

Classification

Michael I. Jordan

University of California, Berkeley

Classification

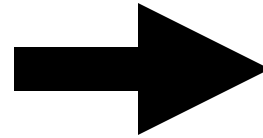
- In classification problems, each entity in some domain can be placed in one of a discrete set of categories: yes/no, friend/foe, good/bad/indifferent, blue/red/green, etc.
- Given a training set of labeled entities, develop a rule for assigning labels to entities in a test set
- Many variations on this theme:
 - binary classification
 - multi-category classification
 - non-exclusive categories
 - ranking
- Many criteria to assess rules and their predictions
 - overall errors
 - costs associated with different kinds of errors
 - operating points

Representation of Objects

- Each object to be classified is represented as a pair (x, y) :
 - where x is a description of the object (see examples of data types in the following slides)
 - where y is a label (assumed binary for now)
- Success or failure of a machine learning classifier often depends on choosing good descriptions of objects
 - the choice of description can also be viewed as a learning problem, and indeed we'll discuss automated procedures for choosing descriptions in a later lecture
 - but good human intuitions are often needed here

Data Types

- Vectorial data:
 - physical attributes
 - behavioral attributes
 - context
 - history
 - etc



- We'll assume for now that such vectors are explicitly represented in a table, but later (cf. kernel methods) we'll relax that assumption

Data Types

- text and hypertext

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.0 Transitional//EN">
<html>
<head>
    <meta http-equiv="Content-Type" content="text/html; charset=utf-8">
    <title>Welcome to FairmontNET</title>
</head>
<STYLE type="text/css">
.stdtext {font-family: Verdana, Arial, Helvetica, sans-serif; font-size: 11px; color: #1F3D4E;}
.stdtext_wh {font-family: Verdana, Arial, Helvetica, sans-serif; font-size: 11px; color: WHITE;}
</STYLE>

<body leftmargin="0" topmargin="0" marginwidth="0" marginheight="0" bgcolor="BLACK">
<TABLE cellpadding="0" cellspacing="0" width="100%" border="0">
  <TR>
    <TD width=50% background="/TFN/en/CDA/Images/common/labels/decorative_2px_blk.gif">&nbsp;</TD>
    <TD></td>
    <TD width=50% background="/TFN/en/CDA/Images/common/labels/decorative_2px_blk.gif">&nbsp;</TD>
  </TR>
</TABLE>
<tr>
<td align="right" valign="middle"><IMG src="/TFN/en/CDA/Images/common/labels/centrino_logo_blk.gif"></td>
</tr>
</body>
</html>
```

Data Types

- email

Return-path <bmill@eecs.berkeley.edu>Received from relay2.EECS.Berkeley.EDU (relay2.EECS.Berkeley.EDU [169.229.60.28]) by imap4.CS.Berkeley.EDU (iPlanet Messaging Server 5.2 HotFix 1.16 (built May 14 2003)) with ESMTP id <0HZ000F506JV5S@imap4.CS.Berkeley.EDU>; Tue, 08 Jun 2004 11:40:43 -0700 (PDT)Received from relay3.EECS.Berkeley.EDU (localhost [127.0.0.1]) by relay2.EECS.Berkeley.EDU (8.12.10/8.9.3) with ESMTP id i58Ieg3N000927; Tue, 08 Jun 2004 11:40:43 -0700 (PDT)Received from redbirds (dhcp-168-35.EECS.Berkeley.EDU [128.32.168.35]) by relay3.EECS.Berkeley.EDU (8.12.10/8.9.3) with ESMTP id i58IegFp007613; Tue, 08 Jun 2004 11:40:42 -0700 (PDT)Date Tue, 08 Jun 2004 11:40:42 -0700From Robert Miller <bmill@eecs.berkeley.edu>Subject RE: SLT headcount = 25In-reply-to <6.1.1.1.0.20040607101523.02623298@imap.eecs.Berkeley.edu>To 'Randy Katz' <randy@eecs.berkeley.edu>Cc "'Glenda J. Smith'" <glendajs@eecs.berkeley.edu>, 'Gert Lanckriet' <gert@eecs.berkeley.edu>Message-id <200406081840.i58IegFp007613@relay3.EECS.Berkeley.EDU>MIME-version 1.0X-MIMEOLE Produced By Microsoft MimeOLE V6.00.2800.1409X-Mailer Microsoft Office Outlook, Build 11.0.5510Content-type multipart/alternative; boundary="-----_NextPart_000_0033_01C44D4D.6DD93AF0"Thread-index AcRMtQRp+R26IVFaRiuz4BfImikTRAA0wf3Qthe headcount is now 32. -----
----- Robert Miller, Administrative Specialist University of California, Berkeley Electronics Research Lab 634 Soda Hall #1776 Berkeley, CA 94720-1776 Phone: 510-642-6037 fax: 510-643-1289

Data Types

- protein sequences

QFDACCFIDDVSKIYG-DYGPI
QFDACCFIDDVSKIYG-DHGPI
QFGACCFIDDVSKTFRLEDGPI
QFDAC-FIDDVSKIIFRLEDGPI
RFDASCFIDDVSKIIFRLEDGPI
QFSVYCLIDDVSKIYR-HDGPM
QFPVCSIIDDL SKMYR-HDSPV
QFPVFCLIDDL SKIYR-DDGLI
QFDARCFIDDL SKIYR-HDGQV
QFDARCFIDDL SKIYR-HDGQV
QFDARCFIDDL SKIYR-HDGPI
RFDACCFIDDVSKICK-HDGPV
QFDACCFIDDVSKICK-HDGPV

Data Types

- sequences of Unix system calls

Process Management

<code>pid = fork()</code>	Create a child process
<code>s=waitpid(pid, &status, opts)</code>	Wait for a child to terminate
<code>s=execve(name, argv, envp)</code>	Replace a process' core image
<code>exit(status)</code>	Terminate execution
<code>s=signal(sig, &act, &oact)</code>	Specify action to take for a signal
<code>s=kill(pid, sig)</code>	Send a signal to process
<code>residual=alarm(seconds)</code>	Schedule a SIGALRM signal later
<code>pause()</code>	Suspend the caller until next signal

Memory Management

<code>size=brk(addr)</code>	Set the size of data segment
-----------------------------	------------------------------

Input/Output Management

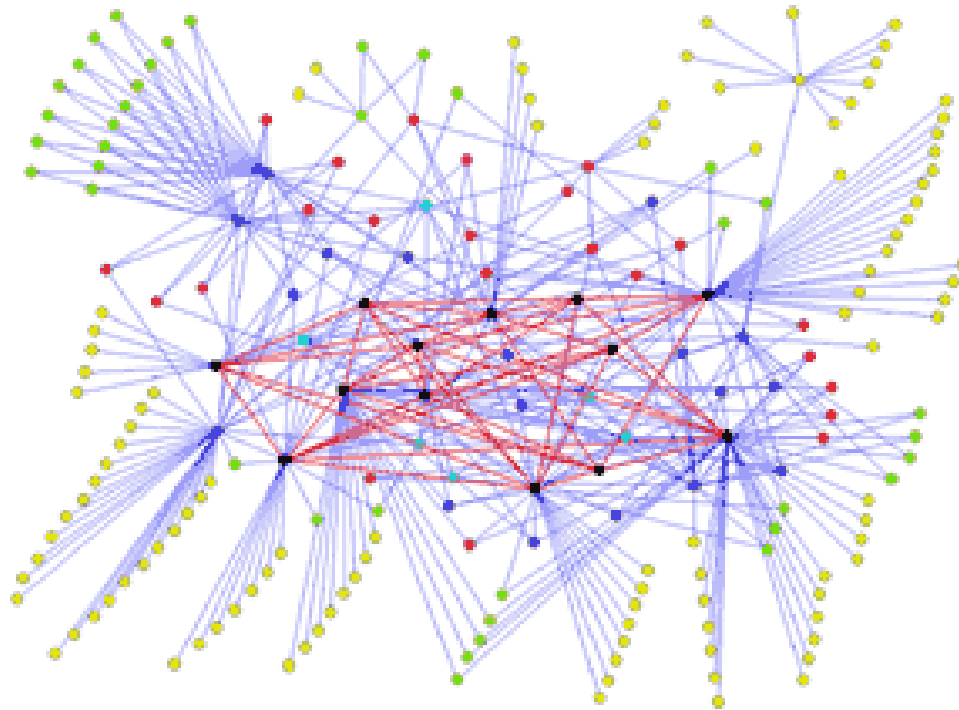
<code>s=cfsetospeed(&termios, speed)</code>	Set the output speed
<code>s=cfsetispeed(&termios, speed)</code>	Set the input speed
<code>s=cfgetospeed(&termios, speed)</code>	Get the output speed
<code>s=cfgetispeed(&termios, speed)</code>	Get the input speed
<code>s=tcsetattr(fd, opt, &termios)</code>	Set terminal attributes
<code>s=tcgetattr(fd, &termios)</code>	Get terminal attributes

Files and Directories Management

<code>fd=create(name, mode)</code>	Create a new file
<code>fd=open(name, how)</code>	Open a file for reading or writing
<code>s=close(fd)</code>	Close an open file
<code>n=read(fd, buffer, nbytes)</code>	Read data from file into a buffer
<code>n=write(fd, buffer, nbytes)</code>	Write data from buffer to file
<code>pos=lseek(fd, offset, whence)</code>	Move the file pointer somewhere
<code>s=stat(name, &buf)</code>	Read and return info. about file
<code>s=mkdir(name, mode)</code>	Create a new directory
<code>s=rmdir(name)</code>	Delete an empty directory
<code>s=link(name 1, name 2)</code>	Create a new directory entry for an old file
<code>s=unlink(name)</code>	Remove a directory entry
<code>s=chdir(dirname)</code>	Change the working directory
<code>s=chmod(name, mode)</code>	Change a file's protection bits

Data Types

- network layout: graph



Data Types

- images



Example: Spam Filter

- Input: email
- Output: spam/ham
- Setup:
 - Get a large collection of example emails, each labeled “spam” or “ham”
 - Note: someone has to hand label all this data
 - Want to learn to predict labels of new, future emails
- Features: The attributes used to make the ham / spam decision
 - Words: FREE!
 - Text Patterns: \$dd, CAPS
 - Non-text: SenderInContacts
 - ...



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Example: Digit Recognition

- Input: images / pixel grids
- Output: a digit 0-9
- Setup:
 - Get a large collection of example images, each labeled with a digit
 - Note: someone has to hand label all this data
 - Want to learn to predict labels of new, future digit images
- Features: The attributes used to make the digit decision
 - Pixels: (6,8)=ON
 - Shape Patterns: NumComponents, AspectRatio, NumLoops
 - ...
- Current state-of-the-art: Human-level performance



0



1



2



1



??

