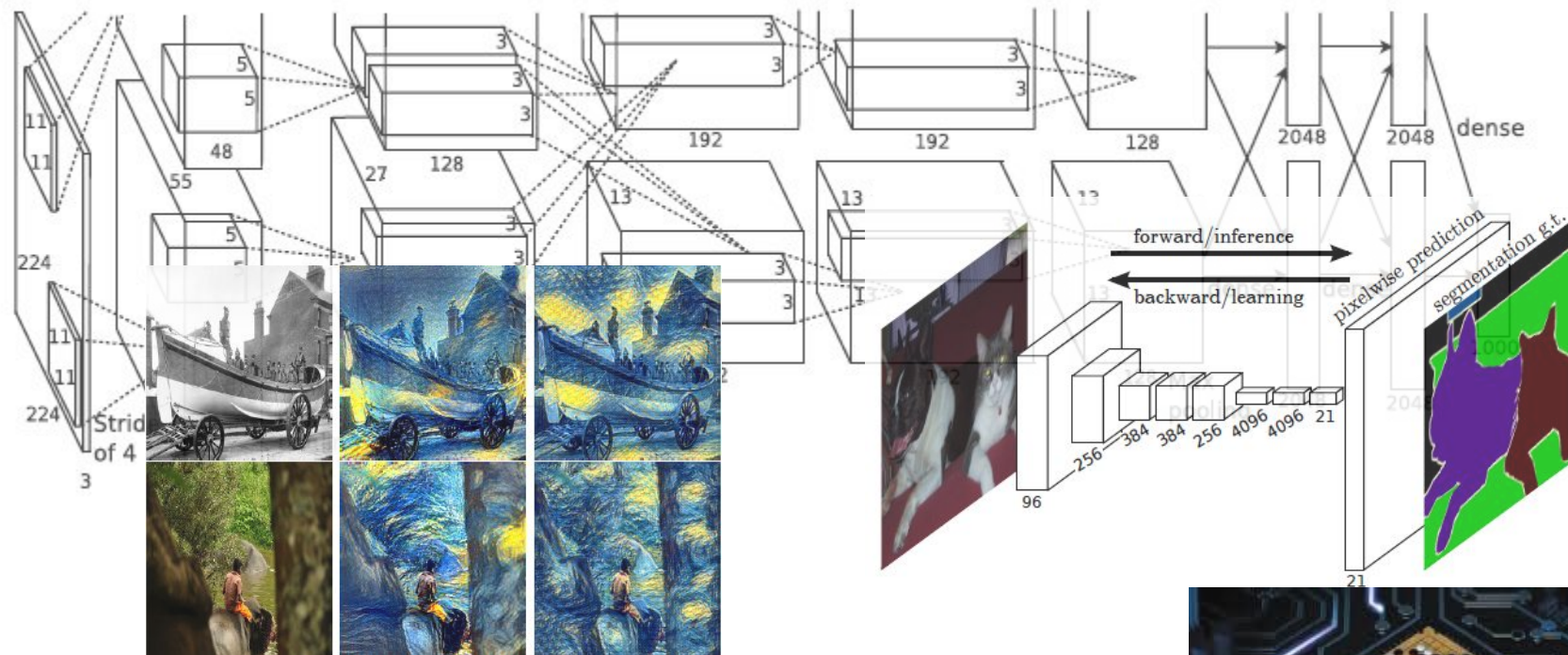


Lecture:

Introduction to Learning Based Vision

The modern world of machine learning




Figures from Krizhevsky et al., Shelhamer et al, Johnson et al, van den Oord et al, Silver et al.

Acquisitions

Google snaps up object recognition startup DNNresearch

Google has acquired a research startup founded within the University of Toronto, whose work includes object recognition.

by **Josh Lowensohn**  **@Josh** / 13 March 2013, 9:22 am AEDT

 2 /  0 /  0 /  0 /  0 /  more +

Google has acquired a three-person Canadian research company that specializes in voice and image recognition.

DNNresearch, which was founded last year within the the University of Toronto's computer science department, specializes in object recognition and now belongs to Google.



From left: Ilya Sutskever, Alex Krizhevsky and University Professor Geoffrey Hinton of the University of Toronto's Department of Computer Science. *(photo by John Guatto, University of Toronto)*

Acquisitions

Google snaps up object recognition startup DNNresearch

Google has acquired a research startup founded within the University of Toronto, whose work includes object recognition.

by Josh Lowensohn @Josh / 13 March 2013, 9:22 am AEDT

2 / f 0 / t 0 / in 0 / g+ / ... more +

Google has acquired a three-person Canadian research company that specializes in voice and image recognition.

DNNresearch, which was founded last year within the the University of Toronto's computer science department, specializes in object recognition and now belongs to Google.



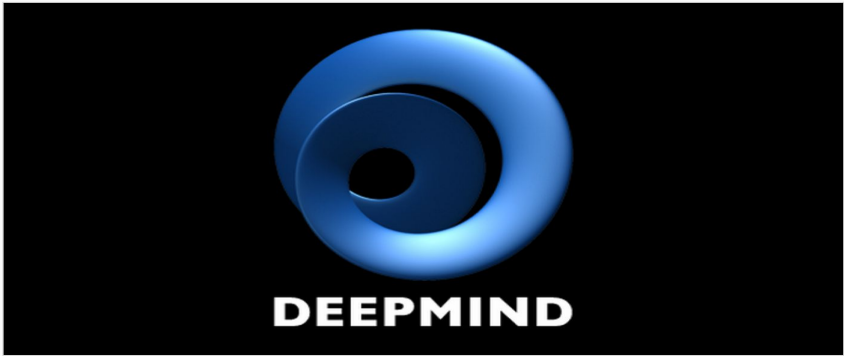
8TH ANNUAL CRUNCHIES AWARDS Celebrate the Best of Tech in 2014 Get Your Tickets Now ▶

Google Acquires Artificial Intelligence Startup DeepMind For More Than \$500M

Posted Jan 26, 2014 by Catherine Shu (@catherineshu)

11.3k SHARES [social sharing icons]

Next Story ▶




CrunchBase

DeepMind	
FOUNDED	2011
OVERVIEW	
DeepMind is a cutting edge artificial intelligence company. We combine the best techniques from machine learning and systems neuroscience to build powerful general-purpose learning algorithms. Founded by Demis Hassabis, Shane Legg and Mustafa Suleyman, the company is based in London and supported by some of the most iconic technology entrepreneurs and investors of the past decade. Our first commercial ...	

Acquisitions

Google snaps up object recognition startup DNNresearch

Google has acquired a research startup founded within the University of Toronto, whose work includes object recognition.

by **Josh Lowensohn**  **@Josh** / 13 March 2013, 9:22 am AEDT

 2 /  0 /  0 /  0 /  0 /  more +

Google has acquired a three-person Canadian research company that specializes in voice and image recognition.

