

Machine Learning

PHYS 453

Dr. Daugherty



ABILENE
CHRISTIAN
UNIVERSITY

Welcome to ML

Today's goal:

- what topics this class covers (i.e. what is ML?)
- what it will be like
- why am I teaching it

Part 0

WHY IS A PHYSICIST TEACHING THIS?

About Me

Bio:

- Grew up in Oklahoma City
- **2002 ACU: Physics and CS**
- 2008 PhD in Nuclear Physics from UT
- Married in 2000, two kids

Research:

- Particle and Nuclear Physics
- Atom smashers
- Radiation Detectors
- Other interests: artificial intelligence, cosmology, amateur theology

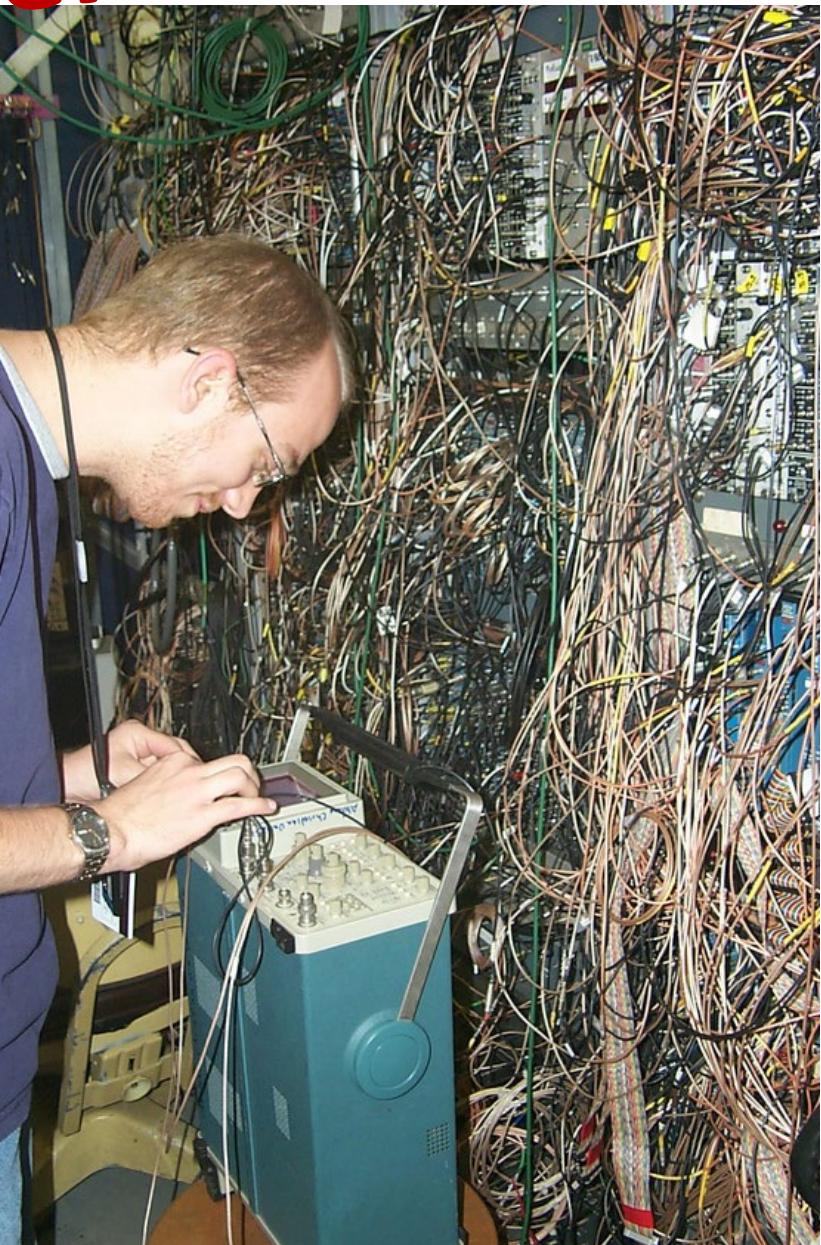
Fun Facts:

- Have been paid as a professional model
- Once burned off an eyebrow in class
- Plays drums in Abilene's best cover band

Hobbies:

- Eating
- Putting things in a laser
- Building dangerous things

First Experiment



PhD Research

Two-Particle Correlations in Ultra Relativistic Heavy
Ion Collisions

by

Michael Scott Daugherty, B.S.; B.S.

