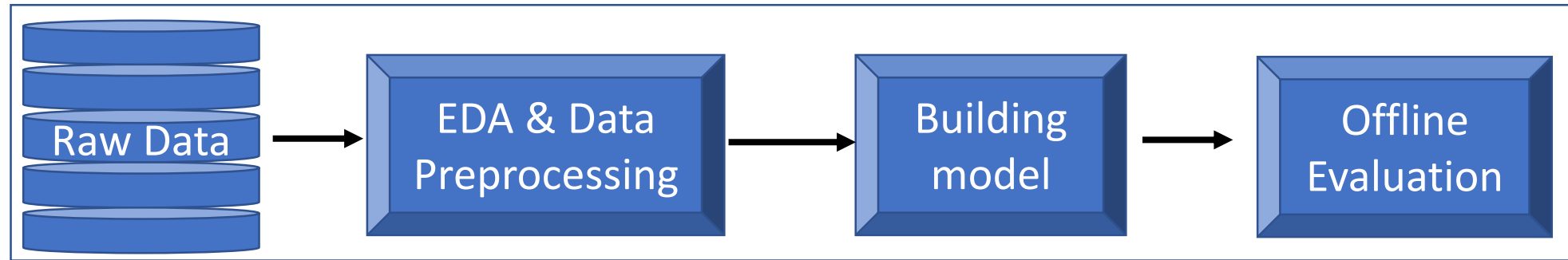
Source: www.sitescout.com

Audience look-alike modeling

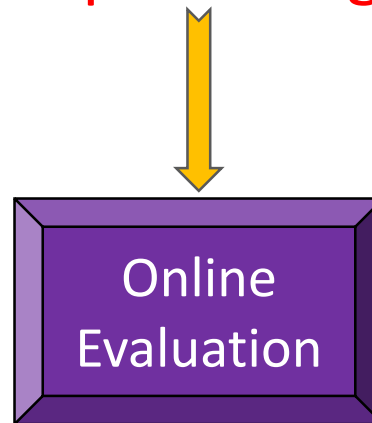
Monitoring performance,
strategic troubleshooting

Testing a new idea in production – Online experimentation

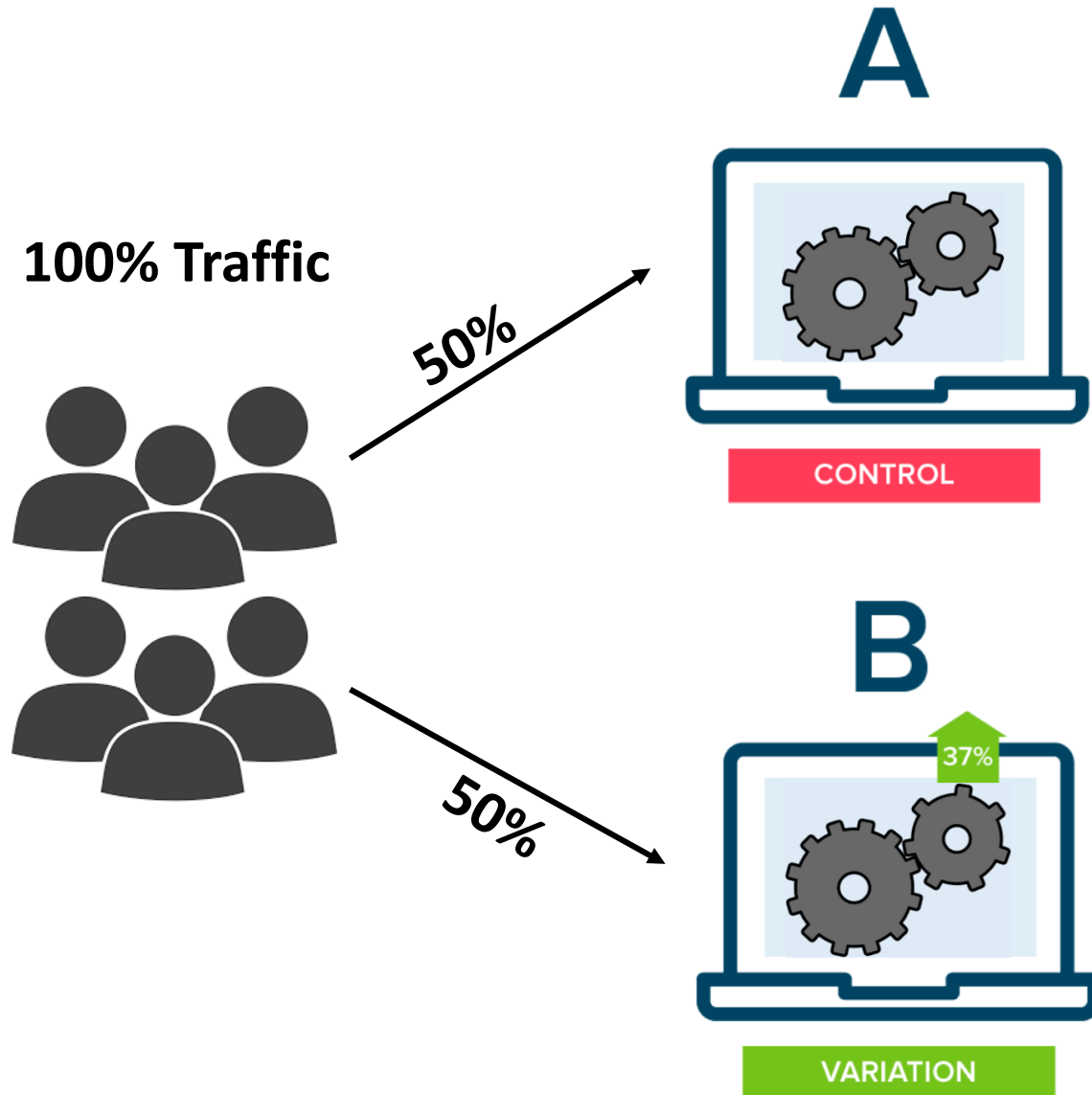
- Ad creative
- Bidding strategy
- Machine learning model



