

## Decision Tree Induction

### Top-Down Induction of Decision Tree — TDIDT

#### ▷ Error Measures

**Classification:**  $Y$  is discrete, a small finite, unordered set of classes

$$\text{error}(h(x), f(x)) = 0 \quad \text{if } h(x) = f(x) \\ \text{else} \quad 1 \quad (0-1 \text{ loss})$$

**Regression:**  $Y$  is continuous, a numeric set (typically real numbers)

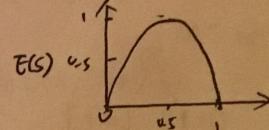
$$\text{error}(h(x), f(x)) = (h(x) - f(x))^2 \quad (\text{squared error})$$

**training error:**  $\text{error}_E(h) := \frac{1}{|E|} * \sum_{(x,y) \in E} \text{error}(y, h(x))$   
( $E \subseteq D$ ,  $E$  training error)

minimal training err doesn't entail minimal test err.

#### ▷ TDIDT. Decision Tree

$$\text{Entropy}(S) = -P_0 \log_2 P_0 - P_1 \log_2 P_1 \\ I(S) \text{ measures the impurity of } S$$



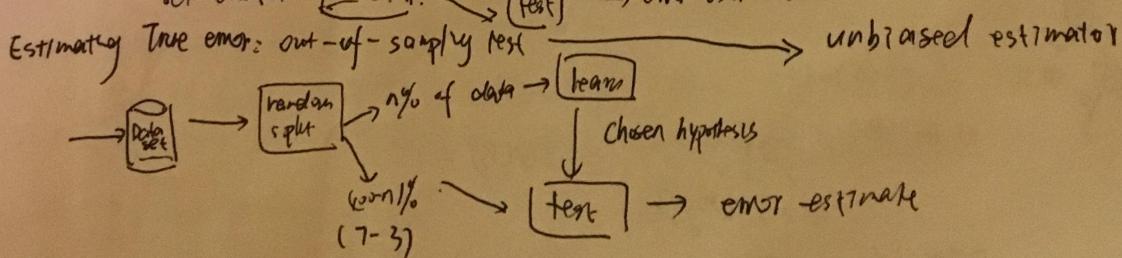
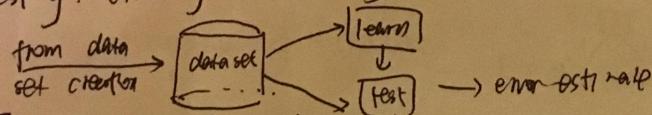
$$\text{Gain}(S, A) = \text{expected reduction in entropy due to sorting on } A \\ = \text{Entropy}(S) - \sum_{\text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Handling numeric data by:

- ① expert preprocessing (特征值划分)
- ② automatic preclustering (用 cluster 分组; 预聚类)
- ③ search at split time  
(sort attributes, create many binary attributes; compare and select best split using entropy/redundancy)

overfitting: error is lower than it would be on independent test set

▷ Full set testing: (testing on training set, training error)

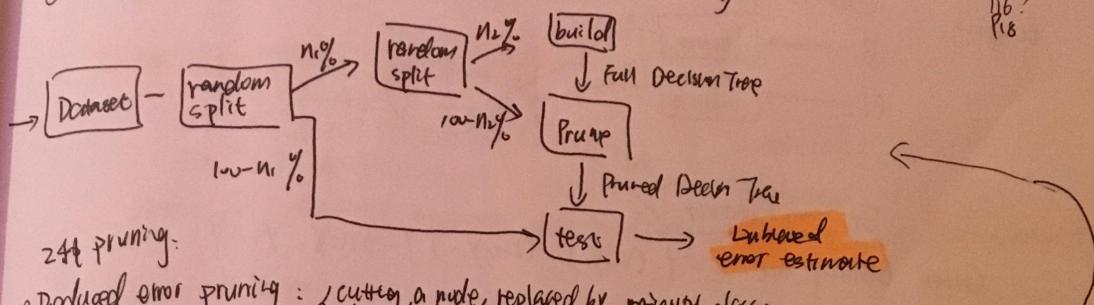


## Out-of-Sample Testing for Decision Tree Pruning

11-13 PPT

P6?

P18



Iterative Workflow: Example set  $E$ )

Sample := random sample ( $E, N\%$ )

$T = TDZDT(\text{sample})$ ,  $Err = \text{errors}(T, E)$

while  $Err \neq 0$ , And "Progress"

- Sample := Sample  $\sqcup$  random sample ( $E_{\text{err}}, M\%$ )

- $T = TDZDT(\text{Sample})$

- $Err = \text{errors}(T, E)$

return  $T$

Boosting: - The basic algorithm

- . Given training set  $E = \{(x_1, y_1), \dots, (x_m, y_m)\}$

- .  $D_i(x_i) = 1/m$ ,  $1 \leq i \leq m$

- . Repeat for  $t=1, \dots, T$ :

- . Learn hypothesis  $h_t$  on current weighted training set

- . Modify  $D_t$  so that incorrectly classified examples get emphasized

- $y_i: h_t(x_i) \neq y_i$ :  $D_{t+1}(x_i) = D_t(x_i) + \frac{\alpha}{m}$

- $y_i: h_t(x_i) = y_i$ :  $D_{t+1}(x_i) = D_t(x_i) - \frac{\alpha}{m}$

- . Combine hypotheses by weighted majority vote into one single accurate rule.

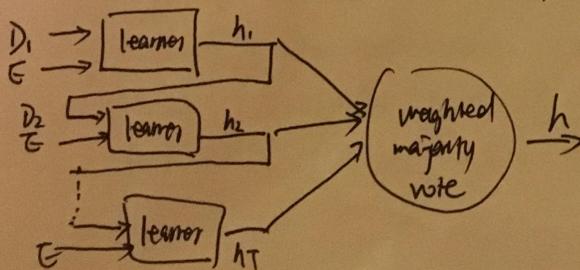
$$h(x) = \text{argmax}_y \sum_{t=1}^T h_t(x) = y \text{ (if } h_t(x) = y \text{ for all } t)$$

Bagging (Bootstrap aggregating): 随机选取  $n$  例样本来训练多个模型  $h_1, \dots, h_n$ .  
then. 各个  $H$  对分不采用投票方式.

boosting: 初始时对每一个训练例赋初权值  $\frac{1}{n}$ .

