# Statistical Methods for NLP

## Text Similarity, Text Categorization, Linear Methods of Classification

Sameer Maskey

# Announcement

- **Reading Assignments**
  - Bishop Book – 6.1, 6.2, 7.1 (upto 7.1.1 only)
  - J&M Book – 3.1, 4.5.1
  - Journal paper on ML for Text Categorization assigned
- **Teaching Assistant**
  - Kapil Thadani, PhD student, Computer Science
  - kapil@cs.columbia.edu
  - Office Hour: 3 to 4pm, Thursdays, NLP Lab
- **February 1st, Monday, 11am at Davis Auditorium**
  - Distinguished lecture by John Lafferty

# Class Information from Feedback Survey

- **Students**
  - PhD, Masters, Undergraduate
  - Computer Science, Applied Math, Statistics, Bioinformatics
  - Quite a few have taken NLP
  - Some have taken ML or NLP or Both
  - Wide variation in background

- Course load (3 HW, 1 project, 1 Final):  Most said ok
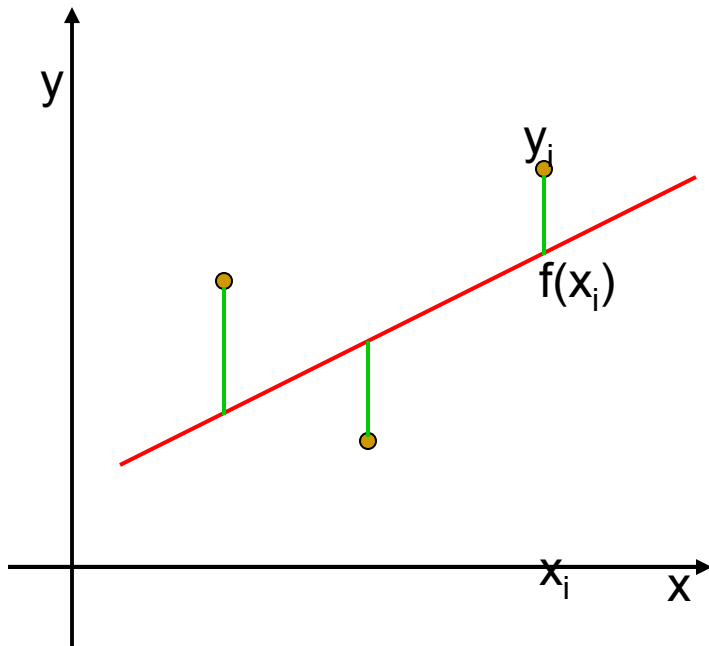- Style:  More wanted mixed

# Project

- You should start thinking about the kind of project you want to do
  - This is a semester long project
  - Goal is for you to understand one NLP-ML problem of your interest in-depth
  - Depending on your expertise and passion you may try to get a research paper out of it
- Some Examples
  - Review to Ratings: Automatic Reviewer
  - Build a search engine writing your own ML algorithm
  - Information Extraction for Relations
  - Statistical Parser with your statistical algorithm for decoding

# Topics for Today

- Text Categorization/Classification
- Linear Discriminant Functions
- Perceptron
- Naïve Bayes Classifier for Text

# Review: Linear Regression

$$J(\theta) = \frac{1}{N} \sum_{i=1}^{N} (y_i - (\theta_0 + \theta_1 x_i))^2$$



$$\theta_1 = \frac{\sum_{i=1}^{N} x_i y_i - \frac{1}{N} \sum_{i=1}^{N} x_i \sum_{i=1}^{N} y_i}{\sum_{i=1}^{N} x_i^2 - \frac{1}{N} \sum_{i=1}^{N} x_i \sum_{i=1}^{N} x_i}$$

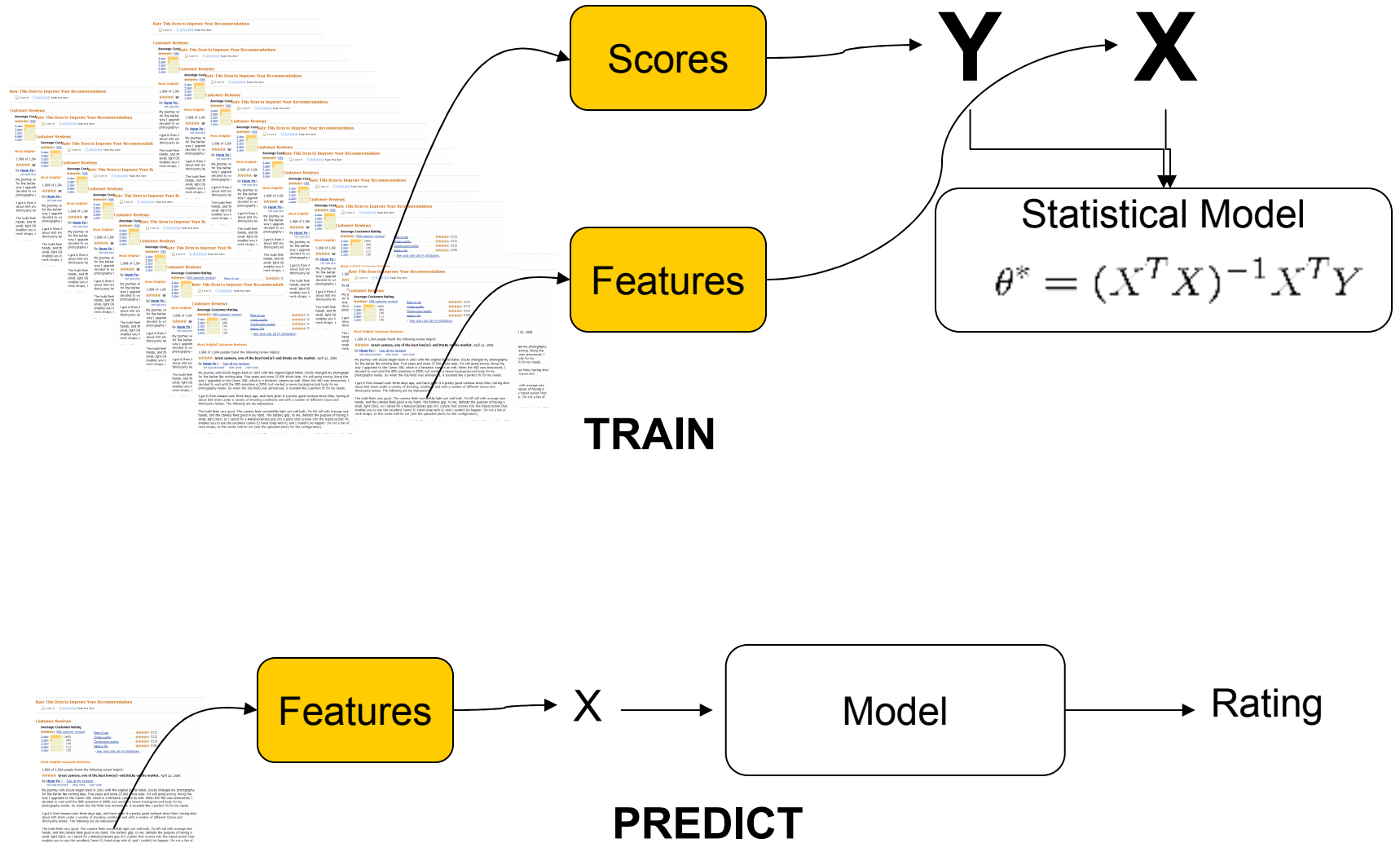$$\theta_0 = \frac{1}{N} \sum_{i=1}^{N} y_i - \frac{1}{N} \theta_1 \sum_{i=1}^{N} x_i$$

# Review: Multiple Linear Regression

$$J(\theta) = \frac{1}{2N} \left\| \begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ \cdot \\ Y_N \end{bmatrix} - \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1K} \\ 1 & x_{21} & x_{22} & \cdots & x_{2K} \\ & & \cdots & & \\ & & \cdots & & \\ & & \cdots & & \\ 1 & x_{N1} & x_{N2} & \cdots & x_{NK} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \cdot \\ \cdot \\ \cdot \\ \theta_K \end{bmatrix} \right\|^2$$

**Y**        **X**        **θ**

$$\theta^* = (X^T X)^{-1} X^T Y$$

# Reviews to Automatic Ratings
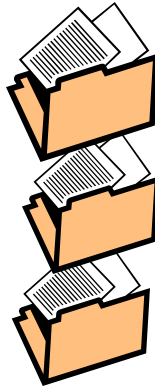
**Scores**

**Features**

**Y** **X**

Statistical Model

$$\theta^* = (X^T X)^{-1} X^T Y$$

**TRAIN**

**Features** **X** Model Rating

**PREDICT**

# Predicted Output May Be Classes

- "…is writing a paper"
- "… has flu ☹"
- "… is happy, yankees won!"

