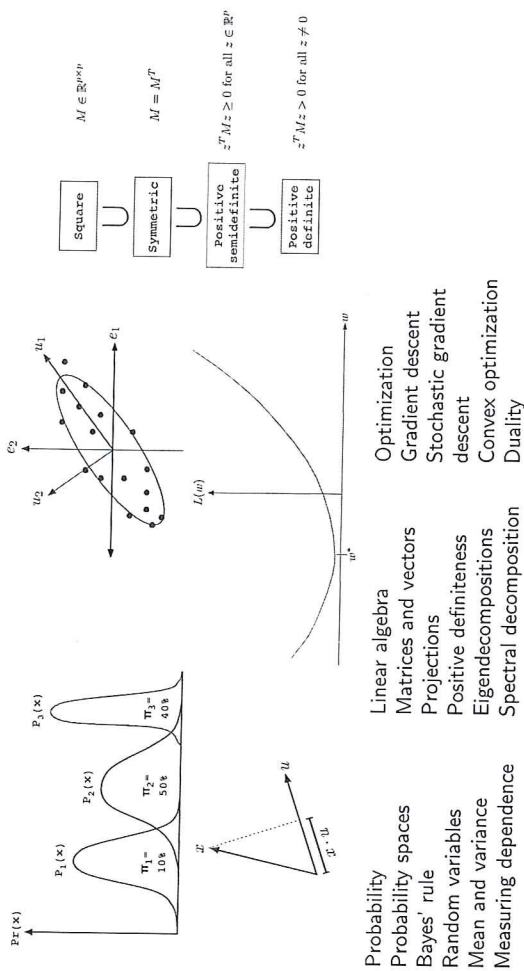


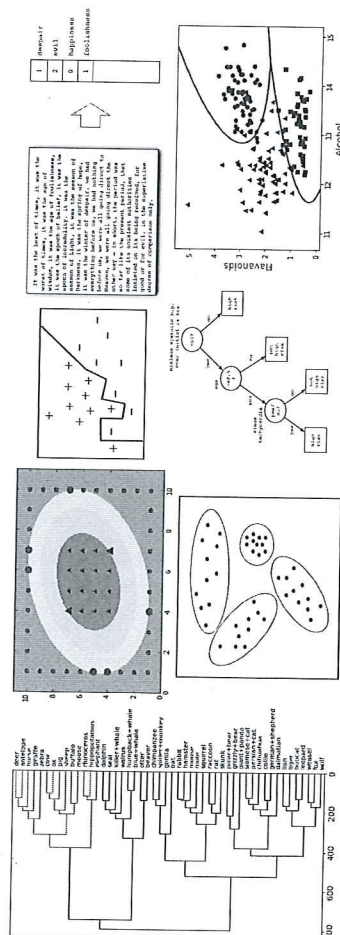
WEEK 11 18 June 2018

DSE 220X: Fundamentals of Machine Learning



Skills you will acquire

- 1 Familiarity with most widely-used ML methods
 - How they work
 - The kinds of data they are suited to
 - Their strengths and weaknesses
- 2 Adapting existing methods to a particular application
- 3 The foundational knowledge to keep pace with a fast-moving field



