

Michael Grossberg

Intro to Data Science CS59969

Classification

Classification

Example Problem MNIST

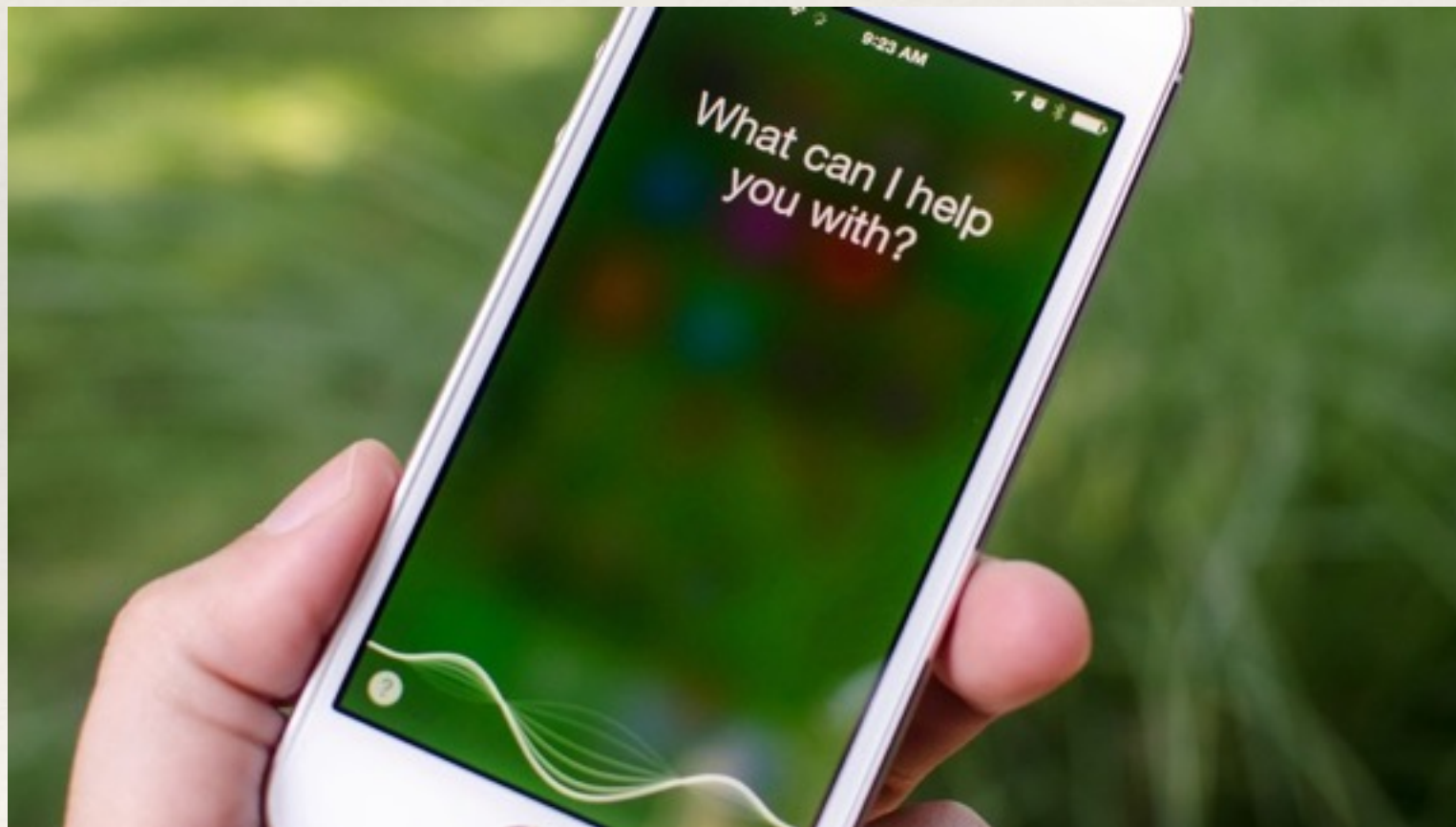


Assign Handwritten Digit
class 0,1,2, ..., 9

0

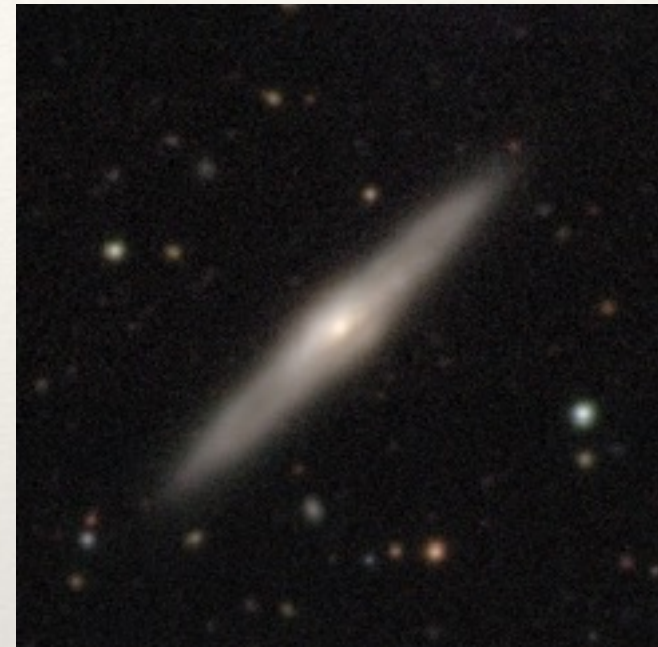
