

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 6th, 2006

Announcements 1

- Midterm on Wednesday
 - ☐ open book, texts, notes,...
 - ☐ no laptops
 - ☐ bring a calculator

Review session today at 5pm
NSH 3305

Announcements 2

Final project details are out!!!

- <http://www.cs.cmu.edu/~guestrin/Class/10701/projects.html>
- Great opportunity to apply ideas from class and learn more
- Example project:
 - Take a dataset
 - Define learning task
 - Apply learning algorithms
 - Design your own extension
 - Evaluate your ideas
- many of suggestions on the webpage, but you can also do your own

talk to us first.

Boring stuff:

- Individually or groups of two students
- It's worth 20% of your final grade
- You need to submit a one page proposal on Wed. 3/22 (just after the break)
- A 5-page initial write-up (milestone) is due on 4/12 (20% of project grade)
- An 8-page final write-up due 5/8 (60% of the grade)
- A poster session for all students will be held on Friday 5/5 2-5pm in NSH atrium (20% of the grade)
- You can use late days on write-ups, each student in team will be charged a late day per day.

page limits are strict...

MOST IMPORTANT:

Have some Fun!!

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

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?

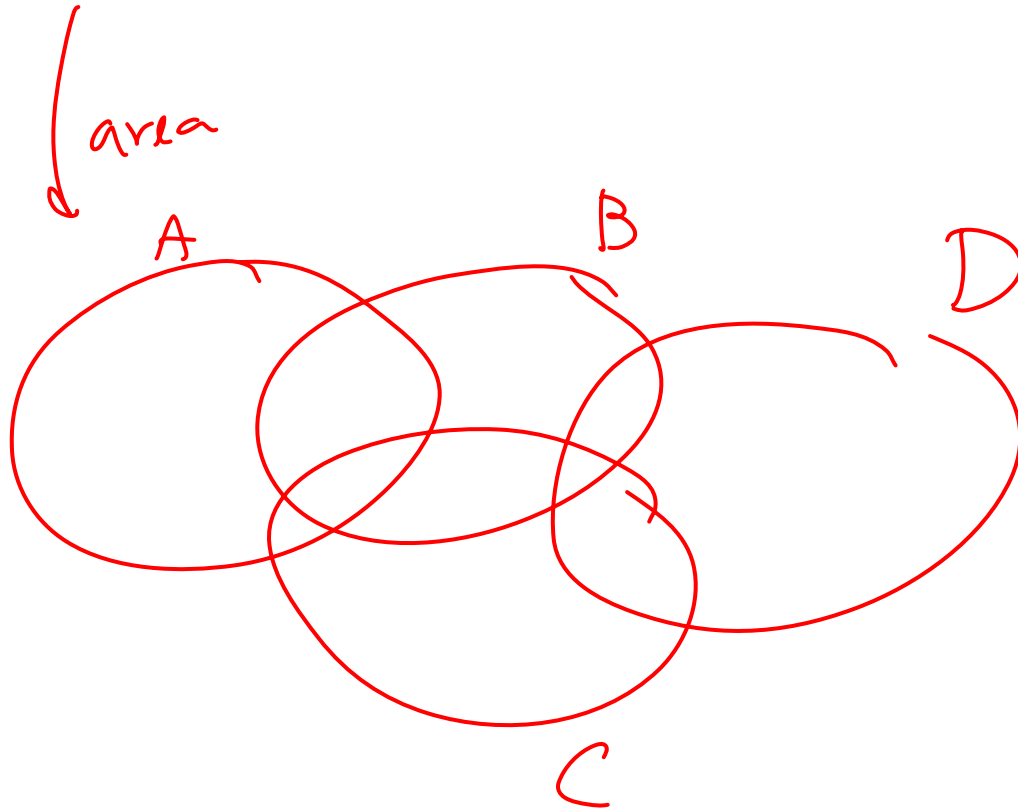
$P(\text{at least one of the } k \text{ was bad and it got lucky}) ?$

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

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?

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

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

how big is K

$$K \leq |H| \quad (\text{loose bound!})$$

$$1-\varepsilon \leq e^{-\varepsilon} \quad \left(\begin{array}{l} \text{make eq.} \\ \text{simpler} \end{array} \right)$$

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

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

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H|e^{-m\epsilon} \leq \sigma_{0.01}$$



Using a PAC bound

■ Typically, 2 use cases:

□ 1: Pick ϵ and δ , give you m

□ 2: Pick m and δ , give you ϵ

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

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

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

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

Smaller δ → need more data

Smaller ϵ more data

not as much data as you may think

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, $0 < \epsilon < 1$: for any learned hypothesis h that is consistent on the training data:

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

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 ^{true} h in class!

- ② ■ Size of hypothesis space
- bound depends on $|H|$
- really really large?
 - infinite?
w continuous

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

- The error of a hypothesis is like estimating the parameter of a coin!

