

SIGNAL PROCESSING & ML IN UBICOMP (CONT)

CSE 590 Ubiquitous Computing | Lecture 4 | April 19

Jon Froehlich • Liang He (TA)

SCHEDULE TODAY: 6:30-9:20

06:35-06:55: Daisy Isibor's discussion of the Saponas et al., CHI'08 reading

06:55-07:05: Apurv Suman's discussion of the Harrison and Hudson Scratch Input paper

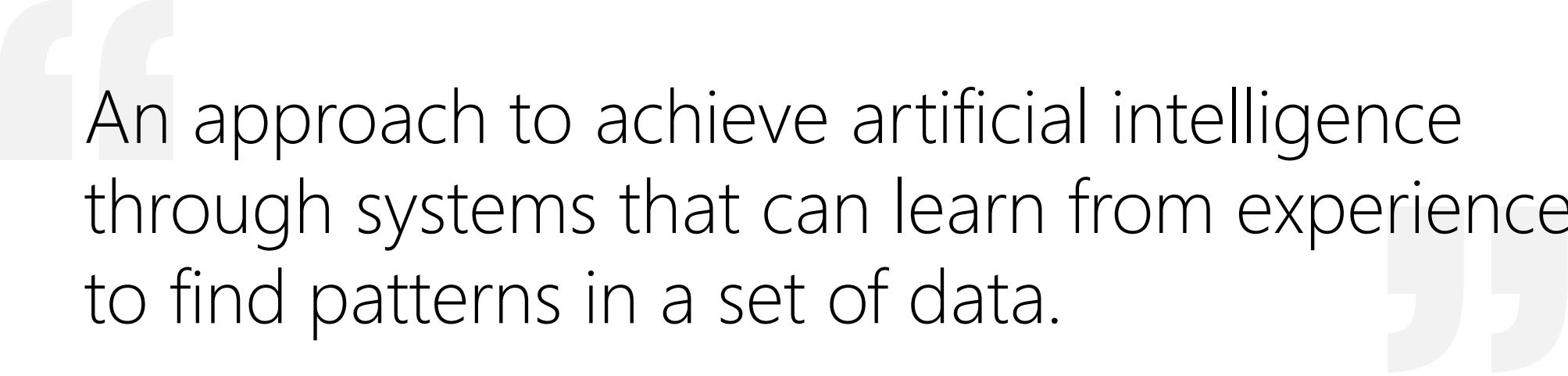
07:05-08:00: Signal processing & ML in ubicomp systems cont (w/ Jupyter Notebook exercises)

07:55-08:05: Break

08:05-08:30: Continue lecture and exercises

08:30-09:20: Work on Part 2 of A2: Gesture Recognizer. Specifically, implementing an SVM.

MACHINE LEARNING



An approach to achieve artificial intelligence through systems that can learn from experience to find patterns in a set of data.

Jason Mayes

Senior Creative Engineering, Google

MACHINE LEARNING

The goal of machine learning is to build computer systems that can adapt and learn from their experience.

Tom Dietterich
Professor, University of Oregon

MACHINE LEARNING

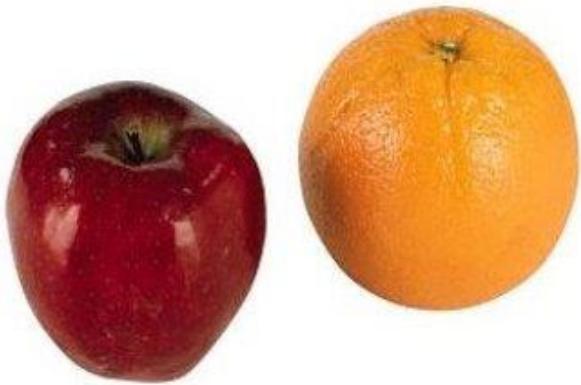
ML involves **teaching a computer to recognize patterns** by example rather than programming a system to recognize such patterns with specific rules.

In other words, **ML is about creating algorithms that learn complex models/functions from data** and then make predictions on similar data.

SUPERVISED LEARNING

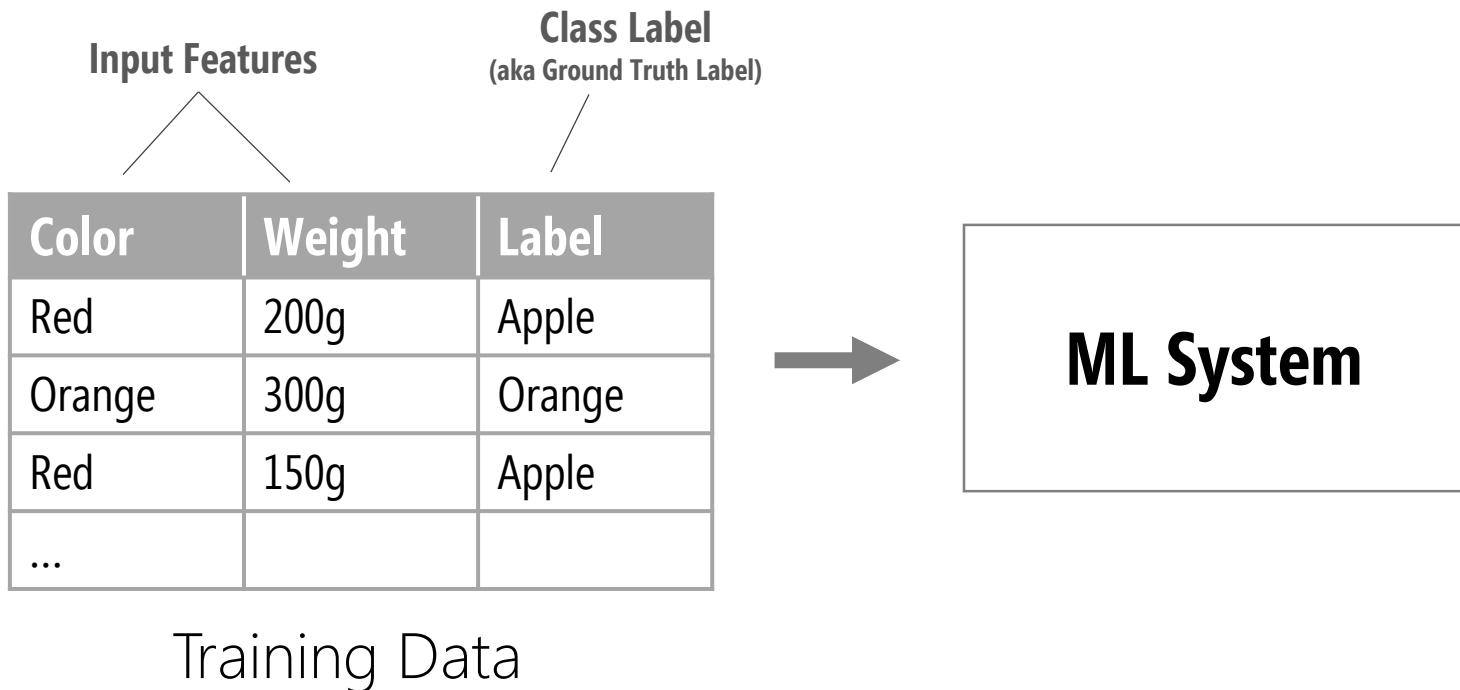
1. Input representative data to train ML system
2. ML system learns patterns from data
3. Then feed in data ML system has never seen before and it returns its best guess for what that data is

SUPERVISED LEARNING



Example: let's build an ML model for
recognizing apples and oranges....

SUPERVISED LEARNING



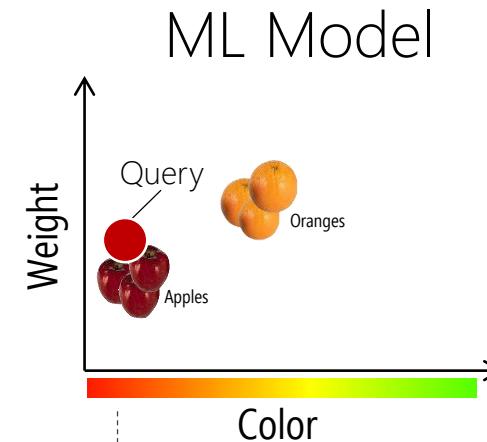
SUPERVISED LEARNING

Input Features		Class Label (aka Ground Truth Label)
Color	Weight	Label
Red	200g	Apple
Orange	300g	Orange
Red	150g	Apple
...		

Training Data

Color	Weight	Label
Red	170g	Unknown

Query Data



ML System

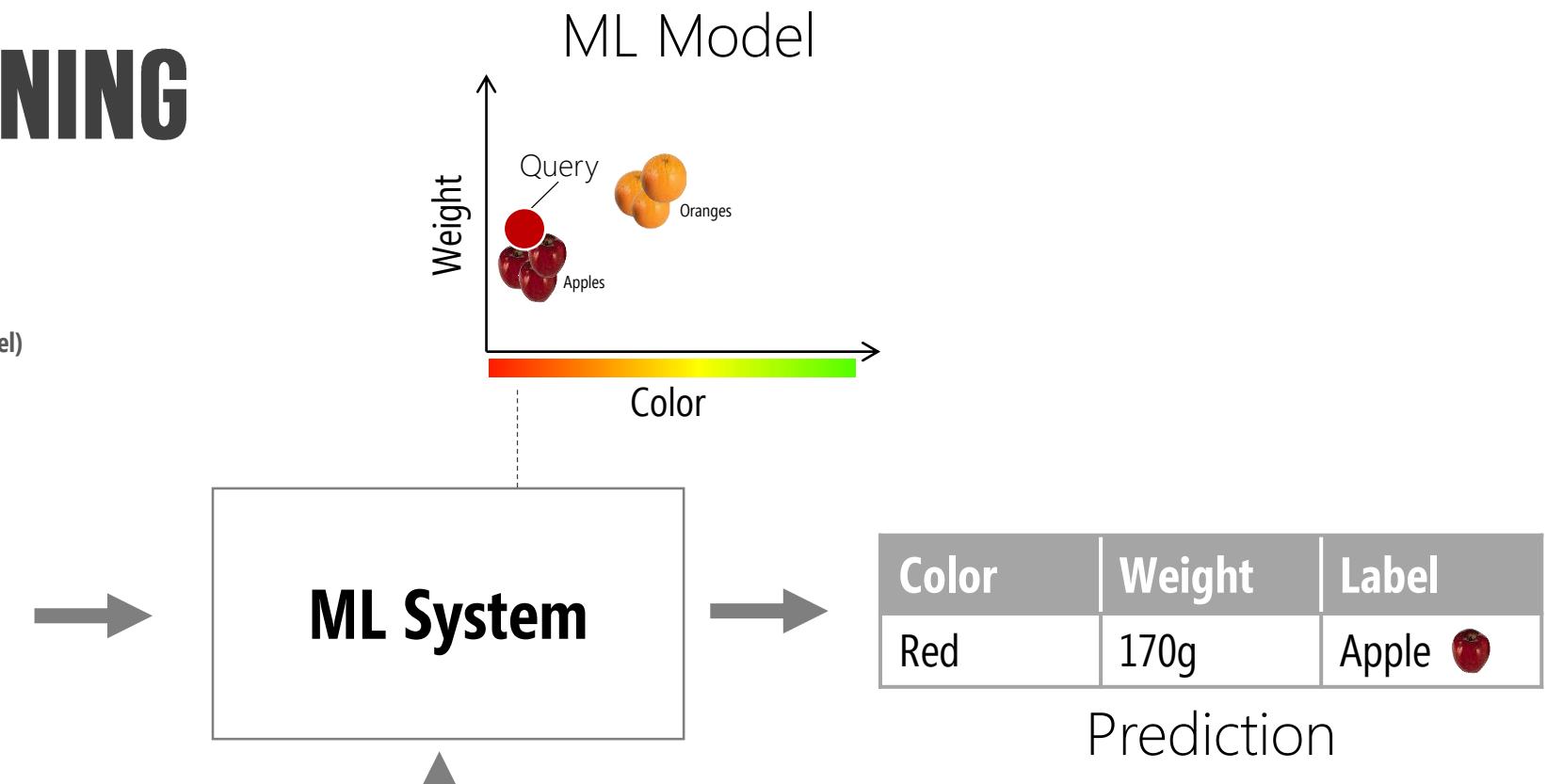
SUPERVISED LEARNING

Input Features		Class Label (aka Ground Truth Label)
Color	Weight	Label
Red	200g	Apple
Orange	300g	Orange
Red	150g	Apple
...		

Training Data

Color	Weight	Label
Red	170g	Unknown

Query Data



ML System

SUPERVISED LEARNING PROBLEMS AND APPROACHES

CLASSIFICATION PROBLEMS

When you need to classify some input vector into a category (*e.g.*, Apple or Orange). Output is some class label C .

Popular approaches:

Support Vector Machines (SVMs)

Random forests

Neural Networks

Deep Learning

REGRESSION PROBLEMS

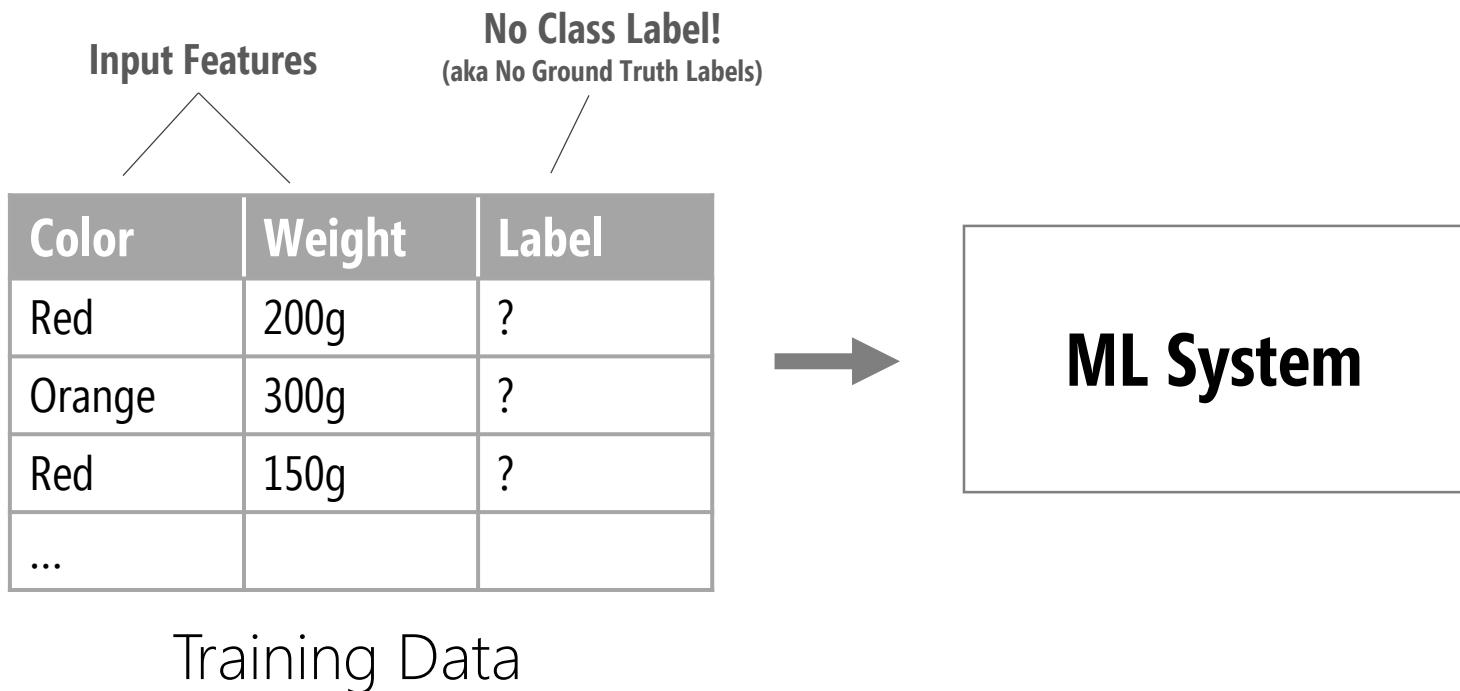
When you want to predict some numeric value (*e.g.*, house prices next year in Seattle). Output has infinitely many values (a number).

Popular approaches:

Linear regression

Random forests

UNSUPERVISED LEARNING



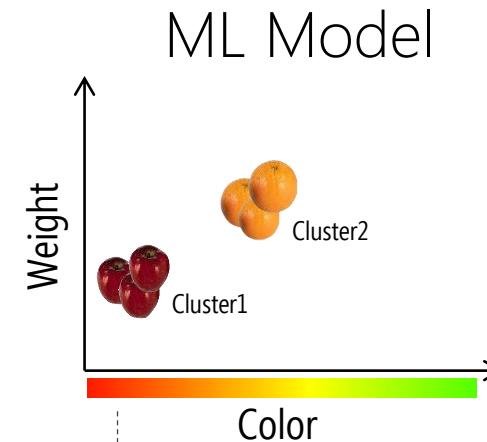
UNSUPERVISED LEARNING

Input Features

No Class Label!
(aka No Ground Truth Labels)

Color	Weight	Label
Red	200g	?
Orange	300g	?
Red	150g	?
...		

Training Data



UNSUPERVISED LEARNING

Input Features

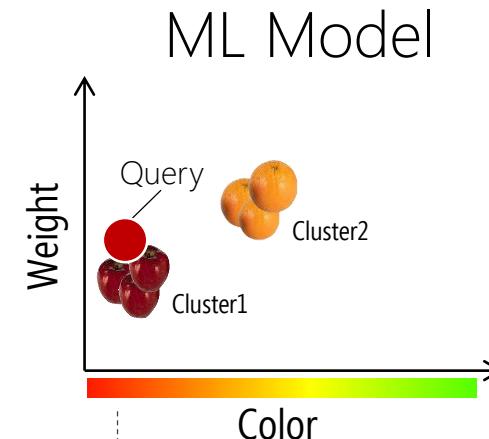
No Class Label!
(aka No Ground Truth Labels)

Color	Weight	Label
Red	200g	?
Orange	300g	?
Red	150g	?
...		

Training Data

Color	Weight	Label
Red	170g	?

Query Data



ML System

Color	Weight	Label
Red	170g	Cluster1

Prediction

UNSUPERVISED LEARNING PROBLEMS AND APPROACHES

CLUSTERING PROBLEMS

When you want to discover inherent groupings in your data (*e.g.*, grouping customers by purchasing behavior)

Popular approaches:

K-means clustering

ASSOCIATION PROBLEMS

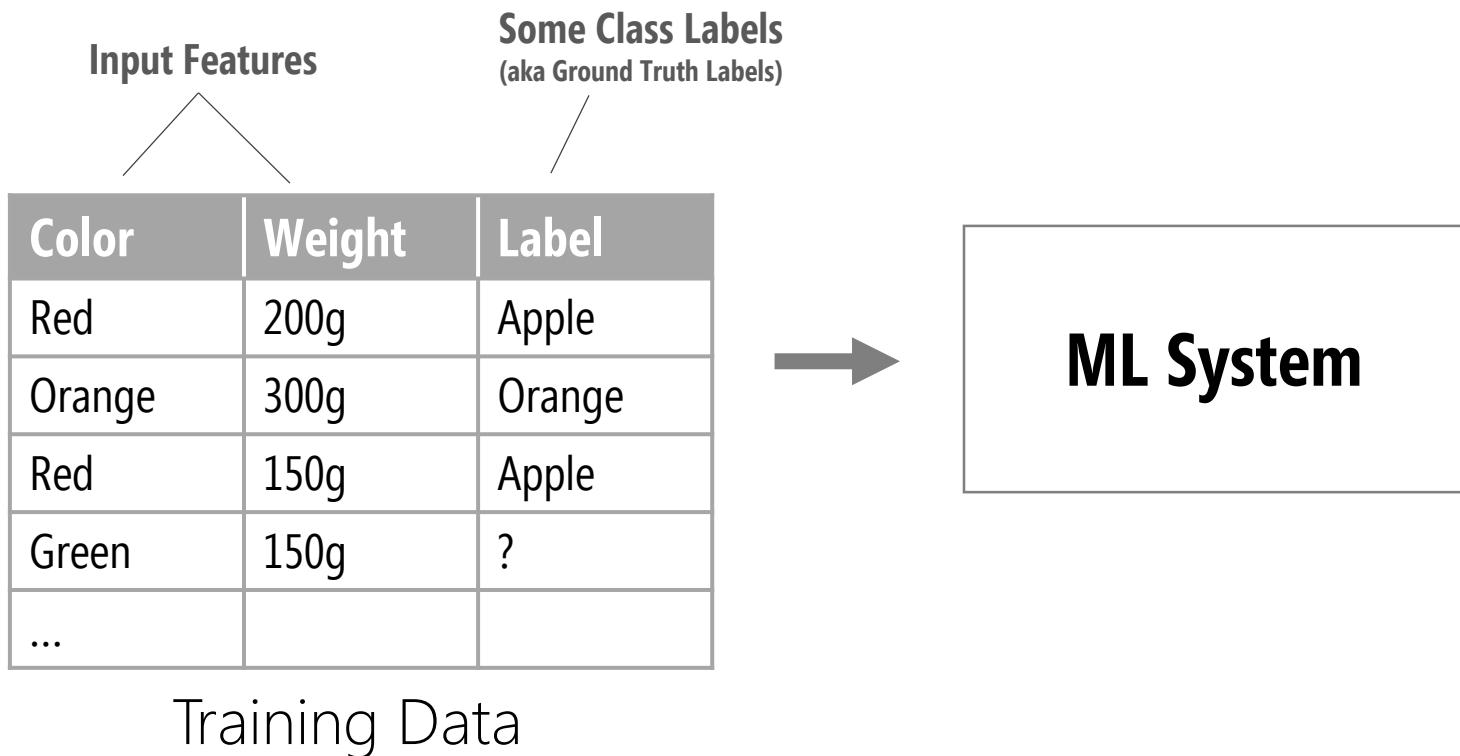
When you want to discover associative rules between large portions of data (*e.g.*, people that buy X also tend to buy Y)

Popular approaches:

Apriori algorithm

Collaborative filtering

SEMISUPERVISED LEARNING



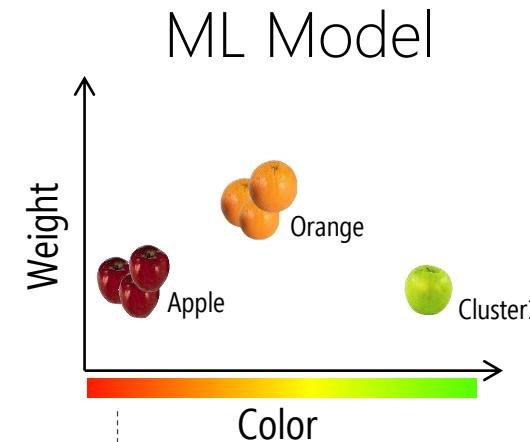
SEMISUPERVISED LEARNING

Input Features

Some Class Labels
(aka Ground Truth Labels)

Color	Weight	Label
Red	200g	Apple
Orange	300g	Orange
Red	150g	Apple
Green	150g	?
...		

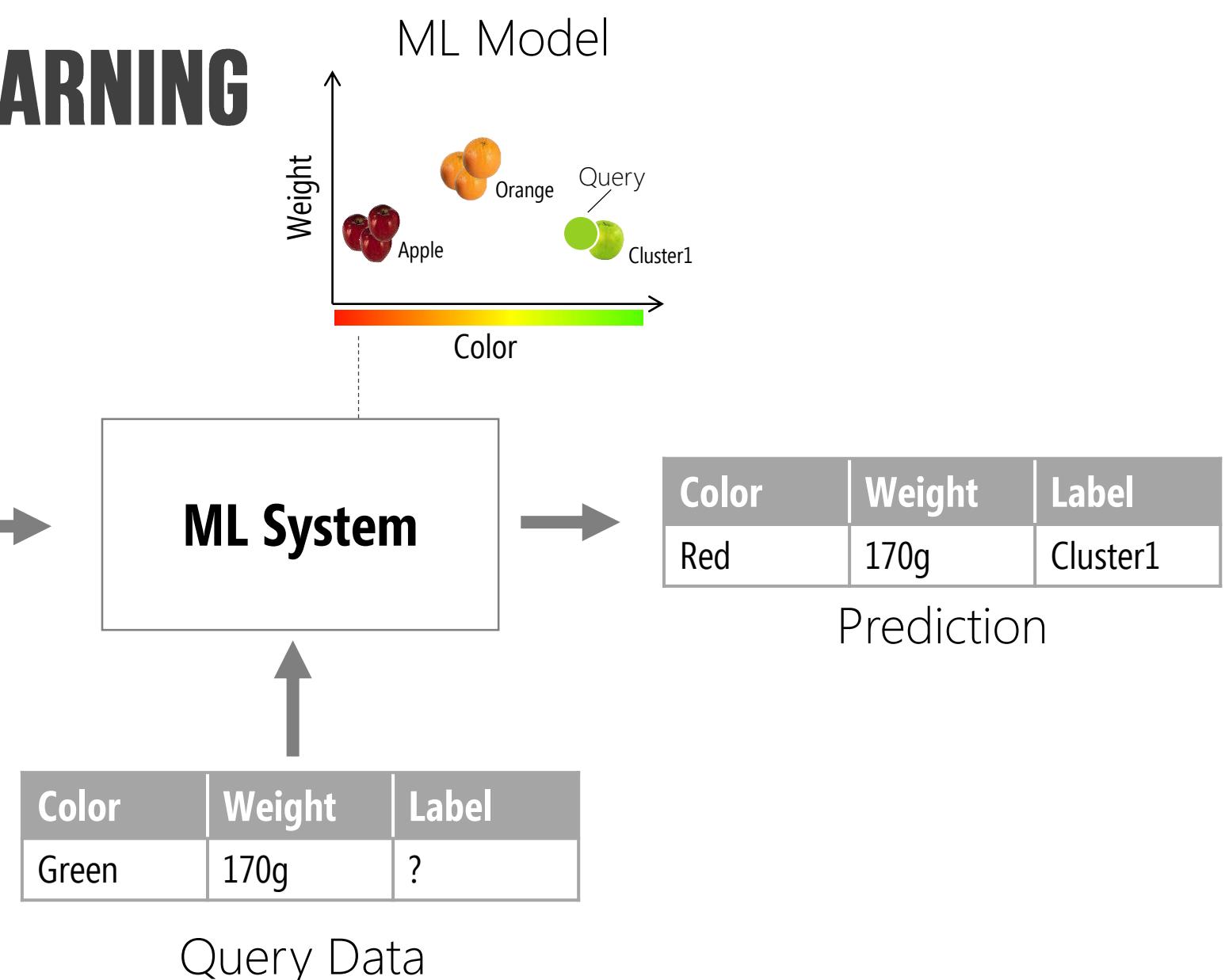
Training Data



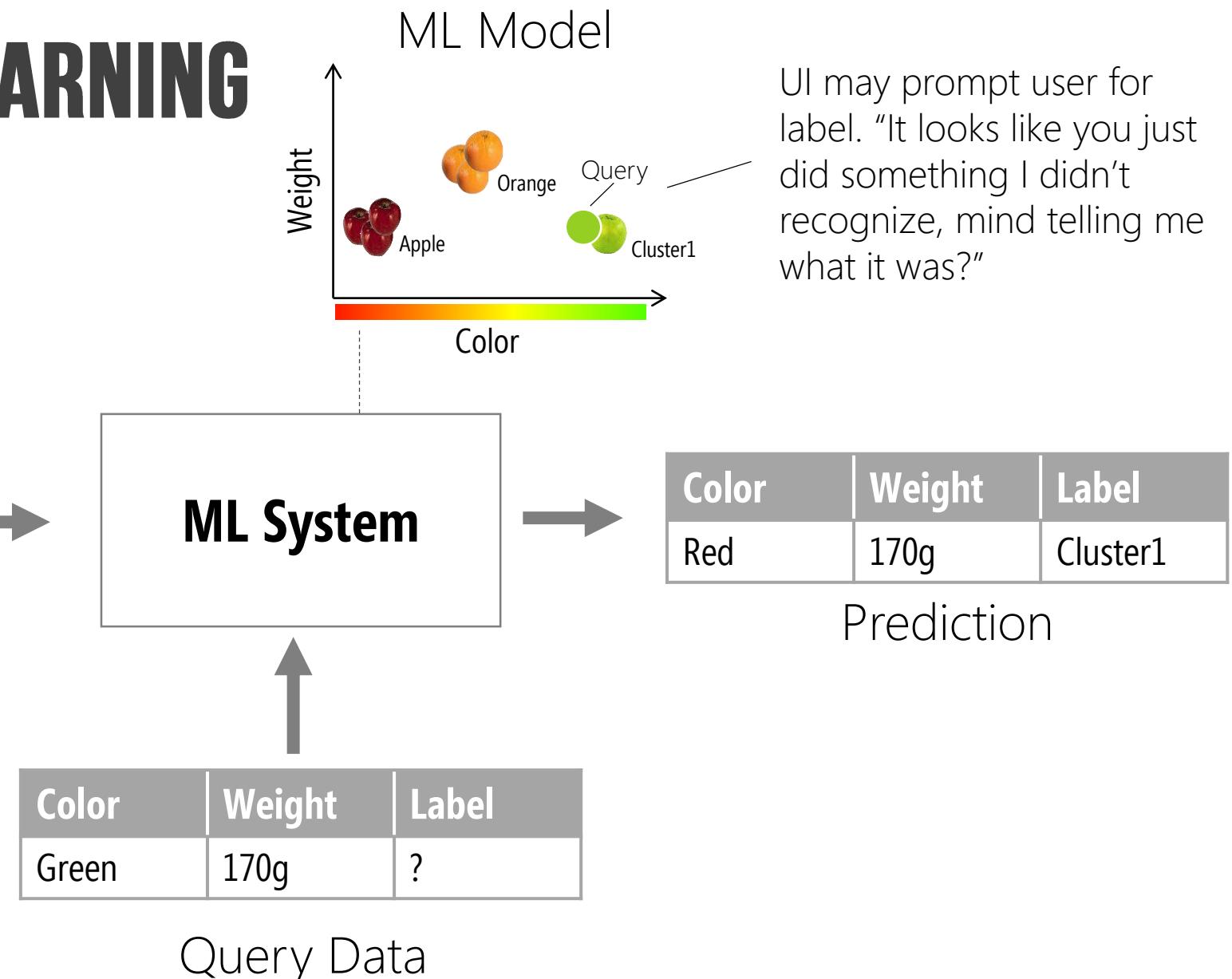
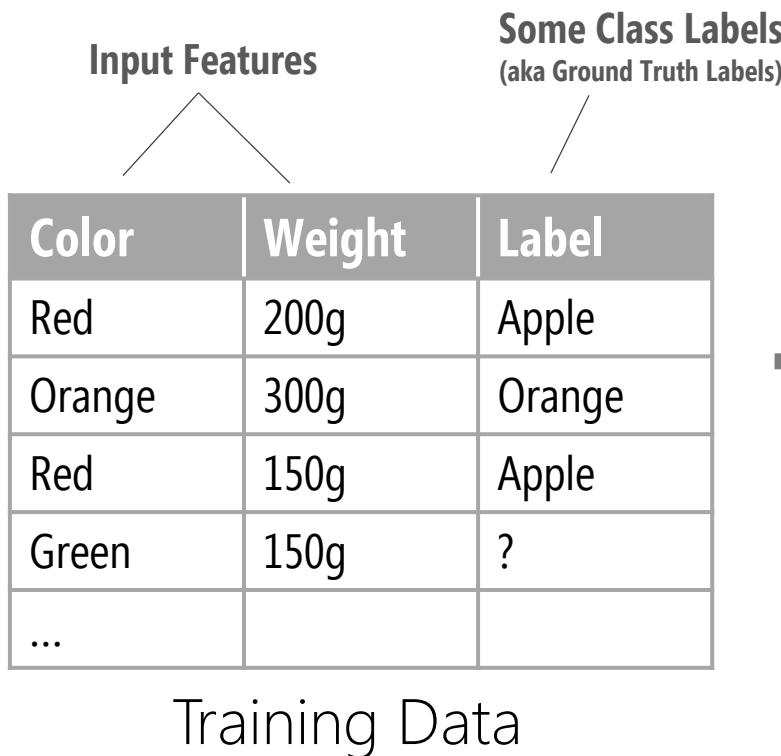
SEMISUPERVISED LEARNING

Input Features		Some Class Labels (aka Ground Truth Labels)
Color	Weight	Label
Red	200g	Apple
Orange	300g	Orange
Red	150g	Apple
Green	150g	?
...		

