

AIN2002 – Introduction to Data Science

ML Overview

Dr. Fatih KAHRAMAN

BAU AI Engineering Department

fatih.kahraman@bau.edu.tr



Machine Learning Problems

Supervised Learning

Unsupervised Learning

Discrete
Continuous

classification or
categorization

clustering

regression

dimensionality
reduction



Clustering Strategies

- K-means
 - Iteratively re-assign points to the nearest cluster center
- Agglomerative clustering
 - Start with each point as its own cluster and iteratively merge the closest clusters
- Mean-shift clustering
 - Estimate modes of pdf
- Spectral clustering
 - Split the nodes in a graph based on assigned links with similarity weights

As we go down this chart, the clustering strategies have more tendency to transitively group points even if they are not nearby in feature space

Machine Learning Problems

Supervised Learning

Unsupervised Learning

Discrete
Continuous

classification or
categorization

clustering

regression

dimensionality
reduction

The machine learning framework

- Apply a prediction function to a feature representation of the image to get the desired output:

$f(\text{ } \img alt="apple" data-bbox="315 510 400 620" \text{ })$ “apple”

$f(\text{ } \img alt="tomato" data-bbox="315 665 400 775" \text{ })$ “tomato”

$f(\text{ } \img alt="cow" data-bbox="315 830 400 942" \text{ })$ “cow”

The machine learning framework

$$y = f(x)$$

A diagram illustrating the machine learning equation $y = f(x)$. The equation is written in large blue font. Below it, three labels are positioned: 'output' under 'y', 'prediction function' under 'f', and 'Image feature' under 'x'. Red arrows point from each label to its corresponding symbol in the equation: an upward arrow from 'output' to 'y', an upward arrow from 'prediction function' to 'f', and a diagonal arrow from 'Image feature' to 'x'.

- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, estimate the prediction function f by minimizing the prediction error on the training set
- **Testing:** apply f to a never before seen *test example* \mathbf{x} and output the predicted value $y = f(\mathbf{x})$

