

# Can a Machine Be Racist or Sexist?

*On Social Bias in Machine Learning*

Renée M. P. Teate

HelioCampus

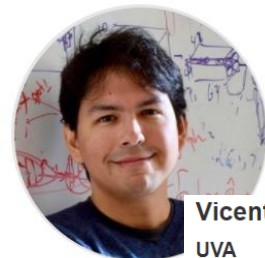
1. This talk: Introducing concepts and examples of social bias in machine learning to get us all on the same page (~30 mins)
2. Panel Discussion w/ Audience Q&A (~20 mins)



**Emily Crose**  
Undisclosed  
Threat Hunter

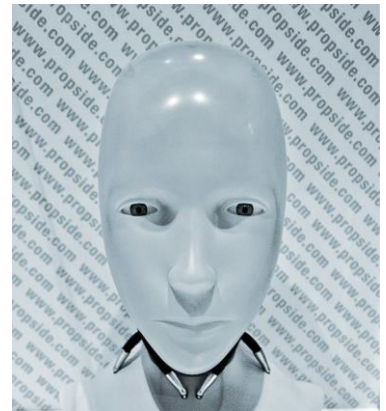


**Ines Montani**  
Explosion AI  
Founder



**Vicente Ordonez**  
UVA  
Assistant Professor

What comes to mind  
for most people  
when they are asked about  
their fears related to  
“Artificial Intelligence”  
or “Machine Learning”?



[Los Exterminadores De Skynet \(Terminator\)](#) by Hersion Piratoba on Flickr  
[Ex Machina](#) from film affinity  
[EXPOSIFY I. Robot](#) by Urko Dorronsoro from Donostia via Wikimedia Commons

So what is ***already*** going on  
with AI and Machine Learning  
that ***should*** concern us?

And are the impacted people and  
communities aware of what's already  
happening?

Are the people who design these systems  
aware of the possible impacts of their work  
on people's lives as they  
design and deploy data products?

“But if we take humans out of the loop and leave decisions up to computers, won’t it reduce the problems inherent in human decision-making?”

Can a Machine Be  
Racist or Sexist?

Can a Machine Learning Model  
Be ***Trained*** to Be  
Racist or Sexist  
(or made biased or unjust  
in other ways --  
***intentionally or not***)?

Let's define

# Machine [Algorithm]

“a step-by-step procedure for solving a problem”

Merriam-Webster Dictionary

# Racism

“racial prejudice or discrimination”

“a belief that race is the primary determinant of human traits and capacities and that **racial differences produce an inherent superiority** of a particular race”

Merriam-Webster Dictionary

# Racism

“a doctrine or political program based on the assumption of racism and *designed to execute its principles*”

[or, not designed NOT to execute its principles!]

# Sexism

“prejudice or discrimination based on sex”

“behavior, conditions, or attitudes that foster **stereotypes of social roles based on sex**”

Merriam-Webster Dictionary

# Institutional or Systemic Racism and Sexism

a **system** that **codifies** and **perpetuates discrimination**  
against individuals or communities  
based on their race or sex

*(note: these systems are designed/engineered by people)*

## Statuses Protected by Laws in the U.S.

- Race
- Sex
- Religion
- National Origin
- Age
- Disability Status
- Pregnancy
- Citizenship
- Familial Status
- Veteran Status
- Genetic Information



## Example: Bank Loans Before Machine Learning

*Bank officer* deciding whether to give a loan, assessing **likelihood to repay**:

- Employment Status and History
- Amount of Debt and Payment History
- Income & Assets
- “Personal Character”
- Co-Signer
- References
- Credit Score
  - Based on amount of debt, credit card payment history, debt-credit ratio etc.
  - May seem fair, but remember things like on-time payment of rent not included
  - Feedback loop - no/bad credit history, can't get credit, high interest, can't improve credit score
  - Is somewhat transparent, and errors can be corrected

## Example: Bank Loans With Machine Learning

*Algorithm* assessing **likelihood to repay**:

- |  |                                      |
|--|--------------------------------------|
| ● Employment Status and History  | ● Where You Live                     |
| ● Amount of Debt and Payment History   | ● Social Media Usage                 |
| ● Income & Assets  | ● Time You Wake Up                   |
| ● <del>“Personal Character”</del>  | ● Workout Consistency                |
| ● Co-Signer  | ● Driving Habits                     |
| ● <del>References</del>  | ● Time Spent Playing Video Games     |
| ● Credit Score   | ● Favorite Music                     |
| ● Detailed Spending Habits <ul style="list-style-type: none"> <li>○ Expenditures per month</li> <li>○ Where you shop</li> <li>○ Bill payment patterns</li> </ul> | ● Browser History                    |
|  | ● Facebook Friends' Financial Status |
|  | ● etc etc etc                        |

Is it **fair** for your interest rate, or whether you even get a loan offer, to be based on the default rates of “similar” people who, for instance, listen to the same kind of music as you?

What does it mean for decisions to become increasingly data-driven and automated?

We're still making the same types of decisions  
*(Who should receive funds from government programs? Who is at risk of dropping out of a university, and how do we intervene? What medical treatment should be applied based on a set of symptoms? Where should we locate the next branch of our business? etc etc),*

but now we're using much more data, and programming computers to help us find patterns from datasets larger than humans could sensibly process.

If designed well,  
machine learning systems  
can improve our world!

Better more targeted answers faster!

*More efficient use of taxpayer dollars, students receiving financial aid and intervention tutoring to help keep them in school, highly customized medical treatments, lower-risk business decisions!*

But we have to keep in mind that now:

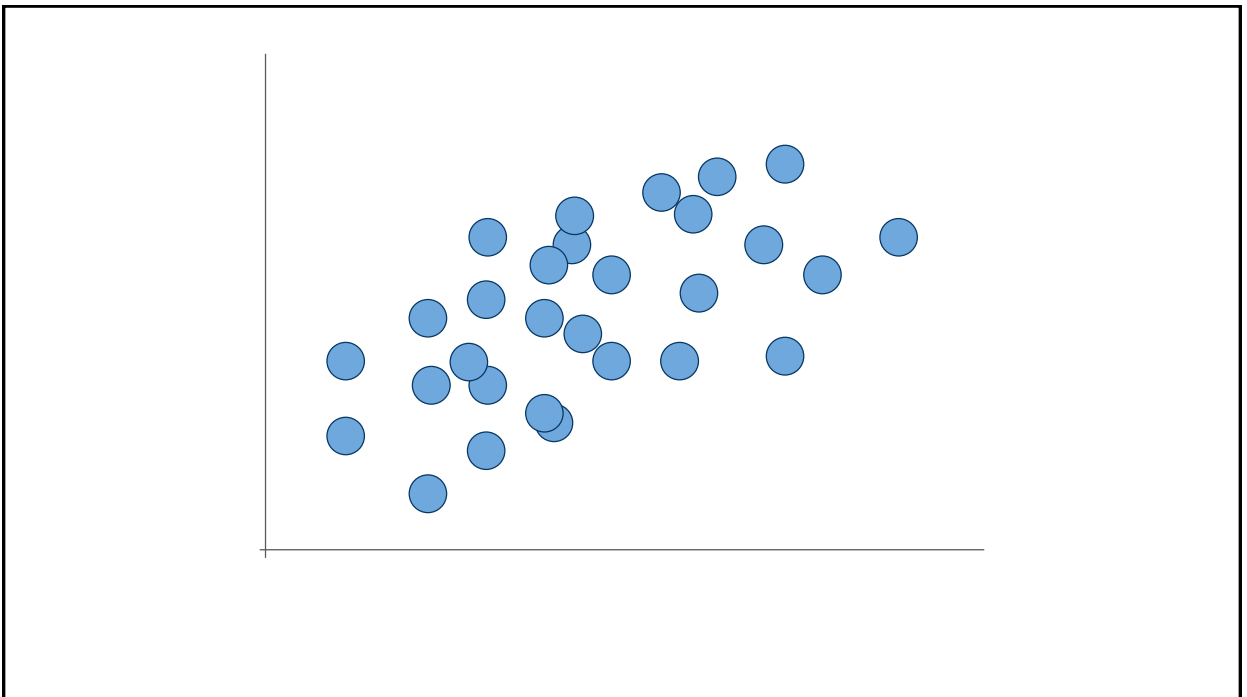
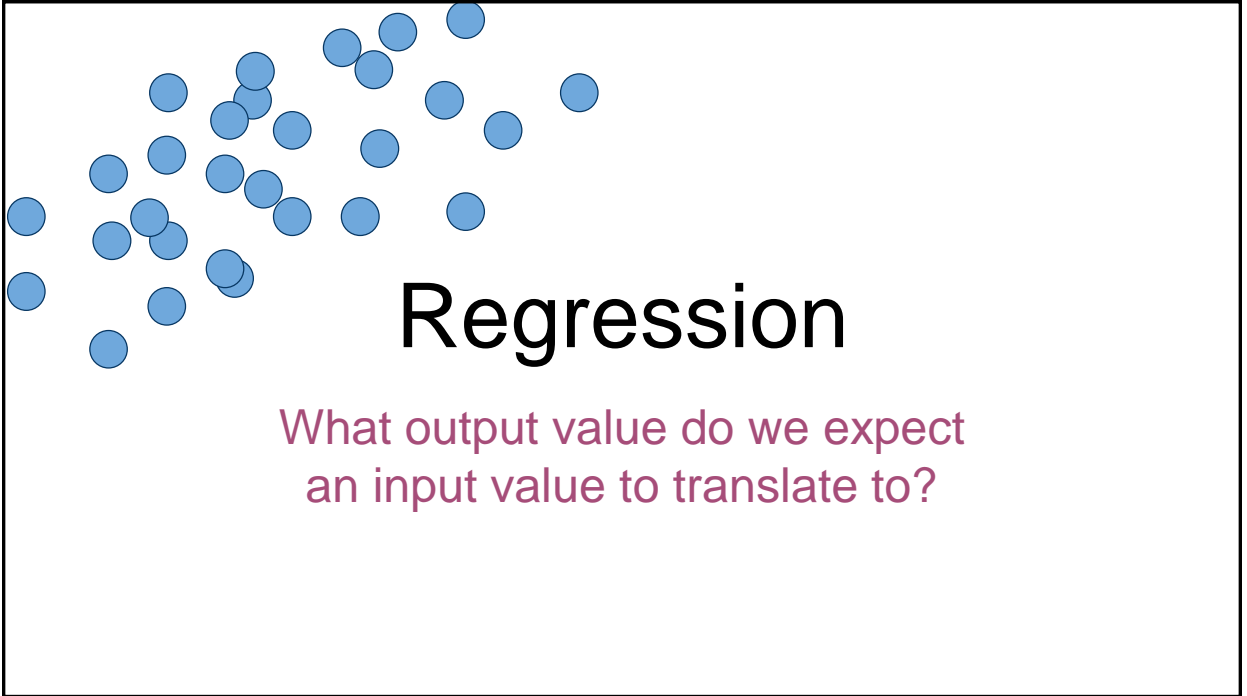
**“...we have the potential to make bad decisions far more quickly, efficiently, and with far greater impact than we did in the past”**

