

A/B Testing for Game Design Iteration: A Bayesian Approach

Steve Collins
CTO / Swrve

Steve@swrve.com
@stevec64

Biography



...



kore

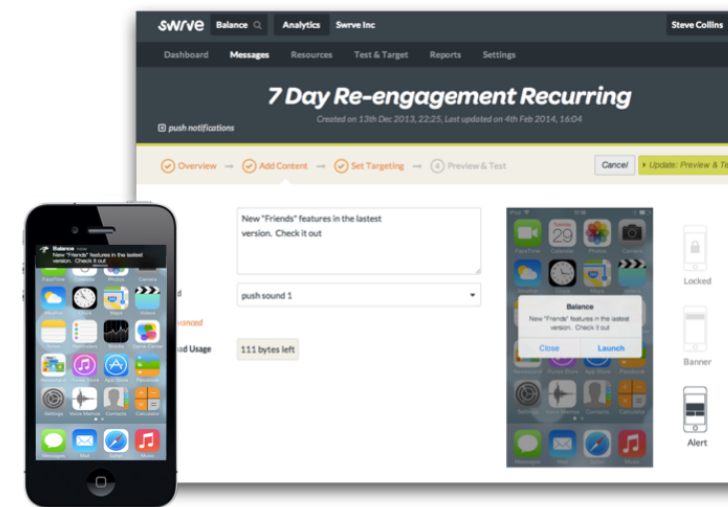




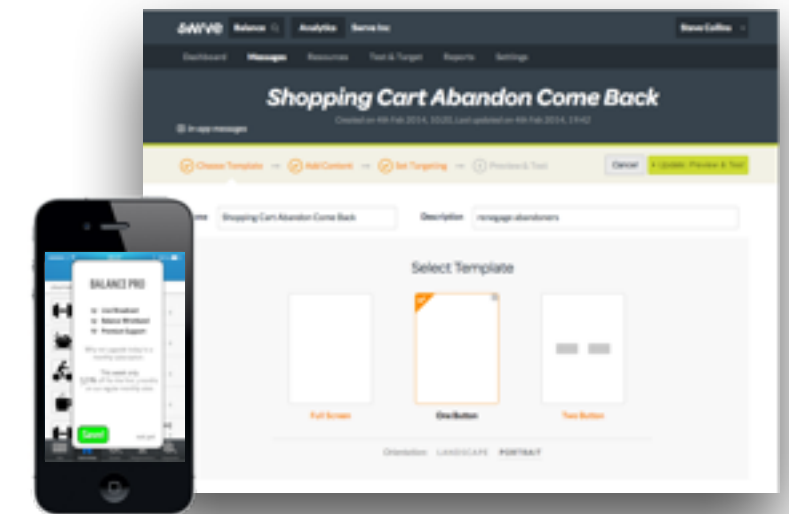
Targeting & Analytics



A/B Testing



Push Notifications

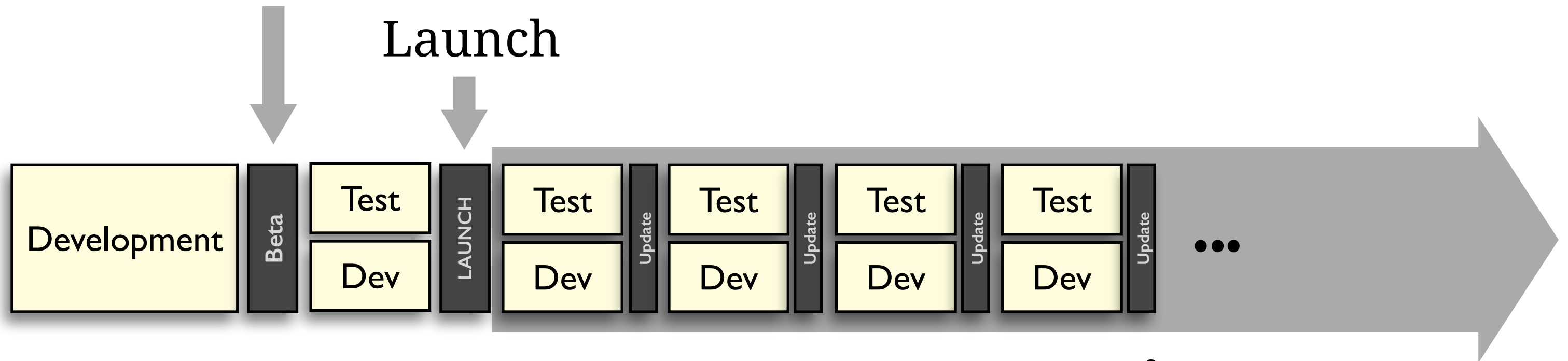


In-App Messaging

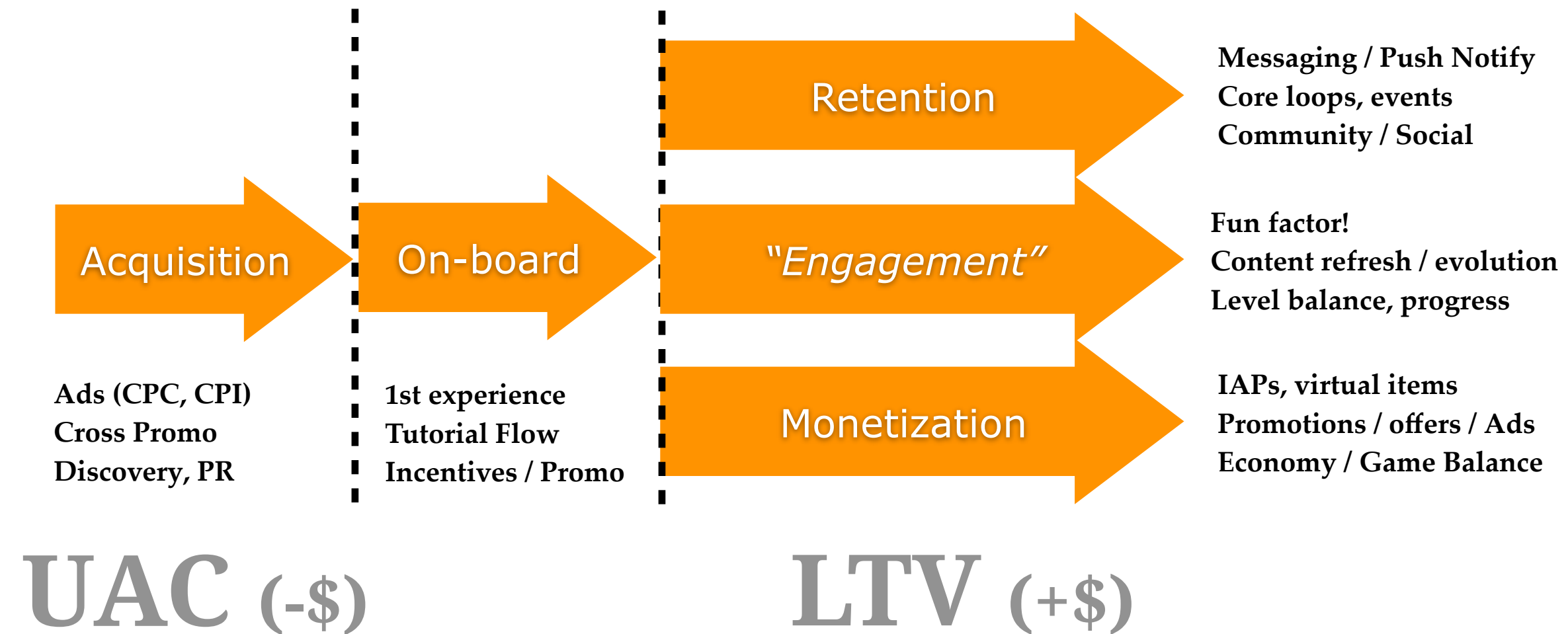
Introduction to A/B Testing

“Soft”
Launch

Launch




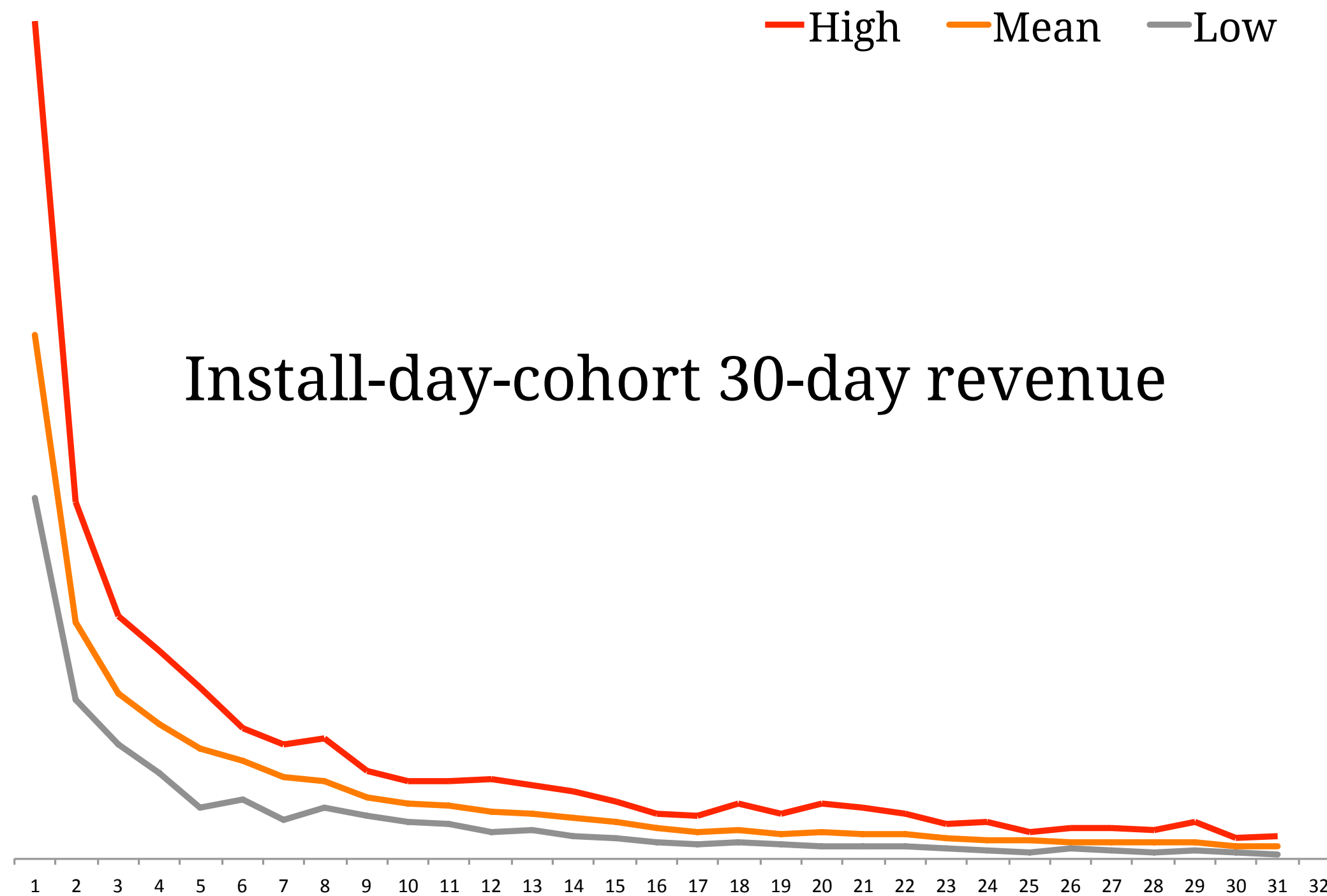
Game Service



$$\text{ROI} = \text{LTV} - \text{UAC}$$

$$\text{ROI} = \sum_{\text{d}=1}^{\text{lifetime?}} \text{ARPU}_{\text{d}} - \text{UAC}$$


LTV



Understand

Metrics -> Analytics -> Insight



Test Hypotheses

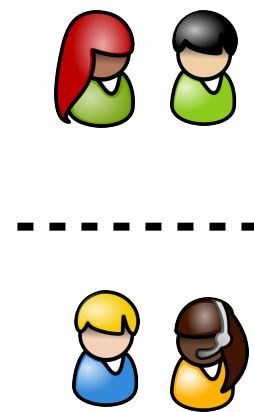
Data driven



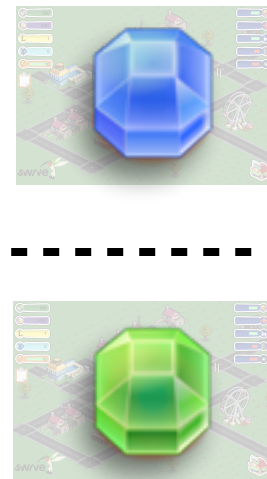
Take action

Iterate, fail-fast

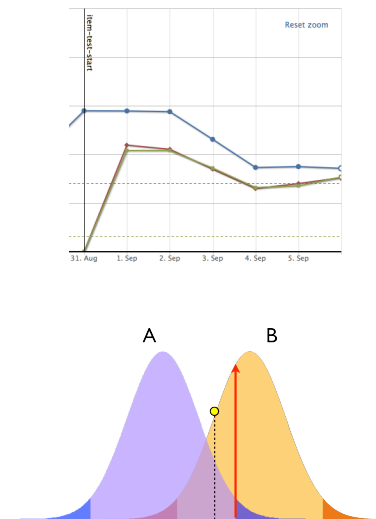




1. Split population



2. Show variations



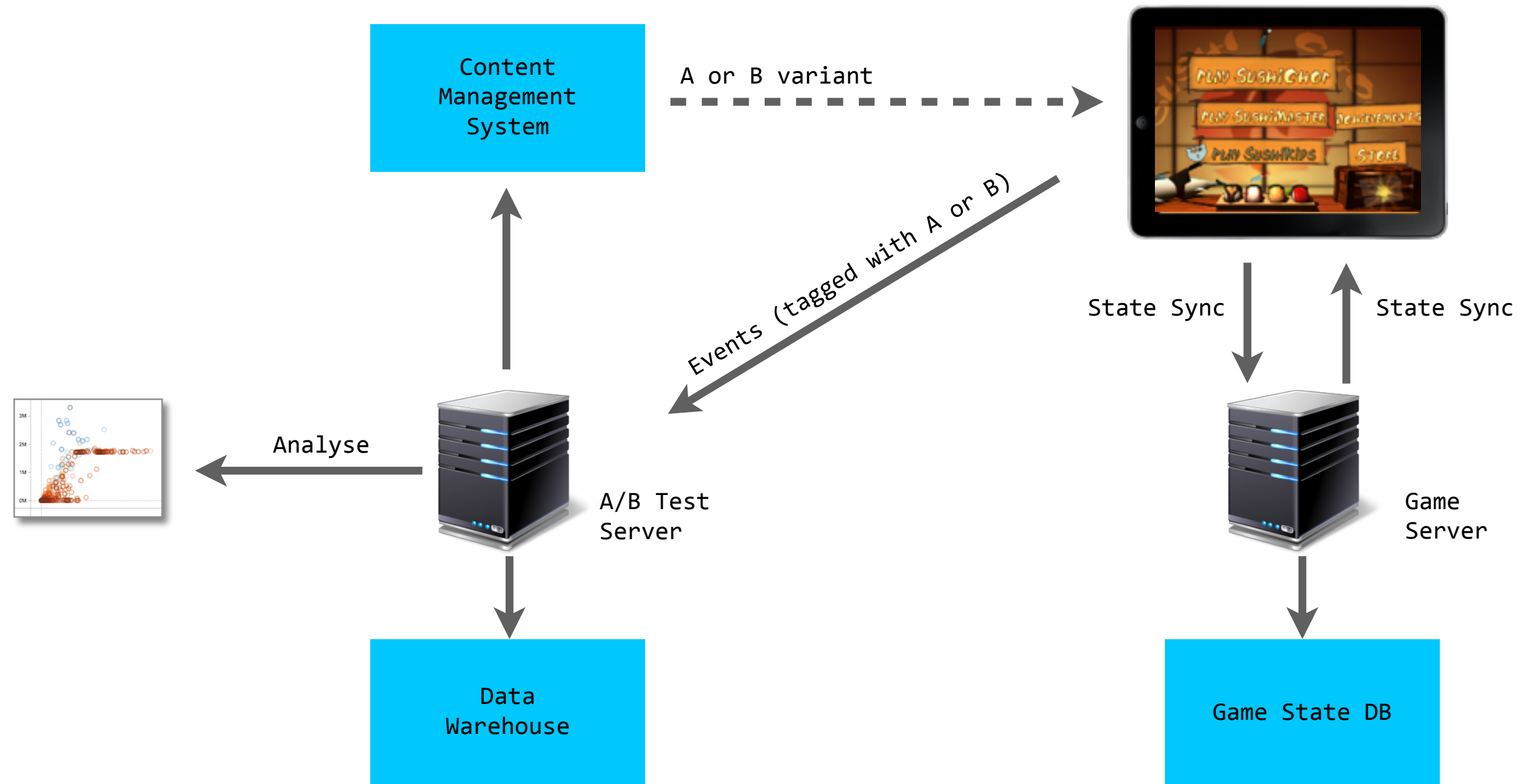
3. Measure response



4. Choose winner



4. Deploy winner



What to test?