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

flip coin m times
(train error) $\hat{\theta}$

θ v. $\hat{\theta}$

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

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

frutti

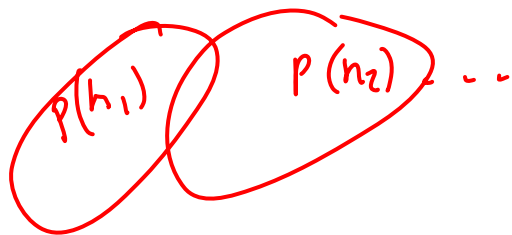
train error sample average

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 ?



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 :

$$P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq \frac{|H|e^{-2m\epsilon^2}}{1}$$

$2m\epsilon^2 = 20$ not as good !!
↓

side note: Haussler's Bound for consistent h :

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

$$\epsilon = 0.1$$

$$m = 1000$$

$$\Rightarrow m \cdot \epsilon = 100$$

PAC bound and Bias-Variance tradeoff

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

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

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

want to
minimize

"bias"		"variance"
↑ larger	↓ smaller	H small
↓ smaller	↑ larger	H large

■ 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? $|H|$

$\ln |H|$? \int

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	:

2^n rows
each
2 possible

conjunctions:

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

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

n attrib..

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

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

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

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

too large

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

"small"

Number of decision trees of depth k

binary features

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

what's $\ln |H|$

Recursive solution

Given n attributes

H_k = Number of decision trees of depth k

$$H_0 = 2$$

$$H_{k+1} = (\text{\#choices of root attribute}) * (\text{\# possible left subtrees}) * (\text{\# possible right subtrees})$$

$$= n * H_k * H_k$$

Write $L_k = \log_2 H_k$

$$L_0 = 1$$

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

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



upper bound ...
(may not use same feature twice)

PAC bound for decision trees of depth k

plug into bound

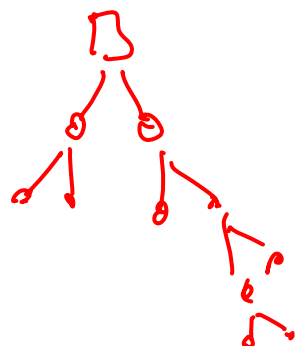
$$m \geq \frac{\ln 2}{2\epsilon^2} \left((2^k - 1)(1 + \log_2 n) + 1 + \ln \frac{1}{\delta} \right)$$

■ Bad!!!

□ Number of points is exponential in depth!

■ But, for m data points, decision tree can't get too big...

DT:



after m leaves
not worth

splitting any more!

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

$$H_0 = 2$$

bound n features for root k ways to split k between left & right

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

i leaves left \leftarrow rest on right.

10 leaves 11 12 \leftarrow $m-10$ leaves $m-11$ $m-12 \dots$

quick bound

Loose bound:

$$\underline{H_k} \leq n^{k-1} (k+1)^{2k-1}$$

Reminder:

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

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

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

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

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

$K = m$

0

> 1 (bad)

$K = \alpha m$
 $\alpha < 1$

↑ up

↓ down

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

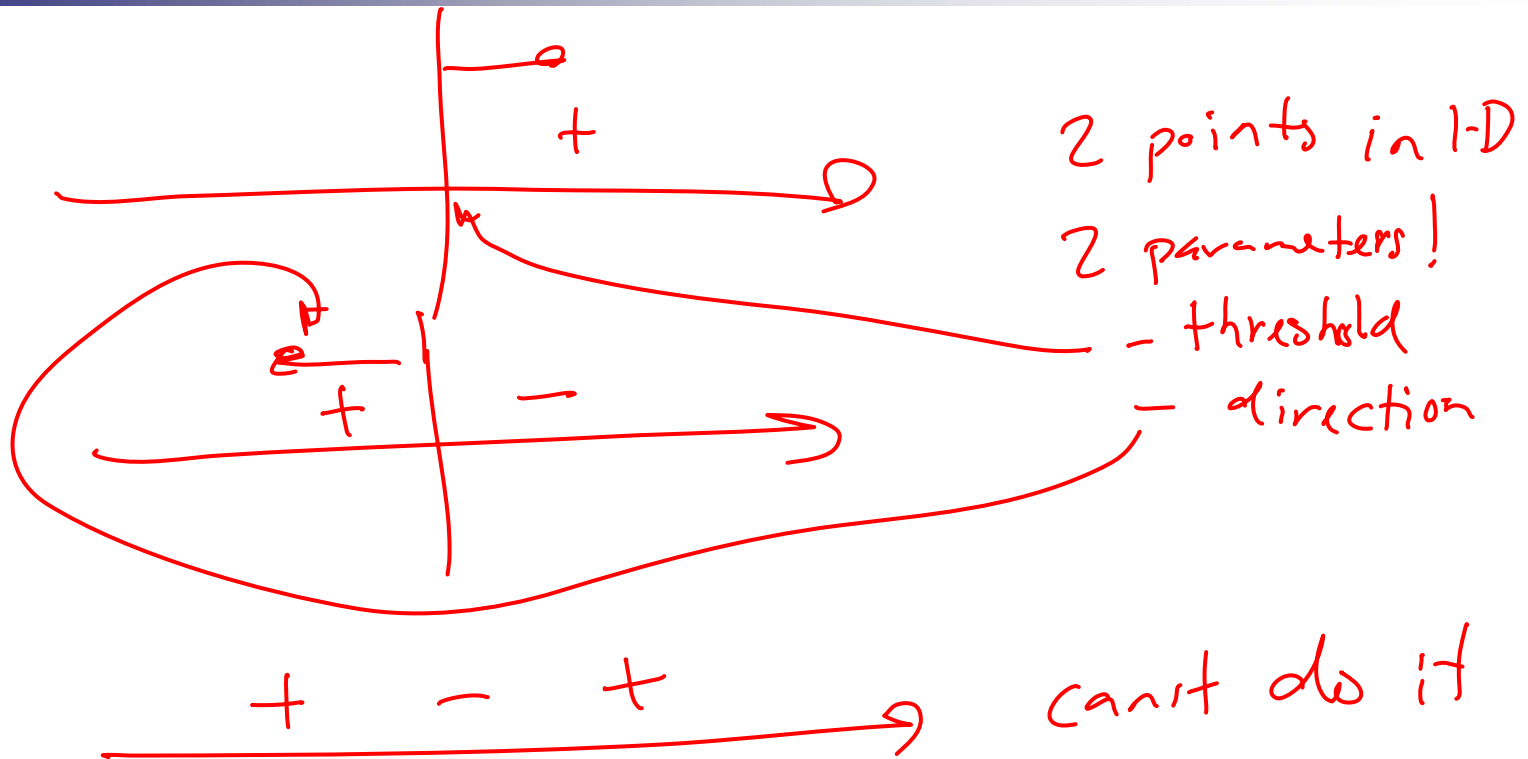
- Complexity m – no bias, lots of variance
- Lower than m – some bias, less variance