Other Examples of Real-World Classification Tasks

- Fraud detection (input: account activity, classes: fraud / no fraud)
 - Web page spam detection (input: HTML/rendered page, classes: spam / ham)
 - Speech recognition and speaker recognition (input: waveform, classes: phonemes or words)
 - Medical diagnosis (input: symptoms, classes: diseases)
 - Automatic essay grader (input: document, classes: grades)
 - Customer service email routing and foldering
 - Link prediction in social networks
 - Catalytic activity in drug design
 - ... many many more
-
- Classification is an important commercial technology

Training and Validation

- Data: labeled instances, e.g. emails marked spam/ham
 - Training set
 - Validation set
 - Test set
- Training
 - Estimate parameters on training set
 - Tune hyperparameters on validation set
 - Report results on test set
 - Anything short of this yields over-optimistic claims
- Evaluation
 - Many different metrics
 - Ideally, the criteria used to train the classifier should be closely related to those used to evaluate the classifier
- Statistical issues
 - Want a classifier which does well on *test* data
 - Overfitting: fitting the training data very closely, but not generalizing well
 - Error bars: want realistic (conservative) estimates of accuracy

Training
Data

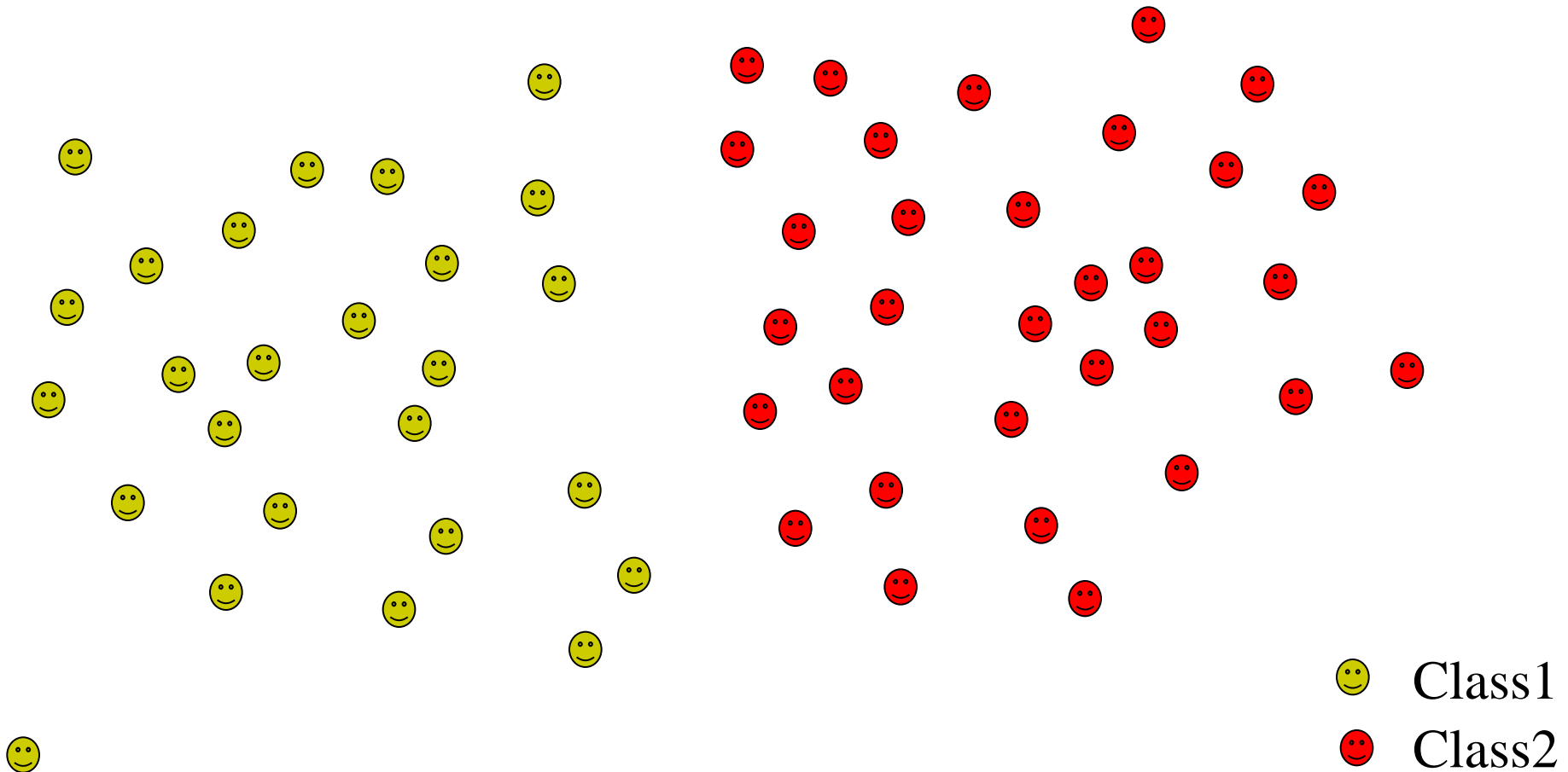
Validation
Data

Test
Data

Some State-of-the-art Classifiers

- Support vector machine
- Random forests
- Kernelized logistic regression
- Kernelized discriminant analysis
- Kernelized perceptron
- Bayesian classifiers
- Boosting and other ensemble methods
- (Nearest neighbor)

Intuitive Picture of the Problem



Some Issues

- There may be a simple separator (e.g., a straight line in 2D or a hyperplane in general) or there may not
- There may be “noise” of various kinds
- There may be “overlap”
- One should not be deceived by one’s low-dimensional geometrical intuition
- Some classifiers explicitly represent separators (e.g., straight lines), while for other classifiers the separation is done implicitly
- Some classifiers just make a decision as to which class an object is in; others estimate class probabilities

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes

- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron

- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees

- 2) Boosted Decision Trees

- 3) Random Forest

.note: $\Theta_{w|c} = 0$ if word w was
not observed in training

$$= \frac{P(\vec{w}|c)P(c)}{\sum_{c'} P(\vec{w}|c')P(c')}$$

\Rightarrow overfitting \Rightarrow need smoothing

Naive Bayes - Summary

pros:

- very simple!
- can scale easily to millions of training examples; just need counts!
- do surprisingly well in practice [in text application for example]

cons:

- bogus independence assumption
- can't include complex features
- not best prediction accuracy

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes

- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

2) Logistic Regression

* assume now \vec{x} is continuous: $\vec{x} = (x_1, \dots, x_d)$

a generative model using NB assumption
could be:

$$\text{let } N(x | \mu, \sigma) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

$$p_{\theta}(\vec{x} | C) = \prod_{i=1}^d N(x_i | \mu_{i|C}, \sigma_i) \quad \text{[NB]}$$

assume
variance doesn't
depend on C

$$p_{\theta}(C) = \pi_C \quad \text{[e.g., Bernoulli if } C \in \{0, 1\}]$$

and parameter $\theta = \left(\begin{matrix} \text{mean} & \text{variance} \\ \mu_{i|C}, \sigma_{i|C}, \pi_C \end{matrix} \right)$

for $i=1, \dots, d$

$C = 0$ to #classes-1

consider decision boundary: [in binary case]

$$\frac{P(C=1|\vec{x})}{P(C=0|\vec{x})} = \frac{P(C=1, \vec{x}) / \cancel{p(\vec{x})}}{P(C=0, \vec{x}) / \cancel{p(\vec{x})}}$$

$$= \frac{\pi_1 \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_{i|1})^2}{2\sigma_i^2}\right)}{\pi_0 \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_{i|0})^2}{2\sigma_i^2}\right)}$$

square terms cancel! ▽

$$= \frac{\pi_1}{\pi_0} \exp\left(\sum_{i=1}^d \left(\frac{\mu_{i|1} - \mu_{i|0}}{\sigma_i^2}\right) x_i - \frac{(\mu_{i|0}^2 + \mu_{i|1}^2)}{2\sigma_i^2}\right)$$

[because equal variance]

$$\Rightarrow \ln \frac{P(C=1|\vec{x})}{P(C=0|\vec{x})} = \sum_{i=1}^d \underbrace{w_i'}_{\text{parameter}} \cdot x_i + w_0'$$

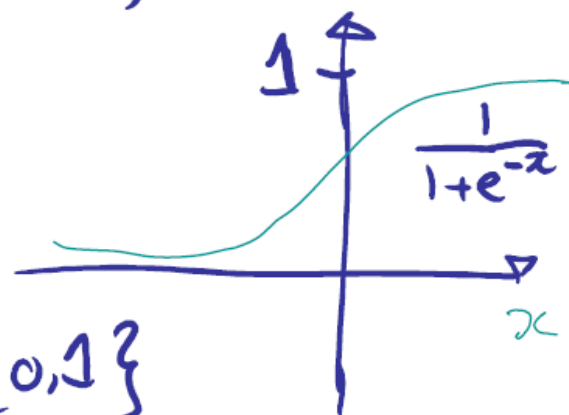
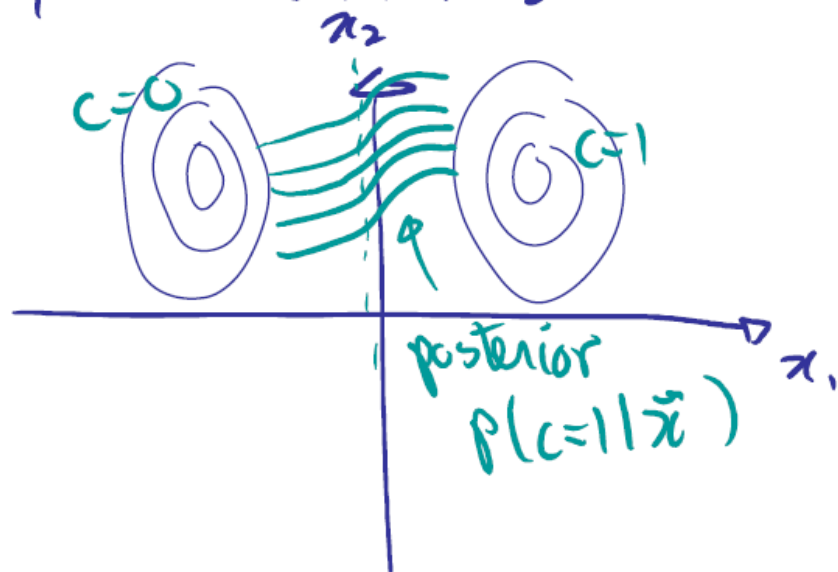
linear decision boundary!

can similarly write:

$$P(c=1 | \vec{x}) = \frac{1}{1 + \exp(\vec{w} \cdot \vec{x} + w_0)}$$

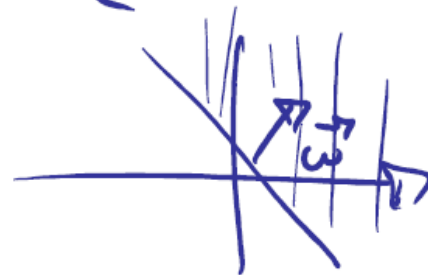
logistic
function

posterior for Gaussian with equal covariance:



$$y_i \in \{0, 1\}$$

$$h(x) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} + w_0 \geq 0 \\ 0 & \text{o.w.} \end{cases}$$



Learning parameters? given (\vec{x}_i, y_i)

• could learn $\mu_{iy}, \sigma_i, \pi_y$ by ML: $\max_{\theta} P_{\theta}(\vec{x}_i, y_i)$

• why not learn \vec{w} directly?

\Rightarrow maximum conditional Likelihood

$$\max_w \prod_{i=1}^n P_w(y_i | \vec{x}_i)$$

• advantage: more robust to model assumptions

[e.g. think of exponential family]

Maximum Conditional likelihood

• unfortunately, no closed formula for \vec{w}^*

→ use gradient ascent or your favourite optimization dg. on conditional likelihood

• log cond. likelihood

$$l(w) = \ln \left(\prod_{i=1}^n p_w(y_i | \vec{x}_i) \right)$$

$$= \sum_{i=1}^n y_i (w_0 + \vec{w} \cdot \vec{x}_i) - \sum_{i=1}^n \ln(1 + \exp(w_0 + \vec{w} \cdot \vec{x}_i))$$

$$\Rightarrow \nabla l = \sum_{i=1}^n (y_i - p_w)$$



Stochastic gradient ascent:

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} + \underbrace{\eta}_{\text{step size}} (y_i - p_{\vec{w}^t})$$

iterate over data until convergence criterion

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

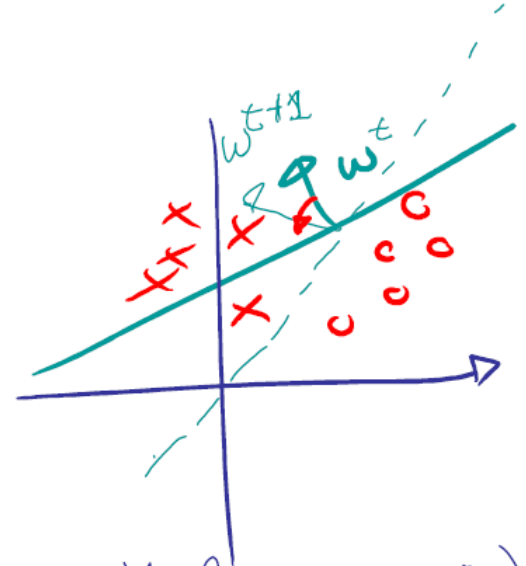
IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

III) Linear Models :

1) Perceptron

$$h(x) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} + b \geq 0 \\ 0 & \text{o.w.} \end{cases}$$



(compare with log. regression)

mistake driven: $\vec{w}^{t+1} = \vec{w}^t + \eta \underbrace{(y_i - h(x_i))}_{\neq 0 \text{ only if mistake}} \vec{x}_i$

2) Average perceptron:

$$\vec{w}_{\text{final}} = \frac{1}{T} \sum_{t=1}^T \vec{w}^{(t)}$$

$$\vec{w}_{\text{final}} \cdot \vec{x} = \frac{1}{T} \sum_{t=1}^T \underbrace{\vec{w}^{(t)} \cdot \vec{x}}_{\text{voting scheme?}}$$

averaged perceptron: practical issues

- work quite well in practice
- but result sensitive to order of data points
⇒ randomize order
- choose stopping criterion by validation set...

