

**Two SVM tutorials linked in class website
(please, read both):**

- High-level presentation with applications (Hearst 1998)
- Detailed tutorial (Burges 1998)

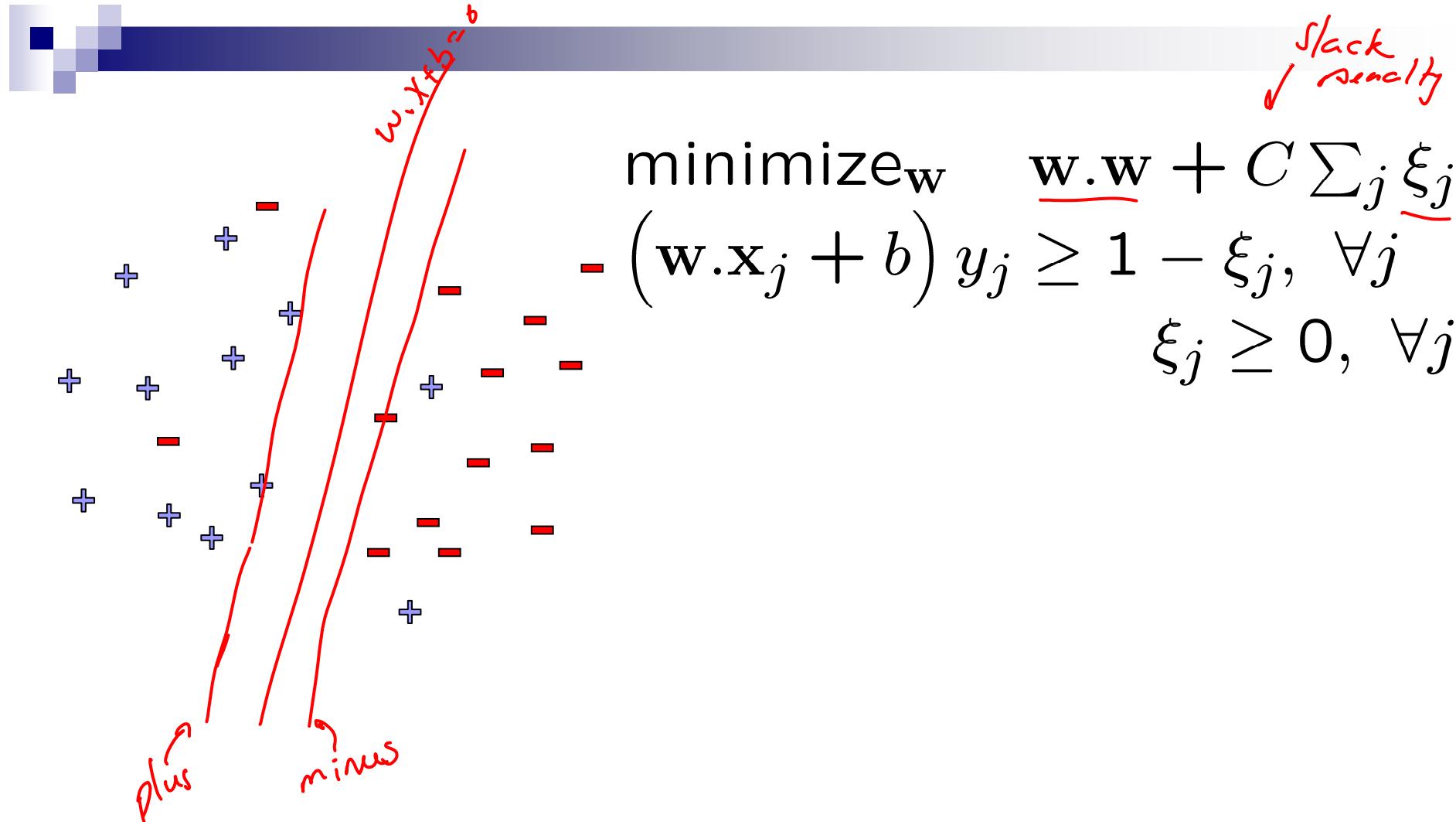
SVMs, Duality and the Kernel Trick (cont.)

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University

March 1st, 2006

©2006 Carlos Guestrin

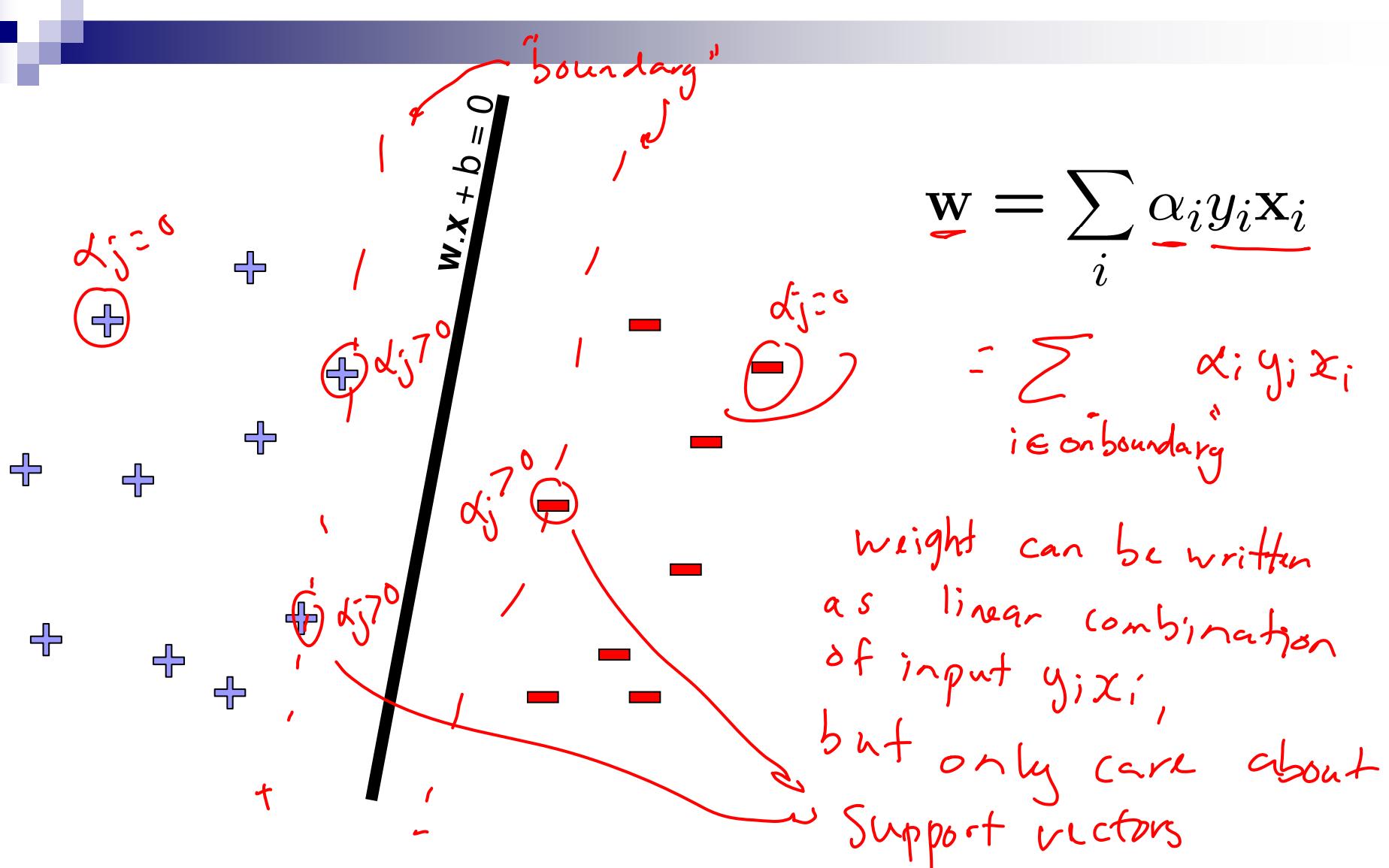
SVMs reminder



Today's lecture

- Learn one of the most interesting and exciting recent advancements in machine learning
 - The “kernel trick”
 - High dimensional feature spaces at no extra cost!
- But first, a detour
 - Constrained optimization!

