

DATA MINING

LECTURE 11

Classification

Support Vector Machines

Logistic Regression

Naïve Bayes Classifier

Supervised Learning

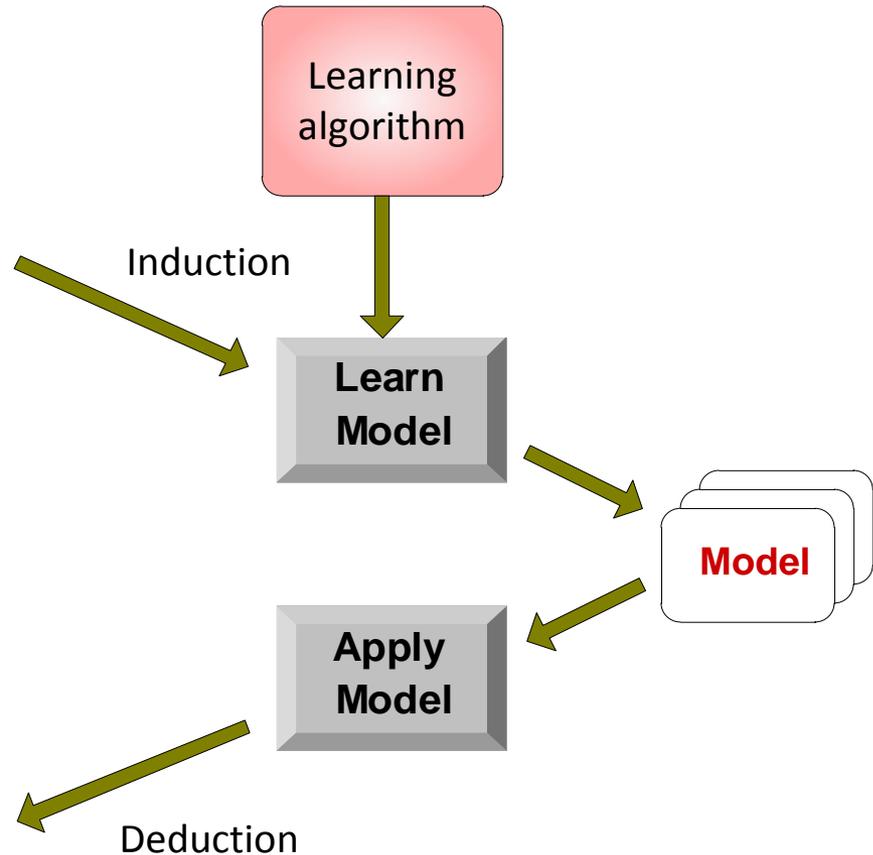
Illustrating Classification Task

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

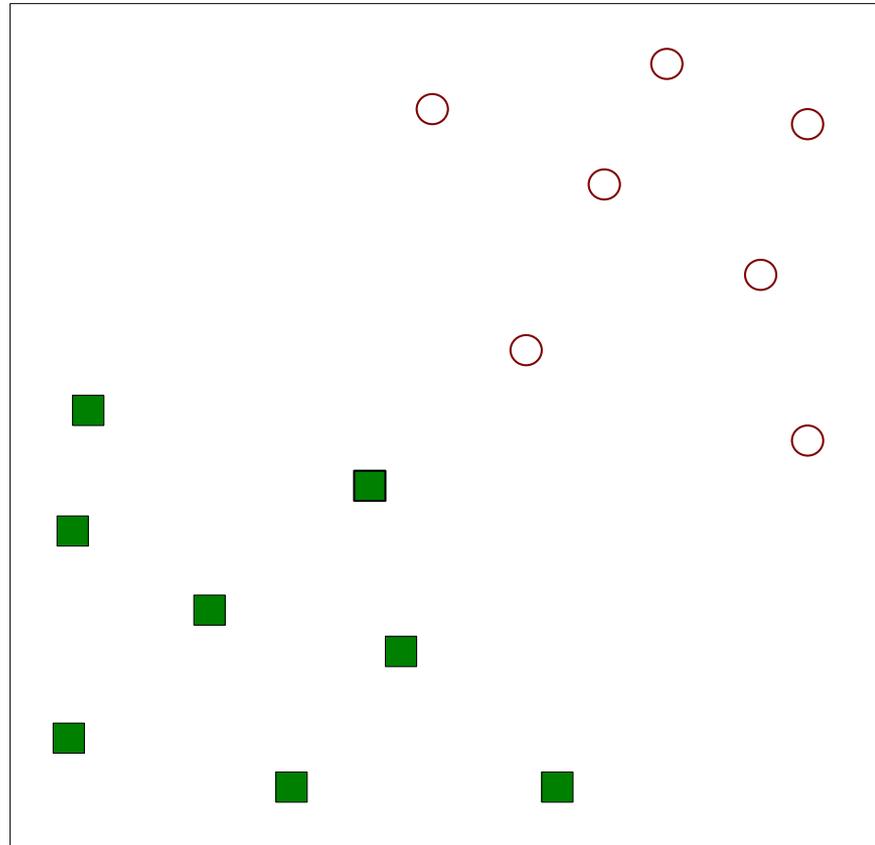
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



