

Introduction to Machine Learning and R

Jonathan Hersh, PhD (Chapman Argyros School of Business)

Belize Central Bank

10/24/2022

Outline

1. What is Machine Learning?

- Machine learning versus econometrics

2. Why Machine Learning for Public Policy

- Big data requires it
- Non-linear relationships
- Better forecasts/econometrics
- Anomaly detection

3. Some Basic Machine Learning Concepts

- Supervised vs Unsupervised learning
- Testing/Training Sets
- Bias-Variance Tradeoff

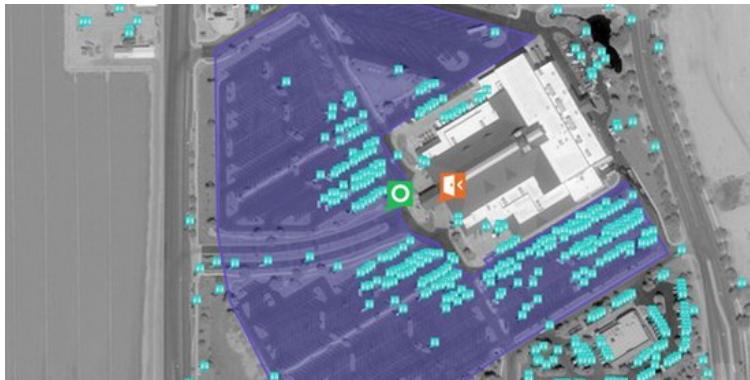
About Me

- Assistant Professor Economics and Management Science Chapman University
- PhD in economics, Boston University
- **Research Fields:**
 - Applications of artificial intelligence (computer vision)
 - Economics of information systems
 - Development economics
 - Digitization strategy
- **Teaching Fields:**
 - Machine learning
 - Applications of artificial intelligence



My Research

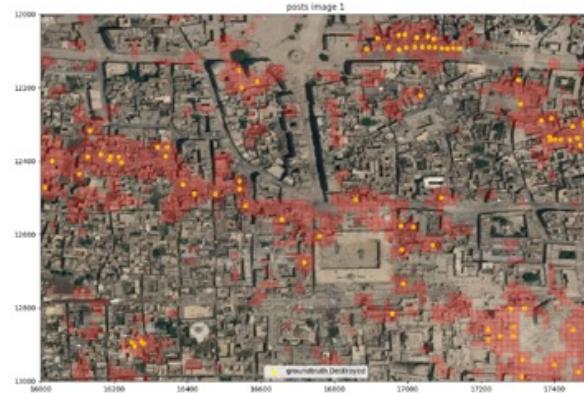
- Satellite Imagery + Computer Vision + Machine Learning



Count cars in parking lots!

Dense Prediction: Scanning Aleppo

Damaged buildings in Syria!



- Advised World Bank/IDB on COVID poverty transfers in Belize, Togo, Guinea



How Satellite Data And Artificial Intelligence Could Help Us Understand Poverty Better

New technology lets computers understand what they see in an image—or a million images.

Economics

Poverty Surveyors in Sri Lanka Get Some Help From Satellites Orbiting the Earth

The World Bank is teaming with a Silicon Valley startup to test whether poverty can be measured using satellite images.

By Adam Satariano

November 6, 2015, 7:00 AM PST Updated on November 6, 2015, 1:57 PM PST

In mountainous areas of Pakistan or far-flung villages in Sri Lanka, finding reliable economic information is extremely difficult. The World Bank's solution has been to send surveyors to study the conditions on the ground, which is an expensive, time-consuming, and imprecise task. The resulting dearth of data leaves governments, aid groups, and researchers unsure of where to put resources that can be critical to helping the world's most impoverished areas.



BY MAYA CRAIG 3 MINUTE READ

Data analytics firm Orbital Insight is partnering with the World Bank to test technology that could help measure global poverty using satellite imagery and artificial intelligence.

More “Business” Research

- Online Media Piracy

Forbes

There's Hope To Combat Piracy If Hollywood, Industry, and Government Unite

 Nelson Granados Contributor @
Hollywood & Entertainment
I cover digital trends in travel, media and entertainment.

⌚ This article is more than 5 years old.

Several studies have shown that piracy hurts the revenues of content owners, and instead pirate sites are reaping hundreds of millions of dollars in online advertising. Yet theft of movies and TV content seems to be as rampant today as ever. The Motion Picture Association of America (MPAA) reports that in 2014, just in the U.S. alone, 710 million movies and TV shows were shared via BitTorrent sites. Extrapolating to a global scale (the U.S. is less than 5% of the world's population) and adding streaming and other piracy methods, losses were likely in the billions of dollars. The staggering order of magnitude may lead some to wonder if it's even worth fighting the battle, or if it has been lost already. Can the battle against piracy be won? If so, how?

- IT Strategy

How APIs Create Growth by Inverting the Firm

Seth G. Benzell*, Jonathan Hersh† Marshall Van Alstyne ‡

This draft: August 7, 2021

Abstract

How might technology increase firm value? One method might be to facilitate more efficient use of internal capital. Another method might be to help the firm tap third party capital. This paper uses four unique data sets to measure growth in firm value based on adoption of Application Programming Interfaces (APIs), a technology that lets firms modularize and reconfigure resources for internal use or expose them to third parties for external use. The latter includes apps and services of the platform economy. We perform difference-in-difference and synthetic control analyses of financial outcomes for public firms and find that adopters of externally facing APIs grew an additional 38% over 16 years relative to non-adopters. Internal use cases were inconclusive. Using proprietary data on private APIs, we find that firms with public APIs grew faster after adoption than firms with private APIs. Then, using a Tobin's Q framework, we measure whether API adopting firms grew by lowering capital adjustment costs. Consistent with an inverted firm hypothesis, where value creation moves from inside to outside, we find that using the technology for external value creation explains more firm growth than using it for internal value creation. Finally, we document an important downside of API adoption: increased risk of data breach. Together these facts lead us to conclude that APIs, as the foundation of digital ecosystems, have a large and positive impact on economic growth and do so primarily by enabling external complementors rather than boosting internal productivity.

Most Proud of: Cited on the Wikipedia Page for “Waffle”



WIKIPEDIA
The Free Encyclopedia

Main page
Contents
Current events
Random article
About Wikipedia
Contact us
Donate

Contribute
Help
Community portal
Recent changes
Upload file

Tools
What links here
Related changes
Special pages
Permanent link
Page information
Cite this page
Wikidata item

Print/export
Download as PDF
Printable version

Not logged in Talk Contributions Create account Log in

Article Talk Read View source View history

Search Wikipedia



Waffle



From Wikipedia, the free encyclopedia

This article is about the batter/dough-based food. For other uses, see Waffle (disambiguation).

A **waffle** is a dish made from leavened batter or dough that is cooked between two plates that are patterned to give a characteristic size, shape, and surface impression. There are many variations based on the type of waffle iron and recipe used. Waffles are eaten throughout the world, particularly in Belgium, which has over a dozen regional varieties.^[1] Waffles may be made fresh or simply heated after having been commercially cooked and frozen.



Waffle

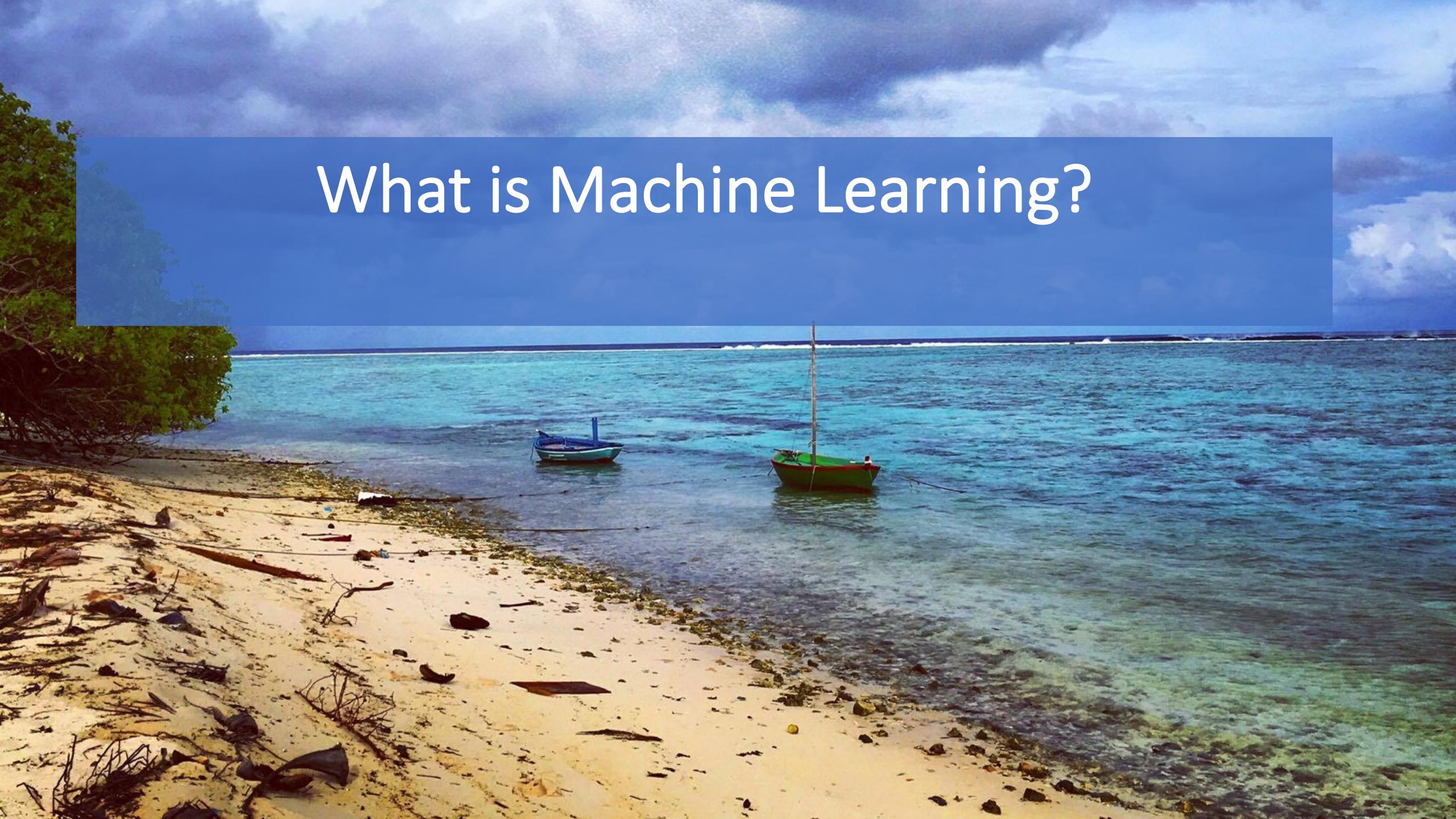
52. ^ a b "Sweet Diversity: Overseas Trade and Gains from Variety after 1492" Archived 2013-07-26 at the Wayback Machine, Jonathan Hersh, Hans-Joachim Voth, Real Sugar Prices and Sugar Consumption Per Capita in England, 1600–1850, p.42

Place of origin	France, Belgium
Main ingredients	Batter or dough
Variations	Liège waffle, Brussels Waffle, Flemish Waffle, Bergische waffle, Stroopwafel and others
Cookbook: Waffle	
Media: Waffle	

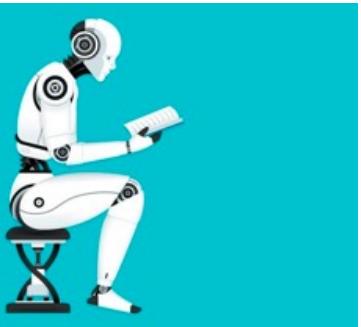
References

1. ^ "Les Gaufres Belges" Archived 2012-08-20 at the Wayback Machine. Gaufresbelges.com. Retrieved on 2013-04-07.
2. ^ Robert Smith (1725). *Court Cookery*. p. 176 .
3. ^ "Waffle" Archived 2013-04-07 at the Wayback Machine, The Merriam-Webster Unabridged Dictionary

What is Machine Learning?