Training Data



SEMISUPERVISED LEARNING



HOW TO CHOOSE A MODEL?

Your decision should be based on the nature of the dataset (e.g., images, videos, text, continuous time-series)

And the question you're trying to answer (e.g., What is likely to happen next? What are the anomalies in this dataset? What objects are in this image? What gesture is the user performing?)

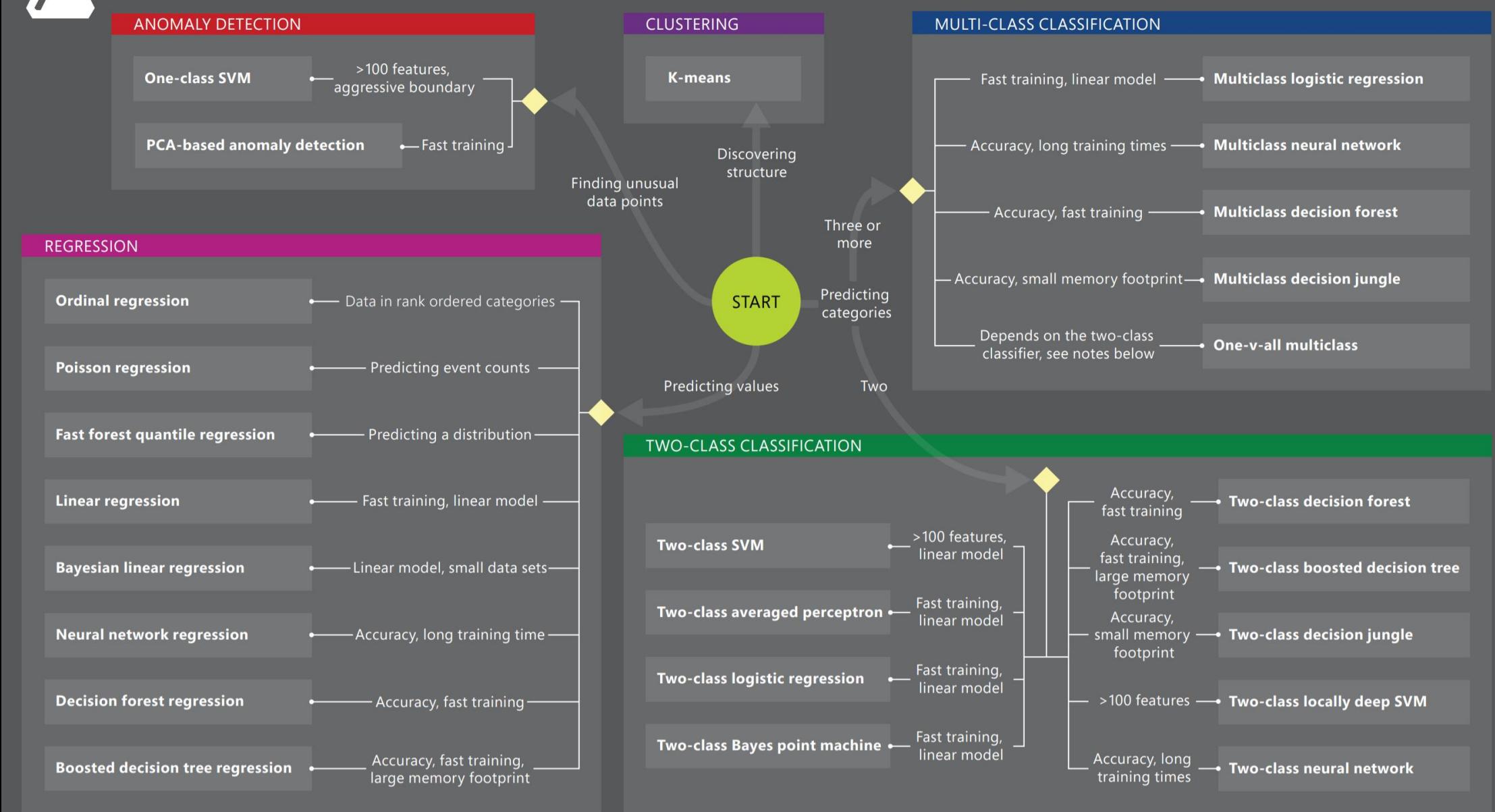
Often, to answer these questions, we explore the literature to see how others have solved similar problems and adapt...

For A2, we strongly recommend an SVM



Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.



MACHINE LEARNING APPROACHES: ORGANIZED BY TYPE



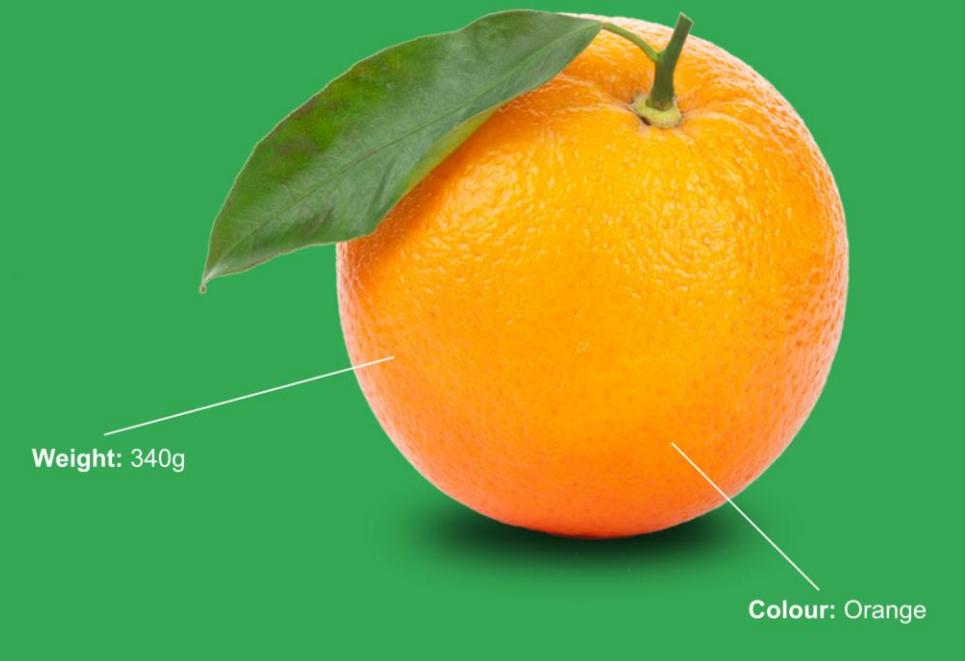
How to select our **input features?**

FEATURES

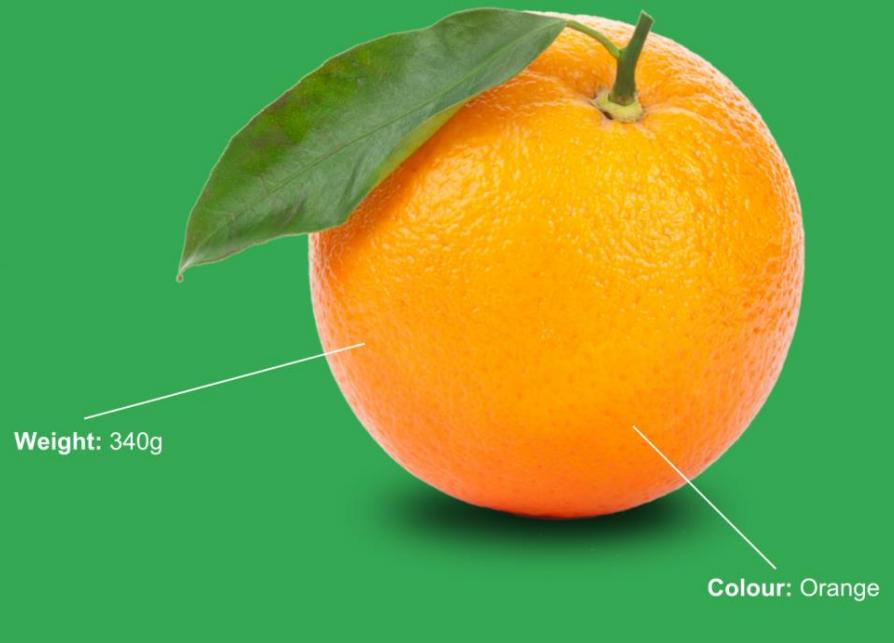
Features are used to **train** the ML system.

You want to **pick features that you think help discriminate** your dataset.

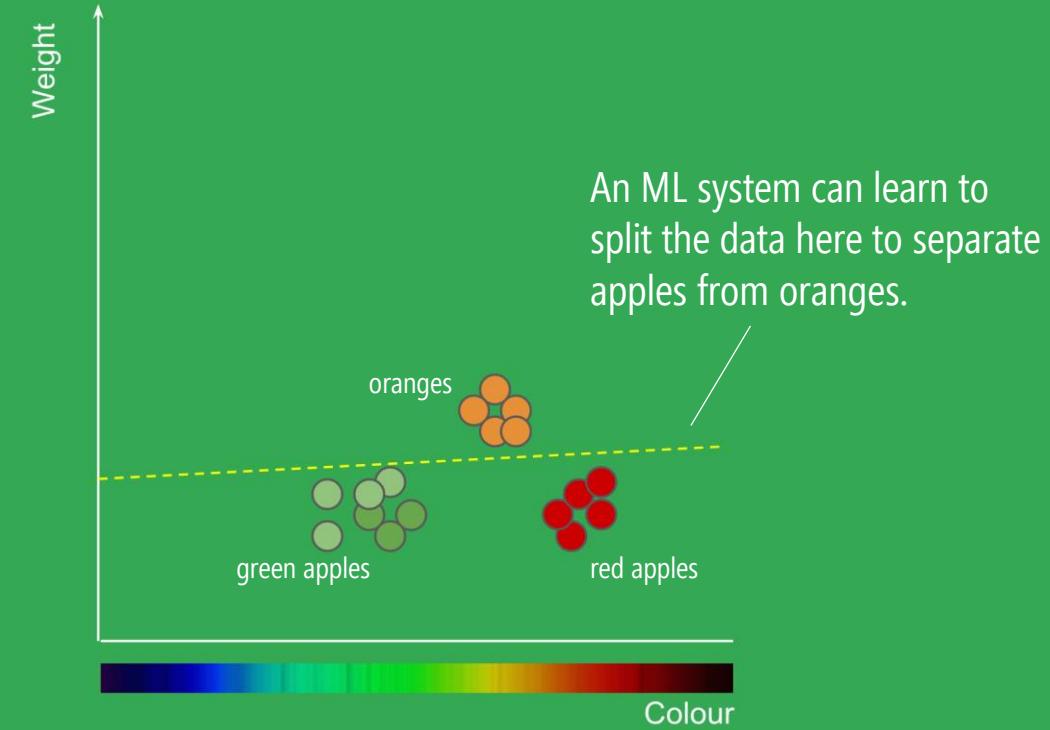
With fruit, for example, you might try **weight** and **color**. **Two features**, which means your featureset is **two-dimensional**



SELECTING FEATURES

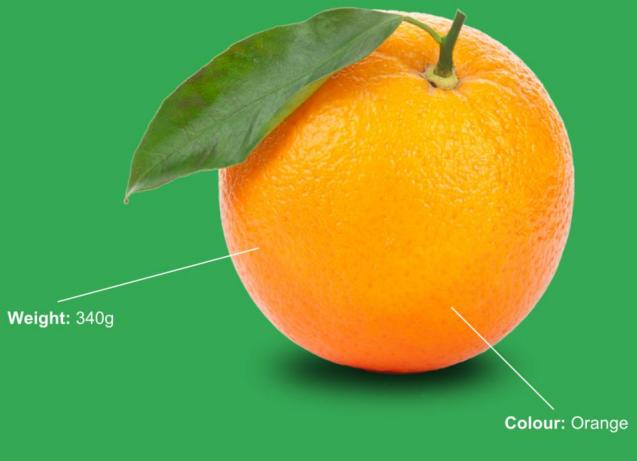


Picking features. If we were building a classifier to automatically distinguish fruits, we could setup a 2D feature vector: weight and color



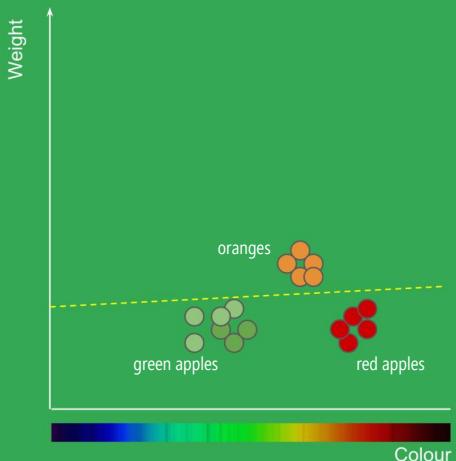
Brainstorm, visualize, investigate, & iterate your feature space. When exploring features, it's a good idea to plot them on a graph (when possible) to develop intuition & assess discriminability.

SELECTING FEATURES

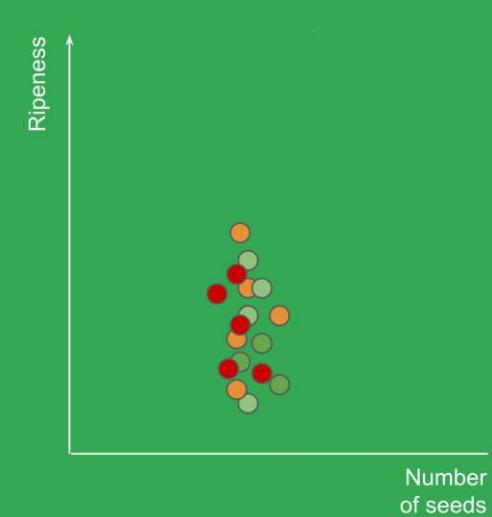


Picking features. if we were building a classifier to automatically distinguish fruits, we could setup a 2D feature vector: weight and color

Brainstorm, visualize, investigate, & iterate your feature space. When exploring features, it's a good idea to plot them on a graph (when possible) to develop intuition and understand their discriminability.



Ripeness and # of seeds are poor features for distinguishing oranges and apples



Choose your features carefully. Some features may not help you distinguish data. For example, ripeness and number of seeds are bad features.

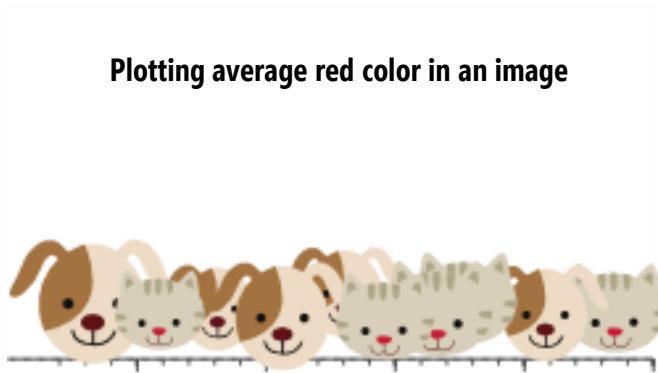
SELECTING FEATURES: HIGHER DIMENSIONALITY

Generally, ML systems are trained with (super!) high dimensional feature vectors that are hard (impossible) to plot and develop an intuition for.

We can plot maybe ~4 dimensions at a time (x, y, z-axis + some visual attribute of the marker).

SELECTING FEATURES: HIGHER DIMENSIONALITY

Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



1D Feature Vector

Plotting average red color in image.

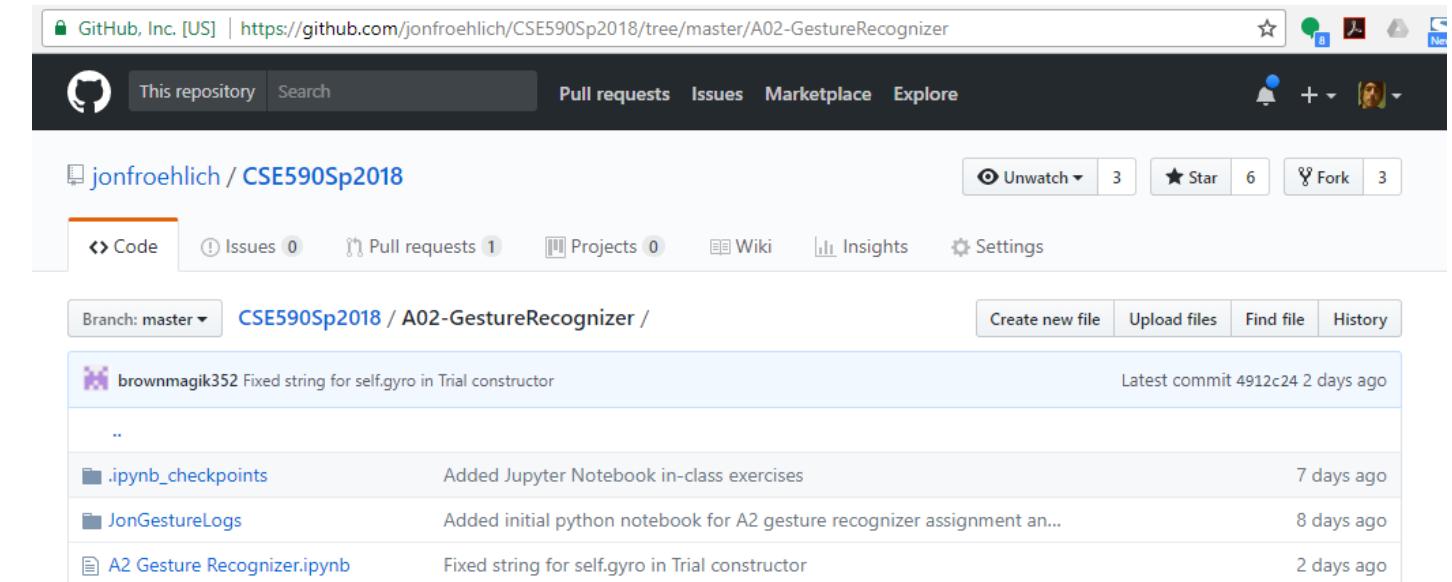
Notice how a single feature does not result in a perfect separation of our training data.

EXERCISE: BRAINSTORM & PLOT FEATURES

We are going to be doing exercises in Jupyter Notebook again.

Please clone/pull a fresh copy of the Android Gesture Recognizer.ipynb

Scroll to “Model-Based ML Exercises”



Feel free to work with a partner or alone!

EXERCISE: BRAINSTORM & PLOT FEATURES

What's the best 1D feature vector you can think of for recognizing gestures?

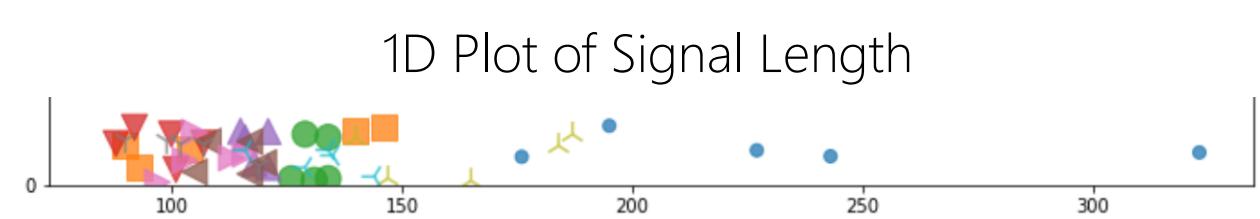
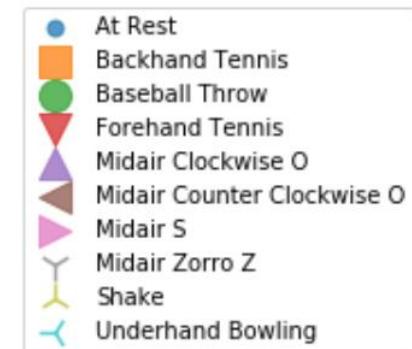
Length of signal

Max accel magnitude

Fundamental frequency of FFT

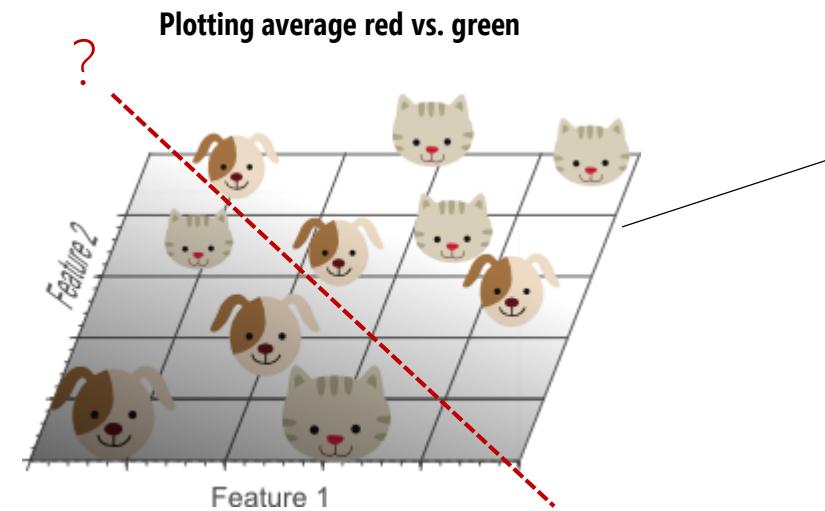
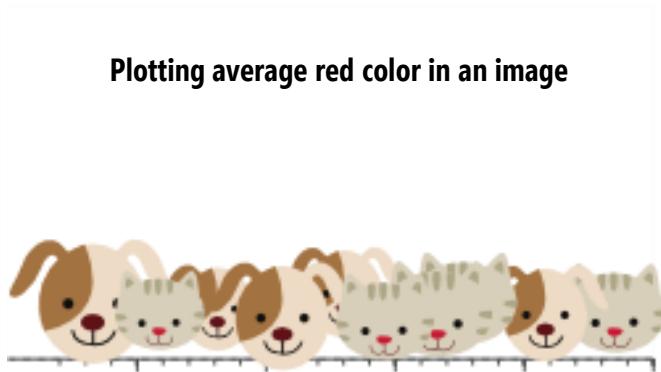
...

Plot your ideas on individual scatterplot graphs; color the gesture classes uniquely (you'll make a 2D scatter graph where the y-values are just small random perturbations to view your data better). We want to see **classes grouped together** and a **clear separation between them**.



SELECTING FEATURES: HIGHER DIMENSIONALITY

Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



1D Feature Vector

Plotting average red color in image.
Notice how a single feature does not result in a perfect separation of our training data.

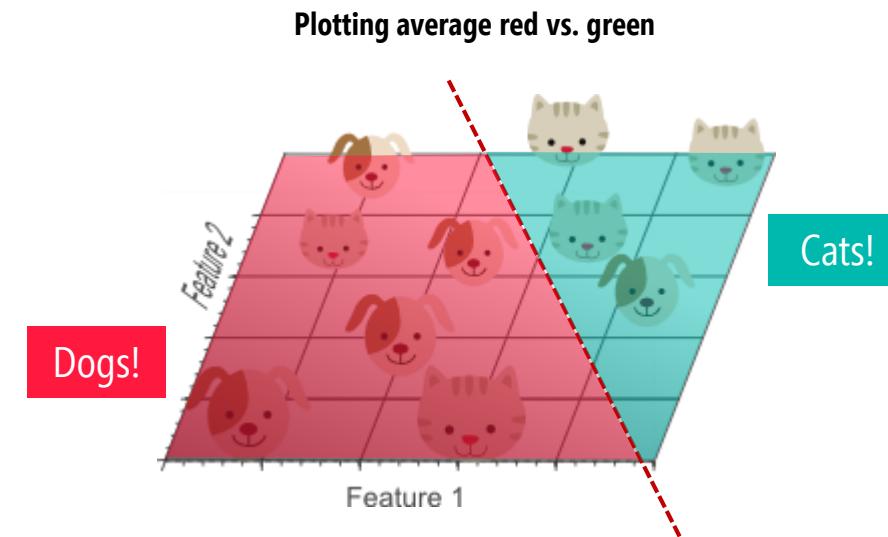
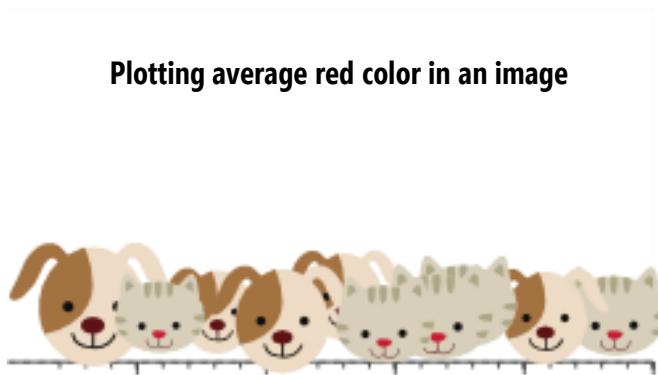
2D Feature Vector

Plotting average red and green colors in image. Now we begin to see some more separation but not perfect.

Is there any straight line that we can draw through our data that separates cats from dogs?

SELECTING FEATURES: HIGHER DIMENSIONALITY

Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



Overall Accuracy: 7/10 (70%)

True	Cat	Predicted		Overall Accuracy
		Cat	Dog	
Cat	3	2		3/5 (60%)
Dog	1	4		4/5 (80%)

1D Feature Vector

Plotting average red color in image.
Notice how a single feature does not result in a perfect separation of our training data.

2D Feature Vector

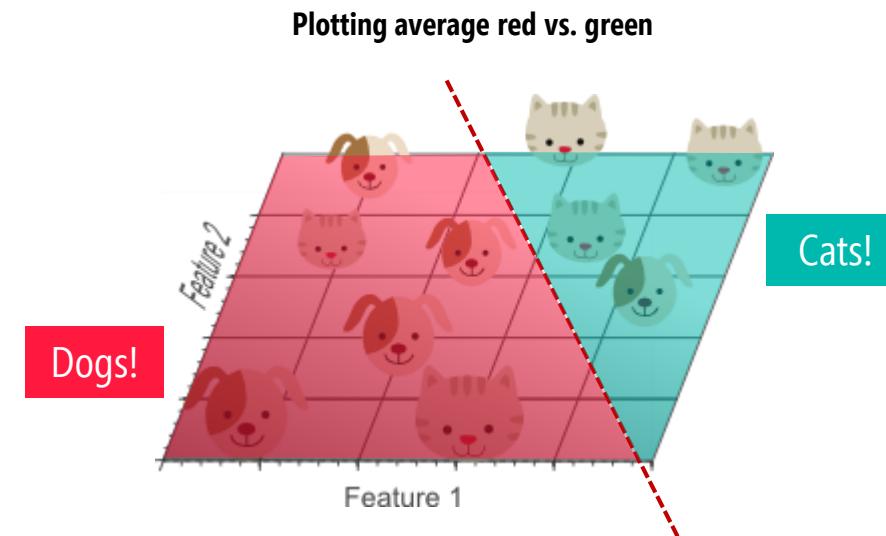
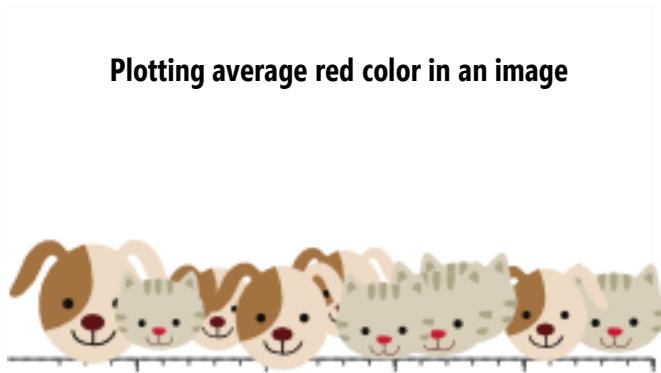
Plotting average red and green colors in image. Now we begin to see some more separation but not perfect.

Results

Overall accuracy is 70% with more cats being misclassified as dogs.