Dissertation

Presented to the Faculty of the Graduate School of

The University of Texas at Austin

in Partial Fulfillment

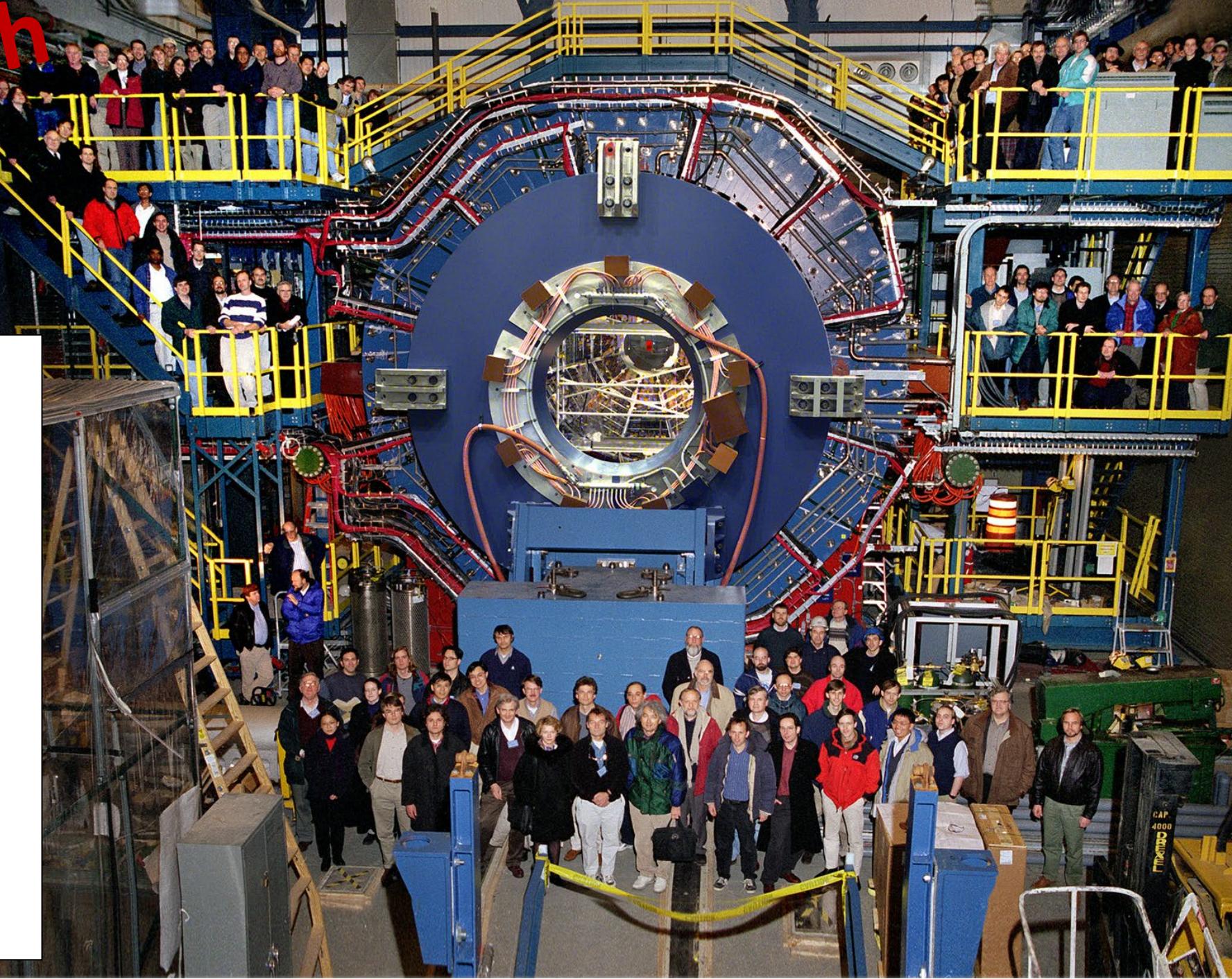
of the Requirements

for the Degree of

Doctor of Philosophy

The University of Texas at Austin

August 2008



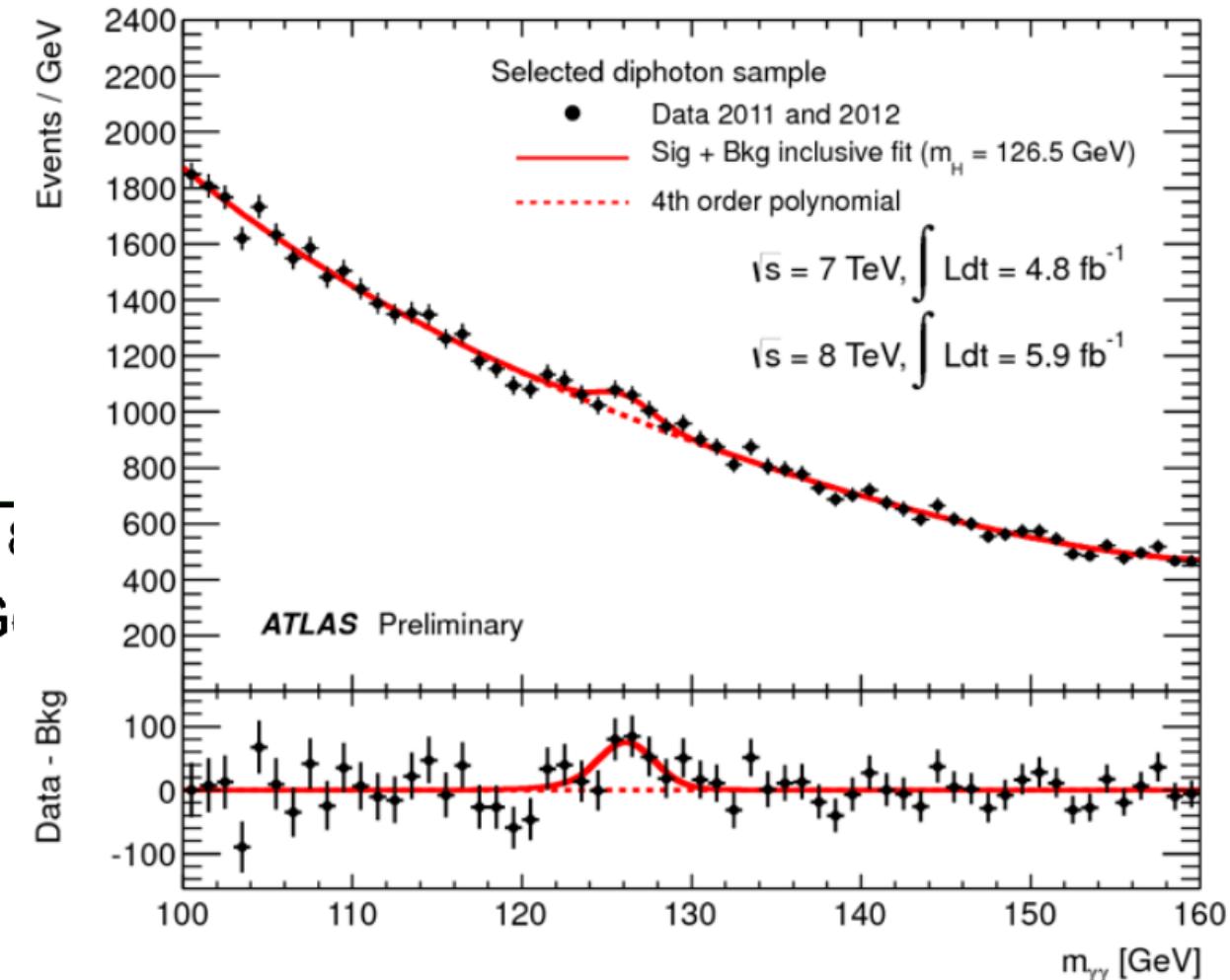
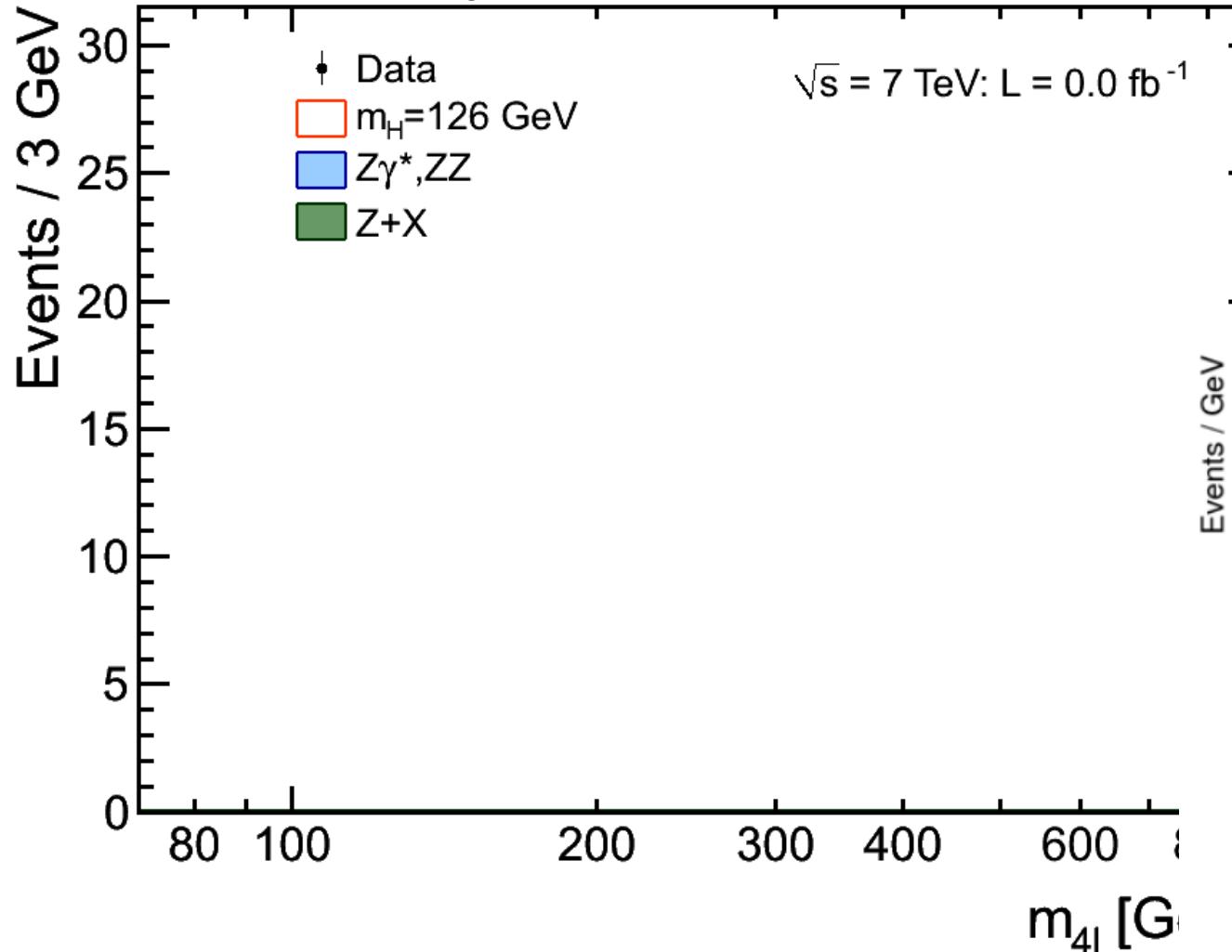
Before “Big Data” was a thing, there was high-energy physics...

]] 16 Mar 2010

A Search for the Higgs Boson Using Neural Networks in Events with Missing Energy and b -quark Jets in $p\bar{p}$ Collisions at $\sqrt{s} = 1.96$ TeV

T. Aaltonen,²⁴ J. Adelman,¹⁴ B. Álvarez González^v,¹² S. Amerio^{dd},⁴⁴ D. Amidei,³⁵ A. Anastassov,³⁹ A. Annovi,²⁰ J. Antos,¹⁵ G. Apollinari,¹⁸ A. Apresyan,⁴⁹ T. Arisawa,⁵⁸ A. Artikov,¹⁶ J. Asaadi,⁵⁴ W. Ashmanskas,¹⁸ A. Attal,⁴ A. Aurisano,⁵⁴ F. Azfar,⁴³ W. Badgett,¹⁸ A. Barbaro-Galtieri,²⁹ V.E. Barnes,⁴⁹ B.A. Barnett,²⁶ P. Barria^{ff},⁴⁷ P. Bartos,¹⁵ G. Bauer,³³ P.-H. Beauchemin,³⁴ F. Bedeschi,⁴⁷ D. Beecher,³¹ S. Behari,²⁶ G. Bellettini^{ee},⁴⁷ J. Bellinger,⁶⁰ D. Benjamin,¹⁷ A. Beretvas,¹⁸ A. Bhatti,⁵¹ M. Binkley,¹⁸ D. Bisello^{dd},⁴⁴ I. Bizjak^{jj},³¹ R.E. Blair,² C. Blocker,⁷ B. Blumenfeld,²⁶ A. Bocci,¹⁷ A. Bodek,⁵⁰ V. Boisvert,⁵⁰ D. Bortoletto,⁴⁹ J. Boudreau,⁴⁸ A. Boveia,¹¹ B. Brau^a,¹¹ A. Bridgeman,²⁵ L. Brigliadori^{cc},⁶ C. Bromberg,³⁶ E. Brubaker,¹⁴ J. Budagov,¹⁶ H.S. Budd,⁵⁰ S. Budd,²⁵ K. Burkett,¹⁸ G. Busetto^{dd},⁴⁴ P. Bussey,²² A. Buzatu,³⁴ K. L. Byrum,² S. Cabrera^x,¹⁷ C. Calancha,³² S. Camarda,⁴ M. Campanelli,³¹ M. Campbell,³⁵ F. Canelli¹⁴,¹⁸ A. Canepa,⁴⁶ B. Carls,²⁵ D. Carlsmith,⁶⁰ R. Carosi,⁴⁷ S. Carrilloⁿ,¹⁹ S. Carron,¹⁸ B. Casal,¹² M. Casarsa,¹⁸ A. Castro^{cc},⁶ P. Catastini^{ff},⁴⁷ D. Cauz,⁵⁵ V. Cavaliere^{ff},⁴⁷ M. Cavalli-Sforza,⁴ A. Cerri,²⁹ L. Cerrito^q,³¹ S.H. Chang,²⁸ Y.C. Chen,¹ M. Chertok,⁸ G. Chiarelli,⁴⁷ G. Chlachidze,¹⁸ F. Chlebana,¹⁸ K. Cho,²⁸ D. Chokheli,¹⁶ J.P. Chou,²³ K. Chung^o,¹⁸ W.H. Chung,⁶⁰ Y.S. Chung,⁵⁰ T. Chwalek,²⁷ C.I. Ciobanu,⁴⁵ M.A. Ciocci^{ff},⁴⁷ A. Clark,²¹ D. Clark,⁷ G. Compostella,⁴⁴ M.E. Convery,¹⁸ J. Conway,⁸ M. Corbo,⁴⁵ M. Cordelli,²⁰ C.A. Cox,⁸ D.J. Cox,⁸ F. Crescioli^{ee},⁴⁷ C. Cuenca Almenar,⁶¹ J. Cuevas^v,¹² R. Culbertson,¹⁸ J.C. Cully,³⁵ D. Dagenhart,¹⁸ M. Datta,¹⁸ T. Davies,²² P. de Barbaro,⁵⁰ S. De Cecco,⁵² A. Deisher,²⁹ G. De Lorenzo,⁴ M. Dell'Orso^{ee},⁴⁷ C. Deluca,⁴ L. Demortier,⁵¹ J. Deng^f,¹⁷ M. Deninno,⁶ M. d'Errico^{dd},⁴⁴ A. Di Canto^{ee},⁴⁷ G.P. di Giovanni,⁴⁵ B. Di Ruzza,⁴⁷ J.R. Dittmann,⁵ M. D'Onofrio,⁴ S. Donati^{ee},⁴⁷ P. Dong,¹⁸ T. Dorigo,⁴⁴ S. Dube,⁵³ K. Ebina,⁵⁸ A. Elagin,⁵⁴ R. Erbacher,⁸ D. Errede,²⁵ S. Errede,²⁵ N. Ershaidat^{bb},⁴⁵ R. Eusebi,⁵⁴ H.C. Fang,²⁹ S. Farrington,⁴³ W.T. Fedorko,¹⁴ R.G. Feild,⁶¹