- Nearest neighbor
- Generative models
- Least-squares regression
- Ridge regression, Lasso
- Logistic regression
- Support vector machines
- Kernel methods
- Decision trees
- Boosting and bagging
- Random forests
- k-means
- Mixtures of Gaussians
- Hierarchical clustering
- Principal component analysis
- Singular value decomposition
- Autoencoders
- Deep learning

```

1  # Import the necessary modules
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from sklearn.preprocessing import StandardScaler
6  from sklearn.metrics import mean_squared_error
7  from sklearn.linear_model import LinearRegression
8  from sklearn.ensemble import RandomForestRegressor
9  from sklearn.svm import SVR
10 from sklearn.neural_network import MLPRegressor
11 from sklearn.metrics import r2_score
12
13 # Load the data
14 data = pd.read_csv('data.csv')
15
16 # Split the data into training and testing sets
17 X_train = data[['X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X7', 'X8', 'X9', 'X10']]
18 y_train = data['Y']
19 X_test = data[['X1', 'X2', 'X3', 'X4', 'X5', 'X6', 'X7', 'X8', 'X9', 'X10']]
20 y_test = data['Y']
21
22 # Standardize the features
23 scaler = StandardScaler()
24 X_train = scaler.fit_transform(X_train)
25 X_test = scaler.transform(X_test)
26
27 # Train the Linear Regression model
28 lr = LinearRegression()
29 lr.fit(X_train, y_train)
30
31 # Predict the values using the Linear Regression model
32 y_pred_lr = lr.predict(X_test)
33
34 # Calculate the Mean Squared Error (MSE) for the Linear Regression model
35 mse_lr = mean_squared_error(y_test, y_pred_lr)
36
37 # Train the Random Forest model
38 rf = RandomForestRegressor()
39 rf.fit(X_train, y_train)
40
41 # Predict the values using the Random Forest model
42 y_pred_rf = rf.predict(X_test)
43
44 # Calculate the Mean Squared Error (MSE) for the Random Forest model
45 mse_rf = mean_squared_error(y_test, y_pred_rf)
46
47 # Train the Support Vector Regression (SVR) model
48 svr = SVR()
49 svr.fit(X_train, y_train)
50
51 # Predict the values using the SVR model
52 y_pred_svr = svr.predict(X_test)
53
54 # Calculate the Mean Squared Error (MSE) for the SVR model
55 mse_svr = mean_squared_error(y_test, y_pred_svr)
56
57 # Train the Multi-Layer Perceptron (MLP) model
58 mlp = MLPRegressor()
59 mlp.fit(X_train, y_train)
60
61 # Predict the values using the MLP model
62 y_pred_mlp = mlp.predict(X_test)
63
64 # Calculate the Mean Squared Error (MSE) for the MLP model
65 mse_mlp = mean_squared_error(y_test, y_pred_mlp)
66
67 # Compare the results
68 print('MSE for Linear Regression: ', mse_lr)
69 print('MSE for Random Forest: ', mse_rf)
70 print('MSE for SVR: ', mse_svr)
71 print('MSE for MLP: ', mse_mlp)
72
73 # Plot the results
74 plt.figure(figsize=(10, 10))
75 plt.scatter(X_test, y_test, label='Actual Values')
76 plt.plot(X_test, y_pred_lr, label='Linear Regression')
77 plt.plot(X_test, y_pred_rf, label='Random Forest')
78 plt.plot(X_test, y_pred_svr, label='SVR')
79 plt.plot(X_test, y_pred_mlp, label='MLP')
80 plt.legend()
81 plt.show()

```

Week 1-2

12 June 2018

Nearest neighbor classification

The problem we'll solve today

Given an image of a handwritten digit, say which digit it is.



\Rightarrow 3

The number of loops / straight line & their positions to decide which digit it is.



Some more examples:

Problems. Handwriting are super noisy / a lot of variability in the way people write 4, 7, 9 etc.

Topics we'll cover

- 1 What is a classification problem?
- 2 The training set and test set
- 3 Representing data as vectors
- 4 Distance in Euclidean space
- 5 The 1-NN classifier
- 6 Training error versus test error
- 7 The error of a random classifier

The machine learning approach

Assemble a data set:

i machine needs a huge amount of data (MNIST dataset handwritten image + correct digit)

14161191548577268032264141
8663597202942947225100467
0130841145910106154061036
311044111030473262097799
66891208670855571214279554
6020177801871129930879709
8401097075973319720155190
5610755182554603546035460
6317875416554603546035460
5518255108503047520439401

An entirely different approach. - ML.

Rather than trying to figure out the underlying patterns

ourselves, we'll just let the machine do it for us.

The MNIST data set of handwritten digits:

- ii • Training set of 60,000 images and their labels. to learn a classifier
- iii • Test set of 10,000 images and their labels. to test errors

And let the machine figure out the underlying patterns.

a function that

takes an image and then outputs

what digit is

Nearest neighbor classification

Training images $x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(60000)}$
 Labels $y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(60000)}$ are numbers in the range 0 - 9

```

1416111915485732680322644141
86235972029929472223100467
01308444459101061540610366
3110641110304752220099799
6689180767285571214279554
60201872801871129930899709
84010755182551825518255190
4317875416554603546035460
55182551085030475220439401
    
```

Input: New Image x . How to classify a new image x ?