➡

$$\begin{bmatrix} \text{SAD} \\ \text{SAD} \\ \text{HAPPY} \end{bmatrix}$$
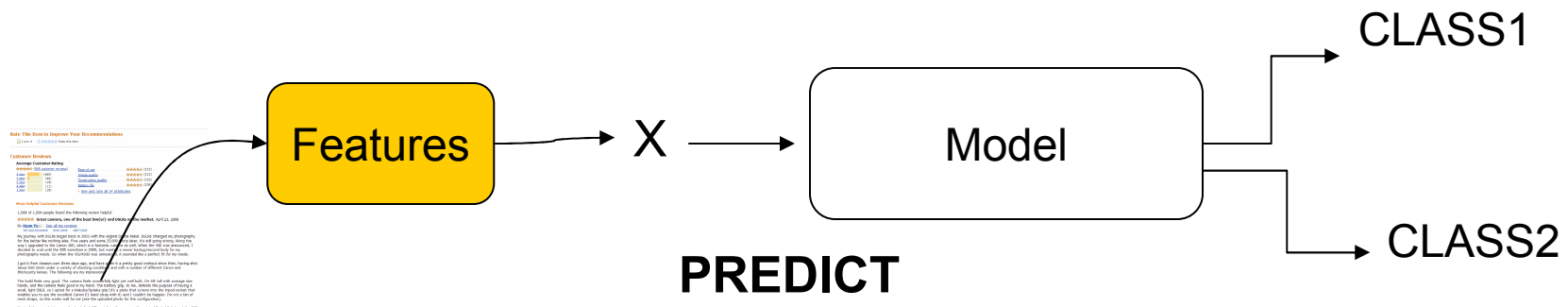


➡

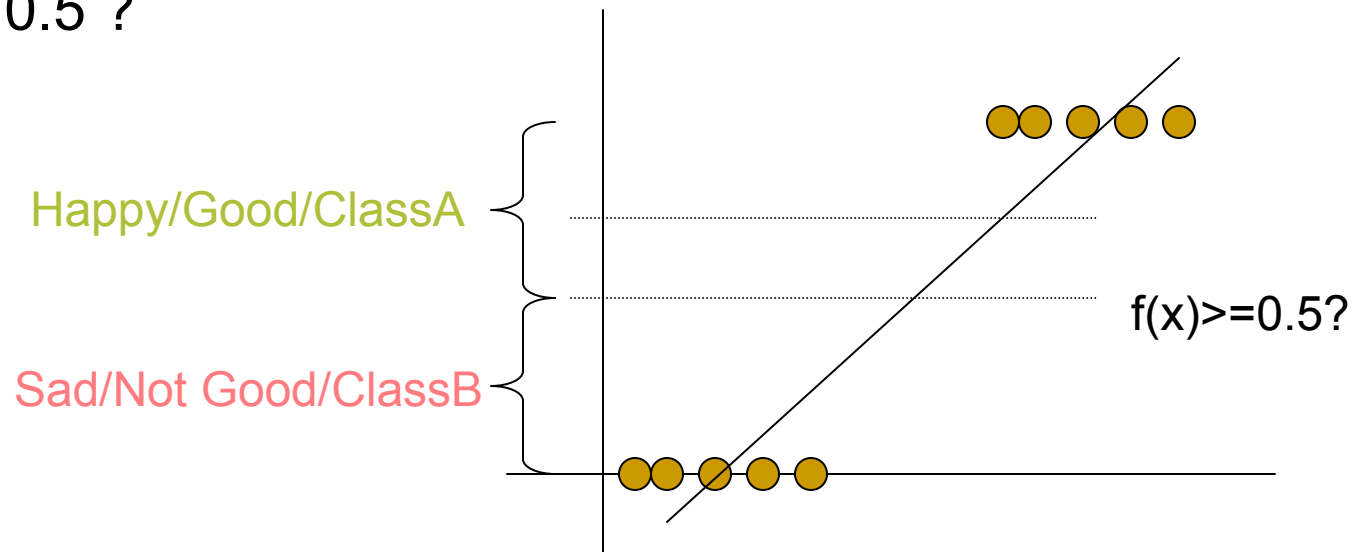$$\begin{bmatrix} \text{DIABETES} \\ \text{HEPATITIS} \\ \text{HEPATITIS} \end{bmatrix}$$

# Text Categorization/Classification

- Given any text (sentence, document, stories, etc), we want to classify it into some predefined class set



CLASS1

Features → X → Model

**PREDICT**

CLASS2

- Training Data consists of Y values that are 0 and 1
  - Review is good or bad
  - Status update is happy or sad

# Regression to Classification
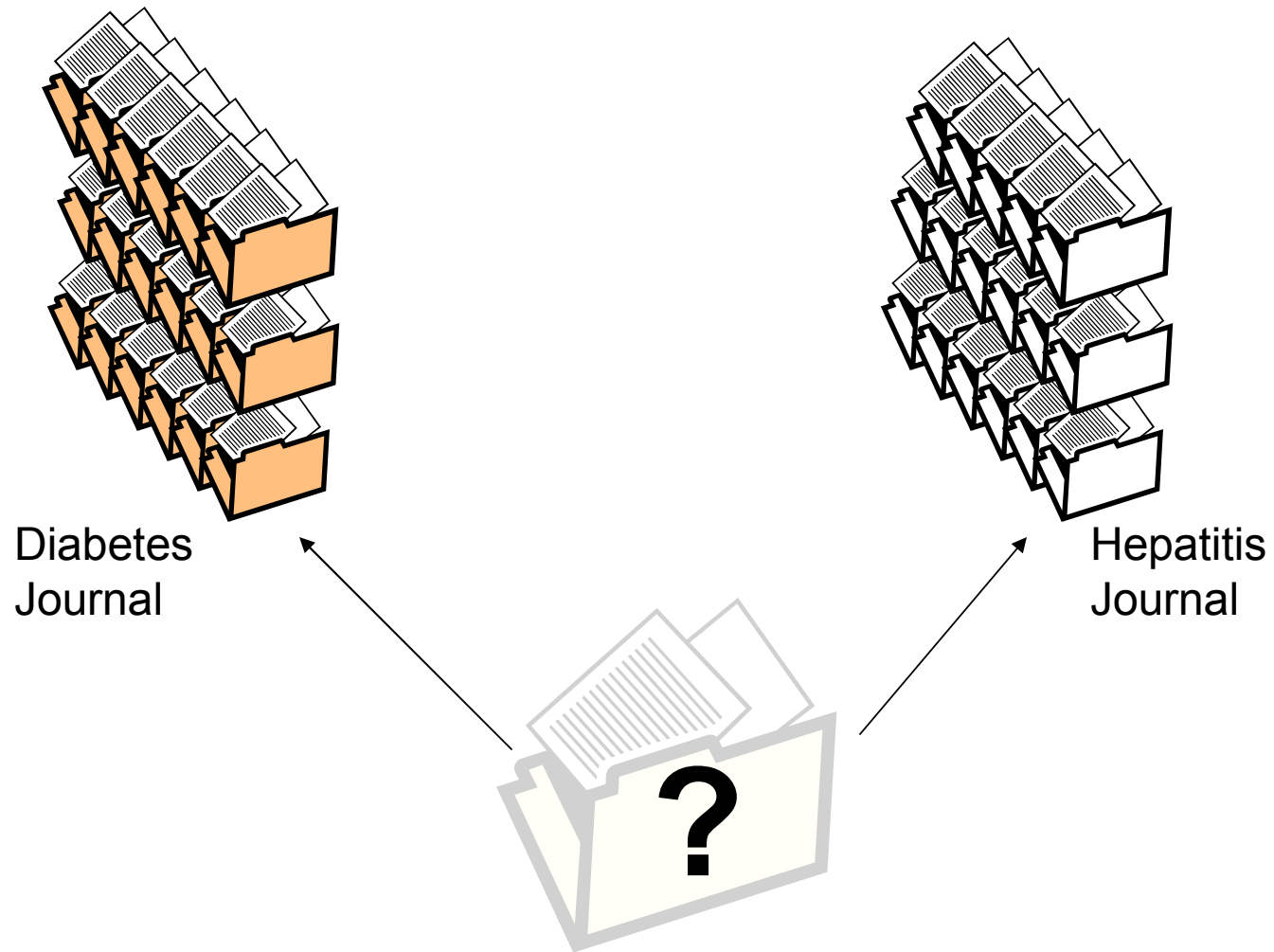
- Can we build a regression model to model such binary classes?

- Train Regression and threshold the output
  - If $f(x) >= 0.7$ CLASS1
  - If $f(x) < 0.7$ CLASS2
  - $f(x) >= 0.5$ ?

Happy/Good/ClassA

Sad/Not Good/ClassB

$f(x)>=0.5?$

# Regression to Classification

- Thresholding on regression function does not always work

- Gaussian assumption on noise

- When the output is binary class, we may want to dry a different technique of modeling than regression

- Many modeling techniques that will better produce class category values we want for Y
  - Linear Classifiers is one such method

# Text Classification
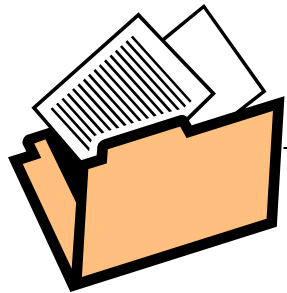


Diabetes Journal

Hepatitis Journal

?

# Text Similarity

- To classify a new journal paper into either diabetes group of hepatitis group we could probably compute similarity of this document with the two groups

- How do we compute similarity between text or group of documents?

