

JSC 270 - LECTURE 4 DATA SIMULATION

<https://jsc270.github.io/>



ANNOUNCEMENTS

Today: Visualization Talk by Dr Fanny Chevalier 2-3pm

Next Monday: Talk on Reproducibility by Dr Benjamin Haibe-Kains also 2-3pm

Assignment grades will be available by Lab time

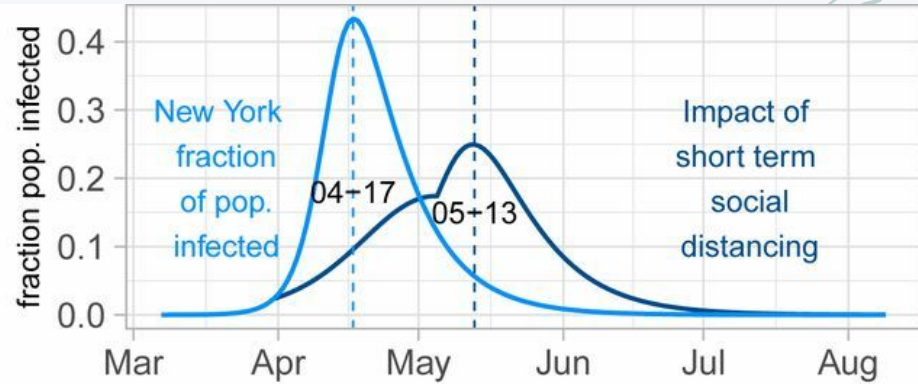
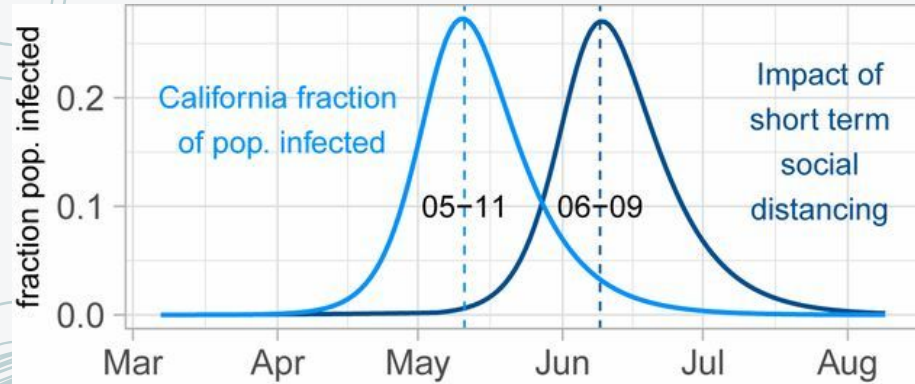
Fantastic job on reflection quizzes and perusall!
Grades for these will be available after class today



SIMULATIONS

Can you think of any examples yourselves?

EXAMPLE 1. COVID19 SPREAD

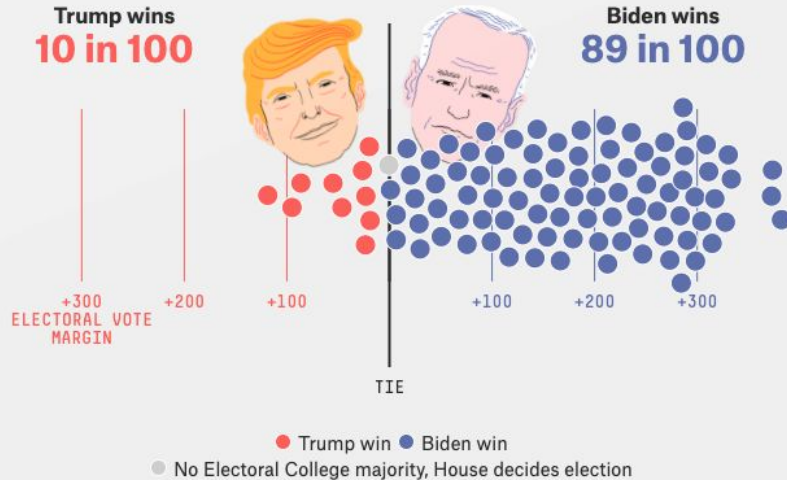


SIR Model. (Left) short term distancing (Right) longer term distancing

EXAMPLE 2. ELECTIONS

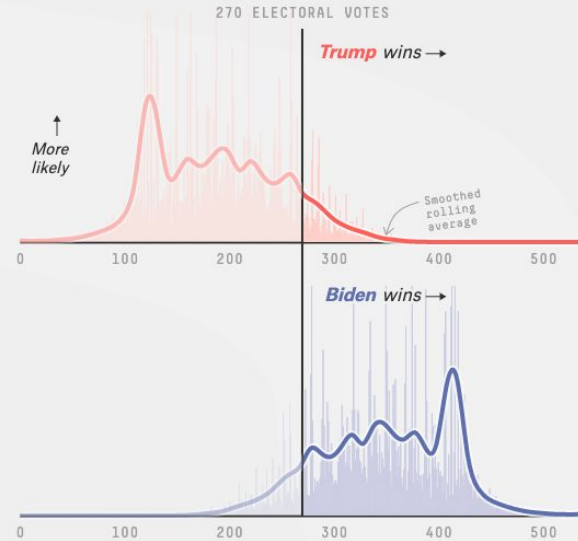
Biden is *favored* to win the election

We simulate the election 40,000 times to see who wins most often. The sample of 100 outcomes below gives you a good idea of the range of scenarios our model thinks is possible.