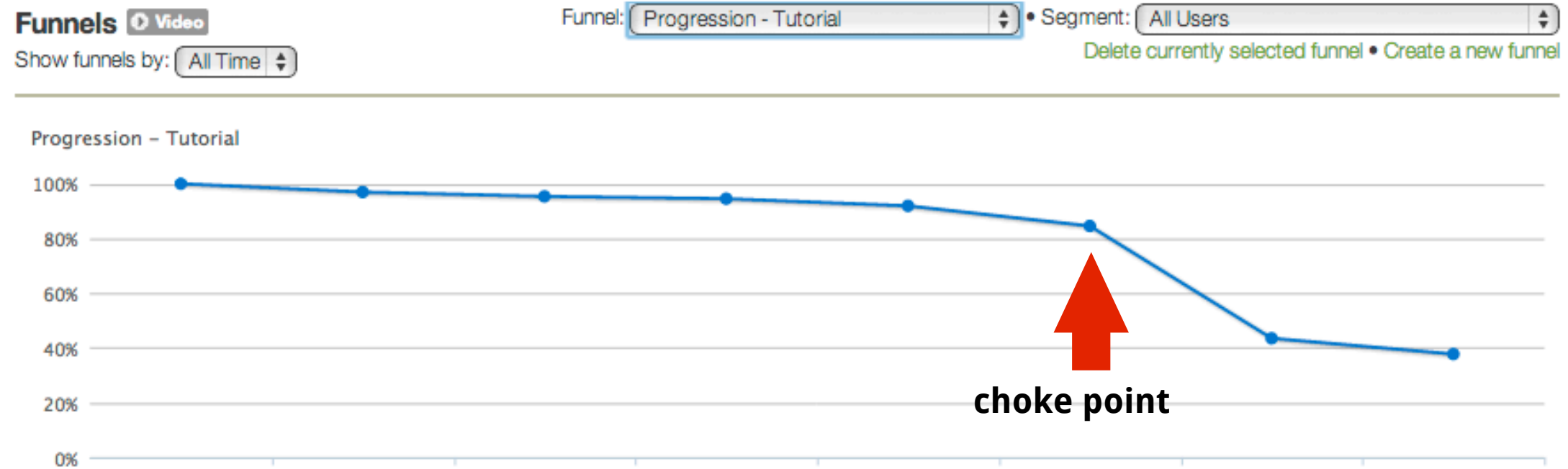
Message layouts / content



VS



Tutorial Flow



Promotion Discounts



vs



Elasticity testing: exchange rate



\$4.99

A

VS



\$4.99

B

Store Inventory

Price set A

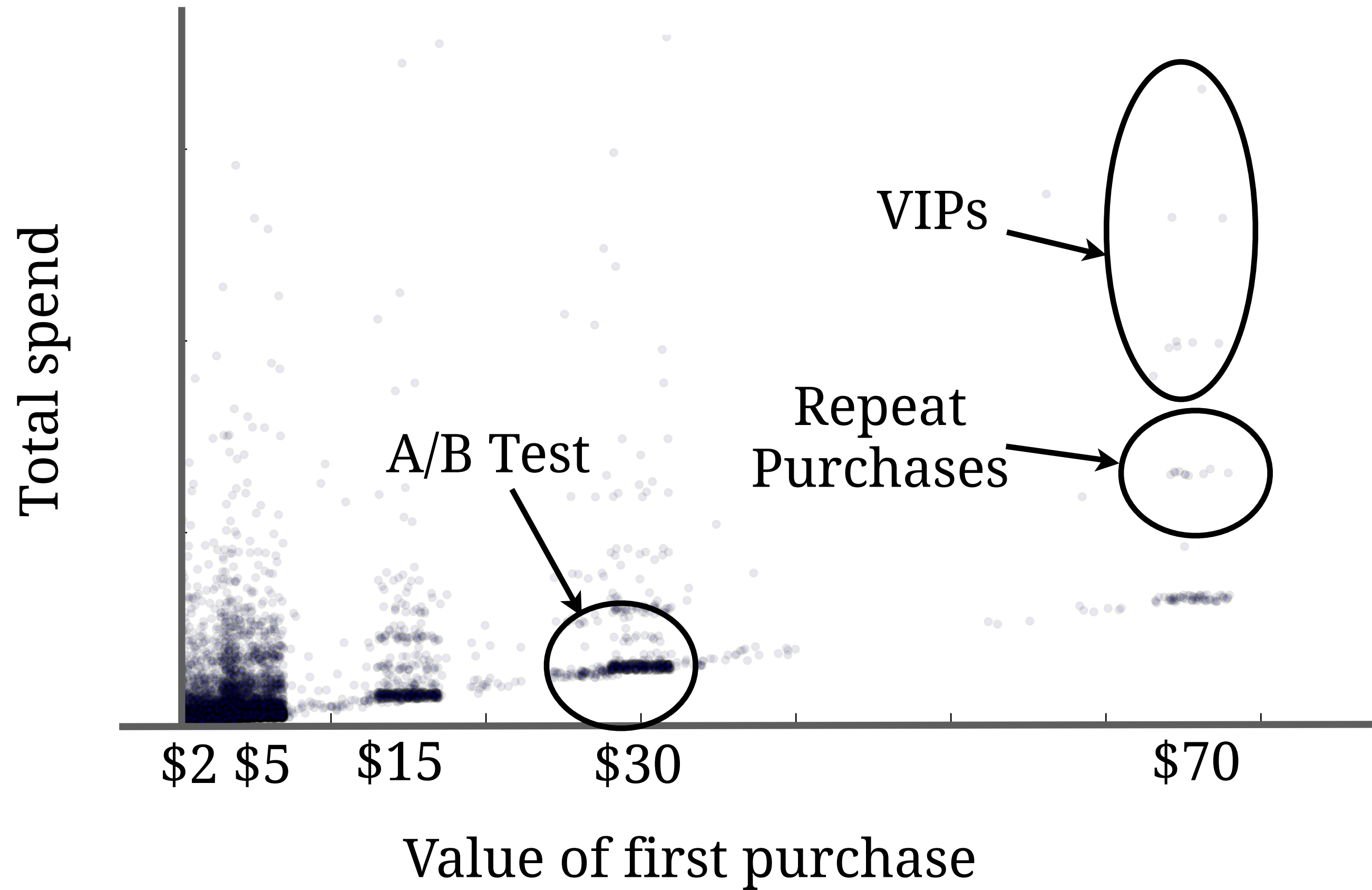


Price set B

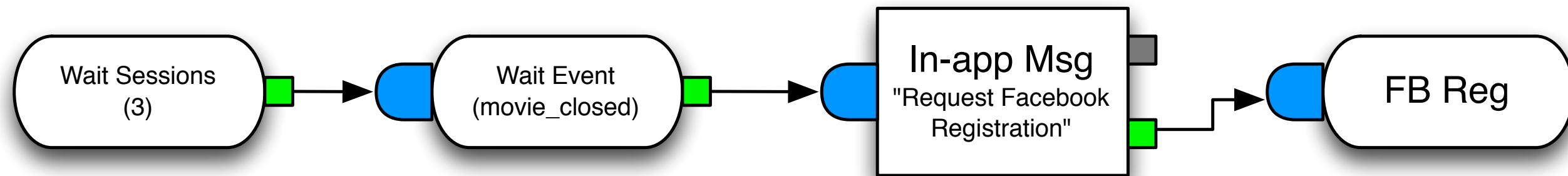


Price set C



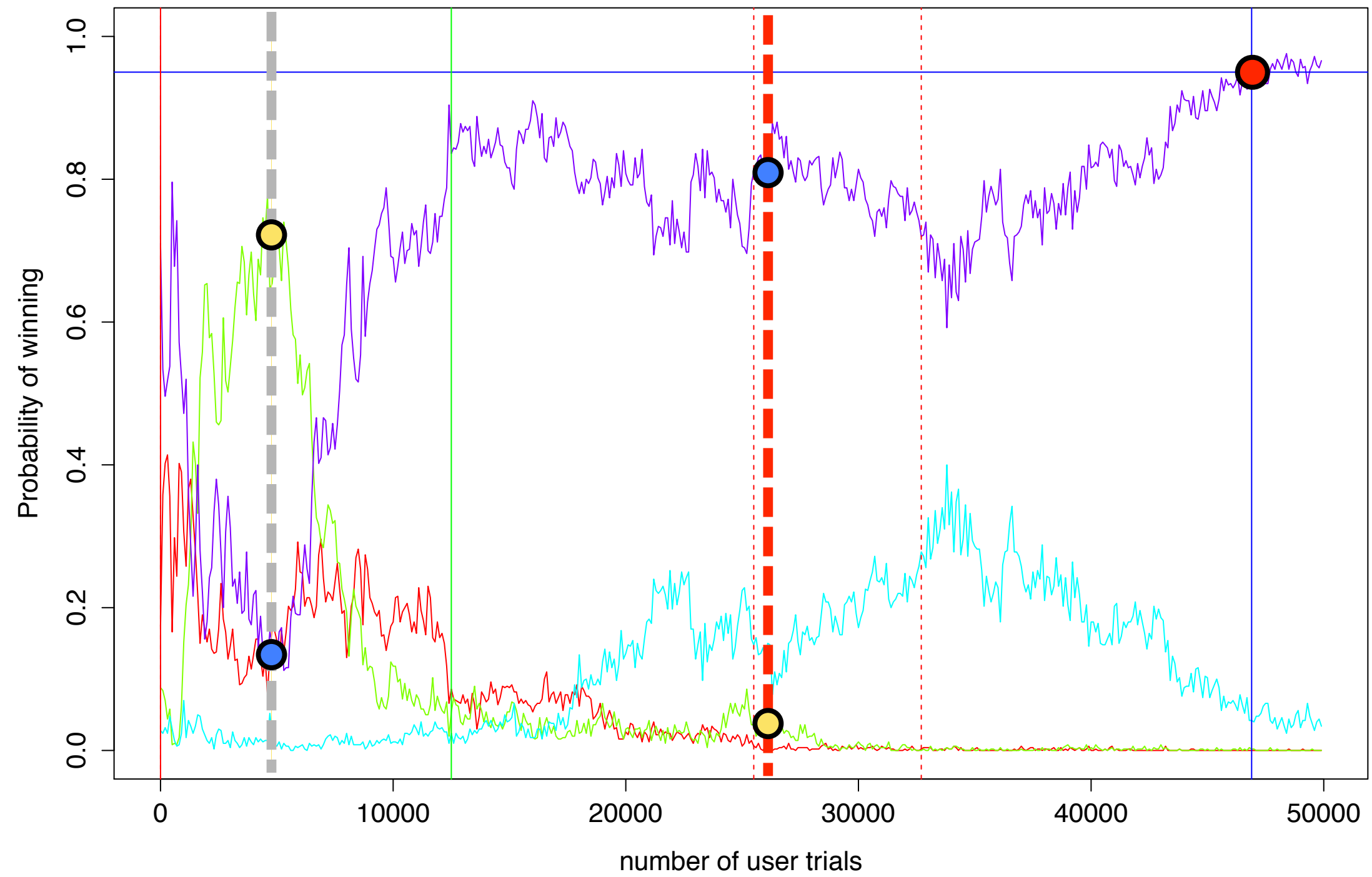


Timing

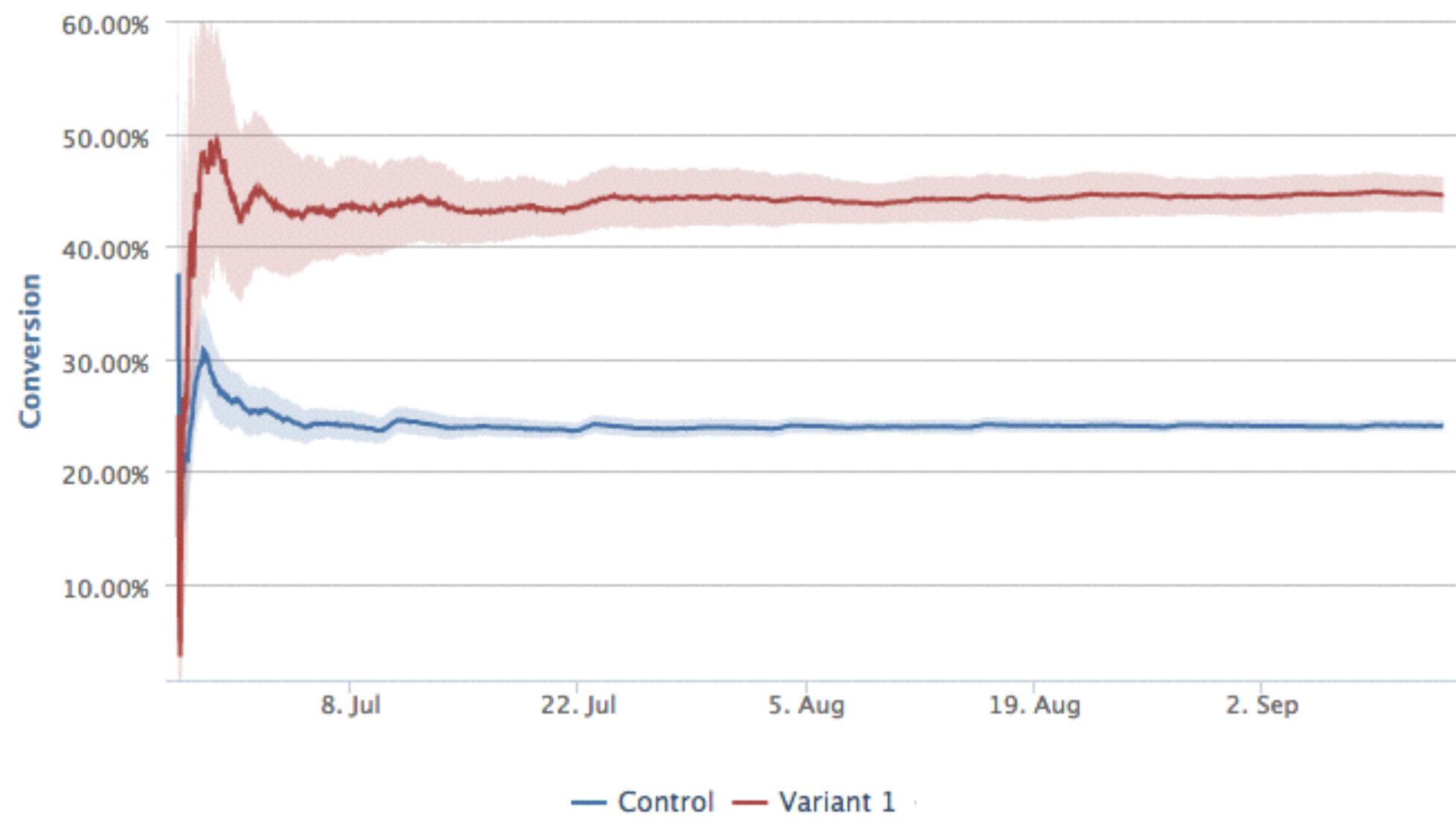


when

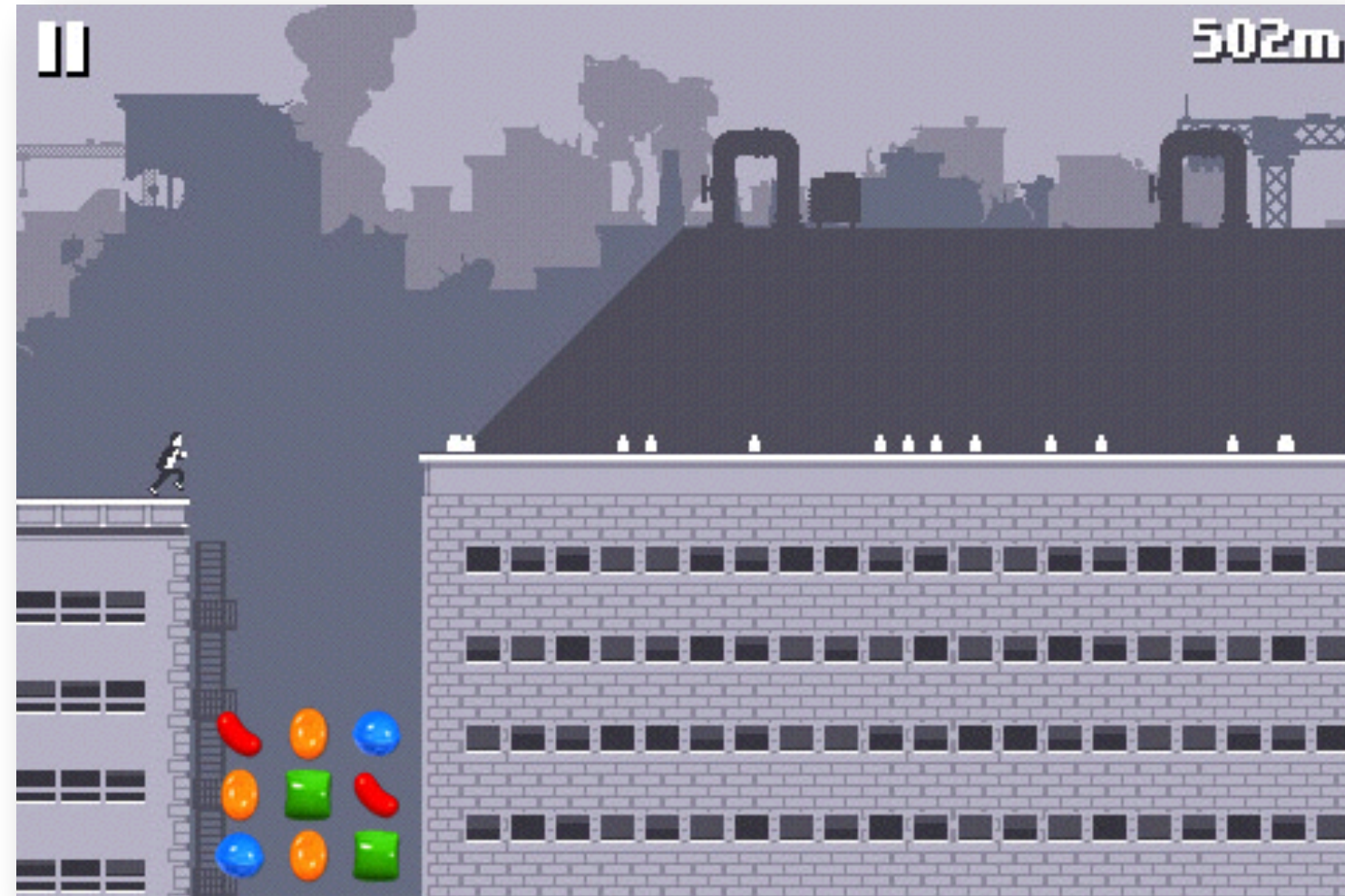
where



Conversion over time

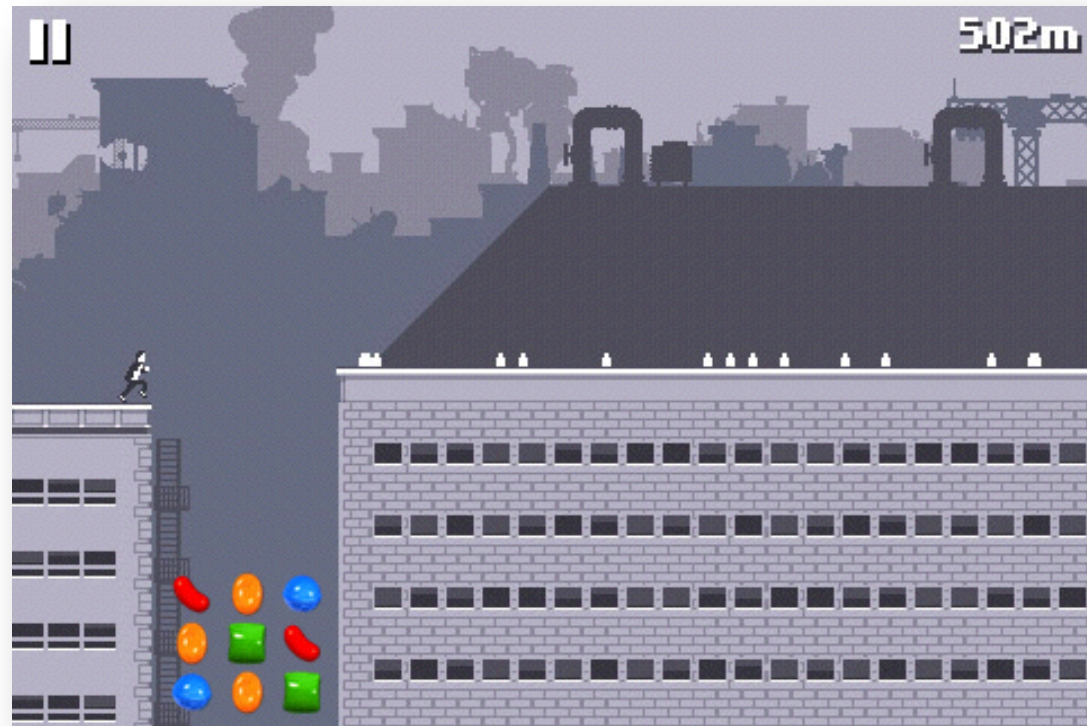


Canacandycrushbalt*

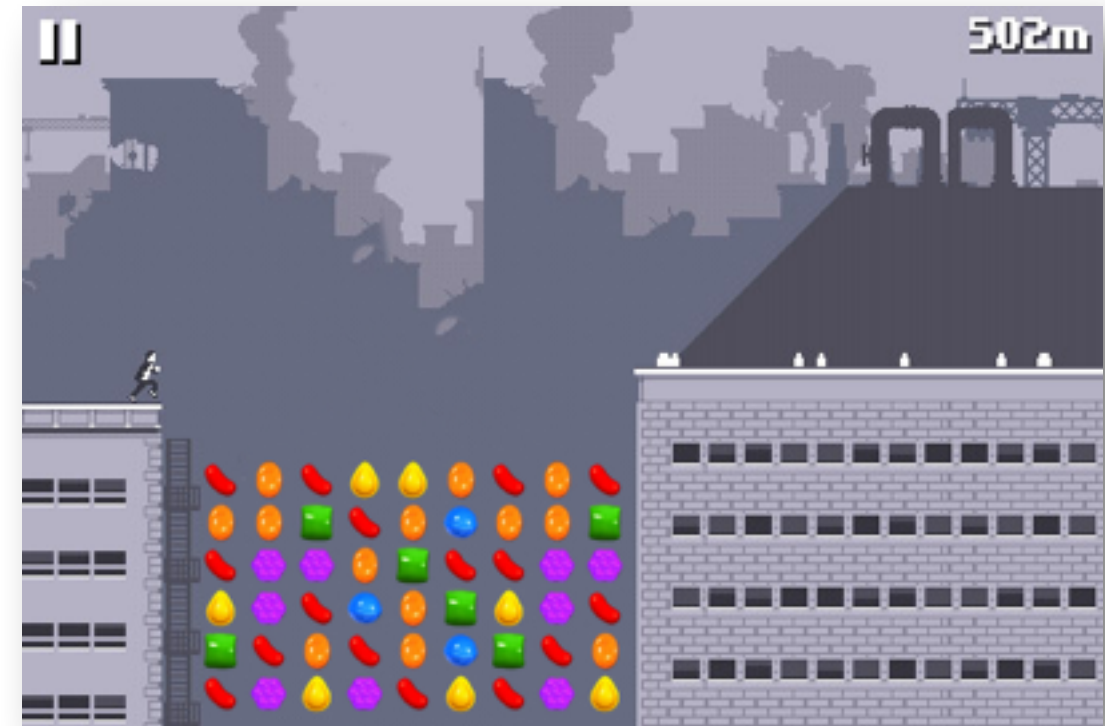


Day 1 retention = 30%

Beta Test



Original



New Version

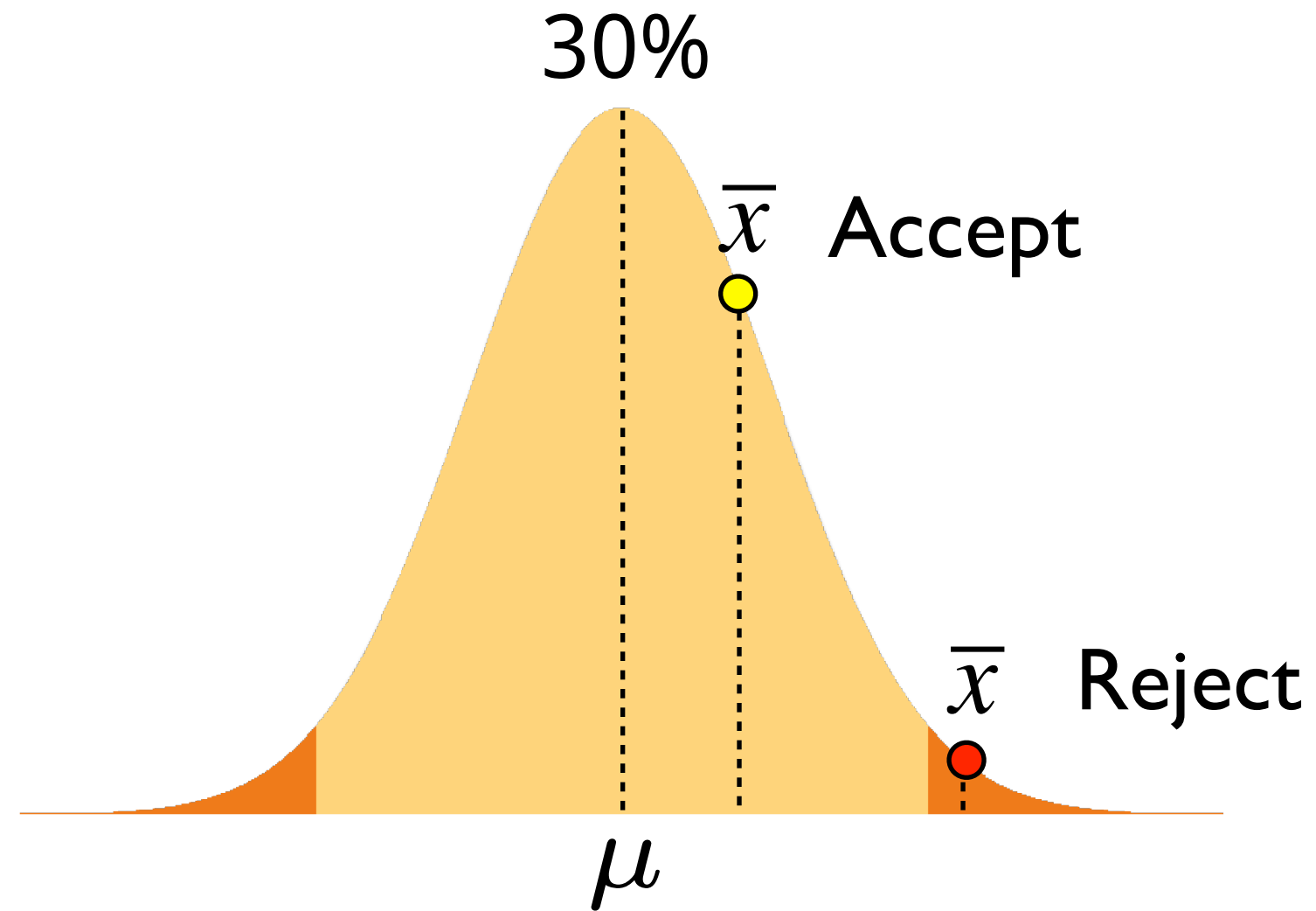
Expecting 30% day-1 retention

After 50 users, we see 0%

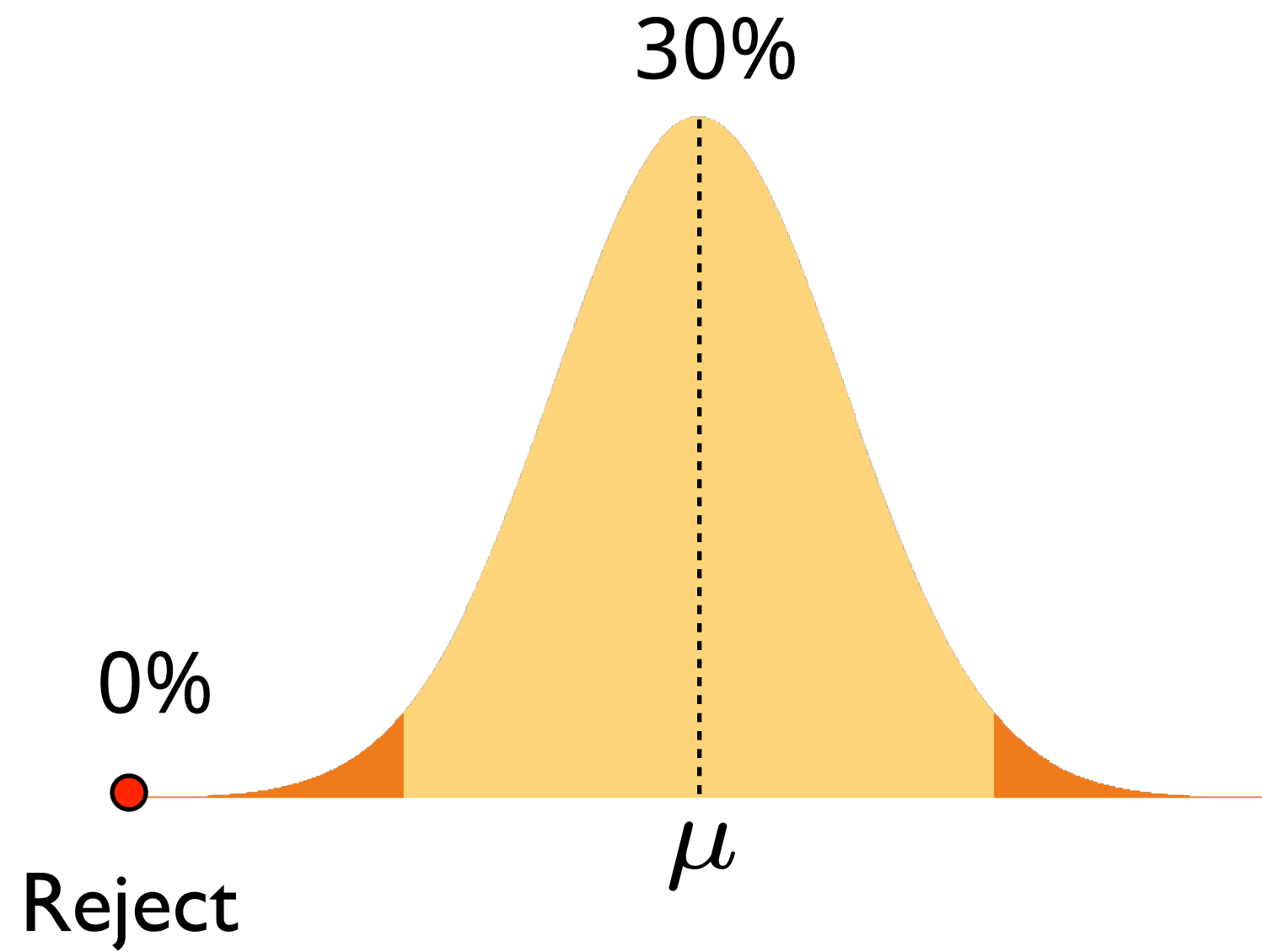
Is this bad?

