

Machine Learning Jam

A gentle introduction
to Machine Learning

The goal

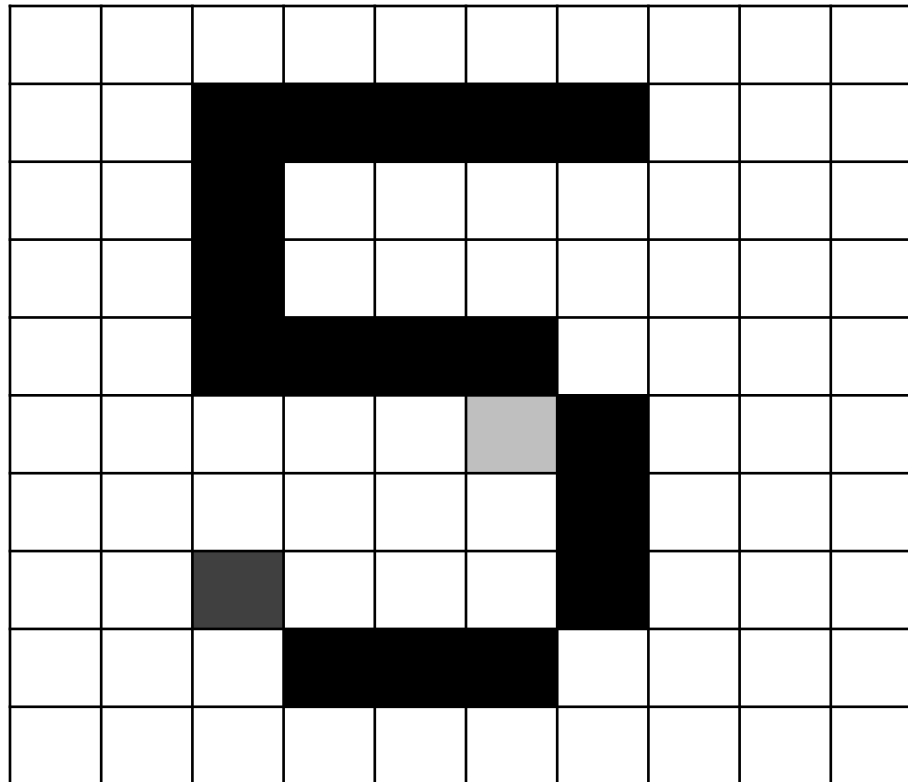
- Take a real Kaggle data science contest
- Write some code and have fun
- Write a classifier, from scratch
- Compare & contrast functional languages
- Learn some Machine Learning concepts
- Bonus goal: send results to Kaggle contest?

WHAT YOU MAY NEED TO KNOW

Kaggle Digit Recognizer contest

- Full description on [Kaggle.com](https://www.kaggle.com/c/digit-recognizer)
- Dataset: hand-written digits (0, 1, ... , 9)
- Goal = automatically recognize digits
- Training sample = 50,000 examples
- Contest: predict 20,000 “unknown” digits