Public Conception of Machine Learning



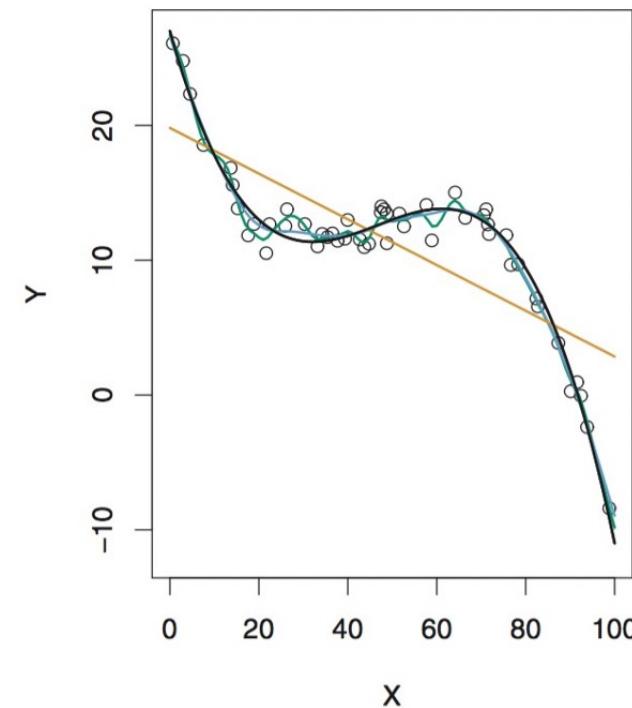
SkyNet



Reality (90% of the time)

Target or
Output

$$\hat{y} = \hat{f}(x)$$



Machine Learning Versus Econometrics

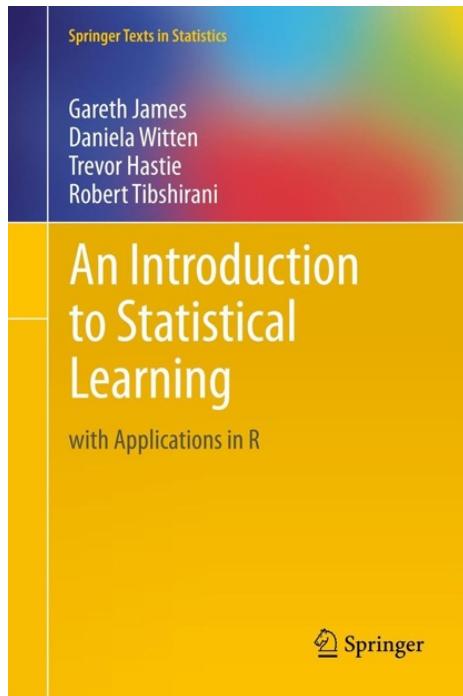
- **Machine Learning**

- Developed to solve problems in computer science
- Prediction/classification
- Desire: goodness of fit
- Huge Datasets! (Terabytes)
Thousands of variables!
- Whatever works

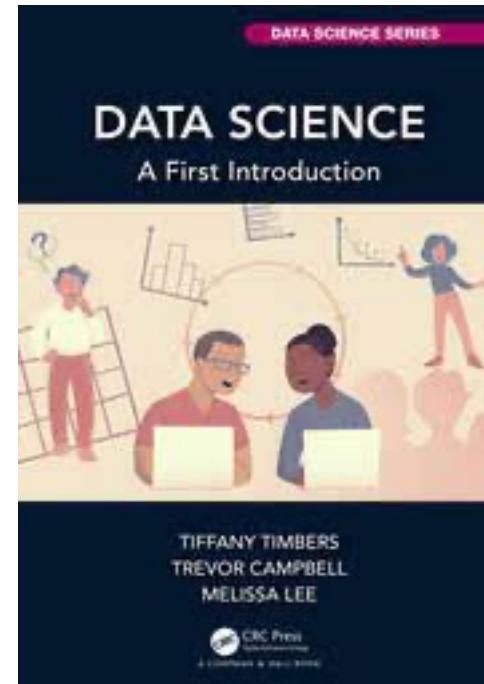
- **Econometrics**

- Developed to solve problems in economics
- Explicitly testing a theory
- “Statistical significance” more important than model fit
- Small datasets
Few dozen variables
- “It works in practice, but what about theory?”

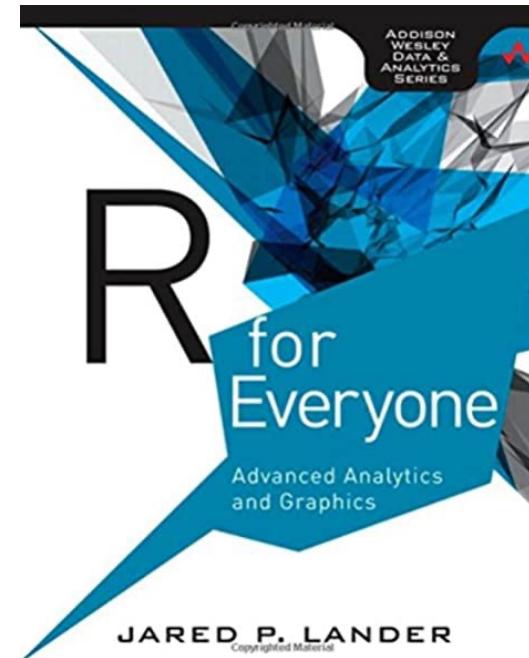
This Class – Introduction to Machine Learning with R



<https://www.statlearning.com/>

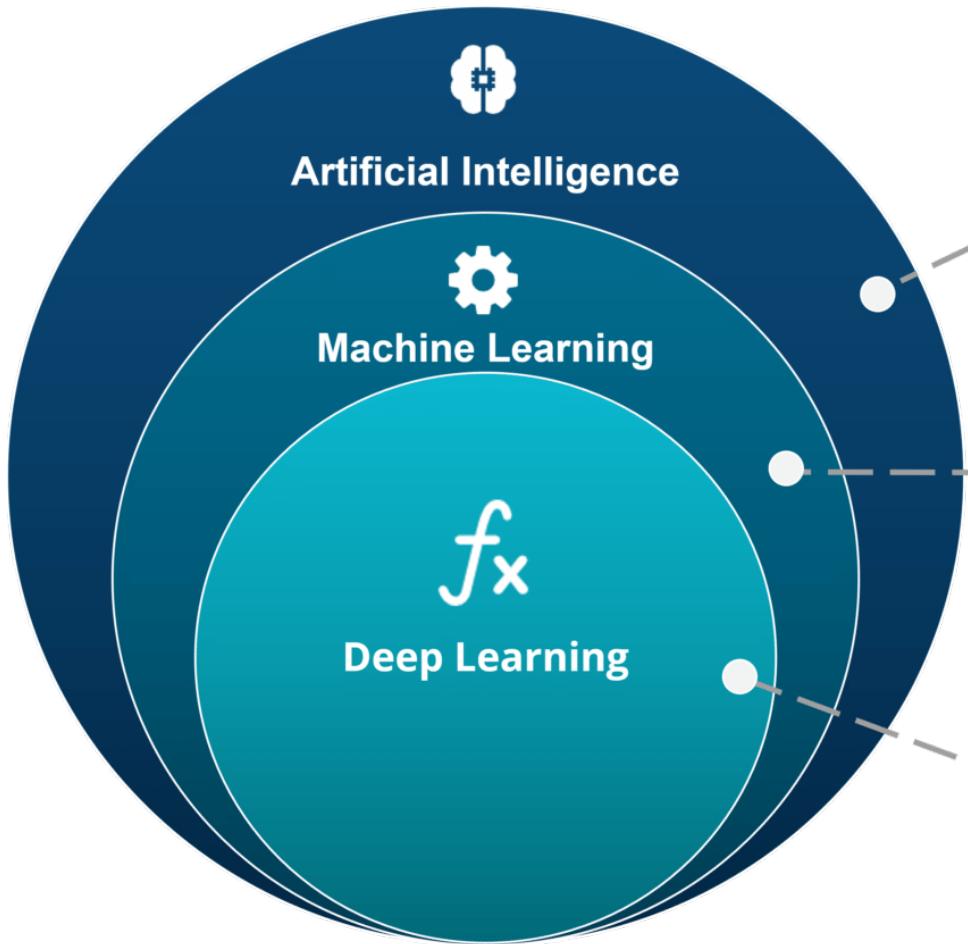


<https://datasciencebook.ca>



<https://www.jaredlander.com/r-for-everyone/>

Machine Learning Versus Artificial Intelligence



ARTIFICIAL INTELLIGENCE

A technique which enables machines to mimic human behaviour

MACHINE LEARNING

Subset of AI technique which use statistical methods to enable machines to improve with experience

DEEP LEARNING

Subset of ML which make the computation of multi-layer neural network feasible

Why Machine Learning?