SELECTING FEATURES: HIGHER DIMENSIONALITY

Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



1D Feature Vector

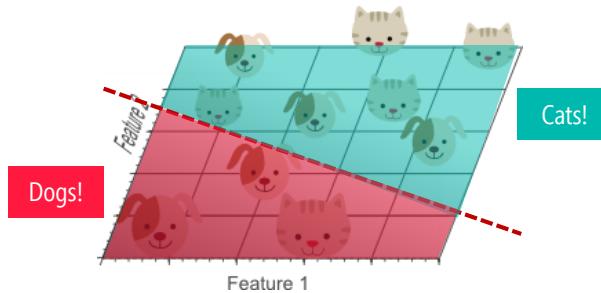
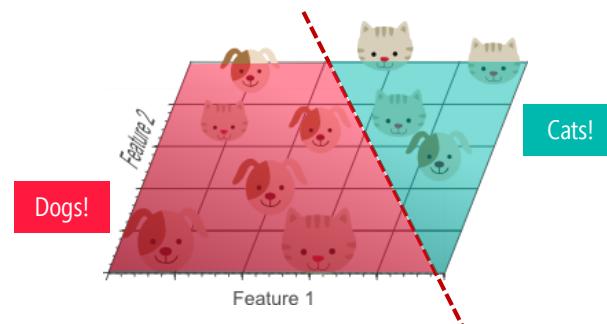
Plotting average red color in image.
Notice how a single feature does not result in a perfect separation of our training data.

2D Feature Vector

Plotting average red and green colors in image. Now we begin to see some more separation but not perfect.

SELECTING FEATURES: HIGHER DIMENSIONALITY

And, of course, you can find different splits of the data. But we want to minimize error—an optimization problem—to find best split.



Overall Accuracy: 7/10 (70%)

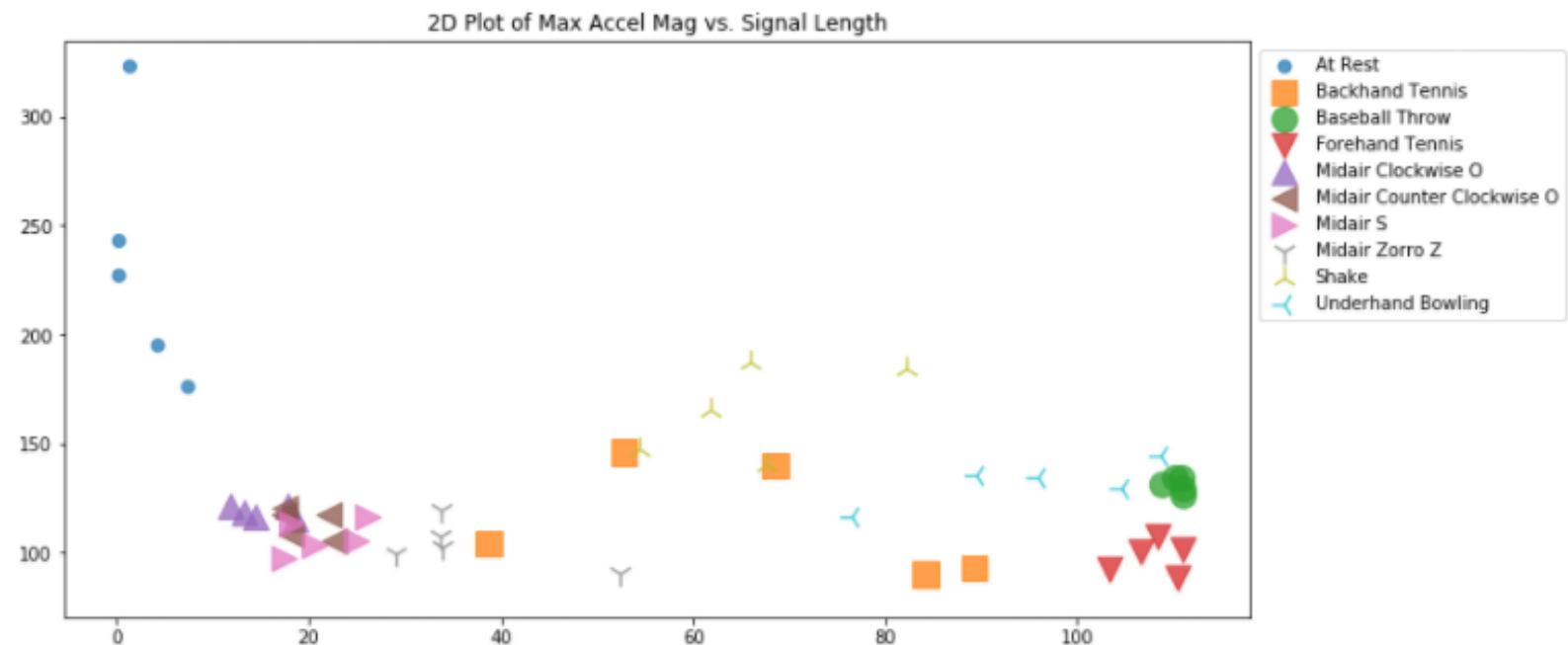
True			Overall Accuracy: 7/10 (70%)
	Cat	Dog	
Cat	3	2	3/5 (60%)
Dog	1	4	4/5 (80%)
	Cat	Dog	
	Predicted		

Overall Accuracy: 6/10 (60%)

True			Overall Accuracy: 6/10 (60%)
	Cat	Dog	
Cat	4	1	4/5 (80%)
Dog	3	2	2/5 (40%)
	Cat	Dog	
	Predicted		

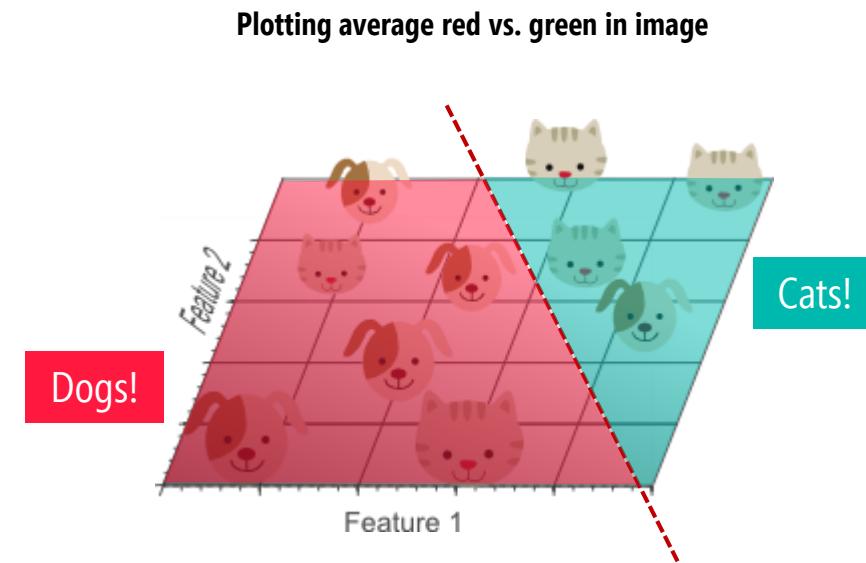
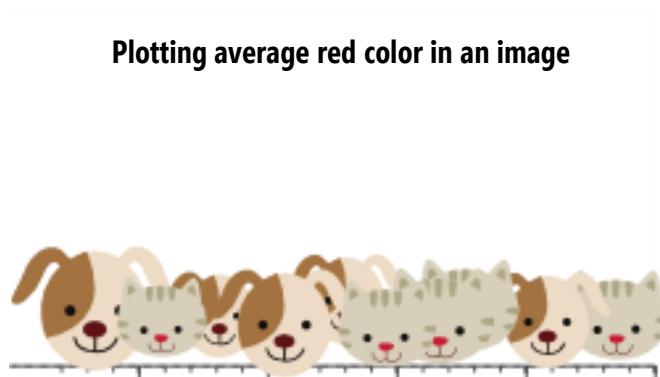
EXERCISE: PLOT FEATURES IN 2D SPACE

Now in a new cell,
start creating 2D
scatter plots by
combining your 1D
features (assign one
feature to x, another
to y and graph!).



SELECTING FEATURES: HIGHER DIMENSIONALITY

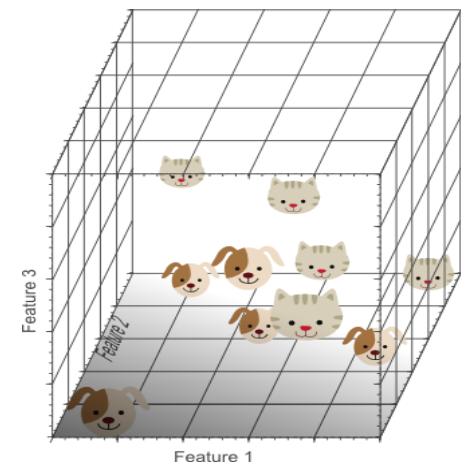
Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



1D Feature Vector

Plotting average red color in image.
Notice how a single feature does not result in a perfect separation of our training data.

Plotting average red, green, blue in an image



2D Feature Vector

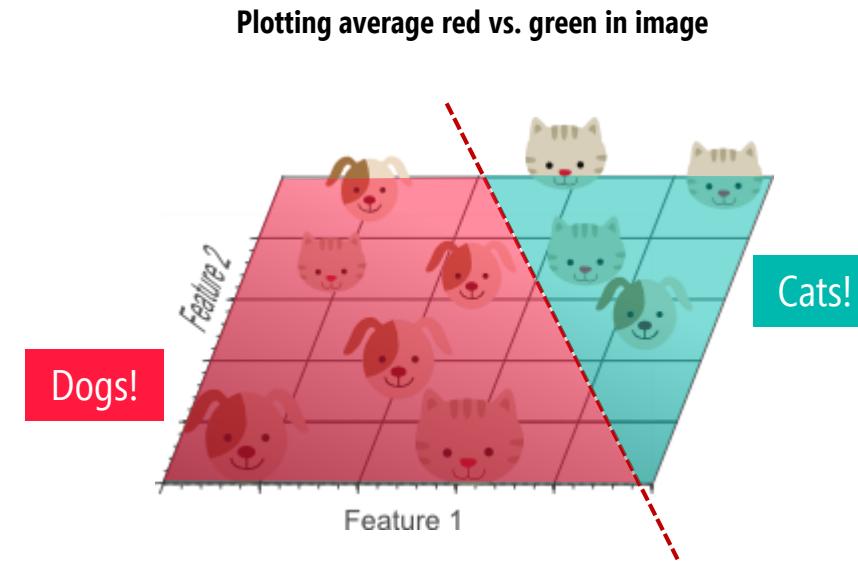
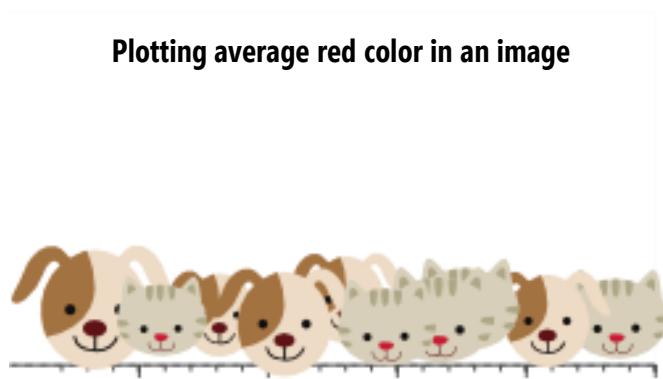
Plotting average red and green colors in image. Now we begin to see some more separation but not perfect.

3D Feature Vector

Using a 3D-input feature vector (R,G,B), we begin to see more separation. However, this is harder to see in the graph.

SELECTING FEATURES: HIGHER DIMENSIONALITY

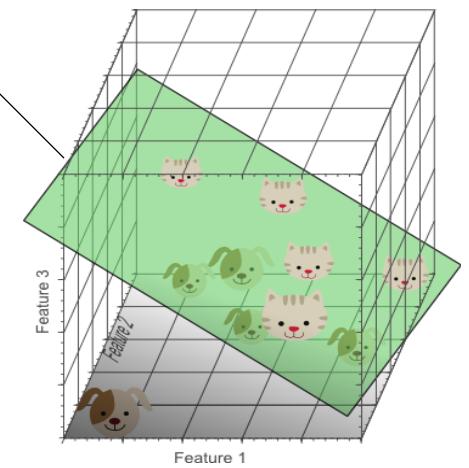
Let's look at another example: this time trying to separate cats and dogs using average pixel color in image.



1D Feature Vector

Plotting average red color in image. Notice how a single feature does not result in a perfect separation of our training data.

Now, in a three-dimensional feature space, we can find a plane that perfectly splits dogs from cats.



2D Feature Vector

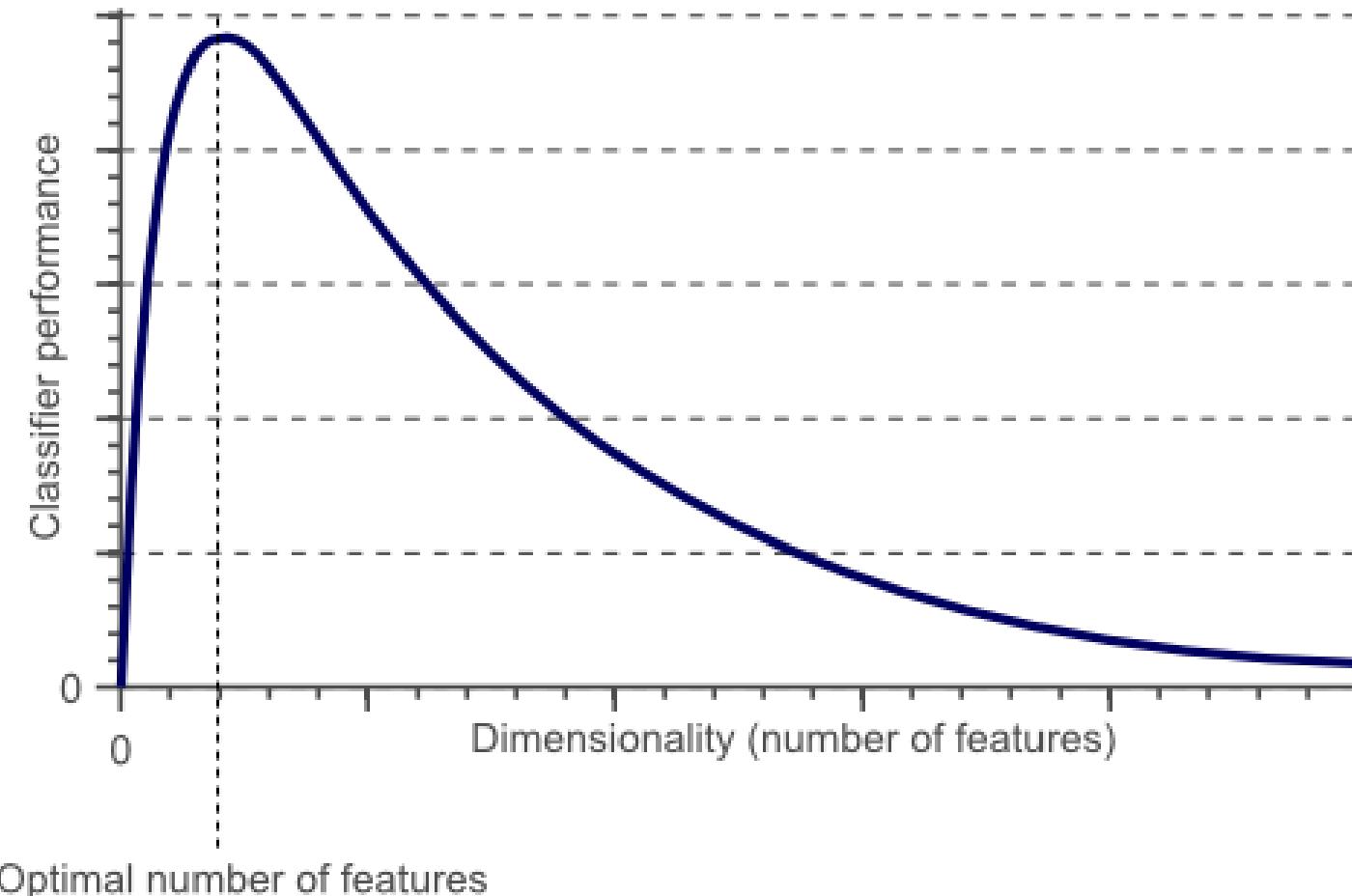
Plotting average red and green colors in image. Now we begin to see some more separation but not perfect.

3D Feature Vector

Using a 3D-input feature vector (R,G,B), we begin to see more separation. (Of course, plotting in 3D often requires interactivity to tilt, rotate, explore data)

THE CURSE OF DIMENSIONALITY: OVERFITTING

You might think that higher dimensionality always results in a better classifier...



DIMENSIONALITY REDUCTION

Dimensionality reduction techniques transform feature vectors from high-dimensional space to fewer dimensions.

The most common technique is Principal Component Analysis, which performs a linear mapping to a lower-dimensional space such that the variance in the low-dimensional representation is maximized

e.g., You might find that those using HMMs in gesture recognition systems reduce high dimensional data using knn clustering or PCA

HOW DO WE SELECT FEATURES?

Based on **domain expertise** (e.g., a medical doctor helping select features for automatic melanoma recognition in images)

Based on **intuitions, experimentations, and plotting the data** (e.g., like we did in the in-class exercises)

Feature selection algorithms!

FEATURE SELECTION ALGORITHMS

Given a giant list of input features, a feature selection algorithm can be seen as a search technique for finding the best combination of feature subsets

Simplest approach: test each possible subset of features finding the one that minimizes the error rate. This is computationally intractable.

Sci-kit learn and other toolkits have support for this. Note: we are not expecting you to do this for A2

FEATURE SELECTION

FEATURE SELECTION EXAMPLE

Classification Features for Touchscreen Typing

The above subsections defined three dimensions of adaptation for personalized touchscreen keyboards. Here we focus on *Dimension 1*, key-press classification models, and consider what features may be most useful for classification (Figure 4). Ultimately, the set of possible features will be determined by the input hardware used (*e.g.*, a vision system may provide more features than a capacitive touchscreen). While our focus is on ten-finger typing, many features listed here are also applicable to smaller devices. Our goal is to highlight the breadth of features that exist for personalizing keyboards as well as to introduce the particular features that we later test.

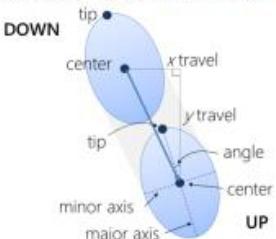


Figure 4. A subset of classification features based on a single key press event. This subset was tested in our simulations.

Touch location features are based on finger-down and finger-up locations. These features can include: (1) *finger center* (x,y) and (2) *fingertip* (x,y) locations. On large touchscreens, the keyboard can move, so these coordinates must be relative to some point(s) of reference. For adaptation to the left and right hands, *two* relative center locations can be calculated, *e.g.*, one as an offset from the F key and one as an offset from the J key (the *home* keys). Absolute values can be used for smaller screens.

Touch area features describe the size and angle of the finger on the screen; these features greatly depend on the hardware device. We used a Microsoft Surface, which provides ellipses fitted to the touch area and allows for (3) *major axis*, (4) *minor axis* and (5) *orientation*.

Movement features describe the movement of the finger during the key-press action. Different keys may result in different movement signatures based on the finger used and how the user reaches for that key. Specific features include (6) *travel* and (7) *elapsed time* between finger down and up. Travel can be calculated as the distance in x and y directions for either the fingertip or the finger center.

Features relative to previous key presses incorporate knowledge about previous key presses. These include (8) *inter-key-press timing*, or the elapsed time between key presses, and (9) *(x,y) distance offset from previous key*. The latter may be useful if, for example, the current key-press location is affected by whether the last key was to its left or right. A simulation study by Rashid and Smith [28] showed low accuracies for a touch model based only on this feature.

Hand or arm features are based on the position of the hand and/or arm. Again, these features are highly dependent on the underlying touchscreen system in that only certain systems can provide this information.

To select which subset of classification features to include in the J48 classification model, we used a wrapper-based feature selection technique [18]. That is, we incrementally added features to the model until no improvement in performance occurred. This process resulted in all features in Table 1 marked with a *. Additionally, we included the finger center y relative to the F and J keys for continuity, since the finger center x was already in the feature set; this inclusion did not reduce classification accuracy. With the

Feature Set	Classification Accuracy
Tip x relative to F and J keys*	53.3% ($SD = 3.3$)
Center x relative to F and J keys*	43.4% ($SD = 6.7$)
Tip y relative to F and J keys*	35.4% ($SD = 2.0$)
Center y relative to F and J keys*	29.7% ($SD = 1.5$)
Angle	28.0% ($SD = 3.9$)
Tip travel in x and y directions*	27.4% ($SD = 1.9$)
Center travel in x and y directions	24.8% ($SD = 1.4$)
Major axis	22.4% ($SD = 1.5$)
Minor axis	19.0% ($SD = 1.9$)
Elapsed time	14.2% ($SD = 1.0$)
Inter-key-press time	11.7% ($SD = 2.3$)
Majority classifier (right space key)	11.8%

* Used in feature set for final implementation.

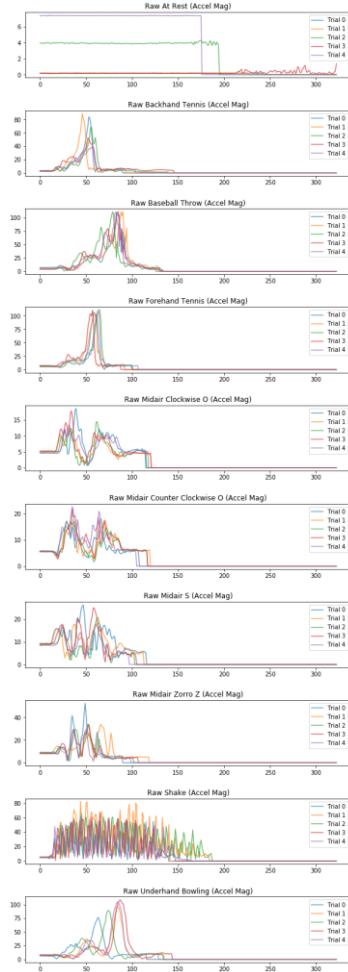
Table 1. Explanatory power of classification features based on the pilot study data: mean 10-fold cross-validation results from a J48 classifier trained on individual features or pairs of related features. The majority classifier is the accuracy rate if we always classify as the most frequent key (right space key).

Where do we get our labeled training data?

THINKING ABOUT LABELING GESTURES IN A2

I labeled pre-segmented streams. What about a more “realistic” approach where you label continuous streams (*i.e.*, every datapoint has a label)

Name	Date modified
At Rest_Accelerometer_1523388944961.csv	4/10/2018 12:37 PM
At Rest_Accelerometer_1523388953684.csv	4/10/2018 12:37 PM
At Rest_Accelerometer_1523388964725.csv	4/10/2018 12:37 PM
At Rest_Accelerometer_1523388984213.csv	4/10/2018 12:37 PM
At Rest_Accelerometer_1523389002850.csv	4/10/2018 12:37 PM
At Rest_Gyroscope_1523388944961.csv	4/10/2018 12:37 PM
At Rest_Gyroscope_1523388953684.csv	4/10/2018 12:37 PM
At Rest_Gyroscope_1523388964725.csv	4/10/2018 12:37 PM
At Rest_Gyroscope_1523388984213.csv	4/10/2018 12:37 PM
At Rest_Gyroscope_1523389002850.csv	4/10/2018 12:37 PM
Backhand Tennis_Accelerometer_1523343191085.csv	4/9/2018 11:54 PM
Backhand Tennis_Accelerometer_1523343197400.csv	4/9/2018 11:54 PM
Backhand Tennis_Accelerometer_1523343205168.csv	4/9/2018 11:54 PM
Backhand Tennis_Accelerometer_1523343214494.csv	4/9/2018 11:54 PM
Backhand Tennis_Accelerometer_1523343221881.csv	4/9/2018 11:54 PM
Backhand Tennis_Gyroscope_1523343191085.csv	4/9/2018 11:54 PM
Backhand Tennis_Gyroscope_1523343197400.csv	4/9/2018 11:54 PM
Backhand Tennis_Gyroscope_1523343205168.csv	4/9/2018 11:54 PM
Backhand Tennis_Gyroscope_1523343214494.csv	4/9/2018 11:54 PM
Backhand Tennis_Gyroscope_1523343221881.csv	4/9/2018 11:54 PM
Baseball Throw_Accelerometer_1523388888168.csv	4/10/2018 12:37 PM
Baseball Throw_Accelerometer_152338895884.csv	4/10/2018 12:37 PM
Baseball Throw_Accelerometer_1523388904126.csv	4/10/2018 12:37 PM
Baseball Throw_Accelerometer_1523388911971.csv	4/10/2018 12:37 PM
Baseball Throw_Accelerometer_1523388920465.csv	4/10/2018 12:37 PM
Baseball Throw_Gyroscope_1523388888168.csv	4/10/2018 12:37 PM
Baseball Throw_Gyroscope_152338895884.csv	4/10/2018 12:37 PM
Baseball Throw_Gyroscope_1523388904126.csv	4/10/2018 12:37 PM
Baseball Throw_Gyroscope_1523388911971.csv	4/10/2018 12:37 PM
Baseball Throw_Gyroscope_1523388920465.csv	4/10/2018 12:37 PM



Gestures are pre-segmented

Labels are at gesture level
(rather than data point level)

GETTING TRAINING DATA

HOW SHOULD WE COLLECT & LABEL CONTINUOUS SENSOR STREAM DATA?

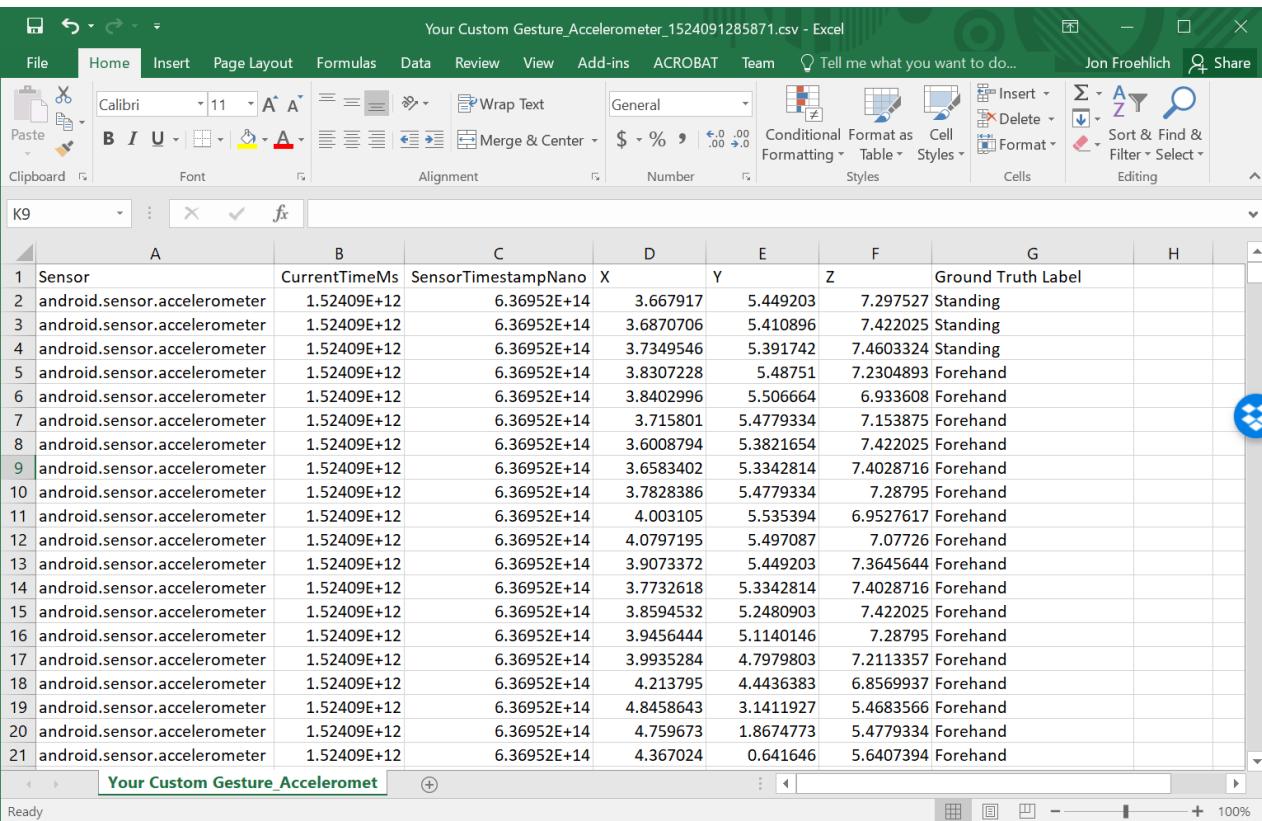


Accelerometer sensor stream when multiple gestures performed over time

GETTING TRAINING DATA

HOW SHOULD WE COLLECT & LABEL CONTINUOUS SENSOR STREAM DATA?

Every timestamped record (i.e., row) in gesture log is labeled with a ground truth activity label



1	Sensor	CurrentTimeMs	SensorTimestampNano	X	Y	Z	Ground Truth Label
2	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.667917	5.449203	7.297527	Standing
3	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6870706	5.410896	7.422025	Standing
4	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7349546	5.391742	7.4603324	Standing
5	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8307228	5.48751	7.2304893	Forehand
6	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8402996	5.506664	6.933608	Forehand
7	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.715801	5.4779334	7.153875	Forehand
8	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6008794	5.3821654	7.422025	Forehand
9	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6583402	5.3342814	7.4028716	Forehand
10	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7828386	5.4779334	7.28795	Forehand
11	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.003105	5.535394	6.9527617	Forehand
12	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.0797195	5.497087	7.07726	Forehand
13	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.9073372	5.449203	7.3645644	Forehand
14	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7732618	5.3342814	7.4028716	Forehand
15	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8594532	5.2480903	7.422025	Forehand
16	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.9456444	5.1140146	7.28795	Forehand
17	android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.9935284	4.7979803	7.2113357	Forehand
18	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.213795	4.4436383	6.8569937	Forehand
19	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.8458643	3.1411927	5.4683566	Forehand
20	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.759673	1.8674773	5.4779334	Forehand
21	android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.367024	0.641646	5.6407394	Forehand



Think/Pair/Share

How would you design a study and analysis approach to get this granular ground truth data?

GETTING TRAINING DATA

HOW SHOULD WE COLLECT & LABEL CONTINUOUS SENSOR STREAM DATA?