Every outcome in our simulations

All possible Electoral College outcomes for each candidate, with higher bars showing outcomes that appeared more often in our 40,000 simulations



EXAMPLE 3. DALL-E (OPEN AI, JAN 5, 2021)

GPT-3 based model

Generates images from English text from text-image pairs

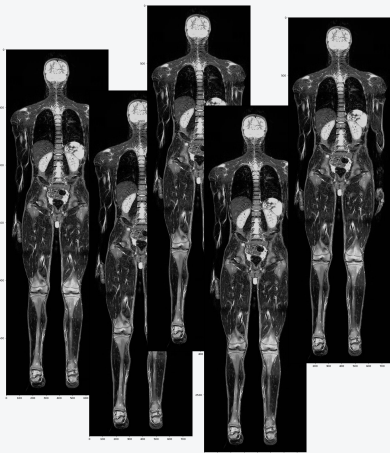
TEXT PROMPT

an armchair in the shape of an avocado [...]

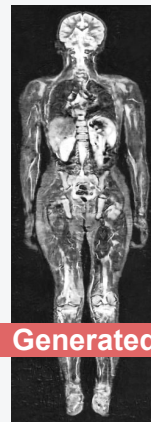
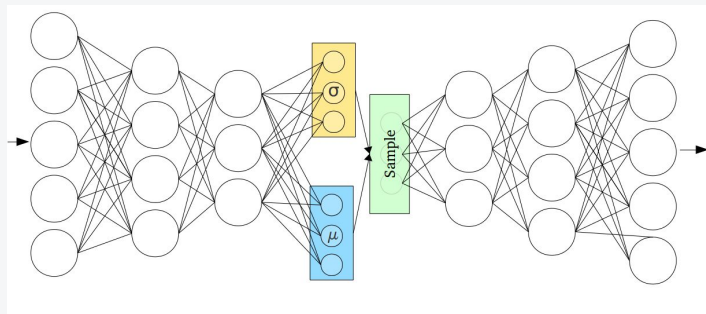
AI-GENERATED IMAGES



EXAMPLE 4. MEDICAL IMAGING



Healthy wbMRI



Generated

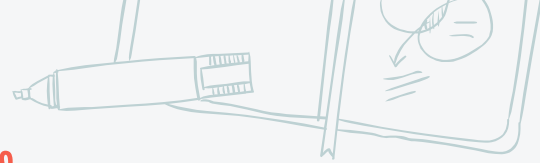
EXAMPLE 5. SOCIAL NETWORKS



<https://arxiv.org/pdf/0912.5410.pdf>

<http://reports-archive.adm.cs.cmu.edu/anon/ml/CMU-ML-06-107.pdf>

WHY WE DO SIMULATIONS?





WHY WE DO SIMULATIONS

- To think through the design of a study
- To test/verify a hypothesis
- To study the system (complexity and dynamics) in a controlled environment
- To augment data
- To generate new knowledge/discovery
- Educational/training purposes

WHEN ARE SIMULATIONS NOT APPROPRIATE?





WHEN ARE SIMULATIONS NOT APPROPRIATE?

- ✗ When the system is deterministic
- ✗ When not enough is known about the system such that the assumptions you would have to make are unlikely to be accurate



SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?



SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

What do you expect?



SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

How do we simulate this to study the outcome?



SIMULATING COIN FLIPS

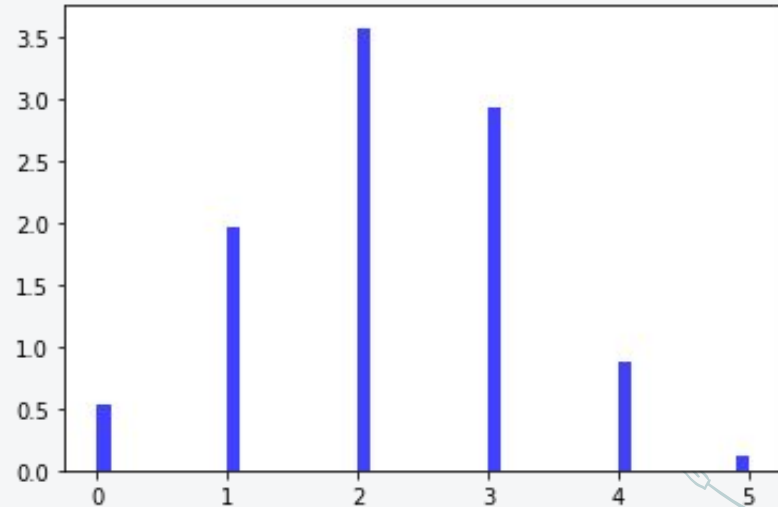
How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