Arguments for Using Machine Learning for Public Policy

1. Needed ML Big Data (models with 100+ variables)

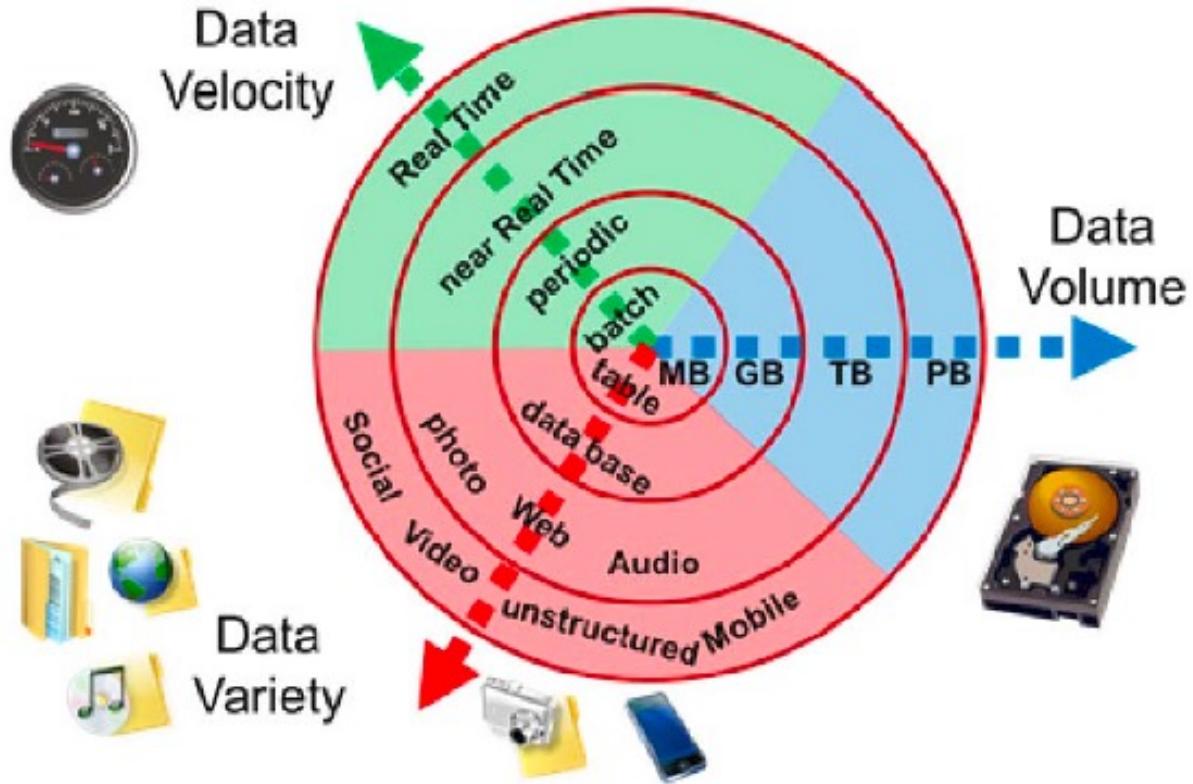
- “Unstructured” data e.g. satellite imagery, text

2. Can learn non-linear relationships

3. Better forecasts / econometrics

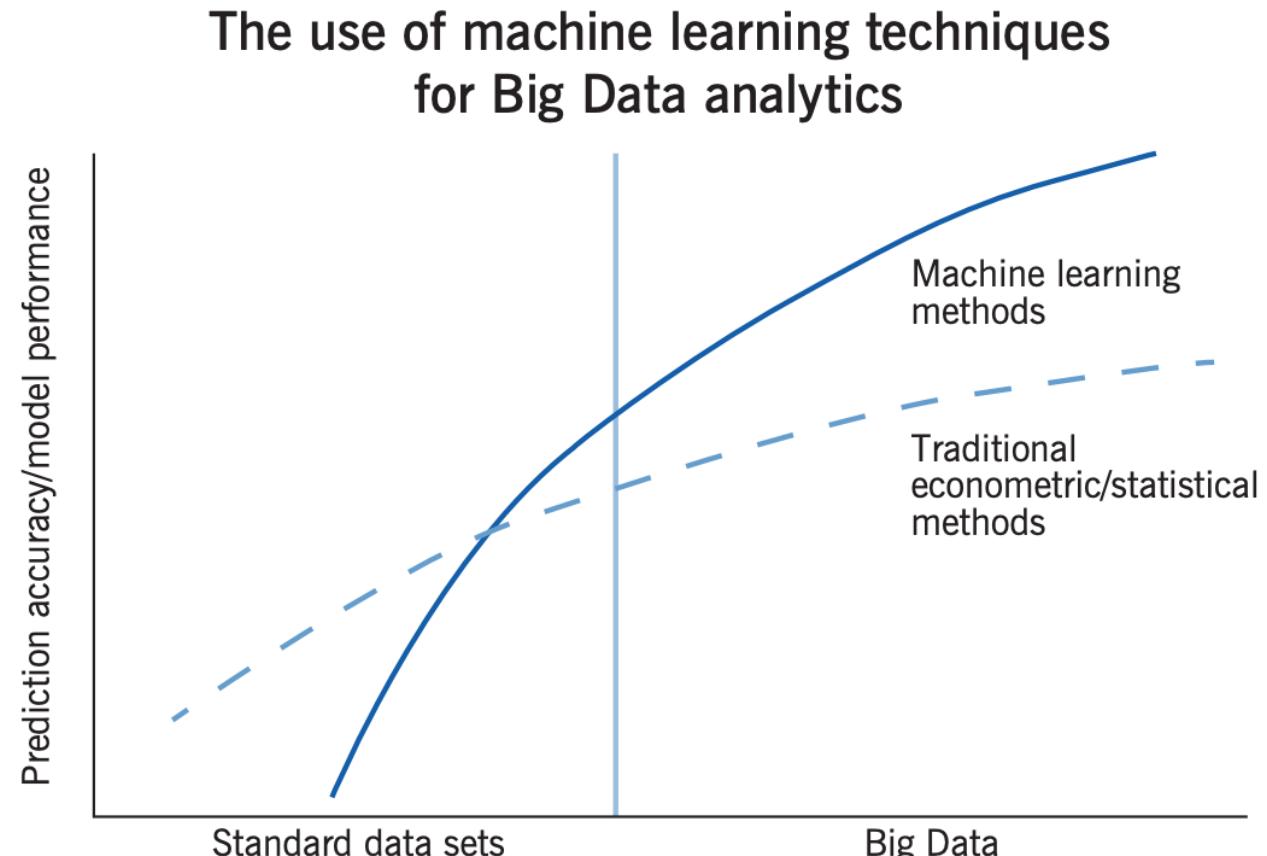
4. Anomaly detection (for fraud detection)

What is Big Data?



- **Big data is Data with Three “v’s”**
 - High volume
 - High variety
 - High velocity

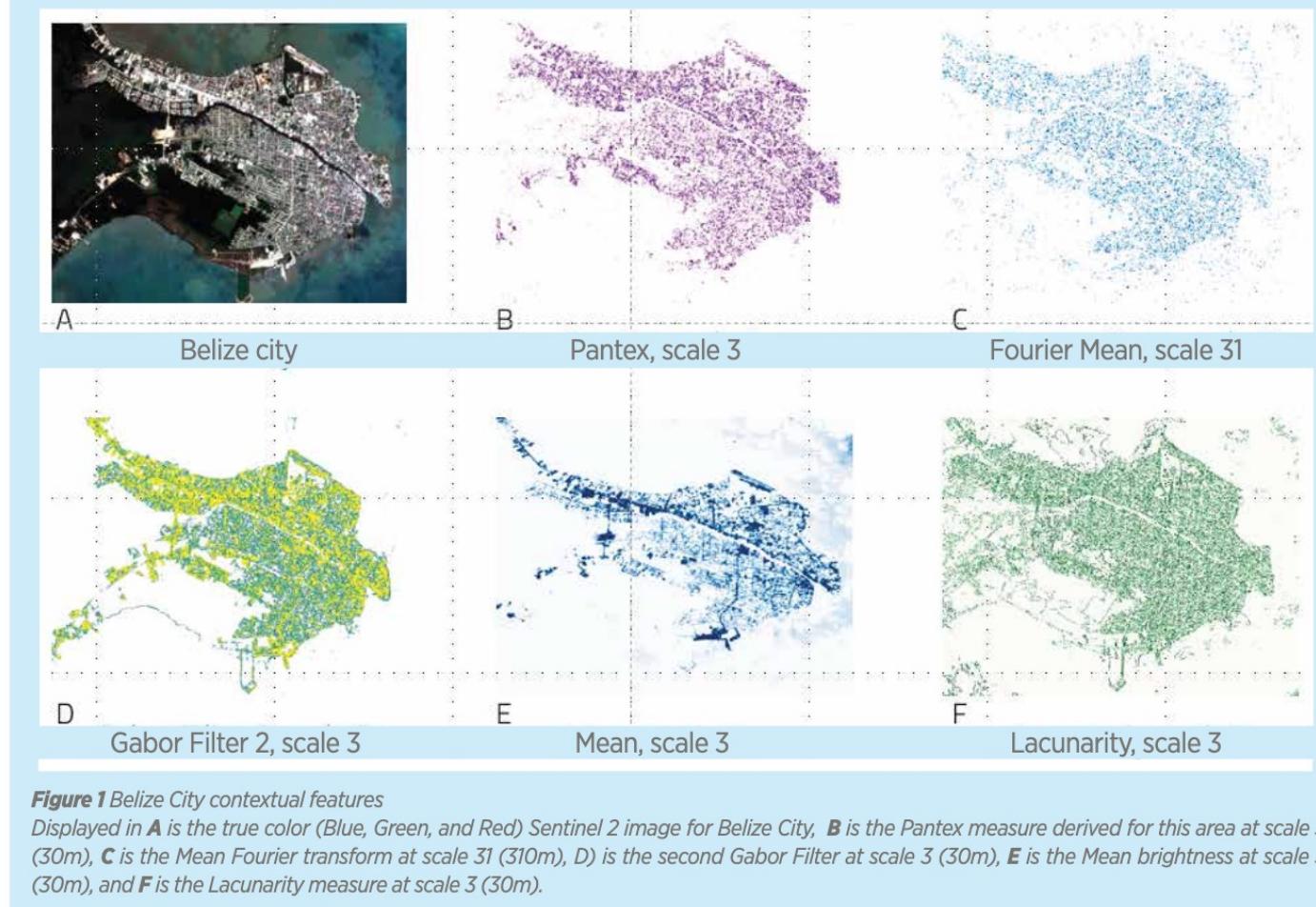
Why ML? (Big Data Needs Machine Learning)



Source: Author's own compilation.

- Machine Learning models continue to improve given more data (both # of variables and # of observations)
- Bigger datasets: bigger gain from machine learning vs econometrics

Why ML? (To Use Satellite Imagery “Big Data”)