What about continuous hypothesis spaces?

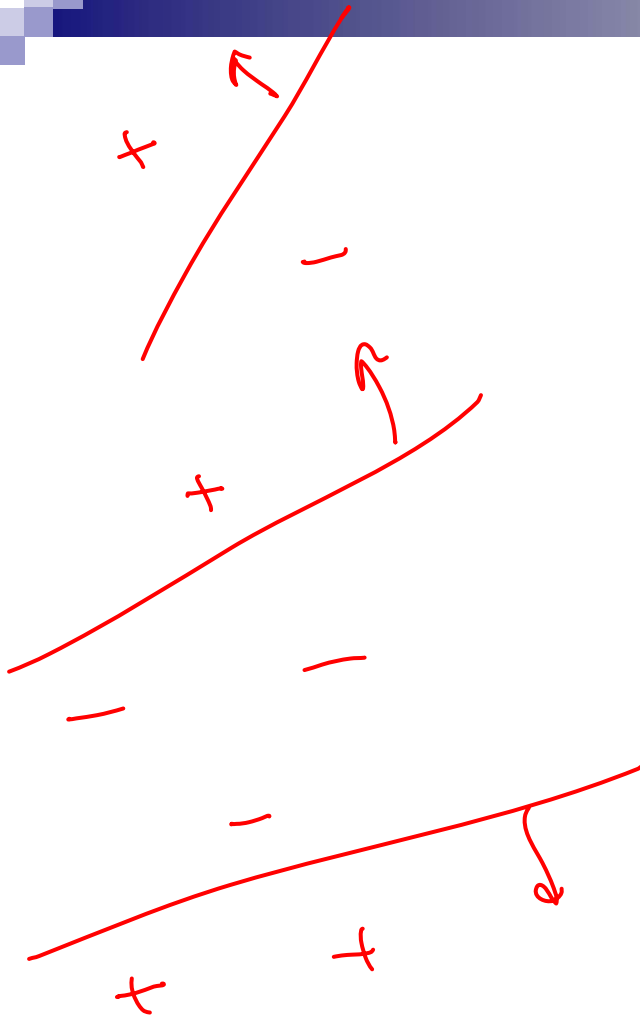
$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

- Continuous hypothesis space: *linear classifiers...*
 - $|H| = \infty$
 - 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)

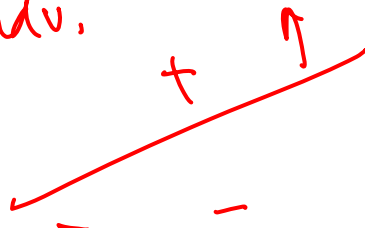


How many points can a linear boundary classify exactly? (2-D)



- you pick locations
- adversary picks labels

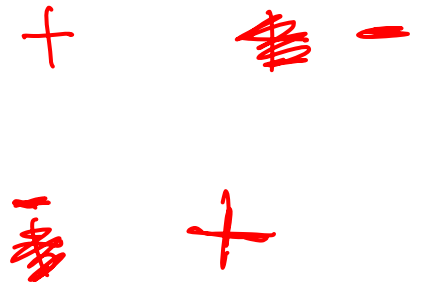
you
adv.