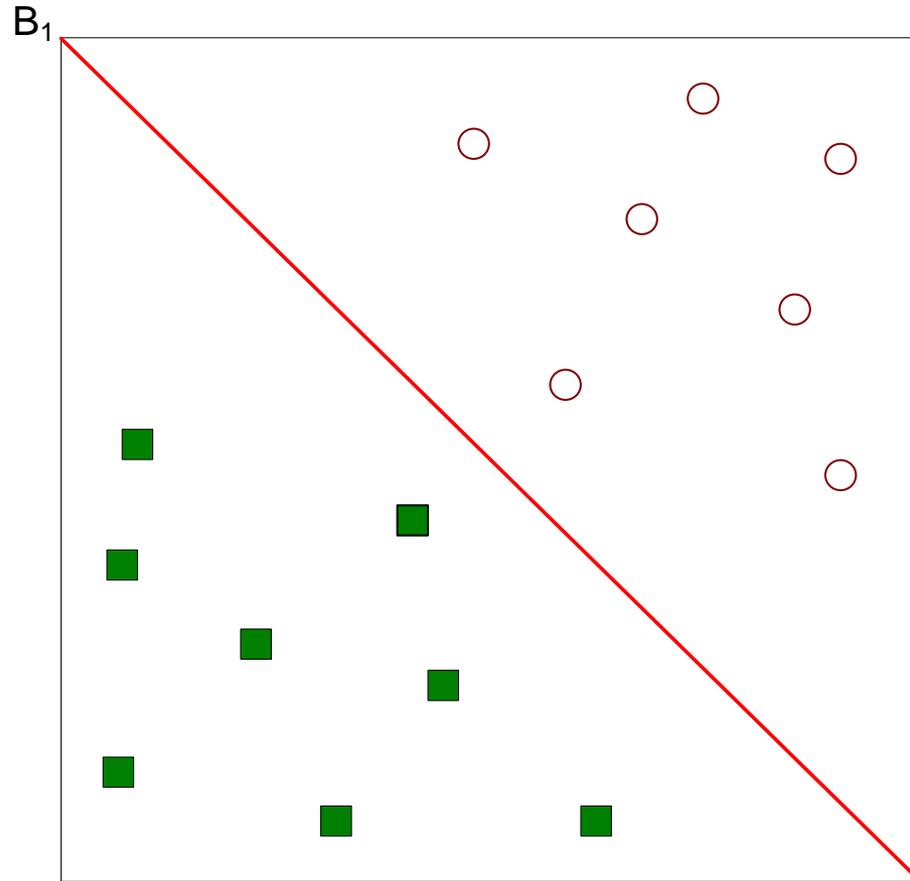
SUPPORT VECTOR MACHINES

Support Vector Machines



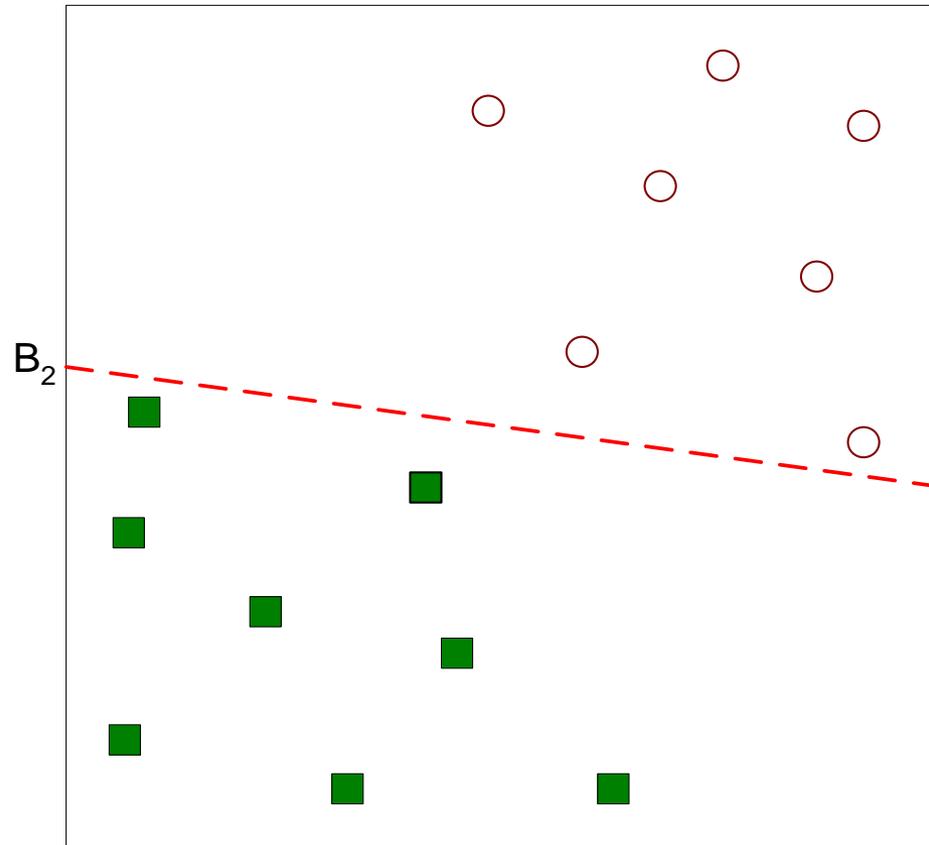
- Find a linear hyperplane (decision boundary) that will separate the data