CMS Preliminary



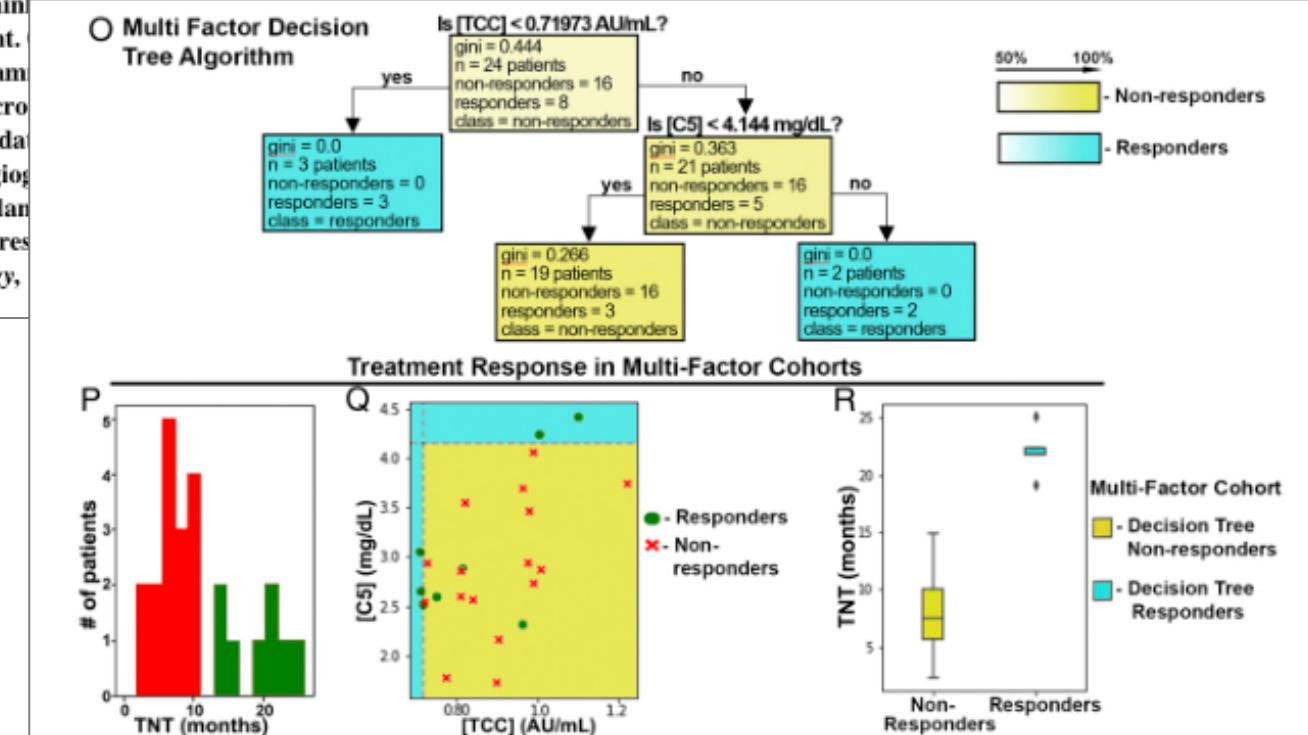
One More Shameless Example

<https://www.jimmunol.org/content/early/2020/11/05/jimmunol.2000511>

Complement as Prognostic Biomarker and Potential Therapeutic Target in Renal Cell Carcinoma

Britney Reese,*¹ Ashok Silwal,*¹ Elizabeth Daugherty,* Michael Daugherty,[†] Mahshid Arabi,* Pierce Daly,* Yvonne Paterson,[‡] Layton Woolford,^{§,¶} Alana Christie,^{§,¶} Roy Elias,^{§,¶} James Brugarolas,^{§,¶} Tao Wang,^{¶,||} Magdalena Karbowniczek,* and Maciej M. Markiewski*

Preclinical studies demonstrated that complement promotes tumor growth. Therefore, we sought to determine the best target for complement-based therapy among common human malignancies. High expression of 11 complement genes was linked to unfavorable prognosis in renal cell carcinoma. Complement protein expression or deposition was observed mainly in tumor vasculature, corresponding to a role of complement in regulating the tumor microenvironment. In tumors correlated with a high nuclear grade. Complement genes clustered within an aggressive inflammatory cancer characterized by poor prognosis, markers of T cell dysfunction, and alternatively activated macrophages. Complement proteins correlated with response to immune checkpoint inhibitors. Corroborating human data and blockade reduced tumor growth by enhancing antitumor immunity and seemingly reducing angiogenesis of kidney cancer resistant to PD-1 blockade. Overall, this study implicates complement in the immune landscape of renal cell carcinoma, and notwithstanding cohort size and preclinical model limitations, the data suggest that tumors resistant to immunotherapy might be suitable targets for complement-based therapy. *The Journal of Immunology*,



Part 1

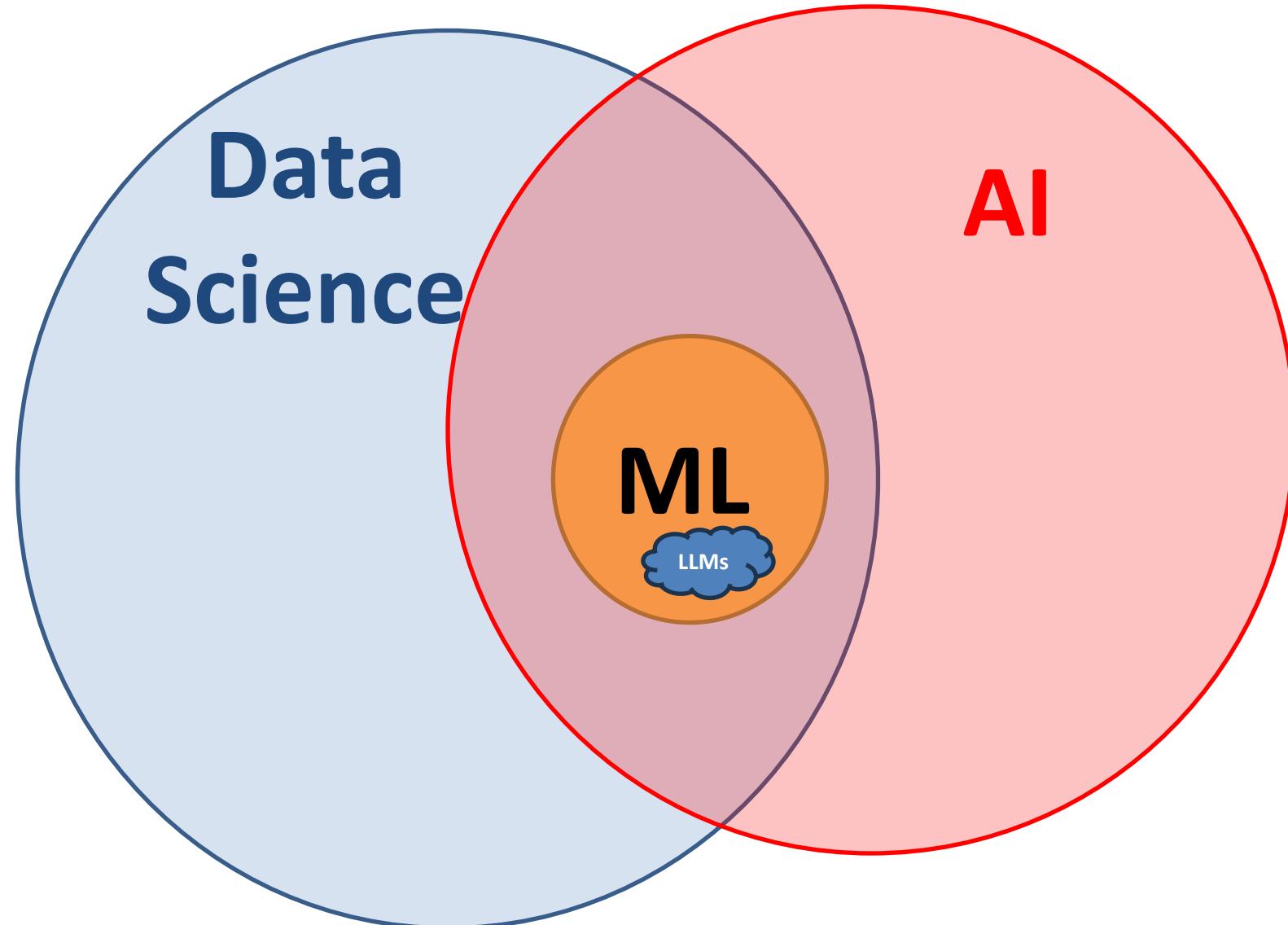
WHAT IS MACHINE LEARNING: A TOO SHORT INTRO

Data Science



- 1) Having data is like sitting on top of a gold mine
- 2) Everyone has data...

Data Science is learning how to dig through the mountain and get the gold out



* not to scale

We will see some of the math behind
LLMs and generative AI models

This class is about how they **work**,
not how to use them

What is “learning”?

The goal:

Make a decision based on data

i.e. find “patterns” in the data. (Besides, if you aren’t making your decisions based on data, then what are you basing them on?)

To get a good answer you **MUST**:

- ask a good question
- have good data

Just the Basics

CONCEPTS

High-Level Overview

Input: represent one sample as numbers (called **features**) in an efficient way that is hopefully useful for solving the problem

Output: our answer depends on our question, common examples:

- Classification: what discrete category does the sample belong to?
- Regression: predict a real number
- Reinforcement: optimize a certain outcome

Model: how we calculate the output from the input. Usually has lots of adjustable parameters that need to be set somehow