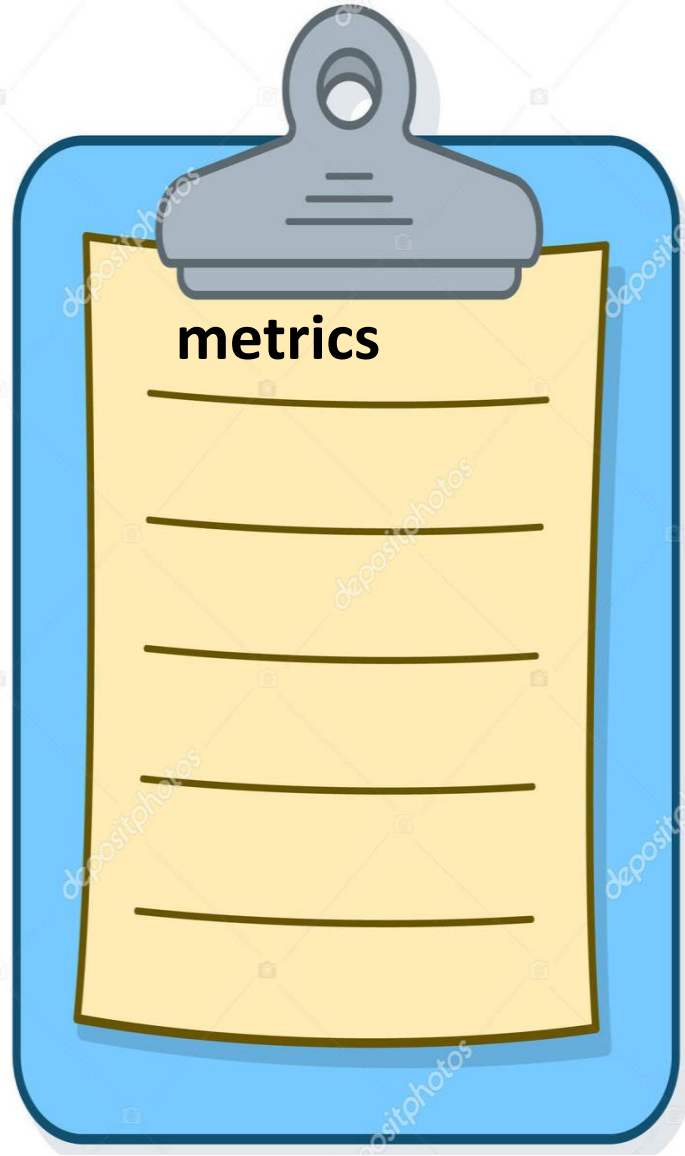
Arriving at a promising idea offline



A/B testing – The conventional method

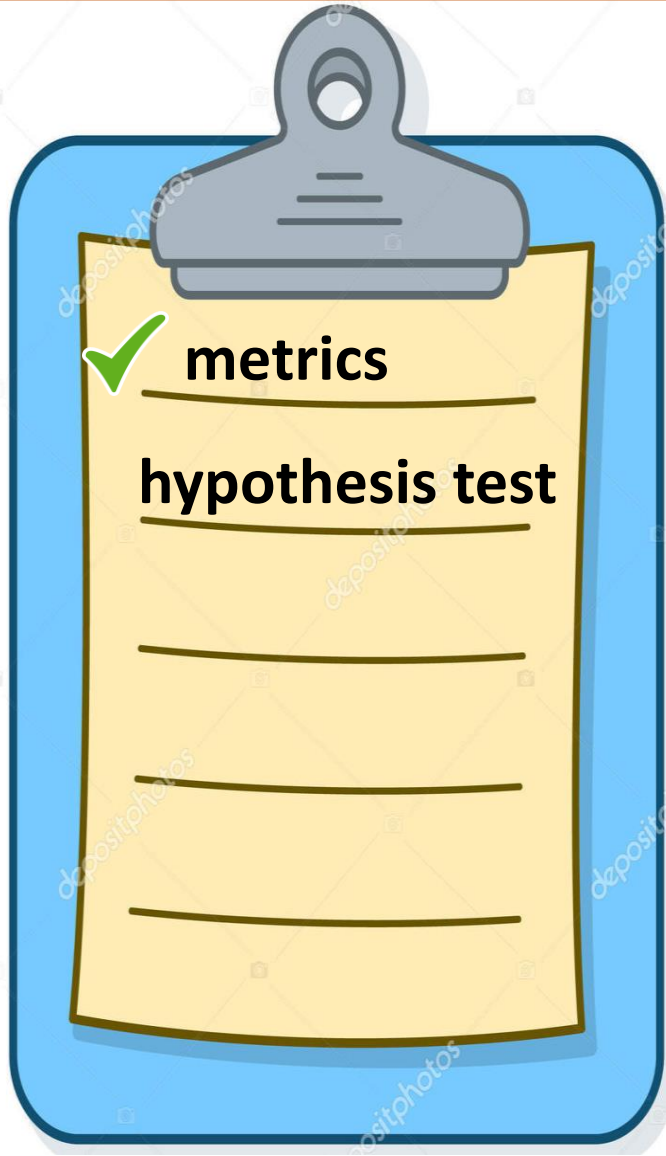


- straight-forward, easy to implement
- 100 years ago: Sir Ronald Fisher – tested different fertilizers on crop yields
- Google ran the first test in 2000. In 2011: they ran > 7000 A/B tests



- E.g, CPC, CPA, CTR
- Choosing the right metrics is critical
- Avoiding confounding variables

A/B testing – Is test result statistically significant?

