



@dhianadeva

MACHINE LEARNING FOR EVERYONE

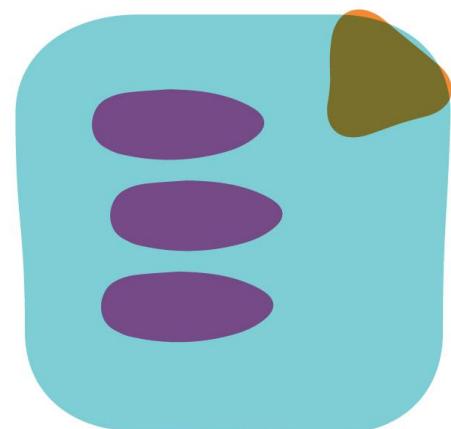
Demystifying machine learning!

AGENDA

Goal:

Encourage you to start a machine learning project. Today!

- About me
- About you
- Machine Learning
- Problems
- Design
- Algorithms
- Evaluation
- Code snippets
- Pay-as-you-go
- Competitions



ThoughtWorks®

ABOUT ME

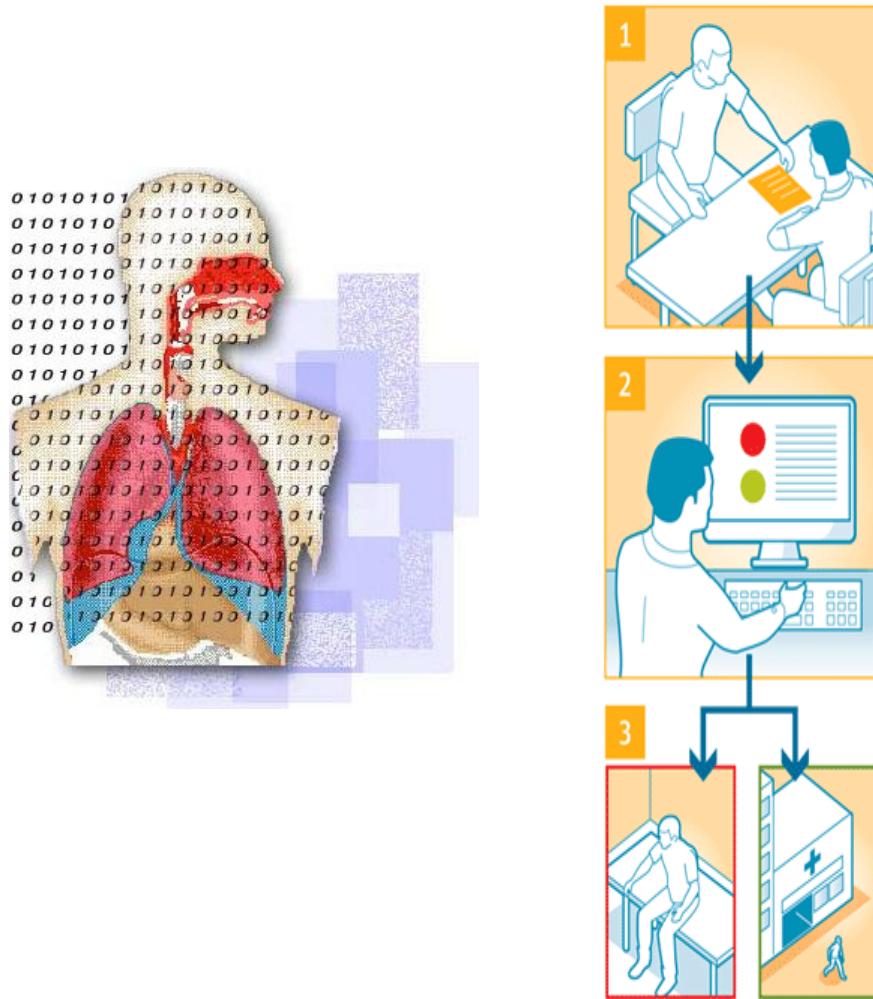
Electronics Engineering, Software Development and Data Science... Why not?

ThoughtWorks®

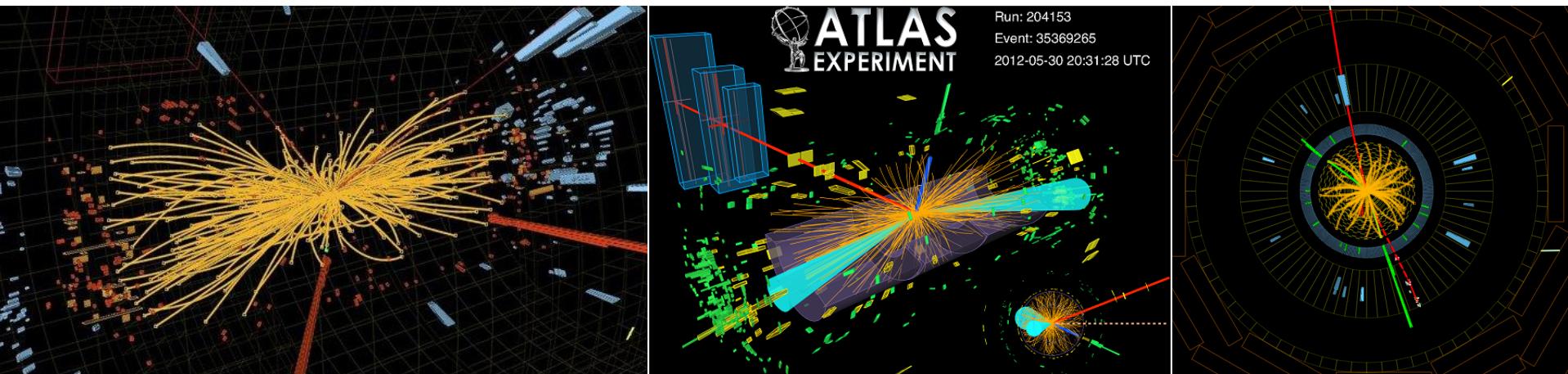
A photograph of a woman with long brown hair, wearing sunglasses and a dark jacket, sitting on a low stone wall. She is smiling and has her right arm raised in a wave. Behind her are the large stone pyramids of Teotihuacan under a dramatic, cloudy sky.

DHIANA DEVA

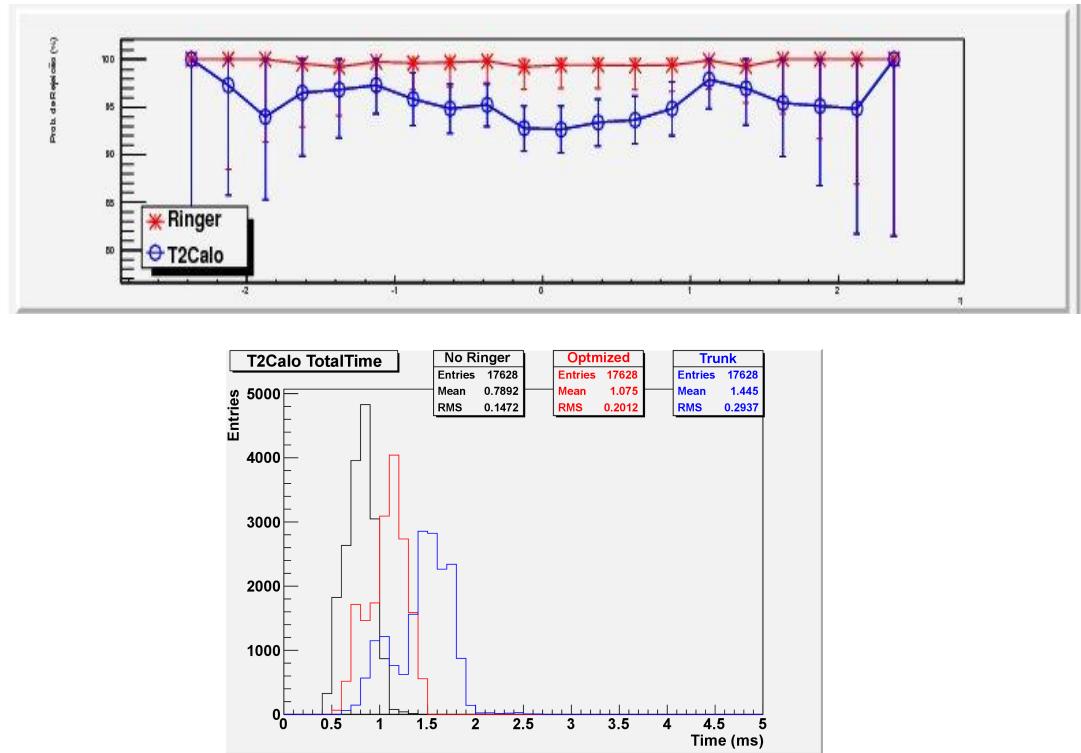
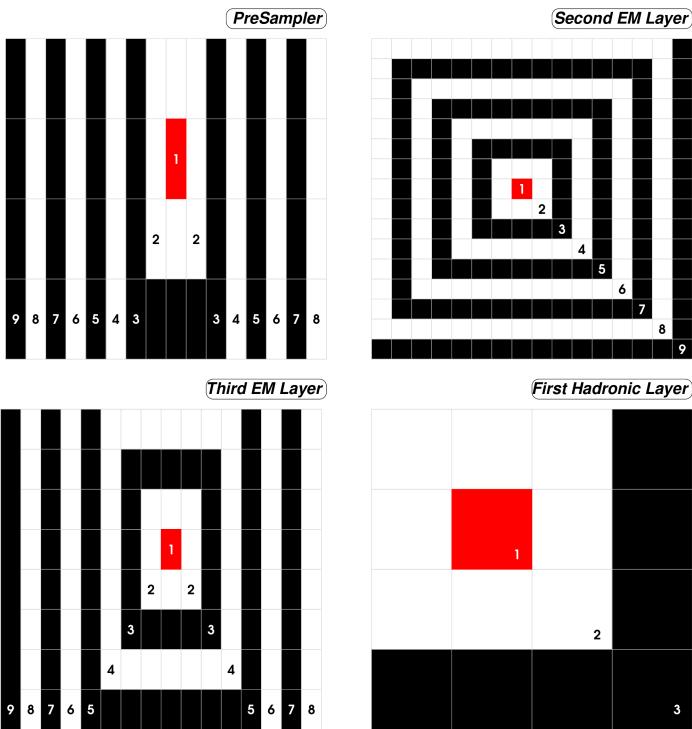
NEURALTB



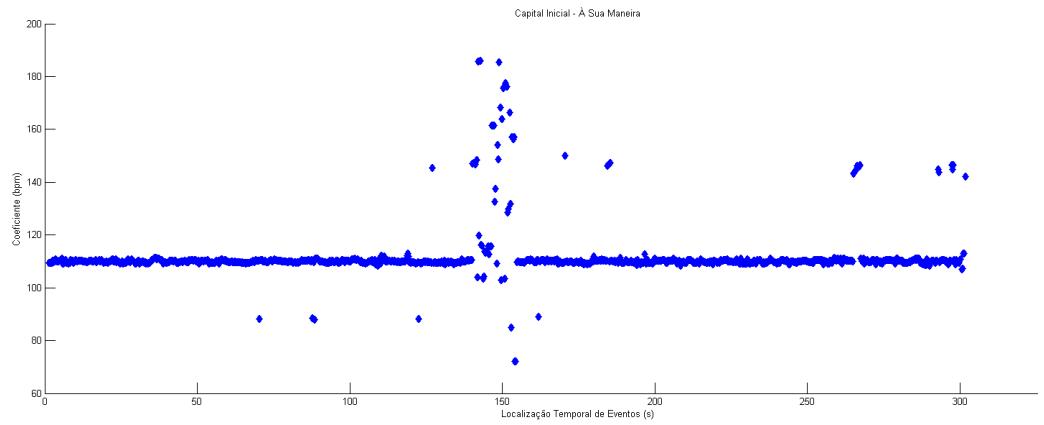
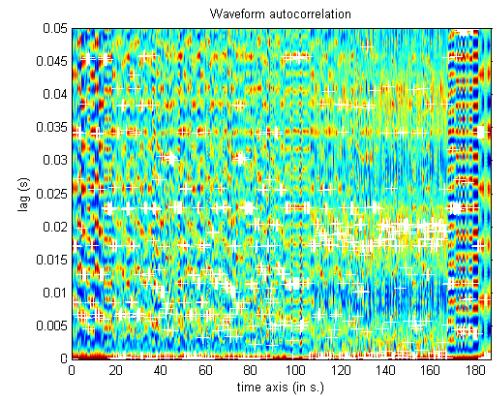
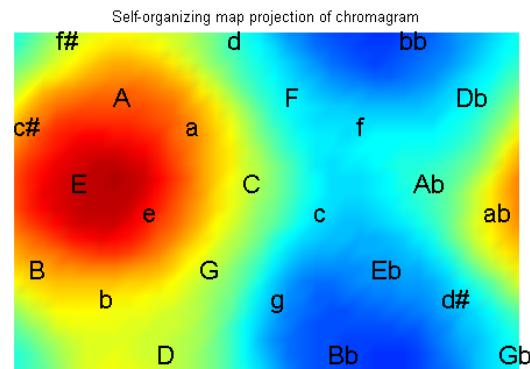
CERN



NEURALRINGER



DJBRAZIL



HIGGS CHALLENGE

Winner announcement

The Higgs Machine Learning Challenge has completed the 15th September 2014, gathering 1785 teams and 1942 participants.

After due verifications, we are pleased to announce the three winners, with the three best scores on the private leaderboard when disclosed.

1 : Gabor Melis : 7000 dollars

2 : Tim Salimans : 4000 dollars

3 : Pierre Courtiol (nhlx5haze) : 2000 dollars

All three have been invited (at the HiggsML organisation expense) to NIPS conference at Montreal, where a special workshop is organised the 13th December 2014, to discuss machine learning techniques applications to high energy physics, and specific developments made for this challenge.

<https://nips.cc/Conferences/2014/Program/event.php?ID=4292>

In addition, documented software was scrutinized, and the special HEP meets ML award is given to :

crowwork (Tianqi Chen and Tong He)

They have developed XGBoost <https://github.com/tqchen/xgboost> and made it available to other participants early, and it was indeed used by many of them ; while not giving the very best score, it appears to be an excellent compromise between performance and simplicity, which makes it a promising improvements to tools currently used by high energy physicists.

The team will be invited at CERN in 2015 for a workshop (being organised) where machine learning techniques application to high energy physics, in particular as they emerged in this challenge, will be discussed further.

The organizers would like to make special mention of CSE_TEAM_0 (Chamila Wijayarathna, Dimuthu Upeksha, Maduranga Siriwardena, Sachith Withana) for the detailed documentation of their optimisation, and Dhiana Deva an enthusiastic undergraduate.

ThoughtWorks®

ABOUT YOU

You can do it!