Null Hypothesis Testing (NHT) View

The Null Hypothesis



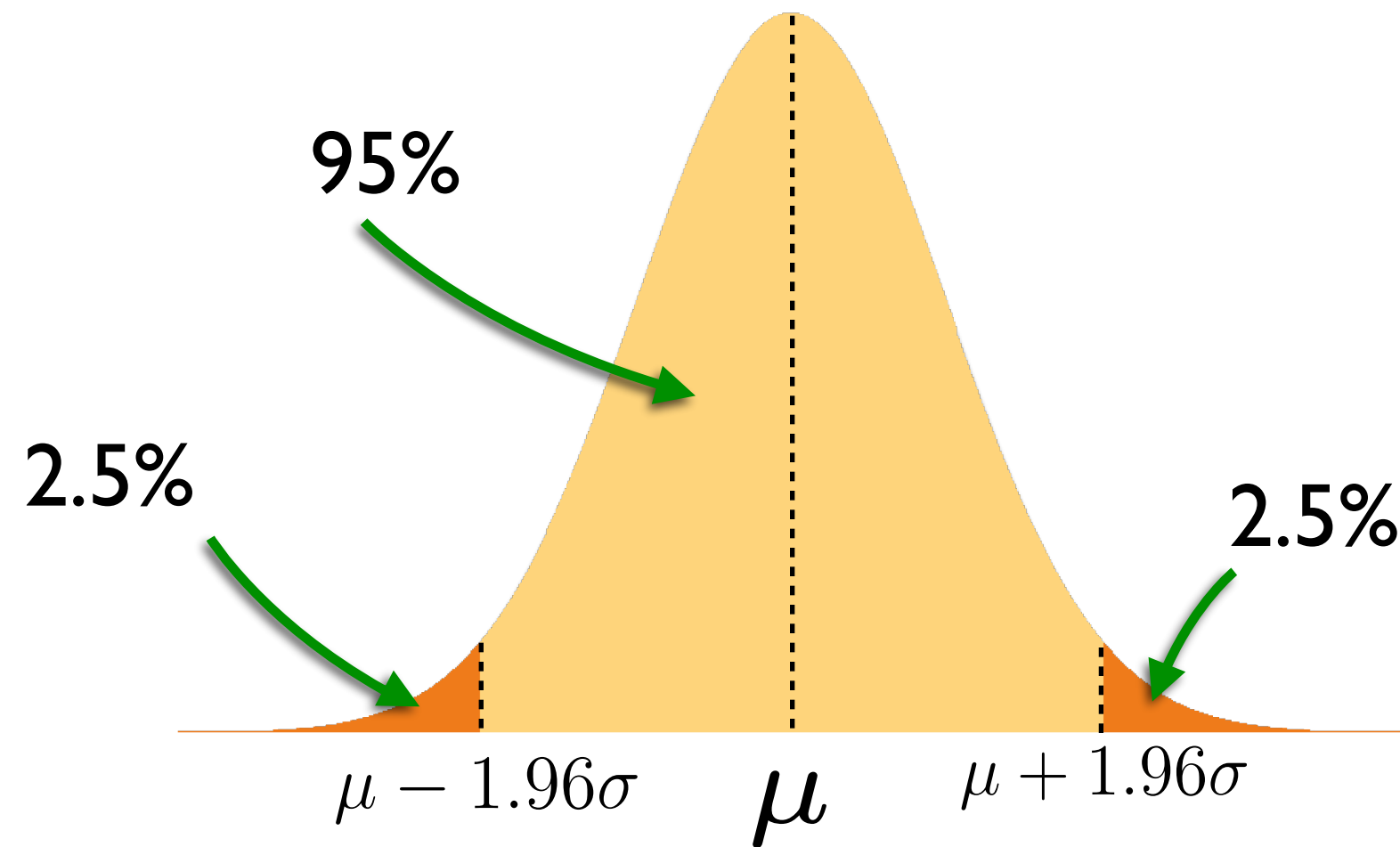
$$H_0 : \bar{x} = \mu$$



*Conclusion: 30% is unlikely to be the
retention rate*

Issue #1
p-value

p-Value



***p*-Value**

*The probability of observing as extreme a result
assuming the null hypothesis is true*

OR

The probability of the data given the model

Null Hypothesis: H_0

Truth

H

H

Observation

Accept H



Type-II Error



Reject H



False Positive



***p*-Value**

$$p < 0.05$$

All we can ever say is either

- *not enough evidence that retention rates are the same*
- *the retention rates are different, 95% of the time*

actually...

$$p < 0.05$$

*The evidence supports a rejection of the null hypothesis, i.e.
the probability of seeing a result as extreme as this, assuming
the retention rate is actually 30%, is less than 5%.*

Issue #2

“Peeking”

Number of participants per group

Standard deviation

$$n \approx \frac{16\sigma^2}{\delta^2}$$

Required change

The diagram illustrates the formula for determining the number of participants per group (n) based on the standard deviation (sigma) and the required change (delta). The formula is $n \approx \frac{16\sigma^2}{\delta^2}$. A green arrow points from the text 'Number of participants per group' to the variable 'n'. Another green arrow points from the text 'Standard deviation' to the variable 'sigma'. A third green arrow points from the text 'Required change' to the variable 'delta'.

To get 5% false positive rate you need...

Peeks	5% Equivalent
1	2.9%
2	2.2%
3	1.8%
5	1.4%
10	1%

i.e. 5 times

Issue #3

Family-wise Error

Family-wise Error

$$p = P(\text{Type-I Error}) = 0.05$$

$$P(\text{no Type-I Error}) = 0.95$$

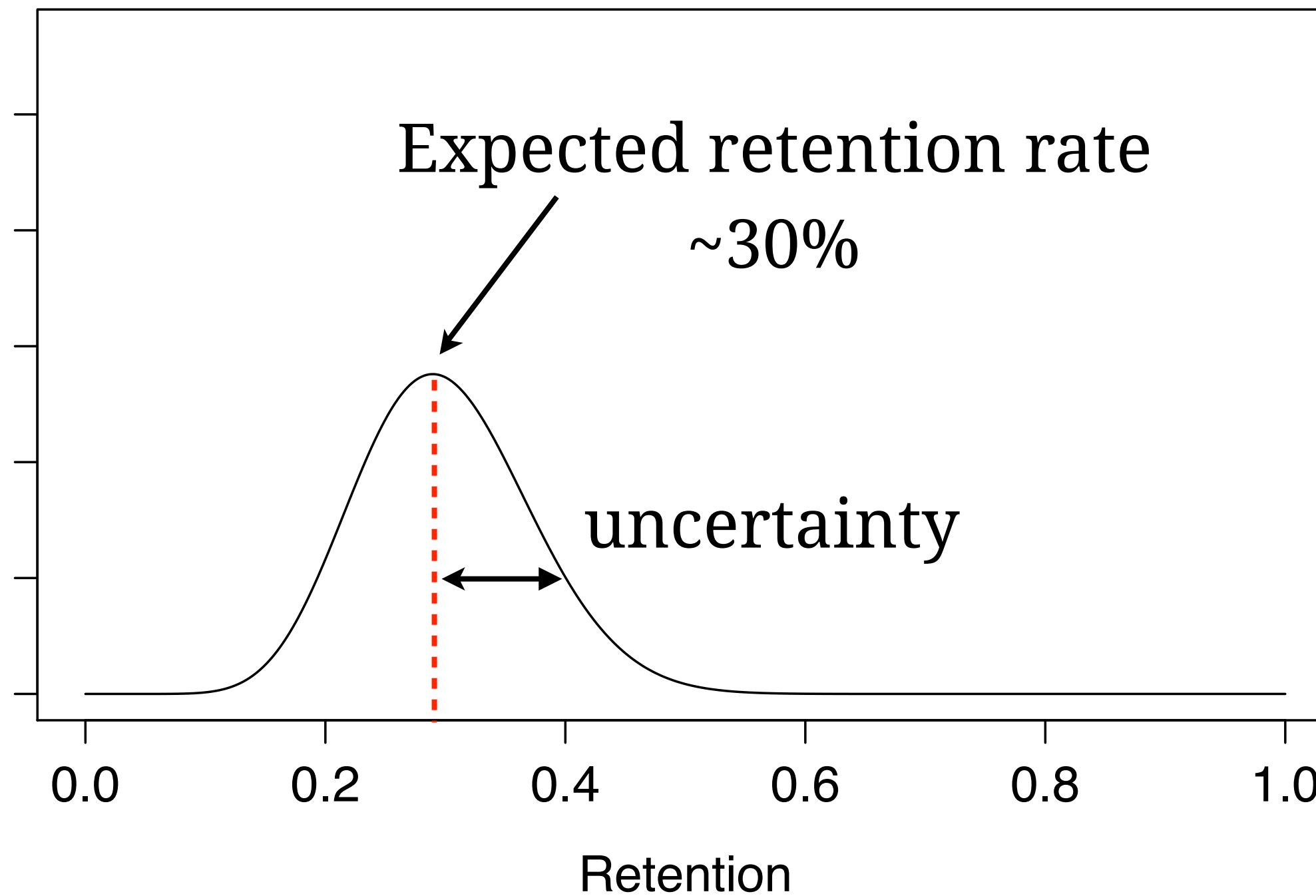
5% of the time we will get a false positive - for **one** treatment

$$P(\text{no Type-I Error for 2 treatments}) = (0.95)(0.95) = 0.9025$$

$$P(\text{at least 1 Type-I Error for 2 treatments}) = (1 - 0.9025) = 0.0975$$

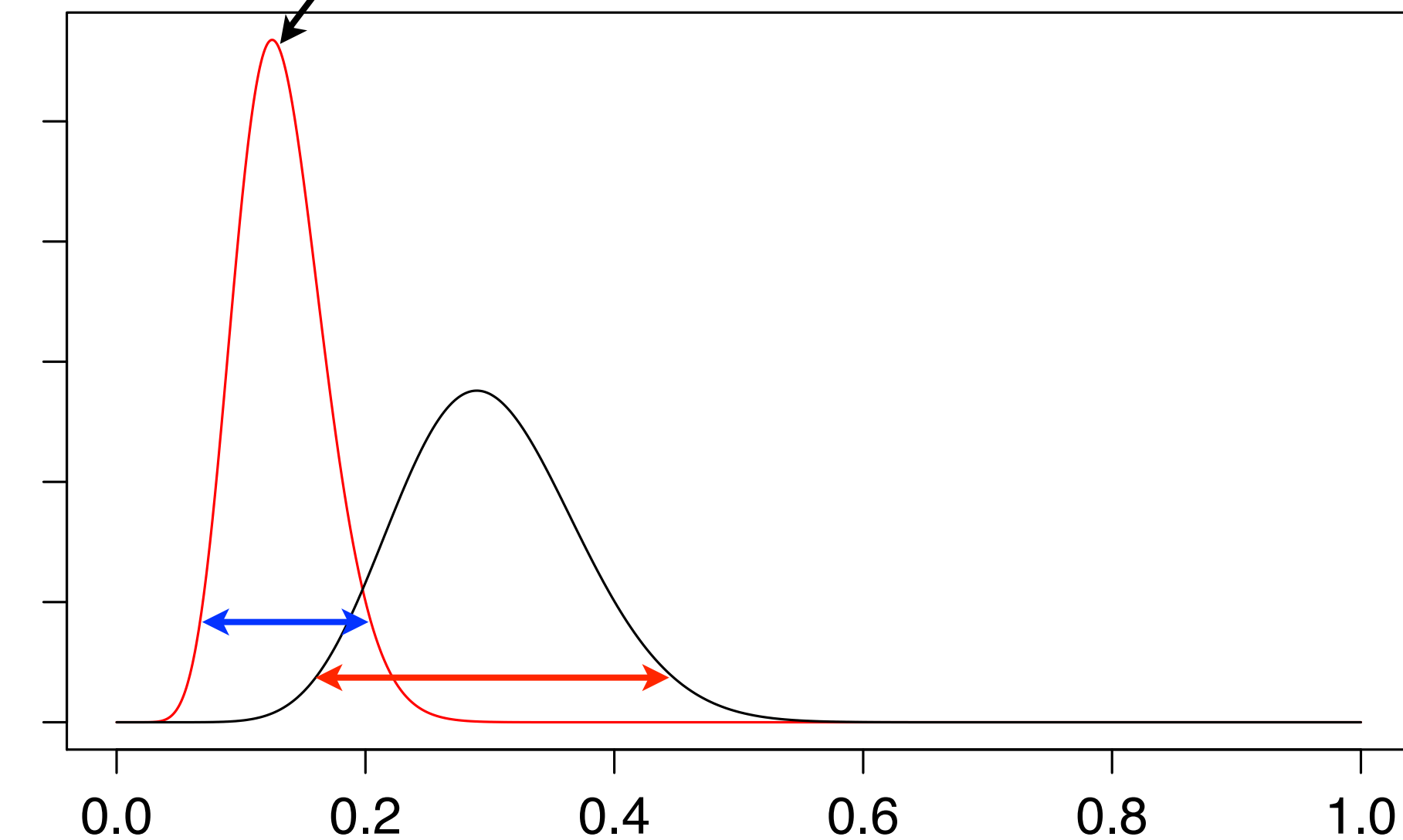
Bayesian View

“Belief”



New “belief” after 0 retained users

~15%



Stronger belief after
observing evidence

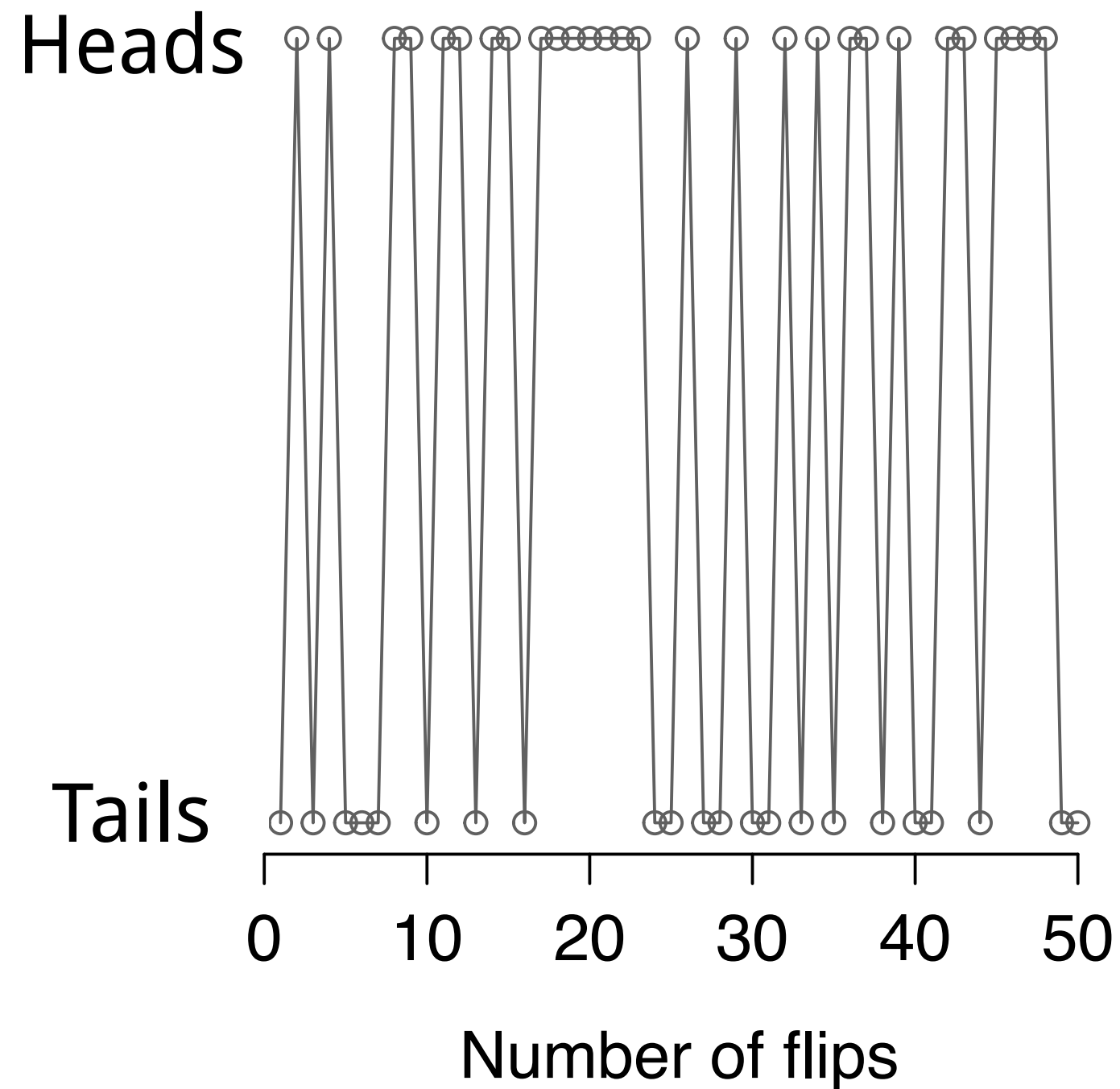
Retention

*The probability of the **model** given the **data***



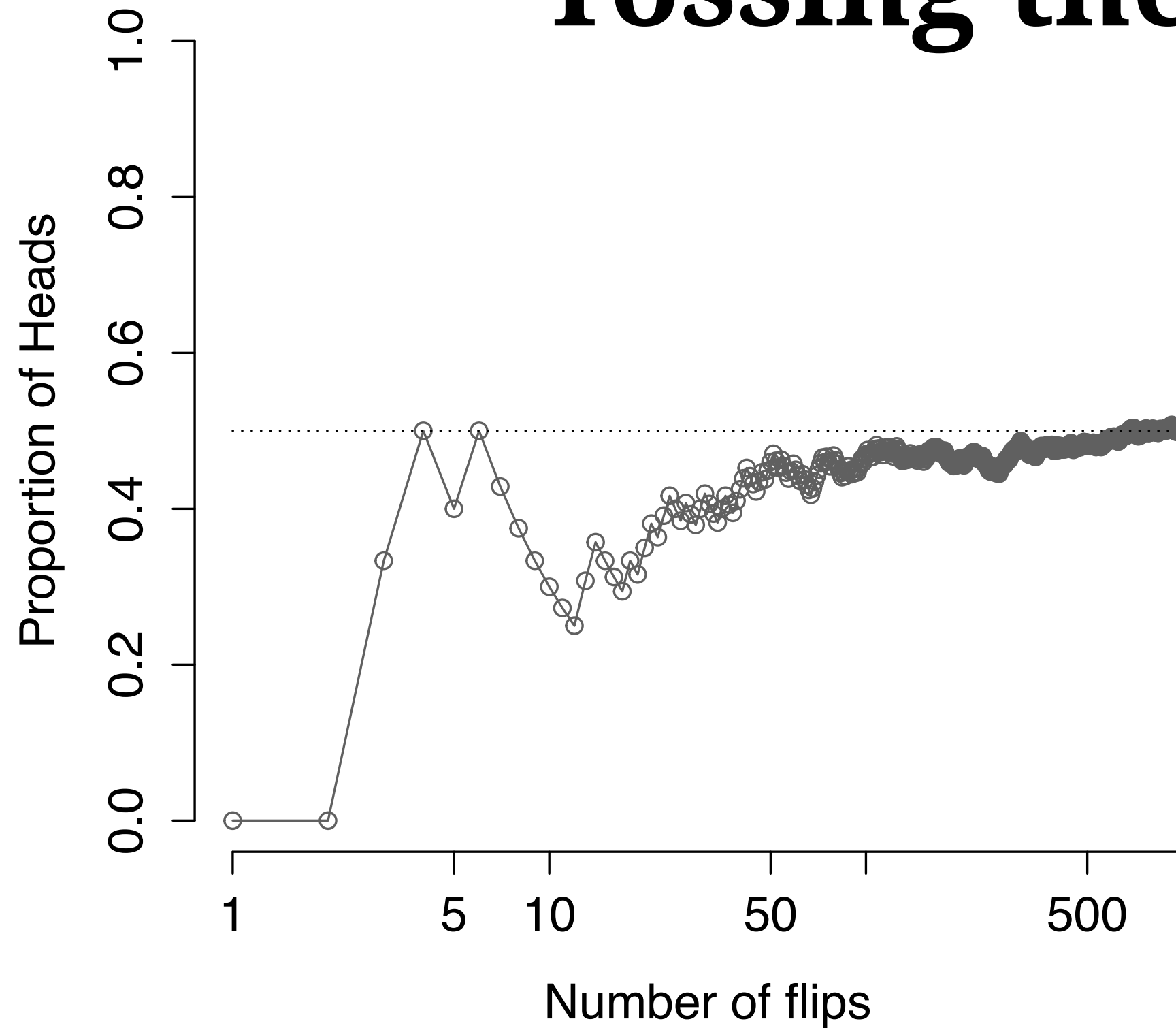
$$p(\text{heads}) = p(\text{tails}) = 0.5$$

Tossing the Coin



THTHTTTHTHTH...