$$X \sim \text{Bin}(5, 0.45)$$

SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

$$X \sim \text{Bin}(5, 0.45)$$

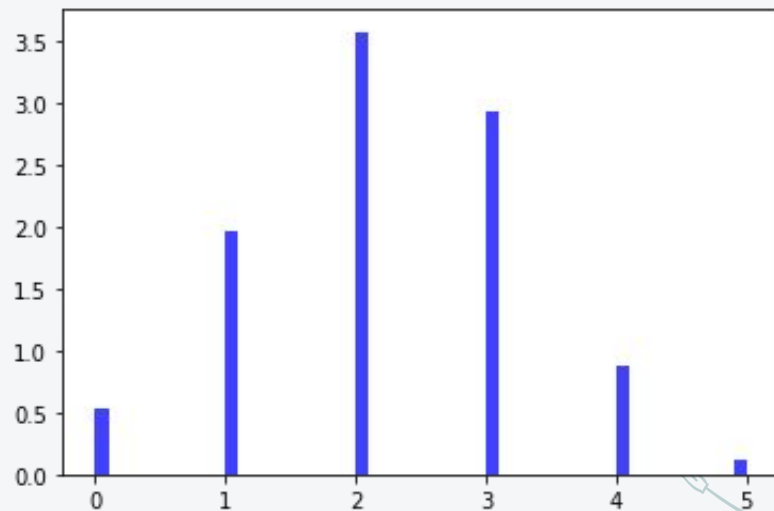


SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

$$X \sim \text{Bin}(5, 0.45)$$

How did we do this?



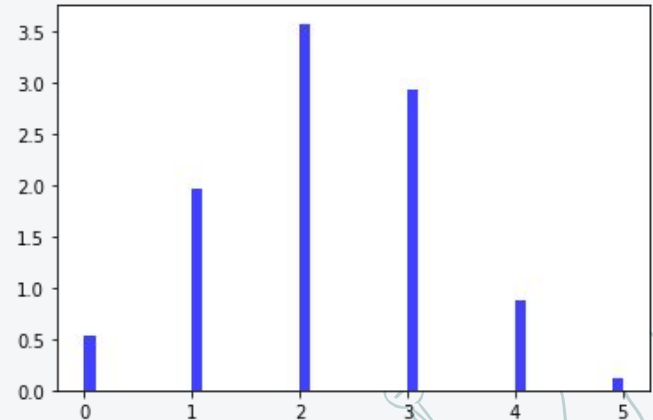
SIMULATING COIN FLIPS

How many heads are we likely to see if we throw 5 slightly unfair (0.45) coins?

$$X \sim \text{Bin}(5, 0.45)$$

How did we do this?

Sample from the binomial!





HOW EXACTLY DO WE SAMPLE FROM A GIVEN BINOMIAL DISTRIBUTION?

Assumption – there exists a uniform random number generator

Step 1. Generate 5 numbers from 0 to 1

Step 2. Record how many numbers fall below 0.45

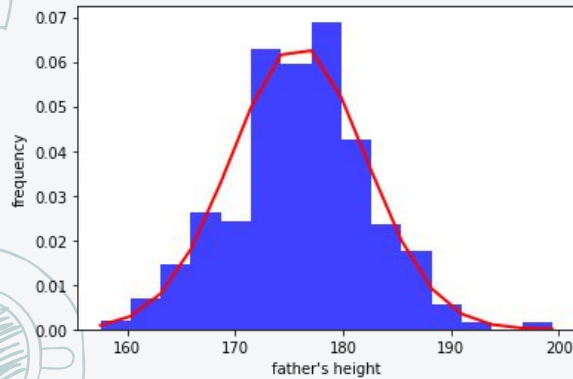
Step 3. Repeat (100/1000/... times)

GALTON'S DATA - AN EXAMPLE OF SAMPLING FROM A GAUSSIAN

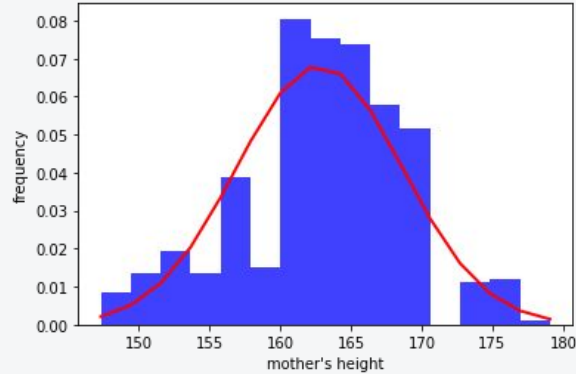


GENERATING NEW HEIGHT DATA BASED ON GALTON'S DATA

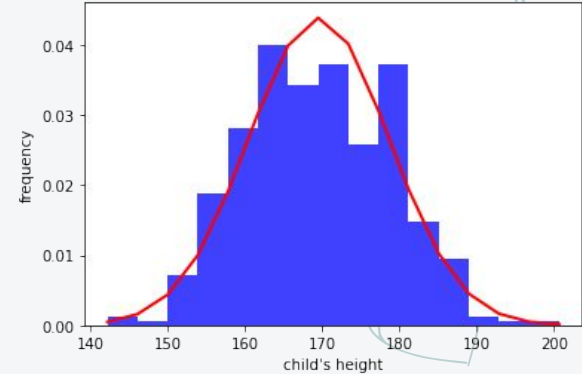
We have a set of heights from the population