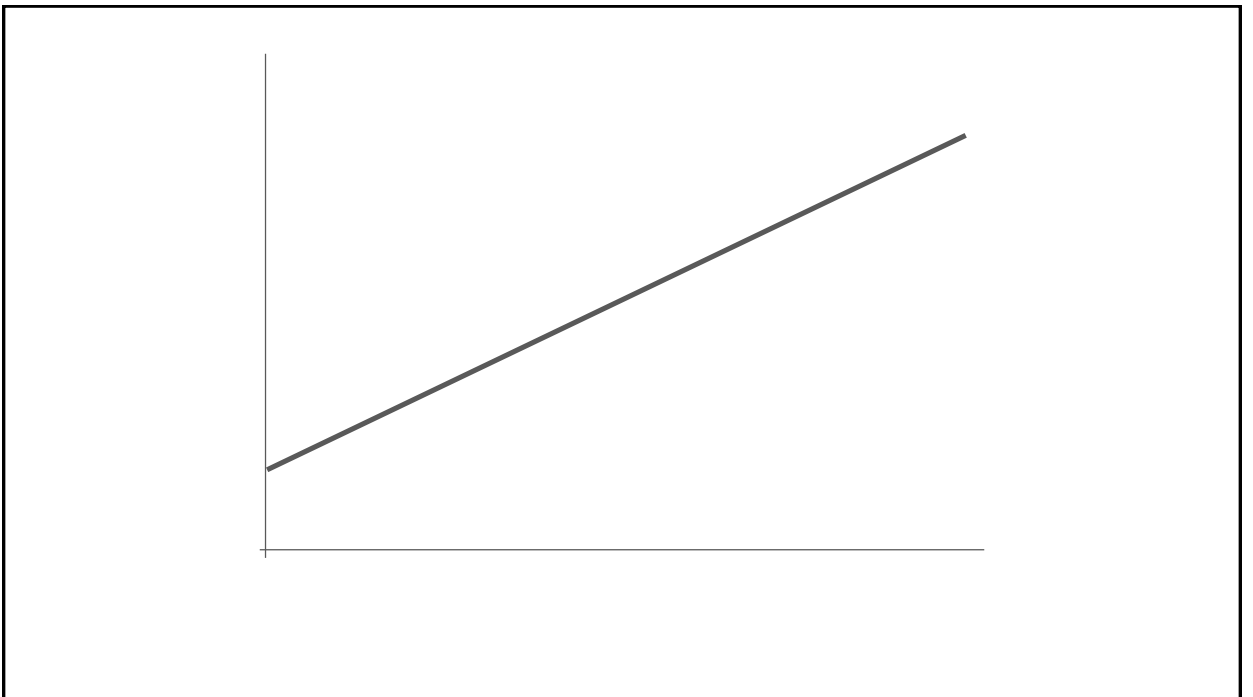
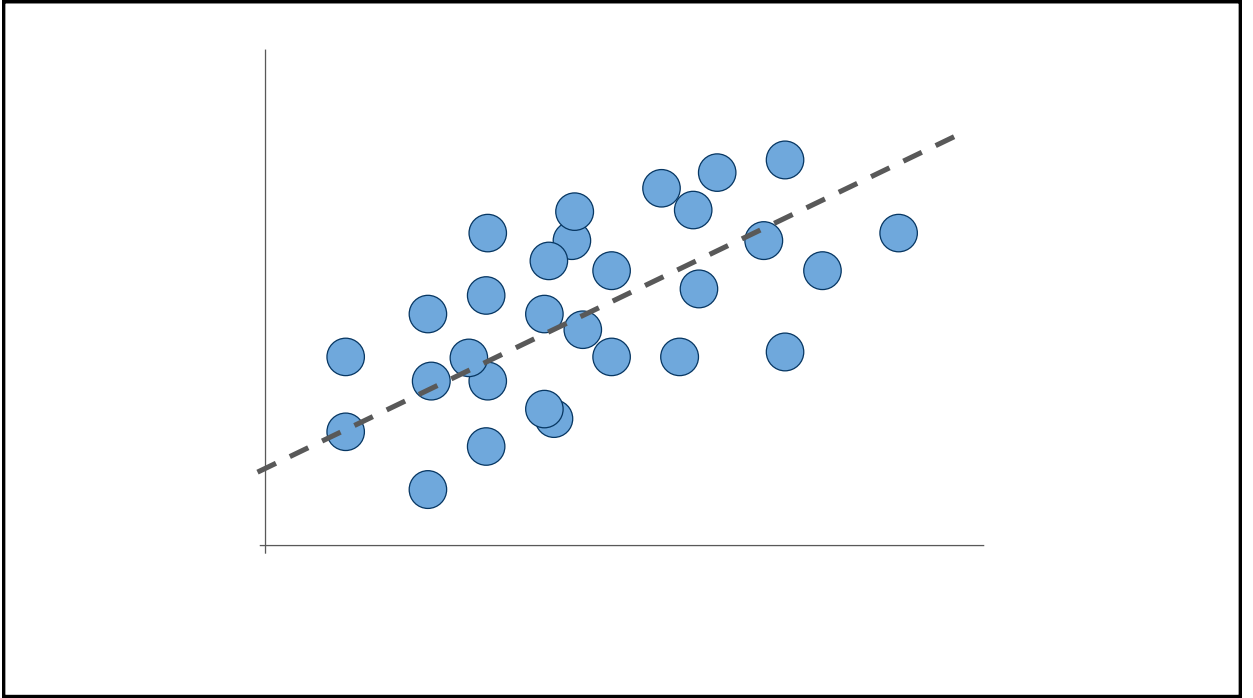
-Susan Etlinger, 2014 TED Talk

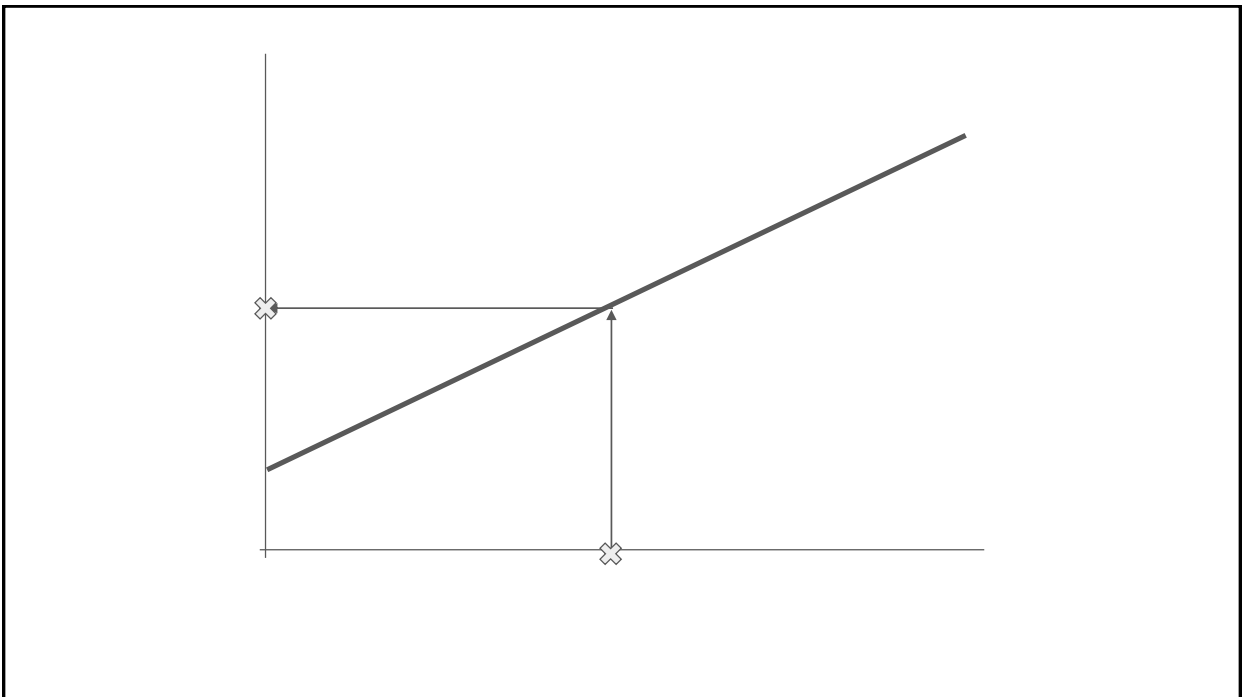
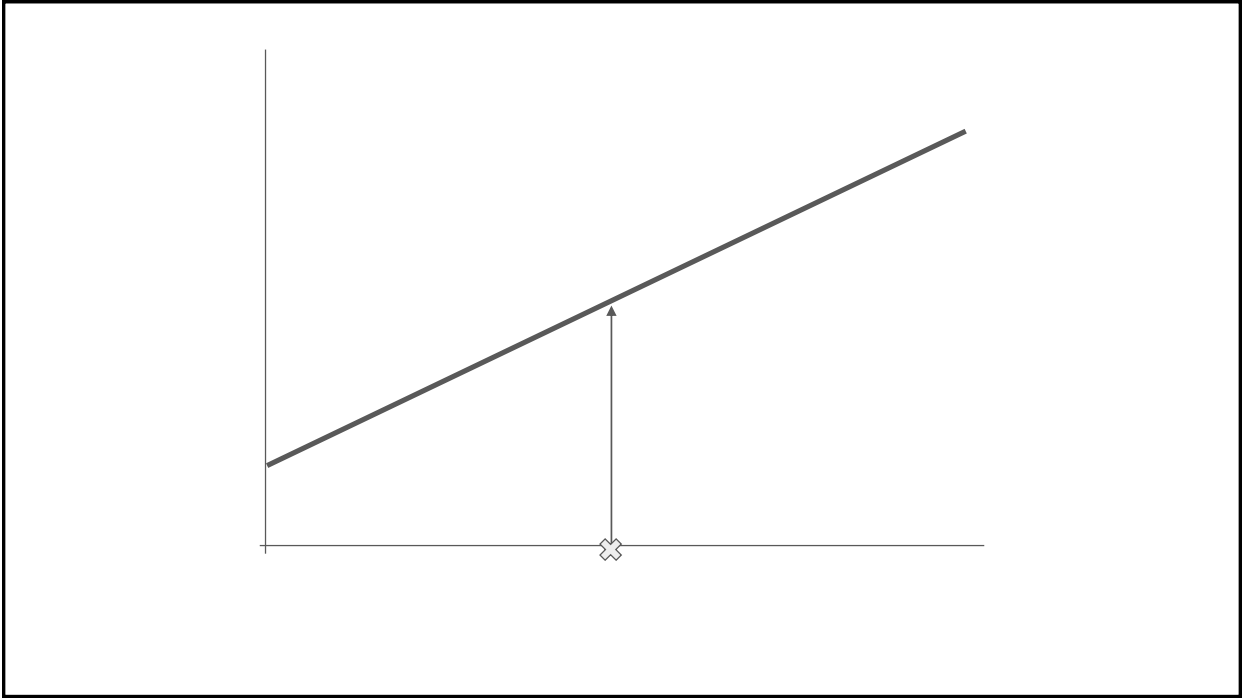
How can human biases get into a machine learning model?