DNNresearch, which was founded last year within



8TH ANNUAL CRUNCHIES AWARDS Celebrate the Best of Tech in 2014 [Get Your Tickets Now](#) ▶

Google Acquires Artificial Intelligence Startup DeepMind



Yann LeCun

December 9, 2013 · 

Big news today!

Facebook has created a new research laboratory with the ambitious, long-term goal of bringing about major advances in Artificial Intelligence.

Slide by Dhruv Batra

Acquisitions

Google snaps up object recognition startup

DNNr

Google has ac
Toronto, who

by Josh Lowensohn

2 / f o

Google has acqui
research compan
image recognition

DNNresearch. wh



Big ne

Facebook h
long-term go
Intelligence.

Slide by Dhruv Batia

« Search needs a shake-up

Songbirds use grammar rules »

Machine Learning Startup Acquired by ai-one

Press Release

For Immediate Release: August 4, 2011

San Diego artificial intelligence startup acquired by leading

IBM acquires deep learning startup AlchemyAPI

by Derrick Harris Mar. 4, 2015 - 8:15 AM PDT

1 Comment

San
prov
late
lead

up
the

IBM Watson. Photo by Clockready/Wikimedia Commons

What is ~~Machine~~ Learning?

- “the acquisition of knowledge or skills through experience, study, or by being taught.”
- Can be (almost) mapped to reinforcement, unsupervised and supervised machine learning.

What is Machine Learning?

- [Arthur Samuel, 1959]
 - Field of study that gives computers
 - the ability to learn without being explicitly programmed
- [Kevin Murphy] algorithms that
 - automatically detect patterns in data
 - use the uncovered patterns to predict future data or other outcomes of interest
- [Tom Mitchell] algorithms that
 - improve their performance (P)
 - at some task (T)
 - with experience (E)

What is Machine Learning?



ML in Nutshell

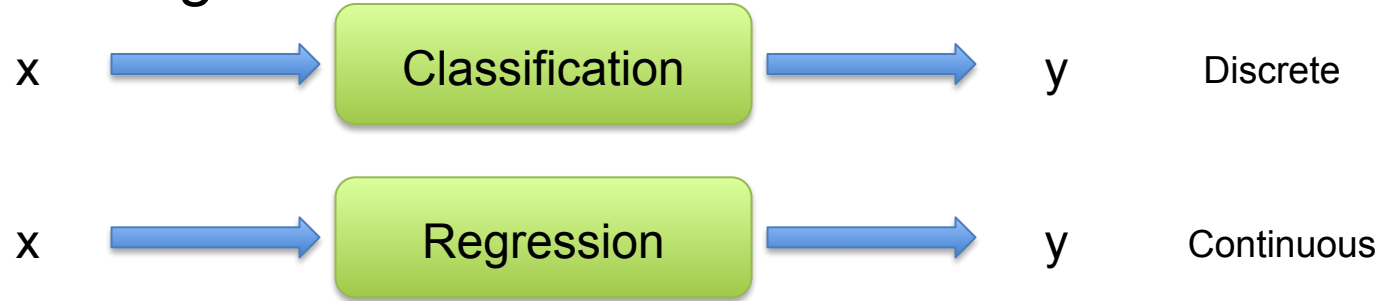
- Tens of thousands of machine learning algorithms
 - Hundreds new every year
- Decades of ML research oversimplified:
 - All of Machine Learning:
 - Learn a mapping from input to output $f: X \rightarrow Y$
 - e.g. X : emails, Y : {spam, notspam}

Types of Learning

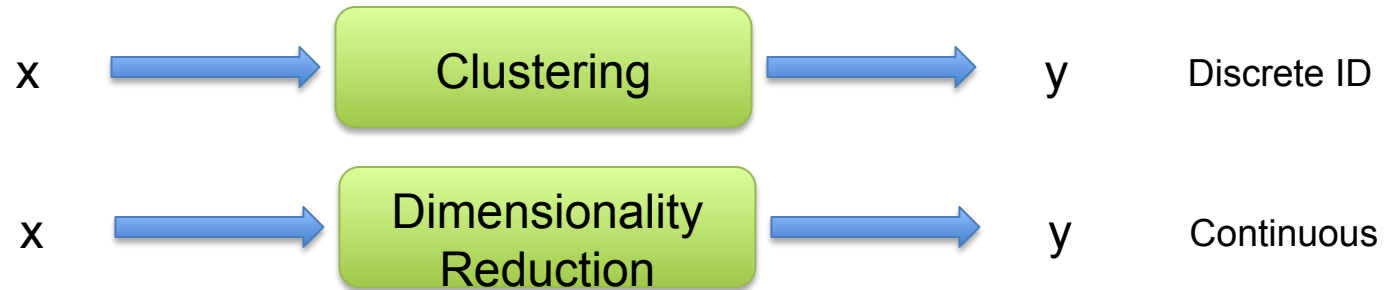
- Supervised learning
 - Training data includes desired outputs
- Unsupervised learning
 - Training data does not include desired outputs
- Weakly or Semi-supervised learning
 - Training data includes a few desired outputs
- Reinforcement learning
 - Rewards from sequence of actions