$N(175.85, 6.3)$



$N(162.8, 5.9)$

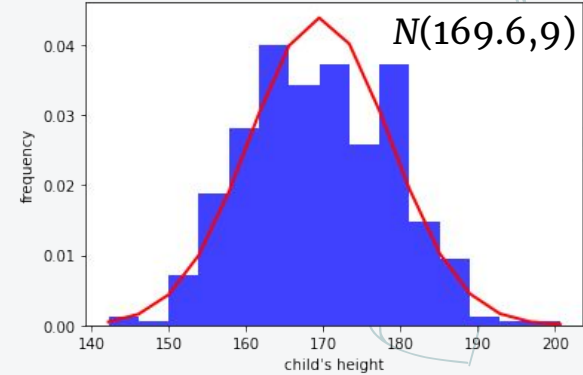
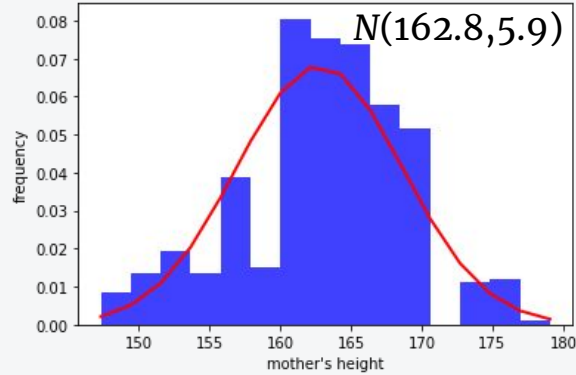
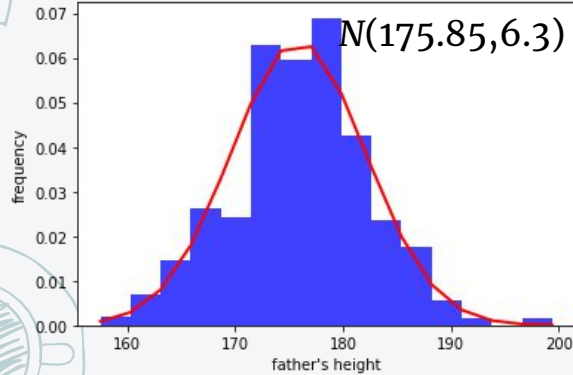


$N(169.6, 9)$

We want to generate another sample from this population (a new family). Can we just sample from these three distributions?

GENERATING NEW HEIGHT DATA BASED ON GALTON'S DATA

We want to generate another sample from this population (a new family).
Can we just sample from these three distributions? Let's try....



Correlation
matrix

	father	mother	child
father	1	0	0
mother	0	1	0
child	0	0	1

GALTON'S DATA - AN EXAMPLE OF SAMPLING FROM A GAUSSIAN

Child's
height

$$Y = 56.67 + 0.38 \cdot FH + 0.28 \cdot MH$$

What is the right algorithm for generating info about a set of new families?

Step 1.

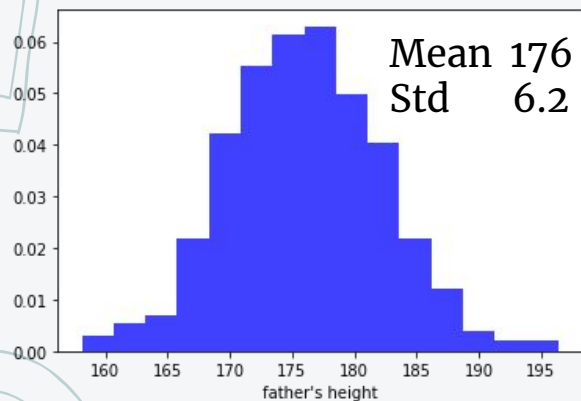
Step 2.

Step 3.

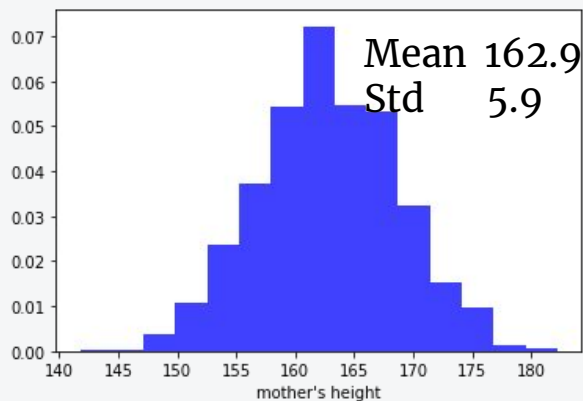
Step 4.