FOR ALL

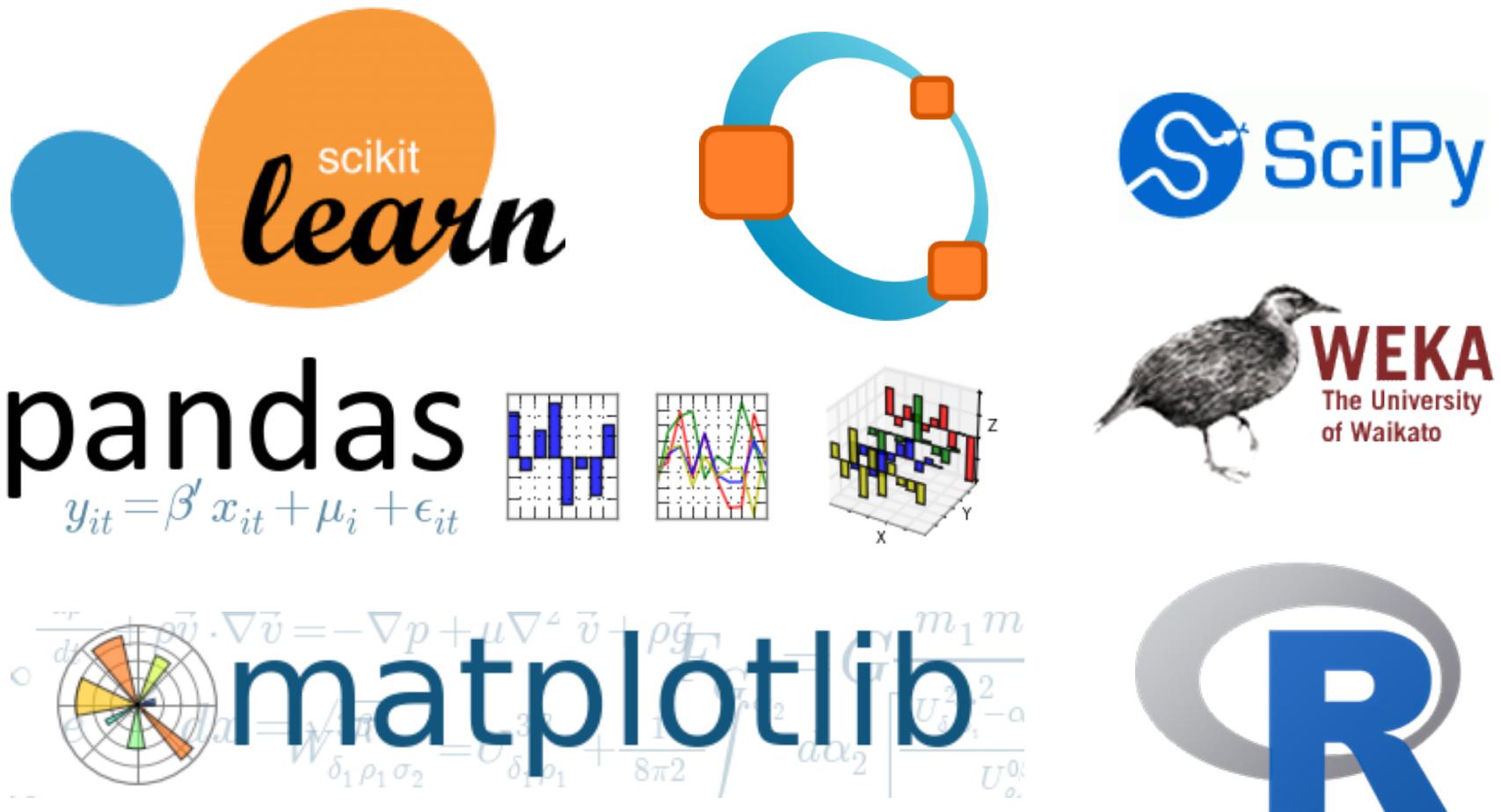


MASSIVE ONLINE OPEN COURSES

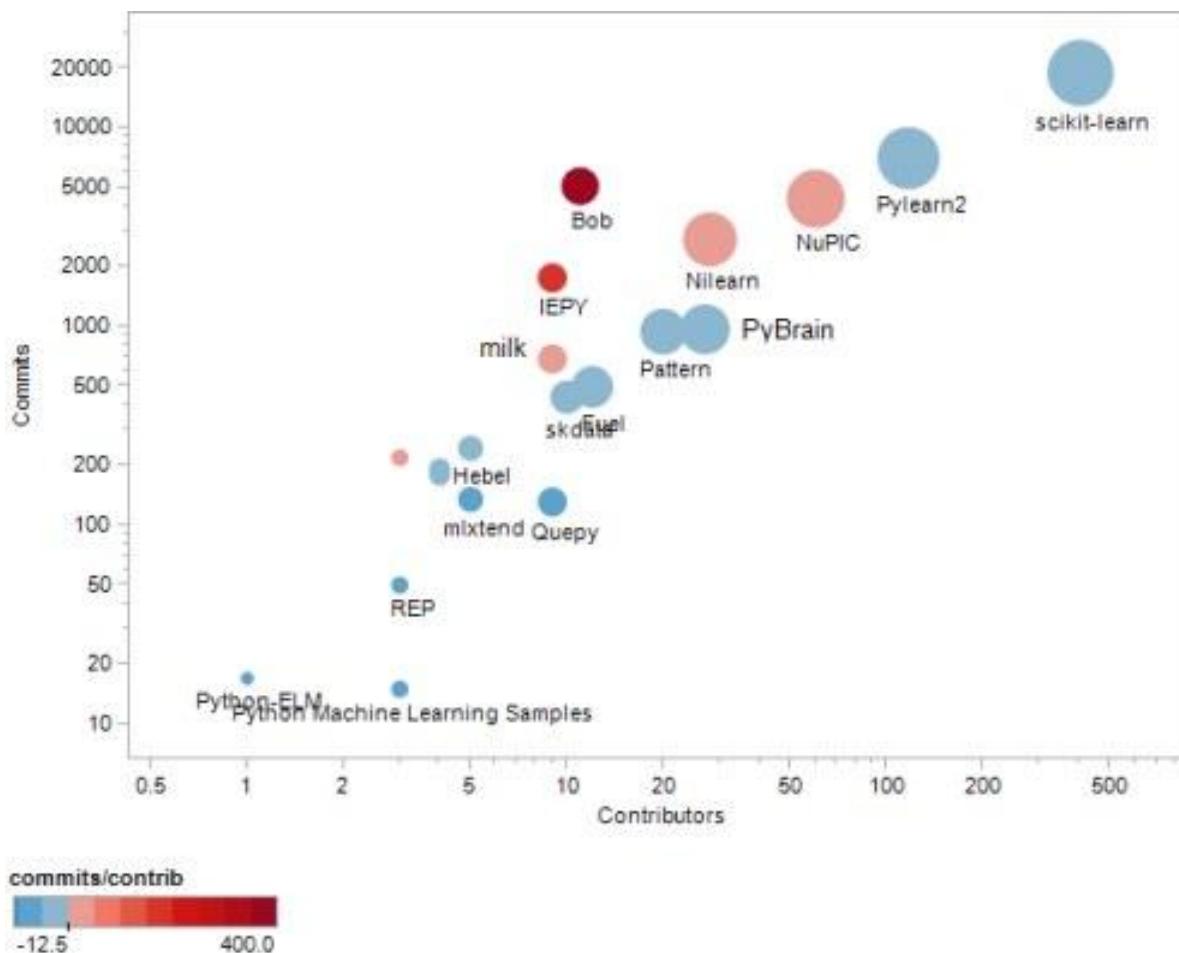
U
UDACITY

coursera
edX

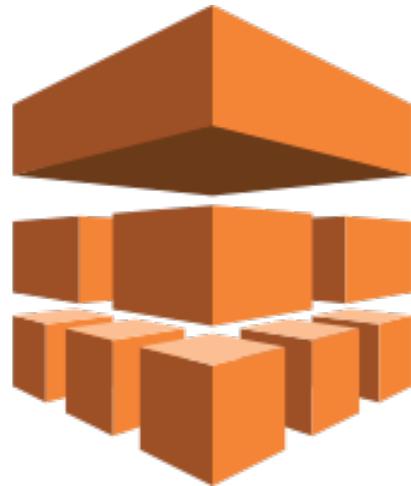
OPEN SOURCE TOOLS



OPEN SOURCE PYTHON TOOLS



PAY-AS-YOU-GO SERVICES



Google Cloud Platform

ThoughtWorks®

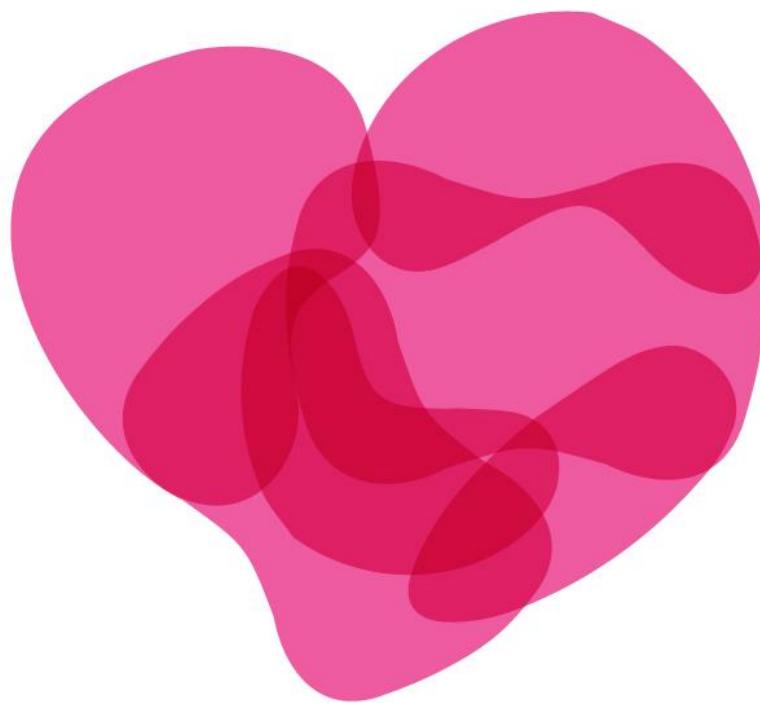
MACHINE LEARNING

Learning, machine learning!

EXPECTATIONS



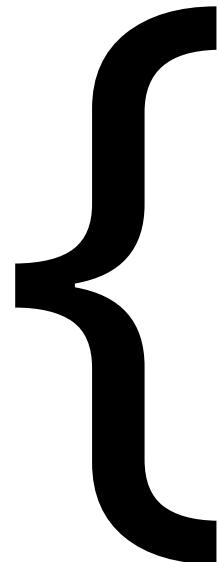
REALITY



FEATURE EXTRACTION

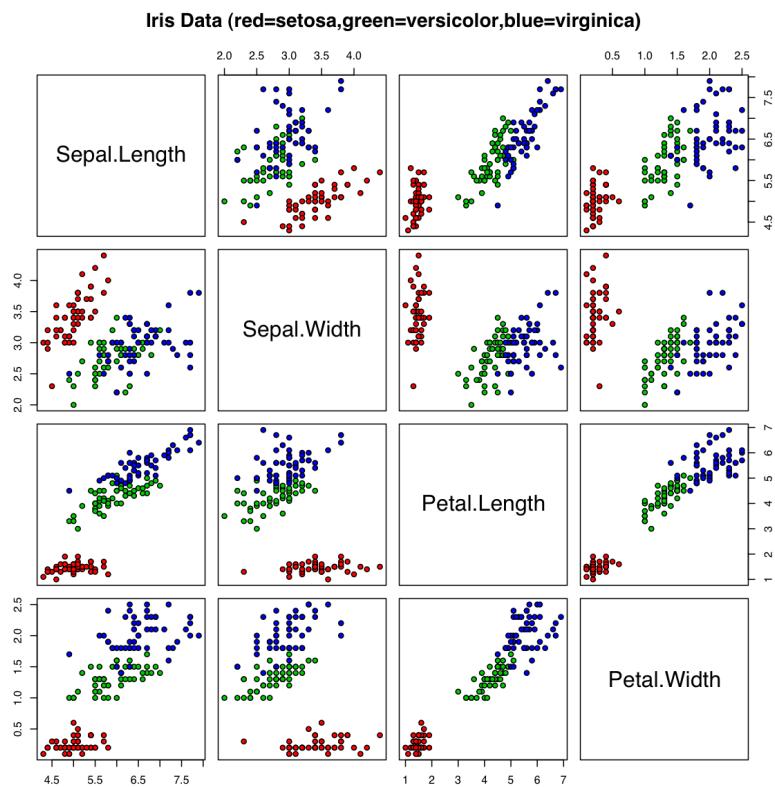
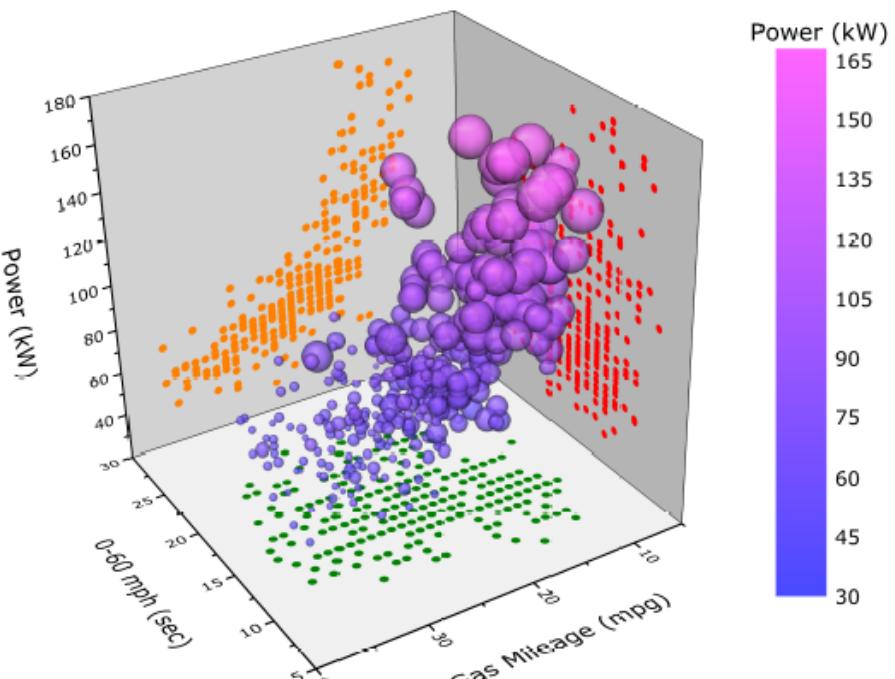


