

UCLA



CS145 Discussion

Week 5

Junheng, Shengming, Yunsheng
11/2/2018



-
- Announcement
 - K-NN
 - Similarity Metrics
 - ROC
 - K-Means
 - Homework 3
 - Course Project Midterm Report



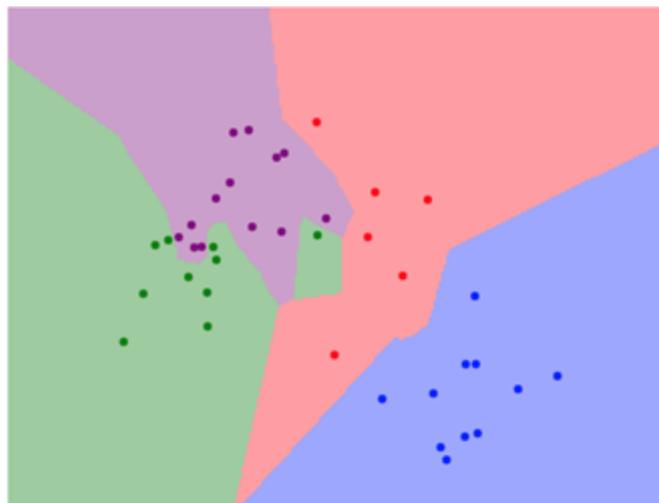
- Homework 3 out
 - Deadline: 11/09 Friday 11:59 pm
 - KNN(30%) and Neural Network(70%)
 - Environment Requirement: Jupyter + Python 3
- Course Project Midterm Report (about to) out
 - Deadline: 11/12 Monday 11:59 pm
 - According to the guideline file
- Midterm date out
 - 11/14 Wednesday 12:00-1:50 pm(in class)
 - Remember to carry: one-page reference paper(letter size), simple calculator



KNN



- Demo: <http://vision.stanford.edu/teaching/cs231n-demos/knn/>



Metric

L1 L2

Num classes

2 3 4 5

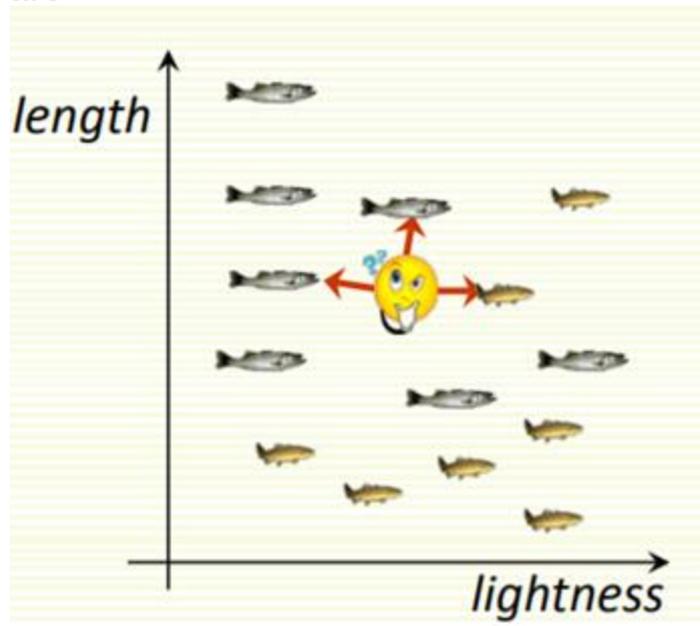
Num Neighbors (K)

1 2 3 4 5 6 7

Num points

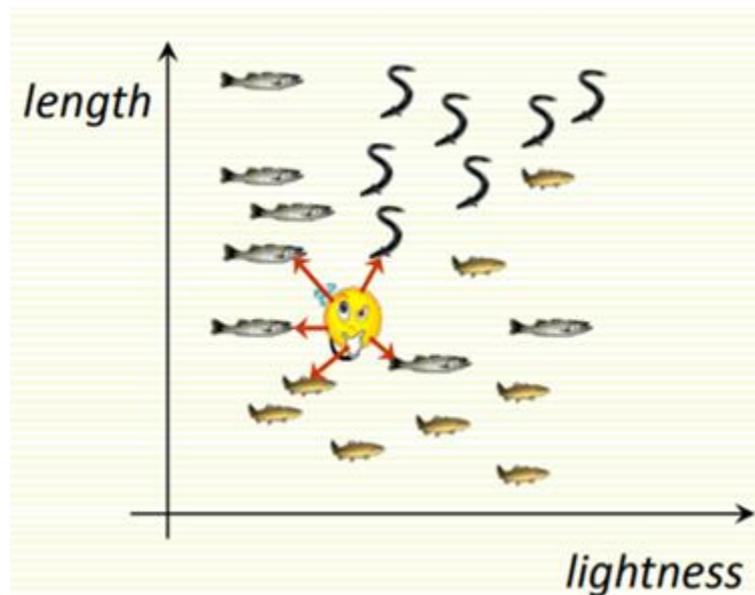
20 30 40 50 60

- Classify an unknown example with the most common class among K nearest examples
 - “Tell me who your neighbors are, and I’ll tell you who you are”
- Example
 - $K = 3$
 - 2 sea bass, 1 salmon
 - Classify as sea bass