Support Vector Machines



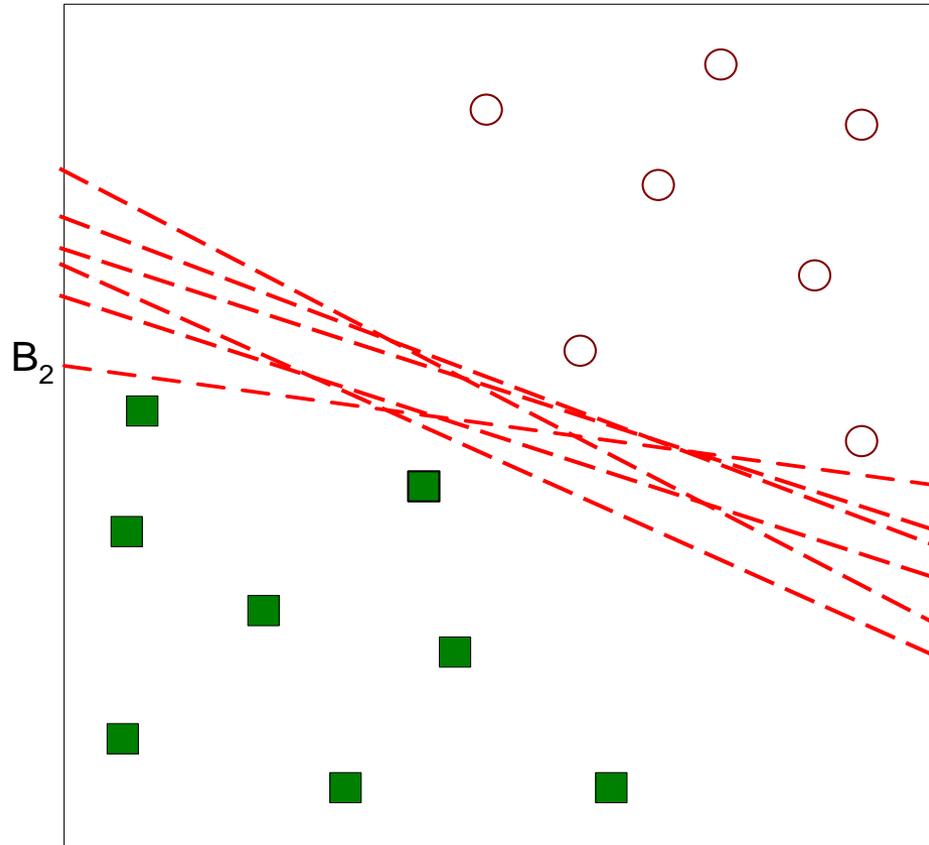
- One Possible Solution

Support Vector Machines



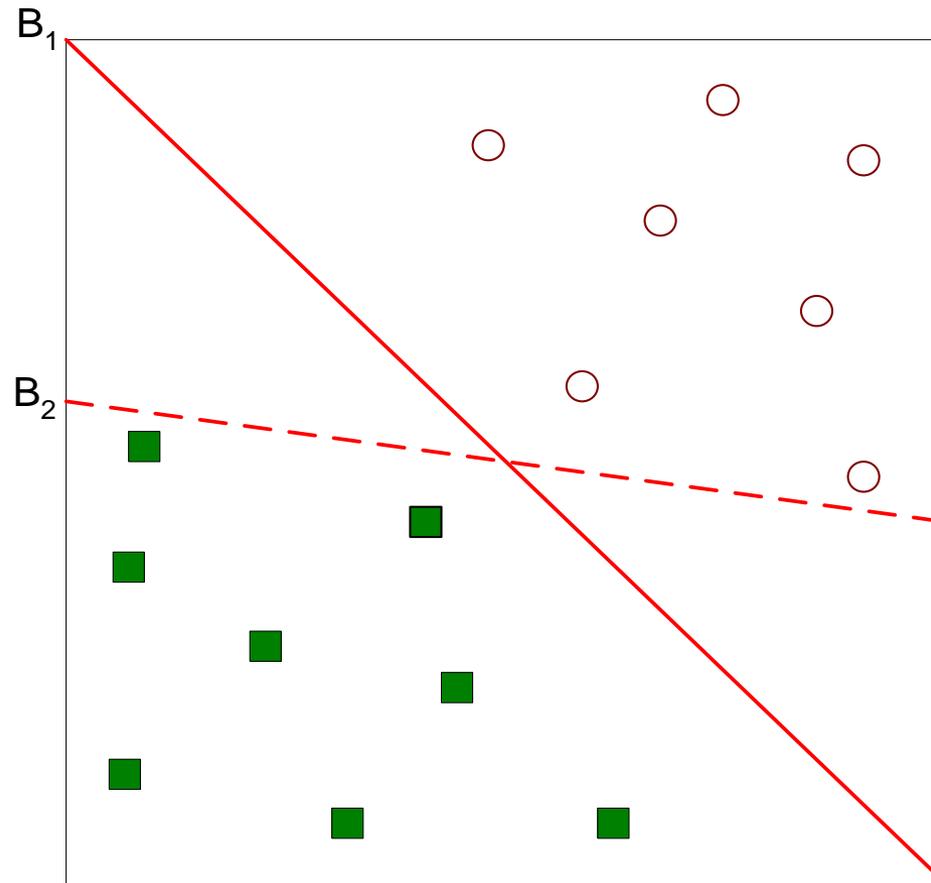
- Another possible solution

Support Vector Machines



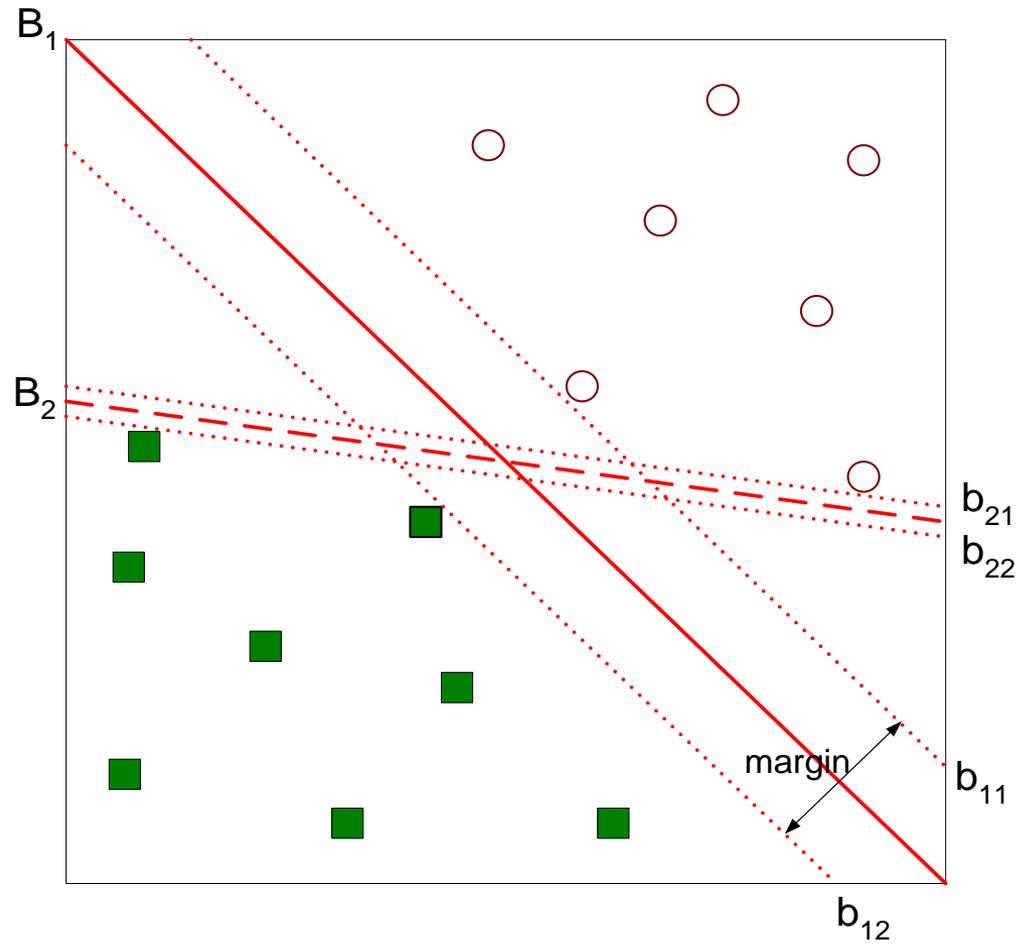
- Other possible solutions

Support Vector Machines



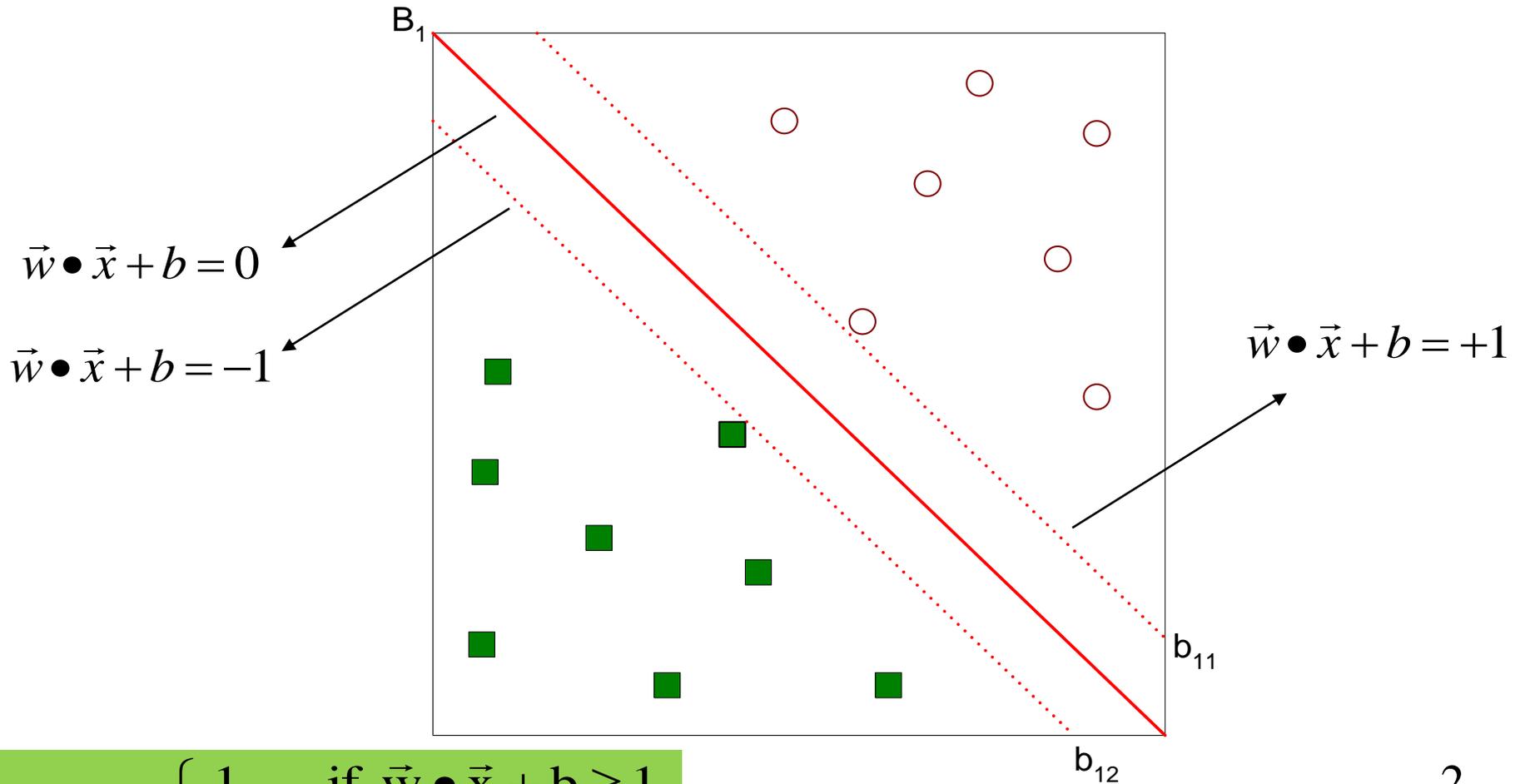
- Which one is better? B_1 or B_2 ?
- How do you define better?

Support Vector Machines



- Find hyperplane **maximizes** the margin => B1 is better than B2

Support Vector Machines



$$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x} + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x} + b \leq -1 \end{cases}$$

$$\text{Margin} = \frac{2}{\|\vec{w}\|}$$

Support Vector Machines

- We want to **maximize**: $\text{Margin} = \frac{2}{\|\vec{w}\|}$
- Which is equivalent to **minimizing**: $L(w) = \frac{\|\vec{w}\|^2}{2}$
- But subjected to the following **constraints**:

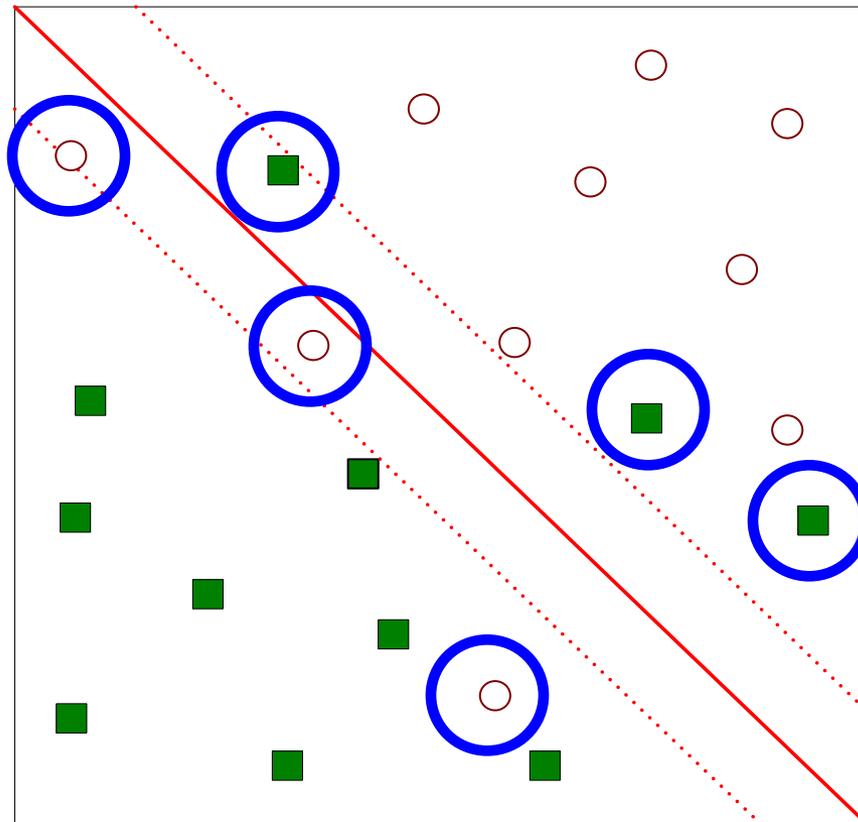
$$\vec{w} \cdot \vec{x}_i + b \geq 1 \text{ if } y_i = 1$$

$$\vec{w} \cdot \vec{x}_i + b \leq -1 \text{ if } y_i = -1$$

- This is a **constrained optimization problem**
 - Numerical approaches to solve it (e.g., **quadratic programming**)