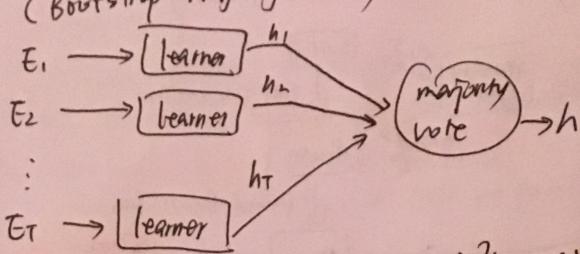
每次训练后，对训练失败的训练例赋予较大的权重，得到  $h_1, \dots, h_m$ .

预测时，权重大，反 $\downarrow$ 小. 最终  $H$  对分用 有放回投票.



$$\text{if } \left( y_i (\tilde{w}, \tilde{x}) + b \right) \leq 0$$

### Boosting. (Bootstrap Aggregation)



Given training set  $E = \{(x_1, y_1), \dots, (x_m, y_m)\} \quad x_i \in X, y_i \in Y$

Repeat for  $i=1 \dots T$ :

- Produce bootstrap replicates  $E_i$  of training set  $E$  by sampling with replacement s.t.  $|E_i| = |E|$

- Learn classifier  $h_i$  from replicate  $E_i$ .

Output majority vote hypothesis  $H$ .

$$H(x) = \text{argmax}_{y \in Y} \left| \{i \in \{1 \dots T\} | h_i(x) = y\} \right|$$

### Predictive Learning:

#### Function Approximation

Assume:  $X$ : instance space (fixed distribution  $D_X$ )

$Y$ : target space

function  $f: X \rightarrow Y$

A set of allowed hypotheses  $H$

Given (Input):

examples  $E \subseteq X \times Y$  such that

- $X$  drawn i.i.d. from  $D_X$
- $\forall (x, y) \in E, f(x) = y$

Find (Output):

A hypothesis  $h \in H$  such that the true error of  $h$

$$\text{error}_{D_X}(h) := E_{D_X} [\text{error}(f(x), h(x))]$$

↳ expected (avg) classification error on instances drawn according to  $D_X$  is minimal

$$\begin{matrix} x \\ f(x) \\ \vdots \\ x \\ f(x) \end{matrix}$$

$$\text{error}_{D_X}(h) = E_{D_X} [\text{error}(f(x), h(x))]$$

??

P17.18

## Concept learning Generality

### • Consistent hypothesis finding problem

For all  $e^+ = (x, y) \in E^+$ ,  $h(x) = +1$  - "completeness" of  $h$

For all  $e^- = (x, y) \in E^-$ ,  $h(x) = -1$  - "correctness" of  $h$

Finding a consistent hypothesis means finding a hypothesis with training error zero

### • Covers Relation

For an example  $e = (x, y)$ , we say

-  $h$  covers  $e$  iff  $h(x) = +1$  ( $h$  classified  $e$  as positive)

- if  $y = +1$ , — true pos

- if  $y = -1$ , — false pos

-  $h$  doesn't cover  $e$  iff  $h(x) = -1$  ( $h$  classified  $e$  as negative)

- if  $y = -1$ , — true neg

- if  $y = +1$ , — false neg

### • Generality.

A hypothesis (concept)  $h_1$  is more general (less specific) than hypothesis (concept)  $h_2$  iff  $h_1$  covers a superset of instances of  $h_2$ .

$$h_1 \leq h_2 \text{ iff } \{x \in X | h_1 \text{ covers } x\} \supseteq \{x \in X | h_2 \text{ covers } x\}$$

Most specific hypothesis  $h_{ms} = <\varphi, \varphi, \varphi, \varphi>$

Most general hypothesis  $h_{mg} = <?, ?, ?, ?, ?>$

$$\therefore h_1 \leq_{func} h_2 \text{ iff } h_1 \leq_{rule} \left[ \begin{array}{l} h_1 \leq_{rule} h_2 \\ \geq \end{array} \right] \text{ iff } h_1 \leq_{rule} \left[ \begin{array}{l} h_2 \leq_{rule} \\ \geq \end{array} \right] \text{ iff } h_2 \leq_{pred} \left[ \begin{array}{l} h_1 \leq_{pred} \\ \geq \end{array} \right]$$

P17.18 - 11-20

## Find-S Algorithm.

1. Initialize  $h$  to the most specific hypothesis in  $H$ .

2. For each positive training instance  $x$ .

    • For each attribute constraint  $a_i$  in  $h$

        if the constraint  $a_i$  in  $h$  is satisfied by  $x$

            Then do nothing

        Else replace  $a_i$  in  $h$  by the next more general constraint  
            that is satisfied by  $x$

3. Output hypothesis  $h$

## Find-S Disadvantages:

- Can't tell whether it has learned concept
- Can't tell when training data is inconsistent
- Picks a maximally specific  $h$ .
- Depending on  $H$ , there might be several

minimal generalization

"Progress  
random sample"

## - Version Spaces

$D$ : training examples      target concept  $\rightarrow c$   
 $\text{Consistent}(h, D) := h(x) = c(x) \text{ for all } (x, c(x)) \text{ in } D$

Version Space  $V_{H,D}$ .

$$V_{H,D} = \{h \in H \mid \text{Consistent}(h, D)\}$$

↳ with respect to hypothesis space  $H$  and training examples  $D$ , is no subset of hypothesis from  $H$  consistent with all training examples in  $D$

## - List-Then-Eliminate Algorithm

1. Version Space  $\leftarrow$  a list containing hypotheses in  $H$
2. For each training example,  $(x, c(x))$   
remove from Version Space any hypothesis for which  $h(x) \neq c(x)$
3. Output the list of hypotheses in Version Space.

General bound  $G$  of  $V_{H,D}$  is the set of its maximally general numbers

">": general       $G = \{g \in H \mid \text{Consistent}(g, D) \wedge (\forall g' \in H) [ (g > g') \wedge \text{Consistent}(g', D)]\}$

Specific bound  $S$ ,

$$S = \{s \in H \mid \text{consistent}(s, D) \wedge (\exists s' \in H) [(s > s') \wedge \text{consistent}(s', D)]\}$$

Every member of  $V_{H,D}$  lies between these boundaries

$$V_{H,D} = \{h \in H \mid \exists s \in S \ (\exists g \in G), g \leq h \leq s\}$$

## - Candidate Elimination Algorithm.

$G, S$ .