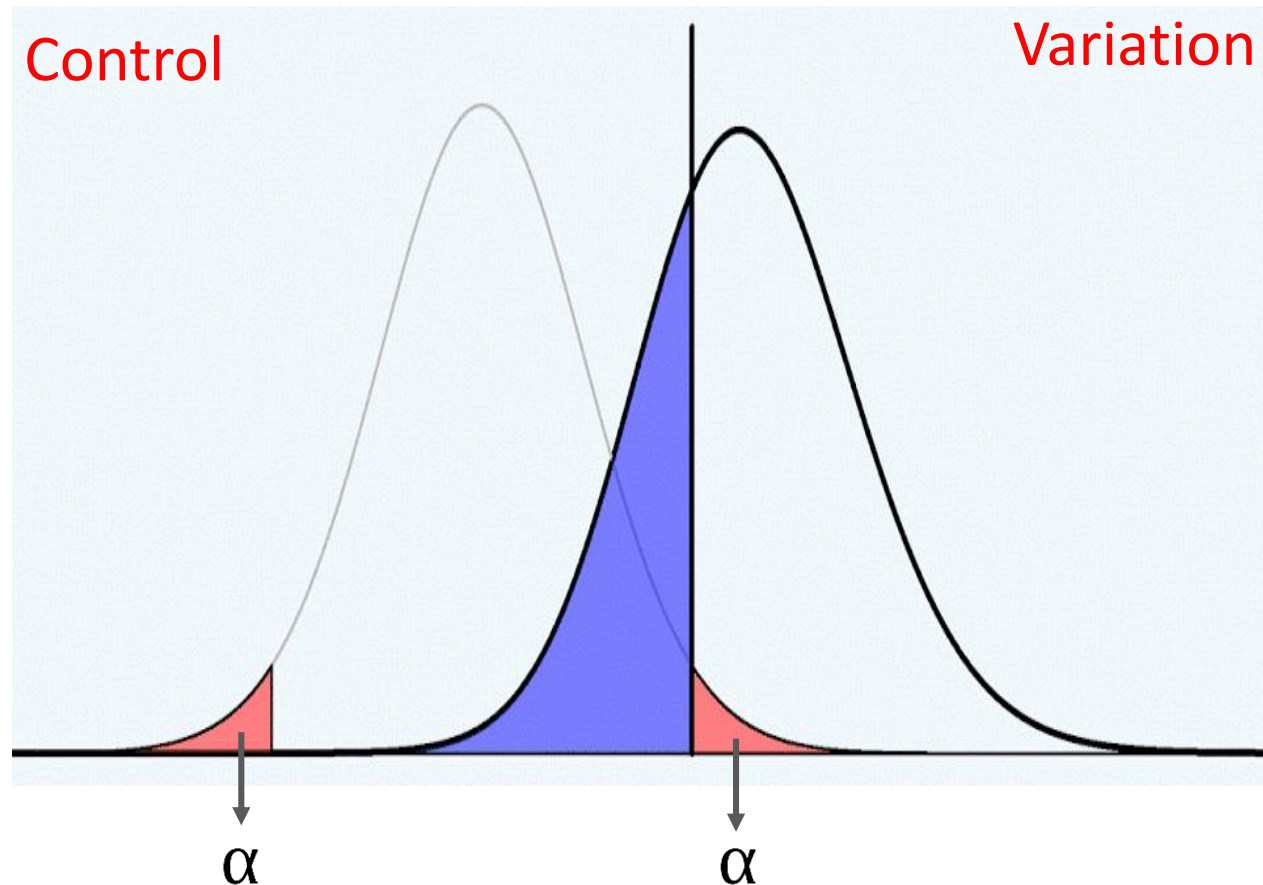


Null hypothesis (H_0): $A = B$

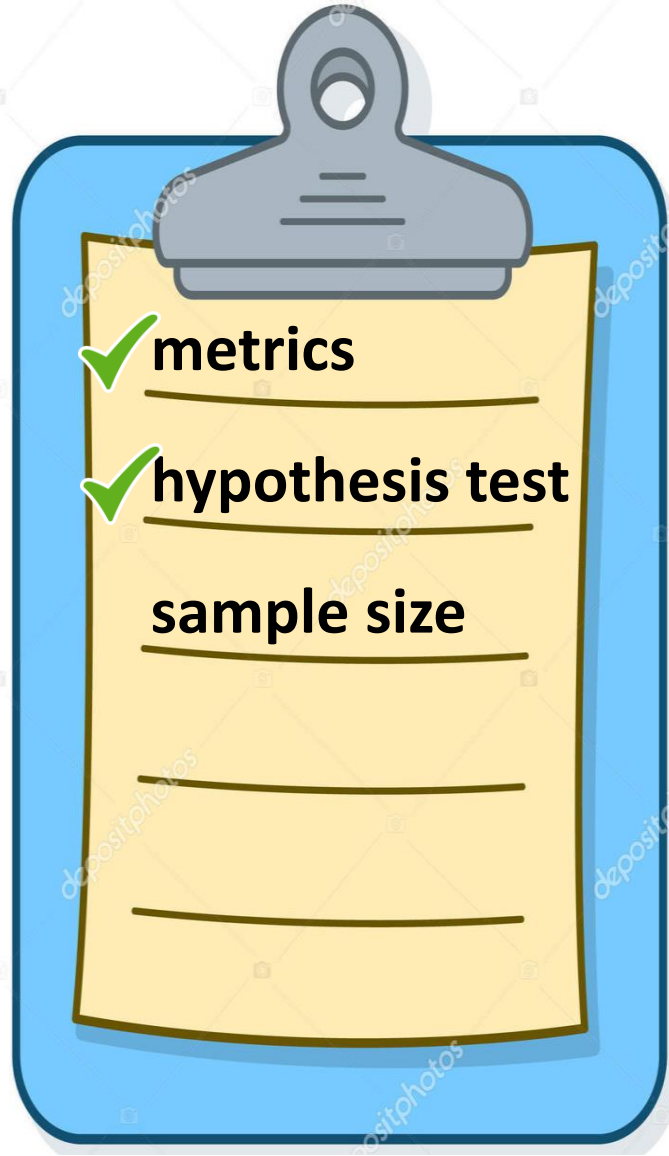
Alternative hypothesis (H_a): $A \neq B$

α : significant level

p-value $< \alpha$: reject H_0



A/B testing – How long to run the test?

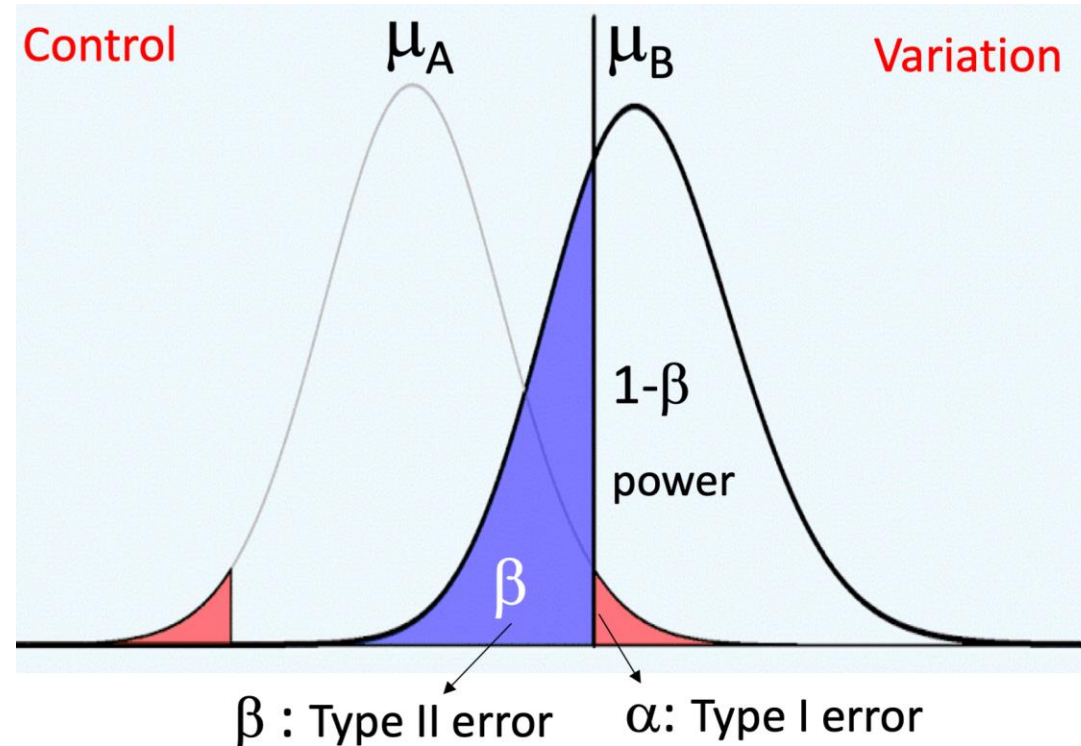


$$\left. \begin{array}{l} \text{effect size} = \left| \frac{\mu_a - \mu_o}{\sigma} \right| \\ \text{significant level} = \alpha \\ \text{power} = 1 - \beta \end{array} \right\}$$