Let's explore how machine learning systems are designed and developed

Some Types of  
Machine Learning Models



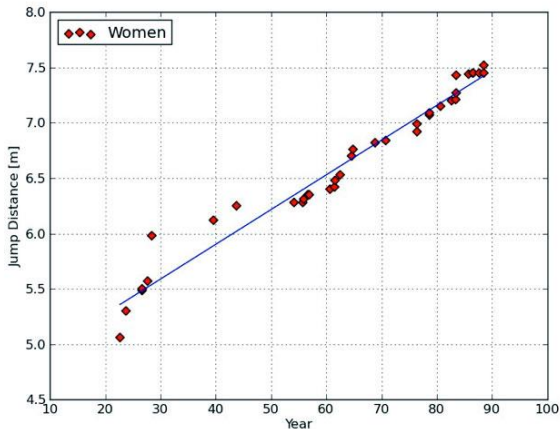




## Forecasting Record-Breaking Long Jump Distance by Year

“Olympics Physics: The Long Jump and Linear Regression”

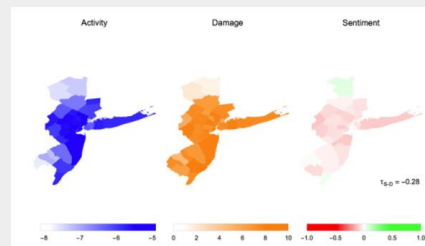
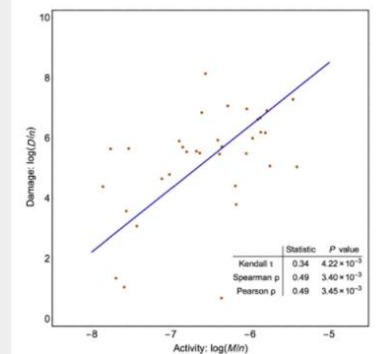
<https://www.wired.com/2012/08/physics-long-jump-linear-regression/>



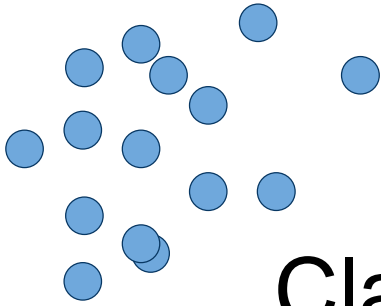
## Predicting Natural Disaster Damage by Counting Relevant Social Media Posts

“Rapid assessment of disaster damage using social media activity”

<http://advances.sciencemag.org/content/2/3/e1500779/tab-figures-data>

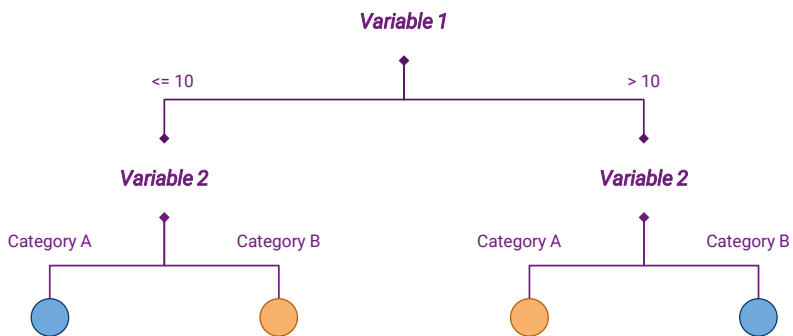
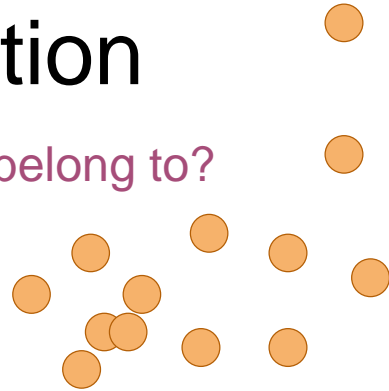


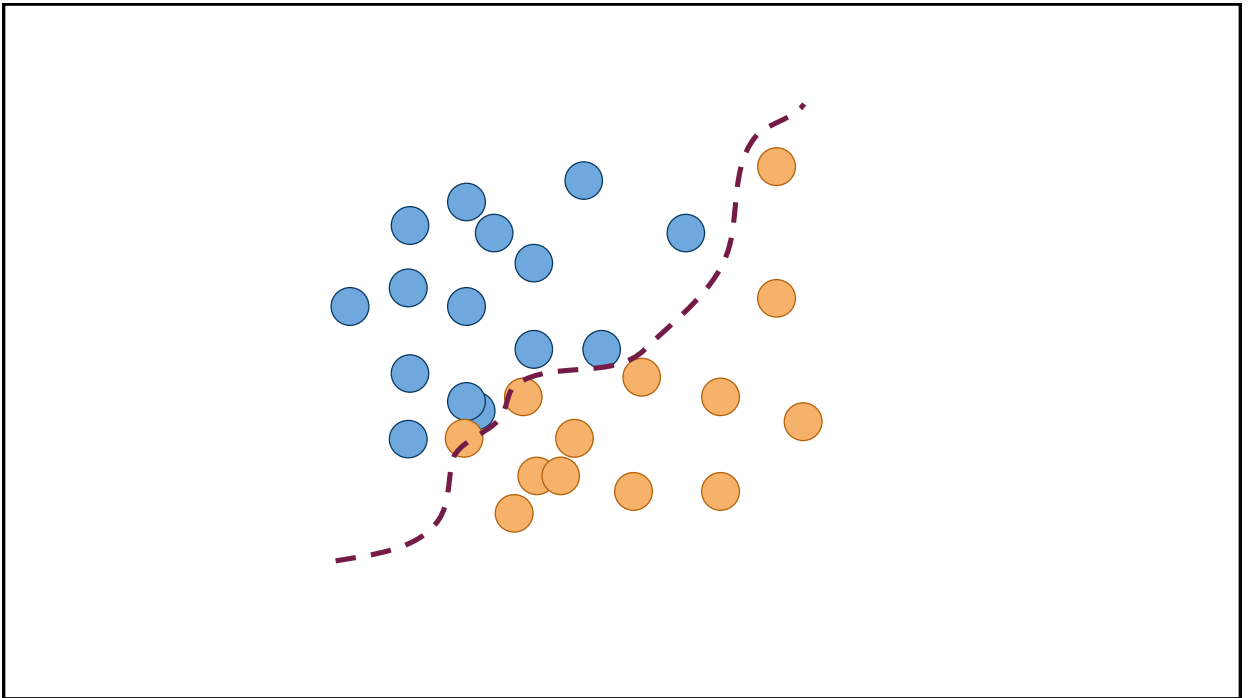
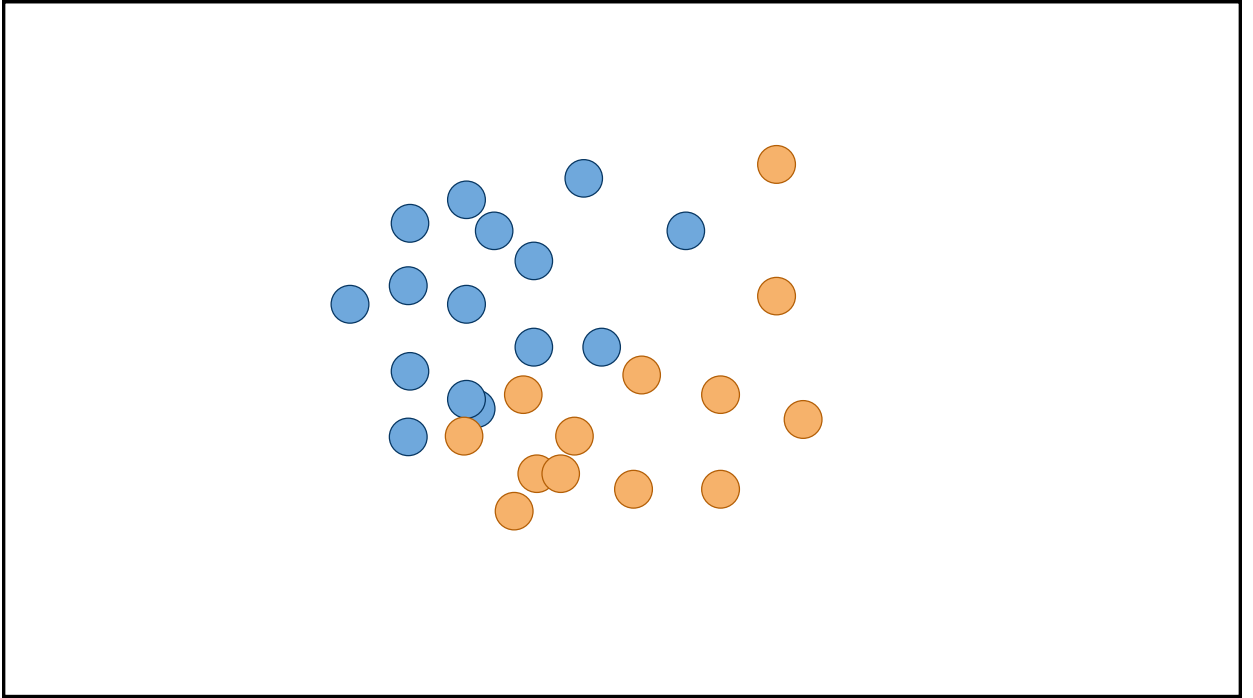