For each training example  $d$ , do

- if  $d$  is "+"
  - remove from  $G$  any hypothesis inconsistent with  $d$
  - For each hypothesis  $s$  in  $S$  that is not consistent with  $d$ 
    - remove  $s$  from  $S$
    - Add to  $S$  all minimal generalization  $h$  of  $s$  such that
      1.  $h$  is consistent with  $d$
      2. some members of  $G$  is more general than  $h$
    - remove from  $S$  any hypothesis that is more general than  $h$
- if  $d$  is "-"
  - remove from  $S$  any hypothesis inconsistent with  $d$
  - For each hypothesis  $g$  in  $G$  that is not consistent with  $d$ 
    - remove  $g$  from  $G$
    - Add to  $G$  all minimal specialization  $h$  of  $g$  such that
      1.  $h$  is consistent with  $d$
      2. some members of  $S$  is more specific than  $h$
    - Remove from  $G$  any hypothesis that is less general than another hypothesis in  $G$

Def. Let  $h_g, h_h$  be boolean-valued function defined over  $X$ .

Then,  $h_g$  is more-general-than-or-equal-to  $h_h$  ( $h_g \leq h_h$ ) if and only if

$$(\forall x \in X) [h_h(x)=1 \rightarrow (h_g(x)=1)]$$

Q1:

### Inductive Bias

concept learning algorithm  $L$

instances  $X$ , target concept  $c$ .

training examples  $D_c = \{(x_i, c(x_i))\}$

$L(X_i, D_c)$  denotes the classifier assigned to the instance  $x_i$  by

$L$  training on data  $D_c$

Def. the inductive bias of  $L$  is any minimal set of assertions  $B$  such that for any target concept  $c$  and corresponding training examples  $D$

$$(\forall x_i \in X) [(B \wedge D_c \wedge x_i) \vdash L(x_i, D_c)]$$

$A \vdash B$  means  $A$  logically entails  $B$

$B$  is a logical consequence (follow deductively, 即  $B$  可以由  $A$  推理得出)

$A > B$  表示  $B$  是  $A$  的逻辑推论 (包含在  $A$  的逻辑推论集中)

有偏概念 Rule learner < version space candidate.. < Tools

Q1:

Candidate  $\hookrightarrow$  目标概念 (包含在候选的假设空间  $H$  中)

Tiny PIB ICF  
用小的伪代码

correctness:  $G, S \quad VS_{H,D} == \text{Output}$

1) soundness  $VS_{H,D} \subseteq \text{Output}$

2) completeness  $\text{Output} \subseteq VS_{H,D}$

Proof  $D = \{e_1, \dots, e_n\}$

Math induction  $\star$  从  $e_1$  到  $e_n$

1)  $D = \emptyset \quad G_0 = \{\text{most general}\} \quad S_0 = \{\text{most specific}\}$

$$VS_{H,D} = H = \emptyset$$

2)  $D = \{e_1, \dots, e_l\} \quad VS_{H,D} = \text{output}$

$$3) D + \{e_{l+1}\}$$

3. 1)  $e_{l+1}$  is positive

$G_i$  and  $S_i$

check  $G_i$  for consistency

check  $S_i$  for consistency

3. 2)  $e_{l+1}$  is negative  $\#$

$G_{i+1}$  and  $S_{i+1} == VS_{H,D_{i+1}}$

Suppose that are not consistent with  $e_{l+1}$ ,  
for all  $s \in S_i$ , that don't cover  $e_{l+1}$ ,  
remove  $s$  from  $S_i$ ,

add to  $S_i$  all generalizations h.s.t.

- a)  $h$  covers  $e_{l+1}$
- b)  $\exists g \in G_{i+1}, g \subset h$

c)  $\nexists h' \quad h < h' \leq S$

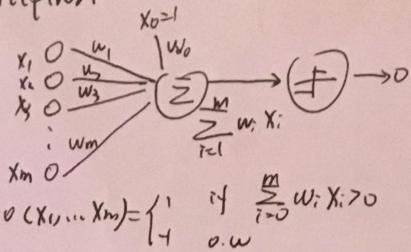
not satisfy  $e_{l+1}$

$$\sqrt{\frac{S}{G}}$$

## Neural Network

P9.

- Perception



perception training rule

$$w_i \leftarrow w_i + \alpha w_i$$

$$\Delta w_i = \eta (t - o) x_i$$

learning rate

- Gradient Descent:

$$O = w_0 + w_1 x_1 + \dots + w_m x_m$$

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \quad D: \text{set of examples}$$

$$\text{Gradient: } \nabla E(\vec{w}) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_m} \right]$$

$$\text{training rule: } \Delta \vec{w} = -\eta \nabla E(\vec{w}) \quad (\Delta w_i = -\eta \frac{\partial E}{\partial w_i})$$

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 = \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \end{aligned}$$

$$\frac{\partial E}{\partial w_i} = \sum_d (t_d - o_d) (-X_i, d)$$

- Sigma function: sigmoid unit

$$\sigma(x) = \frac{1}{1 + e^{-x}} \Rightarrow (0, 1) \quad \left| \begin{array}{l} \tanh h = \frac{e^h - e^{-h}}{e^h + e^{-h}} \Rightarrow (-1, 1) \\ \frac{d}{dx} \tanh h = 1 - \tanh^2 h \end{array} \right.$$

with multiple layers: ① arbitrary decision boundaries  
② including non-convex, disconnected areas

error gradient:

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_{d \in D} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \left( -\frac{\partial o_d}{\partial w_i} \right) \\ &= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial \text{net}_{d,i}} \frac{\partial \text{net}_{d,i}}{\partial w_i} \end{aligned}$$

$$\textcircled{1}: \frac{\partial o_d}{\partial \text{net}_{d,i}} = \frac{\partial \sigma(\text{net}_{d,i})}{\partial \text{net}_{d,i}} = o_d(1 - o_d)$$

$$\textcircled{2}: \frac{\partial \text{net}_{d,i}}{\partial w_i} = X_i, d$$

$$\therefore \frac{\partial E}{\partial w_i} = -\sum_{d \in D} (t_d - o_d) o_d(1 - o_d) X_i, d$$

