

# CSC 2515: Machine Learning

## Lecture 1 - Introduction and Nearest Neighbours

Roger Grosse

University of Toronto

# This course

- Broad introduction to machine learning
  - ▶ First half: algorithms and principles for supervised learning
    - ▶ nearest neighbors, decision trees, ensembles, linear regression, logistic regression, SVMs
    - ▶ neural nets!
  - ▶ Unsupervised learning: PCA, K-means, mixture models
  - ▶ Basics of reinforcement learning
- This course is taught as a stand-alone grad course for the first time.
  - ▶ But the structure and difficulty will be similar to past years, when it was cross-listed as an undergrad course.
  - ▶ The majority of students are from outside Computer Science.

# Course Information

Course Website:

[https://www.cs.toronto.edu/~rgrosse/courses/csc2515\\_2019/](https://www.cs.toronto.edu/~rgrosse/courses/csc2515_2019/)

Slides will be posted to web page in advance of lecture, but I'll continue to make edits up to Thursday night. So please re-download!

We will use Piazza for **discussions**.

- URL to be sent out
- Your grade **does not depend on your participation on Piazza**. It's just a good way for asking questions, discussing with your instructor, TAs and your peers

# Course Information

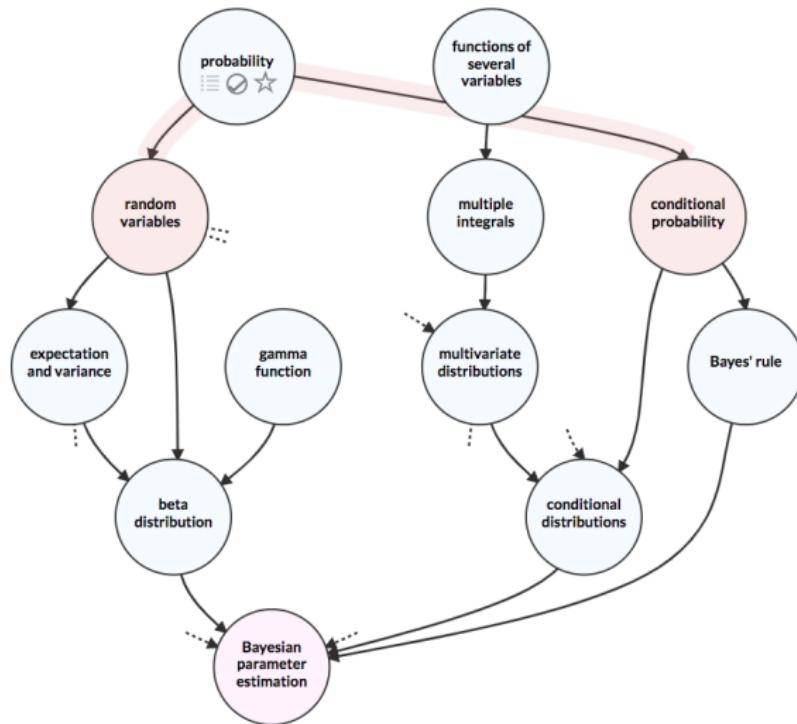
Recommended readings will be given for each lecture. But the following will be useful throughout the course:

- Hastie, Tibshirani, and Friedman: “The Elements of Statistical Learning”
- Christopher Bishop: “Pattern Recognition and Machine Learning”, 2006.
- Kevin Murphy: “Machine Learning: a Probabilistic Perspective”, 2012.
- David Mackay: “Information Theory, Inference, and Learning Algorithms”, 2003.
- Shai Shalev-Shwartz & Shai Ben-David: “Understanding Machine Learning: From Theory to Algorithms”, 2014.

There are lots of freely available, high-quality ML resources.

# Course Information

See Metacademy (<https://metacademy.org>) for additional background, and to help review prerequisites.



# Requirements and Marking

- 5 written homeworks, due roughly every other week.
  - ▶ Combination of pencil & paper derivations and short programming exercises
  - ▶ Each counts for 10%, except that the lowest mark counts for 5%.
  - ▶ Worth 45% in total.
- Read some classic papers.
  - ▶ Worth 5%, honor system.
- Midterm
  - ▶ Oct. 30, 4-6pm
  - ▶ Worth 15% of course mark
- Final Exam
  - ▶ Dec. 17, 3-6pm
  - ▶ Worth 35% of course mark

# More on Assignments

**Collaboration** on the assignments is not allowed. Each student is responsible for his/her own work. Discussion of assignments should be limited to clarification of the handout itself, and should not involve any sharing of pseudocode or code or simulation results. Violation of this policy is grounds for a semester grade of F, in accordance with university regulations.

The schedule of assignments will be posted on the course web page.

Assignments should be handed in by 11:59pm; a late penalty of 10% per day will be assessed thereafter (up to 3 days, then submission is blocked).

Extensions will be granted only in special situations, and you will need a Student Medical Certificate or a written request approved by the course coordinator at least one week before the due date.

What is learning?