Every timestamped record
(i.e., row) in gesture log is
labeled with a ground truth
activity label



Sensor	CurrentTimeMs	SensorTimestampNano	X	Y	Z	Ground Truth Label
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.667917	5.449203	7.297527	Standing
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6870706	5.410896	7.422025	Standing
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7349546	5.391742	7.4603324	Standing
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8307228	5.48751	7.2304893	Standing
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8402996	5.50664	6.933608	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.715801	5.4779334	7.153875	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6807974	5.3821054	7.422025	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6581342	5.3342814	7.4028716	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7823886	5.4779334	7.28795	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.001105	5.535394	6.9527617	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.091795	5.4849397	6.95276126	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8070742	5.449303	7.3645644	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.7733618	5.3342814	7.4028716	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6591532	5.2489093	7.422025	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.6456444	5.1140146	7.3879357	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	3.8935284	4.7979803	7.2113357	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.211795	4.4436383	6.8569937	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.8459643	3.1411927	5.4683566	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.759673	1.8674773	5.4779334	Forehand
android.sensor.accelerometer	1.52409E+12	6.36952E+14	4.367024	0.641546	5.6407394	Forehand



Some Potential Approaches

Have users perform a set of gestures in a prescribed order (follow a script)

Build an interactive system (like a game) that has user perform gestures in reaction to prompts

Record video of gesture study session and build a tool that aligns video with sensor streams then manually label data in tool

Have users wear motion capture suits and use Vicon motion tracking to automatically measure movement and label data

... other ideas?

GETTING TRAINING DATA

EXAMPLE FROM UBICOMP: HYDROSENSE



Evaluating HydroSense

Recall that HydroSense is a single-sensor based approach for automatically identifying water usage activities using pressure signatures.



Initial Study (UbiComp'09)

Proof-of-concept study used controlled water use trials that were manually labeled *in situ* using a custom tool.



Deployment Study (Pervasive'11)

Five-week real-world deployment study that analyzed real water usage data. Ground truth data was collected using custom embedded hardware kit.

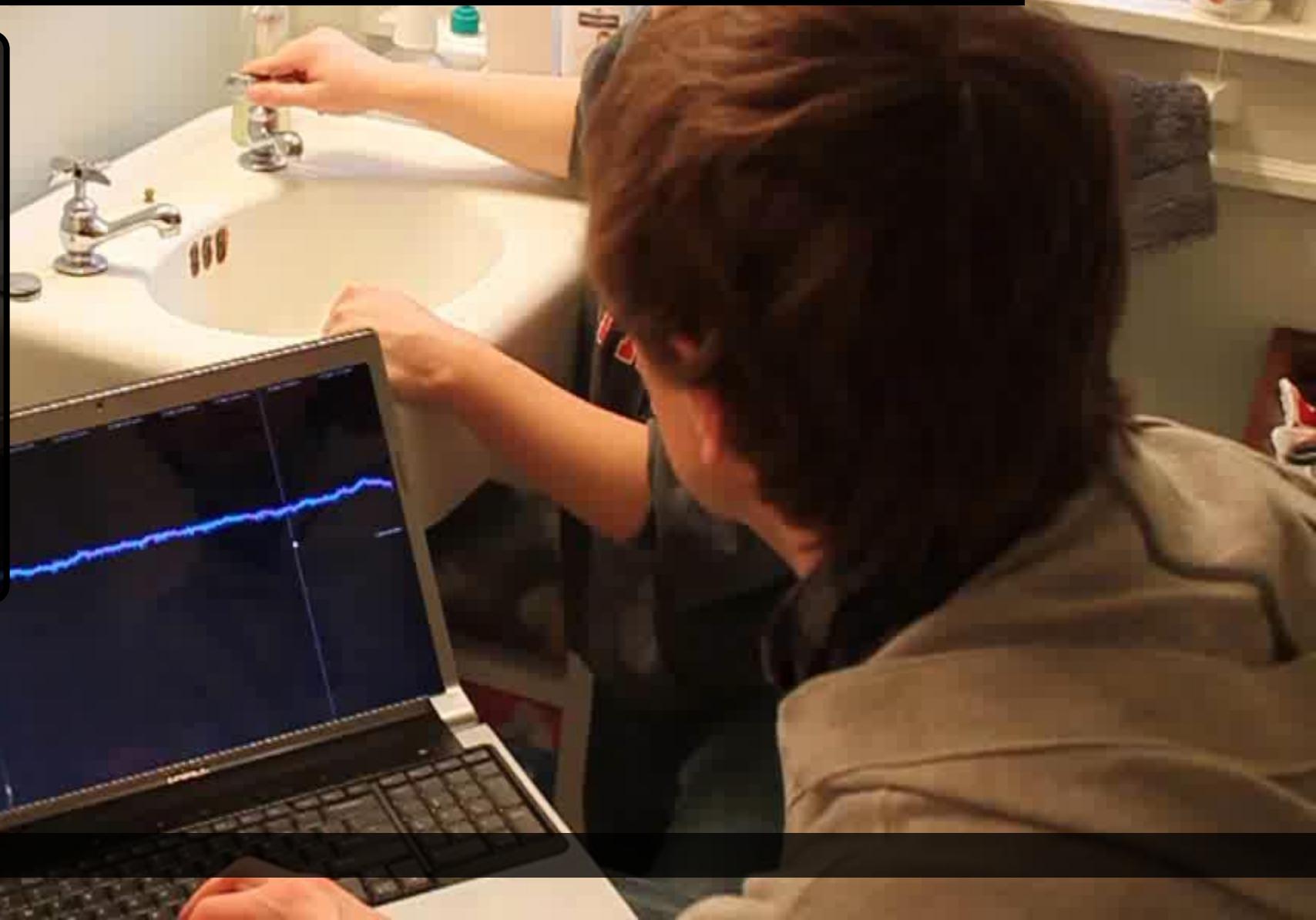
COLLECTING GROUND TRUTH DATA FOR HYDROSENSE 1

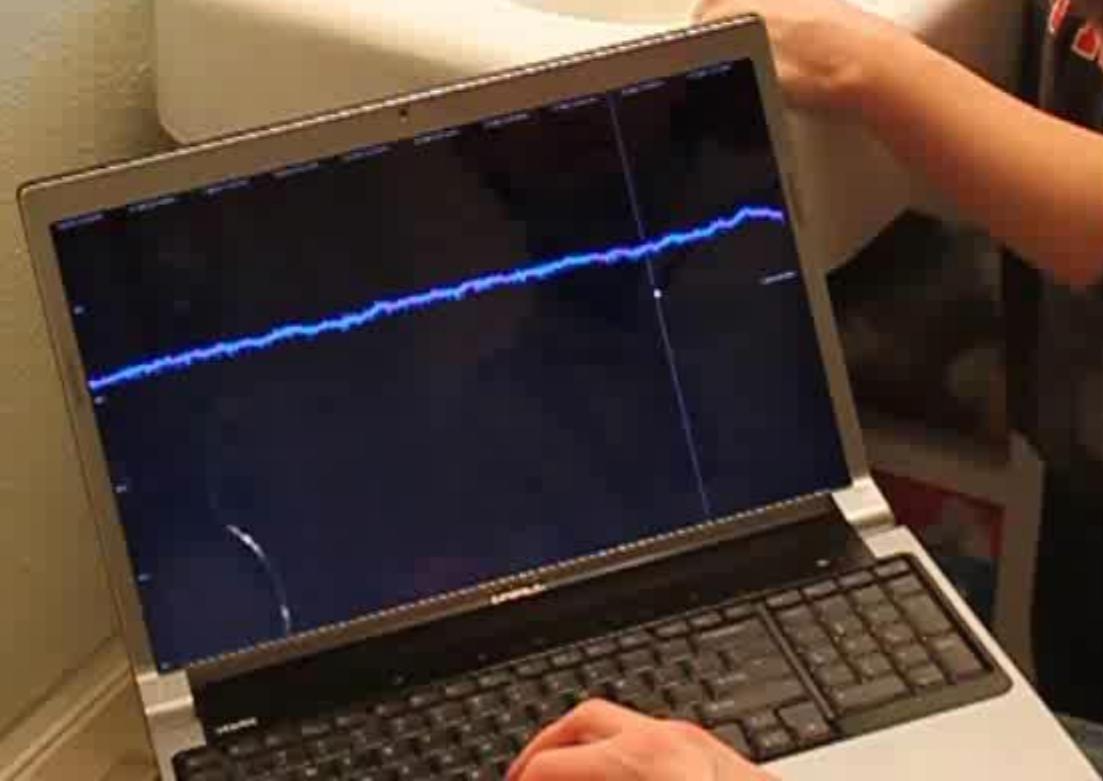
Controlled experiment

- 10 homes
- 2 researchers per site
- 5 trials per valve

Experimental script

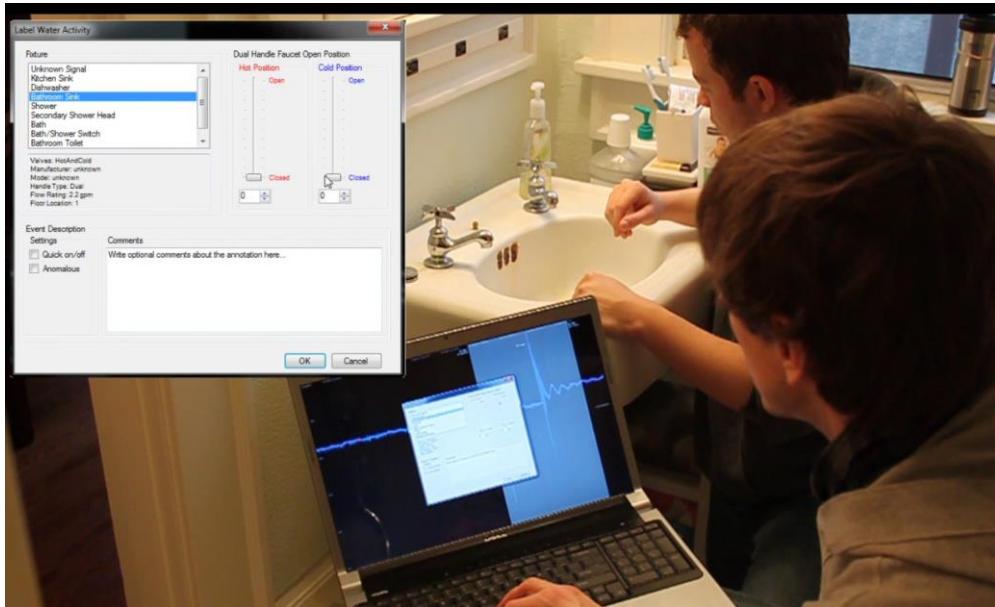
- valve opened full stop
- pause for ~5 seconds
- valve closed





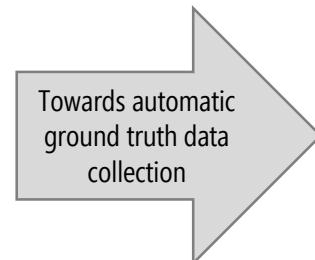
GETTING TRAINING DATA

HOW CAN WE COLLECT GROUND TRUTH LABELS DURING DEPLOYMENT?



How can we automate this?

How can we automate ground truth data collection so that we can get ground truth labels during our field deployment of HydroSense?



Hardware Capabilities

1. Wireless communication
2. Low-power
3. Water resistant

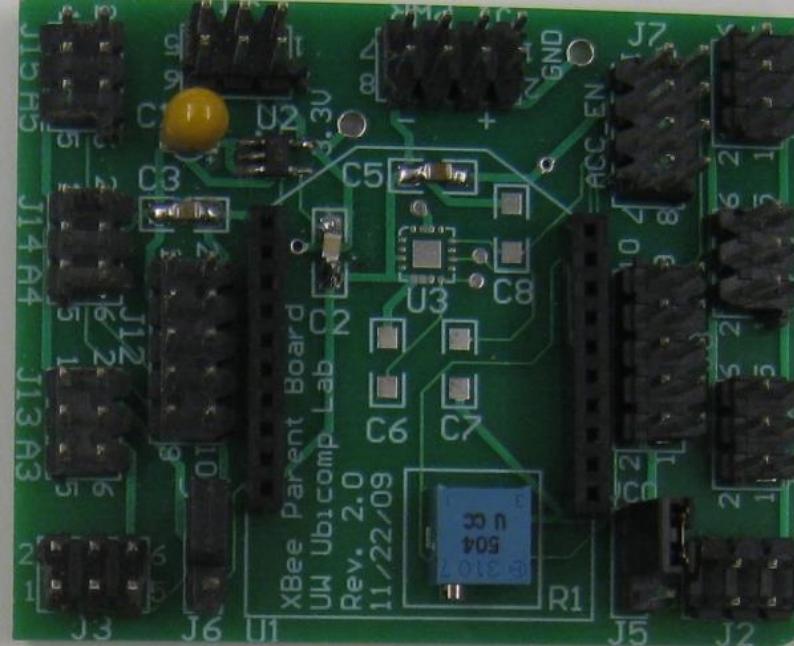
Sensing Capabilities

1. Work across fixtures/appliances
2. Detect opens/closes
3. Discriminate hot/cold/mixed

GETTING TRAINING DATA CUSTOM GROUND TRUTH DATA COLLECTION KIT HYDROSENSE 2.0



xbee wireless modem



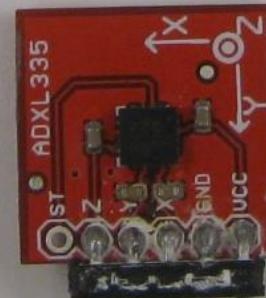
fixture usage sensor board



hall
effect



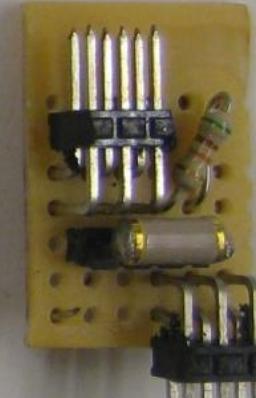
reed
switch



3-axis
accelerometer



unidirectional ball
switch



omnidirectional ball
switch



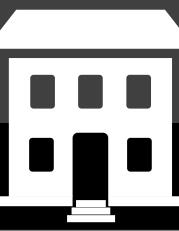
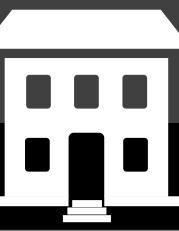
accelerometer



Accelerometer
& Ball Switch
Taped on

GETTING TRAINING DATA

HYDROSENSE 2.0: 5-WEEK FIELD DEPLOYMENT

						Total
residents	2	2	4	2	2	12
size	3000 sqft	750 sqft	1200 sqft	700 sqft	750 sqft	N/A
floors	3	2	2	3 rd flr	6 th flr	N/A
fixtures	17	8	13	8	8	54
valves	28	13	21	13	13	88

GETTING TRAINING DATA

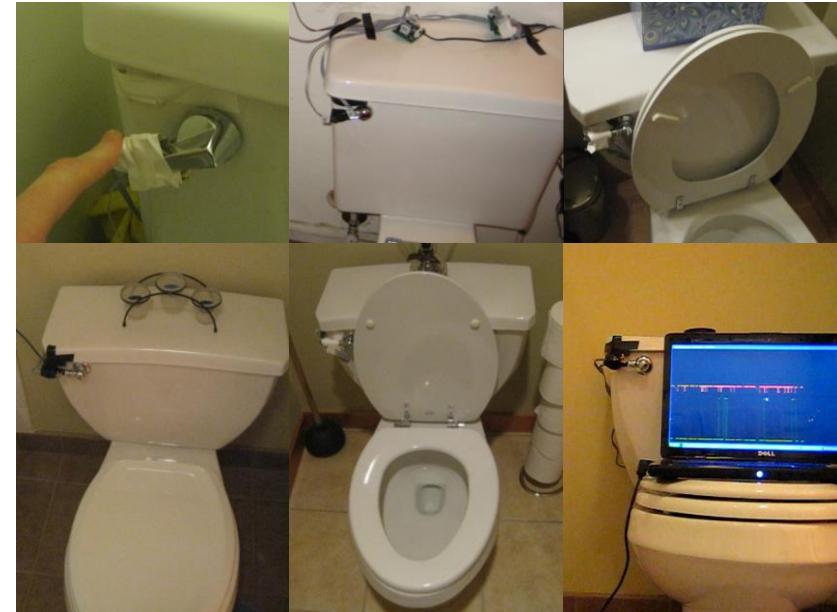
HYDROSENSE 2.0: INSTALLING GROUND TRUTH KIT



Sinks



Showers/Baths



Toilets

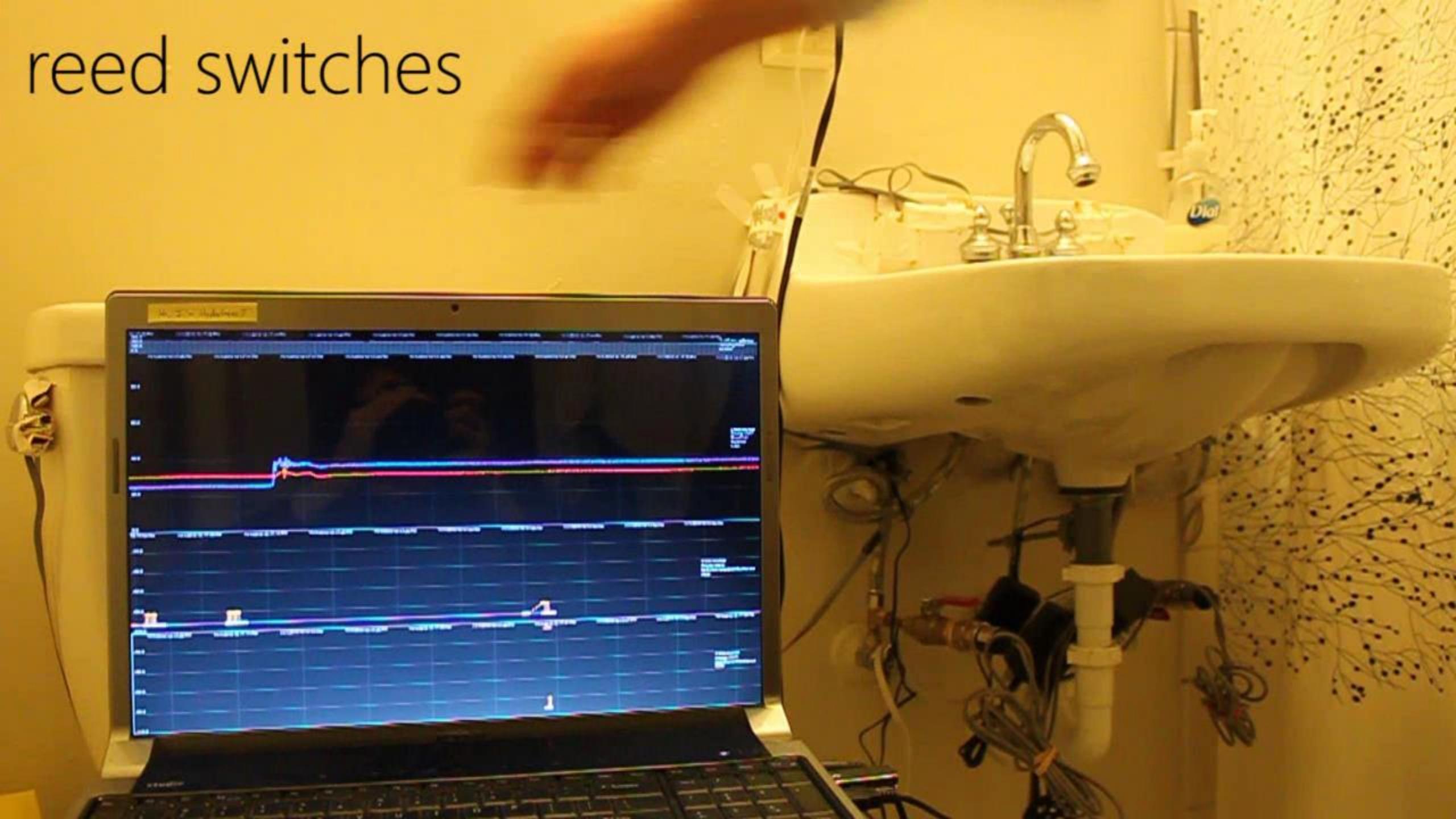


Thermistor to measure laundry cycle (e.g., warm/cold)

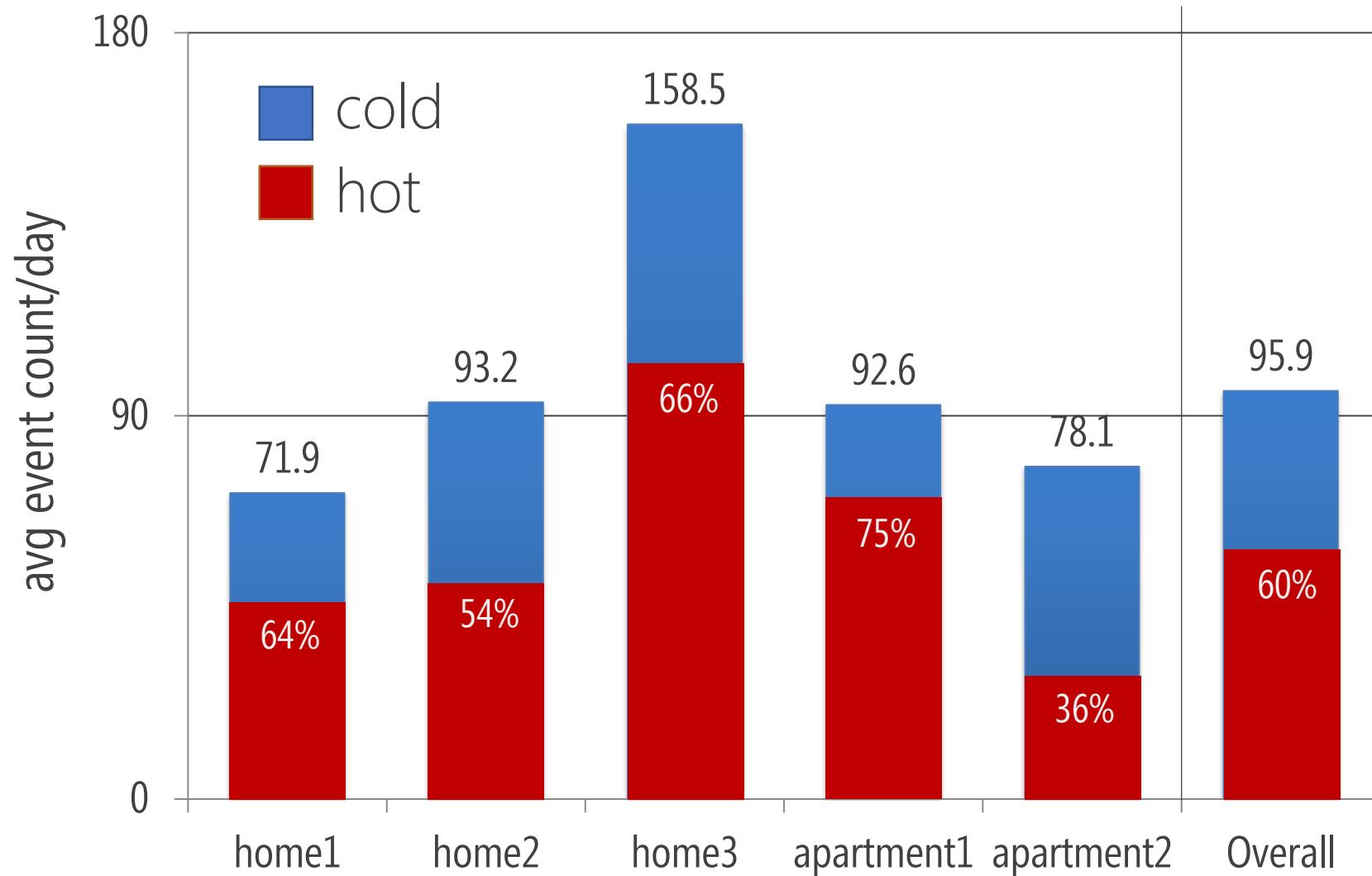




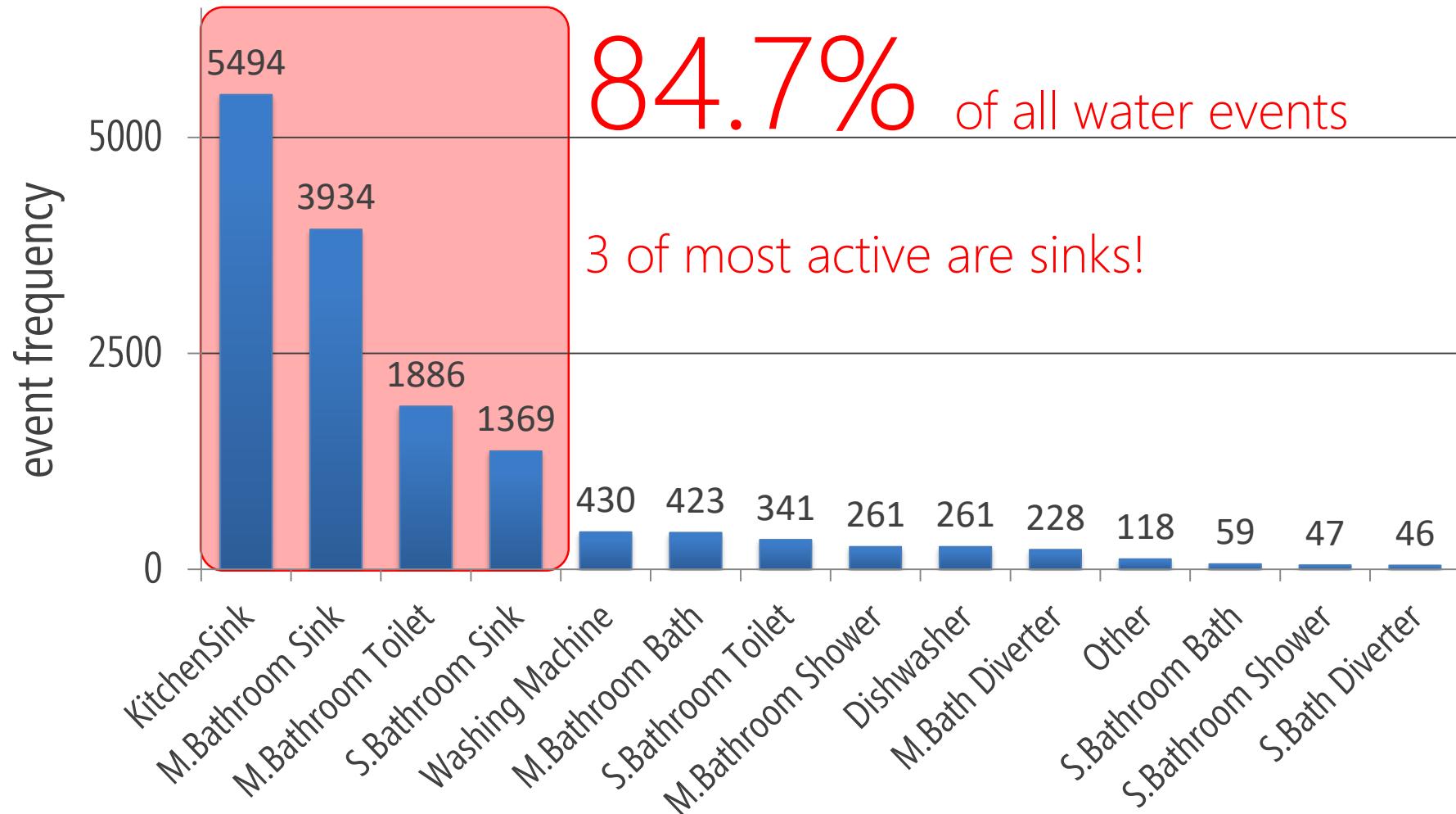
reed switches



AVG WATER EVENTS/DAY



Fixture Activity Frequency



HOW MANY EVENTS WERE COMPOUND?



How many events happened simultaneously with another water usage event?