Speech Recognition

Sound \longrightarrow Word



Automatic Identification Galaxies

Smooth

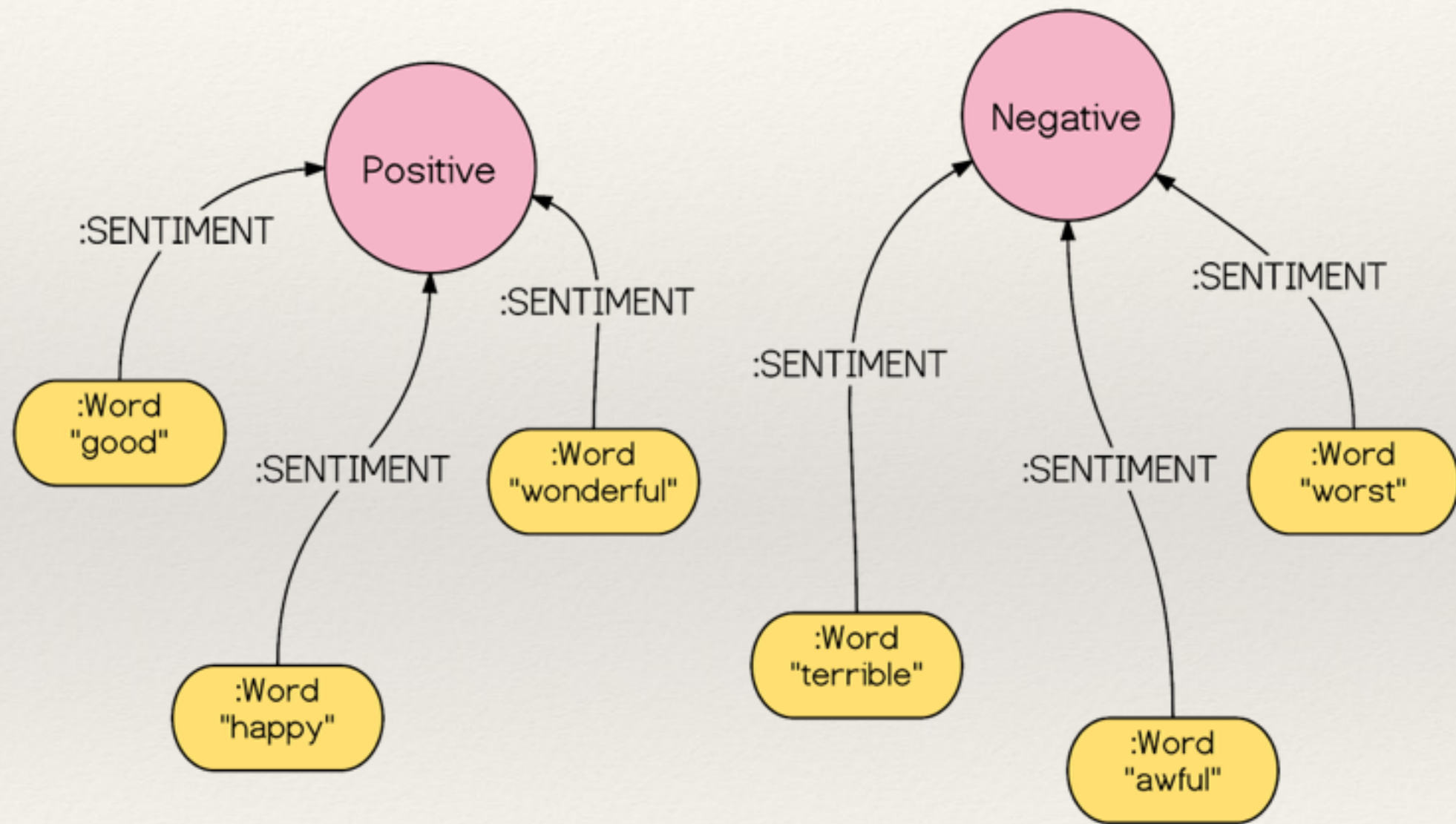


Disk

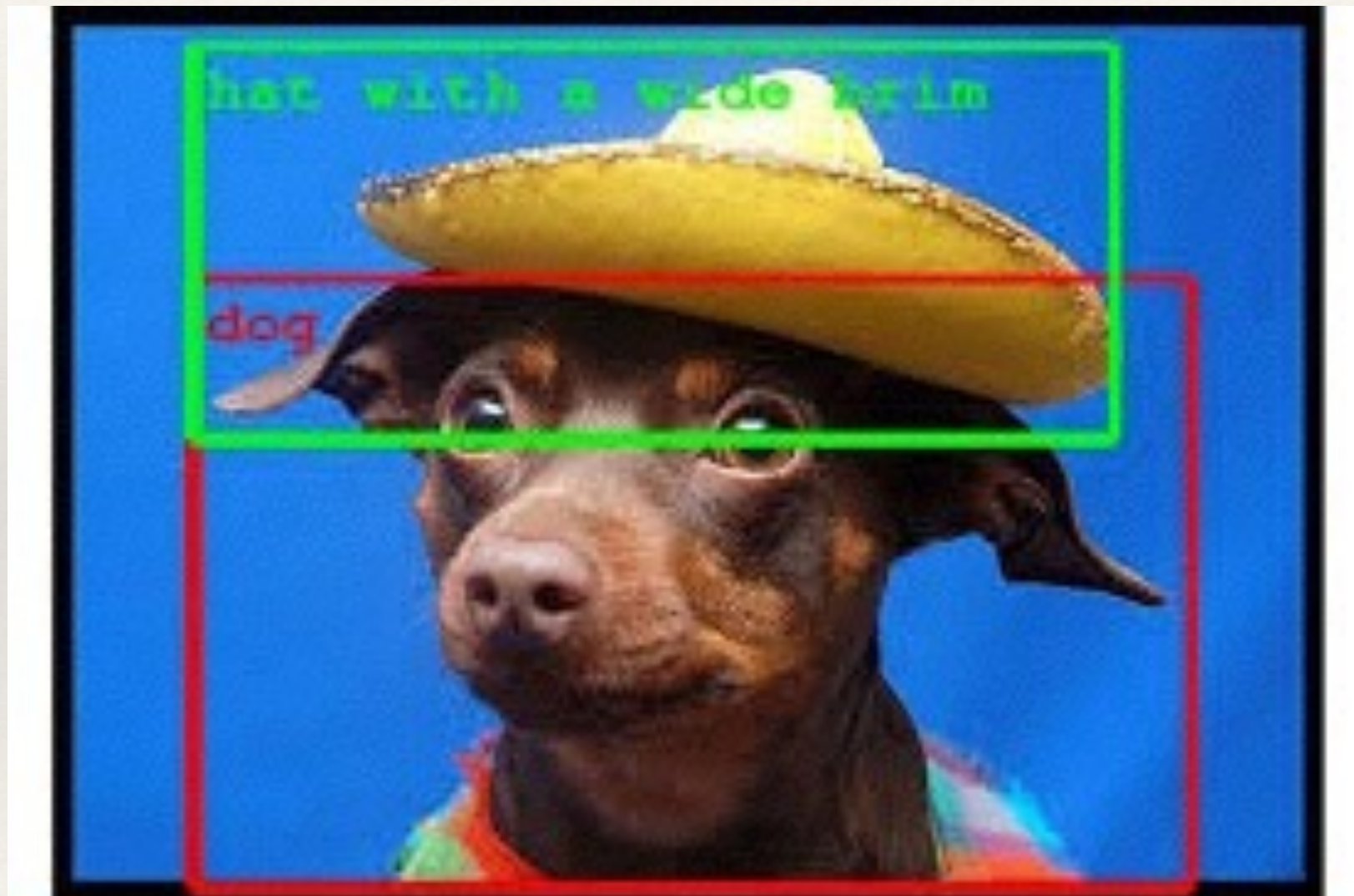
Spiral



Sentiment Analysis