Or: actual  
a: target

X\_i, d: input  
to

BP

① output

Δ w

- Error function with multiple output units

$$E(\vec{w}) = \frac{1}{2} \sum_{\text{outputs}} \sum_{\text{outputs}} (t_{k\text{out}} - o_{k\text{out}})^2$$

function is smooth  
smooth interpolation between data points

### Backpropagation Algorithm ①

Initialize all weights to small random numbers

Until satisfied, Do

- for each training example, do

1. Input the training example to the network and compute the network outputs

2. For each output unit  $k$ :  $\sigma'(x) = f(x) - f'(x)$

$$\delta_k \leftarrow o_k(1-o_k)(t_k - o_k)$$

3. For each hidden unit  $h$ :

$$\delta_h \leftarrow o_h(1-o_h) \sum_{\text{outputs}} w_{hk} \delta_k$$

$$\frac{\partial}{\partial x} f(x) = f'(x)$$

4. Update each network weight  $w_{i,j}$ :

$$w_{i,j} \leftarrow w_{i,j} + \alpha w_{i,j}, \text{ where } \alpha w_{i,j} = \eta \delta_j x_{i,j}$$

$o_k$ : actual output of a node  $k$

$x_{i,j}$ : input from node  $i$  to node  $j$

- Every boolean function can be represented by network with a single hidden layer, but require exponential hidden units
- Every bounded continuous function can be approximated with arbitrary small error with one hidden layer.
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers.

Random weight initialization ✓✓

① necessary

② reduce the risk of the network getting trapped in a local minimum

③ not too large to avoid hitting the saturation region of sigma-function

not too small to not reduce the network to a linear network

### Gradient descent Algorithm

Gradient-descent (training-examples,  $\eta$ )

### \* linear sigma network perceptron

- Initialize each  $w_i$  to some small random value

Until the termination condition is met, Do

• Initialize each  $\Delta w_i$  to zero

• For each  $(\vec{x}, t)$  in training-examples, Do

• Input the instance  $\vec{x}$  to the unit and compute the output  $o$

• For each linear unit weight  $w_i$ , Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

• For each linear unit weight  $w_i$ , Do

$$w_i \leftarrow w_i + \Delta w_i$$

BP at output layer  
hidden layers off

④ output

$\Delta w$

Bayes Optimal classifier: 俗名叫朴素贝叶斯，对新数据分类时只看父节点

$$\arg \max_{V \in V} \sum_{h_i \in H} P(V_j | h_i) P(h_i | D)$$

$$P(V_j | D) = \sum_{h_i \in H} P(V_j | h_i) P(h_i | D)$$

V: 所有分类集合

Naive Bayes classifier:

Assume target function  $f: X \rightarrow V$ , where each instance  $x$  described by attributes  $\langle a_1, a_2, \dots, a_n \rangle$ .

Most probable value of  $f(x)$  is:

$$\begin{aligned} V_{MAP} &= \arg \max_{V \in V} P(V_j | a_1, \dots, a_n) \\ &= \arg \max_{V \in V} \frac{P(a_1, \dots, a_n | V_j) P(V_j)}{P(a_1, \dots, a_n)} \\ &= \arg \max_{V \in V} \frac{\underbrace{P(a_1, \dots, a_n | V_j)}_{\prod P(a_i | V_j)} P(V_j)}{\underbrace{P(a_1, \dots, a_n)}_{\prod P(a_i)}} \end{aligned}$$

Naive Bayes classifier:

$$V_{NB} = \arg \max_{V \in V} P(V_j) \prod_{i=1}^n P(a_i | V_j)$$

#从 example 到 NDA

NB Algorithm:

1. Naive Bayes Learn (examples)

1. For each target value  $V_j$

1.  $\tilde{P}(V_j) \leftarrow \text{estimate } P(V_j)$

1. For each attribute value  $a_i$  of each attribute  $a$

$\tilde{P}(a_i | V_j) \leftarrow \text{estimate } P(a_i | V_j)$

Classify New instance ( $x$ )

$$V_{NB} = \arg \max_{V \in V} \tilde{P}(V_j) \prod_{i=1}^n \tilde{P}(a_i | V_j)$$

If none of training instances with target value  $V_j$  have attribute  $a_i$ ? Then

so, we use  $\tilde{P}(a_i | V_j) = 0 \dots P(V_j) \prod_i P(a_i | V_j) = 0$

$$\tilde{P}(a_i | V_j) \leftarrow \frac{n_c + m_p}{n + m}$$

n: number of training examples for which  $V = V_j$

$n_c$ : number of examples for which  $V = V_j$  and  $a = a_i$

P: Prior estimate for  $P(a_i | V_j) = \frac{1}{m}$

m: weight given to prior (i.e. number of virtual examples)

## Bayesian Belief Networks

描述一组变量之间的联合分布  
 ↳ describe conditioned independence among subsets of variables  
 → allow combining prior knowledge about independencies  
 among variables with observed training data

Def:  $X$  is conditionally independent of  $Y$  given  $Z$  if the probability distribution governing  $X$  is independent of the value of  $Y$  given the value of  $Z$ : that is, if

$$(\forall x_i, y_j, z_k) P(X=x_i | Y=y_j, Z=z_k) = P(X=x_i | Z=z_k)$$

$$\quad \text{① } P(X|Y, Z) = P(X|Z)$$

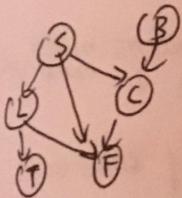
$\Leftarrow$  DNB use:  $P(X, Y | Z) = P(X|Y, Z) \cdot P(Y|Z)$

↓↑↑↑↑↑

$$= P(X|Z) \cdot P(Y|Z)$$

$$P(Y_1, Y_2, \dots, Y_n) = \prod_{i=1}^n P(Y_i | \text{Parents}(Y_i))$$

where  $\text{Parents}(Y_i)$  denotes immediate predecessors of  $Y_i$  in graph.



Kernel Method: a class of pattern detecting algorithms

domain independence: generic algorithms parameterized by kernels

computational efficiency: runtime - polynomial

robustness: handle noise data, identify good approximations

statistical stability:

- high-dimension
- (i) Feature space (data is embedded into an inner product space)
  - (ii) Patterns: linear functions in the feature space
  - (iii)  $\sim$  pairwise inner product
  - (iv)  $\sim$  can be computed by a function  $\rightarrow$  Kernel

Kernel methods: a class of pattern detecting algorithms

Key Aspects

- (i) data is embedded into a inner product space  $\Rightarrow$  feature space (high-dim)
- (ii) Patterns: linear functions in the feature space  
(usually correspond to nonlinear functions in the original space)

Def. A dichotomy  $(S^+, S^-)$  of  $S \subseteq \mathbb{R}^d$  is linearly separable if there exists a  $(d-1)$ -dim affine hyperplane separating  $S^+$  and  $S^-$ .