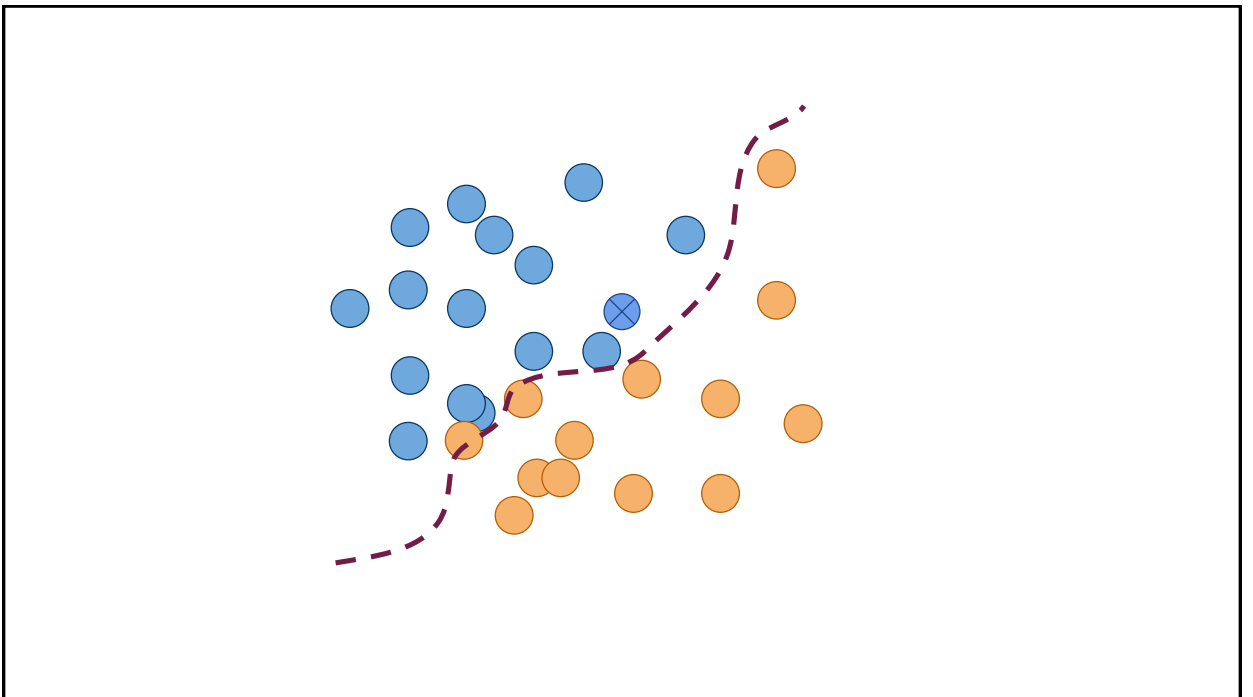
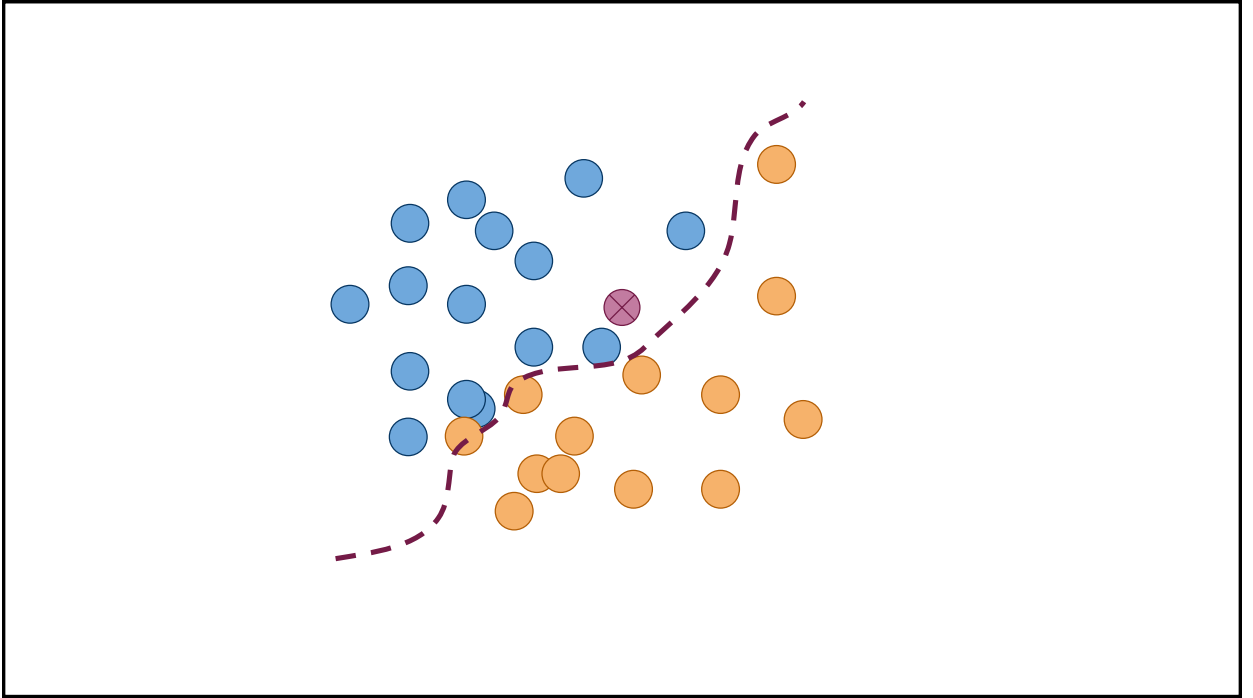


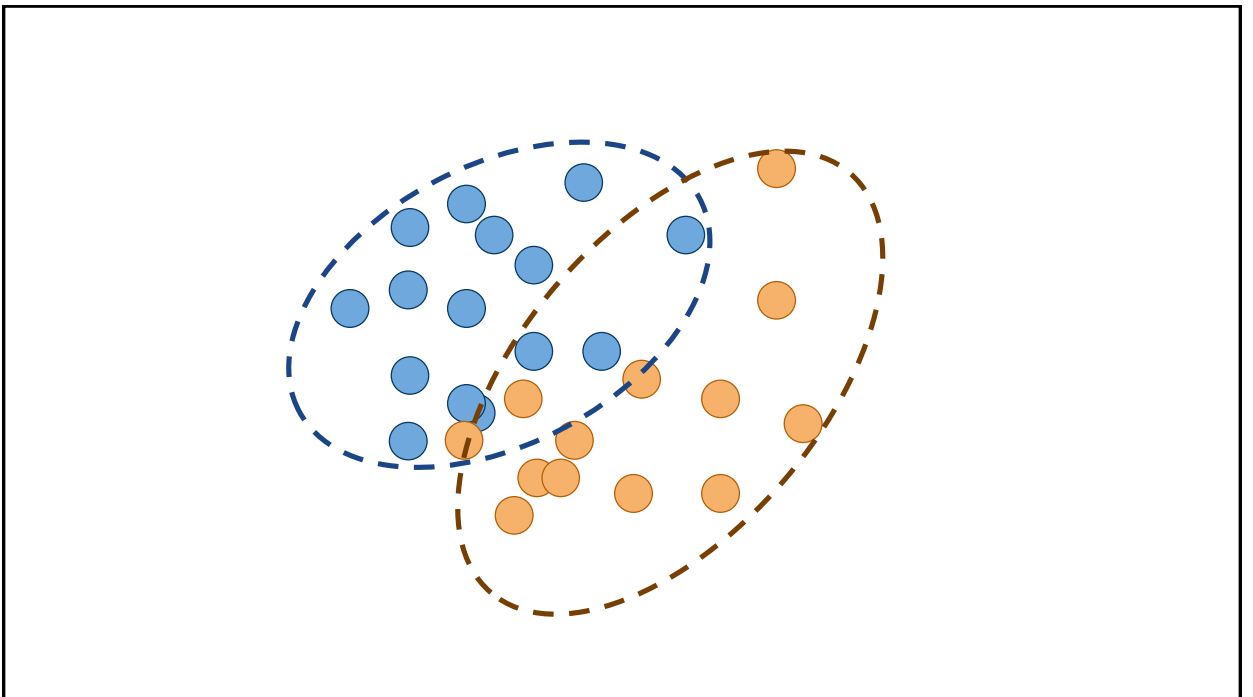
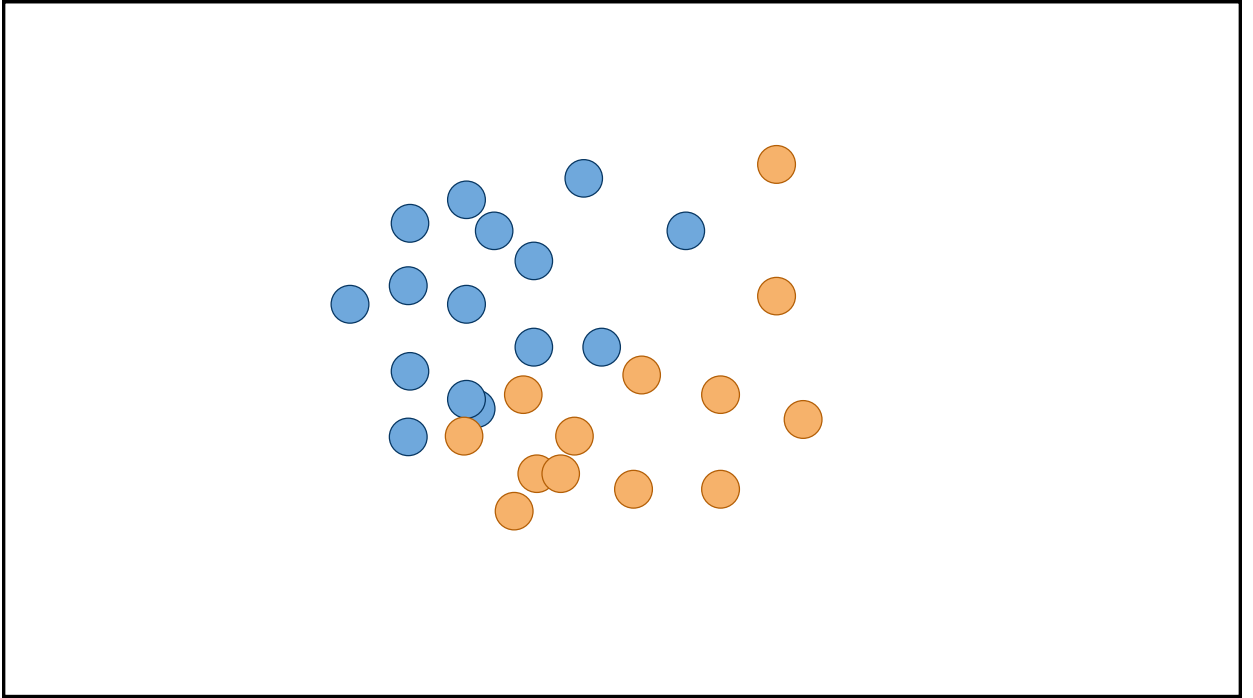
# Classification