Tossing the Coin



Long run average = 0.5

Terminology

$p(x)$ Probability of x

$p(x, y)$ Probability of x and y
(conjoint)

$p(x|y)$ Probability of x given y
(conditional)

The Bernoulli distribution

Head (H) = 1, Tails (T) = 0

A single toss: $p(x|\theta) = \theta^x (1 - \theta)^{(1-x)}$

For a “fair” coin, $\theta = 0.5$

$$p(heads = 1|0.5) = 0.5^1(1-0.5)^{(1-1)} = 0.5$$

$$p(tails = 0|0.5) = 0.5^0(1-0.5)^{(1-0)} = 0.5$$

The Binomial

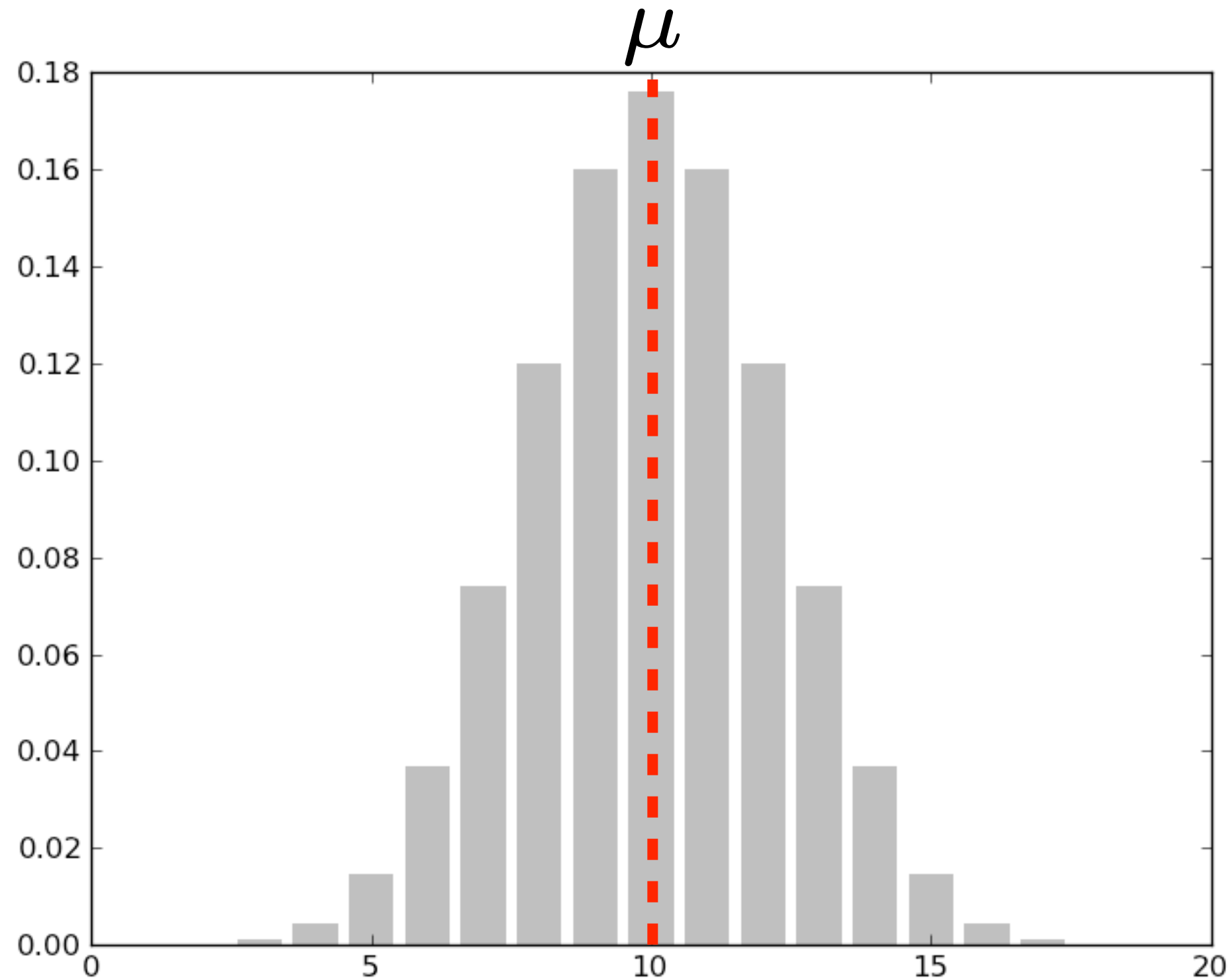
Probability of heads in *a single* throw:

$$p(x|\theta) = \theta^x (1 - \theta)^{(1-x)}$$

Probability of x heads in n throws:

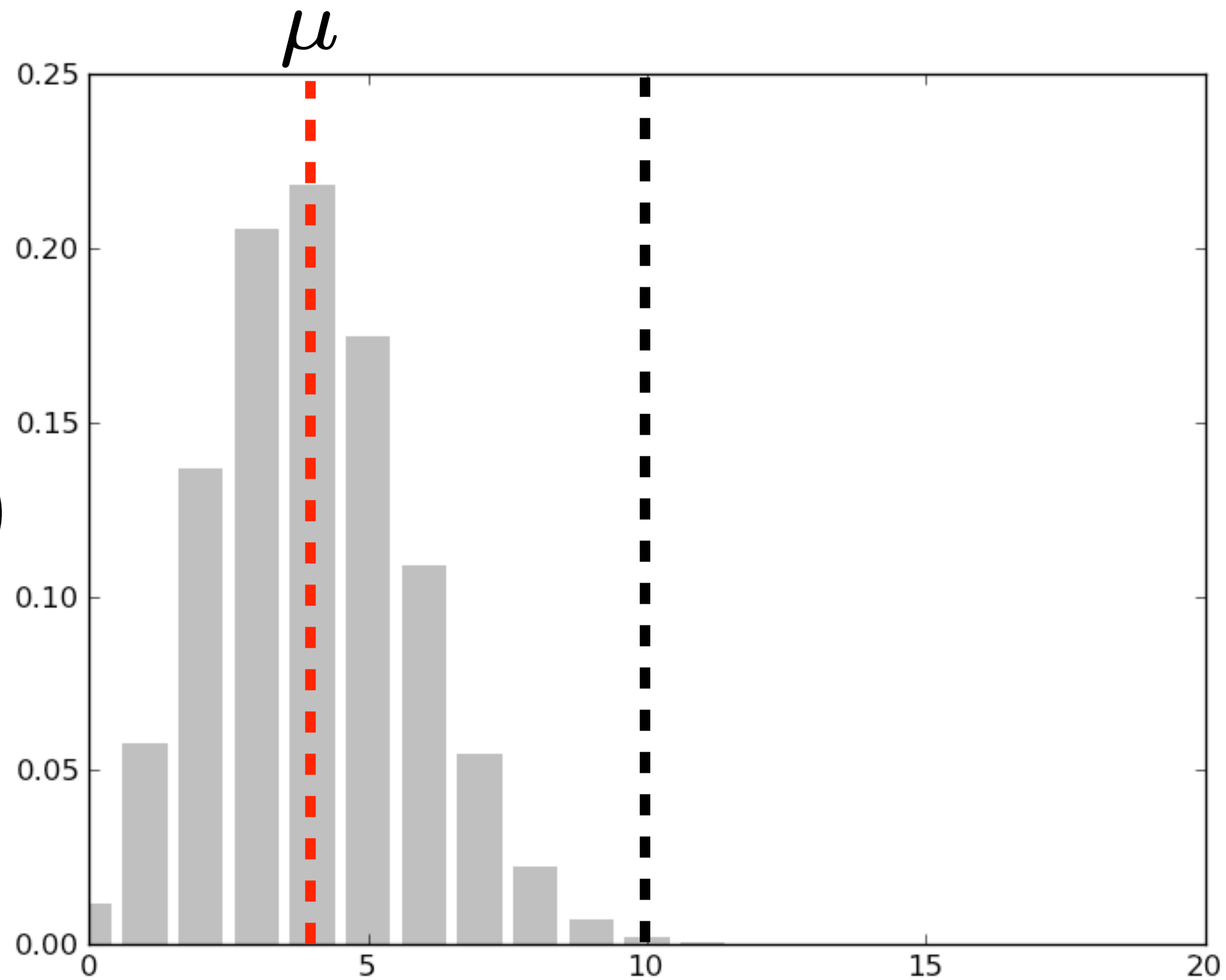
$$p(x|\theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{(n-x)}$$

$$p(x|0.5, 20)$$



20 tosses of a fair coin

$$p(x|0.2, 20)$$



20 tosses of an “un-fair” coin

$$p(x|\theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{(n-x)}$$

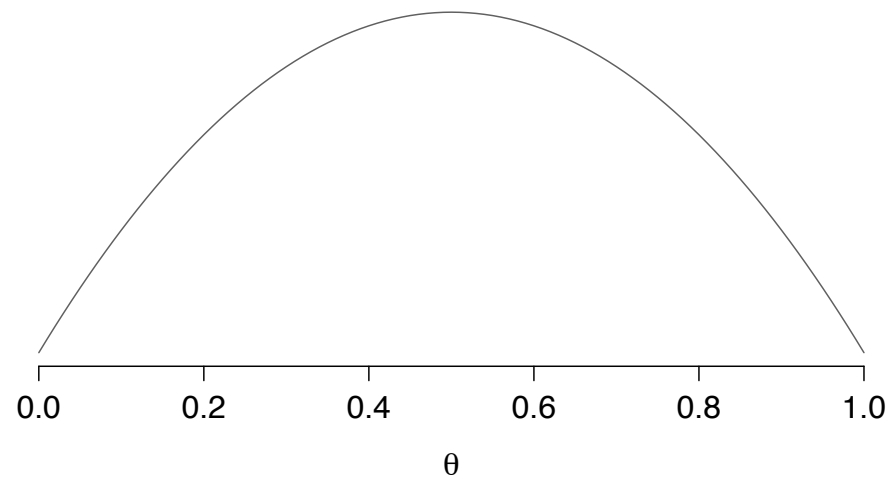
Likelihood of θ given observation i of x heads in n throws:

$$L(\theta|x_i, n_i) = \binom{n_i}{x_i} \theta^{x_i} (1 - \theta)^{(n_i-x_i)}$$

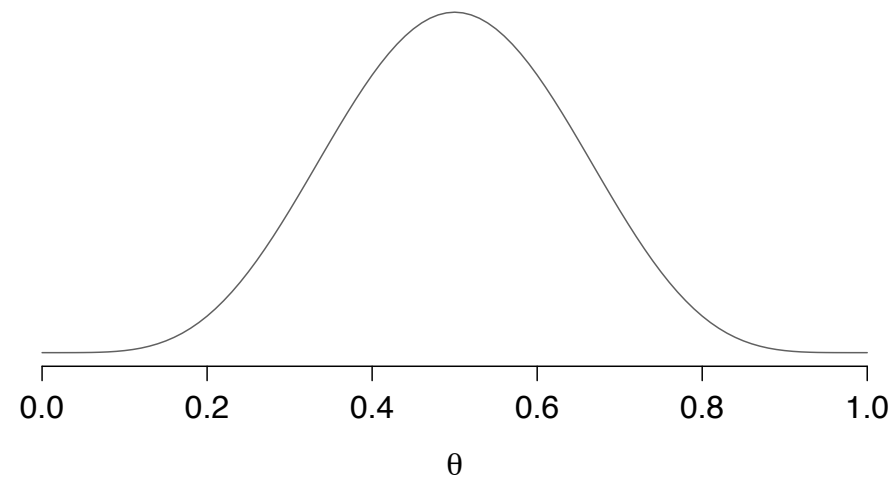
“Binomial Likelihood”

The Likelihood

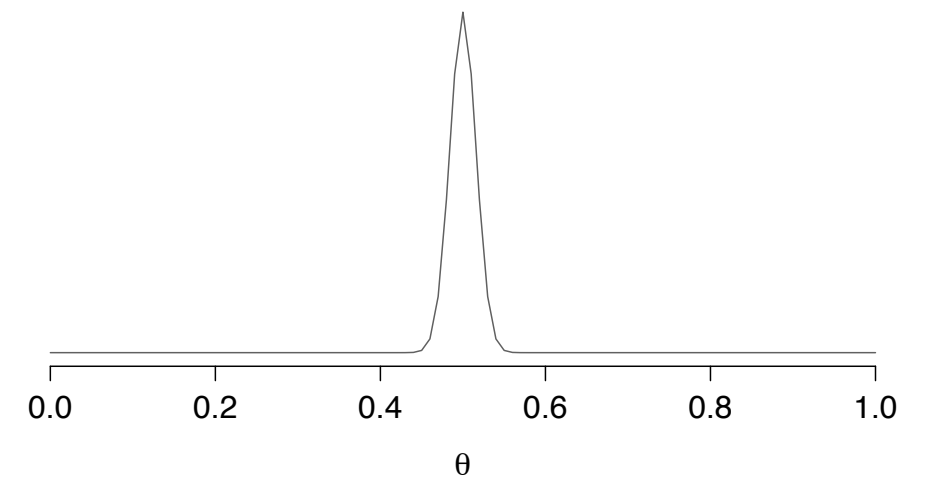
1 heads in 2 throws.



5 heads in 10 throws.



500 heads in 1000 throws.



Increasing likelihood of θ with more observations...

A recap...

x The observations (#heads)

θ The model parameter (e.g. fair coin)

$p(x|\theta)$ Probability of data given model

$p(\theta|x)$ We want to know this

Note that $p(x|\theta) \neq p(\theta|x)$

$$p(\textit{cloudy}|\textit{raining}) \neq p(\textit{raining}|\textit{cloudy})$$

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x)$$

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} \qquad p(x) = \sum_y p(x|y)p(y)$$

Bayes' Rule

$$p(y|x) = \frac{p(x|y)p(y)}{\sum_y p(x|y)p(y)}$$

discrete form

$$p(y|x) = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

continuous form

prob #heads given model

prob model given #heads

$$p(\theta|x)$$

$$p(x|\theta)$$

prob of model

$$p(\theta)$$

$$p(y|x)$$

=

$$\frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

$$p(x)$$

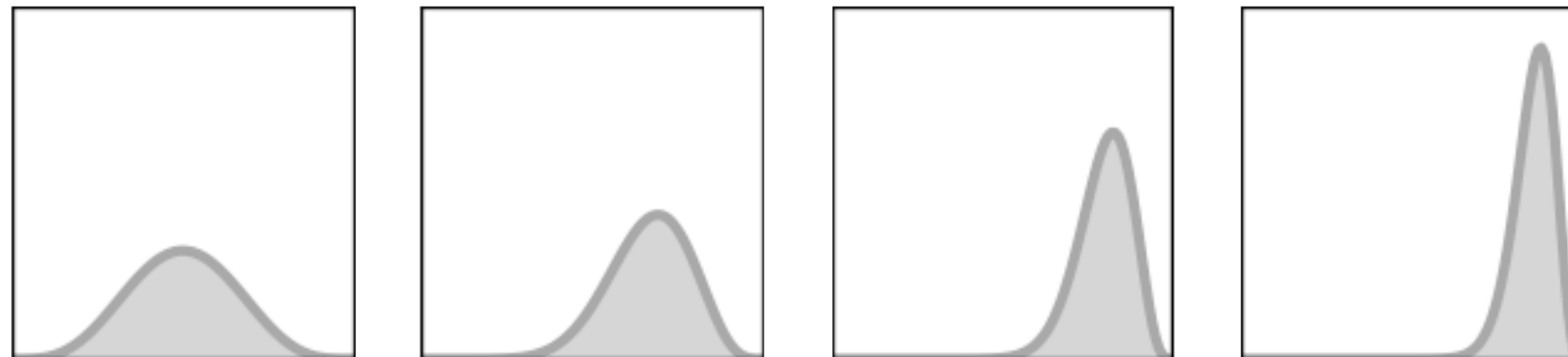
$$\underbrace{p(\theta|x)}_{\text{posterior}} = \underbrace{p(x|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}} / \underbrace{p(x)}_{\substack{\text{normalizing} \\ \text{factor}}}$$

$$\text{normalizing factor } p(x) = \int p(x|\theta)p(\theta)d\theta$$

The prior

$$p(\theta)$$

Captures our “belief” in the model
based on prior experience, observations or knowledge



$$p(\theta|x) = p(x|\theta) p(\theta) / p(x)$$



$$\hat{p}_0(\theta|x_0) = p(x_0|\theta) p(\theta) / p(x_0)$$



$$\hat{p}_1(\theta|x_1) = p(x_1|\theta) \hat{p}_0(\theta) / \hat{p}_0(x_1)$$



⋮

$$\hat{p}_n(\theta|x_n) = p(x_n|\theta) \hat{p}_{n-1}(\theta) / \hat{p}_{n-1}(x_n)$$

Best estimate so far

Iterations with more data...

Selecting a prior

$$p(x|\theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{(n-x)}$$

$$p(\theta|x) = \frac{\theta^x (1 - \theta)^{(n-x)} p(\theta)}{\int \theta^x (1 - \theta)^{(n-x)} p(\theta) d\theta}$$

We'd like the product of prior and likelihood
to be “like” the likelihood

We'd like the integral to be easily evaluated

“Conjugate prior”

$$p(\bar{\theta}) = p(x|\theta)p(\theta)$$

Beta distribution

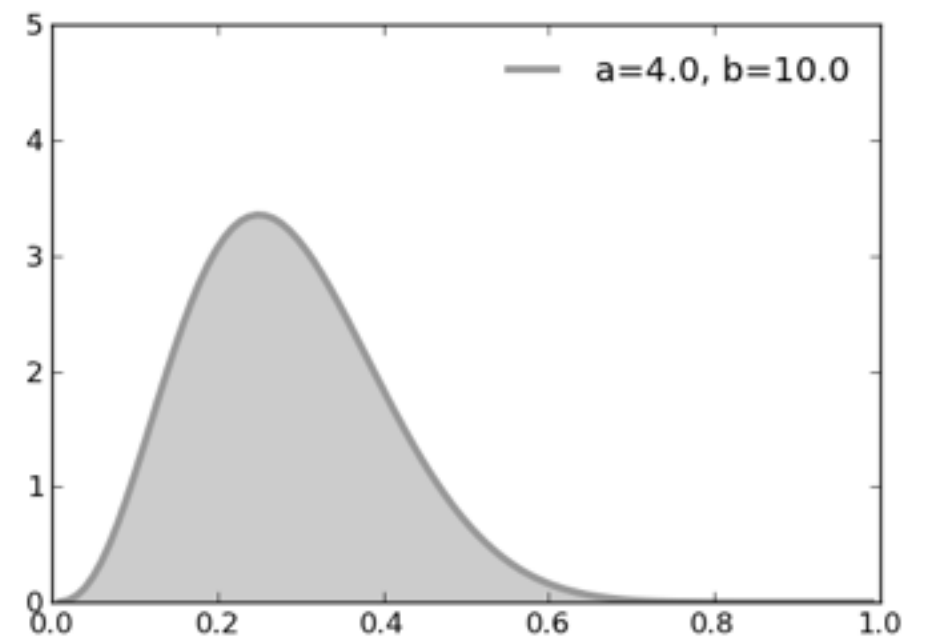
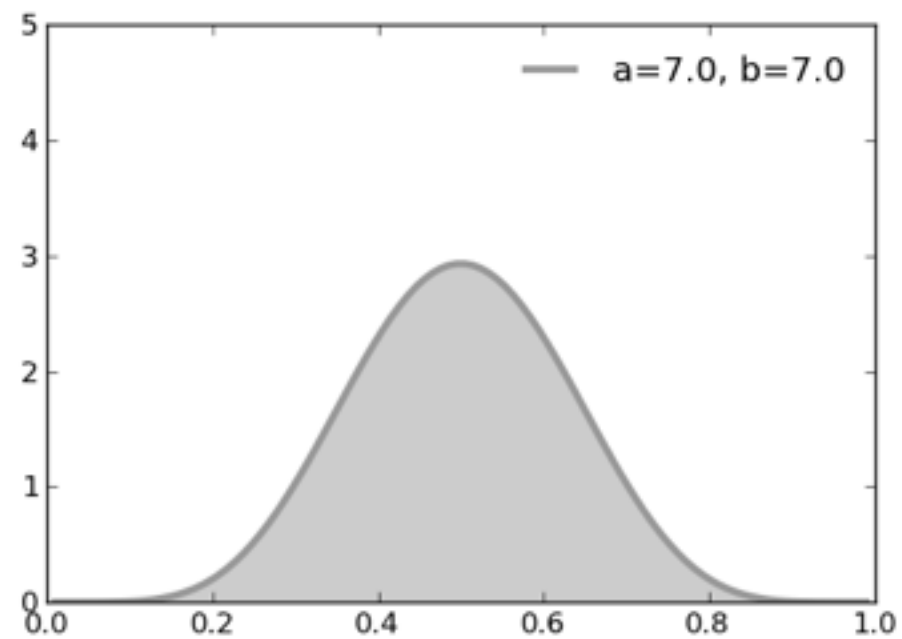
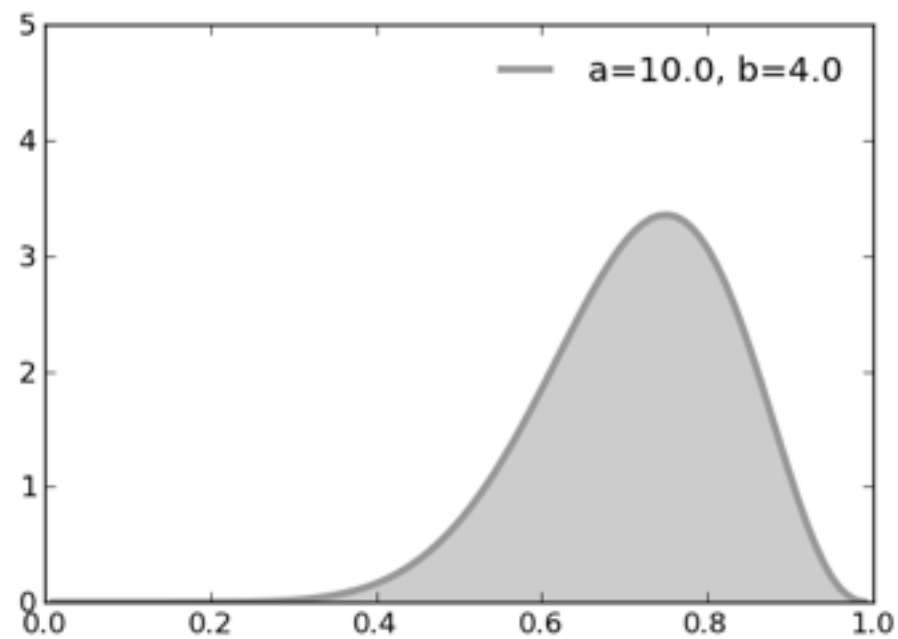
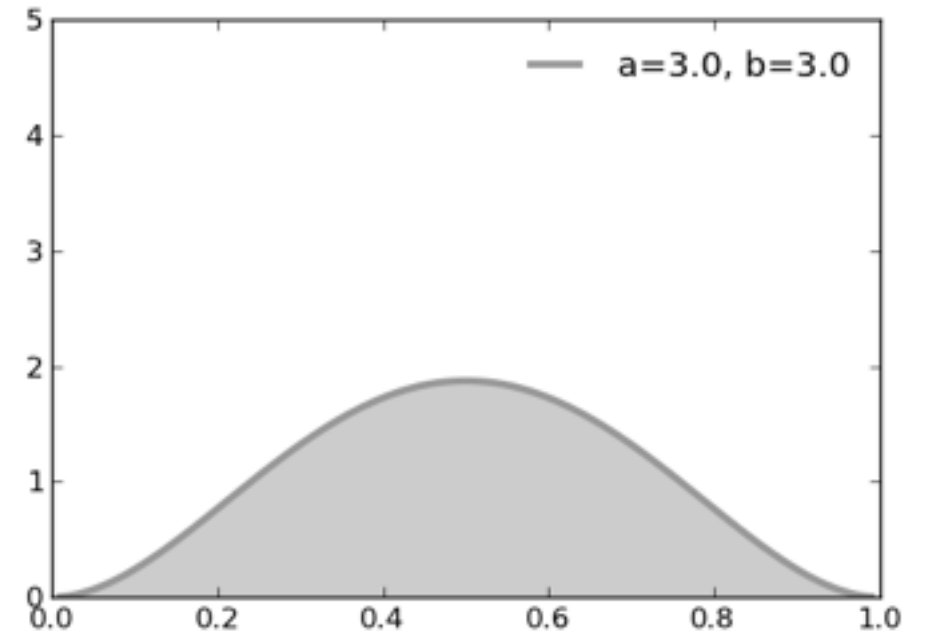
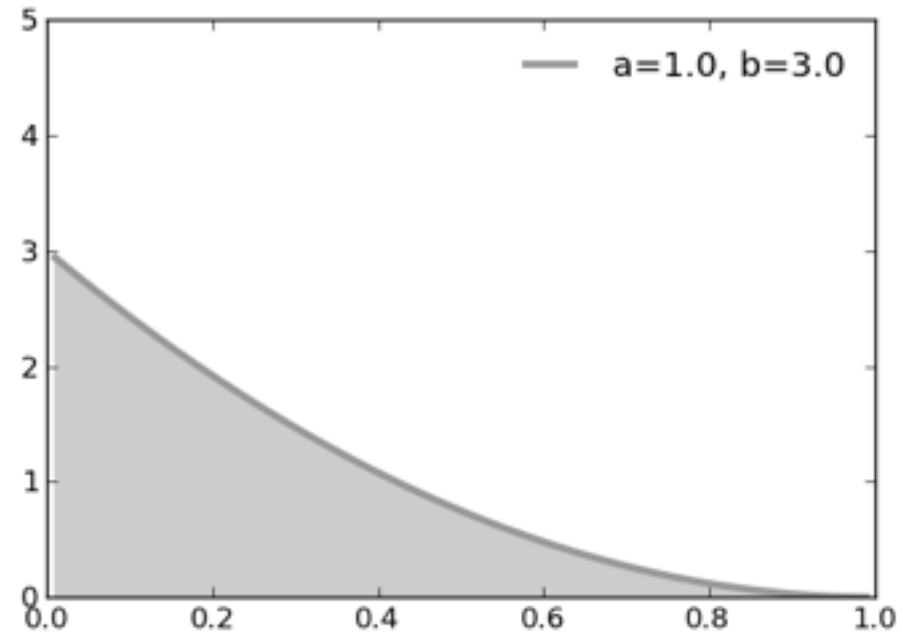
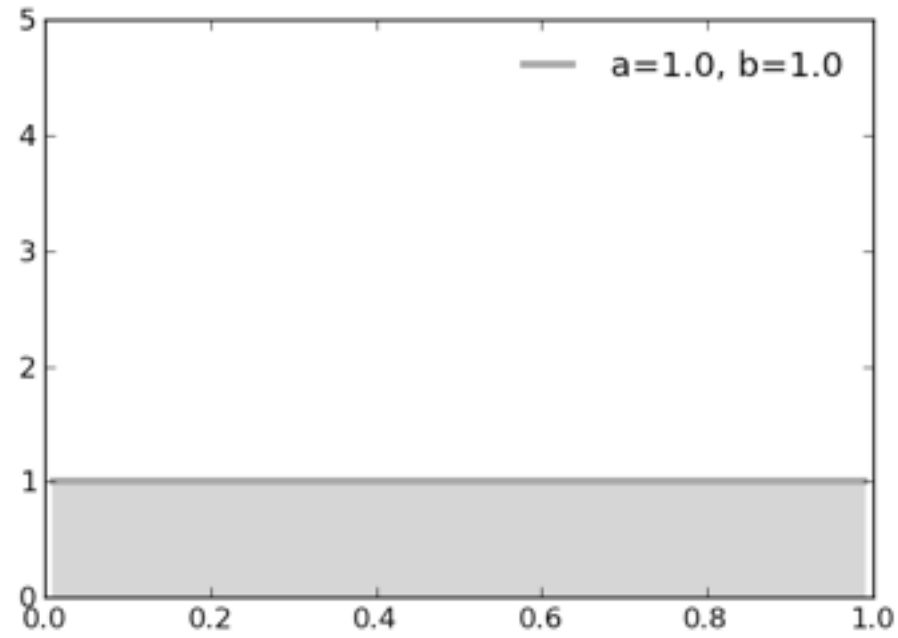
$$\text{beta}(\theta|a, b) = \theta^{(a-1)} (1 - \theta)^{(b-1)} / B(a, b)$$

number of heads + 1

number of tails + 1

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} = \frac{(a-1)!(b-1)!}{(a+b-1)!}$$

