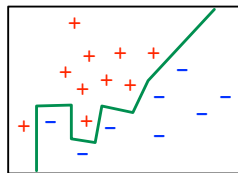
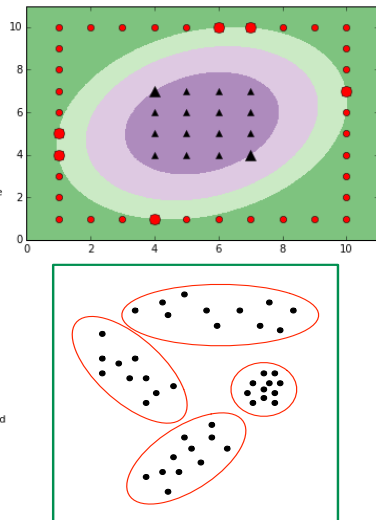
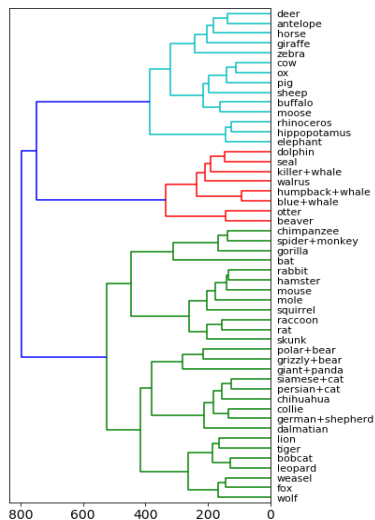


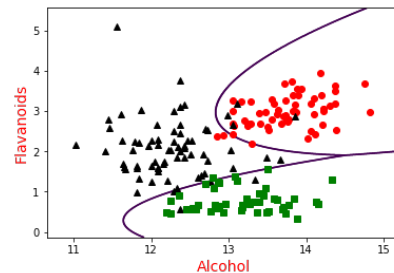
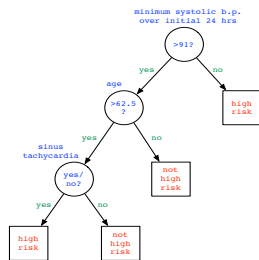
Fundamentals of Machine Learning



It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way – in short, the period was so far like the present period, that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.



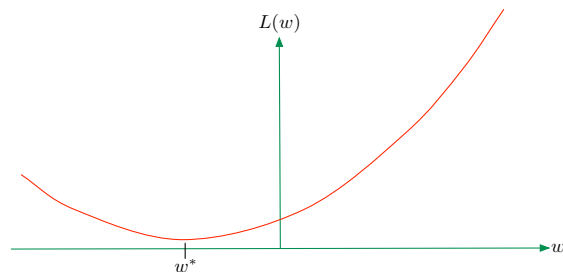
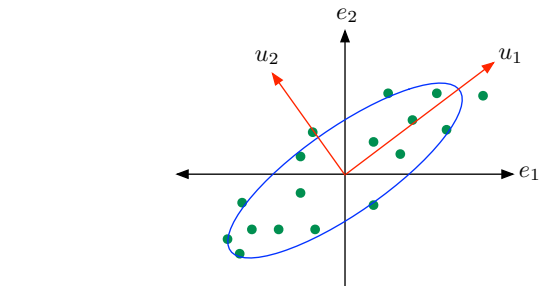
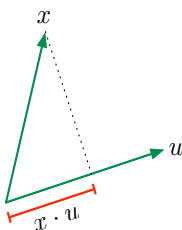
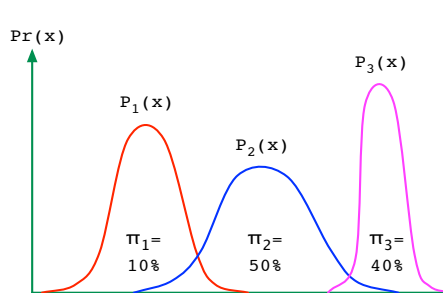
1	despair
2	evil
0	happiness
1	foolishness