LET'S SIMULATE DATA ACCORDING TO THIS ALGORITHM

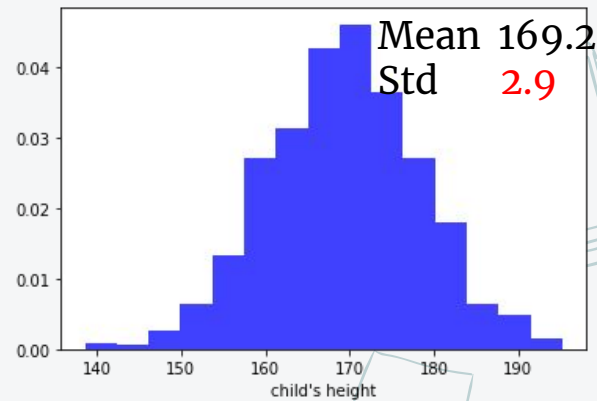
$N(175.85, 6.3)$



$N(162.8, 5.9)$



$N(169.6, 9)$



What could be the reason(s) for the distribution of children heights to deviate?

GALTON'S DATA - AN EXAMPLE OF SAMPLING FROM A GAUSSIAN

Child's
height

$$Y = 56.67 + 0.38 \cdot FH + 0.28 \cdot MH + \epsilon$$

$$\epsilon \sim N(0, \sigma^2_{(y-\hat{y})}) - \text{Variance of the residuals of the model!}$$

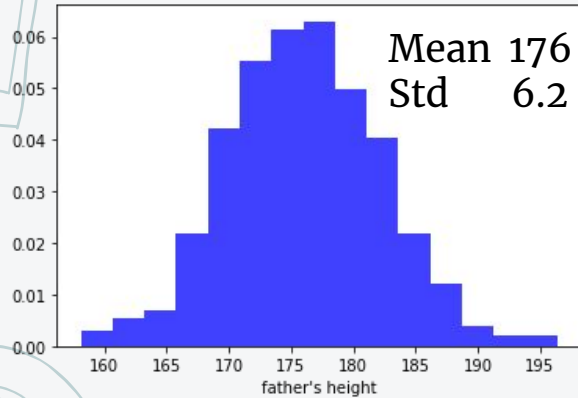
8.8

Don't forget that our model was not perfect!

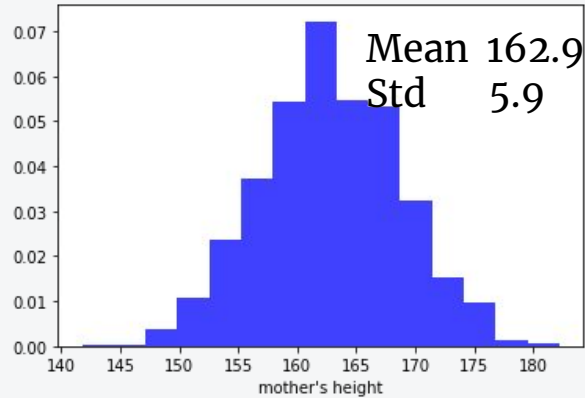
As models almost never are...

A NEW SAMPLE FROM GALTON'S POPULATION

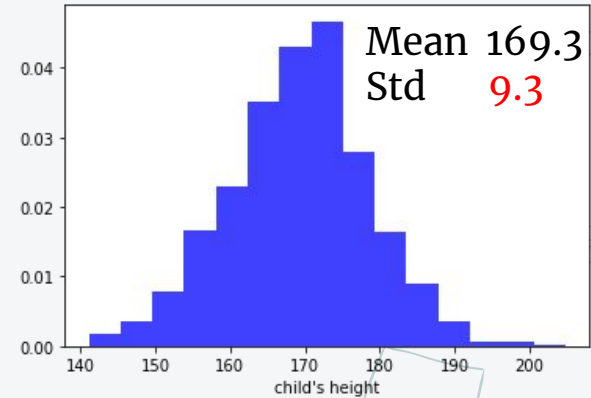
$N(175.85, 6.3)$



$N(162.8, 5.9)$



$N(169.6, 9)$



Don't forget to simulate realistic data (i.e. including the noise) using regression!



COVID19 EXAMPLE

Suppose we know that people are at different levels of immunity, in fact we have four subpopulations with immunities in the 0 to 1 range

How would you simulate data from this population?

COVID19 EXAMPLE

Suppose we know that people are at different levels of immunity, in fact we have four subpopulations with immunities are in the 0 to 1 range

How would you simulate data from this population?

1. Average Immunity per population $\sim \text{Uniform}(0,1)$
2. For simplicity, let's fix the standard dev to be 1
3. Sample each subpopulation (1,000 people) from a $\text{Beta}(\alpha, \beta)$ with the params indicated above

