

# Classification

- Can a computer learn to recognize objects?
- Shown 10,000 flowers, can a computer “understand” flowers? Can it say if the new photograph shown is a flower?



Iris Setosa



Iris Versicolor



Iris Virginica

# Let's try our brain's algorithm!



Iris Versicolor

Iris Setosa

Iris Virginica

???

# What is Similarity?

The quality or state of being similar; likeness; resemblance; as, a similarity of features.

Webster's Dictionary

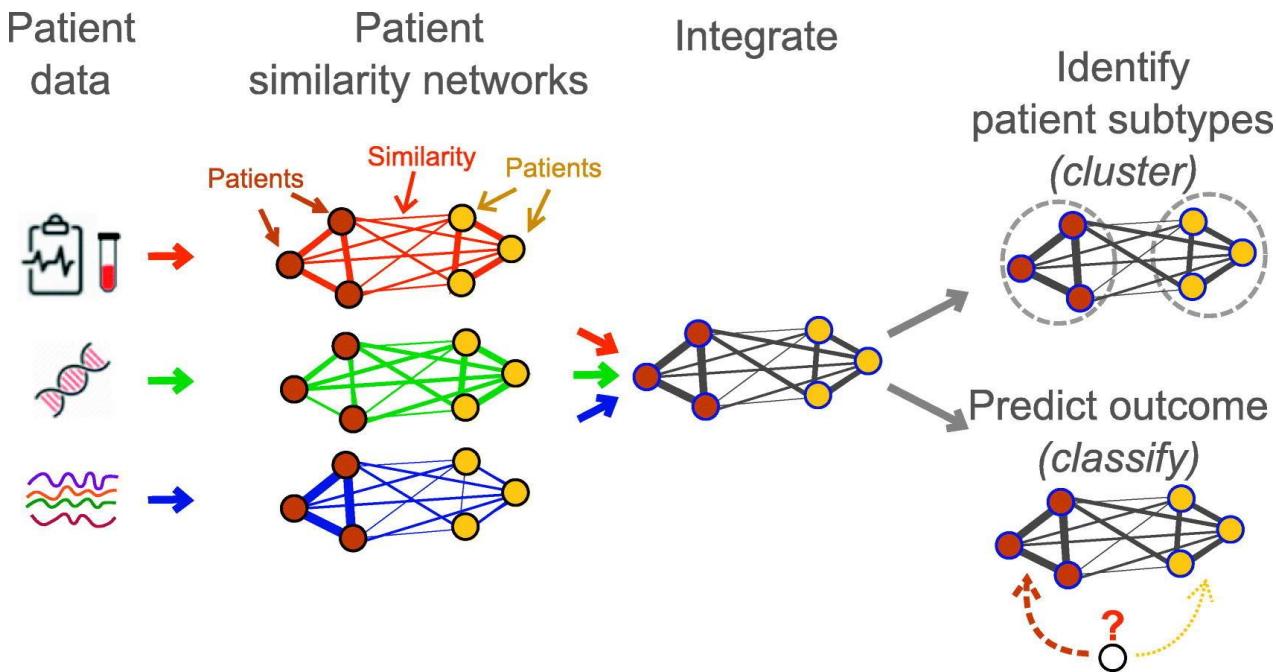


For example, for someone who is writing a software for healthcare industry, They may have to deal with the questions of “how similar are two patients.”

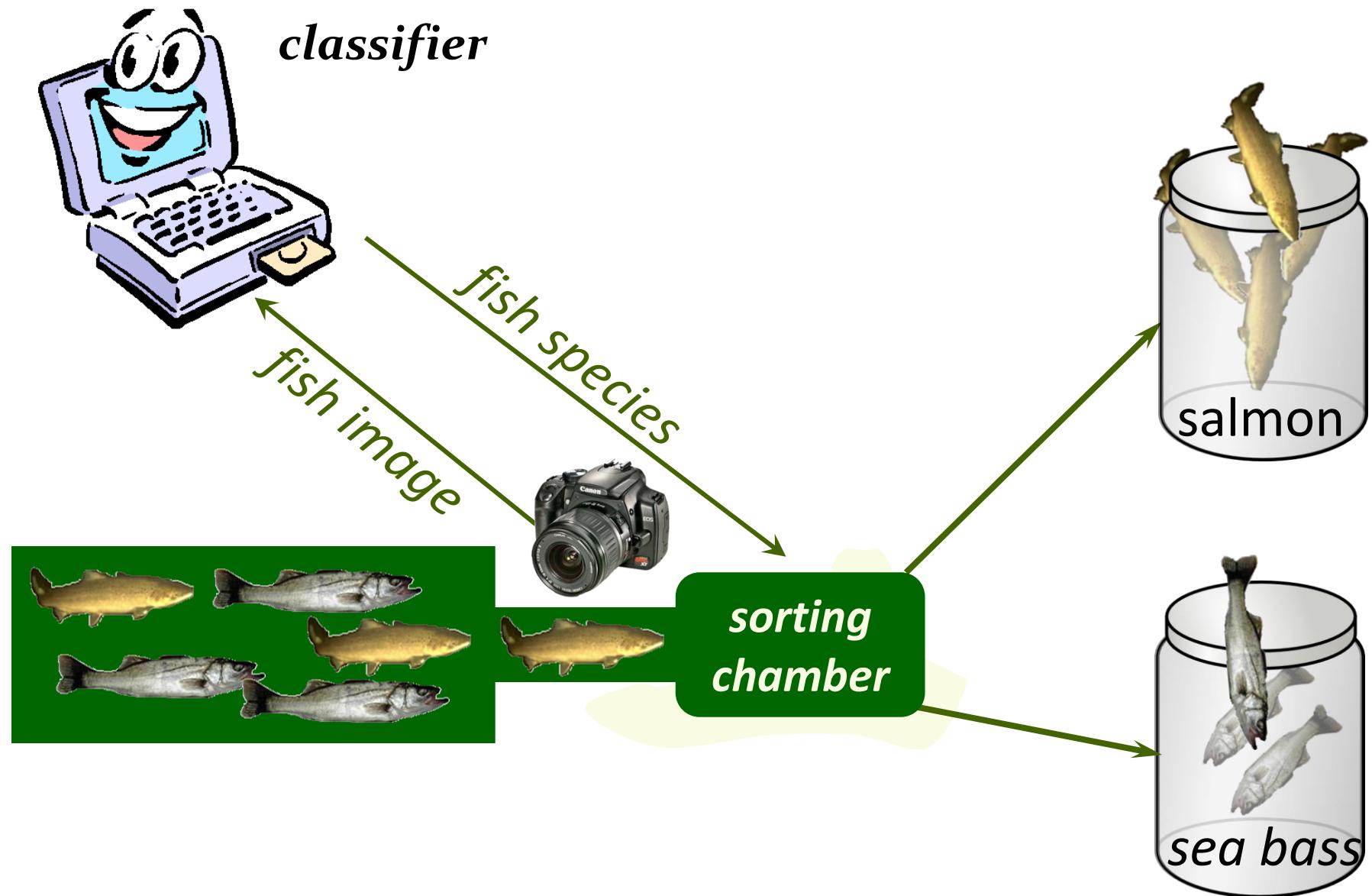
It depends on what you are comparing the two objects for.

Whole lot of research and Ph.D. thesis, just on the concept of similarity.

1. Patient Similarity Networks for Precision Medicine
2. Patient Similarity: Emerging Concepts in Systems and Precision Medicine
3. Machine learning of patient similarity: A case study on predicting survival in cancer patient after locoregional chemotherapy