- First, we need representation of text that takes account of all the information in it?
  - Count of +ve words may not be enough

# Text/Document Vectors

- Document Vectors
  - Documents can be represented in different types of vectors: binary vector, multinomial vector, feature vector

- Binary Vector: For each dimension, 1 if the word type is in the document and 0 otherwise

- Multinomial Vector: For each dimension, count # of times word type appears in the document

- Feature Vector: Extract various features from the document and represent them in a vector. Dimension equals the number of features
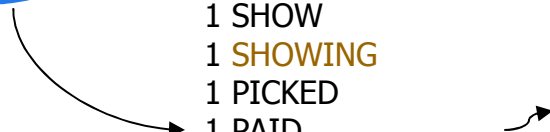
# Example of a Multinomial Document Vector

Screening of the critically acclaimed film NumaFung Reserved tickets can be picked up on the day of the show at the box office at Arledge Cinema. Tickets will not be reserved if not paid for in advance.

| Count | Word | |
|---|---|---|
| 4 | THE | 4 |
| 2 | TICKETS | 2 |
| 2 | RESERVED | 2 |
| 2 | OF | 2 |
| 2 | NOT | 2 |
| 2 | BE | 2 |
| 2 | AT | 2 |
| 1 | WILL | 1 |
| 1 | UP | 1 |
| 1 | SHOW | 1 |
| 1 | SCREENING | 1 |
| 1 | PICKED | 1 |
| 1 | PAID | 1 |
| 1 | ON | 1 |
| 1 | OFFICE | 1 |
| 1 | NUMAFUNG | 1 |
| 1 | IN | 1 |
| 1 | IF | 1 |
| 1 | FOR | 1 |
| 1 | FILM | 1 |
| 1 | DAY | 1 |
| 1 | CRITICALLY | 1 |
| 1 | CINEMA | 1 |
| 1 | CAN | 1 |
| 1 | BOX | 1 |
| 1 | ARLEDGE | 1 |
| 1 | ADVANCE | 1 |
| 1 | ACCLAIMED | |

# Example of a Multinomial Document Vector

| | |
|---|---|
| 4 THE | 4 |
| 2 SEATS | 2 |
| 2 RESERVED | 2 |
| 2 OF | 2 |
| 2 NOT | 2 |
| 2 BE | 2 |
| 2 AT | 2 |
| 1 WILL | 1 |
| 1 UP | 1 |
| 1 SHOW | 1 |
| 1 SHOWING | 1 |
| 1 PICKED | 1 |
| 1 PAID | 1 |
| 1 ON | 1 |
| 1 OFFICE | 1 |
| 1 VOLCANO | 1 |
| 1 IN | 1 |
| 1 IF | 1 |
| 1 FOR | 1 |
| 1 FILM | 1 |
| 1 DAY | 1 |
| 1 CRITICALLY | 1 |
| 1 CINEMA | 1 |
| 1 CAN | 1 |

Find out how similar two text vectors are to find document similarity?

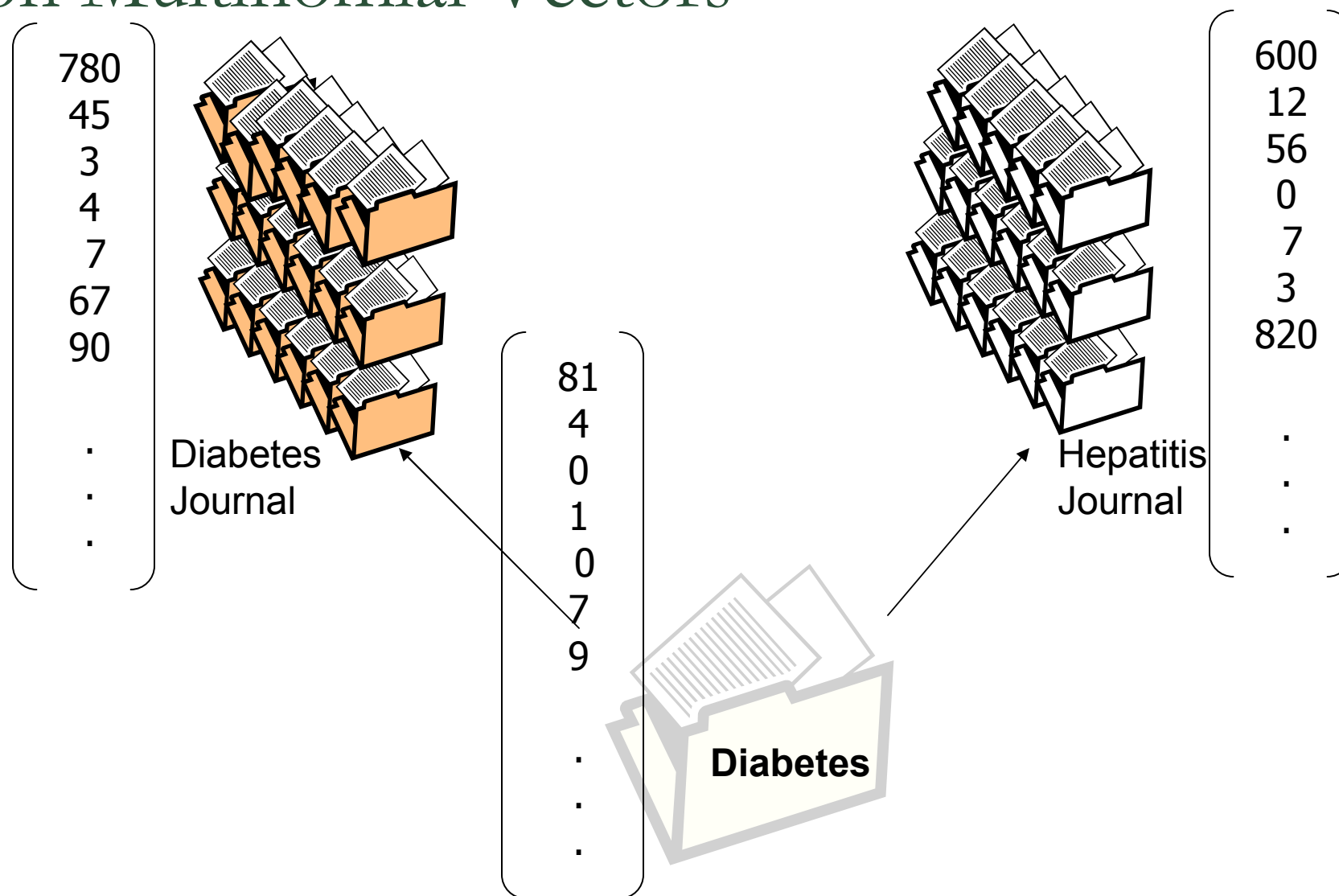# Text Similarity : Cosine of Text Vectors

- Given a set of vectors we can find the cosine similarity between them
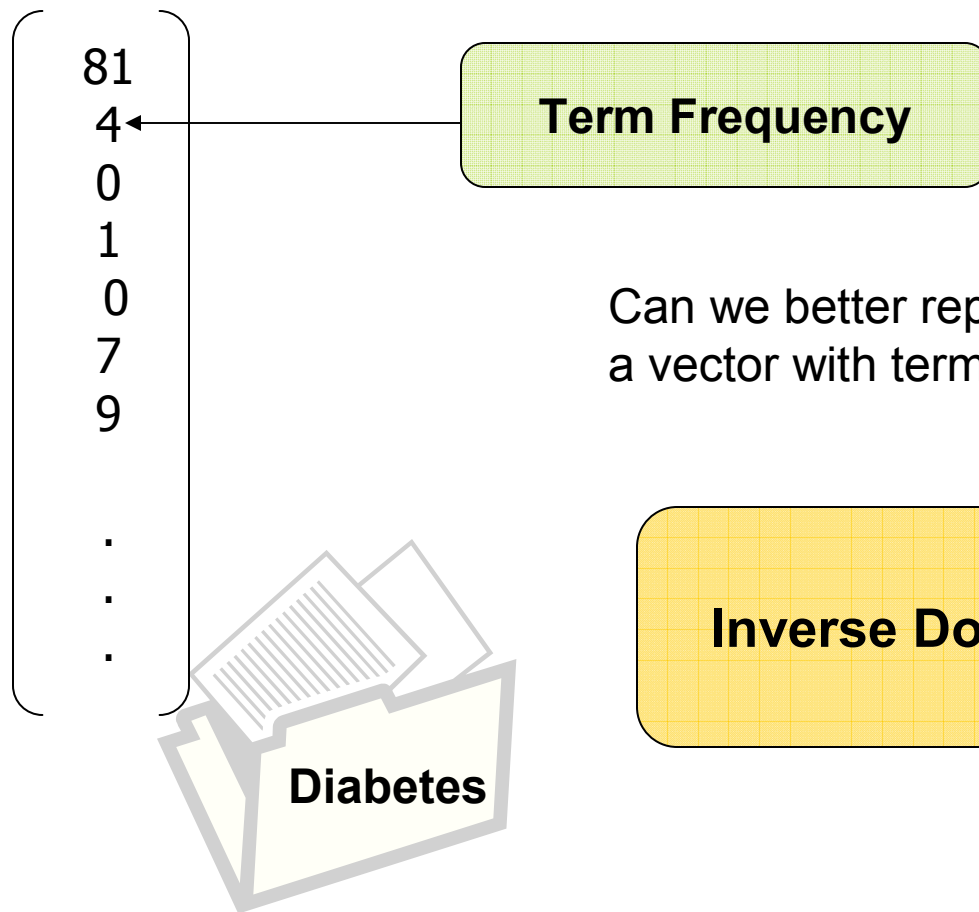
$$Cos\theta = \frac{v_1 . v_2}{\|v_1\| \|v_2\|}$$