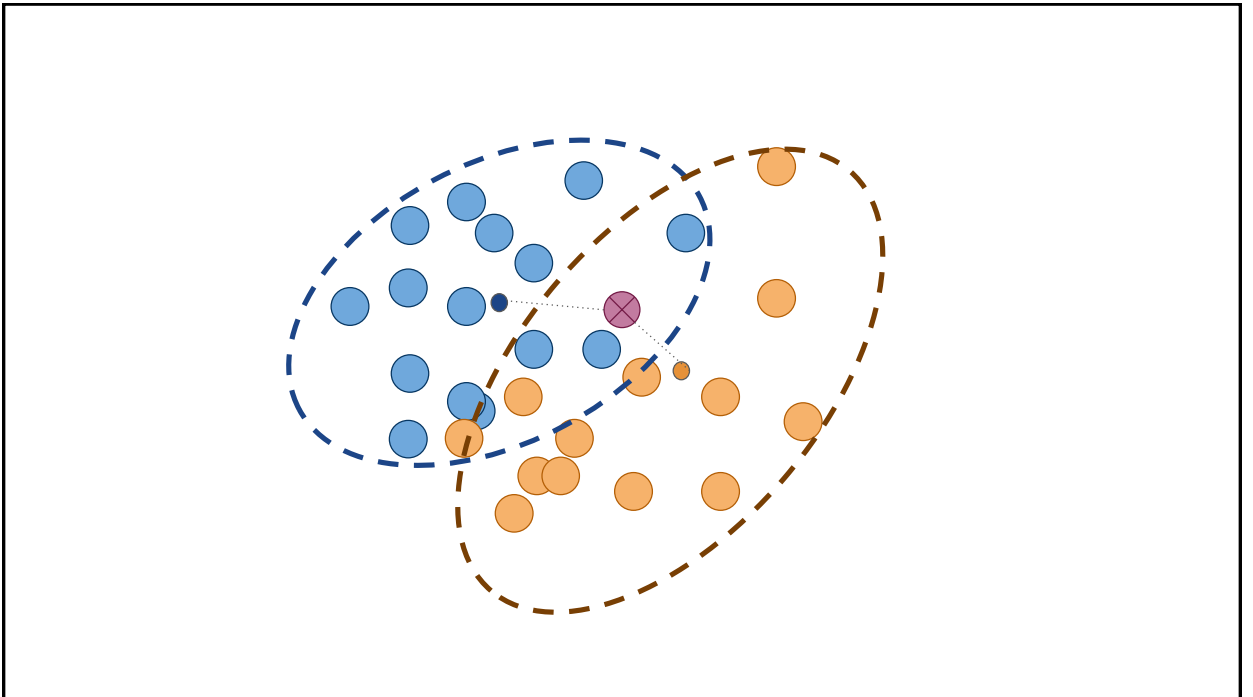
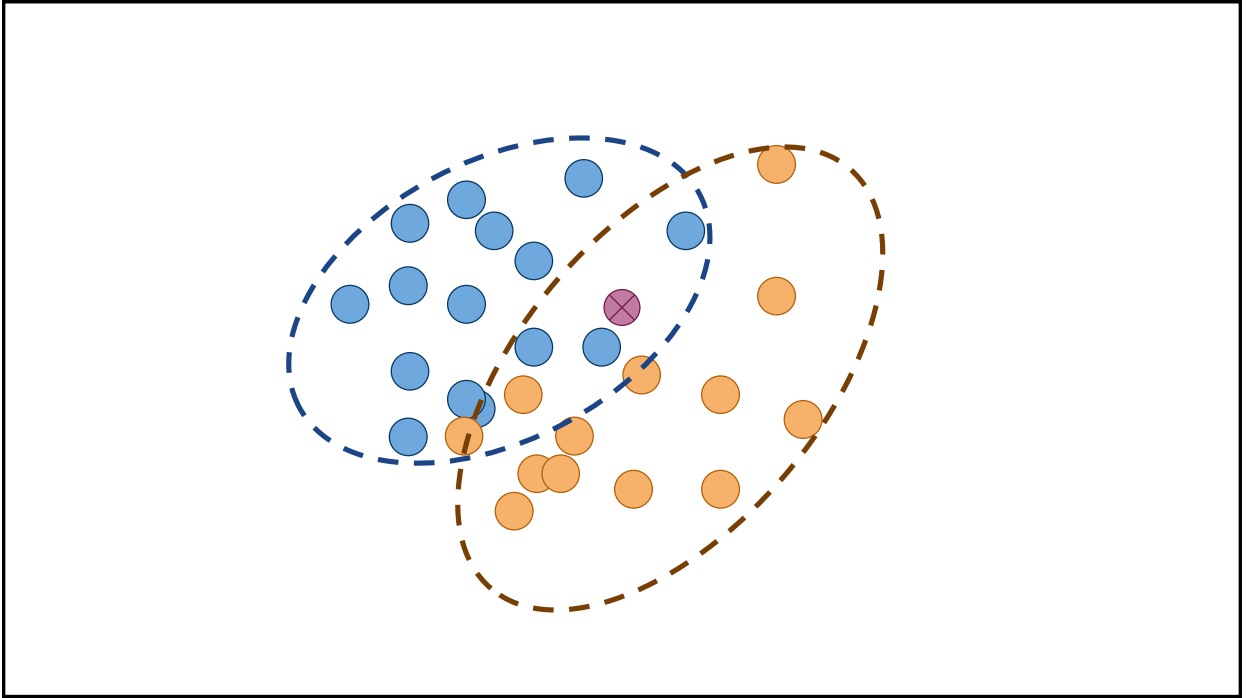
Which group does X belong to?

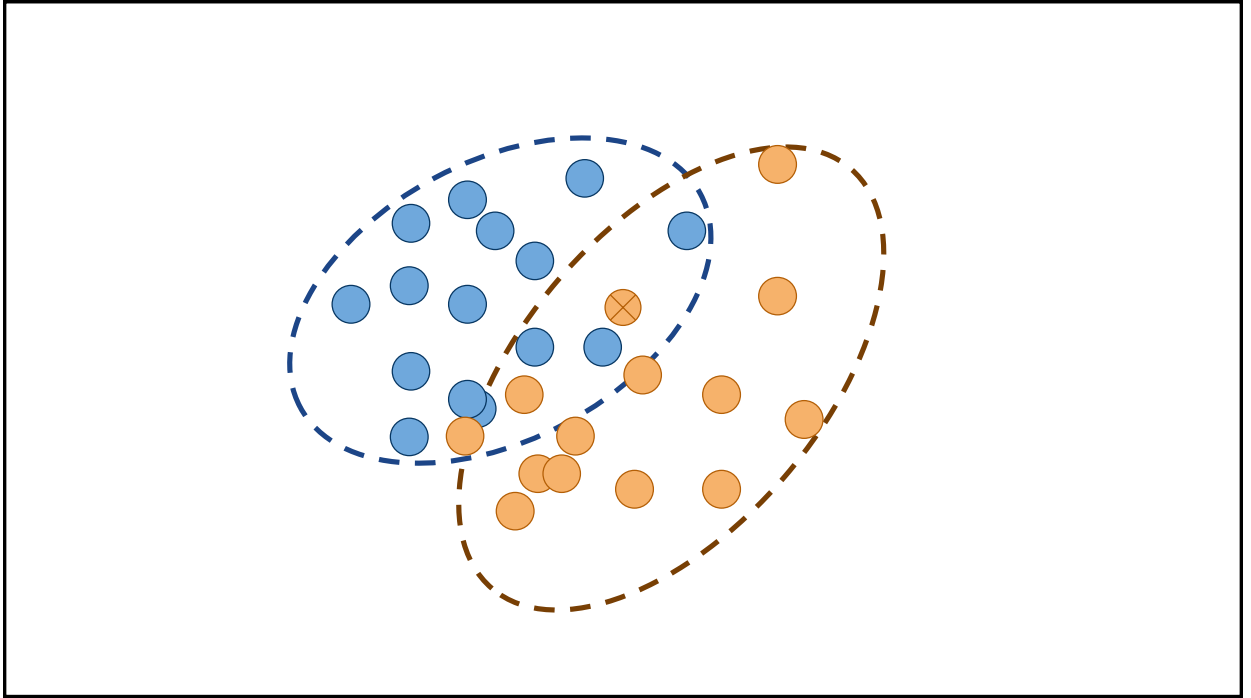












Is there an animal  
in this camera trap  
image?

<https://creativecommons.org/licenses/by-nc-sa/3.0/>



“Deep learning tells giraffes  
from gazelles in the  
Serengeti”

<https://www.newscientist.com/article/2127541-deep-learning-tells-giraffes-from-gazelles-in-the-serengeti/>

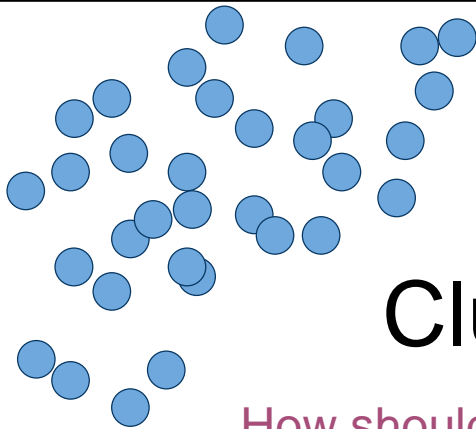
## Is a crime scene gang-related?

“Artificial intelligence could identify gang crimes—and ignite an ethical firestorm”

<http://www.sciencemag.org/news/2018/02/artificial-intelligence-could-identify-gang-crimes-and-ignite-ethical-firestorm>

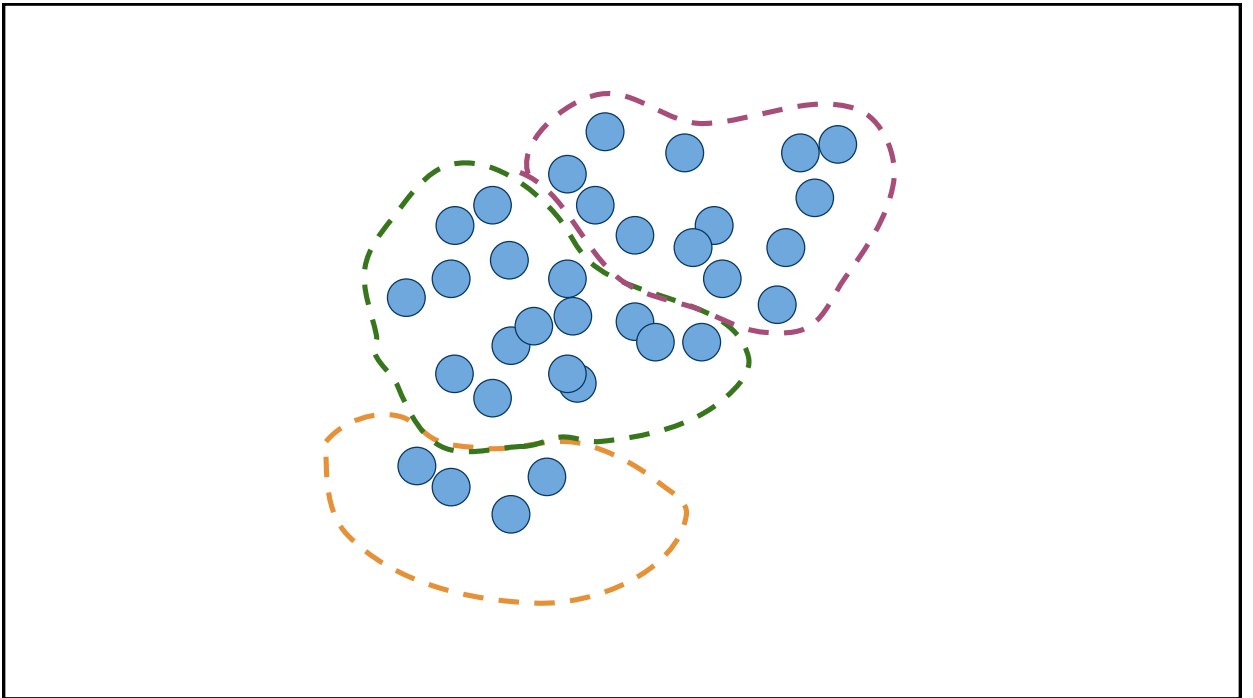
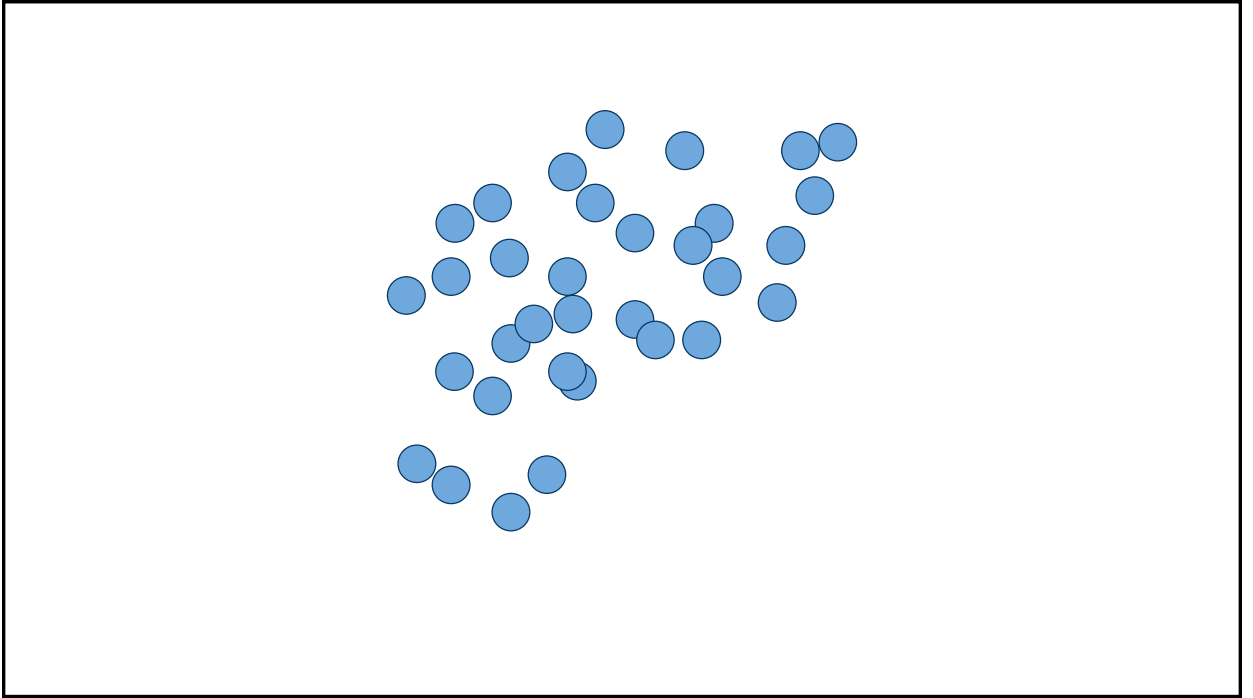