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

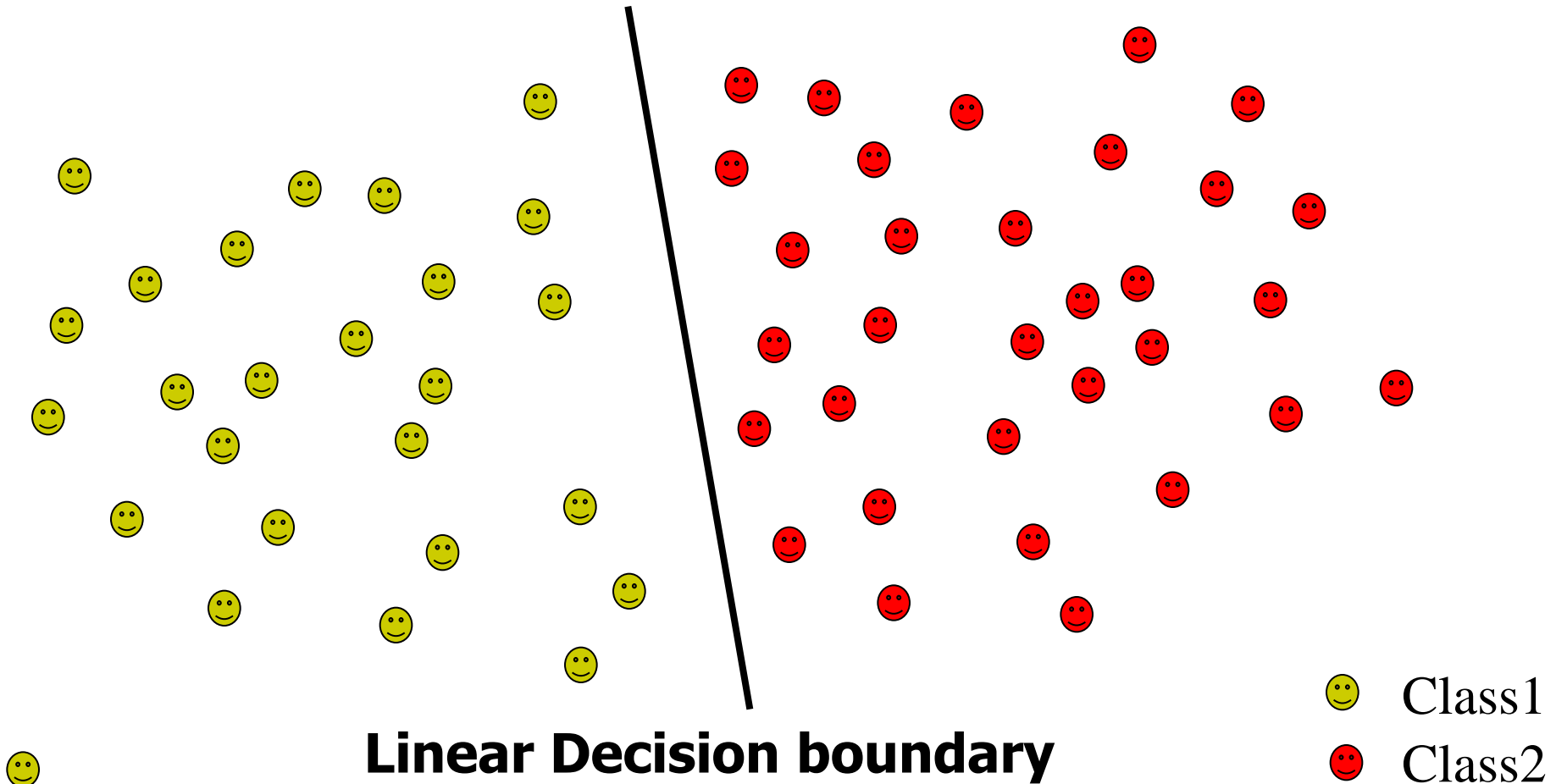
III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

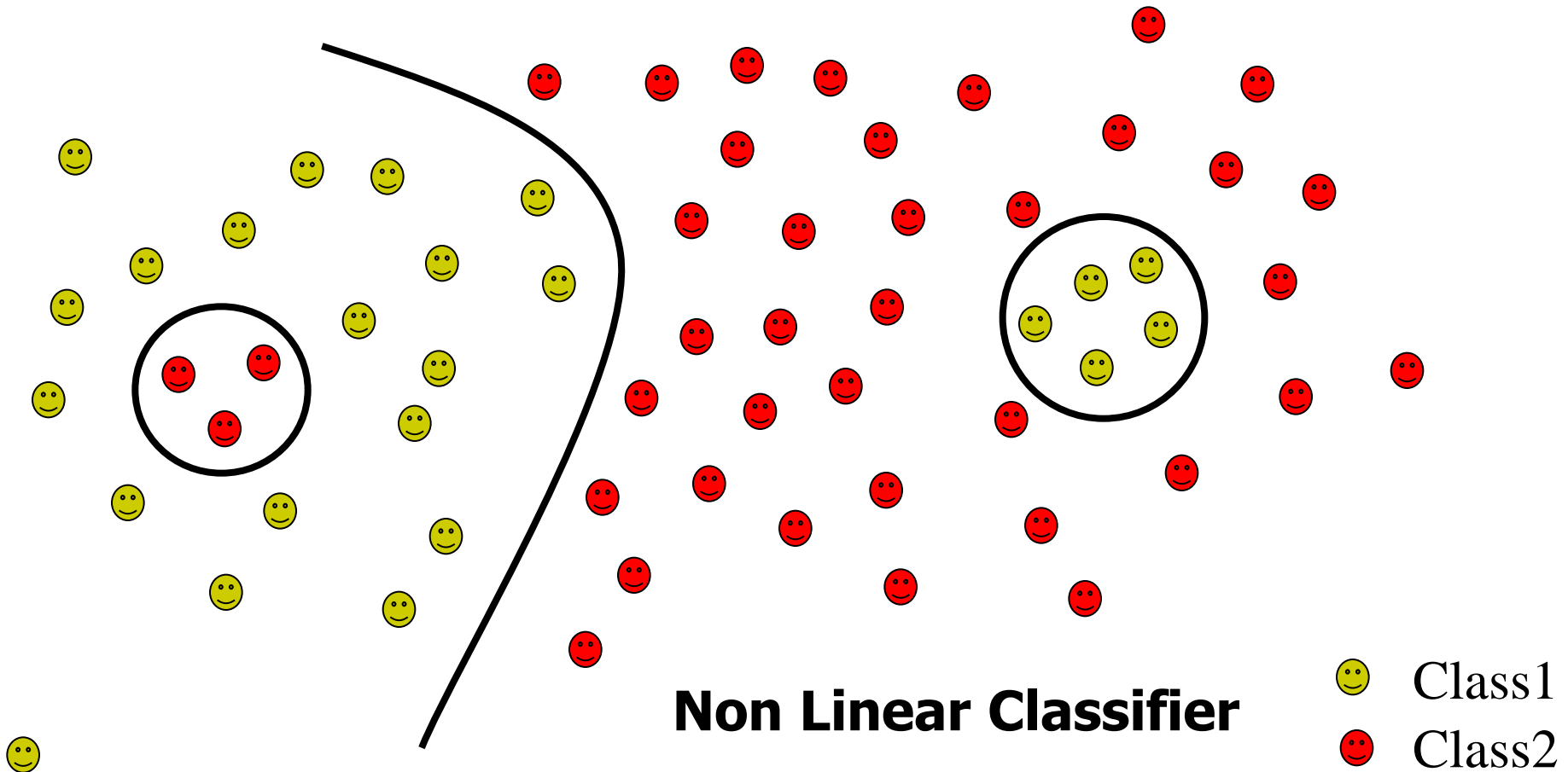
IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

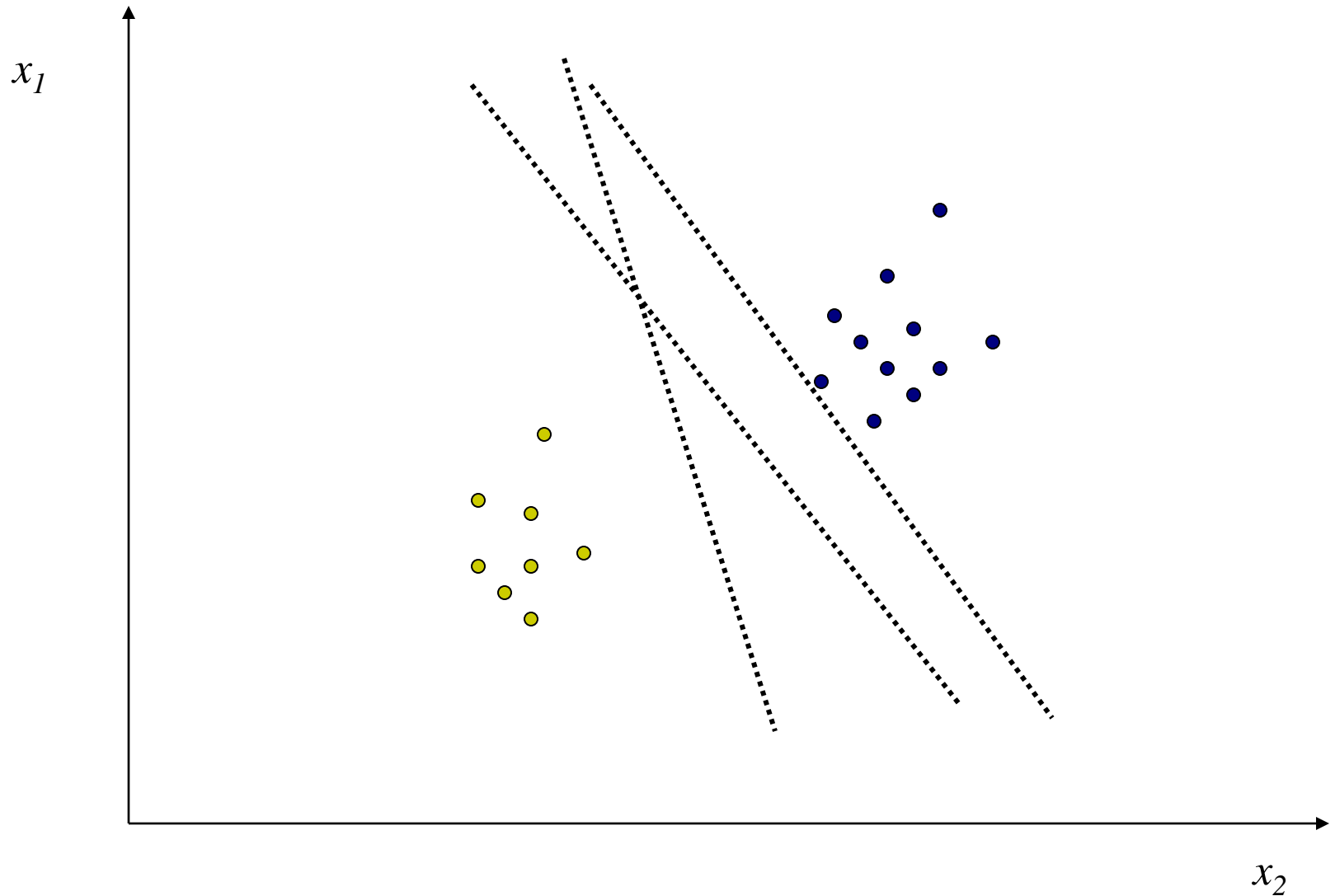
Linearly Separable Data



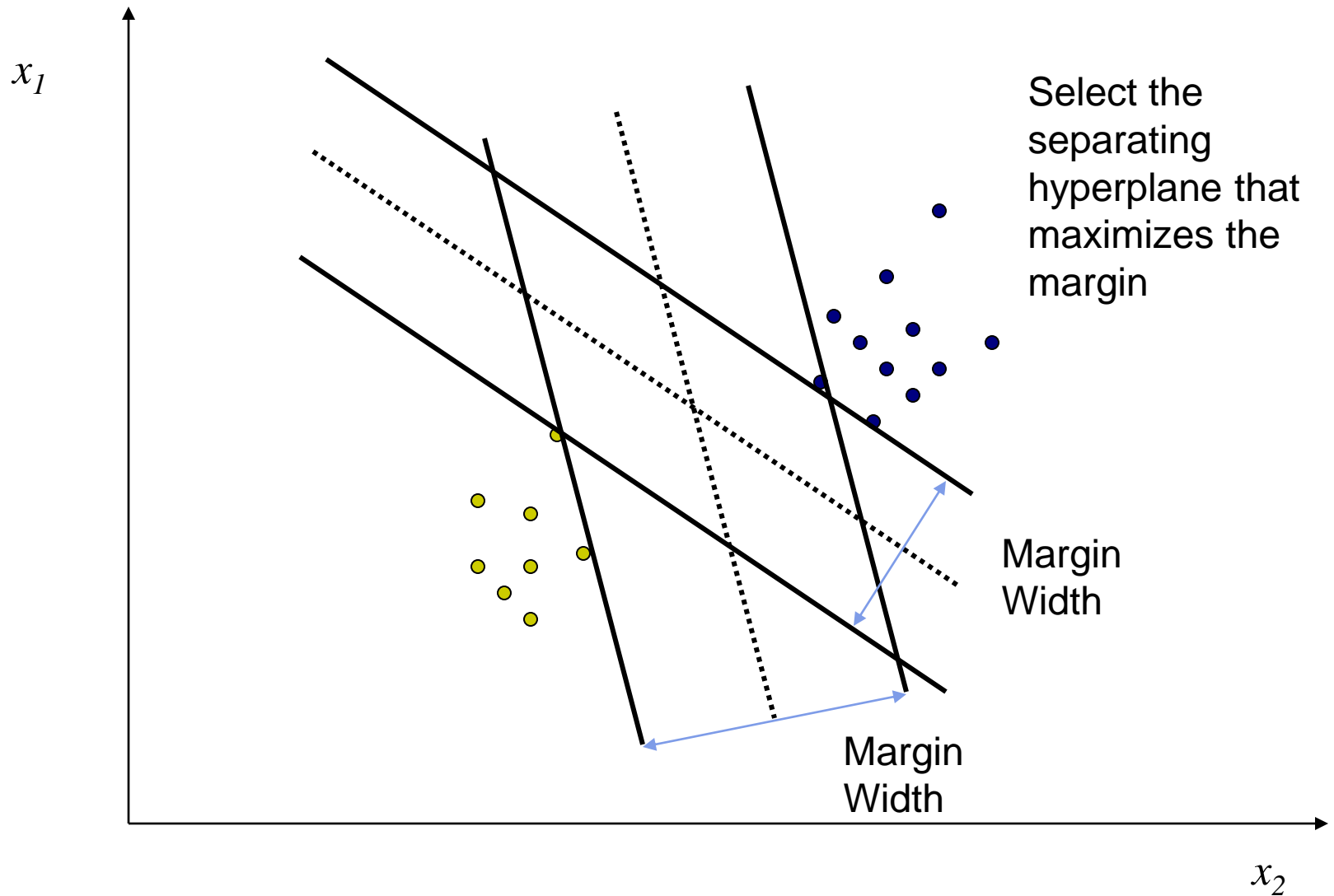
Nonlinearly Separable Data



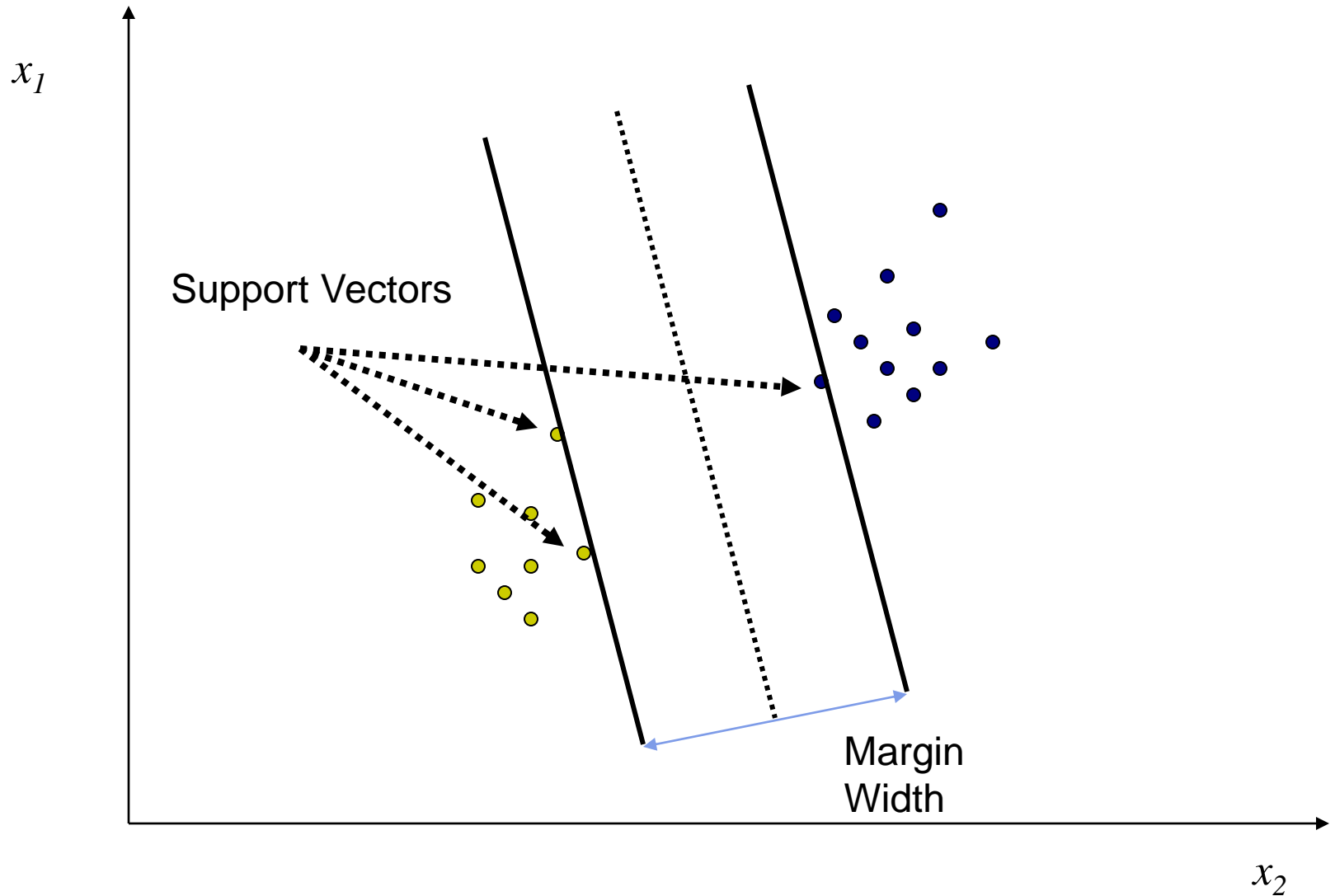
Which Separating Hyperplane to Use?



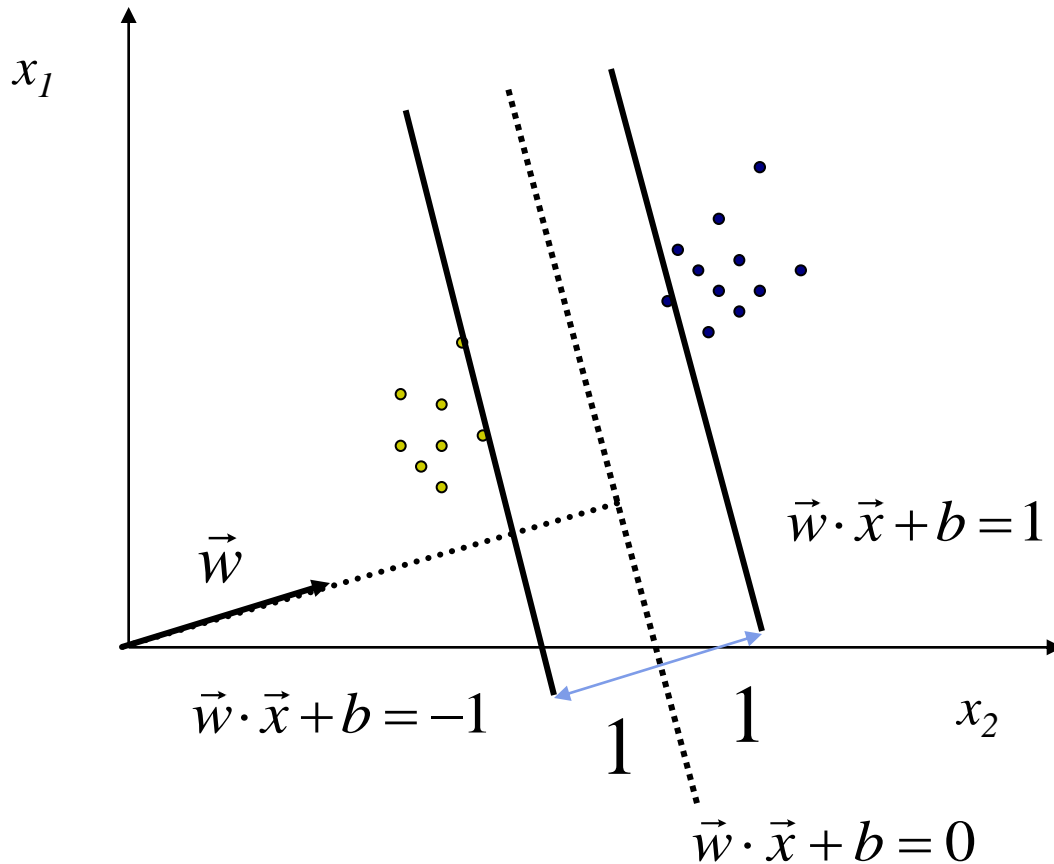
Maximizing the Margin



Support Vectors



Setting Up the Optimization Problem



The maximum margin can be characterized as a solution to an optimization problem:

$$\max \frac{2}{\|\vec{w}\|}$$

s.t. $(\vec{w} \cdot \vec{x} + b) \geq 1, \forall x$ of class 1

$(\vec{w} \cdot \vec{x} + b) \leq -1, \forall x$ of class 2

Setting Up the Optimization Problem

- If class 1 corresponds to 1 and class 2 corresponds to -1, we can rewrite

$$(w \cdot x_i + b) \geq 1, \quad \forall x_i \text{ with } y_i = 1$$

$$(w \cdot x_i + b) \leq -1, \quad \forall x_i \text{ with } y_i = -1$$

- as

$$y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

- So the problem becomes:

$$\max \frac{2}{\|w\|}$$

or

$$s.t. \ y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

$$\min \frac{1}{2} \|w\|^2$$

$$s.t. \ y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

Linear, Hard-Margin SVM Formulation

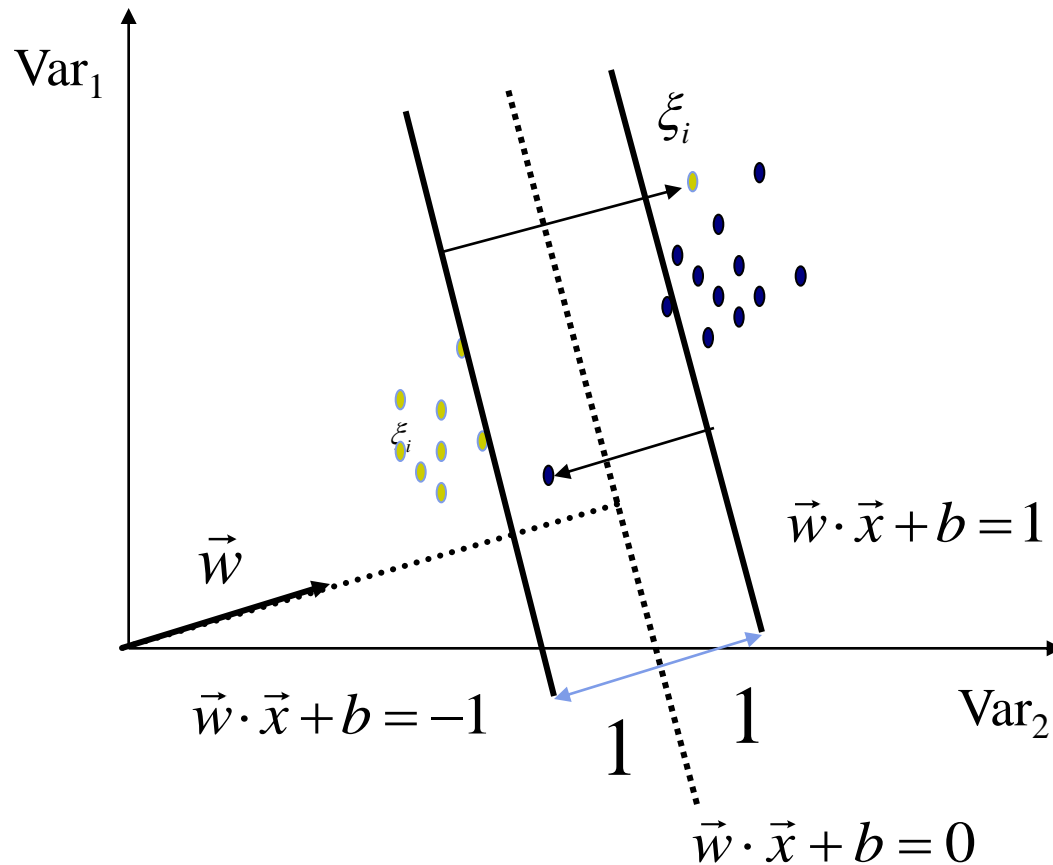
- Find w, b that solves

$$\min \frac{1}{2} \|w\|^2$$

$$s.t. \ y_i (w \cdot x_i + b) \geq 1, \ \forall x_i$$

- Problem is convex so, there is a unique global minimum value (when feasible)
- There is also a unique minimizer, i.e. weight and b value that provides the minimum
- Quadratic Programming
 - very efficient computationally with procedures that take advantage of the special structure