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



Probability

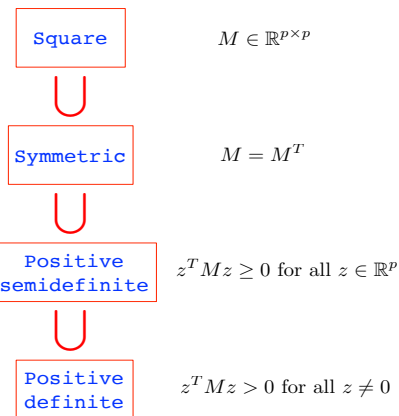
Probability spaces
Bayes' rule
Random variables
Mean and variance
Measuring dependence

Linear algebra

Matrices and vectors
Projections
Positive definiteness
Eigendecompositions
Spectral decomposition

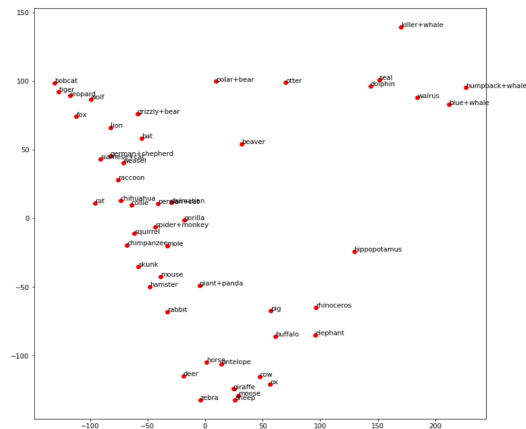
Optimization

Gradient descent
Stochastic gradient descent
Convex optimization
Duality



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

[illegible]

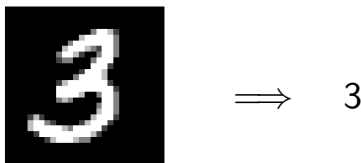
Nearest neighbor classification

Topics we'll cover

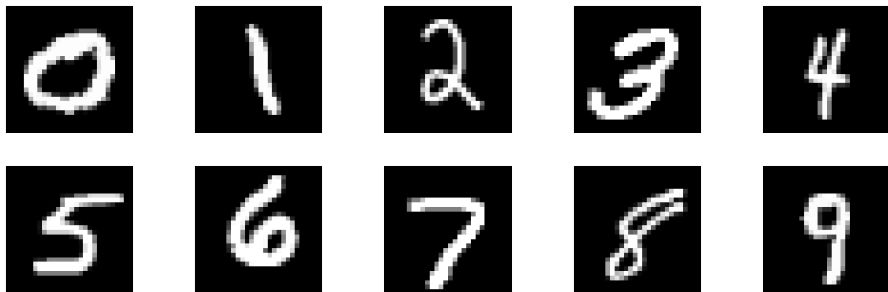
- ① What is a classification problem?
- ② The training set and test set
- ③ Representing data as vectors
- ④ Distance in Euclidean space
- ⑤ The 1-NN classifier
- ⑥ Training error versus test error
- ⑦ The error of a random classifier

The problem we'll solve today

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



Some more examples:



The machine learning approach

Assemble a data set:



The MNIST data set of handwritten digits:

- **Training set** of 60,000 images and their labels.
- **Test set** of 10,000 images and their labels.

And let the machine figure out the underlying patterns.

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

1416119134857268032264141
8663597202992997225100467
0130841115910106154061036
3110641110304752620099799
6689120867085571314279554
6010187801871129930899709
8401097075973319720155190
6610755182551828143580909
4317875216554605546035460
5518255108503047520439401



How to **classify** a new image x ?

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

The data space

How to measure the distance between images?



MNIST images:

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

Stretch each image into a vector with 784 coordinates:



- Data space $\mathcal{X} = \mathbb{R}^{784}$
- Label space $\mathcal{Y} = \{0, 1, \dots, 9\}$

The distance function

Remember Euclidean distance in two dimensions?