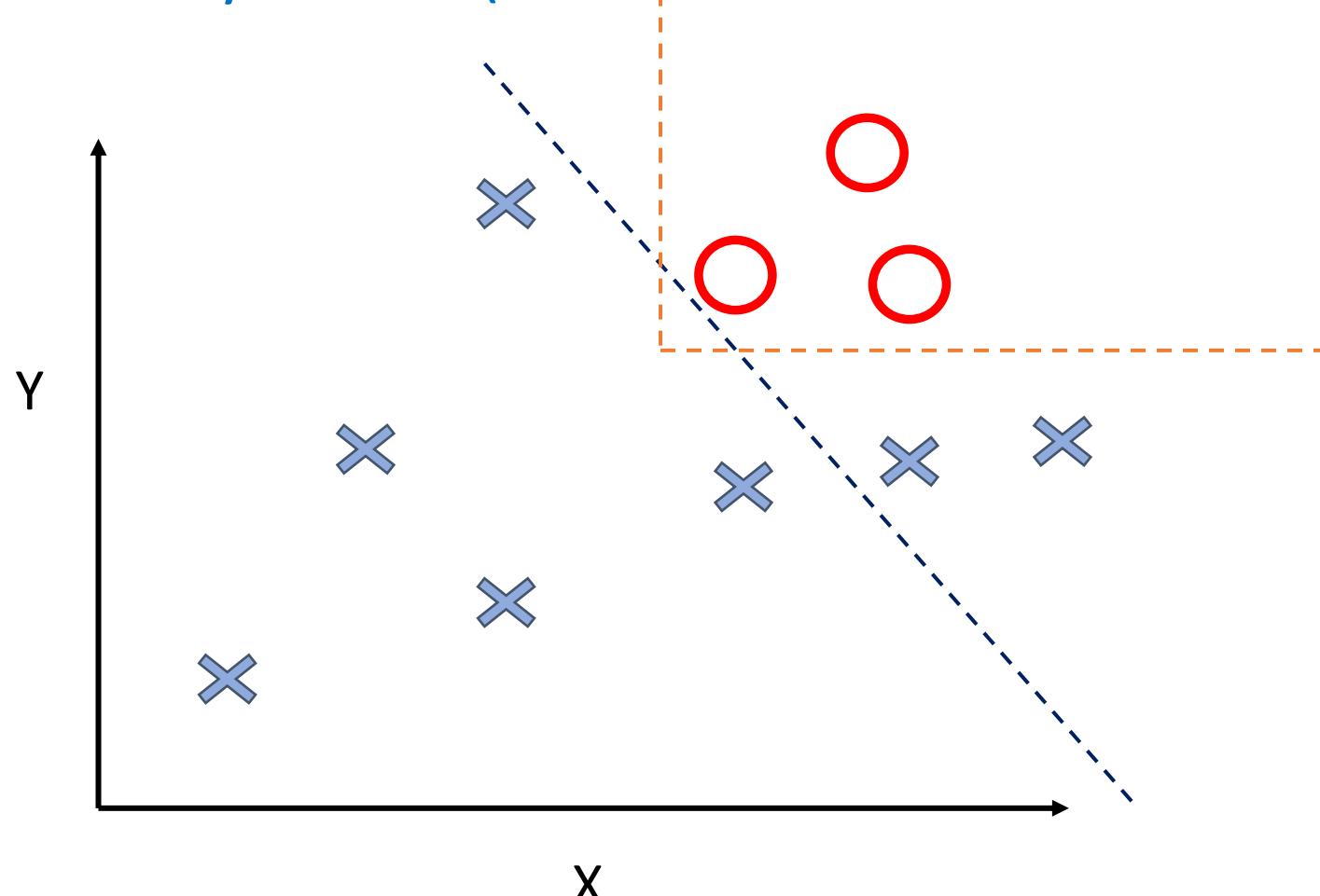


- Satellite Imagery variables too high dimensional for traditional econometric models

- Mapping Poverty in Belize Using Satellite Imagery

<https://publications.iadb.org/publications/english/document/Mapping-Income-Poverty-in-Belize-Using-Satellite-Features-and-Machine-Learning.pdf>

Why ML? (Can learn Nonlinear relationships)



- Example: classify “O”s separate from X’s

Econometrics: $y = X * \beta$

Machine Learning: regression tree

Model	Accuracy
Econometrics	80%
Machine Learning	100%

Why ML? (Better Forecasts For Fiscal Crises)

Predicting Fiscal Crises: A Machine Learning Approach

Klaus-Peter Hellwig¹

International Monetary Fund, Asia Pacific Department

This version: July 27, 2020

Abstract: This paper assesses the ability of econometric and machine learning techniques to predict fiscal crises out of sample. We show that the standard econometric approach used in policy applications cannot outperform a heuristic rule of thumb derived from unconditional historical averages. Elastic net and tree ensemble methods (random forest, gradient boosted trees) deliver significant improvements in accuracy. Performance of machine learning techniques improves, particularly for developing countries, when expanding the set of potential predictors from a small set, preselected manually from the literature, to a large set (748 variables) and relying on algorithmic variable selection techniques. There is considerable agreement across learning algorithms in the set of selected predictors: Results confirm the importance of external sector stock and flow variables found in the literature but also point to demographics and the quality of governance as important predictors of fiscal crises. Fiscal variables appear to have less predictive value, and public debt matters only to the extent that it is owed to external creditors.

Why ML? (Better Forecasts of Inflation)

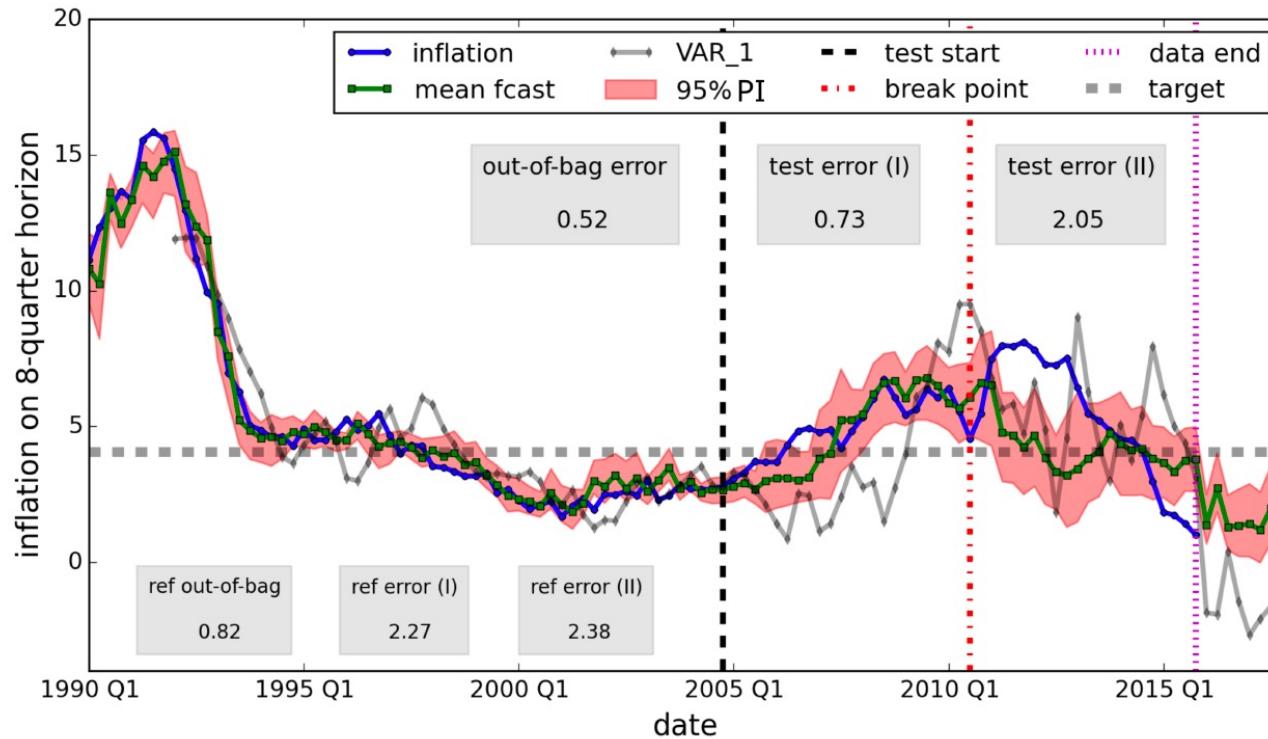
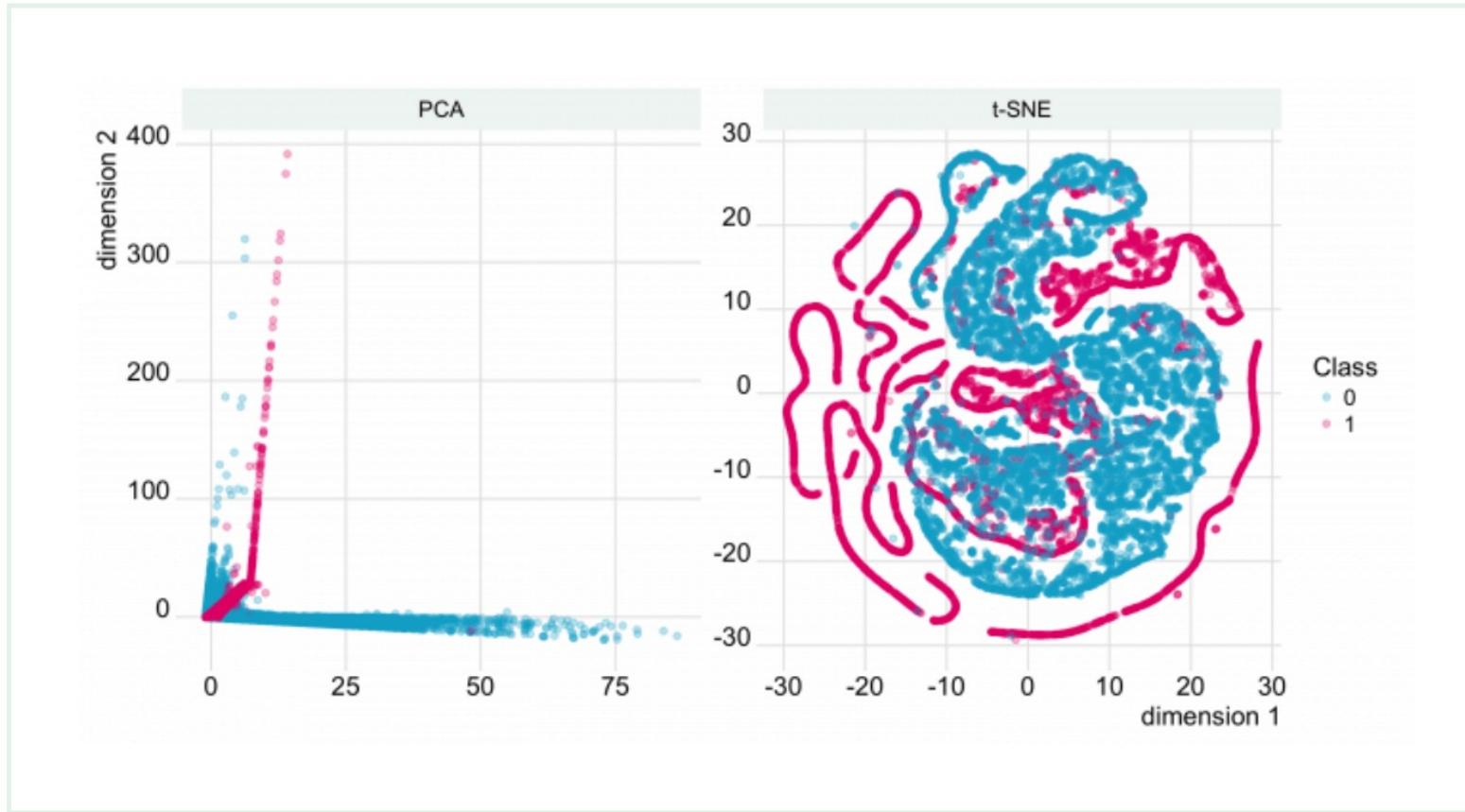


Figure 16: Averaged bootstrapped SVM-dFFANN projections (green line) for two-year changes in CPI (blue line). The shaded band indicates the 95% prediction interval (PI) across bootstrapped models. The vertical dashed lines separates the initial training, pre-crisis, post-crisis and post-data periods, respectively. Model and VAR₁ reference errors are given in the boxes. Sources: BoE, ONS, BIS, World Bank and authors' calculations.