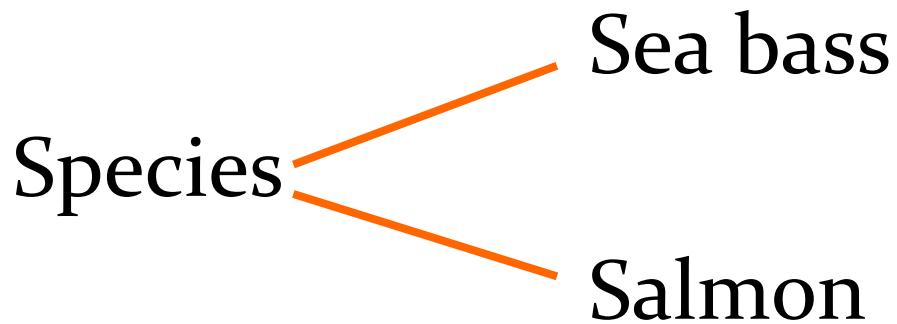


# Fish Sorting: For Packaging



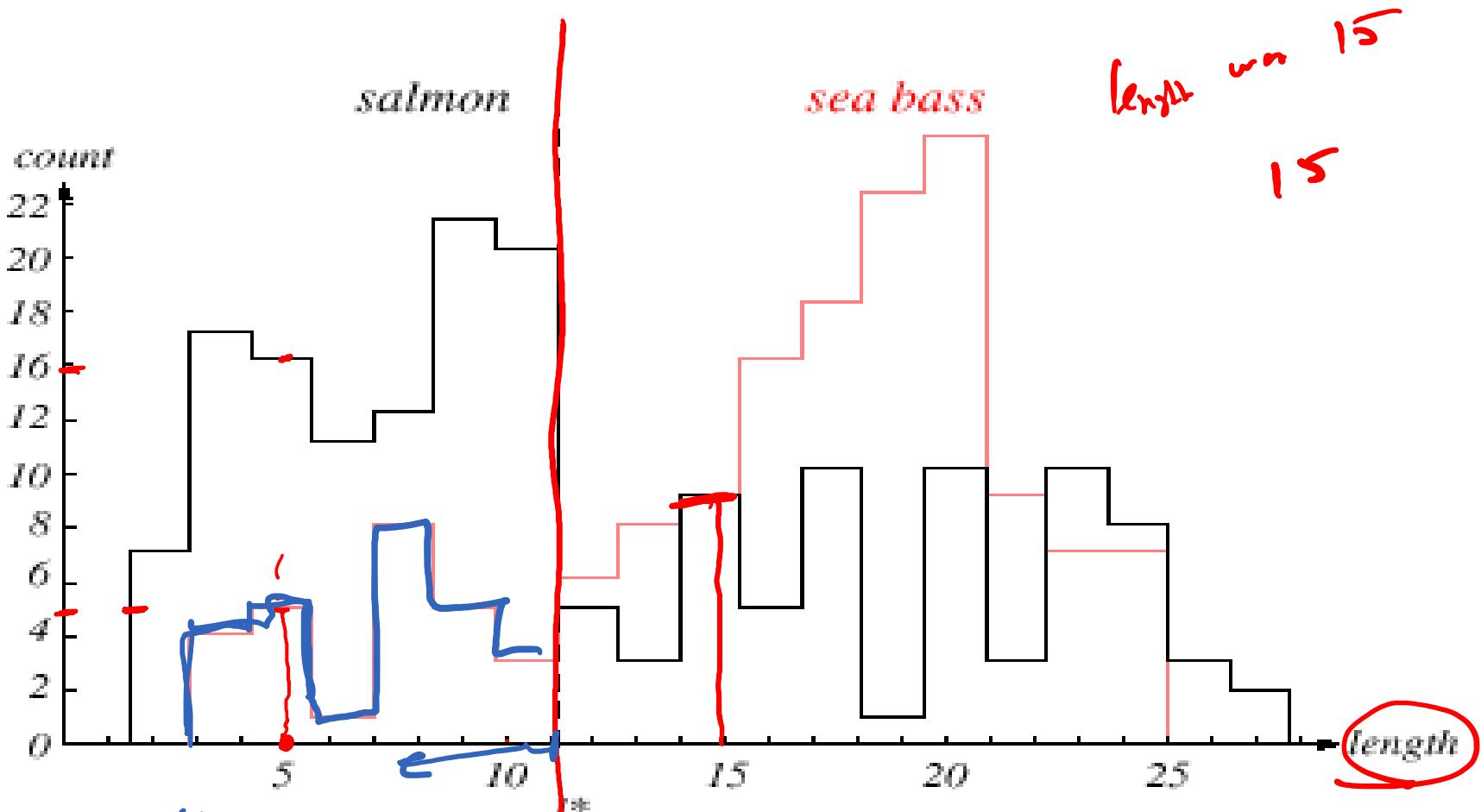
# An Example

- “Sorting incoming Fish on a conveyor according to species using optical sensing”



- Problem Analysis
  - Set up a camera and take some sample images to extract features
    - Length
    - Lightness
    - Width
    - Number and shape of fins
    - Position of the mouth, etc...
  - This is the set of all suggested features to explore for use in our classifier!

- Classification
  - Select the length of the fish as a possible feature for discrimination

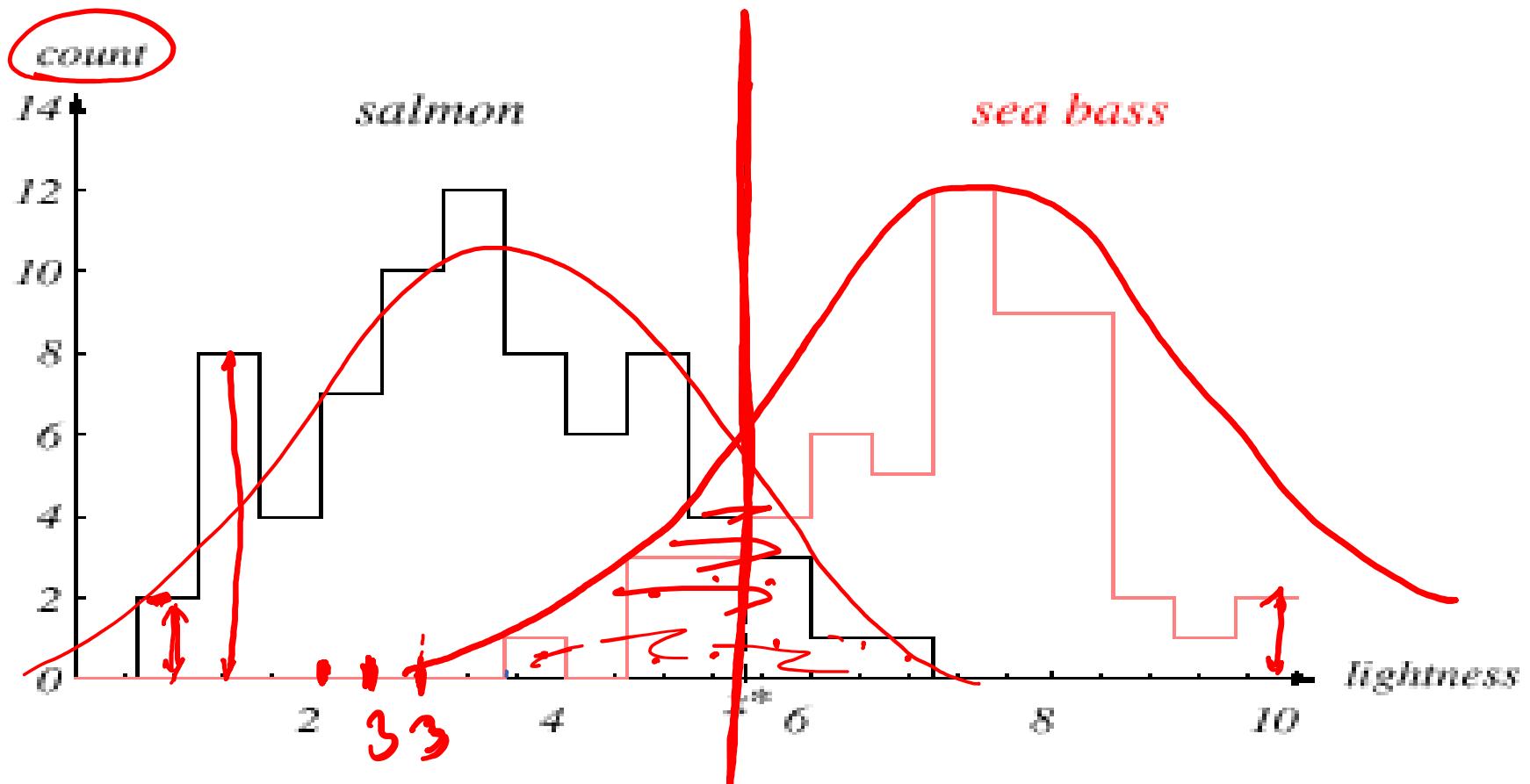


The **length** is a poor feature alone!

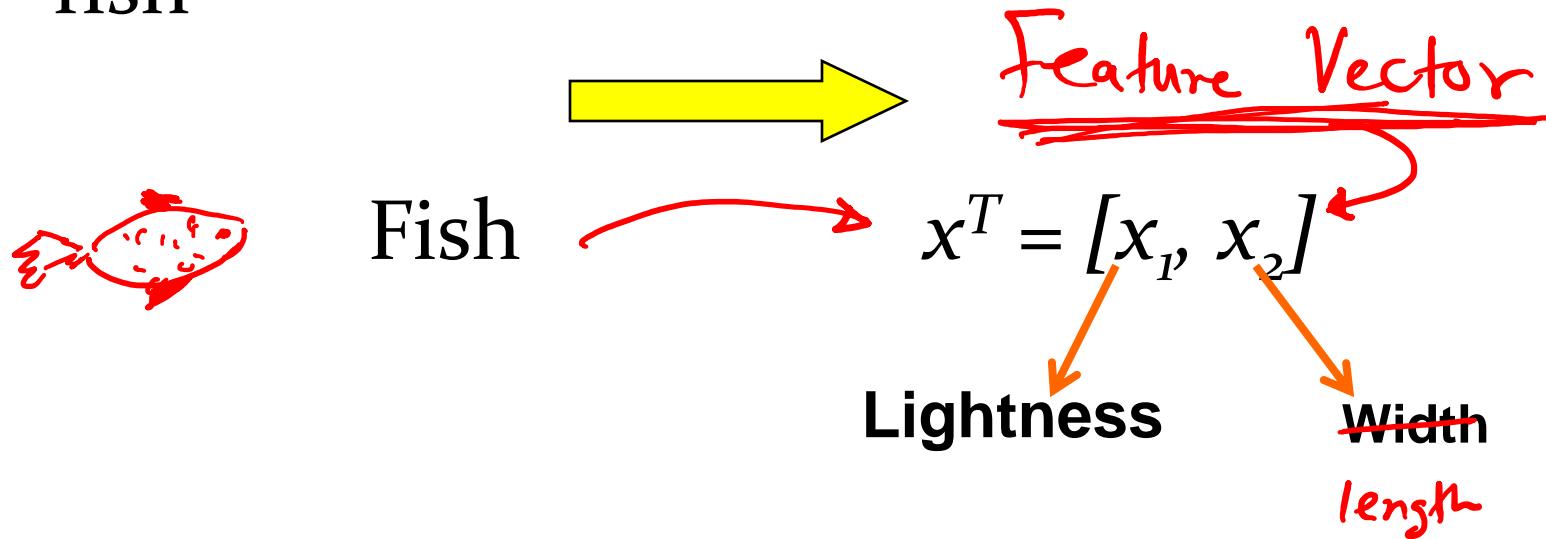
Select the lightness as a possible feature.

6.