Item

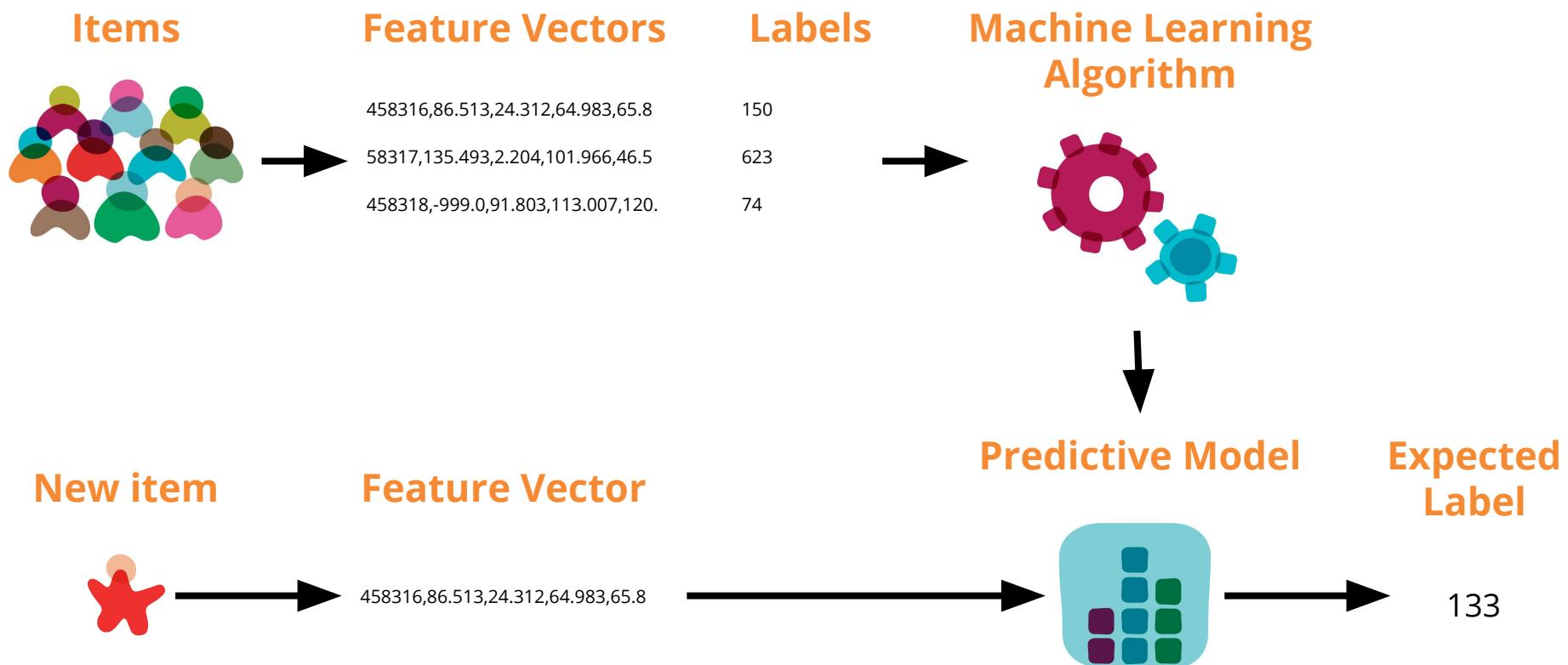


Feature 1
Feature 2
...
Feature N

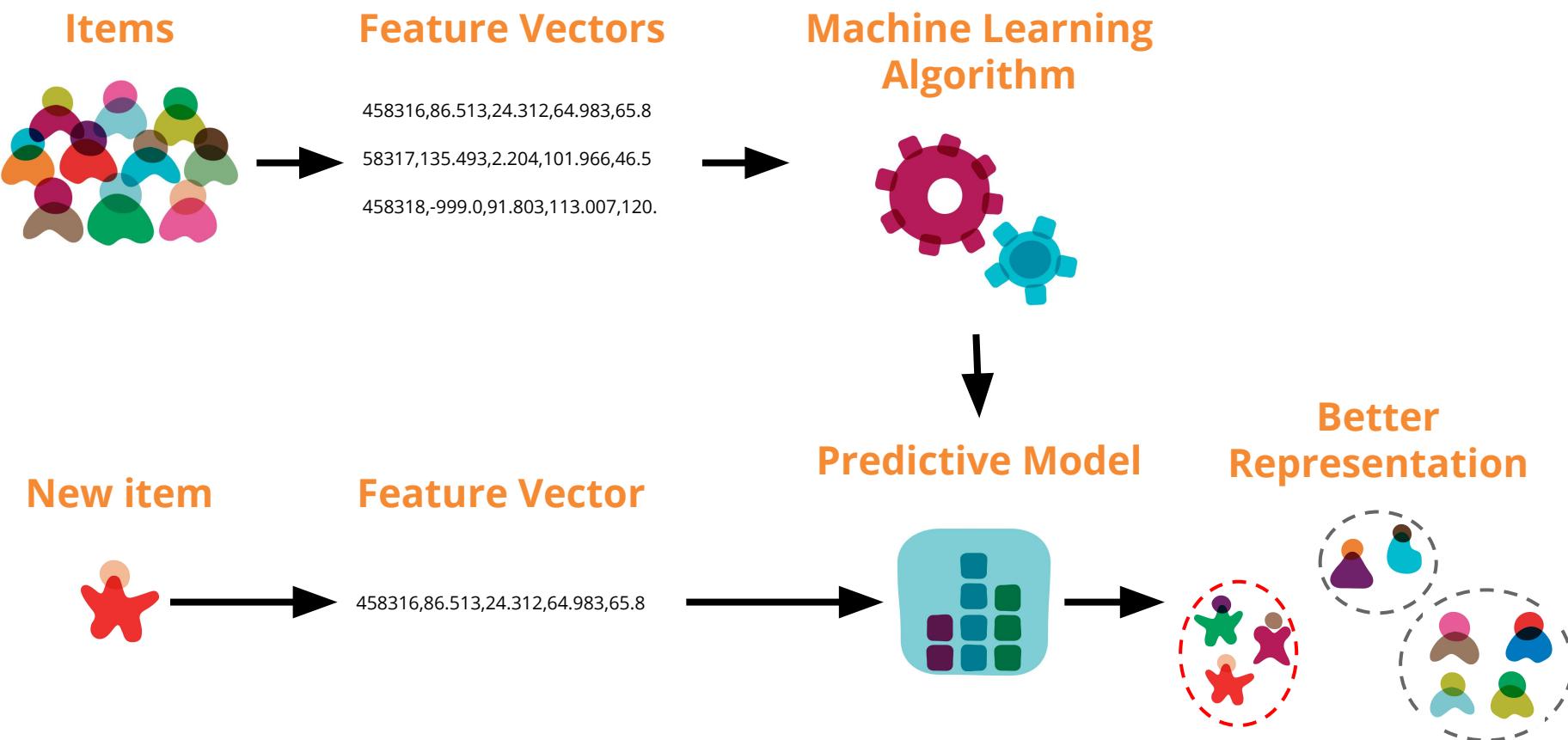
FEATURE SPACE



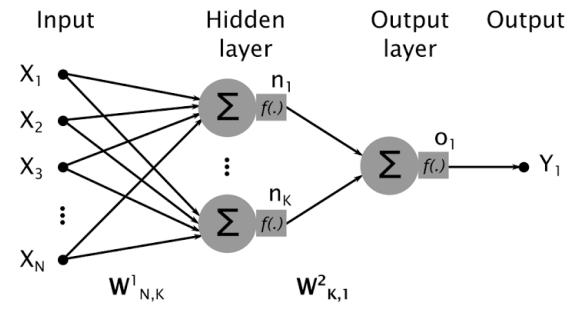
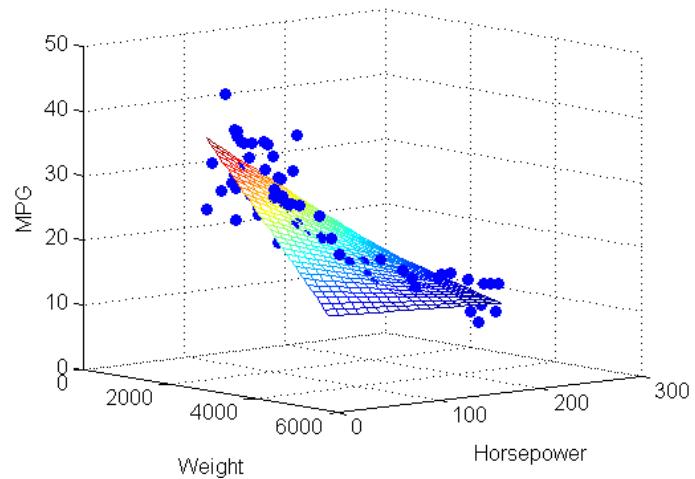
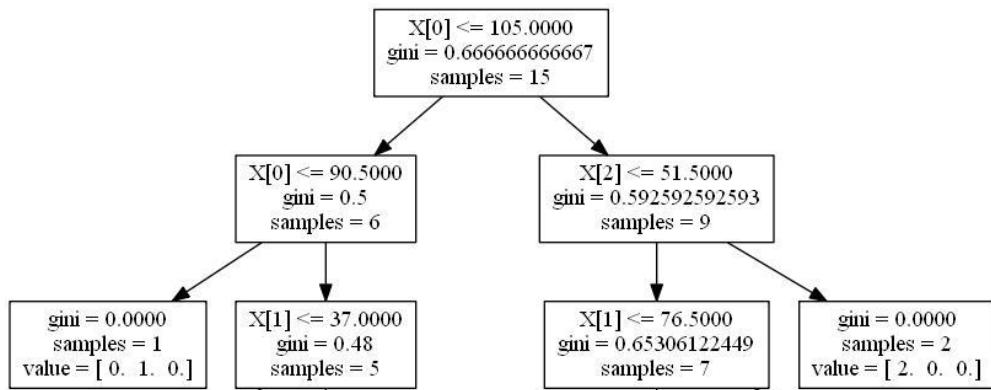
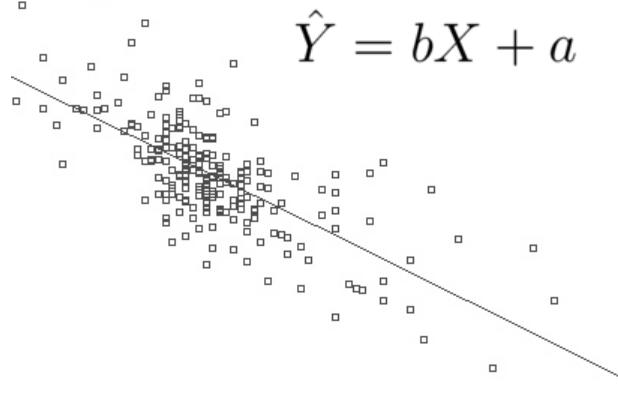
SUPERVISED LEARNING



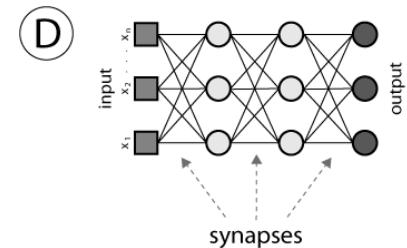
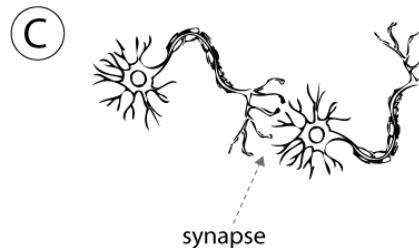
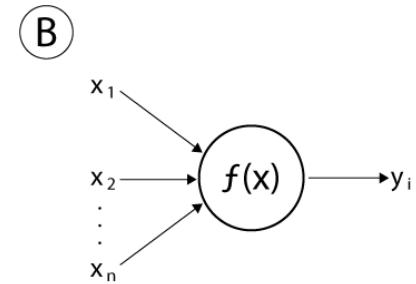
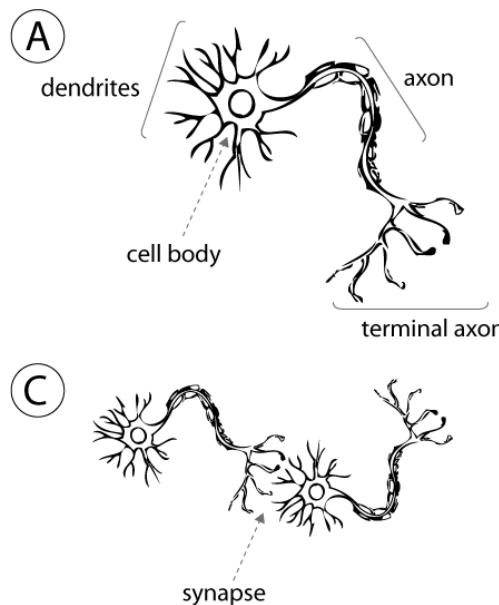
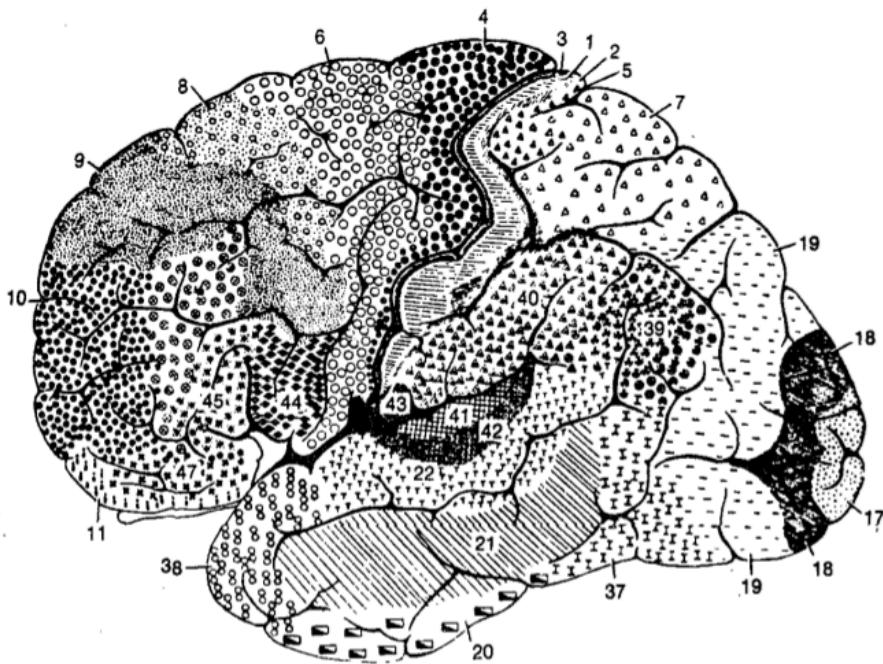
UNSUPERVISED LEARNING



MODELS



BIOLOGICAL MOTIVATION





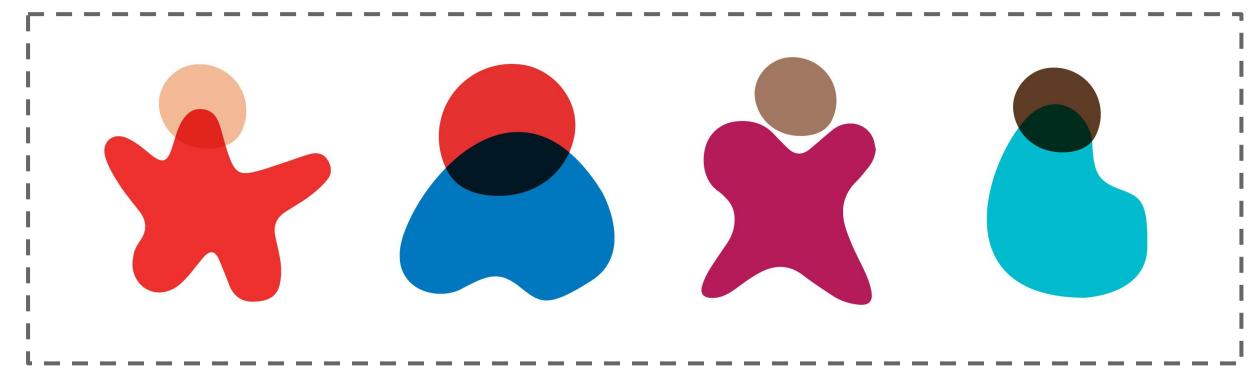
PROBLEMS

I've got 99 problems, but machine learning ain't one!