Beta distribution



Putting it together...

$$p(\theta|x, n) = \overset{\text{binomial likelihood}}{\theta^x (1 - \theta)^{(n-x)}} \overset{\text{beta prior}}{\theta^{(a-1)} (1 - \theta)^{(b-1)}} / B(a, b) p(x)$$

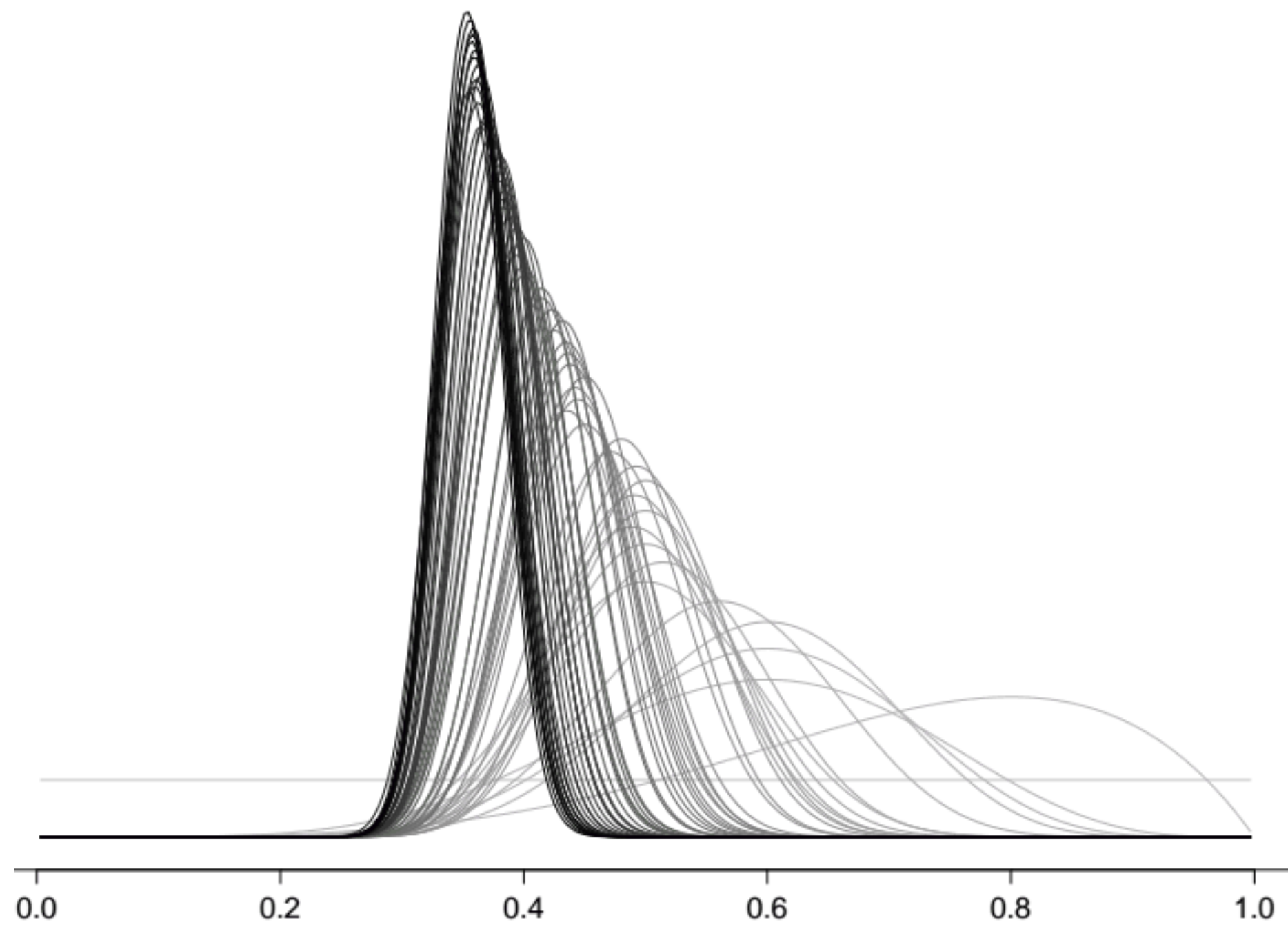
$$= \theta^{x+a-1} (1 - \theta)^{(n-x+b-1)} / B(x + a, n - x + b)$$

number of heads x

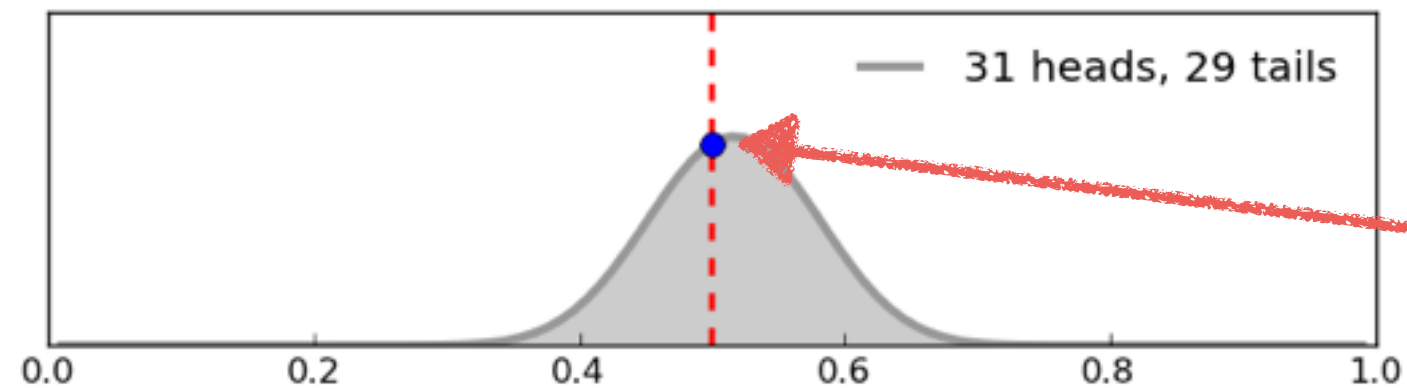
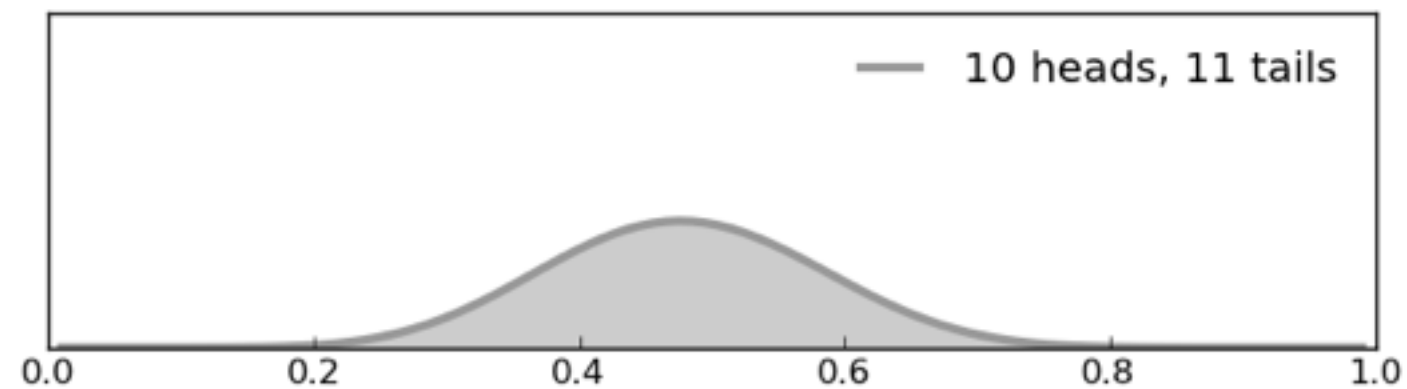
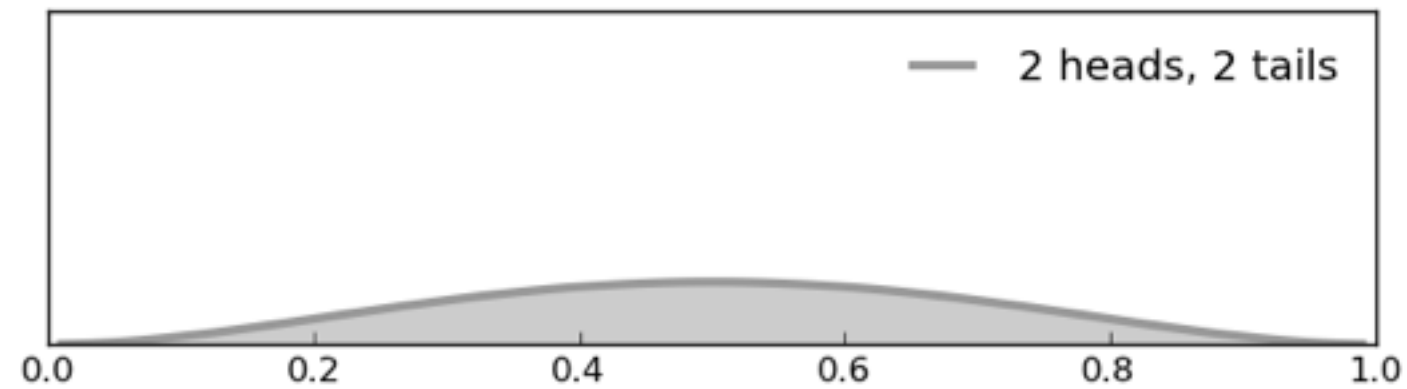
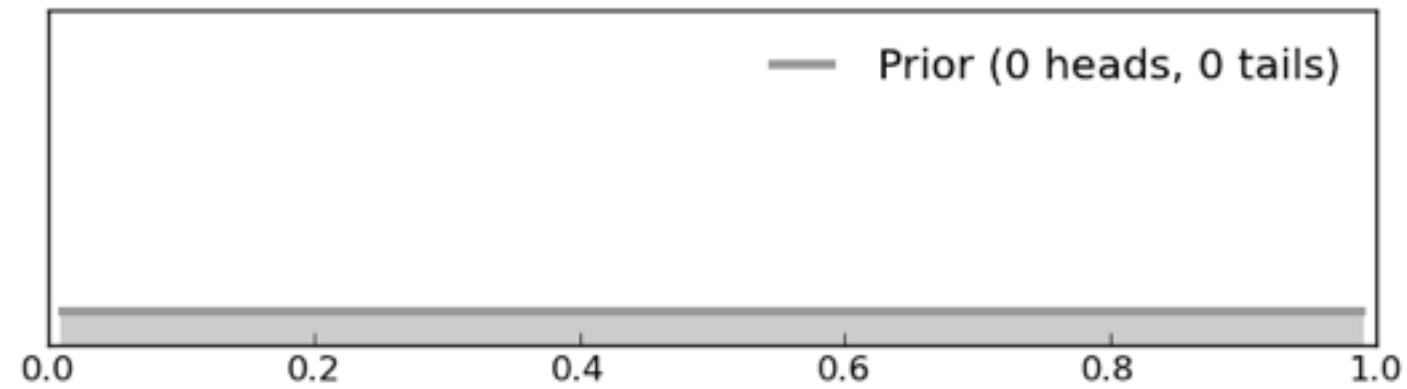
number of tails $(n-x)$

Putting it together...

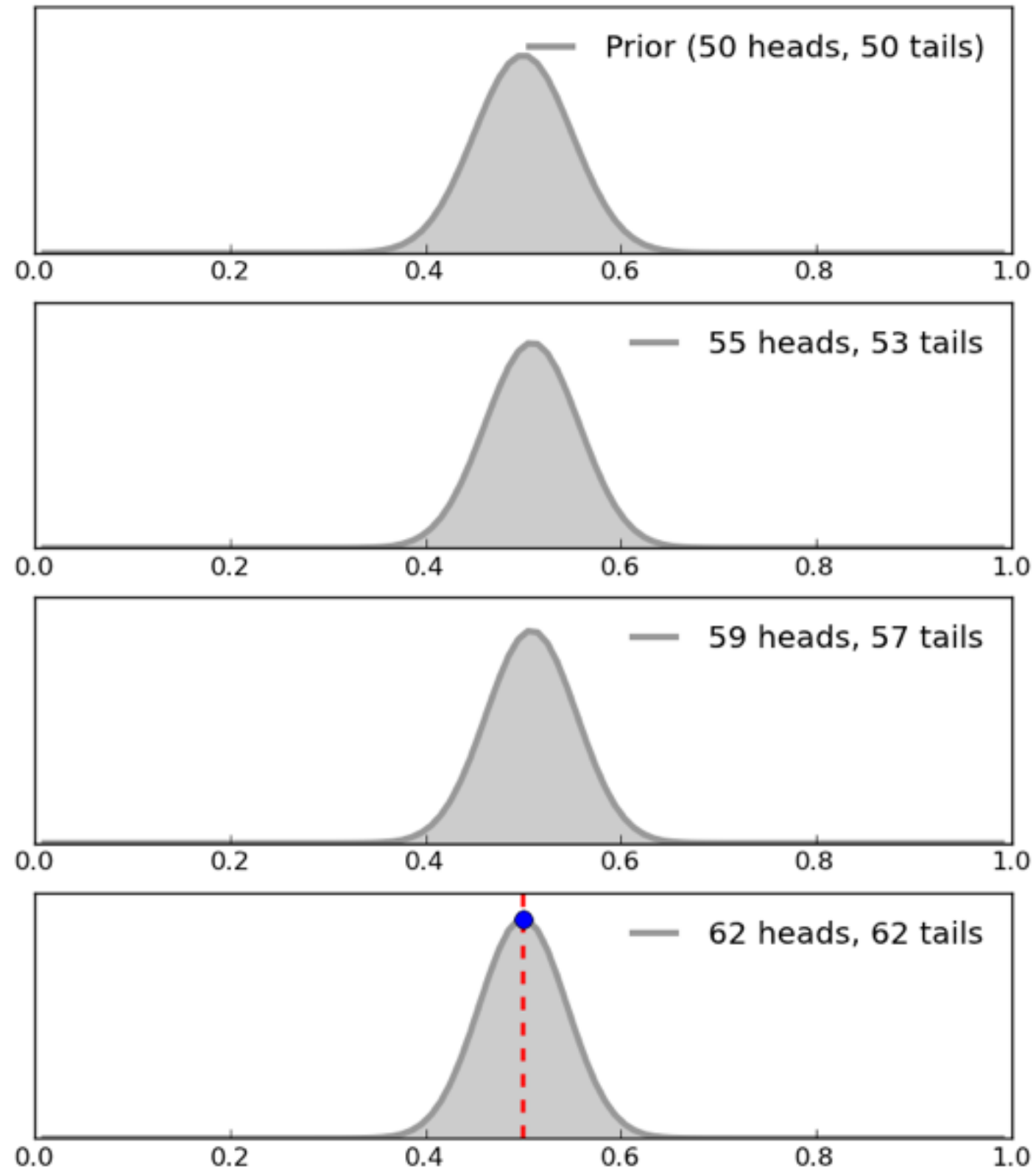
1. Decide on a **prior**, which captures your **belief**
2. Run experiment and **observe data** for heads, tails
3. Determine your **posterior** based on the data
4. Use posterior as your **new belief** and re-run experiment
5. Rinse, repeat until you hit an actionable **certainty**