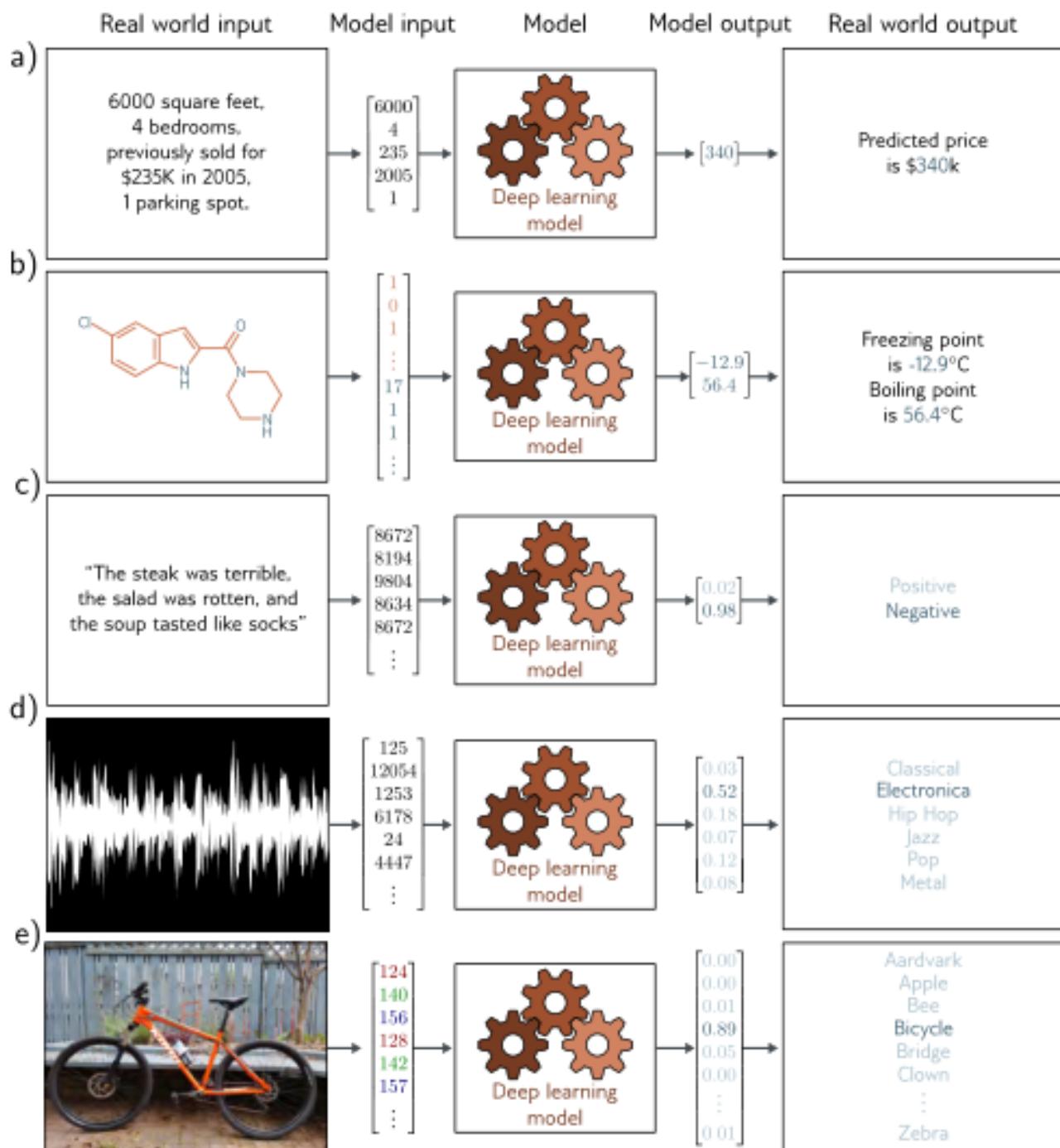
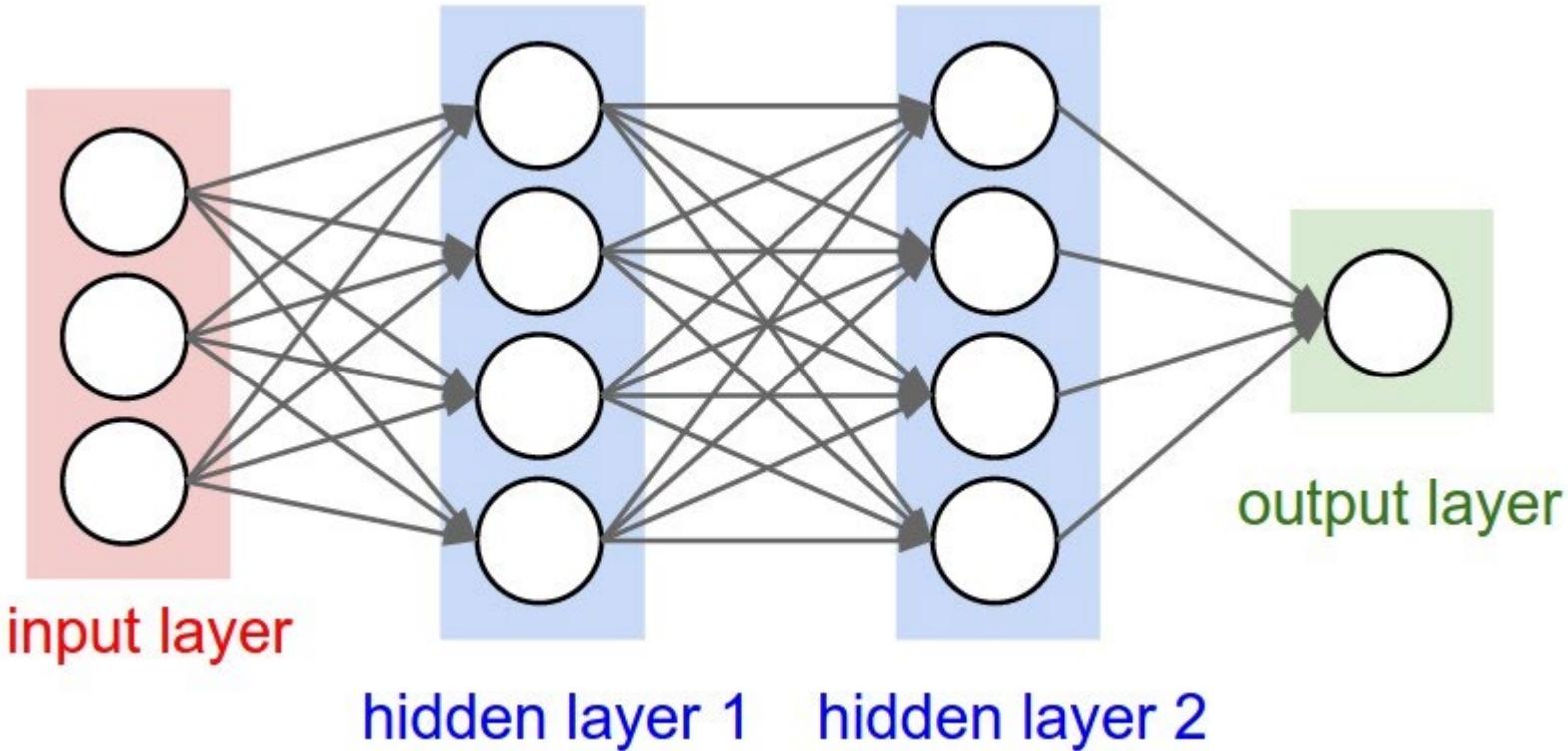


Figure 1.2 Regression and classification problems. a) This *regression* model takes a vector of numbers that characterize a property and predicts its price. b) This *multivariate regression* model takes the structure of a chemical molecule and predicts its melting and boiling points. c) This *binary classification* model takes a restaurant review and classifies it as either positive or negative. d) This *multiclass classification* problem assigns a snippet of audio to one of N genres. e) A second multiclass classification problem in which the model classifies an image according to which of N possible objects it might contain.

Models

- Neural networks are the ones everyone has heard of, but there are lots and lots and lots of other choices
- Neural networks are extremely well-suited to some problems like 2D images, can be complicated and finicky to work with so not always the best choice
- You don't know for certain which model will work best on a given problem ahead of time, so we have to try several
- Most of this class is devoted to learning about different models

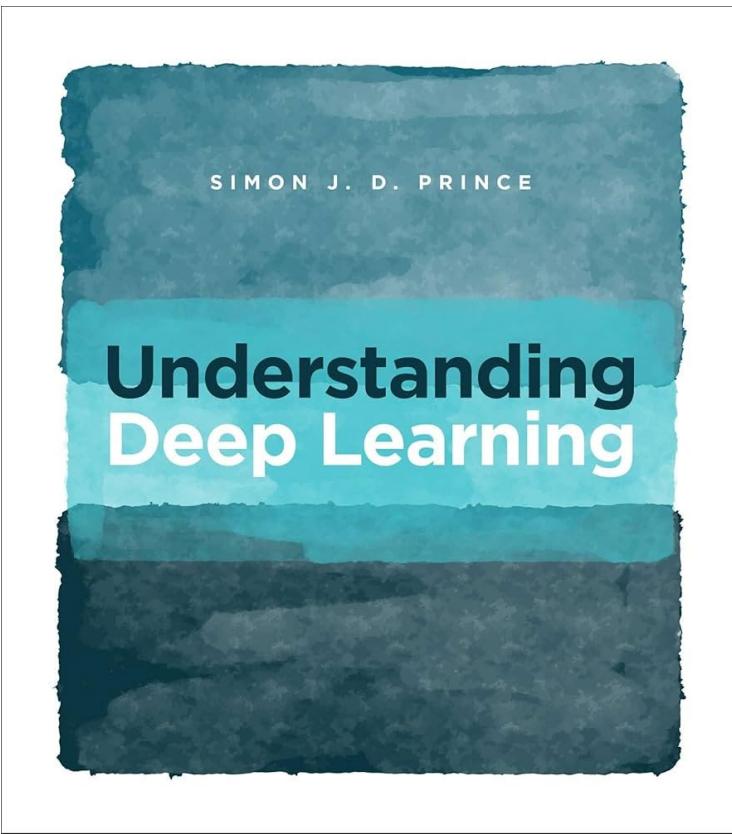
Neural Network



- Each **input** is a SINGLE NUMBER
- **Hidden** and **output** values are directly calculated from inputs and **weights**
- **Weights (arrows)** are adjusted in training to make output match known right answer

Deep Learning

- Deep Learning = Neural networks with lots of layers
- They work extremely well for many hard problems and we honestly aren't really sure why



The title is also partly a joke — *no-one* really understands deep learning at the time of writing. Modern deep networks learn piecewise linear functions with more regions than there are atoms in the universe and can be trained with fewer data examples than model parameters. It is neither obvious that we should be able to fit these functions reliably nor that they should generalize well to new data. The penultimate chapter addresses

Seriously Though

https://scikit-learn.org/stable/user_guide.html

1. Supervised learning

- 1.1. Linear Models
- 1.2. Linear and Quadratic Discriminant Analysis
- 1.3. Kernel ridge regression
- 1.4. Support Vector Machines
- 1.5. Stochastic Gradient Descent
- 1.6. Nearest Neighbors
- 1.7. Gaussian Processes
- 1.8. Cross decomposition
- 1.9. Naive Bayes
- 1.10. Decision Trees
- 1.11. Ensemble methods
- 1.12. Multiclass and multioutput algorithms
- 1.13. Feature selection
- 1.14. Semi-supervised learning
- 1.15. Isotonic regression
- 1.16. Probability calibration
- 1.17. Neural network models (supervised)

2. Unsupervised learning

- 2.1. Gaussian mixture models
- 2.2. Manifold learning
- 2.3. Clustering
- 2.4. Biclustering
- 2.5. Decomposing signals in components (matrix factorization problems)
- 2.6. Covariance estimation
- 2.7. Novelty and Outlier Detection
- 2.8. Density Estimation
- 2.9. Neural network models (unsupervised)

Part 2

TOUR OF ML TOPICS

Classification

- **Classification** – decide which **category** a sample belongs to
- Examples:
 - Email or spam
 - Speech recognition
 - OCR (Optical Character Recognition)
- *i.e. write a program that can recognize patterns*

Spam Filter

Assassinating spam e-mail

SpamAssassin is a widely used open-source spam filter. It calculates a score for an incoming e-mail, based on a number of built-in rules or ‘tests’ in SpamAssassin’s terminology, and adds a ‘junk’ flag and a summary report to the e-mail’s headers if the score is 5 or more.