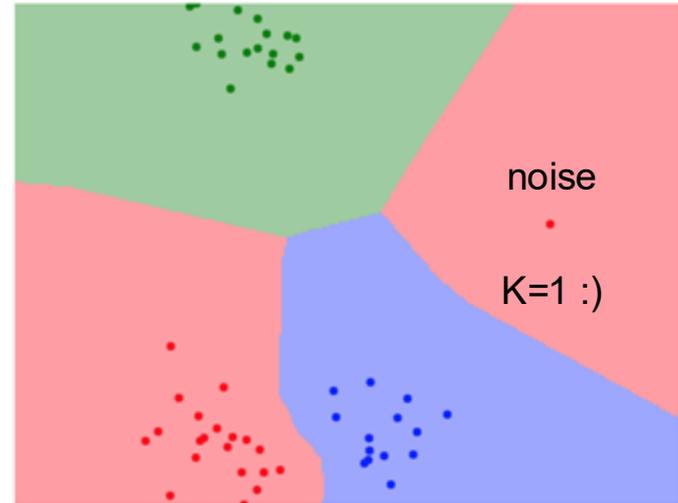


- Easy to implement for multiple classes
- Example for $K = 5$
 - 3 fish species: salmon, sea bass, eel
 - 3 sea bass, 1 eel, 1 salmon \rightarrow classify as sea bass



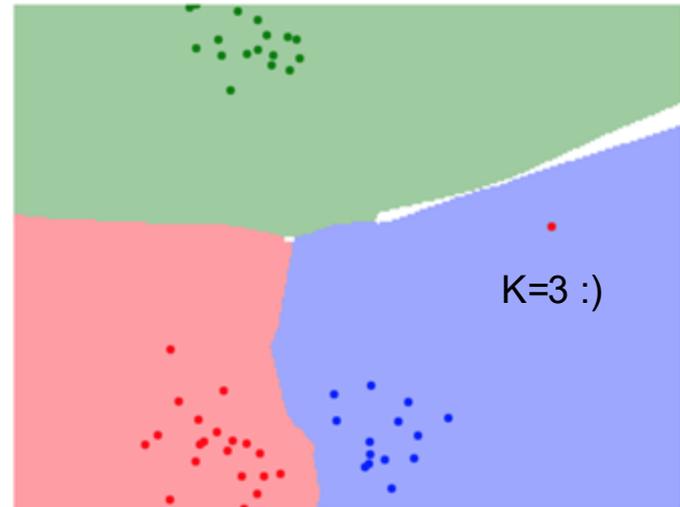


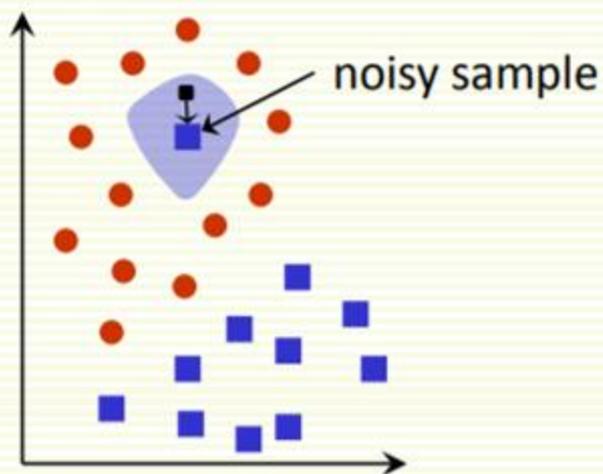
- In theory, if infinite number of samples available, the larger k , the better classification result you'll get.
- Caveat: all K neighbors have to be close
 - Possible when infinite # samples available
 - Impossible in practice since # samples is finite
- Should we “tune” K on training data?
 - Overfitting
- $K = 1 \rightarrow$ sensitive to “noise” (e.g. see right)



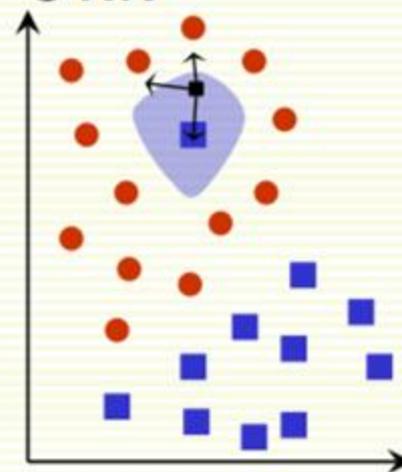


- In theory, if infinite number of samples available, the larger k , the better classification result you'll get.
- Caveat: all K neighbors have to be close
 - Possible when infinite # samples available
 - Impossible in practice since # samples is finite
- Should we “tune” K on training data?
 - Overfitting
- $K = 1 \rightarrow$ sensitive to “noise” (e.g. see right)