Def. A set  $S \subseteq \mathbb{R}^d$  of  $n$  points are in general position if

- (i)  $n > d$  and  $S$  is not contained by an  $(n-2)$ -dim affine hyperplane
- (ii)  $n > d$  and  $S$  has no subset of  $d+1$  points belonging to a  $(d-1)$ -dim affine hyperplane.

Cover's (Function Counting) Theorem:

The number of linearly separable dichotomies of  $n$  points in general position in  $\mathbb{R}^d$  is

$$C(n, d) = 2 \sum_{k=0}^d \binom{n-1}{k} \quad \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Probability of Linear Separability:

$P(n, d)$ : Prob of (A dichotomy selected uniformly at random from a set of  $n$  points in general position in  $\mathbb{R}^d$  is linearly separable)

$$P(n, d) = \frac{1}{2^n} C(n, d) = \frac{1}{2^n} 2 \sum_{k=0}^d \binom{n-1}{k}$$

Implications of Cover's Theorem:

If training data is not linearly separable, with high probability it can be embedded into a higher dimension space

via some non-linear transformation in which it becomes linearly separable.

$$\left. \begin{array}{l} \text{if } n \leq d+1: C(n, d) = 2 \sum_{k=0}^{d+1} \binom{n-1}{k} = 2^n \Rightarrow P(n, d) = 1 \\ \text{if } n = 2(d+1), P(n, d) = \frac{1}{2^{2d+2}} \cdot 2 \sum_{k=0}^d \binom{2d+1}{k} = \frac{1}{2} \\ \text{if } n \gg d, P(n, d) \approx \frac{C(n,d)}{2^n} \text{ for some constant } C > 0 \end{array} \right) P(n, d) = \begin{cases} 1 & \text{if } n \leq d+1 \\ \frac{1}{2} & \text{if } n = 2(d+1) \\ 0 & \text{if } n \geq d \end{cases}$$

So, for any  $\epsilon > 0$

$$\lim_{d \rightarrow \infty} P(2(d+1)(1-\epsilon), d) = 1 \quad \text{phase transition at } n = 2(d+1)$$
$$\lim_{d \rightarrow \infty} P(2(d+1)(1+\epsilon), d) = 0$$

The

Empirical

A prob

Regula

Tikhonov

Ridge

"b

X

Cure of

Kernel meth

The learning problem:

Given: • a set of training data  $S = \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq X \times Y$  generated by a probability distribution  $P(X, Y)$ , and

Find a function  $f \in H$ , minimizing the true error.

$$R[f] = \int_{X \times Y} V(f(x), y) dP(x, y)$$

loss function

Empirical Risk Minimization:

$$\text{empirical error: } R_{\text{emp}}[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i)$$

empirical risk minimization:  $\arg \min_{f \in H} R_{\text{emp}}[f]$  (select a  $h$  that minimizes the training error)

- A problem is well-posed if it satisfies
- 1. existence: the problem has a solution
  - 2. uniqueness: the solution must be unique
  - 3. stability: the solution depends continuously on the data  
delta  $\Rightarrow$  unstable
- now: is ill-posed

- Regularization theory:  
↳ method improving stability of solutions of ill-posed inverse problems

Tikhonov regularization: transforms the original problem into a conditional-optimization problem into a conditional optimization problem that penalizes undesirable solutions by adding a term, called regularizer.

regularizer: enforce the solution to be smooth.

$$\underset{\lambda}{\arg \min} R_{\text{emp}}[f] + \lambda S_2[f]$$

$\lambda \in \mathbb{R}$ : regularization parameter controlling trade-off between

- (i) minimizing the empirical risk
- (ii) minimizing the model complexity

### Ridge Regression:

Given a sample  $S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\} \subseteq \mathbb{R}^d \times \mathbb{R}$ , find a linear function

$$f(\vec{x}) = \vec{x}^\top \vec{w} = \sum_{i=1}^d x_i w_i$$

+ best interpolates  $S$ .

equals: find a vector  $\vec{w} \in \mathbb{R}^d$   $f(\vec{x}) = \langle \vec{x}, \vec{w} \rangle = \vec{x}^\top \vec{w} = \sum_{i=1}^d x_i w_i$

"best interpolates": for  $(\vec{x}, y)$ ,  $V(f(\vec{x}), y) = (f(\vec{x}) - y)^2$

$$\begin{aligned} \text{empirical risk: } R_{\text{emp}}[f] &= \frac{1}{n} \sum_{i=1}^n (f(\vec{x}_i) - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (\langle \vec{x}_i, \vec{w} \rangle - y_i)^2 \\ &= \frac{1}{n} \langle \mathbf{X} \vec{w} - \vec{y} \rangle^\top \langle \mathbf{X} \vec{w} - \vec{y} \rangle = \frac{1}{n} \| \mathbf{X} \vec{w} - \vec{y} \|^2 \end{aligned}$$

$$\min_{\vec{w}} \frac{1}{n} \| \mathbf{X} \vec{w} - \vec{y} \|^2 + \lambda \| \vec{w} \|^2$$

$$= \min_{\vec{w}} R_\lambda(\vec{w}, S) = \min_{\vec{w}} \| \mathbf{X} \vec{w} - \vec{y} \|^2 + \lambda \| \vec{w} \|^2 \text{ is convex in } \vec{w}$$

Curse of dimensionality:

(i) potentially enormous time complexity of computing inner products

(ii) overfitting ↗

Kernel methods solve

(i) Kernel trick

(inner product can be calculated without computing the embedding)

(ii) Regularization

(automatic control of overfitting)

## Ridge Regression

primal

$$\vec{w} = (X^T X + \lambda I_d)^{-1} X^T \vec{y}$$

$$f(\vec{x}) = \vec{x}^T \vec{w}$$

Complexity:

computation of  $\vec{w}$ :  $O(d^3)$

evaluation of  $f$ :  $O(d)$

dual is faster for the first step if  $d > n$ , but evaluation is always slower.

crucial obs

About the dual, matrix  $G = X X^T \Rightarrow$  Gram matrix, is given by  
inner products of the training points,