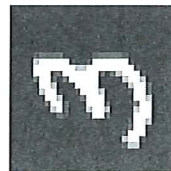
- Find its nearest neighbor amongst the $x^{(i)}$
- Output: Return $y^{(i)}$



Questions: Distance checking for the one that's closest to x , so we have some notion of distance between images.

The data space

How to measure the distance between images?



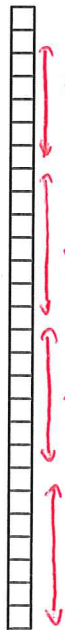
i. represent image as vector.

MNIST images:

- Size 28×28 (total: 784 pixels)
- Each pixel is grayscale: 0-255

Black white.

ii Stretch each image into a vector with 784 coordinates: **Dimensional Vector.**



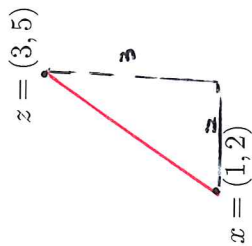
iii Data space $\mathcal{X} = \mathbb{R}^{784}$ Euclidean space.

Label space $\mathcal{Y} = \{0, 1, \dots, 9\}$

Questions: How to compute the distance between vectors?

The distance function (the most common (default)).

Remember Euclidean distance in two dimensions?



distance m, z = length of the line.

$$= \sqrt{2^2 + 3^2}$$

$$= \sqrt{13}$$

Euclidean Distance in 2 Dimensions.

Compare the distance between 2 vectors x, z ,

Simply find out how much differ on each individual coordinates.

Euclidean distance in higher dimension

Euclidean distance between 784-dimensional vectors x, z is

$$\|x - z\| = \sqrt{\sum_{i=1}^{784} (x_i - z_i)^2}$$

Here x_i is the i th coordinate of x .

Nearest neighbor classification

Training images $x^{(1)}, \dots, x^{(60000)}$, labels $y^{(1)}, \dots, y^{(60000)}$

1 4 1 6 1 1 9 1 5 4 8 5 7 3 6 8 0 3 2 2 6 4 1 4 1
8 6 4 3 5 9 7 2 0 2 9 9 2 9 9 7 2 2 5 1 0 0 4 6 7
0 1 3 0 8 4 1 1 1 0 3 0 4 7 5 2 6 2 0 9 7 7 9 9
6 6 8 9 1 2 0 8 6 7 0 8 5 5 7 1 2 1 4 2 7 9 5 5 4
6 0 2 0 1 3 7 3 0 1 8 7 1 1 2 9 7 3 0 1 5 5 1 9 0
8 4 0 1 0 9 7 0 7 5 9 7 3 3 1 9 7 3 0 1 5 5 1 9 0
6 5 1 0 7 5 5 1 8 2 5 5 1 8 2 5 5 1 8 2 5 5 1 9 0
4 3 1 7 8 7 3 4 1 6 5 5 4 6 6 6 6 6 6 6 6 6 6 6
5 5 1 8 2 5 5 1 0 8 5 0 3 0 4 7 5 2 0 4 3 4 4 0 1

To classify a new image x :

- Find its nearest neighbor amongst the $x^{(i)}$ using Euclidean distance in \mathbb{R}^{784}
- Return $y^{(i)}$



How accurate is this classifier?

Accuracy of nearest neighbor on MNIST

Training set of 60,000 points.