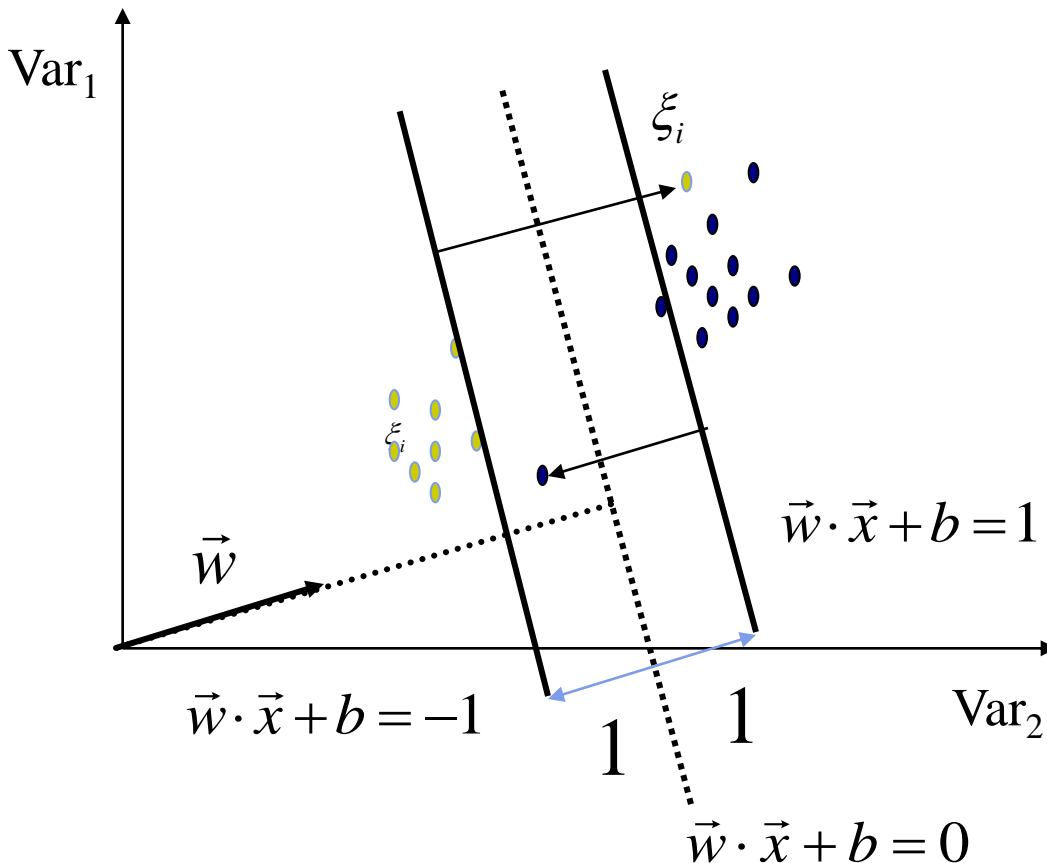
Nonlinearly Separable Data



Introduce slack variables ξ_i

Allow some instances to fall within the margin, but penalize them

Formulating the Optimization Problem



Constraints becomes :

$$y_i (w \cdot x_i + b) \geq 1 - \xi_i, \quad \forall x_i$$
$$\xi_i \geq 0$$

Objective function
penalizes for
misclassified instances
and those within the
margin

$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

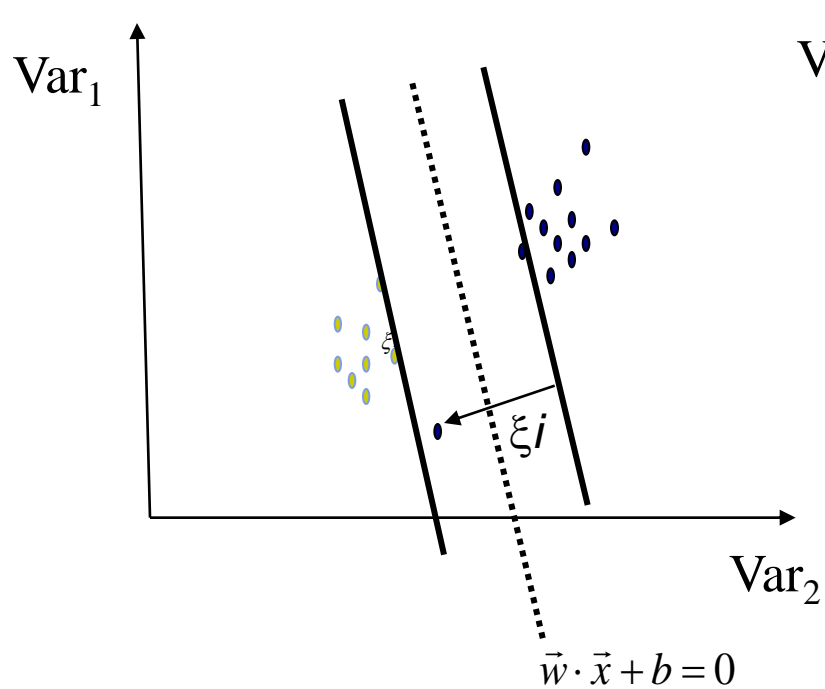
C trades-off margin width
and misclassifications

Linear, Soft-Margin SVMs

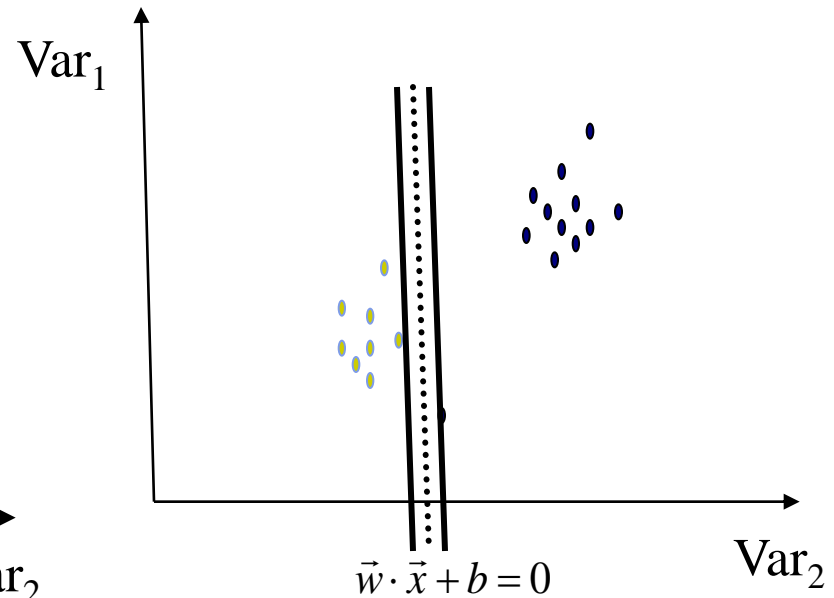
$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \quad \begin{array}{l} y_i(w \cdot x_i + b) \geq 1 - \xi_i, \quad \forall x_i \\ \xi_i \geq 0 \end{array}$$

- Algorithm tries to maintain ξ_i to zero while maximizing margin
- Notice: algorithm does not minimize the *number* of misclassifications (NP-complete problem) but the sum of distances from the margin hyperplanes
- Other formulations use ξ_i^2 instead
- As $C \rightarrow 0$, we get the hard-margin solution

Robustness of Soft vs Hard Margin SVMs

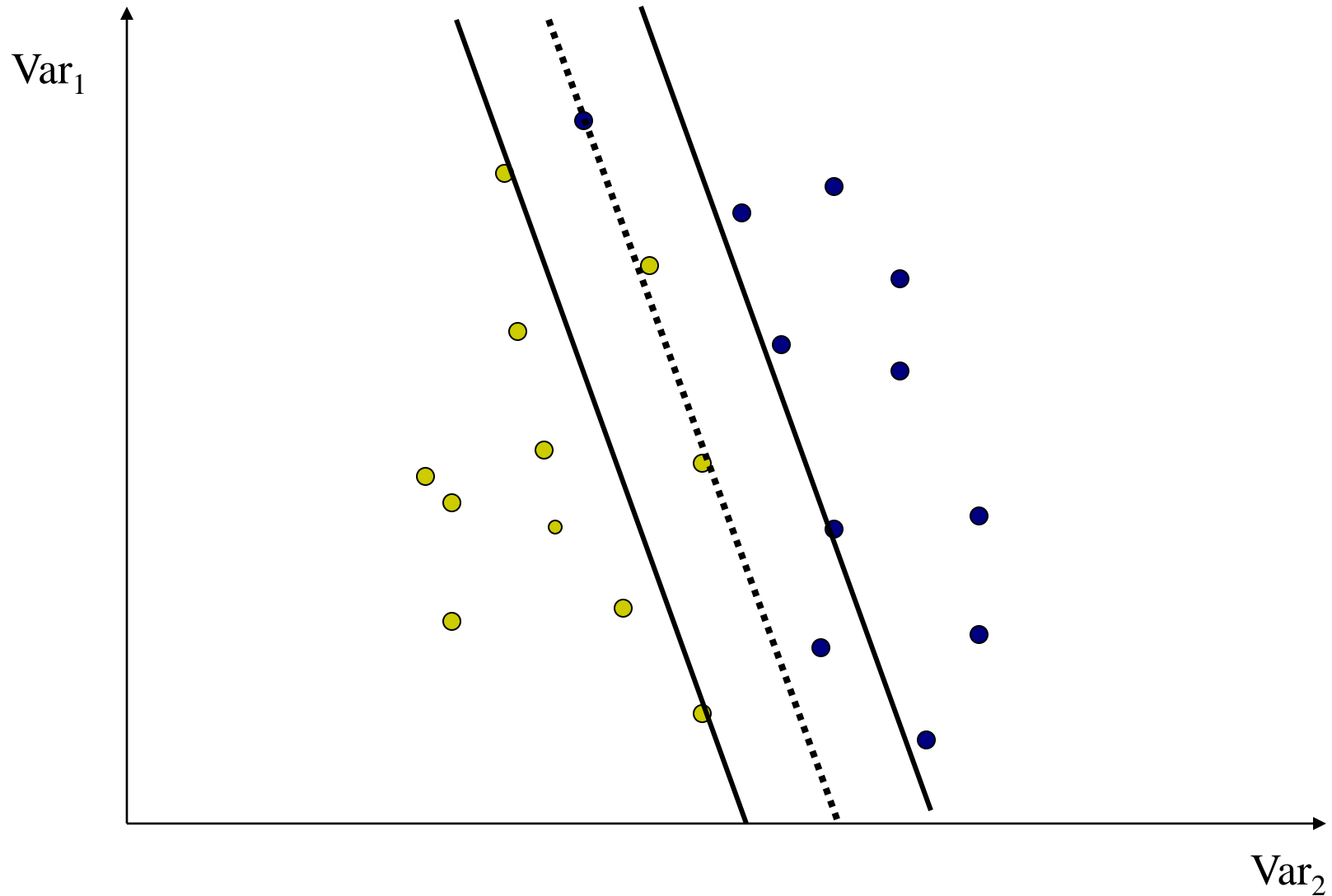


Soft Margin SVN

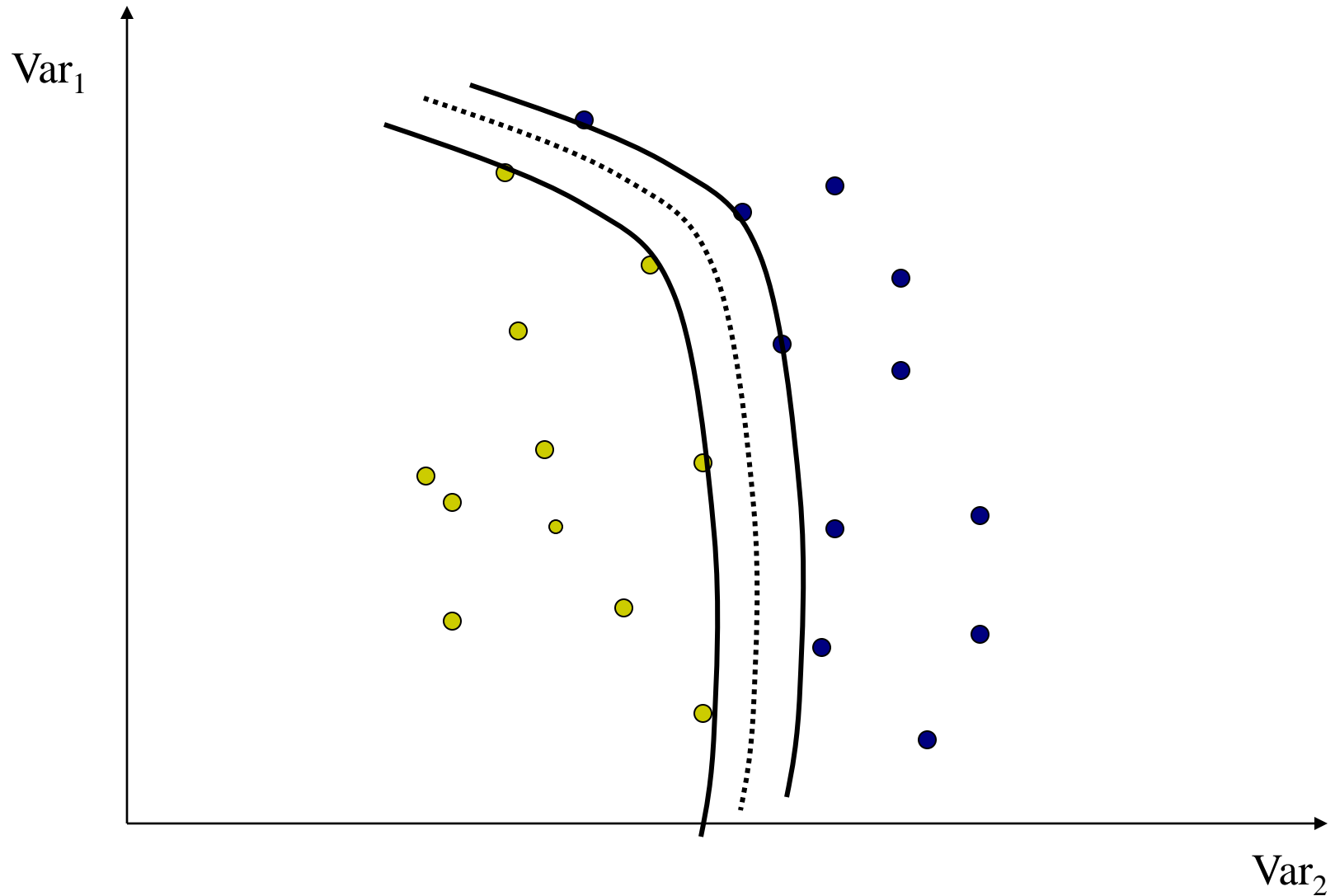


Hard Margin SVN

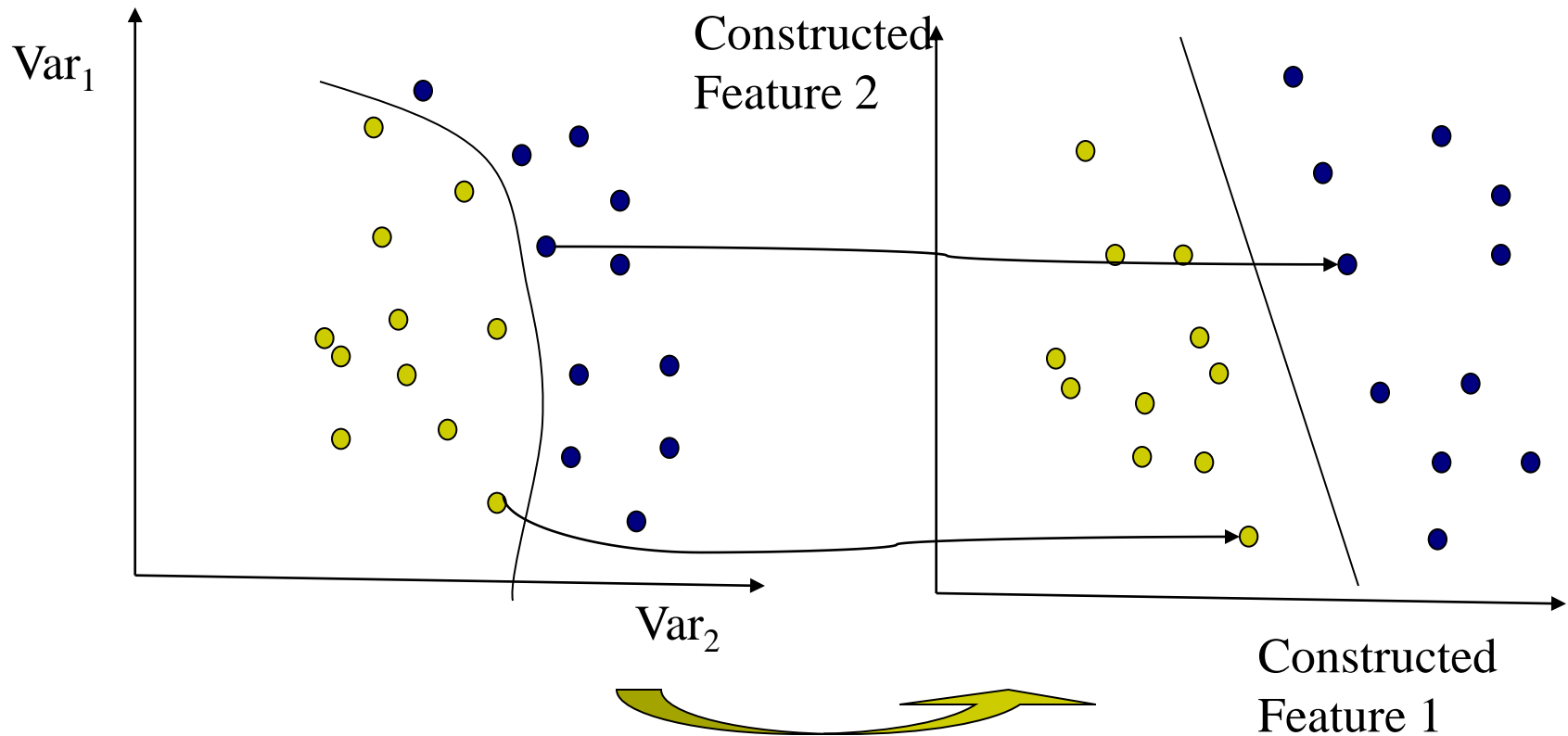
Disadvantages of Linear Decision Surfaces



Advantages of Nonlinear Surfaces



Linear Classifiers in High-Dimensional Spaces



Find function $\Phi(x)$ to map to
a different space

Mapping Data to a High-Dimensional Space

- Find function $\Phi(x)$ to map to a different space, then SVM formulation becomes:

$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \quad \text{s.t. } y_i(w \cdot \Phi(x) + b) \geq 1 - \xi_i, \forall x_i$$
$$\xi_i \geq 0$$

- Data appear as $\Phi(x)$, weights w are now weights in the new space
