

Big Data

Michael Masterson

4/8/19

What is Machine Learning?

Supervised Learning

In supervised learning the goal is to predict an outcome based on data. If the outcome is a category, then you are doing classification. If the outcome is a number, then you are doing regression. The components of the data used to produce predictions are called features.

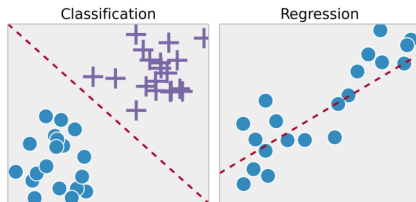


Figure 1:

Examples Of Machine Learning Applications:

- Predict whether a search result is appropriate based on whether people click on it
- Predict whether a tumor is cancerous based on size, shape, and color
- Classify whether a social media post is political
- Predict where political violence within a country is likely to reoccur based on past violence and information about the districts in that country

Why use Machine Learning?

- Statistical models can often outperform experts
- Models can do tasks no human can do

It's the Future

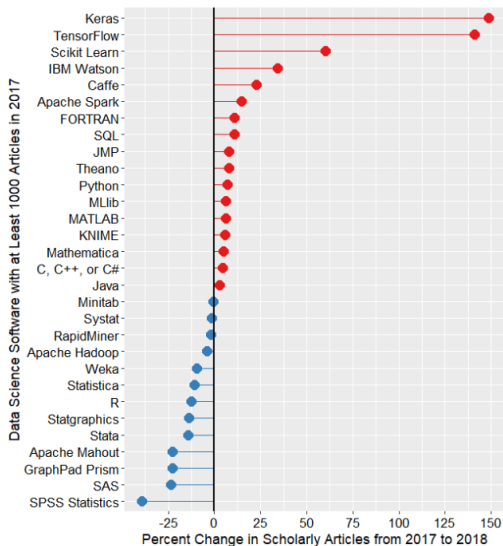


Figure 2:

But We Will Never Use This Right?

You already are!

OLS regression analysis is a type of machine learning!

The main difference is that we are focused on causal inference rather than prediction.

How Does Supervised Learning Work?

- Get set of labeled data
- Divide into a training set and a test set
- Provide training data with labels to the model
- Train the model (more on this later)
- Predict on the test set to evaluate the model

How Does Training Work?

The model minimizes a loss function

For example, we have learned that our linear models minimize squared error. That is a loss function called mean squared error!

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

There are other loss functions for categorical classification.

How Does Testing Work

Usually the model will be tested with some form of *cross-validation*.

Figure 3: 5-fold Cross-validation Example



What is Big Data

- Lots of observations/rows
- Usually must be stored in a database

Why Does Big Data Matter?

Improves prediction/classification

- Represent more variation in train and test set
- Help combat overfitting
- Allows more complicated models
- Allows inclusion of more features

Statistical Issues with Machine Learning

Watch Out for Overfitting

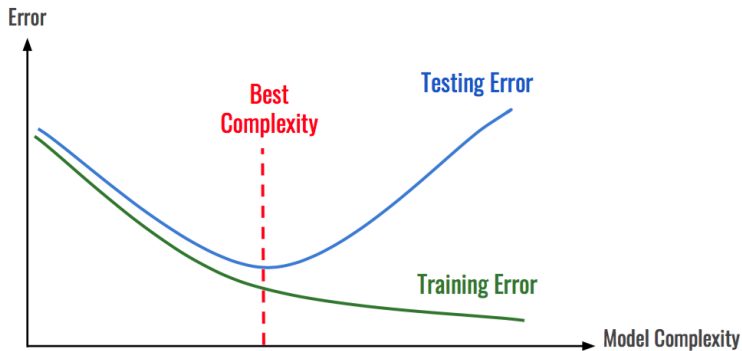


Figure 4: Overfitting

Limits of Big Data

- Correlation not causation
- Big data does not solve causal inference problems
- Typically not a random sample
- Human response
- Comparability

Interpretability

- Some models more interpretable than others.
- Parameters not always directly interpretable by humans.
- Can have accuracy and interpretability tradeoff

Brief Example

Weibo posts

5-fold cross-validation of ML models on a dataset of about 11,000 Chinese social media posts from Weibo (similar to twitter).