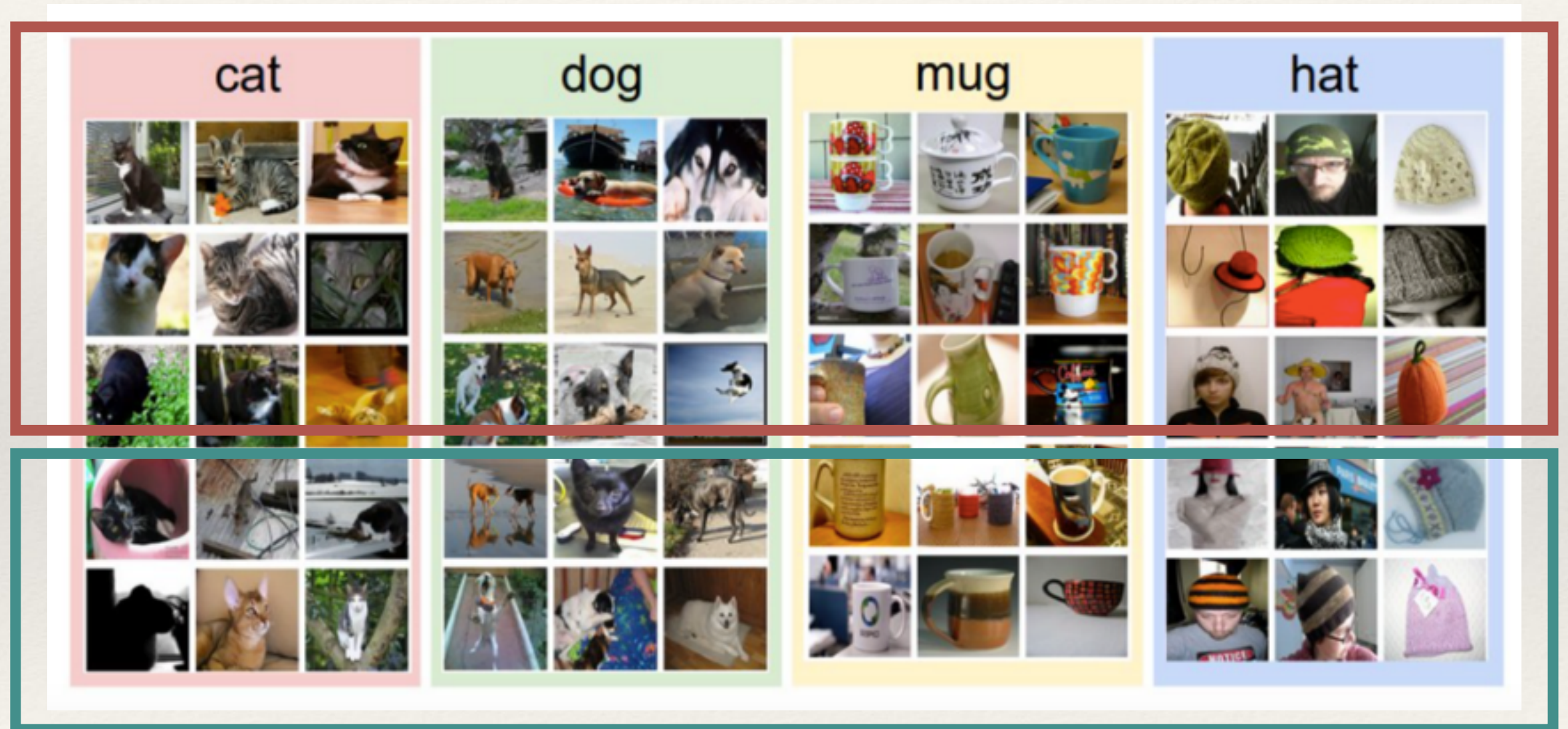


Object Recognition



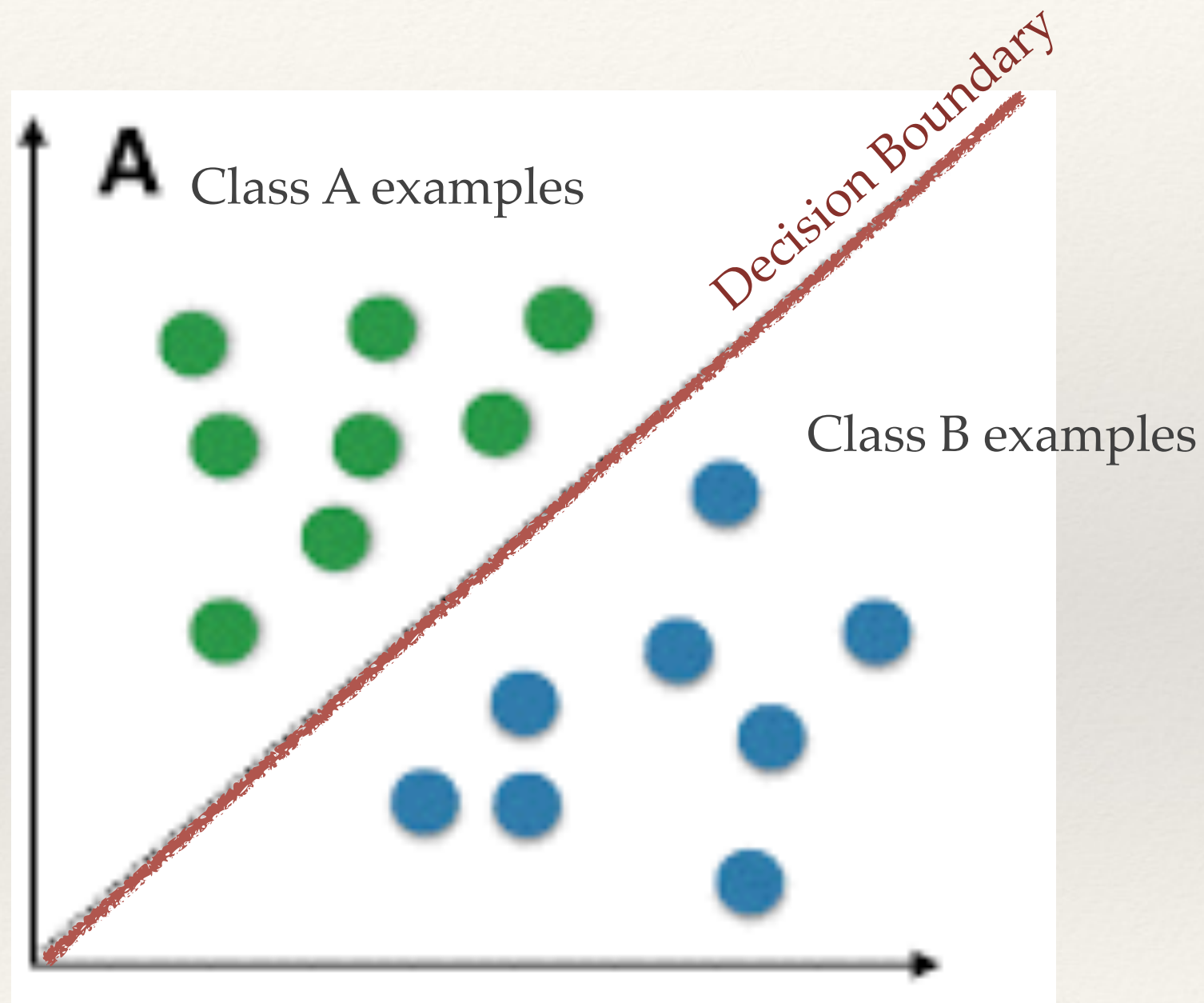
Object Recognition

Train



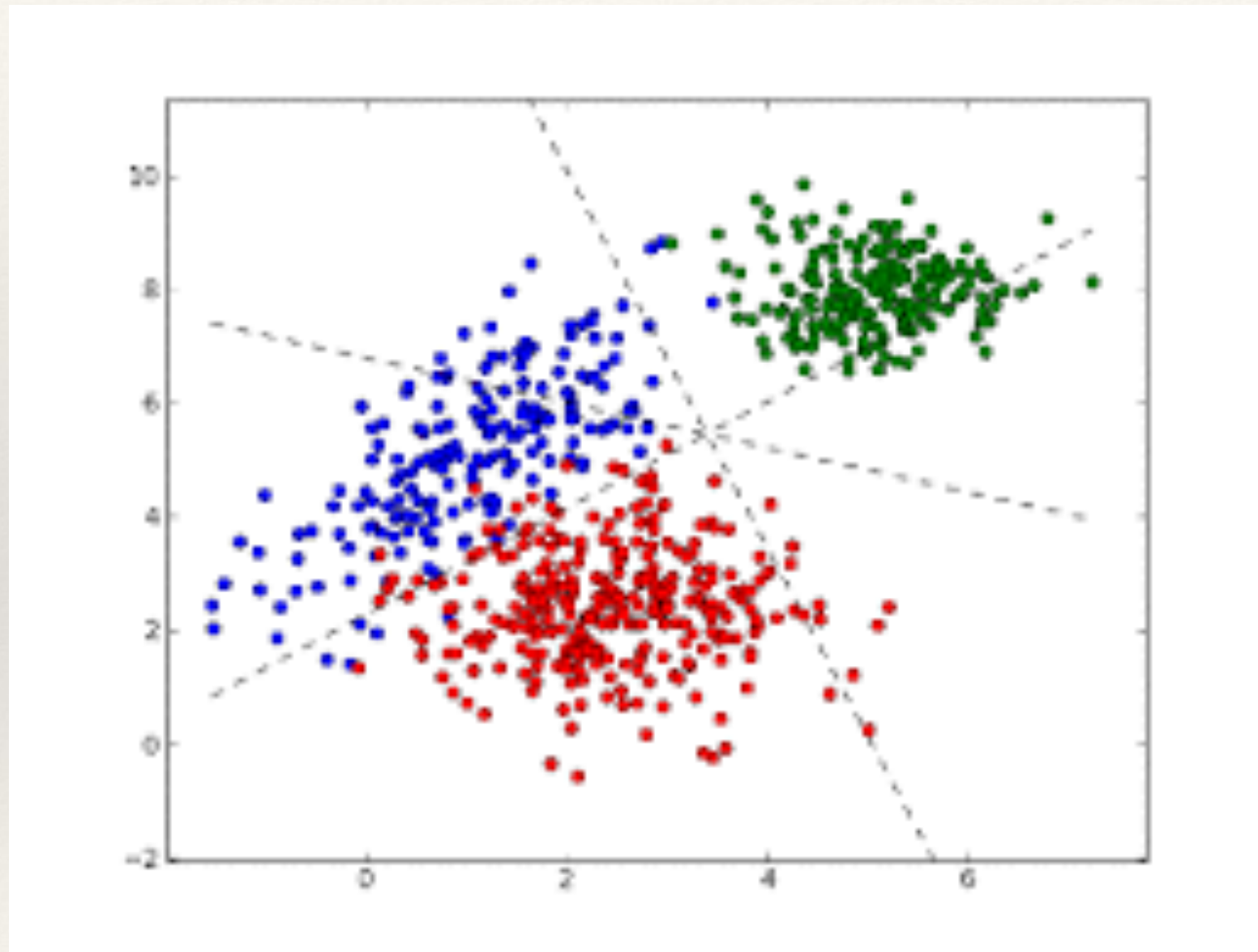
Test

Classification Part (2 class)