**1 NN**

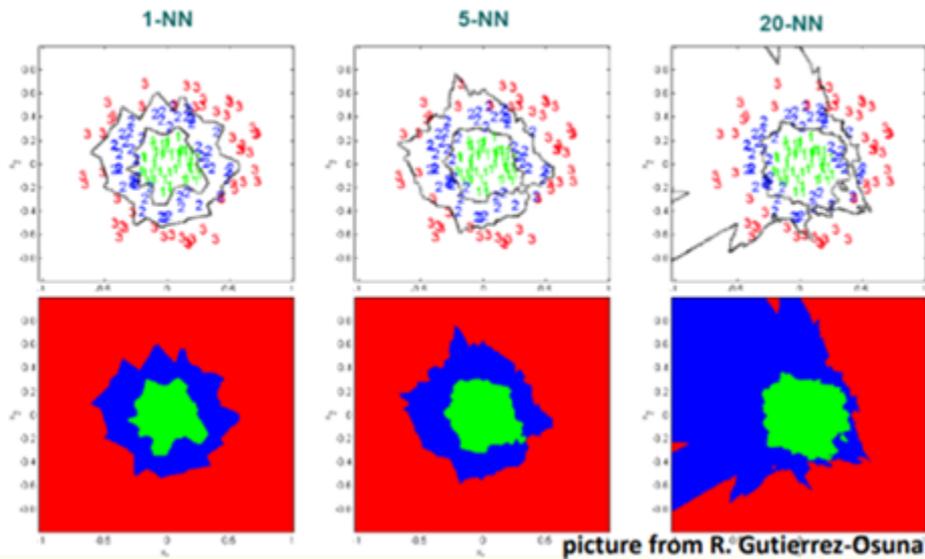
every example in the blue shaded area will be misclassified as the **blue** class

3 NN

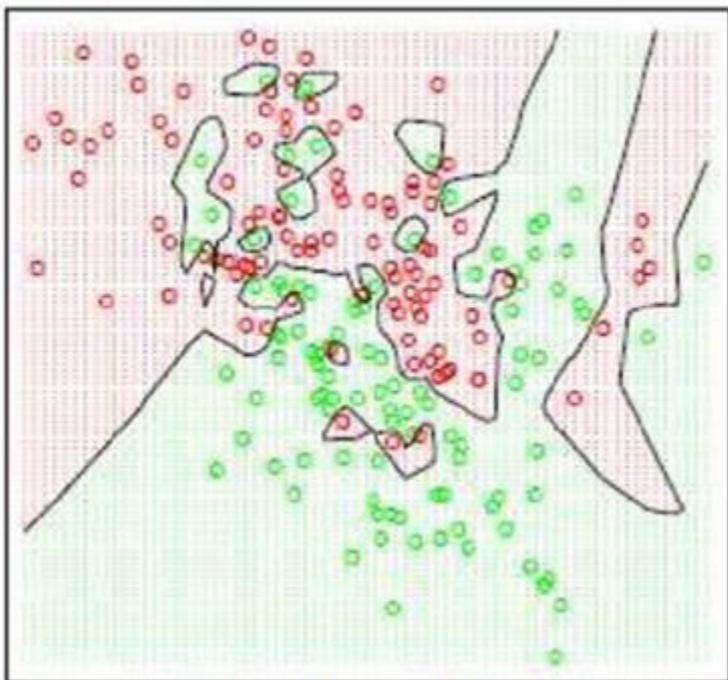
every example in the blue shaded area will be classified correctly as the **red** class



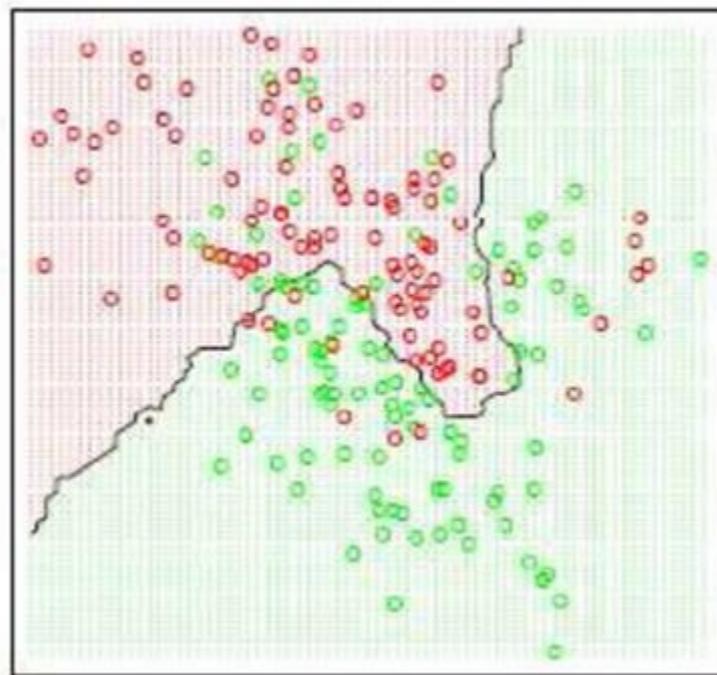
- Larger K gives smoother boundaries, better for generalization
 - Only if locality is preserved
 - K too large → looking at samples too far away that are not from the same class
- Can choose K through cross-validation