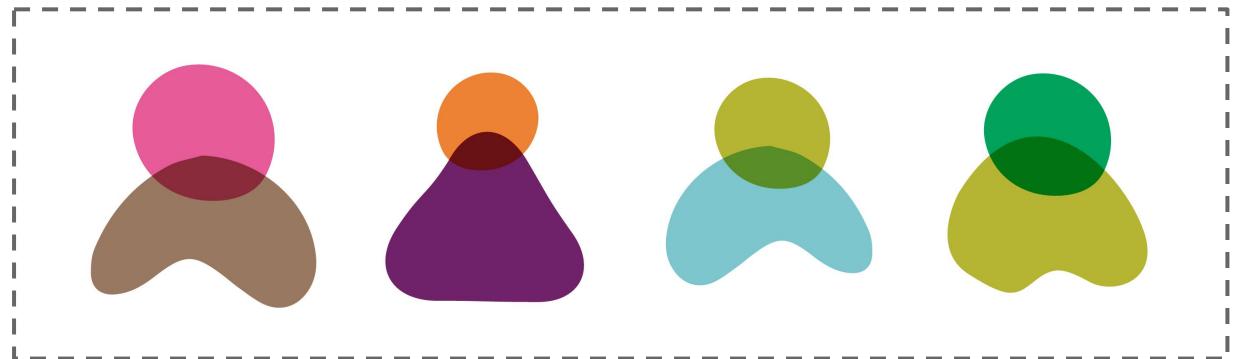
CLASSIFICATION



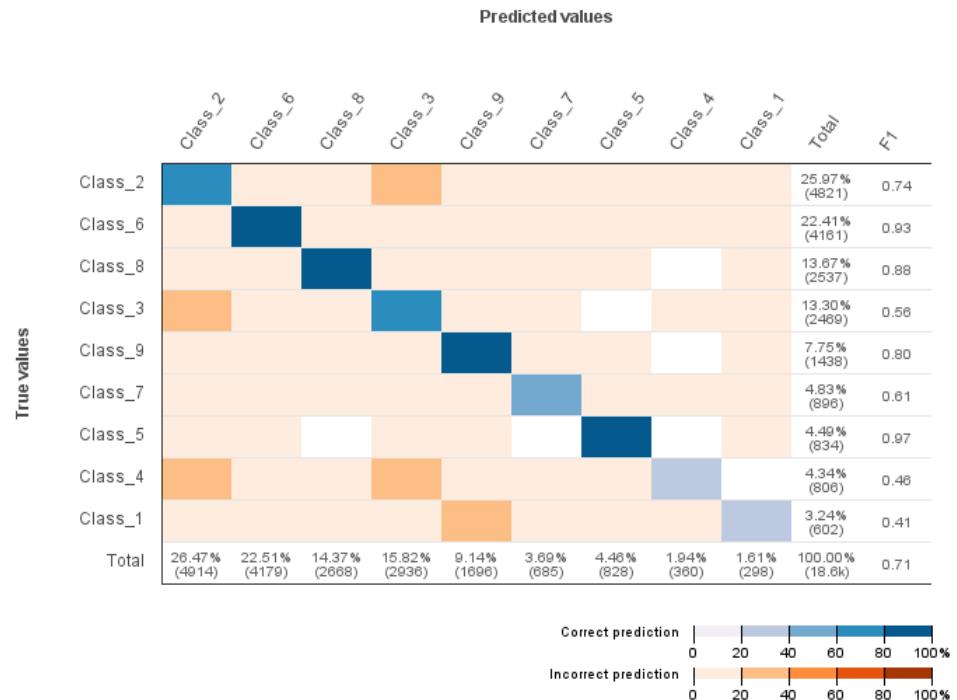
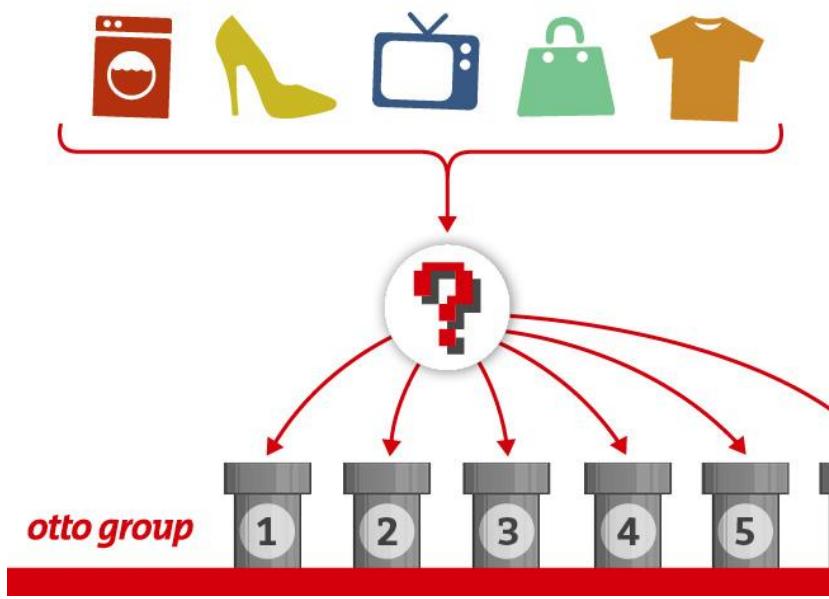
A



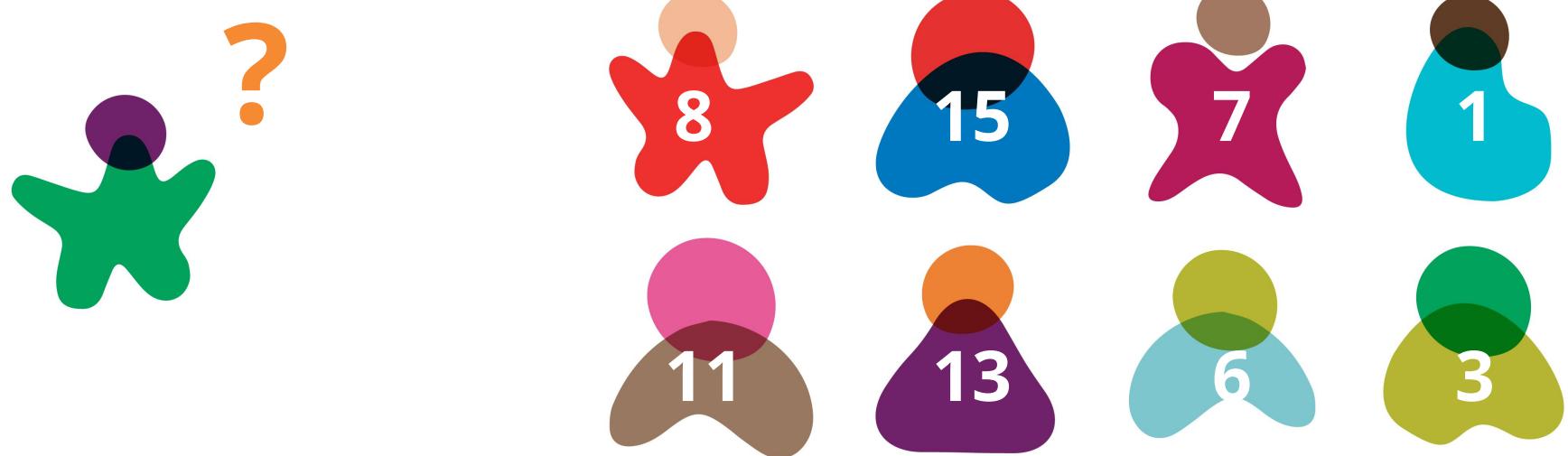
B



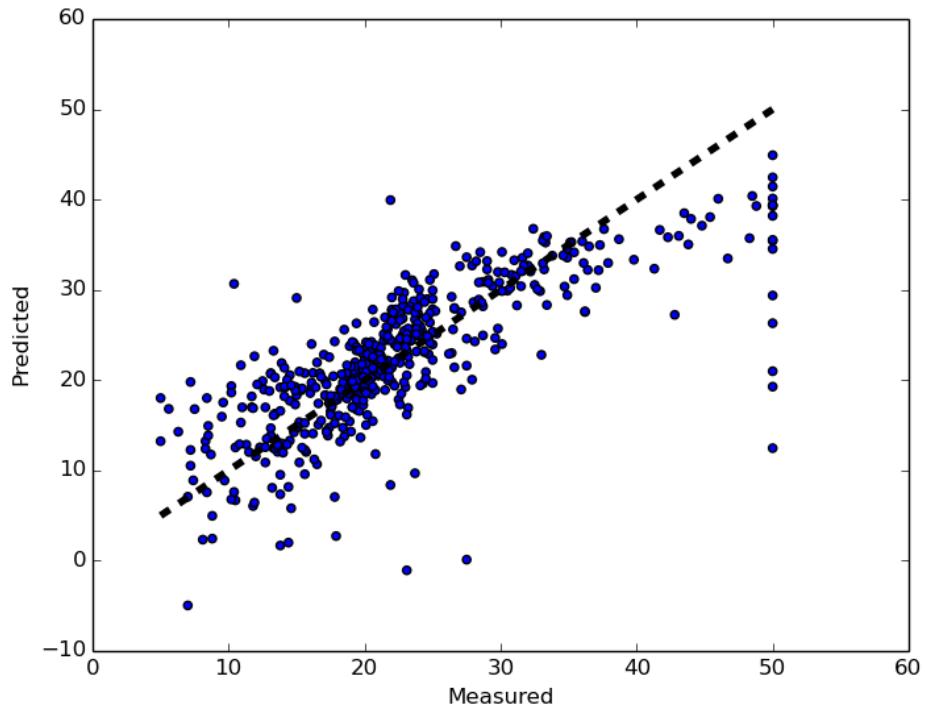
CLASSIFICATION



REGRESSION

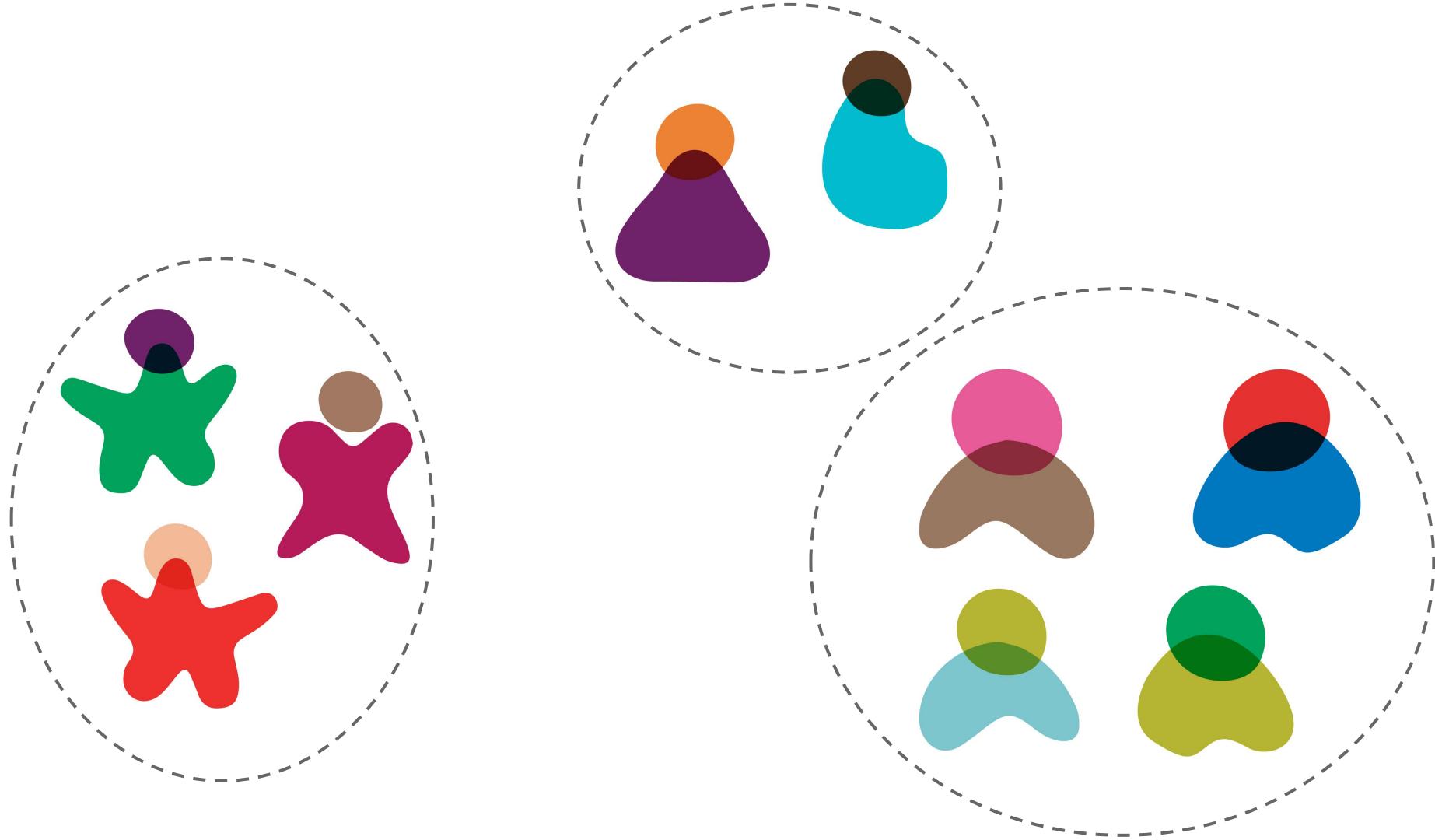


REGRESSION

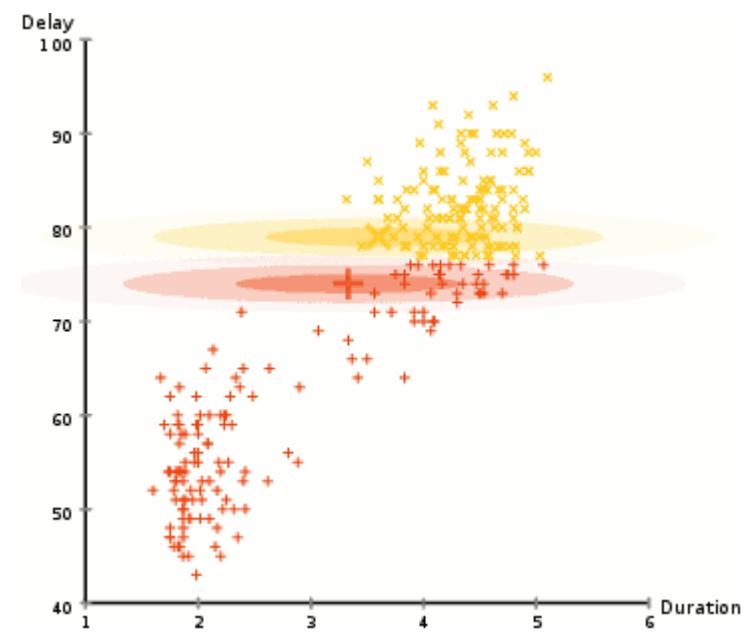
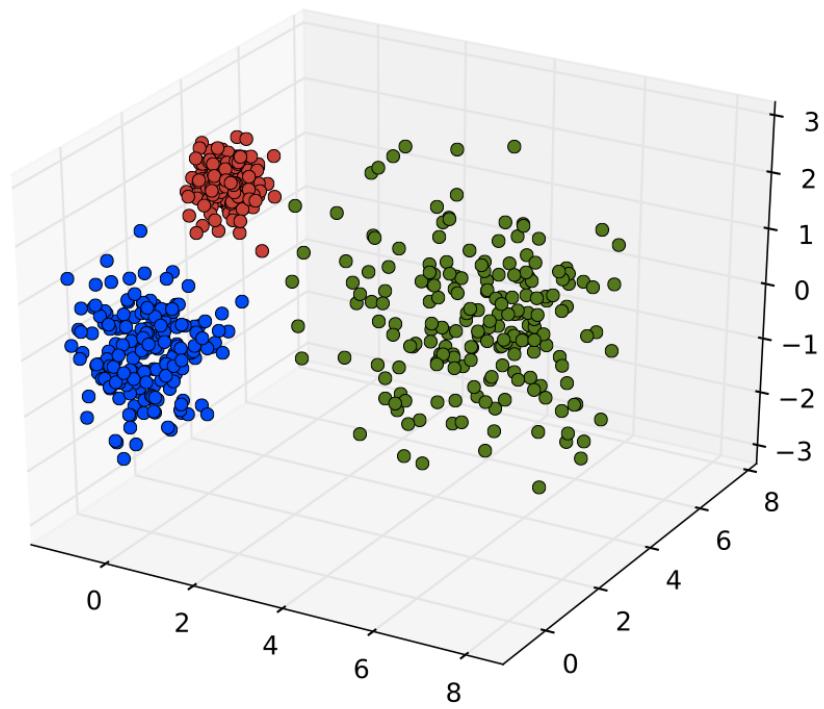