- Explicit mapping expensive if $\Phi(x)$ is very high dimensional
- Solving the problem without explicitly mapping the data is desirable

The Dual of the SVM Formulation

- Original SVM formulation

- n inequality constraints
- n positivity constraints
- n number of ξ variables

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

$$\begin{aligned} s.t. \quad & y_i (w \cdot \Phi(x) + b) \geq 1 - \xi_i, \forall x_i \\ & \xi_i \geq 0 \end{aligned}$$

- The (Wolfe) dual of this problem

- one equality constraint
- n positivity constraints
- n number of α variables (Lagrange multipliers)
- Objective function more complicated

$$\min_{w,b} \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\Phi(x_i) \cdot \Phi(x_j)) - \sum_i \alpha_i$$

$$s.t. \quad C \geq \alpha_i \geq 0, \forall x_i$$

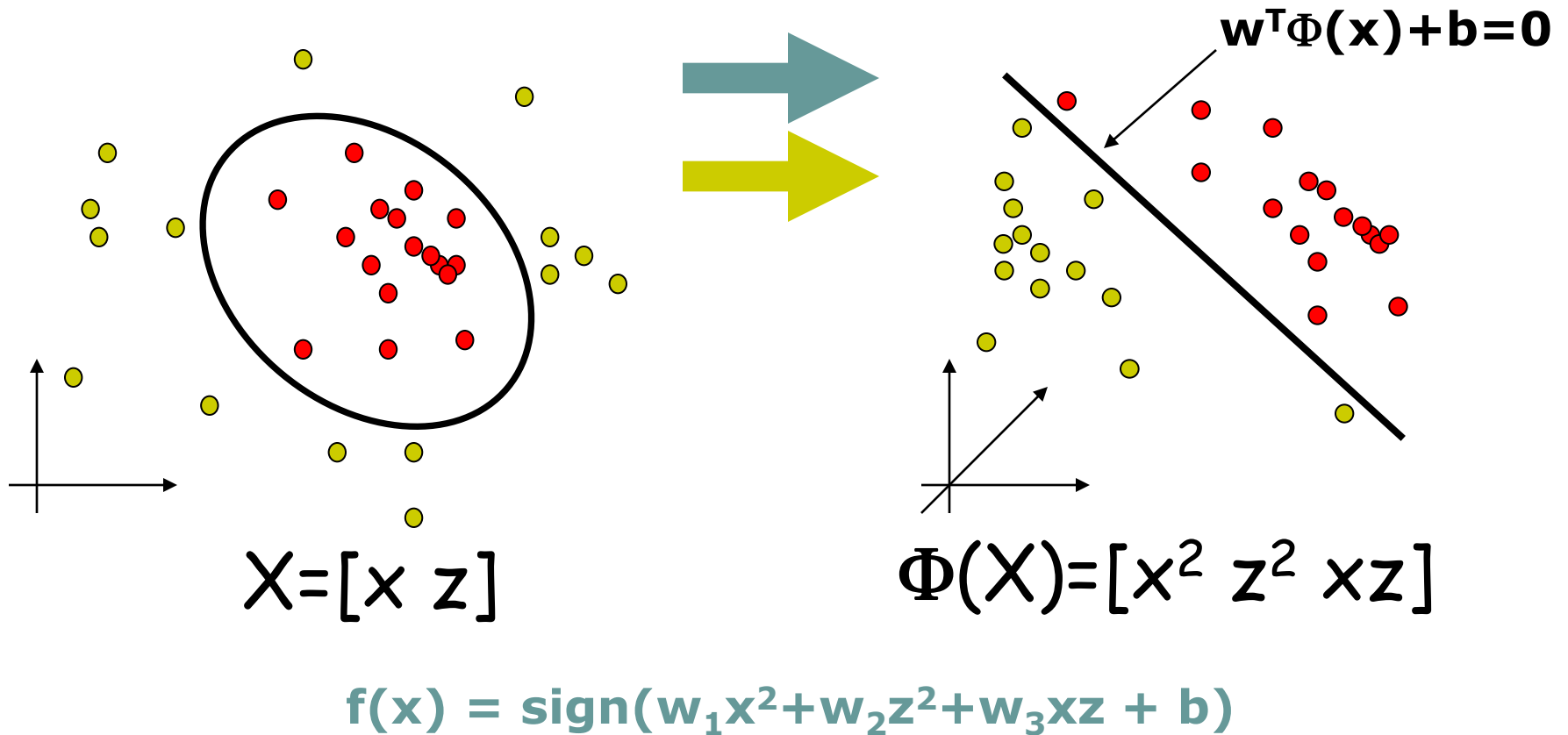
$$\sum_i \alpha_i y_i = 0$$

- NOTE: Data only appear as $\Phi(x_i) \cdot \Phi(x_j)$

The Kernel Trick

- $\Phi(x_i) \cdot \Phi(x_j)$: means, map data into new space, then take the inner product of the new vectors
- We can find a function such that: $K(x_i \cdot x_j) = \Phi(x_i) \cdot \Phi(x_j)$, i.e., the image of the inner product of the data is the inner product of the images of the data
- Then, we do not need to explicitly map the data into the high-dimensional space to solve the optimization problem

Example



Example

$$\begin{array}{ll} \mathbf{X}_1 = [\mathbf{x}_1 \ \mathbf{z}_1] & \xrightarrow{\text{teal}} \Phi(\mathbf{X}_1) = [\mathbf{x}_1^2 \ \mathbf{z}_1^2 \ 2^{1/2}\mathbf{x}_1\mathbf{z}_1] \\ \mathbf{X}_2 = [\mathbf{x}_2 \ \mathbf{z}_2] & \xrightarrow{\text{yellow}} \Phi(\mathbf{X}_2) = [\mathbf{x}_2^2 \ \mathbf{z}_2^2 \ 2^{1/2}\mathbf{x}_2\mathbf{z}_2] \end{array}$$