Brightness



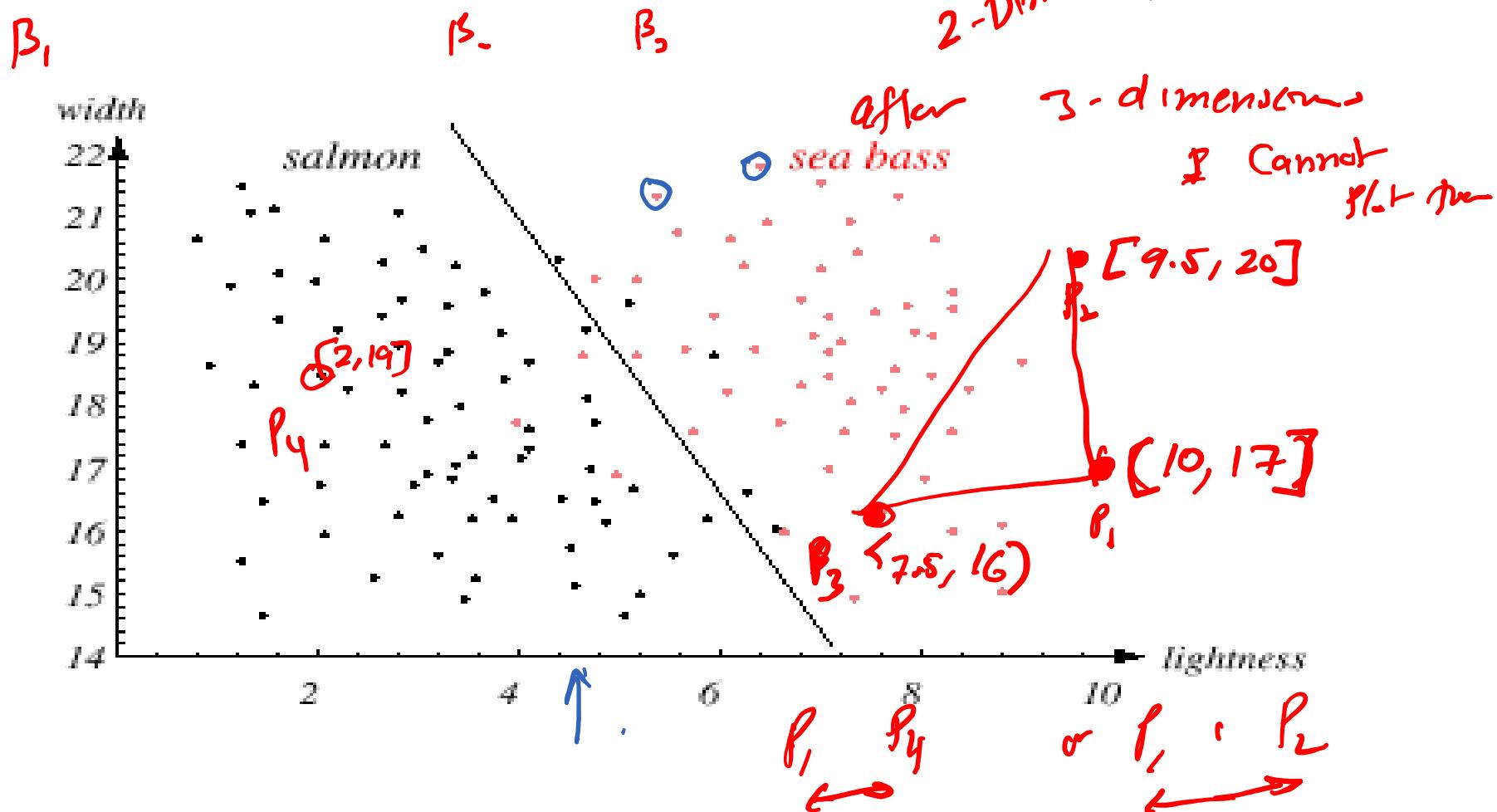
- Adopt the lightness and add the width of the fish



## 2-dimensional Vectors in 2-d space

C

- Plotting Salmon and Seabass based on two-dimensional feature vector.

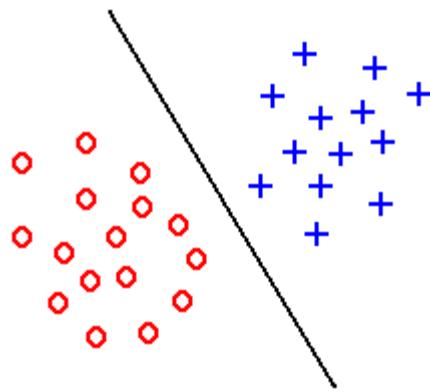


# Feature extraction

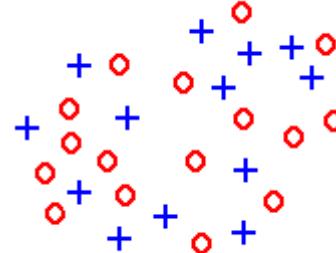
Task: to extract features which are good for classification.

Good features:

- Objects from the same class have similar feature values.
- Objects from different classes have different values.

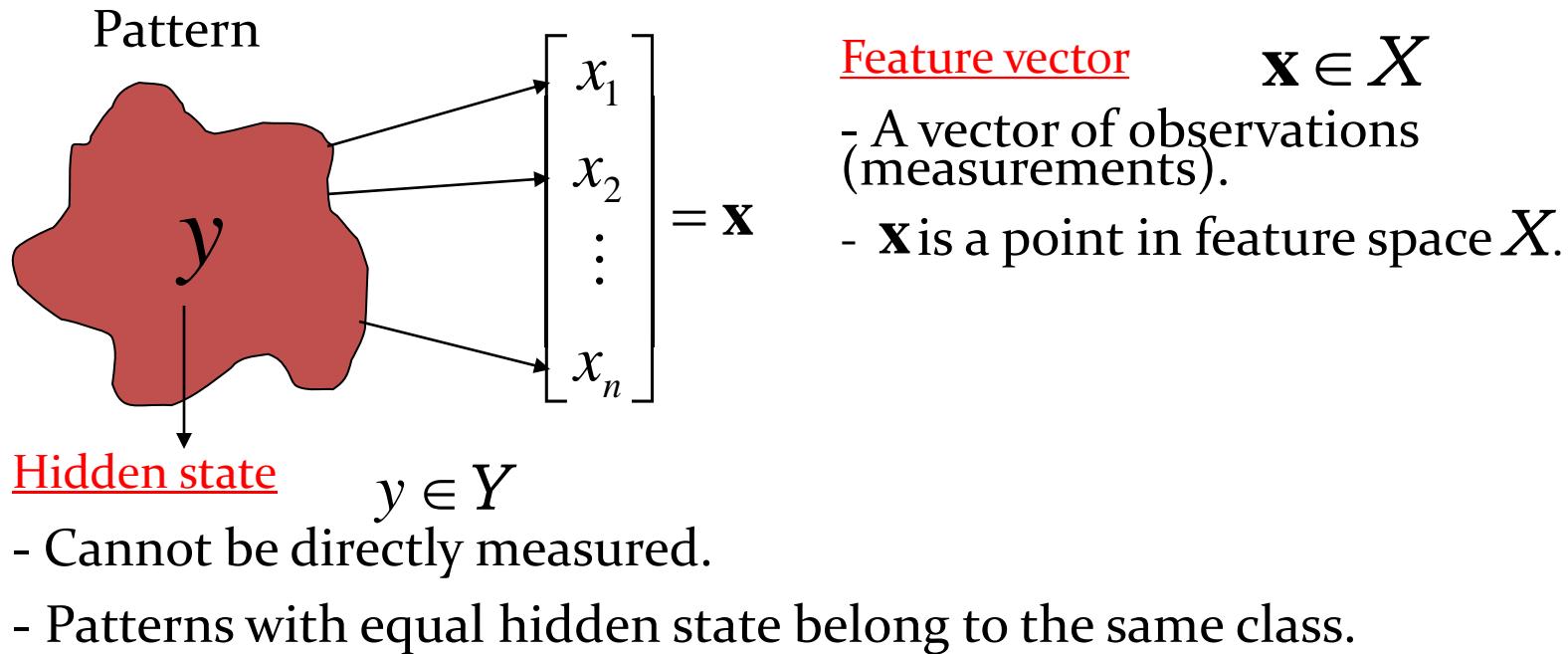


“Good” features



“Bad” features

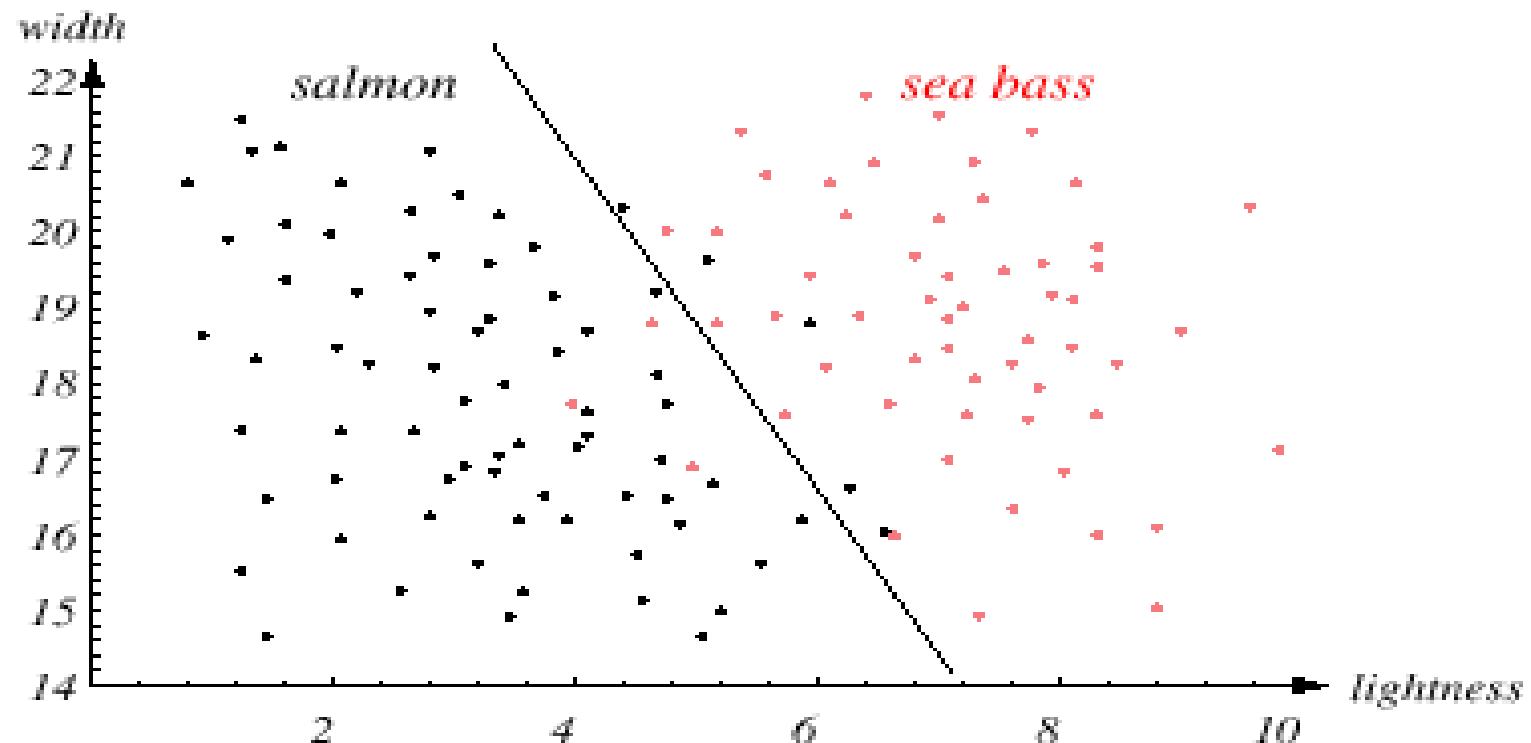
# Basic concepts



## Task

- To design a classifier (decision rule)  $q : X \rightarrow Y$  which decides about a hidden state based on an observation.