Uniform prior “Fair” coin



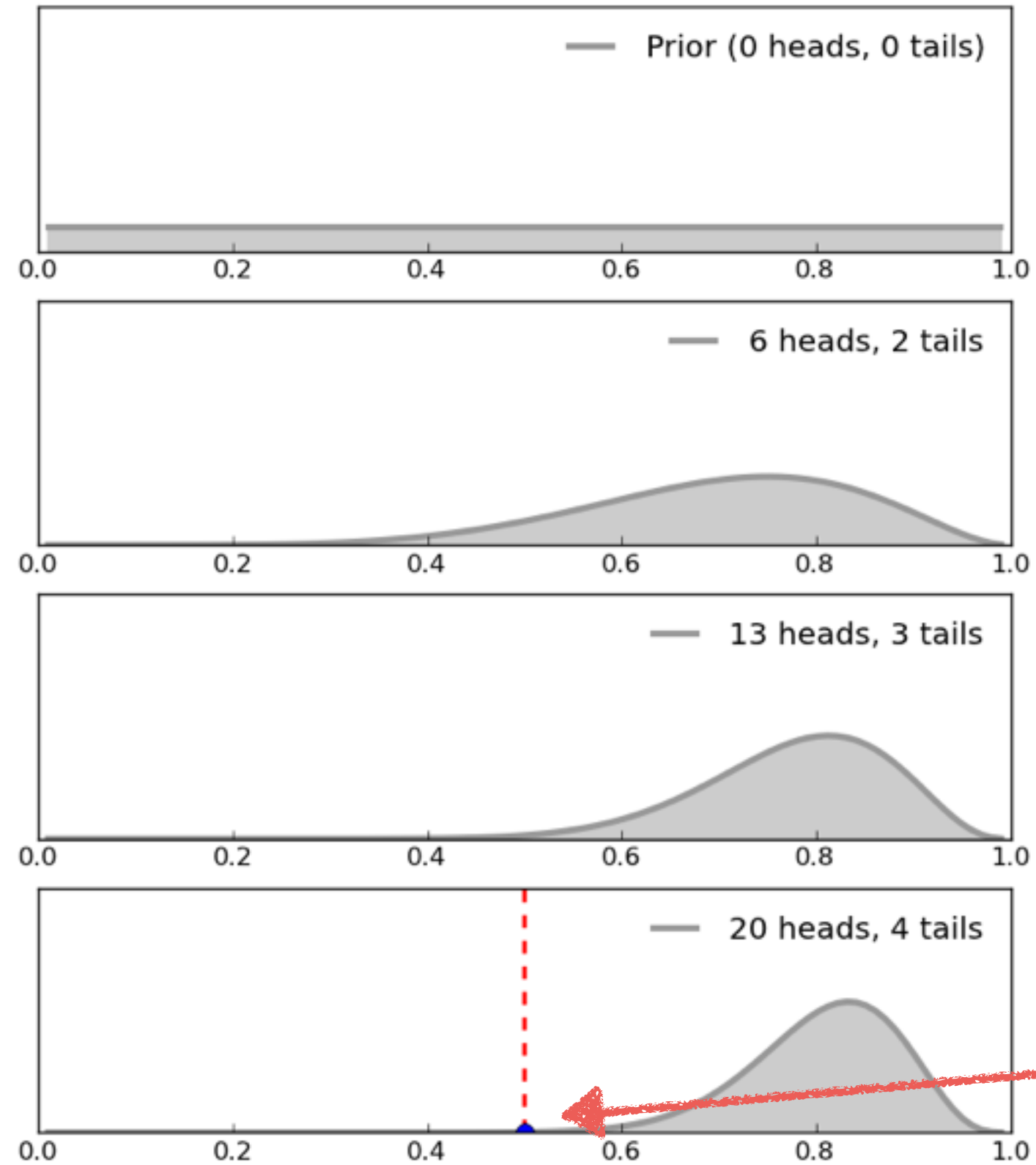
Pretty sure coin is fair



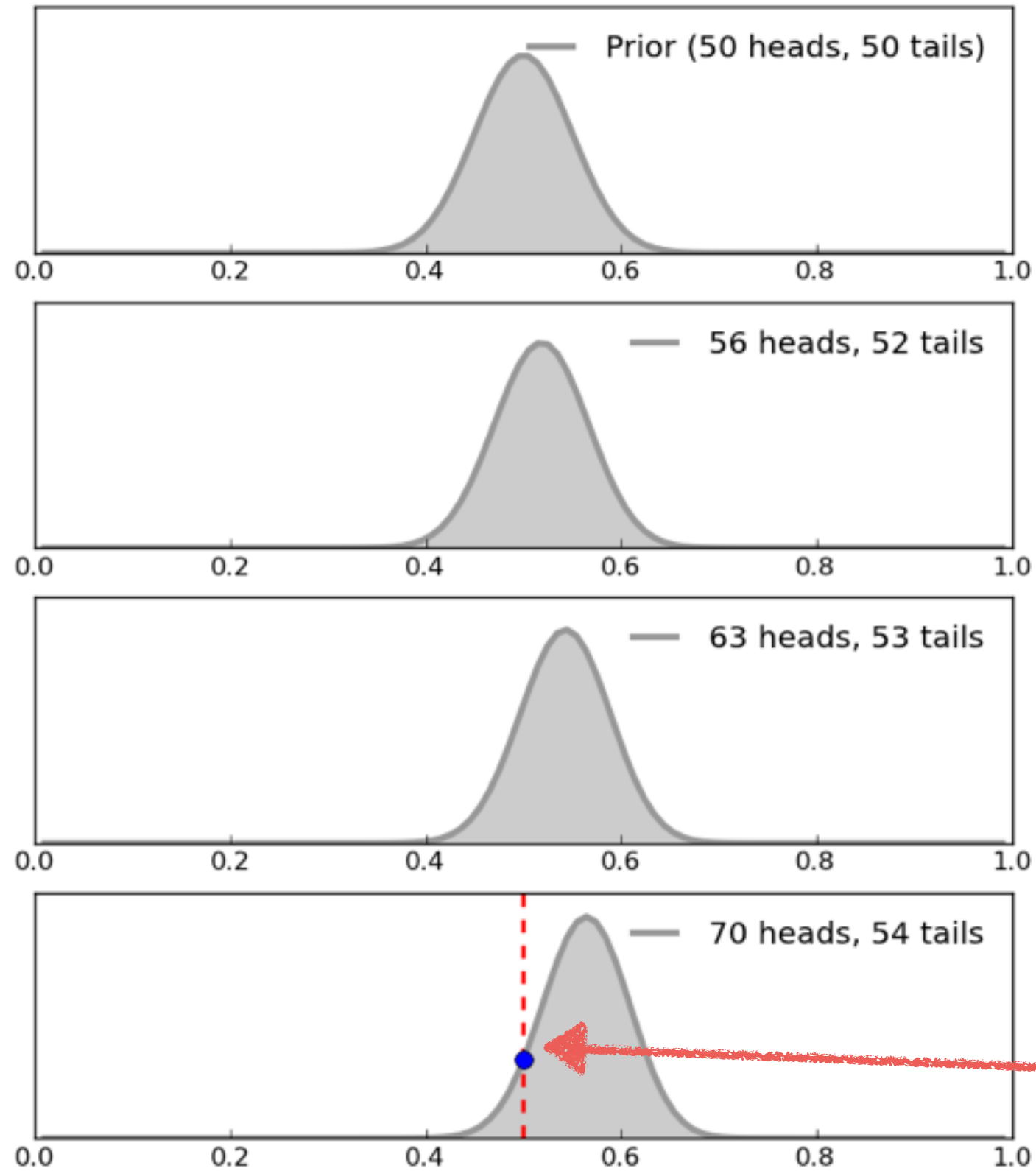
“Coin is fair” prior
“Fair” coin

Very sure coin is fair

Uniform prior “Biased” coin



Pretty sure coin is unfair

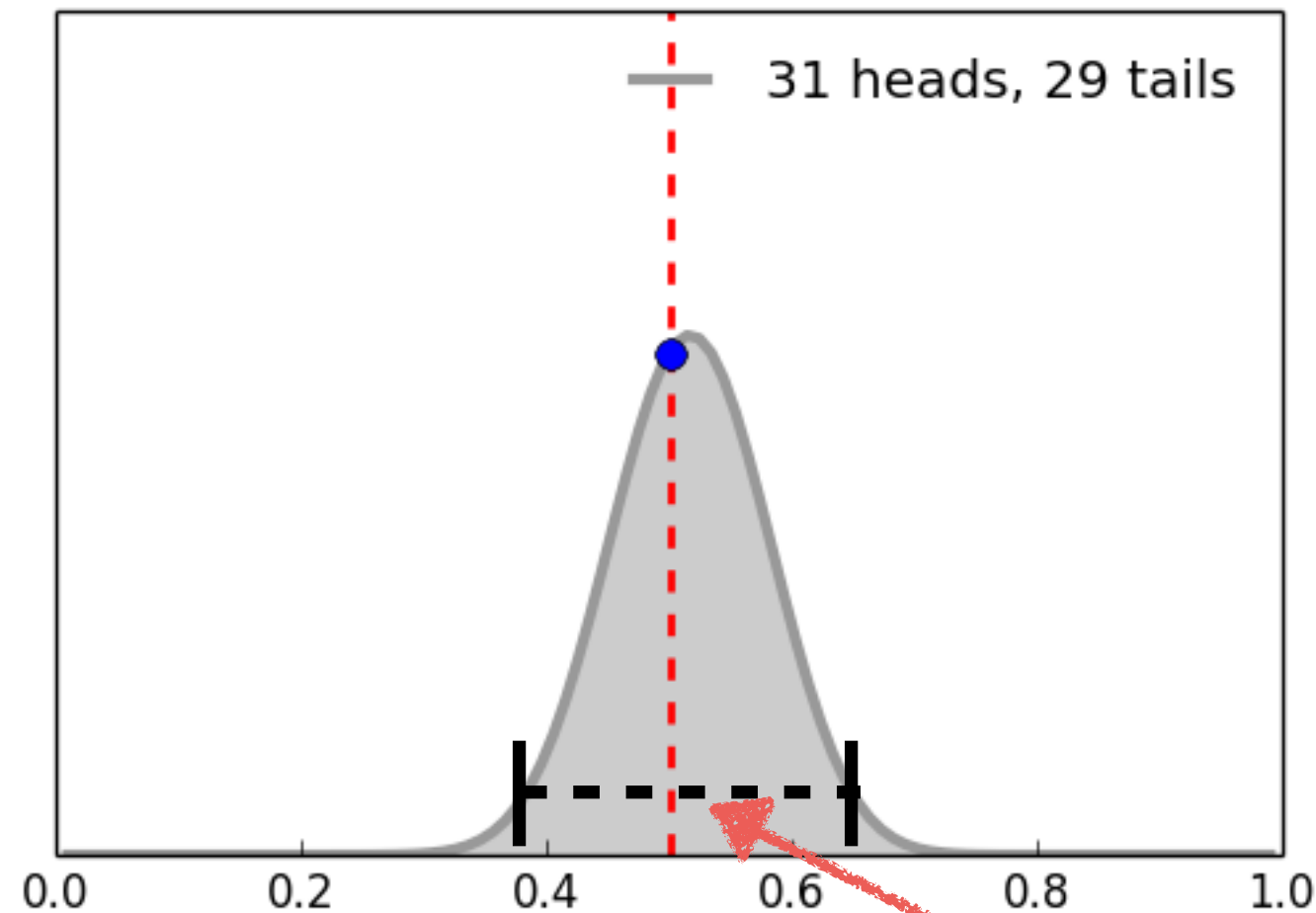


“Coin is fair” prior
“biased” coin

Not sure of anything yet!

When to reject?

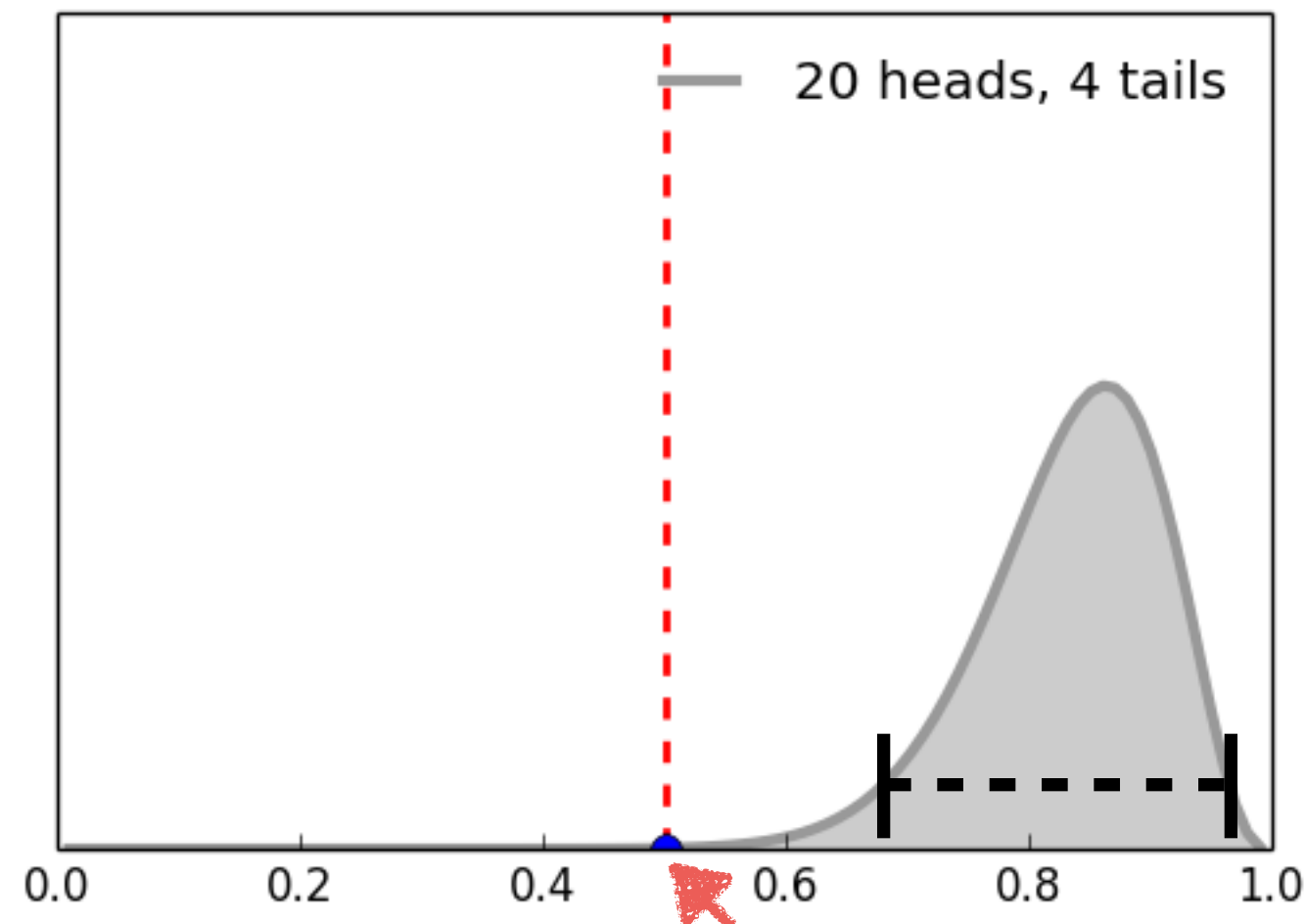
The Credible Interval



Uniform prior
“Fair” coin

95% credible interval

The Credible Interval



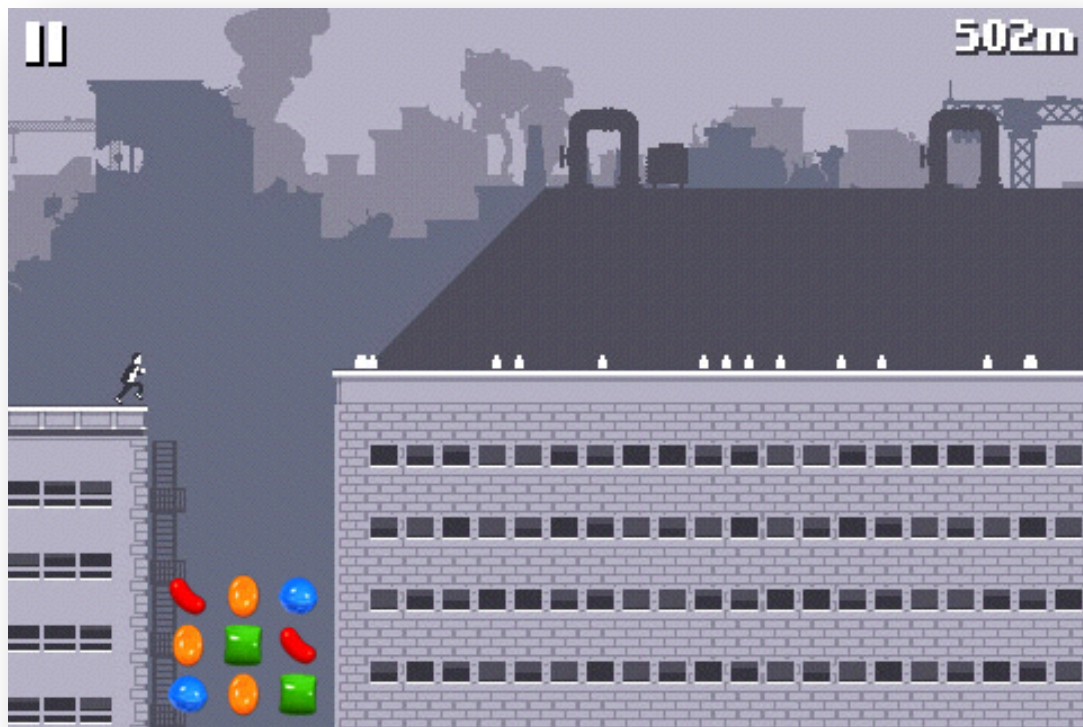
Uniform prior
“Biased” coin

Outside credible interval

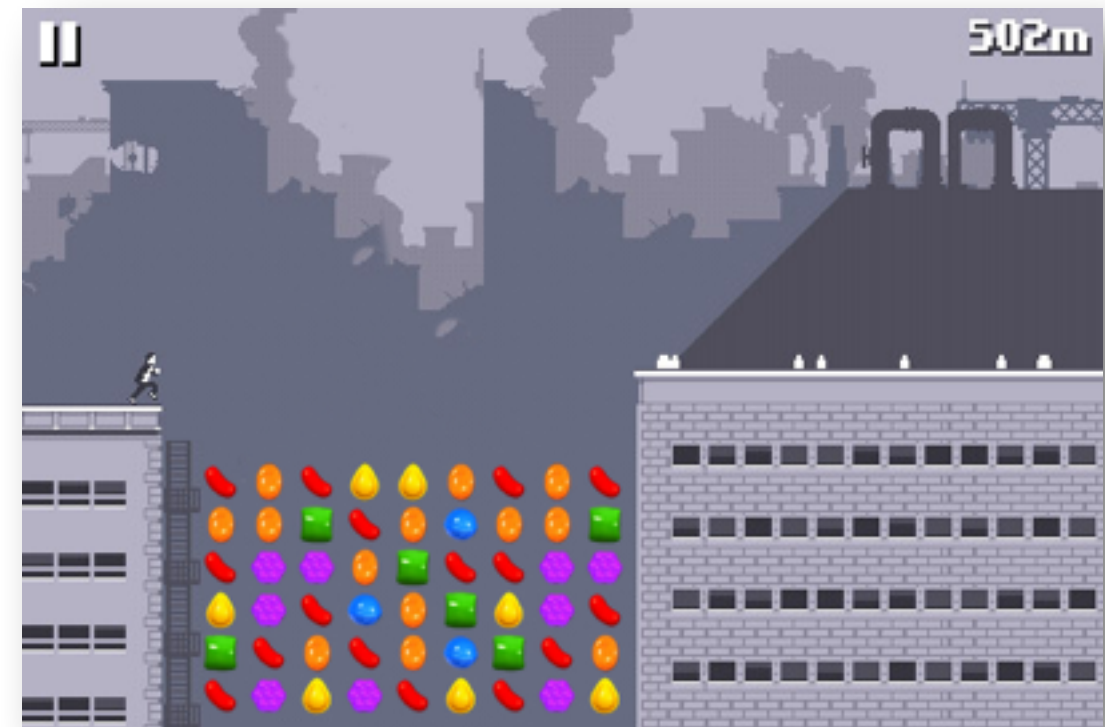
The Prior

- Captures our prior belief, expertise, opinion
- Strong prior belief means:
 - ▶ we need lots of evidence to contradict
 - ▶ results converge more quickly (if prior is “relevant”)
- Provides inertia
- With enough samples, prior’s impact diminishes, rapidly

Running a test...



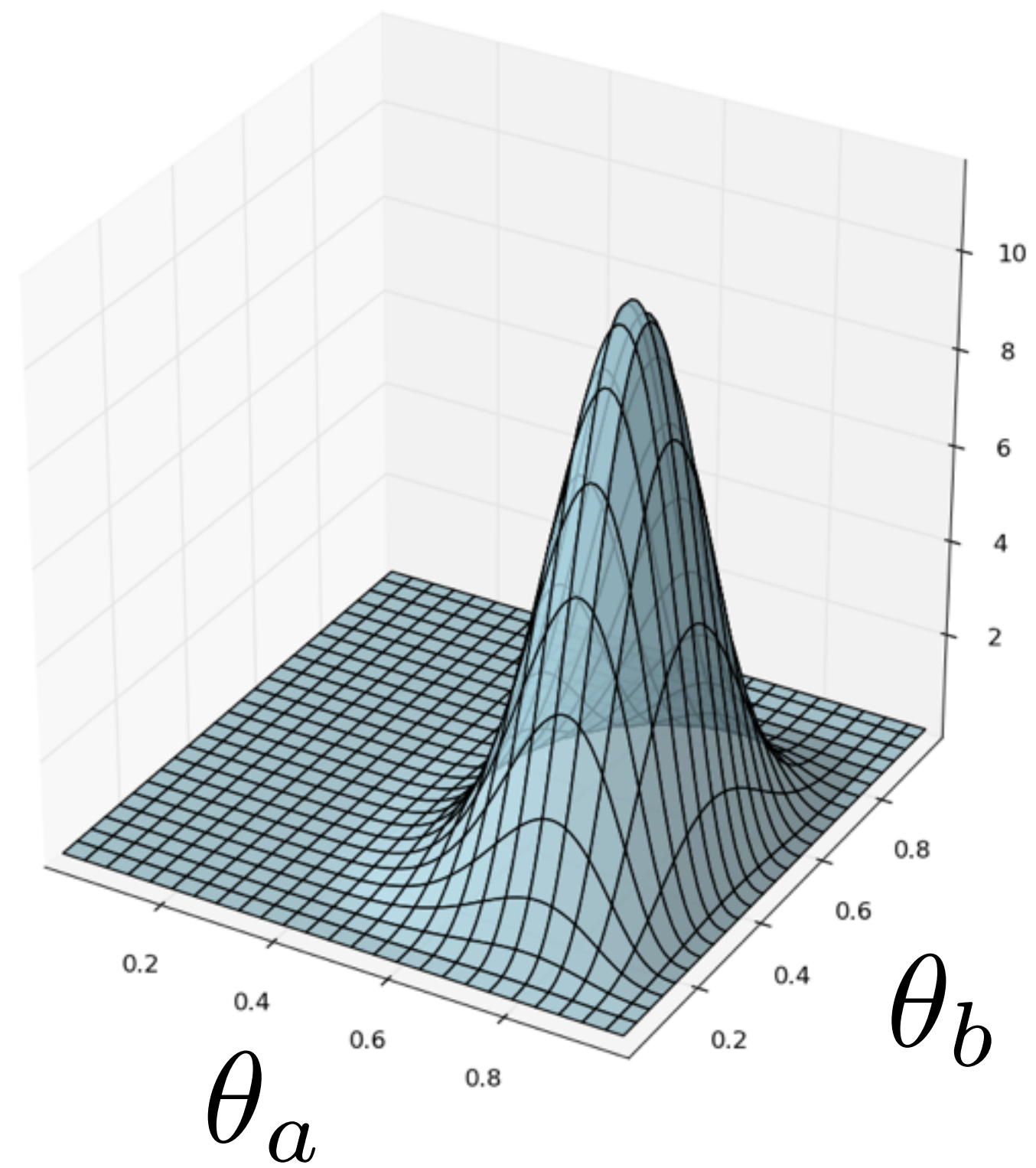
A



B

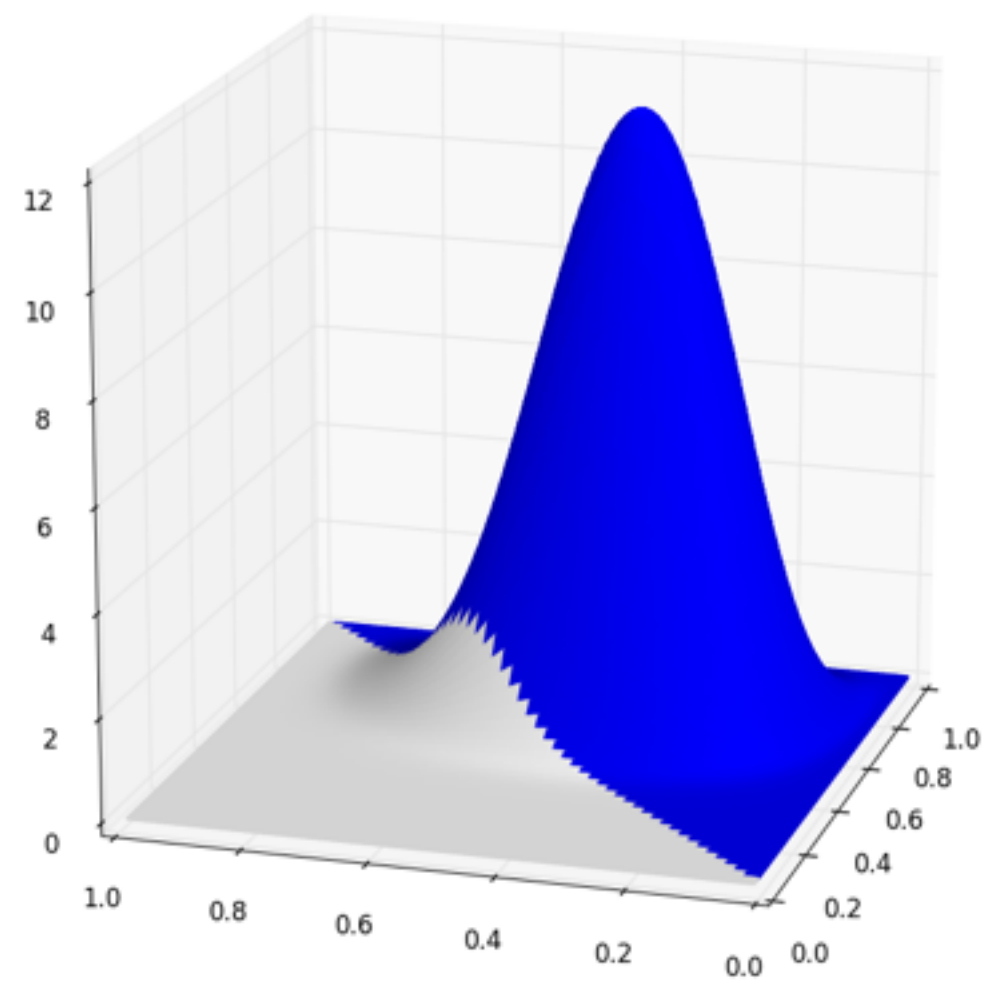
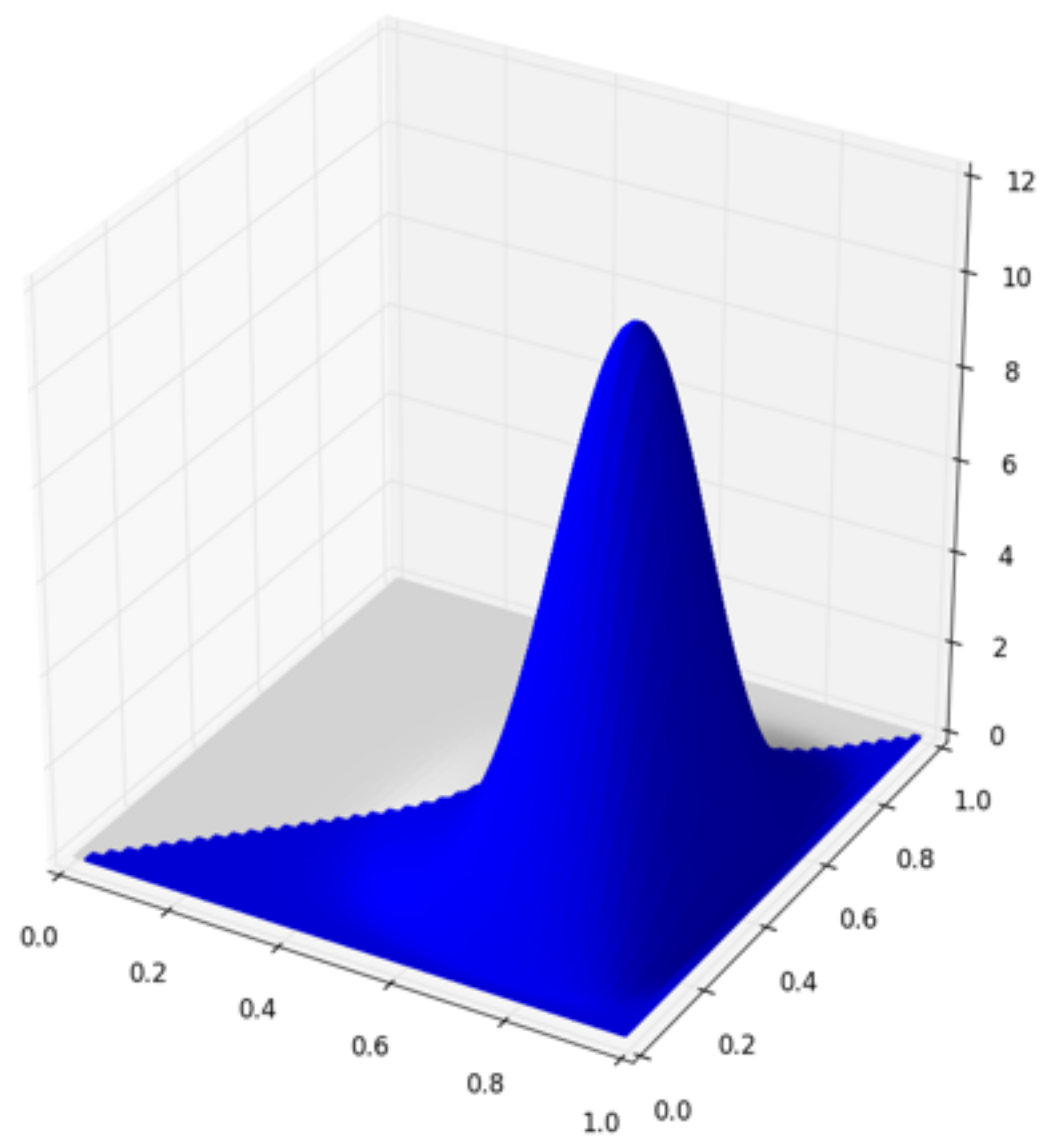
Multiple variant tests