$$\begin{aligned} \Phi(\mathbf{X}_1)^T \Phi(\mathbf{X}_2) &= [\mathbf{x}_1^2 \ \mathbf{z}_1^2 \ 2^{1/2}\mathbf{x}_1\mathbf{z}_1] [\mathbf{x}_2^2 \ \mathbf{z}_2^2 \ 2^{1/2}\mathbf{x}_2\mathbf{z}_2]^T \\ &= \mathbf{x}_1^2\mathbf{z}_1^2 + \mathbf{x}_2^2 \ \mathbf{z}_2^2 + 2 \ \mathbf{x}_1 \ \mathbf{z}_1 \ \mathbf{x}_2 \ \mathbf{z}_2 \quad \text{Expensive!} \quad \mathcal{O}(d^2) \\ &= (\mathbf{x}_1\mathbf{z}_1 + \mathbf{x}_2\mathbf{z}_2)^2 \\ &= (\mathbf{X}_1^T \mathbf{X}_2)^2 \quad \text{Efficient!} \quad \mathcal{O}(d) \end{aligned}$$

Kernel Trick

- **Kernel function**: a symmetric function

$$\mathbf{k} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

- **Inner product kernels**: additionally,

$$\mathbf{k}(\mathbf{x}, \mathbf{z}) = \Phi(\mathbf{x})^\top \Phi(\mathbf{z})$$

- Example:

$\mathcal{O}(d^2)$

$$\Phi(\mathbf{x})^\top \Phi(\mathbf{z}) = \sum_{i,j=1}^{d,d} (x_i x_j)(z_i z_j) = \left(\sum_{i=1}^d x_i z_i \right)^2 = (x^\top z)^2 = K(\mathbf{x}, \mathbf{z})$$

$\mathcal{O}(d)$

Kernel Trick

- Implement an infinite-dimensional mapping **implicitly**
- Only inner products **explicitly** needed for training and evaluation
- Inner products computed **efficiently**, in finite dimensions
- The underlying mathematical theory is that of **reproducing kernel Hilbert space** from functional analysis

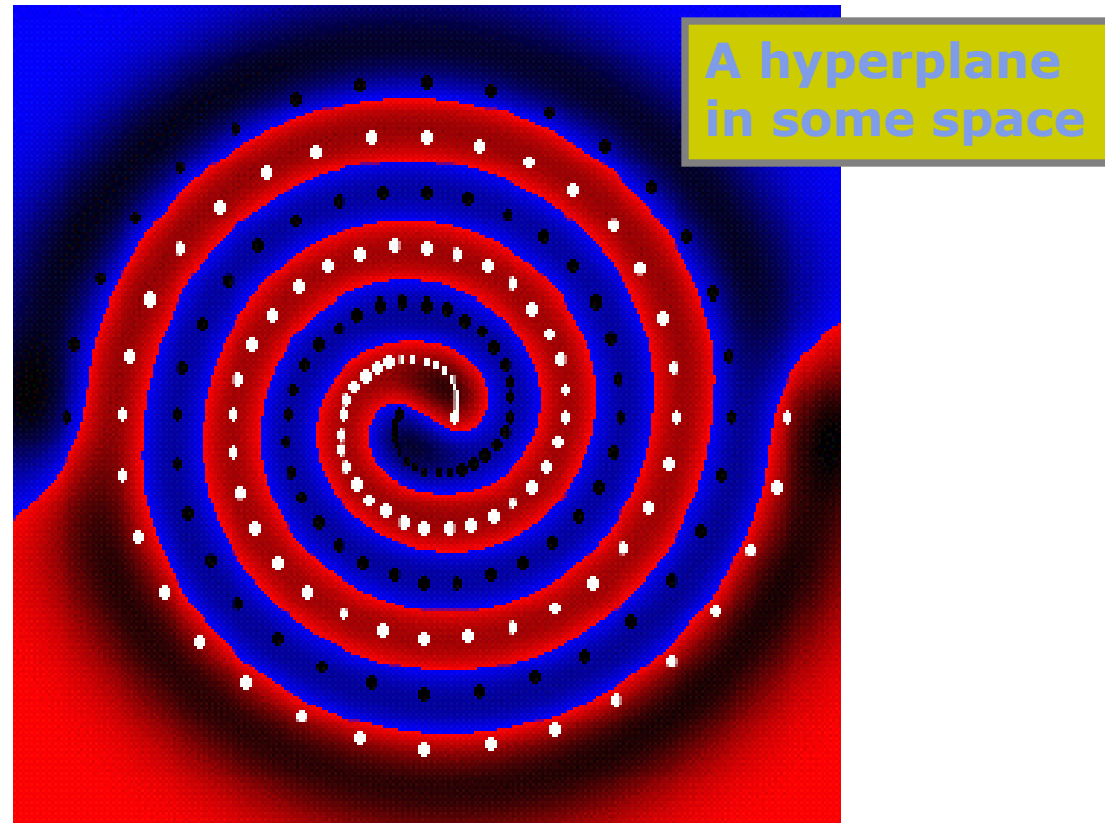
Kernel Methods

- If a linear algorithm can be expressed only in terms of inner products
 - it can be “kernelized”
 - find linear pattern in high-dimensional space
 - nonlinear relation in original space
- Specific kernel function determines nonlinearity

Kernels

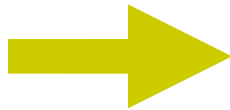
- Some simple kernels
 - Linear kernel: $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$
→ equivalent to linear algorithm
 - Polynomial kernel: $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^d$
→ polynomial decision rules
