

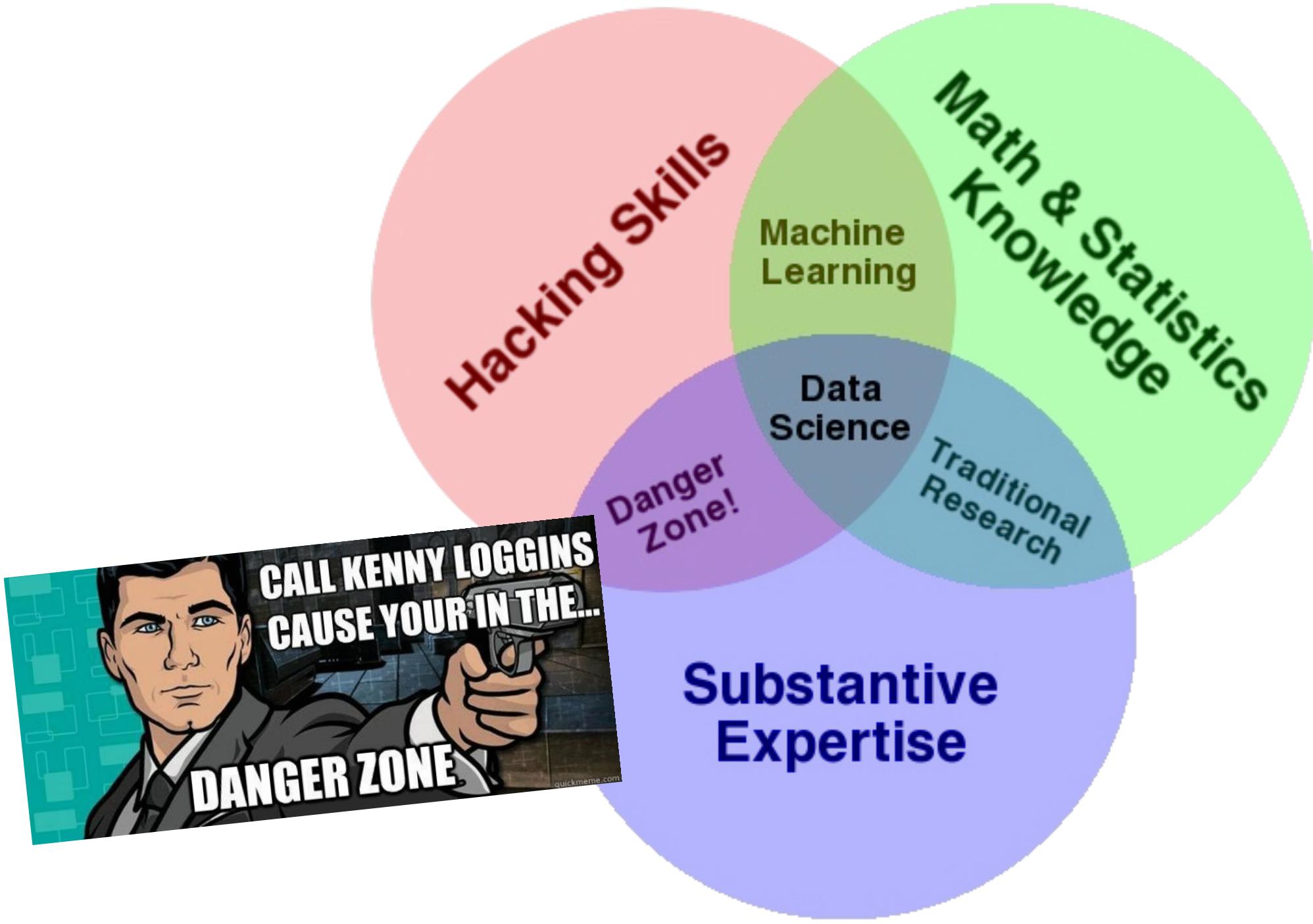
Machine Learning

PHYS 453 – Spring 2023

Dr. Daugherty



ABILENE
CHRISTIAN
UNIVERSITY

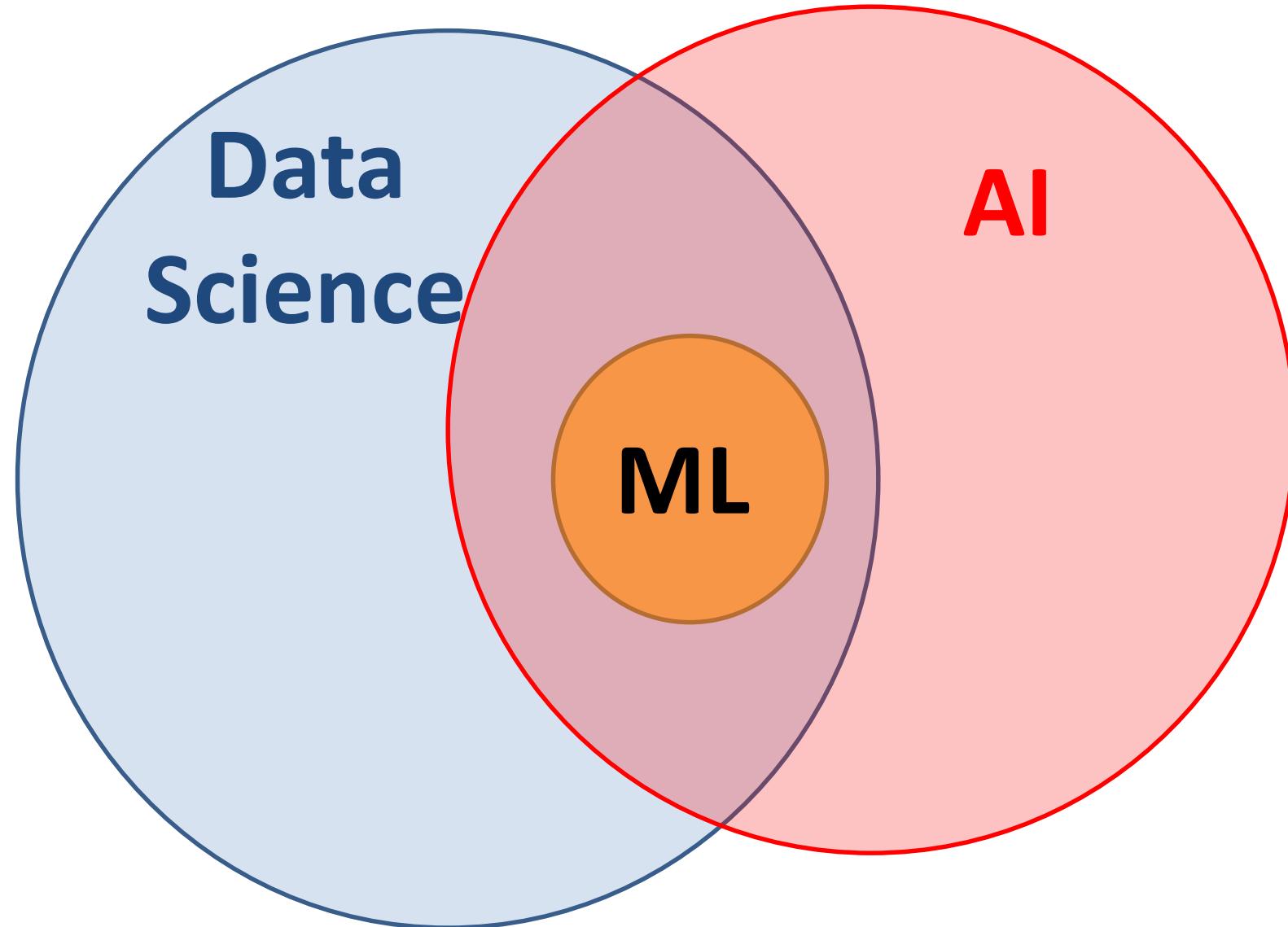


“Data science, as it’s practiced, is a blend of Red-Bull-fueled hacking and espresso-inspired statistics.

But data science is not merely hacking, because when hackers finish debugging their Bash one-liners and Pig scripts, few care about non-Euclidean distance metrics.

And data science is not merely statistics, because when statisticians finish theorizing the perfect model, few could read a tab delimited file into R if their job depended on it.

Data science is the civil engineering of data. Its acolytes possess a practical knowledge of tools & materials, coupled with a theoretical understanding of what’s possible.”



* not to scale

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:

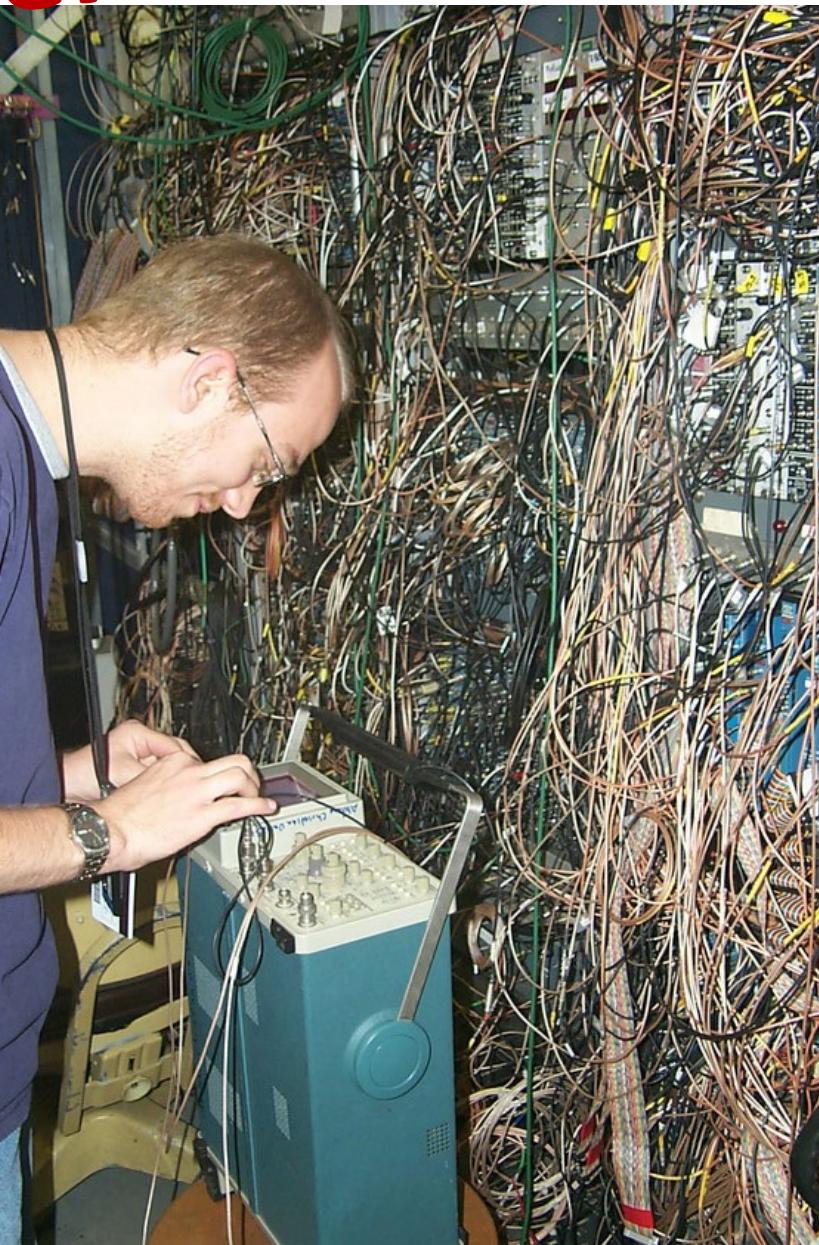
- Have been paid as a professional model
- Once burned off an eyebrow in class
- Plays drums in Abilene's best cover band

Hobbies:

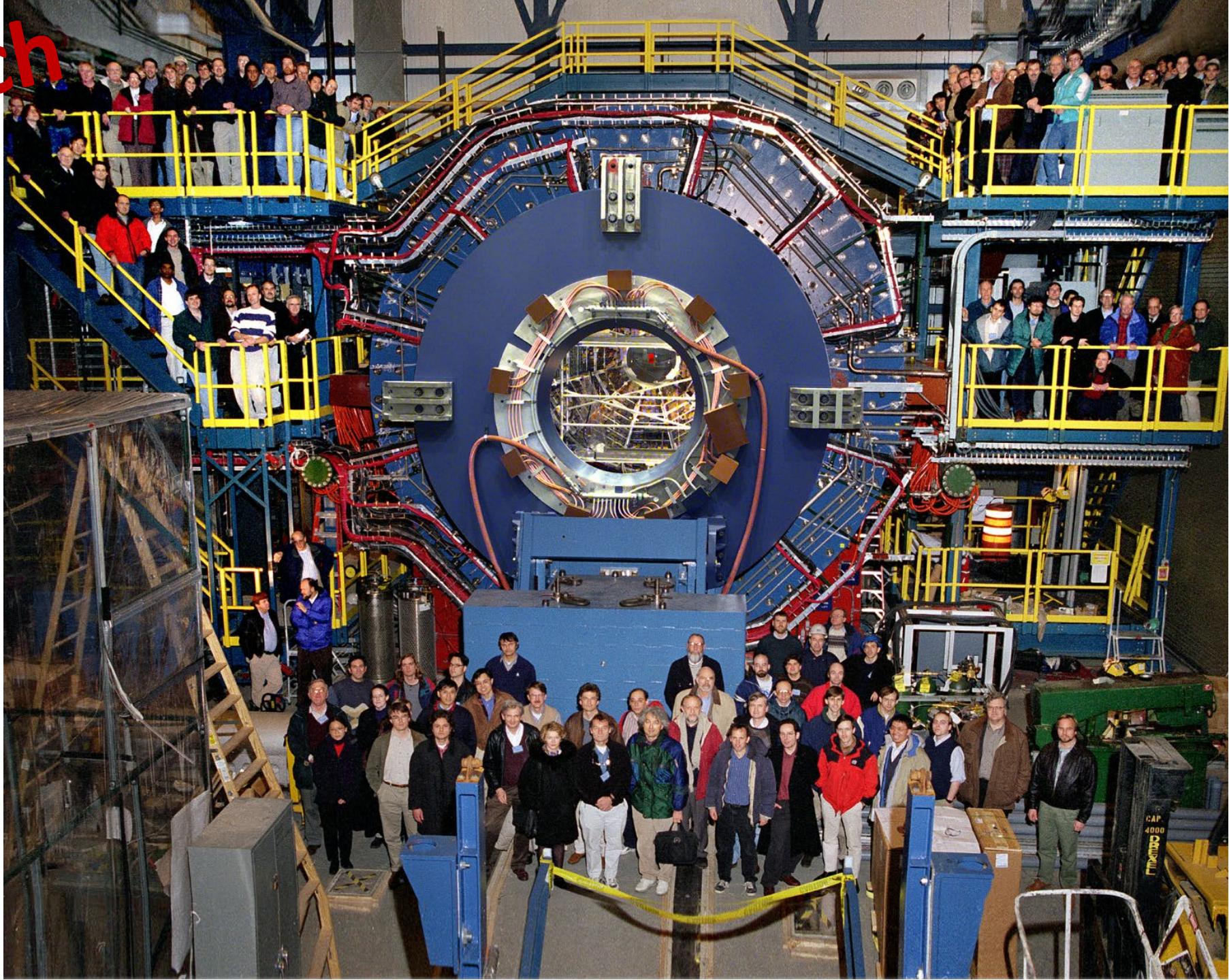
- Eating
- Putting things in a laser
- Building dangerous things



First Experiment

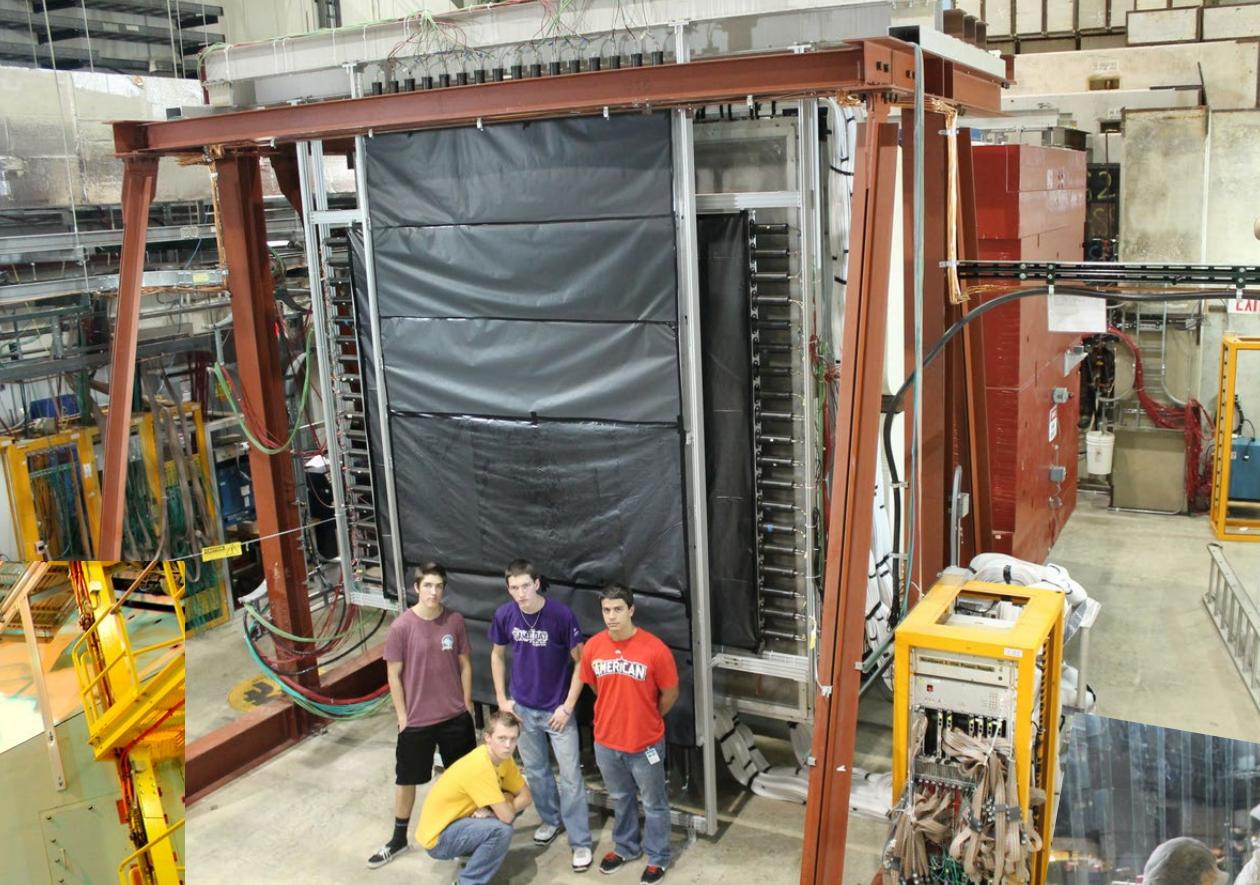
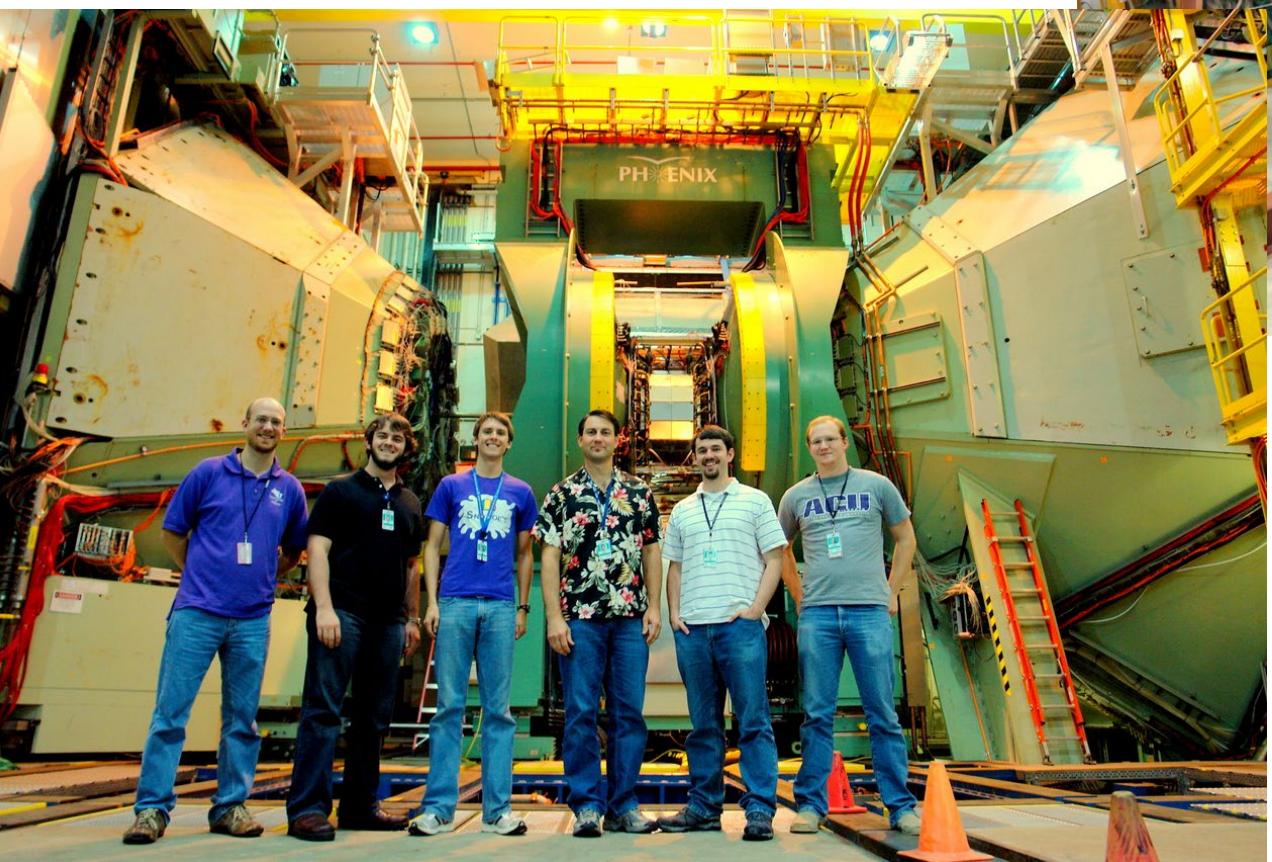


PhD Research



Mike Daugherty

ACU Research





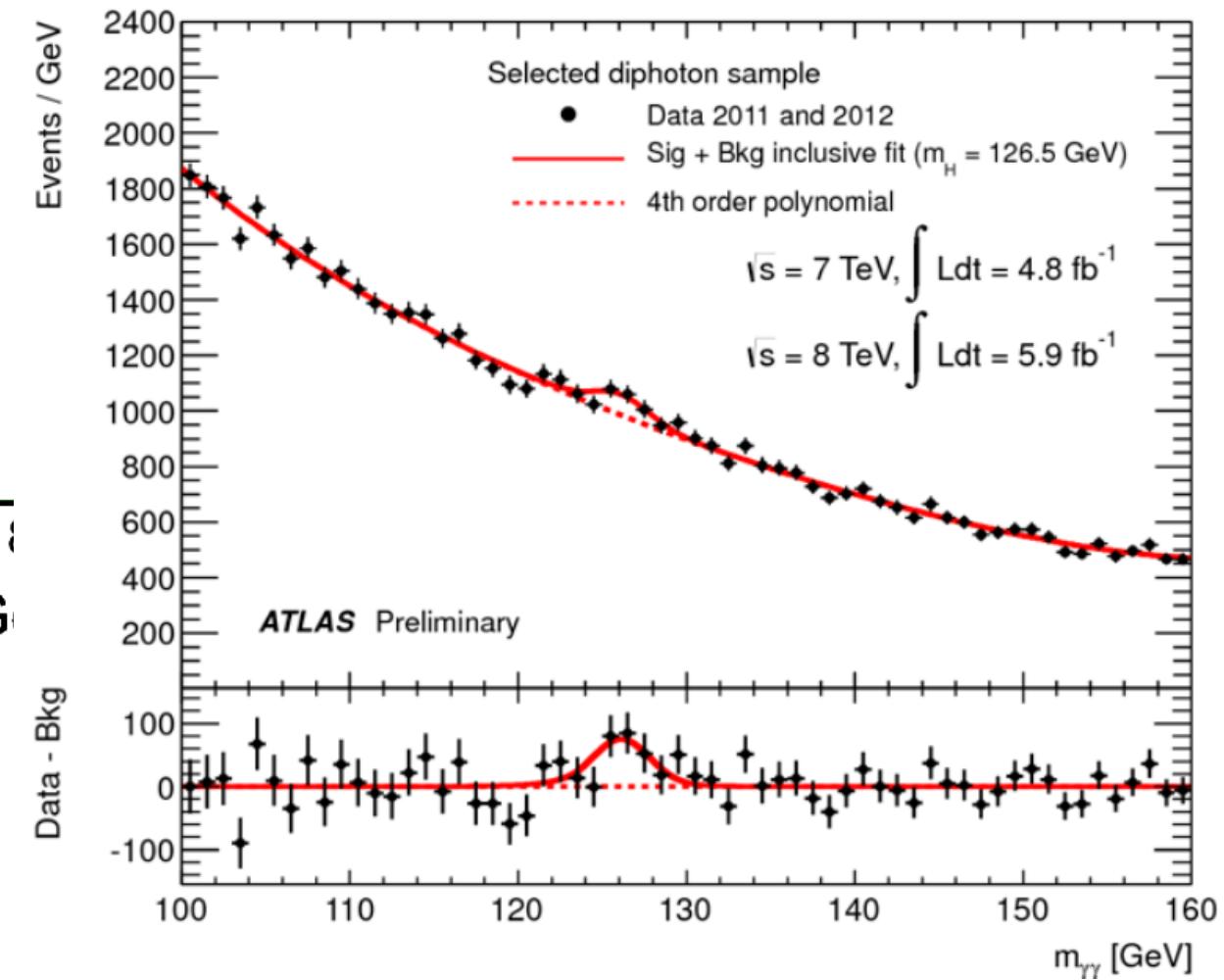
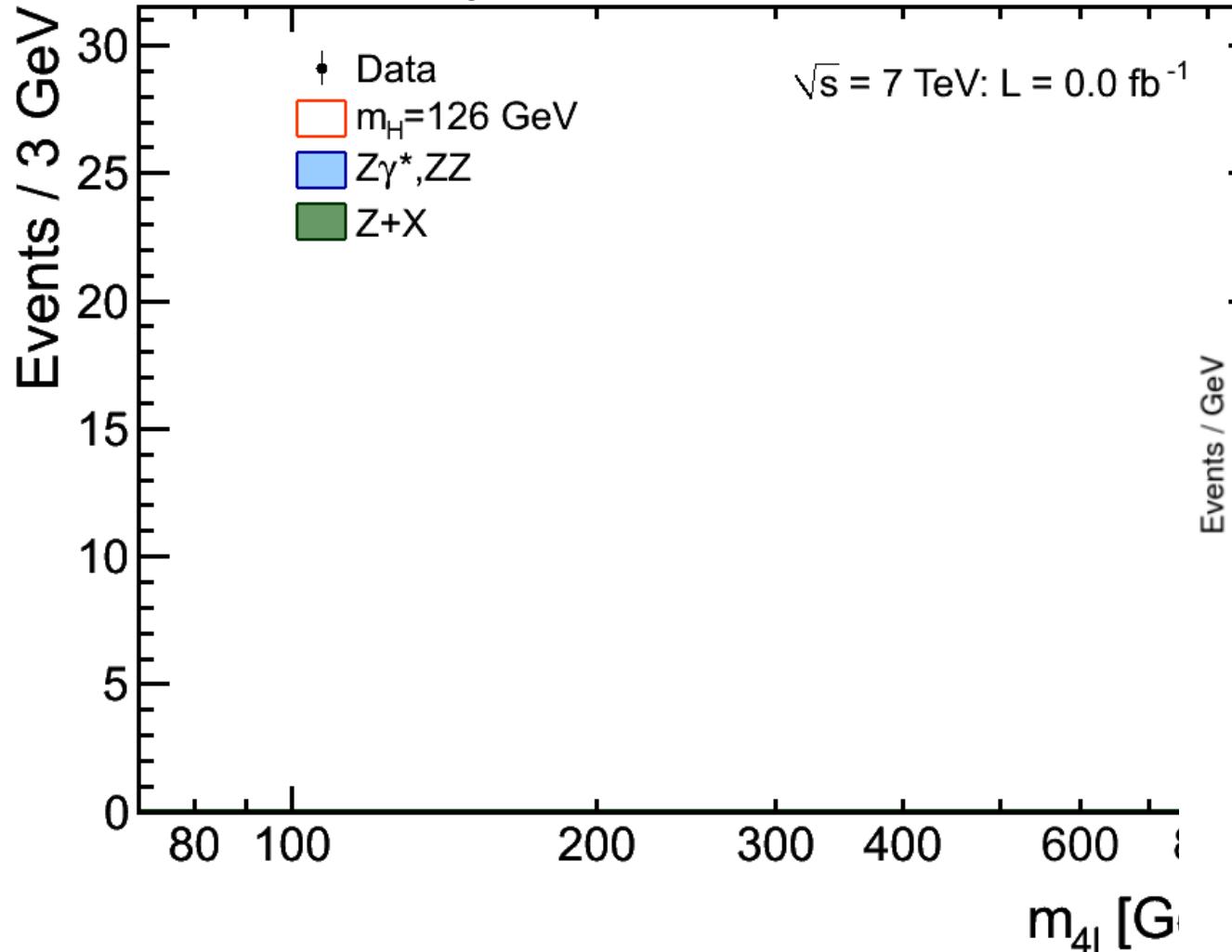
Before “Big Data” was a thing, there was high-energy physics...