Tasks

Supervised Learning



Unsupervised Learning



Examples for supervised learning

Vision: Image Classification

- <http://cloudcv.org/classify/>

x



y

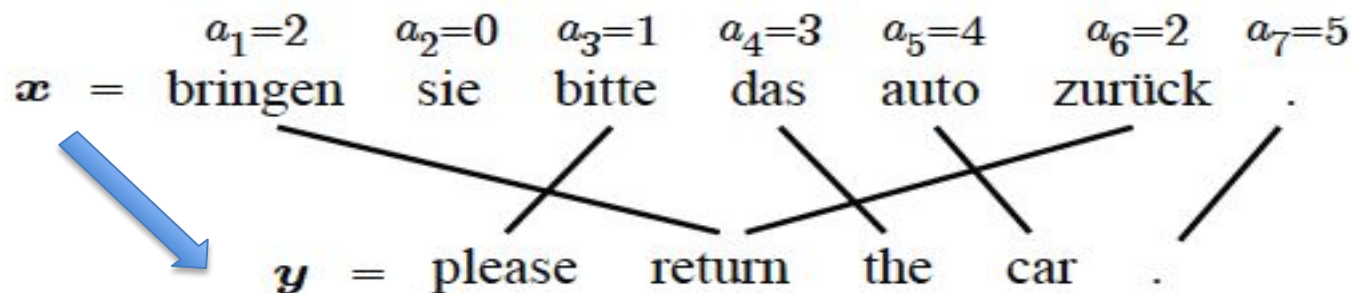


scuba diver

tiger shark

hammerhead
shark

NLP: Machine Translation



Speech: Speech2Text



Slide by Dhruv Batra

Slide Credit: Carlos Guestrin

Image captioning



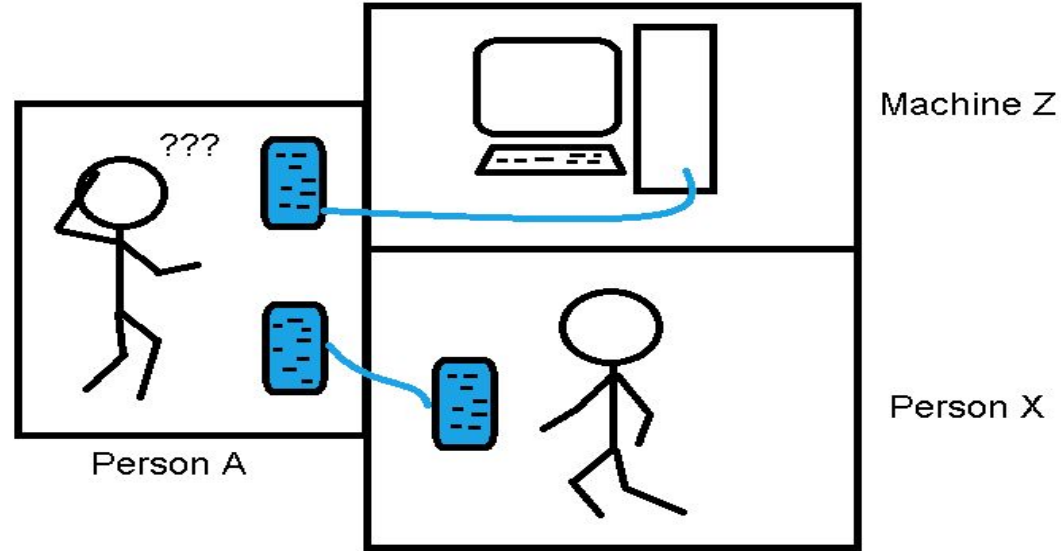
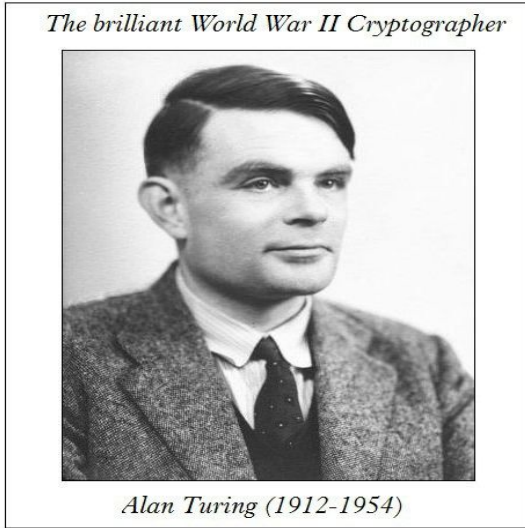
"woman is holding bunch of
bananas."



"black cat is sitting on top of
suitcase."

AI: Turing Test

“Can machines think”



Q: Please write me a sonnet on the subject of the Forth Bridge.

A: Count me out on this one. I never could write poetry.

Q: Add 34957 to 70764.

A: (Pause about 30 seconds and then give as answer) 105621.

AI: Visual Turing Test

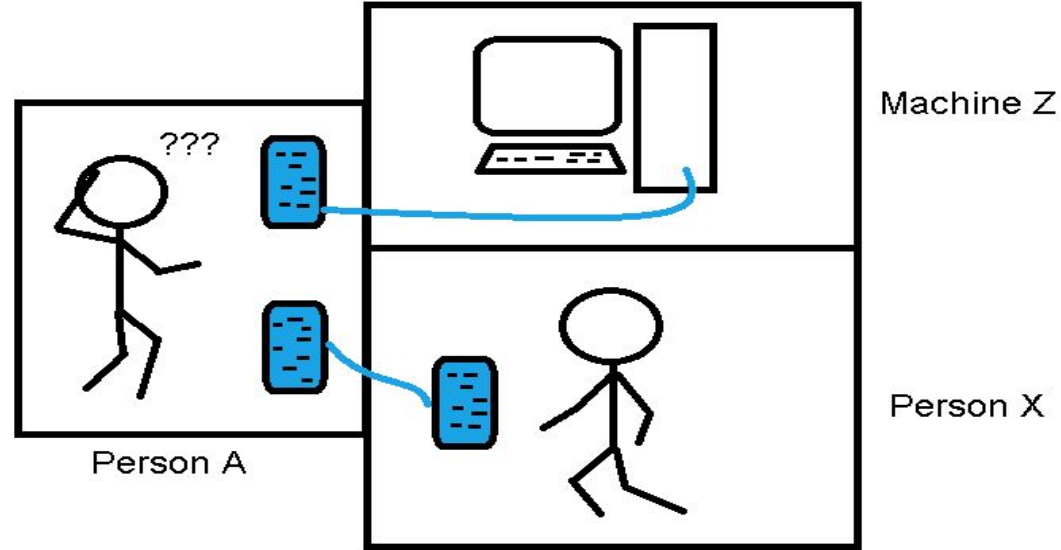


Q: How many slices
of pizza are there?



A: 6

Slide by Dhruv Batra



AI: **Visual** Turing Test

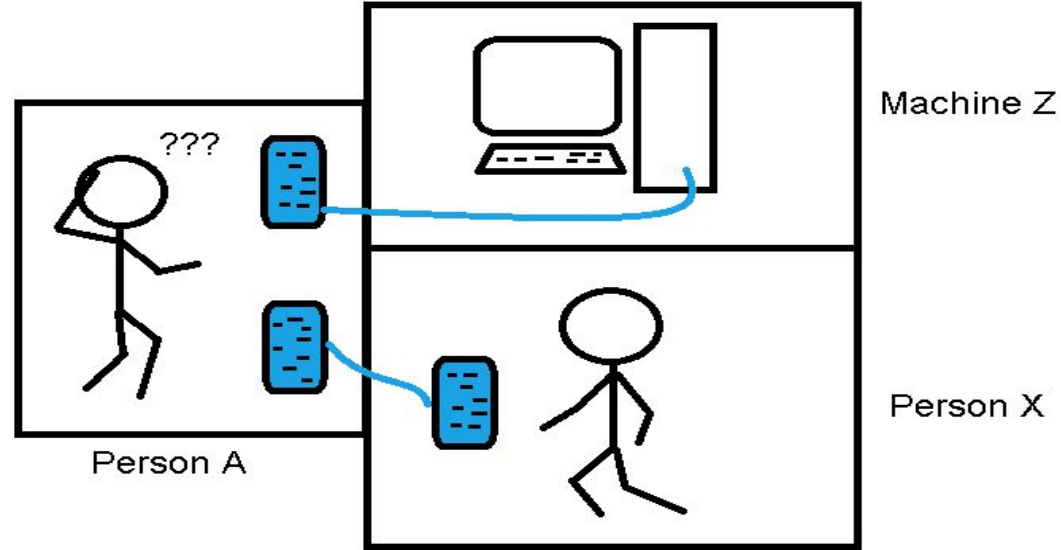


Q: How many slices
of pizza are there?

x
↓
y

A: 6

Slide by Dhruv Batra



Demo: <http://cloudcv.org/vqa/>

Supervised Learning

- Input: x (images, text, emails...)
- Output: y (spam or non-spam...)