-0.1 RCVD_IN_MXRATE_WL	RBL: MXRate recommends allowing [123.45.6.789 listed in sub.mxrate.net]
0.6 HTML_IMAGE_RATIO_02	BODY: HTML has a low ratio of text to image area
1.2 TVD_FW_GRAPHIC_NAME_MID	BODY: TVD_FW_GRAPHIC_NAME_MID
0.0 HTML_MESSAGE	BODY: HTML included in message
0.6 HTML_FONx_FACE_BAD	BODY: HTML font face is not a word
1.4 SARE_GIF_ATTACH	FULL: Email has a inline gif
0.1 BOUNCE_MESSAGE	MTA bounce message
0.1 ANY_BOUNCE_MESSAGE	Message is some kind of bounce message
1.4 AWL	AWL: From: address is in the auto white-list

From left to right you see the score attached to a particular test, the test identifier, and a short description including a reference to the relevant part of the e-mail. As you see, scores for individual tests can be negative (indicating evidence suggesting the e-mail is ham rather than spam) as well as positive. The overall score of 5.3 suggests the e-mail might be spam.

Spam Filter

SEND

Save Now

Discard

To

jenny [REDACTED]@g

Add Cc Add Bo

Subject

<3 <3 <3

Attach a file

B I U T + T + A

Dearest Jenny,

I can't stop thinking about you. I had that dream about you again last night. Y'know that dream where I am Captain Picard and you are Dr. Crusher. I miss you so much. Why didn't you return my calls last night? I tried calling you several times last night but I kept getting your voicemail. I know you were home last night because of your facebook and twitter posts. Are you upset with me or something? I hope I didn't do anything to mess up our relationship because I think we have something really special like we were destined to be



« Plain Text

Training

Supervised Learning:

- we are given training data where we know the right answers
- pattern recognition finds hidden trends
- training data can be a huge bottleneck
- examples: ChatGPT, facial recognition, OCR

Unsupervised Learning:

- we don't know the right answer, so optimize for some metric
- can I find natural clusters or groups in the data?
- stable diffusion: adjust these pixels to make the picture look like a duck
- social media/YouTube “algorithm”

Supervised learning example:

Interviewer: What's your biggest strength?

Me: I'm an expert in machine learning.

Interviewer: What's $9 + 10$?

Me: Its 3.

Interviewer: Not even close. It's 19.

Me: It's 16.

Interviewer: Wrong. Its still 19.

Me: It's 18.

Interviewer: No, it's 19.

Me: it's 19.

Interviewer: You're hired

[Post](#)

Felix

@felix_red_panda

[Follow](#)

...

Microsoft paper claims ChatGPT 3.5 has ~20 billion parameters
arxiv.org/abs/2310.17680

System description		Python (0)		
System	Model	#P	top-1	tokens
T5	t5-large	770M	80.4	
CodeT5	codet5-large	770M	80.5	
GPT-3	text-davinci-003	175B	82.5	
ChatGPT	gpt-3.5-turbo	20B	80.6	
StarCoder	starcoder	15.5B	79.2	
CodeT5+	codet5p-16b	16B	79.6	
CodeGen	codegen-350m	350M	80.1	
Diffusion-LM	Custom	50M	70.4	
GENIE	Custom	93M	73.2	
CODEFUSION	Custom	75M	80.7	

8:03 PM · Oct 30, 2023 · 1M Views

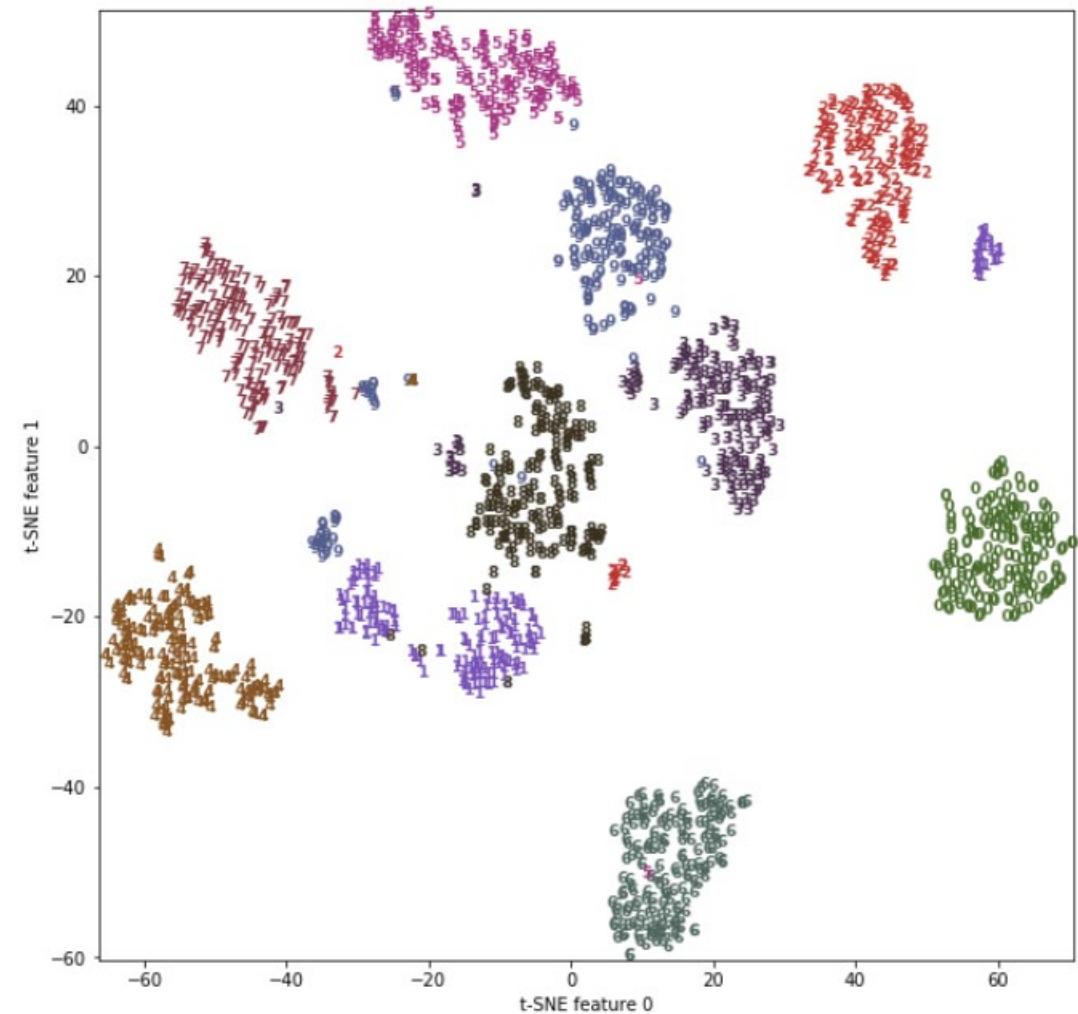
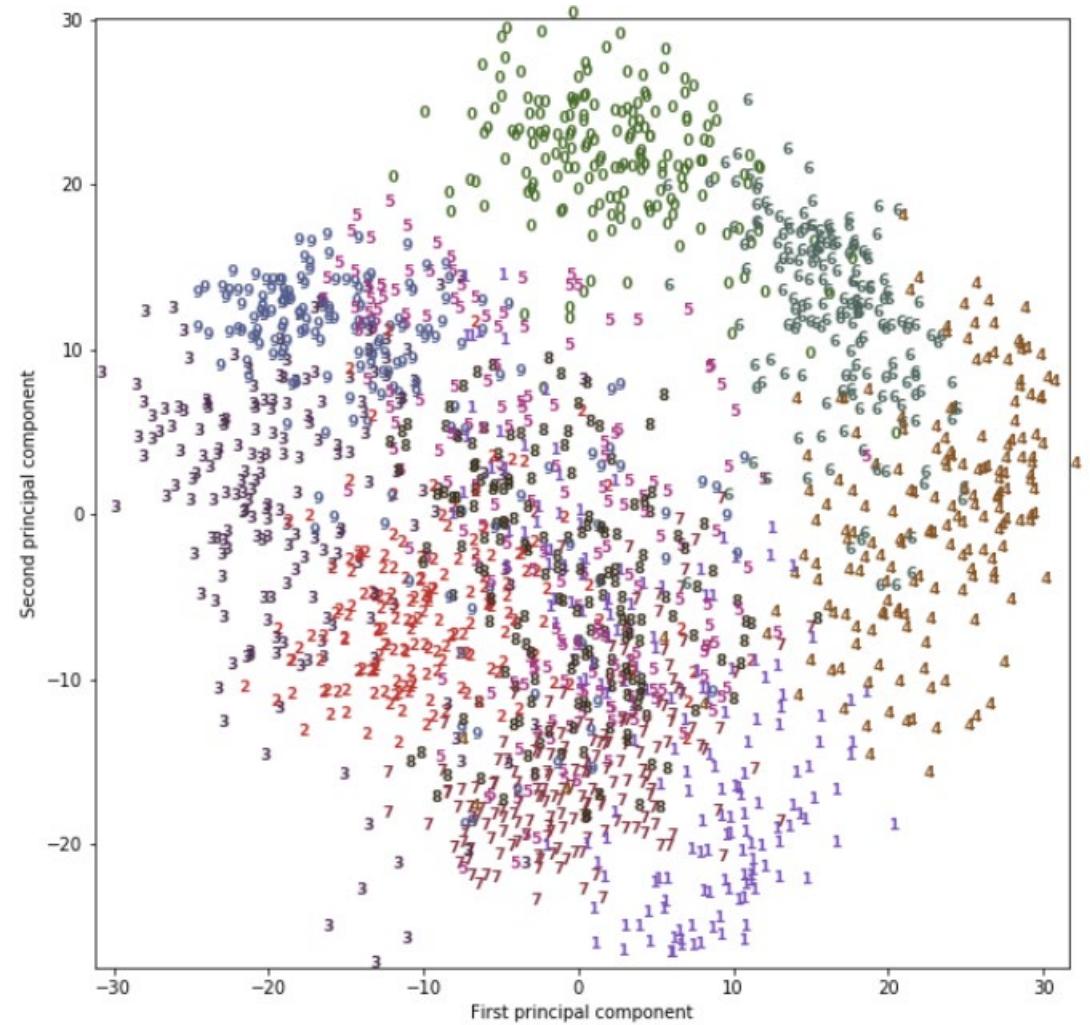
68

417

2.3K

676

Training data is the “fossil fuel” of LLMs. We have essentially scraped the entire internet for data of human-written text and we are reaching the upper limit of parameters we can train.