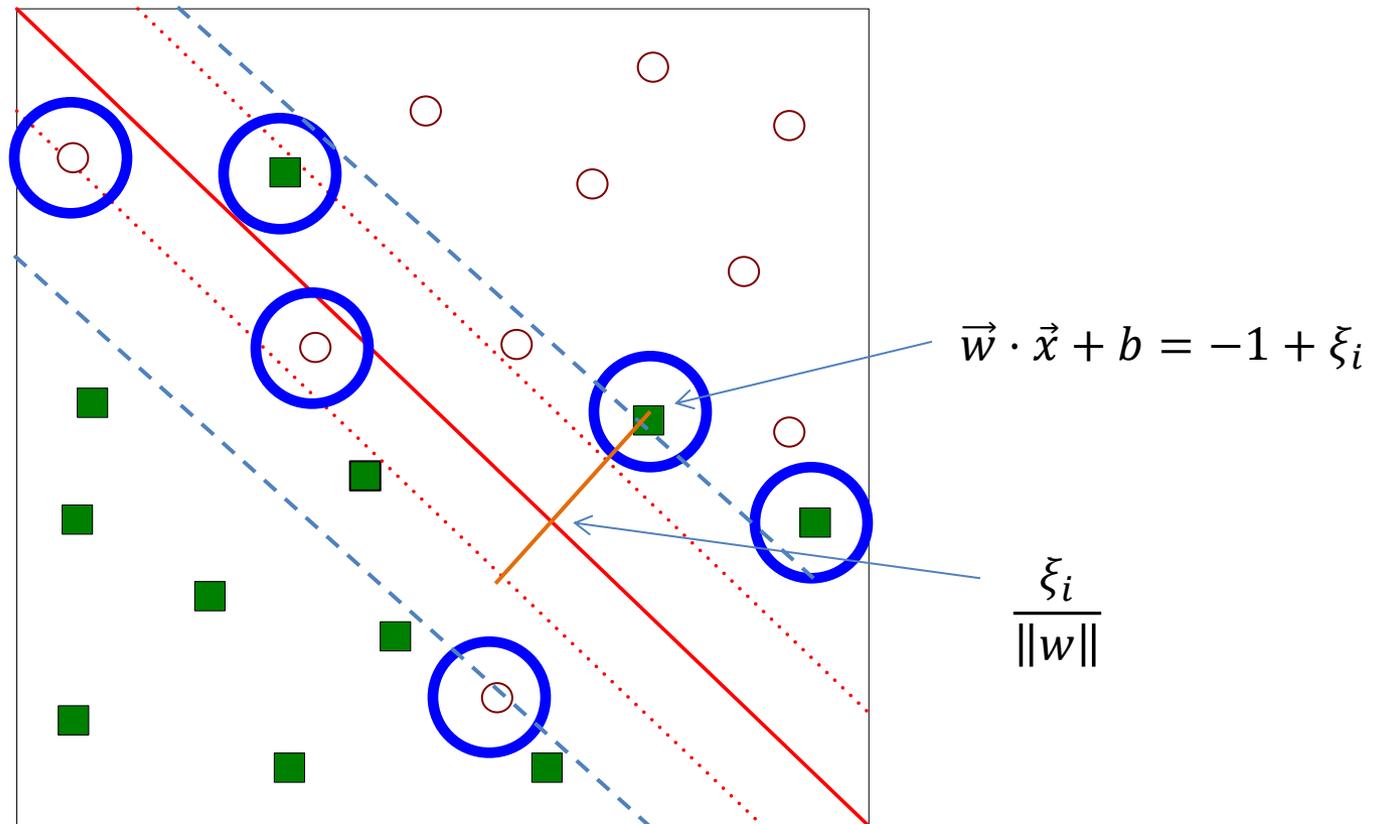
Support Vector Machines

- What if the problem is **not linearly separable**?



Support Vector Machines

- What if the problem is not linearly separable?



Support Vector Machines

- What if the problem is not linearly separable?
 - Introduce slack variables

- Need to minimize:

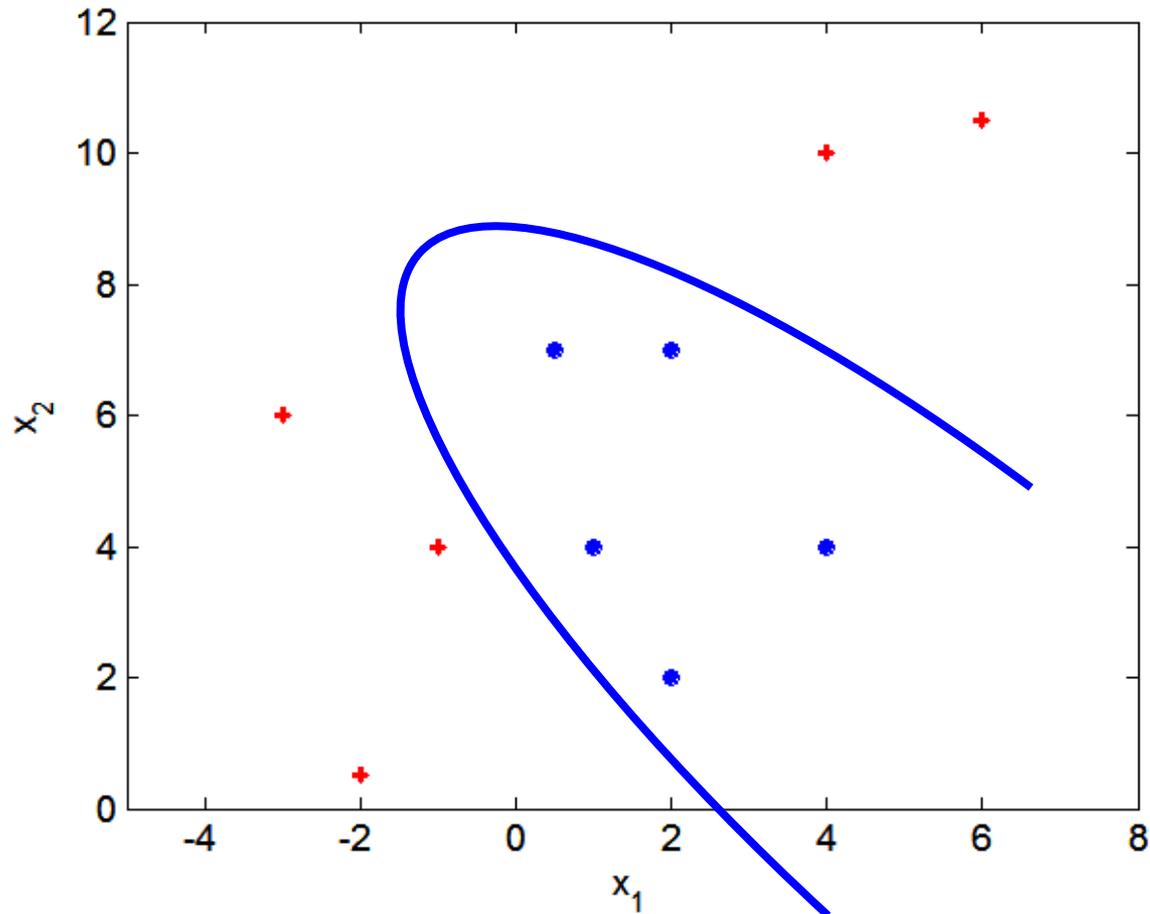
$$L(w) = \frac{\|\vec{w}\|^2}{2} + C \left(\sum_{i=1}^N \xi_i \right)$$