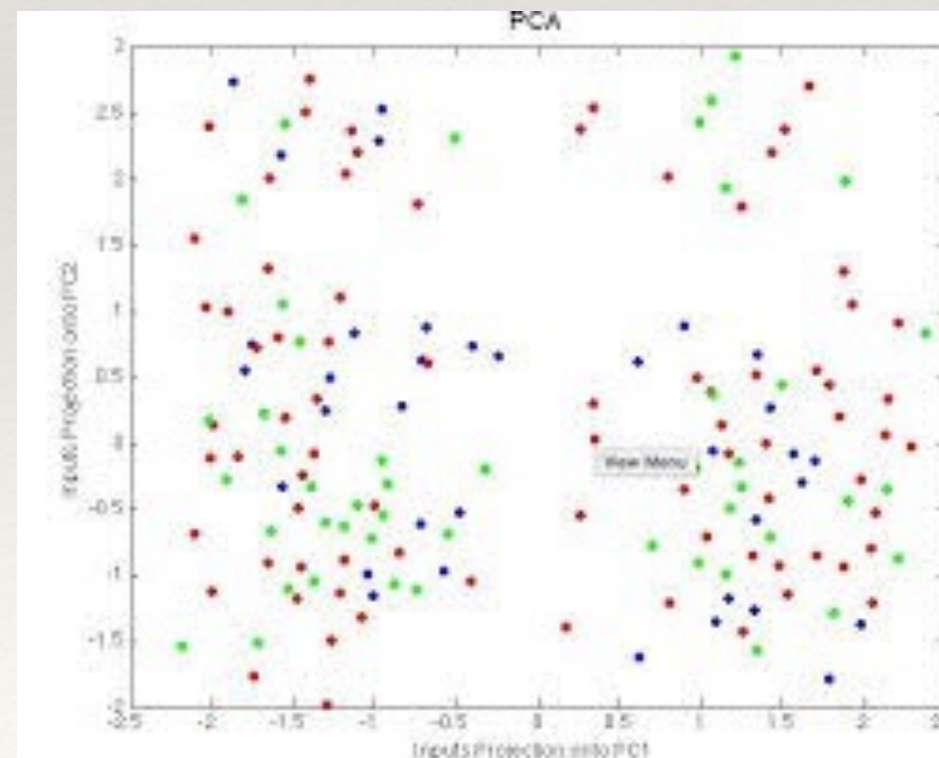
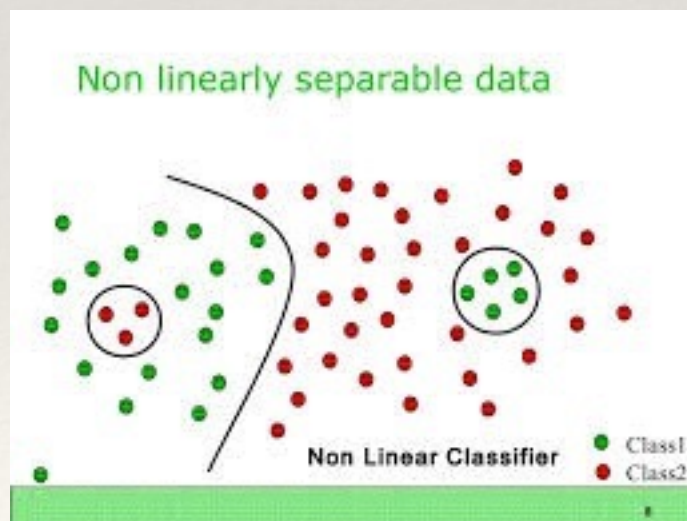
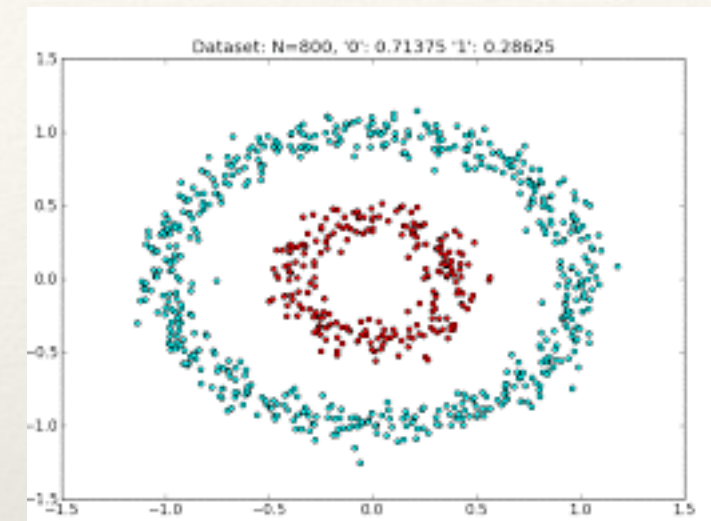
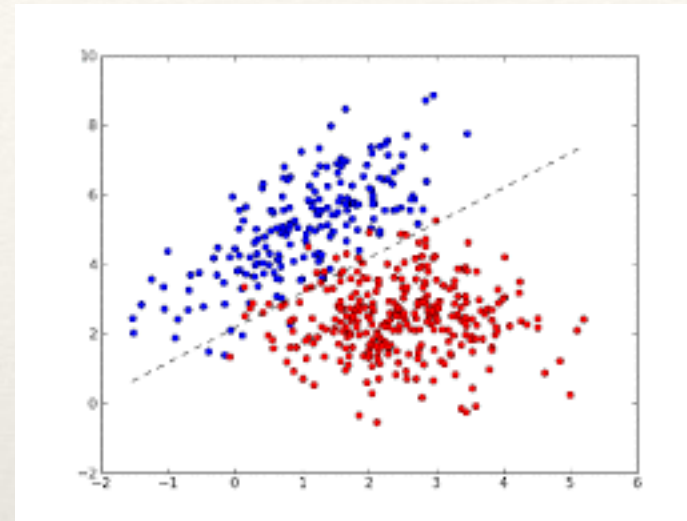
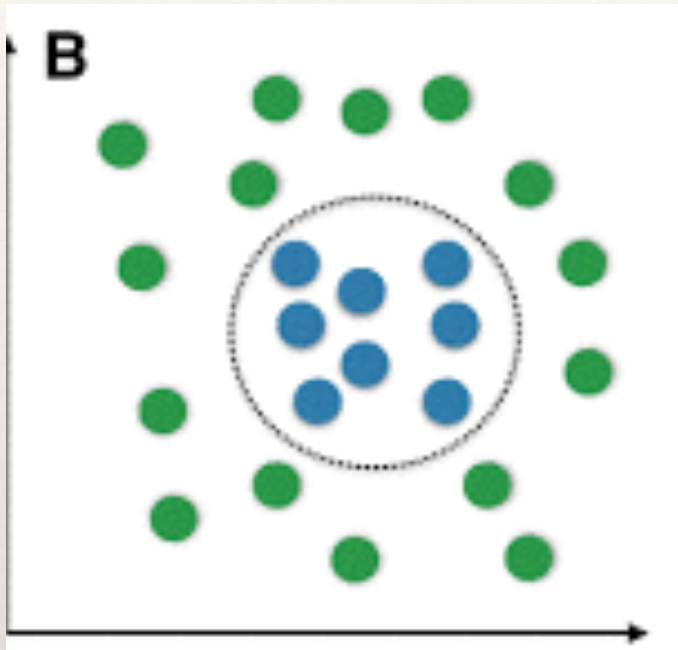