Other Tricks

REGRESSION

Regression

Instead of classifying into categories, our job is now to predict a real-valued number based on training data

Examples:

- predict the price of a house in Boston
- predict how much money a movie will make
- guess the value of a stock
- Diabetes dataset: predict measure of disease based on patient data

Other Tricks

REINFORCEMENT

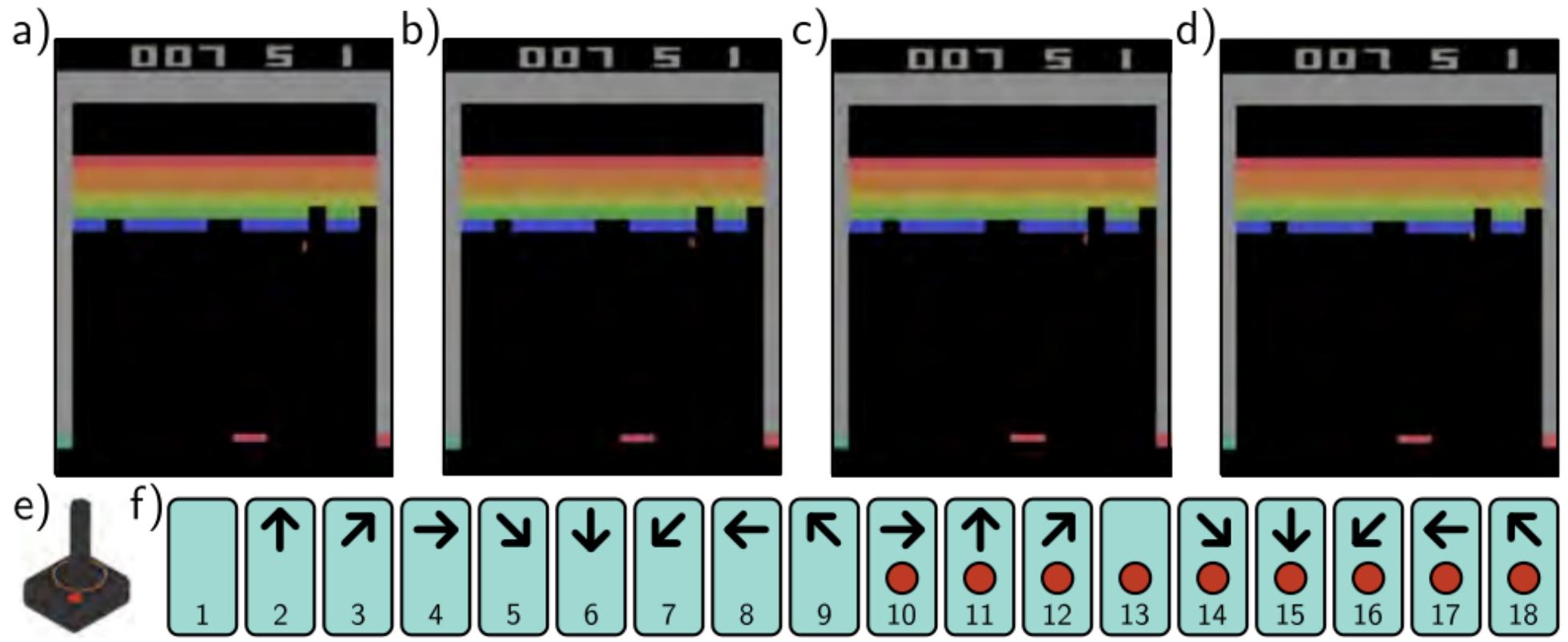
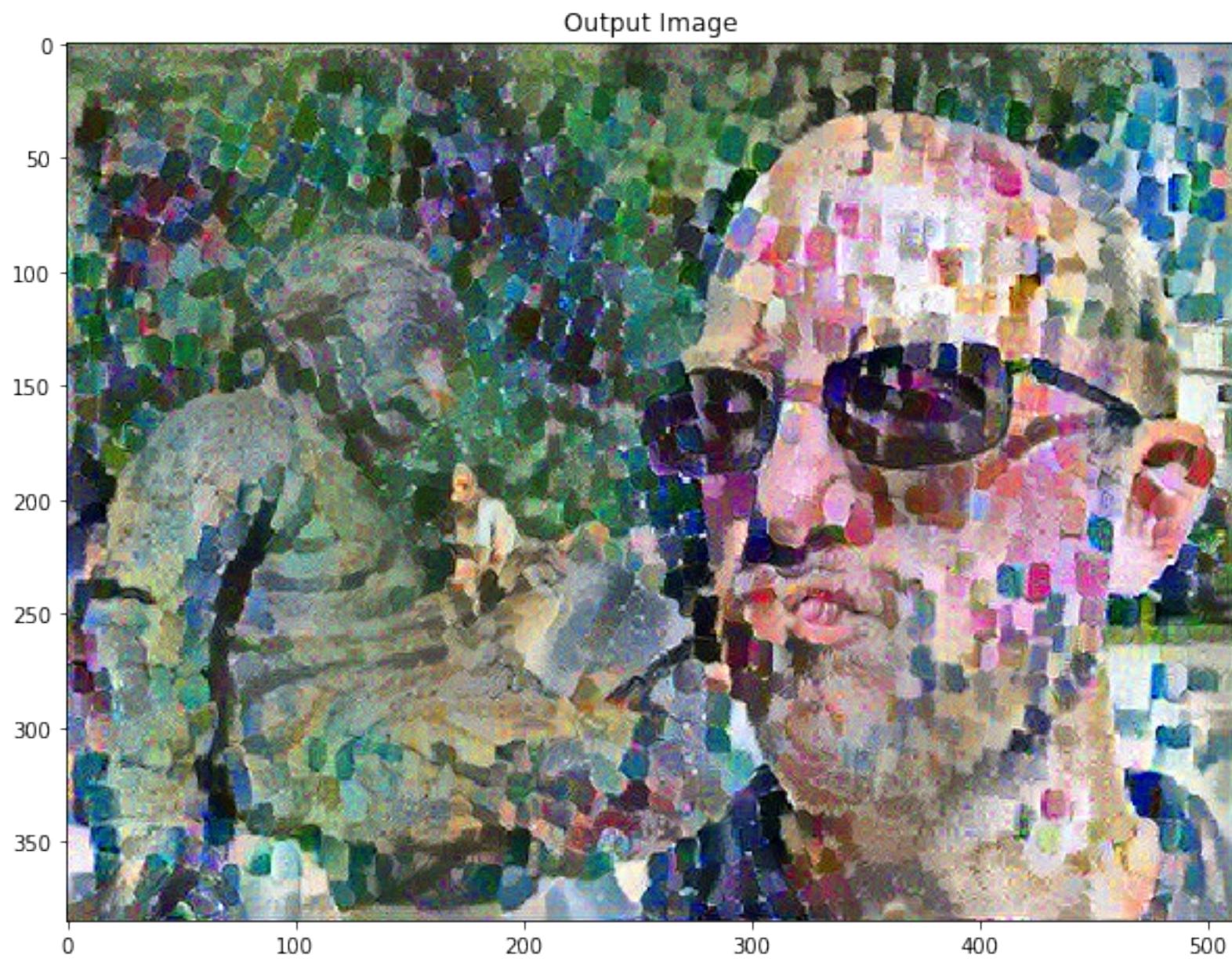


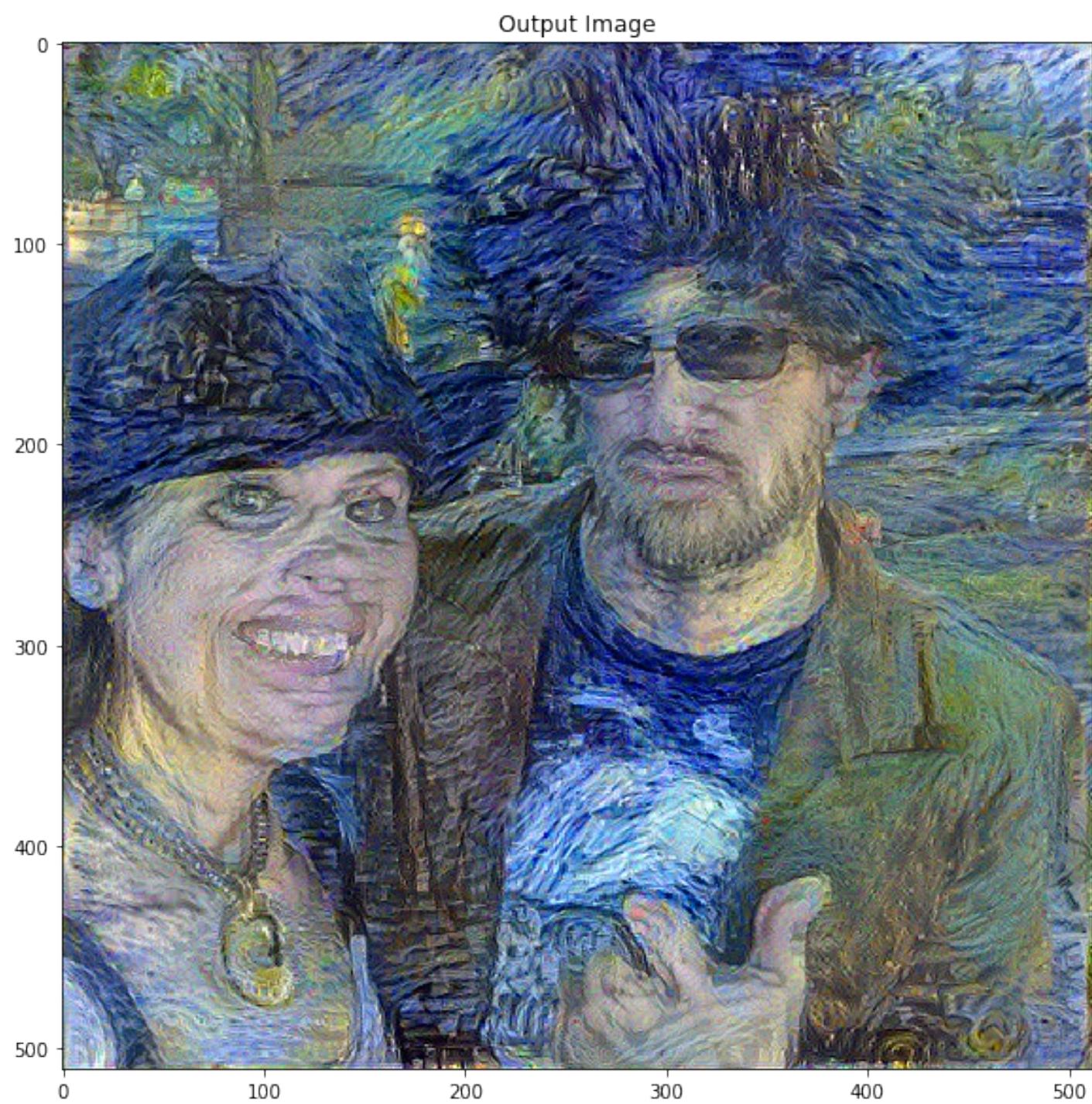
Figure 19.13 Atari Benchmark. The Atari benchmark consists of 49 Atari 2600 games, including Breakout (pictured), Pong, and various shoot-em-up, platform, and other types of games. a-d) Even for games with a single screen, the state is not fully observable from a single frame because the velocity of the objects is unknown. Consequently, it is usual to use several adjacent frames (here, four) to represent the state. e) The action simulates the user input via a joystick. f) There are eighteen actions corresponding to eight directions of movement or no movement, and for each of these nine cases, the button being pressed or not.