 - RBF kernel: $k(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|^2 / 2\sigma)$
→ highly nonlinear decisions

Gaussian Kernel: Example

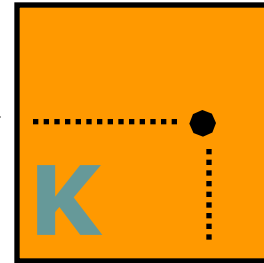


Kernel Matrix

$k(x, y)$



i

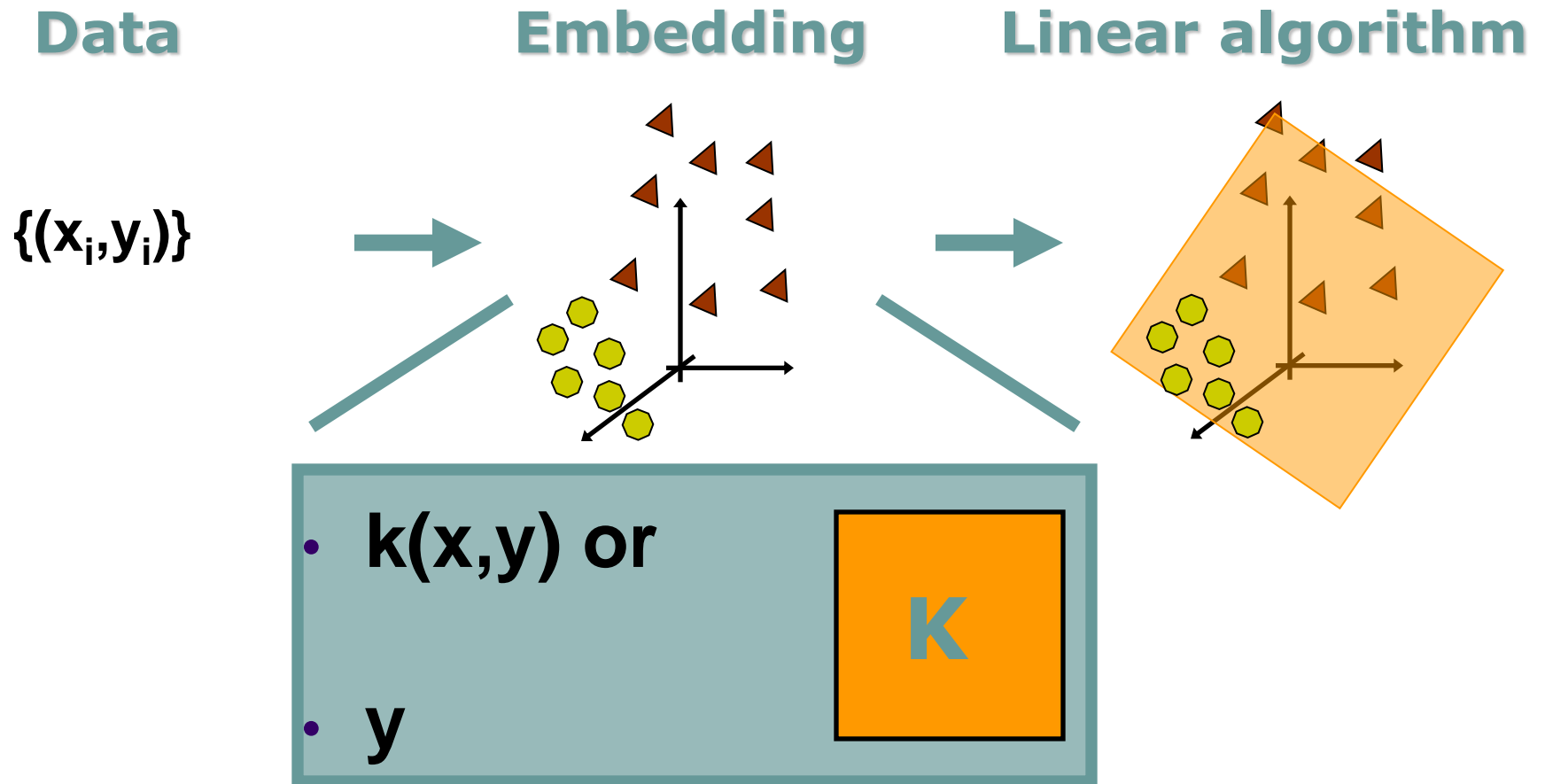


j

- Kernel matrix K defines all pairwise inner products
- Mercer theorem: K positive semidefinite
- Any symmetric positive semidefinite matrix can be regarded as an inner product matrix in some space

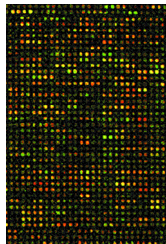
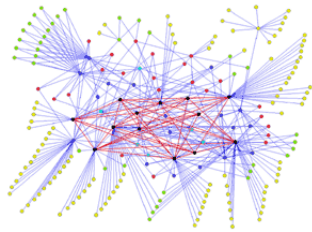
$$K_{ij} = k(x_i, x_j)$$

Kernel-Based Learning

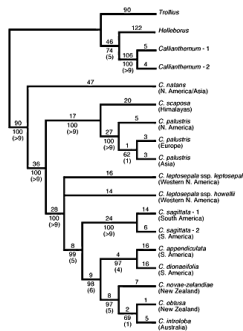


Kernel-Based Learning

Data



QFDACCFIDDVSKIYG-DYGP
QFDACCFIDDVSKIYG-DHGP
QFGACCFIDDVSKIFRLHDGP
QFDAC-FIDDVSKIFRLHDGP
RFDACCFIDDVSKIFRLHDGP
QFSVYCLIDDVSKIYR-HDGP
QFPVCSIIDDLKMYR-HDSP
QFPVFCLIDDLSKIYR-DDGL
QFDARCFIDDLSKIYR-HDGP
QFDARCFIDDLSKIYR-HDGP
QFDARCFIDDLSKIYR-HDGP
RFDACCFIDDVSKICK-HDGP
QFDACCFIDDVSKICK-HDGP

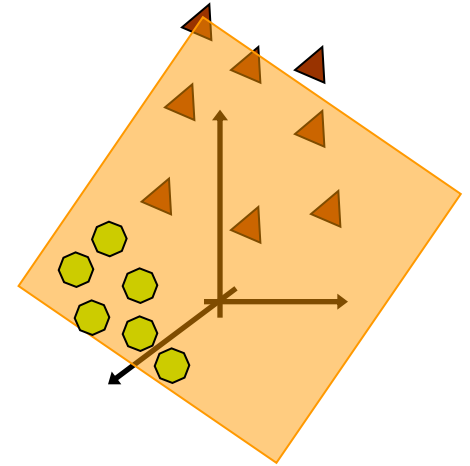


Embedding



Kernel design

Linear algorithm



Kernel algorithm

Kernel Design

- Simple kernels on vector data
- More advanced
 - string kernel
 - diffusion kernel
 - kernels over general structures (sets, trees, graphs...)
 - kernels derived from graphical models
 - empirical kernel map

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

Top-Down Induction of Decision Trees

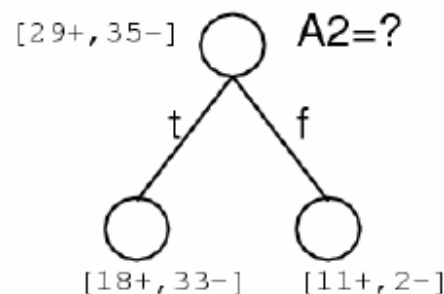
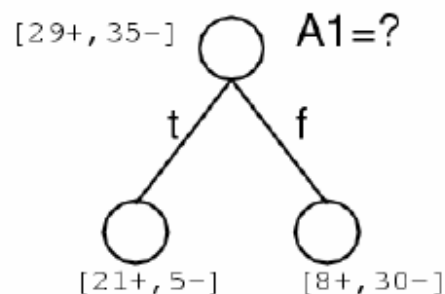
[ID3, C4.5, ...]

node = Root

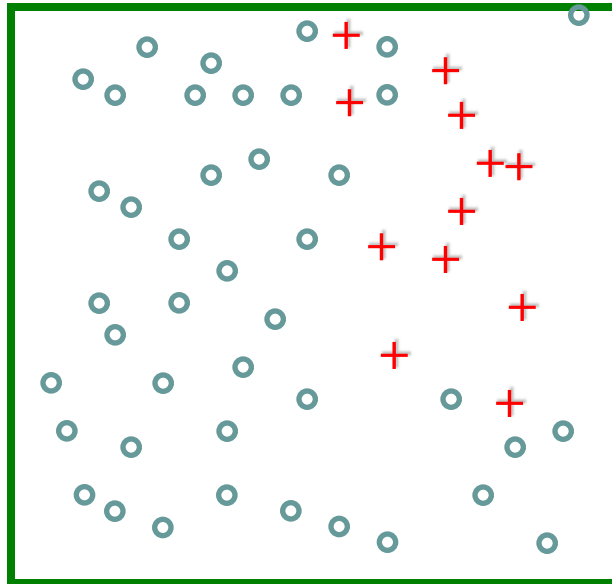
Main loop:

1. $A \leftarrow$ the “best” decision attribute for next *node*
2. Assign A as decision attribute for *node*
3. For each value of A , create new descendant of *node*
4. Sort training examples to leaf nodes
5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

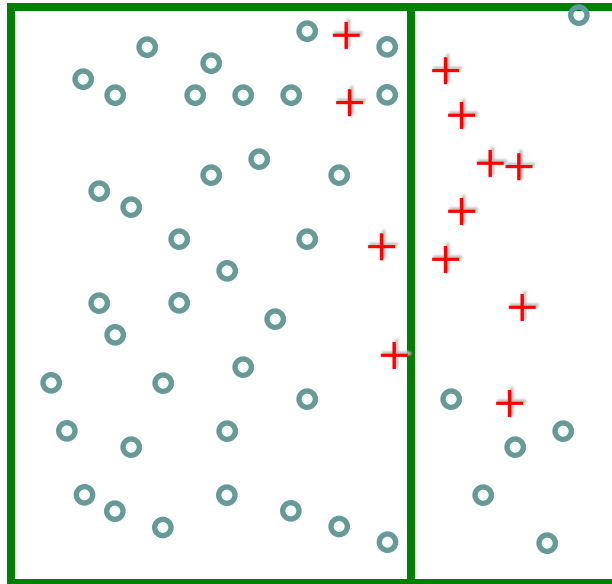
Which attribute is best?