- Cos 90 = 0, vectors are not similar
- Higher cosine value = higher similarity

# Text Classification with Cosine Similarity on Multinomial Vectors

# TF.IDF (Term Frequency and Inverse Document Frequency)

$$
\begin{bmatrix}
81 \\
4 \\
0 \\
1 \\
0 \\
7 \\
9 \\
. \\
. \\
.
\end{bmatrix}
$$

**Term Frequency**

**Diabetes**

Can we better represent the text besides a vector with term frequency for classification?
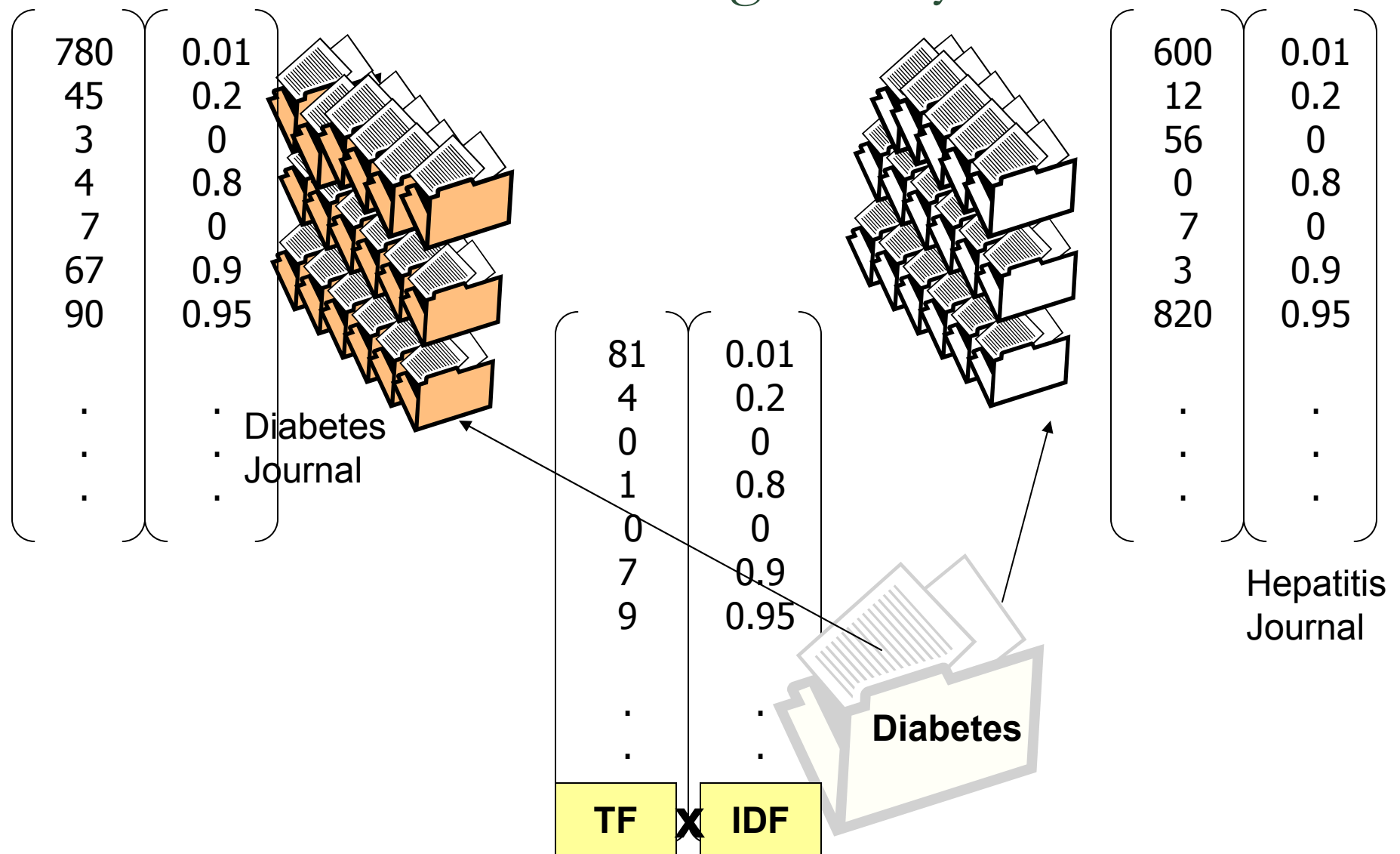
**Inverse Document Frequency?**

# Inverse Document Frequency

- Some words occur a lot no matter what the class of the document is
  - 'the, of, an, …"

- Find words that occur in a fewer number of documents
  - These words could be more discriminating across documents

# Inverse Document Frequency Weighting

- ## $IDF = \log(N/n_i)$
  - Where N is the total number documents
  - $n_i$ is the total number of documents the word occur in
  - Fewer the documents word occur in higher the IDF value
- Words such as 'the, a, on, … ' will occur in many document so will have lower IDF value
- Multiply TF with IDF to get TF.IDF weighting of words in our multinomial vector

# Text Classification with Cosine Similarity on Multinomial Vectors Weighted by TF.IDF



| 780 | 0.01 |
| 45 | 0.2 |
| 3 | 0 |
| 4 | 0.8 |
| 7 | 0 |
| 67 | 0.9 |
| 90 | 0.95 |
| . | . |
| . | . |
| . | . |

Diabetes
Journal

| 81 | 0.01 |
| 4 | 0.2 |
| 0 | 0 |
| 1 | 0.8 |
| 0 | 0 |
| 7 | 0.9 |
| 9 | 0.95 |
| . | . |
| . | . |