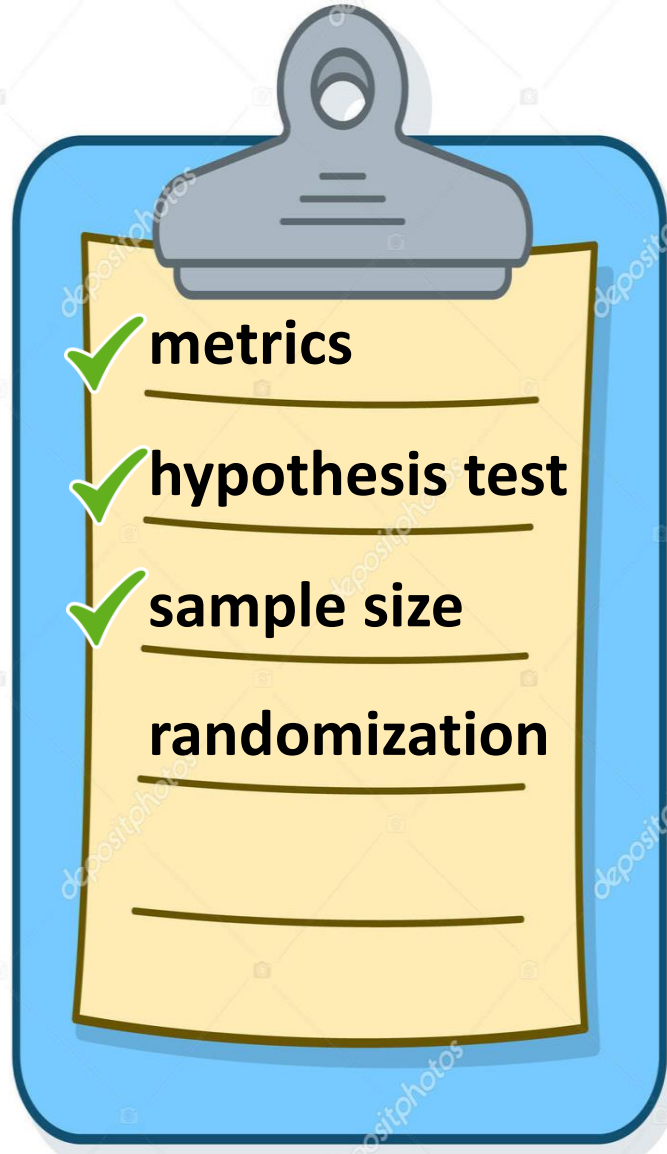
Power
analysis

→ Sample size = n

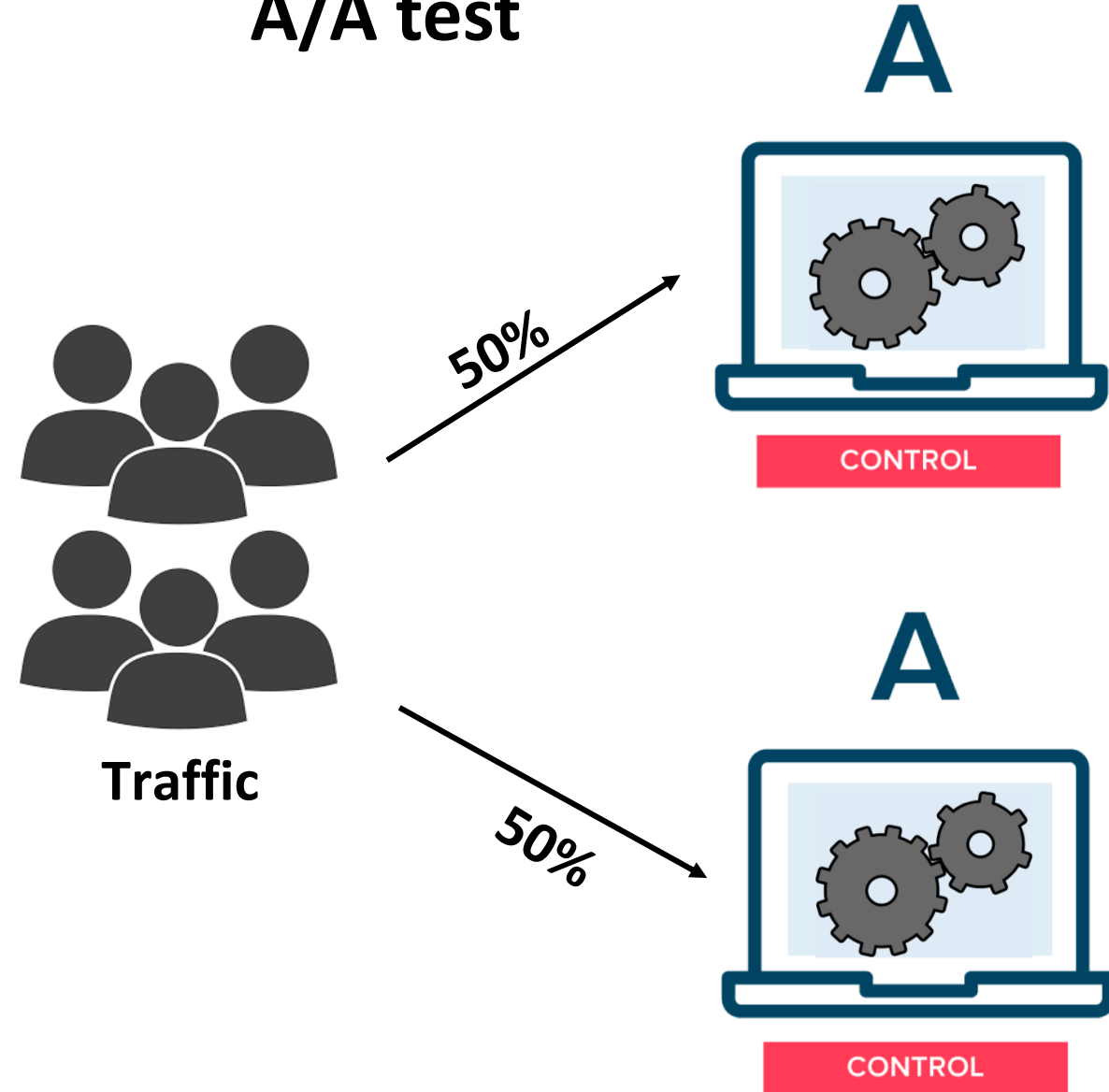
↓
Duration of test



A/B testing – Is randomization guaranteed?