- Plotting Salmon and Seabass based on two-dimensional feature vector.

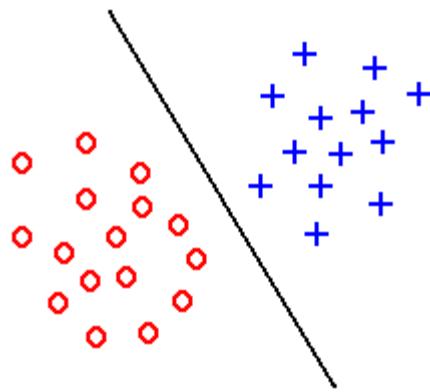


# Feature extraction

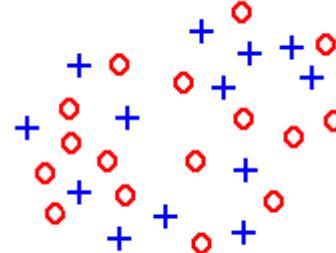
Task: to extract features which are good for classification.

Good features:

- Objects from the same class have similar feature values.
- Objects from different classes have different values.



“Good” features



“Bad” features

# Text Classification



→ Converts  $\text{TFIDF}$  (<sup>Text Document, <sup>Count</sup> words in the dataset</sup>)

From: [soccer@csail.mit.edu](#)  
Newsgroups: comp.graphics  
Subject: Need specs on Apple Q/T  
  
I need to get the specs, or at least a very verbose interpretation of the specs, for QuickTime. Technical articles from magazines like references to books would be nice, too.  
  
I also need the specs in a format usable on a Unix or MS-Dos system. I can't do much with the QuickTime stuff they have on ...

0	baseball
3	specs
0	graphics
1	references
0	hockey
0	car
0	clinton
.	
.	
.	
1	unix
0	space
2	quicktime
0	computer

Term Frequency / Inverse Document Frequency (TFIDF)

- Representing Text as a Vector.
- Stem words used, such that “computer, computes ..” all get noted under “compute.”
- The number in the vector is actually divided by the number of documents that number appears in. “Inverse Document Frequency”

Divide by # of document

Remove "stop words" e.g. is, the, a, an

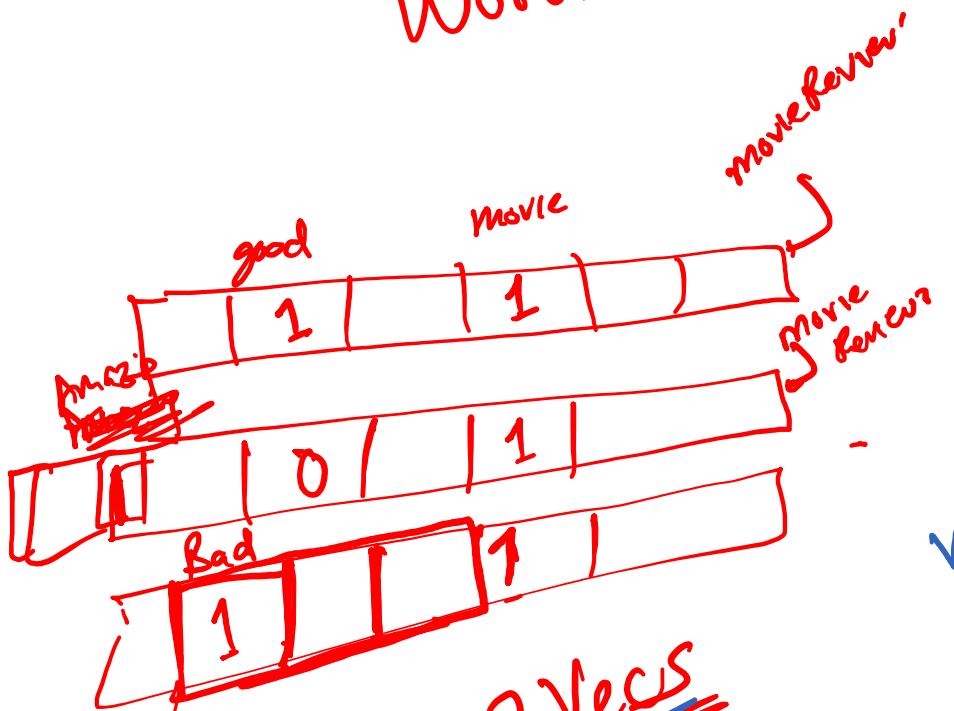
'Stemming' ← keep only the stem words

Compute, Computer, Computing.

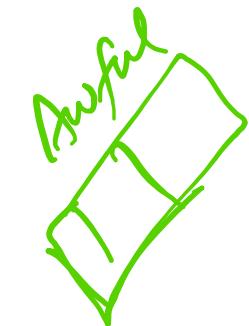
!, ?,

Sentiment detection  
for movie reviews

## Word 2 Vecs



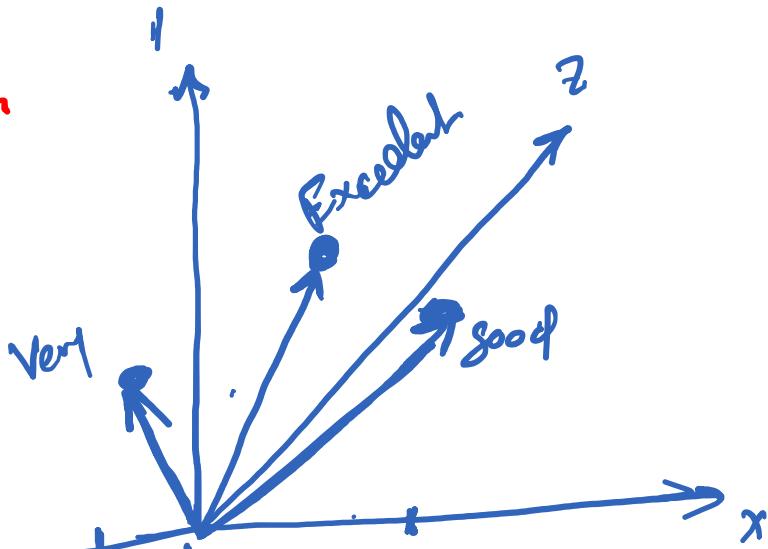
## Word 2 Vecs



"awful"  
"Very"  
"excellent"

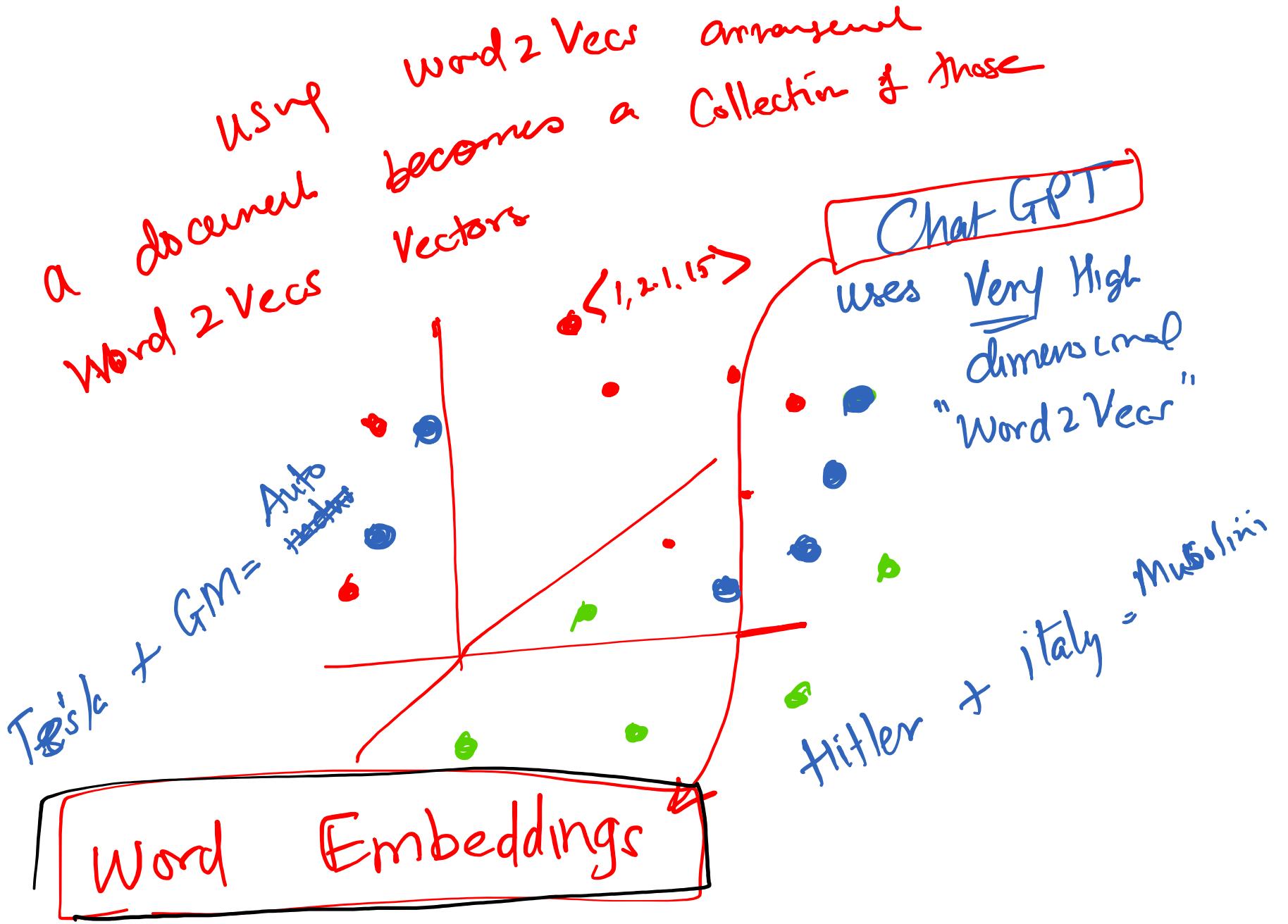
Word 1 "good"	[2 3.1 1.2]
"bad"	[-1 1 1.5]
"excellent"	[1 4.1 2.7]