Supervised Learning

- Input: x (images, text, emails...)
- Output: y (spam or non-spam...)
- (Unknown) Target Function
 - $f: X \rightarrow Y$ (the “true” mapping / reality)

Supervised Learning

- Input: x (images, text, emails...)
- Output: y (spam or non-spam...)
- (Unknown) Target Function
 - $f: X \rightarrow Y$ (the “true” mapping / reality)
- Data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Supervised Learning

- Input: x (images, text, emails...)
- Output: y (spam or non-spam...)
- (Unknown) Target Function
 - $f: X \rightarrow Y$ (the “true” mapping / reality)
- Data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
- Model / Hypothesis Class
 - $g: X \rightarrow Y$
 - $\hat{y} = g(x) = \text{sign}(w^T x)$
- Learning = Search in hypothesis space
 - Find best g in model class.

Supervised Learning

- Input: x (images, text, emails...)
- Output: y (spam or non-spam...)
- (Unknown) Target Function
 - $f: X \rightarrow Y$ (the “true” mapping / reality)
- Data: $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
- Model / Hypothesis Class
 - $g: X \rightarrow Y$
 - $\hat{y} = g(x) = \text{sign}(w^T x)$
- Learning = Search in hypothesis space
 - Find best g in model class.

Synonyms

- Representation Learning
- Deep (Machine) Learning
- Deep Neural Networks
- Deep Unsupervised Learning
- Simply: Deep Learning

So what *is* Deep (Machine) Learning?

- A few different ideas:
- (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
- End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extraction
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons work together

Incoming slides and lectures

- Some fundamentals of machine learning
 - A traditional image classification pipeline (Use raw pixels as features and nearest neighbor classifier)
 - Linear classification
 - Loss functions
 - Training by numerical optimization
- Then, deep learning will be a natural extension..

A simple image classification pipeline

Image Classification: a core task in Computer Vision



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

→ cat

The problem: *semantic gap*

Images are represented as
3D arrays of numbers, with
integers between $[0, 255]$.

E.g.

$300 \times 100 \times 3$

(3 for 3 color channels RGB)

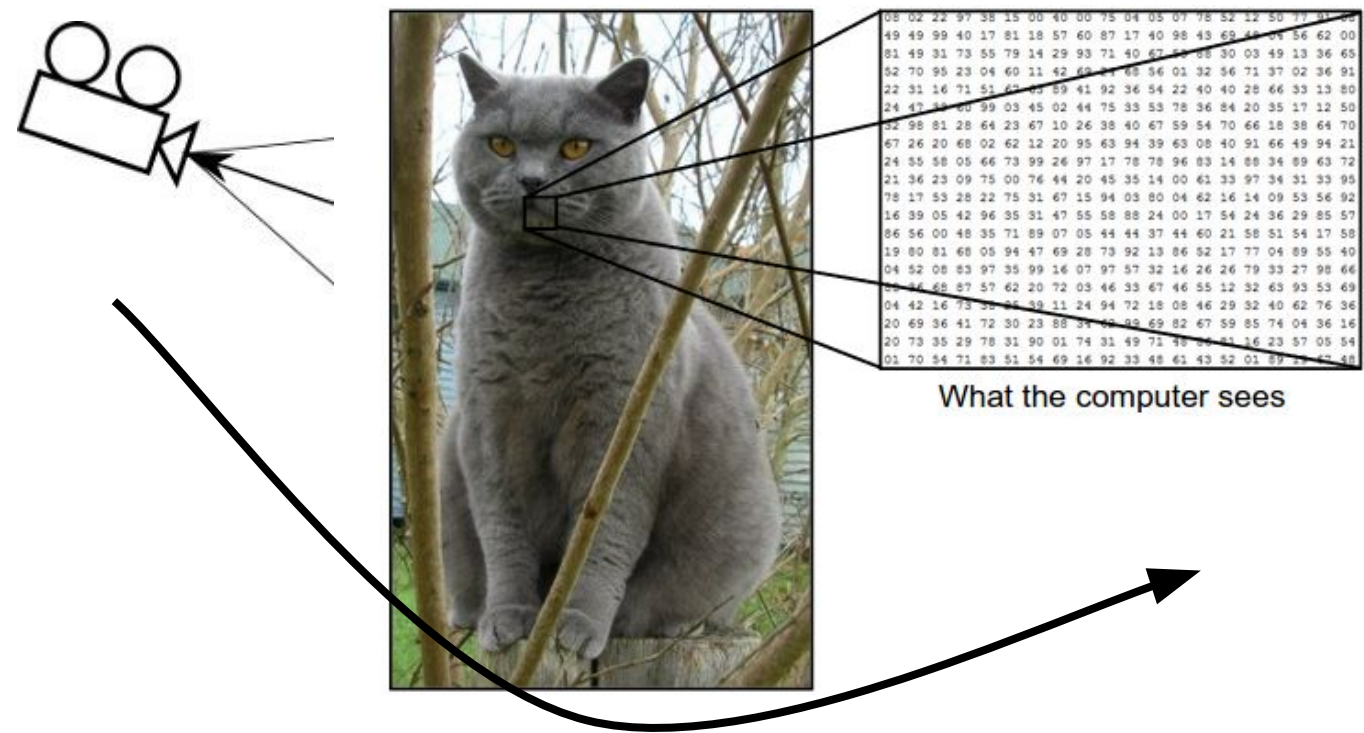


08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	01	29
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	53	88	30	03	49	13	36	65
32	70	95	23	04	60	11	42	69	14	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	03	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	33	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
25	46	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	58	39	11	24	94	72	18	08	46	29	32	40	62	76	36	
20	69	36	41	72	30	23	88	34	62	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	38	85	81	16	23	57	05	54
03	70	54	71	83	51	54	69	16	92	33	48	61	43	82	01	89	11	67	48