ISTOCK.COM/DENISTANGNEYJR

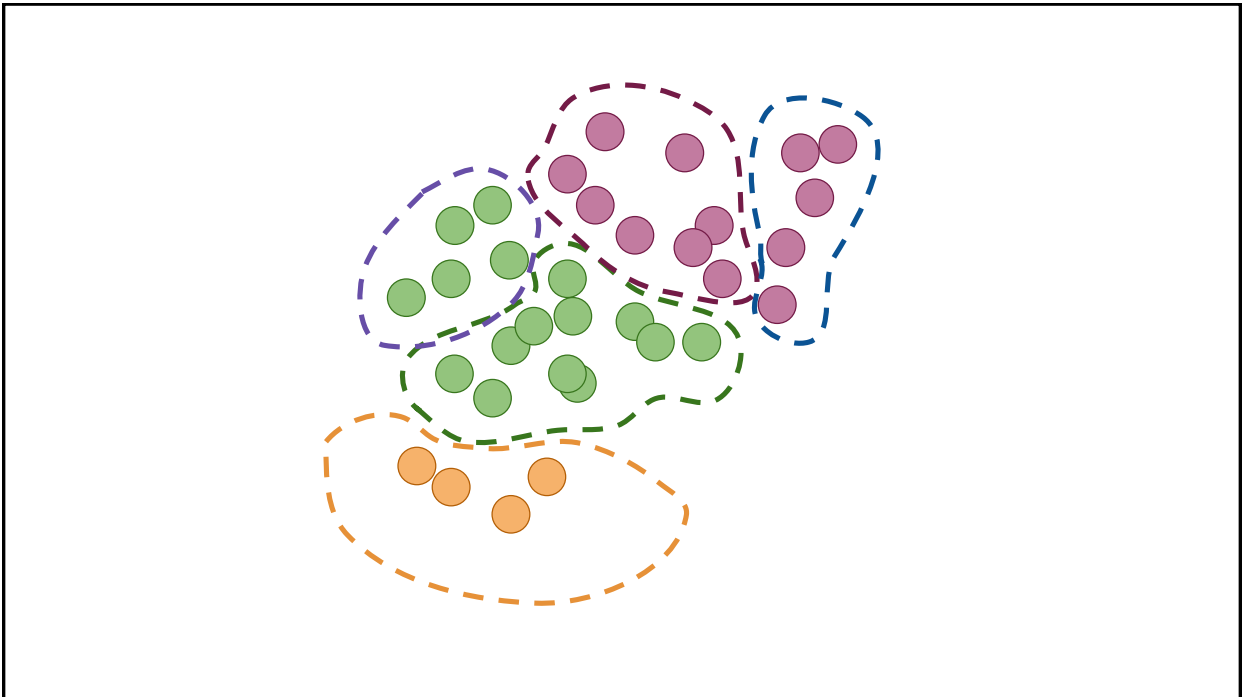
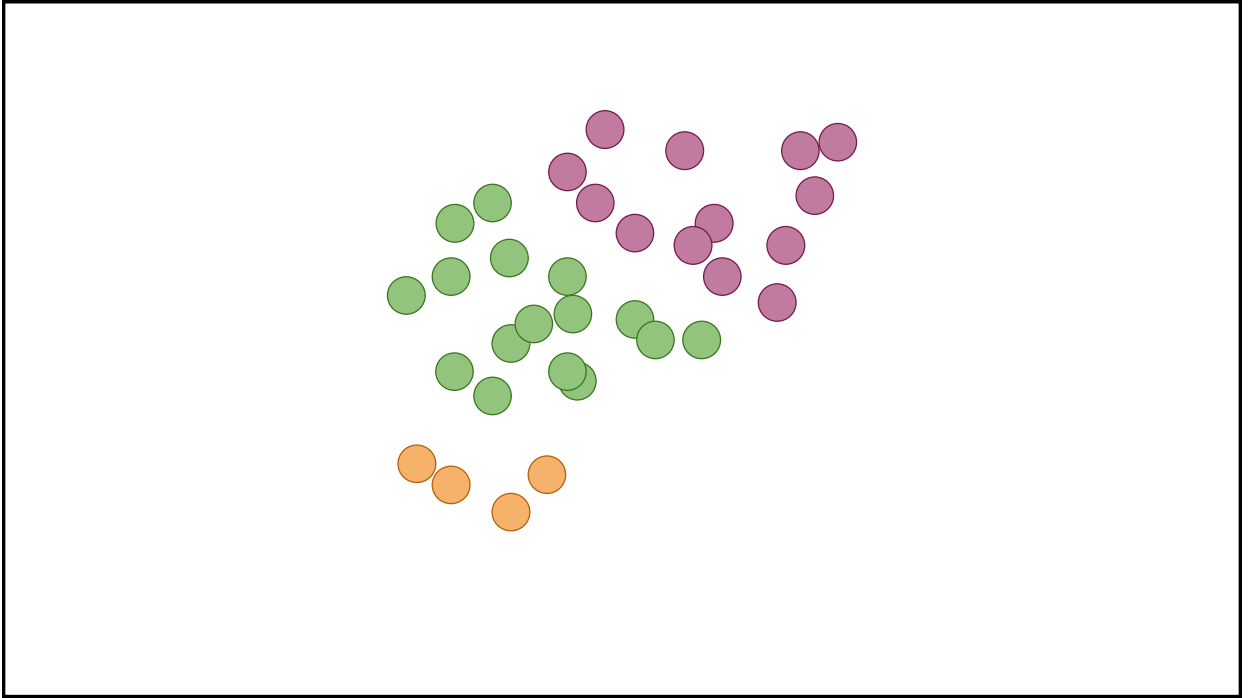


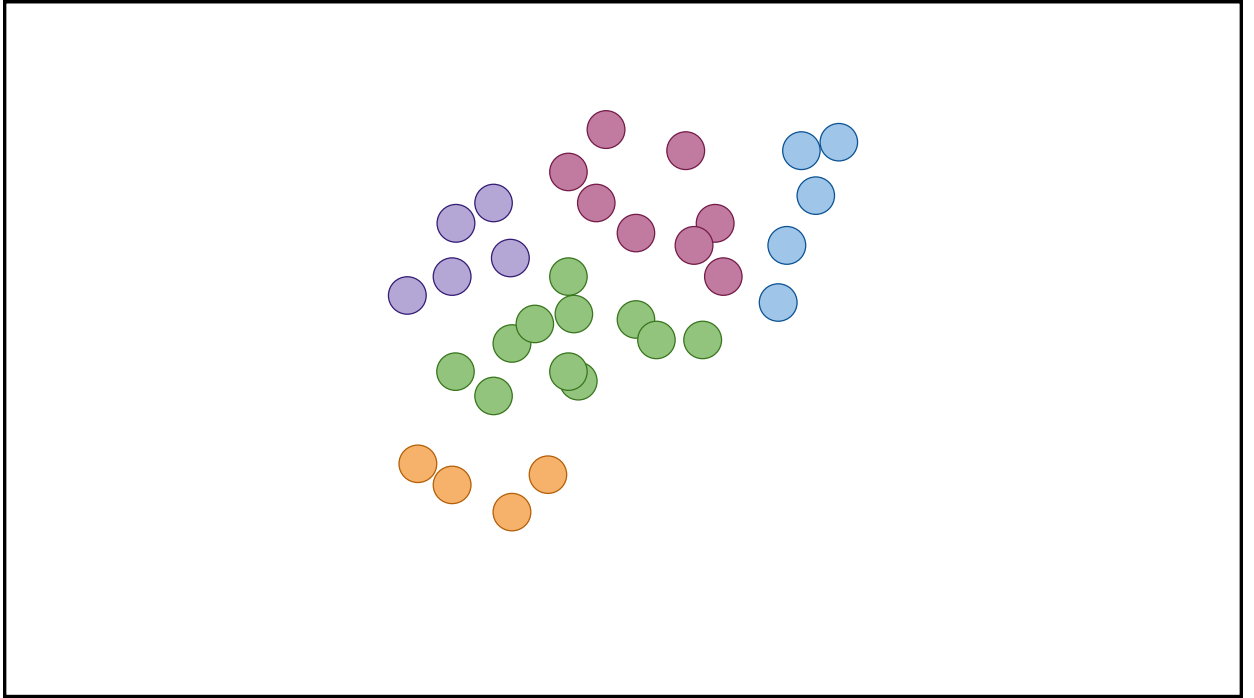
## Clustering

How should we group this data?









## How might we segment our customers?

“Understanding, Analyzing, and Retrieving Knowledge from Social Media”

<http://cucis.ece.northwestern.edu/projects/Social/>

Which neighborhoods  
are most likely to  
experience home  
burglaries this month?