**TF** **X** **IDF**

**Diabetes**

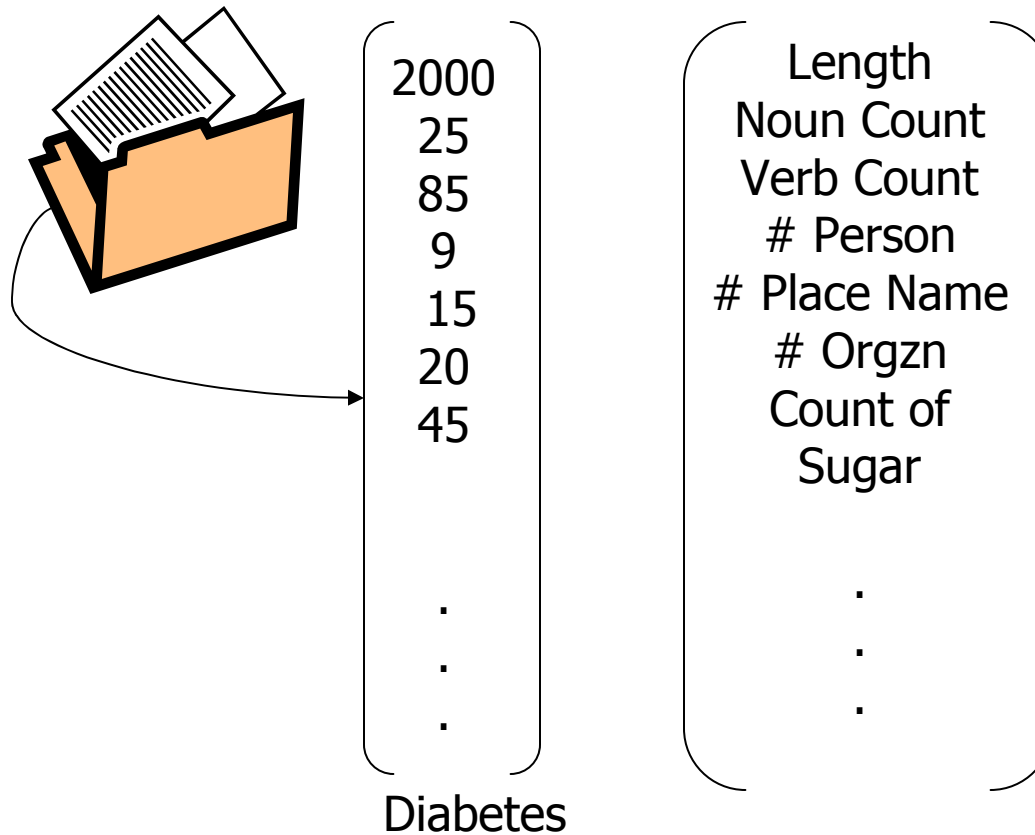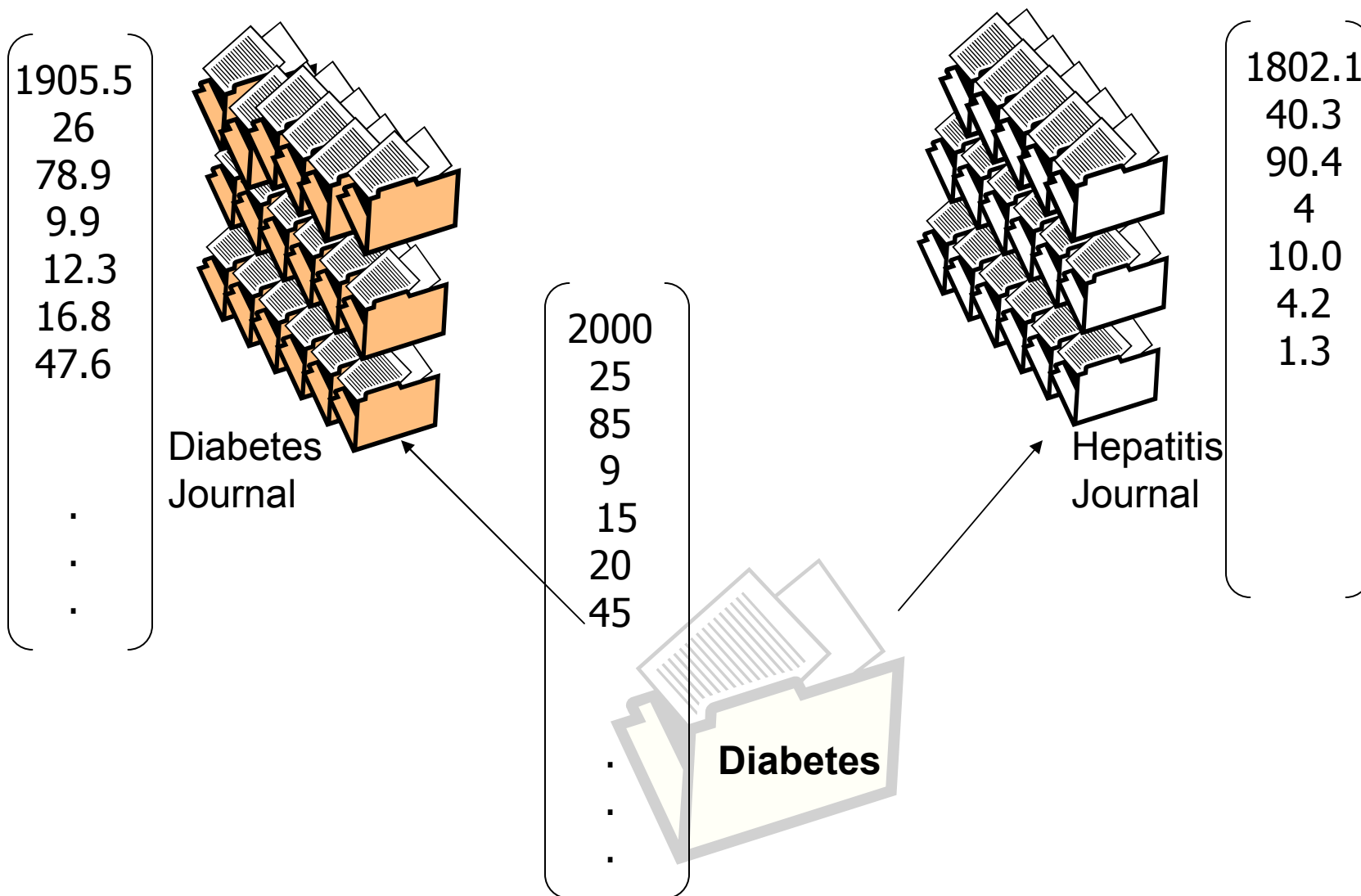| 600 | 0.01 |
| 12 | 0.2 |
| 56 | 0 |
| 0 | 0.8 |
| 7 | 0 |
| 3 | 0.9 |
| 820 | 0.95 |
| . | . |
| . | . |
| . | . |

Hepatitis
Journal

# Feature Vectors

- Instead of just using words to represent documents, we can also extract features and use them to represent the document

- We can extract features like document length (LN), number of nouns (NN), number of verbs (VB), number of person names (PN), number of place (CN) names, number of organization names (ON), number of sentences (NS), number of pronouns (PNN)

# Feature Vectors

- Extracting such features you get a feature vector of length 'K' where 'K' is the number of dimensions (features) for each document



| | |
|---|---|
| 2000 | Length |
| 25 | Noun Count |
| 85 | Verb Count |
| 9 | # Person |
| 15 | # Place Name |
| 20 | # Orgzn |
| 45 | Count of Sugar |
| . | . |
| . | . |
| . | . |

Diabetes

# Text Classification with Cosine Similarity on Feature Vectors



1905.5
26
78.9
9.9
12.3
16.8
47.6

.
.
.

Diabetes Journal

2000
25
85
9
15
20
45

.
.
.

**Diabetes**

Hepatitis Journal

1802.1
40.3
90.4
4
10.0
4.2
1.3

# Cosine Similarity Based Text Classifier

- Build multinomial vectors or feature vectors for each document in the given class
  - Each dimension represent count of the given word or feature
  - Can take average of the vectors to represent a corpus
- 'N' averaged vectors would be the model for 'N' classes of the documents
- For any new document compute similarity of its multinomial vector to the 'N' class vectors
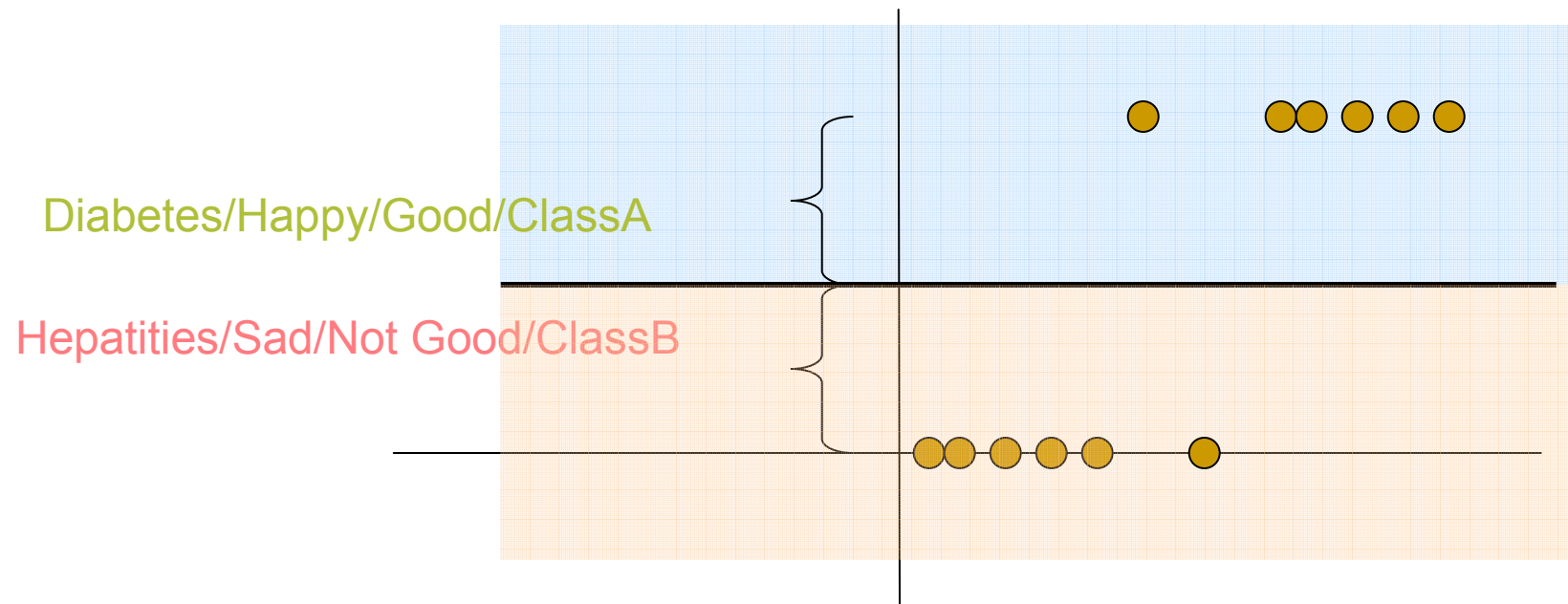- Highest Cosine Similarity represents the class

# Text Classification using Regression?

- **Regression for binary text classes not ideal**
  - 'Diabetes vs. Hepatitis'
- **We want better modeling technique**



Diabetes/Happy/Good/ClassA

Hepatities/Sad/Not Good/ClassB

f(x)>=0.5?