CRIM	Per capita crime rate by town
ZN	Proportion of residential land zoned for lots over 25,000 sq.ft.
INDUS	Proportion of non-retail business acres per town
CHAS	Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
NOX	Nitric oxides concentration (parts per 10 million)
RM	Average number of rooms per dwelling
AGE	Proportion of owner-occupied units built prior to 1940
DIS	Weighted distances to five Boston employment centers
RAD	Index of accessibility to radial highways
TAX	Full-value property-tax rate per \$10,000
PTRATIO	Pupil-teacher ratio by town
B	$1000(Bk - 0.63)^2$ where Bk is the proportion of African-Americans by town
LSTAT	% Lower status of the population
MEDV	Median value of owner-occupied homes in \$1000's
CAT.MEDV	Binary variable that indicates based on the MEDV variable. If MEDV > 30, CAT.MEDV = 1

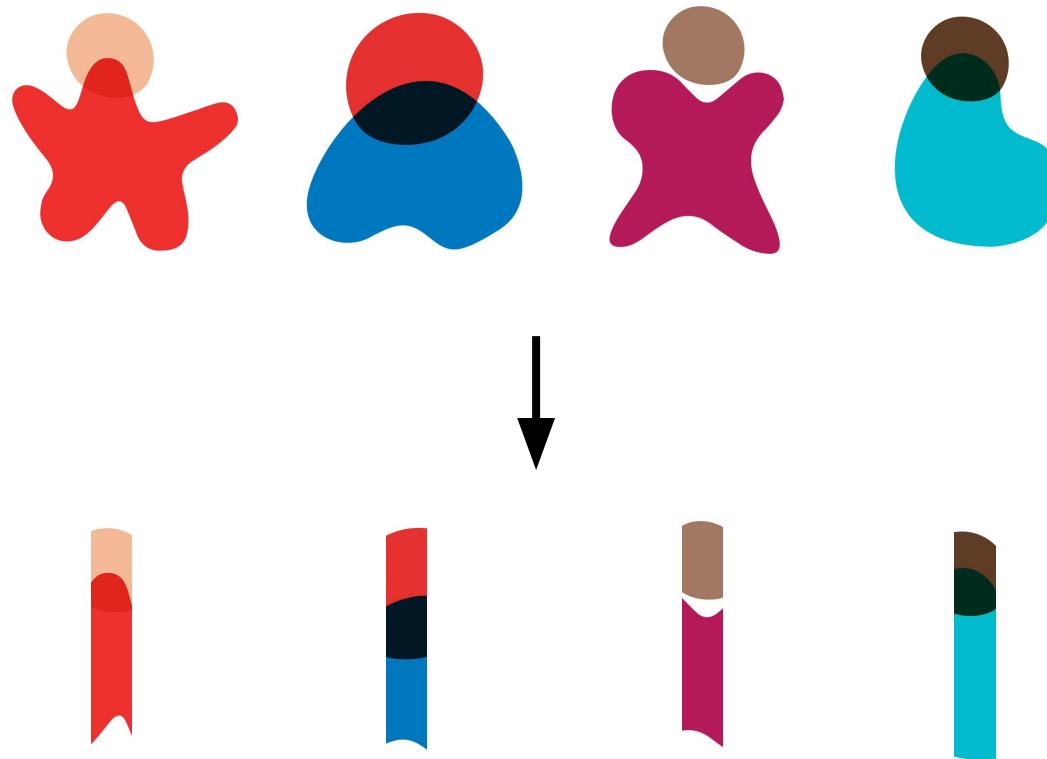
CLUSTERING



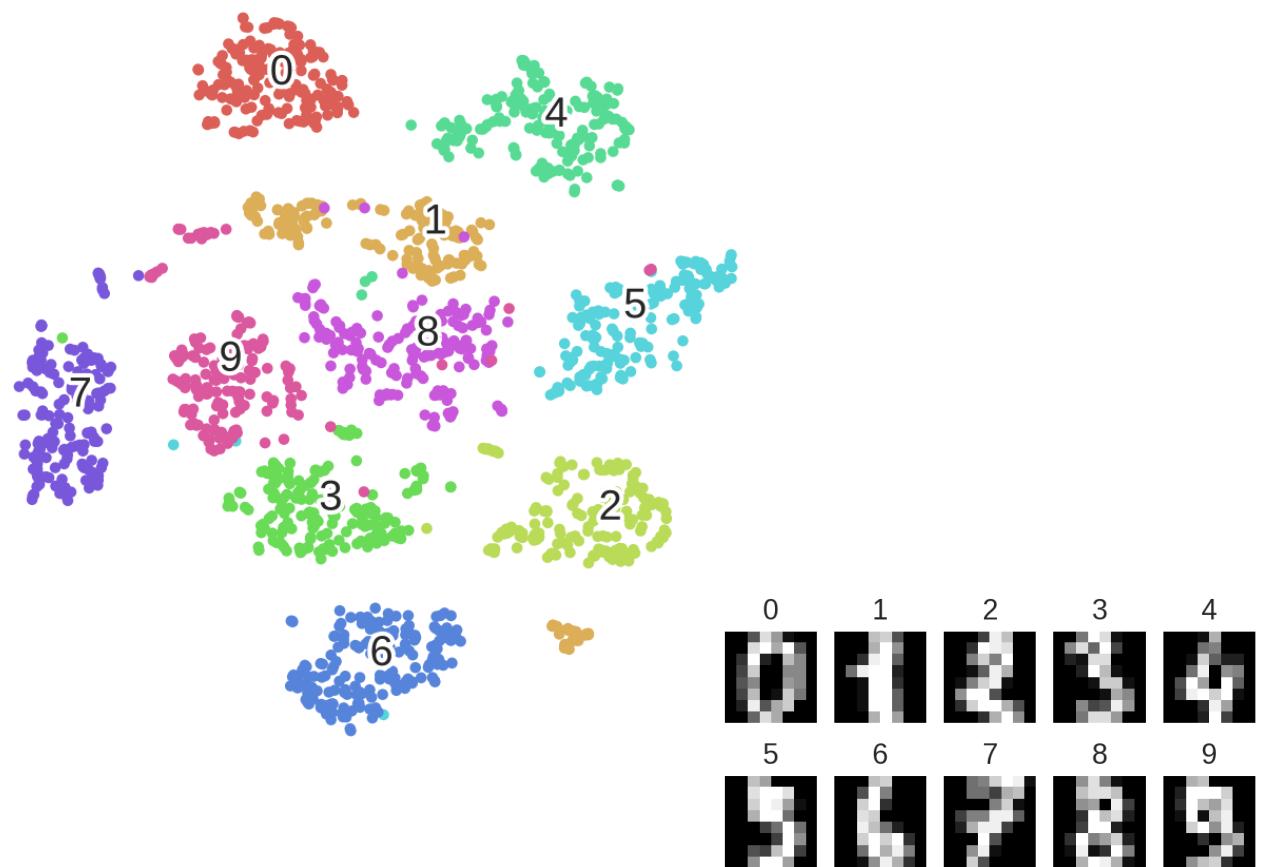
CLUSTERING



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION

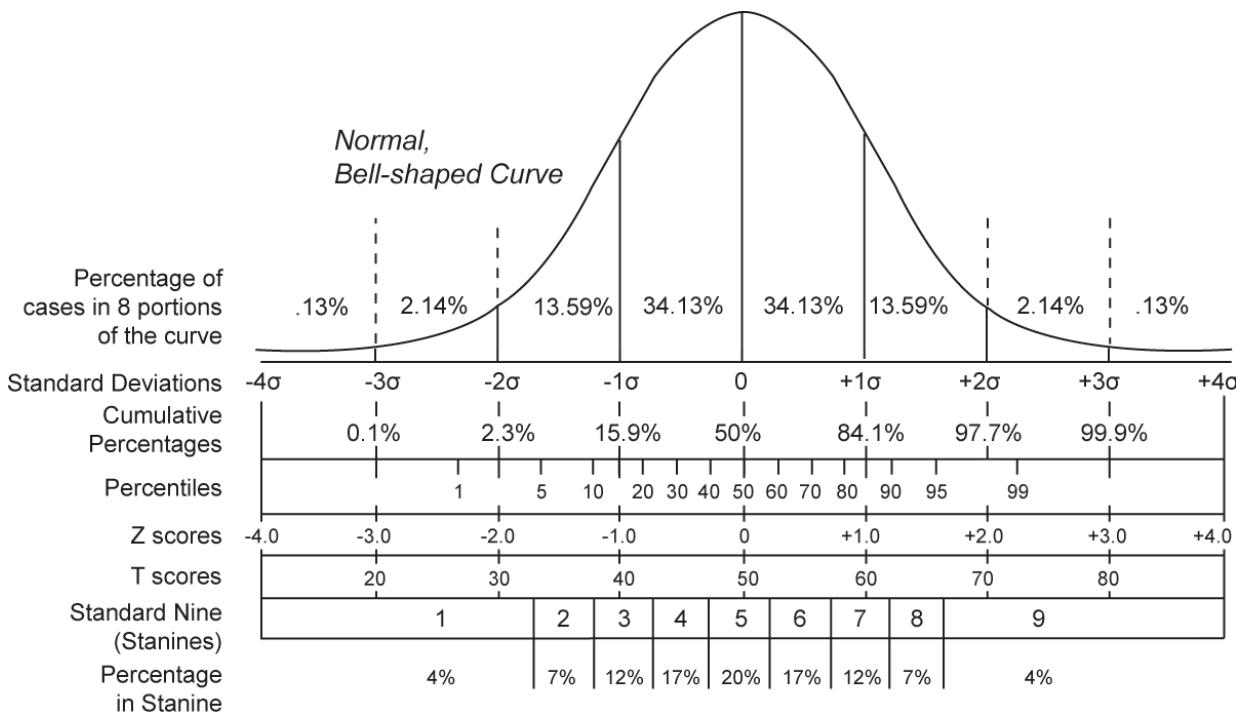


ThoughtWorks®

DESIGN DECISIONS

1, 2 steps!

NORMALIZATION



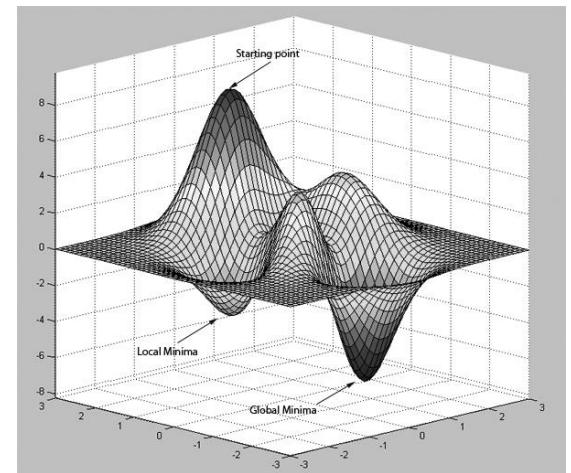
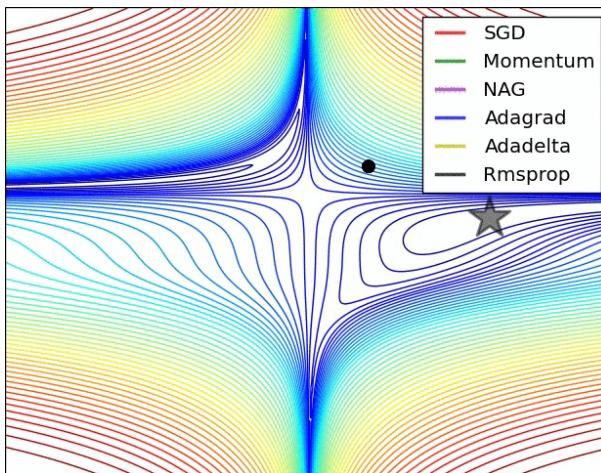
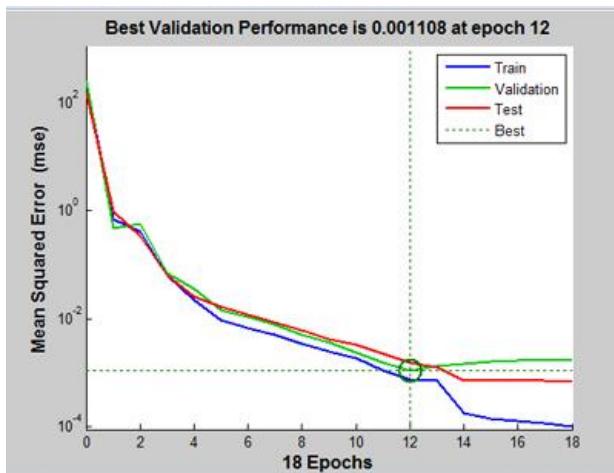
- z-score

