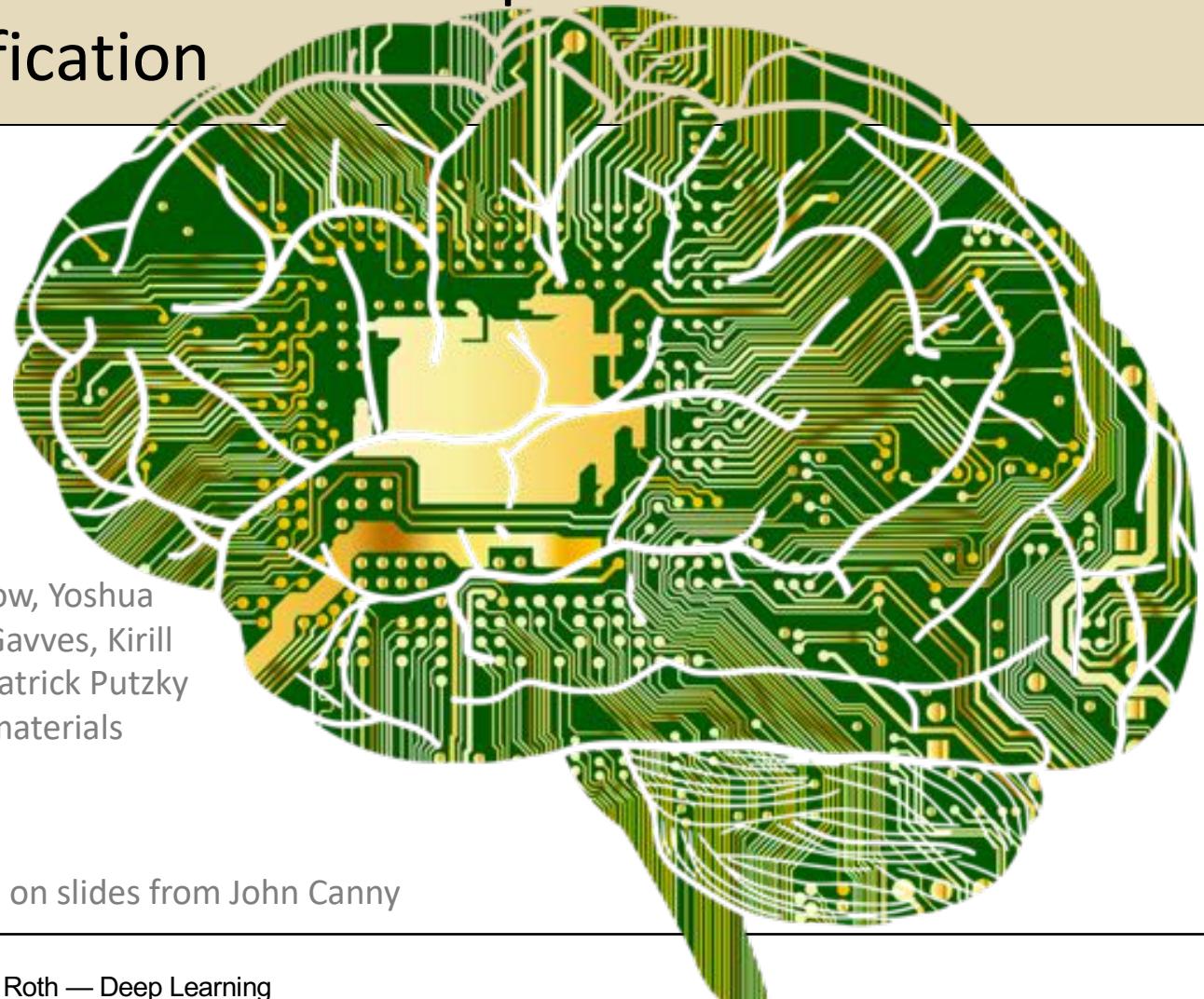


Deep Learning

Architectures and Methods: Computer Vision and Image Classification



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Thanks to John Canny, Ian Goodfellow, Yoshua Bengio, Aaron Courville, Efstratios Gavves, Kirill Gavrilyuk, Berkay Kicanaoglu, and Patrick Putzky and many others for making their materials publically available.

The present slides are mainly based on slides from John Canny

Computer Vision – a brief history



MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".



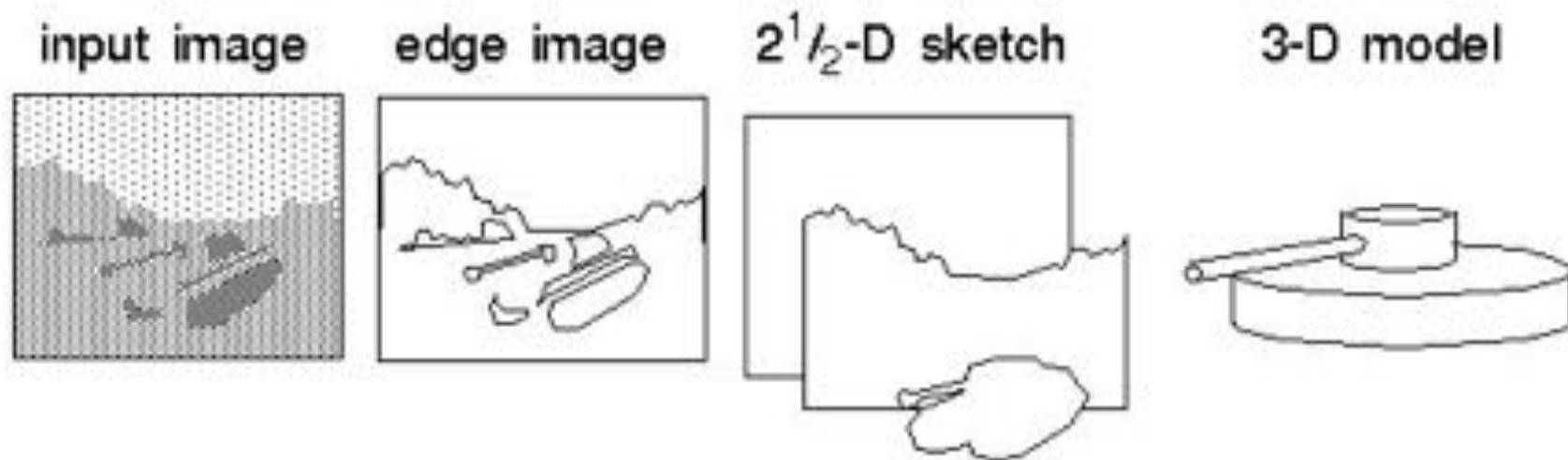
Computer Vision – a brief history

David Marr's approach 1970-80s

- Computation
- Algorithm
- Implementation

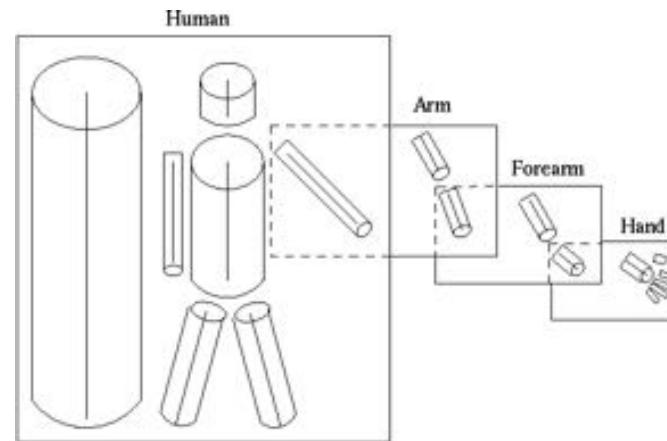


primal sketch and $2\frac{1}{2}$ D sketch



Computer Vision - Shape description

Skeletons and Cylinders: Marr and Nishihara 1978:

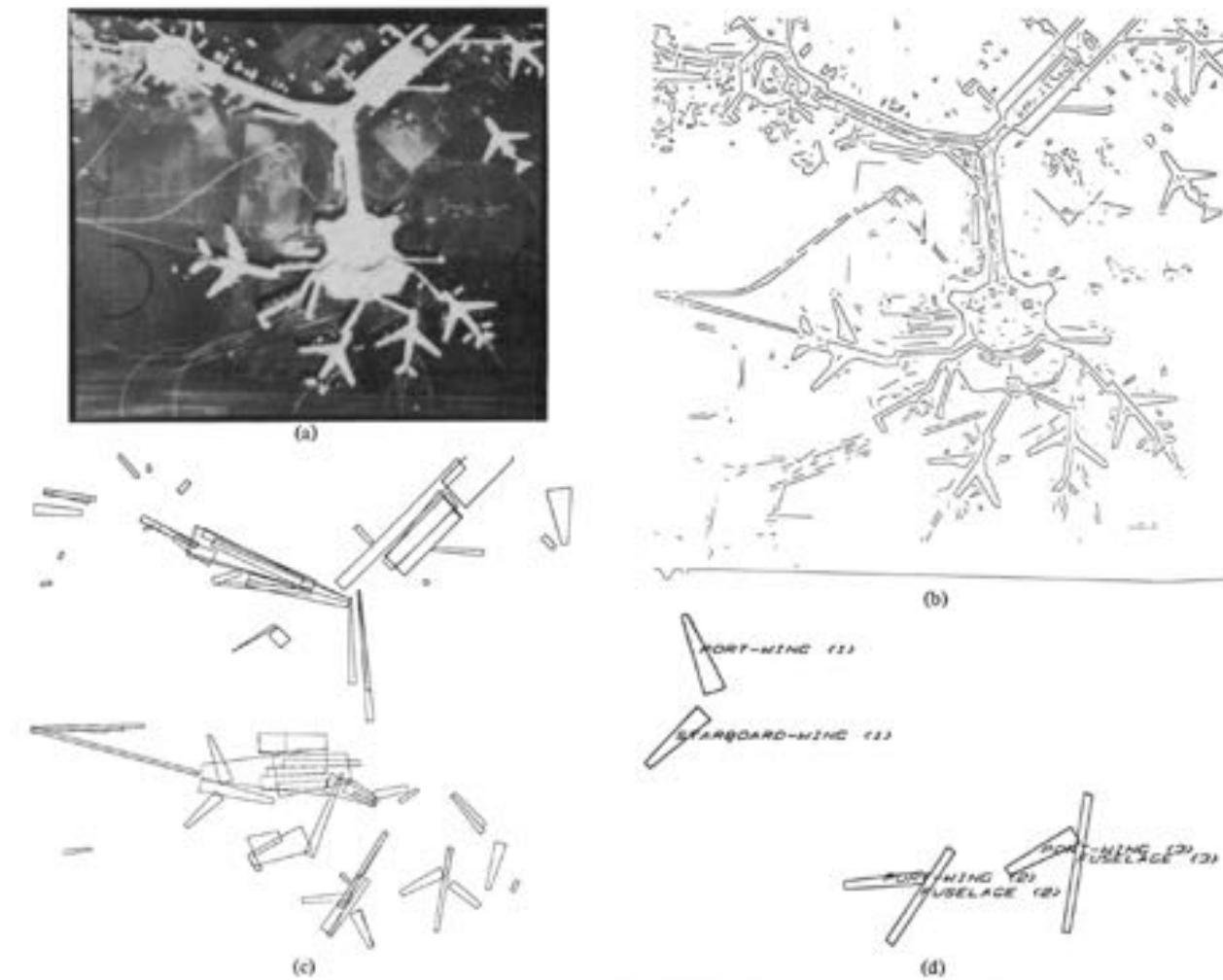


Tom Binford, generalized
cylinders 1971



Computer Vision - Shape matching

ACRONYM (Brooks 1982) matching with generalized cylinders + search!:



Computer Vision - Eigenfaces

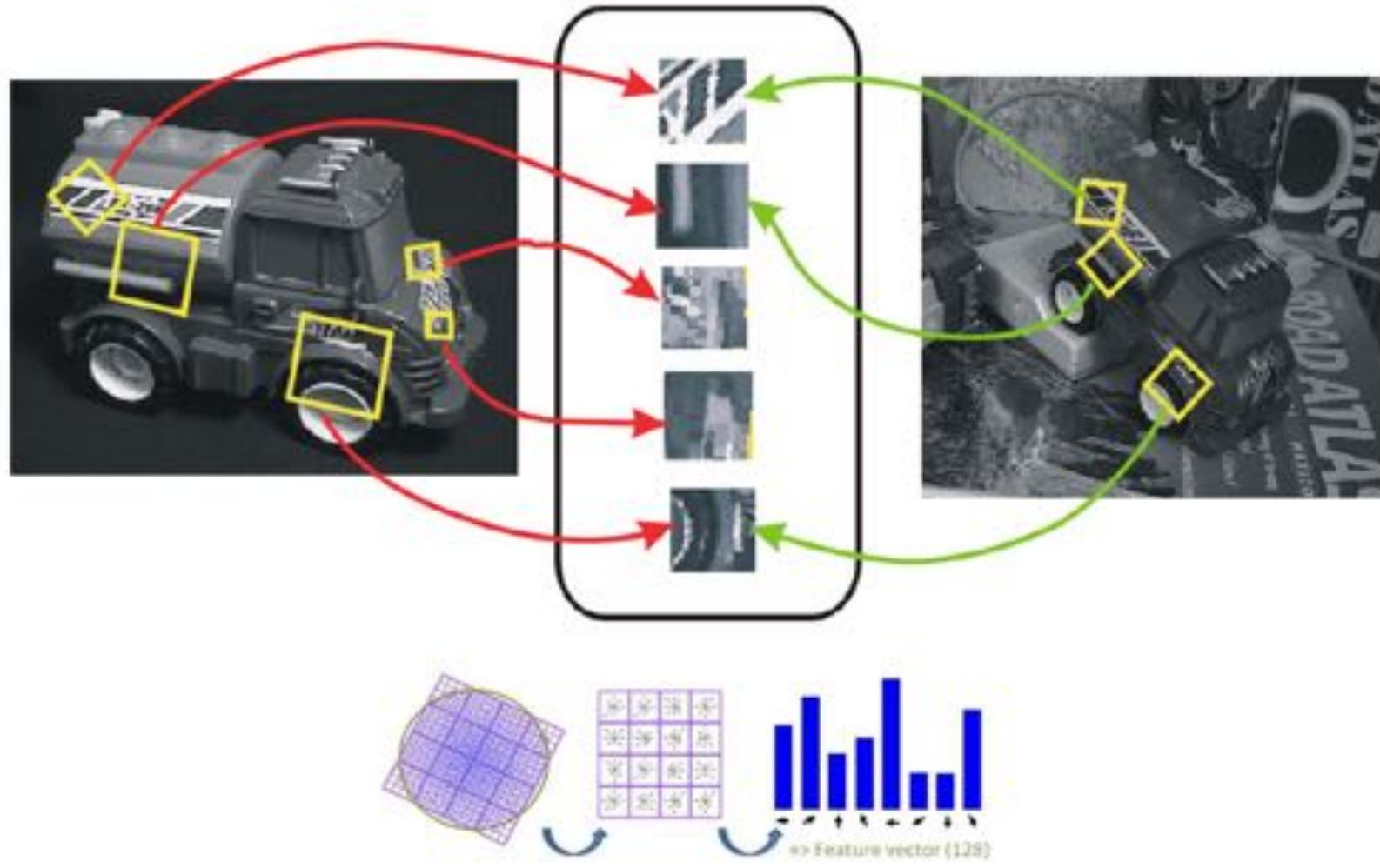
Turk and Pentland 1991:



A data-driven approach?



Computer Vision – Invariant Features

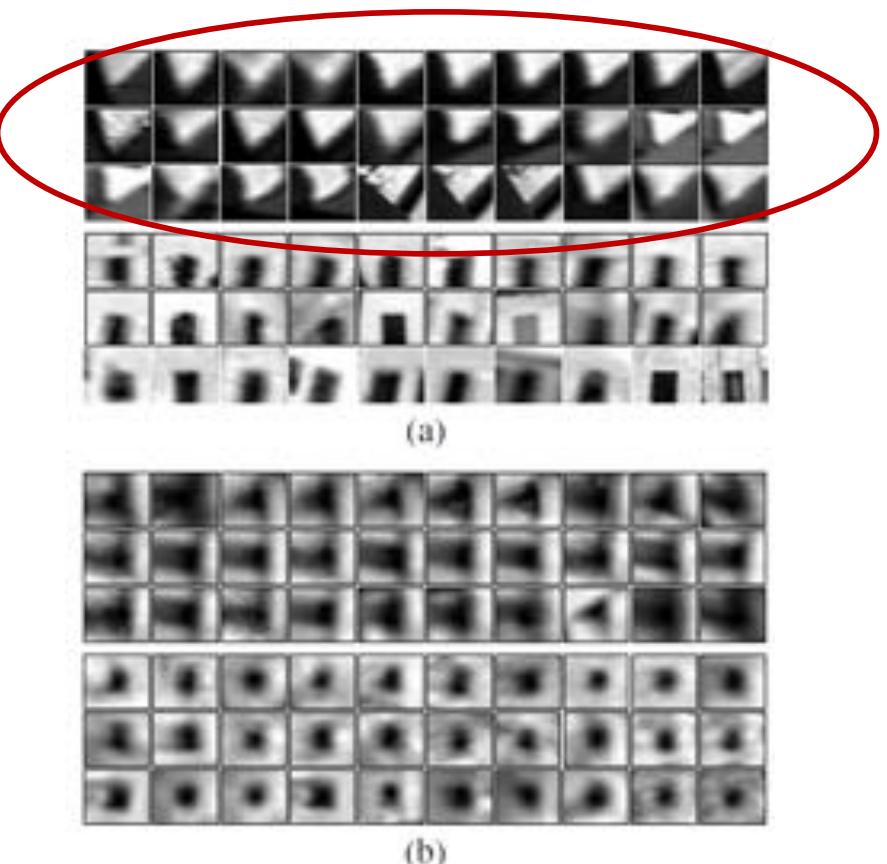


"SIFT" & Object Recognition, David Lowe, 1999



Computer Vision – Bag-of-words

Sivic and Zisserman (2003) “Video Google”

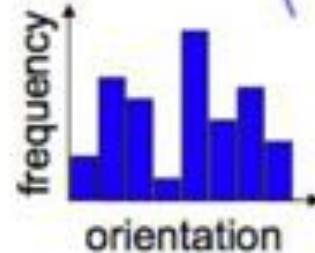
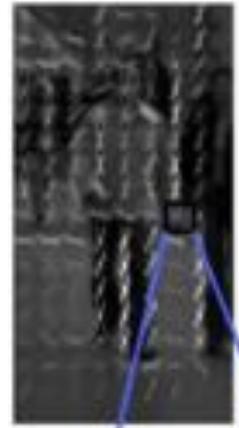