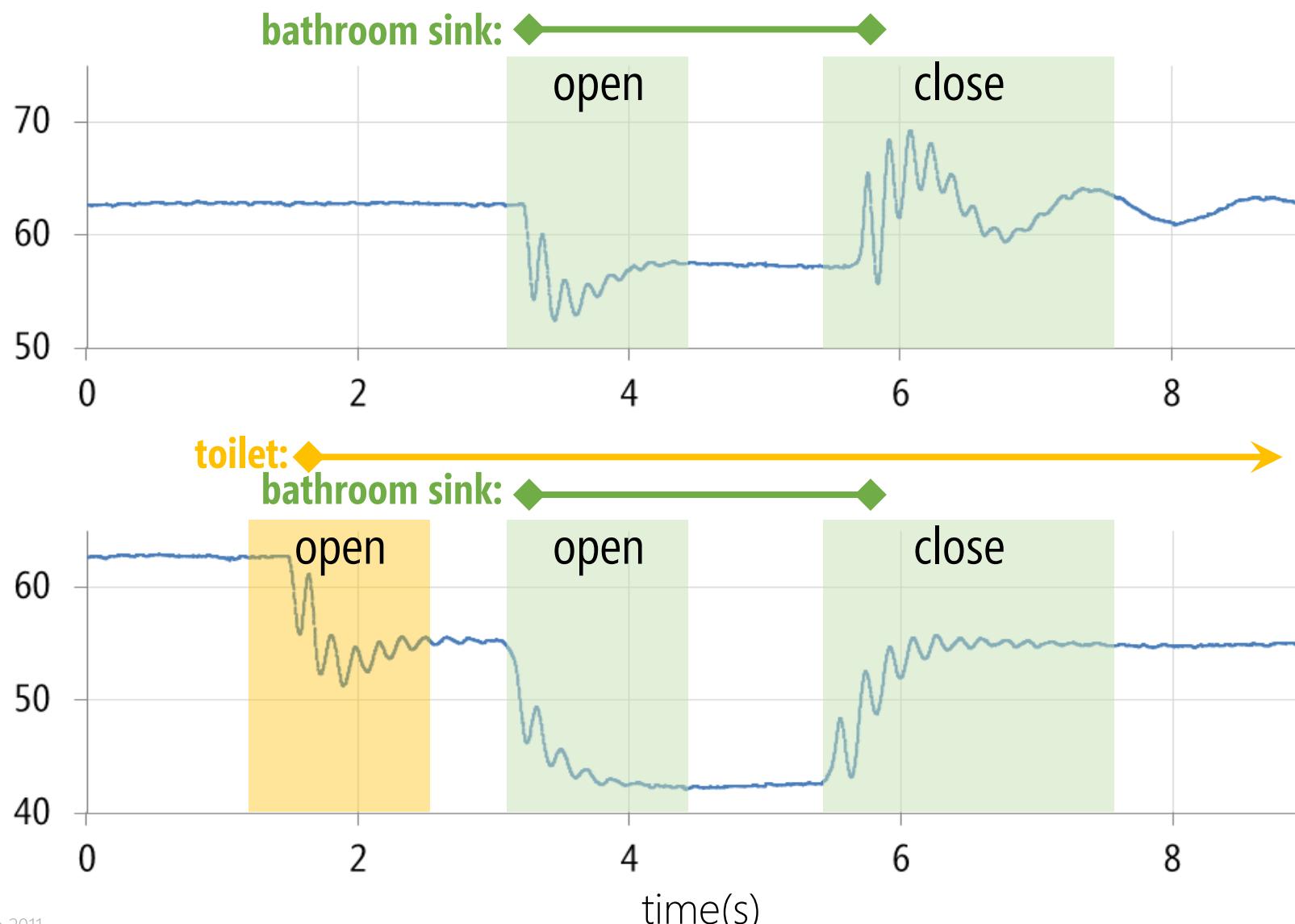
22%

of all **water** events were compound

41.8%

of all **bathroom sink** events were compound

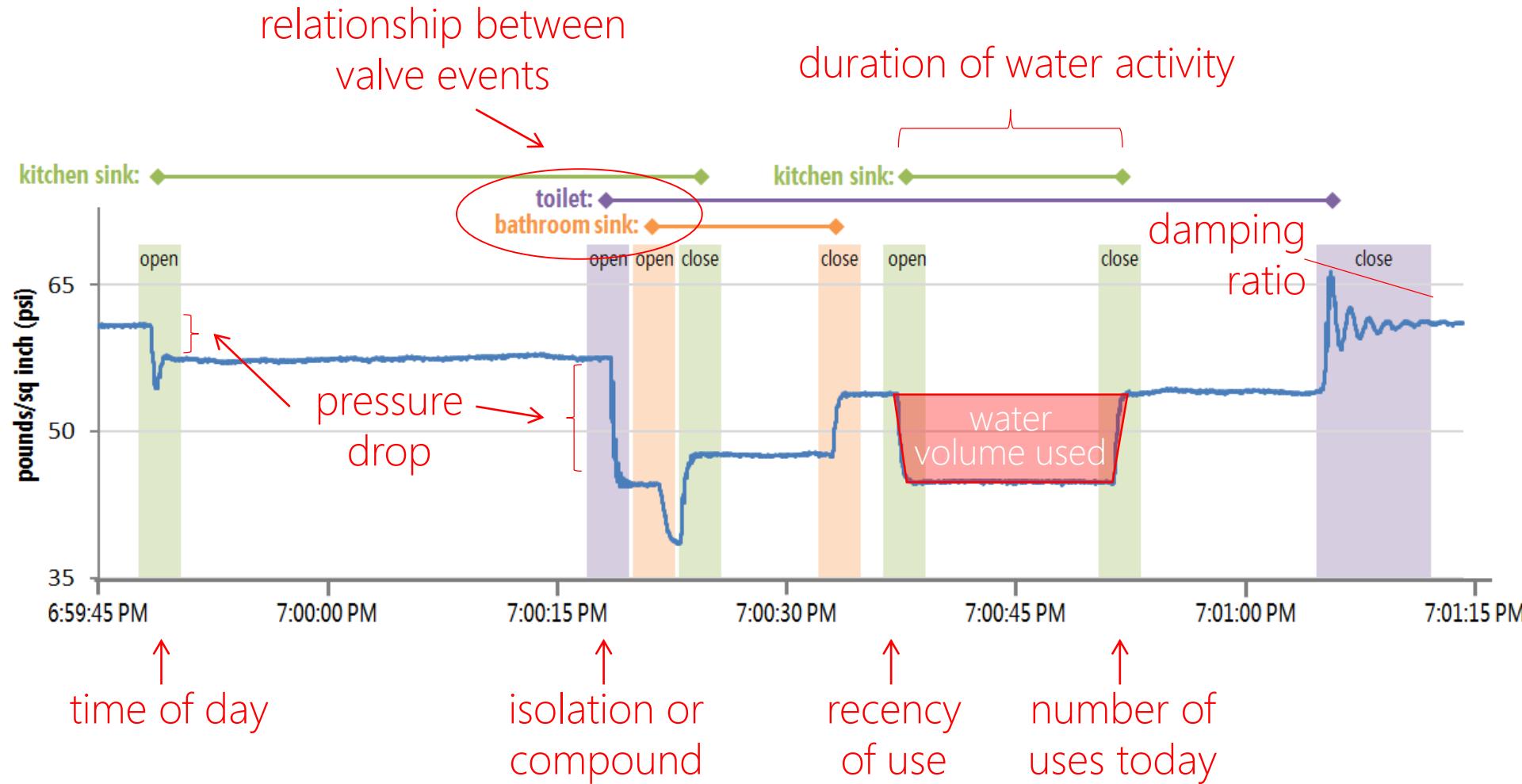
COMPOUND EVENT EXAMPLE



HYDROSENSE 2.0 FIELD DEPLOYMENT RESULTS

MOVING BEYOND TEMPLATE MATCHING

There's lots of information beyond just the pressure signal that we weren't taking advantage of...



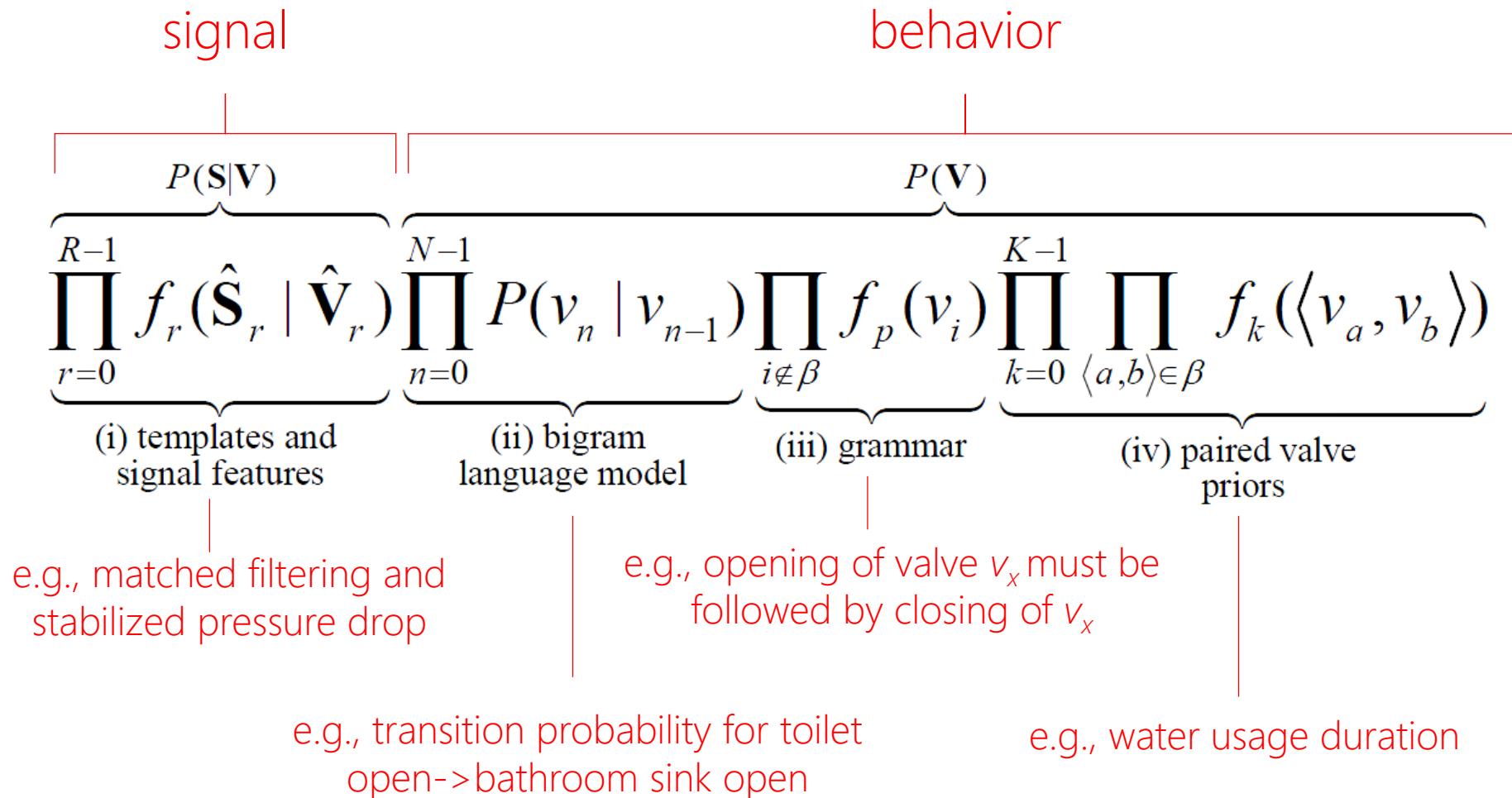
bayesian approach

New algorithm borrows from Bayesian inference in speech recognition—a probabilistic approach (assumes each input feature contributes independent information)

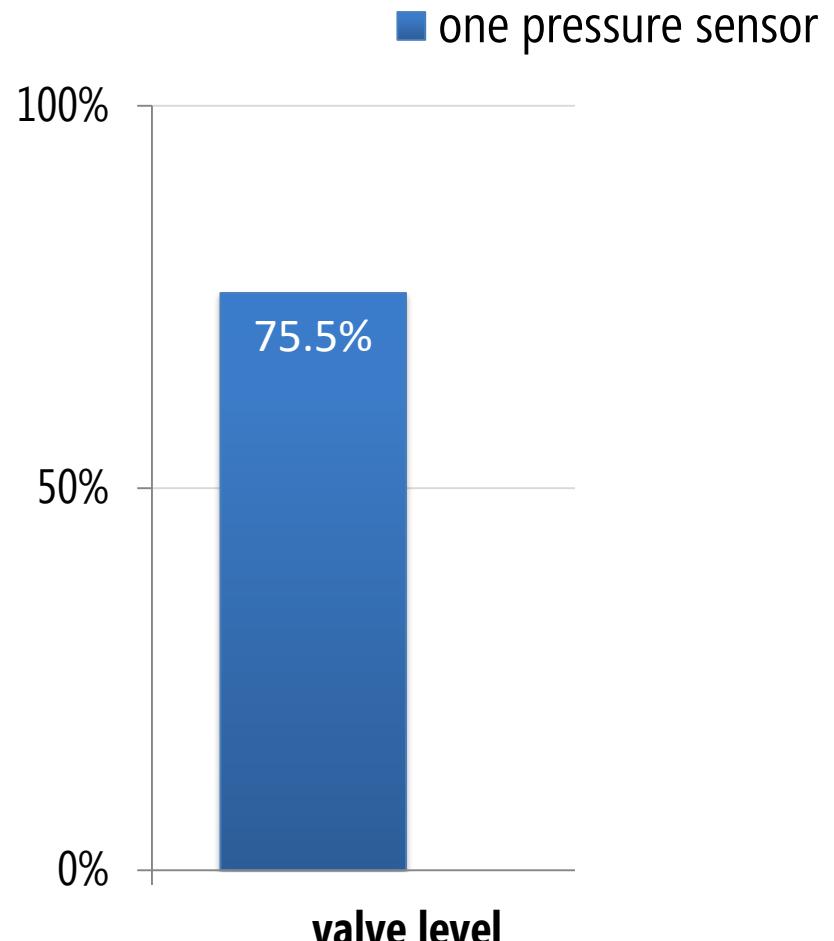
$$\begin{array}{c} \text{signal} \qquad \qquad \qquad \text{behavior} \\ \qquad \qquad \qquad | \\ P(\mathbf{S}|\mathbf{V}) \qquad \qquad \qquad P(\mathbf{V}) \\ \overbrace{\prod_{r=0}^{R-1} f_r(\hat{\mathbf{S}}_r \mid \hat{\mathbf{V}}_r)}^{\text{(i) templates and signal features}} \prod_{n=0}^{N-1} P(v_n \mid v_{n-1}) \prod_{i \notin \beta} f_p(v_i) \prod_{k=0}^{K-1} \prod_{\langle a,b \rangle \in \beta} f_k(\langle v_a, v_b \rangle) \\ \text{(ii) bigram language model} \qquad \qquad \qquad \text{(iii) grammar} \qquad \qquad \qquad \text{(iv) paired value priors} \end{array}$$

bayesian approach

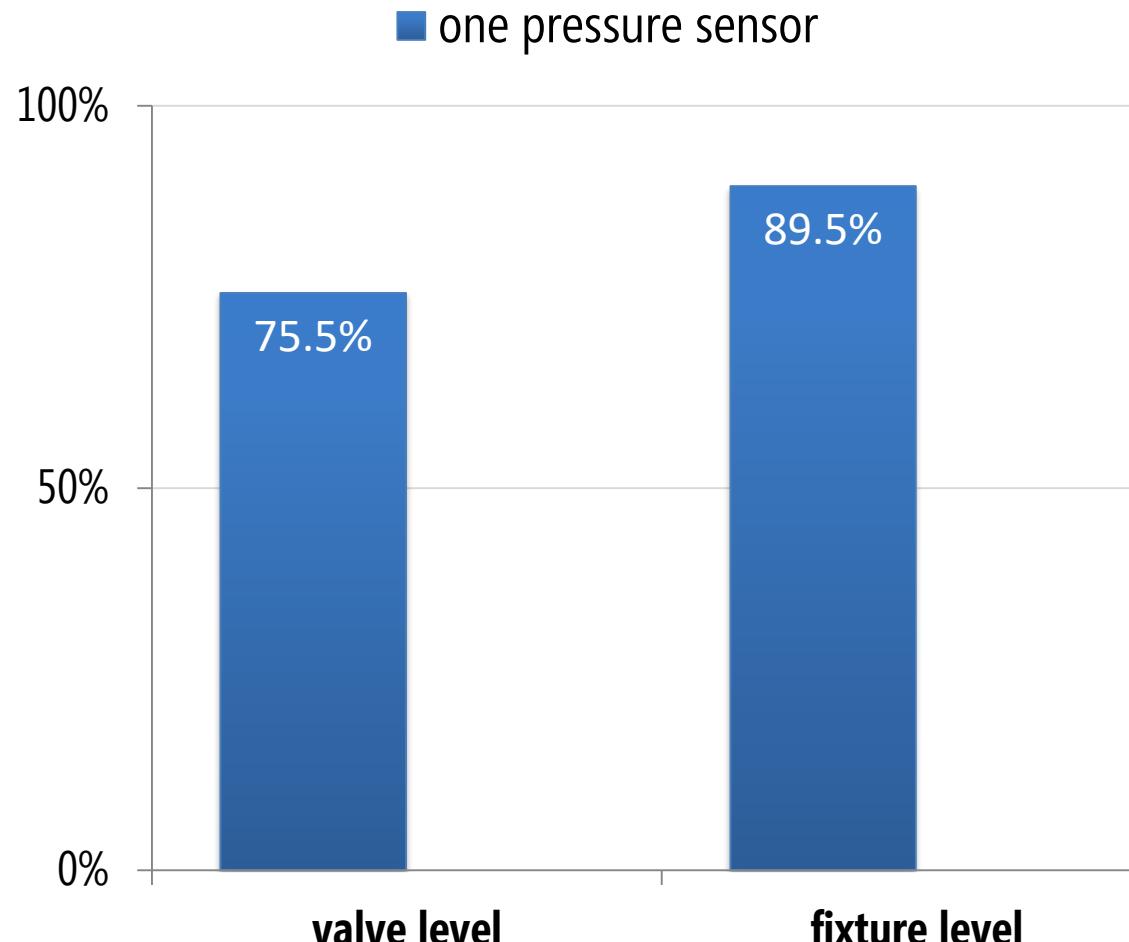
New algorithm borrows from Bayesian inference in speech recognition



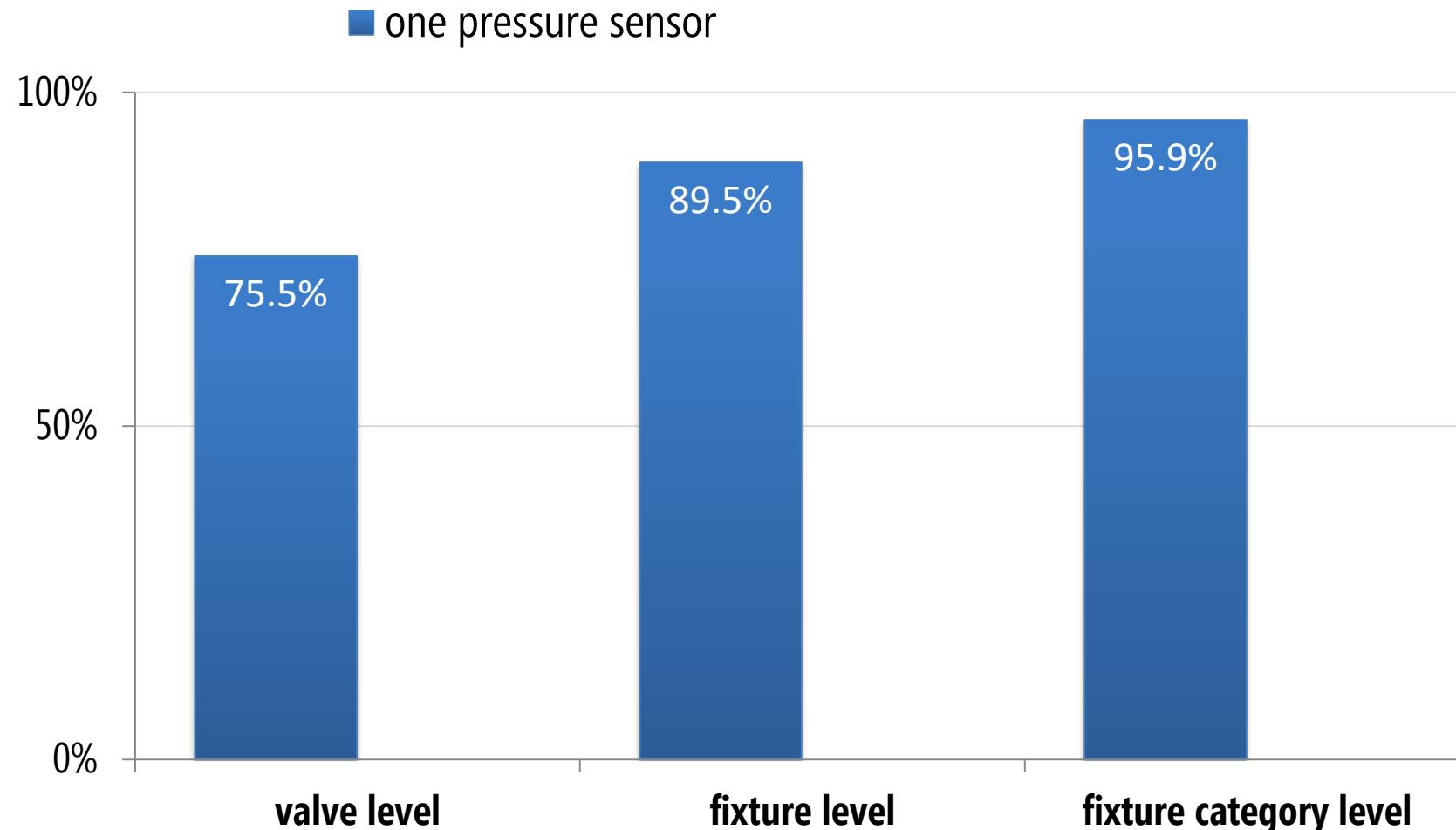
BAYESIAN CLASSIFIER RESULTS



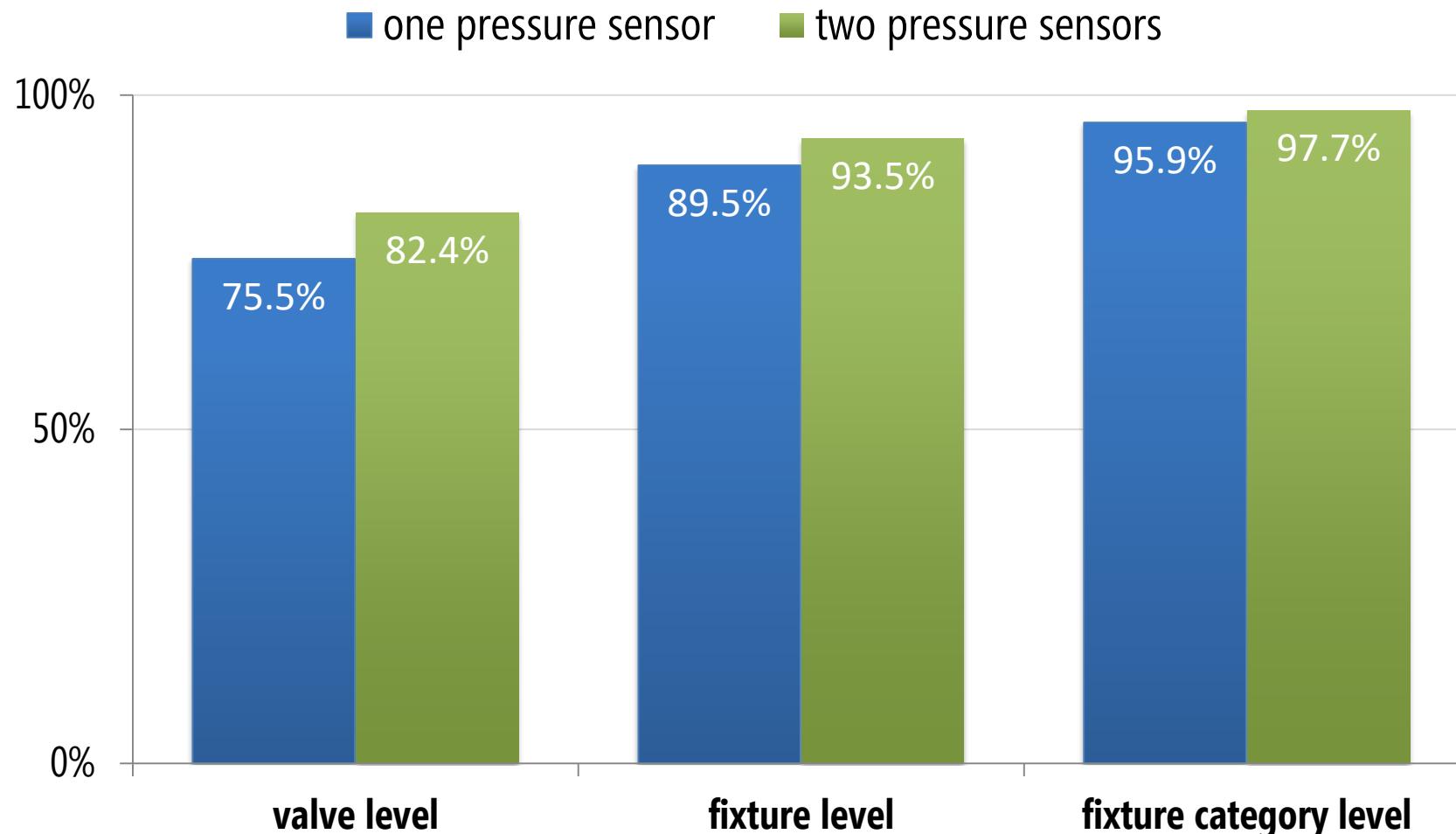
BAYESIAN CLASSIFIER RESULTS



BAYESIAN CLASSIFIER RESULTS

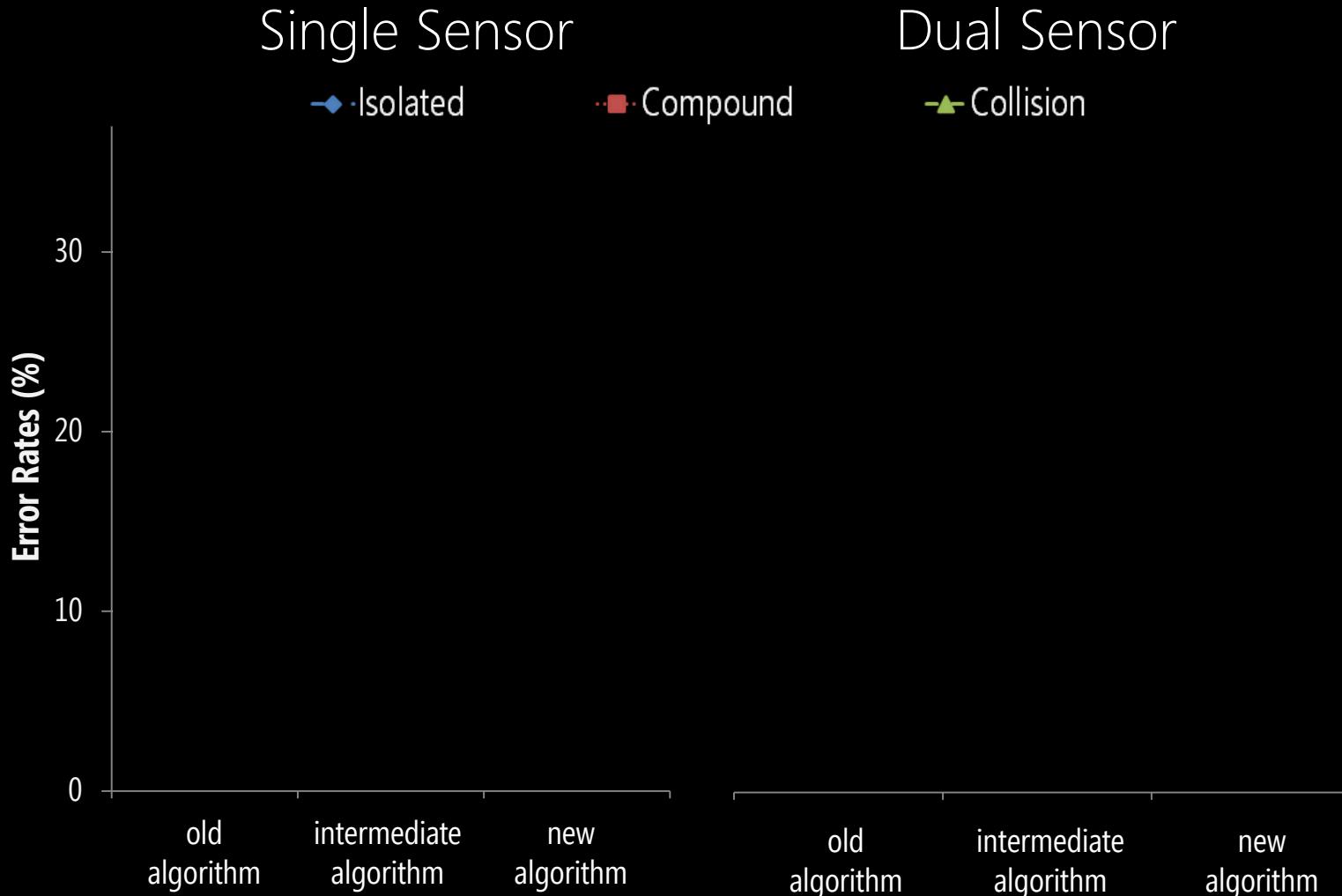


BAYESIAN CLASSIFIER RESULTS



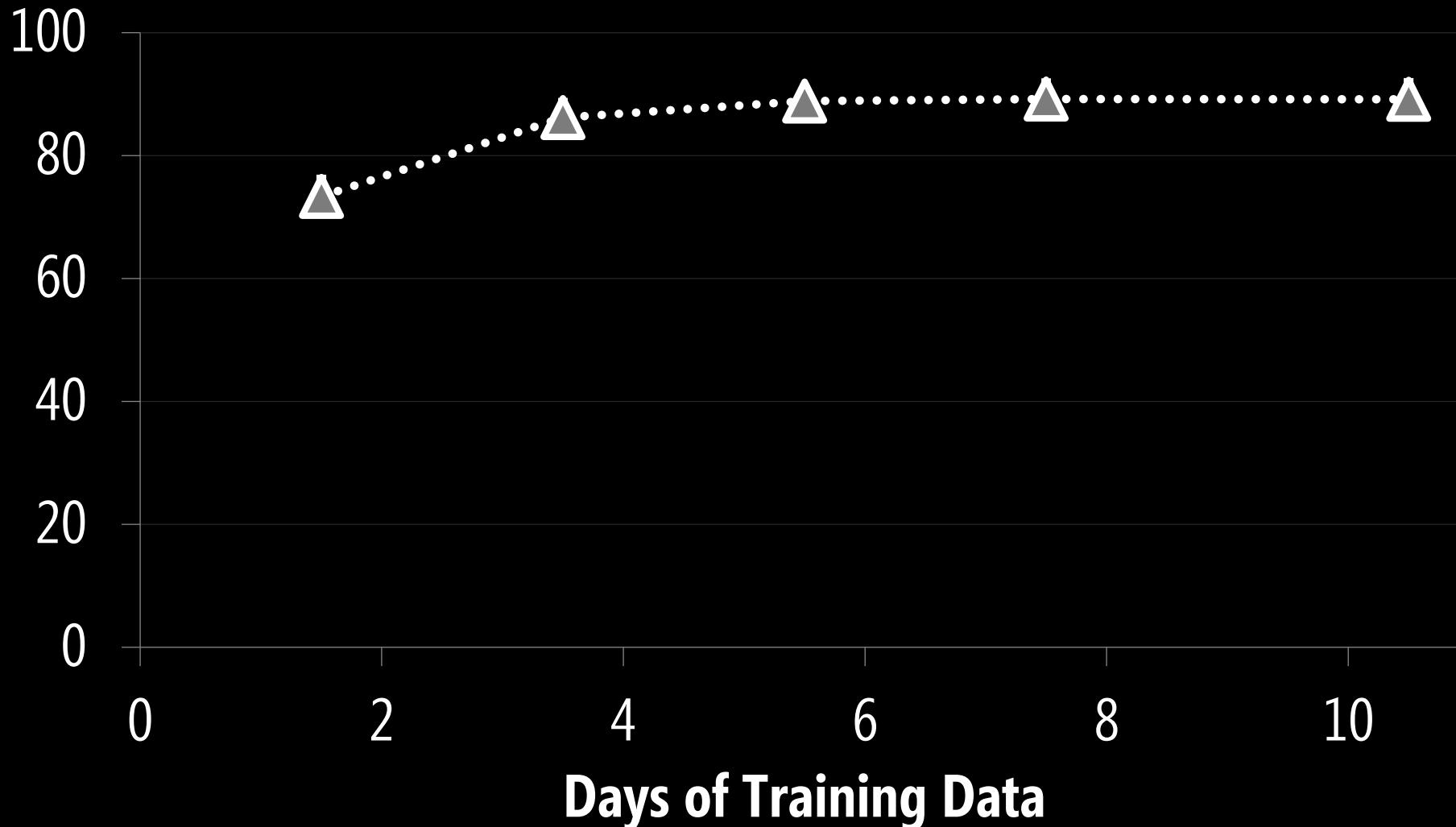
compound events results

real-world water usage data



hydro**sense** training results

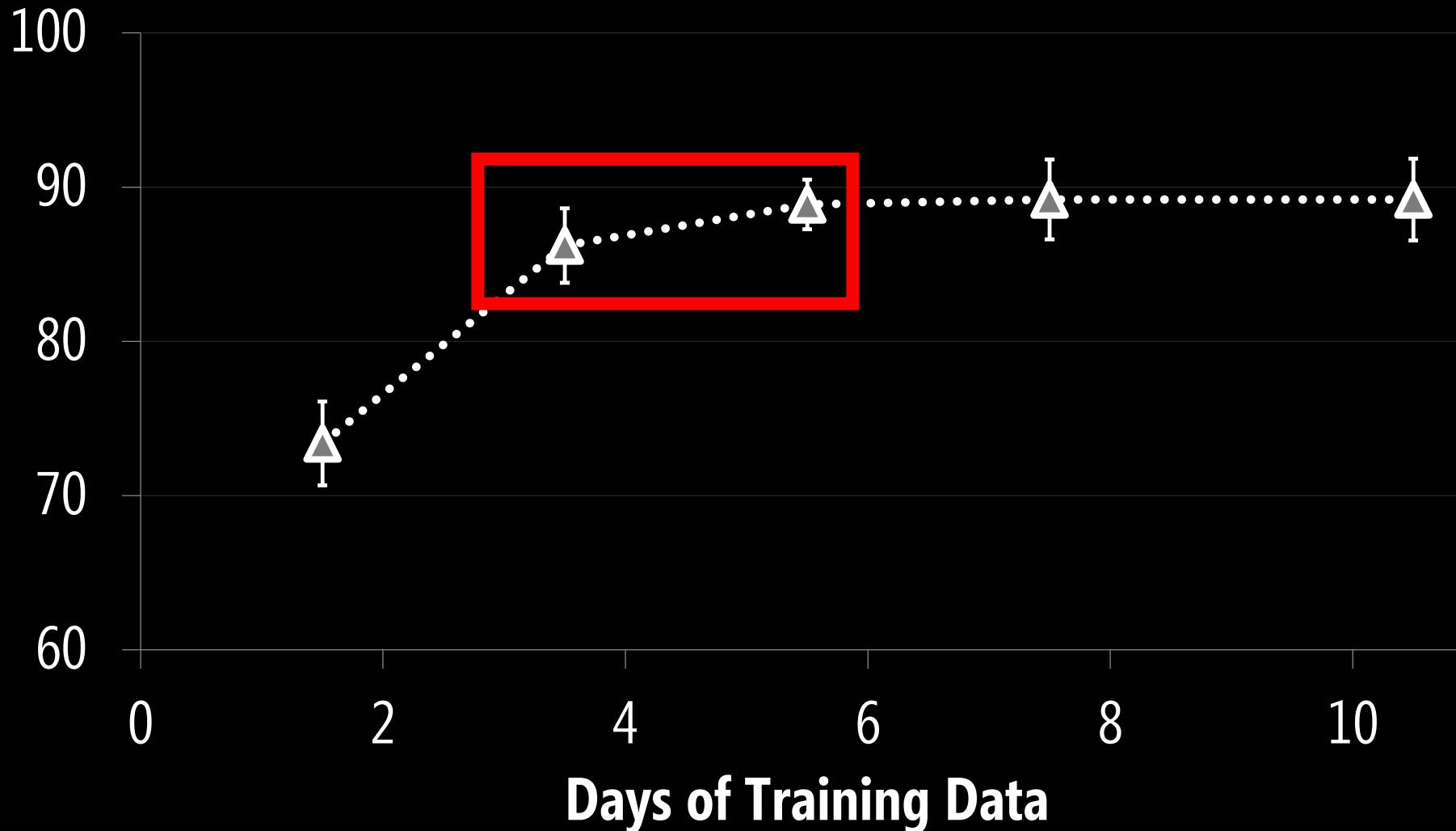
real-world water usage data



*error bars = std error

hydro**sense** training results

real-world water usage data



*error bars = std error

GETTING TRAINING DATA

GETTING HIGH-QUALITY TRAINING DATA CAN BE A SIGNIFICANT BARRIER TO ML IN ACADEMIA AND INDUSTRY

TRANSPORTATION

HOW TESLA AND WAYMO ARE TACKLING A MAJOR PROBLEM FOR SELF-DRIVING CARS: DATA

Autonomous cars won't happen without tons of data, but Tesla and Waymo have a big head start