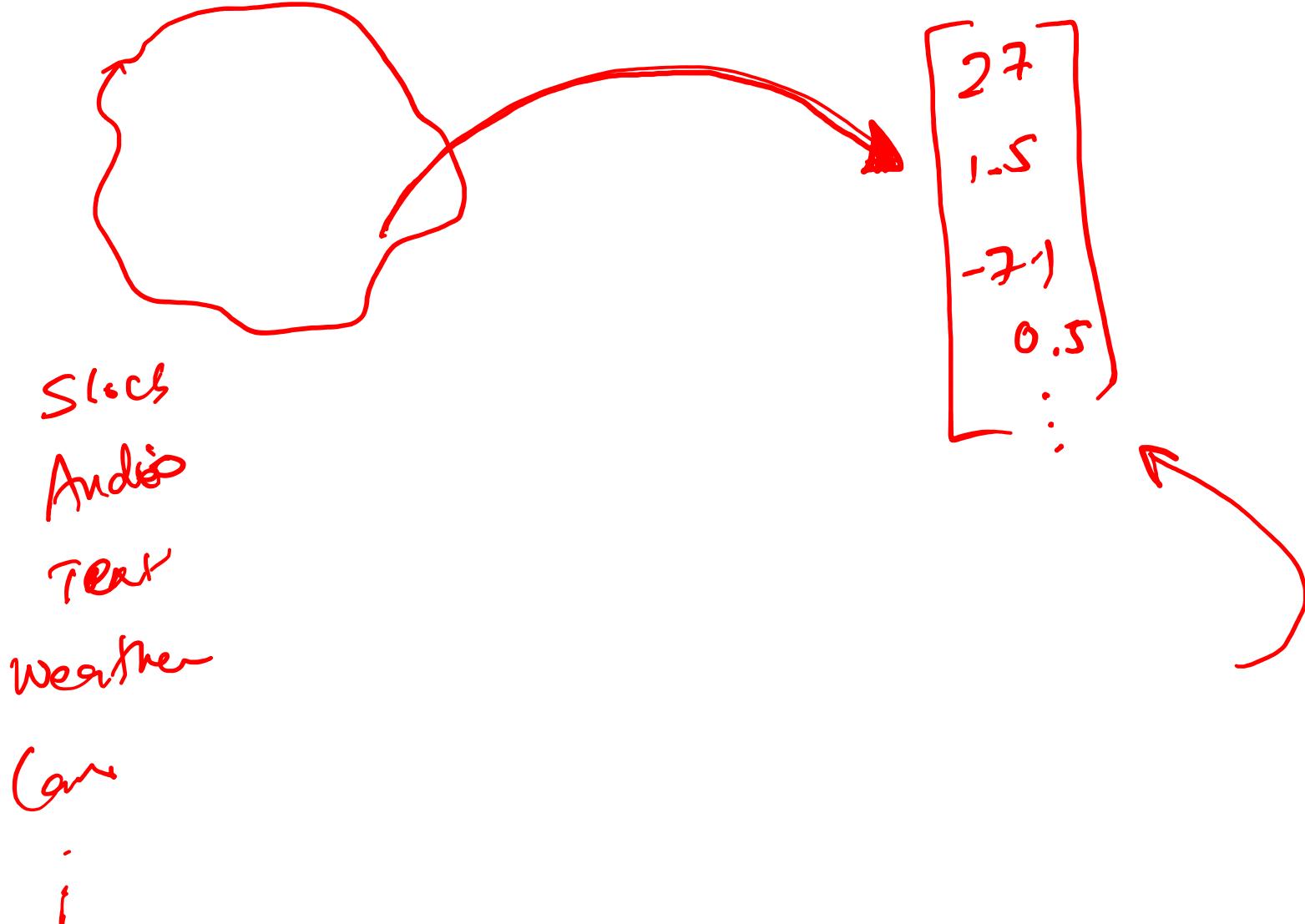
One document is one feature vector  
What should be the vector for the word "bad"



Document	[1 1 1]
Document	[2 3.1 1.2]

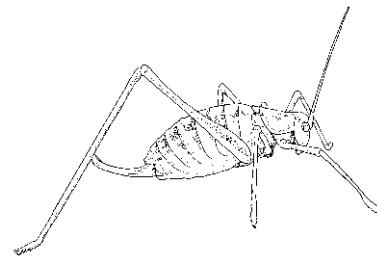
will be collected  
of all word vecs





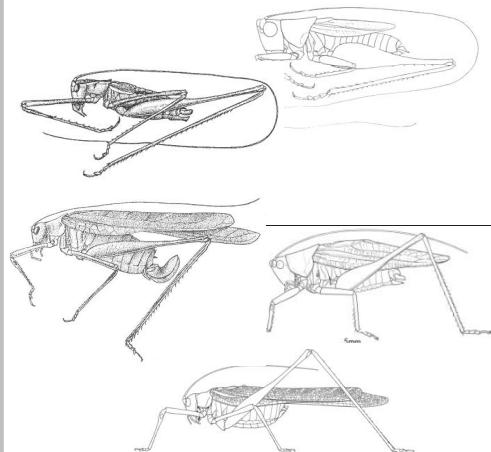
# Let's go back to agriculture!

Given a collection of annotated data. In this case 5 instances **Katydid**s and five of **Grasshoppers**, decide what type of insect the unlabeled example is.