- With 1 or more variants we have a multi-dimensional problem
- Need to evaluate volumes under the posterior
- In general requires numerical quadrature = Markov Chain Monte-Carlo (MCMC)



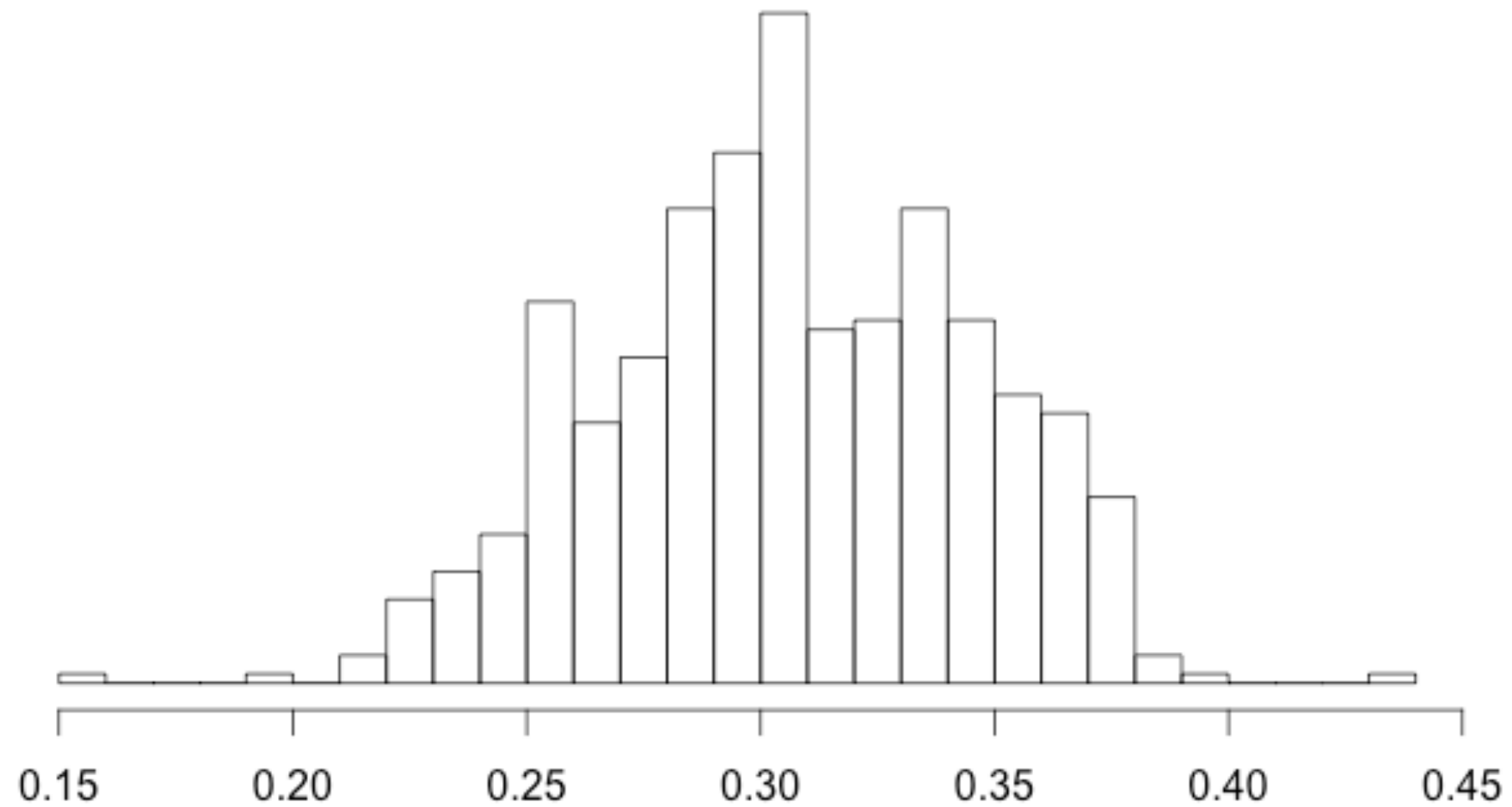
Probability of Winning

$$p(\theta_a > \theta_b) = \int_{\theta_a > \theta_b} p(\theta_a | x_a) p(\theta_b | x_b) d\theta_a d\theta_b$$

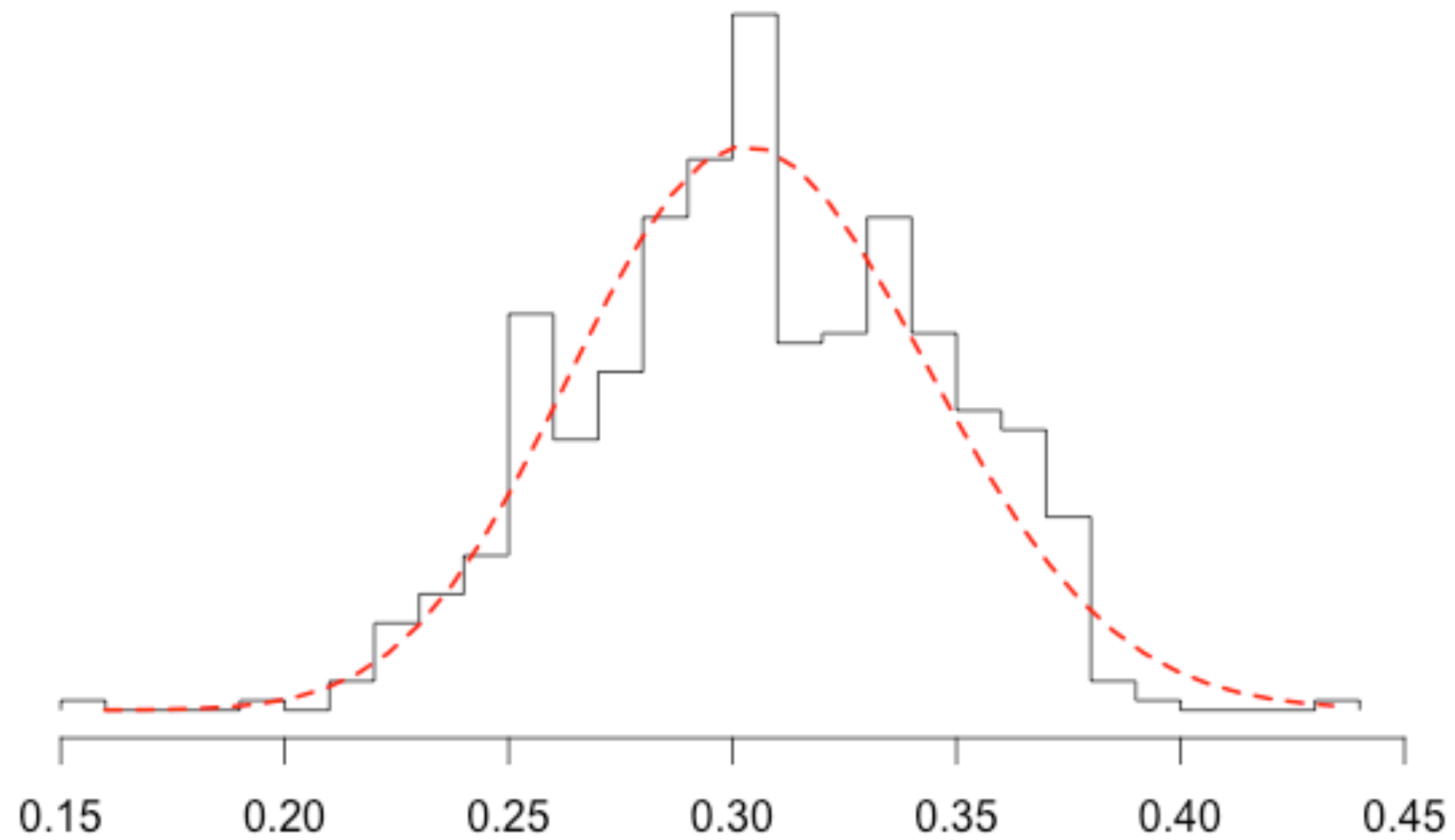


■ $\theta_a > \theta_b$

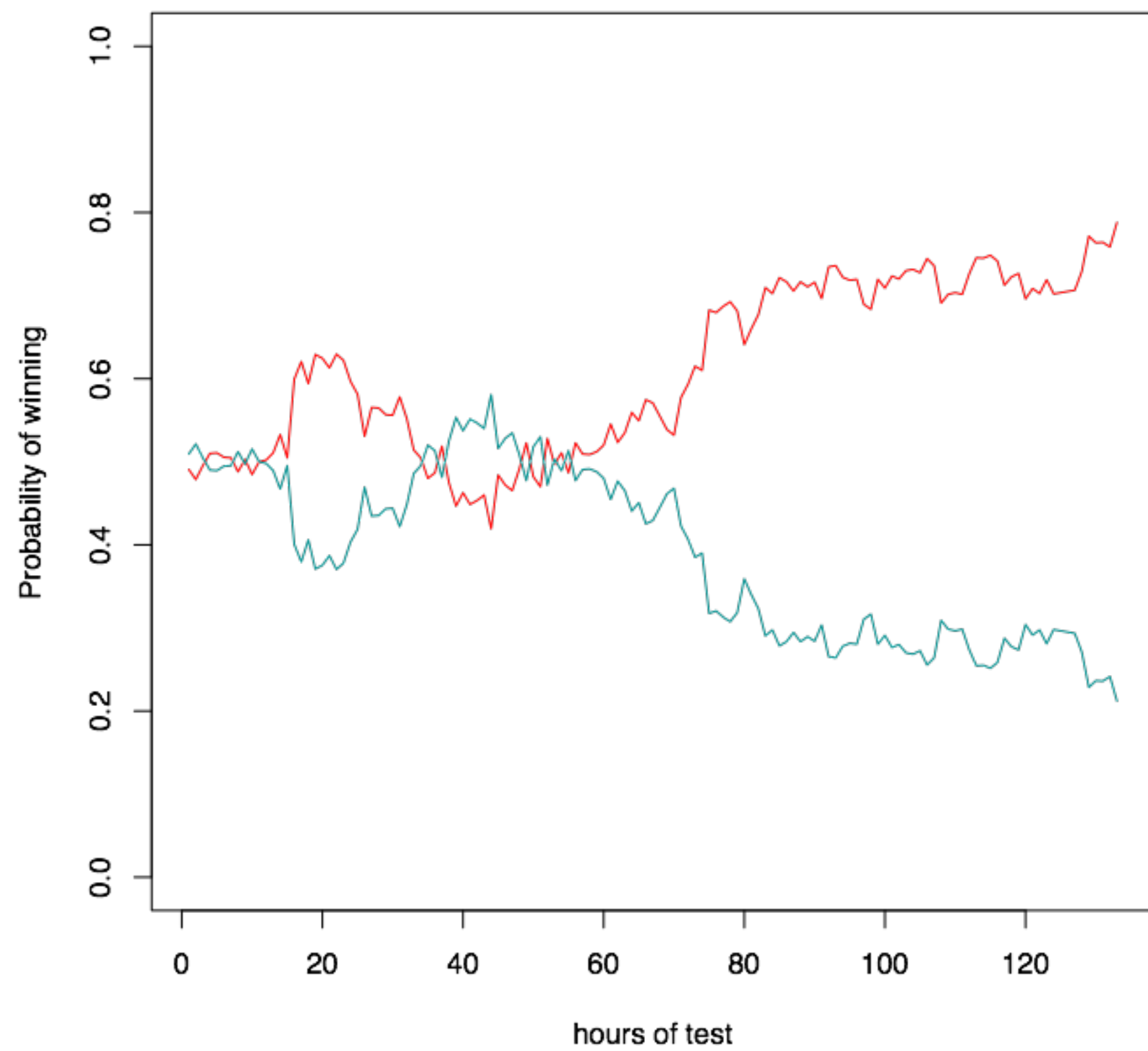
What's the prior?



Fit a beta



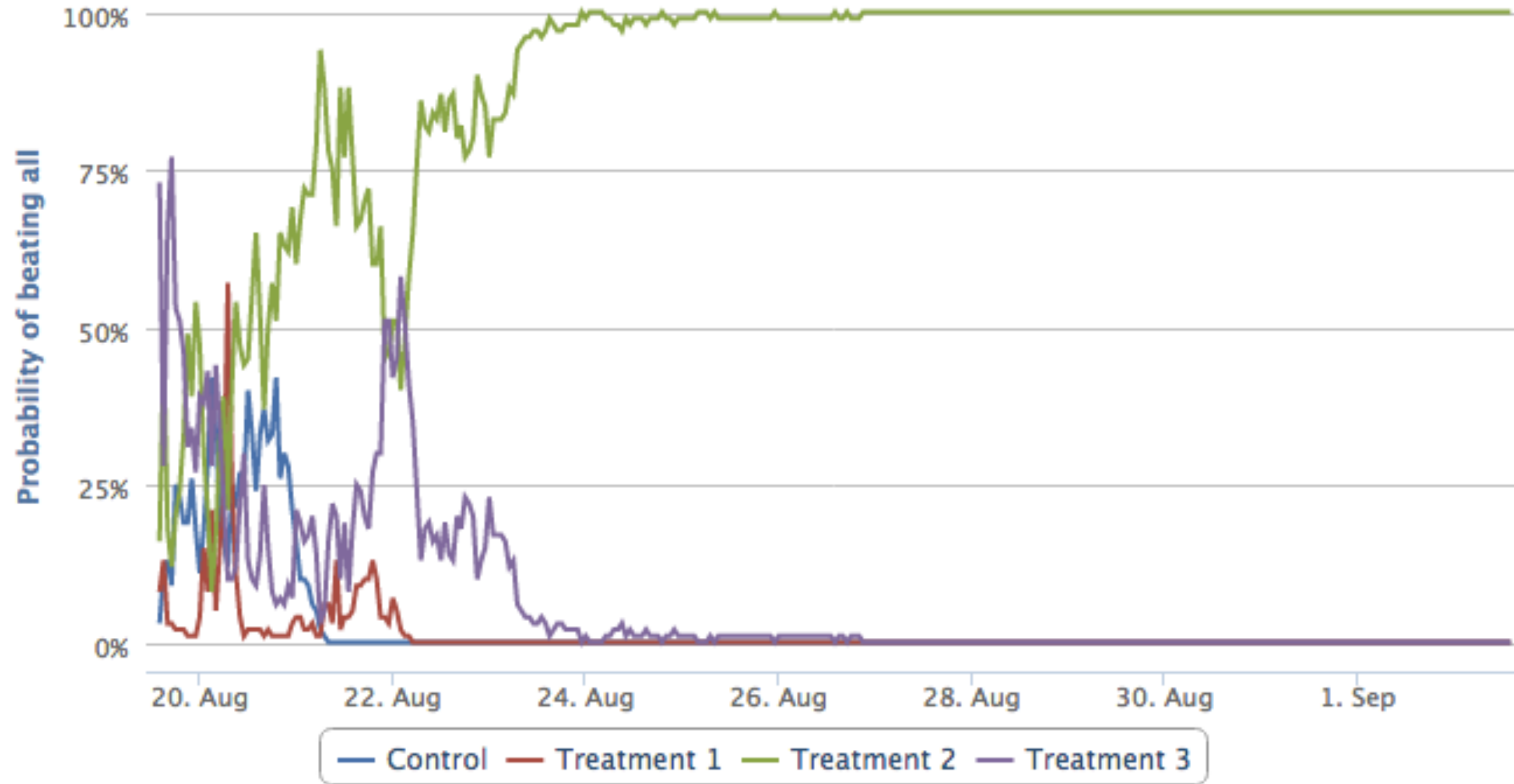
$a=42, b=94$



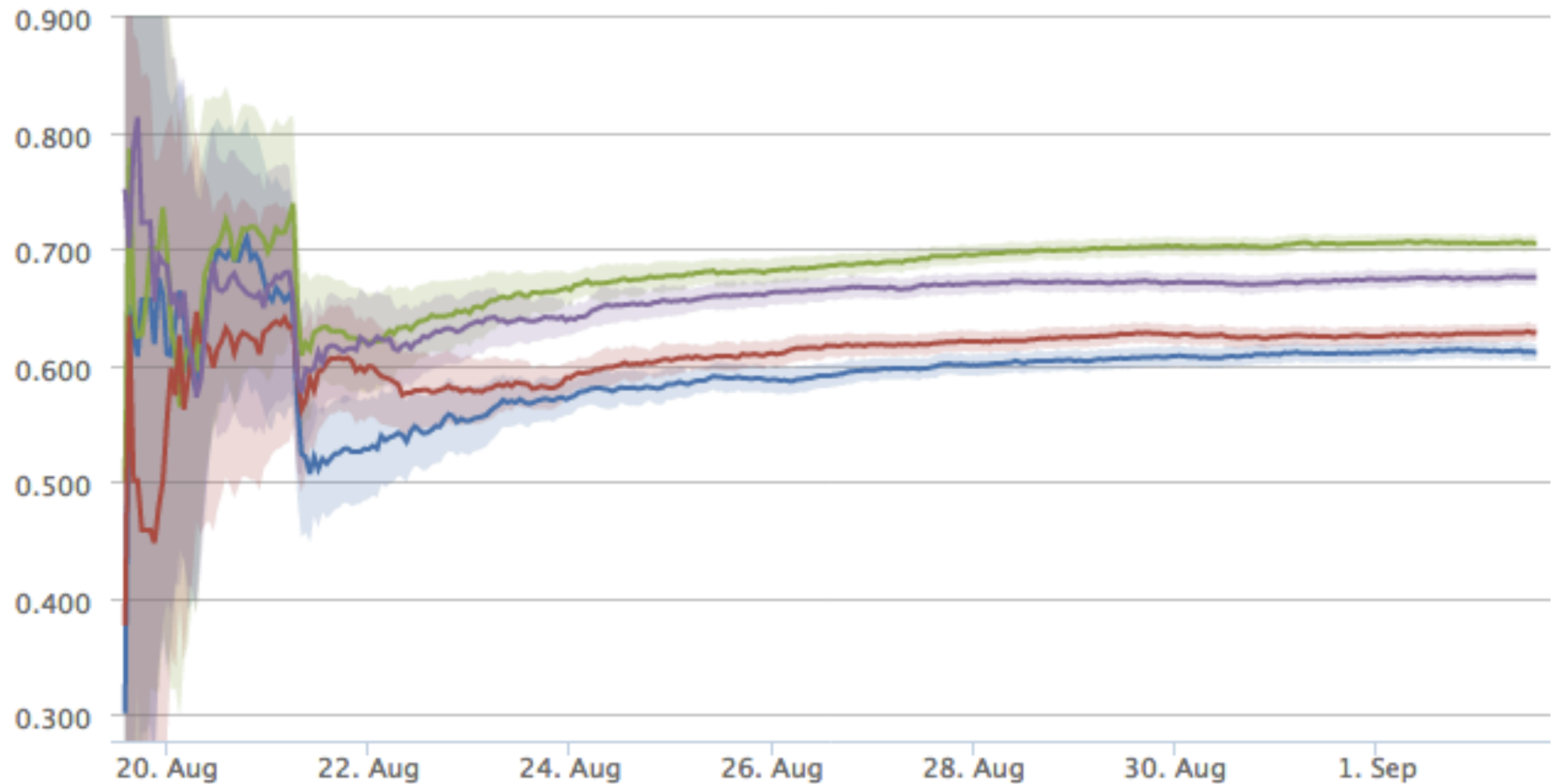
Some examples...