What the computer sees

82% cat

Challenges: Viewpoint Variation



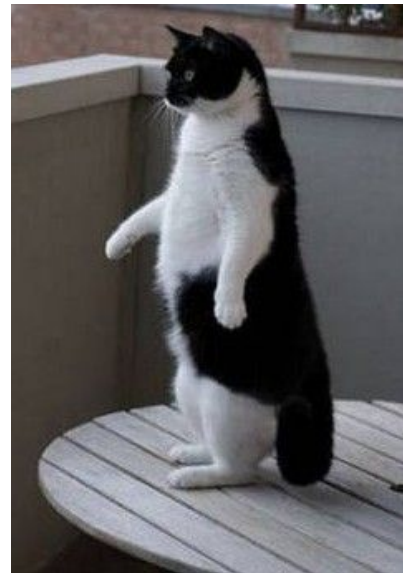
Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Challenges: Illumination



Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Challenges: Deformation



Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Challenges: Occlusion



Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Challenges: Background clutter



Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Challenges: Intraclass variation



Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

An image classifier

```
def predict(image):  
    # ????  
    return class_label
```

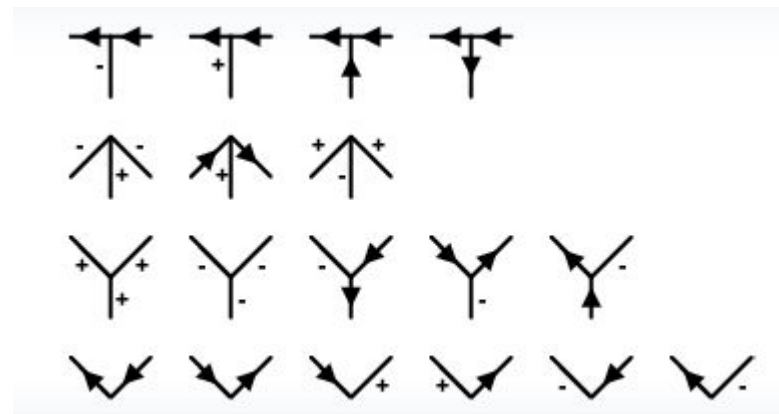
Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.

Attempts have been made



???

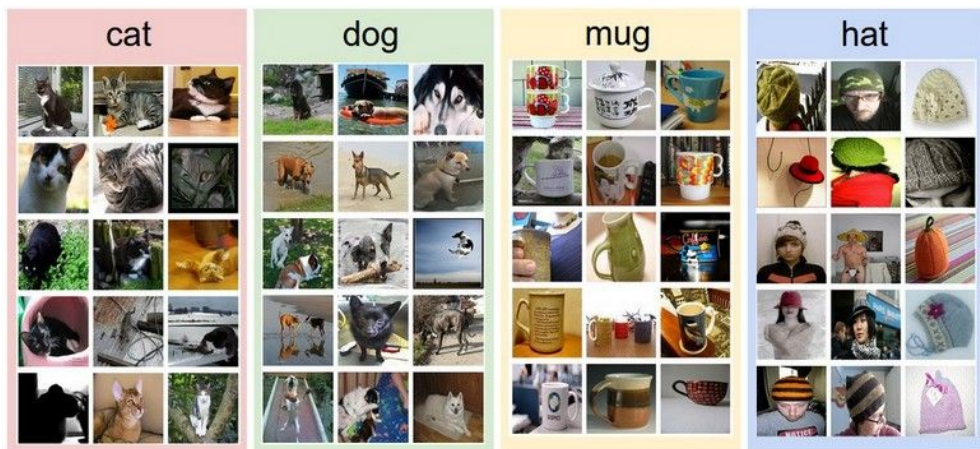


Data-driven approach:

1. Collect a dataset of images and labels
2. Use Machine Learning to train an image classifier
3. Evaluate the classifier on a withheld set of test images

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Example training set



First classifier: **Nearest Neighbor Classifier**

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Two arrows originate from the right side of the code block. The top arrow points from the `train` function to the text 'Remember all training images and their labels'. The bottom arrow points from the `predict` function to the text 'Predict the label of the most similar training image'.

Remember all training images and their labels

Predict the label of the most similar training image

Example dataset: **CIFAR-10**

10 labels

50,000 training images, each image is tiny: 32x32

10,000 test images.

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



Example dataset: **CIFAR-10**

10 labels

50,000 training images

10,000 test images.

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



For every test image (first column),
examples of nearest neighbors in rows



How do we compare the images? What is the **distance metric**?

L1 distance:
$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

test image					training image					pixel-wise absolute value differences				
56	32	10	18		10	20	24	17		46	12	14	1	
90	23	128	133		8	10	89	100		82	13	39	33	
24	26	178	200	-	12	16	178	170	=	12	10	0	30	add
2	0	255	220		4	32	233	112		2	32	22	108	456

Nearest Neighbor classifier

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

```
import numpy as np
```

```
class NearestNeighbor:
```

```
    def __init__(self):  
        pass
```

```
    def train(self, X, y):
```

```
        """ X is N x D where each row is an example. Y is 1-dimension of size N """  
        # the nearest neighbor classifier simply remembers all the training data  
        self.Xtr = X  
        self.ytr = y
```

```
    def predict(self, X):
```

```
        """ X is N x D where each row is an example we wish to predict label for """  
        num_test = X.shape[0]  
        # lets make sure that the output type matches the input type  
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)  
  
        # loop over all test rows  
        for i in xrange(num_test):  
            # find the nearest training image to the i'th test image  
            # using the L1 distance (sum of absolute value differences)  
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)  
            min_index = np.argmin(distances) # get the index with smallest distance  
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

```
    return Ypred
```

remember the training data

Nearest Neighbor classifier

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

for every test image:

- find nearest train image with L1 distance
- predict the label of nearest training image

Nearest Neighbor classifier

Q: how does the classification speed depend on the size of the training data?

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred

```

Nearest Neighbor classifier

Q: how does the classification speed depend on the size of the training data?
linearly :(

This is **backwards**:

- test time performance is usually much more important in practice.
- CNNs flip this: expensive training, cheap test evaluation

Aside: Approximate Nearest Neighbor

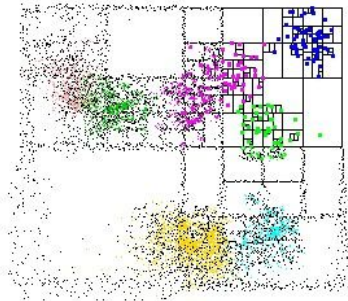
find approximate nearest neighbors quickly

ANN: A Library for Approximate Nearest Neighbor Searching

David M. Mount and Sunil Arya

Version 1.1.2

Release Date: Jan 27, 2010



What is ANN?

ANN is a library written in C++, which supports data structures and algorithms for both exact and approximate nearest neighbor searching in arbitrarily high dimensions.