- Subject to:

$$\begin{aligned} \vec{w} \cdot \vec{x}_i + b &\geq 1 - \xi_i \text{ if } y_i = 1 \\ \vec{w} \cdot \vec{x}_i + b &\leq -1 + \xi_i \text{ if } y_i = -1 \end{aligned}$$

Nonlinear Support Vector Machines

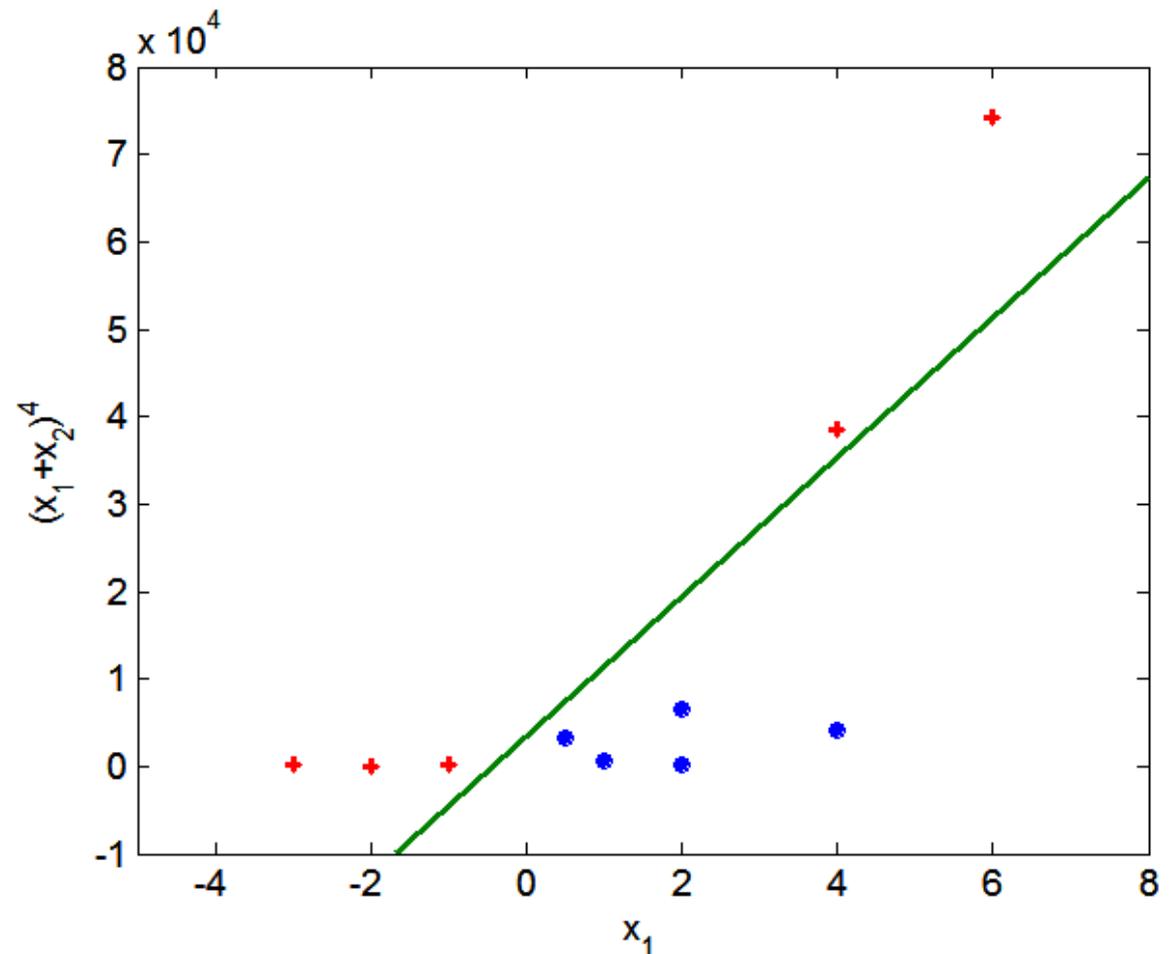
- What if decision boundary is not linear?



Nonlinear Support Vector Machines

- Transform data into higher dimensional space

Use the **Kernel Trick**



LOGISTIC REGRESSION

Classification via regression

- Instead of predicting the **class** of an record we want to **predict the probability of the class** given the record
- The problem of **predicting continuous values** is called **regression** problem
- General approach: find a continuous function that models the continuous points.

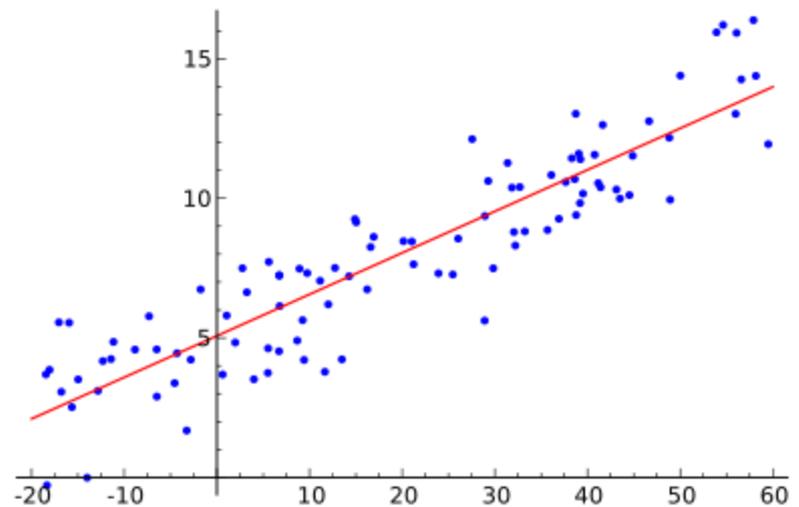
Example: Linear regression

- Given a dataset of the form $\{(x_1, y_1), \dots, (x_n, y_n)\}$ find a linear function that given the vector x_i predicts the y_i value as $y'_i = w^T x_i$

- Find a vector of weights w that **minimizes the sum of square errors**

$$\sum_i (y'_i - y_i)^2$$

- Several techniques for solving the problem.



Classification via regression

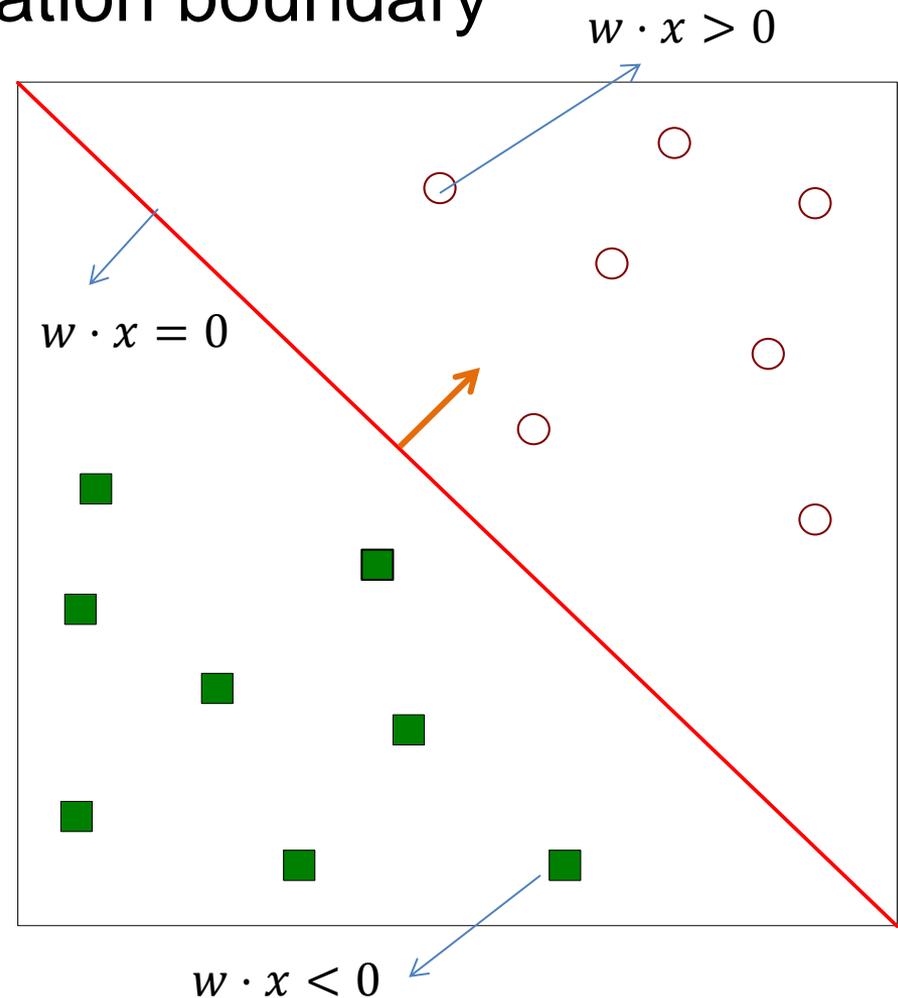
- Assume a linear classification boundary