Dual SVM interpretation



Dual SVM formulation – the linearly separable case

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

$$\begin{aligned} \sum_i \alpha_i y_i &= 0 \\ \alpha_i &\geq 0 \end{aligned}$$

dual program, ^{solve} SVM :
→ solve dual
obtain the α

get w, b

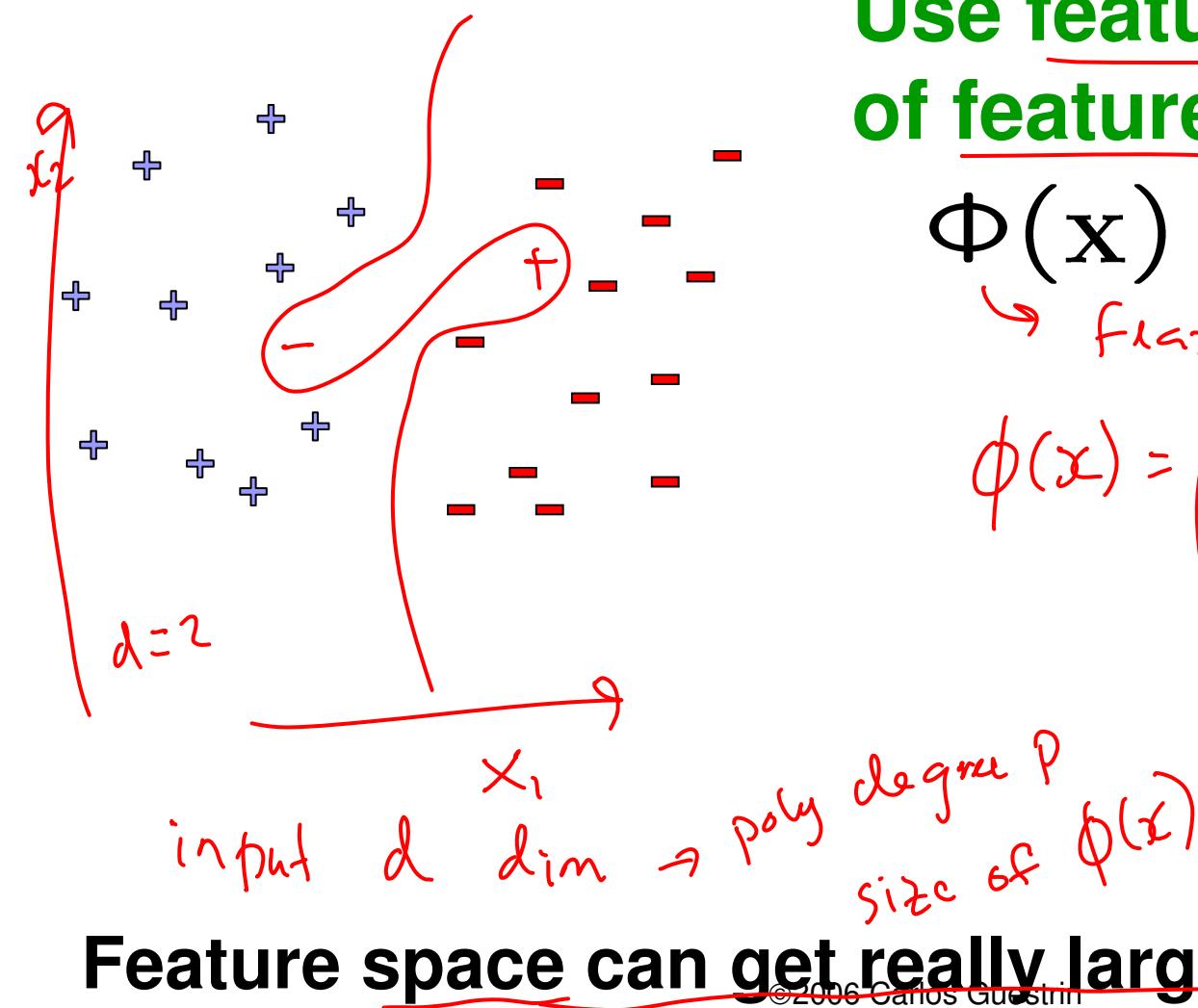
$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

for any k where $\alpha_k > 0$

obj function dual → quadratic → dual quadratic program

Reminder from last time: What if the data is not linearly separable?



Use features of features
of features of features....

$$\Phi(x) : \mathbb{R}^m \mapsto F$$

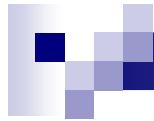
feature mapping $x = (x_1, x_2)$

$$\Phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_1^2 \\ x_1^3 \\ \vdots \\ x_1^P \end{pmatrix}$$

$$\Phi(x) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ x_1^3 \\ x_2^3 \\ x_1 x_2^2 \\ x_1^2 x_2 \\ \vdots \end{pmatrix}$$

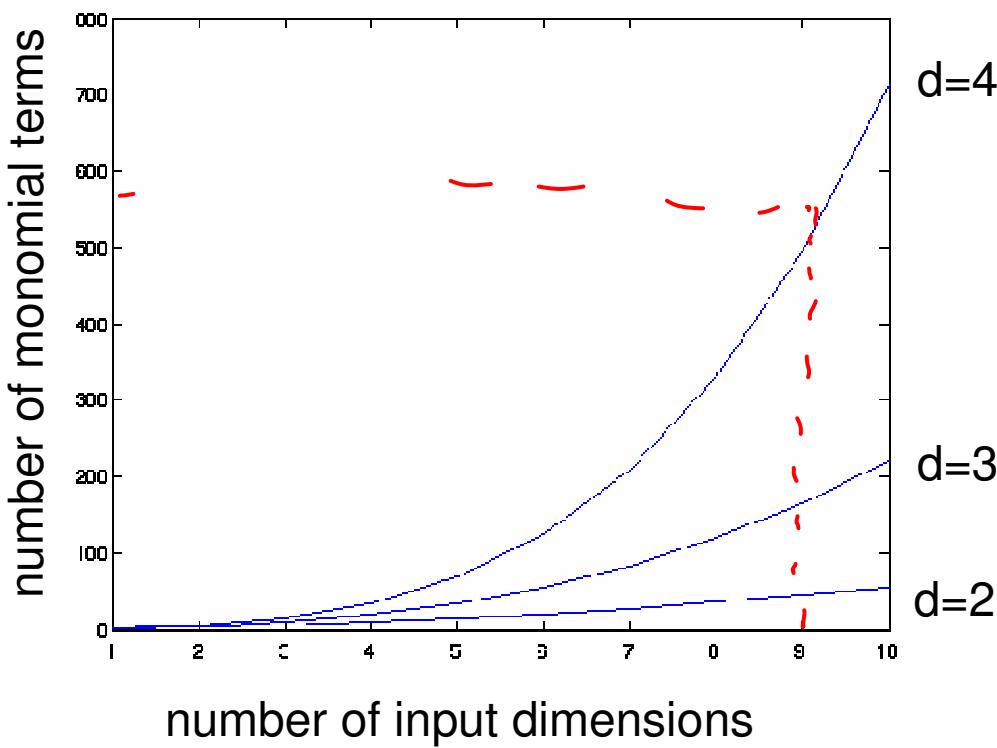
Feature space can get really large really quickly!