SIFT → SA and MS regions

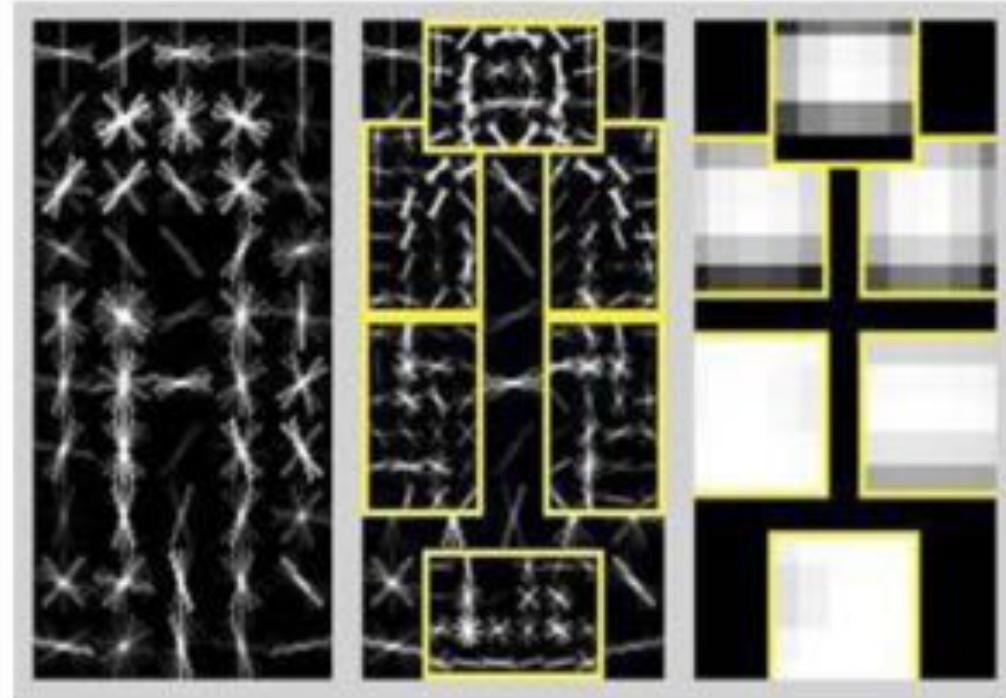
code words



Computer Vision - Shape matching



Histogram of Gradients (HoG)
Dalal & Triggs, 2005



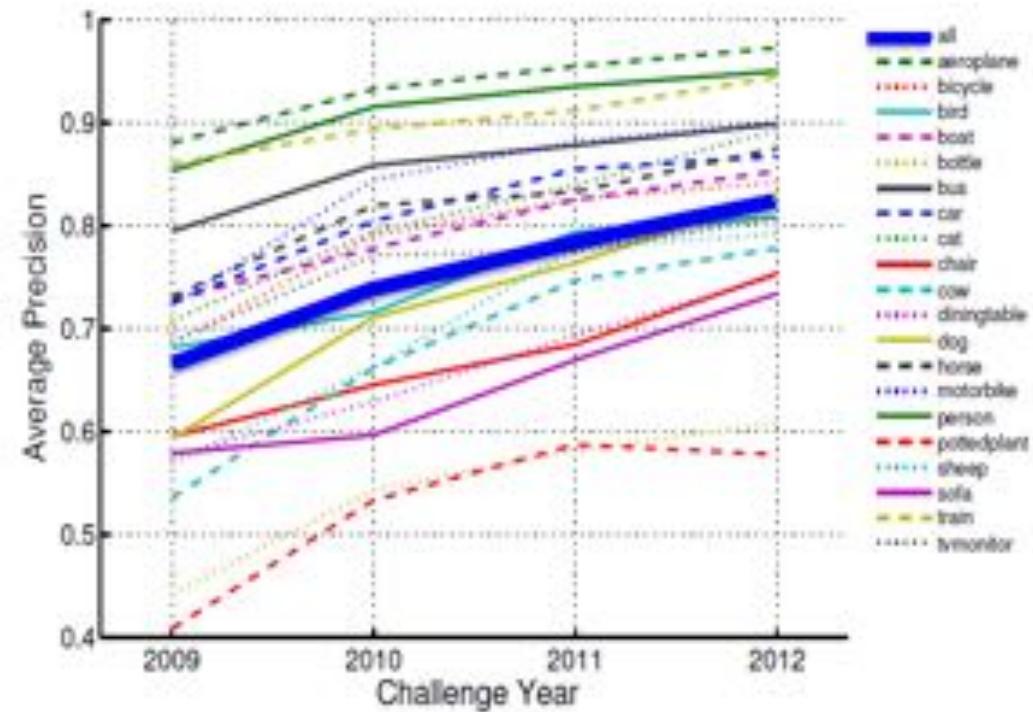
Deformable Part Model
Felzenswalb, McAllester, Ramanan,
2009



Computer Vision – Challenge Datasets

PASCAL Visual Object Challenge (20 object categories)

[Everingham et al. 2006-2012]





www.image-net.org

22K categories and **14M** images

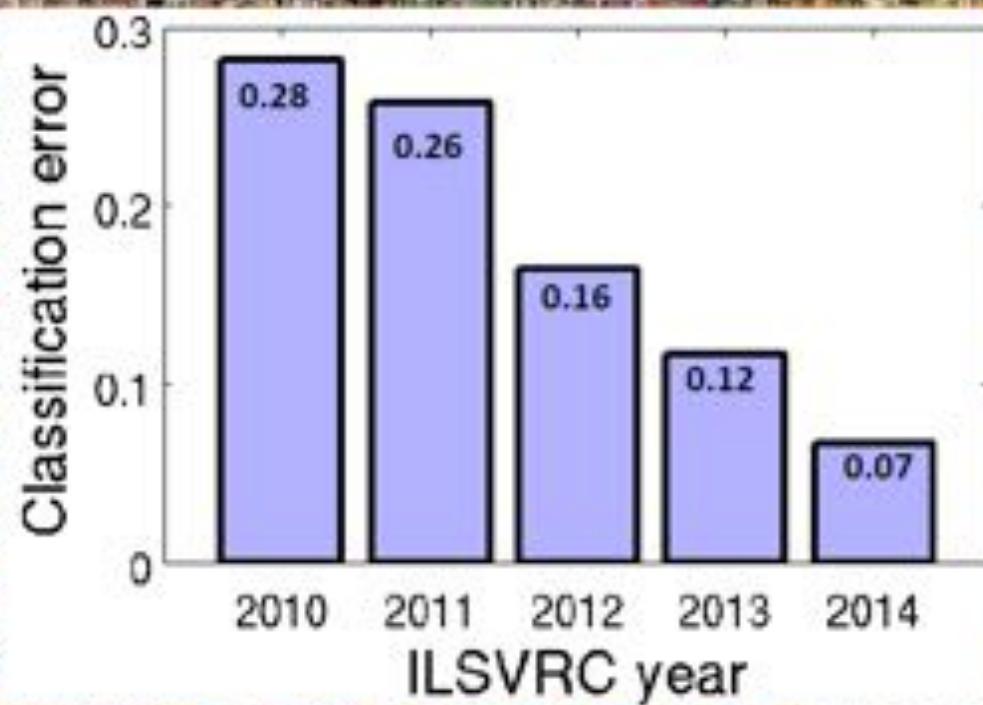
- Animals
 - Bird
 - Fish
 - Mammal
 - Invertebrate
- Plants
 - Tree
 - Flower
 - Food
 - Materials
- Structures
 - Artifact
 - Tools
 - Appliances
 - Structures
- Person
- Scenes
- Indoor
- Geological Formations
- Sport Activities

Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009



IMAGENET Large Scale Visual Recognition Challenge

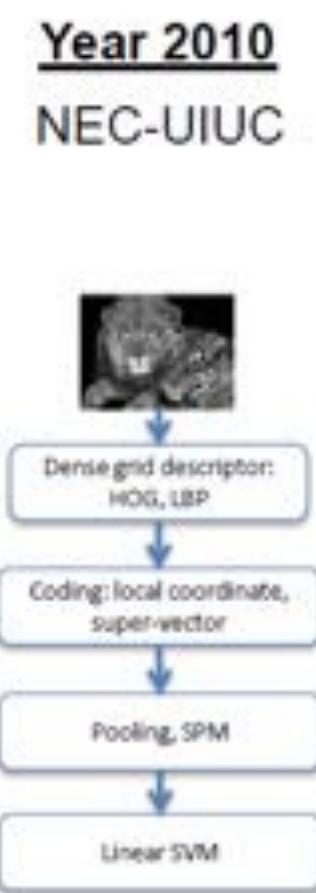
The Image Classification Challenge:
1,000 object classes
1,431,167 images