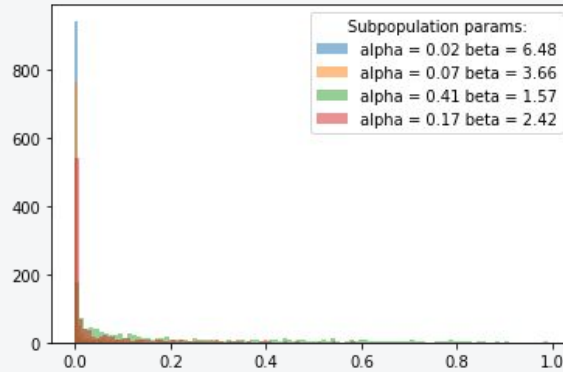
COVID19 IMMUNITY SAMPLED DATA

$\mu \sim \text{Uni}(0,1,4) : [0.15 \ 0.27 \ 0.64 \ 0.41], \ \sigma = 1$

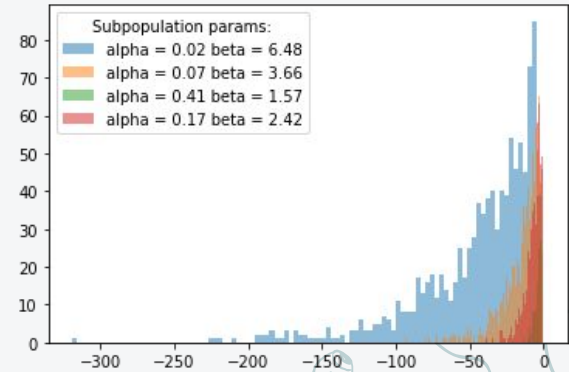
Sample each subpopulation j from Beta:

$$\alpha_j = \mu_j^2 / \sigma_j^2$$

$$\beta_j = \sigma_j^2 / \mu$$



Immunity in various populations



Log(immunity)

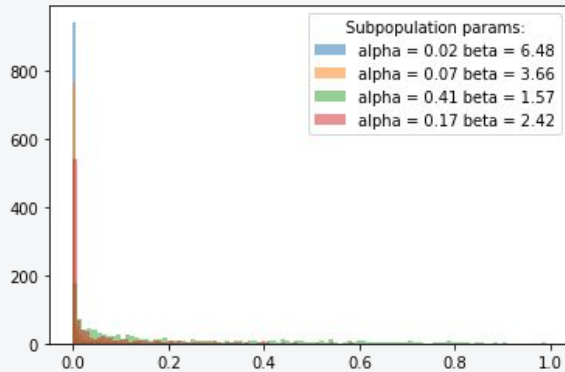
COVID19 IMMUNITY SAMPLED DATA - HIERARCHICAL MODEL

$\mu \sim \text{Uni}(0,1,4) : [0.15 \ 0.27 \ 0.64 \ 0.41], \ \sigma = 1$

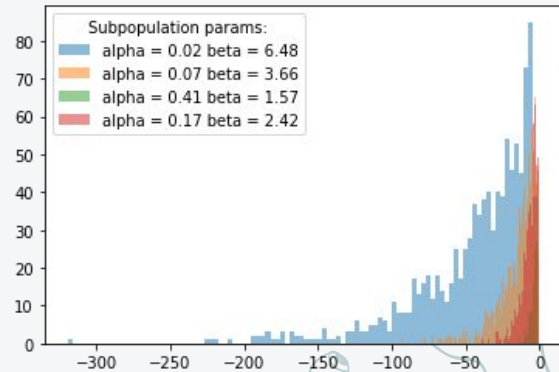
Sample each subpopulation j from Beta:

$$\alpha_j = \mu_j^2 / \sigma_j^2$$

$$\beta_j = \sigma_j^2 / \mu$$



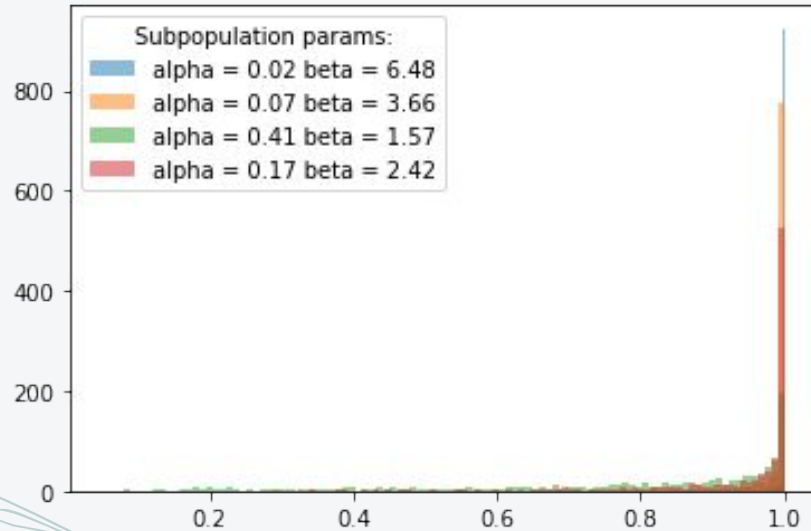
Immunity in various populations

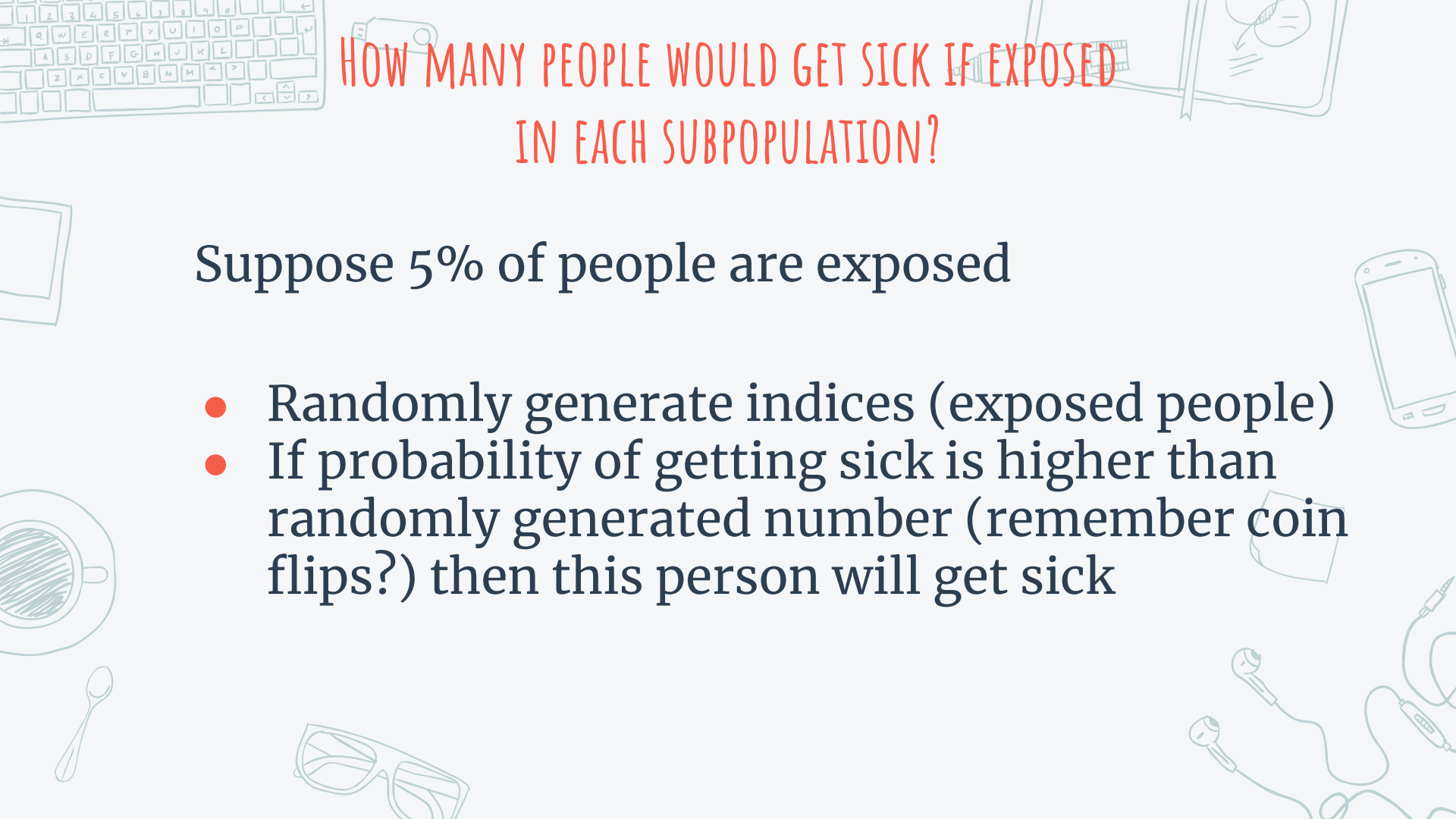


Log(immunity)

COVID19 THE PROBABILITY OF A PERSON BEING INFECTED

The probability of a person getting covid in various populations is $1 - \text{immunity}$





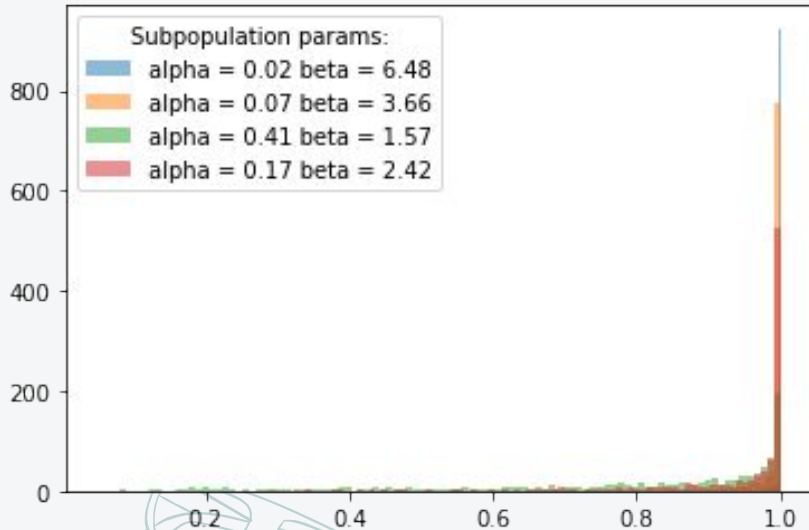
HOW MANY PEOPLE WOULD GET SICK IF EXPOSED IN EACH SUBPOPULATION?