*"The activity or process of gaining knowledge or skill by studying, practicing, being taught, or experiencing something."*

*Merriam Webster dictionary*

*"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E."*

*Tom Mitchell*

# What is machine learning?

- For many problems, it's difficult to program the correct behavior by hand
  - ▶ recognizing people and objects
  - ▶ understanding human speech
- Machine learning approach: program an algorithm to automatically learn from data, or from experience
- Why might you want to use a learning algorithm?
  - ▶ hard to code up a solution by hand (e.g. vision, speech)
  - ▶ system needs to adapt to a changing environment (e.g. spam detection)
  - ▶ want the system to perform *better* than the human programmers
  - ▶ privacy/fairness (e.g. ranking search results)

# What is machine learning?

- It's similar to statistics...
  - ▶ Both fields try to uncover patterns in data
  - ▶ Both fields draw heavily on calculus, probability, and linear algebra, and share many of the same core algorithms
- But it's not statistics!
  - ▶ Stats is more concerned with helping scientists and policymakers draw good conclusions; ML is more concerned with building autonomous agents
  - ▶ Stats puts more emphasis on interpretability and mathematical rigor; ML puts more emphasis on predictive performance, scalability, and autonomy

# What is machine learning?

- Types of machine learning
  - ▶ **Supervised learning:** have labeled examples of the correct behavior
  - ▶ **Reinforcement learning:** learning system receives a reward signal, tries to learn to maximize the reward signal
  - ▶ **Unsupervised learning:** no labeled examples – instead, looking for interesting patterns in the data

# History of machine learning

- 1957 — Perceptron algorithm (implemented as a circuit!)
- 1959 — Arthur Samuel wrote a learning-based checkers program that could defeat him
- 1969 — Minsky and Papert's book *Perceptrons* (limitations of linear models)
- 1980s — Some foundational ideas
  - ▶ Connectionist psychologists explored neural models of cognition
  - ▶ 1984 — Leslie Valiant formalized the problem of learning as PAC learning
  - ▶ 1988 — Backpropagation (re-)discovered by Geoffrey Hinton and colleagues
  - ▶ 1988 — Judea Pearl's book *Probabilistic Reasoning in Intelligent Systems* introduced Bayesian networks

# History of machine learning

- 1990s — the “AI Winter”, a time of pessimism and low funding
- But looking back, the ’90s were also sort of a golden age for ML research
  - ▶ Markov chain Monte Carlo
  - ▶ variational inference
  - ▶ kernels and support vector machines
  - ▶ boosting
  - ▶ convolutional networks
- 2000s — applied AI fields (vision, NLP, etc.) adopted ML
- 2010s — deep learning
  - ▶ 2010–2012 — neural nets smashed previous records in speech-to-text and object recognition
  - ▶ increasing adoption by the tech industry
  - ▶ 2016 — AlphaGo defeated the human Go champion

Computer vision: Object detection, semantic segmentation, pose estimation, and almost every other task is done with ML.



Figure 4. More results of Mask R-CNN on COCO test images, using ResNet-101-FPN and running at 5 fps, with 35.7 mask AP (Table 1).



Source: <https://www.youtube.com/watch?v=YUOOWAqgjg>

Instance segmentation - [Link](#)



DAQUAR 1553  
What is there in front of the sofa?  
Ground truth: table  
IMG+BOW: table (0.74)  
2-VIS+BLSTM: table (0.88)  
LSTM: chair (0.47)



COCOQA 5078  
How many leftover donuts is the red bicycle holding?  
Ground truth: three  
IMG+BOW: two (0.51)  
2-VIS+BLSTM: three (0.27)  
BOW: one (0.29)

Speech: Speech to text, personal assistants, speaker identification...



NLP: Machine translation, sentiment analysis, topic modeling, spam filtering.

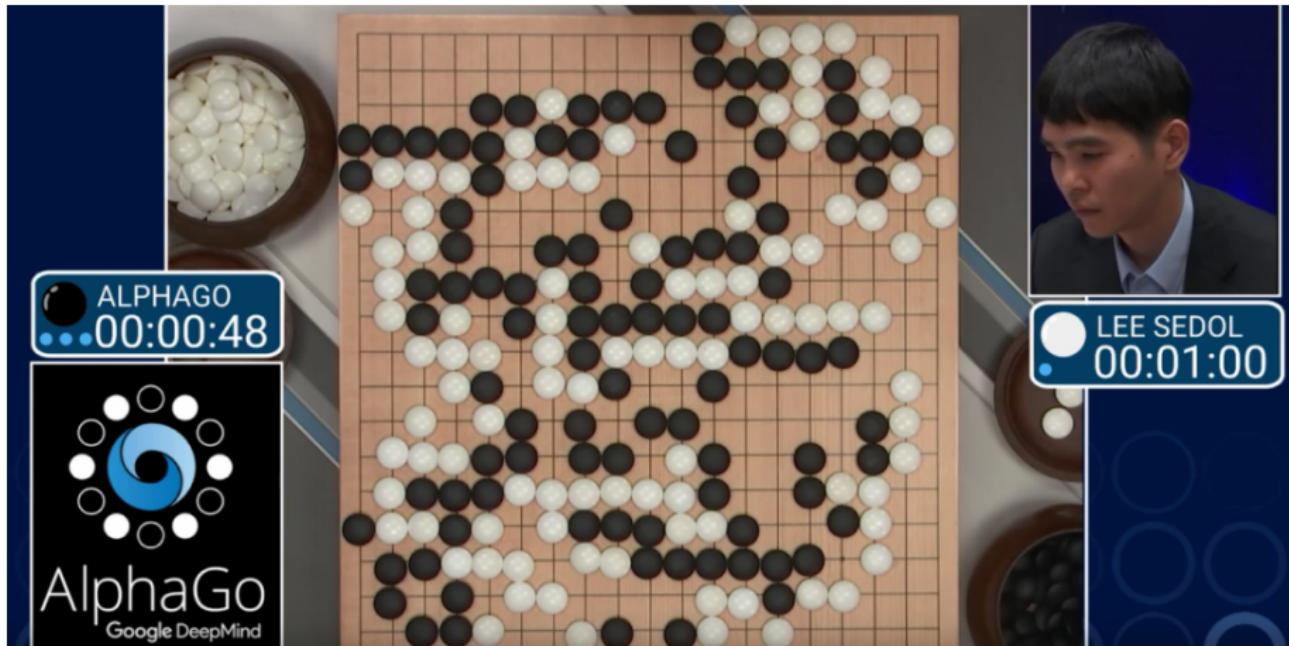
## Real world example:

The New York Times

LDA analysis of 1.8M New York Times articles:

music band songs rock album jazz pop song singer night	book life novel story books man stories love children family	art museum show exhibition artist artists paintings painting century works	game Knicks nets points team season play games night coach	show film television movie series says life man character know
theater play production show stage street broadway director musical directed	clinton bush campaign gore political republican dole presidential senator house	stock market percent fund investors funds companies stocks investment trading	restaurant sauce menu food dishes street dining dinner chicken served	budget tax governor county mayor billion taxes plan legislature fiscal

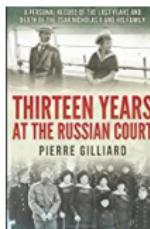
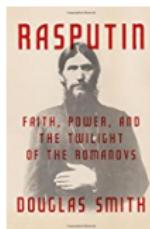
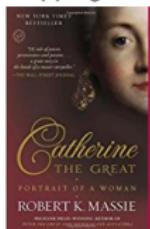
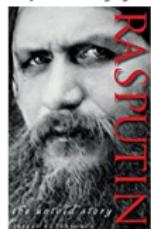
# Playing Games



DOTA2 - [Link](#)

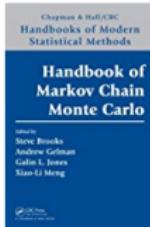
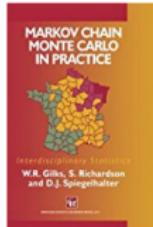
# E-commerce & Recommender Systems : Amazon, netflix, ...

Inspired by your shopping trends



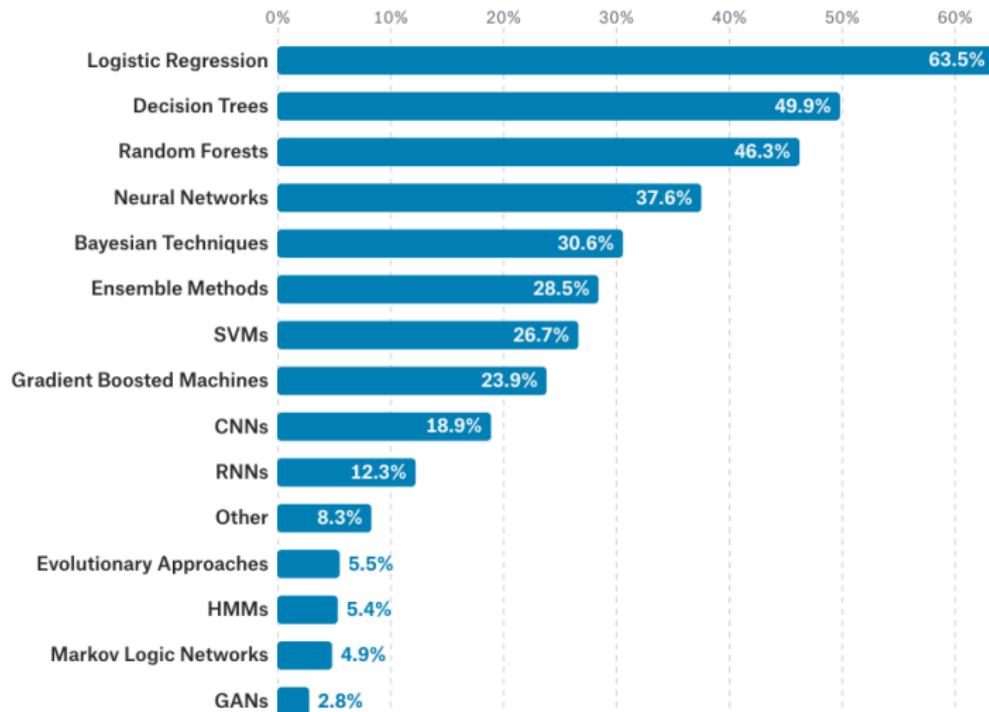
Related to items you've viewed

[See more](#)



# Why this class?

2017 Kaggle survey of data science and ML practitioners: what data science methods do you use at work?



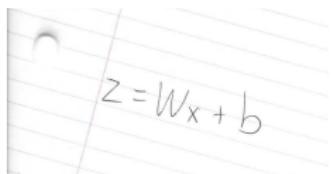
# ML Workflow

ML workflow sketch:

1. Should I use ML on this problem?
  - ▶ Is there a pattern to detect?
  - ▶ Can I solve it analytically?
  - ▶ Do I have data?
2. Gather and organize data.
3. Preprocessing, cleaning, visualizing.
4. Establishing a baseline.
5. Choosing a model, loss, regularization, ...
6. Optimization (could be simple, could be a Phd...).
7. Hyperparameter search.
8. Analyze performance and mistakes, and iterate back to step 5 (or 3).