$$z = \underset{\bullet}{(3, 5)}$$

$$x = \overset{\bullet}{(1, 2)}$$

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 0 1 1 9 1 5 4 8 5 7 2 6 8 0 3 2 2 6 4 1 4 1
8 6 6 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 5 9 1 0 1 0 6 1 5 4 0 6 1 0 3 6
3 1 1 0 6 4 1 1 1 0 3 0 4 7 5 2 6 2 0 0 9 9 7 9 9
6 6 8 9 1 2 0 8 6 7 8 8 5 5 7 1 3 1 4 2 7 9 5 5 4
6 0 6 0 1 8 7 8 0 1 8 7 1 1 2 9 9 3 0 8 9 9 7 0 9
8 4 0 1 0 9 7 0 7 5 9 7 3 3 1 9 7 2 0 1 5 5 1 9 0
6 5 1 0 7 5 5 1 8 2 5 5 1 8 2 8 1 4 3 5 8 0 9 0 9
4 3 1 7 8 7 5 2 1 6 5 5 4 6 0 5 5 4 6 0 3 5 4 6 0
5 5 1 8 2 5 5 1 0 8 5 0 3 0 4 7 5 2 0 4 3 9 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%.**

Examples of errors

Test set of 10,000 points:

- 309 are misclassified
- Error rate 3.09%

Examples of errors:

Query



NN



Improving the performance of nearest neighbor

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.

K -nearest neighbor classification

To classify a new point:

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

MNIST:	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?

Cross-validation

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

10-fold cross-validation

- Divide the training set into 10 equal pieces.
Training set (call it S): 60,000 points
Call the pieces S_1, S_2, \dots, S_{10} : 6,000 points each.
- 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
- Take the average of these 10 numbers:

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

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:

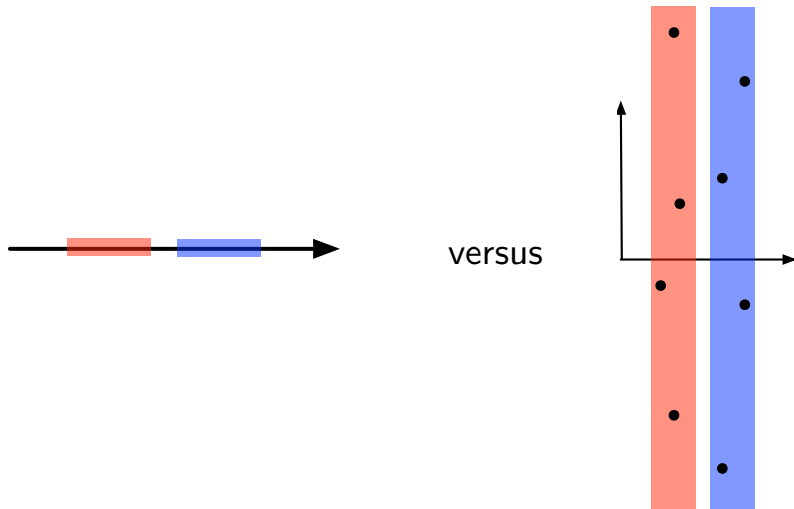
- Small translations and rotations. e.g. *tangent distance*.
- A broader family of natural deformations. e.g. *shape context*.

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

Related problem: feature selection

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.

Useful distance functions for machine learning

Topics we'll cover

① L_p norms

② Metric spaces

Measuring distance in \mathbb{R}^m

Usual choice: **Euclidean distance**:

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

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
- ℓ_1 distance: $\|x - z\|_1 = \sum_{i=1}^m |x_i - z_i|$
- ℓ_∞ distance: $\|x - z\|_\infty = \max_i |x_i - z_i|$

Example 1

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

Example 2

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

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

Metric spaces

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)

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)

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:

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

A host of prediction problems

Topics we'll cover

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

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.

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

E.g., Spam detection

$\mathcal{X} = \{\text{email messages}\}$

$\mathcal{Y} = \{\text{spam, not spam}\}$

Multiclass

E.g., News article classification

$\mathcal{X} = \{\text{news articles}\}$

$\mathcal{Y} = \{\text{politics, business, sports, ...}\}$

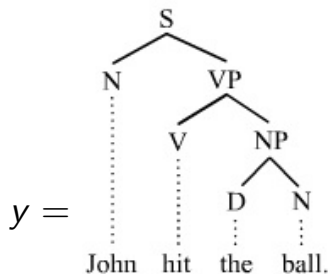
Structured outputs

E.g., Parsing

$\mathcal{X} = \{\text{sentences}\}$

$\mathcal{Y} = \{\text{parse trees}\}$

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



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?

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?

Roadmap for the course

① Solving prediction problems

Classification, regression, probability estimation

② Representation learning

Clustering, projection, dictionary learning, autoencoders

③ Deep learning