In the nearest neighbor problem a set of data points in d -dimensional space is given. These points are preprocessed into a data structure, so that given any query point q , the nearest or generally k nearest points of P to q can be reported efficiently. The distance between two points can be defined in many ways. ANN assumes that distances are measured using any class of distance functions called Minkowski metrics. These include the well known Euclidean distance, Manhattan distance, and max distance.

Based on our own experience, ANN performs quite efficiently for point sets ranging in size from thousands to hundreds of thousands, and in dimensions as high as 20. (For applications in significantly higher dimensions, the results are rather spotty, but you might try it anyway.)

The library implements a number of different data structures, based on kd-trees and box-decomposition trees, and employs a couple of different search strategies.

The library also comes with test programs for measuring the quality of performance of ANN on any particular data sets, as well as programs for visualizing the structure of the geometric data structures.

FLANN - Fast Library for Approximate Nearest Neighbors

- Home
- News
- Publications
- Download
- Changelog
- Repository

What is FLANN?

FLANN is a library for performing fast approximate nearest neighbor searches in high dimensional spaces. It contains a collection of algorithms we found to work best for nearest neighbor search and a system for automatically choosing the best algorithm and optimum parameters depending on the dataset.

FLANN is written in C++ and contains bindings for the following languages: C, MATLAB and Python.

News

- (14 December 2012) Version 1.8.0 is out bringing incremental addition/removal of points to/from indexes
- (20 December 2011) Version 1.7.0 is out bringing two new index types and several other improvements.
- You can find binary installers for FLANN on the [Point Cloud Library](#) project page. Thanks to the PCL developers!
- Mac OS X users can install flann through MacPorts (thanks to Mark Moll for maintaining the Portfile)
- New release introducing an easier way to use custom distances, kd-tree implementation optimized for low dimensionality search and experimental MPI support
- New release introducing new C++ templated API, thread-safe search, save/load of indexes and more.
- The FLANN license was changed from LGPL to BSD.

How fast is it?

In our experiments we have found FLANN to be about one order of magnitude faster on many datasets (in query time), than previously available approximate nearest neighbor search software.

Publications

More information and experimental results can be found in the following papers:

- Marius Muja and David G. Lowe: "Scalable Nearest Neighbor Algorithms for High Dimensional Data", Pattern Analysis and Machine Intelligence (PAMI), Vol. 36, 2014. [\[PDF\]](#) [\[BibTeX\]](#)
- Marius Muja and David G. Lowe: "Fast Matching of Binary Features", Conference on Computer and Robot Vision (CRV) 2012. [\[PDF\]](#) [\[BibTeX\]](#)
- Marius Muja and David G. Lowe: "Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration", In International Conference on Computer Vision Theory and Applications (VISAPP'09), 2009 [\[PDF\]](#) [\[BibTeX\]](#)

The choice of distance is a **hyperparameter**
common choices:

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

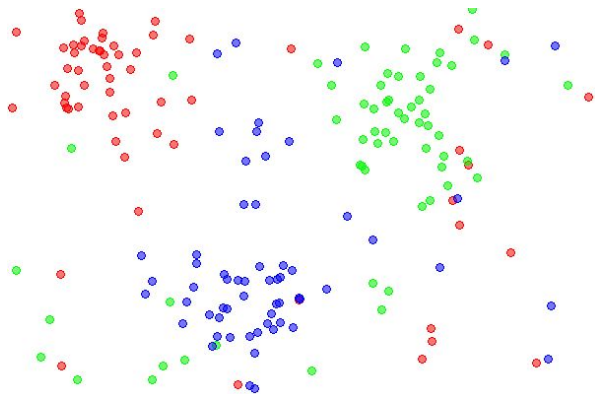
L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

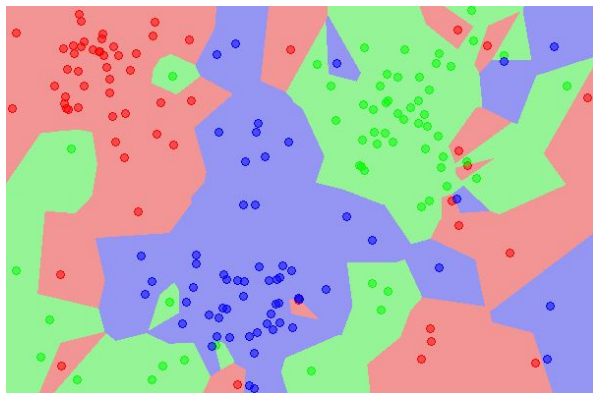
k-Nearest Neighbor

find the k nearest images, have them vote on the label

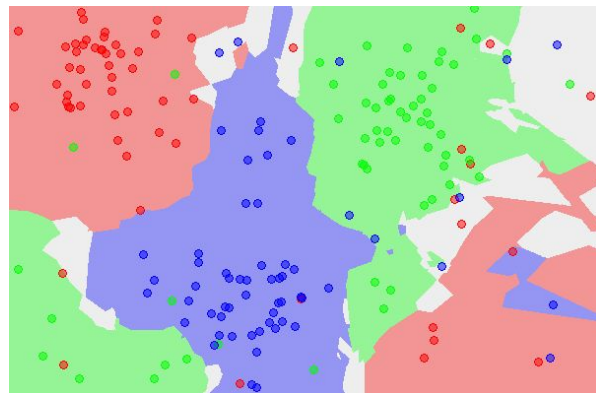
the data



NN classifier



5-NN classifier



http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

Slide by Fei-Fei Li, Andrej Karpathy & Justin Johnson

Example dataset: **CIFAR-10**

10 labels

50,000 training images

10,000 test images.

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



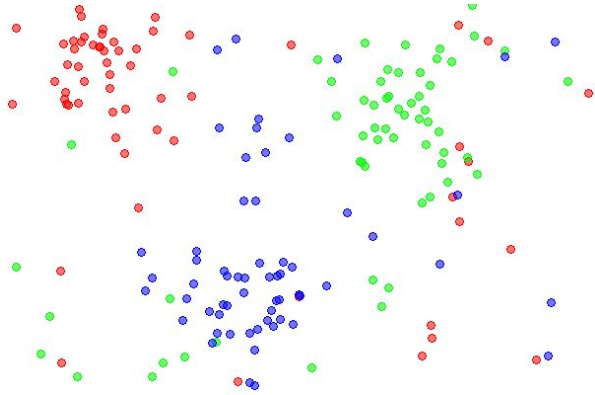
truck



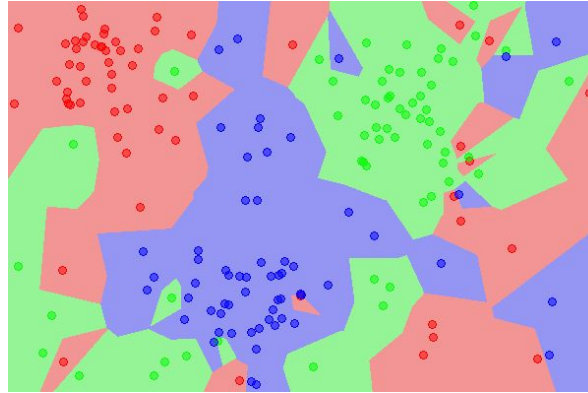
For every test image (first column),
examples of nearest neighbors in rows