	Variant ?	Score ?	Change ?	Probability of beating control ?	Probability of beating all ?	Conversions / Participants ?
	Control	0.611			0% 🚫	6,870 / 11,243
	Treatment 1	0.6276	+2.71%	100% ✅	0% 🚫	7,037 / 11,212
🏆	Treatment 2	0.7044	+15.27%	100% ✅	100% ✅	7,955 / 11,294
	Treatment 3	0.6755	+10.55%	100% ✅	0% 🚫	7,616 / 11,274

A successful test



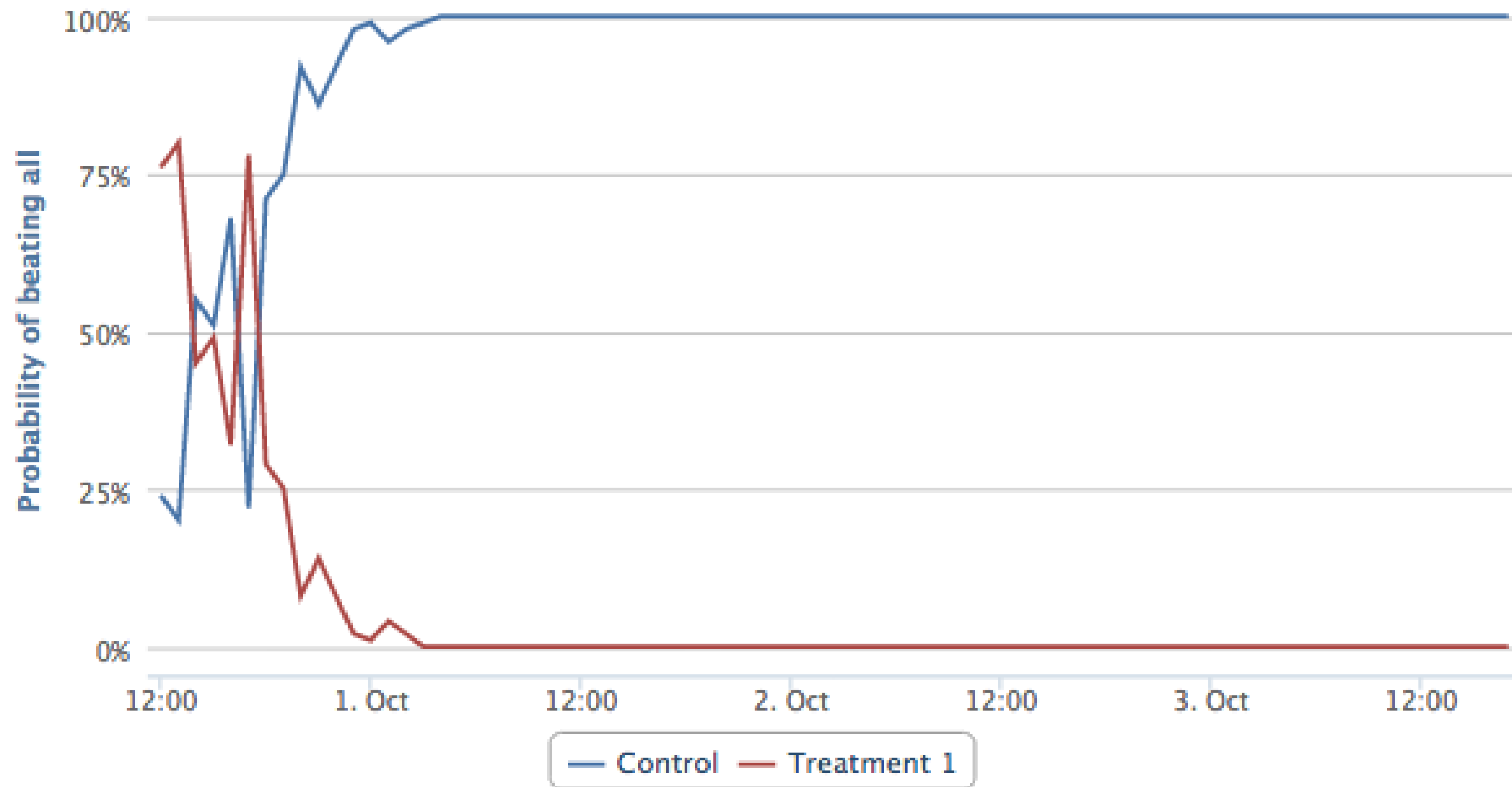
Probability of beating all



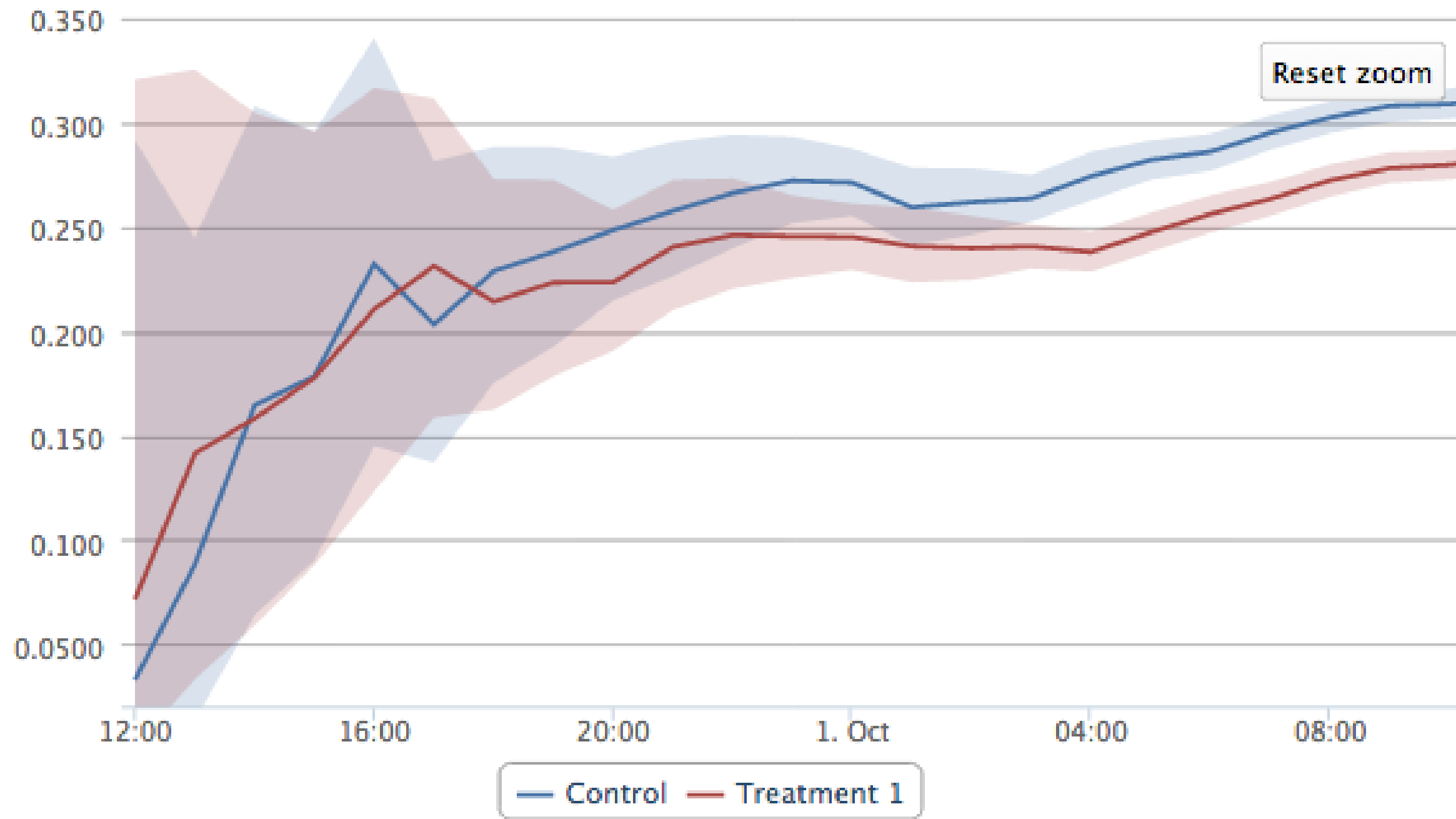
Observed conversion rates (with CI bounds)

	Variant ?	Score ?	Change ?	Probability of beating control ?	Probability of beating all ?	Conversions / Participants ?
★	Control	0.3774			100% ✓	15,567 / 41,244
	Treatment 1	0.3477	-7.88%	0% 🚫	0% 🚫	14,385 / 41,372

A ~~not~~ so successful test...



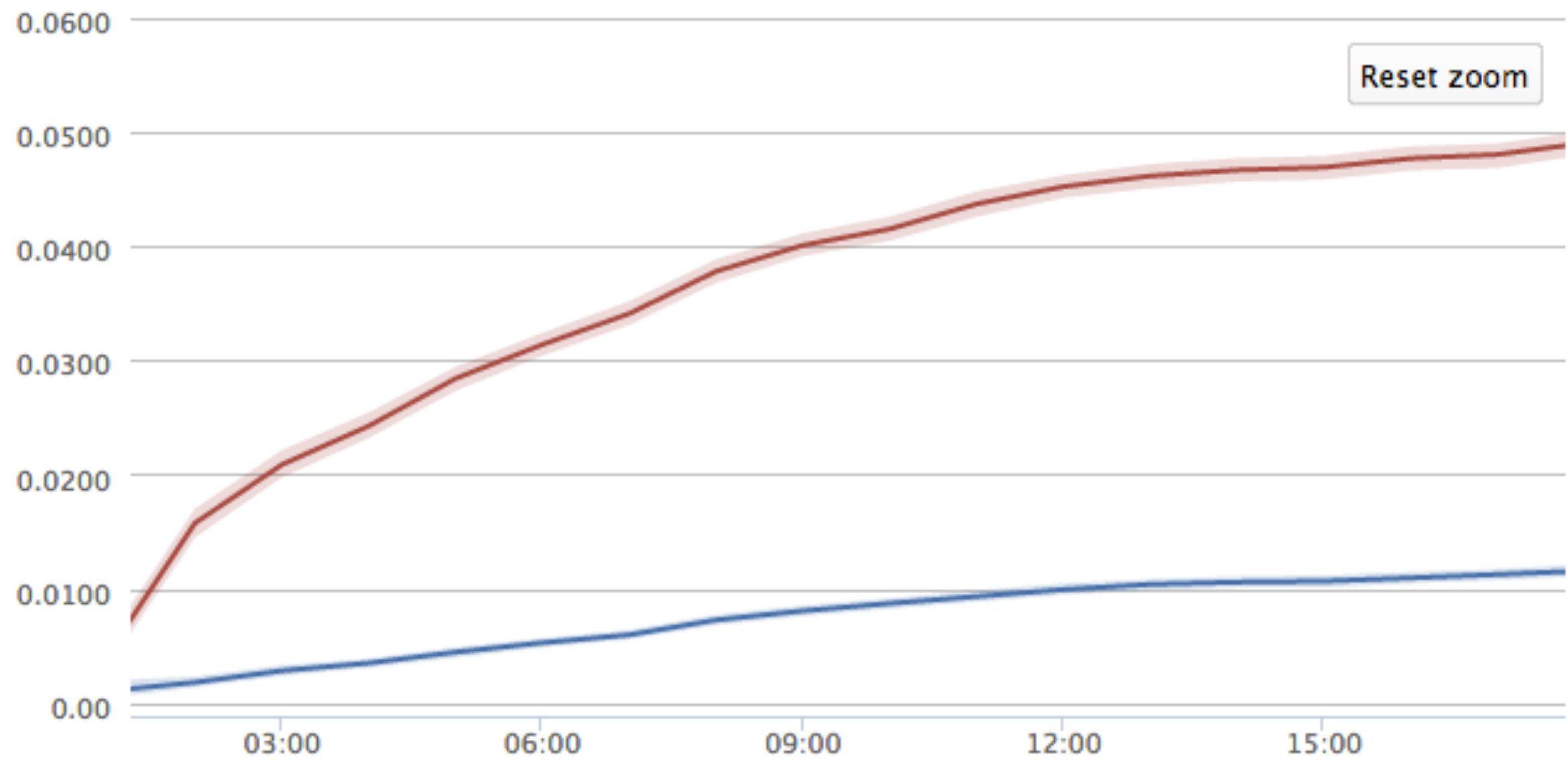
Probability of beating all



Observed conversion rate (posterior)

Assumptions

- Users are independent
- User's convert quickly (immediately)
- Probability of conversion is independent of time



Un-converged conversion rate

Benefits / Features

- Continuously observable
- No need to fix population size in advance
- Incorporate prior knowledge / expertise
- Result is a “true” probability
- A measure of the difference magnitude is given
- Consistent framework for lots of different scenarios

Useful Links

- <https://github.com/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers>
- “Doing Bayesian Data Analysis: A Tutorial with R and Bugs”, John K. Kruschke
- <http://www.evanmiller.org/how-not-to-run-an-ab-test.html>
- <http://www.kaushik.net> - Occam's Razor Blog
- <http://exp-platform.com> - Ron Kovahi et al.

Thanks

Steve@swrve.com
@stevec64

Multi-arm bandits

5%



6%



10%



4%



Pull 1



Pull 2

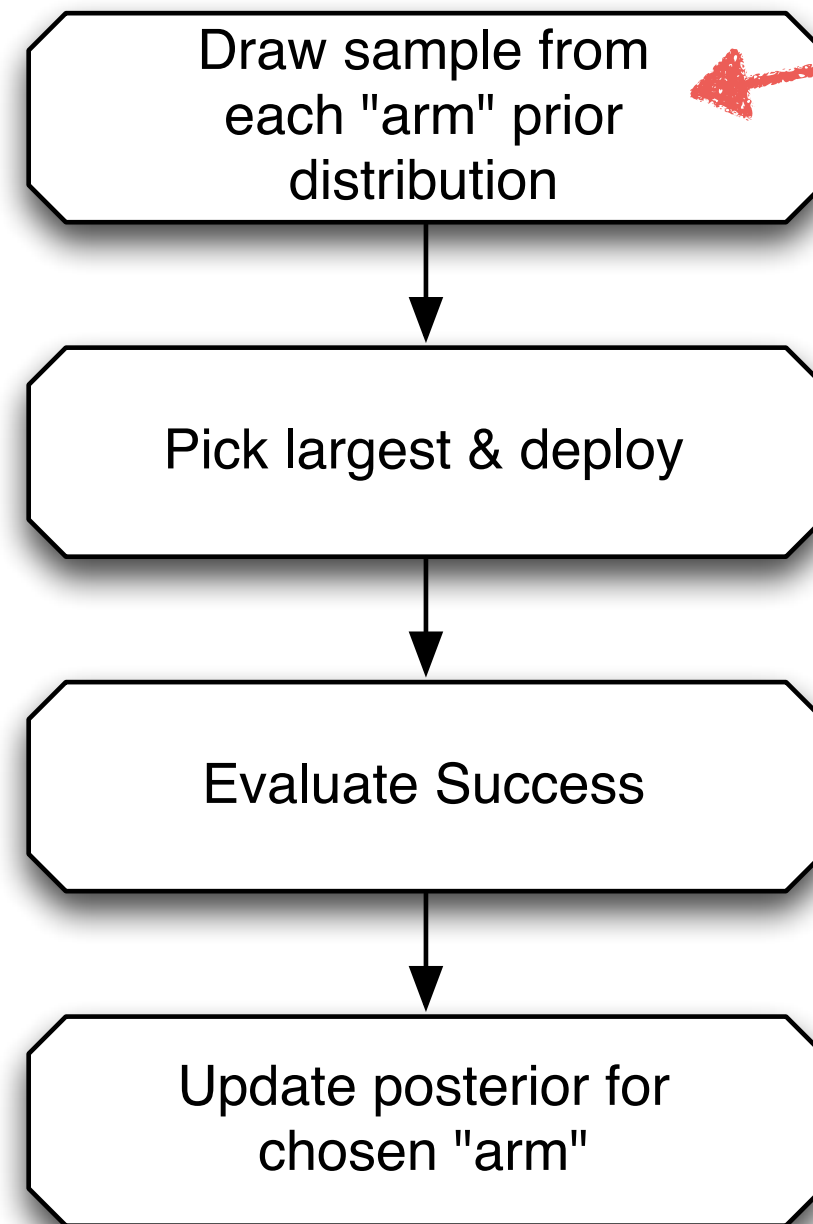


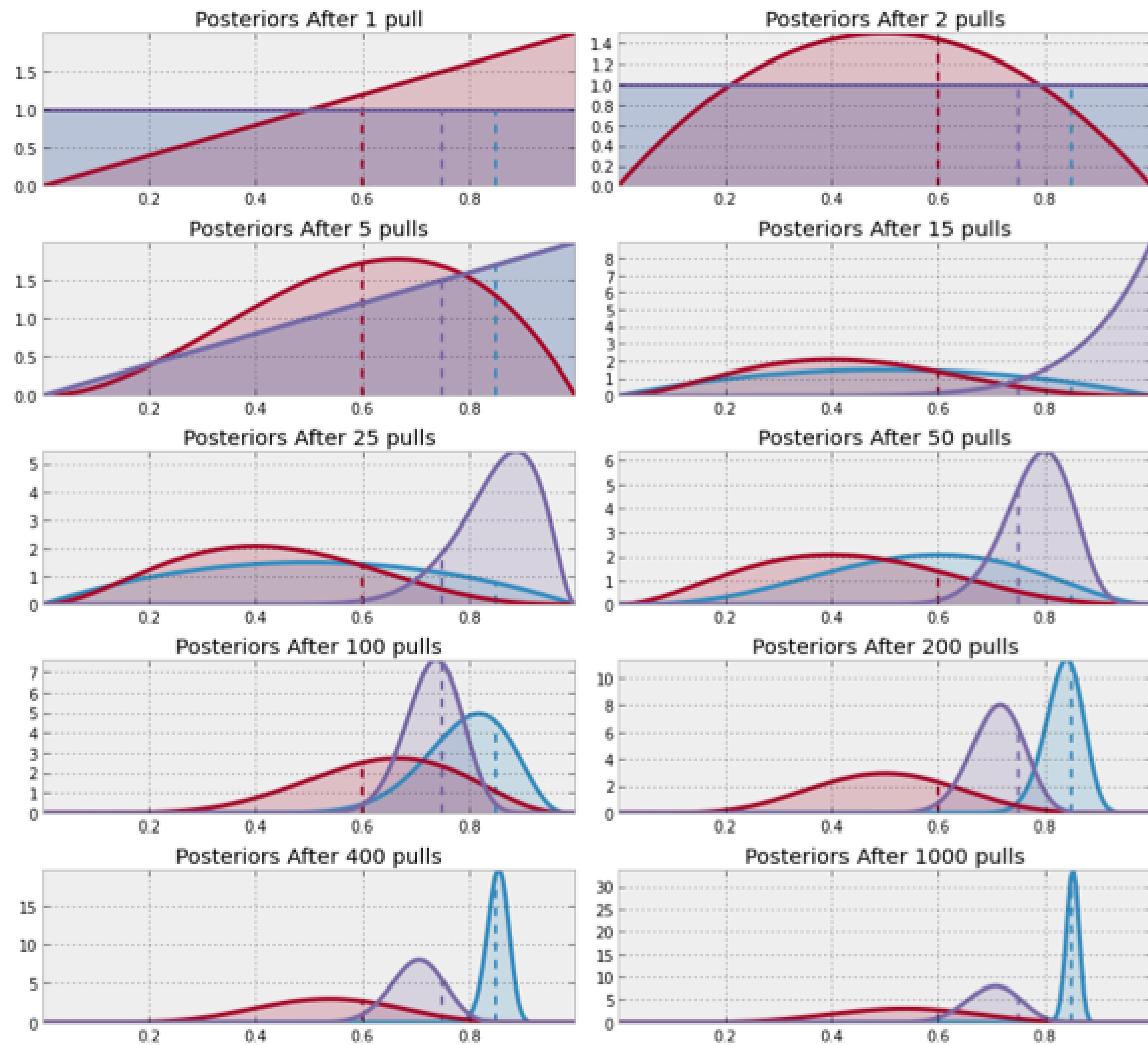
Pull 3



Multi-arm bandits

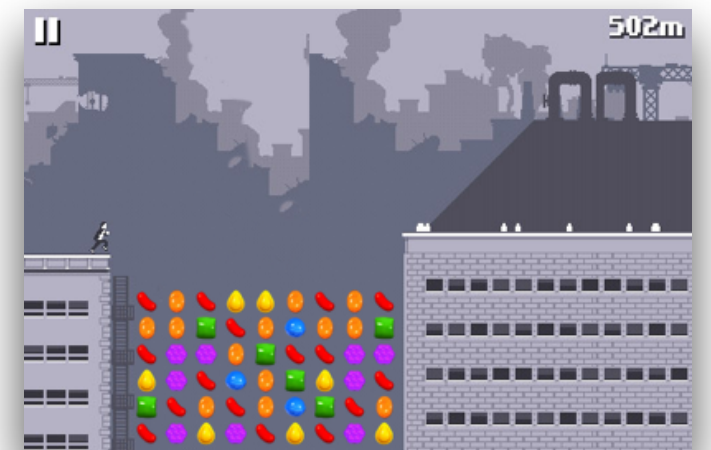
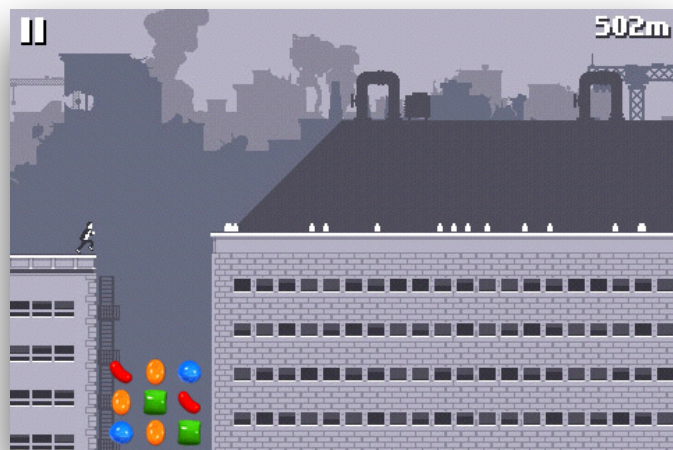
Thompson Sampling



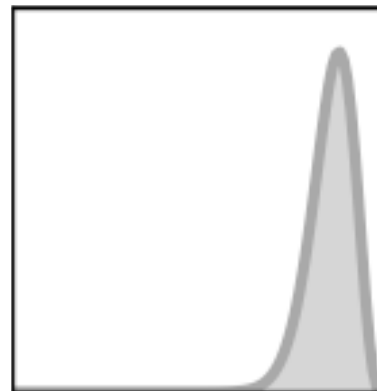
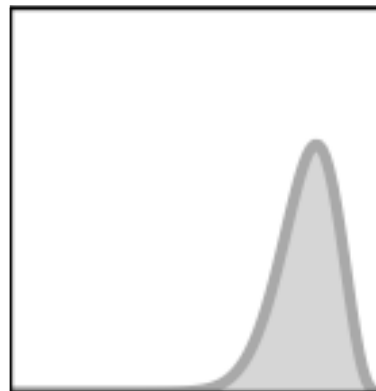
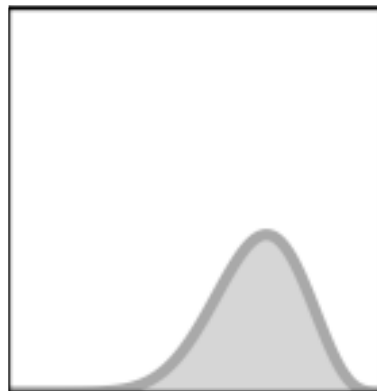
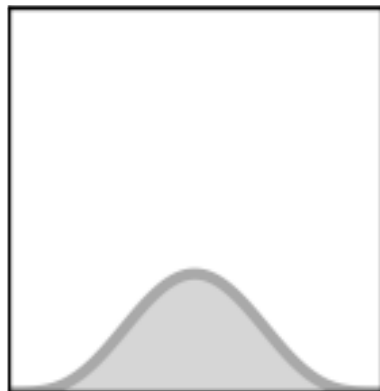
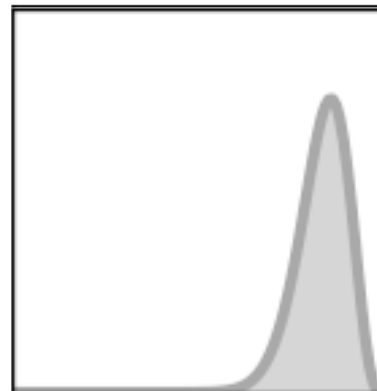
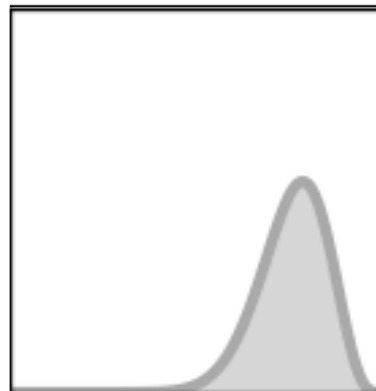
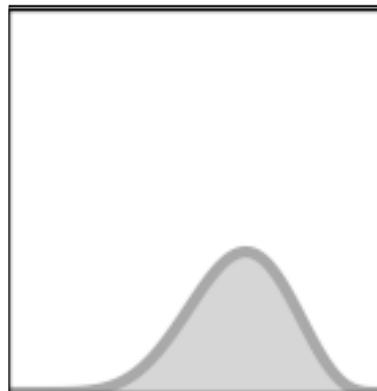
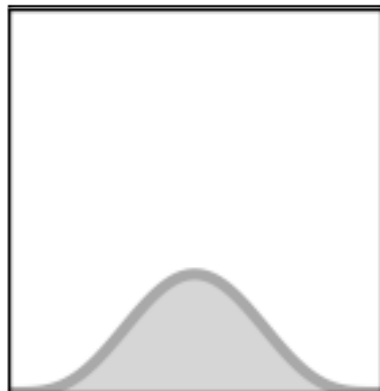
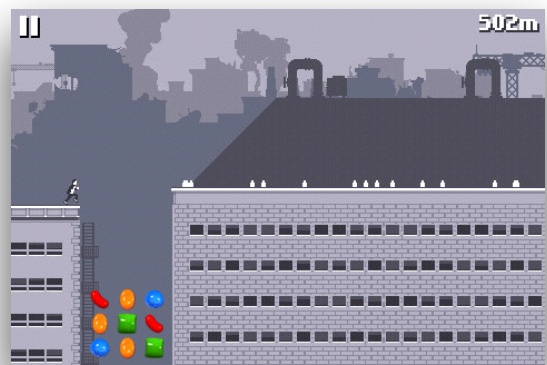


Difficulty tuning

a slight silly example...



Canacandycrushbalt



Winner