# Implementing machine learning systems

- You will often need to derive an algorithm (with pencil and paper), and then translate the math into code.
- Array processing (NumPy)
  - ▶ **vectorize** computations (express them in terms of matrix/vector operations) to exploit hardware efficiency
  - ▶ This also makes your code cleaner and more readable!



```
z = np.zeros(m)
for i in range(m):
    for j in range(n):
        z[i] += W[i, j] * x[j]
    z[i] += b[i]
```

$z = \text{np.dot}(W, x) + b$

# Implementing machine learning systems

- Neural net frameworks: PyTorch, TensorFlow, etc.
  - ▶ automatic differentiation
  - ▶ compiling computation graphs
  - ▶ libraries of algorithms and network primitives
  - ▶ support for graphics processing units (GPUs)
- Why take this class if these frameworks do so much for you?
  - ▶ So you know what to do if something goes wrong!
  - ▶ Debugging learning algorithms requires sophisticated detective work, which requires understanding what goes on beneath the hood.
  - ▶ That's why we derive things by hand in this class!

Questions?

?

## Nearest Neighbours

# Introduction

- Today (and for the next 6 weeks) we're focused on **supervised learning**.
- This means we're given a **training set** consisting of **inputs** and corresponding **labels**, e.g.

Task	Inputs	Labels
object recognition	image	object category
image captioning	image	caption
document classification	text	document category
speech-to-text	audio waveform	text
⋮	⋮	⋮

# Input Vectors

What an image looks like to the computer:



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

What the computer sees

→  
image classification  
82% cat  
15% dog  
2% hat  
1% mug

[Image credit: Andrej Karpathy]

# Input Vectors

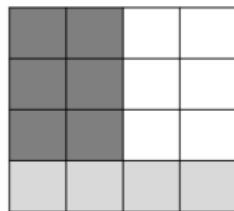
- Machine learning algorithms need to handle lots of types of data: images, text, audio waveforms, credit card transactions, etc.
- Common strategy: represent the input as an **input vector** in  $\mathbb{R}^d$ 
  - ▶ **Representation** = mapping to another space that's easy to manipulate
  - ▶ Vectors are a great representation since we can do linear algebra!



# Input Vectors

Can use raw pixels:

Images  $\leftrightarrow$  Vectors



60	60	255	255
60	60	255	255
60	60	255	255
128	128	128	128

60
60
255
255
60
60
255
255
60
60
255
255
128
128
128
128

Can do much better if you compute a vector of meaningful features.

# Input Vectors

- Mathematically, our training set consists of a collection of pairs of an input vector  $\mathbf{x} \in \mathbb{R}^d$  and its corresponding **target**, or **label**,  $t$ 
  - ▶ **Regression**:  $t$  is a real number (e.g. stock price)
  - ▶ **Classification**:  $t$  is an element of a discrete set  $\{1, \dots, C\}$
  - ▶ These days,  $t$  is often a highly structured object (e.g. image)
- Denote the training set  $\{(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(N)}, t^{(N)})\}$ 
  - ▶ Note: these superscripts have nothing to do with exponentiation!

# Nearest Neighbors

- Suppose we're given a novel input vector  $\mathbf{x}$  we'd like to classify.
- The idea: find the nearest input vector to  $\mathbf{x}$  in the training set and copy its label.
- Can formalize "nearest" in terms of Euclidean distance

$$\|\mathbf{x}^{(a)} - \mathbf{x}^{(b)}\|_2 = \sqrt{\sum_{j=1}^d (x_j^{(a)} - x_j^{(b)})^2}$$

## Algorithm:

1. Find example  $(\mathbf{x}^*, t^*)$  (from the stored training set) closest to  $\mathbf{x}$ . That is:

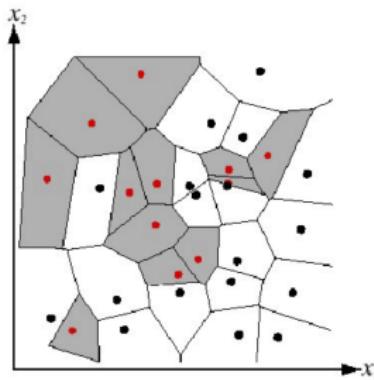
$$\mathbf{x}^* = \underset{\mathbf{x}^{(i)} \in \text{train. set}}{\operatorname{argmin}} \text{distance}(\mathbf{x}^{(i)}, \mathbf{x})$$

2. Output  $y = t^*$

- Note: we don't need to compute the square root. Why?