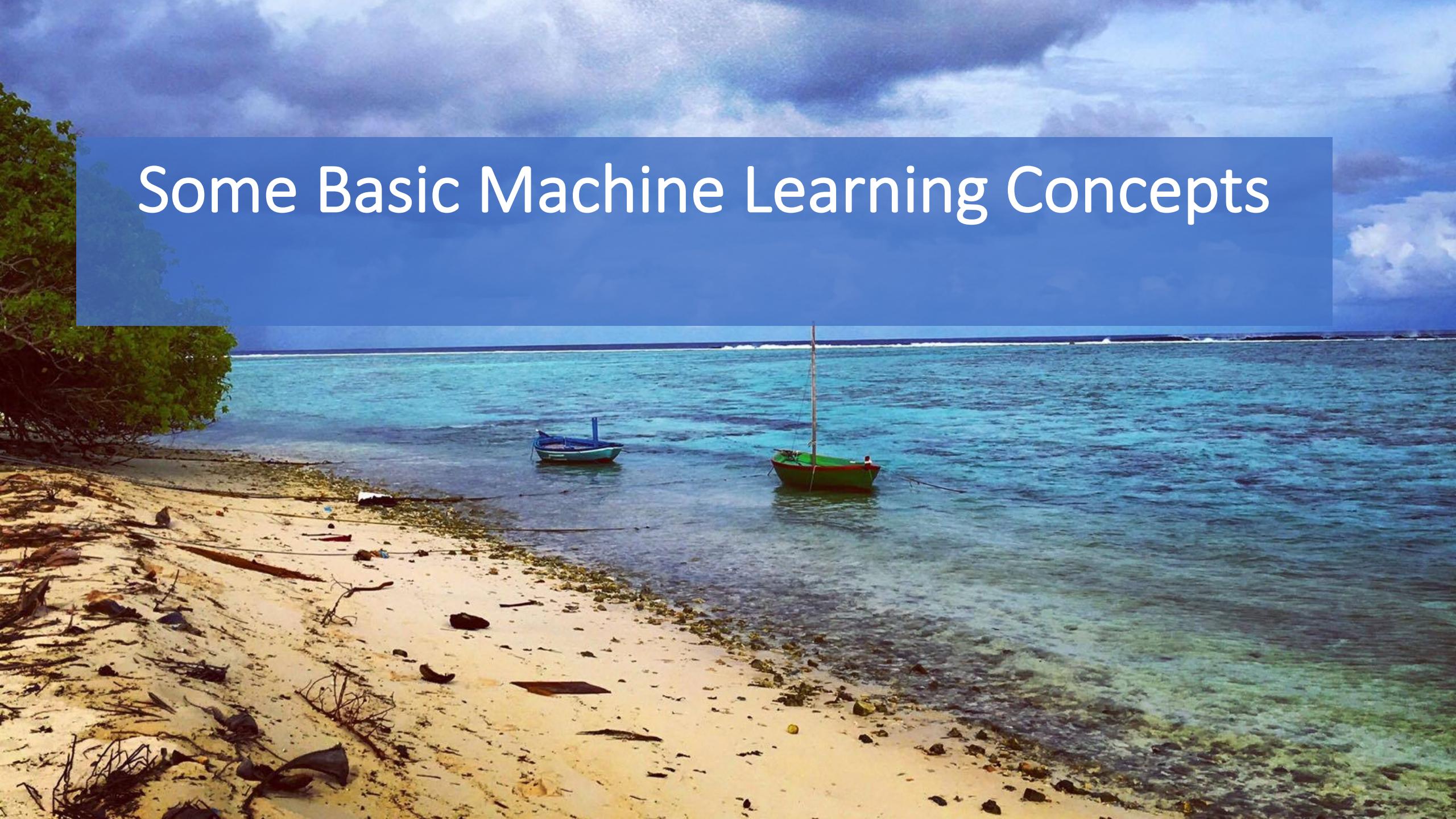
Why ML? (Anomaly Detection Aka Fraud Alerts)



Dimensionality reduction techniques in fraud analytics. The plots show the first two dimensions of PCA (left) and t-SNE (right) for fraudulent (Class = 1) and regular (Class = 0) transactions.

Source: https://shiring.github.io/machine_learning/2017/05/01/fraud

Some Basic Machine Learning Concepts



Supervised vs Unsupervised Learning

Supervised Learning:

- For every x_i we observe some y_i
- Ex: random forests to predict loan default (y_i) based on applicant characteristics (x_i)

Supervised Learning



Unsupervised Learning



[DataEnRanc.wordpress.com](http://dataenranc.wordpress.com)

Unsupervised Learning:

- We only observe x_i
- Ex: clustering loan applicants based on characteristics (x_i)

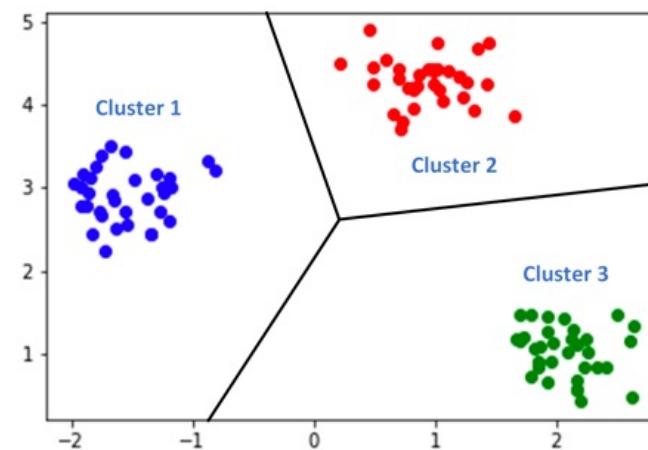
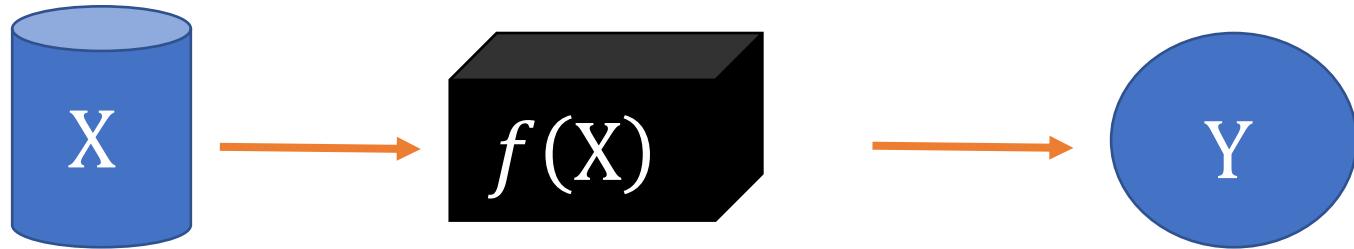


Fig.1. An Example Of Data Clustering

Supervised learning: learning $f(X)$ our predicted out given inputs

$$Y = f(X) + \epsilon$$



ϵ = “epsilon” (unexplained portion)

“Estimating” $\hat{f}(X)$

- $Y = f(X) + \epsilon$ is the true value
- We can only use data to “guess” at $f(X)$
- We call this guess $\hat{f}(X)$

How do we know when we’ve selected a “good” $\hat{f}(X)$?

- We reserve a portion of our data into a “test” set, estimate a model on the other part, and see how our model performs on this test set

Testing Training Data Subsets

Training set: (observation-wise) subset of data used to develop models

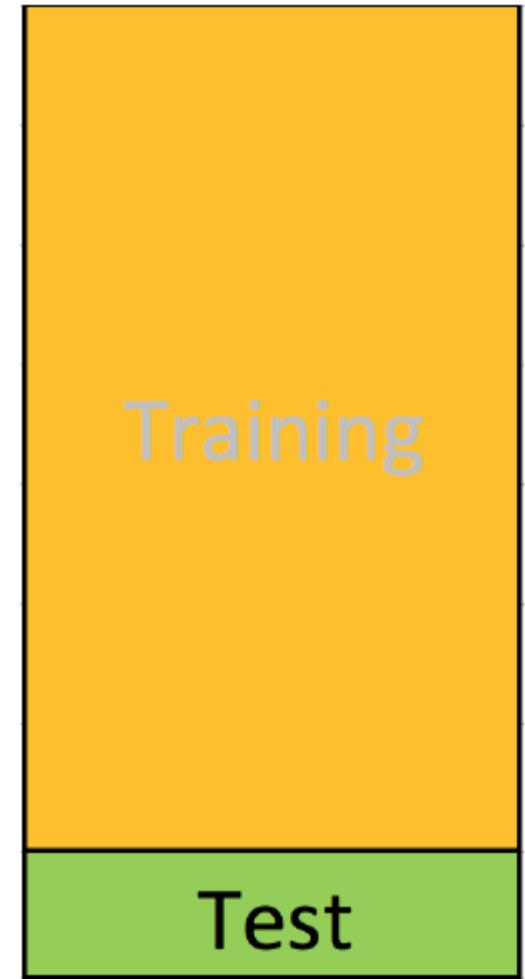


Testing/Training Split

Training set: (observation-wise) subset of data used to develop models

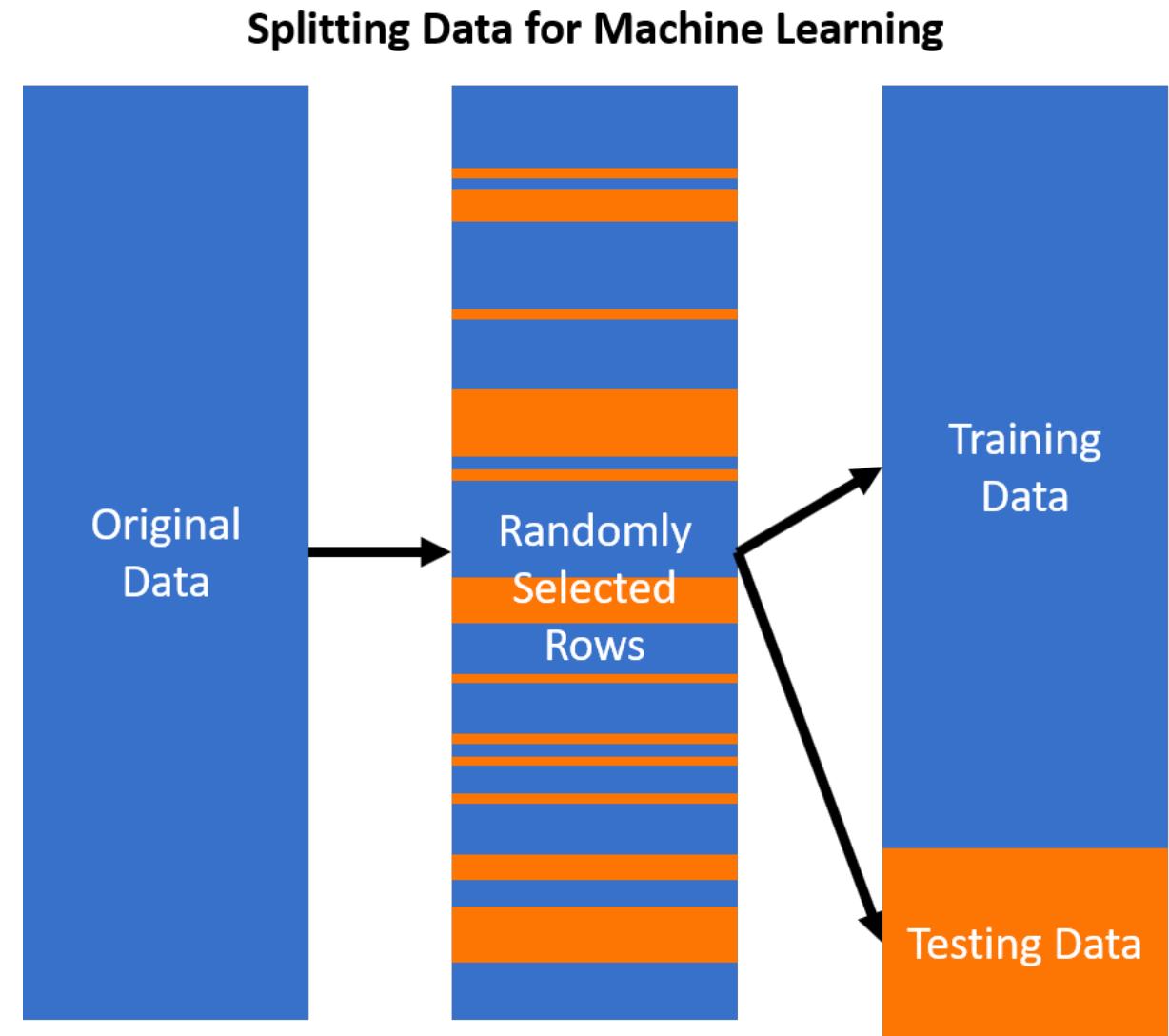
Test or Validation set: subset of data used during intermediate stages to “tune” model parameters