[] 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



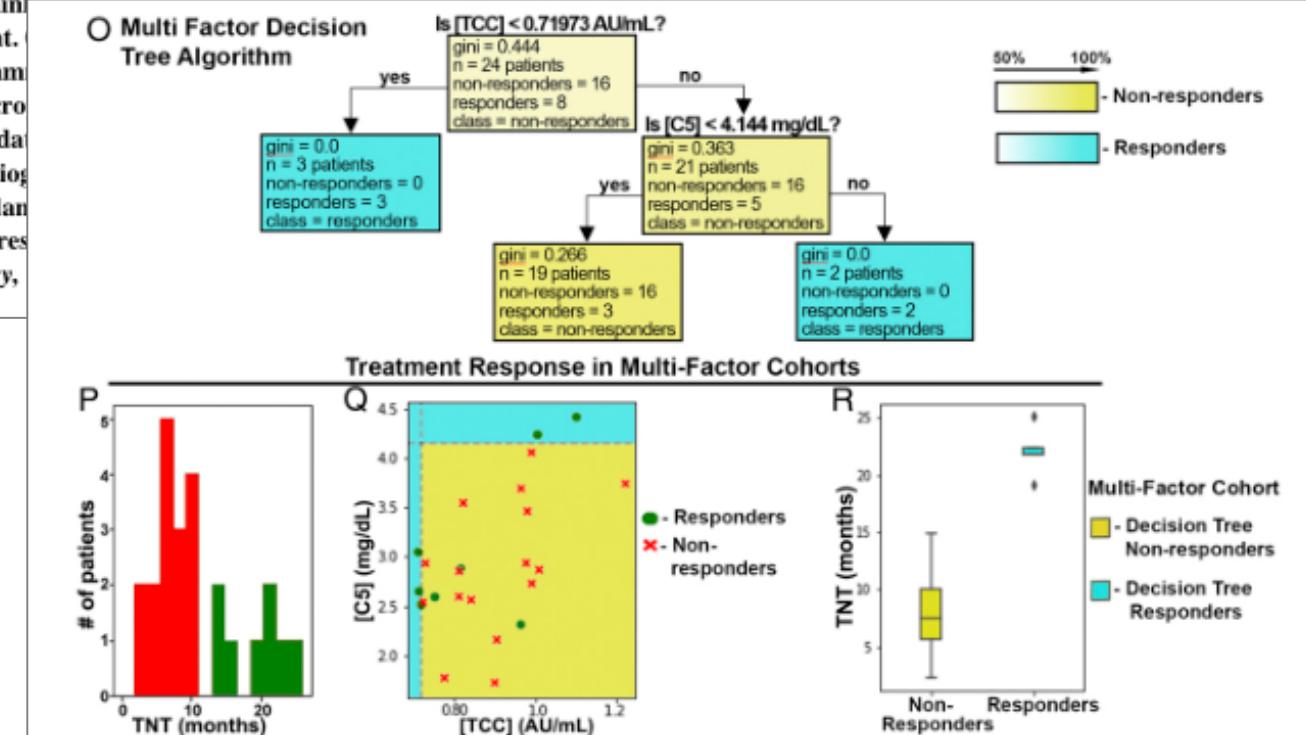
One More Example

<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*, 2020; 195(22): 7333–7344.



Data Science



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

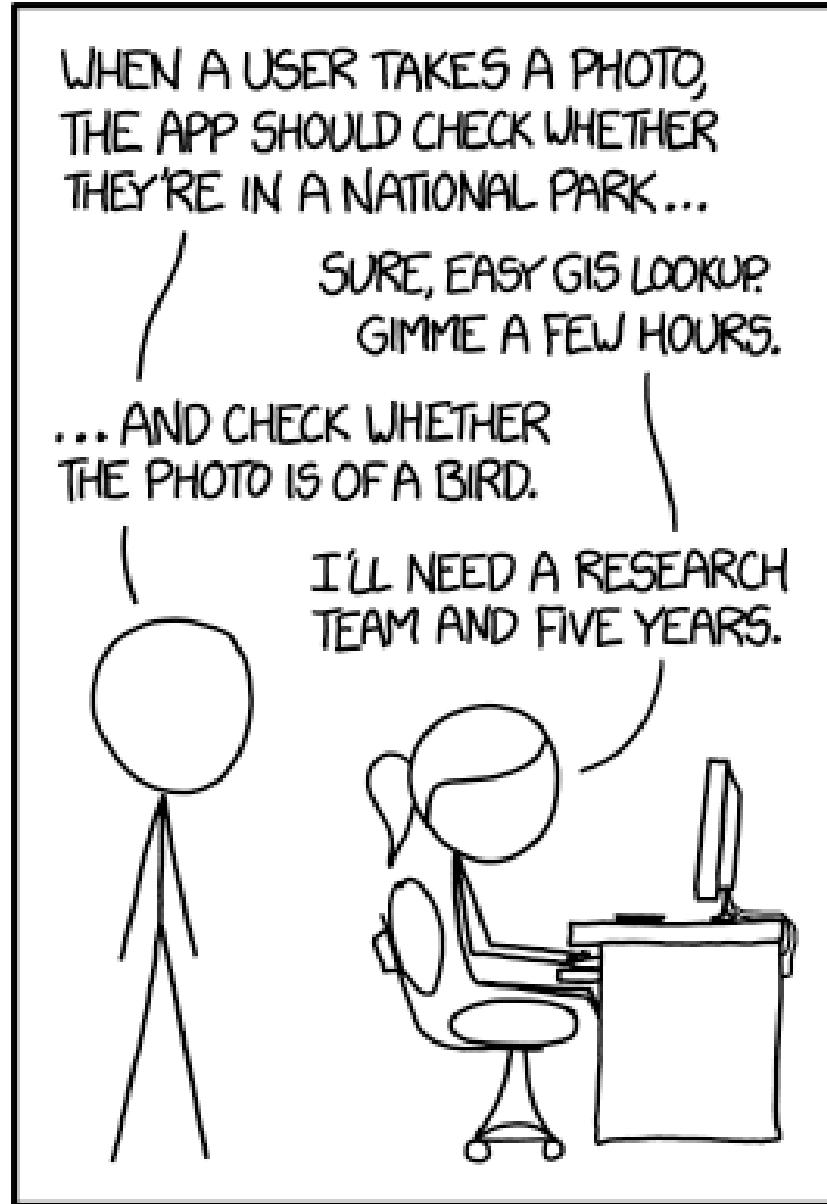
Data Science is learning how to dig through the mountain and get the gold out

Part 2

WHAT IS PATTERN RECOGNITION: A TOO SHORT INTRO

On One Hand

Relevant XKCD



IN CS, IT CAN BE HARD TO EXPLAIN
THE DIFFERENCE BETWEEN THE EASY
AND THE VIRTUALLY IMPOSSIBLE.

On The Other



What is classification?

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**:

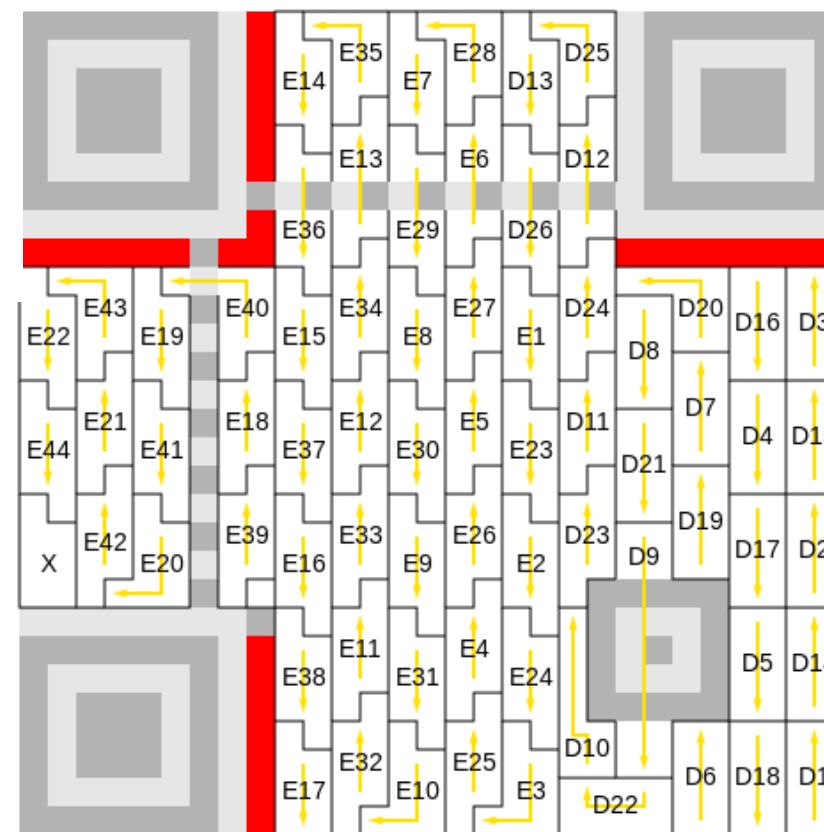
- ask a good question
- have good data

EXAMPLE

QR CODES



- Positioning/Orientation
- Format Information
- Timing marks
- Version Information
- Spacing
- Alignment



■ Fixed Patterns ■ Format Info

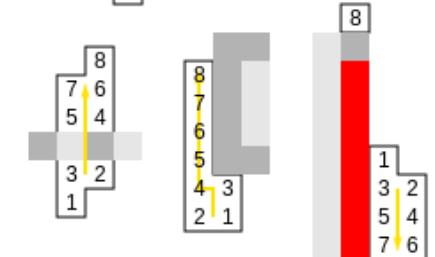
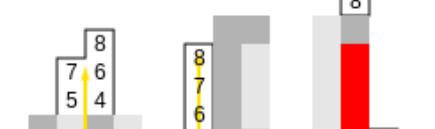
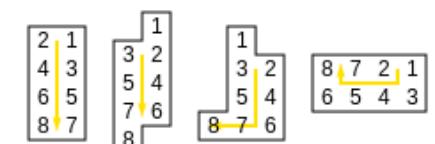
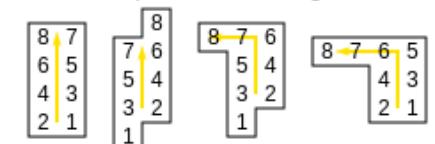
D: Data, E: Error Correction, X: Unused
Error Correction Level H is shown