# Half Plane and Half Spaces

- Half plane is a region on one side of an infinite long line, and does not contain any points from other side

- Half space n-dimensional space obtained by removing points on one side of hyperplane (n-1 dimension)
  - What would it look like for a 3 dimensional space

Diabetes/Happy/Good/ClassA

Hepatities/Sad/Not Good/ClassB

# Linear Discriminant Functions

- A linear discriminant function is defined by

$$f(x) = w^T x + w_0$$

  - where 'w' is the weight vector and $w_0$ is bias

- For a binary classifier

  Decision Surface $f(x) = 0$

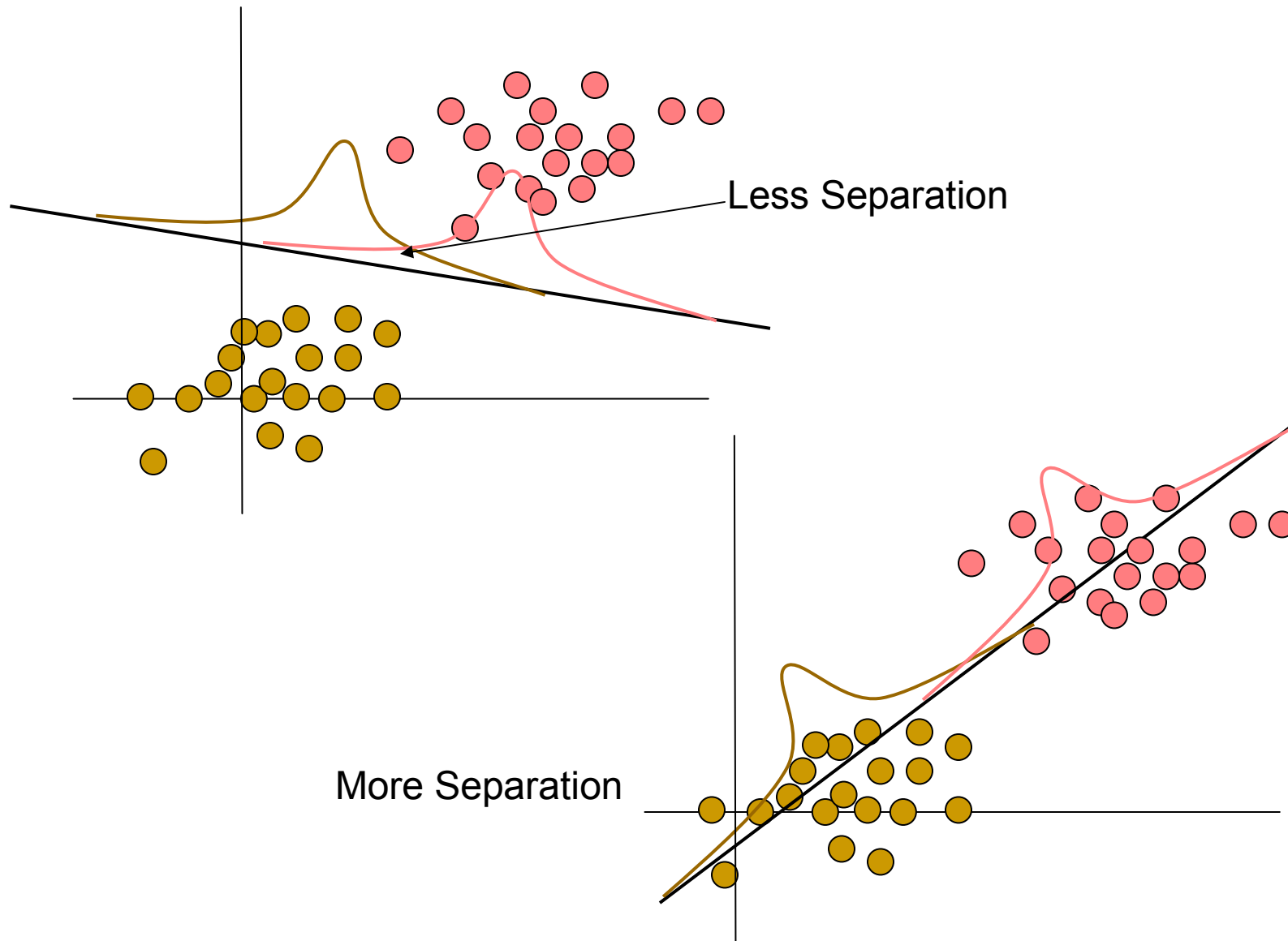  Class $C_0$ if $f(x) > 0$

  $Class C_1$ if $f(x) < 0$

# Decision Surface

→ Linear Discriminant function divides the space features live in
by a hyperplane decision surface
→Orientation is defined by normal vector (**w**) to the decision surface
→Location defined by offset



$f(x)>0$

$f(x)=0$

$f(x)<0$

# How Do We Find the Decision Surface

- We want to find a decision surface that will produce least class error in the training data
  - Fisher's Linear Discriminant
    - Dimensionality reduction, project data on a line and classify
  - Linear Discrimination with a hyperplane in (d-1) dimension
    - Perceptrons
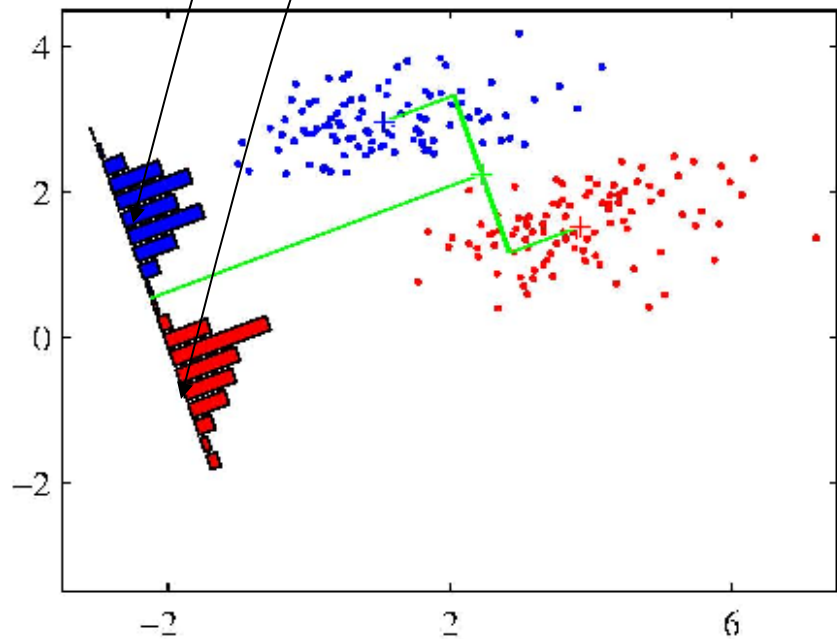
# Varying Levels of Separation



Less Separation

More Separation

# Fisher's Linear Discriminant Function

- We want to find the direction (w) of the decision surface such that points are well separated

  - Project points to a line

  - Compute mean and variances for the classes

  - Maximize

    $$J(w) = \frac{\text{square of separation of projected means}}{\text{Sum of within-class variance}}$$
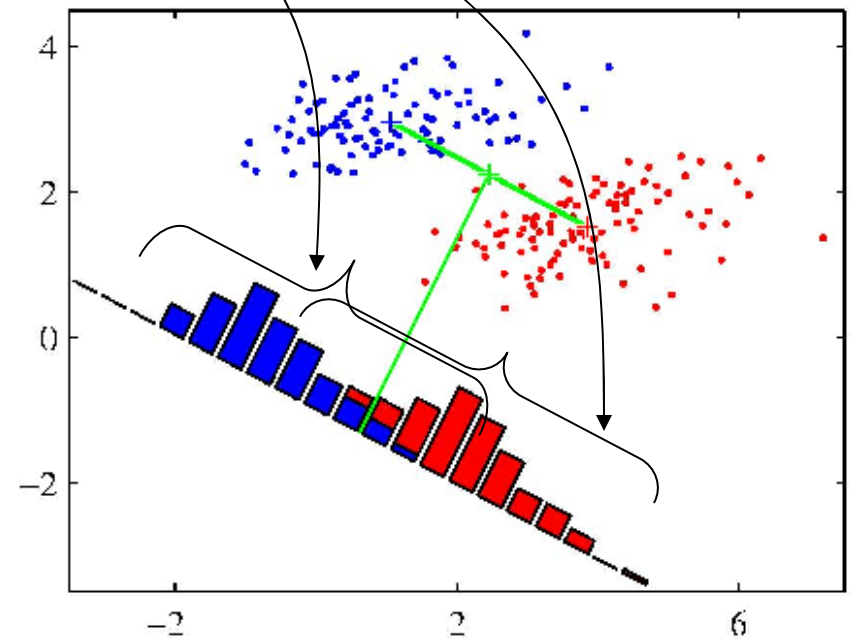
Maximize a function that will produce large separation between class means (projected) and has smaller within-class variance
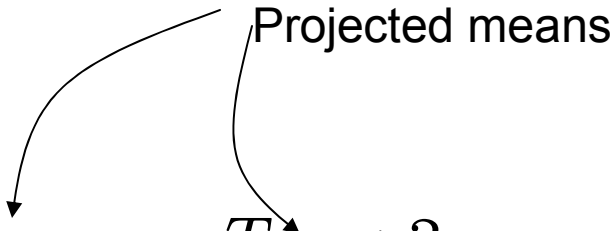
# Why This Criteria Makes Sense?

$$J(w) = \frac{\text{square of separation of projected means}}{\text{Sum of within-class variance}}$$



Projecting points on a line [1]

# Fisher's Linear Discriminant

Projected means

$$J(\mathbf{w}) = \frac{(\mathbf{w}^T \mu_1 - \mathbf{w}^T \mu_0)^2}{\mathbf{w}^T(\sum_1 + \sum_0)\mathbf{w}}$$

$$\mathbf{w} = (\sum_1 + \sum_0)^{-1}(\mu_1 - \mu_0)$$

Higher the interclass difference in means and
lower the in-class variance better
is the model

# Linear Discrimination with a Hyperplane

- Dimensionality reduction is one way of classification

- We can also try to find they discriminating hyperplane by reducing the total error in training
  - Perceptrons is one such algorithm

# Perceptron

- We want to find a function that would produce least training error

$$R_n(w) = \frac{1}{n} \sum_{i=1}^{n} Loss(y_i, f(x_i; w))$$

Can't we just find minima of the loss function by setting derivatives to zero?