$$X' = \frac{X - \mu}{\sigma}$$

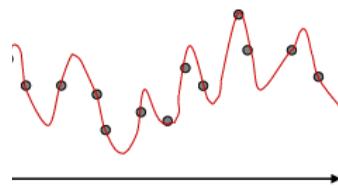
- min-max

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

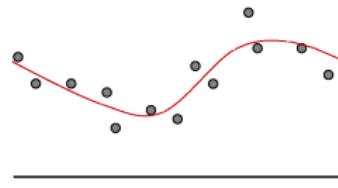
TRAINING



REGULARIZATION

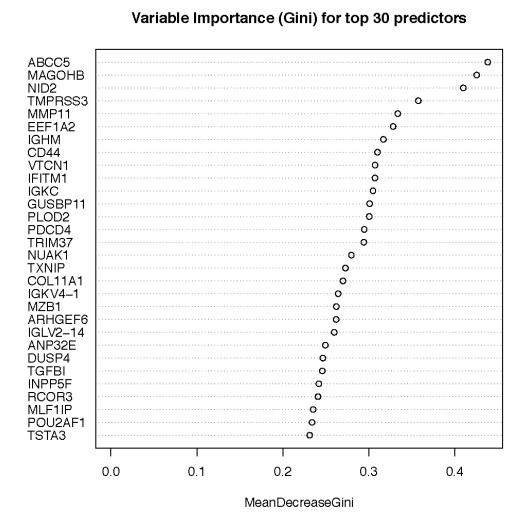
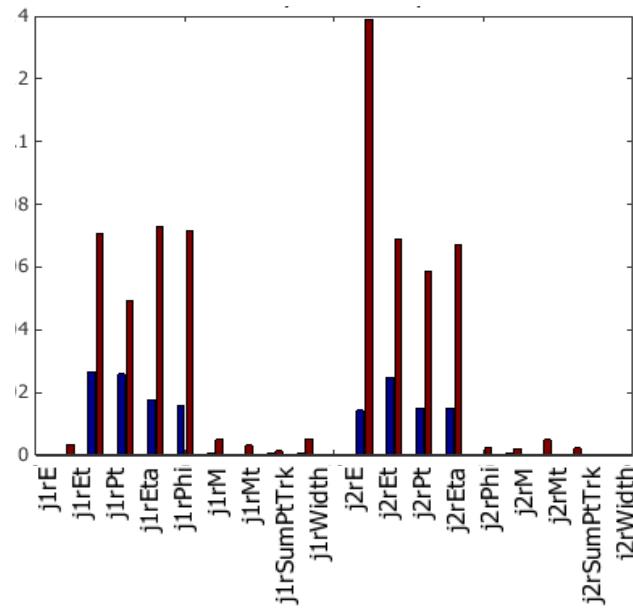
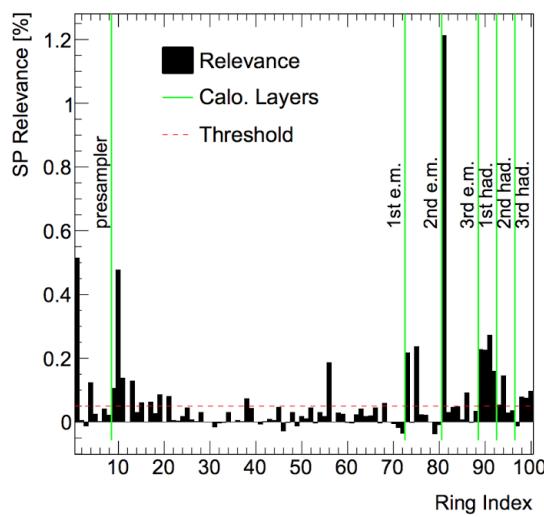


$$\min_f |Y_i - f(X_i)|^2$$

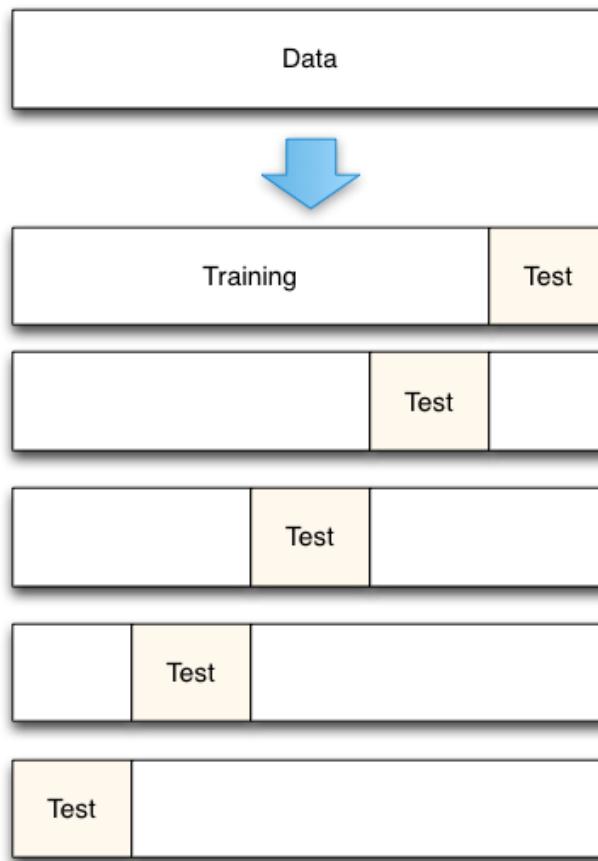


$$\min_{f \in H} \sum_{i=1}^n |Y_i - f(X_i)|^2 + \lambda \|f\|_H^2$$

RELEVANCE ANALYSIS



CROSS VALIDATION

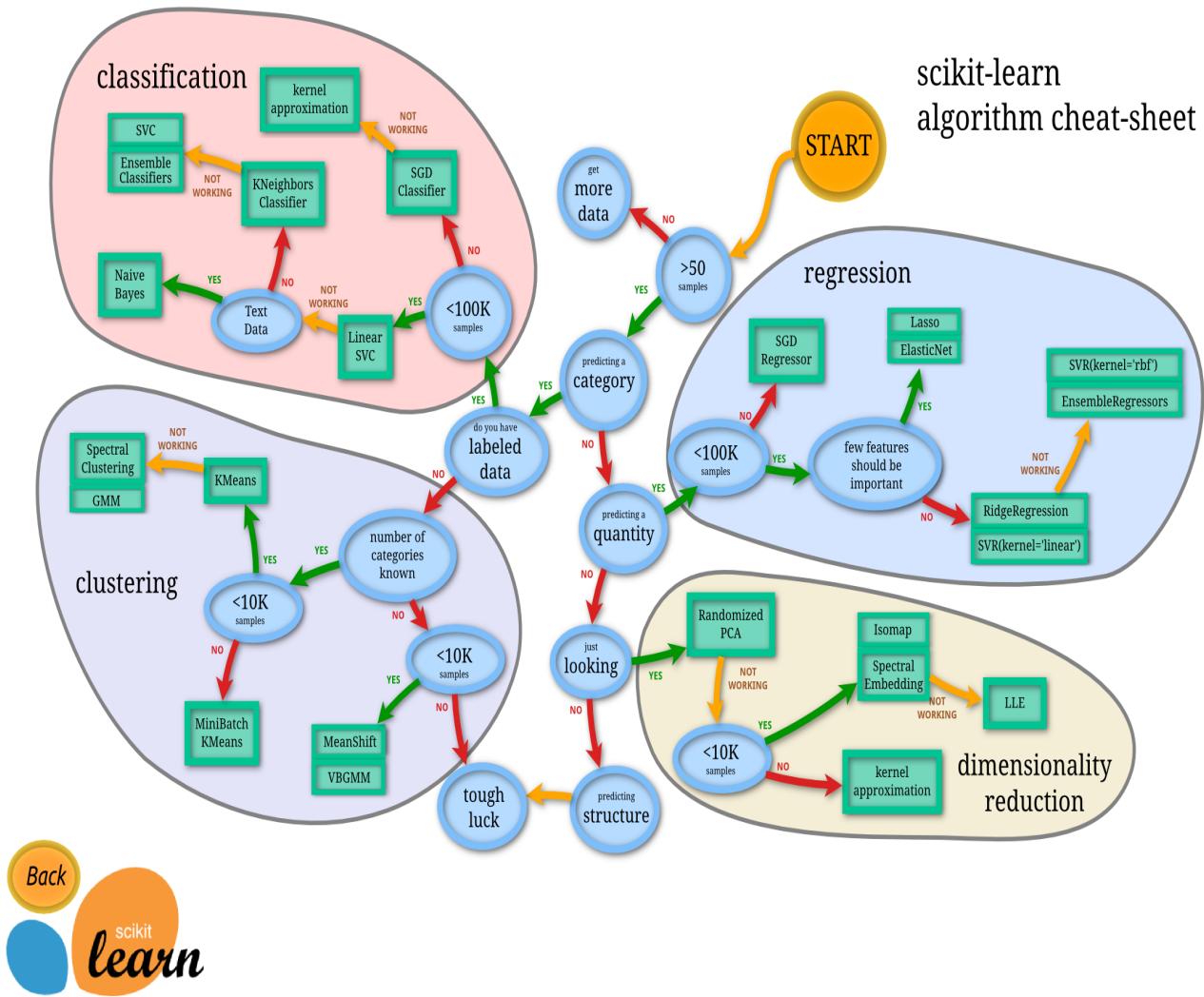


ThoughtWorks®

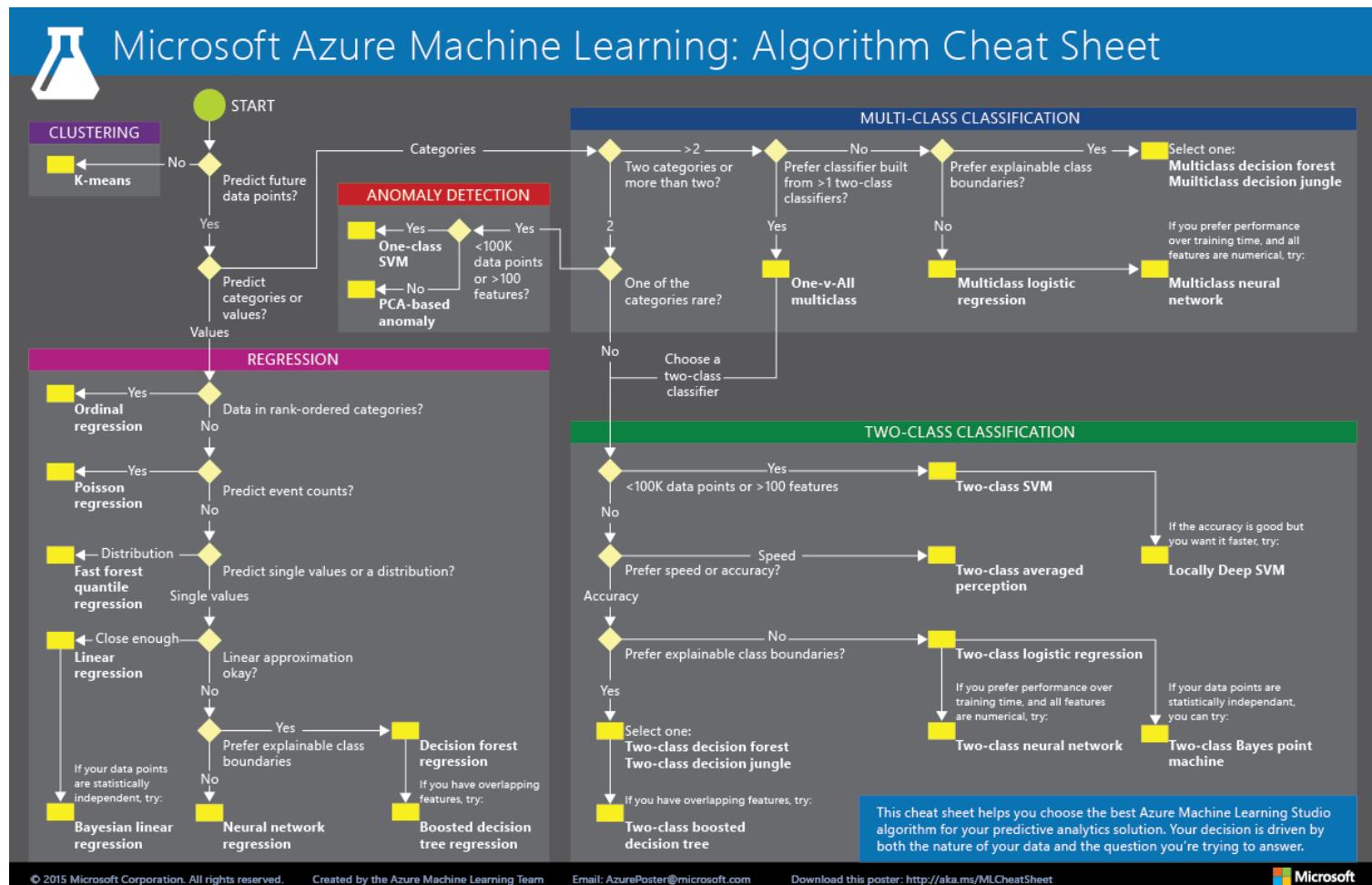
ALGORITHMS

Cheat sheet included!