you pick

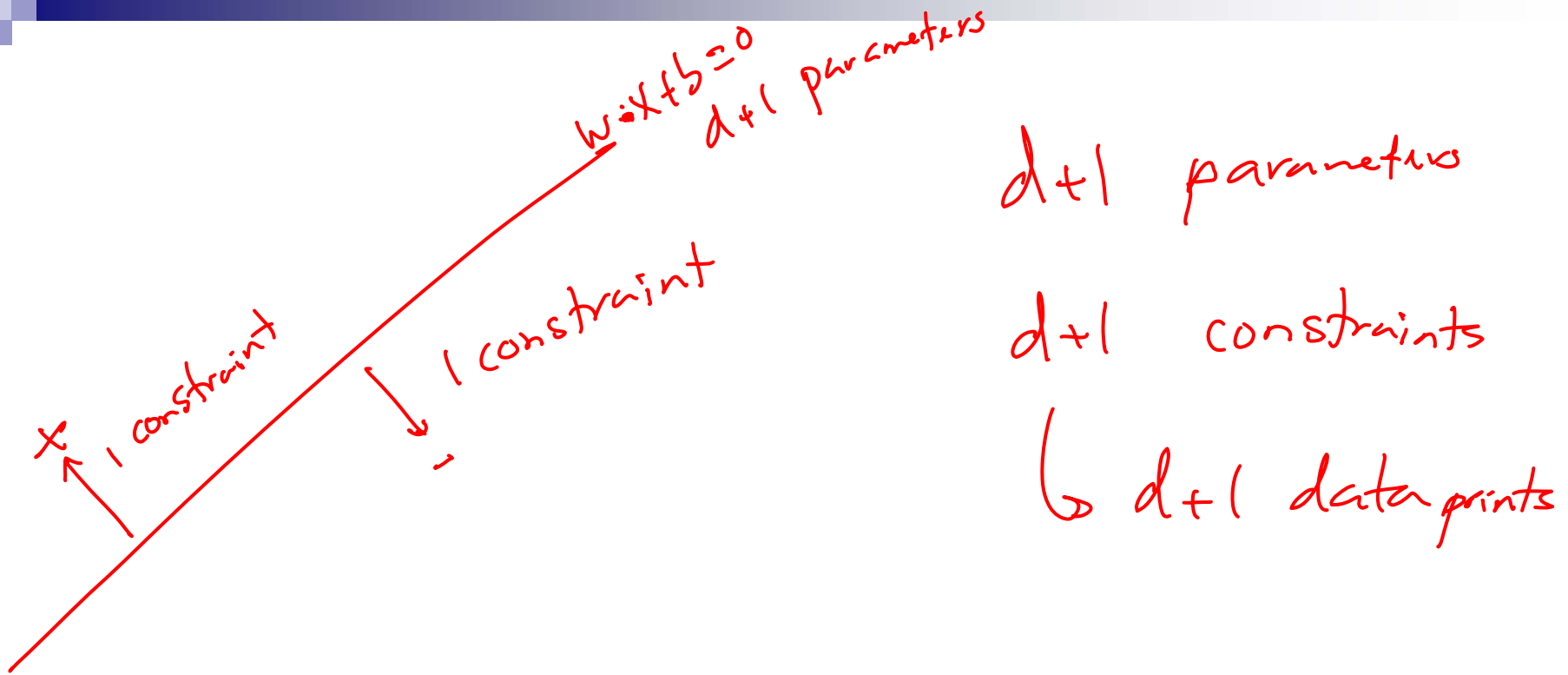
+

adversary:
XOR



Gap ...

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



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.

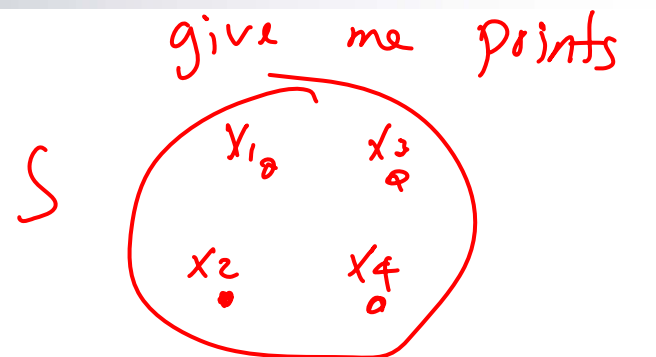
\forall dichotomies

$\exists h$ labels correctly:

$\{x_1, x_2\} \rightarrow + \quad \{x_3, x_4\} \rightarrow - \}$ choose h_7

$\{x_3\} \rightarrow + \quad \{x_1, x_2, x_4\} \rightarrow - \}$ choose h_{52}

\forall splits $\exists h \rightarrow$ shattered !!



dichotomy: set of labels
 $S = S_1 \cup S_2 \quad S_1 \cap S_2 = \emptyset$

$S_1 \rightarrow +$

$S_2 \rightarrow -$

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 $VC(H) \equiv \infty$.