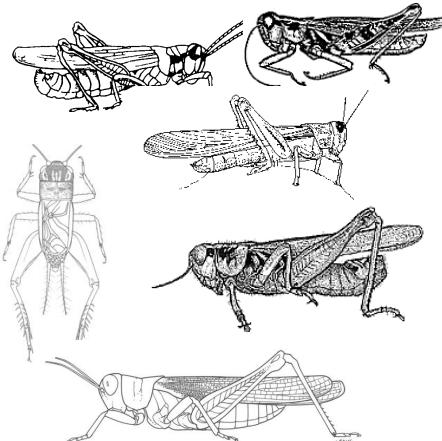


**Katydid or Grasshopper?**

## Katydid



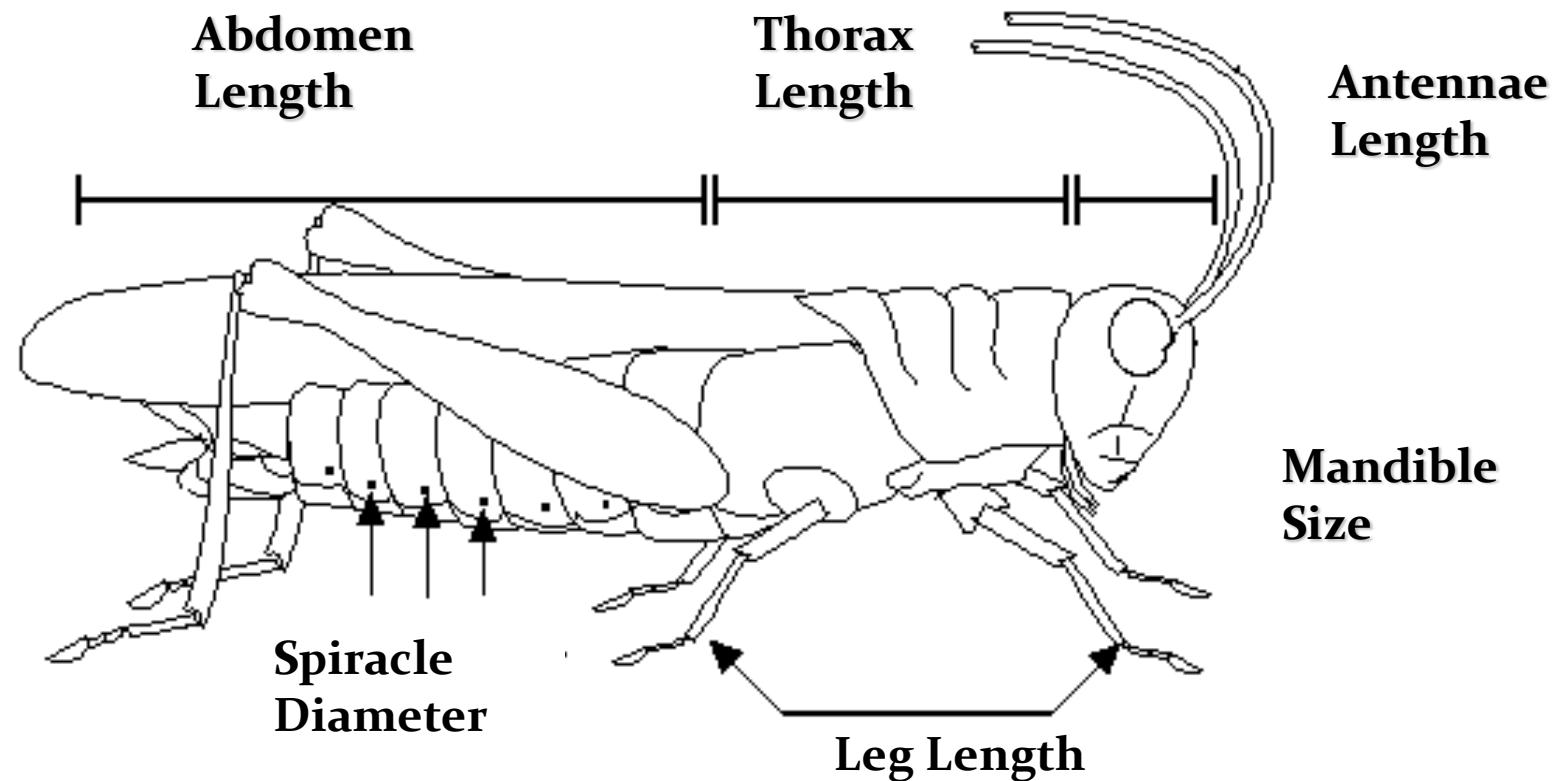
## Grasshoppers



For any domain of interest, we can measure *features*

Color {Green, Brown, Gray, Other}

Has Wings?



We can store features in a database.

The classification problem can now be expressed as:

- Given a training database (**My\_Collection**), predict the **class** label of a **previously unseen** instance

**My\_Collection**

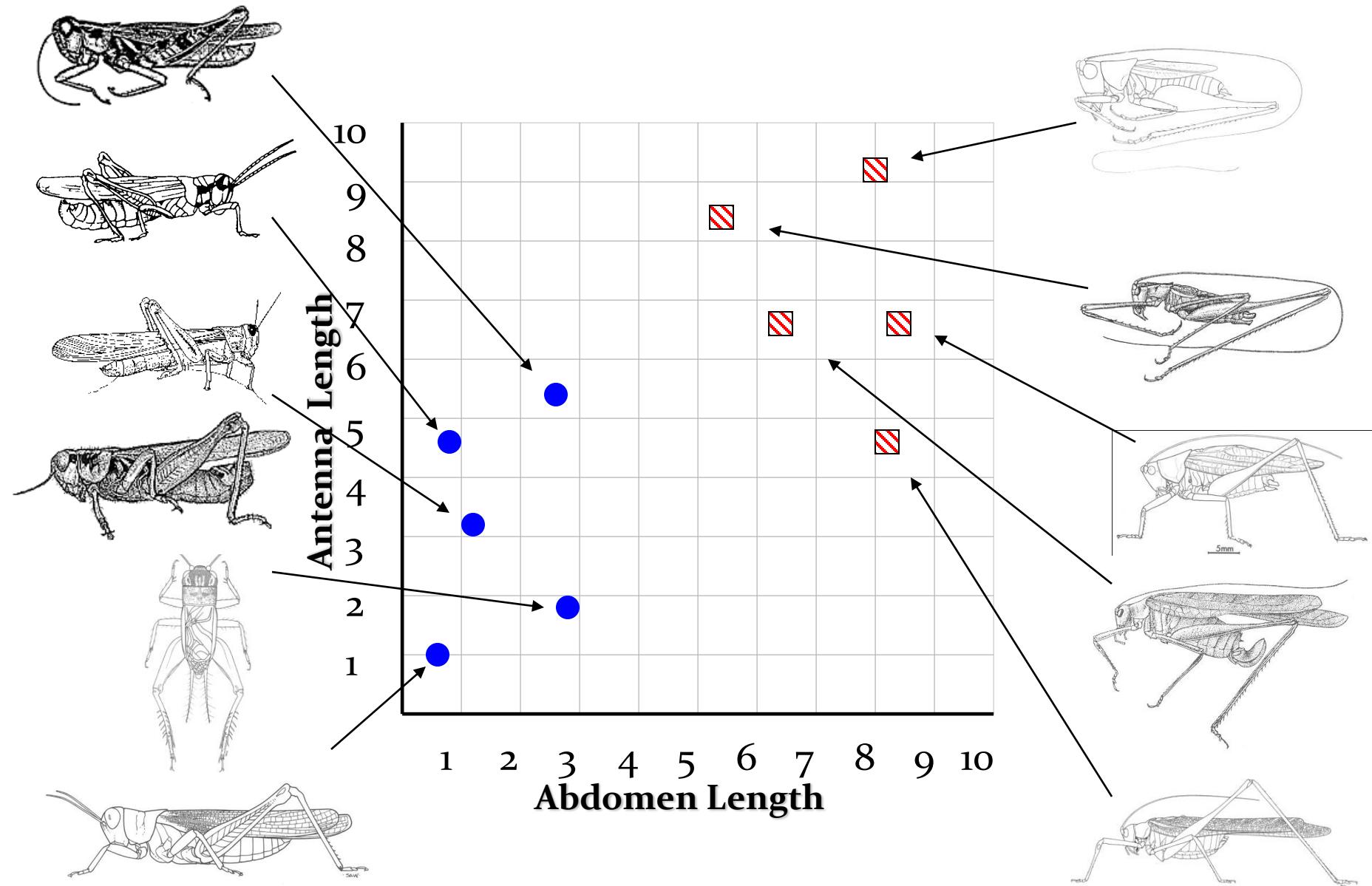
Insect ID	Abdomen Length	Antennae Length	Insect Class
1	2.7	5.5	Grasshopper
2	8.0	9.1	Katydid
3	0.9	4.7	Grasshopper
4	1.1	3.1	Grasshopper
5	5.4	8.5	Katydid
6	2.9	1.9	Grasshopper
7	6.1	6.6	Katydid
8	0.5	1.0	Grasshopper
9	8.3	6.6	Katydid
10	8.1	4.7	Katydids

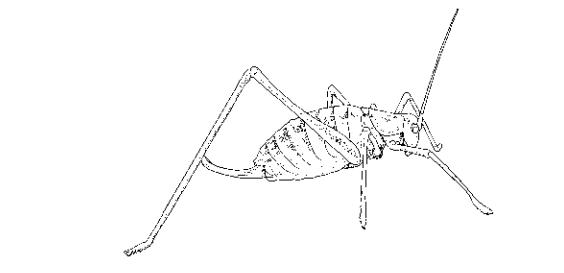
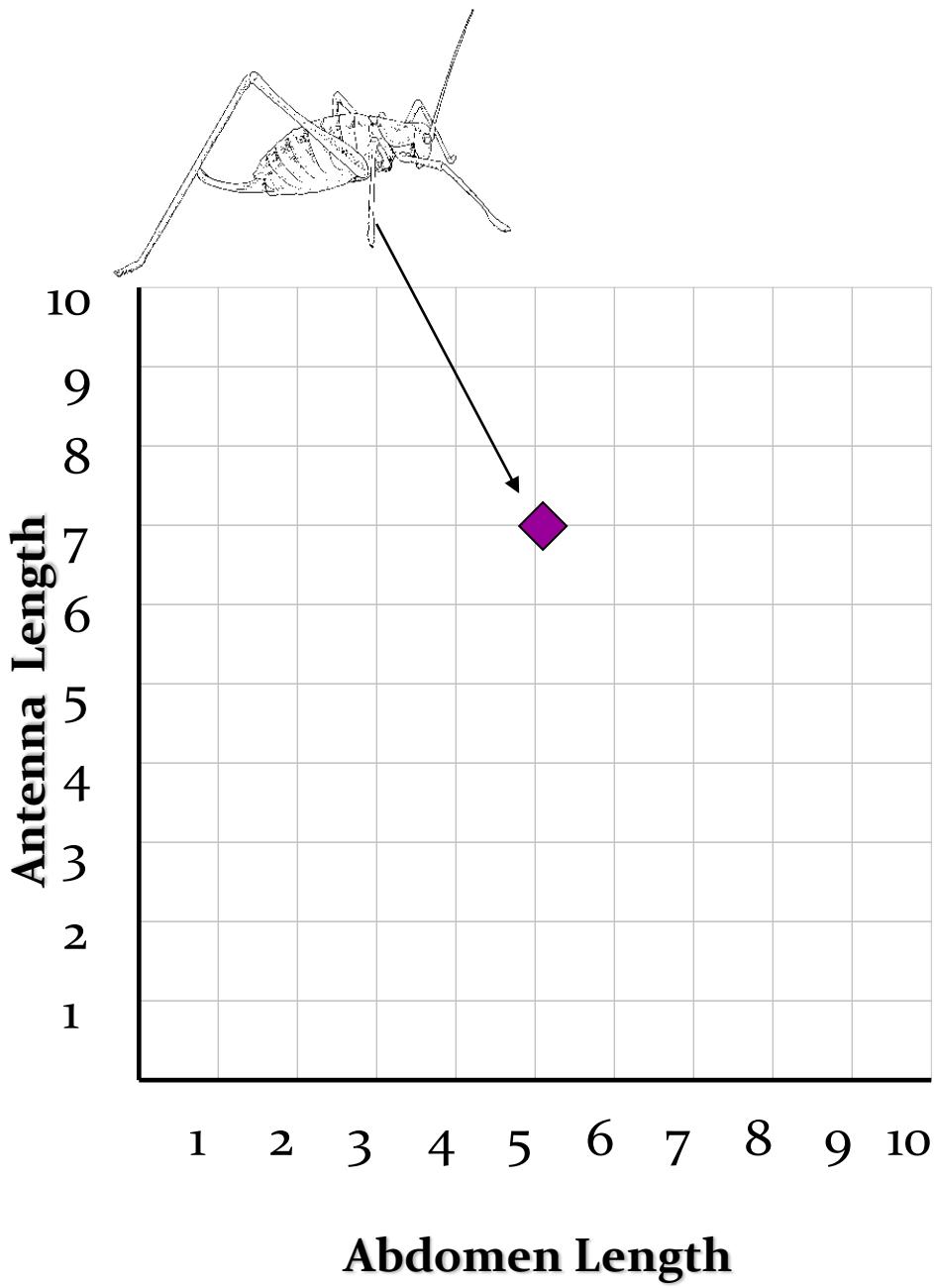
previously unseen instance =

11	5.1	7.0	???????
----	-----	-----	---------

# Grasshoppers

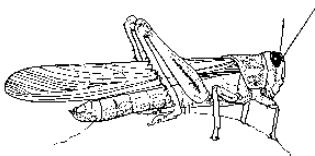
# Katydid



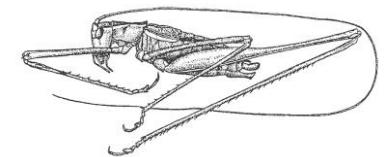


Katydid or Grasshopper?