Higher order polynomials



degree of poly

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!}$$



d=4 input 4

m – input features
d – degree of polynomial

grows fast!
d = 6, m = 100
about 1.6 billion terms

Dual formulation only depends on dot-products, not on w!

$$\text{minimize}_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

no w !

use features $\phi(x)$

all I need is $\phi(x_j) \cdot \phi(x_i)$

$$K(x_j, x_i) = \frac{\phi(x_j) \cdot \phi(x_i)}{\phi(x_j) \cdot \phi(x_i)}$$

$$\text{minimize}_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

Finally: the “kernel trick”!

*i, j all pairs
of data points
including datapoint
with itself.*

$$\text{minimize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

- Never represent features explicitly
 - Compute dot products in closed form
- Constant-time high-dimensional dot-products for many classes of features
- Very interesting theory – Reproducing Kernel Hilbert Spaces
 - Not covered in detail in 10701/15781, more in 10702

$$\mathbf{w} = \sum_i \alpha_i y_i \Phi(\mathbf{x}_i)$$

$$b = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

for any k where $C > \alpha_k > 0$

Common kernels

- Polynomials of degree d

$$K(\mathbf{u}, \mathbf{v}) = \underbrace{(\mathbf{u} \cdot \mathbf{v})^d}_{\text{exactly}}$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = \underbrace{(\mathbf{u} \cdot \mathbf{v} + 1)^d}_{\text{including}}$$

- Gaussian kernels

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|}{2\sigma^2}\right)$$

- Sigmoid

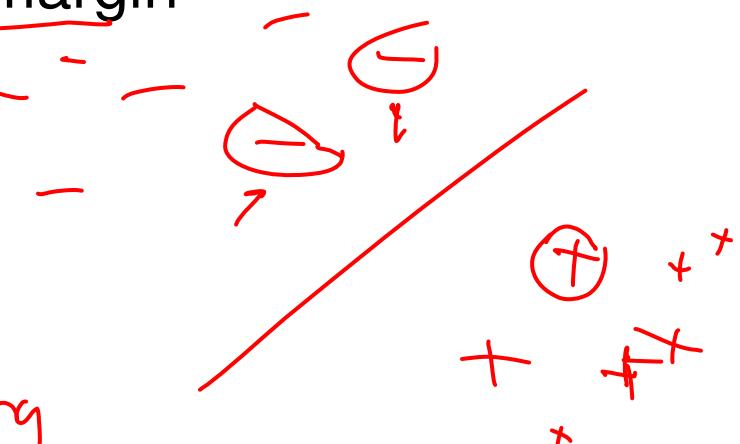
$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

correspond infinite
dimensional
feature space
 $\dim[\phi(x)]$ infinite

Overfitting?

- Huge feature space with kernels, what about overfitting???
- Maximizing margin leads to sparse set of support vectors
- Some interesting theory says that SVMs search for simple hypothesis with large margin
- Often robust to overfitting

*Sparse solutions \rightarrow
a few support vectors
 \rightarrow less overfitting*



What about at classification time

- For a new input \mathbf{x} , if we need to represent $\Phi(\mathbf{x})$, we are in trouble! *if have to write w, b , too large*
- Recall classifier: sign($\mathbf{w} \cdot \Phi(\mathbf{x}) + b$)
- Using kernels we are cool!

$$K(\mathbf{u}, \mathbf{v}) = \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v})$$

$$\underline{\mathbf{w} \cdot \Phi(\mathbf{x})} = \sum_{i \in \text{training data}} \alpha_i y_i \underbrace{\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)}_{\text{easy to compute}}$$

$$\text{new input } \mathbf{x} \rightsquigarrow \Phi(\mathbf{x}), \text{ sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

$$\underline{\mathbf{w}} = \sum_i \alpha_i y_i \Phi(\mathbf{x}_i)$$

$$\underline{b} = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

for any k where $C > \alpha_k > 0$

SVMs with kernels

- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors α_i
- At classification time, compute:

The diagram illustrates the classification process. A black box contains the equations for the learned model. A green arrow labeled "Classify as" points from the box to the classification rule. Red annotations with arrows point to the "new point" \mathbf{x} and the "old data" \mathbf{x}_i in the equations. A red curved arrow points from the "old data" \mathbf{x}_i in the first equation to the "new point" \mathbf{x} in the second equation, indicating that the learned model $\mathbf{w} \cdot \Phi(\mathbf{x})$ is used to predict the value of b for the new point \mathbf{x} .

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i)$$
$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

for any k where $C > \alpha_k > 0$

Classify as $\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$

Remember kernel regression

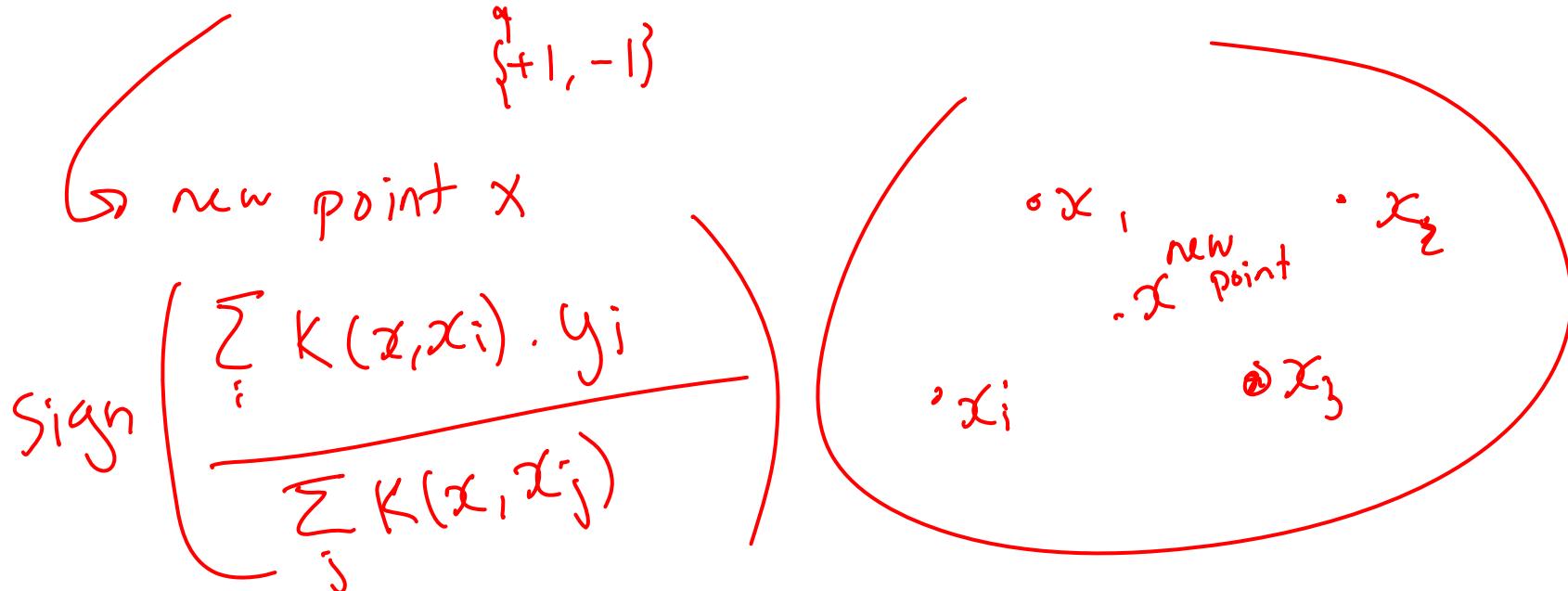
Remember kernel regression???