hyper plane in 2d

• → can't shatter

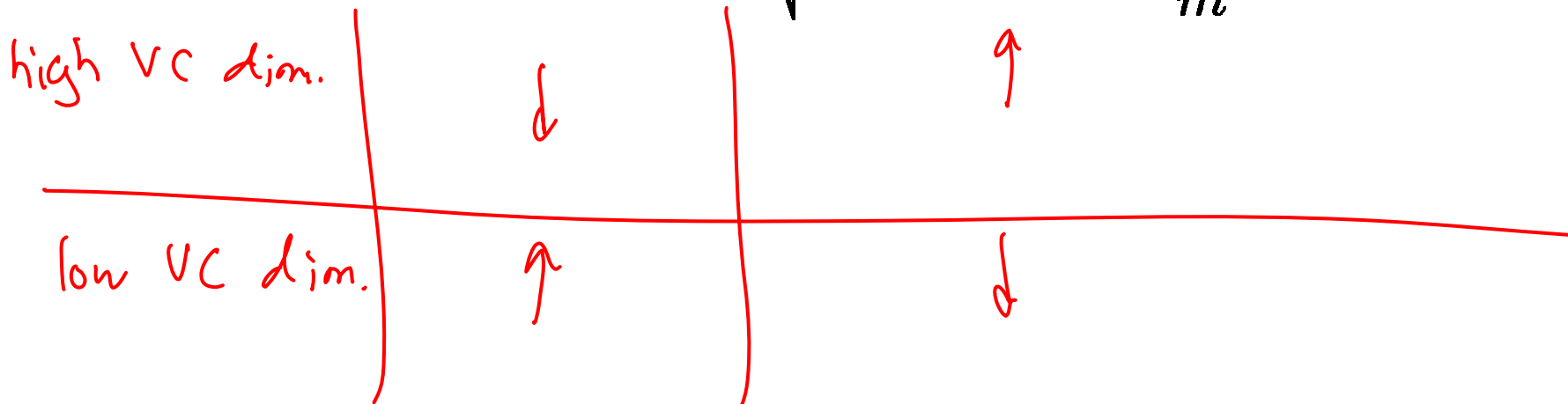
but I don't
care ...

→ I get to pick locations
I pick: •
•

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
 - Bound for infinite dimension hypothesis spaces:

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



Examples of VC dimension

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

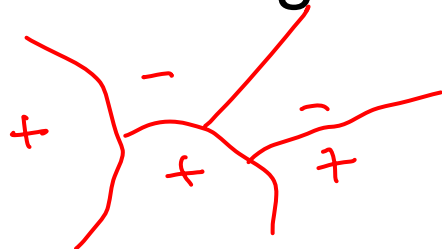
■ Linear classifiers:

- $VC(H) = d+1$, 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?



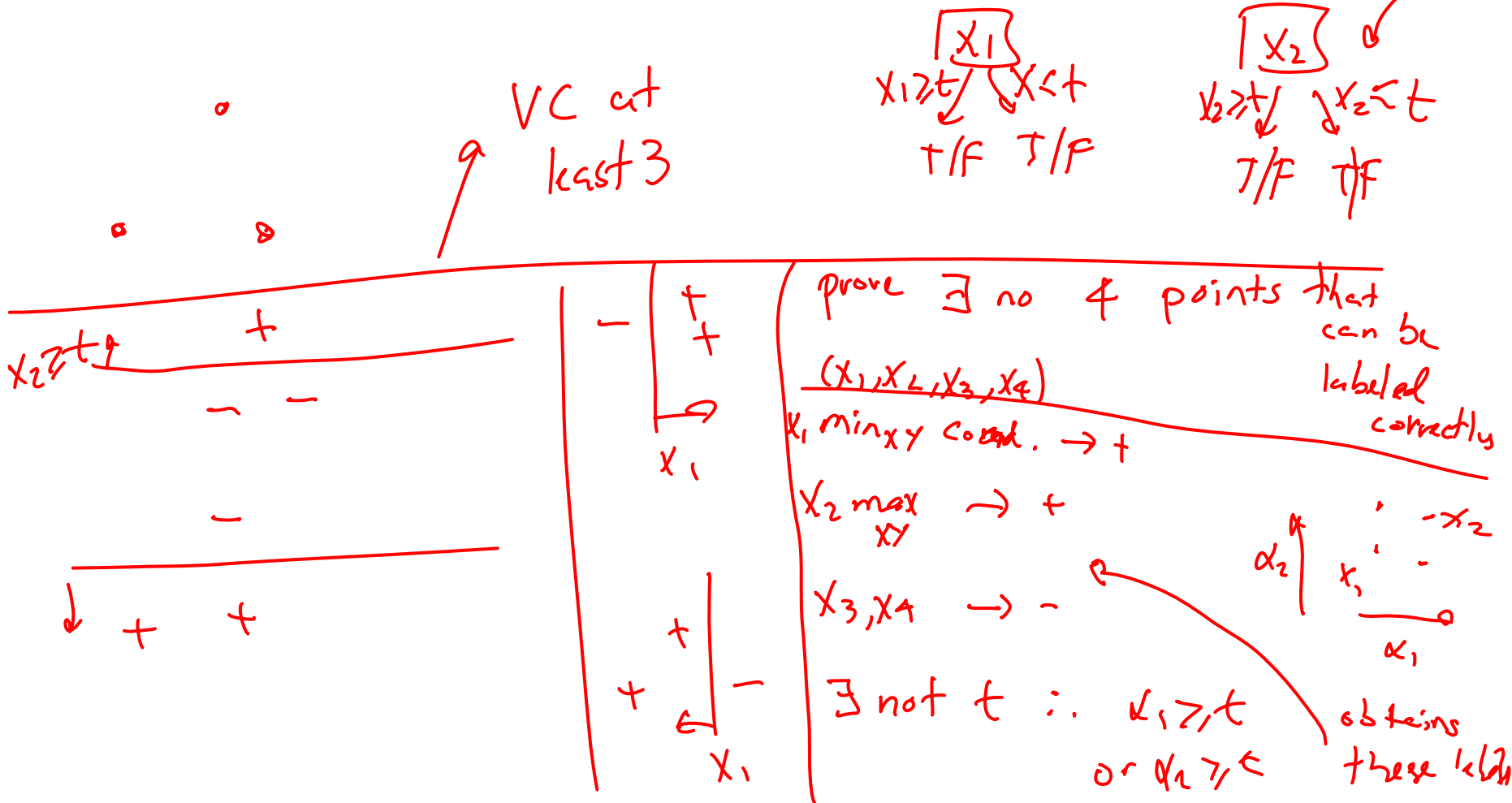
(In this def., an input point is its own NN.)