Block 1 Codewords: D1–D13, E1–E22

Block 2 Codewords: D14–D26, E23–E44

Message Data: D1–D13, D14–D26

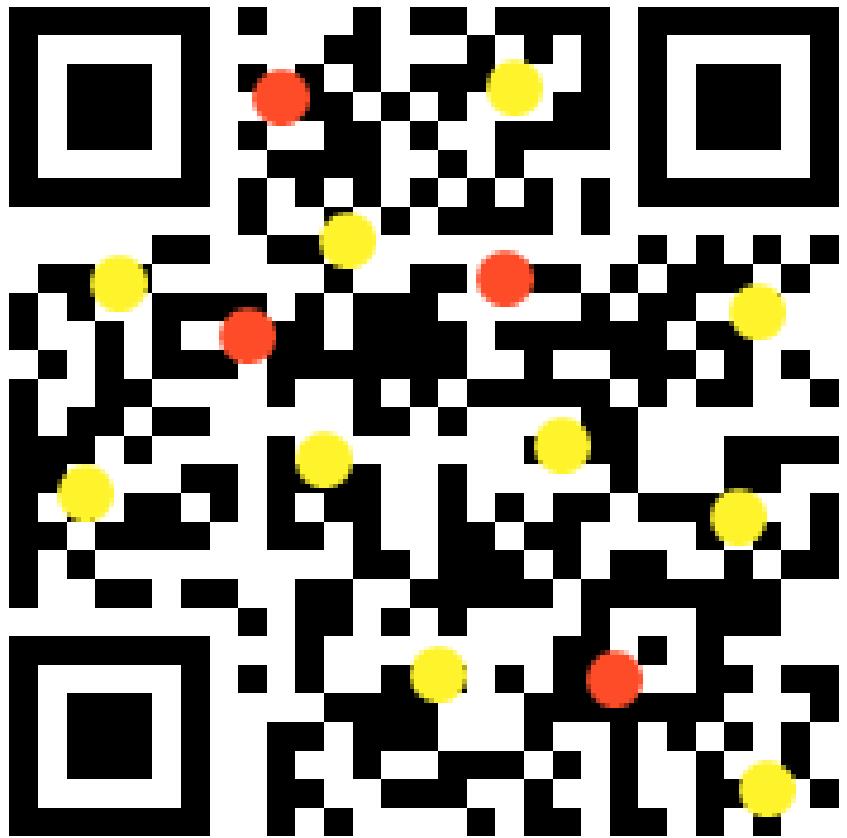
Bit order (1 is the most significant bit):

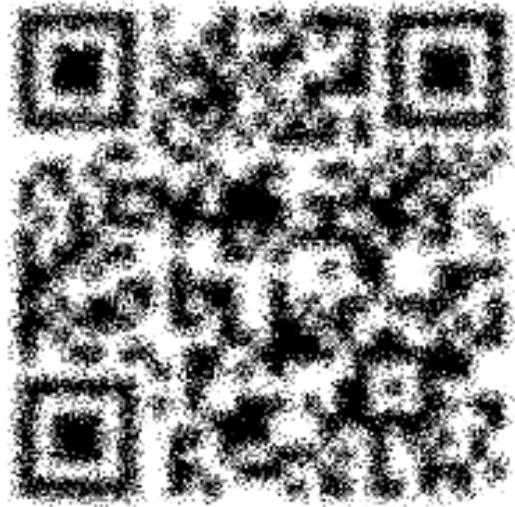












MOAR EXAMPLES

Machine Perception

- Build a machine that can recognize patterns:
 - Email or spam
 - Speech recognition
 - OCR (Optical Character Recognition)
- We often don't care *how* it works!?!
 - Lives or dies based on quality of training data

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_FONT_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

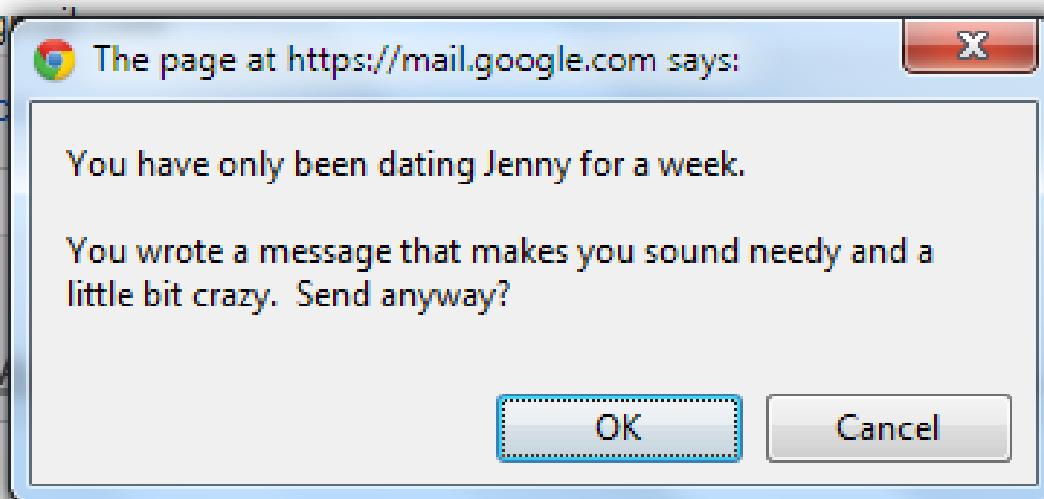
<3 <3 <3

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



Kinect



Real-Time Human Pose Recognition in Parts from Single Depth Images

Jamie Shotton

Andrew Fitzgibbon

Mat Cook

Toby Sharp

Mark Finocchio

Richard Moore

Alex Kipman

Andrew Blake

Microsoft Research Cambridge & Xbox Incubation

Abstract

We propose a new method to quickly and accurately predict 3D positions of body joints from a single depth image, using no temporal information. We take an object recognition approach, designing an intermediate body parts representation that maps the difficult pose estimation problem into a simpler per-pixel classification problem. Our large and highly varied training dataset allows the classifier to estimate body parts invariant to pose, body shape, clothing, etc. Finally we generate confidence-scored 3D proposals of several body joints by reprojecting the classification result and finding local modes.

The system runs at 200 frames per second on consumer hardware. Our evaluation shows high accuracy on both synthetic and real test sets, and investigates the effect of several training parameters. We achieve state of the art accuracy in our comparison with related work and demonstrate improved generalization over exact whole-skeleton nearest neighbor matching.

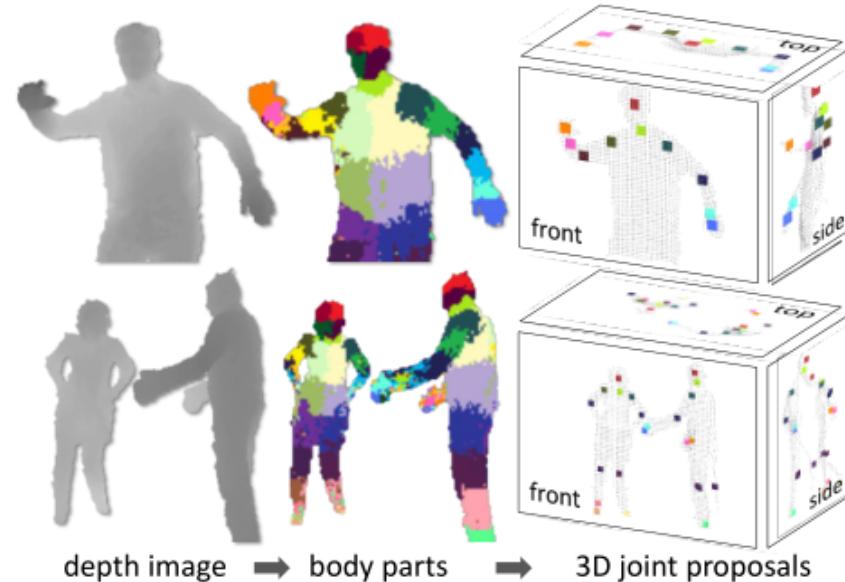
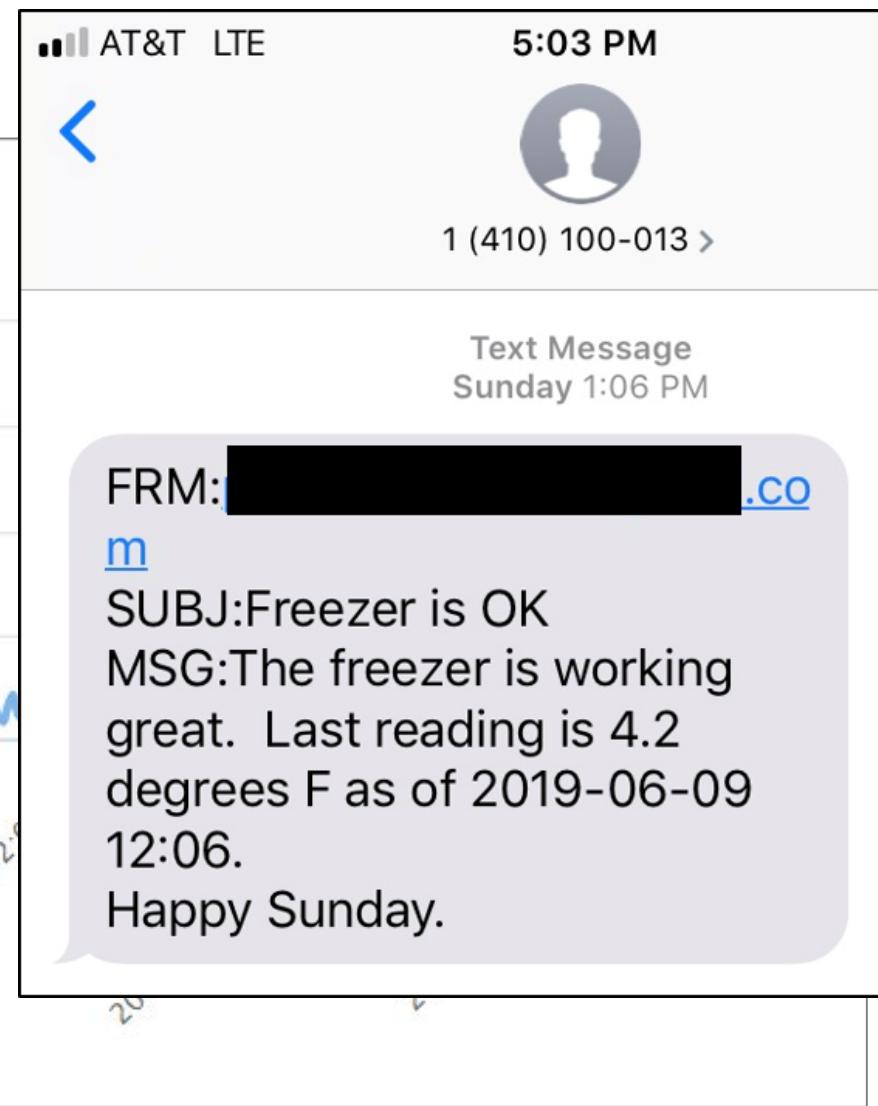
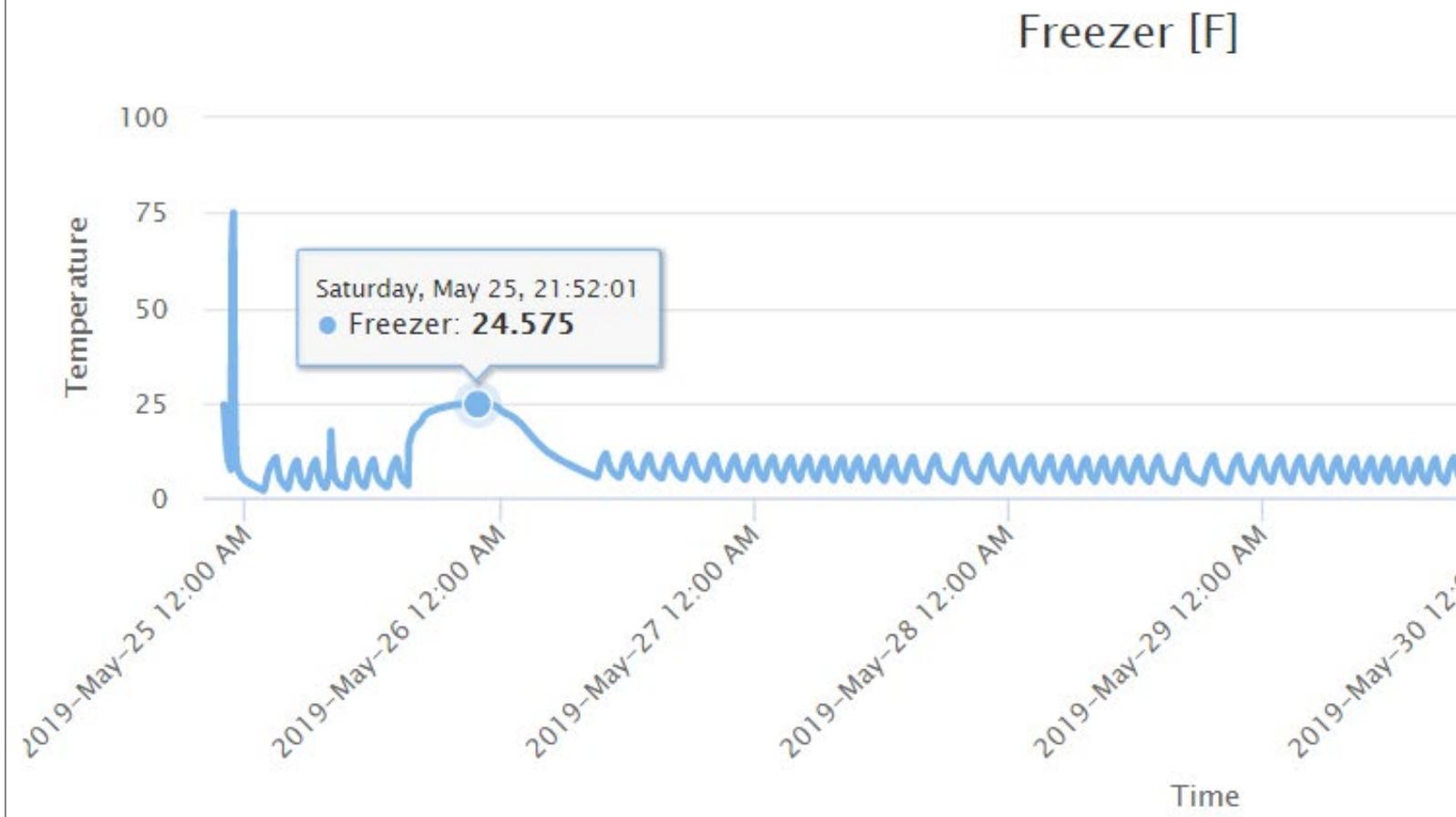


Figure 1. **Overview.** From an single input depth image, a per-pixel body part distribution is inferred. (Colors indicate the most likely part labels at each pixel, and correspond in the joint proposals). Local modes of this signal are estimated to give high-quality proposals for the 3D locations of body joints, even for multiple users.