K=1



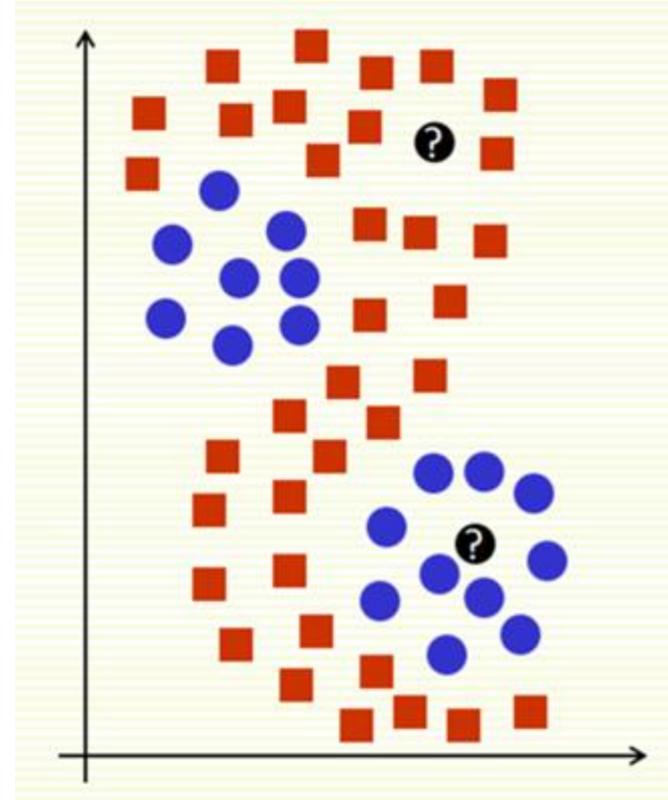
K=15



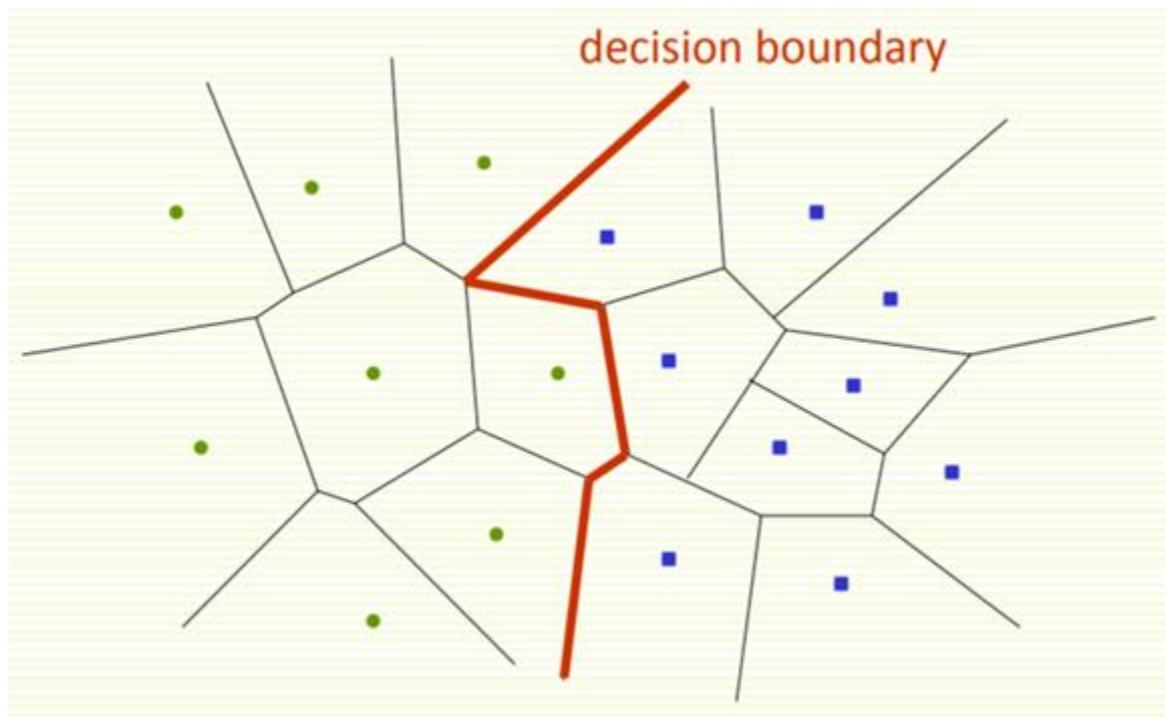
Figures from Hastie, Tibshirani and Friedman (Elements of Statistical Learning)



- Many classification models would not work for this 2-class classification problem
 - Linear-R
 - Logistic-R
 - Decision-Tree
 - SVM
- Nearest neighbors will do reasonably well