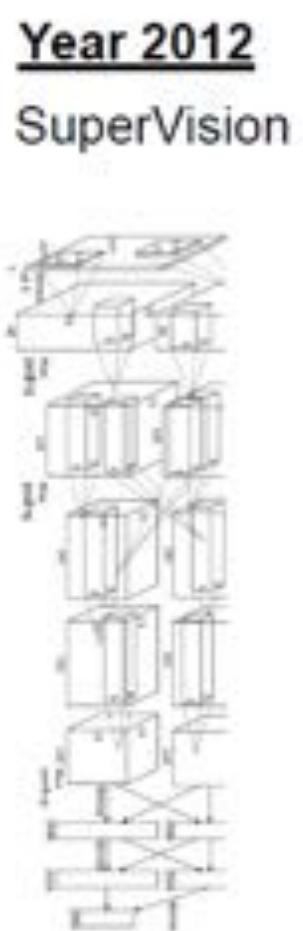
Russakovsky et al. arXiv, 2014



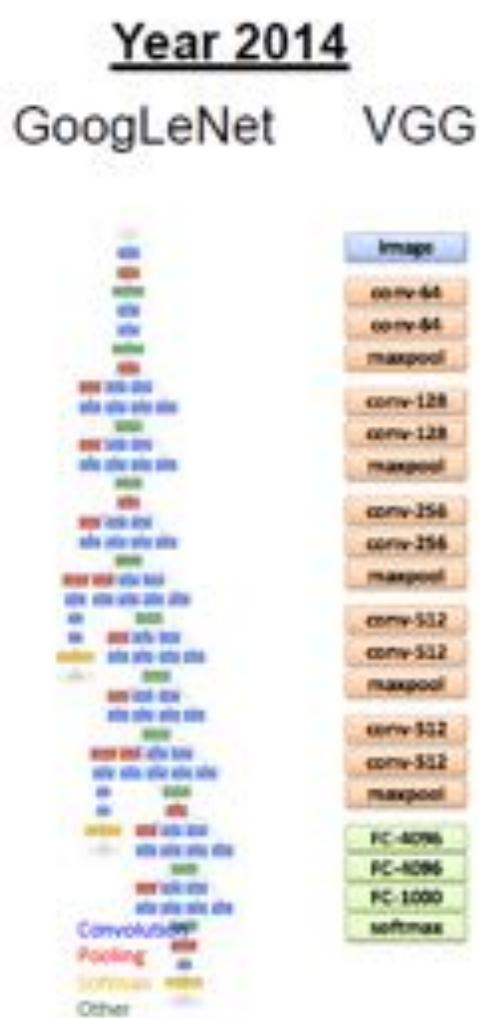
IMAGENET Large Scale Visual Recognition Challenge



[In CVPR 2011]



[Krizhevsky NIPS 2012]



[Szegedy arxiv 2014]

[Simonyan arxiv 2014]



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)

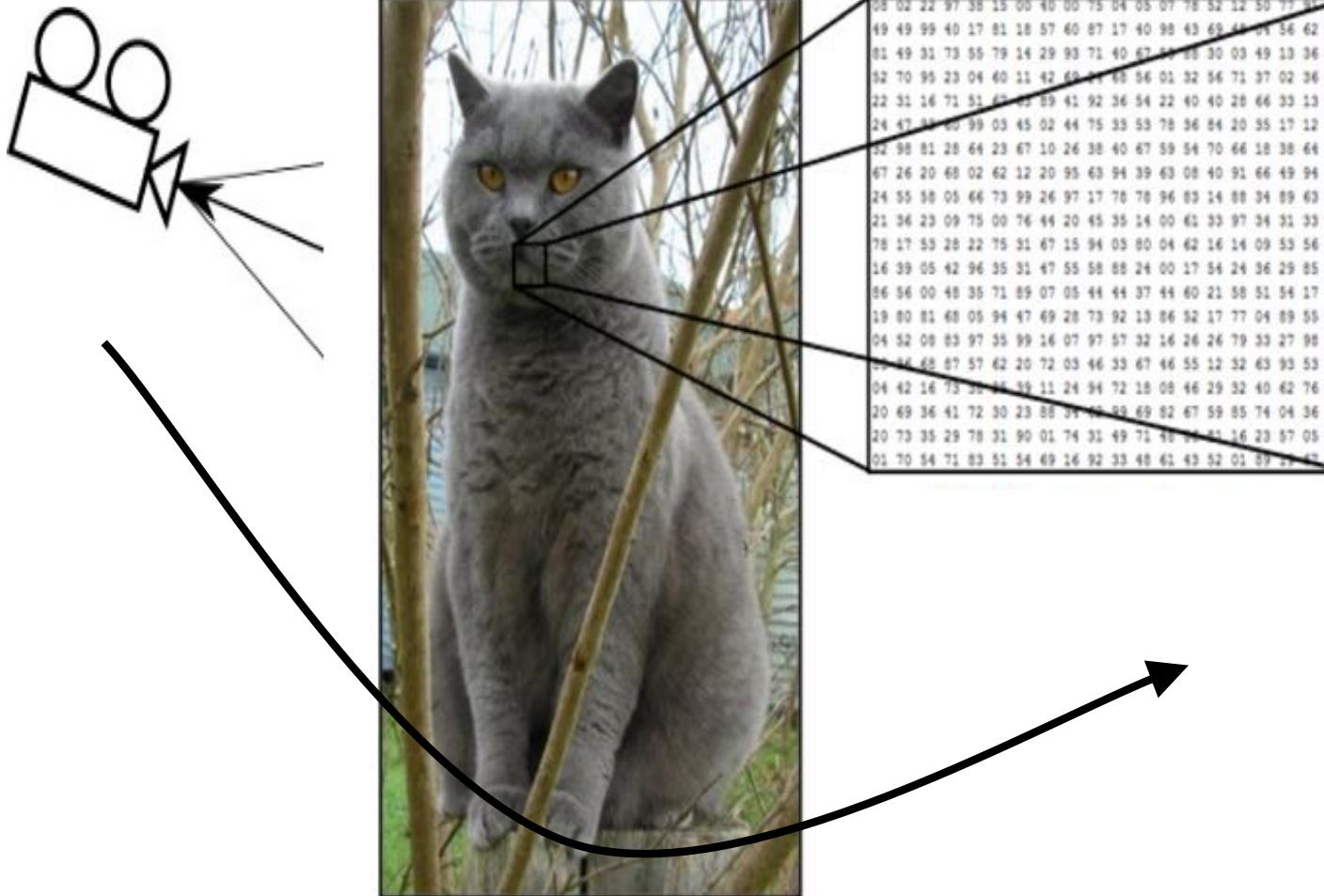


08	02	22	97	38	15	00	40	00	73	04	05	07	78	32	12	30	77	31	60
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	13	30	03	49	13	36	65	
52	70	95	23	04	60	11	42	69	44	69	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	03	59	41	92	36	54	22	40	40	28	66	33	13	80
24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

What the computer sees



Challenges: Viewpoint Variation



Challenges: Illumination



Challenges: Deformation



Challenges: Occlusion



Challenges: Background clutter



Challenges: Intraclass variation



An image classifier

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

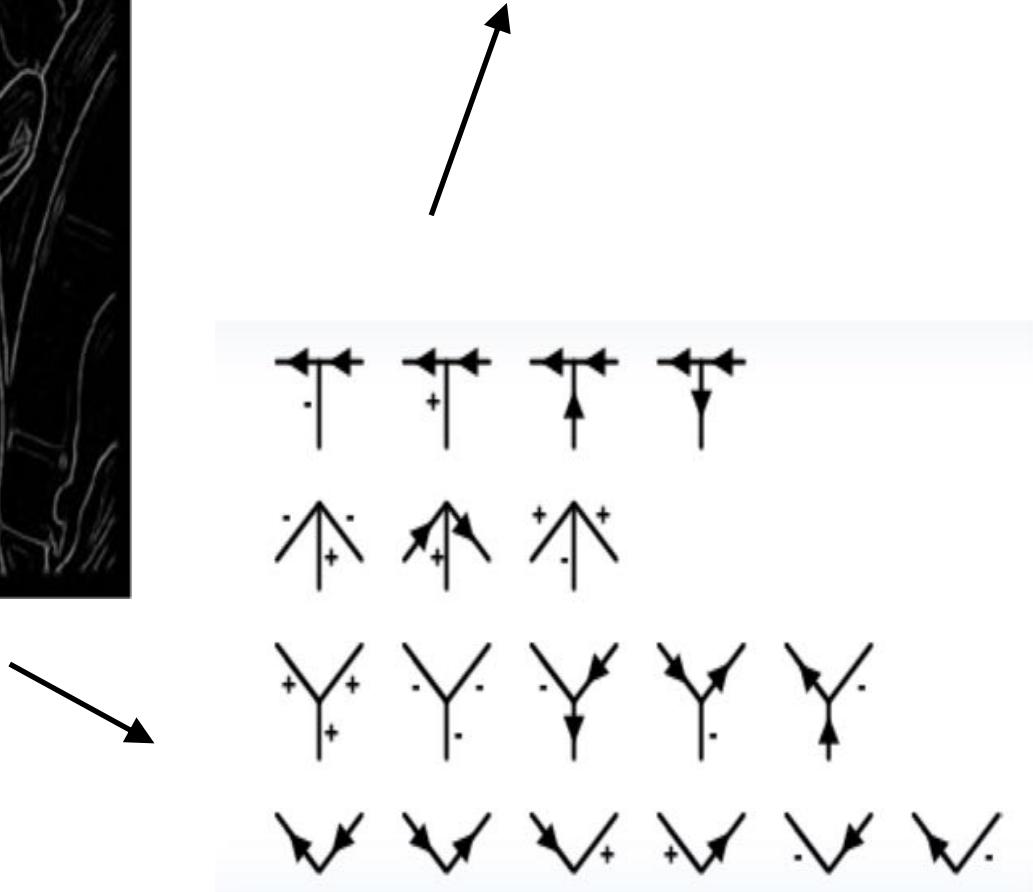
Unlike e.g. sorting a list of numbers,

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



Classifying using constraints? ?

???



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

Example training set

```
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
```



First classifier: Nearest Neighbor Classifier

Remember all training images
and their labels

```
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
```



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.