Rule of thumb 75% training 25% test -ish

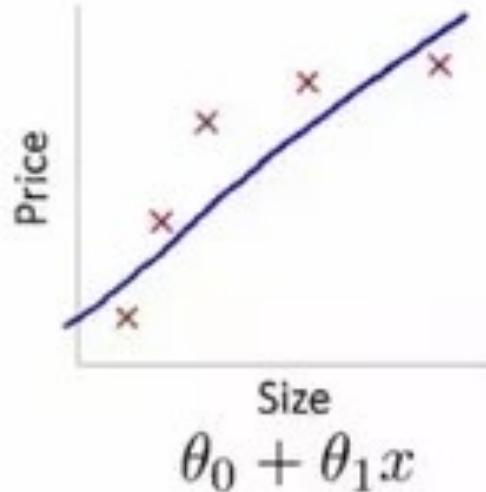


Randomly Selecting Rows for Test or Training Sets

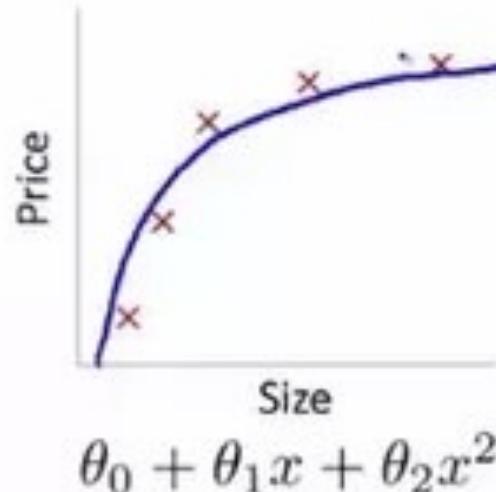
- Observations are randomly selected into either testing or training splits of the data



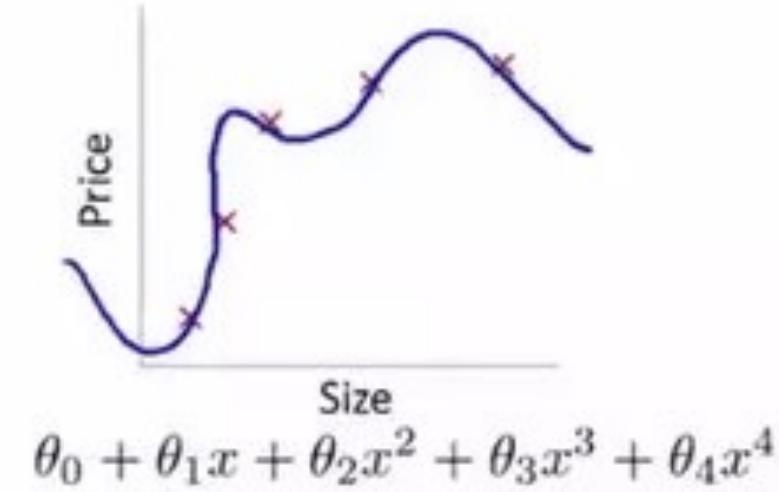
Optimal Model Complexity: Neither Underfit Nor Overfit



High bias
(underfit)

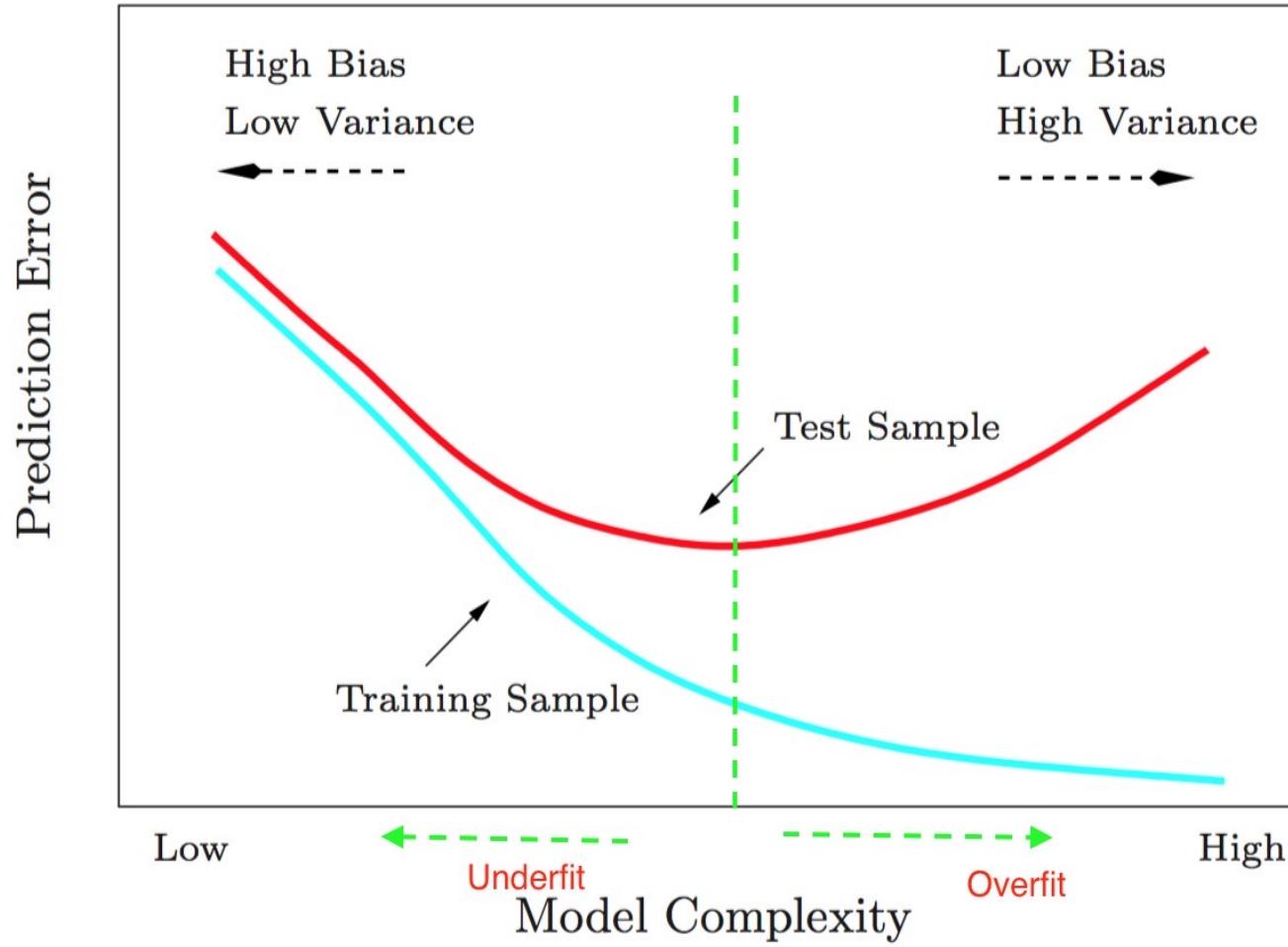


"Just right"



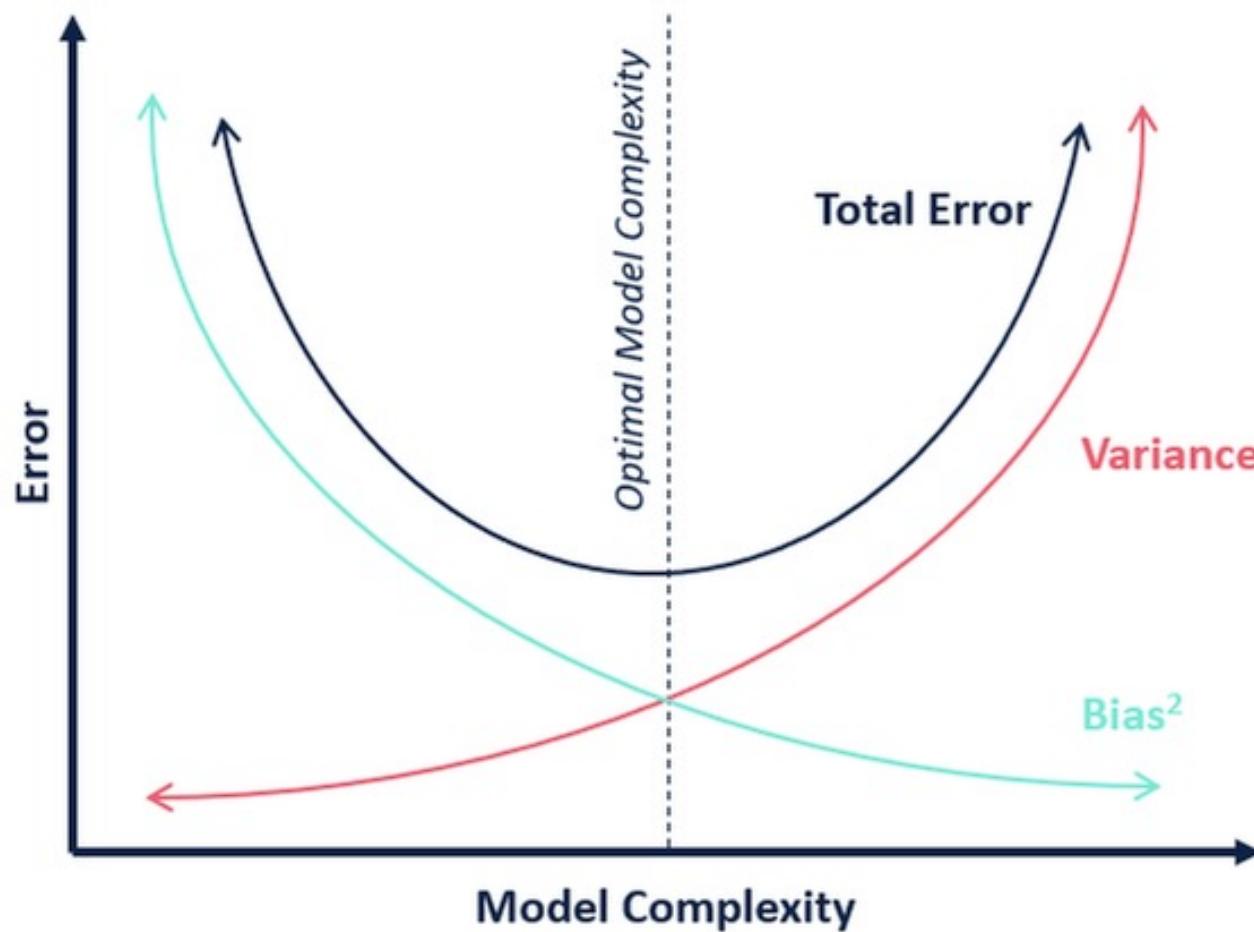
High variance
(overfit)