For the positive class the **bigger** the **value of $w \cdot x$** , the further the point is from the classification boundary, the higher our **certainty** for the membership to the **positive class**

- Define $P(C_+|x)$ as an **increasing** function of $w \cdot x$

For the negative class the **smaller** the **value of $w \cdot x$** , the further the point is from the classification boundary, the higher our **certainty** for the membership to the **negative class**

- Define $P(C_-|x)$ as a **decreasing** function of $w \cdot x$



Logistic Regression

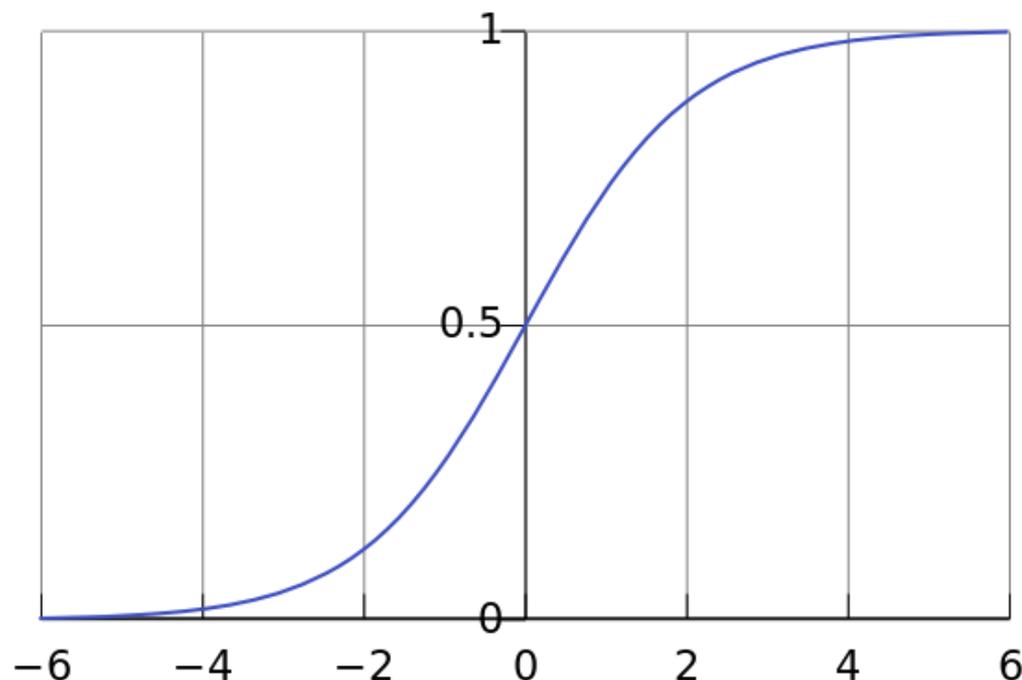
The **logistic function**

$$f(t) = \frac{1}{1 + e^{-t}}$$

$$P(C_+|x) = \frac{1}{1 + e^{-w \cdot x}}$$

$$P(C_-|x) = \frac{e^{-w \cdot x}}{1 + e^{-w \cdot x}}$$

$$\log \frac{P(C_+|x)}{P(C_-|x)} = w \cdot x$$



Logistic Regression: Find the vector w that **maximizes the probability** of the observed data

Linear regression on the **log-odds ratio**

Logistic Regression

- Produces a **probability estimate** for the **class membership** which is often very useful.
- The **weights** can be useful for understanding the **feature importance**.
- Works for relatively large datasets
- Fast to apply.

NAÏVE BAYES CLASSIFIER

Bayes Classifier

- A probabilistic framework for solving classification problems
- **A, C** random variables
- **Joint** probability: **$\Pr(A=a, C=c)$**
- **Conditional** probability: **$\Pr(C=c | A=a)$**
- Relationship between joint and conditional probability distributions

$$\Pr(C, A) = \Pr(C | A) \times \Pr(A) = \Pr(A | C) \times \Pr(C)$$

- **Bayes Theorem:**
$$P(C | A) = \frac{P(A | C)P(C)}{P(A)}$$

Bayesian Classifiers

- How to classify the new record $X = (\text{'Yes'}, \text{'Single'}, 80\text{K})$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Find the class with the highest probability given the vector values.

Maximum A posteriori Probability estimate:

- Find the value c for class C that maximizes $P(C=c | X)$

How do we estimate $P(C|X)$ for the different values of C ?

- We want to estimate $P(C=\text{Yes} | X)$
- and $P(C=\text{No} | X)$

Bayesian Classifiers

- In order for probabilities to be well defined:
 - Consider each attribute and the class label as **random variables**
 - Probabilities are determined from the data

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Evade C

Event space: {Yes, No}
 $P(C) = (0.3, 0.7)$

Refund A_1

Event space: {Yes, No}
 $P(A_1) = (0.3, 0.7)$

Marital Status A_2

Event space: {Single, Married, Divorced}
 $P(A_2) = (0.4, 0.4, 0.2)$

Taxable Income A_3

Event space: R
 $P(A_3) \sim \text{Normal}(\mu, \sigma^2)$
 $\mu = 104$:sample mean, $\sigma^2 = 1874$:sample var

Bayesian Classifiers

- Approach:
 - compute the posterior probability $P(C | A_1, A_2, \dots, A_n)$ using the Bayes theorem

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- Maximizing $P(C | A_1, A_2, \dots, A_n)$ is equivalent to maximizing $P(A_1, A_2, \dots, A_n | C) P(C)$
 - The value $P(A_1, \dots, A_n)$ is the same for all values of C .
- How to estimate $P(A_1, A_2, \dots, A_n | C)$?

Naïve Bayes Classifier

- Assume **conditional independence** among attributes A_i when class C is given:
 - $P(A_1, A_2, \dots, A_n | C) = P(A_1 | C) P(A_2 | C) \cdots P(A_n | C)$
 - We can estimate $P(A_i | C)$ from the data.
 - New point $X = (A_1 = \alpha_1, \dots, A_n = \alpha_n)$ is classified to class c if

$$P(C = c | X) = P(C = c) \prod_i P(A_i = \alpha_i | c)$$

is maximum over all possible values of C .

Example

- Record

$X = (\text{Refund} = \text{Yes}, \text{Status} = \text{Single}, \text{Income} = 80\text{K})$

- For the class $C = \text{'Evade'}$, we want to compute:

$P(C = \text{Yes}|X)$ and $P(C = \text{No}| X)$

- We compute:

- $P(C = \text{Yes}|X) = P(C = \text{Yes}) * P(\text{Refund} = \text{Yes} | C = \text{Yes})$

- * $P(\text{Status} = \text{Single} | C = \text{Yes})$

- * $P(\text{Income} = 80\text{K} | C = \text{Yes})$

- $P(C = \text{No}|X) = P(C = \text{No}) * P(\text{Refund} = \text{Yes} | C = \text{No})$

- * $P(\text{Status} = \text{Single} | C = \text{No})$

- * $P(\text{Income} = 80\text{K} | C = \text{No})$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Class Prior Probability:

$$P(C = c) = \frac{N_c}{N}$$

N_c : Number of records with class c

N = Number of records

$$P(C = \text{No}) = 7/10$$

$$P(C = \text{Yes}) = 3/10$$

How to Estimate Probabilities from Data?