Freezer

Microcontrollers are immanently practical



A too-short introduction

THEORY

Theory

- Features:
 - Turn a sample into some numbers
 - Curse of Dimensionality
 - Symmetries!
- Models:
 - Lots and lots and lots of classifiers
 - “No Free Lunch” Theorem, and no silver bullets either
- Training: “learn” and evaluate
- Repeat?

Feature Choice - Examples

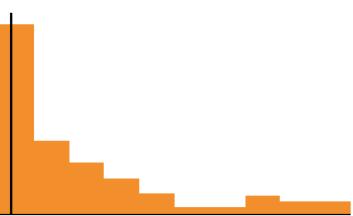
What does Spotify use for songs?

<https://developer.spotify.com/documentation/web-api/reference/tracks/get-audio-features/>

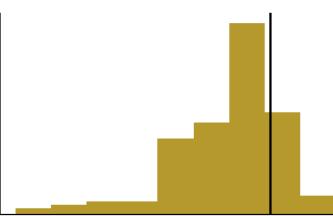
How Songs on Spotify's 'Grime Shutdown' Playlist Compare to my Music Taste

The song '*Strictly Business*' by Shorty, Wiley has an average variance of 0.11 compared with the mean of each audio feature.

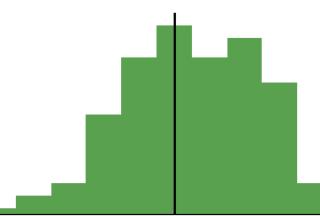
Acousticness



Danceability



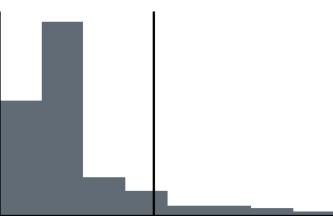
Energy



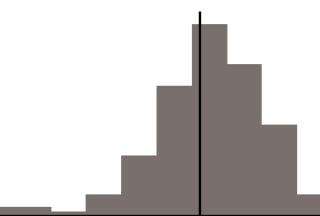
Instrumentalness



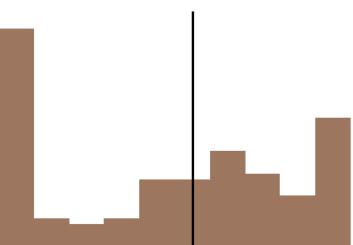
Liveness



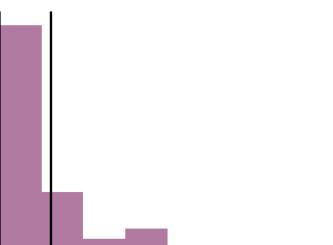
Loudness



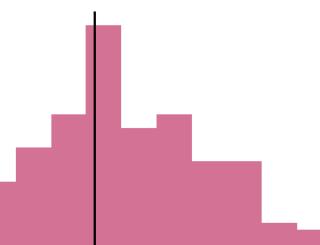
Key



Speechiness



Valence



acousticness number<float>

A confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic.

=> 0 <= 1

analysis_url string

A URL to access the full audio analysis of this track. An access token is required to access this data.

danceability number<float>

Danceability describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. A value of 0.0 is least danceable and 1.0 is most danceable.

duration_ms integer

The duration of the track in milliseconds.

energy number<float>

Energy is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity. Typically, energetic tracks feel fast, loud, and noisy. For example, death metal has high energy, while a Bach prelude scores low on the scale. Perceptual features contributing to this attribute include dynamic range, perceived loudness, timbre, onset rate, and general entropy.

id string

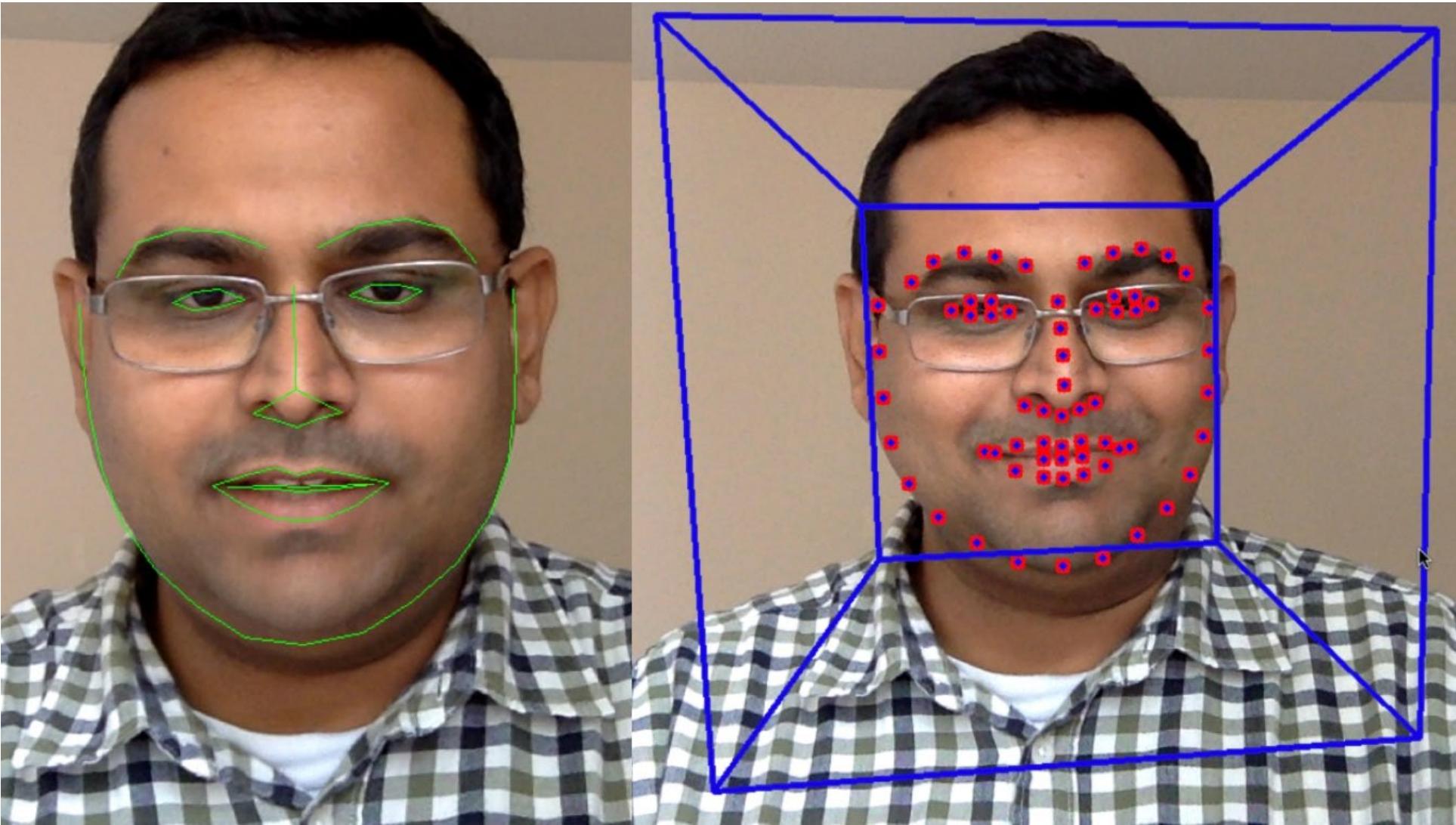
The Spotify ID for the track.

instrumentalness number<float>

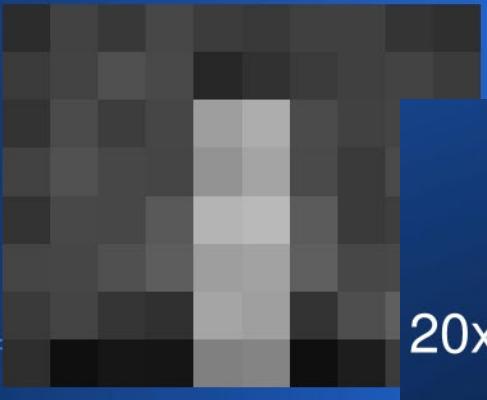
Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Rap or spoken word tracks are clearly "vocal". The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content. Values above 0.5 are intended to represent instrumental tracks, but confidence is higher as the value approaches 1.0.

Feature Choice - Examples

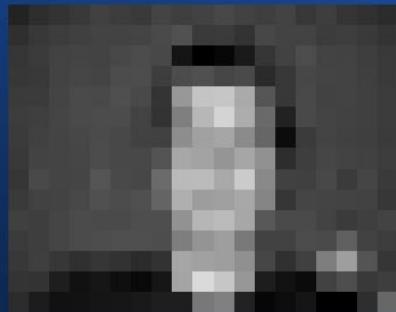
What about faces?



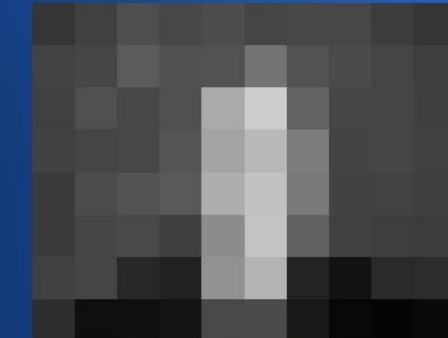
Humans may require more inputs



20x15 pixels



10x7 pixels



<https://www.wired.com/story/best-algorithms-struggle-recognize-black-faces-equally/>

TOM SIMONITE

BUSINESS 07.22.2019 07:00 AM

The Best Algorithms Struggle to Recognize Black Faces Equally

US government tests find even top-performing facial recognition systems misidentify blacks at rates five to 10 times higher than they do whites.

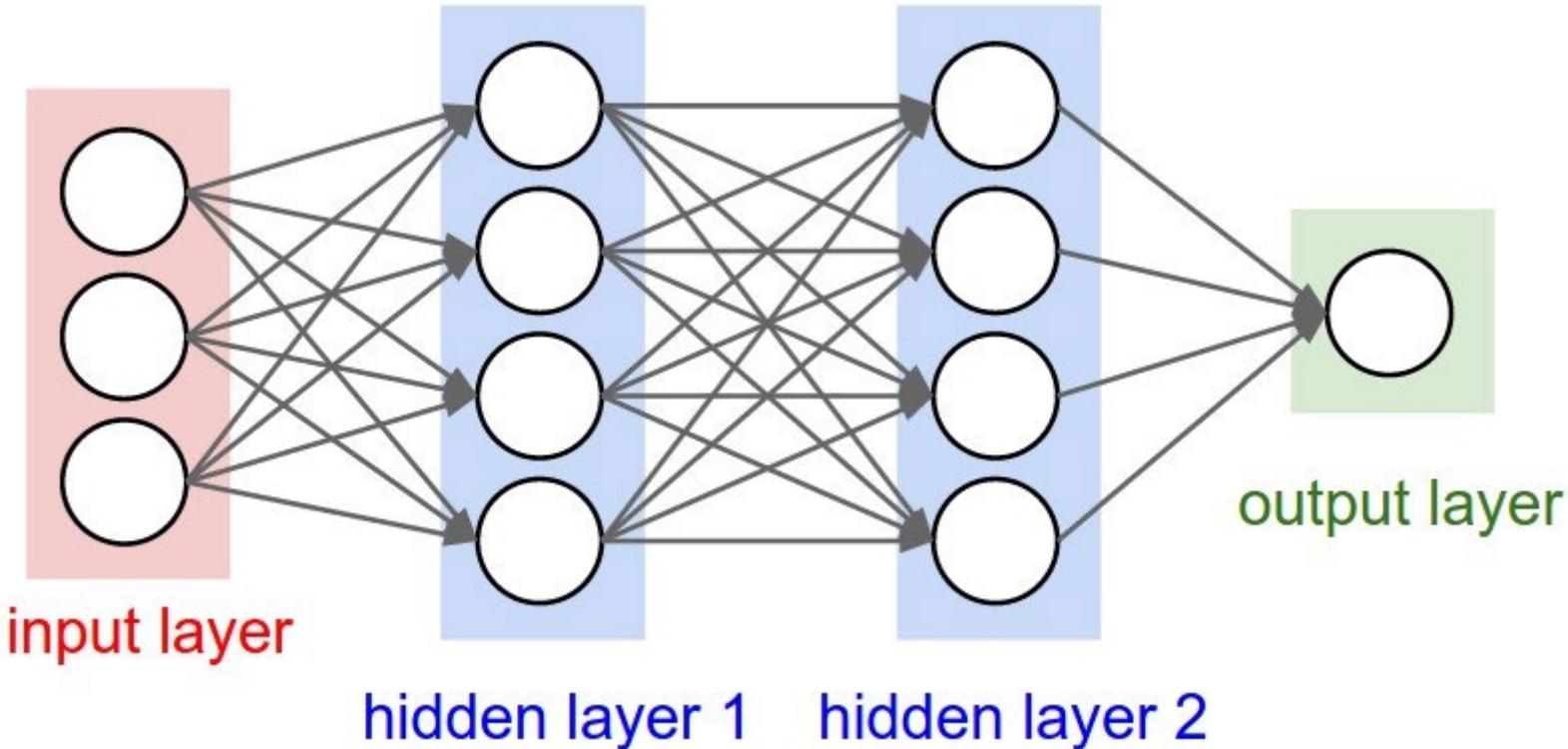
OCTOBER 24, 2020

BLOG, SCIENCE POLICY, SPECIAL EDITION: SCIENCE POLICY AND SOCIAL JUSTICE

Racial Discrimination in Face Recognition Technology

<http://sitn.hms.harvard.edu/flash/2020/racial-discrimination-in-face-recognition-technology/>