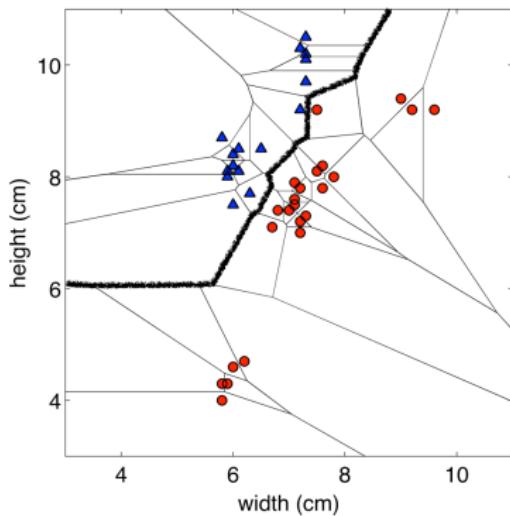
# Nearest Neighbors: Decision Boundaries

We can visualize the behavior in the classification setting using a [Voronoi diagram](#).

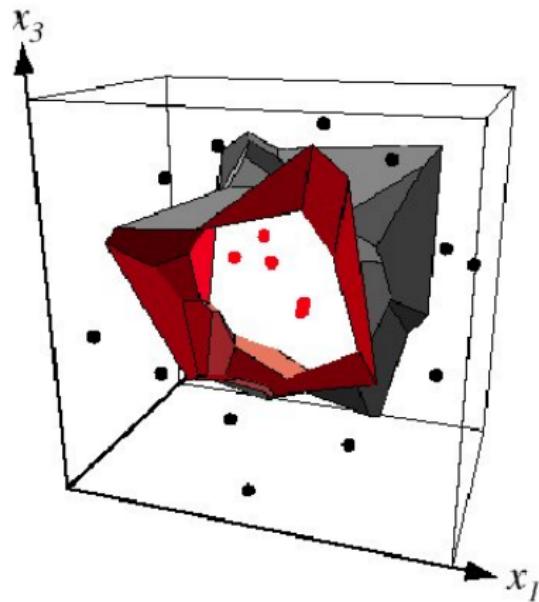


# Nearest Neighbors: Decision Boundaries

**Decision boundary:** the boundary between regions of input space assigned to different categories.



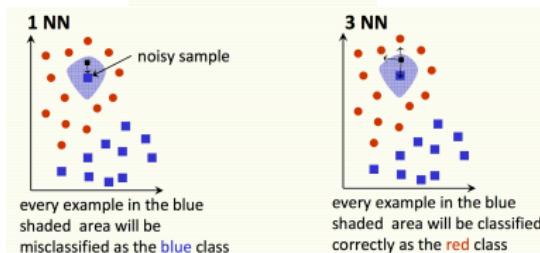
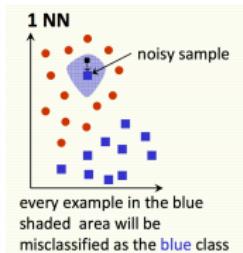
# Nearest Neighbors: Decision Boundaries



Example: 3D decision boundary

# k-Nearest Neighbors

[Pic by Olga Veksler]



- Nearest neighbors **sensitive to noise or mis-labeled data** (“class noise”).  
Solution?
- Smooth by having k nearest neighbors vote

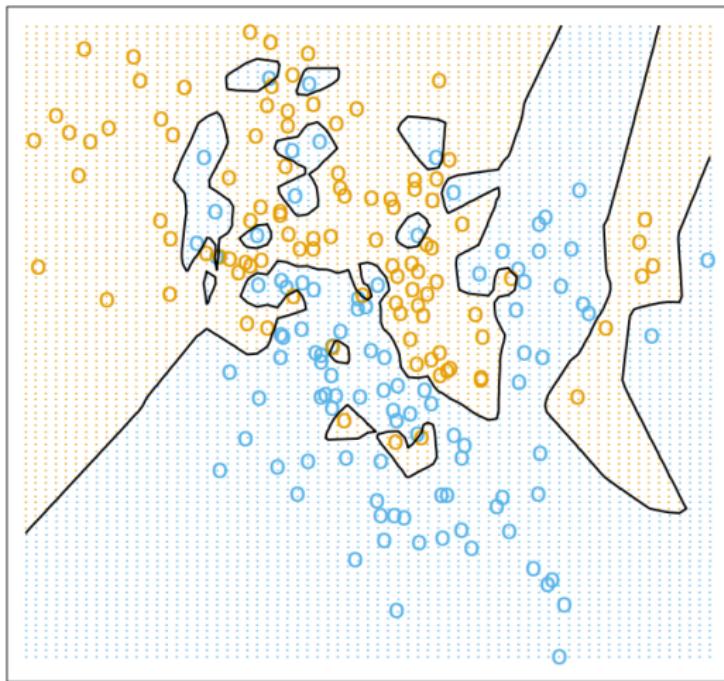
## Algorithm (kNN):

1. Find k examples  $\{\mathbf{x}^{(i)}, t^{(i)}\}$  closest to the test instance  $\mathbf{x}$
2. Classification output is majority class