By Sean O'Kane | @sokane1 | Apr 19, 2018, 8:00am EDT
Illustration by William Joe

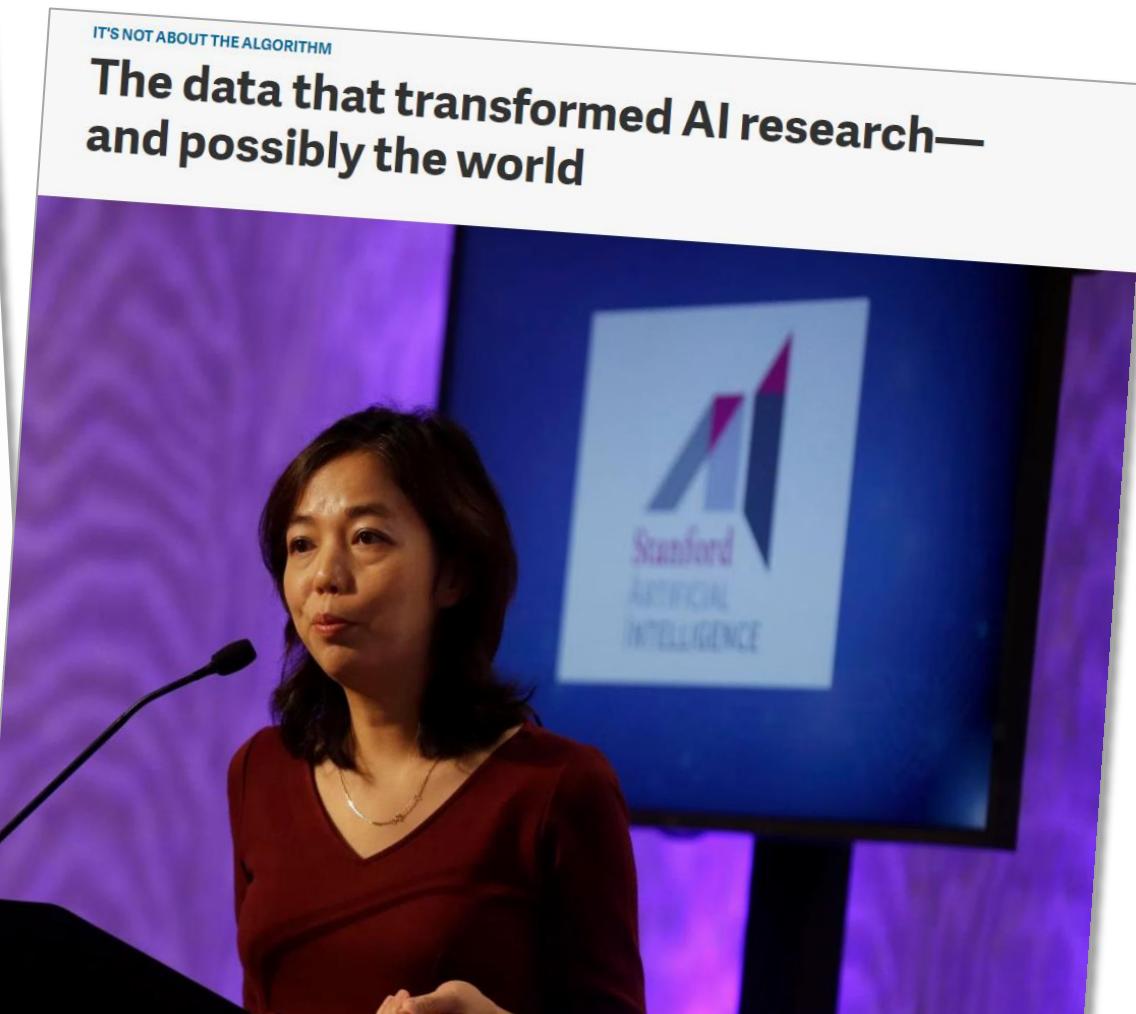
f t SHARE

There's a race happening right now that stretches from Silicon Valley to Detroit and back: who can make a self-driving car that behaves better than a human driver? It's a far harder task than it sounded even a few years ago because human drivers know a lot — not just about their cars but about *how people behave* on the road when they're behind the wheel. To reach that same kind of understanding, computerized cars need lots of data. And the two companies with the most data right now are Tesla and Waymo.

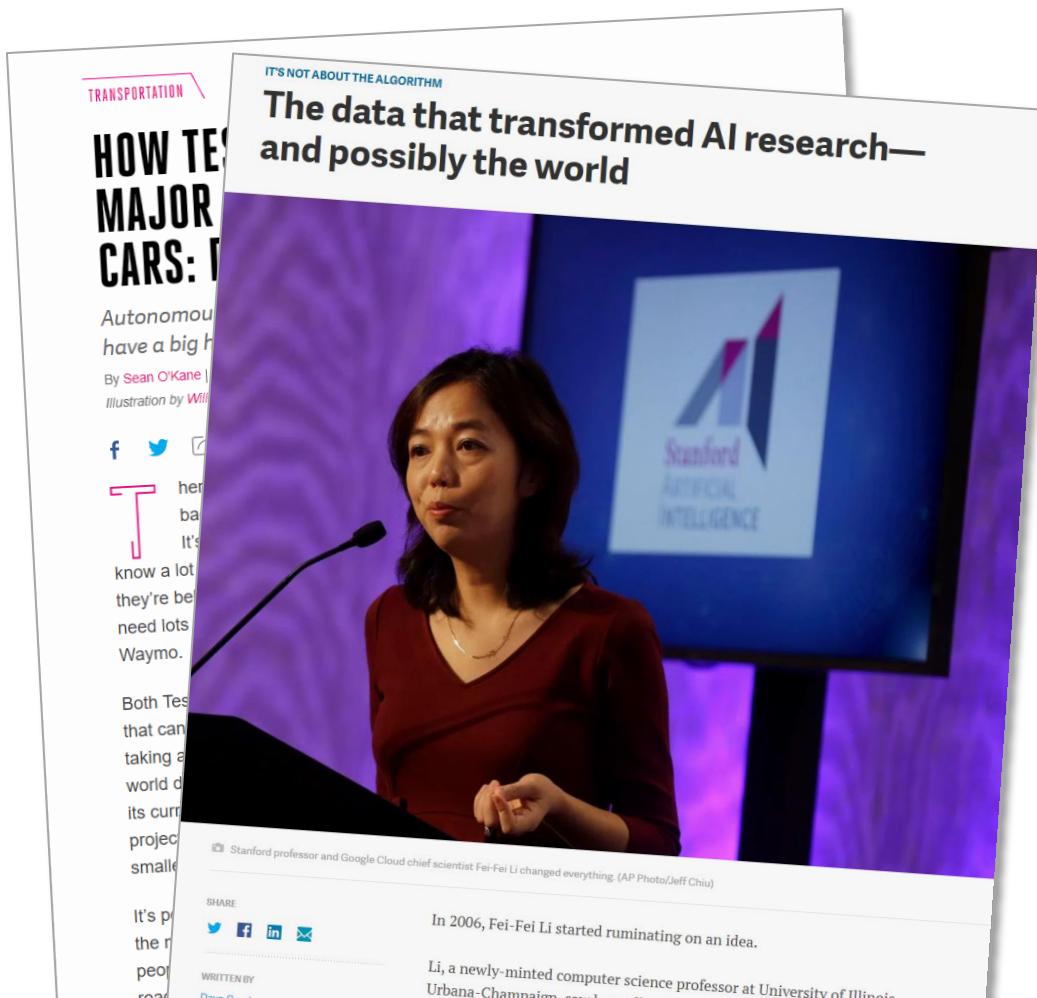
Both Tesla and Waymo are attempting to collect and process enough data to create a car that can drive itself. And they're approaching those problems in very different ways. Tesla is taking advantage of the hundreds of thousands of cars it has on the road by collecting real-world data about how those vehicles perform (and how they *might* perform) with Autopilot, its current semi-autonomous system. Waymo, which started as Google's self-driving car project, uses powerful computer simulations and feeds what it learns from those into a smaller real-world fleet.

IT'S NOT ABOUT THE ALGORITHM

The data that transformed AI research—and possibly the world

A photograph of a woman with dark hair, wearing a red dress, speaking into a microphone at a podium. Behind her is a large screen displaying the Stanford Artificial Intelligence logo, which consists of a stylized 'AI' in blue and pink, with the words 'Stanford' and 'ARTIFICIAL INTELLIGENCE' below it. The background is purple and blue.

GETTING HIGH-QUALITY TRAINING DATA CAN BE A SIGNIFICANT BARRIER TO ML IN ACADEMIA AND INDUSTRY

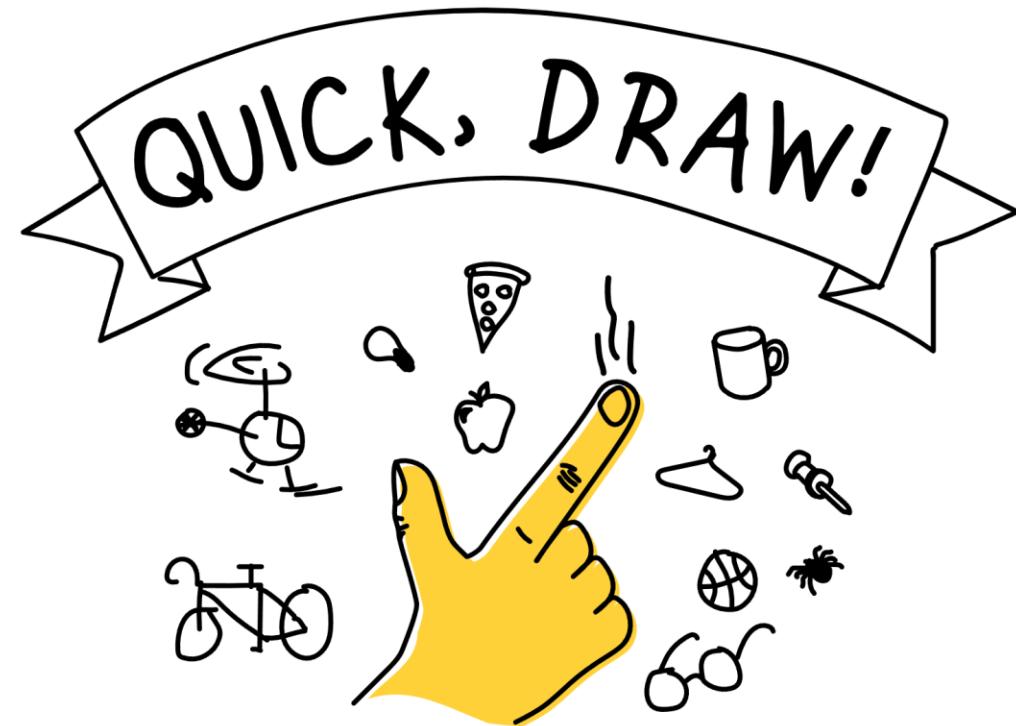


In 2006, Fei-Fei Li started ruminating on an idea.

Li, a newly-minted computer science professor at University of Illinois Urbana-Champaign, saw her colleagues across academia and the AI industry hammering away at the same concept: a better algorithm would make better decisions, regardless of the data.

But she realized a limitation to this approach—the best algorithm wouldn't work well if the data it learned from didn't reflect the real world.

Fei-Fei Li is now a Stanford professor and Chief Scientist at Google Cloud



Can a neural network learn to recognize doodling?

Google for: Google Quick Draw

Let's Draw!

GETTING TRAINING DATA

GAMES WITH A PURPOSE

Labeling Images with a Computer Game

Luis von Ahn and Laura Dabbish
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA, USA
{biglou,dabbish}@cs.cmu.edu

Abstract

We introduce a new interactive system: a game that is fun and can be used to create valuable output. When people play the game they help determine the contents of images by providing useful labels for them. If the game is played as much as popular online games, we estimate that most images on the Web can be labeled in a few months. Having proper labels associated with each image on the Web would allow for more accurate image search, improve the accessibility of sites (by providing descriptions of images to visually impaired individuals), and help users block inappropriate images. Our game has significant commercial benefit because of its valuable output and because of the way it addresses the image-labeling problem. Rather than using computer vision techniques, which don't work well enough, we encourage people to do the work by taking advantage of their desire to be entertained.

Categories & Subject Descriptors: I.2.6 [Learning]: Knowledge acquisition; H.3.m [Information Retrieval]: miscellaneous; H.5.3 [HCI]: Web-based interaction.

General Terms: Design, Human Factors, Languages
Keywords: Distributed knowledge acquisition, image labeling, online games, World Wide Web.

INTRODUCTION

Images on the Web present a major technological challenge. There are millions of them, there are no guidelines about providing appropriate textual descriptions for them, and computer vision hasn't yet produced a program that can determine the contents of an image in a useful way. However, accurate descriptions of images are required by several applications like image search engines and accessibility programs for the visually impaired. Current techniques to categorize images for these applications are insufficient in many ways, mostly because they assume that the contents of images on the Web are related to the text appearing in the page. This is insufficient because the text adjacent to the images is often scarce, and can be misleading or hard to process [4].

Permit to make digital or hard copies of all or part of this work for personal or classroom use without fee and without permission if the copy is not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
CHI 2004, April 24–29, 2004, Vienna, Austria.
Copyright 2004 ACM 1-58113-702-8/04/0004...\$5.00.

Labeling images with a computer game

Von Ahn & Dabbish, CHI'04

2,325 citations



Figure 2. The ESP Game. Players try to “agree” on as many images as they can in 2.5 minutes. The thermometer at the bottom measures how many images partners have agreed on.



Figure 4. An image with all its labels.

GETTING TRAINING DATA

RECAPTCHAS

reCAPTCHA: Human-Based Character Recognition via Web Security Measures

Luis von Ahn,¹ Benjamin Maurer,² Colin McMillen,¹ David Abraham,¹ Manuel Blum¹

CAPTCHAs (Completely Automated Public Turing test to tell Computers and Humans Apart) are widespread security measures on the World Wide Web that prevent automated programs from abusing online services. They do so by asking humans to perform a task that computers cannot yet perform, such as reading distorted characters. Our research shows whether this mental effort can be converted into a useful corpus helping to digitize old printed material by asking users to decipher scanned words from books that computerized optical character recognition failed to recognize. We showed that this method can transcribe text with a word accuracy exceeding 99%, matching the guarantee of professional human transcribers. Our apparatus is deployed in more than 40,000 Web sites and has transcribed over 440 million words.

A CAPTCHA (*1, 2*) is a challenge-response test used on the World Wide Web to determine whether a user is a human or a computer. The acronym stands for Completely Automated Public Turing test to tell Computers and Humans Apart. A typical CAPTCHA is an

Computer Science Department, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, USA.

To whom correspondence should be addressed. E-mail: luis@cs.cmu.edu

random characters rendered by a computer. reCAPTCHA displays words taken from scanned texts. The solutions entered by humans are used to improve the digitization process. To increase efficiency, we ask only two words for which the answer is known. reCAPTCHA programs cannot recognize words that are human. However, to meet the goal of a CAPTCHA (distinguishing between humans and computers), the system needs to be able to verify the user's answer. To do this, reCAPTCHA gives the user two words: the one for which the answer is not known and a second "control" word for which the answer is known. If users correctly type the control word, the system assumes they are human and gains confidence that they also typed the other word correctly (Fig. 1). We describe the exact process below.

We start with an image of a scanned page. Two different OCR programs analyze the image; their respective outputs are then aligned with each other by standard string matching algorithms (*3*). They are compared to each other and to an English dictionary. Any word that is deciphered differently by both OCR programs or that is not in the English dictionary is marked as "suspicious." These are typically the words that the OCR programs failed to decipher correctly. According to our analysis, about 96% of these suspicious words are recognized incorrectly by at least one of the OCR programs; conversely, 99.74% of the words marked as suspicious are deciphered correctly by both programs. Each suspicious word ("morning") was not recognized by OCR, another word for which the answer was known ("overlooks") was also presented to be determined if the user entered the correct answer.

image containing several distorted characters that appears at the bottom of Web registration forms. Users are asked to type the wavy characters to "prove" they are human. Computer programs that guess correctly read distorted words as well as humans can (*3*), so CAPTCHAs act as sentries against automated programs that attempt to abuse online services. Owing to their effectiveness as anti-spam tools, CAPTCHAs are now used to protect many types of Web sites, including free-e-mail providers, ticket sellers, social networks, wikis, guessers match each other, but differ from both of the OCRs' guesses, then (and only then) the word becomes a control word in other challenges. In case of a discrepancy, the program sends the word to humans as an "unknown word" and picks the answer with the highest number of "votes," where each human answer counts as one vote and each OCR guess counts as one half of a vote. The word that has been selected is then previously processed by OCR). In practice, these weights seem to yield the best results, though our accuracy is not very sensitive to them (as long as more weight is given to human guesses

The Norwich line steamboat train, from New-London for Boston, this morning ran off the track seven miles north of New-London.



1466

12 SEPTEMBER 2008 VOL 321 SCIENCE www.sciencemag.org

reCAPTCHA: Human-Based Character Recognition via Web Security Measures

Von Ahn *et al.*, CHI'08

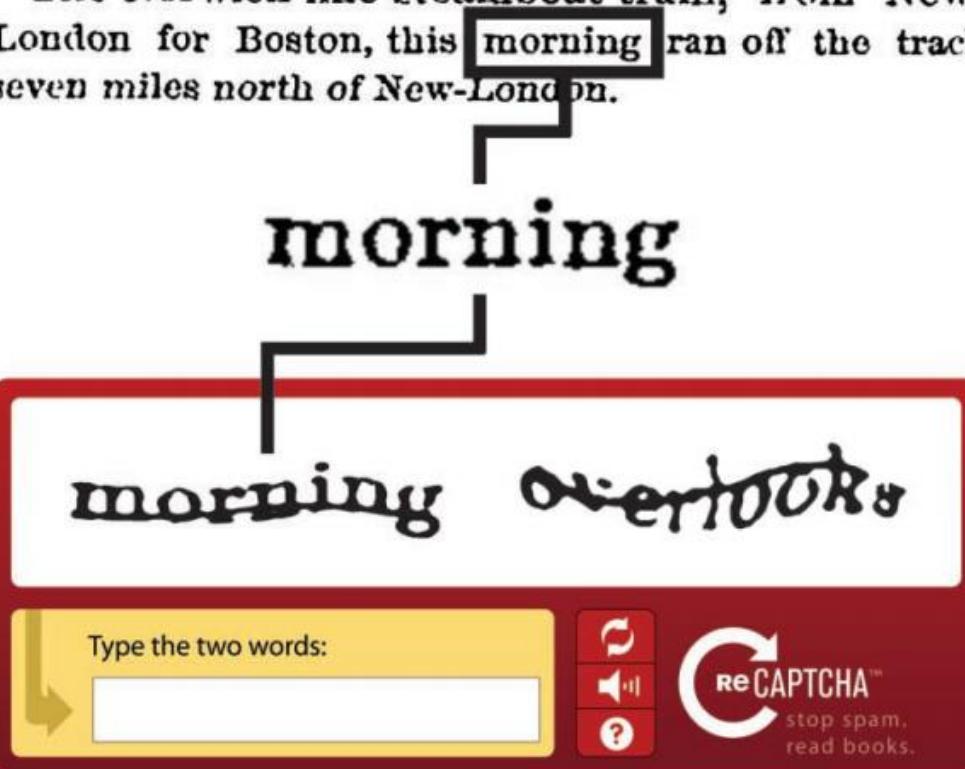
1,043 citations

http://science.sciencemag.org/10.1126/science.1154410

http://science.sciencemag.org/10.1126/science.1154410

Fig. 1. The reCAPTCHA system displays words from scanned texts to humans on the World Wide Web. In this example, the word "morning" was unrecognizable by OCR. reCAPTCHA isolated the word, distorted it using random transformations including adding a line through it, and then presented it as a challenge to a user. Because the original word ("morning") was not recognized by OCR, another word for which the answer was known ("overlooks") was also presented to determine if the user entered the correct answer.

The Norwich line steamboat train, from New-London for Boston, this morning ran off the track seven miles north of New-London.



reCAPTCHA™
stop spam.
read books.

GETTING TRAINING DATA **RECAPTCHAS**

reCAPTCHA: Human-Based Character Recognition via Web Security Measures

Luis von Ahn,* Benjamin Maurer, Colin McMillen, David Abraham, Manuel Blum

CAPTCHAs (Completely Automated Public Turing test to tell Computers and Humans Apart) are widespread security measures on the World Wide Web that prevent automated programs from abusing online services. They do so by asking humans to perform a task that computers cannot yet perfectly execute. Our research group has developed a new type of CAPTCHA that we believe can help combat spam. This effort can be characterized as a useful corporate resource: helping to digitize old printed material by asking users to decipher scanned words from books that comprised optical character recognition failed to recognize. We showed that this method can transcribe text with a word accuracy exceeding 99%, matching the guarantee of professional human transcribers. Our apparatus is deployed in more than 40,000 Web sites and has transcribed over 440 million words.

ACAPTCHA (*I*, 2) is a challenge-response test used on the World Wide Web to determine whether a user is a human or a computer. The acronym stands for Completely Automated Public Turing test to tell Computers and Humans Apart. A typical CAPTCHA is an

Computer Science Department, Carnegie Mellon University,
5000 Forbes Avenue, Pittsburgh, PA 15213, USA.

of random characters rendered by a computer, reCAPTCHA displays words taken from texts, such as books and movies, which are used to improve the digitization process. To increase efficiency and security, only the words that automated OCR programs cannot recognize are sent to humans. However, to meet the goal of a CAPTCHA (differentiating between humans and computers), the system needs to be able to verify the user's answer. To do this, it asks the user a second question, one for which the answer is not known and a second "control" word for which the answer is known. If users correctly type the control word, the system assumes they are human and gains confidence that they also typed the other word correctly (Fig. 1). We describe the exact process below.

We start an analysis of a scanned page by two OCR programs. The results of their respective outputs are then aligned with each other by standard string matching algorithms (7) and compared to each other and to an English dictionary. Any word that is deciphered differently by both OCR programs or that is not in the English dictionary is marked as "suspicous". These are typically words that the OCR programs failed to decipher correctly. According to our tests, about 96% of the words that are marked as suspicious are in fact at least one of the words the OCR programs, conversely, 97.4% of the words marked as suspicious are deciphered correctly by both programs. Each suspicious word is then placed in an image along with another word for which the answer is already known, the two words are distorted further to ensure that automated programs cannot decipher them, and the resulting image is used as a CAPTCHA. Users are asked to

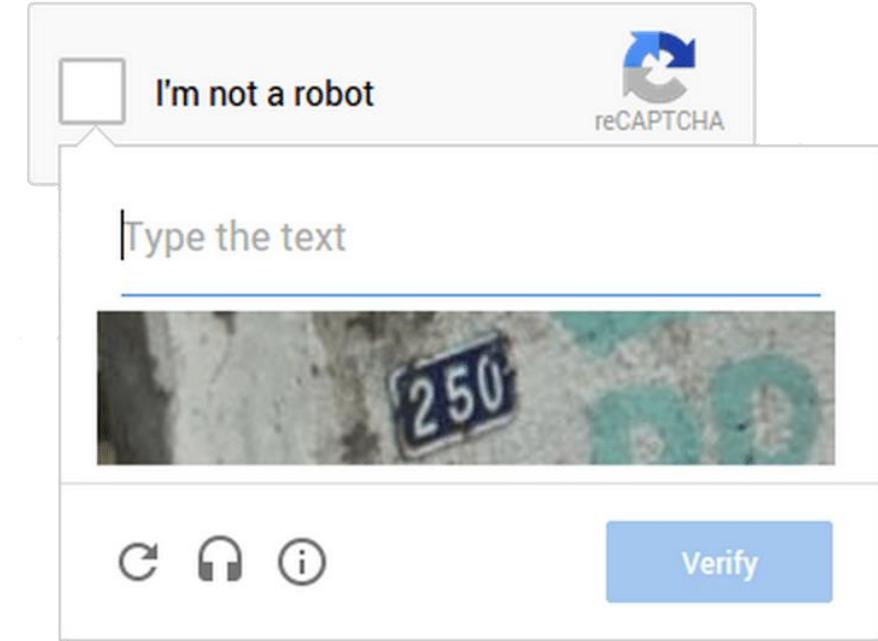
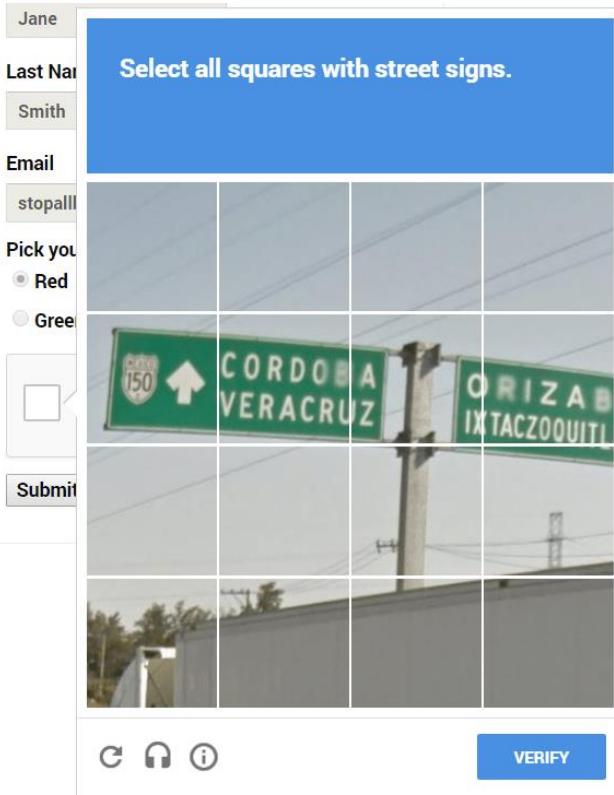
The Norwich line steamboat train, from New-London for Boston, this morning ran off the track seven miles north of New-London.