$$\downarrow$$

$$VC \equiv \infty$$

Another VC dim. example


- What's the VC dim. of decision stumps in 2d?



PAC bound for SVMs

- SVMs use a linear classifier

- For d features, $VC(H) = d+1$:


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

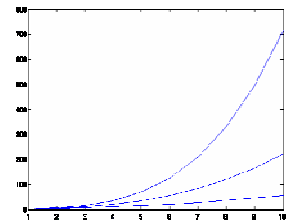
VC dimension and SVMs: Problems!!!

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

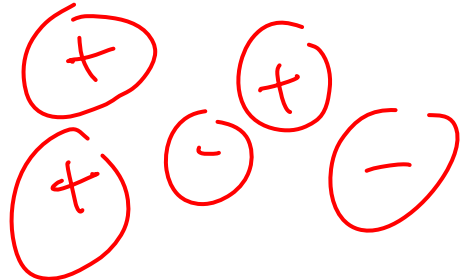


$$\text{num. terms} = \binom{p+n-1}{p} = \frac{(p+n-1)!}{p!(n-1)!}$$

n – input features

p – degree of polynomial

□ Gaussian kernels can classify any set of points exactly

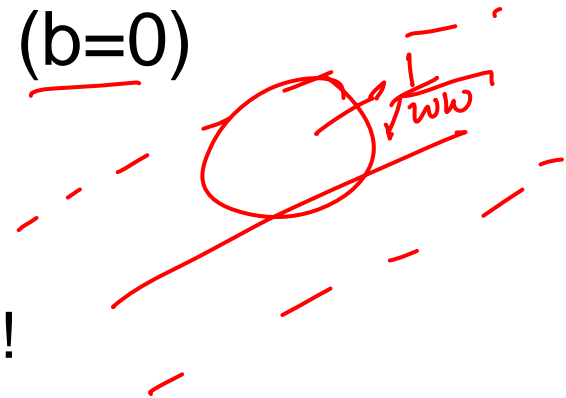


bound bad!!

"VC = ∞"

Margin-based VC dimension

- H: Class of linear classifiers: $\mathbf{w} \cdot \Phi(\mathbf{x})$ ($b=0$)
 - Canonical form: $\min_j |\mathbf{w} \cdot \Phi(\mathbf{x}_j)| = 1$
training example
- $VC(H) = R^2 \mathbf{w} \cdot \mathbf{w}$ $= \frac{R^2}{\text{margin}^2}$
 - 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 → bounded VC dimension
- Large margin, low $\mathbf{w} \cdot \mathbf{w}$, low VC dimension – Very cool!



Applying margin VC to SVMs?

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

- $VC(H) = R^2 \mathbf{w} \cdot \mathbf{w}$ *k*
 - $R^2 = \max_j \Phi(\mathbf{x}_j) \cdot \Phi(\mathbf{x}_j)$ – 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? *Smallest VC dim. hyp. consistent with data*
- **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