CHEAT SHEET

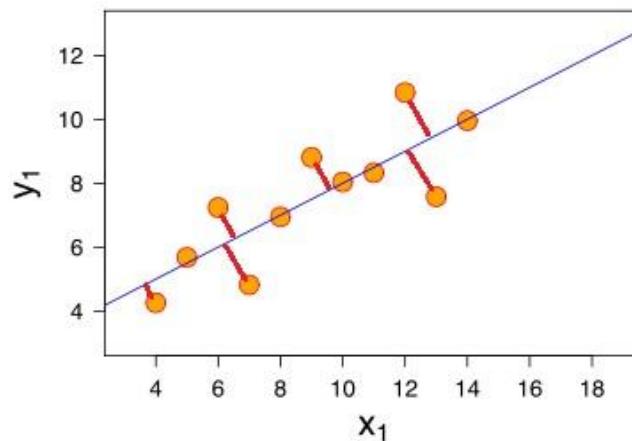


CHEAT SHEET

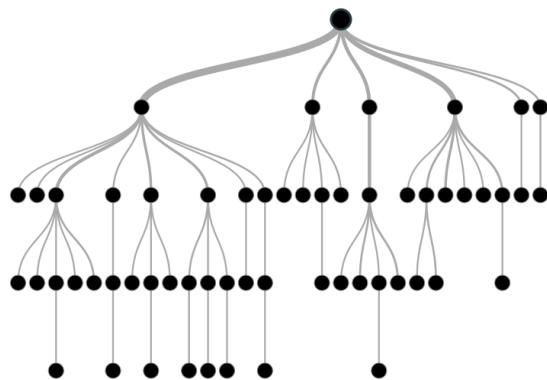


ALGORITHMS PT. I

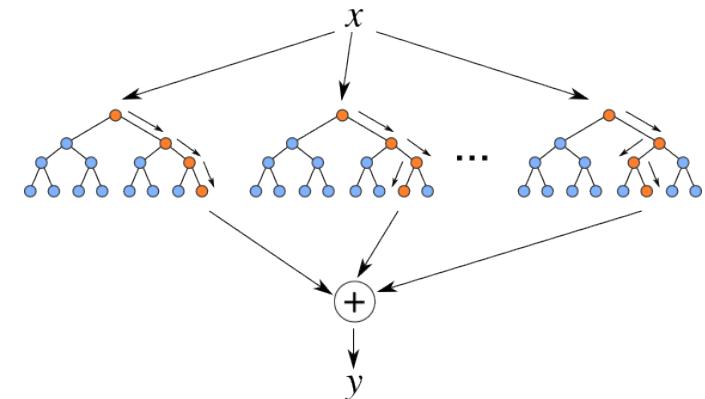
Linear Regression



Decision Trees

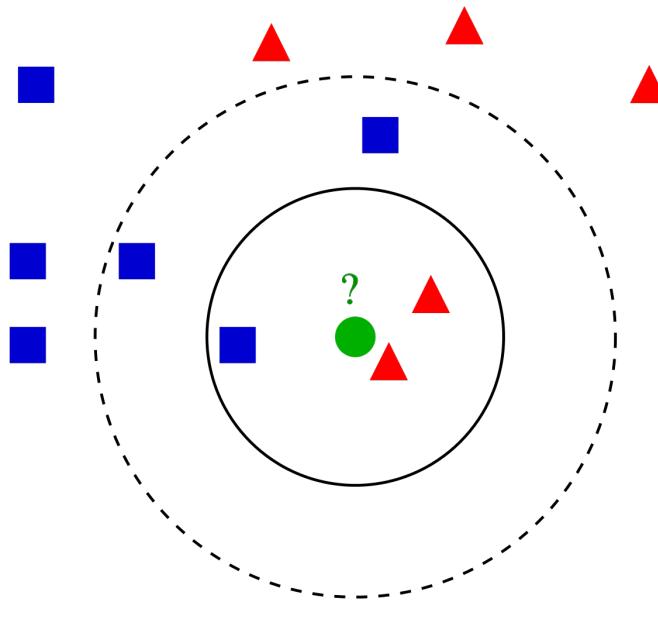


Random Forest

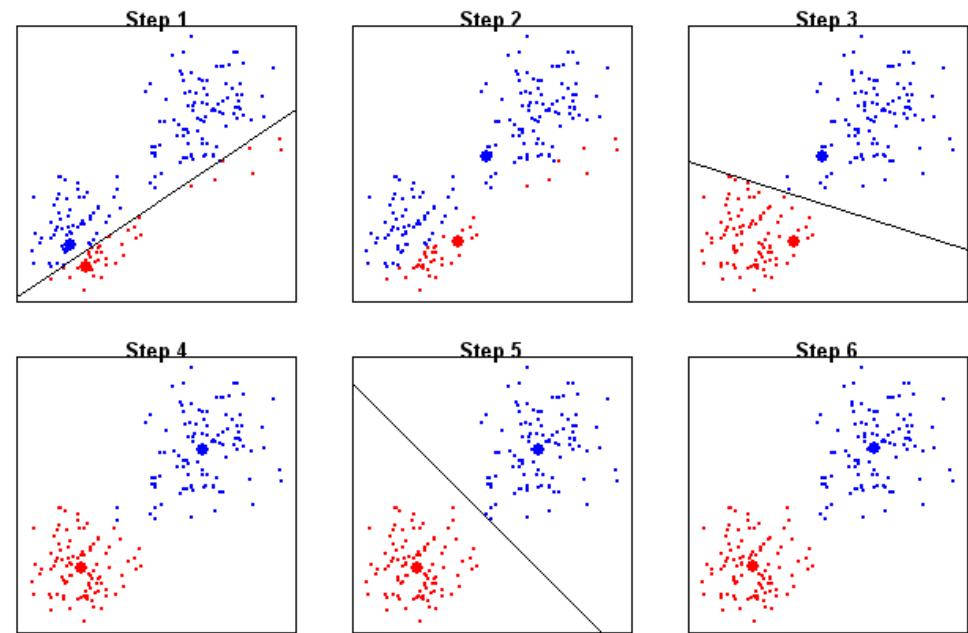


ALGORITHMS PT. II

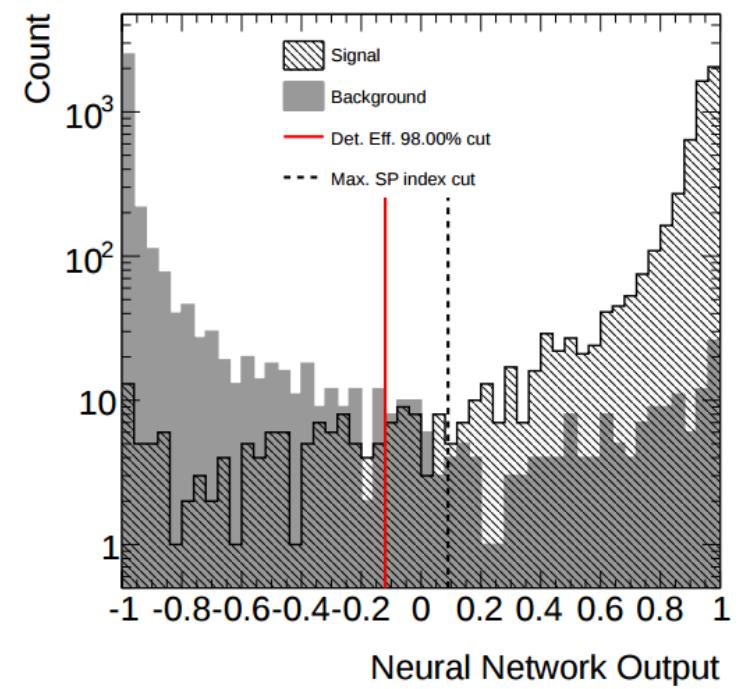
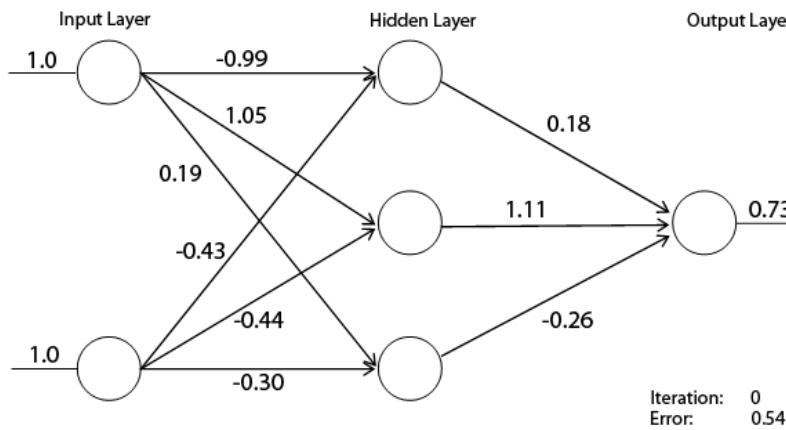
K-Nearest Neighbors



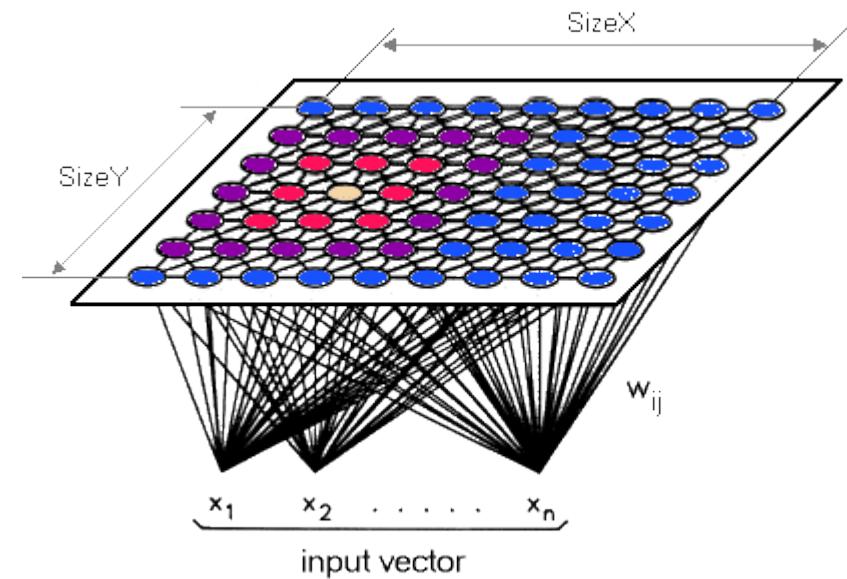
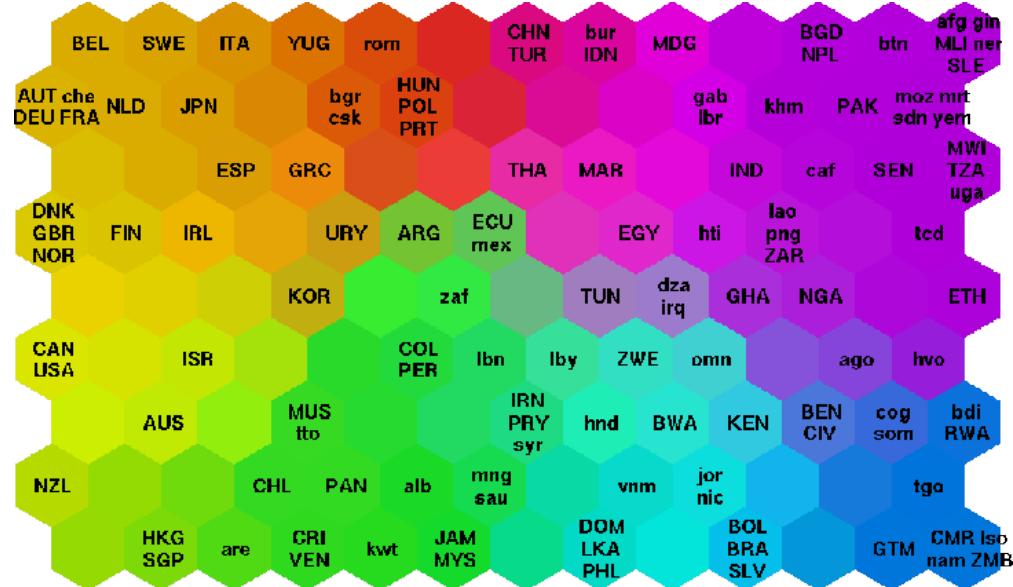
K-Means



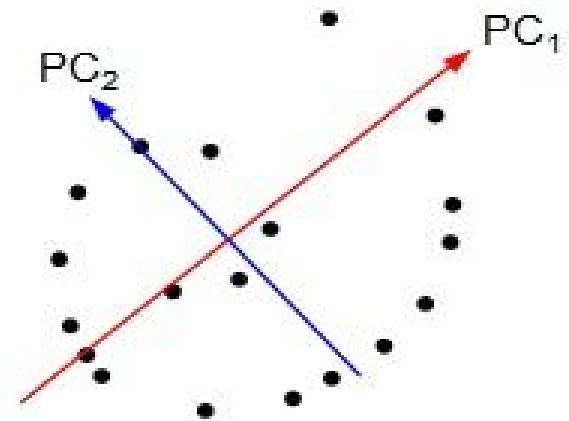
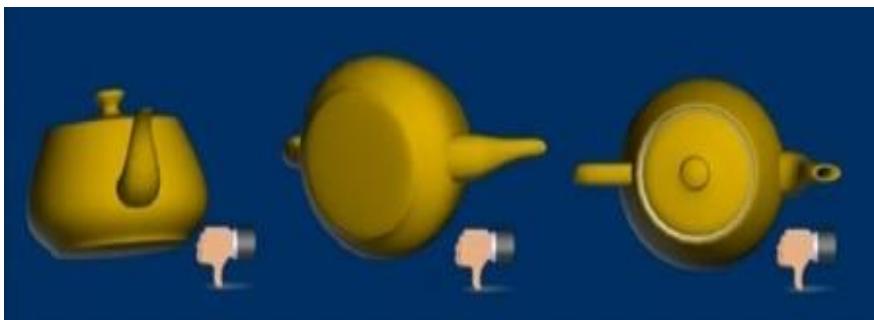
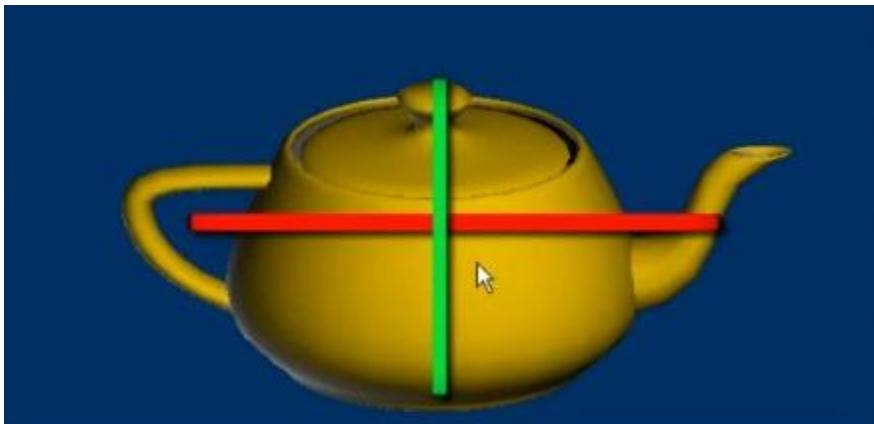
NEURAL NETWORKS



SELF-ORGANIZING MAPS



PRINCIPAL COMPONENTS ANALYSIS



T-SNE

9

ThoughtWorks®

EVALUATION

How you doin'?

PRECISION AND ACCURACY



CONFUSION MATRIX

TP = True Positives

TN = True Negatives

FP = False Positives

FN = False Negatives

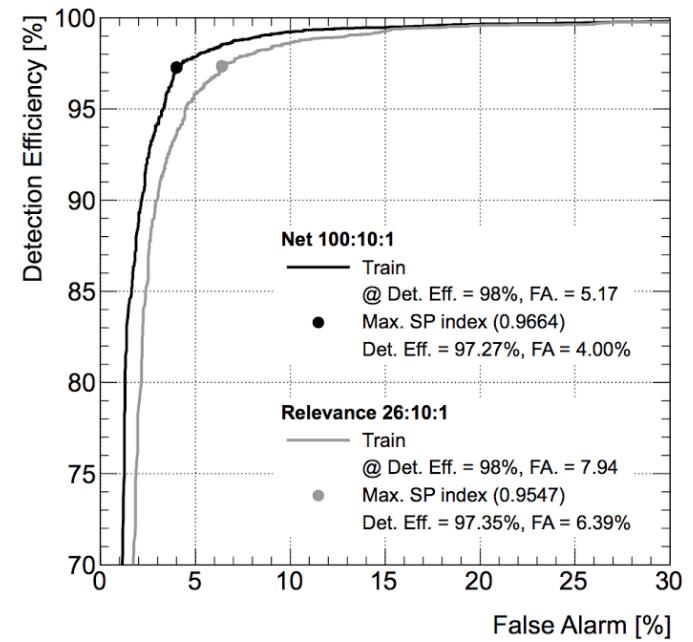
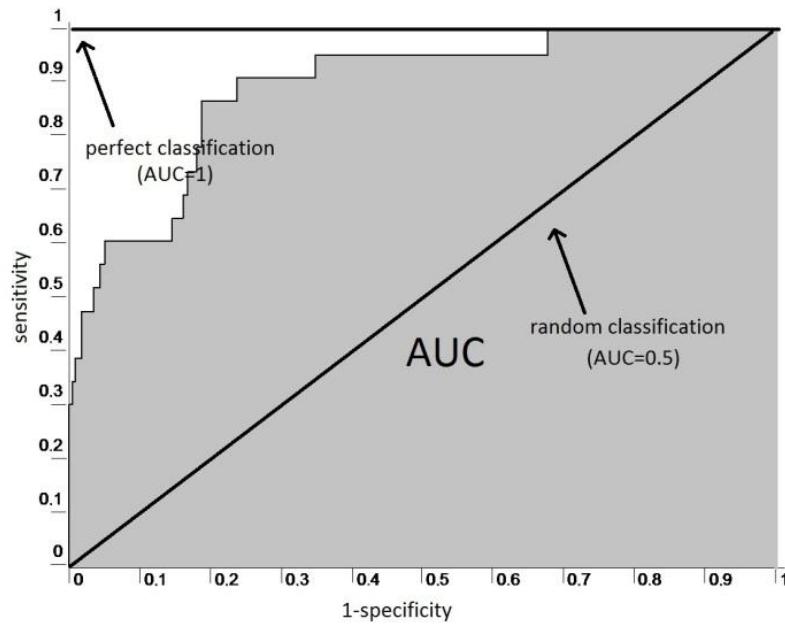
	p' (Predicted)	n' (Predicted)
P (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