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

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)

Other Links

Word Lens

- <https://www.youtube.com/watch?v=h2OfQdYrHRs>
- <https://www.youtube.com/watch?v=zSOlYdlqyTQ>

Other Stuff

- <https://www.netsafe.org.nz/rescam/>
- <http://detexify.kirelabs.org/classify.html>
- <https://www.youtube.com/watch?v=qv6UVOQ0F44> – MARI/O
- <https://www.youtube.com/watch?v=DMyJ6fGUqRI&t=46s> – Rock Wall Pong
- <https://xkcd.com/2236/> - Is it Christmas?

What is Pattern Recognition?

Your turn:

My turn:

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

Other Tricks

REGRESSION

Regression

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

Regression

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

Basically all of the techniques we've learned can be modified to work for regression

http://scikit-learn.org/stable/auto_examples/plot_cv_predict.html

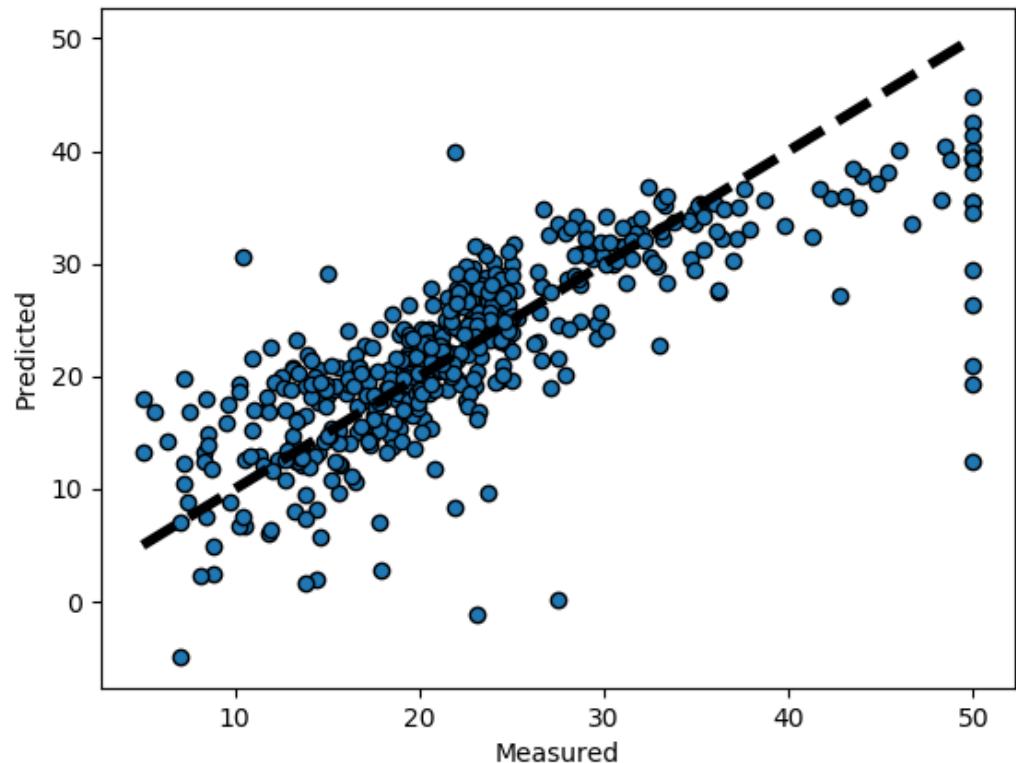
Predicting the price of a house based on features like square footage, # bedrooms, age, etc.

```
from sklearn import datasets
from sklearn.model_selection import cross_val_predict
from sklearn import linear_model
import matplotlib.pyplot as plt

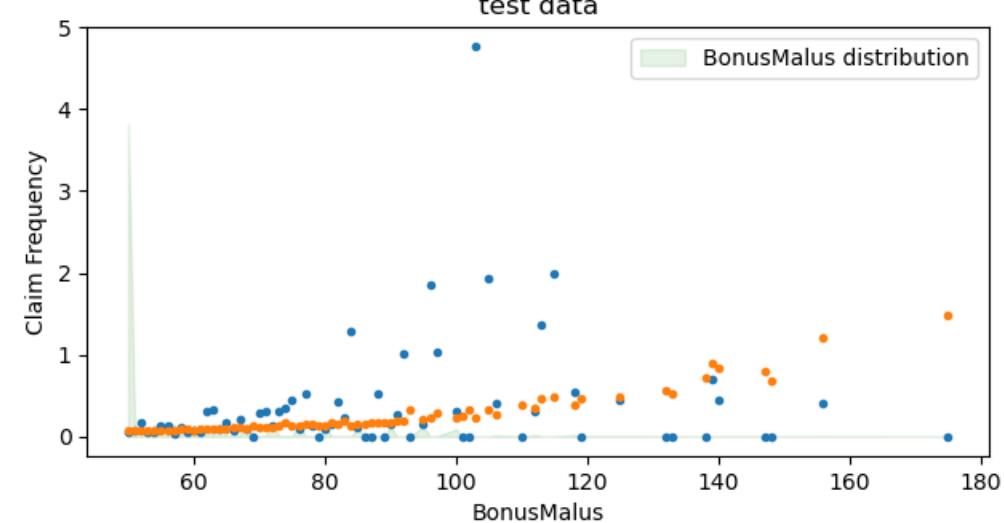
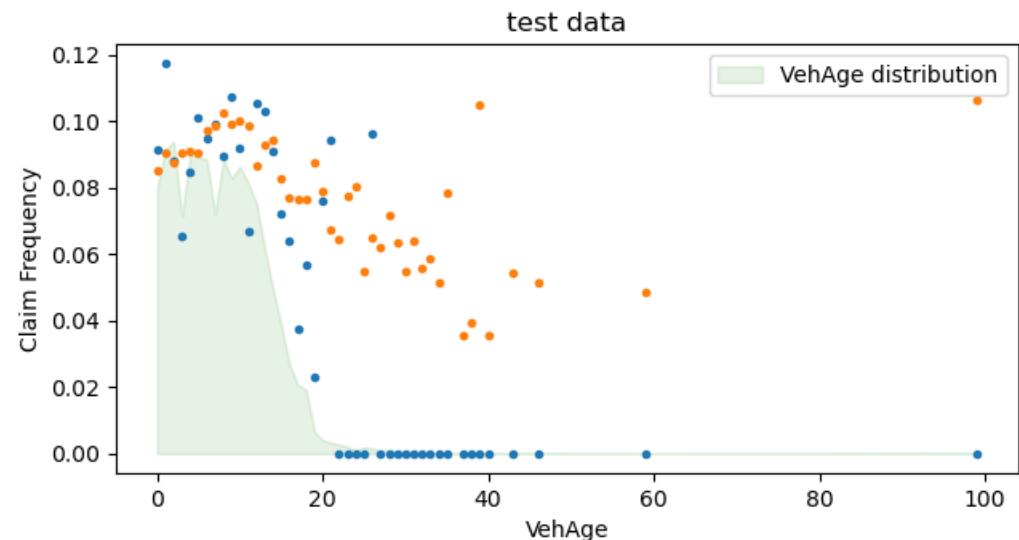
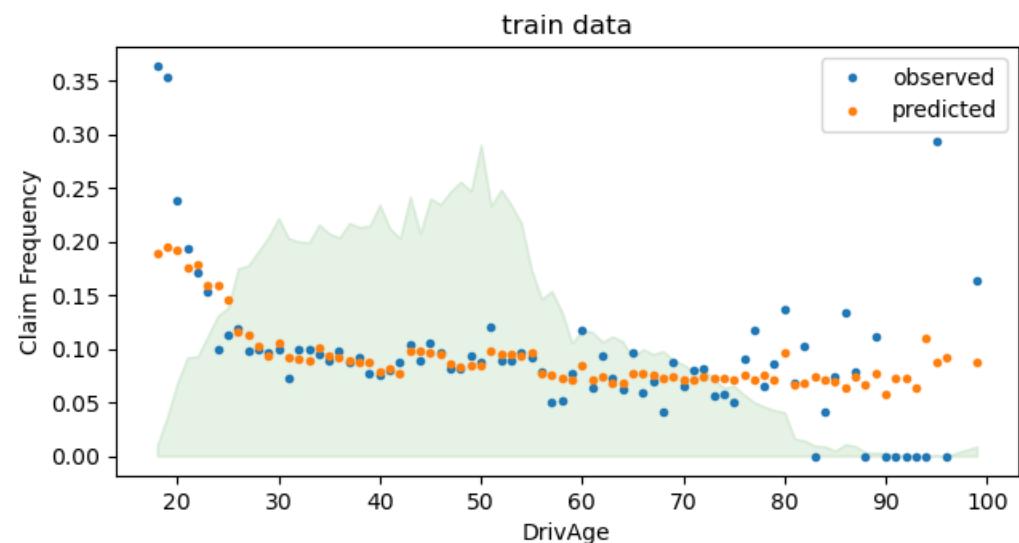
lr = linear_model.LinearRegression()
boston = datasets.load_boston()
y = boston.target

# cross_val_predict returns an array of the same size as `y` where each entry
# is a prediction obtained by cross validation:
predicted = cross_val_predict(lr, boston.data, y, cv=10)

fig, ax = plt.subplots()
ax.scatter(y, predicted, edgecolors=(0, 0, 0))
ax.plot([y.min(), y.max()], [y.min(), y.max()], 'k--', lw=4)
ax.set_xlabel('Measured')
ax.set_ylabel('Predicted')
plt.show()
```



Predicting car insurance claims for drivers in France

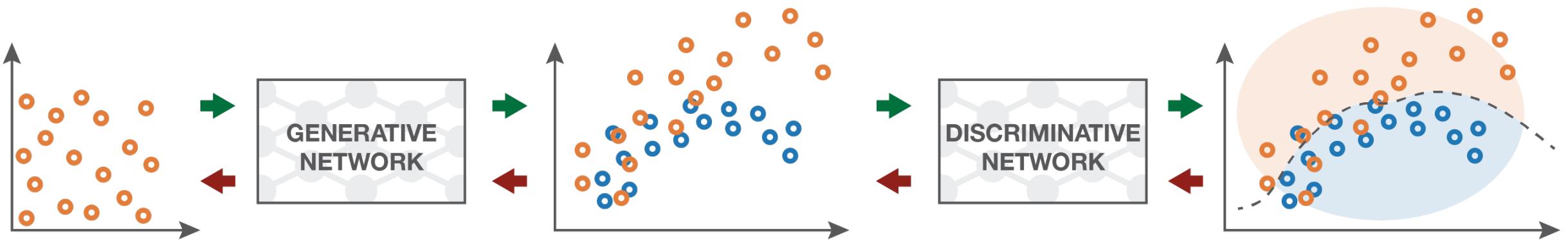


Bonus-Malus (good-bad) is like a good driver score based on how many previous claims a driver has

Other Tricks

GENERATION

■ Forward propagation (generation and classification) ■ Backward propagation (adversarial training)



Input random variables.