Suppose 5% of people are exposed

- Randomly generate indices (exposed people)
- If probability of getting sick is higher than randomly generated number (remember coin flips?) then this person will get sick

COVID19 NUMBER OF PEOPLE IN EACH SUBPOPULATION IF 5% EXPOSED (TOTAL 1000 PER SUBPOPULATION)

The probability of a person getting covid in various populations is $1 - \text{immunity}$



Subpopulation 1: 50/50

Subpopulation 2: 48/50

Subpopulation 3: 42/50

Subpopulation 4: 47/50



SUMMARY

- Saw a few examples of data simulations
- Should be able to articulate advantages and disadvantages of data simulations
- Should be able to simulate data from
 - Binomial (e.g. coin flips)
 - Gaussian
 - Regression
 - Hierarchical model



READING

Original paper introducing Transformer architecture - foundation for AlphaFold2 and GPT-3/DALL-E

<https://arxiv.org/pdf/1706.03762.pdf>

Nice example of Poisson modeled bus waiting times:

<https://jakevdp.github.io/blog/2018/09/13/waiting-time-paradox/>



Extra Slides



SUPPOSE YOU DIDN'T KNOW THE MEAN

Example: We need to simulate a dataset of heights y across various places on earth, but we only know world average height and world standard deviation

What would you do?



SUPPOSE YOU DIDN'T KNOW THE MEAN

What would you do?

1. Make it up based on what you are trying to model
2. Be a little more Bayesian...

BEING BAYESIAN - HIERARCHICAL MODEL

l - location

w - world

$$h_l \sim N(\mu_l, \sigma_l^2)$$

$$\mu_l \sim N(\mu_w, \sigma_w^2) \quad - \quad \text{prior}$$

How to generate data?

1. Sample mean for a location
2. Given the mean, sample subpopulation



BEING BAYESIAN ABOUT HEIGHTS

World height average for women is 165cm
World height std for women is 8.9cm

$$h_l \sim N(\mu_l, \sigma_l^2)$$
$$\mu_l \sim N(\mu_w, \sigma_w^2)$$



BEING BAYESIAN ABOUT HEIGHTS

$$h_l \sim N(\mu_l, \sigma_l^2)$$
$$\mu_l \sim N(\mu_w, \sigma_w^2)$$

World height average for women is 165cm

World height std for women is 8.9cm

Height averages: [178.37 173.76 156 169.69 155.38]

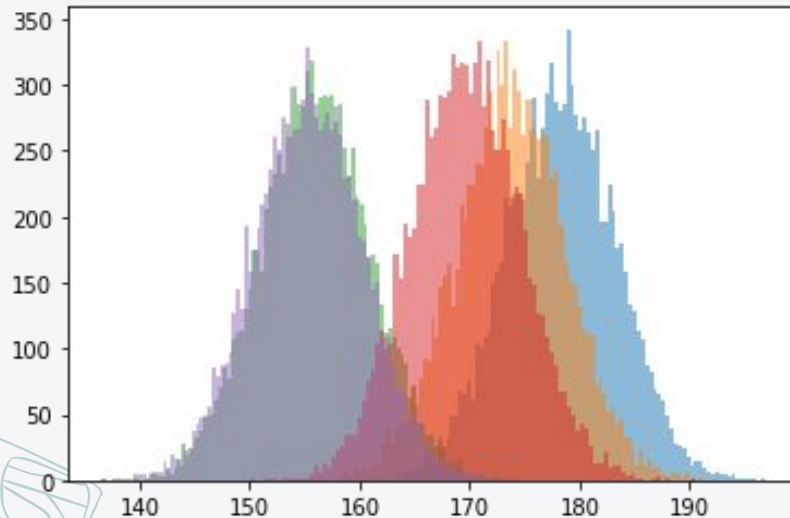
BEING BAYESIAN ABOUT HEIGHTS

$$h_l \sim N(\mu_l, \sigma_l^2)$$
$$\mu_l \sim N(\mu_w, \sigma_w^2)$$

World height average for women is 165cm

World height std for women is 8.9cm

Height averages: [178.37 173.76 156 169.69 155.38]



BEING BAYESIAN

l - location

w - world

$$h_l \sim N(\mu_l, \sigma_l^2)$$

$$\mu_l \sim N(\mu_w, \sigma_w^2) \quad - \quad \text{prior}$$

How to generate data?

1. Sample mean for a location
2. Given the mean, sample subpopulation

Graphical Model
notation



SHOULD BE ABLE TO SAMPLE FROM MANY MODELS NOW!

Example: Latent Dirichlet Allocation (LDA) – inferring topics in a set of publications

We generate the set of publications as follows:

1. Generate the topic distribution for document i $\theta_i \sim \text{Dir}(\alpha)$
2. For each topic k , sample word distribution $\phi_k \sim \text{Dir}(\beta)$
3. For each of the positions for the word j in document i
 - a. Choose topic $z_{ij} \sim \text{Multinomial}(\theta_i)$
 - b. Choose a word $w_{ij} \sim \text{Multinomial}(\phi_{z_{ij}})$