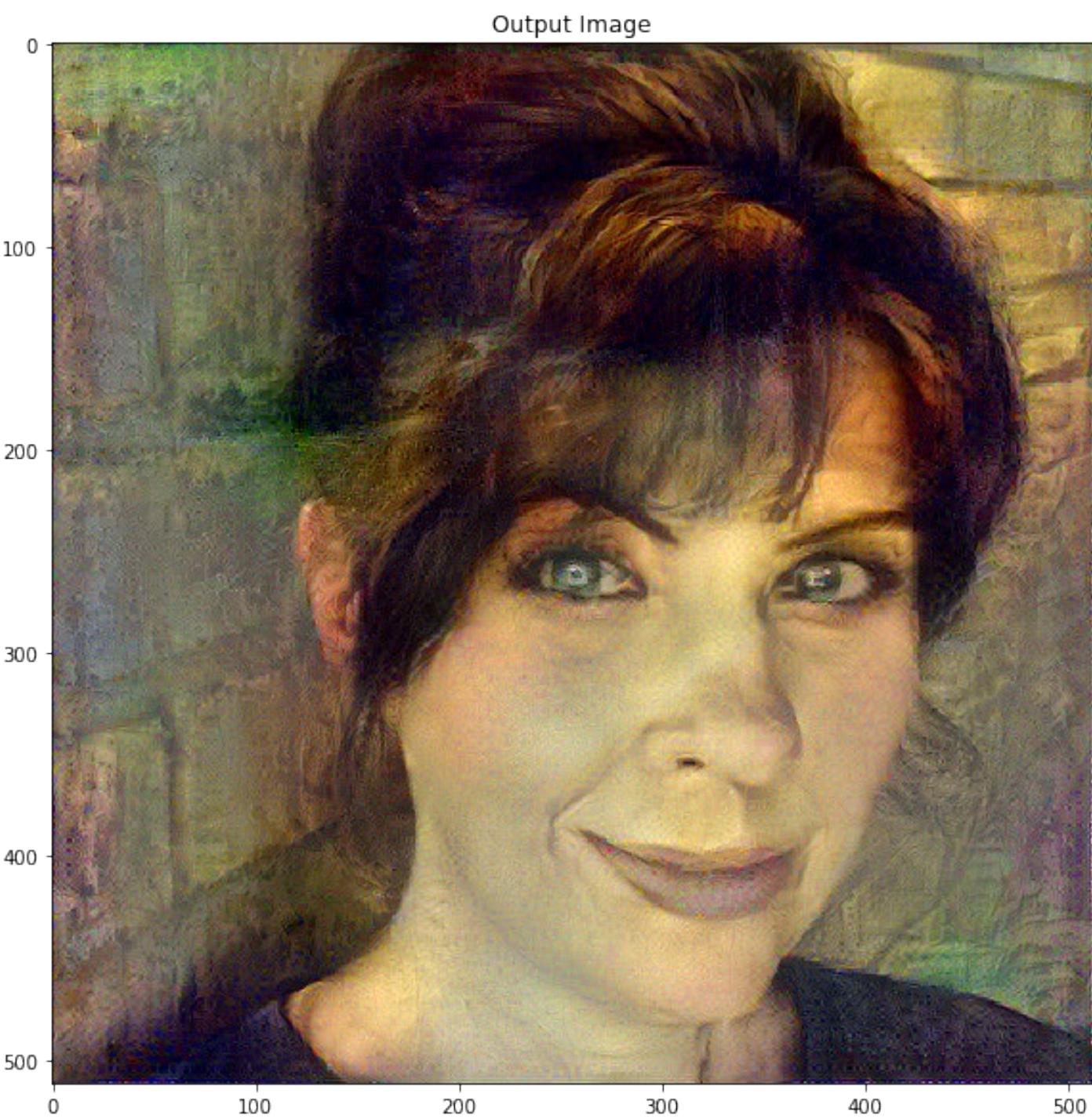
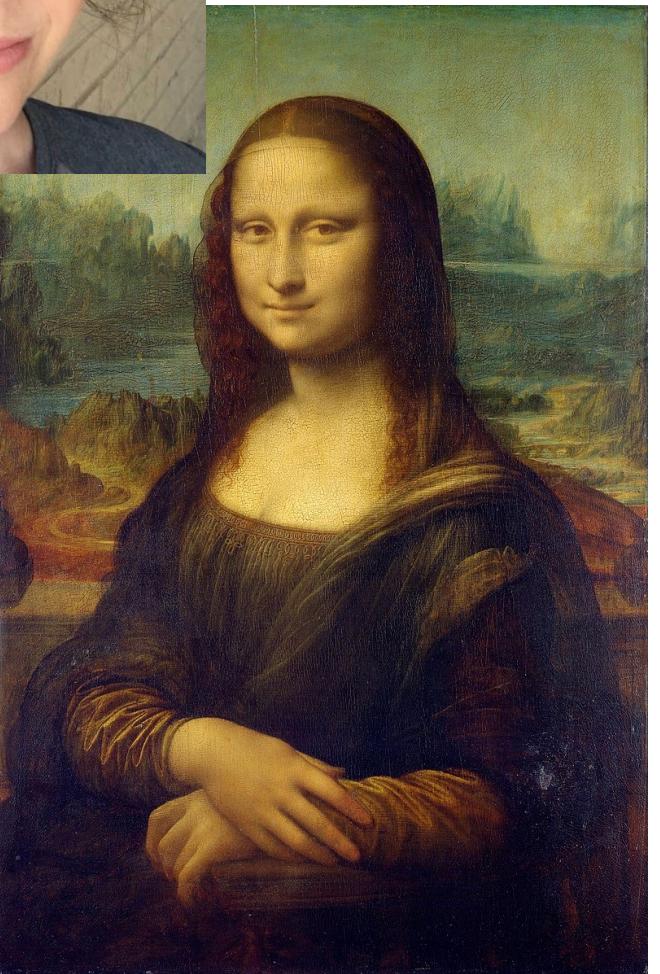
Other Tricks

GENERATION





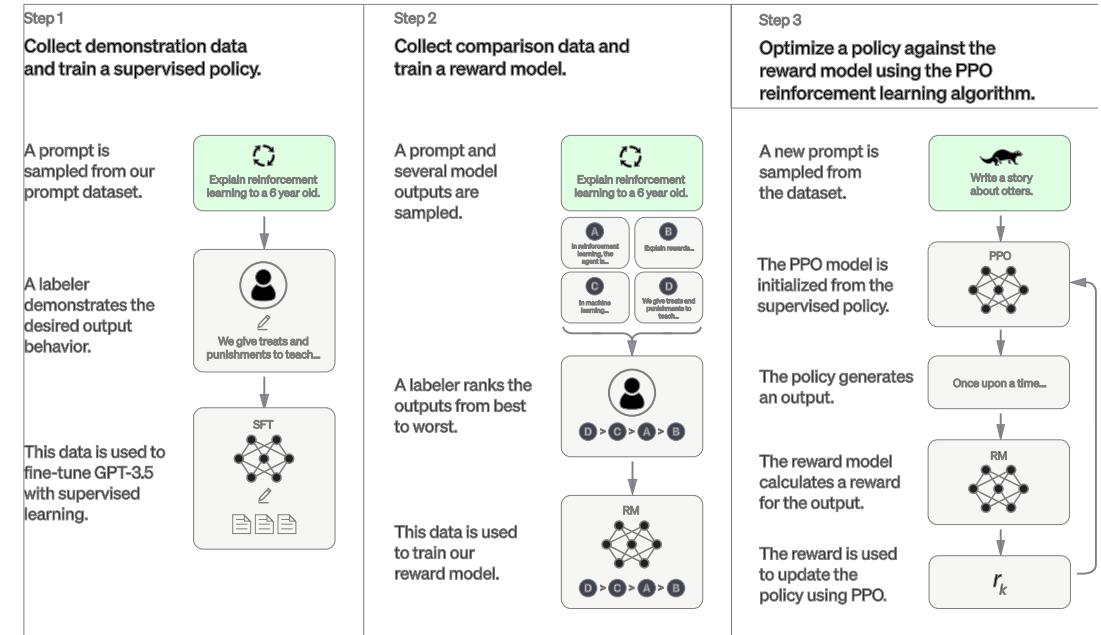






More Examples

- DALL-E: generating images based on input text description
<https://labs.openai.com/>
- ChatGPT: generating text based on input text description
<https://openai.com/blog/chatgpt/>