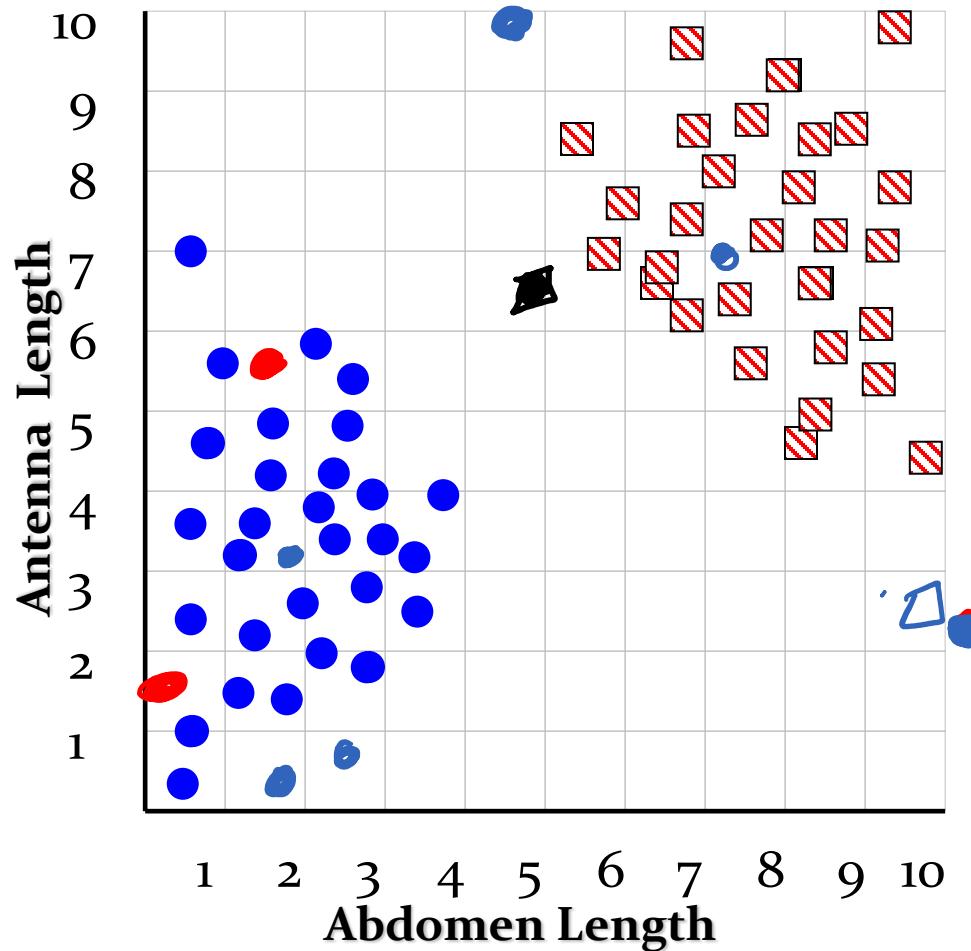
# Grasshoppers



# Katydid



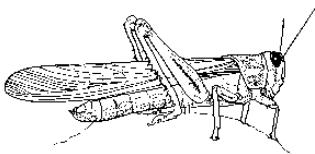
We will also use this larger dataset as a motivating example...



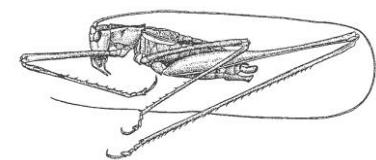
Each of these data objects are called...

- exemplars
- (training) examples
- instances
- tuples

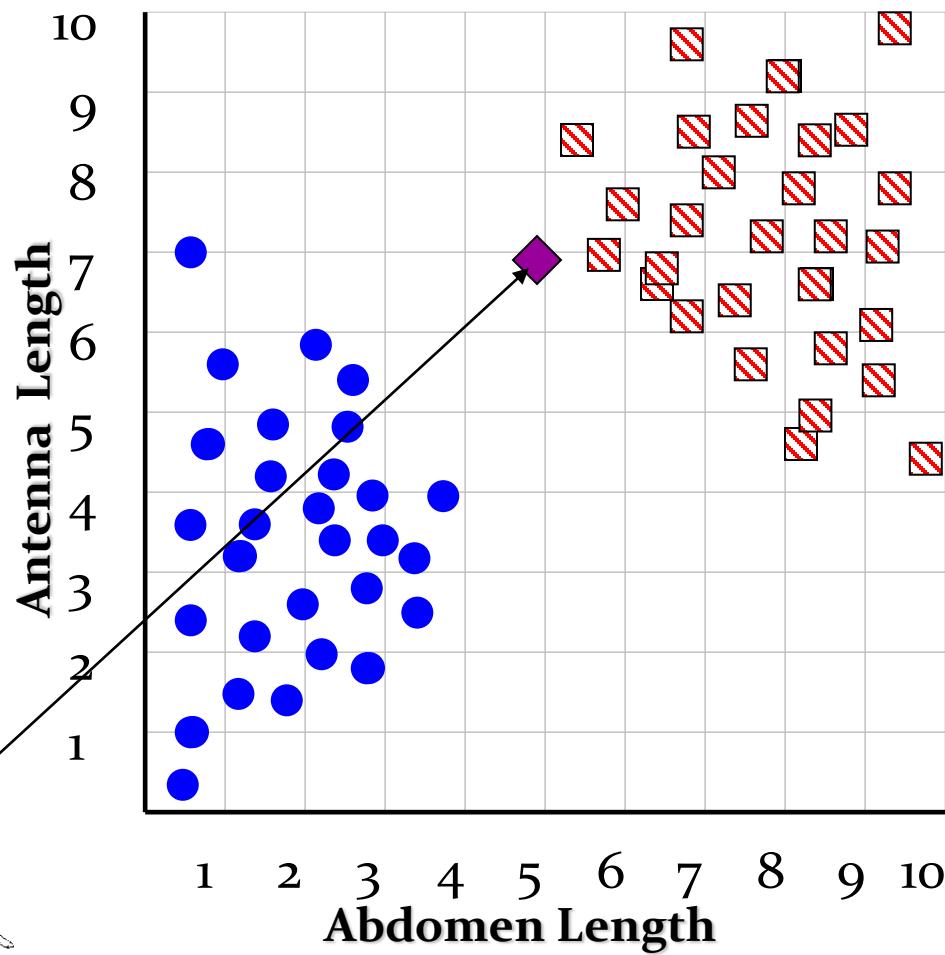
# Grasshoppers



# Katydid



We will also use this larger dataset as a motivating example...

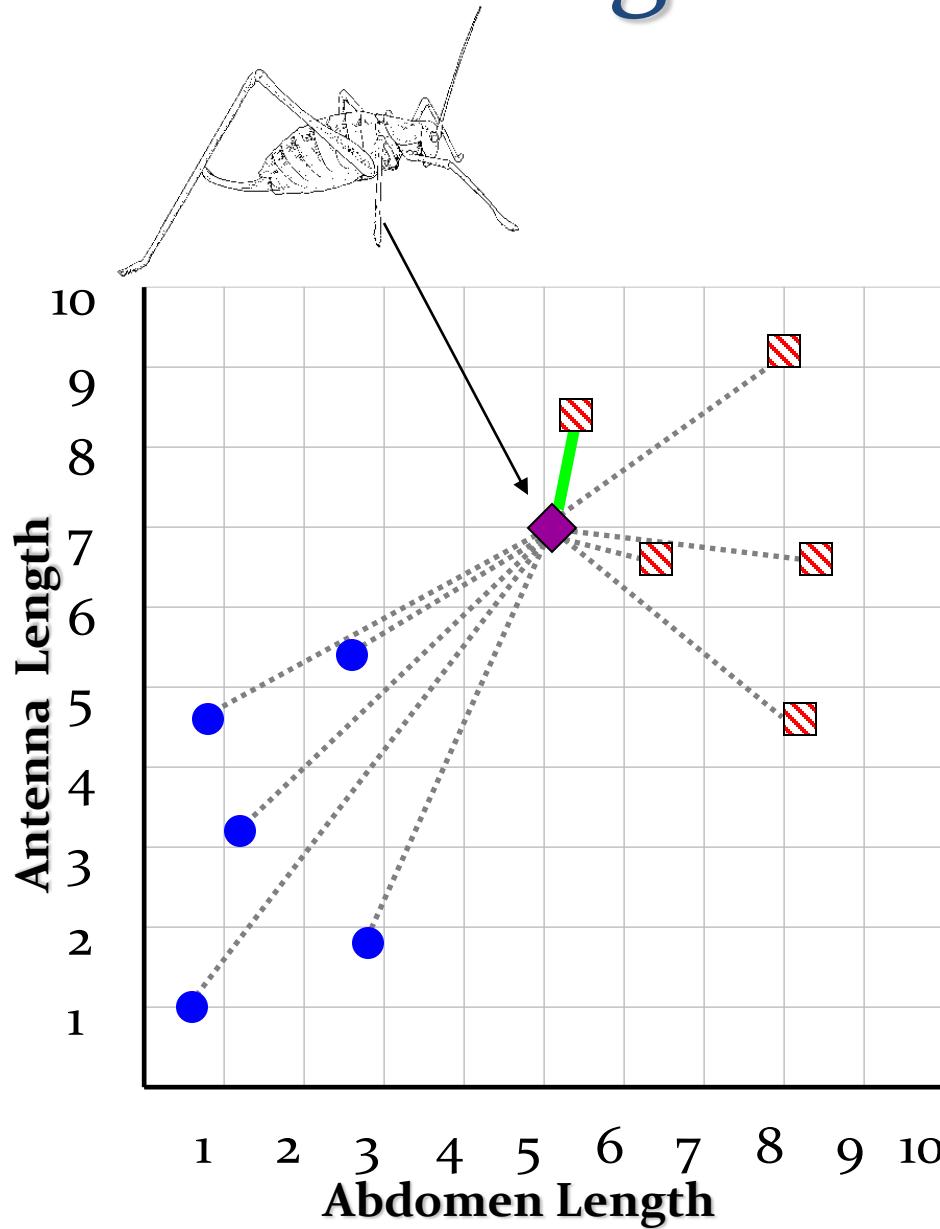


Each of these data objects are called...