1. $w_i = \exp(-D(x_i, \text{query})^2 / K_w^2)$
2. How to fit with the local points?

Gaussian kernel $K(x, x_i)$

Predict the weighted average of the outputs:

$$\text{predict} = \sum w_i y_i / \sum w_i$$



SVMs v. Kernel Regression

SVMs

$$\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

or

$$\text{sign} \left(\sum_i \underbrace{\alpha_i y_i K(\mathbf{x}, \mathbf{x}_i)}_{\text{weights}} + b \right)$$

↑ offset

Kernel Regression

$$\Pi = \sum_j K(\mathbf{x}, \mathbf{x}_j)$$

$$\text{sign} \left(\frac{\sum_i y_i K(\mathbf{x}, \mathbf{x}_i)}{\sum_j K(\mathbf{x}, \mathbf{x}_j)} \right)$$

$$\hookrightarrow \text{sign} \left(\sum_i \frac{1}{\Pi} y_i K(\mathbf{x}, \mathbf{x}_i) \right)$$

SVMs v. Kernel Regression

SVMs

$$\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

or

Differences:

■ SVMs:

- Learn weights α_i (and bandwidth)
- Often sparse solution

of Kernel

■ KR:

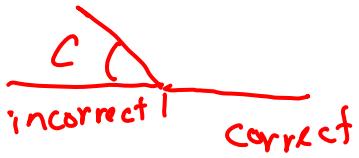
- Fixed “weights”, learn bandwidth
- Solution may not be sparse
- Much simpler to implement

Kernel Regression

$$\text{sign}\left(\frac{\sum_i y_i K(\mathbf{x}, \mathbf{x}_i)}{\sum_i K(\mathbf{x}, \mathbf{x}_i)}\right)$$

sign

What's the difference between SVMs and Logistic Regression?

	SVMs	Logistic Regression
Loss function	Hinge loss  A graph of the hinge loss function. The x-axis is labeled 'incorrect' and the y-axis is labeled 'correct'. The curve is zero for 'correct' points and increases linearly for 'incorrect' points. $\max(0, 1 - \mathbf{y} \cdot \mathbf{w})$	Log-loss  A graph of the log-loss function, which is convex and decreases as the predicted value approaches the true value. $\mathbf{y} \ln(\mathbf{p}) + (1 - \mathbf{y}) \ln(1 - \mathbf{p})$
High dimensional features with kernels	Yes! <i>kernel trick!!</i>	No <i>yes!</i>

Kernels in logistic regression

$$P(Y = 1 | x, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \phi(x) + b)}}$$

$\phi(x)$ high dim.
features

- Define weights in terms of support vectors:

"Representer Theorem":

$$\underline{\mathbf{w}} = \sum_i \underline{\alpha_i} \underline{\phi(\mathbf{x}_i)}$$

$$\begin{aligned} P(Y = 1 | x, \mathbf{w}) &= \frac{1}{1 + e^{-(\sum_i \underline{\alpha_i} \underline{\phi(\mathbf{x}_i)} \cdot \underline{\phi(x)} + b)}} \\ &= \frac{1}{1 + e^{-(\sum_i \underline{\alpha_i} K(\mathbf{x}, \mathbf{x}_i) + b)}} \end{aligned}$$

- Derive simple gradient descent rule on α_i 's using
learn α_i 's using
grad. descent

What's the difference between SVMs and Logistic Regression? (Revisited)

	SVMs	Logistic Regression
Loss function	Hinge loss 	Log-loss 
High dimensional features with kernels	Yes!	Yes!
Solution sparse	Often yes! <i>because hinge loss</i>	Almost always no! <i>because of log-loss</i>
Semantics of output	“margin”	“Real” probabilities

What you need to know

- Dual SVM formulation
 - How it's derived
- The kernel trick
- Derive polynomial kernel
- Common kernels
- Kernelized logistic regression
- Differences between SVMs and logistic regression

playing a ppflet...

option 1: learn $wx+b$

option 2: $\begin{pmatrix} x \\ 1 \end{pmatrix} = \bar{x}$

learn $\begin{pmatrix} w \\ w_0 \end{pmatrix} = \bar{w}$

$\rightarrow \bar{w} \cdot \bar{x}$ write SVM with nob

option 1
 $\min w \cdot w$

option
 $\min \bar{w} \cdot \bar{w}$

$\equiv \min w \cdot w + b^2$

option 1 is sometimes better

Acknowledgment

- SVM applet:
 - <http://www.site.uottawa.ca/~gcaron/applets.htm>

More details:

General: <http://www.learning-with-kernels.org/>

Example of more complex bounds:

http://www.research.ibm.com/people/t/tzhang/papers/jmlr02_cover.ps.gz

PAC-learning, VC Dimension and Margin-based Bounds

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

March 1st, 2005

©2006 Carlos Guestrin

What now...

- We have explored **many** ways of learning from data
- But...
 - How good is our classifier, really?
 - How much data do I need to make it “good enough”?

Learning Theory

A simple setting...

- Classification
 - m data points
 - Finite number of possible hypothesis (e.g., dec. trees of depth d)
- A learner finds a hypothesis h that is consistent with training data
 - Gets zero error in training – $\text{error}_{\text{train}}(h) = 0$
- What is the probability that h has more than ϵ true error?
 - $\text{error}_{\text{true}}(h) \geq \epsilon$

$$P(\text{error}_{\text{train}}(h) \geq \epsilon) \leq \gamma$$

e.g.) $\gamma = 0.01$

How likely is a bad hypothesis to get m data points right?