(2) August 2009

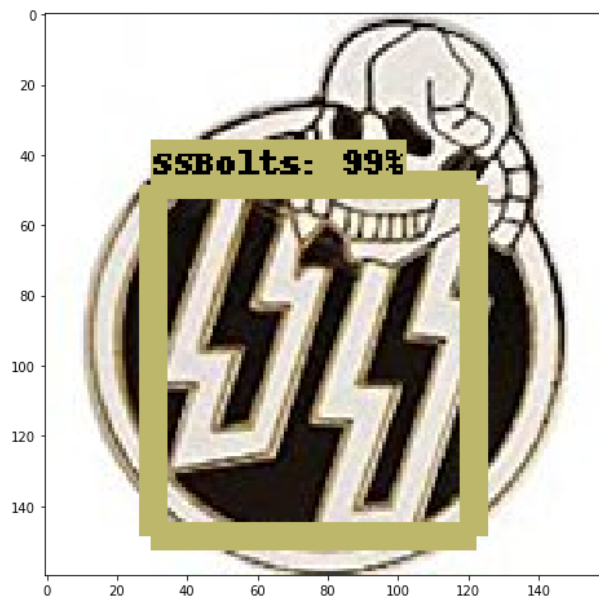
“Crime Forecasting Using  
Spatio-Temporal  
Pattern with Ensemble  
Learning”

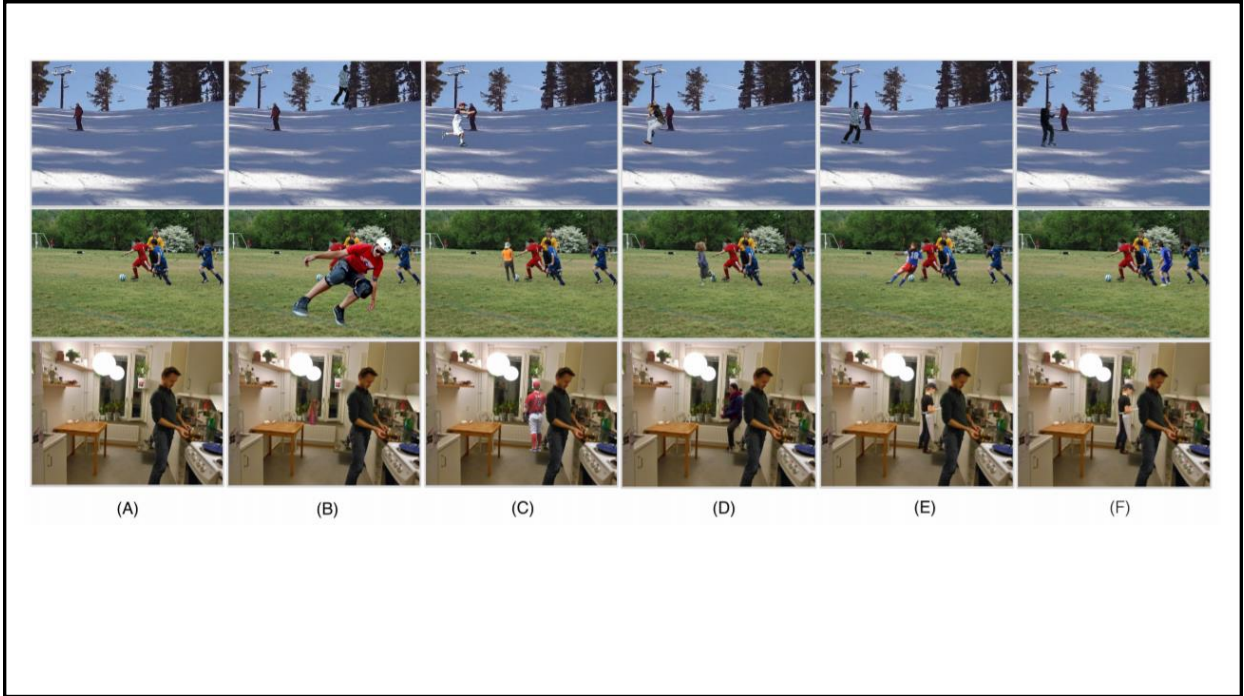
[https://www.cs.umb.edu/~csyu/YU\\_resume%202016\\_01\\_08\\_files/yuPAKDD2014.pdf](https://www.cs.umb.edu/~csyu/YU_resume%202016_01_08_files/yuPAKDD2014.pdf)

Artificial Neural Networks  
Reinforcement Learning  
Collaborative Filtering  
etc.

The purpose of most of these algorithms is to find patterns, trends, group things that are similar...

In other words, we're basically asking the computer to use lots of information to make generalizations, or stereotype.

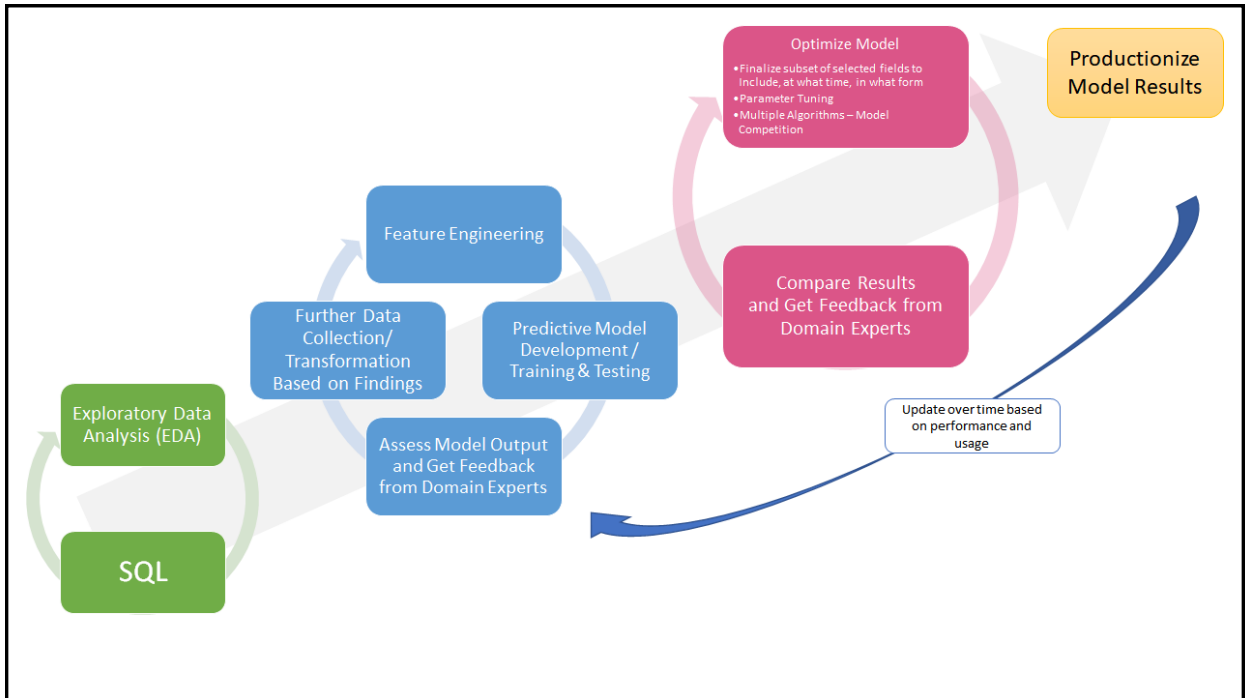




# Predictive Model Development Process

## Predictive Model Development

- Deciding what you're predicting / optimizing for
- Data collection and storage
- Data cleansing/preparation
- Feature selection & engineering
- Importing data into different algorithmic models
- Training & Testing
- Model evaluation & competition; Deciding what qualifies as a “good model”
  - Parameter Tuning, Cost function, Selecting cutoff values or stopping conditions, etc
- “Productionizing” - Applying to live data, building interactive reports for end-users, explaining what the results mean and how to use them to make decisions
- Monitoring, Improving, and Re-training over time



Where within this process can social biases be introduced?

## Data Collection: Incorrectly Recorded

### Hearing focuses on Texas troopers wrongly recording drivers' race

By: Claire Ricke 

Updated: Sep 20, 2016 01:57 AM CDT

## Data Collection: Manipulated

### Prison time for some Atlanta school educators in cheating scandal



By **Ashley Fantz, CNN**

🕒 Updated 7:03 AM ET, Wed April 15, 2015

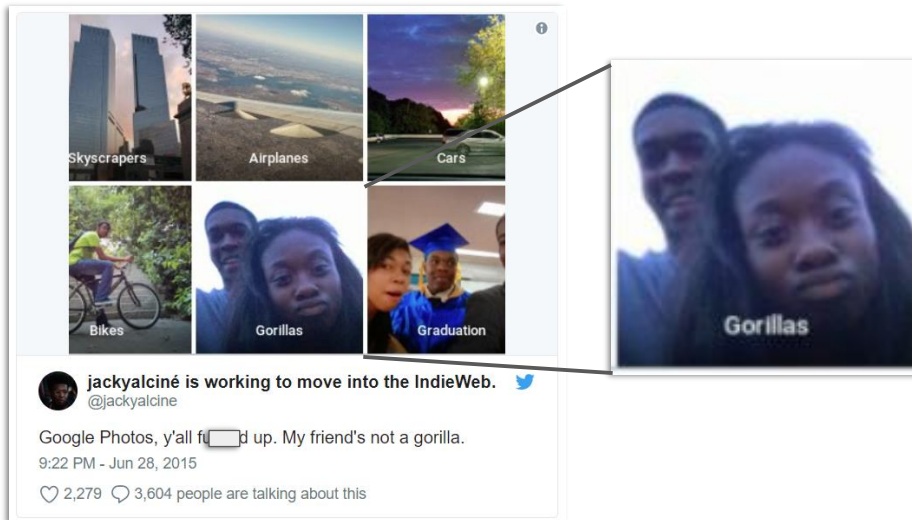
The cheating is believed to date back to 2001, when scores on statewide aptitude tests improved greatly, according to a 2013 indictment. The indictment also states that for at least four years, between 2005 and 2009, test answers were altered, fabricated or falsely certified.

A review that former Gov. Sonny Perdue ordered, determined that some cheating had occurred in more than half the district's elementary and middle schools.

Michael Bowers, a former Georgia attorney general who investigated the cheating scandal, said in 2013 that there were "cheating parties," erasures in and out of classrooms, and teachers were told to make changes to student answers on tests.



## Data Collection: Not Representative



## Update...

### Google 'fixed' its racist algorithm by removing gorillas from its image-labeling tech

*Nearly three years after the company was called out, it hasn't gone beyond a quick workaround*

By James Vincent | @jjvincent | Jan 12, 2018, 10:35am EST

## Data Collection: Contains Historic Biases

### Is Your Computer Sexist?

It may say "boss" is a man's job, BU and Microsoft researchers discover

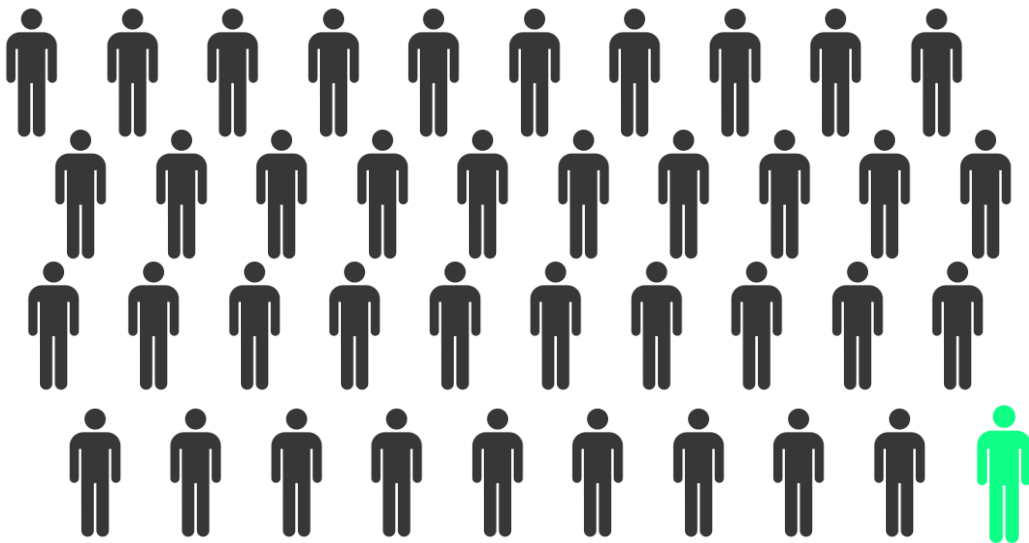
12.06.2016

By [Rich Barlow](#)

But word embeddings can recognize word relationships only by studying batches of writing. The researchers particularly focused on word2vec, a publicly accessible embedding nourished on texts from [Google News](#), an aggregator of journalism articles. Turns out that those articles contain gender stereotypes, as the researchers found when they asked the embedding to find analogies similar to "he/she."