Example dataset: CIFAR-10

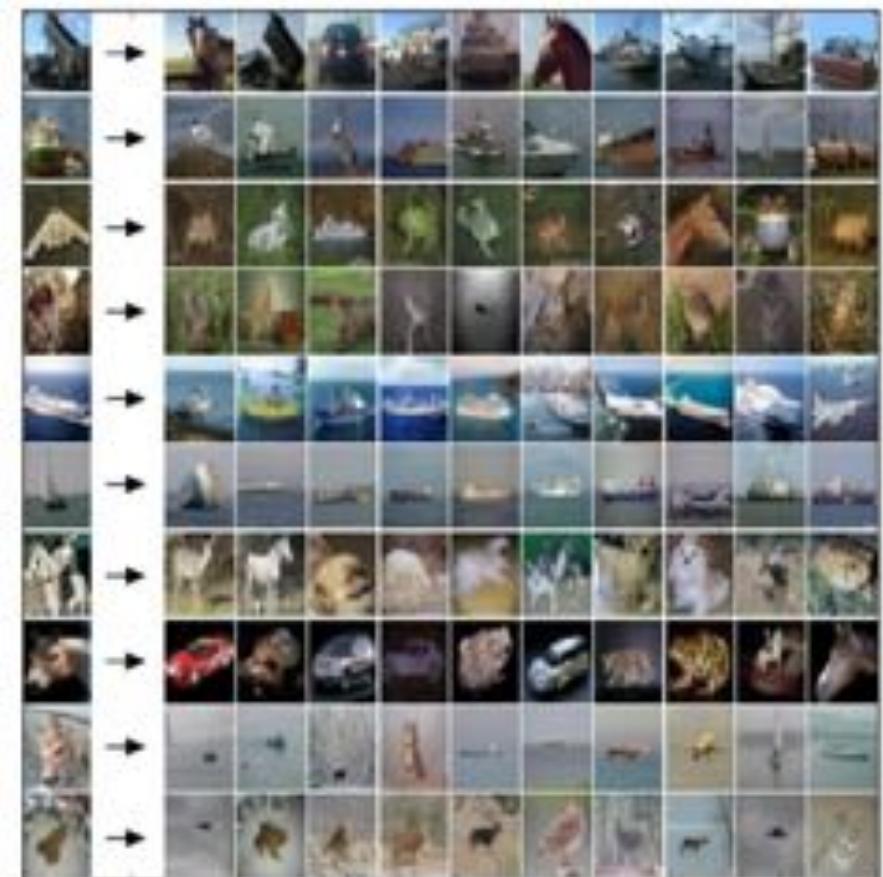
10 labels

50,000 training images

10,000 test images.



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
2	0	255	220	4	32	233	112	2	32	22	108

- = → 456

add



```
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
```

Nearest Neighbor classifier

remember 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

```

Nearest Neighbor classifier

for every test image:

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