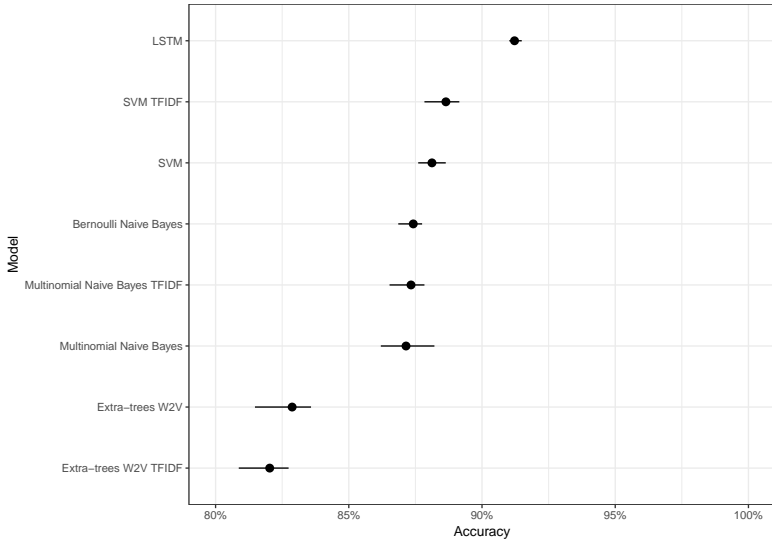
The models must predict whether the posts contain political content or not. About half of the posts do.

Here are Two (Translated) Posts

- *Political*: 'APEC' has gone, and our old friends [air pollution] have returned to us. [Our old friends] never abandon us or give us up. Such friends are absolutely loyal.
- *Not Political*: A while ago, I suddenly wanted to go to the Internet cafe to go online. I went to the school wall and climbed onto the wall. Surprisingly, just when I climbed up, a school security guard came and called to me in the distance. I said, "Is it possible to come in to the school and find someone?" Security: "No! Get out right away!" Then I got out. . .

How do you think the models did?

Result



Learning to Walk

Video

Ethical Issues

Omission

Women and minority groups are often excluded from studies. This can have major consequences:

- Women 47% more likely to have serious injuries in car crash
- Some blood pressure drugs backfire for women
- Common hyper-tension drug less effective for black patients
- Hispanic and Black patients make up $< 1.3\%$ of patients in cancer clinical treatment trials

Discriminatory Classification/Prediction

Models may use race, gender, or religion as features, leading to discriminatory results.

Examples:

- Who should be hired?
- How is likely to reoffend?
- Is this person a civilian or a terrorist?

If human decisions are used as training data, then the model will replicate any discriminatory practices in those decisions.

Crime Prediction

Pacific Standard

[HOME](#) > [SOCIAL JUSTICE](#)

CAN RACIAL BIAS EVER BE REMOVED FROM CRIMINAL JUSTICE ALGORITHMS?

A recent vote over a proposed tool to predict the risk that a person would pose a threat to public safety in Pennsylvania stirred a debate over its unintended consequences.

STEPHANIE WYKSTRA · JUL 12, 2018

Figure 5:

Microsoft Chat bot

We use cookies to offer you a better browsing experience, analyze site traffic, personalize content, and serve targeted advertisements. Read about how we use cookies by clicking "Cookie Information." If you continue to use this site, you consent to our use of cookies.

[Cookie Information](#)

intelligent machines

Microsoft's neo-Nazi sexbot was a great lesson for makers of AI assistants

Yandex's head of machine intelligence says Microsoft's Tay showed how important it is to fix AI problems fast.

by Rachel Metz March 27, 2018



JEREMY PORTJE

R

Remember Tay, the chatbot Microsoft unleashed on Twitter and other social platforms two years ago that quickly turned into a racist, sex-crazed neo-Nazi?