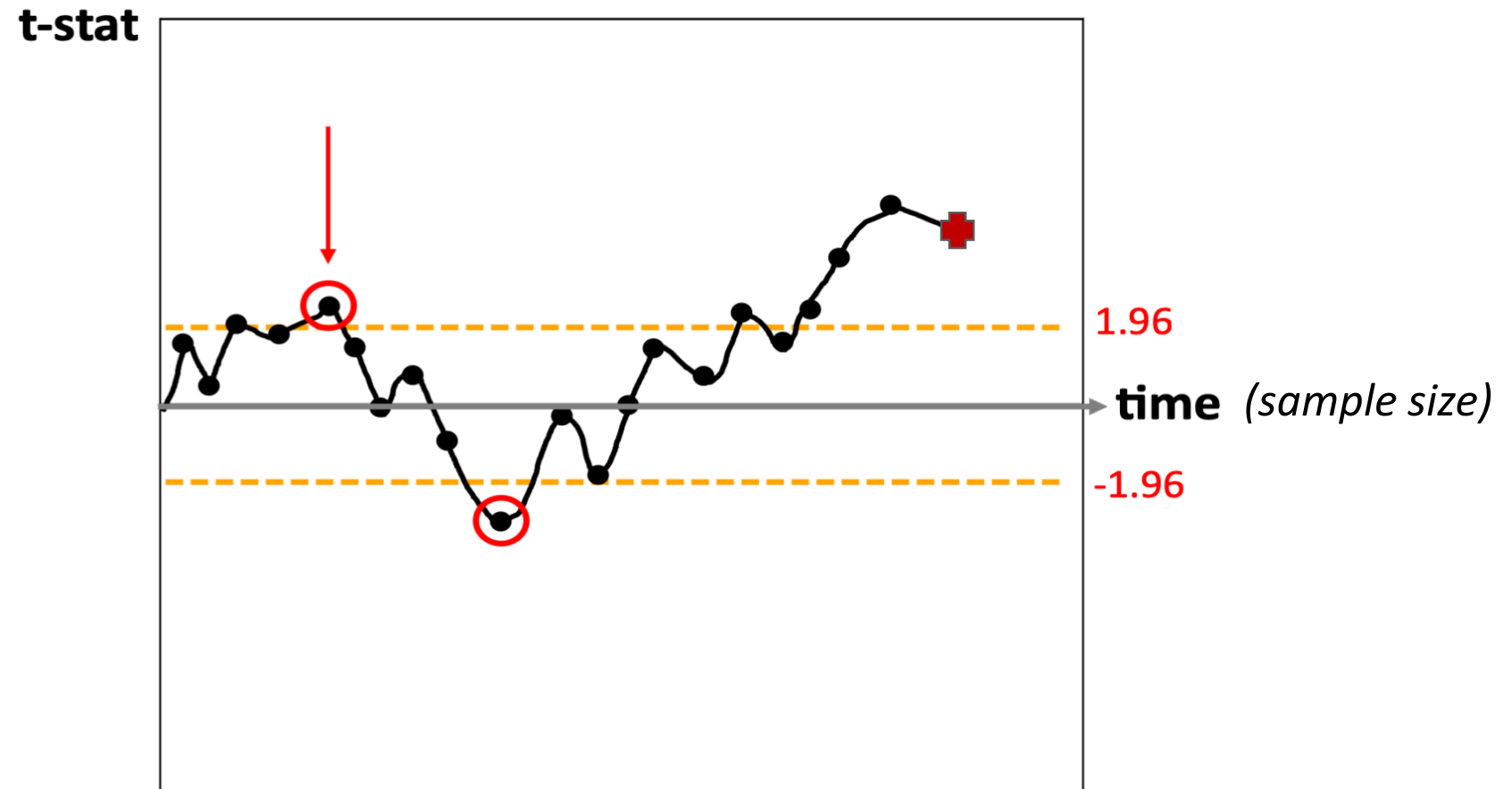
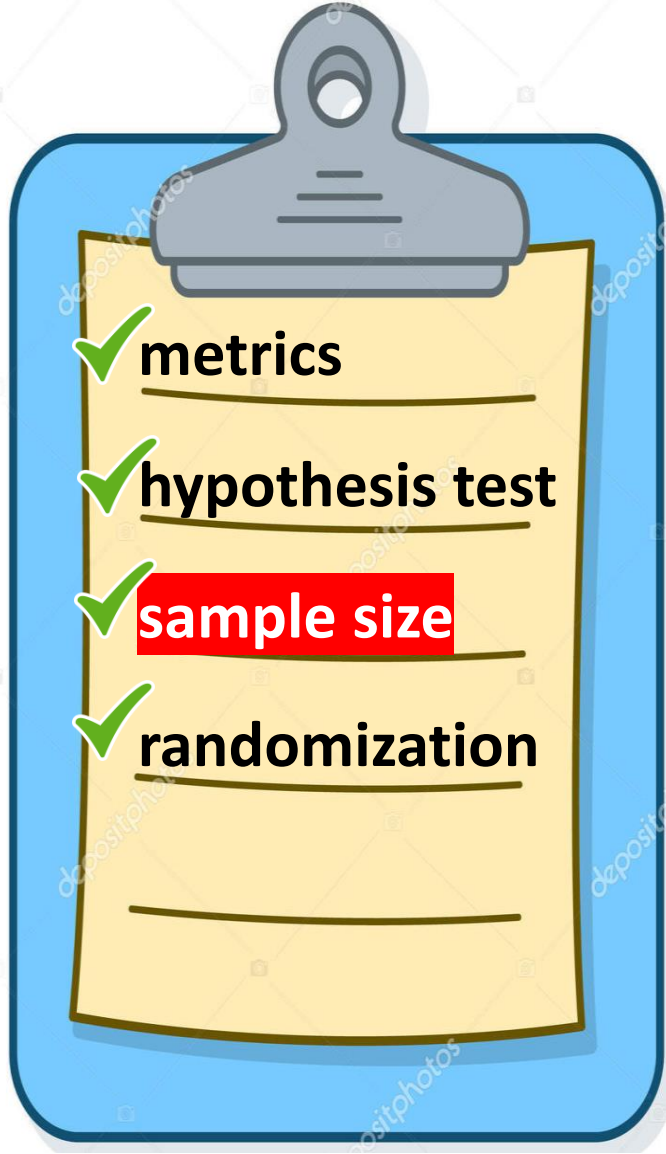


A/A test



Peeking at A/B Tests and early stopping – Why it matters?