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

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

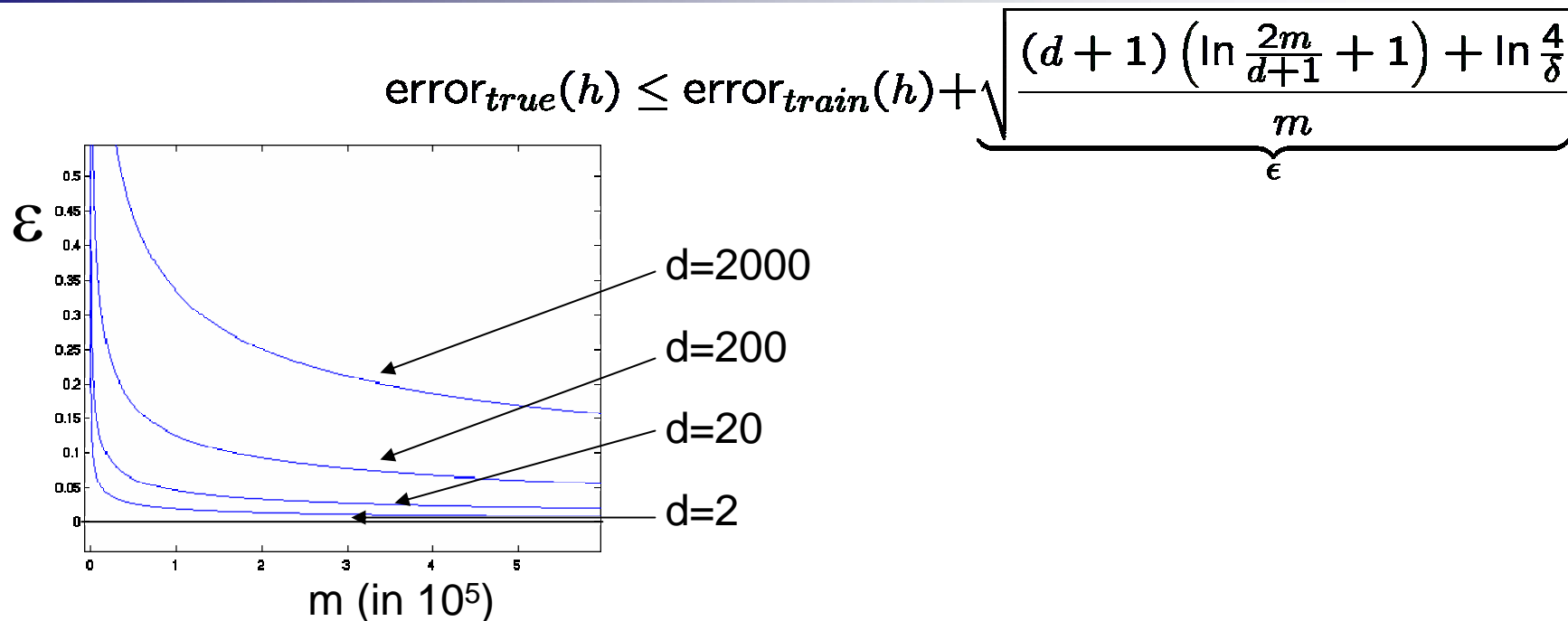
↪ fix margin γ

how much of data do I get right.

Simplest
Structure
→ minimum
error.

- 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 gives 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?



Big Picture

Machine Learning – 10701/15781

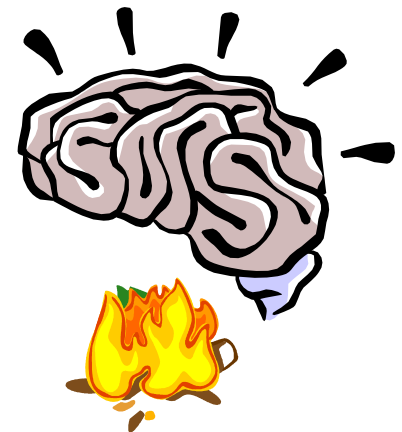
Carlos Guestrin

Carnegie Mellon University

March 6th, 2006

What you have learned thus far

- Learning is function approximation
- Point estimation
- Regression
- Naïve Bayes
- Logistic regression
- Bias-Variance tradeoff
- Neural nets
- Decision trees
- Cross validation
- Boosting
- Instance-based learning
- SVMs
- Kernel trick
- PAC learning
- VC dimension
- Margin bounds
- Mistake bounds



Review material in terms of...



- Types of learning problems
- Hypothesis spaces
- Loss functions
- Optimization algorithms