Goal Find a Decision Boundary

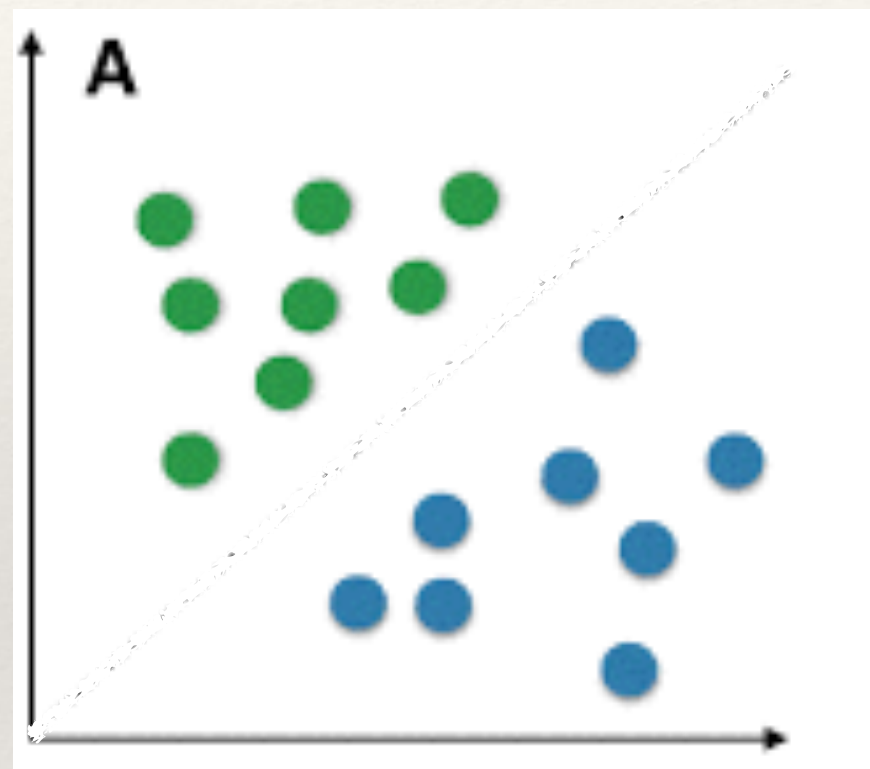
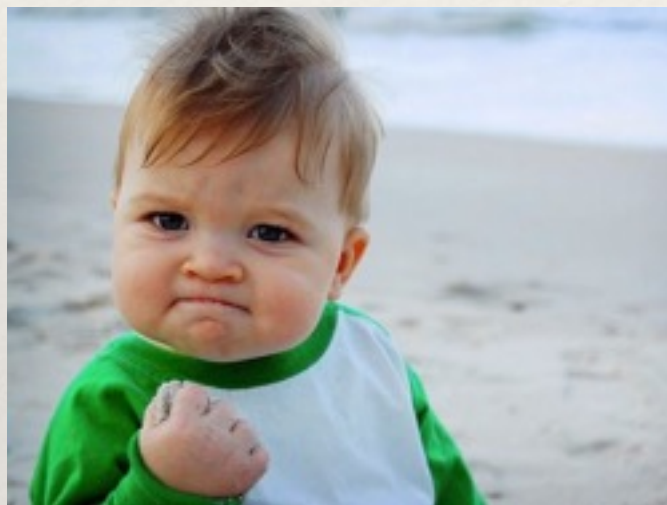
K-Class Similar



Not Always (Usually Not) So Easy



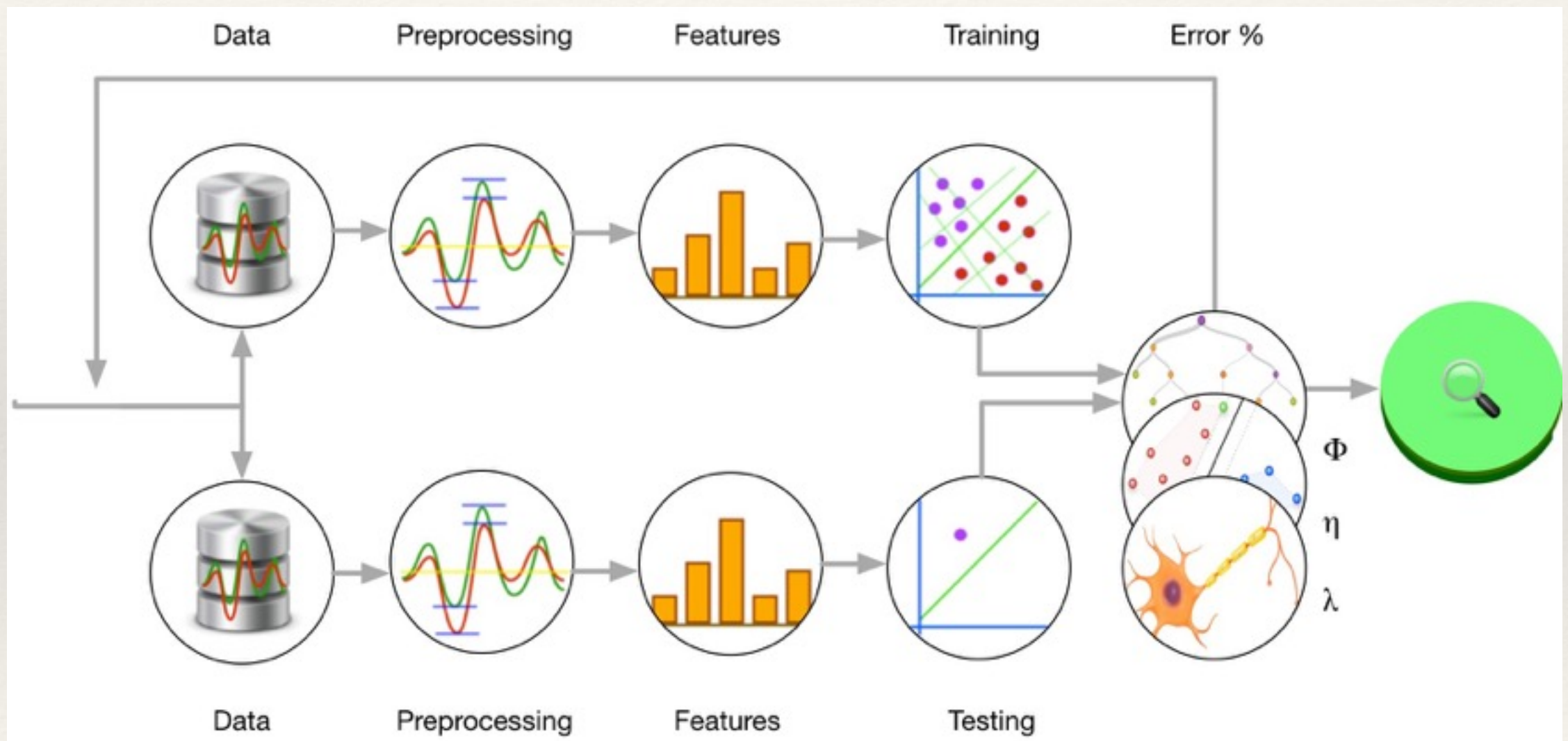
Feature Problem



How do you go from this

to this?

Usual Steps

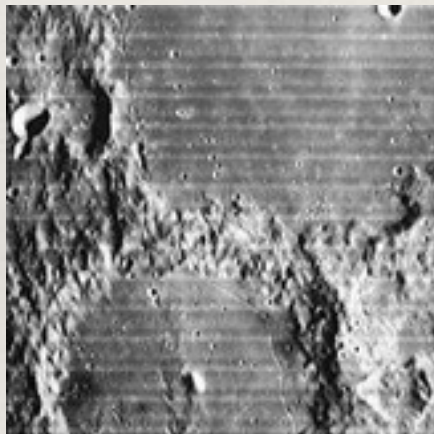


Preprocessing

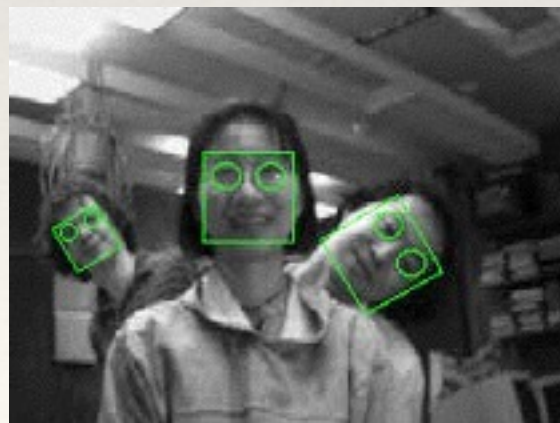


Data Cleaning

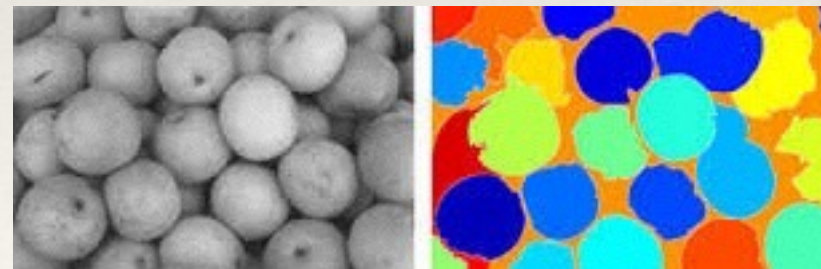
Drop bad rows?
Find bad columns?
Fill in missing data?