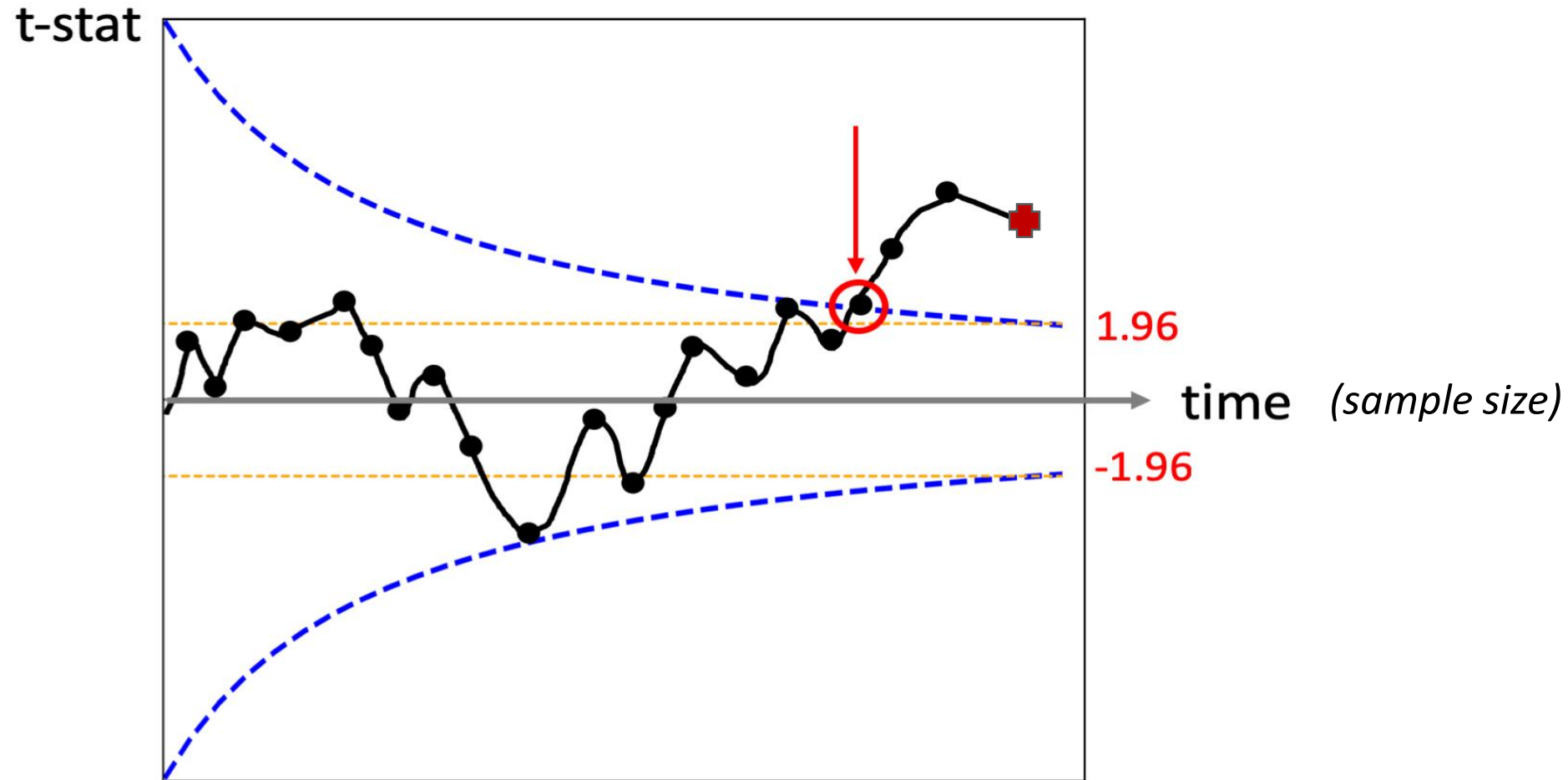
- Impatient experimenters: peeking at the test and early stopping



- But time is money!

Peeking at A/B Tests and early stopping – What to do about it?

Sequential A/B testing (frequentist): Sequential probability ratio test



Peeking at A/B Tests and early stopping – What to do about it?

The slightly better model

Scenario: version B (CTR = 0.142%) is slightly better than version A (CTR = 0.14%),
p-value = 0.11 > significant level α (0.05)

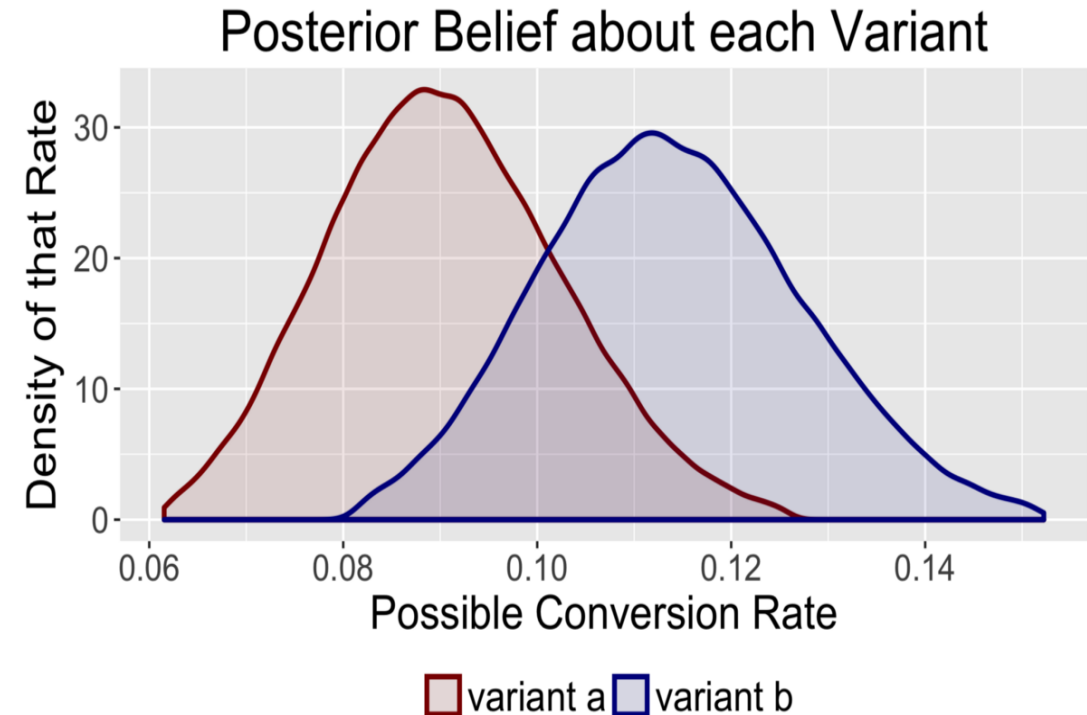
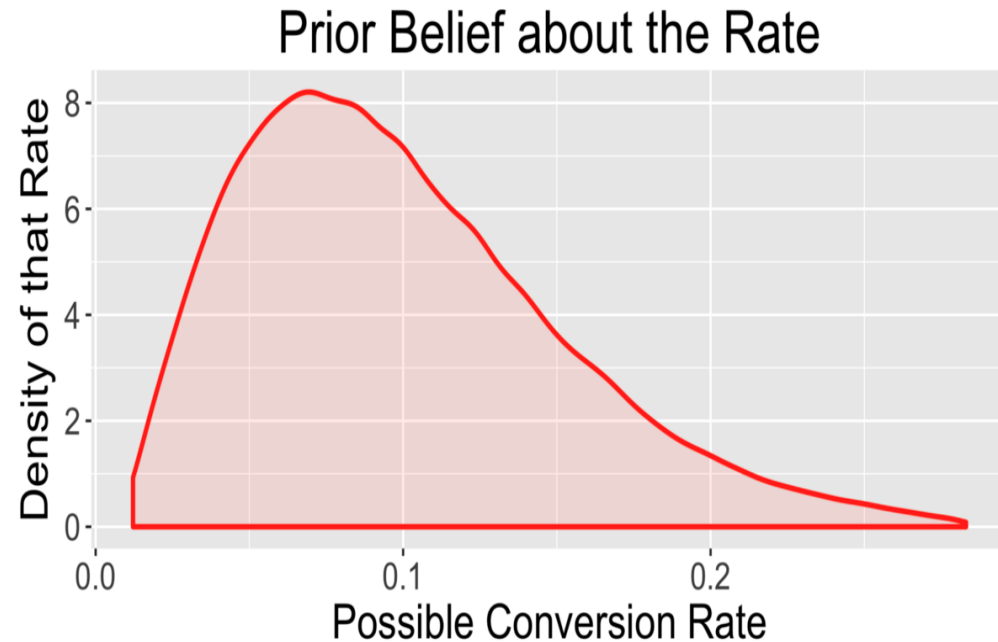
--> Keep A or use B?