- What is the error rate on training points? Zero.
In general, training error is an overly optimistic predictor of future performance.
- A better gauge: separate test set of 10,000 points.
Test error = fraction of test points incorrectly classified.
- What test error would we expect for a *random classifier*?
(One that picks a label 0 — 9 at random?) **90%.**
- Test error of nearest neighbor: **3.09%.** — *Out of 10,000 points, it gets 309 of them wrong.*

Examples of errors

Test set of 10,000 points:

- 309 are misclassified
- Error rate 3.09%

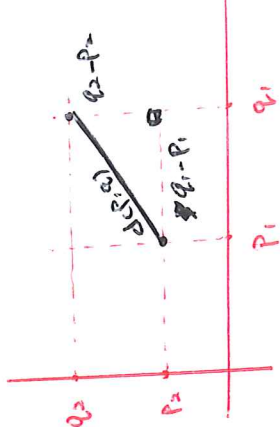
Examples of errors:

Query    
NN    

More: 2-dimensional Euclidean Distance.

if $P = (P_1, P_2)$ and $Q = (Q_1, Q_2)$ the the distance is

$$d(P, Q) = \sqrt{(Q_1 - P_1)^2 + (Q_2 - P_2)^2}$$



$$d(P, Q)^2 = (Q_1 - P_1)^2 + (Q_2 - P_2)^2$$

Pythagorean theorem.

Week 1-3. 12/June 20:30

Improving the performance of nearest neighbor

K-nearest neighbor classification

To classify a new point:

- i Find the k nearest neighbors in the training set.
- ii Return the most common label amongst them.

Error Rate. stop goes down.

k	1	3	5	7	9	11
Test error (%)	3.09	2.94	3.13	3.10	3.43	3.34

In real life, there's no test set. How to decide which k is best?

finding 3 closest images in the train set and return the majority label.

Recall: nearest neighbor on MNIST

- Images of handwritten digits, represented as vectors in \mathbb{R}^{784} .
- Labels 0 - 9
- Training set: 60,000 points; test set: 10,000 points

Test error of nearest neighbor using Euclidean distance: 3.09%

Examples of errors:

Query					
NN					

Ideas for improvement: (1) k -NN (2) better distance function. To improve the performance.

Cross-validation - to select the value of k .

How to estimate the error of k -NN for a particular k ?

10-fold cross-validation k -fold Cross Validation.

- Input i • Divide the training set into 10 equal pieces. Randomly divide the data into k equal sized 'folds'.

Training set (call it S): 60,000 points

Call the pieces S_1, S_2, \dots, S_{10} : 6,000 points each.

- Input ii • For each piece S_i :

- Classify each point in S_i using k -NN with training set $S - S_i$
- Let ϵ_i = fraction of S_i that is incorrectly classified

- Output iii • Take the average of these 10 numbers: Report the classification error as the average of the misclassification rates

$$\text{estimated error with } k\text{-NN} = \frac{\epsilon_1 + \dots + \epsilon_{10}}{10}$$

Another improvement: better distance functions

The Euclidean (ℓ_2) distance between these two images is very high!



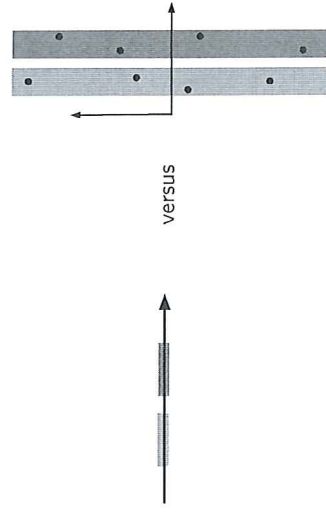
Much better idea: distance measures that are invariant under:

- Small translations and rotations. e.g. tangent distance.
- A broader family of natural deformations. e.g. shape context.

Test error rates:	ℓ_2	tangent distance	shape context
	3.09	1.10	0.63

Related problem: feature selection Prior - using KNN/NN

Feature selection/reweighting is part of picking a distance function. And, one noisy feature can wreak havoc with nearest neighbor!



Algorithmic issue: speeding up NN search

Naive search takes time $O(n)$ for training set of size n : slow!

Luckily there are data structures for speeding up nearest neighbor search, like:

- 1 Locality sensitive hashing
- 2 Ball trees
- 3 K-d trees

These are part of standard Python libraries for NN, and help a lot.

Intro. KNN Classification algorithm.