De-noise
De-Stripe



Localize
Rotate
Segment



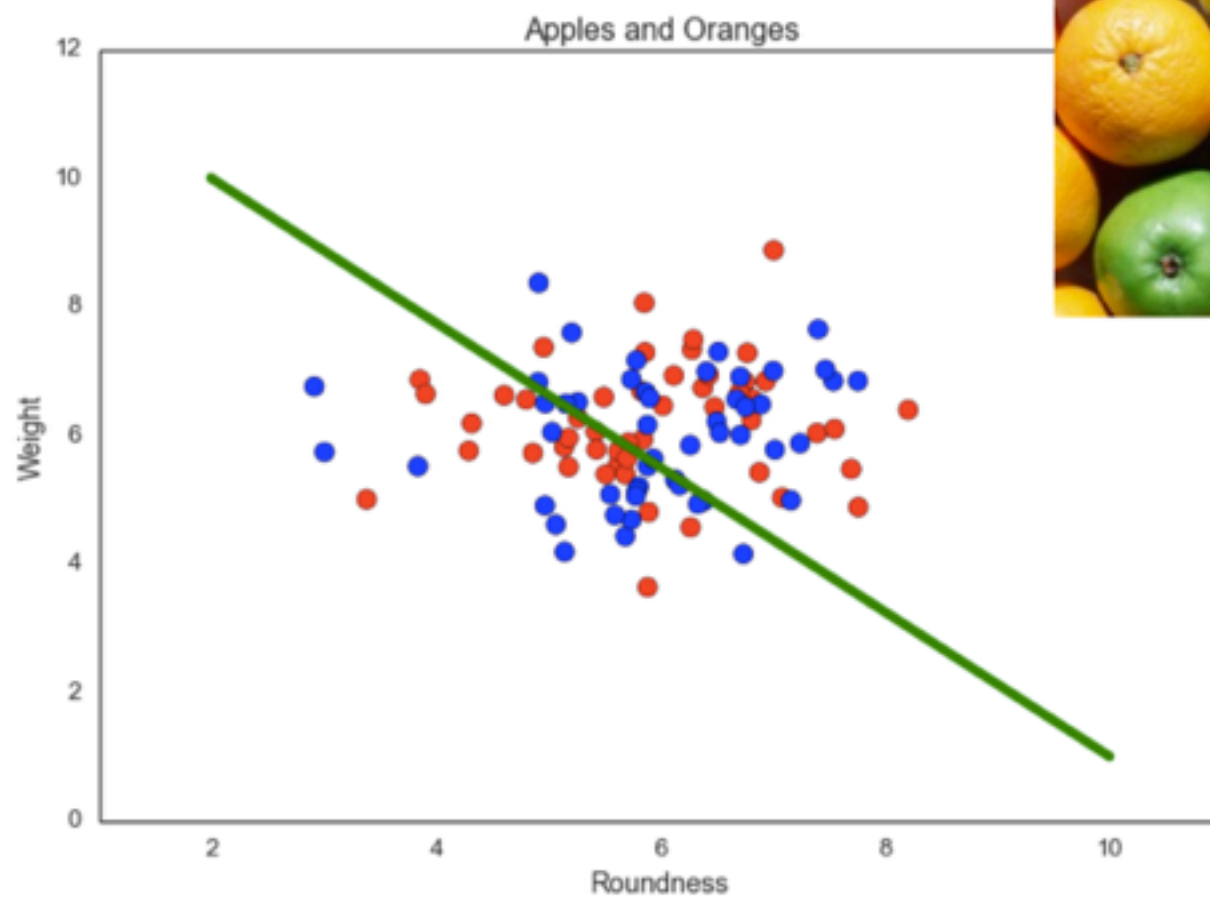
Discovering Good Preprocessing
Easily 80% of Work

Finding Features (preprocessing?)

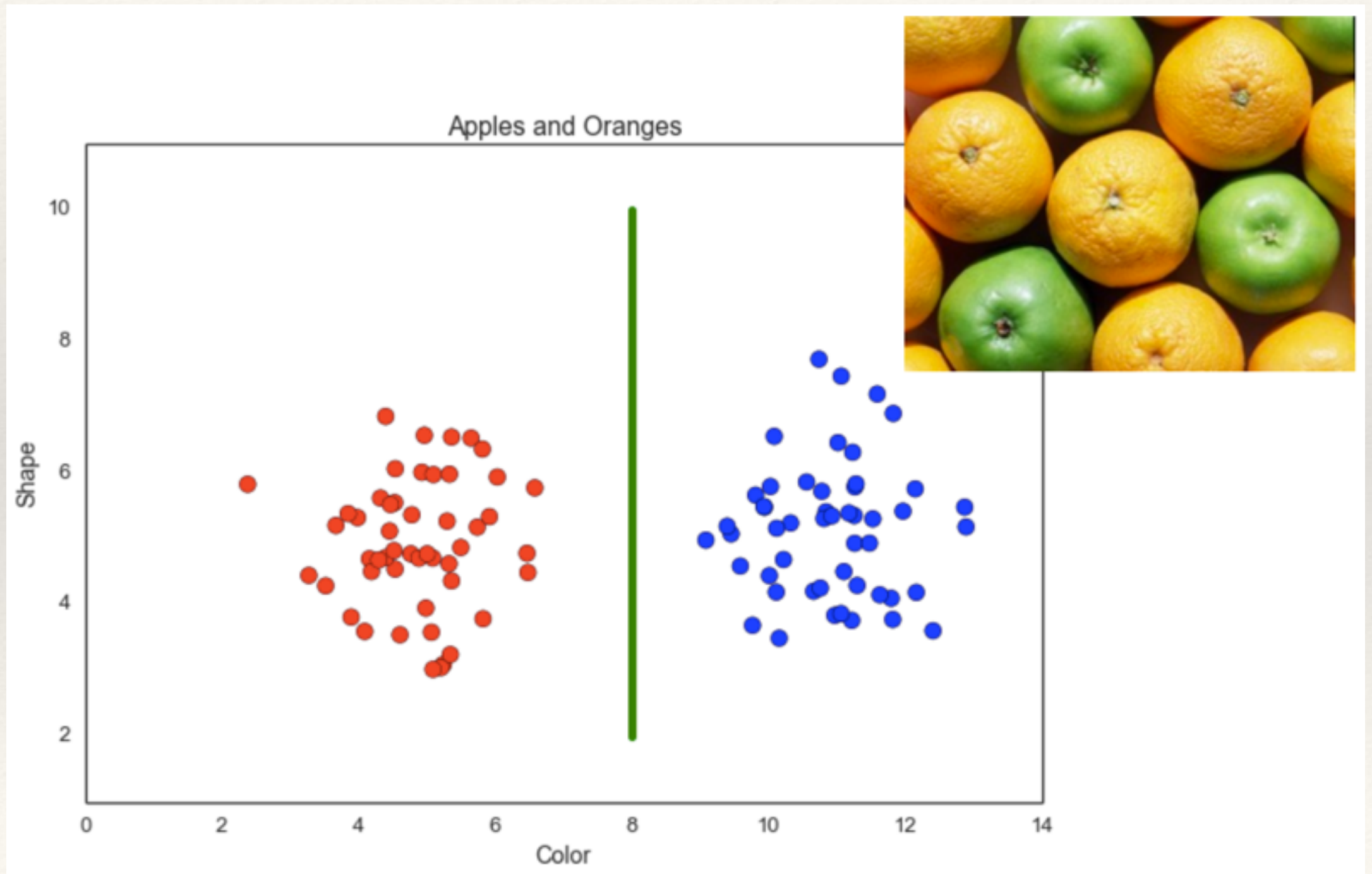


Bad Features

Feature Selection



Good Features



Finding Good Features

Still Can be Alchemy....