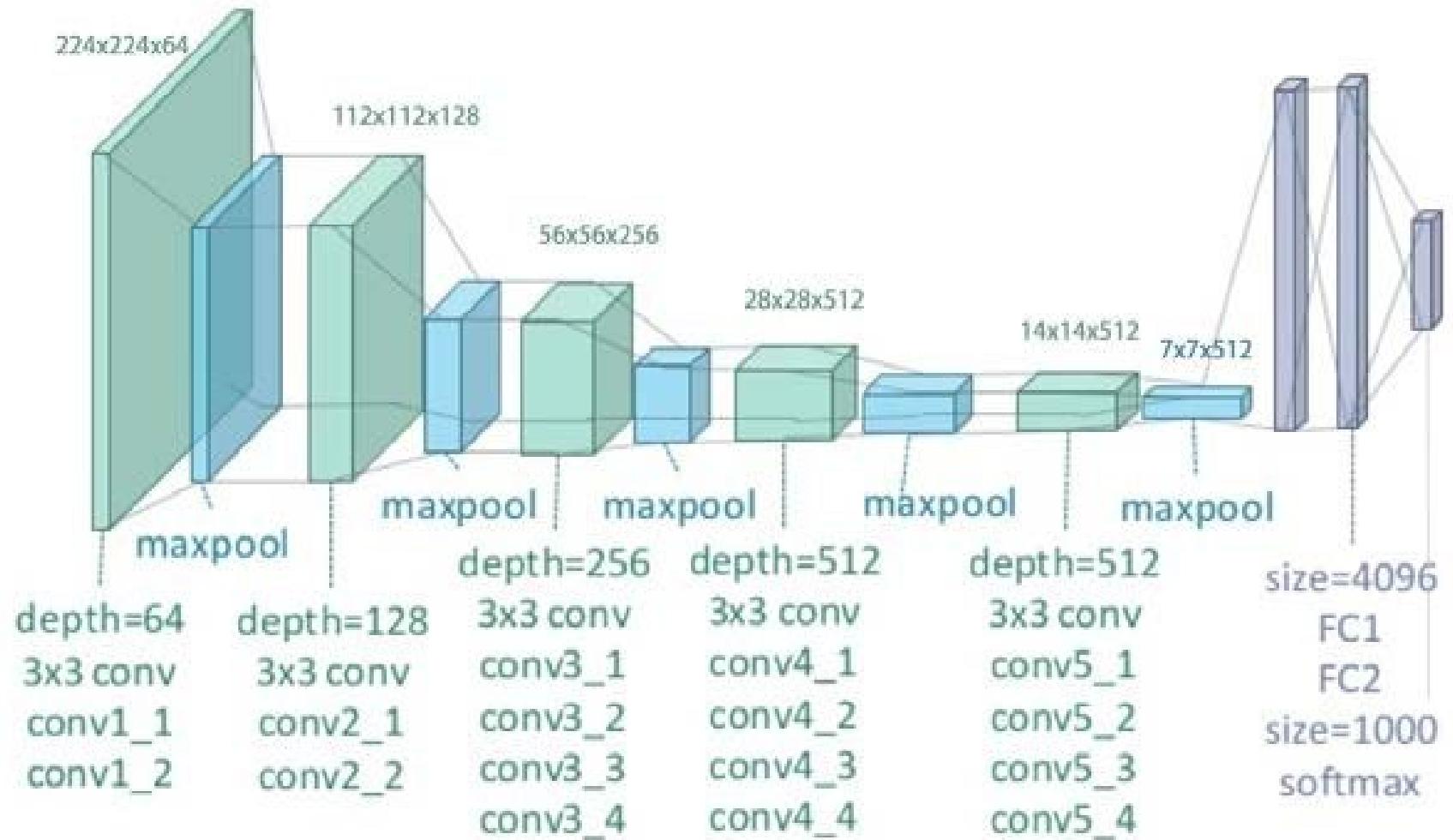
The generative network is trained to **maximise** the final classification error.

The **generated distribution** and the **true distribution** are not compared directly.

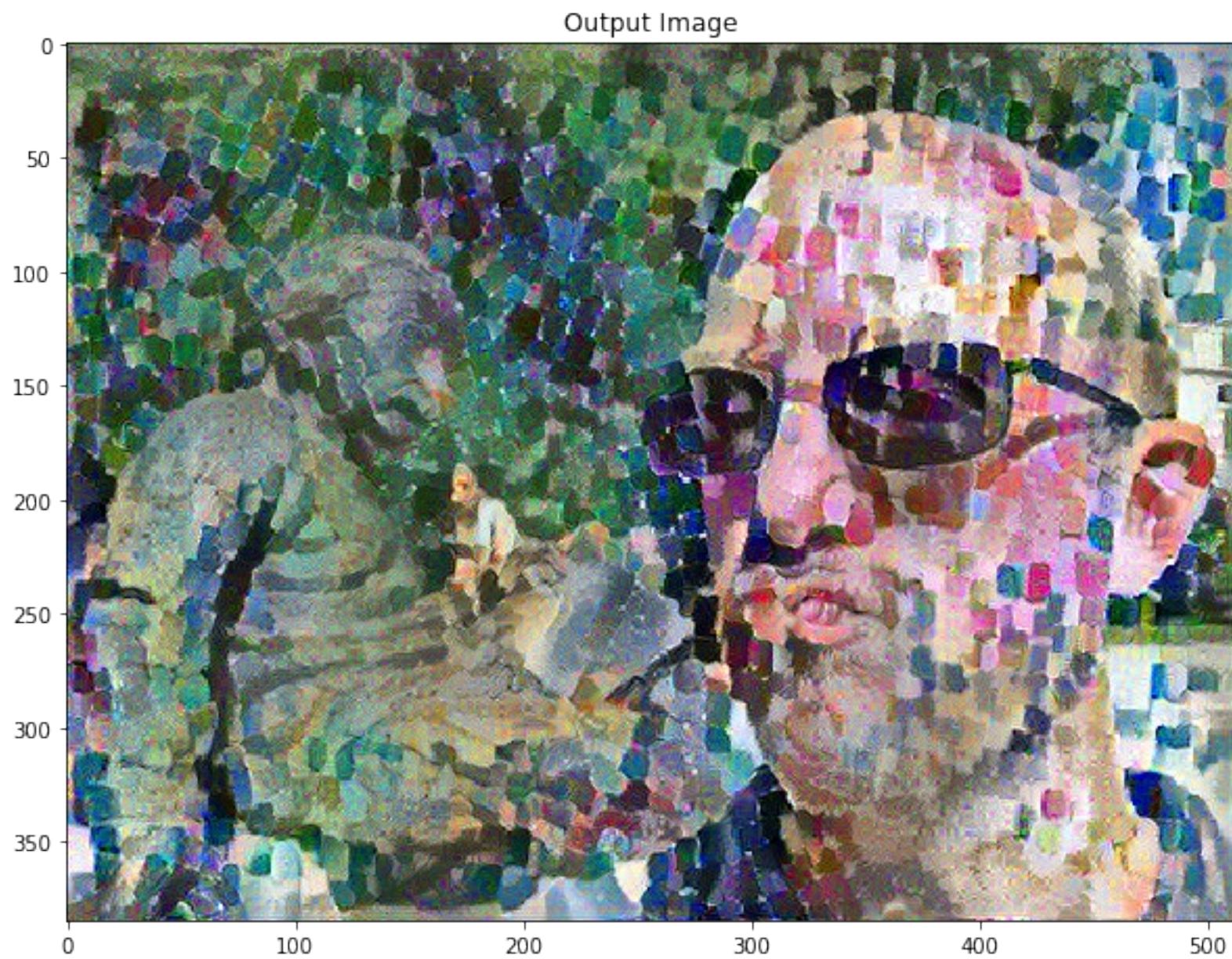
The discriminative network is trained to **minimise** the final classification error.

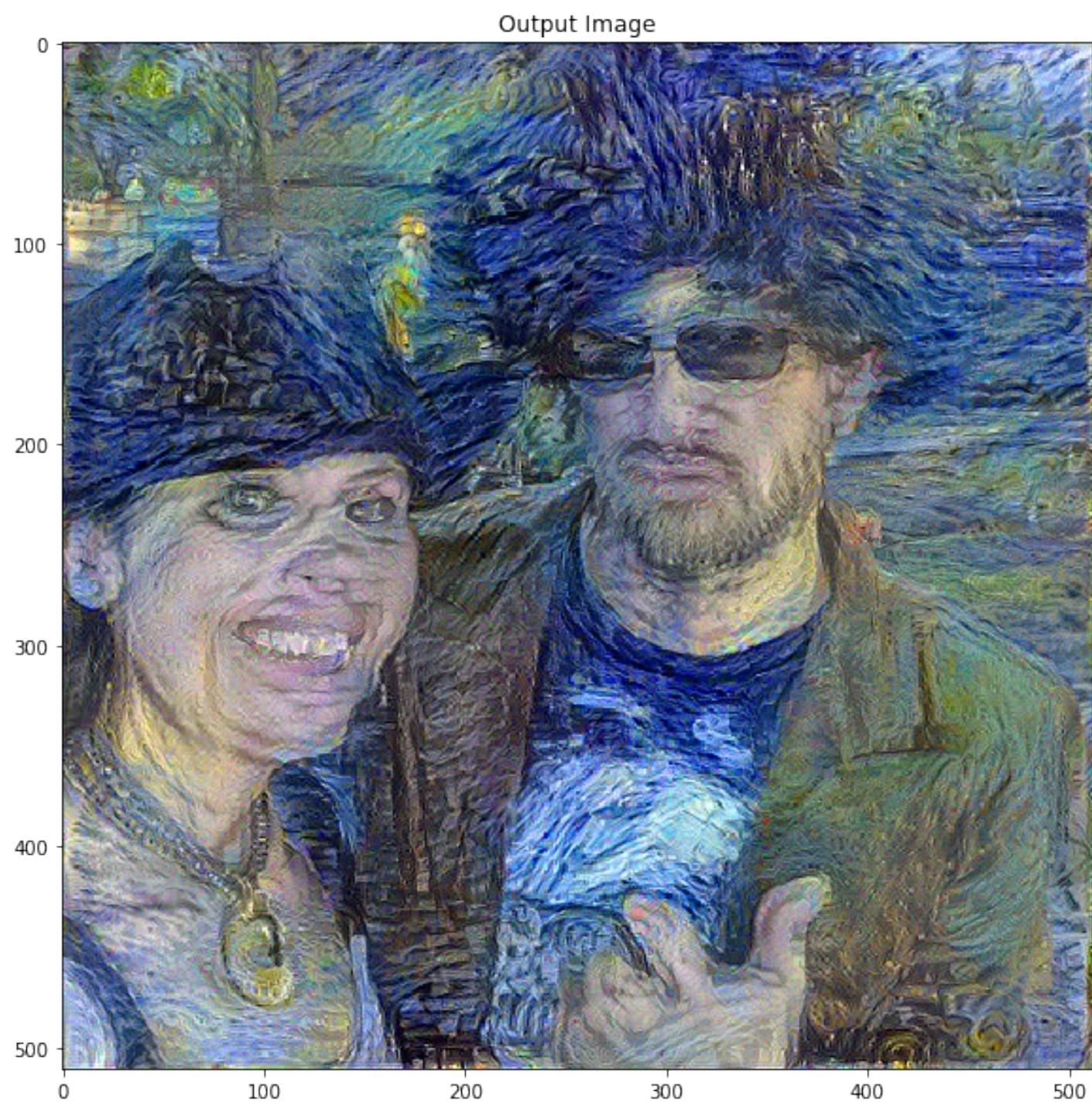
The classification error is the basis metric for the training of both networks.