- exemplars
- (training) examples
- instances
- tuples

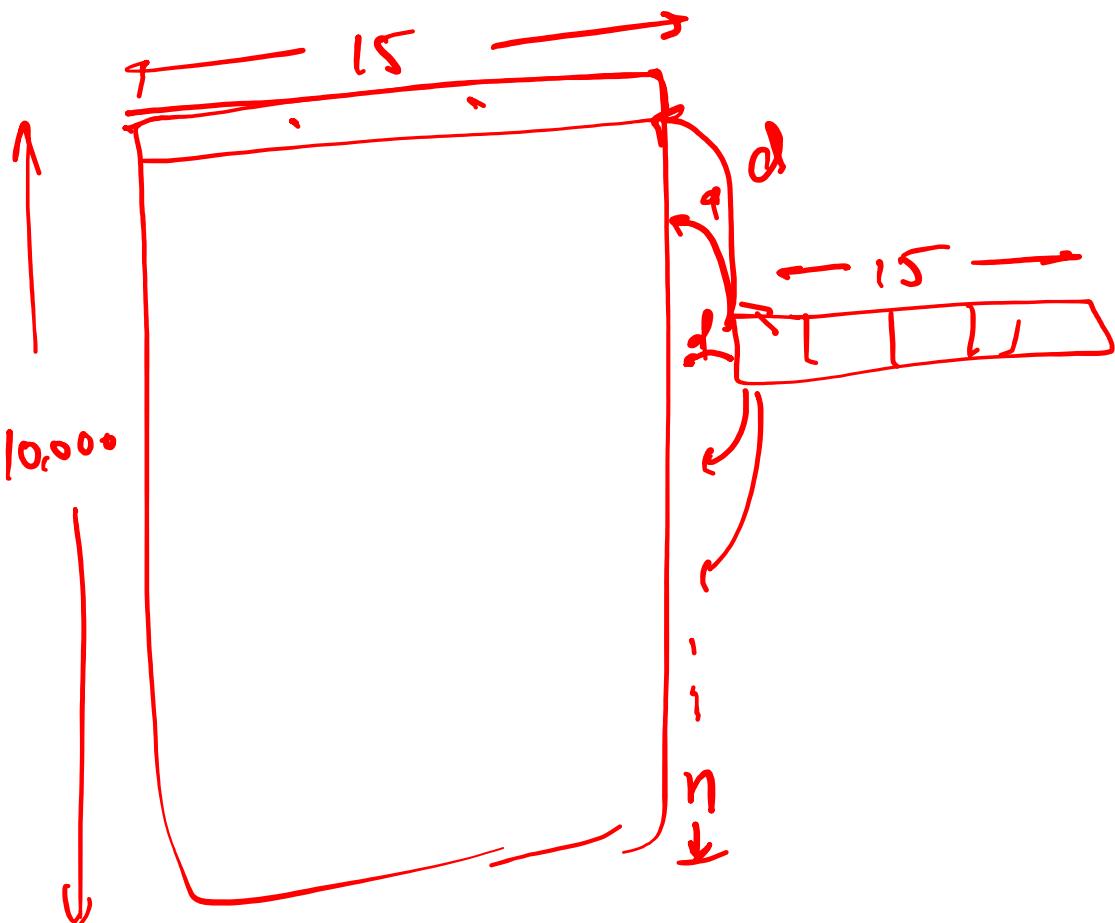
????

# Nearest Neighbor Classifier



If the **nearest** instance to the **previously unseen instance** is a **Katydid**  
class is **Katydid**  
else  
class is **Grasshopper**

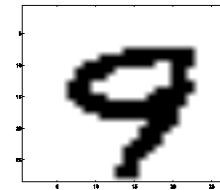
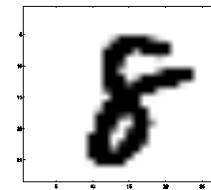
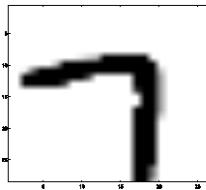
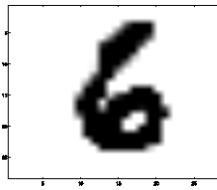
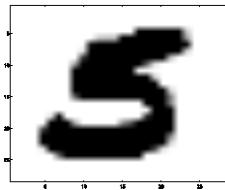
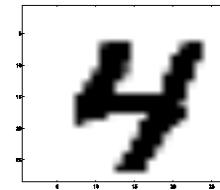
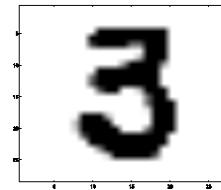
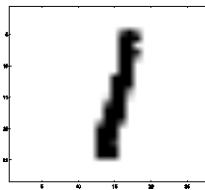
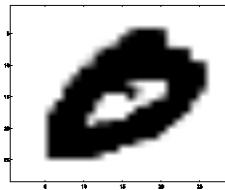
■ **Katydid**  
● **Grasshopper**



Model, optimales  $f$ .

# Hand Written Digits example

Database of 20,000 images of handwritten digits, each labeled by a human (Supervised Learning)



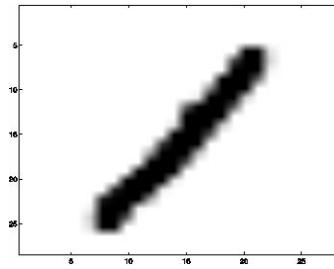
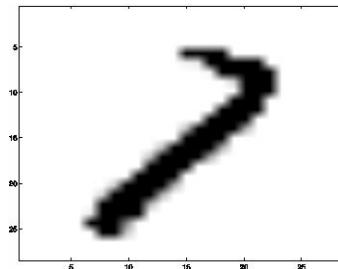
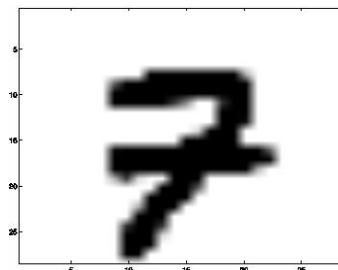
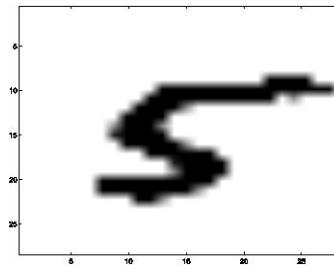
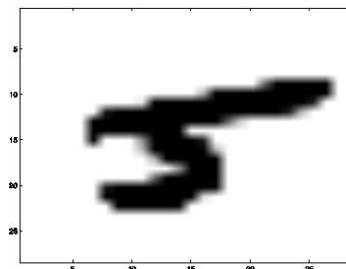
[28 x 28 greyscale; pixel values 0-255; labels 0-9]

Use these to learn a classifier which will label digit-images automatically...

# Nearest neighbor

Image to label

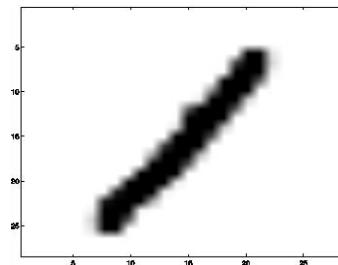
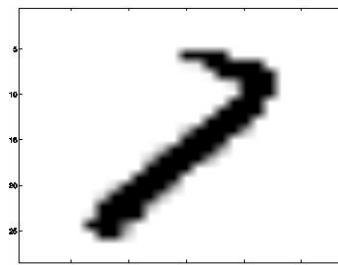
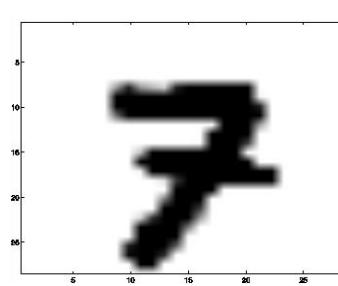
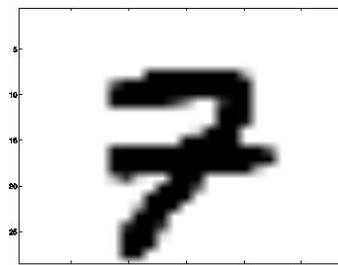
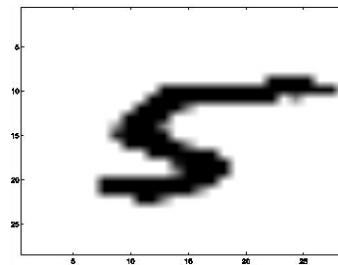
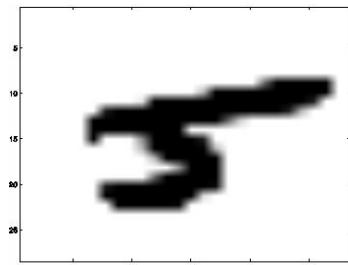
Nearest neighbor



# Nearest neighbor

Image to label

Nearest neighbor



Overall:  
error rate = 6%  
(on test set)