Bayesian A/B testing

Given some prior knowledge, what is the probability that the metric under variant B is larger than the metric under variant A?

Peeking at A/B Tests and early stopping – What to do about it?

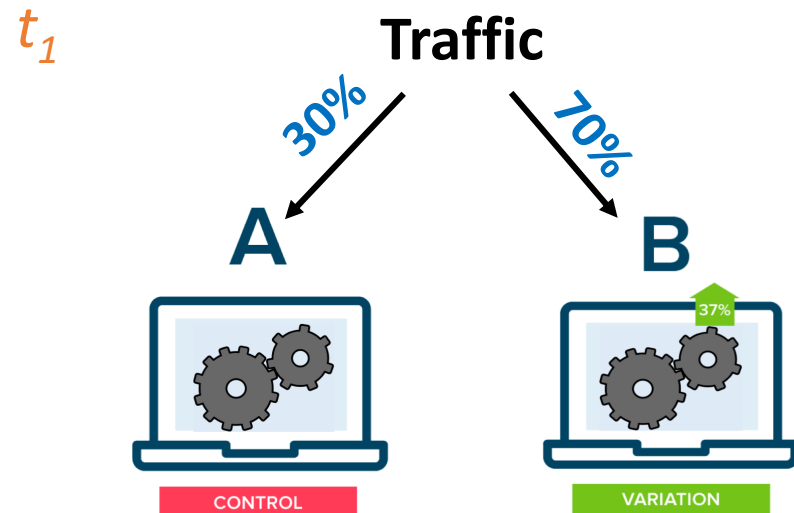
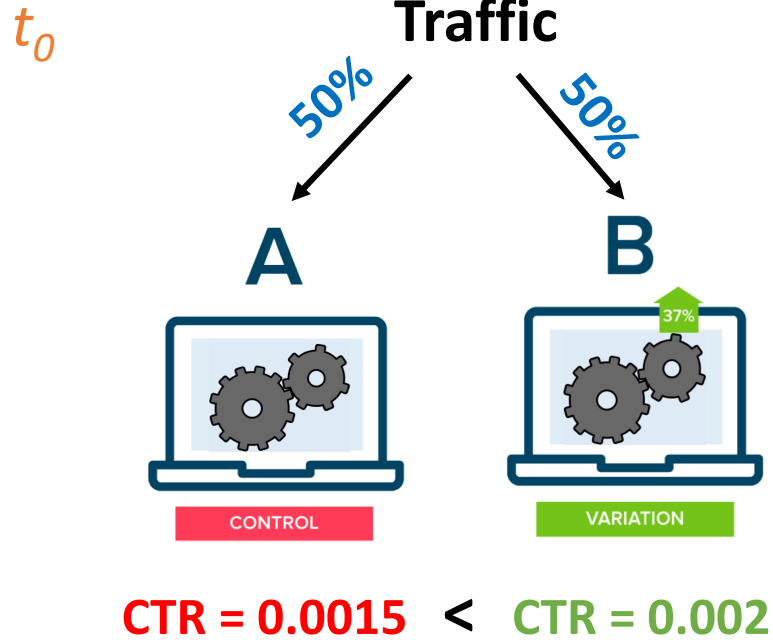
Bayesian A/B testing



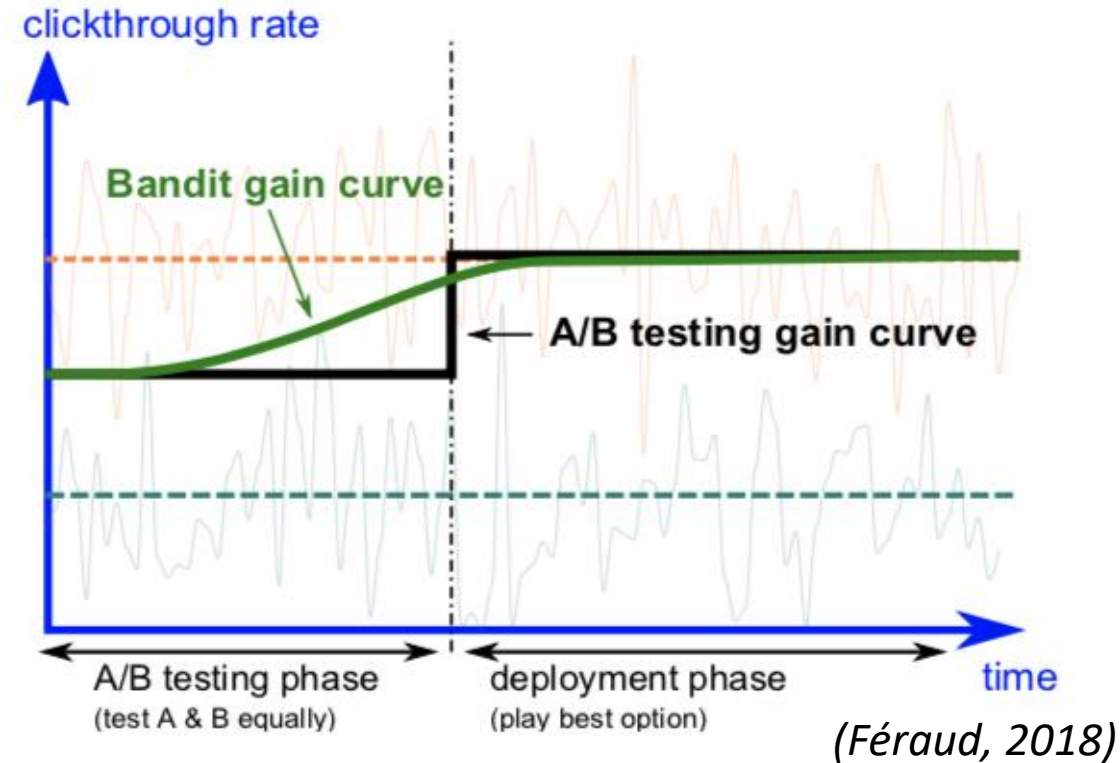
Source: [Medium](#)