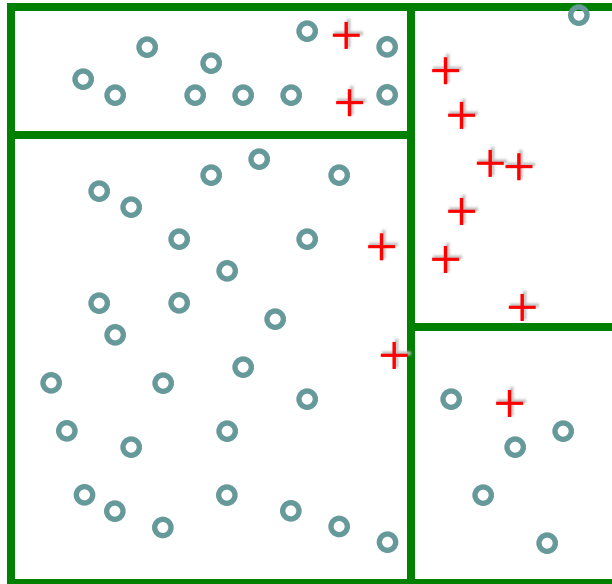
Spatial example: recursive binary splits



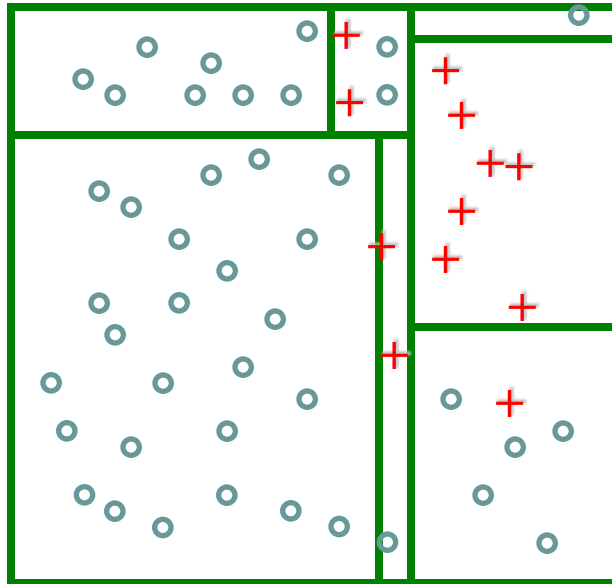
Spatial example: recursive binary splits



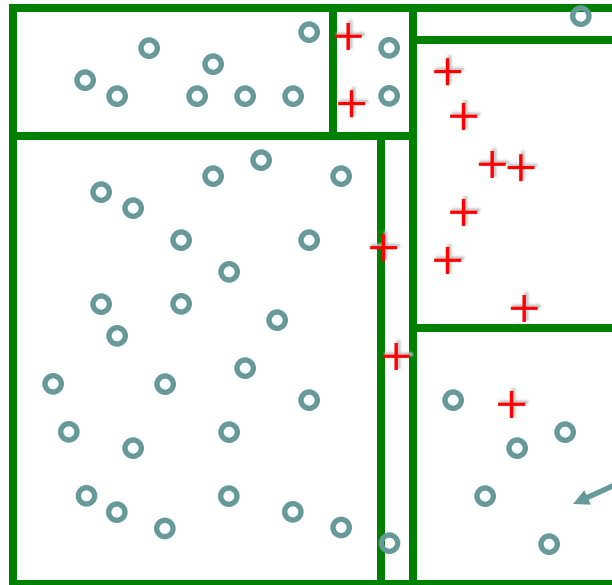
Spatial example: recursive binary splits



Spatial example: recursive binary splits



Spatial example: recursive binary splits

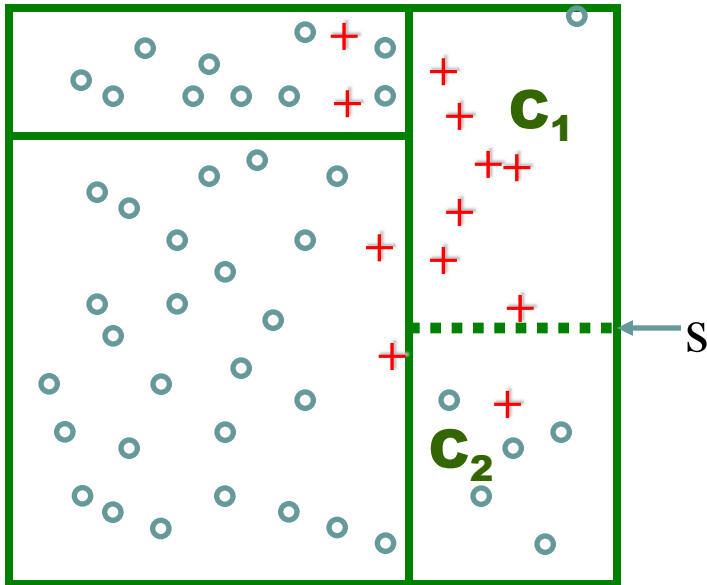


Once regions are
chosen class
probabilities are easy
to calculate

$$p_m = 5/6$$

How to choose a split

$$N_1=9$$
$$p_1=8/9$$



$$p_2=5/6$$
$$N_2=6$$

Impurity measures: $L(p)$

- Information gain (entropy):
- $p \log p - (1-p) \log(1-p)$
- Gini index: $2 p (1-p)$
- (0-1 error: $1-\max(p,1-p)$)

$$\min_s N_1 L(p_1) + N_2 L(p_2)$$

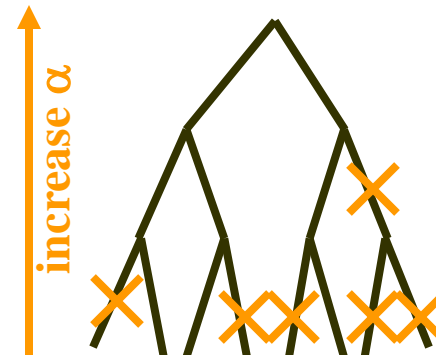
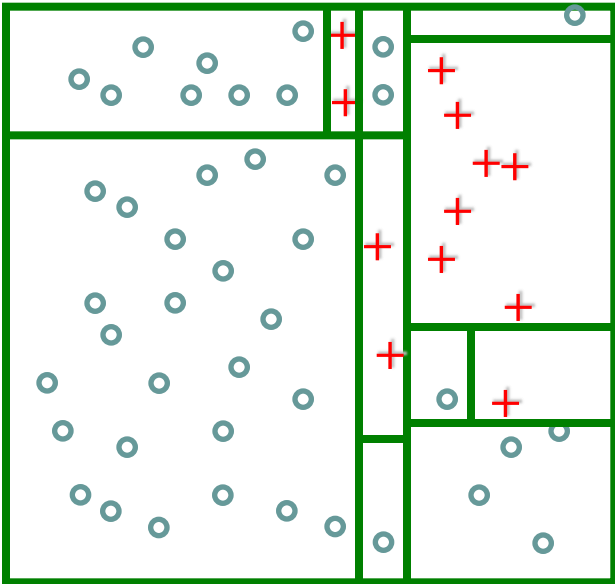
Then choose the region that has the best split

Overfitting and pruning

L: 0-1 loss

$$\min_T \sum_i L(x_i) + \alpha |T|$$

then choose α with CV



Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

Methods

I) Instance-based methods:

- 1) Nearest neighbor

II) Probabilistic models:

- 1) Naïve Bayes
- 2) Logistic Regression

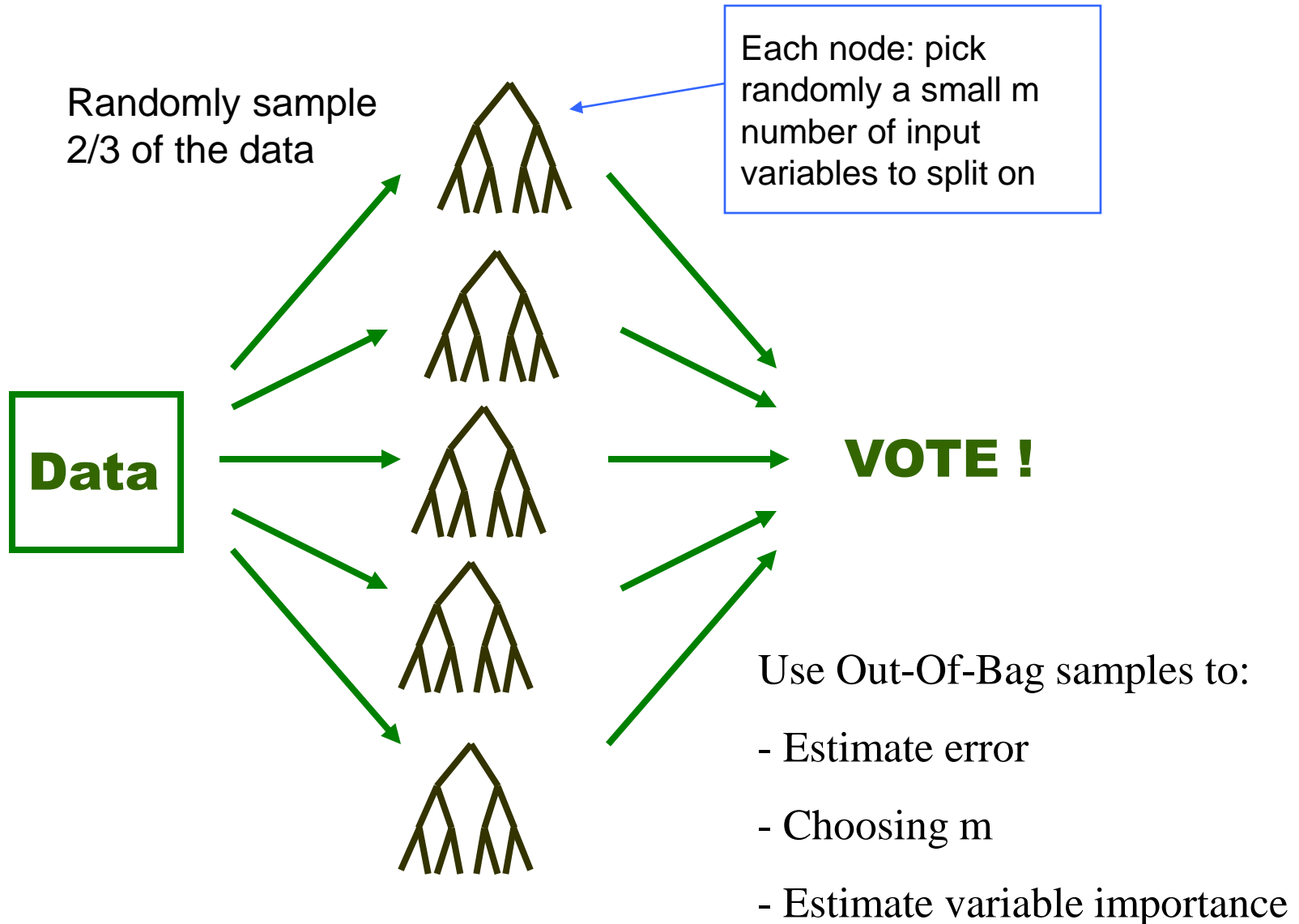
III) Linear Models:

- 1) Perceptron
- 2) Support Vector Machine

IV) Decision Models:

- 1) Decision Trees
- 2) Boosted Decision Trees
- 3) Random Forest

Random Forest



Summary of algorithms

 = state of the art

- 1) KNN : good first try if have "meaningful" $d(x, x')$
- 2) NB :
- 3) avg. perceptron / logistic regression / SVM
(in increasing order of complexity, memory requirement and prediction accuracy)
 - they can all be kernelized \Rightarrow can handle any data for which you can define meaningful kernel
"similarity measure"
 - are all vector space methods
 \Rightarrow deal with nominal data in unnatural way

- log. reg. output prob. [can be useful]
→ can get prob. for SVM using Platt's trick

4) trees methods [decision tree / boosted DT / random forest]

- handle mixed attributes naturally (e.g. nominal + numeric)
- interpretable
- more scalable than SVM

TABLE 10.1. *Some characteristics of different learning methods. Key: ● = good, ● = fair, and ● = poor.*

Characteristic	Neural nets	SVM	Trees	MARS	k-NN, kernels
Natural handling of data of "mixed" type	●	●	●	●	●
Handling of missing values	●	●	●	●	●
Robustness to outliers in input space	●	●	●	●	●
Insensitive to monotone transformations of inputs	●	●	●	●	●
Computational scalability (large N)	●	●	●	●	●
Ability to deal with irrel- evant inputs	●	●	●	●	●
Ability to extract linear combinations of features	●	●	●	●	●
Interpretability	●	●	●	●	●
Predictive power	●	●	●	●	●

Reading

- All of the methods that we have discussed are presented in the following book:

Hastie, T., Tibshirani, R. & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (Second Edition), NY: Springer.

- We haven't discussed theory, but if you're interested in the theory of (binary) classification, here's a pointer to get started:

Bartlett, P., Jordan, M. I., & McAuliffe, J. D. (2006). Convexity, classification and risk bounds. *Journal of the American Statistical Association*, 101, 138-156.