166

12 SEPTEMBER 2008 VOL 321 SCIENCE www.sciencemag.org

reCAPTCHA: Human-Based Character Recognition via Web Security Measures

Von Ahn *et al.*, CHI'08
1,043 citations



Google purchased von Ahn's CAPTCHA company and has integrated it into its products

PARAMETER TUNING

Most ML algorithms have parameters to set

When you train a model, you must provide values for those parameters and the efficacy of the model depends on those parameters

The optimal set of parameters is known as *model selection*

(Yes, this is similar to feature selection and the approaches are generally similar too!)

PARAMETER TUNING PROCESS

1. Define the parameter space. For a given ML model, decide what range of parameter values you want to consider.
2. Define the evaluation method (*e.g.*, decide how to choose cross-validation folds for the dataset)
3. Define your key metrics (*e.g.*, accuracy, precision, recall, f-measure, or root-mean squared error)
4. Train and evaluate. For each unique combination of the parameter values, cross-validation performed and your key metrics produced.
5. Compare output and choose the best-performing parameters for your model.

PARAMETER TUNING: DECISION TREE EXAMPLE ON AZURE

Most ML toolkits offer parameter tuning functionality. This example is from Microsoft Azure's ML toolkit.

Two-Class Boosted Decision Tree

Create trainer mode
Single Parameter

Maximum number of leaves per tree
20

Minimum number of samples per leaf node
10

Learning rate
0.2

Number of trees constructed
100

Random number seed

Allow unknown categorical levels

Two-Class Boosted Decision Tree

Create trainer mode
Parameter Range

Maximum number of leaves per tree
Use Range Builder
Parameter Range : 101 - 900

Number of points : 3
Log Scale

Minimum number of samples per leaf n...
Use Range Builder
1, 10, 50

Learning rate
Use Range Builder
0.025, 0.05, 0.1, 0.2, 0.4

Number of trees constructed
Use Range Builder
1-5, 20, 100, 500

Random number seed

Allow unknown categorical levels

Partition and Sample

Partition or sample mode
Assign to Folds

Use replacement in the partitioning

Randomized split

Random seed
0

Specify the partitioner method
Partition evenly

Specify number of folds to split evenly i...
5

Stratified split
False

Sweep Parameters

Specify parameter sweeping mode
Entire grid

Label column

Selected columns:
Column names: income

Launch column selector

Metric for measuring performance for classification
Accuracy

Metric for measuring performance for regression
Mean absolute error

1. Choose and Setup Model

Select an ML model to tune

2. Parameter Ranges

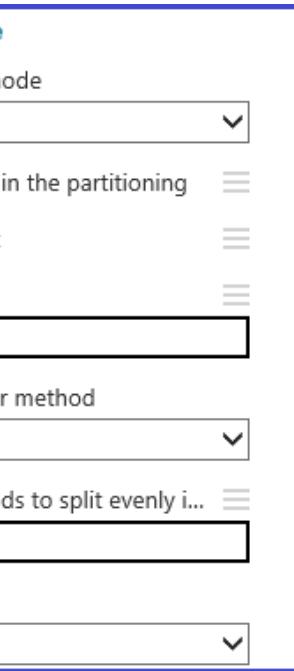
Decide what range of values to input and compare

3. Setup Evaluation

Decide how many folds in cross-validation and key error metrics.

PARAMETER TUNING: DECISION TREE EXAMPLE ON AZURE

Most ML toolkits offer parameter tuning functionality. This example is from Microsoft Azure's ML toolkit.



Iteration

5 folds in cross-validation and key error metrics.

rows	columns											
180	11	Number of leaves	Minimum leaf instances	Learning rate	Number of trees	Accuracy	Precision	Recall	F-Score	AUC	Average Log Loss	Training Log Loss
		8	10	0.1	500	0.866128	0.75785	0.652595	0.701295	0.922837	0.290561	47.363286
		8	1	0.05	500	0.866097	0.762796	0.644306	0.698562	0.922672	0.289389	47.575474
		8	1	0.2	100	0.865944	0.763813	0.641755	0.697484	0.921469	0.291591	47.176649
		32	1	0.05	100	0.865913	0.759367	0.648769	0.699725	0.921503	0.29273	46.970281
		8	10	0.2	100	0.865882	0.763821	0.641372	0.697262	0.922052	0.290024	47.46048
		32	1	0.025	500	0.865668	0.749604	0.663946	0.70418	0.923187	0.292484	47.014813
		8	10	0.05	500	0.865514	0.761322	0.643158	0.697269	0.922803	0.288929	47.658898
		32	1	0.1	100	0.865453	0.751234	0.659737	0.702519	0.922703	0.292779	46.961499
		8	1	0.025	500	0.865391	0.764778	0.636909	0.695011	0.921386	0.29142	47.207519
		8	50	0.05	500	0.865115	0.760932	0.641372	0.696055	0.92261	0.288991	47.647701
		32	10	0.025	500	0.865115	0.748308	0.662798	0.702962	0.923435	0.29145	47.202173
		8	50	0.025	500	0.865084	0.764336	0.635761	0.694145	0.921226	0.291554	47.183315
		32	10	0.05	100	0.865084	0.756395	0.648642	0.698387	0.92155	0.29239	47.031913
		8	10	0.025	500	0.864992	0.764066	0.635633	0.693957	0.921674	0.290844	47.312033
		32	10	0.2	20	0.864961	0.757939	0.645326	0.697114	0.920488	0.294732	46.607674
		8	1	0.1	500	0.864808	0.752756	0.653105	0.699399	0.922545	0.291764	47.145278
		32	10	0.1	100	0.864777	0.748626	0.66012	0.701593	0.92295	0.291274	47.234005
		8	1	0.1	100	0.864715	0.763255	0.635251	0.693395	0.920514	0.293092	46.904731
		8	10	0.1	100	0.864715	0.763416	0.634996	0.693309	0.920639	0.292815	46.954929
		8	10	0.4	100	0.864623	0.754863	0.648387	0.697585	0.921302	0.292915	46.936722
		32	1	0.2	20	0.864439	0.757089	0.64354	0.695712	0.920124	0.295428	46.481506

4. Train, Evaluate, and compare

You can see the parameters in left-hand columns and the output produced.

EVALUATING ML ALGORITHMS

Splitting up training and testing data (most common approach: k-fold cross validation)

Evaluating performance as a function of training data (e.g., train ML model with minimum dataset, evaluate, and repeat with incrementally more data)

Evaluating performance as a function of training data type (e.g., articulation speed of gesture)

Evaluating performance with N-Best lists

Evaluating performance with confusion matrices

KEY ML EVALUATION METRICS

The simplest evaluation metric is **accuracy**

$$\text{Accuracy} = \frac{\text{Number of correct classifications}}{\text{Total number of classifications}}$$

KEY ML EVALUATION METRICS

We can reframe many classification problems as retrieval problems. That is, rather than saying "classify gesture X," we say, "find all instances of gesture X in the dataset." For example, let's imagine we are searching for the "Forehand Tennis" gesture. (Note: for this to work, we would have to have either a learned or manually set "correct match" score threshold)

True positive = # of correctly returned objects

e.g., # of returned "Forehand Tennis" gestures

False positive = # of incorrectly returned objects

e.g., # of returned objects that were not "Forehand Tennis" gestures—they could be "At Rest" or "Backhand Tennis" gestures, etc.

True negative = # of objects that were correctly not returned

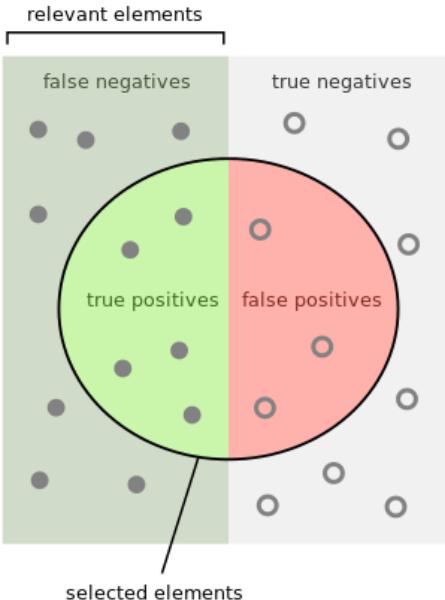
e.g., # of objects that were not "Forehand Tennis" gestures in our dataset and were correctly not returned

False negative = # of objects that should have been found but were not

e.g., # of objects that were "Forehand Tennis" gestures in our dataset but were not actually found and returned.

KEY ML EVALUATION: PRECISION AND RECALL

Recall is the fraction of relevant instances that have been retrieved over the total amount of relevant instances in our dataset and Precision is the fraction of relevant instances among the retrieved instances.



Recall

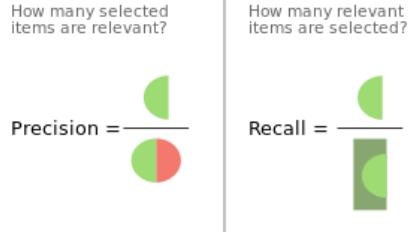
How much of the correct things did we find?

$$\text{Recall} = \frac{\# \text{ of True Positives}}{\# \text{ of True Positives} + \# \text{ of False Negatives}}$$

Precision

Of those things found, how many were correct?

$$\text{Precision} = \frac{\# \text{ of True Positives}}{\# \text{ of True Positives} + \# \text{ of False Positives}}$$



HOW TO EVALUATE ML ALGORITHMS

Input Features				Class Label (aka Ground Truth Label)
Index	Color	Weight	Label	
1	Red	200g	Apple	
2	Orange	300g	Orange	
3	Red	150g	Apple	
N	...			

Total Labeled Dataset

How much of our dataset should we reserve for training and how much for testing?

ML EVALUATION

K-FOLD CROSS VALIDATION

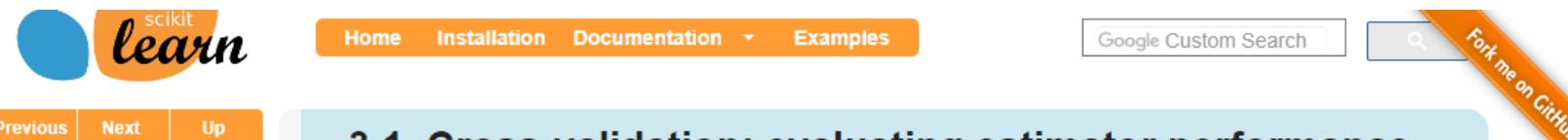


C



ML EVALUATION

SCI-KIT SUPPORT FOR K-FOLD



[Previous](#)
3. Model sele...

[Next](#)
3.2. Tuning t...

[Up](#)
3. Model sele...

scikit-learn v0.19.1
[Other versions](#)

Please [cite us](#) if you use
the software.

3.1. Cross-validation: evaluating estimator performance

3.1.1. Computing cross-validated metrics

- 3.1.1.1. The `cross_validate` function and multiple metric evaluation
- 3.1.1.2. Obtaining predictions by cross-validation
- 3.1.2. Cross validation iterators
- 3.1.2.1. Cross-validation iterators for i.i.d. data
 - 3.1.2.1.1. K-fold
 - 3.1.2.1.2. Repeated K-fold
 - 3.1.2.1.3. Leave One Out (LOO)
 - 3.1.2.1.4. Leave P Out (LPO)
 - 3.1.2.1.5. Random permutations cross-validation a.k.a. Shuffle & Split
- 3.1.2.2. Cross-validation iterators with stratification based on class labels.
 - 3.1.2.2.1. Stratified k-fold
 - 3.1.2.2.2. Stratified Shuffle Split
- 3.1.2.3. Cross-validation iterators for grouped data.
 - 3.1.2.3.1. Group k-fold
 - 3.1.2.3.2. Leave One Group Out
 - 3.1.2.3.3. Leave P Groups Out

3.1. Cross-validation: evaluating estimator performance

Learning the parameters of a prediction function and testing it on the same data is a methodological mistake: a model that would just repeat the labels of the samples that it has just seen would have a perfect score but would fail to predict anything useful on yet-unseen data. This situation is called **overfitting**. To avoid it, it is common practice when performing a (supervised) machine learning experiment to hold out part of the available data as a **test set** `X_test, y_test`. Note that the word "experiment" is not intended to denote academic use only, because even in commercial settings machine learning usually starts out experimentally.

In scikit-learn a random split into training and test sets can be quickly computed with the `train_test_split` helper function.

Let's load the iris data set to fit a linear support vector machine on it:

```
>>> import numpy as np
>>> from sklearn.model_selection import train_test_split
>>> from sklearn import datasets
>>> from sklearn import svm

>>> iris = datasets.load_iris()
>>> iris.data.shape, iris.target.shape
((150, 4), (150,))
```

We can now quickly sample a training set while holding out 40% of the data for testing (evaluating) our classifier:

```
>>> X_train, X_test, y_train, y_test = train_test_split(
...     iris.data, iris.target, test_size=0.4, random_state=0)

>>> X_train.shape, y_train.shape
((90, 4), (90,))
>>> X_test.shape, y_test.shape
((60, 4), (60,))

>>> clf = svm.SVC(kernel='linear', C=1).fit(X_train, y_train)
>>> clf.score(X_test, y_test)
0.96...
```

When evaluating different settings ("hyperparameters") for estimators, such as the `C` setting that must be manually set for an SVM, there is still a risk of overfitting *on the test set* because the parameters can be tweaked until the estimator performs optimally. This way, knowledge about the test set can "leak" into the model and evaluation metrics no longer report the true performance of the model. To avoid this, we need to use a different type of cross-validation fold, called a *cross-validation iterator*.

SCI-KIT SUPPORT FOR K-FOLD

3.1.2.1.3. Leave One Out (LOO)

`LeaveOneOut` (or LOO) is a simple cross-validation. Each learning set is created by taking all the samples except one, the test set being the sample left out. Thus, for n samples, we have n different training sets and n different tests set. This cross-validation procedure does not waste much data as only one sample is removed from the training set:

Don't do this!

```
>>> from sklearn.model_selection import LeaveOneOut  
  
>>> X = [1, 2, 3, 4]  
>>> loo = LeaveOneOut()  
>>> for train, test in loo.split(X):  
...     print("%s %s" % (train, test))  
[1 2 3] [0]  
[0 2 3] [1]  
[0 1 3] [2]  
[0 1 2] [3]
```

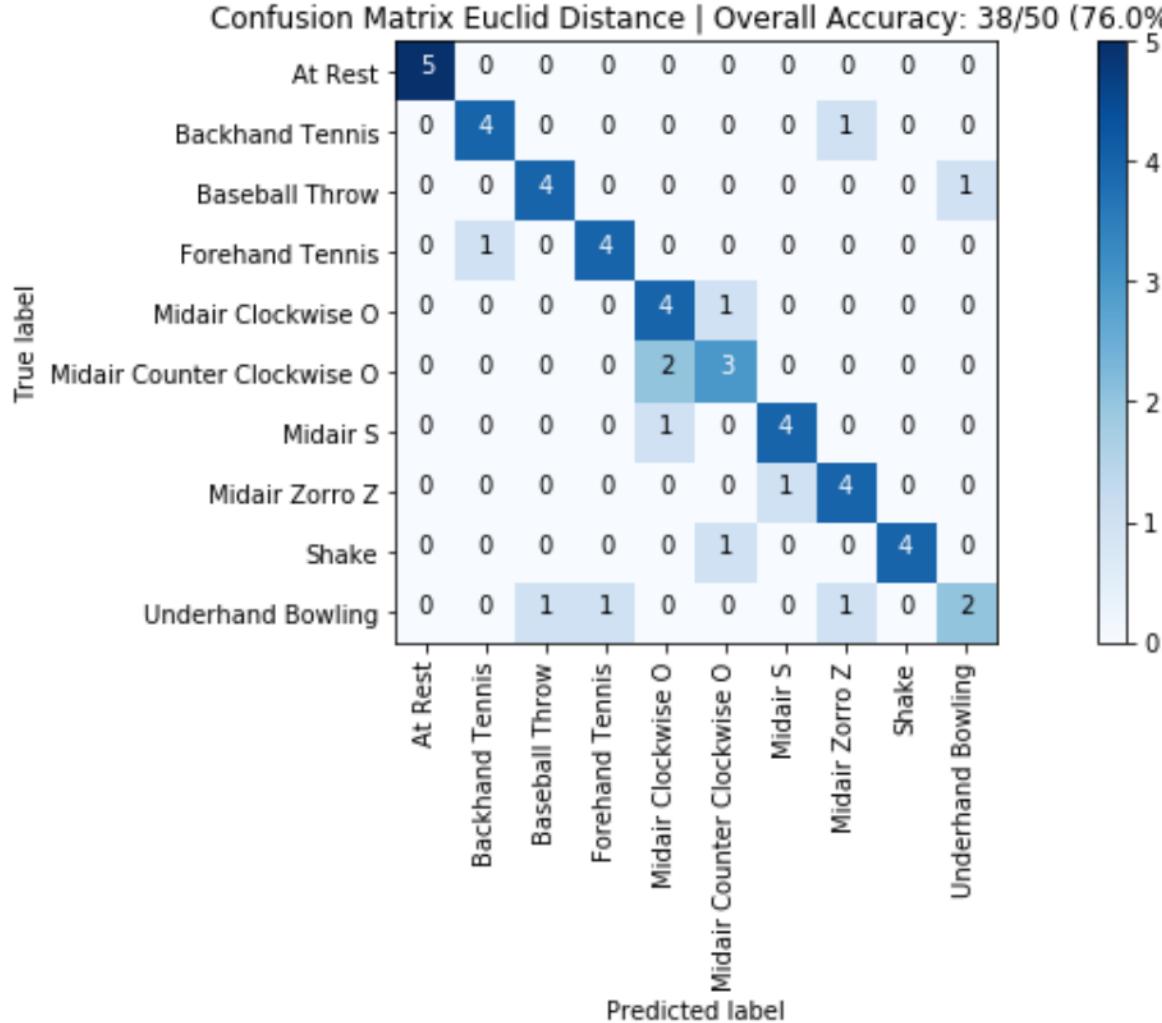
Potential users of LOO for model selection should weigh a few known caveats. When compared with k -fold cross validation, one builds n models from n samples instead of k models, where $n > k$. Moreover, each is trained on $n - 1$ samples rather than $(k - 1)n/k$. In both ways, assuming k is not too large and $k < n$, LOO is more computationally expensive than k -fold cross validation.

In terms of accuracy, LOO often results in high variance as an estimator for the test error. Intuitively, since $n - 1$ of the n samples are used to build each model, models constructed from folds are virtually identical to each other and to the model built from the entire training set.

However, if the learning curve is steep for the training size in question, then 5- or 10- fold cross validation can overestimate the generalization error.

As a general rule, most authors, and empirical evidence, suggest that 5- or 10- fold cross validation should be preferred to LOO.

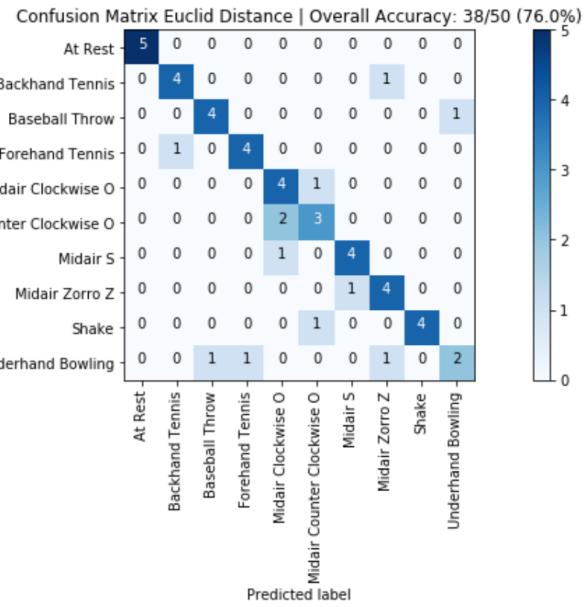
CONFUSION MATRICES



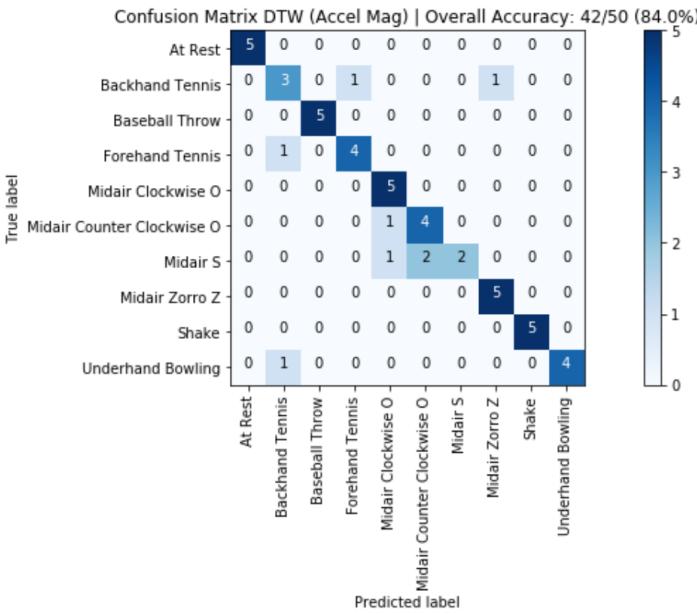
CONFUSION MATRICES

Just Using Accel Magnitude

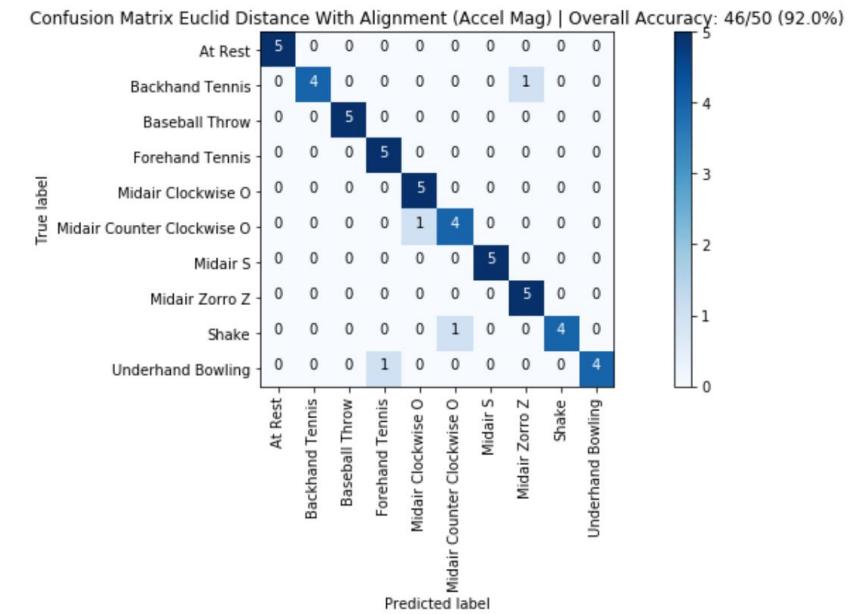
True label



True label



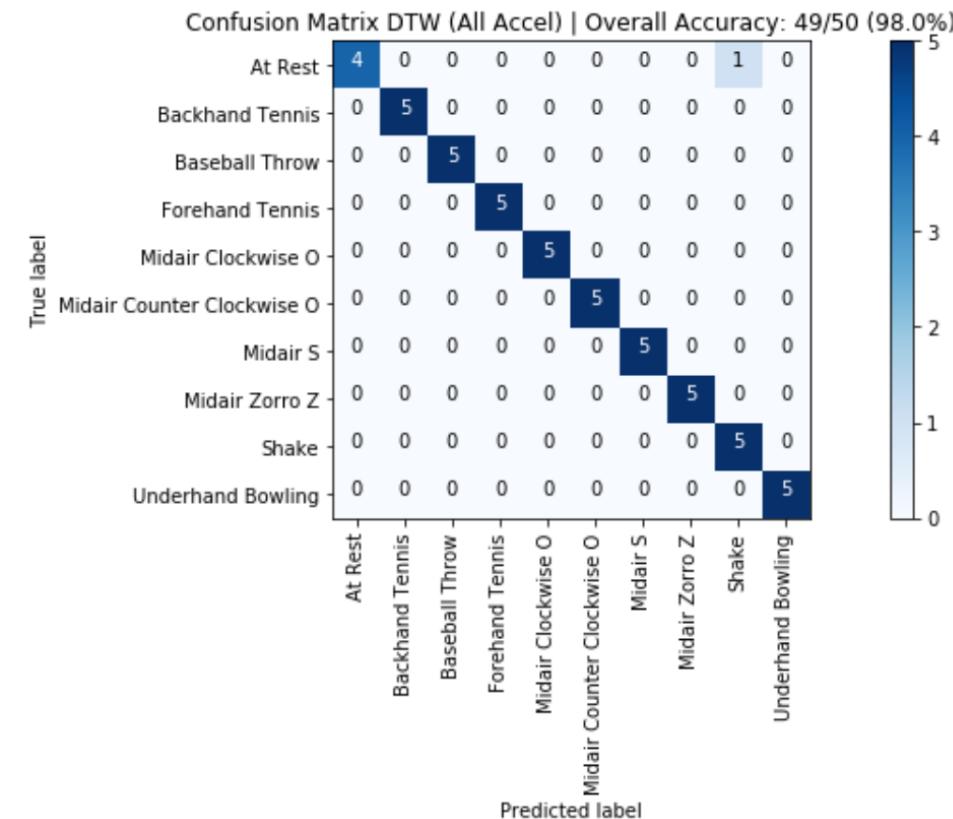
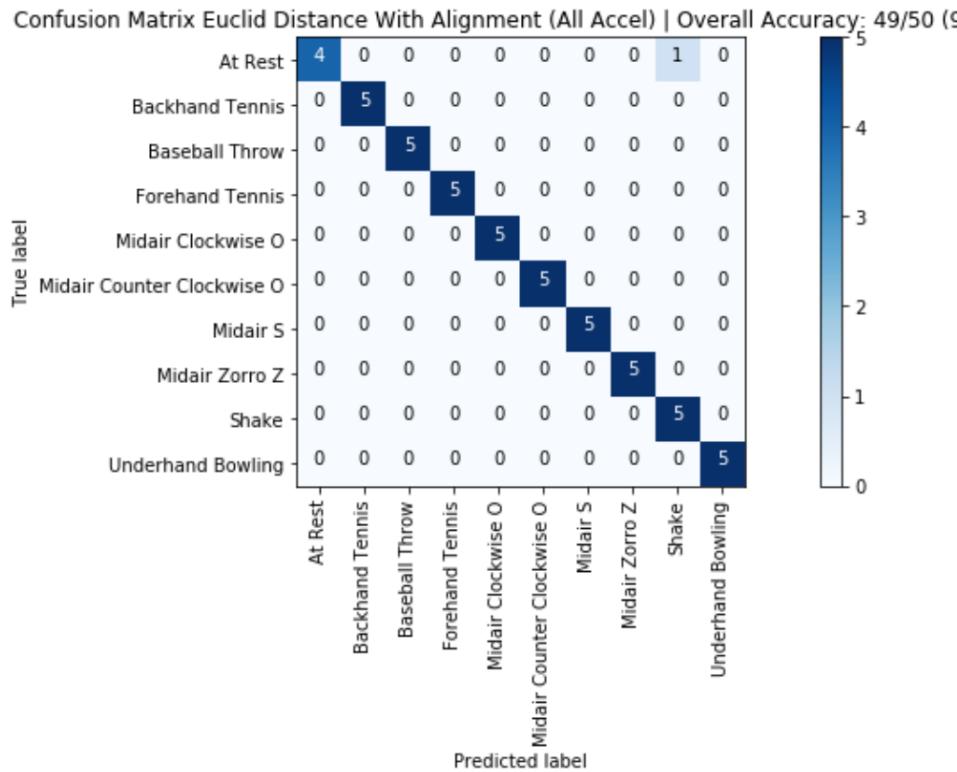
True label



ML EVALUATION

CONFUSION MATRICES

Using Combination of Signals



N-BEST LISTS

For algorithms that provide per item classification scores, we can compute an N-best list

N-best lists depict where the “right” answer is in an ordered list of matching scores between items

Change your matching functions to return a sorted list of (matches, scores) rather than just the (best match trial, score)

N-BEST LIST EXAMPLE CODE

```
# Just returns the best match based on euclidean distance
def find_closest_match_euclid_distance(testTrial, templateFolds, **kwargs):
    minDistance = float('inf')
    bestMatchTrial = None
    for fold in templateFolds:
        for foldGestureName, foldTrial in fold.items():
            # Calculate Euclidean distance
            euclidDistance = distance.euclidean(testTrial.signal, foldTrial.signal)
            if euclidDistance < minDistance:
                minDistance = euclidDistance
                bestMatchTrial = foldTrial

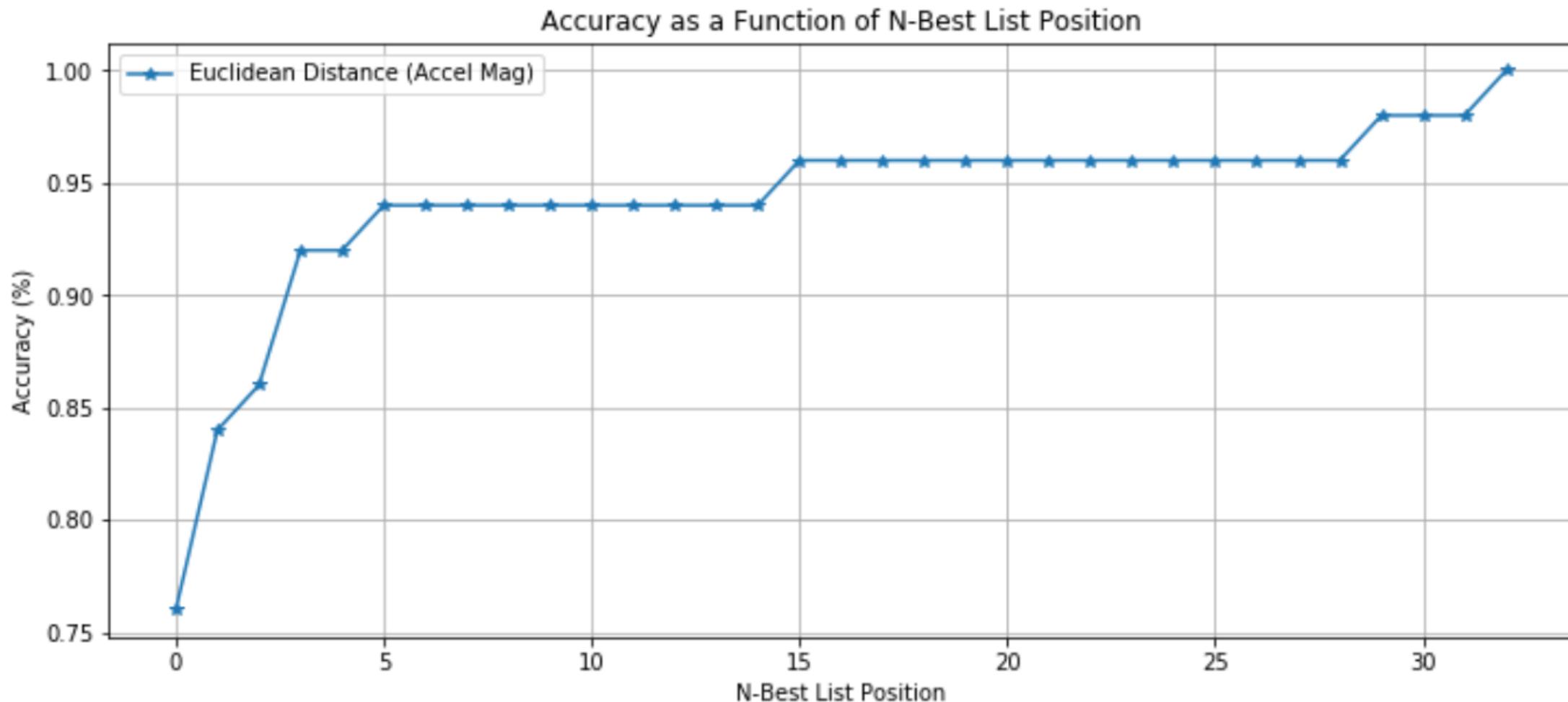
    # Return the best match
    return (bestMatchTrial, minDistance)

# Returns an nBestList tuple of (foldTrial, score) based on euclidean distance
def find_closest_match_euclid_distance(testTrial, templateFolds, **kwargs):
    # Tuple list of trial to score
    nBestList = list()
    for fold in templateFolds:
        for foldGestureName, foldTrial in fold.items():
            # Euclidean distance
            euclidDistance = distance.euclidean(testTrial.signal, foldTrial.signal)
            nBestList.append((foldTrial, euclidDistance))

    # sort the nBestList by score
    nBestList.sort(key=lambda x: x[1])
    return nBestList
```