## Kernel Trick:

Learning in Feature Space

motived example: dual of ridge regression

$$(i) \vec{d} = (G + \lambda I_n)^{-1} \vec{y} \text{ with } G_{ij} = \langle \vec{x}_i, \vec{x}_j \rangle$$

$$(ii) f(\vec{x}) = \sum_{i=1}^n \alpha_i \langle \vec{x}, \vec{x}_i \rangle$$

transformation:  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$

then, work automatically in feature space

- (1) replace all entries  $G_{ij} = \langle \vec{x}_i, \vec{x}_j \rangle$  by  $\langle \phi(\vec{x}_i), \phi(\vec{x}_j) \rangle$
- (2) for (ii), use  $f(\vec{x}) = \sum_{i=1}^n \alpha_i \langle \phi(\vec{x}), \phi(\vec{x}_i) \rangle$

## Inner Product (Pre-Hilbert) Spaces

Def. a vector space  $H$  is an inner product space if it is equipped with a function  $\langle \cdot, \cdot \rangle: H \times H \rightarrow \mathbb{R}$

satisfying

$$\langle \vec{x}, \vec{y} \rangle = \langle \vec{y}, \vec{x} \rangle$$

$$\langle \vec{x}, \vec{x} \rangle \geq 0 \text{ with } \langle \vec{x}, \vec{x} \rangle = 0 \text{ if } \vec{x} = \vec{0}$$

$$\langle \lambda_1 \vec{x} + \lambda_2 \vec{y}, \vec{z} \rangle = \lambda_1 \langle \vec{x}, \vec{z} \rangle + \lambda_2 \langle \vec{y}, \vec{z} \rangle$$

for all  $\vec{x}, \vec{y} \in H$  and scalars  $\lambda_1, \lambda_2$

the standard inner product in  $\mathbb{R}^d$ , i.e.

$$\langle \vec{x}, \vec{y} \rangle = \sum_{i=1}^d x_i y_i$$

$$\begin{aligned} \langle \lambda_1 \vec{x} + \lambda_2 \vec{y}, \vec{z} \rangle &= \sum_{i=1}^d (\lambda_1 x_i + \lambda_2 y_i) z_i \\ &= \sum_{i=1}^d \lambda_1 x_i z_i + \sum_{i=1}^d \lambda_2 y_i z_i \\ &= \lambda_1 \langle \vec{x}, \vec{z} \rangle + \lambda_2 \langle \vec{y}, \vec{z} \rangle \end{aligned}$$

dual solution  $\vec{\alpha} = \frac{1}{\lambda} (\vec{y} - X \vec{w})$   
 $\vec{w} = X^T \vec{\alpha}$

dual

$$\vec{\alpha} = (X X^T + \lambda I_n)^{-1} \vec{y}$$

$$f(\vec{x}) = \sum_{i=1}^n \alpha_i \langle \vec{x}, \vec{x}_i \rangle$$

computation of  $\vec{\alpha} = O(n^3)$

evaluation of  $f$ :  $O(nd)$

The

Kernel

example

Normalized (1/m)

eigenvalues

Properties of

Cauchy-

Properties of

## The Kernel Trick.

Def. a kernel is a function  $X \times X \rightarrow \mathbb{R}$  such that for all  $x, y \in X$

$$k(x, y) = \langle \phi(x), \phi(y) \rangle$$

Kernel trick: substitute all occurrences of  $\langle \cdot, \cdot \rangle$  by a kernel  $k$  with  
 $k(x, y) = \langle \phi(x), \phi(y) \rangle$   
where  $\phi$  is the underlying function mapping the input space into the feature space.

example. let  $k: X \times X \rightarrow \mathbb{R}$  with  $X \subseteq \mathbb{R}^d$  be defined by

$$k(\vec{x}, \vec{y}) = \langle \vec{x}, \vec{y} \rangle^2 \text{ for all } \vec{x}, \vec{y} \in X$$

claim:  $k$  is a kernel corresponding to the feature map  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d^2}$  defined by

$$\phi: \vec{z} \mapsto (z_i z_j)_{i,j=1}^d \text{ for all } \vec{z} \in \mathbb{R}^d$$

$$\langle \phi(\vec{x}), \phi(\vec{y}) \rangle = \langle (x_i x_j)_{i,j=1}^d, (y_i y_j)_{i,j=1}^d \rangle$$

$$= \sum_{i,j=1}^d x_i x_j y_i y_j = \sum_{i=1}^d x_i y_i \sum_{j=1}^d x_j y_j$$

$$= \langle \vec{x}, \vec{y} \rangle^2$$

## Normalized linear space:

Def: a vector space  $H$  is a normalized linear space, if there exists a function  $\|\cdot\|: H \rightarrow \mathbb{R}^+$  satisfying

$$\|\vec{x}\| \geq 0 \text{ with } \|\vec{x}\| = 0 \text{ if } \vec{x} = \vec{0}$$

$$\|\vec{x} + \vec{y}\| \leq \|\vec{x}\| + \|\vec{y}\|$$

$$\|\lambda \vec{x}\| = |\lambda| \cdot \|\vec{x}\| \text{ for all } \vec{x}, \vec{y} \in H \text{ and scalar } \lambda.$$

Def. for a normalized linear space  $H$ , the function  $d(\vec{x}, \vec{y}) = \|\vec{x} - \vec{y}\|$  defines a distance on  $H$ .

eigenvalue 不考, 但用自己解

Properties of Inner Products  $\|\vec{x}\| = \sqrt{\langle \vec{x}, \vec{x} \rangle}$  (inner product space  $H$  is a normalized linear space)

Cauchy-Schwarz Inequality: for any  $\vec{x}, \vec{y} \in H$ .

$$|\langle \vec{x}, \vec{y} \rangle| \leq \|\vec{x}\| \cdot \|\vec{y}\|$$

$$\Rightarrow 1 \leq \frac{\langle \vec{x}, \vec{y} \rangle}{\|\vec{x}\| \|\vec{y}\|} \leq 1$$

$$\text{or } 0 = \frac{\langle \vec{x}, \vec{y} \rangle}{\|\vec{x}\| \|\vec{y}\|}$$

Properties of Kernels: Let  $k: X \times X \rightarrow \mathbb{R}$  be a kernel function with  $k(x, y) = \langle \phi(x), \phi(y) \rangle$ , for  $\phi$ , and ref

$$\alpha(x, y) = \frac{k(x, y)}{\sqrt{k(x, x)} \sqrt{k(y, y)}}$$

$$\text{then } \alpha(x, y) = \frac{\langle \phi(x), \phi(y) \rangle}{\sqrt{\langle \phi(x), \phi(x) \rangle} \sqrt{\langle \phi(y), \phi(y) \rangle}} = \cos(\phi(x), \phi(y))$$