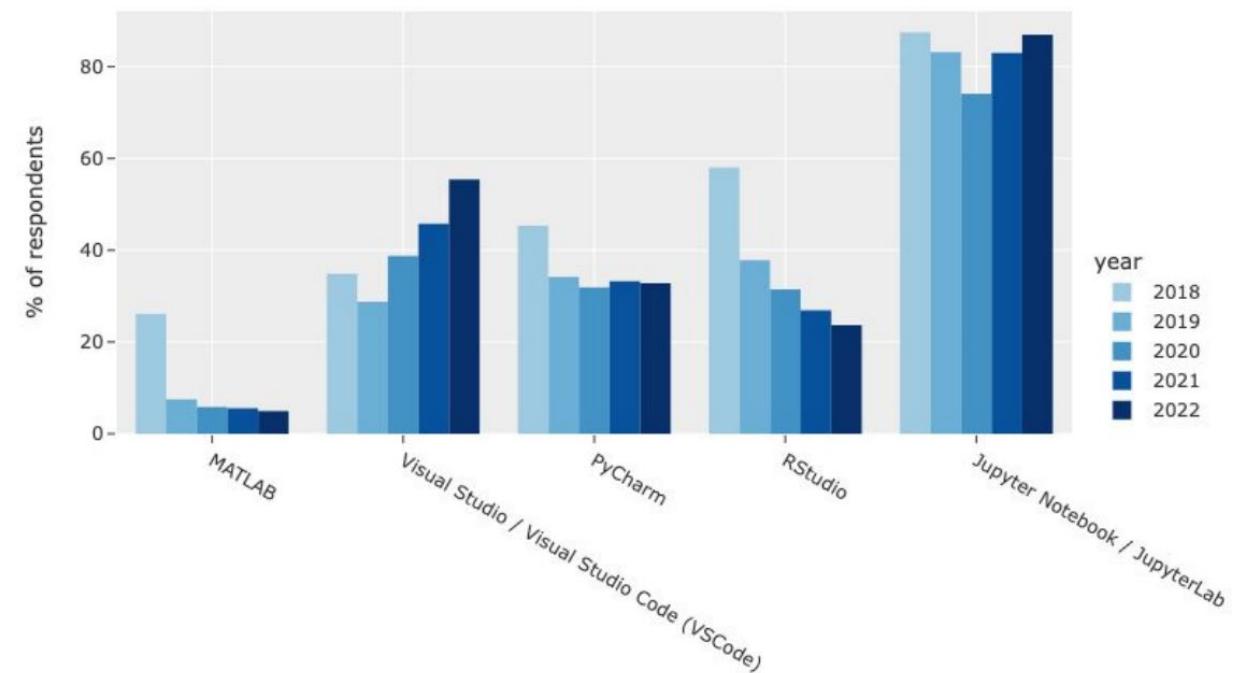
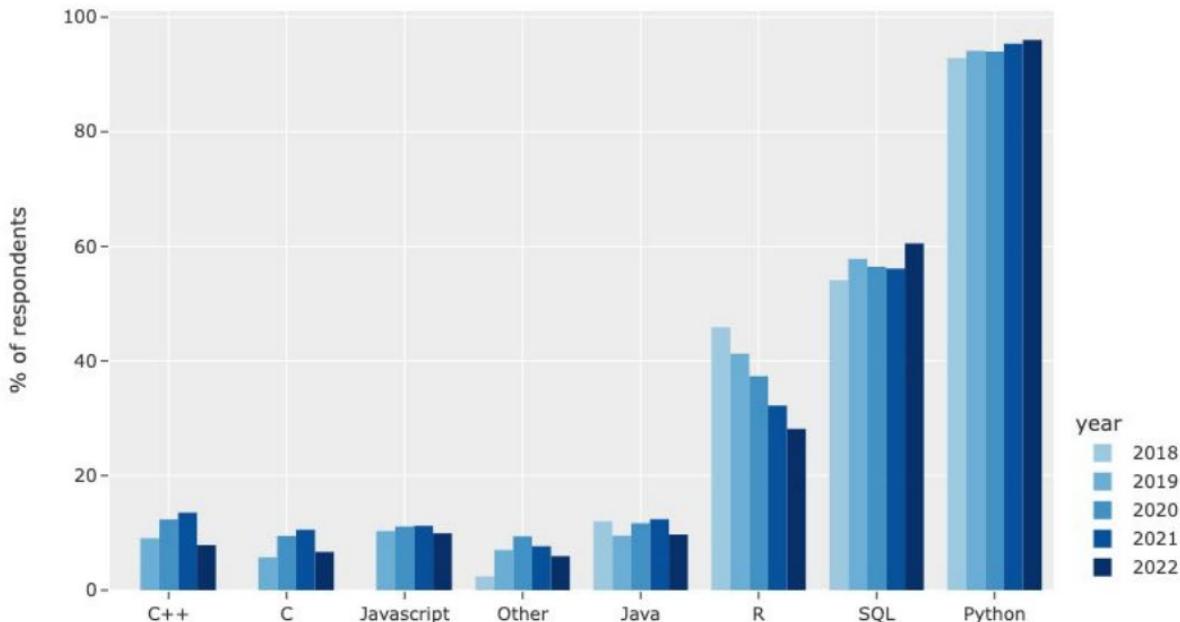


Part 3

TOOLS

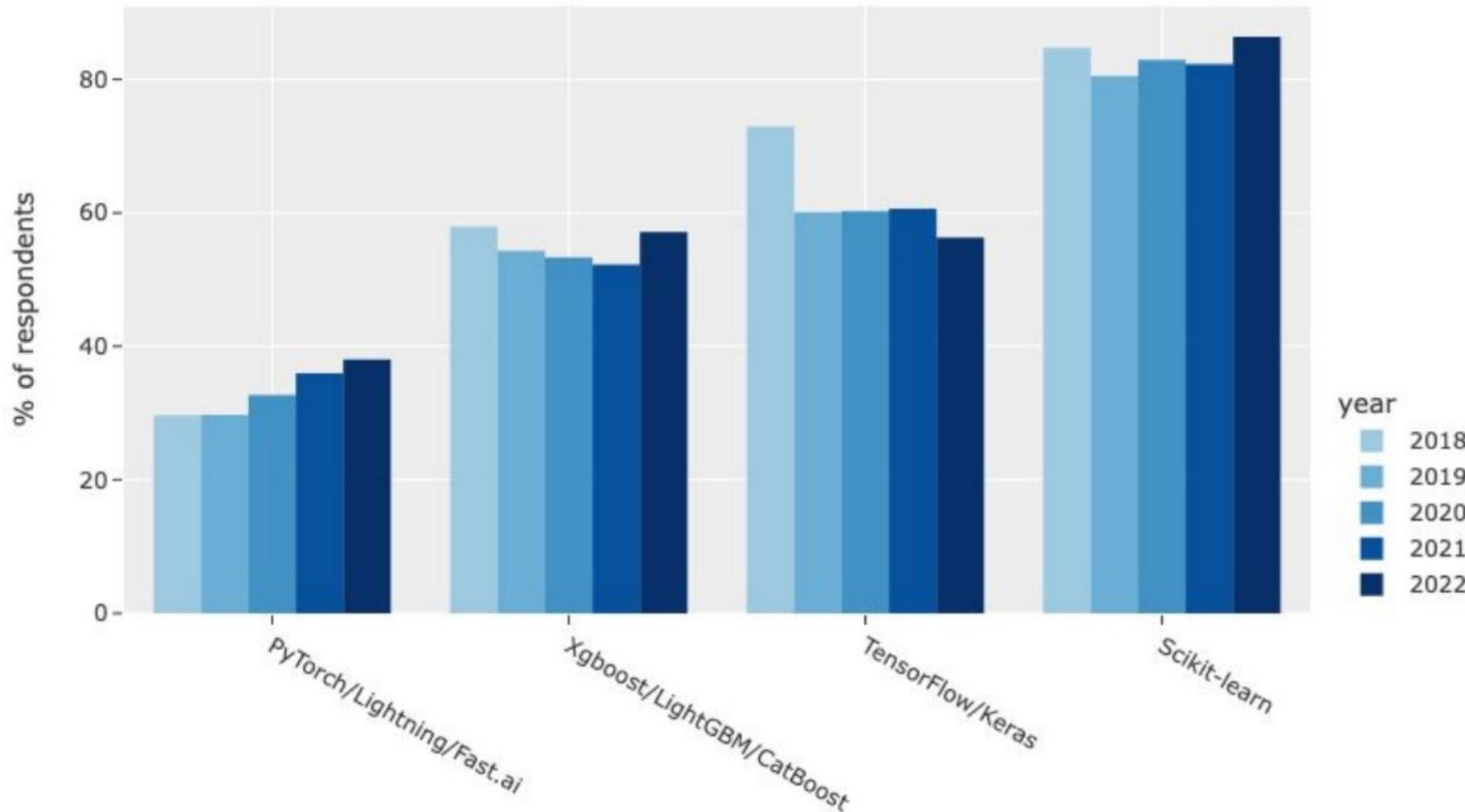
Kaggle: Making Data Science a Sport

<https://www.kaggle.com/kaggle-survey-2022>



Kaggle: Making Data Science a Sport

<https://www.kaggle.com/kaggle-survey-2022>



Getting Started

Great python resources:

A Whirlwind Tour of Python by Jake VanderPlas.

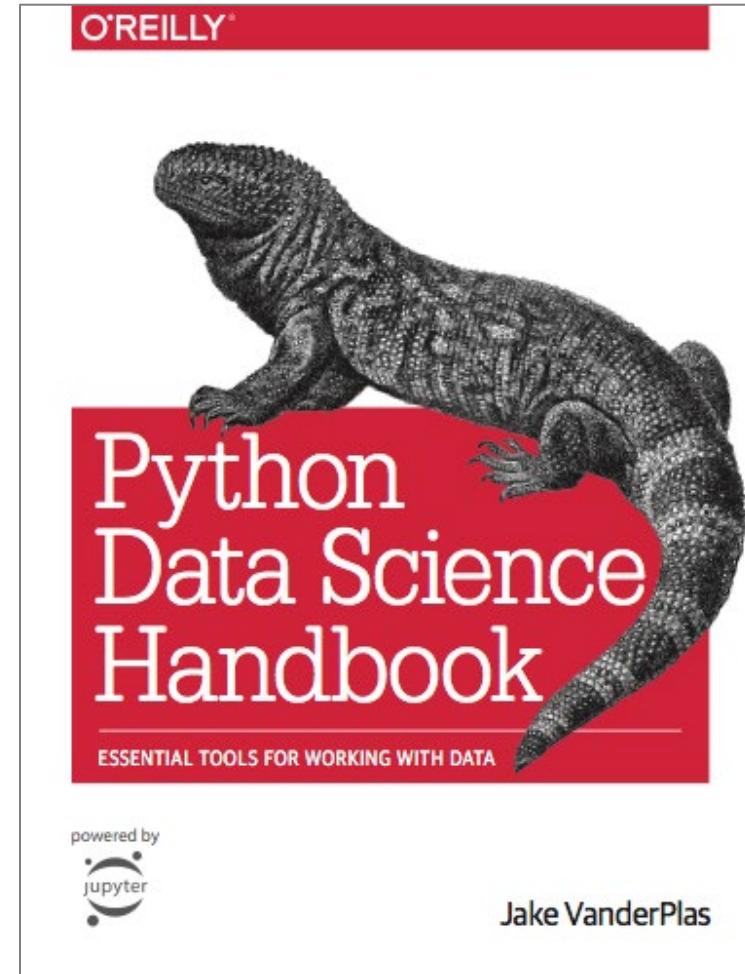
Freely available online at:

<https://jakevdp.github.io/WhirlwindTourOfPython/>

Python Data Science Handbook by Jake VanderPlas.

Freely available online at:

<https://jakevdp.github.io/PythonDataScienceHandbook/>



Part N-1

POLICIES

Doctors: Googling stuff online does not make you a doctor.

Programmers:



Not sharing your source code out of greed



Not sharing your source code out of shame