$$y = \arg \max_{t^{(z)}} \sum_{r=1}^k \delta(t^{(z)}, t^{(r)})$$

# K-Nearest neighbors

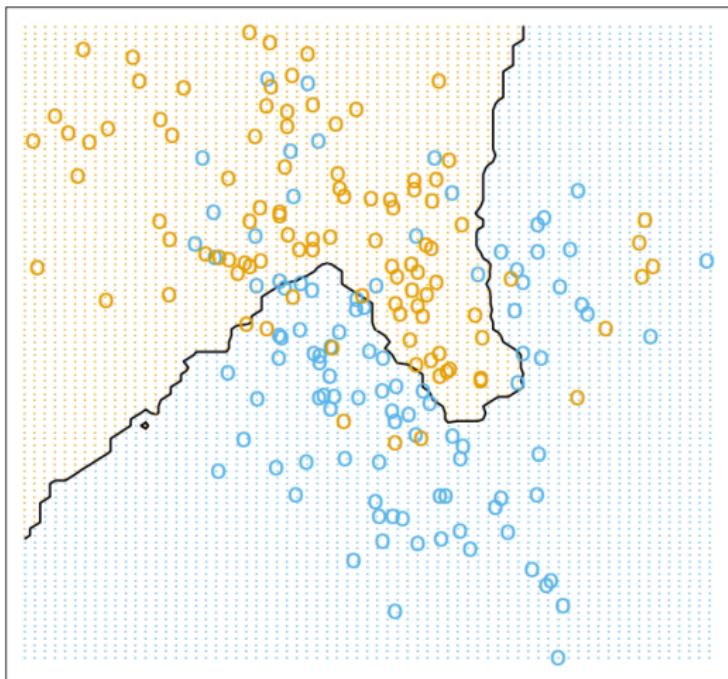
$k=1$



[Image credit: "The Elements of Statistical Learning"]

# K-Nearest neighbors

$k=15$



[Image credit: "The Elements of Statistical Learning"]

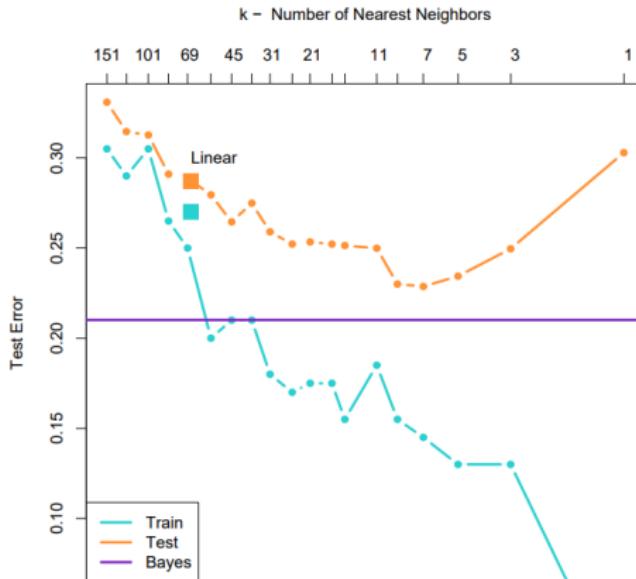
# k-Nearest Neighbors

Tradeoffs in choosing  $k$ ?

- Small  $k$ 
  - ▶ Good at capturing fine-grained patterns
  - ▶ May **overfit**, i.e. be sensitive to random idiosyncrasies in the training data
- Large  $k$ 
  - ▶ Makes stable predictions by averaging over lots of examples
  - ▶ May **underfit**, i.e. fail to capture important regularities
- Rule of thumb:  $k < \sqrt{n}$ , where  $n$  is the number of training examples

# K-Nearest neighbors

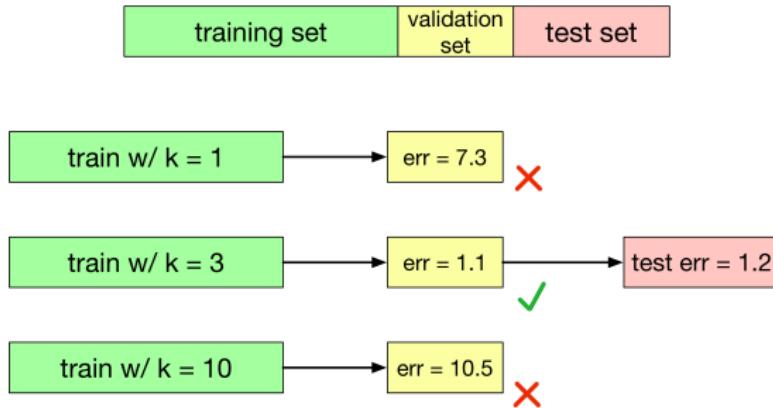
- We would like our algorithm to **generalize** to data it hasn't before.
- We can measure the **generalization error** (error rate on new examples) using a **test set**.



[Image credit: "The Elements of Statistical Learning"]

# Validation and Test Sets

- $k$  is an example of a **hyperparameter**, something we can't fit as part of the learning algorithm itself
- We can tune hyperparameters using a **validation set**:



- The test set is used only at the very end, to measure the generalization performance of the final configuration.

# Consistency

- Is KNN **consistent**? I.e., given enough data, will it give the “right” answer?
- To analyze this, suppose the inputs  $\mathbf{x}$  and targets  $t$  are random variables drawn **independently and identically distributed (i.i.d.)** from a **data generating distribution** with density  $p(\mathbf{x}, t)$ .
- The **Bayes optimal classifier** is the function  $f(\mathbf{x})$  which minimizes the misclassification rate, i.e.

$$f(\mathbf{x}) = y_* = \arg \min_y \Pr(y \neq t \mid \mathbf{x}) = \arg \max_y \Pr(y = t \mid \mathbf{x}).$$

Its error rate is called the **Bayes error**.

- Question: how close does KNN get to the Bayes error in the limit of infinite data?