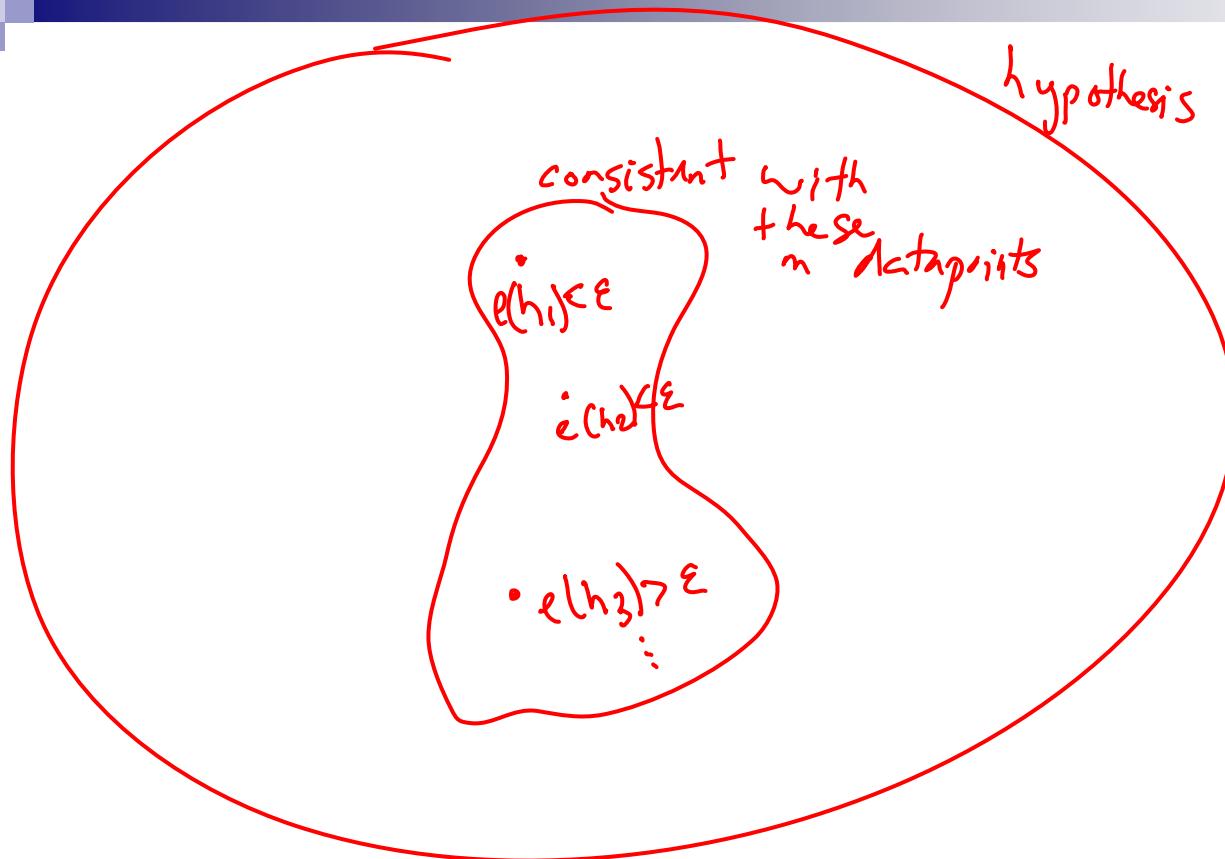
- Hypothesis h that is consistent with training data \rightarrow got m i.i.d. points all right
- Prob. h with $\text{error}_{\text{true}}(h) \geq \varepsilon$ gets one data point right
 $P(\text{error}_{\text{true}}(h) \geq \varepsilon, \text{ and get one data point right}) \leq 1 - \varepsilon$

- Prob. h with $\text{error}_{\text{true}}(h) \geq \varepsilon$ gets m data points right

$$P(\text{bad } h \text{ get lucky}) \leq (1 - \varepsilon)^m$$

more data, ^{exponentially} less likely a bad hyp. gets all right.

But there are many possible hypothesis that are consistent with training data



learner picked one arbitrarily
the worst one!!

How likely is learner to pick a bad hypothesis

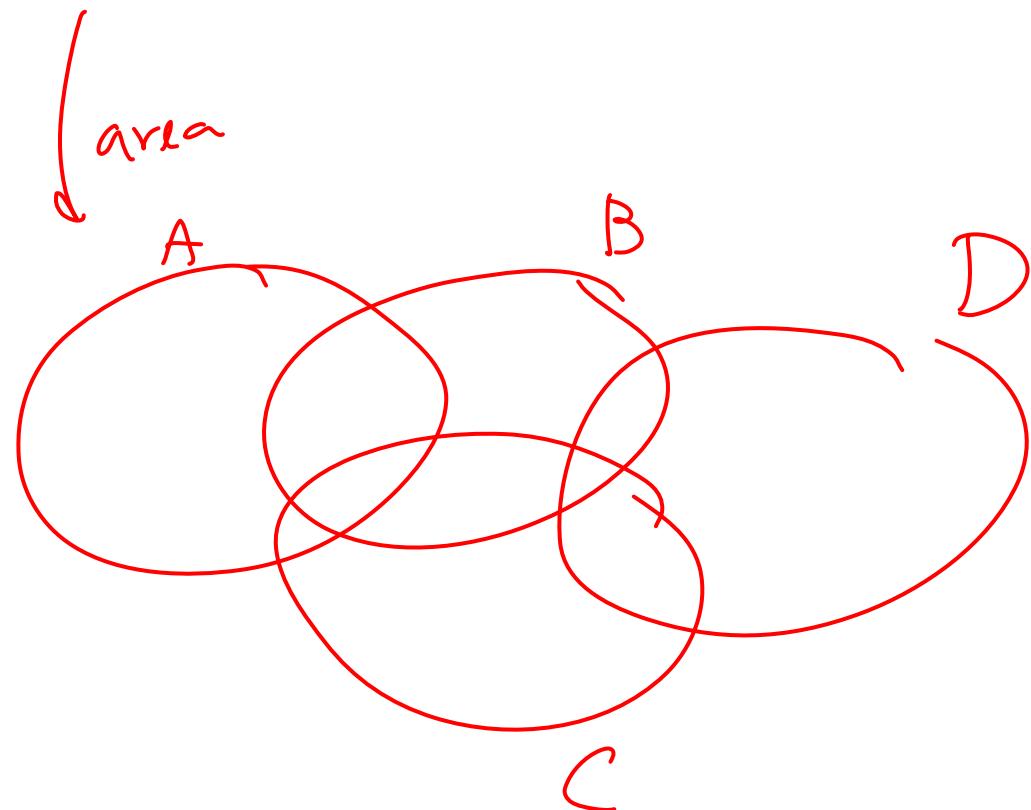
- Prob. h with $\text{error}_{\text{true}}(h) \geq \varepsilon$ gets m data points right
- There are k hypothesis consistent with data
 - How likely is learner to pick a bad one?

$P\left(\text{at least one of the } k \text{ was bad} \text{ and it got lucky}\right) ?$

$$P(\text{one got lucky}) \leq (1-\varepsilon)^m$$

Union bound

- bad hyp. got all m right
- $P(A \text{ or } B \text{ or } C \text{ or } D \text{ or } \dots) \leq P(A) + P(B) + P(C) + \dots$



How likely is learner to pick a bad hypothesis

- Prob. h with $\text{error}_{\text{true}}(h) \geq \varepsilon$ gets m data points right
- There are k hypothesis consistent with data

- How likely is learner to pick a bad one?

$$\begin{aligned} & P(h_1 \text{ bad} \wedge \text{got lucky} \text{ or } h_2 \text{ bad} \wedge \text{got lucky} \text{ or } h_3 \dots) \\ & \leq P(h_1 \text{ bad} \wedge \text{lucky}) + P(h_2 \text{ bad} \wedge \text{lucky}) + P(h_3 \dots) + \dots \\ & \leq (1-\varepsilon)^m \end{aligned}$$

$$\begin{aligned} & \leq K (1-\varepsilon)^m \\ & \leq |H| e^{-\varepsilon m} \end{aligned}$$