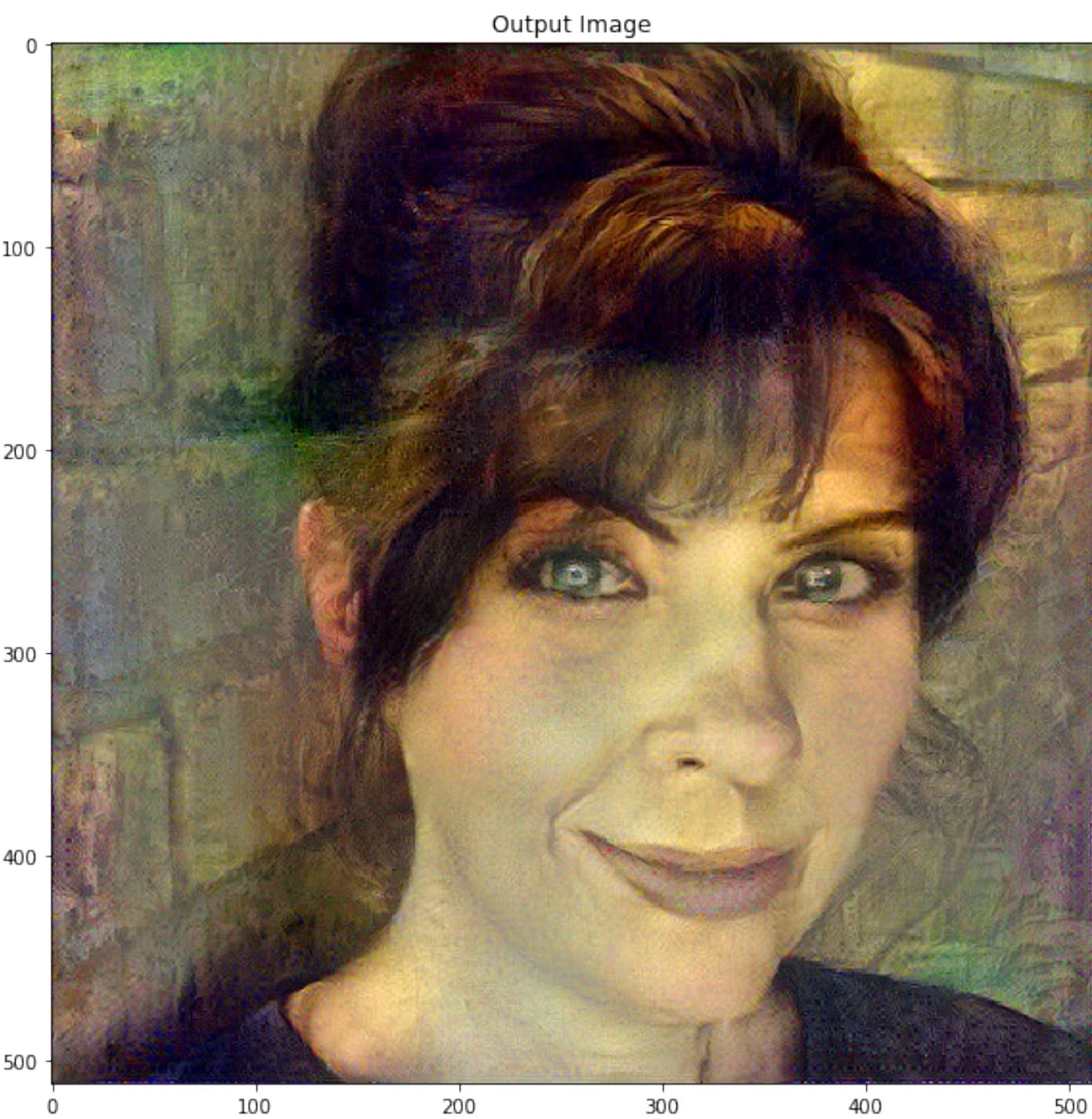
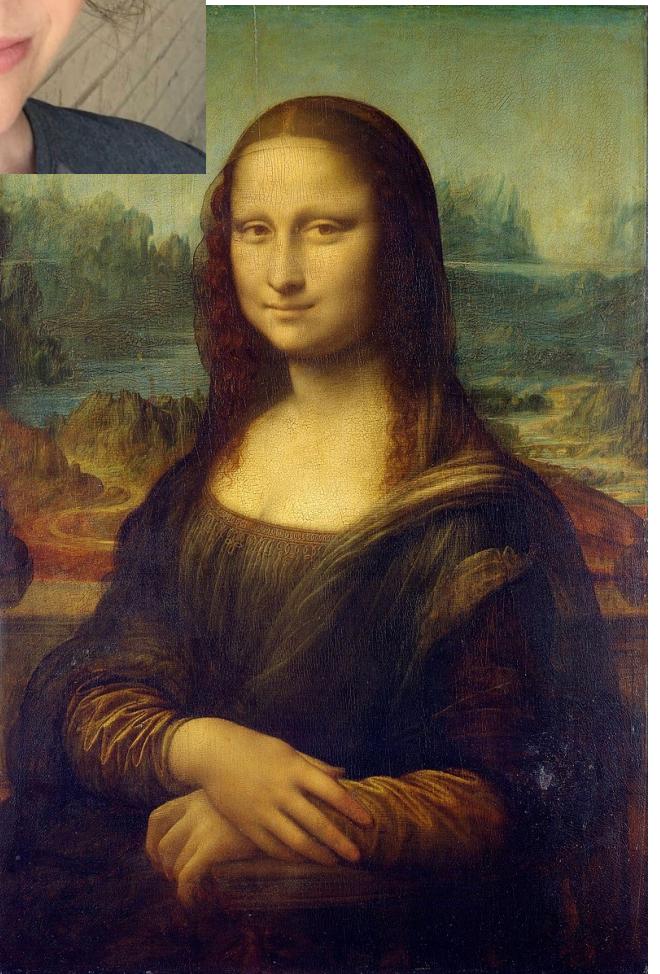
https://colab.research.google.com/github/tensorflow/models/blob/master/research/nst_blogpost/4_Neural%20Style%20Transfer%20with%20Eager%20Execution.ipynb









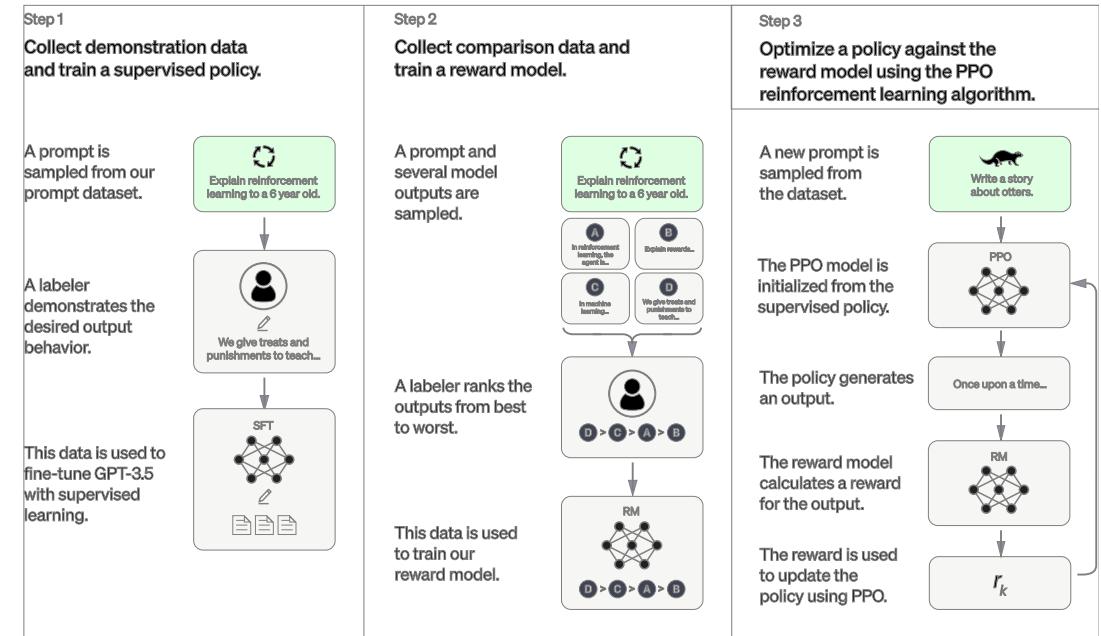






More Examples

- 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/>



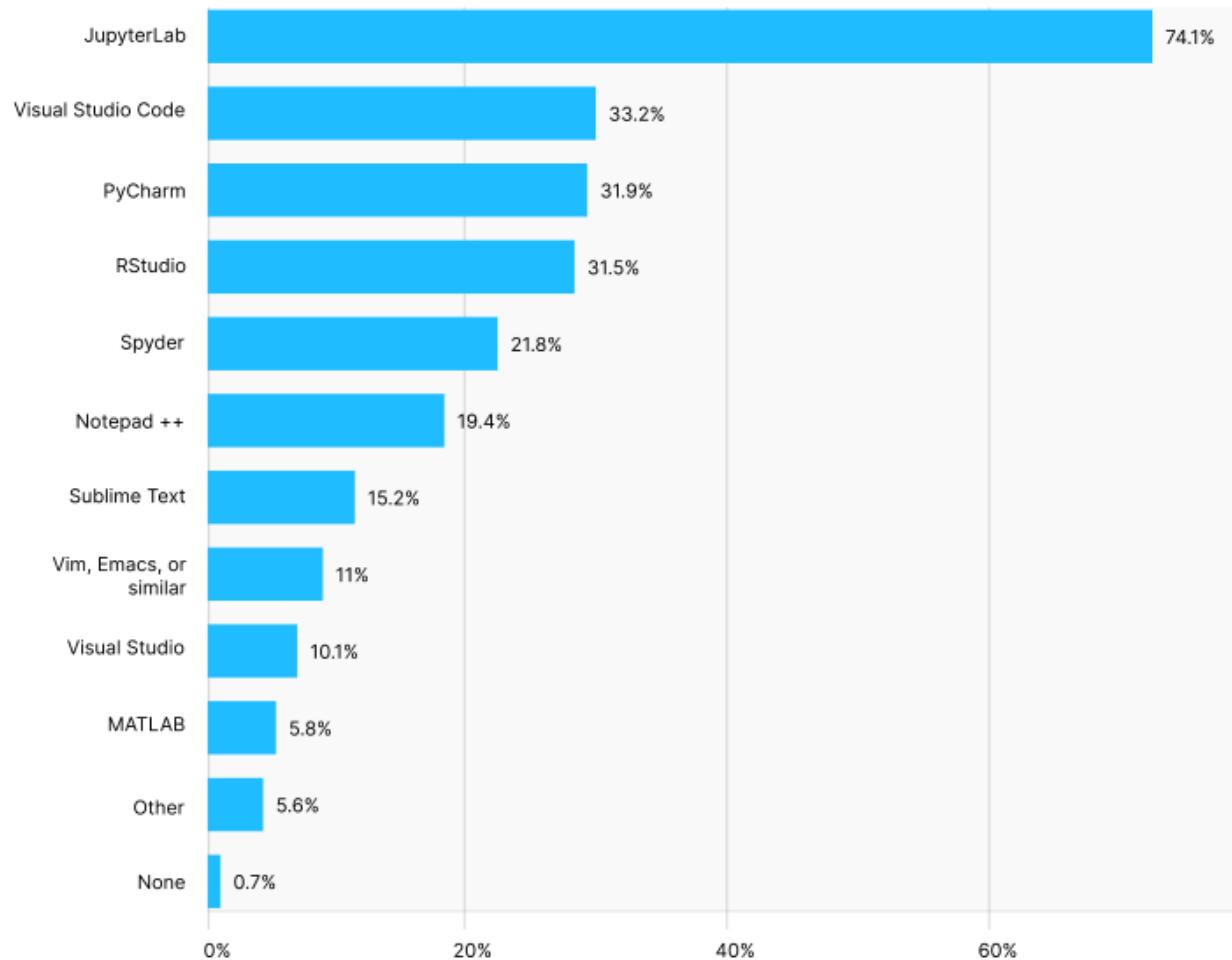
Part 3

TOOLS

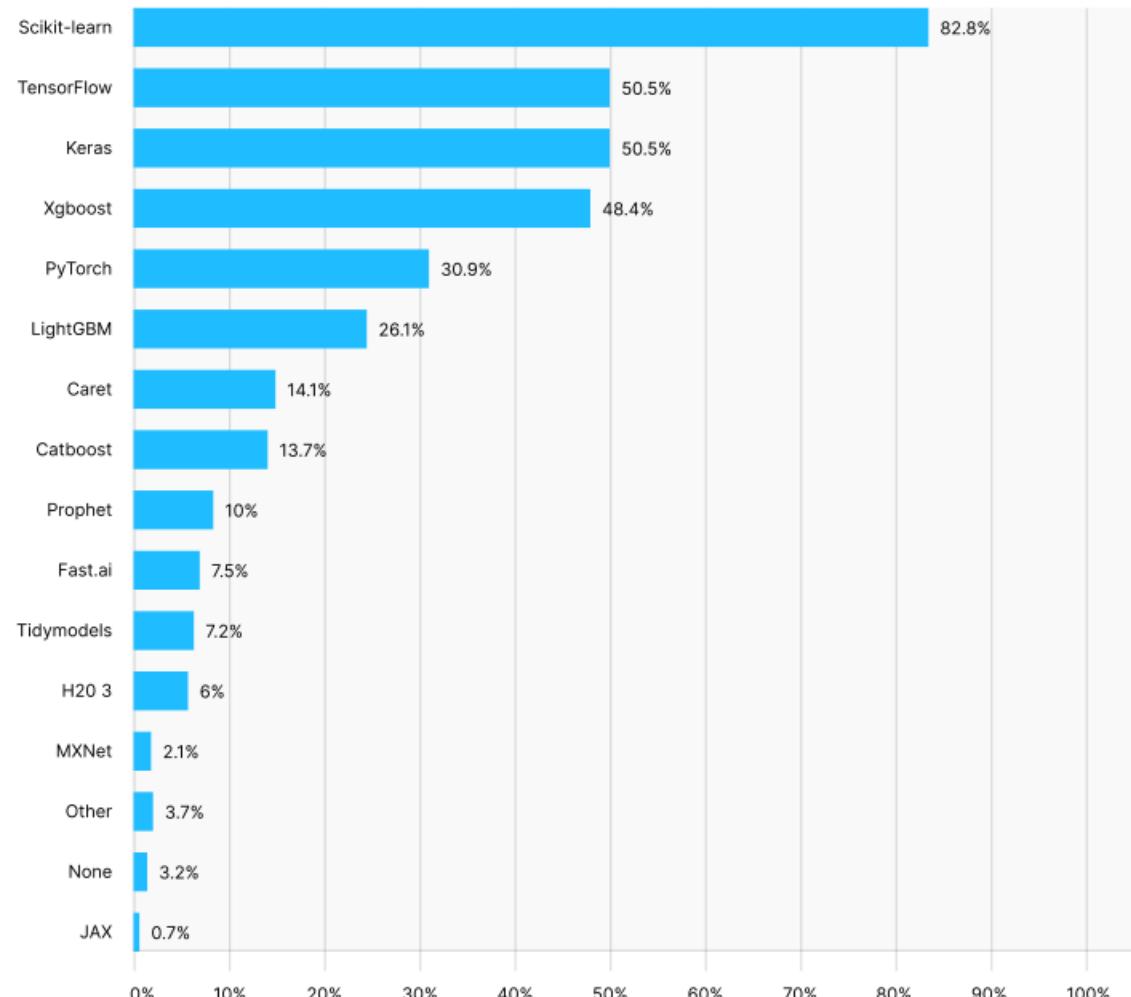
Kaggle: Making Data Science a Sport

<https://www.kaggle.com/kaggle-survey-2020>

POPULAR IDE USAGE



MACHINE LEARNING FRAMEWORK USAGE



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/>