Steps

Training

Training
Images



Image
Features



Training
Labels



Training



Learned
model

Testing



Test Image



Image
Features



Learned
model



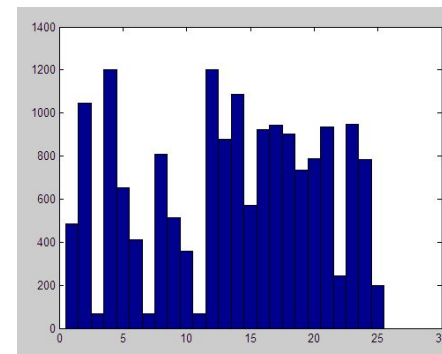
Prediction

Features

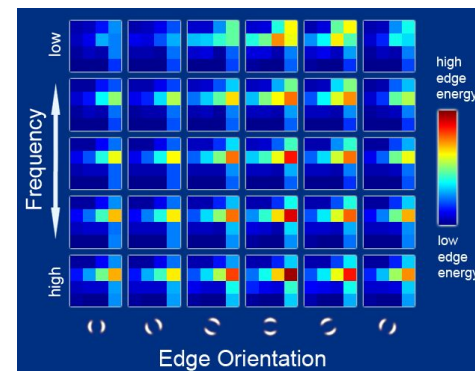
- Raw pixels



- Histograms

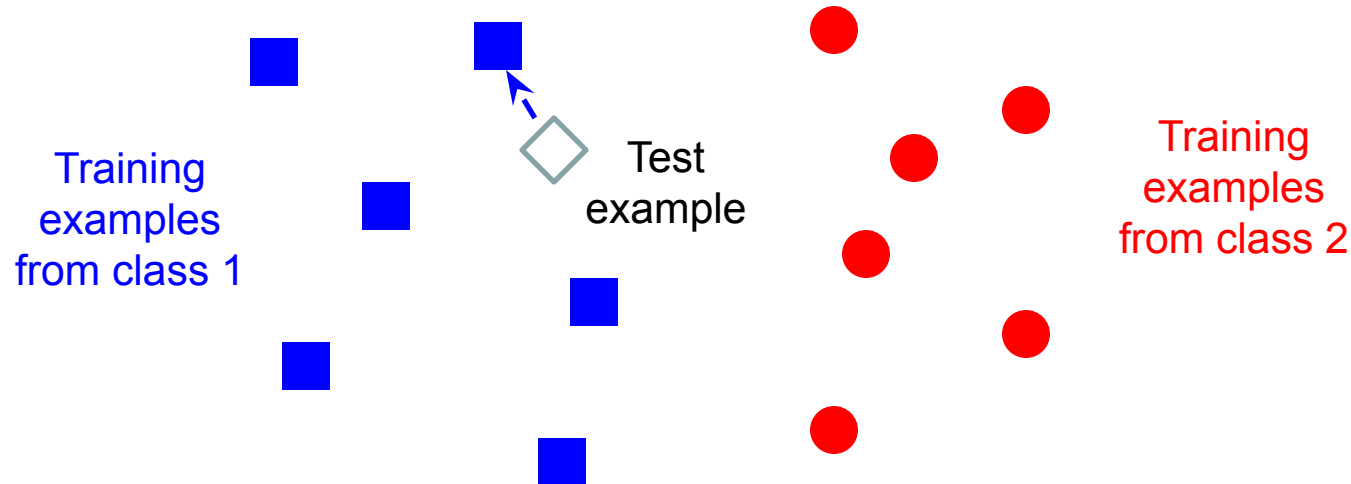


- GIST descriptors



- ...

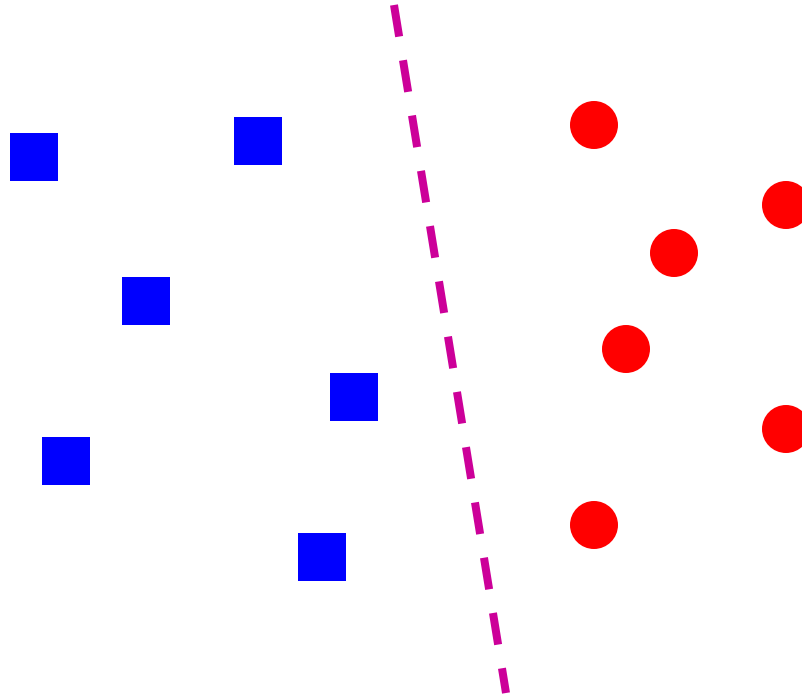
Classifiers: Nearest neighbor



$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

- All we need is a distance function for our inputs
- No training required!