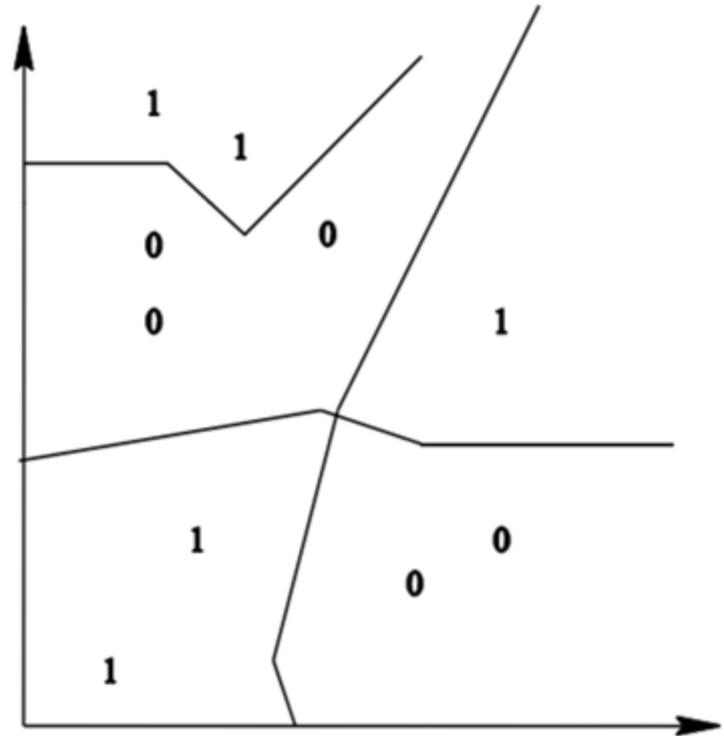


- Voronoi diagram is useful for visualization





- Decision boundaries are formed by a subset of the Voronoi Diagram of the training data
- Each line segment is equidistant between two points of opposite class
- The more examples that are stored, the more fragmented and complex the decision boundaries can be.





- We use Euclidean Distance to find the nearest neighbor:

$$D(a, b) = \sqrt{\sum_k (a_k - b_k)^2}$$

- Euclidean distance treats each feature as equally important
- Sometimes, some features (or dimensions) may be much more discriminative than other features



- Feature 1 gives the correct class: 1 or 2
- Feature 2 gives irrelevant number from 100 to 200
- Dataset: [1, 150], [2, 110]
- Classify [1, 100]

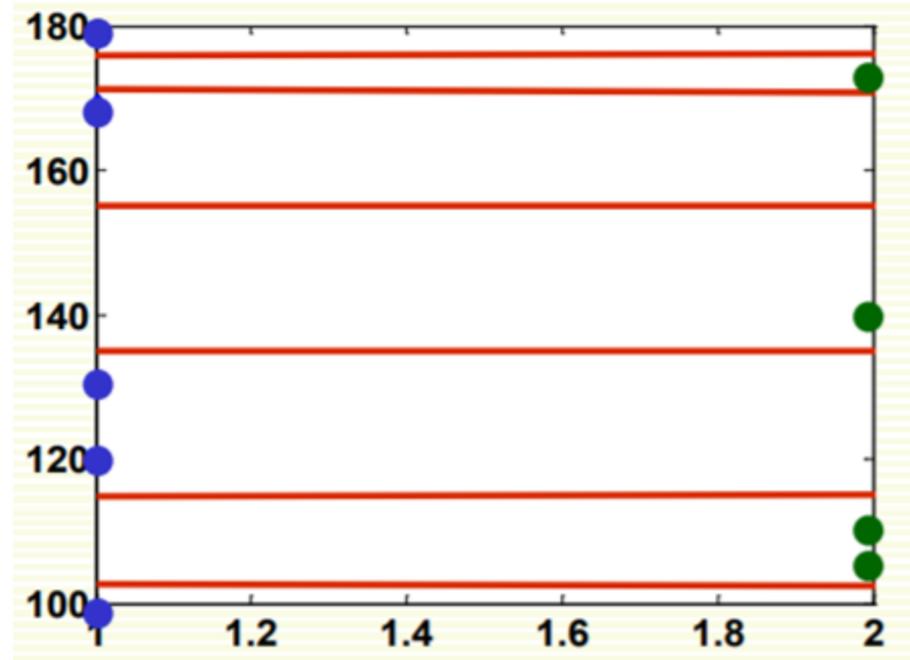
$$D\left(\begin{bmatrix} 1 \\ 100 \end{bmatrix}, \begin{bmatrix} 1 \\ 150 \end{bmatrix}\right) = \sqrt{(1-1)^2 + (100-150)^2} = 50$$

$$D\left(\begin{bmatrix} 1 \\ 100 \end{bmatrix}, \begin{bmatrix} 2 \\ 110 \end{bmatrix}\right) = \sqrt{(1-2)^2 + (100-110)^2} = 10.5$$

- Use Euclidean distance can result in wrong classification
- Dense Example can help solve this problem



- Decision boundary is in red, and is really wrong because:
 - Feature 1 is discriminative, but its scale is small
 - Feature gives no class information but its scale is large, which dominates distance calculation





- Normalize features that makes them be in the same scale
- Different normalization approaches may reflect the result
- Linear scale the feature in range [0,1]:

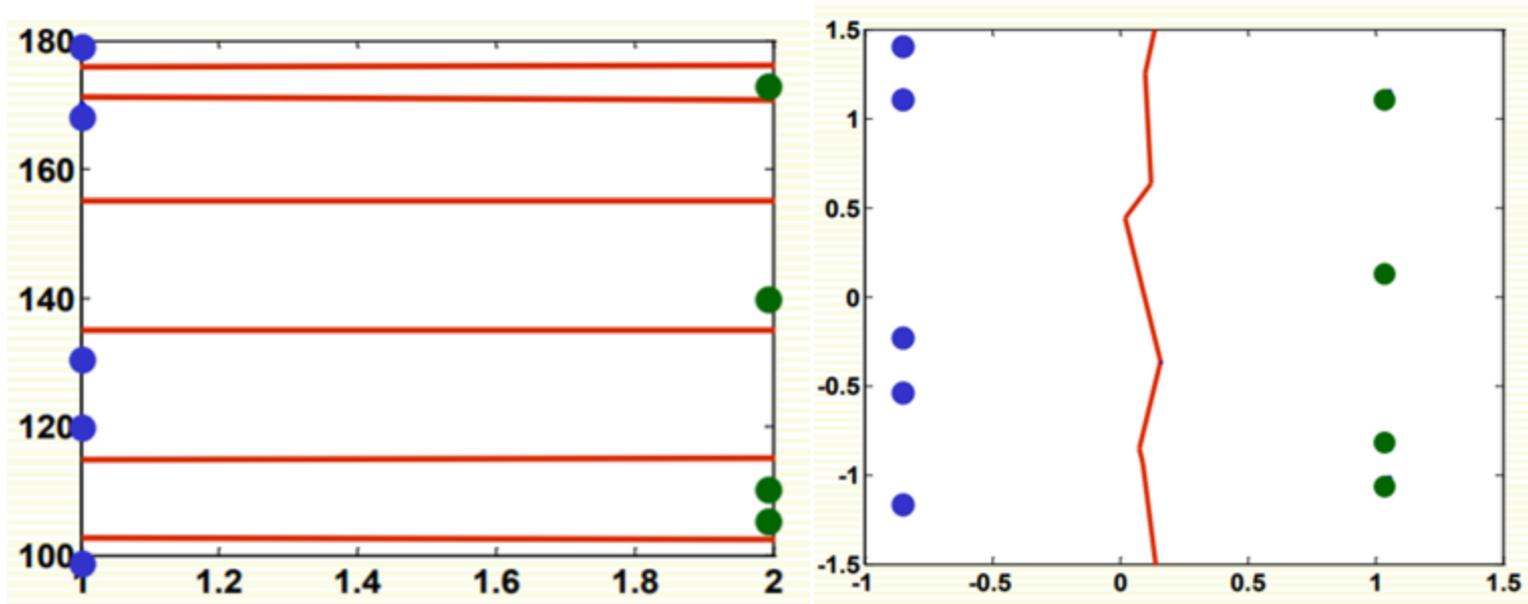
$$f_{new} = \frac{f_{old} - f_{old}^{\min}}{f_{old}^{\max} - f_{old}^{\min}}$$

- Linear scale to 0 mean variance 1(Z-score):

$$f_{new} = \frac{f_{old} - \mu}{\sigma}$$



- Result comparison non-normalized vs normalized

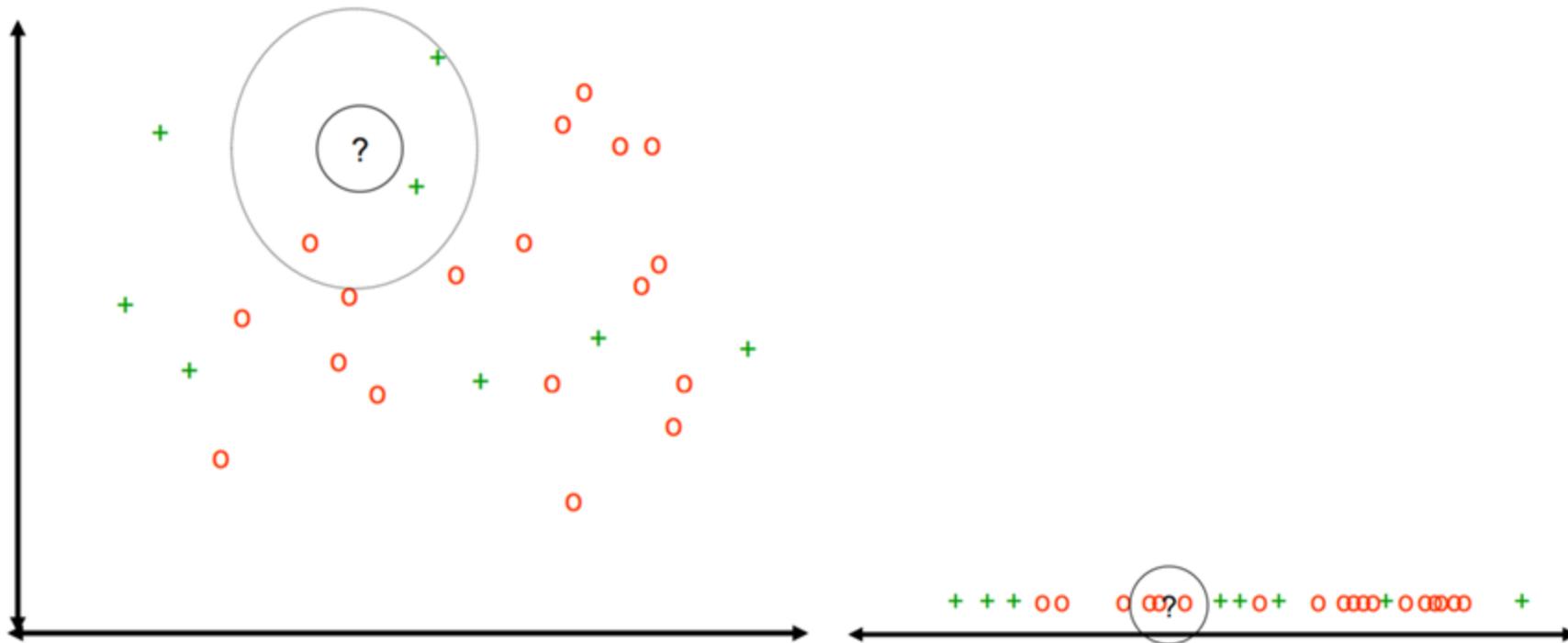




- Feature normalization does not help in high dimensional spaces if most features are irrelevant

$$D(a, b) = \sqrt{\sum_k (a_k - b_k)^2} = \sqrt{\underbrace{\sum_i (a_i - b_i)^2}_{\text{Discriminative features}} + \underbrace{\sum_j (a_j - b_j)^2}_{\text{Noisy features}}}$$

- If the number of useful feature is smaller than the number of noisy features, Euclidean distance is dominated by noise.





