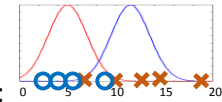


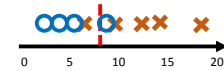
# Logistic Classifier

CISC 5800  
Professor Daniel Leeds

Classification strategy:  
generative vs. discriminative



- Generative, e.g., Bayes/Naïve Bayes:
  - Identify probability distribution for each class
  - Determine class with maximum probability for data example
- Discriminative, e.g., Logistic Regression:
  - Identify boundary between classes
  - Determine which side of boundary new data example exists on



2

## Linear algebra: data features

- Vector – list of numbers:  
each number describes  
a data **feature**

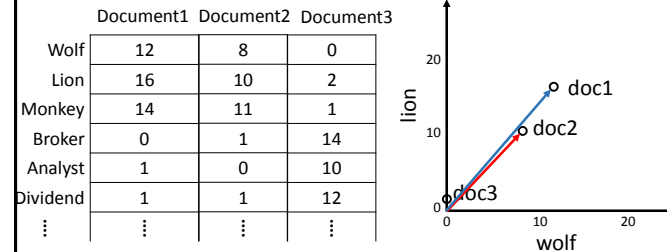
	Document 1	Document 2	Document 3
Wolf	12	8	0
Lion	16	10	2
Monkey	14	11	1
Broker	0	1	14
Analyst	1	0	10
Dividend	1	1	12
⋮	⋮	⋮	⋮

- Matrix – list of lists of numbers:  
features for each data  
point

3

## Feature space

- Each data feature defines a dimension in space

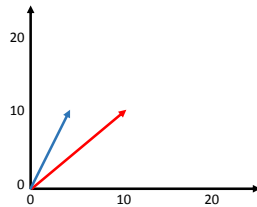


4

## The dot product

The dot product compares two vectors:

$$\mathbf{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}, \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \quad \mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i = \mathbf{a}^T \mathbf{b}$$



$$\begin{bmatrix} 5 \\ 10 \end{bmatrix} \cdot \begin{bmatrix} 10 \\ 10 \end{bmatrix} =$$

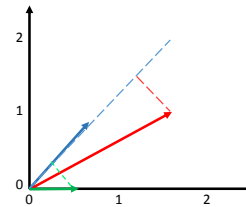
5

The dot product, continued  $\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i$

Magnitude of a vector is sum of squares of the elements

$$|\mathbf{a}| = \sqrt{\sum_i a_i^2}$$

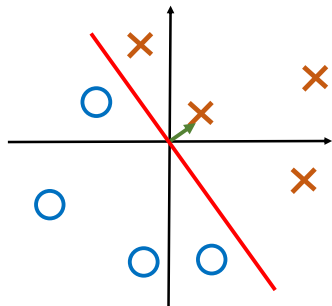
If  $\mathbf{a}$  has unit magnitude,  $\mathbf{a} \cdot \mathbf{b}$  is “projection” of  $\mathbf{b}$  onto  $\mathbf{a}$



$$\begin{bmatrix} 0.6 \\ 0.8 \end{bmatrix} \cdot \begin{bmatrix} 1.5 \\ 1 \end{bmatrix} =$$

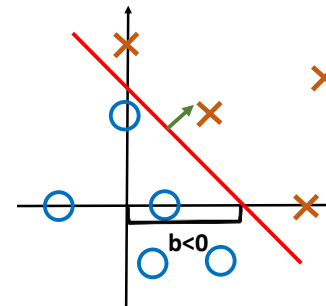
7

## Separating boundary, defined by $\mathbf{w}$



- Separating **hyperplane** splits **class 0** and **class 1**
- Plane is defined by line  $\mathbf{w}$  perpendicular to plane
- Is data point  $\mathbf{x}$  in class 0 or class 1?  $\mathbf{w}^T \mathbf{x} + b > 0$  class **1**  
 $\mathbf{w}^T \mathbf{x} + b < 0$  class **0**

## Separating boundary, defined by $\mathbf{w}$ and $b$



**Example:**

$$\mathbf{w} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b = -4$$

$$\mathbf{x}^1 = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$$

$$\mathbf{x}^2 = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$

11

## Notational simplification

Recall:  $\mathbf{w}^T \mathbf{x} = \mathbf{w} \cdot \mathbf{x} = \sum_{i=1}^n w_i x_i$

Define  $x'_{1:n} = x_{1:n}$  and  $x'_{n+1} = 1$  for all inputs  $\mathbf{x}$  and  $w'_{1:n} = w_{1:n}$  and  $w'_{n+1} = b$

Now  $\mathbf{w}'^T \mathbf{x}' = \mathbf{w}^T \mathbf{x} + b$

Let's assume  $x_{n+1}=1$  always, and  $w_{n+1}=b$  always

## From real-number projection to 0/1 label

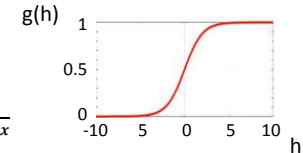
Binary classification: 0 is class A, 1 is class B

Sigmoid function stands in for  $p(\mathbf{x}|y)$

Sigmoid:  $g(h) = \frac{1}{1+e^{-h}}$

$$p(y=0|x; \theta) = 1 - g(\mathbf{w}^T \mathbf{x}) = \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

$$p(y=1|x; \theta) = g(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$



$$\mathbf{w}^T \mathbf{x} = \sum_j w_j x_j + b$$

## Learning parameters for classification

Similar to MLE for Bayes classifier

"Likelihood" for data points  $y^1, \dots, y^n$   
(different from Bayesian likelihood)

- If  $y^i$  in class A,  $y^i=0$ , multiply  $(1-g(\mathbf{x}^i; \mathbf{w}))$
- If  $y^i$  in class B,  $y^i=1$ , multiply  $(g(\mathbf{x}^i; \mathbf{w}))$

$$\operatorname{argmax}_w L(y|x; w) = \prod_i (1 - g(\mathbf{x}^i; \mathbf{w}))^{(1-y^i)} g(\mathbf{x}^i; \mathbf{w})^{y^i}$$

15