$\Rightarrow$  normalized kernels = cosine similarity in the feature space

### Characterization of Kernels: Finite Domain

( $A^{(n \times n)}$  is positive semi-definite, iff all of its eigenvalues are non-negative)  
 Gramm / Kernel matrix  $K_{X,k}$  of  $X$ , wrt  $k$  is  $n \times n$  matrix,  $(K_{X,k})_{ij} = k(x_i, x_j)$   
 $\rightarrow (k(x_i, x_j) = k(y, x_j))$

Theorem (finite case): Let  $k: X \times X \rightarrow \mathbb{R}$  be a symmetric function. Then

$k$  is a kernel  $\Leftrightarrow K_{X,k}$  is positive  $\Rightarrow$  semi-definite

example:  $k = \{0, 1, 2\} \times \{0, 1, 2\} \rightarrow \mathbb{R}$  defined by

$$k(x, y) = \begin{cases} 1 & \text{if } |x-y| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

is not a kernel, because:

$$K_{X,k} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

it has negative eig.,  $\Rightarrow$  is not positive semi-definite

### Characterization of Kernels: Infinite Domain

let  $X$  be a compact subset of  $\mathbb{R}^n$  and let  $k: X \times X \rightarrow \mathbb{R}$  be a symmetric and continuous function, then:

?  $k$  is a kernel  $\Leftrightarrow$  for all finite  $Y \subseteq X$ :  $K_{Y,k}$  is positive semi-definite  
 $k(x,y) = \langle \phi(x), \phi(y) \rangle$

#### Kernel Construction

(i)  $k(x,y) = f(x)f(y)$  is a kernel over  $X \times X$  for all functions  $f: X \rightarrow \mathbb{R}$

Proof. let  $\phi: X \rightarrow \mathbb{R}$  defined by  $\phi = x \mapsto f(x)$  for all  $x \in X$   
 $\Rightarrow k(x,y) = f(x)f(y) = \langle \phi(x), \phi(y) \rangle$

(ii) Let  $A$  be a symmetric positive semi-definite  $n \times n$  matrix, then

$$k(\vec{x}, \vec{y}) = \vec{x}^T A \vec{y} \text{ for all } \vec{x}, \vec{y} \in \mathbb{R}^n \text{ is a kernel over } \mathbb{R}^n \times \mathbb{R}^n$$

Proof. by spectral theorem,  $A = U^T \Lambda U$  ( $U$  orthogonal matrix,  $\Lambda$  diagonal matrix)

- all entries in  $A$  are non-negative, as  $A$  is positive semi-definite
- $\Rightarrow$  we can take  $\sqrt{\Lambda}$  (matrix with the square roots of the eigenvalues)
- $\Rightarrow$  for  $B = \sqrt{\Lambda} U$ ,  $\phi: \mathbb{R} \mapsto B\vec{x}$ ,

$$k(\vec{x}, \vec{y}) = \vec{x}^T A \vec{y} = \vec{x}^T U^T \Lambda U \vec{y} = \vec{x}^T B^T B \vec{y} = \langle B\vec{x}, B\vec{y} \rangle = \langle \phi(\vec{x}), \phi(\vec{y}) \rangle$$

(iii) Let  $k_1, k_2$  be kernels over  $X \times X$ . Then for all  $\alpha, \beta \geq 0$

$k(x, y) = \alpha k_1(x, y) + \beta k_2(x, y)$  is a kernel

Proof. for all finite  $Y = \{x_1, \dots, x_n\} \subseteq X$  and for all  $\vec{a} \in \mathbb{R}^n$

$$\vec{a}^T K_{Y,k_1} \vec{a} \geq 0 \text{ and } \vec{a}^T K_{Y,k_2} \vec{a} \geq 0$$

Mercer's theorem, it suffices to show that  $\vec{a}^T (\alpha K_{Y,k_1} + \beta K_{Y,k_2}) \vec{a} \geq 0$

$$\vec{a}^T (\alpha K_{Y,k_1} + \beta K_{Y,k_2}) \vec{a} = \vec{a}^T (\alpha K_{Y,k_1}) \vec{a} + \vec{a}^T (\beta K_{Y,k_2}) \vec{a}$$

$$= \alpha \underbrace{\vec{a}^T K_{Y,k_1} \vec{a}}_{\geq 0} + \beta \underbrace{\vec{a}^T K_{Y,k_2} \vec{a}}_{\geq 0} \geq 0$$

(iv): ~~let  $k = k_1 + k_2$~~

(iv). Let  $k_1, k_2$  be kernels over  $X \times X$ , then

$$k(x, y) = k_1(x, y) k_2(x, y) \text{ is a kernel}$$

Proof:  $A \otimes B$ : 矩乘

$$A \otimes B = \begin{pmatrix} a_{11}B & \dots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{n1}B & \dots & a_{nn}B \end{pmatrix} \in \mathbb{R}^{n \times n^2}$$

$A \otimes B$  is positive semi-definite  $\Rightarrow A \otimes B$  is also  $\oplus$  semi-definite

Let  $Y = \{x_1, \dots, x_n\} \subseteq X$  be a finite set

$\Rightarrow$  Mercer's theorem: it suffices to show  $K_{Y, k}$  is positive definite

- $K_{Y, k} = K_{Y, k_1} \cup K_{Y, k_2}$

- $k_1, k_2$  are kernels,  $K_{Y, k_1} \cup K_{Y, k_2}$

- $K_{Y, k}$  is a principle matrix of  $K_{Y, k_1} \otimes K_{Y, k_2}$   
 $= K_{Y, k_1} \otimes K_{Y, k_2}$

- $K_{Y, k}$  is positive semi-definite.