- Initiate the experiment and collect data
- Choose a loss threshold **epsilon**
- Periodically, compute expected loss $E[L](A)$ or $E[L](B)$ if we choose version A or B
- Once either $E[L](A)$ or $E[L](B) < \text{epsilon}$: stop the experiment and record the winning variation

Multi-armed bandit test



- Earn while you learn



- Give the poorer performing versions a second chance
- Automation
- >2 arms
- More tricky to implement

- Eg., Implement a new click predictor to optimize for cpc → Metric: ctr, cpc
 - It might take long time to observe a significant difference in the new model
 - Heterogeneity → assumption of 'iid' fails
 - Testing at ad-group level: data size is even smaller
 - Time is money and we don't want to compromise user's experience
- A/B test might not an optimal method
- Multi-armed bandit: the more appropriate approach

- Online evaluation of new model/new idea is important but challenging, especially in RTB
- A/B testing is the most widely used method, though early stopping without a plan is dangerous → sequential A/B test or Bayesian A/B test
- A/B testing: strict experiments where focus is on statistical significance
- Multi-armed bandit (MAB): continuous optimization where focus is on maintaining optimal KPIs

References

Bayesian A/B test:

<https://medium.com/convoy-tech/the-power-of-bayesian-a-b-testing-f859d2219d5>

[https://cdn2.hubspot.net/hubfs/310840/VWO SmartStats technical whitepaper.pdf](https://cdn2.hubspot.net/hubfs/310840/VWO_SmartStats_technical_whitepaper.pdf)

Peeking at A/B test:

<https://www.youtube.com/watch?v=AjX4W3MwKzU>

<https://www.evanmiller.org/how-not-to-run-an-ab-test.html>

Sequential A/B test:

<https://www.evanmiller.org/sequential-ab-testing.html>

<https://www.aarondefazio.com/tangentially/?p=83>

[https://en.wikipedia.org/wiki/Sequential probability ratio test](https://en.wikipedia.org/wiki/Sequential_probability_ratio_test)

Backup slides

Peeking at A/B Tests and early stopping – What to do about it?

Sequential A/B testing (frequentist): Sequential probability ratio test

$$H_o: \theta_o$$

$$H_a: \theta_a$$

$$S_o = 0 \rightarrow S_i = S_{i-1} + \log \Lambda_i \quad \text{for } i \text{ in } 1, 2, \dots, t$$

$$\Lambda_i = \frac{L(\theta_o|x)}{L(\theta_a|x)} : \text{Likelihood ratio}$$

Stopping rule:

$$a < S_i < b$$

$$S_i \geq b: \text{Accept } H_a$$

$$S_i \leq a: \text{Accept } H_o$$

$$a \sim \log\left(\frac{\beta}{1-\alpha}\right) \quad b \sim \log\left(\frac{1-\beta}{\alpha}\right)$$

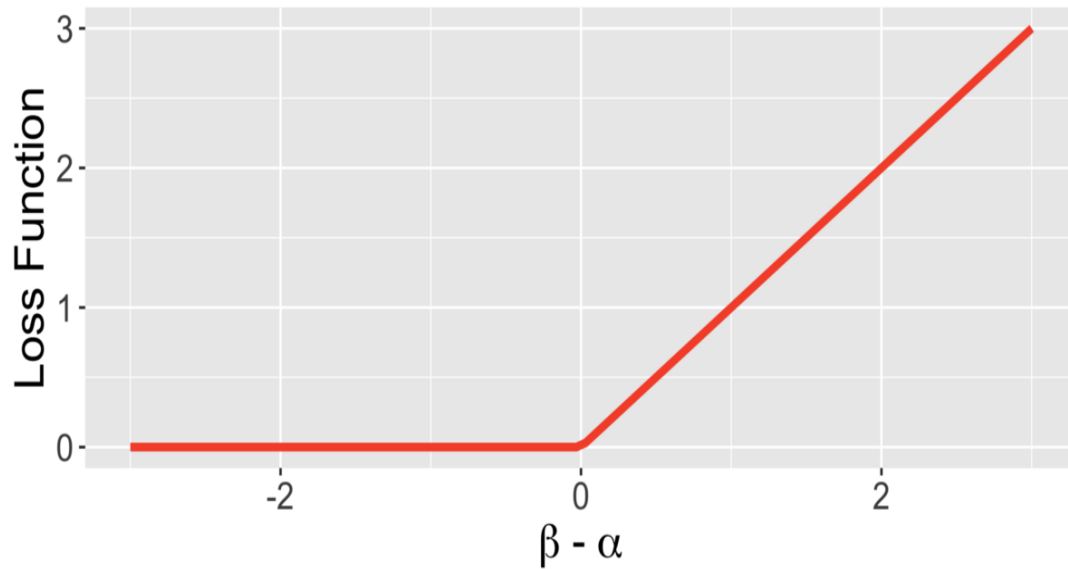
The slightly better model

Scenario: version B (CTR = 0.142%) is slightly better than version A (CTR = 0.14%), p-value = 0.11 > significant level α (0.05)

--> Keep A or use B?

Bayesian A/B testing

Loss Function when Choosing Variant A



- Initiate the experiment and collect data
- Periodically, compute $E[L](A)$ and $E[L](B)$
- If either $E[L](A)$ or $E[L](B) < \text{epsilon}$: stop the experiment and record the winning variation
- Else: continue the experiment

$$L(\alpha, \beta, x) = \begin{cases} \max(\beta - \alpha, 0) & x = a \\ \max(\alpha - \beta, 0) & x = b \end{cases}$$

$$E[L](x) = \int_A \int_B L(\alpha, \beta, x) \underbrace{f(\alpha, \beta)}_{\text{joint posterior density}} d\alpha d\beta$$