# Consistency

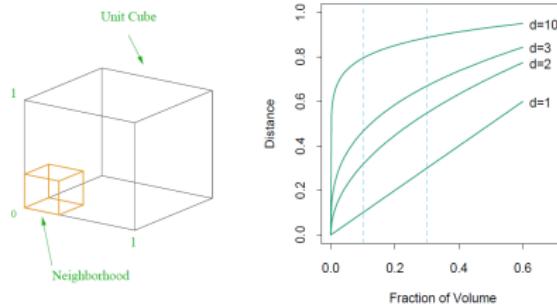
- Assume  $p(\mathbf{x}, t)$  is smooth as a function of  $\mathbf{x}$ .
- Main idea: suppose  $N$  (the number of training examples) is very large, and consider a **query point**  $\mathbf{x}_q$  which we'd like to classify.
  - ▶ By smoothness,  $p(t | \mathbf{x})$  is approximately constant for nearby  $\mathbf{x}$ .
  - ▶ Hence, the labels of the neighbors can be seen as independent random variables with PMF  $p(t | \mathbf{x}_q)$ .

# Consistency

- First consider  $k = 1, N \rightarrow \infty$ .
  - ▶  $y$  (the nearest neighbour prediction) and  $t$  (the true label at  $\mathbf{x}_q$ ) are (approximately) independent random variables with PMF  $p(t | \mathbf{x}_q)$ .
  - ▶ Apply the **Union Bound**:
$$\Pr(t \neq y | \mathbf{x}_q) \leq \Pr(t \neq y_* | \mathbf{x}_q) + \Pr(y_* \neq y | \mathbf{x}_q) = 2\Pr(t \neq y_* | \mathbf{x}_q).$$
- ▶ I.e., the asymptotic error of 1-NN is at most twice the Bayes error.
- Now consider  $k, N \rightarrow \infty$  and  $k/N \rightarrow 0$ .
  - ▶ The counts of neighbors' labels (approximately) follow a multinomial distribution with  $k$  trials.
  - ▶ For large  $k$ , the argmax will agree with the Bayes classifier with high probability. (E.g., apply the Central Limit Theorem.)
  - ▶ Hence, the KNN approaches the Bayes error, i.e. KNN is **Bayes consistent**.
- Bayes consistency is a very special property, and holds for hardly any of the algorithms covered in this course.

# Pitfalls: The Curse of Dimensionality

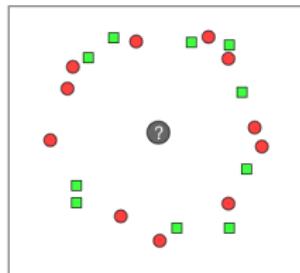
- Consistency is great, but it might take a very large amount of data to get close to the Bayes error.
  - ▶ Especially in high dimensions! KNN suffers from the **Curse of Dimensionality**.
- How large does  $N$  need to be to guarantee we have an  $\epsilon$ -neighbour?
- The volume of a single ball of radius  $\epsilon$  is  $\mathcal{O}(\epsilon^d)$
- The total volume of  $[0, 1]^d$  is 1.
- Therefore  $\mathcal{O}\left((\frac{1}{\epsilon})^d\right)$  balls are needed to cover the volume.



[Image credit: "The Elements of Statistical Learning"]

# Pitfalls: The Curse of Dimensionality

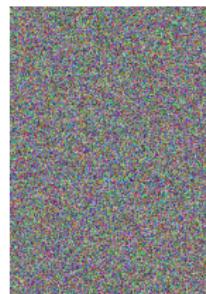
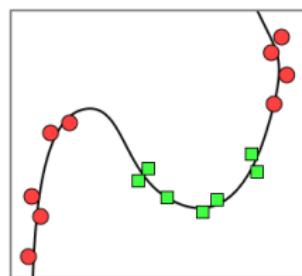
- Another perspective on the Curse of Dimensionality: in high dimensions, “most” points are approximately the same distance. (Homework question coming up...)



- This is just one example of how 2-D visualizations of high-dimensional spaces can be extremely misleading!

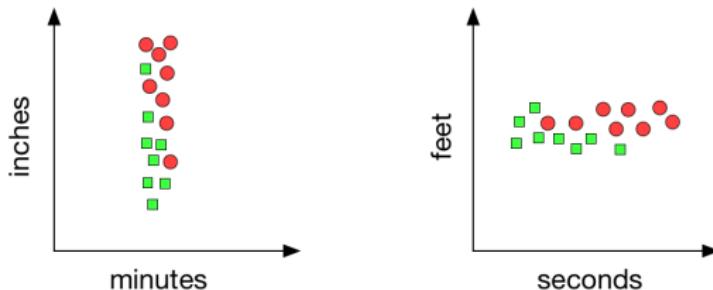
# Pitfalls: The Curse of Dimensionality

- Saving grace: some datasets may have low **intrinsic dimension**, i.e. lie on or near a low-dimensional manifold.
- E.g., natural images have a lot fewer degrees of freedom than the number of pixels in the image.
- The distance to the neighbors depends on the intrinsic dimension, not the dimension of the input space. Hence, KNN can still work in high dimensions, as long as the data are intrinsically low-dimensional.