The embedding spit back worrisome analogies involving jobs. For "he" occupations, it came up with words like "architect," "financier," and "boss," while "she" jobs included "homemaker," "nurse," and "receptionist."

## Data Availability: Imbalanced Dataset



## Model Evaluation: Confusion Matrix & Cost

### Accuracy

If 99 people out of 100 don't have cancer, and there is a test that just always comes back negative (predicts that no one has cancer), that test is still 99% accurate

	Has Cancer	Does Not Have Cancer
Tests Positive for Cancer	TRUE POSITIVE	FALSE POSITIVE
Tests Negative for Cancer	FALSE NEGATIVE	TRUE NEGATIVE

## Model Evaluation: Confusion Matrix & Cost

	Has Cancer	Does Not Have Cancer
Model Predicts Patient Has Cancer	0	0
Model Predicts Patient Does Not Have Cancer	10	990

*This model doesn't have any False Positives, and is "99% Accurate", but also has no Positive Predictive Value.*

*How is it Penalized for that?*

## Model Evaluation: Confusion Matrix & Cost

	Has Cancer	Does Not Have Cancer
Model Predicts Patient Has Cancer	5	90
Model Predicts Patient Does Not Have Cancer	5	900

*What is the “cost” of each type of error?*

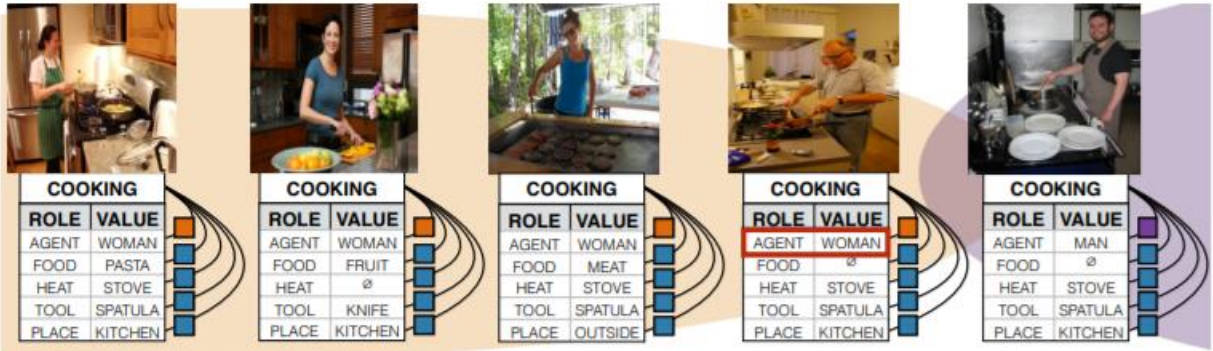
## Data Pre-Processing: Dropping Data

### Remove Missing Data

Now that you know how to mark missing values in your data, you need to learn how to handle them.

A simple way to handle missing data is to remove those instances that have one or more missing values.

## Model Training: Bias Amplification



<http://vicenteordonez.com/files/bias.pdf>

**Feature & Algorithm Selection** - Different algorithms handle different types of data in different ways

**Target Selection/Optimization Goal** - What are you optimizing for? (Technical & Business Decision)

**Model Evaluation** - How good does your model have to be to decide to stop improving it? And how do you define “good”?

Consider “cost” of each type of error.

**Example:** Optimizing for maximum video viewing time may incentivize the display of alarming/intriguing information, whether or not it is true (propaganda)

## YouTube, the Great Radicalizer



By Zeynep Tufekci

March 10, 2018

The New York Times

It seems as if you are never “hard core” enough for YouTube’s recommendation algorithm. It promotes, recommends and disseminates videos in a manner that appears to constantly up the stakes. Given its billion or so users, YouTube may be one of the most powerful radicalizing instruments of the 21st century.

**Implementing the Trained Model** - In what scenarios can your model be applied? How generalizable is it?

**Interpretation** - How do you interpret the results? How do you document and explain to others how to interpret the results?

**Maintenance** - For how long can the current model be applied? When does the model need to be retrained? Does the “ground truth” change?

## Can your model be gamed?

### How to persuade a robot that you should get the job

Do mere human beings stand a chance against software that claims to reveal what a real-life face-to-face chat can't?



**Stephen Buranyi**

Sat 3 Mar 2018 19.05 EST

A fightback against automation has emerged, as applicants search for ways to game the system. On web forums, students trade answers to employers' tests and create fake applications to gauge their processes. One HR employee for a major technology company recommends slipping the words "Oxford" or "Cambridge" into a CV in invisible white text, to pass the automated screening.

Could your model cause harm?

Or perpetuate existing social hierarchies,  
preventing a fair playing field?

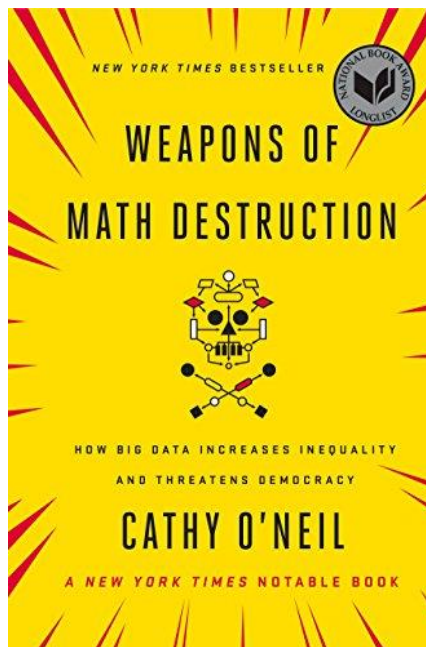
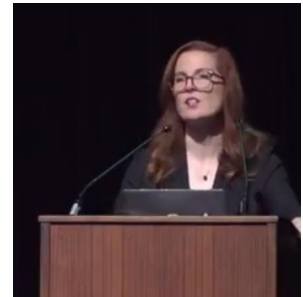
## Some Types of Harm a Model Can Perpetuate

**Allocative harms** - resources are allocated unfairly or withheld (transactional, quantifiable)

**Representational harms** - systems reinforce subordination/perceived inferiority of some groups (cultural, diffuse, can lead to other types of harm)

- stereotyping
- underrepresentation
- recognition
- denigration
- Ex-nomination

*(from Kate Crawford's talk at 2017 NIPS Conference, The Trouble With Bias)*





What makes a model a “Weapon of Math Destruction”?

- **Opacity** - inscrutable “black boxes” (often by design)
- **Scale** - capable of exponentially increasing the number of people impacted

*“The privileged...are processed more by people,  
the masses by machines.”*

- **Damage** - can ruin people’s lives and livelihoods

So, Can a Machine Be  
Racist or Sexist?

# YES

“...there is nothing “artificial” about [Artificial Intelligence] — it is made by humans, intended to behave like humans and affects humans. So if we want it to play a positive role in tomorrow’s world, it must be guided by human concerns.”

“...there are no “machine” values at all, in fact; machine values *are* human values.”

-Fei-Fei Li

How to Make A.I. That’s Good for People,  
New York Times, 3/7/2018



# What can we do about it?

- **Be aware** of the potential for bias and disparate impact machine learning models can perpetuate (and help educate others)
- **Test for bias** in your models & evaluate model performance on minority classes in your dataset
- **Evaluate possible uses** (or misuses) of your model, and “perverse incentives” it may create in the system into which it’s being deployed
- **Improve transparency & explainability** of how your model categorizes and predicts things (LIME)

- **Document Data Source & Transformation Pipeline** to help manage data governance & provenance and allow reproducible research
- **Communicate context**, and explain model generalization limits to end-users
- **Involve Domain Experts** who know the history of the data collection, actual field definitions, data quality issues, system changes over time, how the model will likely be applied, ethical issues and laws in the field of application, etc.
- **Gather Representative Training Data**

- Include **fairness as an optimization objective**, add social bias penalties (research emerging, more needed)
- **Research and Develop** new tools and techniques for detecting bias and reducing disparate impact caused by machine learning models
- **Build adversarial tools** to stress-test your own models, and thwart others that are causing harm
- **Hire Diverse Teams**
- **Demand Accountability and Regulation** in the industry, and in your own organizations

## Organizations working on this

**AI Now Institute** at NYU <https://ainowinstitute.org/>

**Algorithmic Justice League** <https://www.ajlunited.org/>

(Watch Joy Buolamwini's TED Talk!)

Data for Democracy, Bloomberg, BrightHive

**Data Science Code of Ethics**

<http://datafordemocracy.org/projects/ethics.html>



**Fairness, Accountability, and Transparency in Machine Learning**

<https://www.fatml.org/>

## Example research

**Men Also Like Shopping:**

**Reducing Gender Bias Amplification using Corpus-level Constraints**

Jieyu Zhao<sup>§</sup> Tianlu Wang<sup>§</sup> Mark Yatskar<sup>†</sup>  
Vicente Ordonez<sup>§</sup> Kai-Wei Chang<sup>§</sup>

**Certifying and removing disparate impact\***

Michael Feldman Sorelle A. Friedler John Moeller  
Haverford College Haverford College University of Utah  
Carlos Scheidegger Suresh Venkatasubramanian<sup>†</sup>  
University of Arizona University of Utah

**Fairness-Aware Classifier  
with Prejudice Remover Regularizer**

Toshihiro Kamishima<sup>1</sup>, Shotaro Akaho<sup>1</sup>, Hideki Asoh<sup>1</sup>, and Jun Sakuma<sup>2</sup>

**A STUDY OF PRIVACY AND FAIRNESS  
IN SENSITIVE DATA ANALYSIS**

**Learning Fair Representations**

MORITZ A.W. HARDT

Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, Cynthia Dwork ; *Proceedings of the 30th International Conference on Machine Learning, PMLR 28(3):325-333, 2013.*

## More resources

Flipboard Magazine where I collect articles on this topic:

<https://flipboard.com/@becomingdatasci/bias-in-machine-learning-rv7p7r9ry>

<https://www.becomingdatascientist.com/2015/11/22/a-challenge-to-data-scientists/>

<https://developers.google.com/machine-learning/fairness-overview/>

<https://sloanreview.mit.edu/article/the-risk-of-machine-learning-bias-and-how-to-prevent-it/>

<https://www.fatml.org/resources/relevant-scholarship>

[https://twitter.com/random\\_walker/status/961332883343446017](https://twitter.com/random_walker/status/961332883343446017)

“How do we govern ourselves? How do we instill that trust in others that we as stewards of that data understand the power that we have, and want to make sure that we're doing right by the people who are trusting us with that data?”

-Natalie Evans Harris

Data for Good Exchange 2017, Inside Bloomberg





Renée Teate

BecomingADataScientist.com

@becomingdatasci

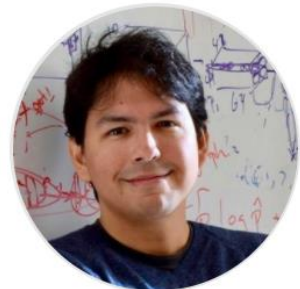
## Social Bias in Machine Learning Panel



**Emily Crose**  
Undisclosed  
Threat Hunter



**Ines Montani**  
Explosion AI  
Founder



**Vicente Ordonez**  
UVA  
Assistant Professor