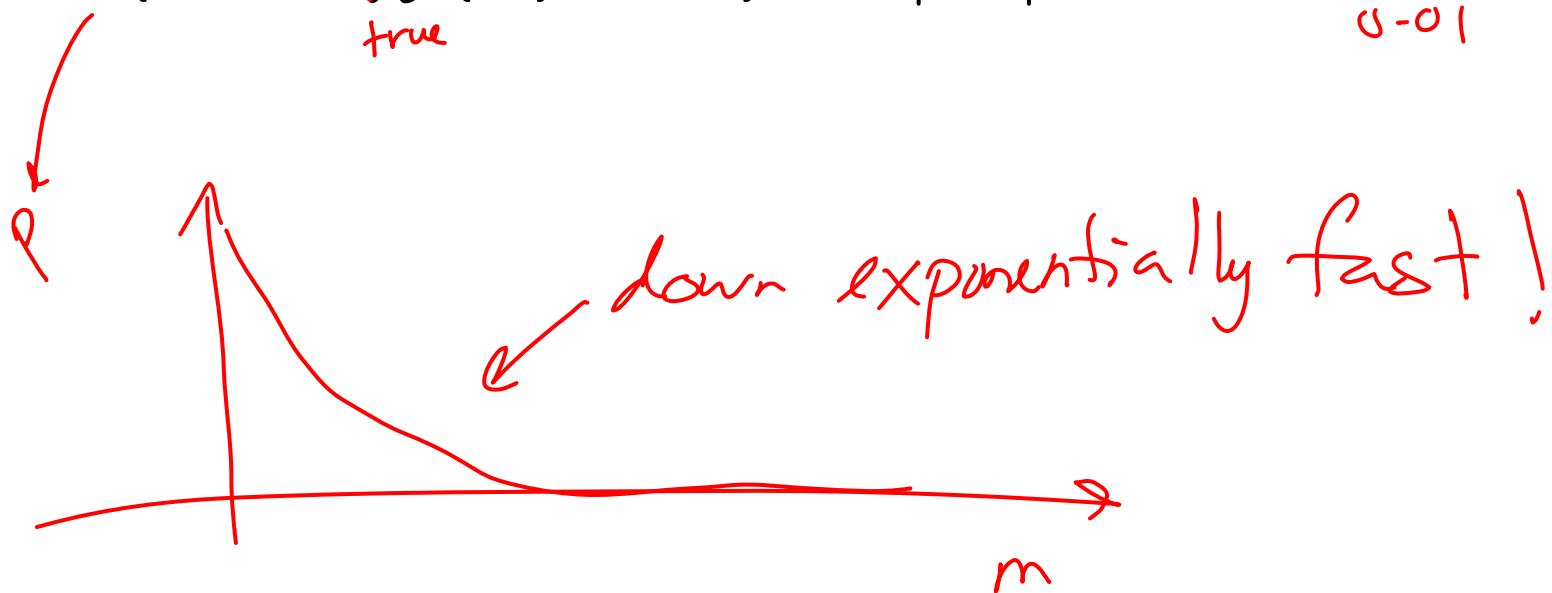
how big is K

$$K \leq |H| \quad (\text{loose bound!})$$
$$1-\varepsilon \leq e^{-\varepsilon} \quad \begin{cases} \text{make eq.} \\ \text{simpler} \end{cases}$$

Review: Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem:** Hypothesis space H finite, dataset D with m i.i.d. samples, $0 < \varepsilon < 1$: for any learned hypothesis h that is consistent on the training data:

$$P(\text{error}_{\mathcal{V}}(h) > \varepsilon) \leq |H| e^{-m\varepsilon} \leq \delta$$



Using a PAC bound

■ Typically, 2 use cases:

- 1: Pick ϵ and δ , give you m
- 2: Pick m and δ , give you ϵ

$$\textcircled{1} \quad P \leq |H| e^{-m\epsilon} \leq \delta$$

$$\ln |H| - m\epsilon \leq \ln \delta$$

$\Rightarrow m \geq \frac{1}{\epsilon} \left(\ln |H| + \ln \frac{1}{\delta} \right)$

smaller ϵ
more data

not as much data as you might think

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H| e^{-m\epsilon}$$

$$\text{Case 2 } \delta \leq |H| e^{-m\epsilon}$$

$$\ln \delta \leq \ln |H| - m\epsilon$$

$$\epsilon \leq \frac{1}{m} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

learn h
true error \leq
before you run the algorithm

Review: Generalization error in finite hypothesis spaces [Haussler '88]

- **Theorem:** Hypothesis space H finite, dataset D with m i.i.d. samples, $0 < \varepsilon < 1$: for any learned hypothesis h that is consistent on the training data:

$$P(\text{error}_{\mathcal{X}}(h) > \varepsilon) \leq |H|e^{-m\varepsilon}$$

if I can always learn a
consistent classifier then

Even if h makes zero errors in training data, may make errors in test

Limitations of Haussler '88 bound

① Consistent classifier

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon}$$

there may ^{not} be such h in class!

② Size of hypothesis space

bound depends on $|H|$

really really large?

infinite?
w continuous

What if our classifier does not have zero error on the training data?

- A learner with zero training errors may make mistakes in test set
- What about a learner with $\overline{\text{error}}_{\text{train}}(h)$ in training set?

h is not
perfect
in training

error
true?

Simpler question: What's the expected error of a hypothesis?

- The error of a hypothesis is like estimating the parameter of a coin! | flip coin m times

true error $\leftarrow \theta$
I don't know it



- Chernoff bound: for m i.i.d. coin flips, x_1, \dots, x_m , where $x_i \in \{0, 1\}$. For $0 < \varepsilon < 1$:

$$P\left(\theta - \frac{1}{m} \sum_i x_i > \epsilon\right) \leq e^{-2m\epsilon^2}$$

frutti

↑
train error
Simple
average

Using Chernoff bound to estimate error of a single hypothesis

$$P \left(\theta - \frac{1}{m} \sum_i x_i > \epsilon \right) \leq e^{-2m\epsilon^2}$$

$\text{error}_{\text{true}}(h) \leftarrow \theta$ # times expect to make a mistake

$x_i \leftarrow$ did I get i th data point wrong
 $x_i = \mathbb{I}(h(i) \neq t(i))$

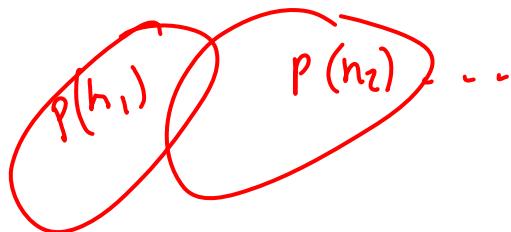
$$\text{error}_{\text{train}} = \frac{1}{n} \sum_i \mathbb{I}(h(i) \neq t(i)) = \frac{1}{n} \sum_i x_i$$

But we are comparing many hypothesis: **Union bound**

For each hypothesis h_i :

$$P(\text{error}_{\text{true}}(h_i) - \text{error}_{\text{train}}(h_i) > \epsilon) \leq e^{-2m\epsilon^2}$$

What if I am comparing two hypothesis, h_1 and h_2 ?


$$P(h_1) \quad P(h_2) \quad \dots$$

learner is going to compare all of h_i 's

$$P(\exists i, \text{error}_{\text{true}}(h_i) - \text{error}_{\text{train}}(h_i) > \epsilon) \leq |H| e^{-2m\epsilon^2}$$