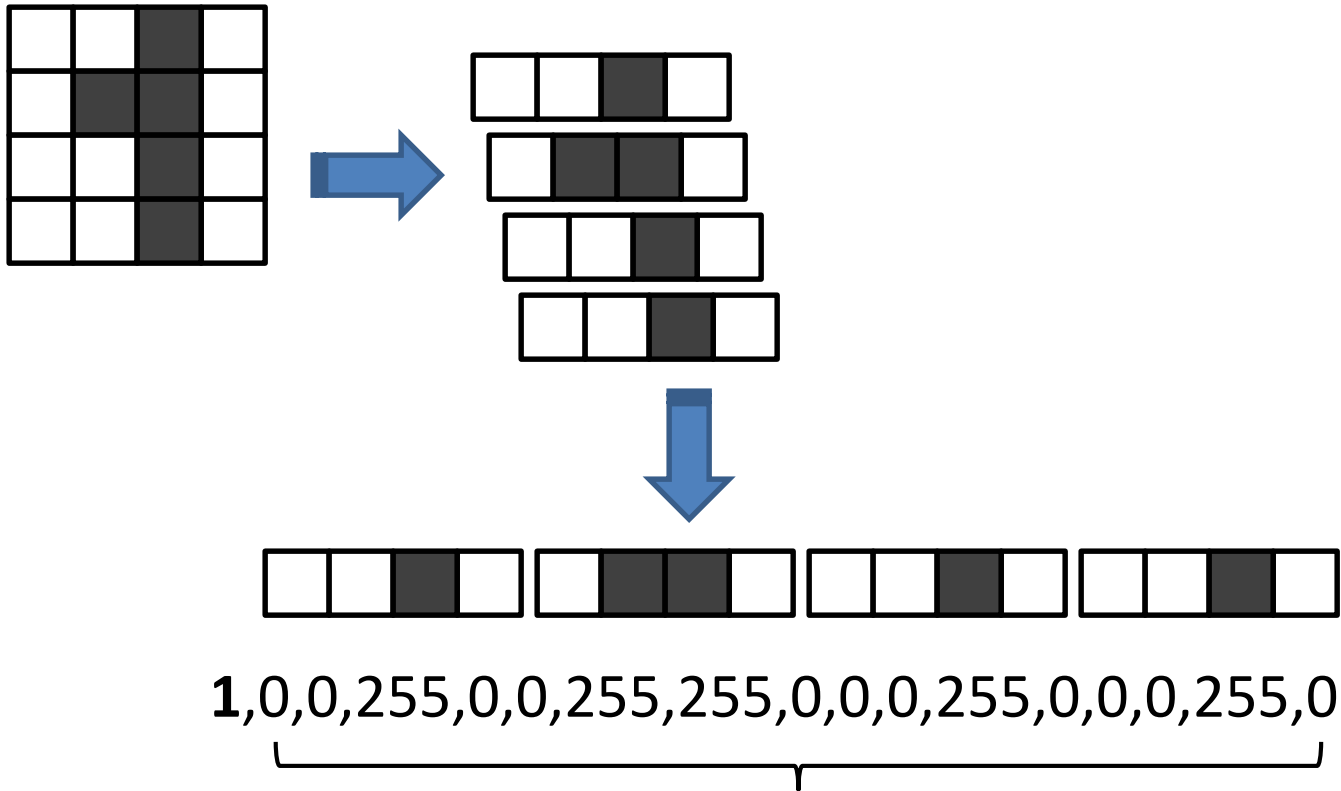
The data “looks like that”



Real data

- 28 x 28 pixels
- Grayscale: each pixel 0 (white) to 255 (black)
- Flattened: one record = Number + 784 Pixels
- CSV file

Illustration (simplified data)



Actual Number Pixels (real: 784 fields, from 0 to 255)

What's a Classifier?

- “Give me an unknown data point, and I will predict what **class** it belongs to”
- In this case, classes = 0, 1, 2, ... 9
- Unknown data point = scanned digit, without the class it belongs to

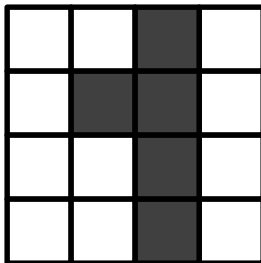
The KNN Classifier

- KNN = K-Nearest-Neighbors algorithm
- Given an unknown subject to classify,
- Look up all the known examples,
- Find the K closest examples,
- Take a majority vote,
- Predict what the majority says

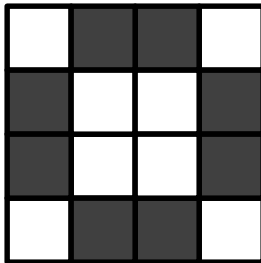
Illustration: 1 nearest neighbor

Sample

1

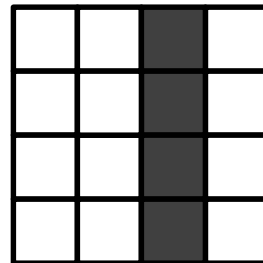


0



Unknown

?



*Which item from the sample
is nearest / closest to the Unknown
item we want to predict?*

What does “close” mean?