# Minimizing Training Error

Given training data $< (x_i, y_i) >$
We want to find $w$ such that
$y_i(w.x_i) > 0$ if $y_i > 0$
$y_i(w.x_i) < 0$ if $y_i < 0$

- We can iterate over all points and adjust the parameters

$$w \leftarrow w + y_i x_i$$

$$\text{if } y \neq f(x_i; w)$$

- Parameters are updated only if the classifier makes a mistake

# Perceptron Algorithm

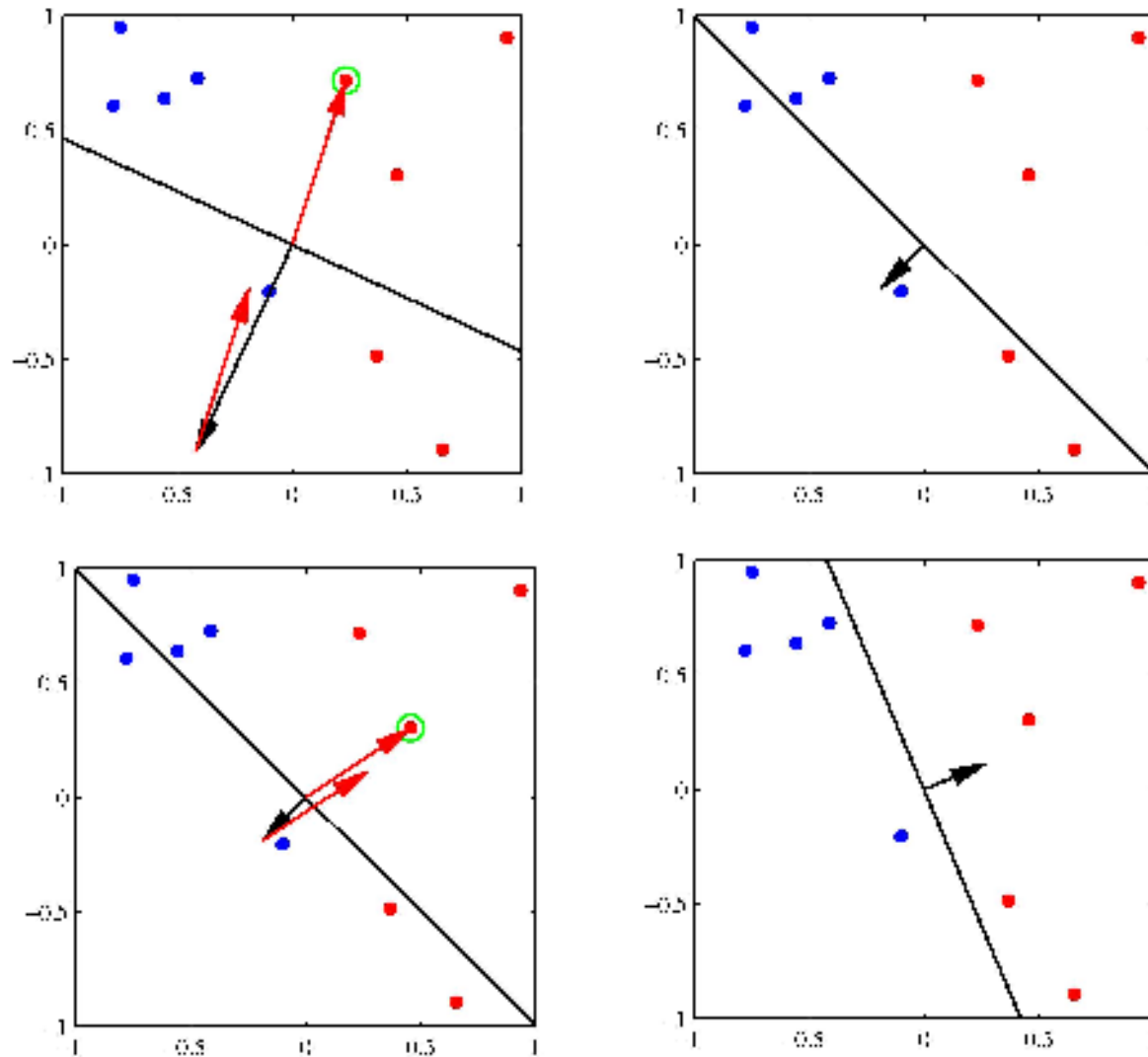We are given $(x_i, y_i)$

Initialize $w$

Do until converged

      if error$((y_i, f(x_i, w)) == TRUE)$
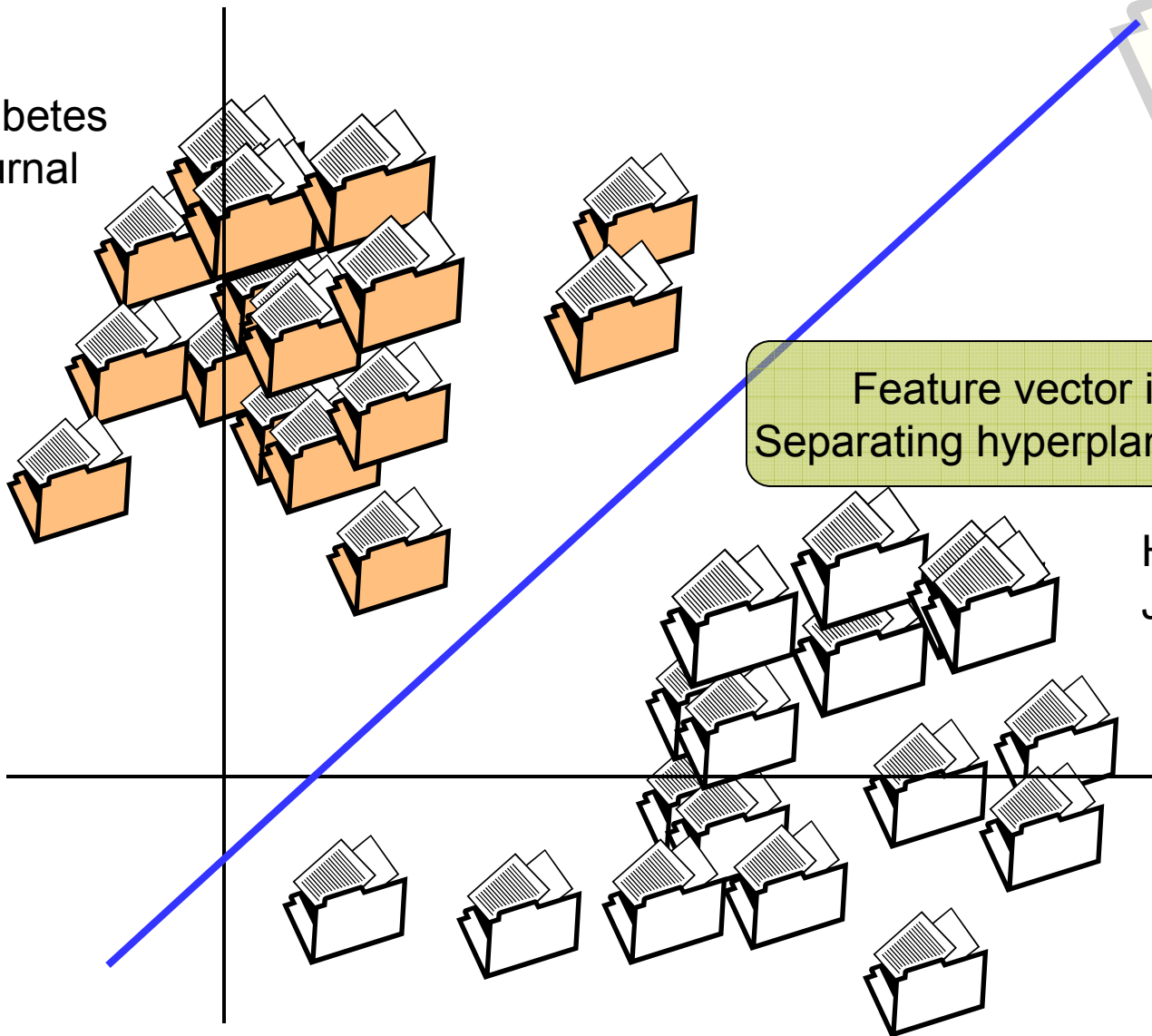
$$w \leftarrow w + y_i x_i$$

    end if

End do

# Why Algorithm Makes Sense?



Convergence illustration [1]

# Text Classification with Perceptron



Diabetes Journal

Hepatitis Journal

**?**

Which side of the hyperplane is this document?