Bias-Variance Tradeoff



- Error in Training sample (~bias) \downarrow as we \uparrow model complexity (e.g. number of variables)
- Error in Test sample (~variance) \uparrow as we \uparrow
- Key: finding optimal model complexity

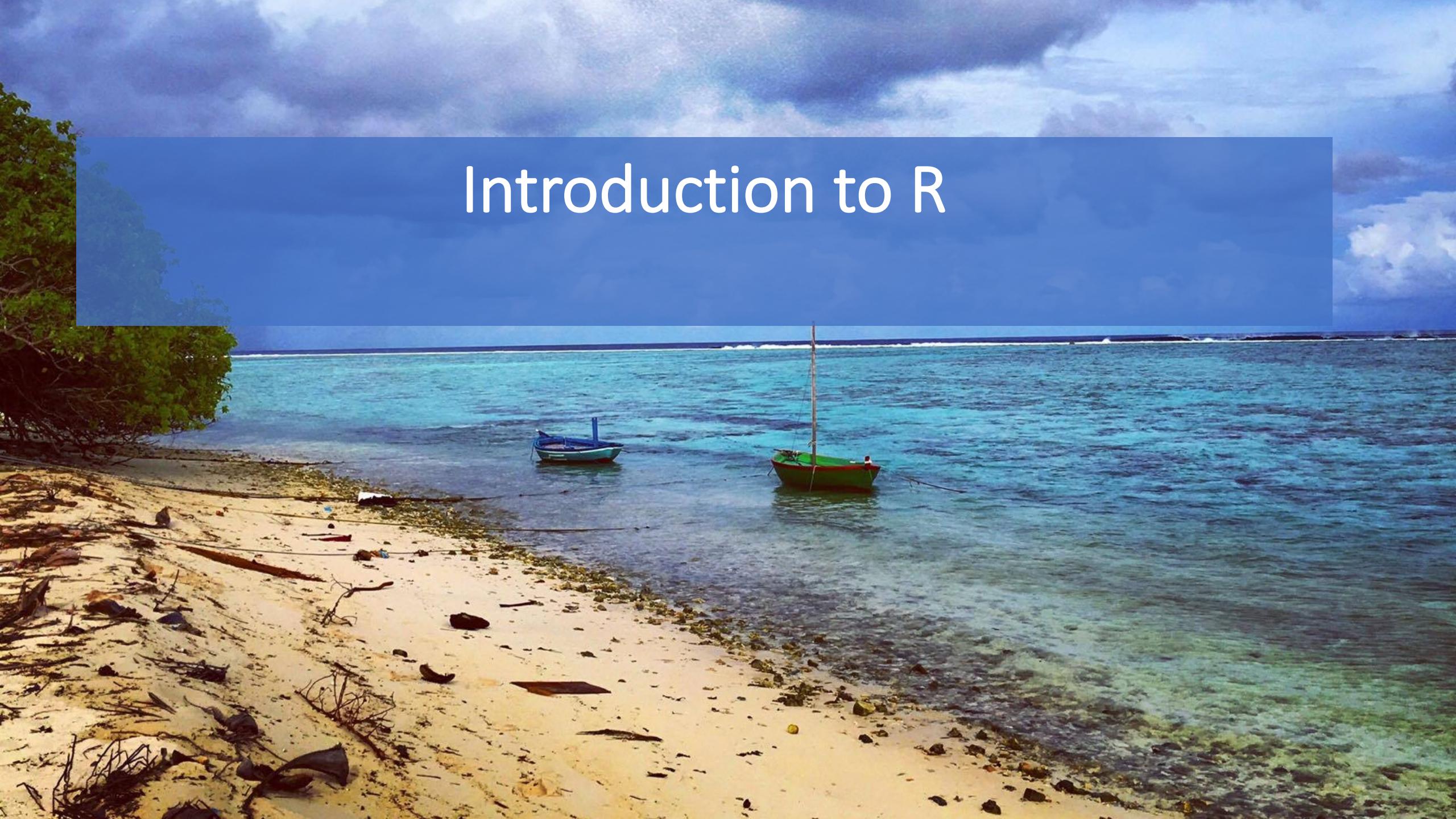
Key: Finding Optimal Model Complexity



Summary – Intro to Machine Learning

- **Machine Learning** is a set of methods developed to find robust patterns across datasets
- **Public Policy can benefit from machine learning.**
 - Big data requires it
 - Non-linear relationships
 - Better forecasts/econometrics
 - Anomaly detection
- **Remember these key concepts**
 - Supervised (Y,X) vs Unsupervised learning (just X)
 - Testing/Training Sets
 - (model -> train, see how it performs on test)
 - Bias-Variance Tradeoff
 - Bias – how far off model from true
 - Variance – precision of estimated model
 - Total error = bias² + variance

Introduction to R



Rstudio Projects

The screenshot shows a GitHub repository page for 'Belize_ML_2022'. The page contains instructions for loading the project into RStudio. It includes a bulleted list of steps and a code block for loading the project using the 'usethis' package.

If you cannot install those programs, please head over to [rstudio.cloud](#).

- Click "GET STARTED FOR FREE"
- Then click "Sign Up".
- You may log in using your email address
- Next click new project.
- You should now see a R Studio session in your browser.

Using Github

If you have never used Github, don't worry. You can either clone the repository, or you may click the "Code" button on the main page, and then "Download Zip" to download all the files. You may also download the files individually, or copy and paste code as needed.

Loading Project in R

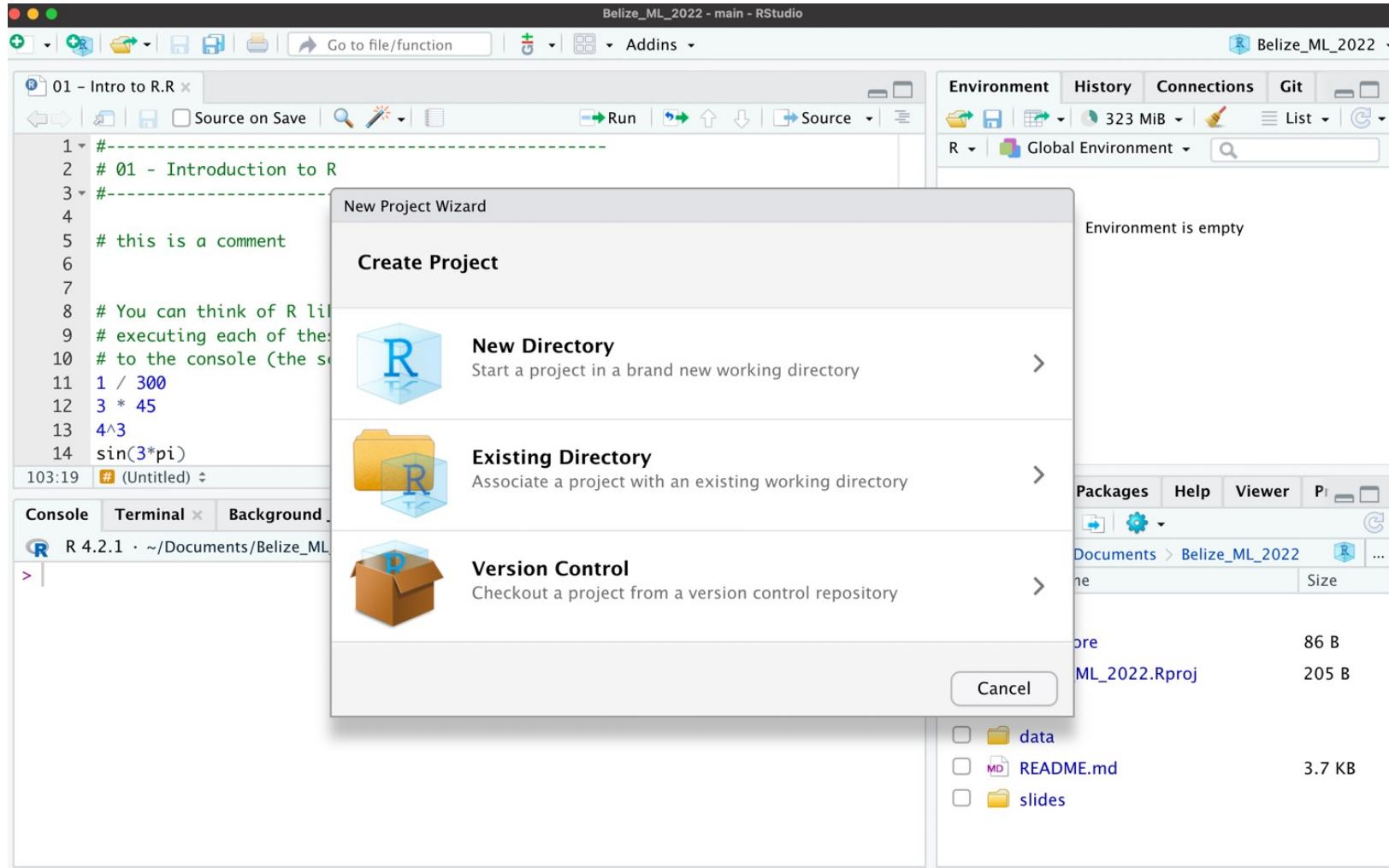
```
install.packages('usethis')
install.packages('tidyverse')

newProject <- usethis::use_course('https://github.com/jonhersh/Belize_ML_2022/archive/main.zip')
```

© 2022 GitHub, Inc. Terms Privacy Security Status Docs Contact GitHub Pricing API Training Blog About