• Given. - training examples $\{(x_i, y_i)\}$

- New point x that we want to classify

• Algorithm. - Compute distance $D(x, x_i)$ to every training example x_i :

- select k closest instances x_{i_1}, \dots, x_{i_k} and their labels

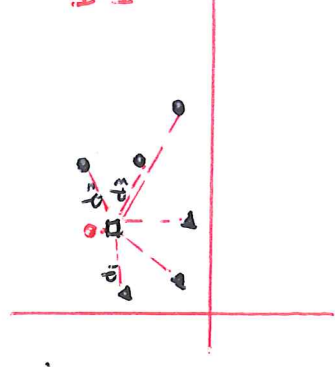
y_{i_1}, \dots, y_{i_k}

- Output the class y^* which is most frequent in

y_{i_1}, \dots, y_{i_k} .

$k=1$

$k=3$



12/ June / 2018.

Useful distance functions for machine learning

Usual choice: Euclidean distance: **1.2**

$$\|x - z\|_2 = \sqrt{\sum_{i=1}^m (x_i - z_i)^2}.$$

For $p \geq 1$, here is ℓ_p distance:

$$\|x - z\|_p = \left(\sum_{i=1}^m |x_i - z_i|^p \right)^{1/p}$$

- $p = 2$: Euclidean distance

- $$\ell_1 \text{ distance: } \|x - z\|_1 = \sum_{i=1}^m |x_i - z_i|$$

- $$\ell_\infty \text{ distance: } \|x - z\|_\infty = \max_i |x_i - z_i|$$

→ finding distance between two coordinates
work at the difference along each coordinate,
and simply add up those differences.

Topics we'll cover

Consider the all-ones vector $(1, 1, \dots, 1)$ in \mathbb{R}^d . What are its ℓ_2 , ℓ_1 , and ℓ_∞ length?

2 families of distance function.

- ① L_p norms
- ② Metric spaces

cell @ points

$$\|x\|_2 = \sqrt{1^2 + 2^2 + \dots + 1^2}$$

$$d_1, d_2, \dots, d_n = |x_1| + \dots + |x_n|$$

$\frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2}$

بسم الله الرحمن الرحيم

diff.

Example 2

In \mathbb{R}^2 , draw all points with:

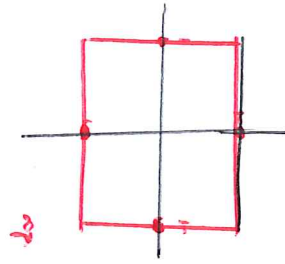
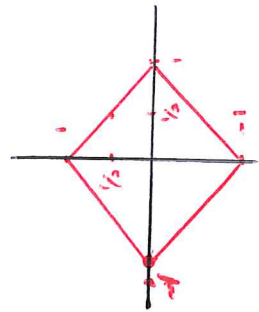
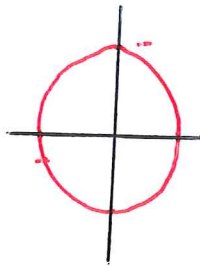
- ① ℓ_2 length 1
- ② ℓ_1 length 1
- ③ ℓ_∞ length 1

- 2-Dimensional plot

- want to draw all points where

$d_2(x_1, x_2)$ length is 1.

$$d_2(x_1, x_2) = \sqrt{x_1^2 + x_2^2} = 1$$



$$d_1(x_1, x_2) = |x_1| + |x_2| = 1$$

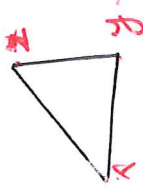
Metric spaces

Other distance function needed, e.g. distance between objects (string graphs)

Let \mathcal{X} be the space in which data lie.

A distance function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a **metric** if it satisfies these properties:

- $d(x, y) \geq 0$ (nonnegativity)
- $d(x, y) = 0$ if and only if $x = y$
- $d(x, y) = d(y, x)$ (symmetry)
- $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality)



Example 1

$$\mathcal{X} = \mathbb{R}^m \text{ and } d(x, y) = \|x - y\|_p$$

Check:

- $d(x, y) \geq 0$ (nonnegativity) ✓
- $d(x, y) = 0$ if and only if $x = y$ ✓
- $d(x, y) = d(y, x)$ (symmetry) ✓
- $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality) ✓

d_1 in Metric Space.

$$d_1(x, y) = \sum_{i=1}^m |x_i - y_i|$$

d_1 distance satisfies all the properties of a metric, and all the CP distances also.