Feature vector in D dimensions
Separating hyperplane in D-1 dimensions

# Text Classification with Perceptron

- Perceptron may not always converge

- Ok for two classes, not trivial to extend it to multiple classes

- Not the optimal hyperplane
  - Many hyperplanes that separates the class
  - Depends on random initialization

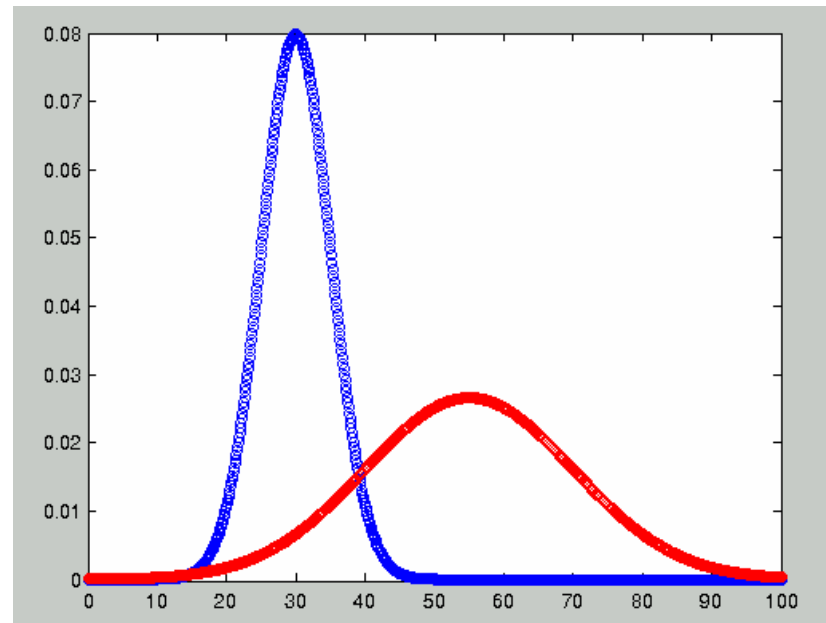# Generative vs. Discriminative

- **Generative Classifier**

  - Model joint probability p(x,y) where x are inputs and y are labels

  - Make prediction using Bayes rule to compute p(y|x)

- **Discriminative Classifier**

  - Try to predict output directly

  - Model p(y|x) directly

# Generative Classifier

- We can model class conditional densities using Gaussian distributions
- If we know class conditional densities
  - $p(x|y=C1)$
  - $p(x|y=C2)$
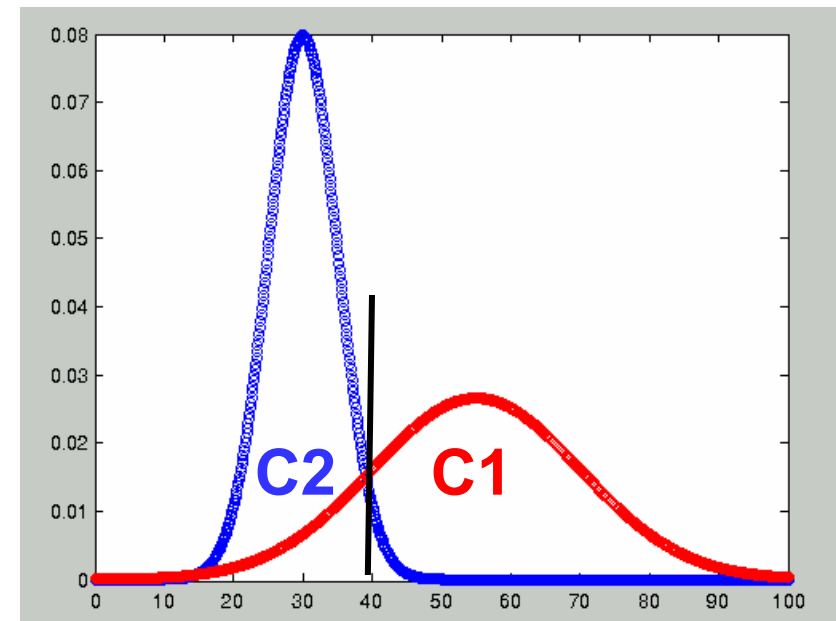- We can find a decision to classify the unseen example

# Bayes Rule

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

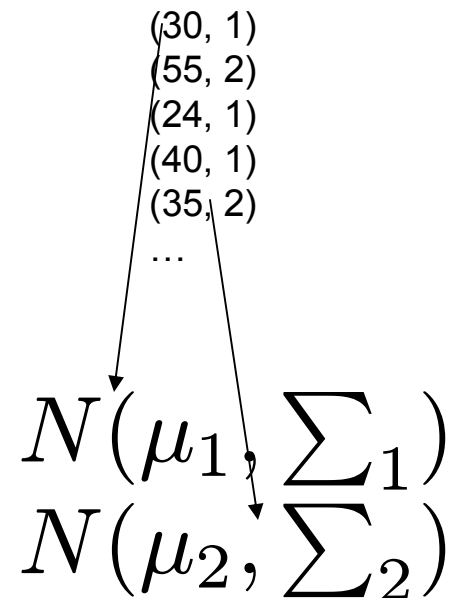→So how would this rule help in classifying text in two different categories; Diabetes vs Hepatitis

→Think about distribution of count of the word diabetes for example
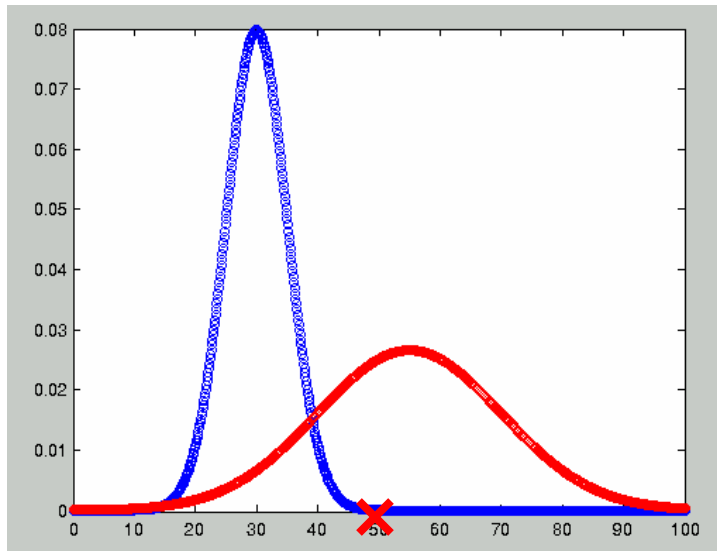
# Generative Classifier

- **If we have two classes C1 and C2**

- **We can estimate Gaussian distribution of the features for both classes**

  - Let's say we have a feature x

    - x = length of a document

  - And class label (y)

    - y = 1 diabetes or 2 hepatitis

(30, 1)
(55, 2)
(24, 1)
(40, 1)
(35, 2)
...

Find out $\mu_i$ and $\sum_i$ from data for both classes

$$N(\mu_1, \sum{}_1)$$
$$N(\mu_2, \sum{}_2)$$

# Generative Classifier

- Given a new data point find out posterior probability from each class and take a log ratio

- If higher posterior probability for C1, it means new x better explained by the Gaussian distribution of C1



$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

$$p(y=1|x) \propto p(x|\mu_1, \textstyle\sum_1)p(y=1)$$

# Naïve Bayes Classifier

- **Naïve Bayes Classifier a type of Generative classifier**
  - Compute class-conditional distribution but with conditional independence assumption
- **Shown to be very useful for text categorization**

# Conditional Independence

- Given random variables X, Y,Z,  X is conditionally independent of Y given Z if and only if

$$P(X|Y, Z) = p(X|Z)$$

$$
\begin{aligned}
P(X|Y) &= P(X_1, X_2|Y) \\
&= P(X_1|X_2, Y)P(X_2|Y) \\
&= P(X_1|Y)P(X_2|Y)
\end{aligned}
$$

# Conditional Independence

- For a feature vector with 'n' features we get

$$P(X_1, X_2, ..., X_N | Y) = \Pi_{i=1}^{N} P(X_i | Y)$$

N features are conditionally independent of one another given Y

Why would this assumption help?

# Naïve Bayes Classifier for Text

$$P(Y_k, X_1, X_2, ..., X_N) = P(Y_k)\Pi_i P(X_i|Y_k)$$

Prior Probability of the Class

Conditional Probability of feature given the Class

Here N is the number of words, not to confuse with the total vocabulary size

# Naïve Bayes Classifier for Text

$$P(Y = y_k | X_1, X_2, ..., X_N) = \frac{P(Y=y_k)P(X_1,X_2,..,X_N|Y=y_k)}{\sum_j P(Y=y_j)P(X_1,X_2,..,X_N|Y=y_j)}$$

$$= \frac{P(Y=y_k)\Pi_i P(X_i|Y=y_k)}{\sum_j P(Y=y_j)\Pi_i P(X_i|Y=y_j)}$$

$$Y \leftarrow argmax_{y_k} P(Y = y_k)\Pi_i P(X_i|Y = y_k)$$

# Naïve Bayes Classifier for Text

- Given the training data what are the parameters to be estimated?

$$P(Y) \qquad P(X|Y_1) \qquad P(X|Y_2)$$

Diabetes : 0.8
Hepatitis : 0.2

the: 0.001
diabetic : 0.02
blood : 0.0015
sugar : 0.02
weight : 0.018
…

the: 0.001
diabetic : 0.0001
water : 0.0118
fever : 0.01
weight : 0.008
…

# References

- [1] Christopher M Bishop, "Pattern Recognition and Machine Learning" 2006