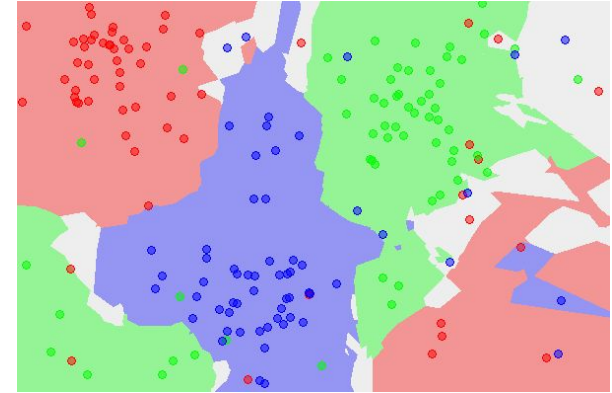
the data



NN classifier

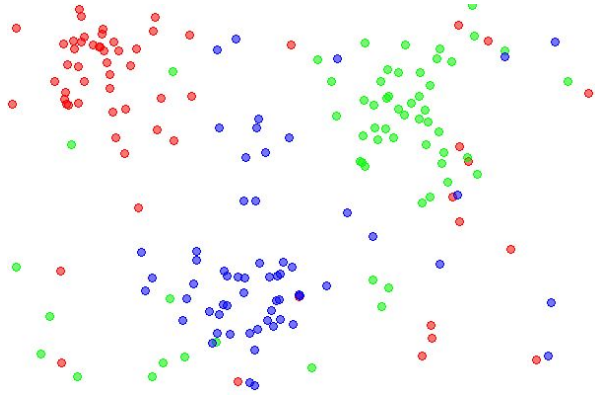


5-NN classifier

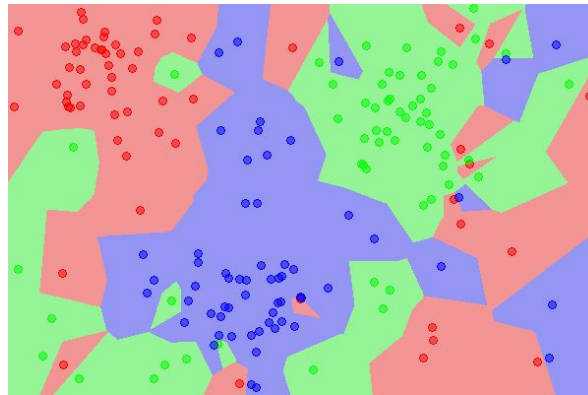


Q: what is the accuracy of the nearest neighbor classifier on the training data, when using the Euclidean distance?

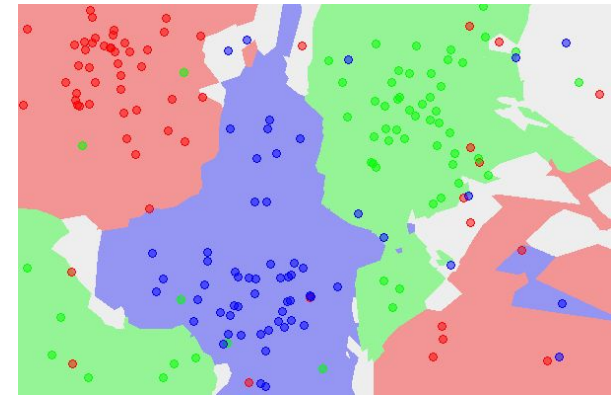
the data



NN classifier



5-NN classifier



Q2: what is the accuracy of the **k**-nearest neighbor classifier on the training data?

What is the best **distance** to use?

What is the best value of **k** to use?

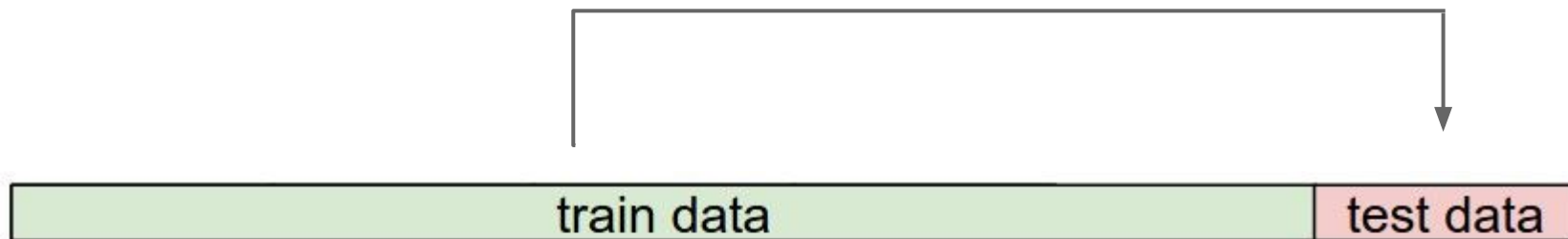
i.e. how do we set the **hyperparameters**?

What is the best **distance** to use?
What is the best value of **k** to use?

i.e. how do we set the **hyperparameters**?

Very problem-dependent.
Must try them all out and see what works best.

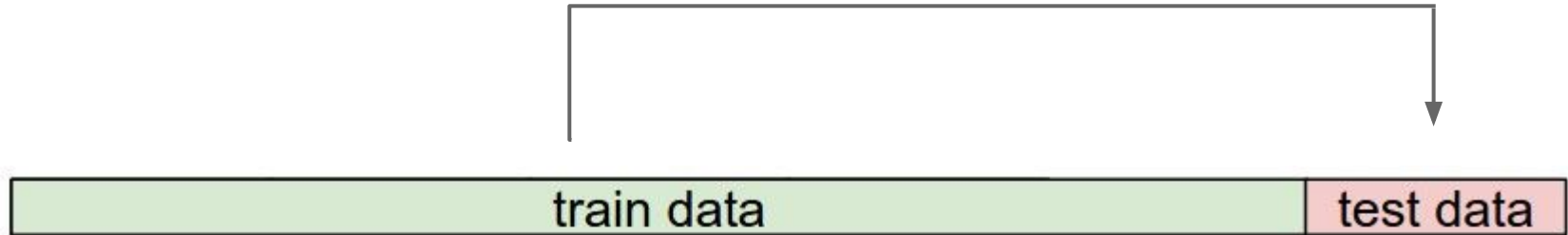
Try out what hyperparameters work best on test set.