Example 2

$\mathcal{X} = \{\text{strings over some alphabet}\}$ and $d = \text{edit distance}$

Check:

- $d(x, y) \geq 0$ (nonnegativity) ✓
- $d(x, y) = 0$ if and only if $x = y$ ✓
- $d(x, y) = d(y, x)$ (symmetry) ✓
- $d(x, z) \leq d(x, y) + d(y, z)$ (triangle inequality) ✓

- Input space doesn't consist of vectors.

• $\mathcal{X} = \{A, C, G, T\}^+$ strings over alphabet. Input.

* means the strings with arbitrary length.

- $x = A C G T$. 2 strings we have.

$y = C C G T$. distance between the 2 strings.

$d(x, y) = \#$ of insertions, deletions, and substitutions from x to y to be a metric.

A non-metric distance function

Let p, q be probability distributions on some set \mathcal{X} .

The Kullback-Leibler divergence or **relative entropy** between p, q is:

$$d(p, q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}.$$

\Downarrow

$$\begin{aligned} p &= (1/2, 1/4, 1/3, 1/3) & d(p, q) &= \frac{1}{2} \log \frac{1/2}{1/6} + \frac{1}{4} \log \frac{1/4}{1/3} + \frac{1}{3} \log \frac{1/3}{1/3} \\ q &= (1/6, 1/3, 1/3, 1/6) & &+ \frac{1}{3} \log \frac{1/3}{1/6} \end{aligned}$$

week 1-5. 12/Jan/2028

A host of prediction problems

Machine learning versus Algorithms

A central goal of both fields:

develop procedures that exhibit a desired input-output behavior.

- Algorithms: input-output mapping can be precisely defined.

Input: Graph G , two nodes u, v in the graph.

Output: Shortest path from u to v in G

- Machine learning: mapping cannot easily be made precise.

Input: Picture of an animal.

Output: Name of the animal.

Instead, provide examples of (input,output) pairs.

Ask the machine to *learn* a suitable mapping itself.

Topics we'll cover

- ① Machine learning versus algorithms
- ② A taxonomy of prediction problems
- ③ Roadmap for the course

Prediction problems: inputs and outputs

Basic terminology:

- The input space, \mathcal{X} .

E.g. 32×32 RGB images of animals.

- The output space, \mathcal{Y} .

E.g. Names of 100 animals.



x :

y : "bear"

After seeing a bunch of examples (x, y) , pick a mapping

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

that accurately recovers the input-output pattern of the examples.

Categorize prediction problems by the type of output space:

(1) discrete, (2) continuous, or (3) probability values

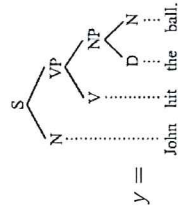
Discrete output space: classification

Binary classification	Multiclass
E.g., Spam detection $\mathcal{X} = \{\text{email messages}\}$ $\mathcal{Y} = \{\text{spam, not spam}\}$	E.g., News article classification $\mathcal{X} = \{\text{news articles}\}$ $\mathcal{Y} = \{\text{politics, business, sports, ...}\}$

Structured outputs

E.g., Parsing
 $\mathcal{X} = \{\text{sentences}\}$
 $\mathcal{Y} = \{\text{parse trees}\}$

$x = \text{"John hit the ball"}$



$y =$

Probability estimation

$\mathcal{Y} = [0, 1]$ represents **probabilities**

Example: Credit card transactions

- x = details of a transaction
- y = probability this transaction is fraudulent

Why not just treat this as a binary classification problem?

Continuous output space: regression

- Pollution level prediction
 Predict tomorrow's air quality index in my neighborhood
 $\mathcal{Y} = [0, \infty)$ (< 100 : okay, > 200 : dangerous)
- Insurance company calculations
 What is the expected life expectancy of this person?
 $\mathcal{Y} = [0, 120]$

What are suitable predictor variables (\mathcal{X}) in each case?

Roadmap for the course

- 1 Solving prediction problems
 Classification, regression, probability estimation
- 2 Representation learning
 Clustering, projection, dictionary learning, autoencoders
- 3 Deep learning