# Pitfalls: Normalization

- Nearest neighbors can be sensitive to the ranges of different features.
- Often, the units are arbitrary:



- Simple fix: **normalize** each dimension to be zero mean and unit variance. I.e., compute the mean  $\mu_j$  and standard deviation  $\sigma_j$ , and take

$$\tilde{x}_j = \frac{x_j - \mu_j}{\sigma_j}$$

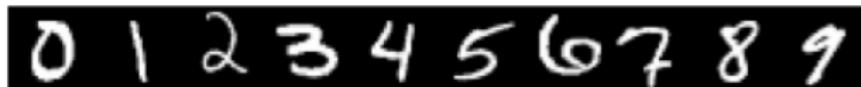
- Caution: depending on the problem, the scale might be important! (Can you think of an example?)

# Pitfalls: Computational Cost

- Number of computations at **training time**: 0
- Number of computations at **test time**, per query (naïve algorithm)
  - ▶ Calculate  $D$ -dimensional Euclidean distances with  $N$  data points:  $\mathcal{O}(ND)$
  - ▶ Sort the distances:  $\mathcal{O}(N \log N)$
- This must be done for *each* query, which is very expensive by the standards of a learning algorithm!
- Need to store the entire dataset in memory!
- Tons of work has gone into algorithms and data structures for efficient nearest neighbors with high dimensions and/or large datasets.

# Example: Digit Classification

- Decent performance when lots of data

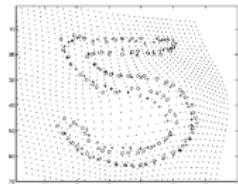
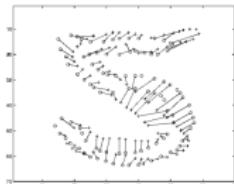
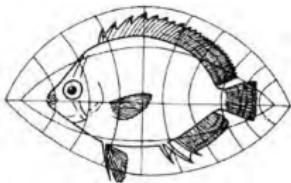
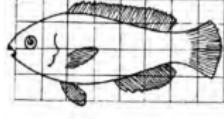


- Yann LeCunn – MNIST Digit Recognition
  - Handwritten digits
  - 28x28 pixel images:  $d = 784$
  - 60,000 training samples
  - 10,000 test samples
- Nearest neighbour is competitive

	Test Error Rate (%)
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
K-NN, Tangent Distance, 16x16	1.1
K-NN, shape context matching	0.67
1000 RBF + linear classifier	3.6
SVM deg 4 polynomial	1.1
2-layer NN, 300 hidden units	4.7
2-layer NN, 300 HU, [deskewing]	1.6
LeNet-5, [distortions]	0.8
Boosted LeNet-4, [distortions]	0.7

# Example: Digit Classification

- KNN can perform a lot better with a good similarity measure.
- Example: shape contexts for object recognition. In order to achieve invariance to image transformations, they tried to warp one image to match the other image.
  - ▶ Distance measure: average distance between corresponding points on *warped* images
- Achieved 0.63% error on MNIST, compared with 3% for Euclidean KNN.
- Competitive with conv nets at the time, but required careful engineering.



[Belongie, Malik, and Puzicha, 2002. Shape matching and object recognition using shape contexts.]

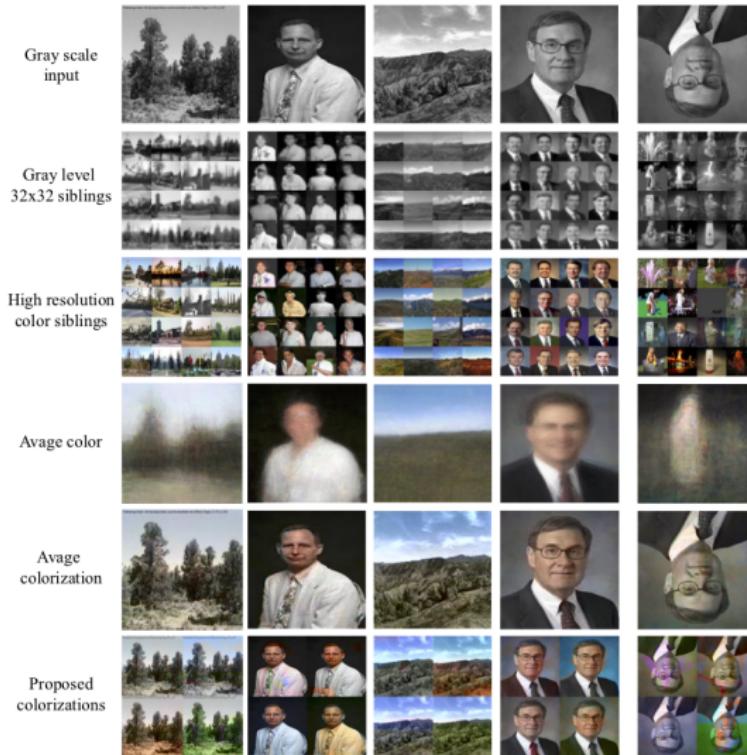
# Example: 80 Million Tiny Images

- 80 Million Tiny Images was the first extremely large image dataset. It consisted of color images scaled down to  $32 \times 32$ .
- With a large dataset, you can find much better semantic matches, and KNN can do some surprising things.
- Note: this required a carefully chosen similarity metric.



[Torralba, Fergus, and Freeman, 2007. 80 Million Tiny Images.]

# Example: 80 Million Tiny Images



[Torralba, Fergus, and Freeman, 2007. 80 Million Tiny Images.]

Questions?

?