Trying out what hyperparameters work best on test set:

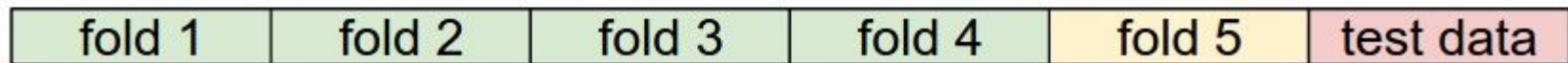
Very bad idea. The test set is a proxy for the generalization performance!

Use only **VERY SPARINGLY**, at the end.



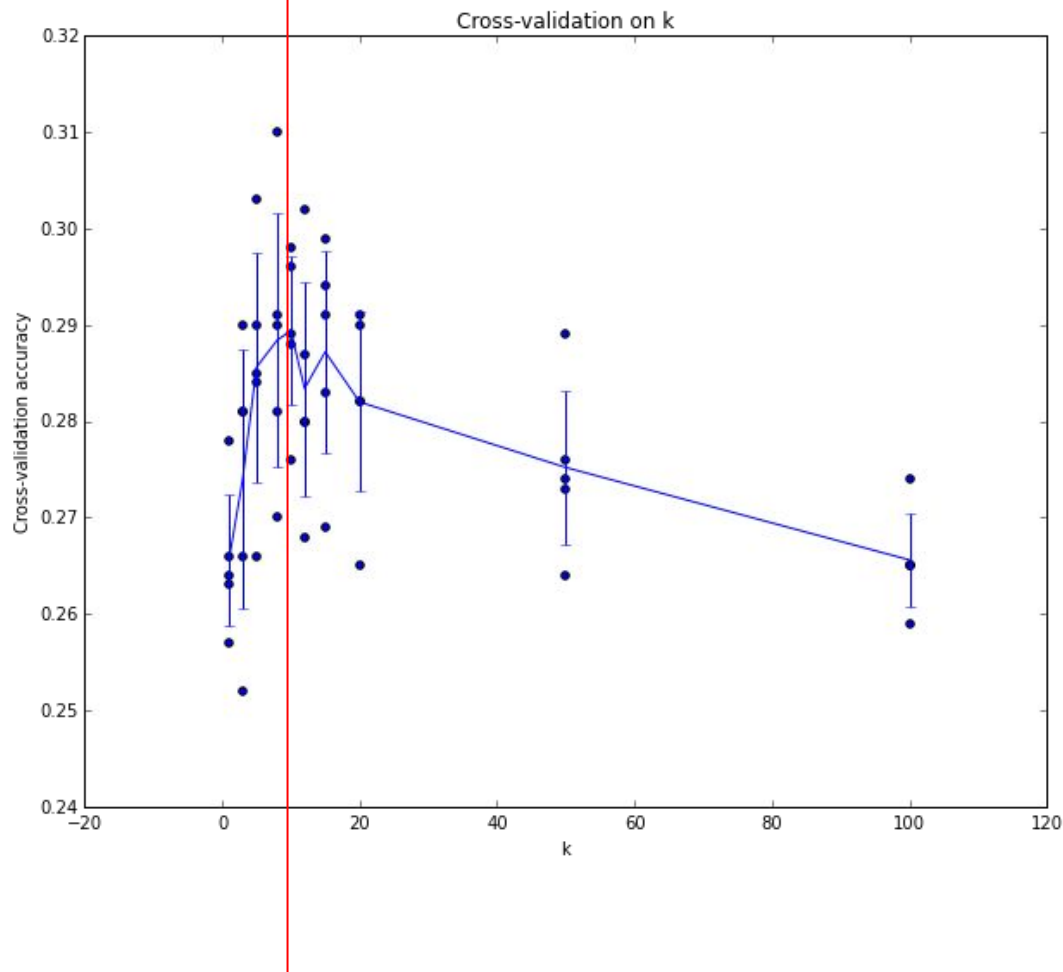


Validation data
use to tune hyperparameters



Cross-validation

cycle through the choice of which fold is the validation fold, average results.



Example of
5-fold cross-validation
for the value of **k**.

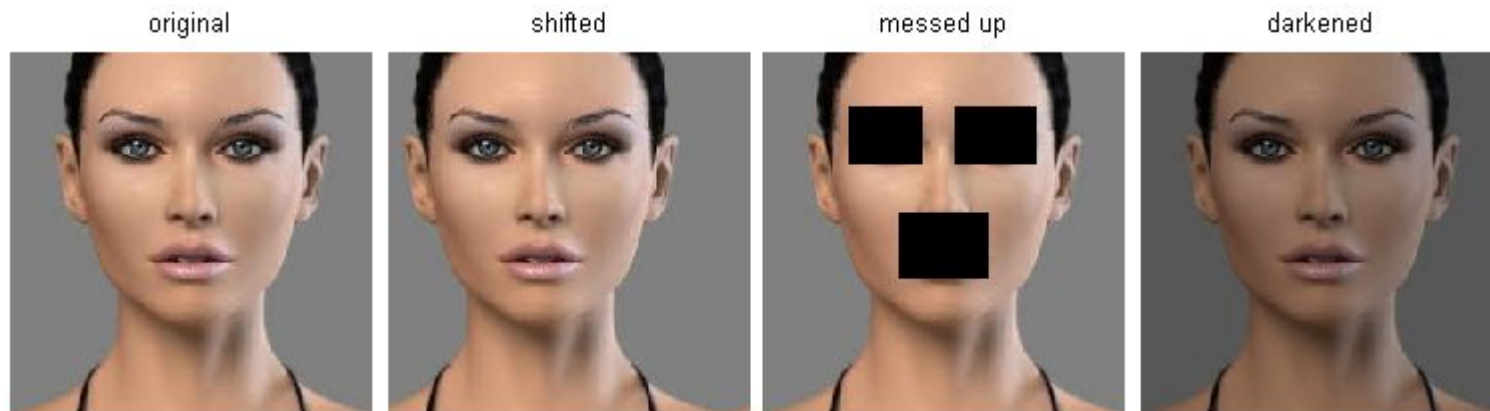
Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)

k-Nearest Neighbor on images is **never used**.

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



(all 3 images have same L2 distance to the one on the left)

Summary

- **Image Classification:** We are given a **Training Set** of labeled images, asked to predict labels on **Test Set**. Common to report the **Accuracy** of predictions (fraction of correctly predicted images)
- We introduced the **k-Nearest Neighbor Classifier**, which predicts the labels based on nearest images in the training set
- We saw that the choice of distance and the value of k are **hyperparameters** that are tuned using a **validation set**, or through **cross-validation** if the size of the data is small.
- Once the best set of hyperparameters is chosen, the classifier is evaluated once on the test set, and reported as the performance of kNN on that data.