- Scale each feature by its importance for classification

$$D(a, b) = \sqrt{\sum_k w_k (a_k - b_k)^2}$$

- Use prior/domain knowledge to set the weight w
- Use cross-validation to learn the weight w



- Suppose n examples with dimension d
- For each point to be classified:
 - $O(d)$ to compute distance to one example
 - $O(nd)$ to compute distances to all examples
 - $O(nk)$ time to find k closest examples
 - Total time: $O(nk + nd)$
- Very expensive for a large number of queries



- Reduce the dimensionality of the data:
 - Find a projection from high dimensional space to a lower dimensional space so that the distance between samples are approximately the same
 - Use Principal Component Analysis(PCA)
- Apply smart data structure, e.v. k-d trees



- Advantages:
 - Can be applied to the data from any distribution
 - The decision boundary is not necessarily to be linear
 - Simple and Intuitive
 - Good Classification with large number of samples
- Disadvantages:
 - Choosing k may be tricky
 - Test stage is computationally expensive
 - No training stage, time-consuming test stage
 - Usually we can afford long training step but fast testing speed
 - Need large number of examples for accuracy



-
- http://www.csd.uwo.ca/courses/CS9840a/Lecture2_knn.pdf
 - <http://classes.engr.oregonstate.edu/eecs/spring2012/cs534/notes/knn.pdf>
 - <http://people.csail.mit.edu/dsontag/courses/ml12/slides/lecture10.pdf>

Similarity Metrics



- Dissimilarity

- $$\begin{bmatrix} 3 & 5 \\ 6 & 9 \\ 11 & 21 \end{bmatrix}$$

- 3 2-d input, x_1, x_2, x_3
- The dissimilarity matrix is a 3*3 lower-triangular matrix:

$$\begin{bmatrix} 0 & 0 & 0 \\ d(2,1) & 0 & 0 \\ d(3,1) & d(3,2) & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ \sqrt{(3-6)^2 + (5-9)^2} & 0 & 0 \\ \sqrt{(3-11)^2 + (5-21)^2} & \sqrt{(11-6)^2 + (21-9)^2} & 0 \end{bmatrix}$$



- Student 1: likes Jazz, eats pizza, roots for the cubs, wears socks
- Student 2: likes Rock, eats pizza, roots for the cubs, goes barefoot
- $d(\text{Student 1}, \text{Student 2})$:
 - m : # of matches $\rightarrow 2$
 - p : total # of variables $\rightarrow 4$
 - $d(\text{Student 1}, \text{Student 2}) = (4-2)/4 = 0.5$



- Symmetric binary attributes:
 - Gender
- Asymmetric attributes:
 - Preference, Character, etc.
- Can be manually defined



- Dissimilarity of Binary Attributes:
 - Define 0 and 1
 - calculate q,s,r,t,p
- Distance measure for symmetric binary variables:

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

		Object <i>j</i>		sum
		1	0	
Object <i>i</i>	1	<i>q</i>	<i>r</i>	<i>q + r</i>
	0	<i>s</i>	<i>t</i>	<i>s + t</i>
sum		<i>q + s</i>	<i>r + t</i>	<i>p</i>

- Distance measure for asymmetric binary variables:

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient (similarity measure for asymmetric binary variables):

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$



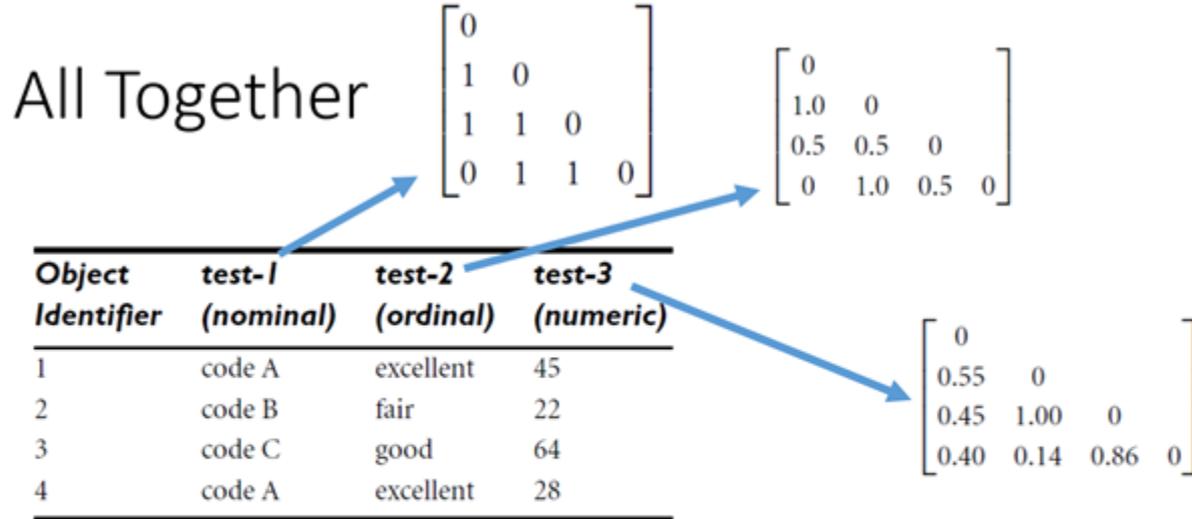
Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- $i = \text{Jack}, j = \text{Mary}$
- Define M,Y,P as 1; Define F,N as 0
- Assume symmetric attributes for all:
 - $r = 0, s = 1, q = 2$
- Assume symmetric for Gender, asymmetric for other attributes
 - For Gender: $r = 1, s = 0, q = 0, t = 0$
 - For others: $r = 0, s = 1, q = 2$