Classifiers: Linear



- Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

Many classifiers to choose from

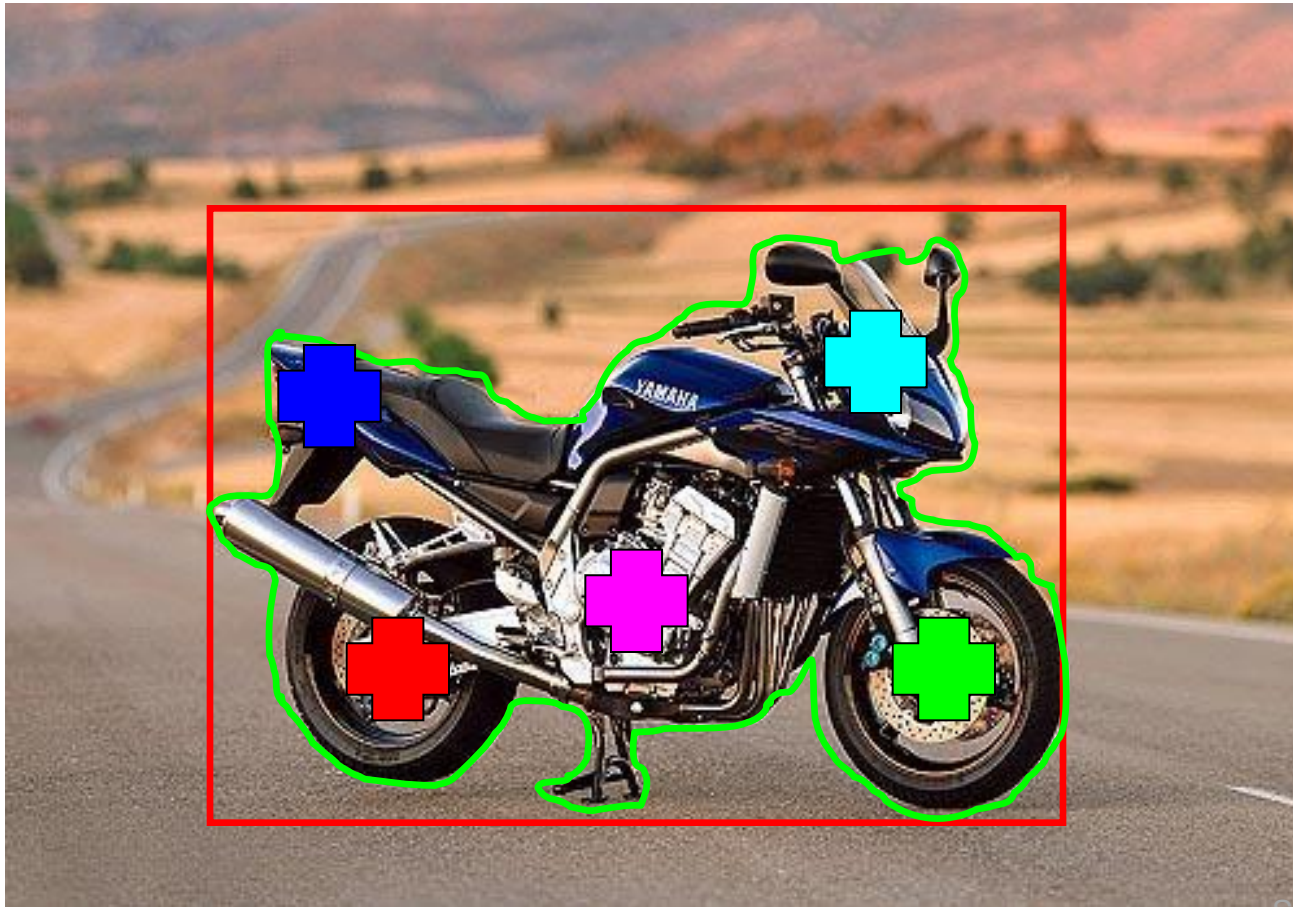
- SVM
- Neural networks
- Naïve Bayes
- Bayesian network
- Logistic regression
- Randomized Forests
- Boosted Decision Trees
- K-nearest neighbor
- RBMs
- Etc.

Which is the best one?

Recognition task and supervision

- Images in the training set must be annotated with the “correct answer” that the model is expected to produce

Contains a motorbike



What to remember about classifiers

- Try simple classifiers first
- Better to have smart features and simple classifiers than simple features and smart classifiers
- Use increasingly powerful classifiers with more training data