```

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?



```

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

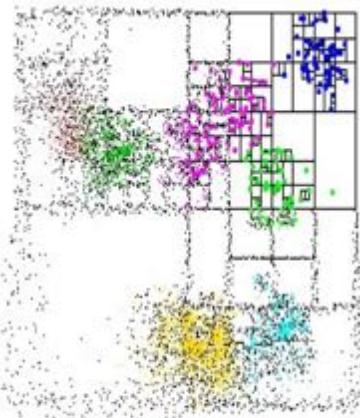
(and you can even learn binary hash codes)

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

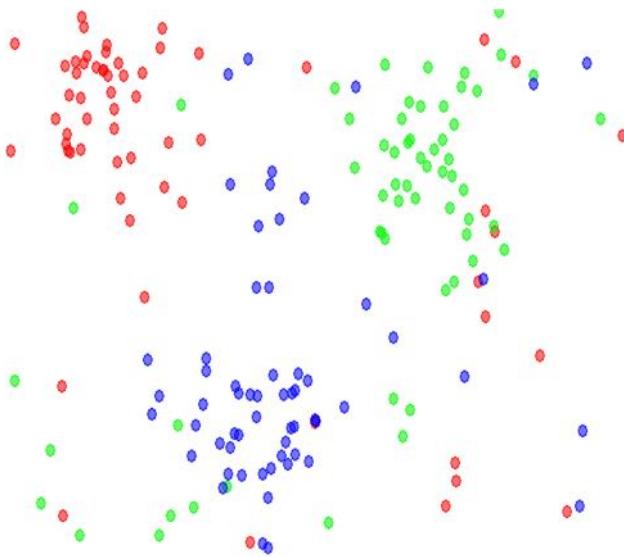
$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



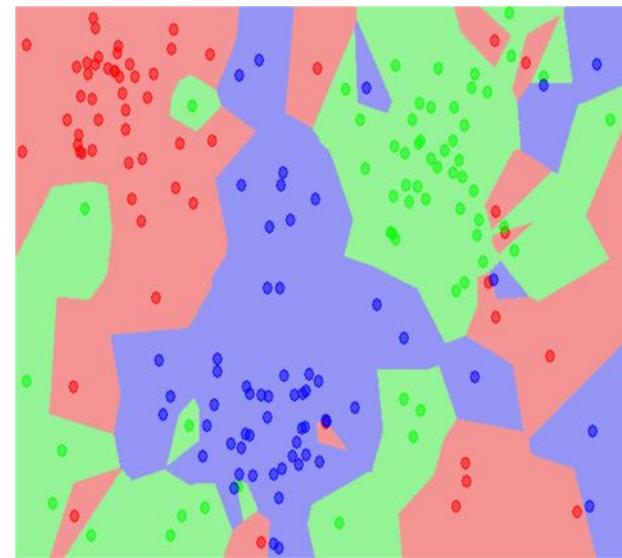
k-Nearest Neighbor

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

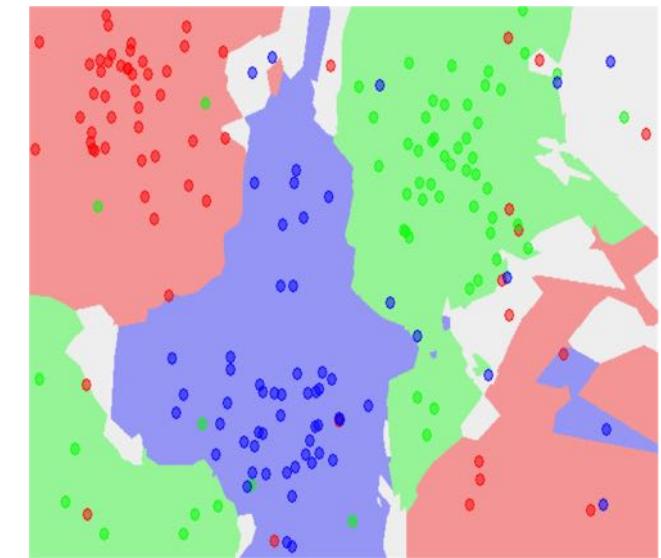
the data



NN classifier



5-NN classifier



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



Example dataset: CIFAR-10

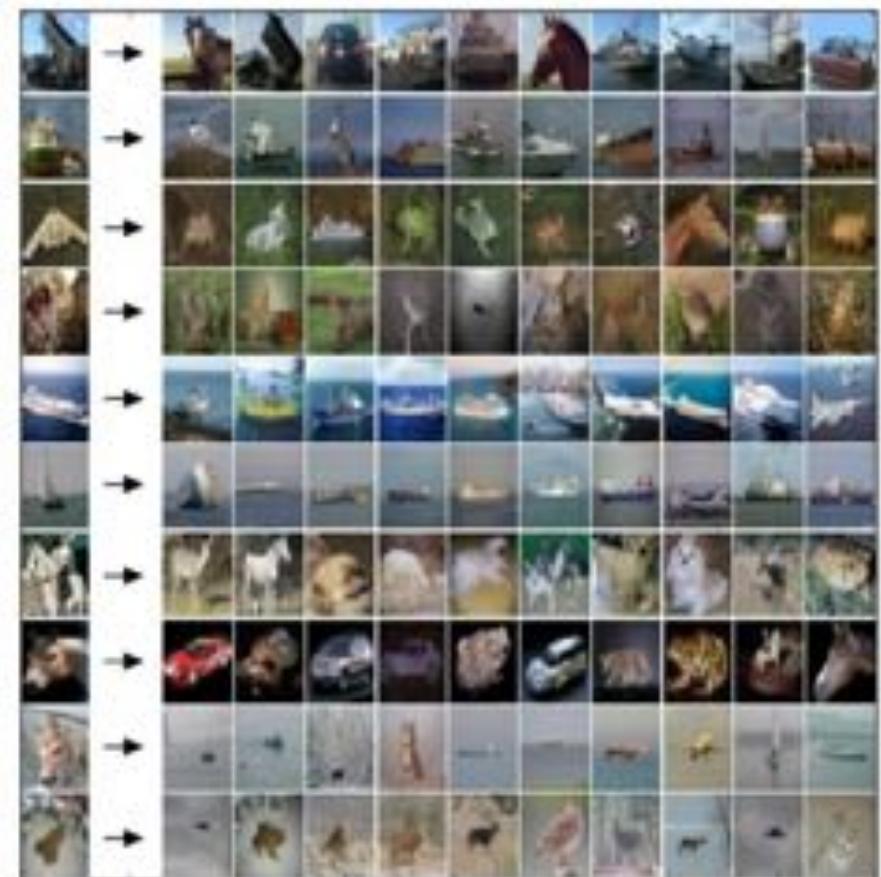
10 labels

50,000 training images

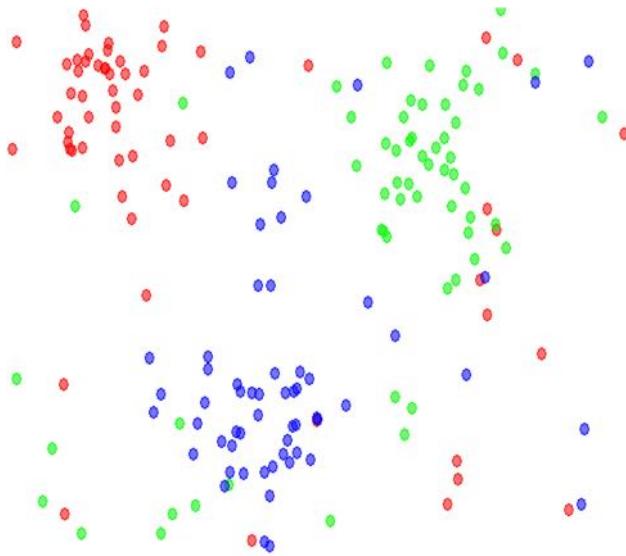
10,000 test images.



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



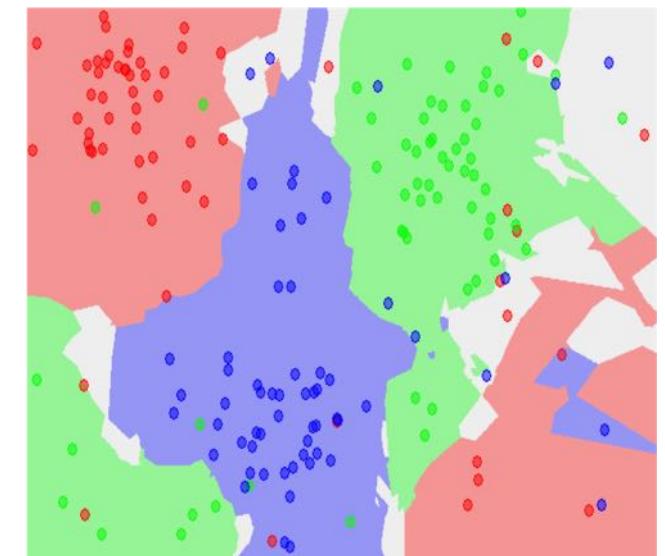
the data



NN classifier



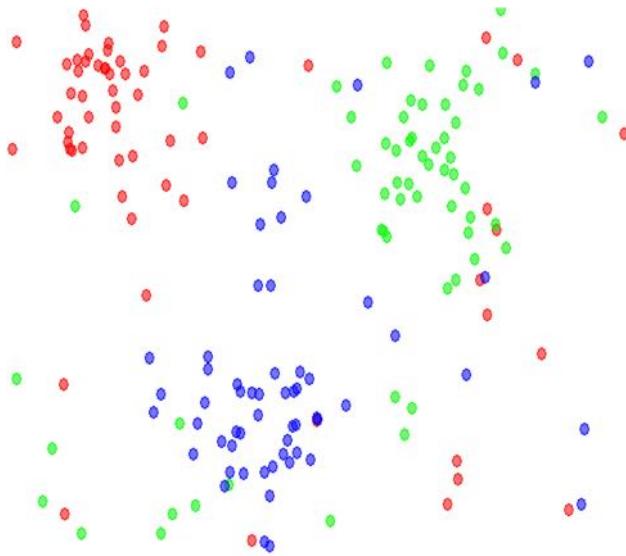
5-NN classifier



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



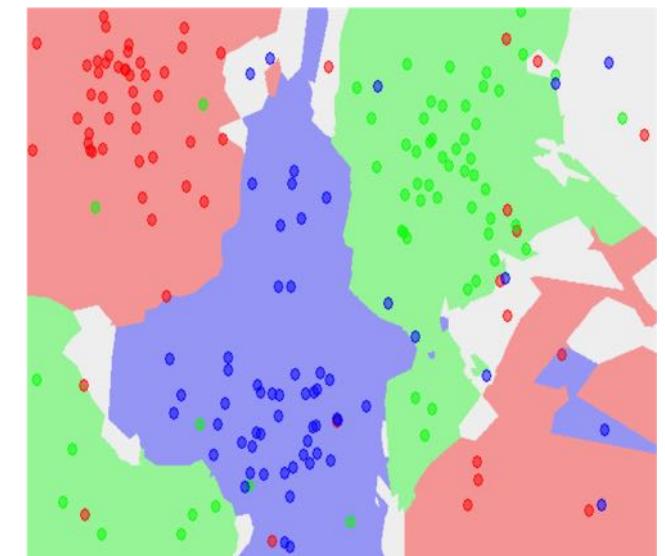
the data



NN classifier



5-NN classifier



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



Hyperparameter tuning:

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

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

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



Bias and Variance

Recall from statistics (or ML), given a true function $y=f(x)$ (mapping an input x to a label y), a machine learning algorithm is a statistical estimator $g_D(x)$. Here D is a data sample.

Estimators have:

- **Bias:** $g_D(x)$ is systematically different from $f(x)$, i.e. $E_D(g_D(x)) \neq f(x)$
- **Variance:** $\text{VAR}_D(g_D(x))$. Even if no bias, $g_D(x)$ may be far from $f(x)$ at most points x .



Bias and Variance

Increasing k in the kNN algorithm should have what effect on:

Bias: ?

Variance: ?



Bias and Variance

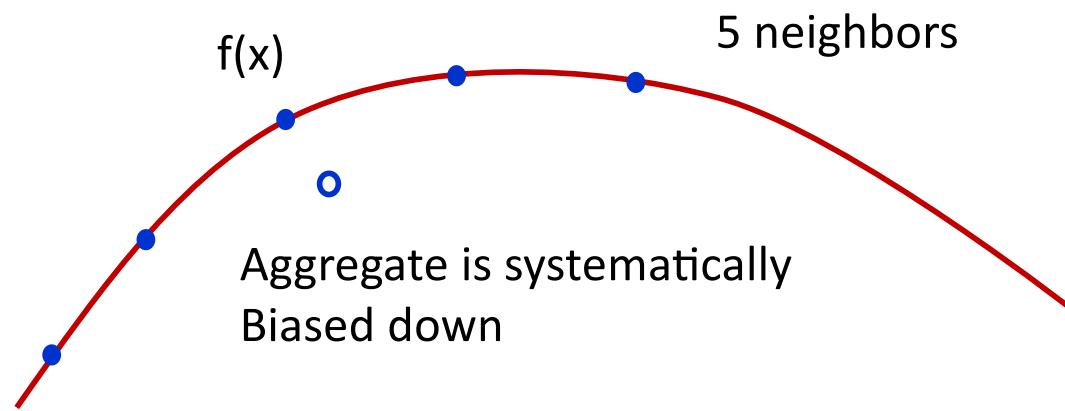
Increasing k in the kNN algorithm should have what effect on:

Bias: Should **increase**. The large k is, the further (on average) are the points being aggregated from x . i.e. the value depends on $f(x')$ for x' further from x .

Variance: ?



Bias



Bias and Variance

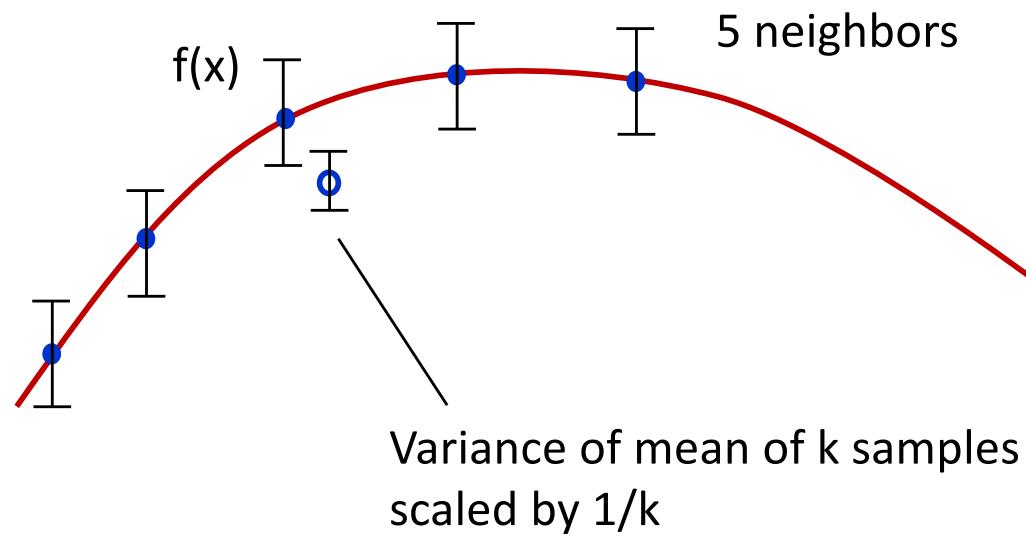
Increasing k in the kNN algorithm should have what effect on:

Bias: Should increase. The large k is, the further (on average) are the points being aggregated from x . i.e. the value depends on $f(x')$ for x' further from x .

Variance: Should **decrease**. The average or majority vote of a set of equal-variance values has lower variance than each value.



Variance



Bias and Variance Rule of Thumb

Compared to simpler (fewer parameters), complex models have what kind of Bias and Variance?:

Bias: ?

Variance:?



Bias and Variance Rule of Thumb

Compared to simpler (fewer parameters), complex models have what kind of Bias and Variance?:

Bias: ? **Lower**, because complex models can better model local variation in $f(x)$.

Variance:?



Bias and Variance Rule of Thumb

Compared to simpler (fewer parameters), complex models have what kind of Bias and Variance?:

Bias: ? Lower, because complex models can better model local variation in $f(x)$.

Variance:? **Higher**, because the parameters are generally more sensitive to few data values.



Bias and Variance Rule of Thumb

Compared to simpler (fewer parameters), complex models have what kind of Bias and Variance?:

Bias: ? **Lower**, because complex models can better model local variation in $f(x)$.

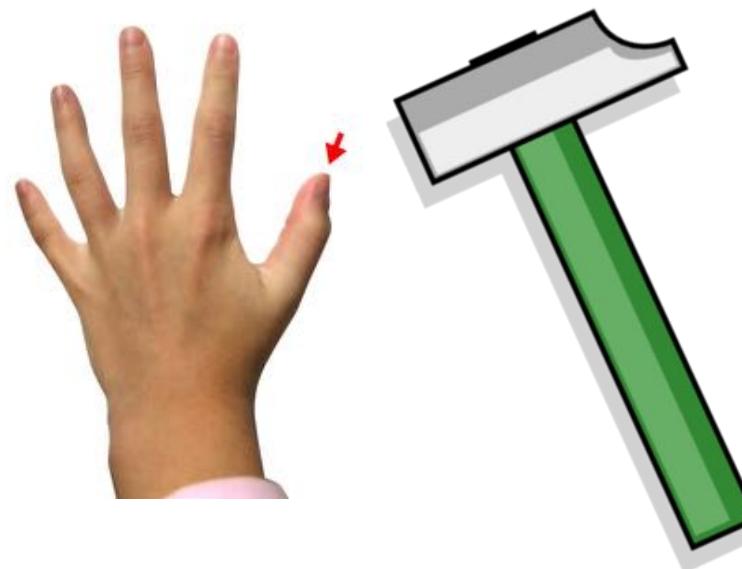
Variance:? **Higher**, because the parameters are generally more sensitive to few data values.