Generalization bound for $|H|$ hypothesis

- **Theorem:** Hypothesis space H finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any learned hypothesis h :

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H| e^{-2m\epsilon^2}$$

$$2m\epsilon^2 = 20 \quad \text{not as good !!}$$

side note: Haussler's bound for consistent h :

$$P \leq |H| e^{-m\epsilon}$$

$$\epsilon = 0.1$$

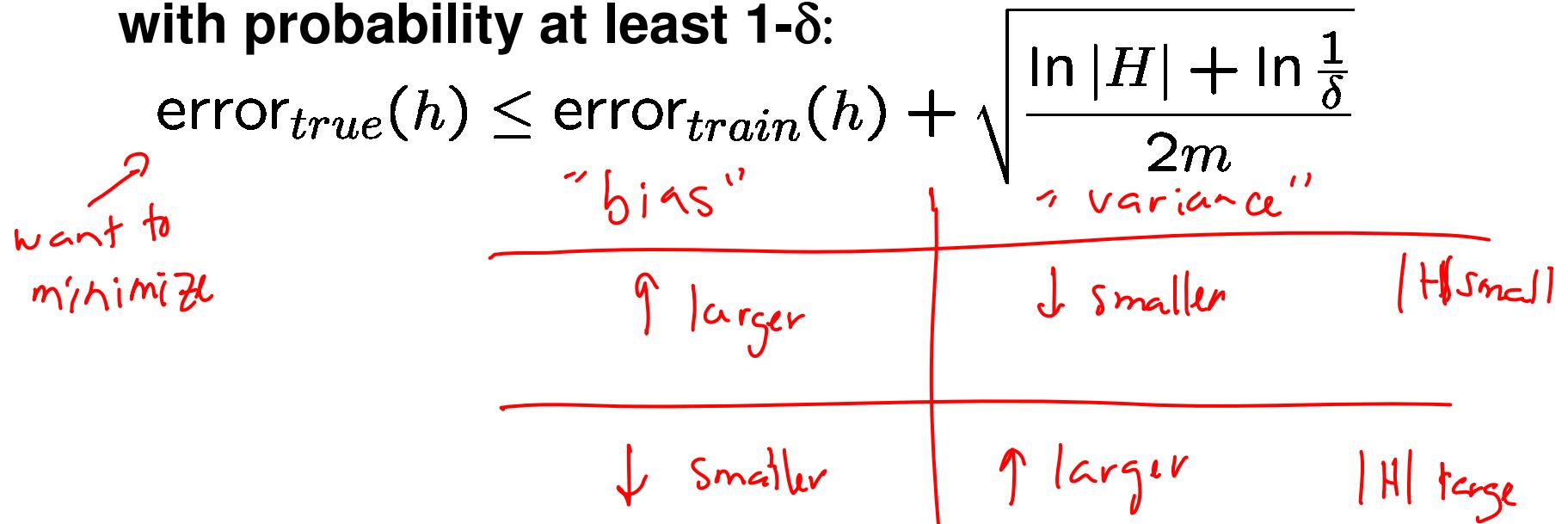
$$\Rightarrow m, \epsilon = 166$$

$$m = 1000$$

PAC bound and Bias-Variance tradeoff

$$\underline{P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H|e^{-2m\epsilon^2}}$$

or, after moving some terms around,
with probability at least $1-\delta$:



- **Important: PAC bound holds for all h , but doesn't guarantee that algorithm finds best h !!!**

What about the size of the hypothesis space?

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

↑
need this amount

- How large is the hypothesis space? $|H|$

$\ln |H|$?

Boolean formulas with n binary features

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

what's $\ln |H|$?

tabular representation

	x_1	x_2	x_3	x_4	y
T	T	T	T	0/1	
T	T	T	F	0/1	
T	T	F	T	0/1	
T	T	F	F	0/1	

2^n rows
each

2 possibl.
y

conjunctions:

$$h_1 = x_1 \wedge x_2 \wedge x_7$$

$$h_2 = x_2 \wedge \neg x_5 \wedge \neg x_8 \dots$$

n attrib..

$$\langle \{\emptyset, 1, 7\}, \{\emptyset, 1, 7\}, \dots \rangle$$

$$|H| = 3^n \text{ (really large)}$$

$$|H| = 2^{2^n} \text{ (really really large)}$$

$$\ln |H| = n \ln 3$$

"small"

$$\ln |H| = 2^n \ln 2$$

too large

Number of decision trees of depth k

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

Recursive solution

Given n attributes

H_k = Number of decision trees of depth k

$H_0 = 2$

$$\begin{aligned} H_{k+1} &= (\text{\#choices of root attribute}) * \\ &\quad (\text{\# possible left subtrees}) * \\ &\quad (\text{\# possible right subtrees}) \\ &= n * H_k * H_k \end{aligned}$$

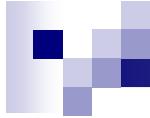
Write $L_k = \log_2 H_k$

$L_0 = 1$

$L_{k+1} = \log_2 n + 2L_k$

So $L_k = (2^k - 1)(1 + \log_2 n) + 1$

PAC bound for decision trees of depth k



$$m \geq \frac{\ln 2}{2\epsilon^2} \left((2^k - 1)(1 + \log_2 n) + 1 + \ln \frac{1}{\delta} \right)$$

- Bad!!!
 - Number of points is exponential in depth!
- But, for m data points, decision tree can't get too big...

Number of leaves never more than number data points

Number of decision trees with k leaves

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{1}{\delta} \right)$$

H_k = Number of decision trees with k leaves

$H_0 = 2$

$$H_{k+1} = n \sum_{i=1}^k H_i H_{k+1-i}$$

Loose bound:

$$H_k = n^{k-1} (k+1)^{2k-1}$$

Reminder:

$$|\text{DTs depth } k| = 2 * (2n)^{2^k - 1}$$

PAC bound for decision trees with k leaves – Bias-Variance revisited

$$H_k = n^{k-1} (k+1)^{2k-1}$$

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(k-1) \ln n + (2k-1) \ln(k+1) + \ln \frac{1}{\delta}}{2m}}$$

What did we learn from decision trees?

- Bias-Variance tradeoff formalized

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(k-1) \ln n + (2k-1) \ln(k+1) + \ln \frac{1}{\delta}}{2m}}$$

- Moral of the story:

Complexity of learning not measured in terms of size hypothesis space, but in maximum *number of points* that allows consistent classification

- Complexity m – no bias, lots of variance
- Lower than m – some bias, less variance

What about continuous hypothesis spaces?


$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

- Continuous hypothesis space:
 - $|H| = \infty$
 - Infinite variance???
- **As with decision trees, only care about the maximum number of points that can be classified exactly!**



How many points can a linear boundary classify exactly? (1-D)



How many points can a linear boundary classify exactly? (2-D)



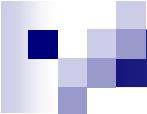
How many points can a linear boundary classify exactly? (d-D)

PAC bound using VC dimension

- Number of training points that can be classified exactly is VC dimension!!!
 - Measures relevant size of hypothesis space, as with decision trees with k leaves

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

Shattering a set of points



Definition: a **dichotomy** of a set S is a partition of S into two disjoint subsets.

Definition: a set of instances S is **shattered** by hypothesis space H if and only if for every dichotomy of S there exists some hypothesis in H consistent with this dichotomy.