Discrete attributes:

$$P(A_i = a | C = c) = \frac{N_{a,c}}{N_c}$$

$N_{a,c}$: number of instances having attribute $A_i = a$ and belong to class c

N_c : number of instances of class c

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Discrete attributes:

$$P(A_i = a | C = c) = \frac{N_{a,c}}{N_c}$$

$N_{a,c}$: number of instances having attribute $A_i = a$ and belong to class c

N_c : number of instances of class c

$$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Discrete attributes:

$$P(A_i = a | C = c) = \frac{N_{a,c}}{N_c}$$

$N_{a,c}$: number of instances having attribute $A_i = a$ and belong to class c

N_c : number of instances of class c

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Discrete attributes:

$$P(A_i = a | C = c) = \frac{N_{a,c}}{N_c}$$

$N_{a,c}$: number of instances having attribute $A_i = a$ and belong to class c

N_c : number of instances of class c

$$P(\text{Status}=\text{Single}|\text{No}) = 2/7$$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Discrete attributes:

$$P(A_i = a | C = c) = \frac{N_{a,c}}{N_c}$$

$N_{a,c}$: number of instances having attribute $A_i = a$ and belong to class c

N_c : number of instances of class c

$$P(\text{Status}=\text{Single}|\text{Yes}) = 2/3$$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i = a | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(a-\mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (a, c) pair
- For **Class=No**
 - sample mean $\mu = 110$
 - sample variance $\sigma^2 = 2975$
- For **Income = 80**

$$P(\text{Income} = 80 | \text{No}) = \frac{1}{\sqrt{2\pi(54.54)}} e^{-\frac{(80-110)^2}{2(2975)}} = 0.0062$$

How to Estimate Probabilities from Data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i = a | c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(a-\mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (a, c) pair
- For **Class=Yes**
 - sample mean $\mu = 90$
 - sample variance $\sigma^2 = 2975$
- For **Income = 80**

$$P(\text{Income} = 80 | \text{Yes}) = \frac{1}{\sqrt{2\pi(5)}} e^{-\frac{(80-90)^2}{2(25)}} = 0.01$$

Example of Naïve Bayes Classifier

- Creating a Naïve Bayes Classifier, essentially means to compute **counts**:

Total number of records: $N = 10$

Class No:

Number of records: 7

Attribute Refund:

Yes: 3

No: 4

Attribute Marital Status:

Single: 2

Divorced: 1

Married: 4

Attribute Income:

mean: 110

variance: 2975

Class Yes:

Number of records: 3

Attribute Refund:

Yes: 0

No: 3

Attribute Marital Status:

Single: 2

Divorced: 1

Married: 0

Attribute Income:

mean: 90

variance: 25

Example of Naïve Bayes Classifier

Given a Test Record:

$X = (\text{Refund} = \text{Yes}, \text{Status} = \text{Single}, \text{Income} = 80\text{K})$

naive Bayes Classifier:

$$P(\text{Refund}=\text{Yes} | \text{No}) = 3/7$$

$$P(\text{Refund}=\text{No} | \text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes} | \text{Yes}) = 0$$

$$P(\text{Refund}=\text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status}=\text{Single} | \text{No}) = 2/7$$

$$P(\text{Marital Status}=\text{Divorced} | \text{No}) = 1/7$$

$$P(\text{Marital Status}=\text{Married} | \text{No}) = 4/7$$

$$P(\text{Marital Status}=\text{Single} | \text{Yes}) = 2/7$$

$$P(\text{Marital Status}=\text{Divorced} | \text{Yes}) = 1/7$$

$$P(\text{Marital Status}=\text{Married} | \text{Yes}) = 0$$

For taxable income:

If class=No: sample mean=110
sample variance=2975

If class=Yes: sample mean=90
sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{Yes}|\text{Class}=\text{No})$
 $\times P(\text{Married} | \text{Class}=\text{No})$
 $\times P(\text{Income}=120\text{K} | \text{Class}=\text{No})$
 $= 3/7 * 2/7 * 0.0062 = 0.00075$
 - $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No} | \text{Class}=\text{Yes})$
 $\times P(\text{Married} | \text{Class}=\text{Yes})$
 $\times P(\text{Income}=120\text{K} | \text{Class}=\text{Yes})$
 $= 0 * 2/3 * 0.01 = 0$
 - $P(\text{No}) = 0.3, P(\text{Yes}) = 0.7$
- Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$
Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
 $\Rightarrow \text{Class} = \text{No}$

Naïve Bayes Classifier

- If one of the conditional probability is **zero**, then the entire expression becomes zero
- **Laplace Smoothing:**

$$P(A_i = a | C = c) = \frac{N_{ac} + 1}{N_c + N_i}$$

- N_i : number of attribute **values** for attribute A_i

Example of Naïve Bayes Classifier

Given a Test Record:

With Laplace Smoothing

$X = (\text{Refund} = \text{Yes}, \text{Status} = \text{Single}, \text{Income} = 80\text{K})$

naive Bayes Classifier:

$$P(\text{Refund}=\text{Yes} | \text{No}) = 4/9$$

$$P(\text{Refund}=\text{No} | \text{No}) = 5/9$$

$$P(\text{Refund}=\text{Yes} | \text{Yes}) = 1/5$$

$$P(\text{Refund}=\text{No} | \text{Yes}) = 4/5$$

$$P(\text{Marital Status}=\text{Single} | \text{No}) = 3/10$$

$$P(\text{Marital Status}=\text{Divorced} | \text{No}) = 2/10$$

$$P(\text{Marital Status}=\text{Married} | \text{No}) = 5/10$$

$$P(\text{Marital Status}=\text{Single} | \text{Yes}) = 3/6$$

$$P(\text{Marital Status}=\text{Divorced} | \text{Yes}) = 2/6$$

$$P(\text{Marital Status}=\text{Married} | \text{Yes}) = 1/6$$

For taxable income:

If class=No: sample mean=110
sample variance=2975

If class=Yes: sample mean=90
sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \times P(\text{Married} | \text{Class}=\text{No}) \times P(\text{Income}=120\text{K} | \text{Class}=\text{No}) = 4/9 \times 3/10 \times 0.0062 = 0.00082$
 - $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No} | \text{Class}=\text{Yes}) \times P(\text{Married} | \text{Class}=\text{Yes}) \times P(\text{Income}=120\text{K} | \text{Class}=\text{Yes}) = 1/5 \times 3/6 \times 0.01 = 0.001$
 - $P(\text{No}) = 0.7, P(\text{Yes}) = 0.3$
 - $P(X|\text{No})P(\text{No}) = 0.0005$
 - $P(X|\text{Yes})P(\text{Yes}) = 0.0003$
- $\Rightarrow \text{Class} = \text{No}$

Implementation details

- Computing the conditional probabilities involves multiplication of many very small numbers
 - Numbers get very close to zero, and there is a danger of numeric instability
- We can deal with this by computing the **logarithm** of the conditional probability

$$\begin{aligned}\log P(C|A) &\sim \log P(A|C) + \log P(C) \\ &= \sum_i \log P(A_i|C) + \log P(C)\end{aligned}$$

Naïve Bayes for Text Classification

- Naïve Bayes is commonly used for **text classification**
- For a document with k terms $d = (t_1, \dots, t_k)$

Fraction of documents in c

$$P(c|d) = P(c)P(d|c) = P(c) \prod_{t_i \in d} P(t_i|c)$$

- $P(t_i|c)$ = Fraction of terms from **all documents** in c that are t_i .

Number of times t_i appears in some document in c

$$P(t_i|c) = \frac{N_{ic} + 1}{N_c + T}$$

Laplace Smoothing

Total number of terms in all documents in c

Number of unique words (vocabulary size)

- Easy to implement and works relatively well
- **Limitation:** Hard to incorporate **additional features** (beyond words).
 - E.g., number of adjectives used.

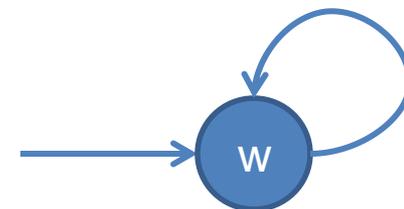
Multinomial document model

- Probability of document $d = (t_1, \dots, t_k)$ in class c :

$$P(d|c) = P(c) \prod_{t_i \in d} P(t_i|c)$$

- This formula assumes a **multinomial distribution** for the document generation:
 - If we have probabilities p_1, \dots, p_T for events t_1, \dots, t_T the probability of a subset of these is

$$P(d) = \frac{N}{N_{t_1}! N_{t_2}! \dots N_{t_T}!} p_1^{N_{t_1}} p_2^{N_{t_2}} \dots p_T^{N_{t_T}}$$



- Equivalently: There is an **automaton** spitting words from the above distribution

```

TRAINMULTINOMIALNB(C, D)
1  V ← EXTRACTVOCABULARY(D)
2  N ← COUNTDOCS(D)
3  for each c ∈ C
4  do Nc ← COUNTDOCSINCLASS(D, c)
5     prior[c] ← Nc/N
6     textc ← CONCATENATETEXTOFALLDOCSINCLASS(D, c)
7     for each t ∈ V
8     do Tct ← COUNTTOKENSOFTERM(textc, t)
9     for each t ∈ V
10    do condprob[t][c] ←  $\frac{T_{ct}+1}{\sum_{t'}(T_{ct'}+1)}$ 
11  return V, prior, condprob

```

```

APPLYMULTINOMIALNB(C, V, prior, condprob, d)
1  W ← EXTRACTTOKENSFROMDOC(V, d)
2  for each c ∈ C
3  do score[c] ← log prior[c]
4     for each t ∈ W
5     do score[c] += log condprob[t][c]
6  return arg maxc∈C score[c]

```

► Figure 13.2 Naive Bayes algorithm (multinomial model): Training and testing.