Complex models are prone to overfitting. They require/benefit from large training datasets in order to reduce variance. This includes most **Deep Networks**.



Bias and Variance Rule of Thumb

This rule of thumb is important to keep in mind, but it is only a **rule of thumb**, not a theorem. It applies to typical ML algorithms and deep networks, but not all.



A mathematical theorem is like a hammer. A rule of thumb is no match for it...



Hyperparameter tuning:

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

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

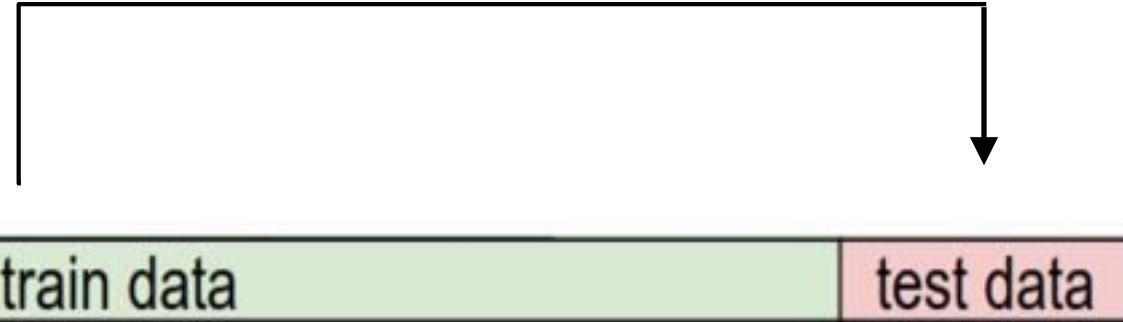
i.e. how do we set the **hyperparameters**?

Very problem-dependent.

Must try them all out and see what works best.



Try out what hyperparameters work best on test set.

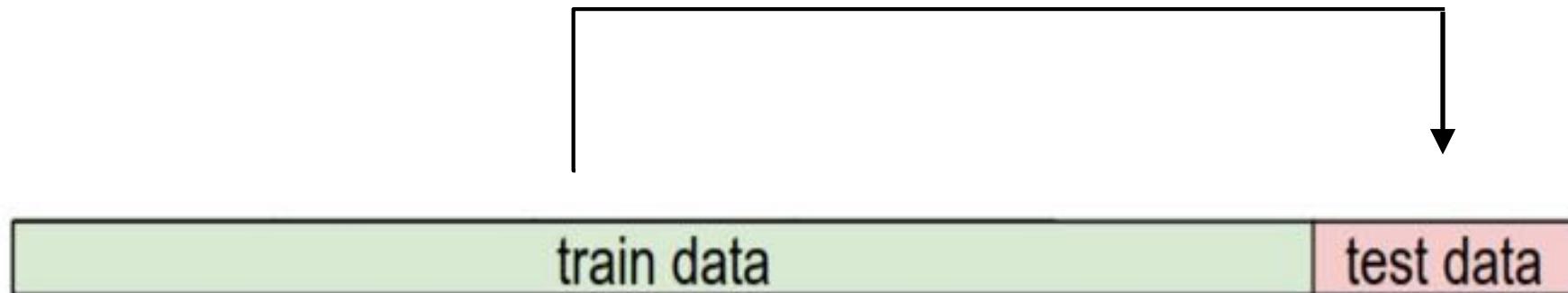


Trying out what hyperparameters work best on test set:

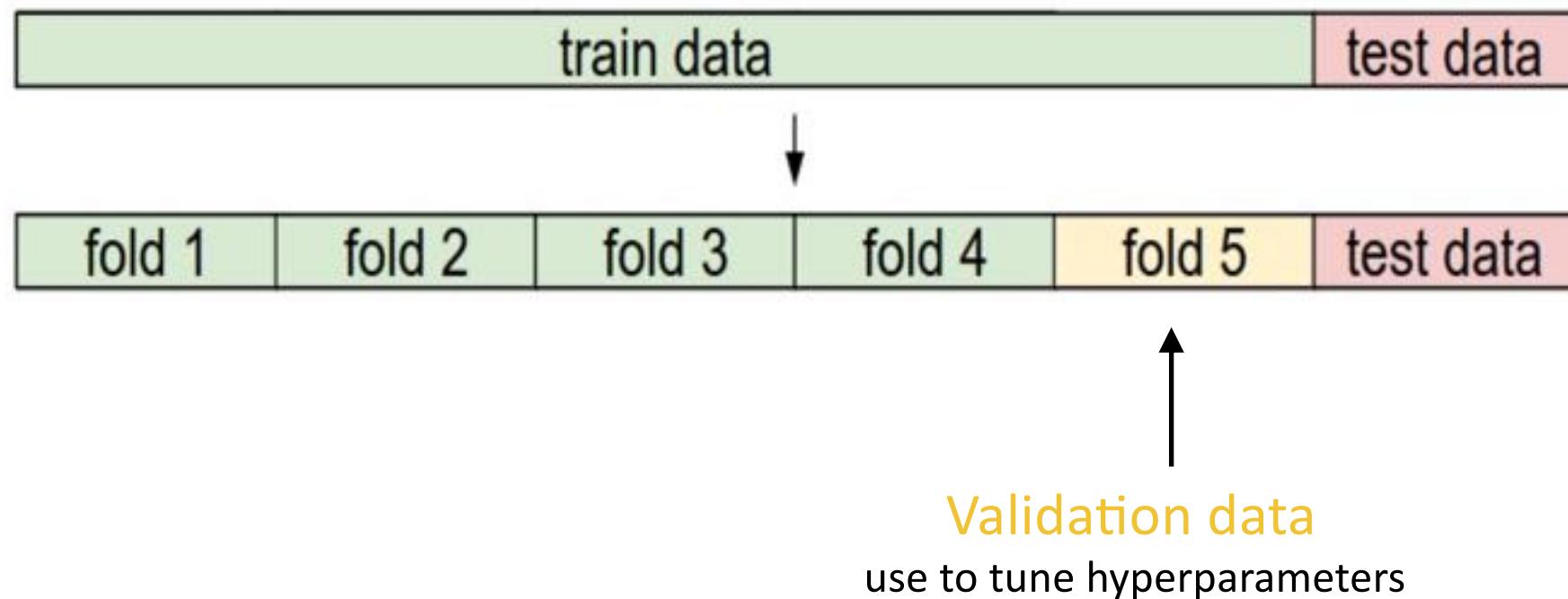


TECHNISCHE
UNIVERSITÄT
DARMSTADT

Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY**, at the end.



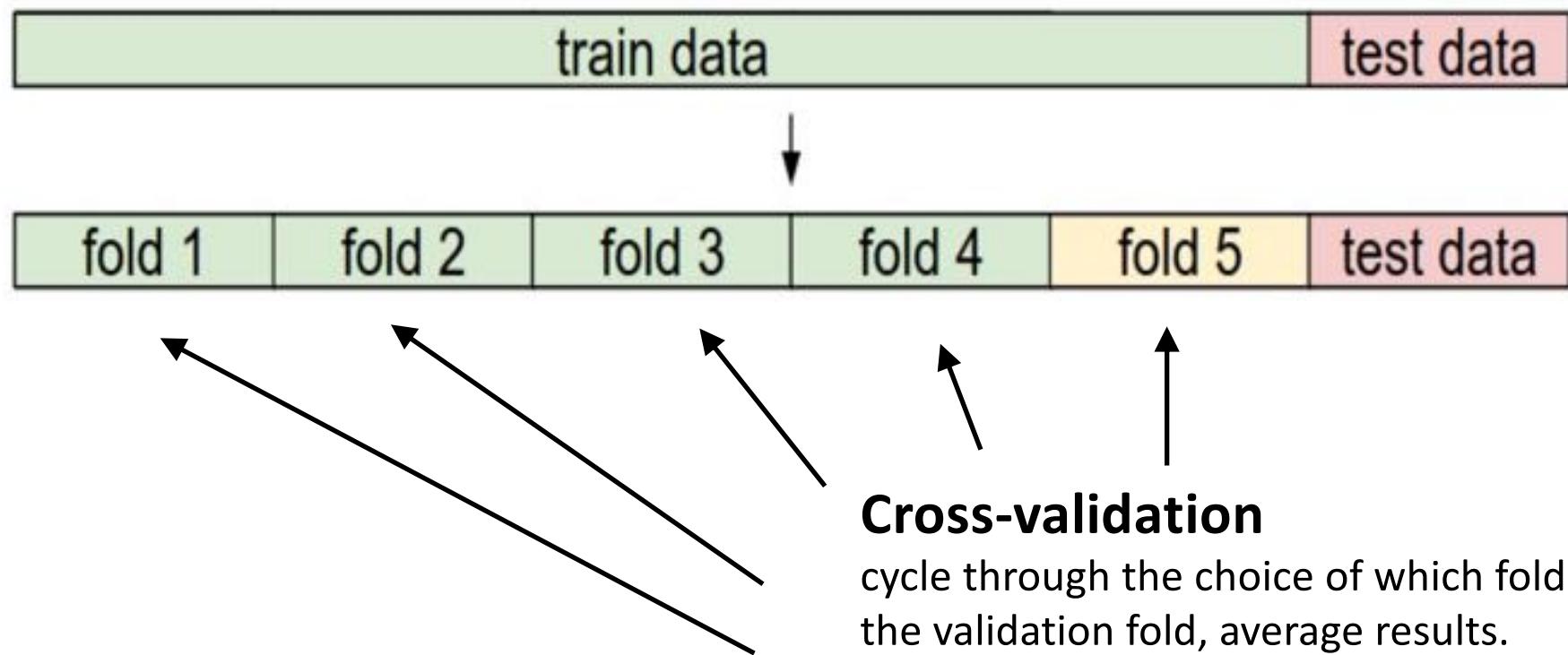
Validation Set:



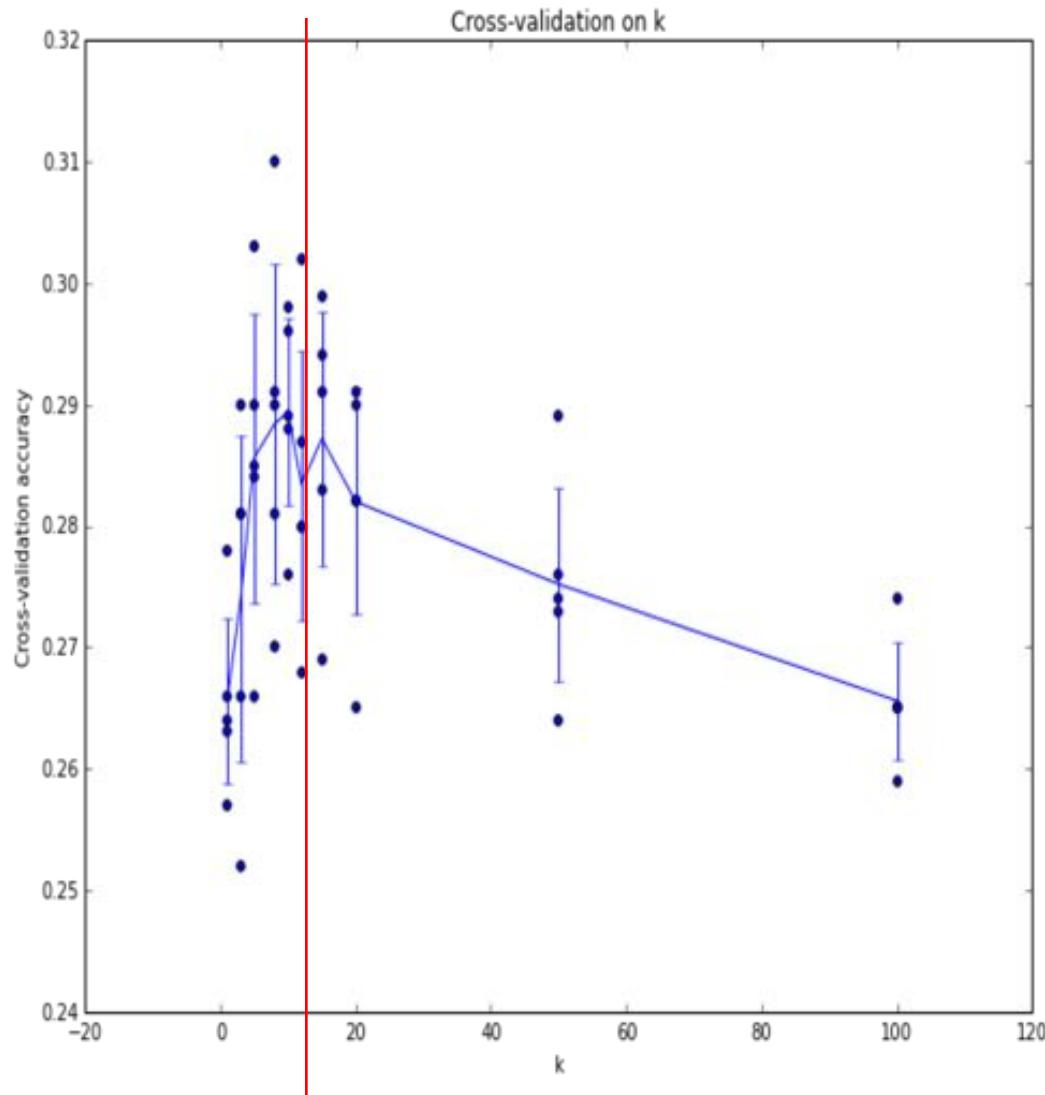
Validation Set:



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Hyperparameter tuning:



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

Each point: single
outcome.

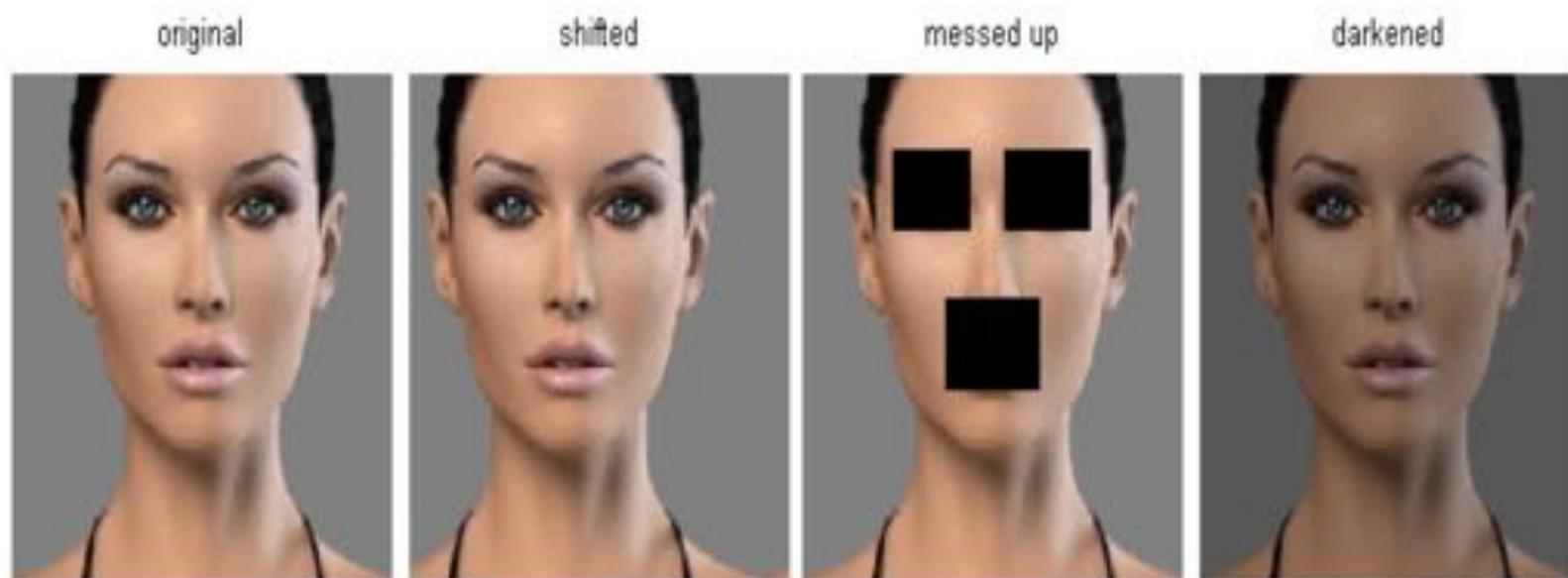
The line goes
through the mean, bars
indicated standard
deviation

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



k-Nearest Neighbor on images **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)



What have you learnt?

- **Image Classification:** given a **Training Set** of labeled images, predict labels on **Test Set**. Common to report the **Accuracy** of predictions (fraction of correct predictions)
- We introduced the **k-Nearest Neighbor Classifier**, which predicts 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.