		Object j		
		1	0	sum
Object i	1	q	r	$q+r$
	0	s	t	$s+t$
sum		$q+s$	$r+t$	p



- Order is important
 - Transfer rank into value
 - Freshman, Sophomore, Junior, Senior
 - 1,2,3,4



- $d(3,1)?$
 - $\frac{1(1)+1(0.5)+1(0.45)}{3}$



- For vector data
- d1: I like to go to the store
- d2: I like the cubs, go cubs go

$$\cos(d_1, d_2) = (d_1 \cdot d_2) / (||d_1|| ||d_2||),$$

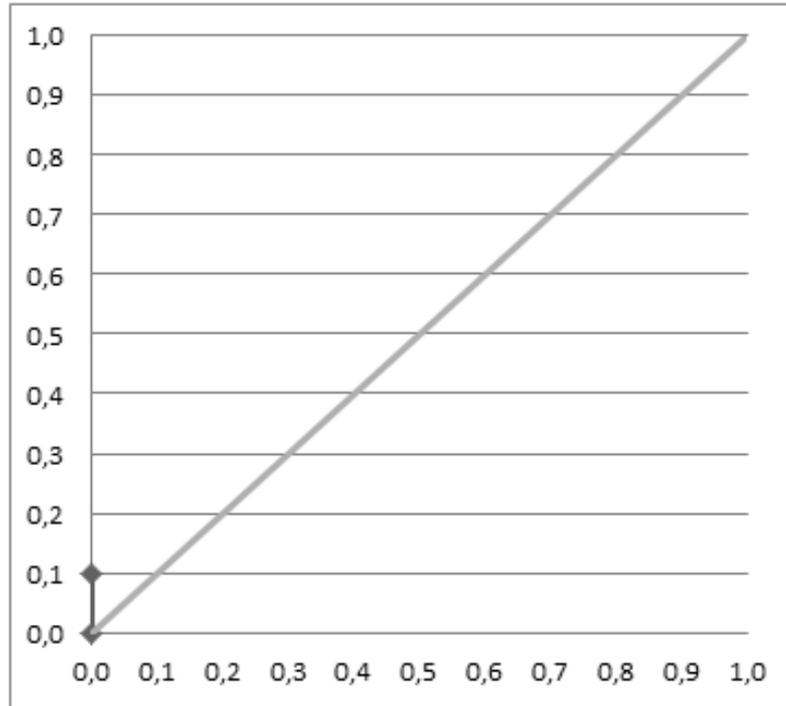
Document	I	like	to	go	the	store	cubs
d1	1	1	2	1	1	1	0
d2	1	1	0	2	1	0	2

- $\cos(d1, d2)$?

- $$\frac{1 \cdot 1 + 1 \cdot 1 + 2 \cdot 0 + 1 \cdot 2 + 1 \cdot 1 + 1 \cdot 0 + 0 \cdot 2}{\sqrt{1^2 + 1^2 + 2^2 + 1^2 + 1^2 + 1^2 + 0^2} \cdot \sqrt{1^2 + 1^2 + 0^2 + 2^2 + 1^2 + 0^2 + 2^2}}$$

ROC

#	C	Score
1	P	0,9
2	P	0,8
3	N	0,7
4	P	0,6
5	P	0,55
6	P	0,54
7	N	0,53
8	N	0,52
9	P	0,51
10	N	0,505
11	P	0,4
12	N	0,39
13	P	0,38
14	N	0,37
15	N	0,36
16	N	0,35
17	P	0,34
18	N	0,33
19	P	0,3
20	N	0,1





-
- Demo 1: <http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>
 - Demo 2: <https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>



- Due Date: Next Friday
- Be prepared (Refer to the doc for hw3)
 - Install Jupyter with Python 3.x(3.6 is preferred)
 - Download the cifar-10 dataset via `get_datasets.sh`
 - For windows user, there might be problems downloading the dataset
- Live demonstration for installation
- Hints on k-fold cross-validation and matrix-level operations involved in NN



- <https://docs.google.com/document/d/1xLeBU-n8nuMT6zhLyLL1NIGuu25SJU1OWI9eybdgjXg/edit?usp=sharing>