Some Approaches:

- Dimension Reduction

 - PCA

 - ICA

 - Matrix Factorization

 - LDA

- Anova

- Lasso

- Deep Learning (feature learning)

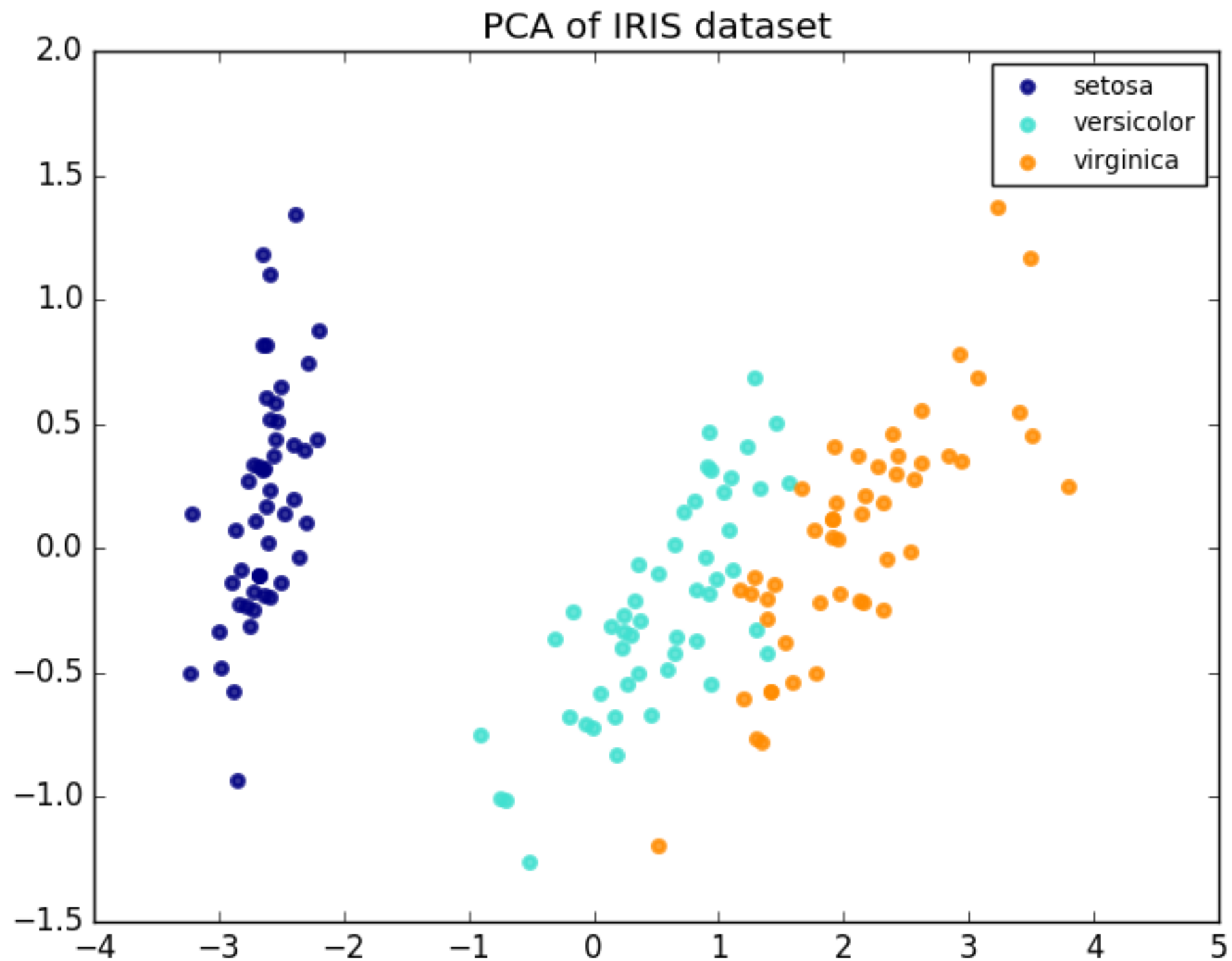
Iris Data Set



Attribute Information:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class:
 - Iris Setosa
 - Iris Versicolour
 - Iris Virginica

PCA



Classification: KNN

Two Simple (“Brain Dead”) algorithms:

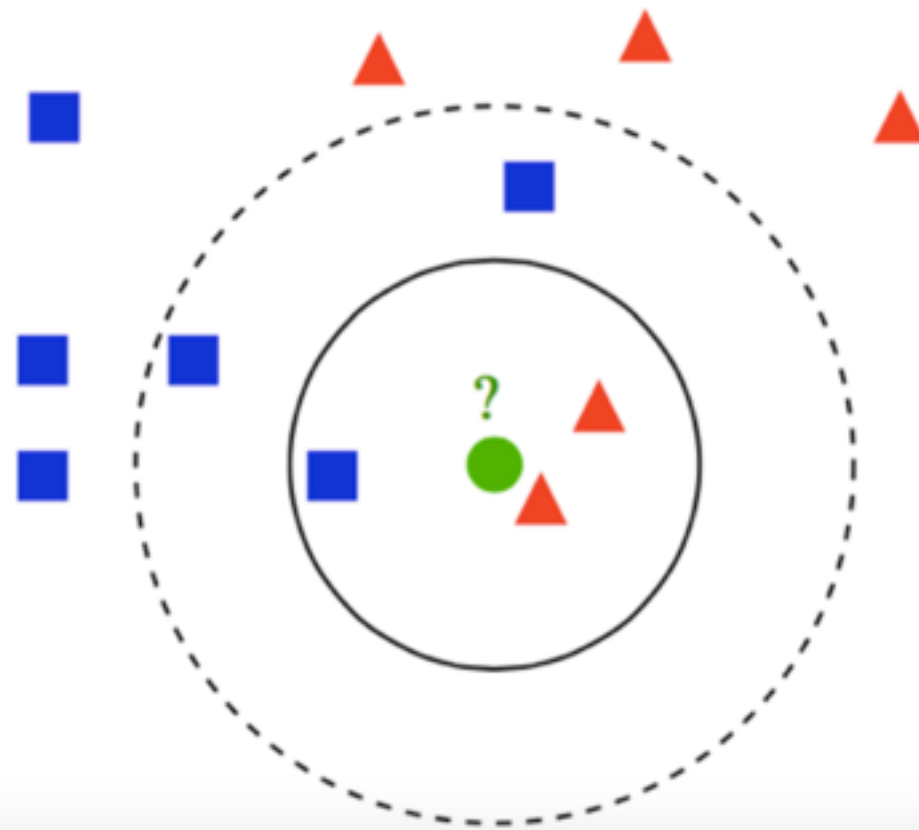
KNN

Linear Discriminant Analysis

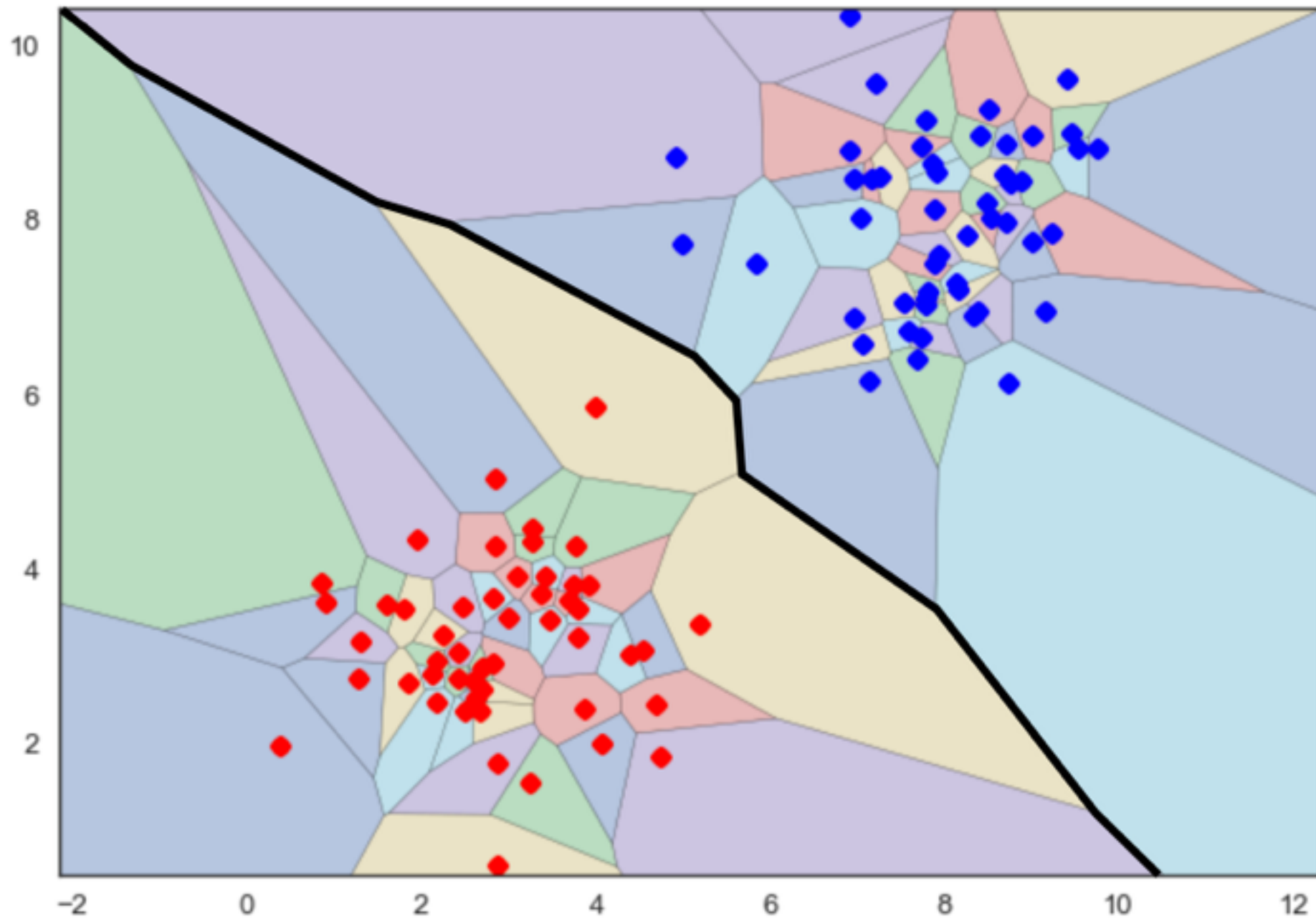
Always Try this first!!!!

KNN

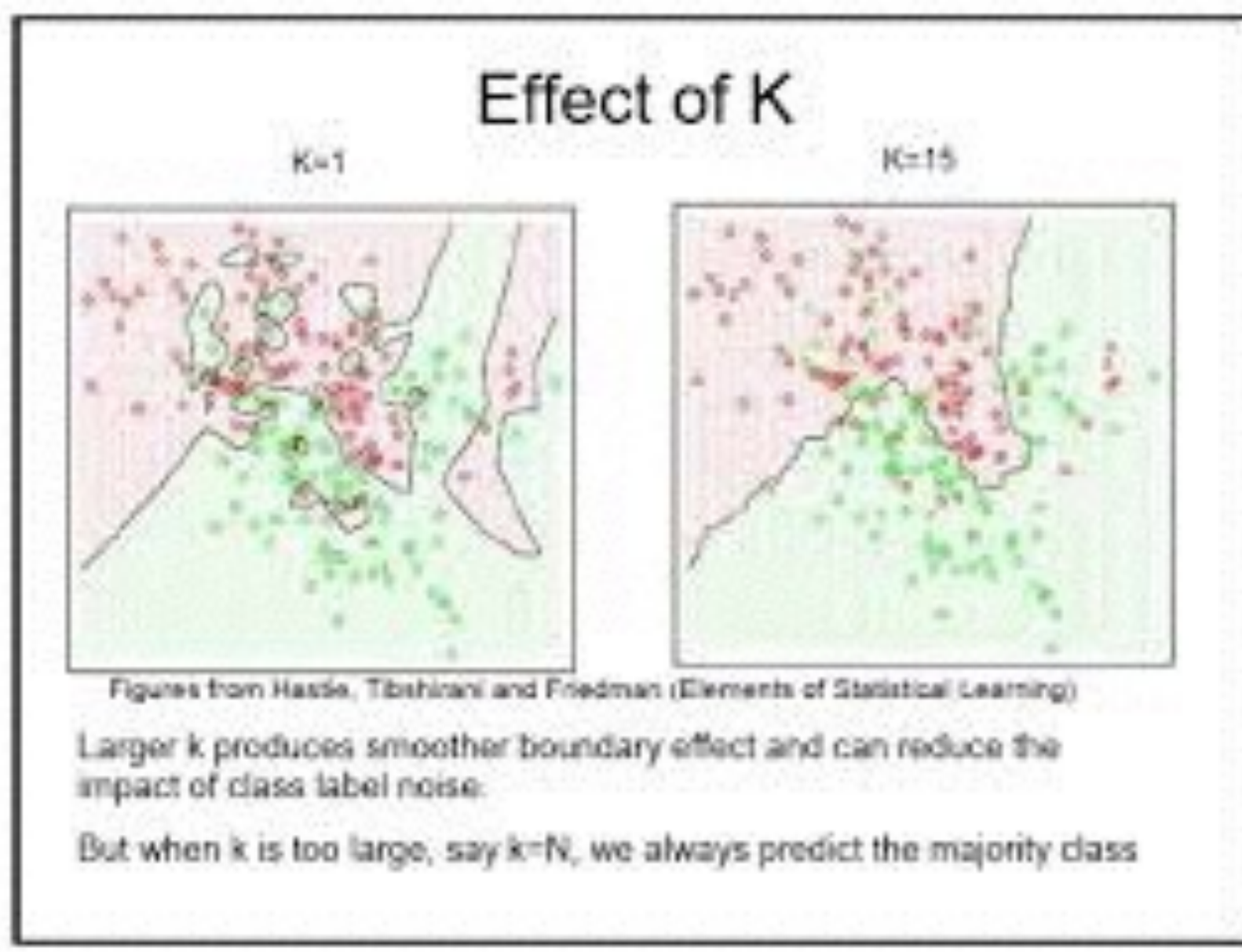
Predict class of new data point by majority vote of K nearest neighbors



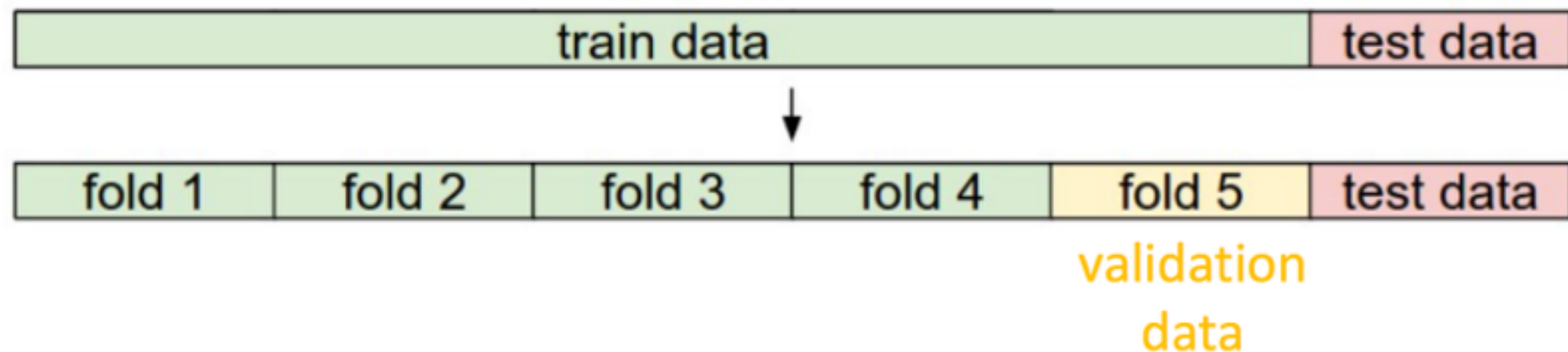
Voronoi Cells



Effect of K



Evaluate With Cross Validation



- **Training data:** train classifier
- **Validation data:** estimate (hyper) parameters (k)
- **Test data:** measure performance