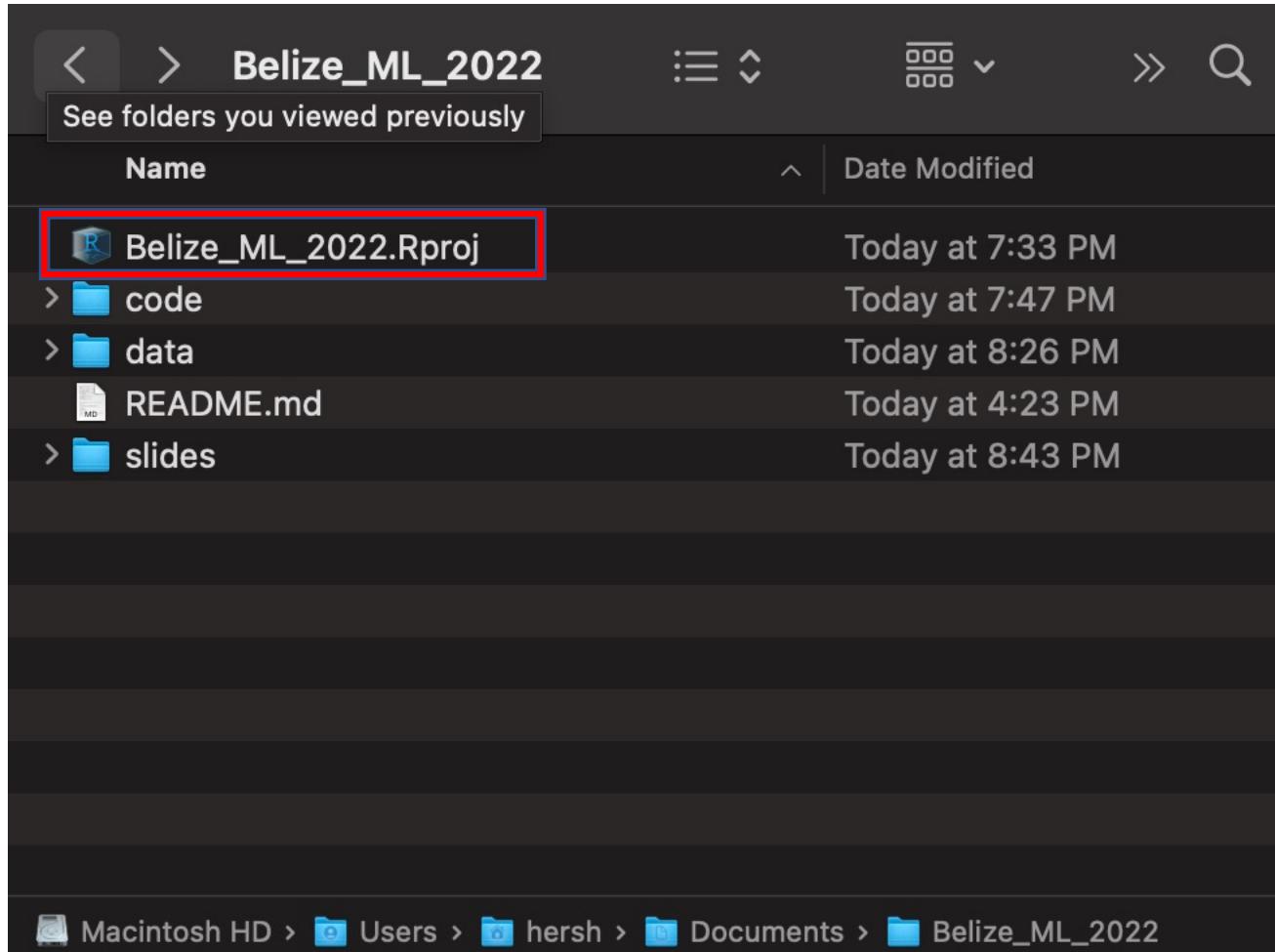
https://github.com/jonhersh/Belize_ML_2022

Rstudio Projects



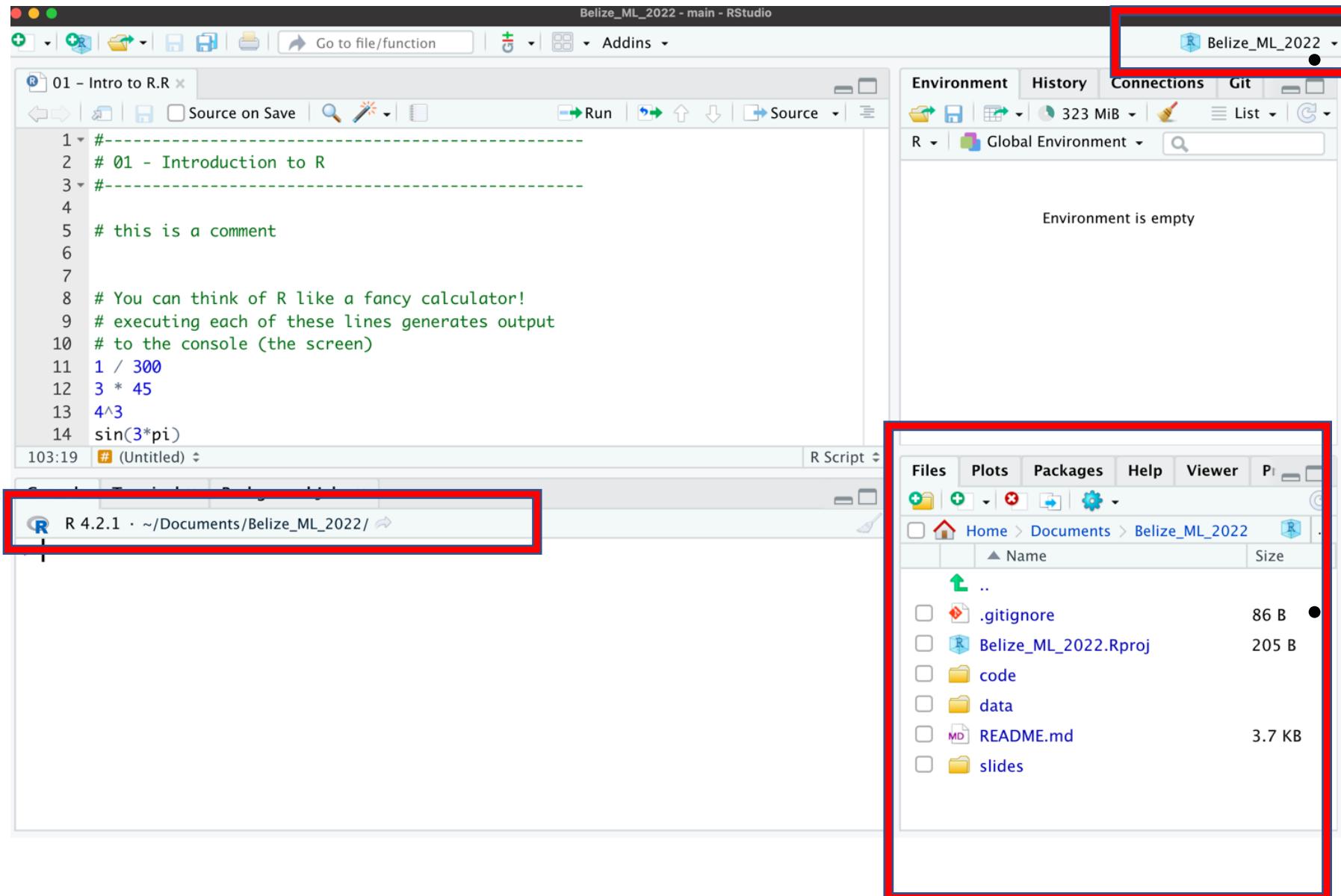
- You can also create new projects by going to file -> new project

Rstudio Projects



- This creates an .Rproj file you should click and open whenever you want to run code for this class

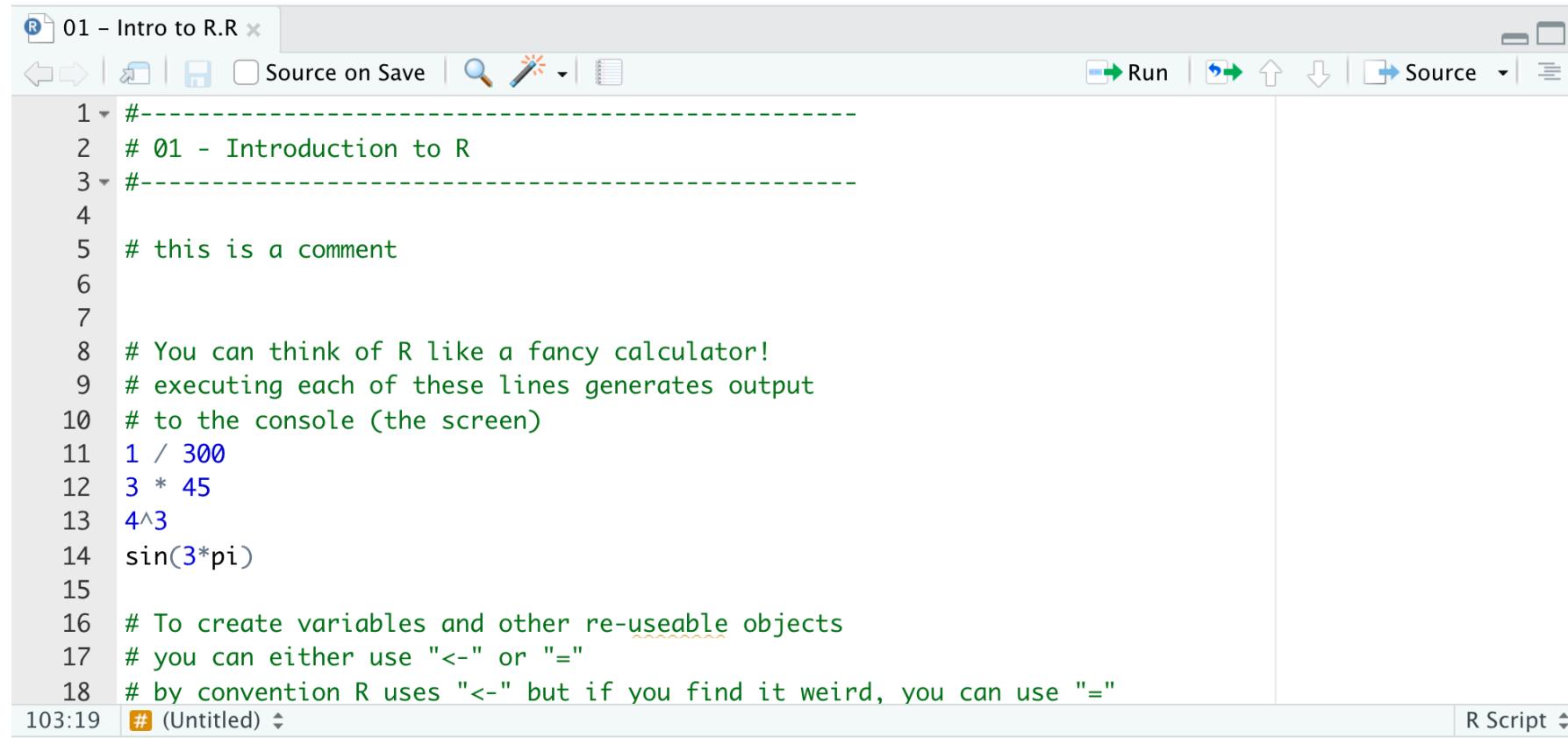
Rstudio Projects



You should see your project name in the top right and the “root directory” of the project in the bottom right

All links to files, datasets are relative to this root folder

Lab 1 – Intro to R



The screenshot shows the RStudio interface with an R script file open. The title bar reads "01 - Intro to R.R". The toolbar includes standard icons for file operations, search, and source code. The main editor area contains the following R code:

```
1 #-----
2 # 01 - Introduction to R
3 #-----
4
5 # this is a comment
6
7
8 # You can think of R like a fancy calculator!
9 # executing each of these lines generates output
10 # to the console (the screen)
11 1 / 300
12 3 * 45
13 4^3
14 sin(3*pi)
15
16 # To create variables and other re-useable objects
17 # you can either use "<->" or "="
18 # by convention R uses "<->" but if you find it weird, you can use "="
```

The status bar at the bottom left shows the line number "103:19" and the file name "# (Untitled)". The status bar at the bottom right shows "R Script".