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_i - q_i)^2 + \dots + (p_n - q_n)^2}.$$

- To define “close” we need a distance
- We can use the distance between images as a measure for “close”
- Other distances can be used as well
- Note: Square root not important here

Illustration: 1 nearest neighbor

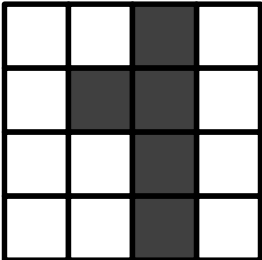
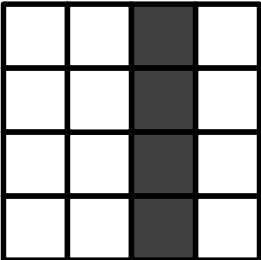
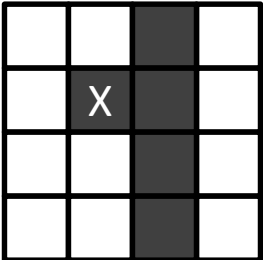
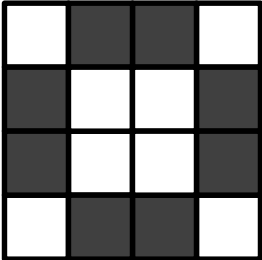
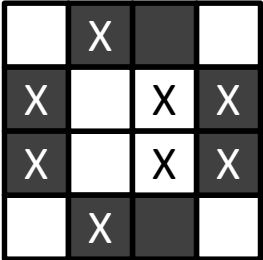
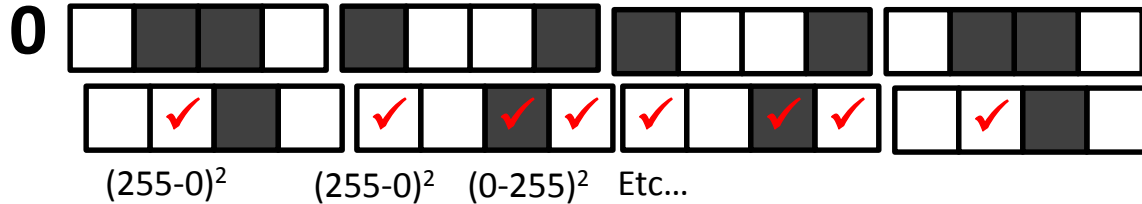
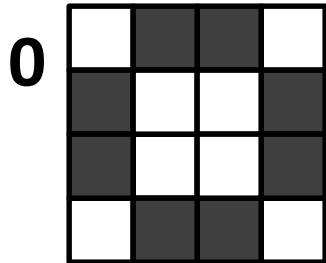
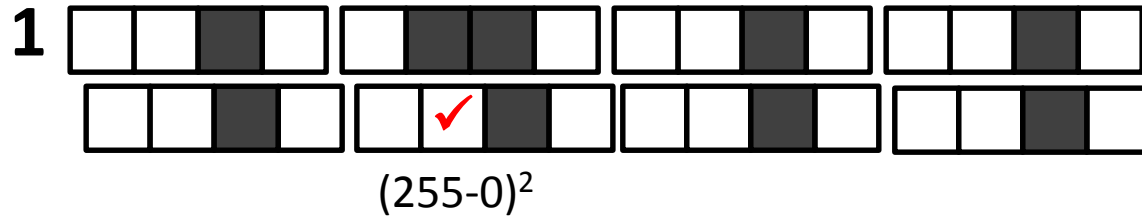
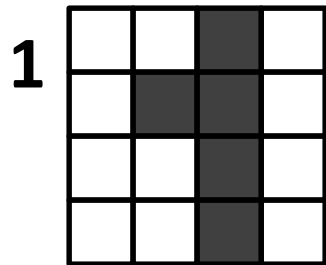
Sample	Unknown	Differences	Distances
1 	? 	1 	$\sqrt{255^2}$
0 		0 	$\sqrt{(255^2+255^2+...+255^2)}$

Illustration: 1 nearest neighbor

Sample



Unknown

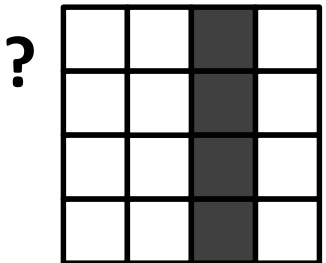


Illustration: 1 nearest neighbor

Sample	Unknown	Differences	Distances
1	?	1	This item is closest: we predict its number, 1
0		0	

Questions?

Organization

- Form teams
- 1:00 – 2:45: code
- 2:45 – 3:00: prepare demo
- 3:00 – 4:00: demos (5 minutes each)
- Slides & “guidance” are on **github.com/strangeloop/lambdajam2013**

Let's start coding!

- Suggested path
 - Use Euclidean distance first
 - Build a 1-neighbor classifier
 - What % of examples in Validation are correctly classified?
 - ... go wild 😊