Figure 6:

Privacy

- Data often collected without people's knowledge/consent
- Predict pregnancy
- Predict sexual orientation
- Can identify “anonymous data”

Deep Fakes

Video

Autonomous Weapons

4/7/2019

Killer Robots | Human Rights Watch



Explore Arms

Killer Robots

Fully autonomous weapons, also known as “killer robots,” would be able to select and engage targets without human intervention. Precursors to these weapons, such as armed drones, are being developed and deployed by nations including China, Israel, South Korea, Russia, the United Kingdom and the United States. It is questionable that fully autonomous weapons would be capable of meeting international humanitarian law standards, including the rules of distinction, proportionality, and military necessity, while they would threaten the fundamental right to life and principle of human dignity. Human Rights Watch calls for a preemptive ban on the development, production, and use of fully autonomous weapons. Human Rights Watch is a founding member and serves as global coordinator of the **Campaign to Stop Killer Robots**.

EUROPE/CENTRAL ASIA



<https://www.hrw.org/topic/armay/killer-robots>

1/6

Figure 7:

Job Loss

4/7/2019



A.I. Expert Kai Fu Lee: 40% of Jobs Will Be Lost to AI, Robots | Fortune

FORTUNE

A.I. Expert Says Automation Could Replace 40% of Jobs in 15 Years

By DON REISINGER January 10, 2019

An artificial intelligence expert and venture capitalist predicts automation will cause major changes in the workforce.

Speaking to [CBS News](#)' Scott Pelley in an [interview for 60 Minutes](#) on Sunday, Kai Fu Lee said that he believes 40% of the world's jobs will be replaced by robots capable of automating tasks. He said that both blue collar and white collar professions will be affected, but he believes those who drive for a living could be most affected.

"Chauffeurs, truck drivers, anyone who does driving for a living—their jobs will be disrupted more in the 15-25 year time frame," he said in the interview. "Many jobs that seem a little bit complex, chef, waiter, a lot of things will become automated."

Lee's comments are not necessarily new. [Many who support artificial intelligence](#) and automation believe that they can fundamentally change the workforce. But many of those people also believe that while some jobs could be affected, humans will find [new opportunities](#) surrounding artificial intelligence and take on new professions.

A growing number of detractors—including Elon Musk, [who has warned](#) about the power of artificial intelligence—worry that automation could disrupt entire communities and [disproportionately affect low-income workers](#).

Still, many, including Lee, believe there's no slowing down artificial intelligence and its impact on society. And he compared artificial intelligence to major innovations in history, like the steam engine and electricity, saying

[fortune.com/2019/01/10/automation-replace-jobs/](#)

1/2

Figure 8:

Solutions

IRB

Institutional Review Board

- Legally required for University Research
- No such requirement for corporations

Principles

- Respect for persons
 - Informed consent
 - Limit deception
- Beneficence
 - Maximize benefits minimize harms
- Justice
 - Risks and benefits should be distributed fairly
 - Protect vulnerable groups

Avoiding Discrimination by Omission

No easy solutions here but some possibilities include:

- Ensure your research is not exclusive
- More female researchers
- More researchers from minority groups

Avoiding Discriminatory Prediction/Classification

Excluding features like race and gender do not solve this because other features could correlate

For applications where discrimination is an issue:

- Models should be publically available (not proprietary)
- Models should be interpretable
- Test trained models to ensure they do not have discriminatory impacts

[Learn More](#)

General Machine learning

- ML intro video playlist (just watch first 3)

Books

- *The Elements of Statistical Learning* by Hastie et al.
- *An Introduction to Statistical Learning* by James et al.

Languages

SQL

- [Video Playlist](#)

Python

- [Video Playlist](#)
- [Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow, 2nd Edition by Sebastian Raschka and Vahid Mirjalili](#)

Classes

Next Semester!

STAT 479 - Special Topics in Statistics

- Section 001 Intro Classification and Regression Trees
- Section 002 Machine Learning