VC dimension



Definition: The **Vapnik-Chervonenkis dimension**, $VC(H)$, of hypothesis space H defined over instance space X is the size of the largest finite subset of X shattered by H . If arbitrarily large finite sets of X can be shattered by H , then $VC(H) \equiv \infty$.

Examples of VC dimension


$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

- Linear classifiers:
 - $VC(H) = d+1$, for d features plus constant term b
- Neural networks
 - $VC(H) = \# \text{parameters}$
 - Local minima means NNs will probably not find best parameters
- 1-Nearest neighbor?

PAC bound for SVMs

- SVMs use a linear classifier
 - For d features, $\text{VC}(H) = d+1$:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(d+1) \left(\ln \frac{2m}{d+1} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

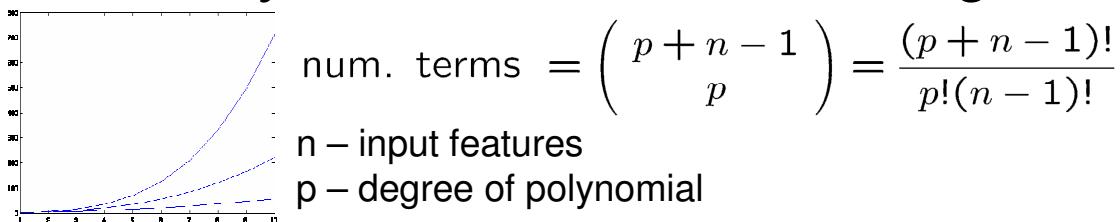
VC dimension and SVMs: Problems!!!

Doesn't take margin into account

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{(d+1) \left(\ln \frac{2m}{d+1} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

■ What about kernels?

- Polynomials: num. features grows really fast = Bad bound



- Gaussian kernels can classify any set of points exactly

Margin-based VC dimension

- H : Class of linear classifiers: $\mathbf{w} \cdot \Phi(\mathbf{x})$ ($b=0$)
 - Canonical form: $\min_j |\mathbf{w} \cdot \Phi(\mathbf{x}_j)| = 1$
- $VC(H) = R^2 \mathbf{w} \cdot \mathbf{w}$
 - Doesn't depend on number of features!!!
 - $R^2 = \max_j \Phi(\mathbf{x}_j) \cdot \Phi(\mathbf{x}_j)$ – magnitude of data
 - R^2 is bounded even for Gaussian kernels \rightarrow bounded VC dimension
- Large margin, low $\mathbf{w} \cdot \mathbf{w}$, low VC dimension – Very cool!

Applying margin VC to SVMs?


$$\text{error}_{true}(h) \leq \text{error}_{train}(h) + \sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

- $VC(H) = R^2 \mathbf{w} \cdot \mathbf{w}$
 - $R^2 = \max_j \Phi(\mathbf{x}_j) \cdot \Phi(\mathbf{x}_j)$ – magnitude of data, doesn't depend on choice of \mathbf{w}
- SVMs minimize $\mathbf{w} \cdot \mathbf{w}$
- SVMs minimize VC dimension to get best bound?
- Not quite right: ☹
 - Bound assumes VC dimension chosen before looking at data
 - Would require union bound over infinite number of possible VC dimensions...
 - But, it can be fixed!

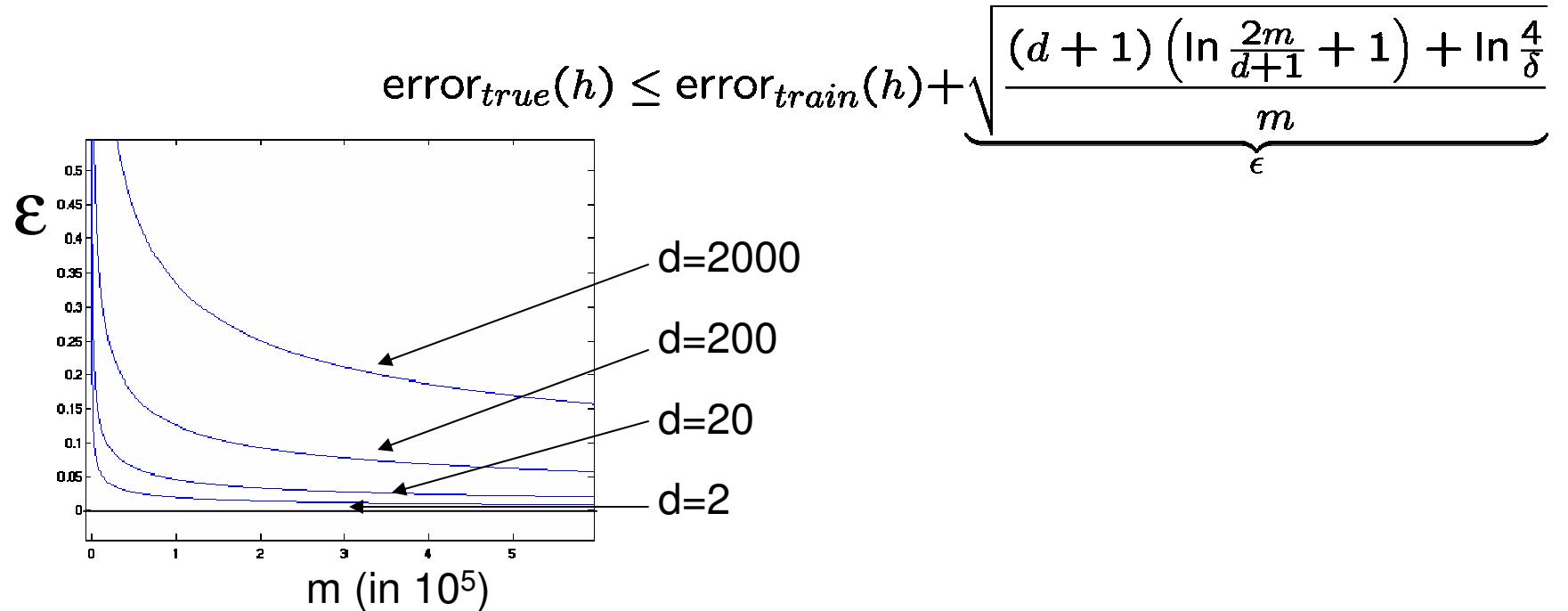
Structural risk minimization theorem


$$\text{error}_{true}(h) \leq \text{error}_{train}^\gamma(h) + C \sqrt{\frac{\frac{R^2}{\gamma^2} \ln m + \ln \frac{1}{\delta}}{m}}$$

$\text{error}_{train}^\gamma(h)$ = num. points with margin < γ

- For a family of hyperplanes with margin $\gamma > 0$
 - $\mathbf{w} \cdot \mathbf{w} \leq 1$
- SVMs maximize margin γ + hinge loss
 - Optimize tradeoff training error (bias) versus margin γ (variance)

Reality check – Bounds are loose



- Bound can be very loose, why should you care?
 - There are tighter, albeit more complicated, bounds
 - Bounds give us formal guarantees that empirical studies can't provide
 - Bounds give us intuition about complexity of problems and convergence rate of algorithms

What you need to know

- Finite hypothesis space
 - Derive results
 - Counting number of hypothesis
 - Mistakes on Training data
- Complexity of the classifier depends on number of points that can be classified exactly
 - Finite case – decision trees
 - Infinite case – VC dimension
- Bias-Variance tradeoff in learning theory
- Margin-based bound for SVM
- Remember: will your algorithm find best classifier?