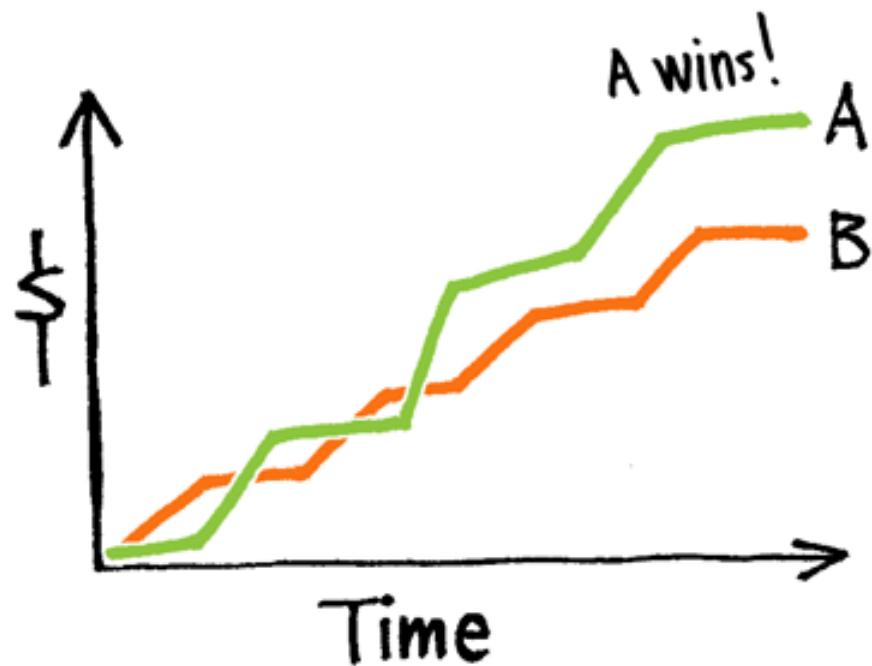
$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F1-score} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

ROC CURVE



A/B TESTS



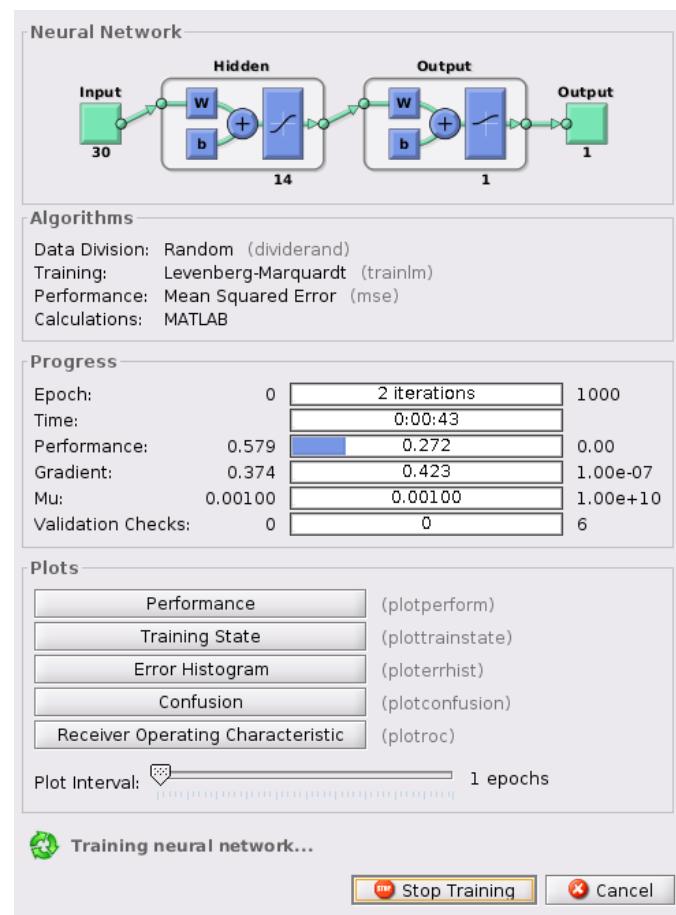
ThoughtWorks®

Code Snippets

"Hello, Machine Learning"

MATLAB 101

```
[x,y] = ovarian_dataset;  
net = patternnet(5);  
[net,tr] = train(net,x,y);  
testX = x(:,tr.testInd);  
testY = net(testX);
```



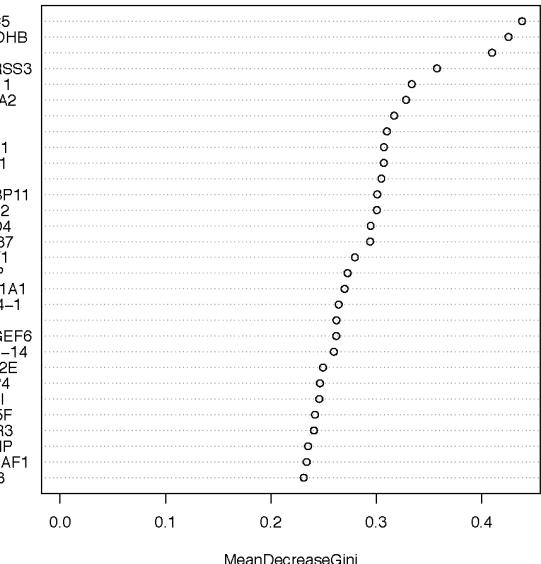
MATLAB 201

```
net = patternnet(14);
net.input.processFcns = {'mapminmax', 'fixunknowns', 'processpca'};
net.inputs{1}.processParams{3}.maxfrac = 0.02;
net.trainFcn = 'trainlm';
net.performFcn = 'mse';
net.divideParam.trainRatio = 70/100;
net.divideParam.valRatio = 15/100;
net.divideParam.testRatio = 15/100;
[net, tr] = train(net_config, test_inputs, train_targets);
outputs = net(test_inputs);
```

R

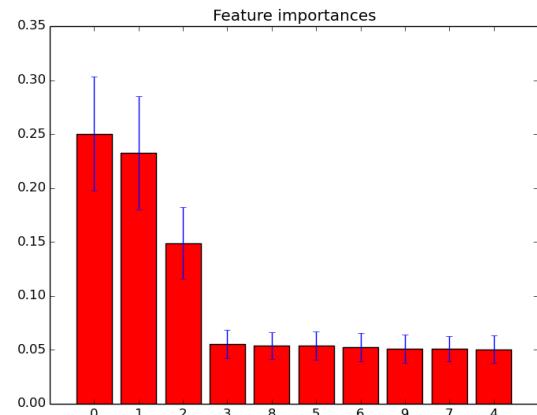
```
library(randomForest)
raw.orig <- read.csv(file="train.txt", header=T, sep="\t")
frmla = Metal ~ OTW + AirDecay + Koc
fit.rf = randomForest(frmla, data=raw)
print(fit.rf)
importance(fit.rf)
```

Variable Importance (Gini) for top 30 predictors



SCIKIT LEARN

```
dataset = pd.read_csv('Data/train.csv')
target = dataset.Activity.values
train = dataset.drop('Activity', axis=1).values
test = pd.read_csv('Data/test.csv').values
rf = RandomForestClassifier(n_estimators=100, n_jobs=-1)
rf.fit(train, target)
predicted_probs = [x[1] for x in rf.predict_proba(test)]
importances = rf.feature_importances_
```



ThoughtWorks®

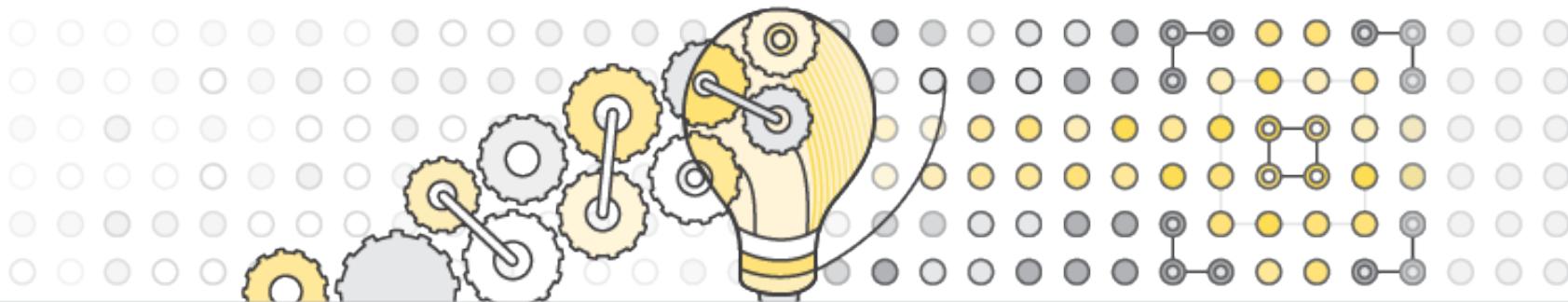
PAY-AS-YOU-GO SERVICES

Amazon Machine Learning

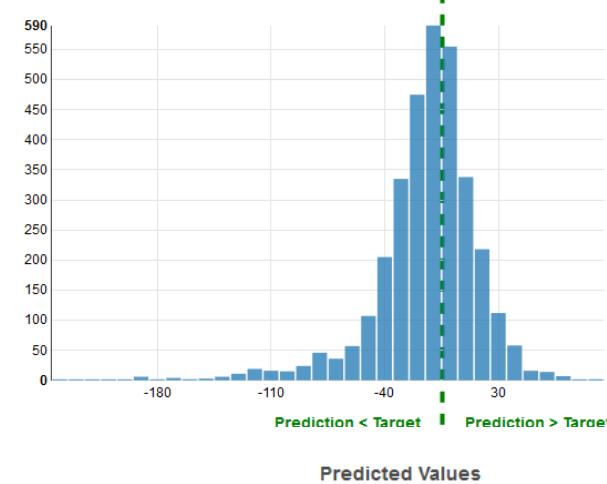
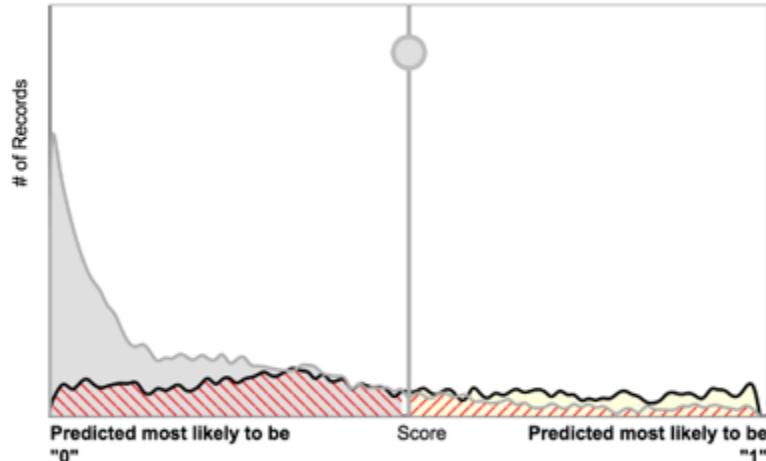
AMAZON MACHINE LEARNING

Five easy steps

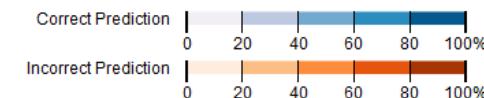
1. Upload csv dataset to Amazon S3
2. Create **Datasource** with metadata about uploaded dataset
3. Create **ML Model** with configurations for model training
4. Create **Evaluation** to analyse and tune model efficiency
5. Create **Prediction** to use trained model with new data



EVALUATION



	Romance	Thriller	Adventure	Total	F1
Romance	57.92% (49.1k)				0.78
Thriller		21.23% (18.0k)			0.33
Adventure			20.85% (17.7k)		0.32
Total	77.56% (65.8k)	9.33% (7910)	13.12% (11.1k)	100.00% (84.8k)	0.47



SDKs



AWS SDK for Java



AWS SDK for .NET



AWS SDK for Python



AWS SDK for PHP



AWS SDK for JavaScript in Node.js



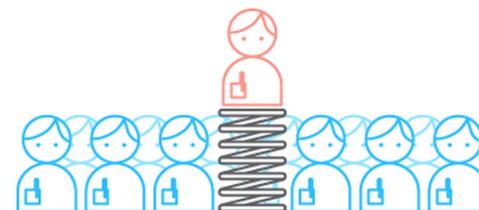
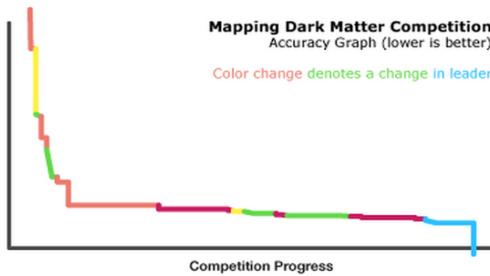
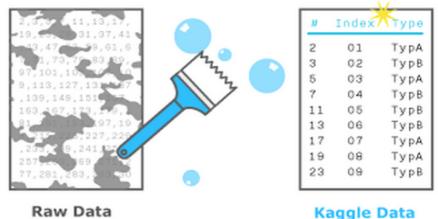
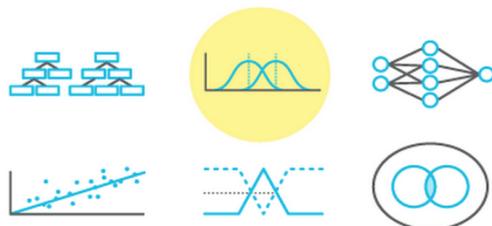
AWS SDK for Ruby

ThoughtWorks®

DATA SCIENCE COMPETITIONS

Challenge accepted!

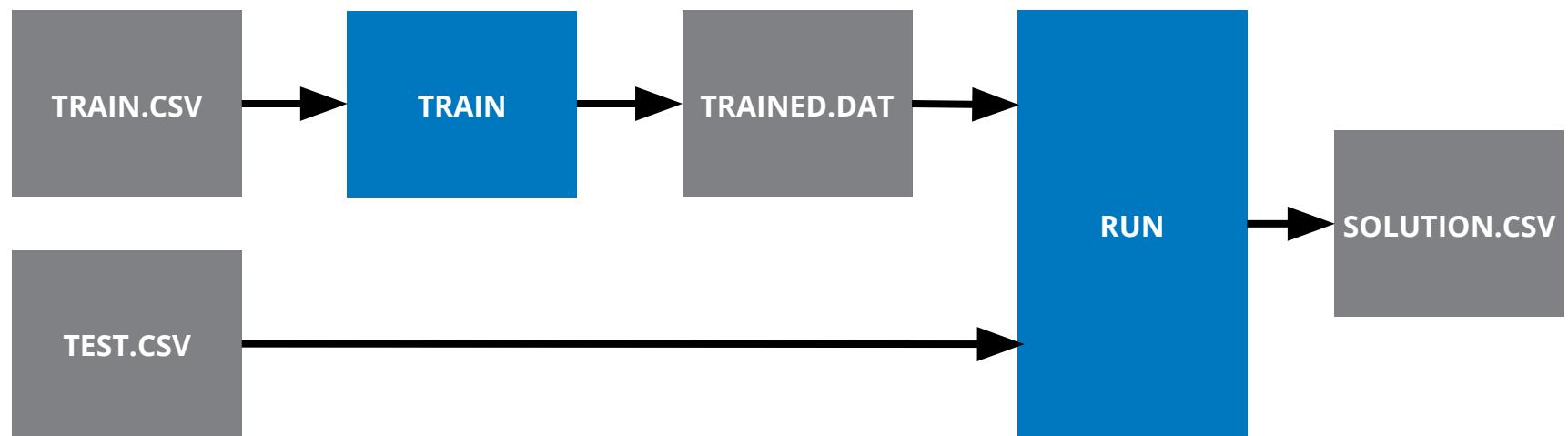
KAGGLE



SPONSORED



END TO END



THANK YOU

Questions?

Dhiana Deva

ddeva@thoughtworks.com

ThoughtWorks®