(v) Let  $k_1$  be a kernel over  $X \times X$ , and  $p$  be a polynomial with positive coefficients. Then

$$k(x, y) = p(k_1(x, y)) \text{ is a kernel.}$$

Proof. It follows from (i), (ii)-(iv)

need (i) to cover the constant term in the polynomial

(vi). Let  $k_1: X \times X \rightarrow \mathbb{R}$  be a kernel and let  $f: \mathbb{R} \rightarrow \mathbb{R}$  be a function expressed as power series with non-negative coefficients only. Then

$$k(x, y) = f(k_1(x, y)) \text{ is a kernel.}$$

Proof.  $f = \lim_{n \rightarrow \infty} p_n$  for polynomials with non-negative coefficients

$\Rightarrow$  for any finite  $Y \subseteq X$ :  $K_{Y, f(k_1)} = \lim_{n \rightarrow \infty} K_{Y, p_n(k_1)}$

by (v)  $K_{Y, p_n(k_1)}$  is positive definite for all  $n$ , and hence for all  $D$ .

$$\nabla^T K_{Y, f(k_1)} \nabla = \nabla^T (\lim_{n \rightarrow \infty} K_{Y, p_n(k_1)}) \nabla = \lim_{n \rightarrow \infty} D^T K_{Y, p_n(k_1)} D \geq 0$$

(vii) Let  $k_1: X \times X \rightarrow \mathbb{R}$  be a kernel. Then

$$k(x, y) = e^{k_1(x, y)} \text{ is a kernel}$$

Proof: immediate from (vi), as all coefficients of  $e^{k_1(x, y)} = \sum_{n=0}^{\infty} \frac{k_1(x, y)^n}{n!}$  are non-negative

(viii) The function  $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  defined by

$$k(\vec{x}, \vec{y}) = \exp\left(-\frac{\|\vec{x} - \vec{y}\|_2^2}{2\sigma^2}\right) \text{ for all } \vec{x}, \vec{y} \in \mathbb{R}^d$$

is a kernel for any  $d$ , and for any  $\sigma > 0$ . It is called Gaussian or radial basis function (RBF) kernel.

Proof.  $k'(\vec{x}, \vec{y}) = \exp\left(-\frac{\langle \vec{x}, \vec{y} \rangle}{\sigma^2}\right)$

$$\frac{k'(\vec{x}, \vec{y})}{k'(\vec{x}, \vec{x}) k'(\vec{y}, \vec{y})} = \underbrace{\exp\left(-\frac{\langle \vec{x}, \vec{y} \rangle}{\sigma^2}\right)}_{k(\vec{x}, \vec{y})} \cdot \underbrace{\frac{\exp(-\frac{\langle \vec{x}, \vec{x} \rangle}{\sigma^2}) \cdot \exp(-\frac{\langle \vec{y}, \vec{y} \rangle}{\sigma^2})}{\exp(-\frac{\langle \vec{x}, \vec{y} \rangle}{\sigma^2})}}_{(i)}$$

common Kernel functions over  $\mathbb{R}^d \times \mathbb{R}^d$

linear kernel:  $k(\vec{x}, \vec{y}) := \vec{x}^T \vec{y}$

poly... kernel:  $k(\vec{x}, \vec{y}) := (\vec{x}^T \vec{y} + c)^k$

Gaussian or RBF:  $k(\vec{x}, \vec{y}) = \exp\left(-\frac{\|\vec{x} - \vec{y}\|_2^2}{2\sigma^2}\right)$

## Support Vector Machine

- Choose hyperplane with the largest margin  $\gamma$ .

- (1) Robust against noise
- (2) Excellent predictive performance in practice
- (3) Separating hyperplane becomes unique.
- (4) Good generalization is at odds with capacity  
minimize the structural risk (high  $\alpha \Rightarrow$  high structural risk)

For any  $\gamma > 0$ , and for any  $S \subset \mathbb{R}^d$ :  $\|\vec{w}\| \leq R/\gamma$  that can be  $\gamma$ -shattered by the hyperplanes.

$$|S| \leq \min\left(\left(\frac{R}{\gamma}\right)^2, d\right) + 1$$

$\Rightarrow$  bound on the VC-dim of hyperplanes separating with margin  $\geq \gamma$   
- VC-dim of half-spaces for "orthogonal" shattering:  $d+1$

$\Rightarrow$  Larger margin implies smaller VC-dim.

$\Rightarrow$  Max margin minimizes the structural risk of separating hyperplanes.

- Hard Margin SVM

hyperplane for  $f(\vec{x}) = \langle \vec{w}, \vec{x} \rangle + b = 0$

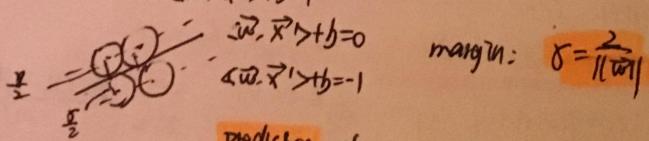
(1)  $\vec{w}$  is orthogonal to the hyperplane

(2) Signed distance of a point  $\vec{x}'$  from the hyperplane

$$\gamma = \frac{|f(\vec{x}')|}{\|\vec{w}\|}$$

- The Maximum Margin hyperplane.

hyperplane  $\langle \vec{w}, \vec{x} \rangle + b = 0$ : scale  $\vec{w}$  and  $b$  such that  $|\langle \vec{w}, \vec{x}' \rangle + b| = 1$  for all points on the closest hyperplanes



Prediction of an unseen instance  $\vec{x}$ :  $\text{sign}(\langle \vec{w}, \vec{x} \rangle + b)$

- Non-convex optimization problem

$$\max_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2$$

subject to  $|\langle \vec{w}, \vec{x}_i \rangle + b| \geq 1$  for  $i=1, \dots, n$

equivalent formulation ( $\max \frac{1}{2} \|\vec{w}\|^2 = \min \frac{1}{2} \sum \|\vec{w}\|^2$ )

Primal optimization problem:  $\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2$

$$\text{s.t. } -[y_i (\langle \vec{w}, \vec{x}_i \rangle + b) - 1] \leq 0 \text{ for } i=1, \dots, n$$

### Dual optimization Problem

$$\max_{\vec{\alpha}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

s.t.  $\sum_{i=1}^n \alpha_i y_i = 0$  and  $\alpha_i \geq 0 \quad i=1, \dots, n$

- Soft Margin SVM. (data is not linearly separable  
notes: allow violations, but penalize them)

### ① Optimization problem with $C > 0$

$$\min_{\vec{w}, b, \vec{s}} C \sum_{i=1}^n s_i + \frac{1}{2} \|\vec{w}\|^2$$

s.t.  $y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - s_i, \quad s_i \geq 0 \quad i=1, \dots, n$

$(\begin{array}{l} s_i = 0: \text{correct classification} \\ 0 < s_i \leq 1: \text{lives inside the margin, but on correct side} \\ s_i > 1: \text{lives on the wrong side} \end{array})$

$\sum_{i=1}^n s_i$  is an upper bound on the # of misclassified points

### ② Dual form:

$$\max_{\vec{\alpha}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

s.t.  $0 \leq \alpha_i \leq C$  and  $\sum_{i=1}^n \alpha_i y_i = 0, \quad i=1, \dots, n$