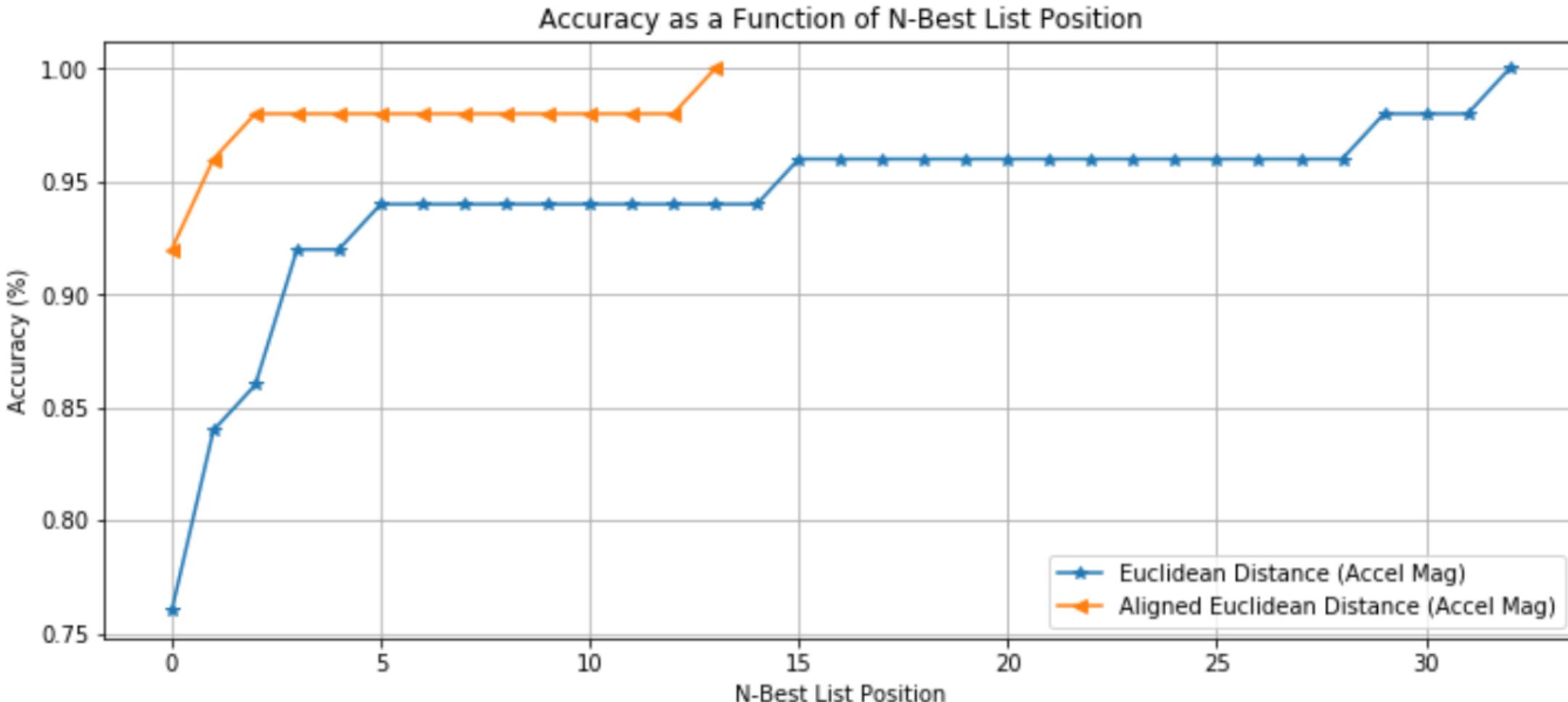
N-BEST LIST: SHAPE-MATCHING GESTURE REC

This is on the JonGestureLogs from <https://github.com/jonfroehlich/CSE590Sp2018/tree/master/A02-GestureRecognizer>



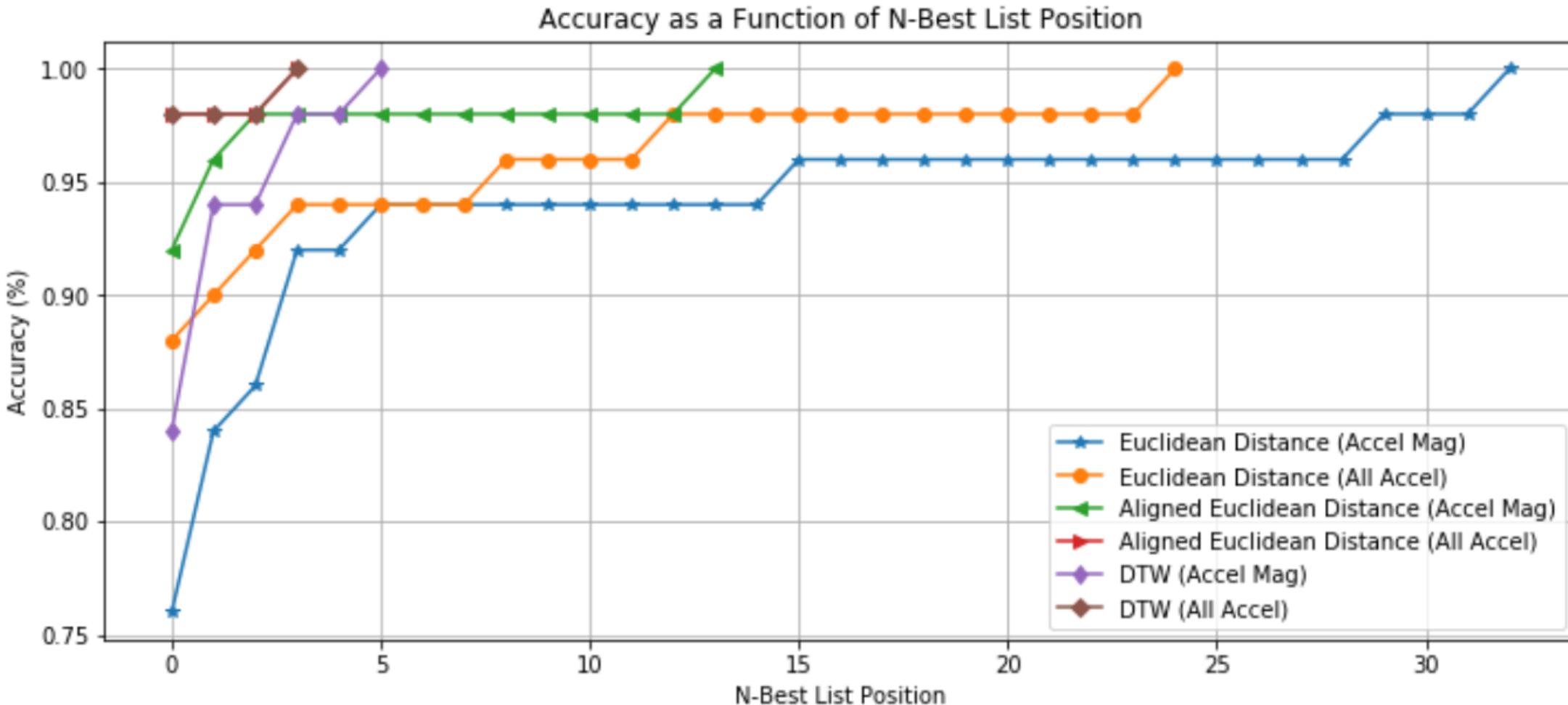
N-BEST LIST: SHAPE-MATCHING GESTURE REC

This is on the JonGestureLogs from <https://github.com/jonfroehlich/CSE590Sp2018/tree/master/A02-GestureRecognizer>



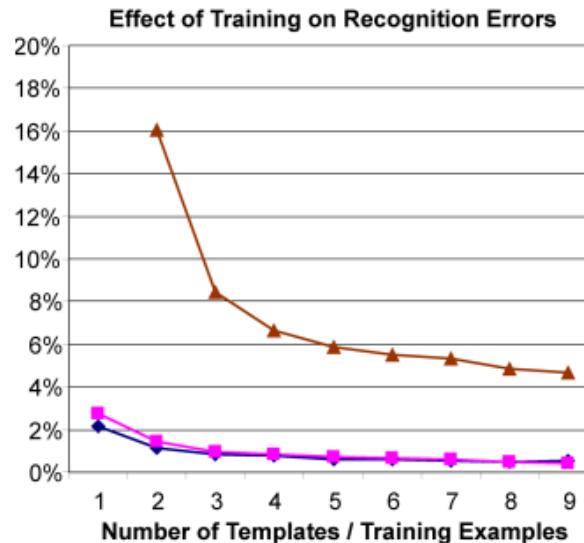
N-BEST LIST: SHAPE-MATCHING GESTURE REC

This is on the JonGestureLogs from <https://github.com/jonfroehlich/CSE590Sp2018/tree/master/A02-GestureRecognizer>

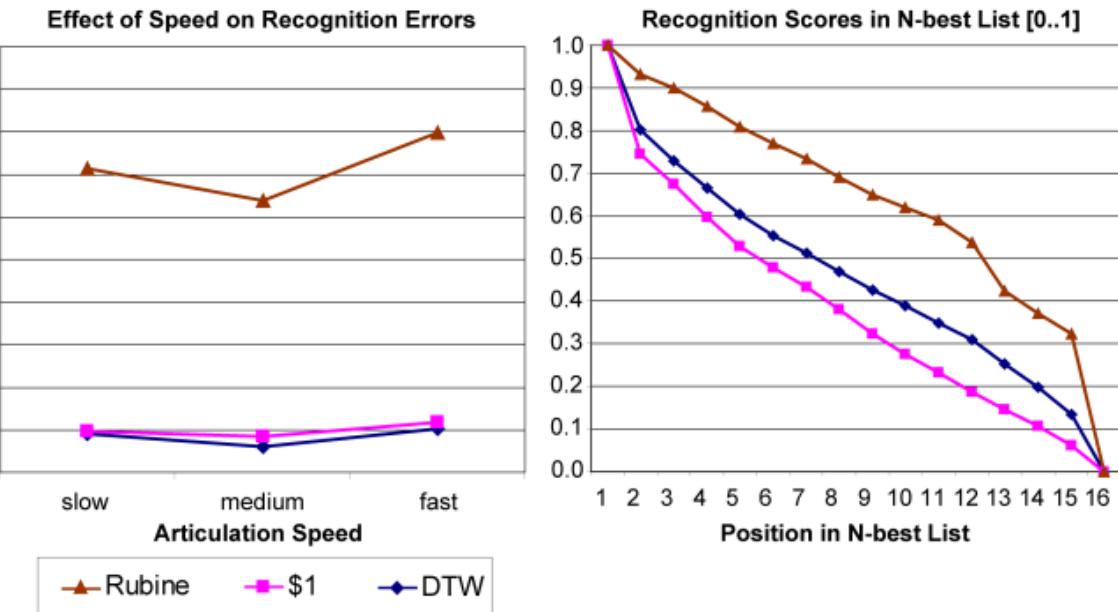
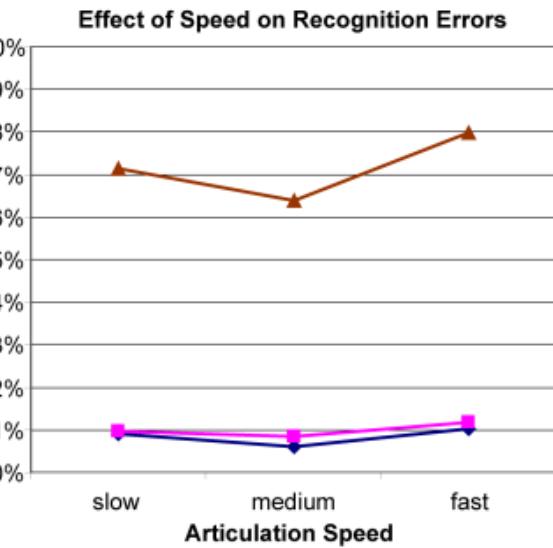


ML EVALUATION

N-BEST LIST EXAMPLE 2



Graphing average recognition scores as a function of N-best list index. Prefer a dramatic dropoff between 1st and 2nd position, indicating a good separation of scores.



Gestures without libraries, toolkits or training: a \$1 recognizer for user interface prototypes
 Wobbrock *et al.*, UIST'07
 744 citations

ML TOOLKITS: LANGUAGES AND LIBRARIES

The screenshot shows the scikit-learn homepage. At the top, there's a navigation bar with links for Home, Installation, Documentation, Examples, Google Custom Search, and a search icon. Below the header, there's a large banner for "scikit-learn Machine Learning in Python". The main content area is divided into several sections: Classification (with examples like SVM, KNN, Random Forest), Regression (with examples like Ridge Regression, Lasso), Clustering (with examples like k-Means, Spectral Clustering), Dimensionality reduction (PCA), Model selection (grid search, cross-validation), and Preprocessing (feature extraction, normalization). Each section includes a brief description, applications, and algorithms.

Sci-Kit Learn

Python-based library for classification, regression, clustering, dimensionality reduction, model selection, preprocessing data, etc.

The screenshot shows the WEKA website. At the top, there's a navigation bar with links for Project, Software (which is currently selected), Book, Courses, Publications, People, and Related. The main content area features a heading "Weka 3: Data Mining Software in Java". It includes a brief introduction about Weka being a collection of machine learning algorithms for data mining tasks, its availability as open source under the GNU General Public License, and its use in several online courses. Below this, there are three columns: Getting started (Requirements, Download, Documentation, FAQ, Getting Help), Further information (Citing Weka, Datasets, Related Projects, Miscellaneous Code, Other Literature), and Developers (Development, History, Subversion, Contributors, Commercial licenses).

Weka

Java-based library for pre-processing, classification, regression, clustering, and visualizing data. Used to be the standard. People now shifting to sci-kit learn, R, and cloud-based APIs.

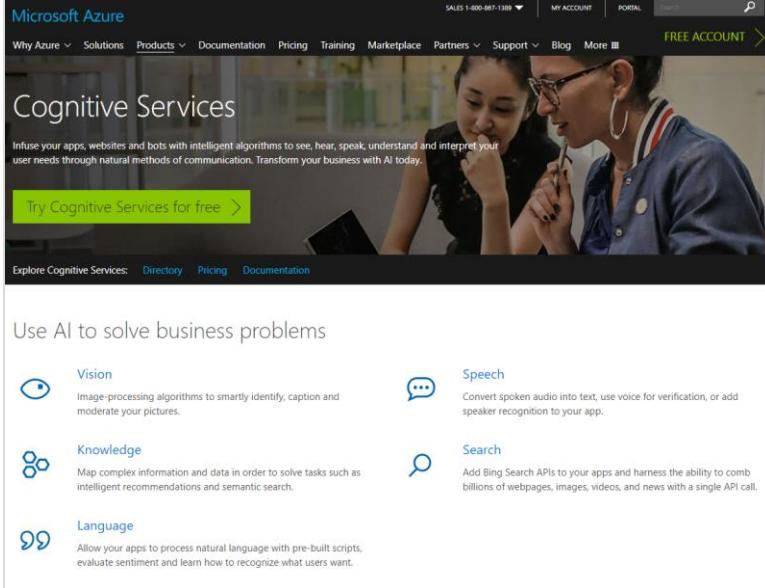
The screenshot shows the R Project for Statistical Computing website. At the top, there's a navigation bar with links for Home (selected), Download, CRAN, Project, About R, Logo, Contributors, What's New?, Reporting Bugs, Development Site, Conferences, Search, Foundation, Board, Members, Donors, Donate, Help With R, Getting Help, Documentation, Manuals, FAQs, The R Journal, and Books. The main content area has a heading "The R Project for Statistical Computing". Below it, there's a "Getting Started" section with a brief introduction to R as a free software environment for statistical computing and graphics. It also includes news items about recent releases and events like the R Conference in London.

R

R arose as a statistical computing language but is now widely used in ML and has a number of well-respected ML packages.

ML TOOLKITS

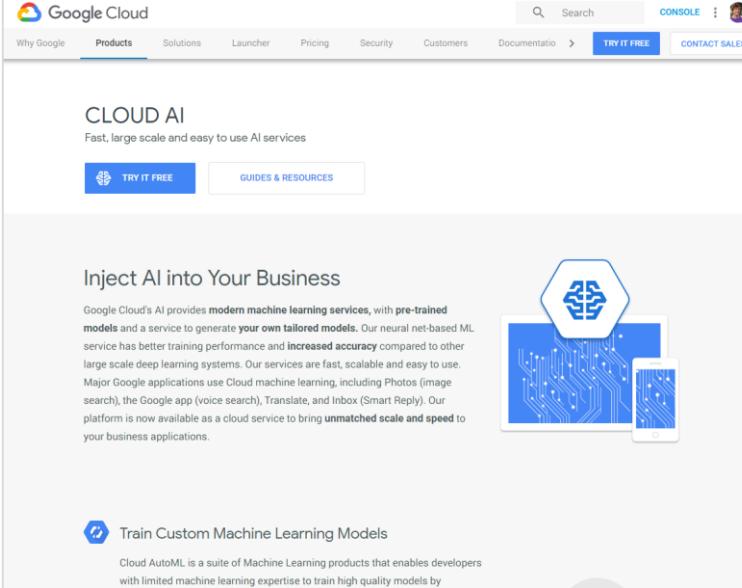
ML TOOLKITS: CLOUD APIs



The screenshot shows the Microsoft Azure Cognitive Services page. At the top, there's a navigation bar with links like 'SALES 1-800-887-1389', 'MY ACCOUNT', 'PORTAL', 'Why Azure', 'Solutions', 'Products', 'Documentation', 'Pricing', 'Training', 'Marketplace', 'Partners', 'Support', 'Blog', 'More', and a 'FREE ACCOUNT' button. Below the navigation is a large image of two women looking at a laptop screen. The text 'Cognitive Services' is prominently displayed. A sub-section titled 'Infuse your apps, websites and bots with intelligent algorithms to see, hear, speak, understand and interpret your user needs through natural methods of communication. Transform your business with AI today.' includes a 'Try Cognitive Services for free' button. Further down, there's a section titled 'Use AI to solve business problems' with four categories: Vision (image processing), Speech (audio-to-text), Knowledge (complex data processing), and Language (natural language processing). Each category has a brief description and a small icon.

Microsoft Cognitive Services

Cloud API for computer vision, NLP, speech recognition, knowledge representation, chatbots, and search.



The screenshot shows the Google Cloud AI page. At the top, there's a navigation bar with links like 'Why Google', 'Products', 'Solutions', 'Launcher', 'Pricing', 'Security', 'Customers', 'Documentation', 'TRY IT FREE', and 'CONTACT SALES'. Below the navigation is a section titled 'CLOUD AI' with the sub-sub-section 'Inject AI into Your Business'. It describes Google Cloud's AI services, mentioning pre-trained models and the ability to generate tailored models. It also highlights the use of Cloud machine learning in various Google applications like Photos, Translate, and Inbox. There's a visual element showing a brain icon inside a hexagon above a smartphone and a laptop. Another section below is 'Train Custom Machine Learning Models' with a brief description of Cloud AutoML.

Google Cloud AI

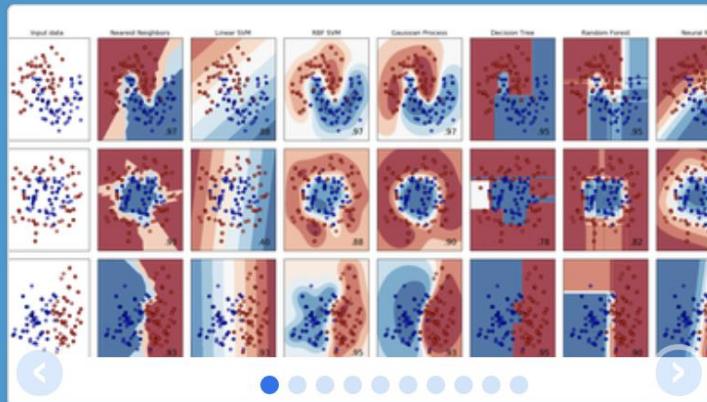
Cloud API for large-scale machine learning, including video and image analysis, speech recognition, speech synthesis, NLP, and language translation



The screenshot shows the AWS Machine Learning on AWS page. At the top, there's a navigation bar with links like 'Menu', 'AWS', 'Contact Sales', 'Products', 'Solutions', 'More', 'English', 'My Account', and 'Create an AWS Account'. Below the navigation is a section titled 'Machine Learning on AWS' with the sub-sub-section 'Why machine learning on AWS?'. It discusses Amazon's long history in AI and its impact across various internal systems and customer experiences. It emphasizes the company's mission to share learnings and ML capabilities as fully managed services. There are four columns under 'Why machine learning on AWS?': 'Machine Learning for everyone', 'API-driven ML services', 'Broad framework support', and 'Breadth of compute options'. Each column has a brief description and a small icon.

Amazon AWS ML Platform

Cloud API for general ML with image and video analysis, chatbots, and language services such as NLP, translation, transcription, and speech synthesis.



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ...

— Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ...

— Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ...

— Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization.

— Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics.

— Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction.

— Examples

SCIKIT-LEARN FEATURES INCLUDE

FEATURE EXTRACTION, SELECTION, & REDUCTION

Feature Extraction: helper methods for extracting features from image and text data.

Feature Selection: for identifying meaningful attributes from which to create supervised models.

Dimensionality Reduction: for reducing the number of features (e.g., Principal component analysis)

ML MODELS

Clustering Algorithms: for grouping unlabeled data such as KMeans.

Supervised Models: a vast array of supervised models, including naive bayes, neural networks, SVMs and decision trees.

Ensemble Methods: for combining the predictions of multiple supervised models.

EVALUATION

Cross Validation: for estimating the performance of supervised models on unseen data.

Parameter Tuning: for getting the most out of supervised models.

Datasets: for test datasets and for generating datasets with specific properties for investigating model behavior.

GETTING STARTED

This screenshot shows the 'An introduction to machine learning with scikit-learn' page. It features a navigation bar at the top with links for Home, Installation, Documentation, Examples, Google Custom Search, and a search icon. Below the navigation is a sidebar with links for Previous, Next, Up, and Other versions. A prominent yellow call-to-action button says 'Please cite us if you use the software.' The main content area has a title 'An introduction to machine learning with scikit-learn' and a 'Section contents' section. It includes a note about the machine learning vocabulary used throughout the library and a simple learning example. A 'Machine learning: the problem setting' section defines learning problems like classification, regression, and clustering. On the left, there's a sidebar with general information about scikit-learn, including its version (v0.19.1), installation, documentation, examples, and a link to the GitHub repository.

This screenshot shows the 'User Guide' page. The navigation bar and sidebar are identical to the first page. The main content area has a title 'User Guide' and a section titled '1. Supervised learning'. This section lists 16 sub-topics: Generalized Linear Models, Linear and Quadratic Discriminant Analysis, Kernel ridge regression, Support Vector Machines, Stochastic Gradient Descent, Nearest Neighbors, Gaussian Processes, Cross decomposition, Naive Bayes, Decision Trees, Ensemble methods, Multiclass and multilabel algorithms, Feature selection, Semi-Supervised, Isotonic regression, and Probability calibration. A 'User Guide' sidebar on the left provides links to other sections of the User Guide.

This screenshot shows the 'Examples' page. The navigation bar and sidebar are identical to the previous pages. The main content area has a title 'Examples' and a section titled 'General examples'. It describes general-purpose and introductory examples for the scikit-learn library. Below this, there are several examples with small images: 'Plotting Cross-Validated Predictions' (a scatter plot with a fitted curve), 'Concatenating multiple feature extraction methods' (three overlapping colored circles), and 'Pipelining: chaining a PCA and a logistic regression' (a graph showing a non-linear decision boundary). The sidebar on the left lists various examples such as Covariance estimation, Cross decomposition, Dataset examples, Decomposition, Ensemble methods, Feature Selection, Gaussian Process for Machine Learning, Generalized Linear Models, Manifold learning, Gaussian Mixture Models, Model Selection, Multicategory methods, Nearest Neighbors, Neural Networks, Preprocessing, Semi Supervised Classification, and Support Vector Machines.

Quick Start Tutorial

Introduces some basic vocabulary and walks through a simple example

User Guide

Goes through each major component of the sci-kit library from supervised and unsupervised models to model selection and evaluation.

Examples

Lots of different examples of using, visualizing, and evaluating models

SVM

SVM is very similar to the Apple /
Oranges and Dog / Cat example

Linear classifier

Can be slow to train but fast to use

We'll watch a short video

Support Vector Machines

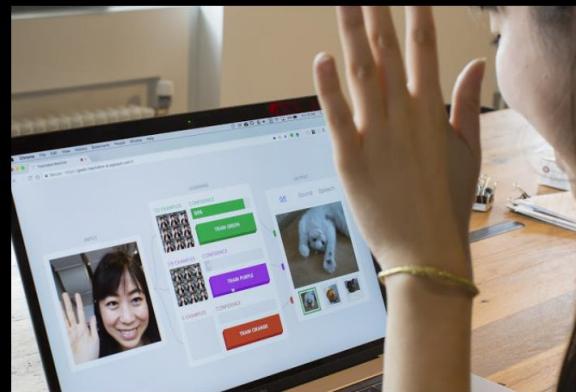
A Visual Explanation + Sample Python Code

<https://github.com/adashofdata>

Explore machine learning in simple, hands-on ways.

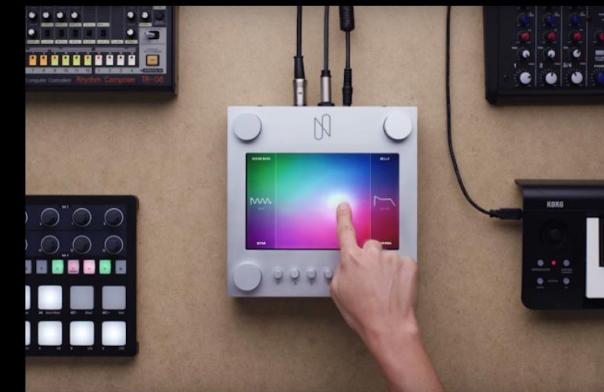
[Read more](#)

FEATURED EXPERIMENTS



TEACHABLE MACHINE

by Google Creative Lab

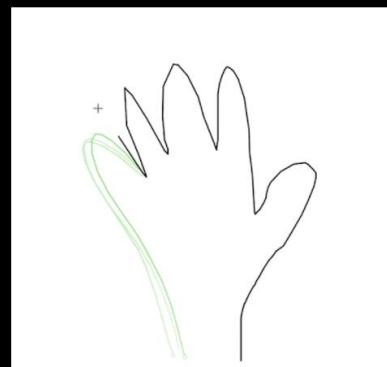


NSYNTH SUPER

by Magenta / Google Creative Lab

BROWSE TAGS ▾

SORT BY: NEWEST ▾



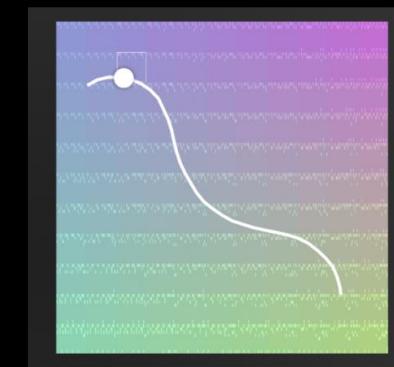
SCRYING PEN

by Andy Matuschak



EMOJI SCAVENGER HUNT

by Google Brand Studio



BEAT BLENDER

by Creative Lab + Magenta

RESOURCES

FIND OUT MORE: ML RESOURCES

BASIC OVERVIEWS

[Machine Learning 101](#)

Jason Mayes, Senior Creative Engineer Google
A intro slidedeck that overviews ML approaches and methods

[The 7 Steps of Machine Learning](#)

Yufeng Guo, Google Cloud
A short article and video about the ML pipeline

COURSES / IN-DEPTH GUIDES

[Machine Learning Mastery](#)

Jason Brownlee, PhD
A free step-by-step guide on many aspects of ML.

[Supervised Learning with Sci-kit Learn](#)

Datacamp.com
A paid course (1st chapter free) on using sci-kit learn

[Coursera's Machine Learning Course](#)

Andrew Ng, former Stanford Prof & Google Brain head
The most popular MOOC on ML of all-time.

MATHEMATICS OF ML

[Machine Learning: A Probabilistic Perspective](#)

Kevin Murphy, former UBC Prof, now at Google
An in-depth introduction to probabilistic and inference models

[Elements of Statistical Learning](#)

Hastie, Tibshirani, & Friedman
Recommended by Machine Learning Mastery