Example

News titles for **Politics** and **Sports**

Politics

Sports

documents

“Obama meets Merkel”
“Obama elected again”
“Merkel visits Greece again”

“OSFP European basketball champion”
“Miami NBA basketball champion”
“Greece basketball coach?”

$P(p) = 0.5$

$P(s) = 0.5$

terms

obama:2, meets:1, merkel:2,
elected:1, again:2, visits:1,
greece:1

OSFP:1, european:1, basketball:3,
champion:2, miami:1, nba:1,
greece:1, coach:1

Vocabulary
size: 14

Total terms: 10

Total terms: 11

New title: **X = “Obama likes basketball”**

$$\begin{aligned} P(\text{Politics}|X) &\sim P(p) \cdot P(\text{obama}|p) \cdot P(\text{likes}|p) \cdot P(\text{basketball}|p) \\ &= 0.5 \cdot 3/(10+14) \cdot 1/(10+14) \cdot 1/(10+14) = 0.000108 \end{aligned}$$

$$\begin{aligned} P(\text{Sports}|X) &\sim P(s) \cdot P(\text{obama}|s) \cdot P(\text{likes}|s) \cdot P(\text{basketball}|s) \\ &= 0.5 \cdot 1/(11+14) \cdot 1/(11+14) \cdot 4/(11+14) = 0.000128 \end{aligned}$$

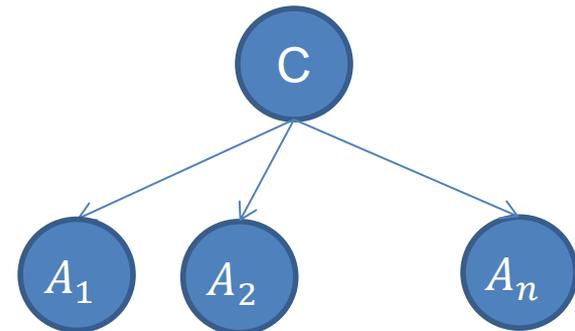
Naïve Bayes (Summary)

- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption may not hold for some attributes
 - Use other techniques such as Bayesian Belief Networks (BBN)
- Naïve Bayes can produce a probability estimate, but it is usually a very biased one
 - Logistic Regression is better for obtaining probabilities.

Generative vs Discriminative models

- Naïve Bayes is a type of a **generative model**
 - Generative process:
 - First pick the category of the record
 - Then given the category, generate the attribute values from the distribution of the category

- Conditional independence given C



- We use the training data to learn the distribution of the values in a class

Generative vs Discriminative models

- Logistic Regression and SVM are **discriminative models**
 - The goal is to find the boundary that discriminates between the two classes from the training data
- In order to classify the language of a document, you can
 - Either learn the two languages and find which is more likely to have generated the words you see
 - Or learn what differentiates the two languages.

SUPERVISED LEARNING

Learning

- **Supervised Learning**: learn a model from the data using **labeled data**.
 - **Classification** and **Regression** are the prototypical examples of supervised learning tasks. Other are possible (e.g., ranking)
- **Unsupervised Learning**: learn a model – extract structure from **unlabeled data**.
 - **Clustering** and **Association Rules** are prototypical examples of unsupervised learning tasks.
- **Semi-supervised Learning**: learn a model for the data using both **labeled and unlabeled** data.

Supervised Learning Steps

- **Model** the problem
 - What is you are trying to predict? What kind of optimization function do you need? Do you need classes or probabilities?
- Extract **Features**
 - How do you find the right features that help to discriminate between the classes?
- Obtain **training data**
 - Obtain a collection of labeled data. Make sure it is large enough, accurate and representative. Ensure that classes are well represented.
- Decide on the **technique**
 - What is the right technique for your problem?
- **Apply** in practice
 - Can the model be trained for very large data? How do you test how you do in practice? How do you improve?

Modeling the problem

- Sometimes it is not obvious. Consider the following three problems
 - Detecting if an email is spam
 - Categorizing the queries in a search engine
 - Ranking the results of a web search

Feature extraction

- Feature extraction, or **feature engineering** is the most tedious but also the most important step
 - How do you separate the players of the Greek national team from those of the Swedish national team?
- One line of thought: throw features to the classifier and the classifier will figure out which ones are important
 - **More features**, means that you need **more training data**
- Another line of thought: **Feature Selection**: Select carefully the features using various functions and techniques
 - Computationally intensive

Training data

- An overlooked problem: How do you get **labeled data** for training your model?
 - E.g., how do you get training data for ranking?
- Usually requires a lot of manual effort and domain expertise and carefully planned labeling
 - Results are not always of high quality (lack of expertise)
 - And they are not sufficient (low coverage of the space)
- Recent trends:
 - Find a **source** that generates the labeled data for you.
 - **Crowd-sourcing** techniques

Dealing with small amount of labeled data

- **Semi-supervised learning** techniques have been developed for this purpose.
- **Self-training**: Train a classifier on the data, and then feed back the high-confidence output of the classifier as input
- **Co-training**: train two “independent” classifiers and feed the output of one classifier as input to the other.
- **Regularization**: Treat learning as an optimization problem where you define relationships between the objects you want to classify, and you exploit these relationships
 - Example: Image restoration

Technique

- The choice of technique depends on the problem requirements (do we need a probability estimate?) and the problem specifics (does independence assumption hold? do we think classes are linearly separable?)
- For many cases finding the right technique may be trial and error
- For many cases the exact technique does not matter.

Big Data Trumps Better Algorithms

- If you have enough data then the algorithms are not so important
- The web has made this possible.
 - Especially for text-related tasks
 - Search engine uses the **collective human intelligence**

Google lecture:
[Theorizing from the Data](#)

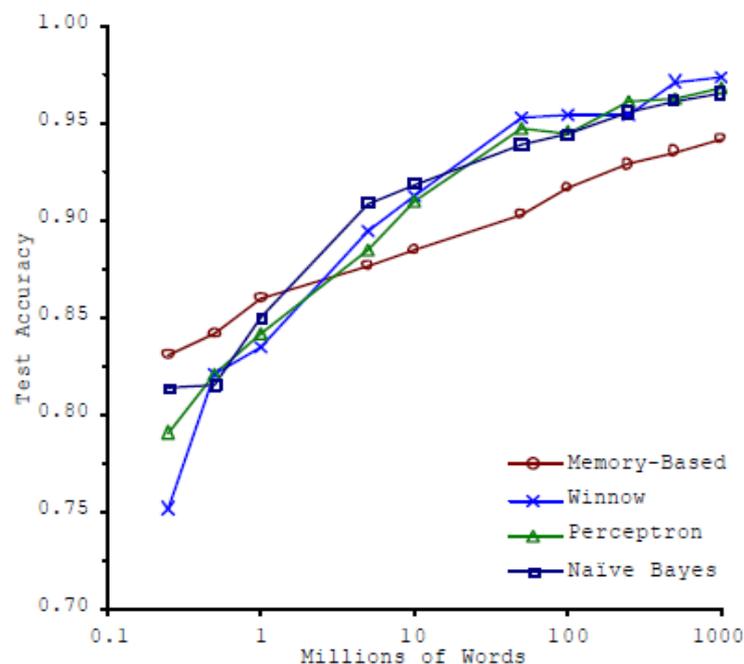


Figure 1. Learning Curves for Confusion Set Disambiguation

Apply-Test

- How do you **scale** to very large datasets?
 - Distributed computing – **map-reduce** implementations of machine learning algorithms (Mahut, over Hadoop)
- How do you test something that is running online?
 - You cannot get labeled data in this case
 - **A/B testing**
- How do you deal with changes in data?
 - **Active learning**