Text Classification



→ Company home page

VS

Personal home page

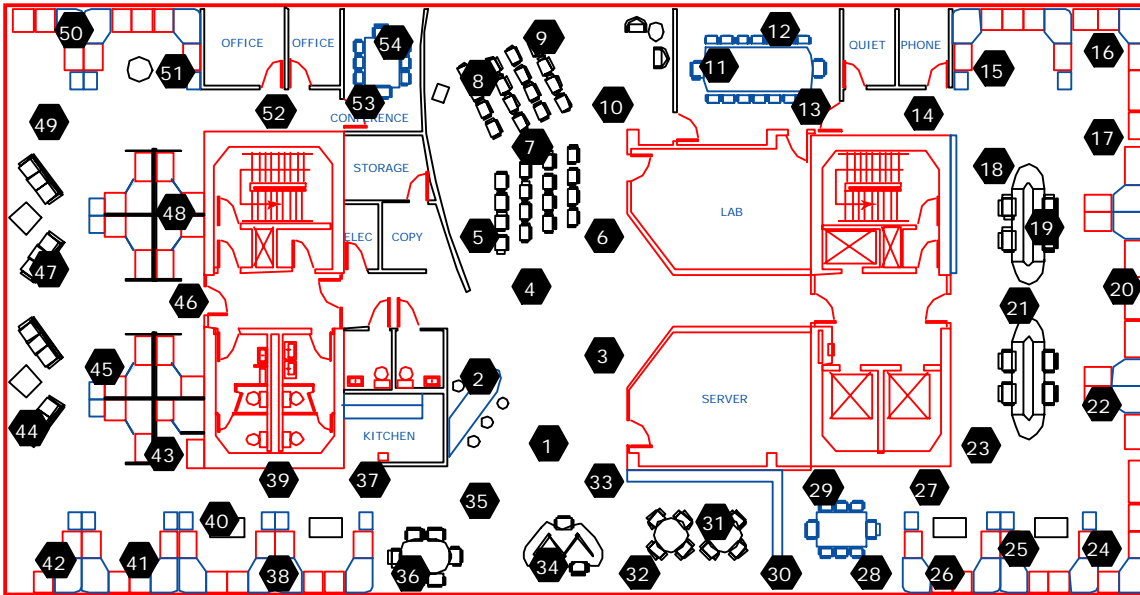
VS

Univeristy home page

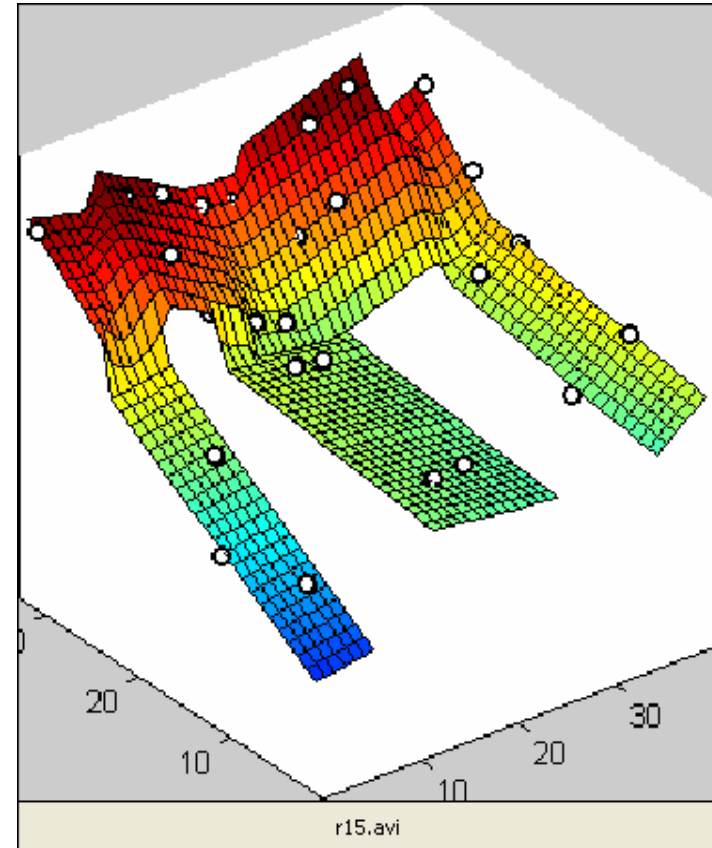
VS

...

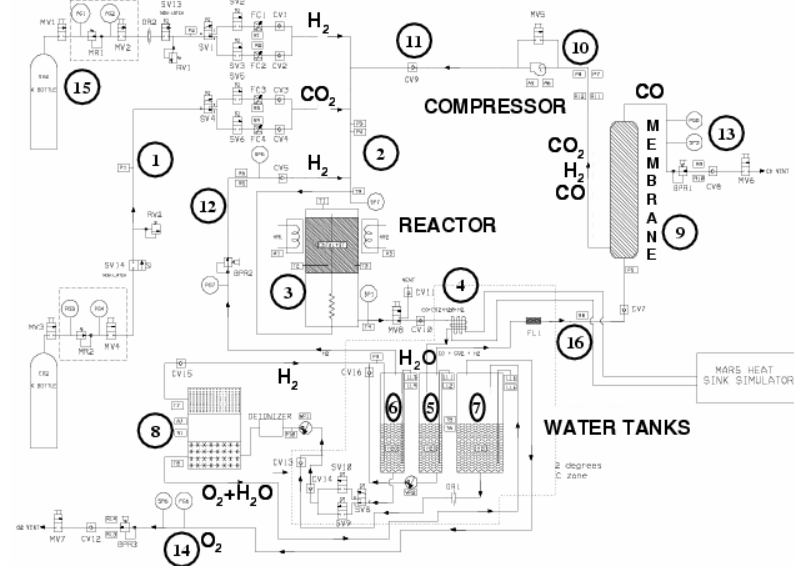
Function fitting



Temperature data



Monitoring a complex system



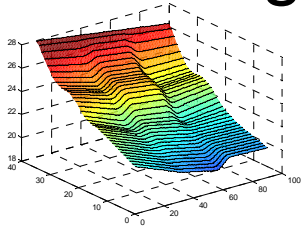
- Reverse water gas shift system (RWGS)
- Learn model of system from data
- Use model to predict behavior and detect faults

Types of learning problems

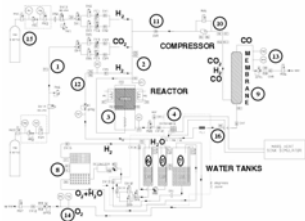
■ Classification



■ Regression



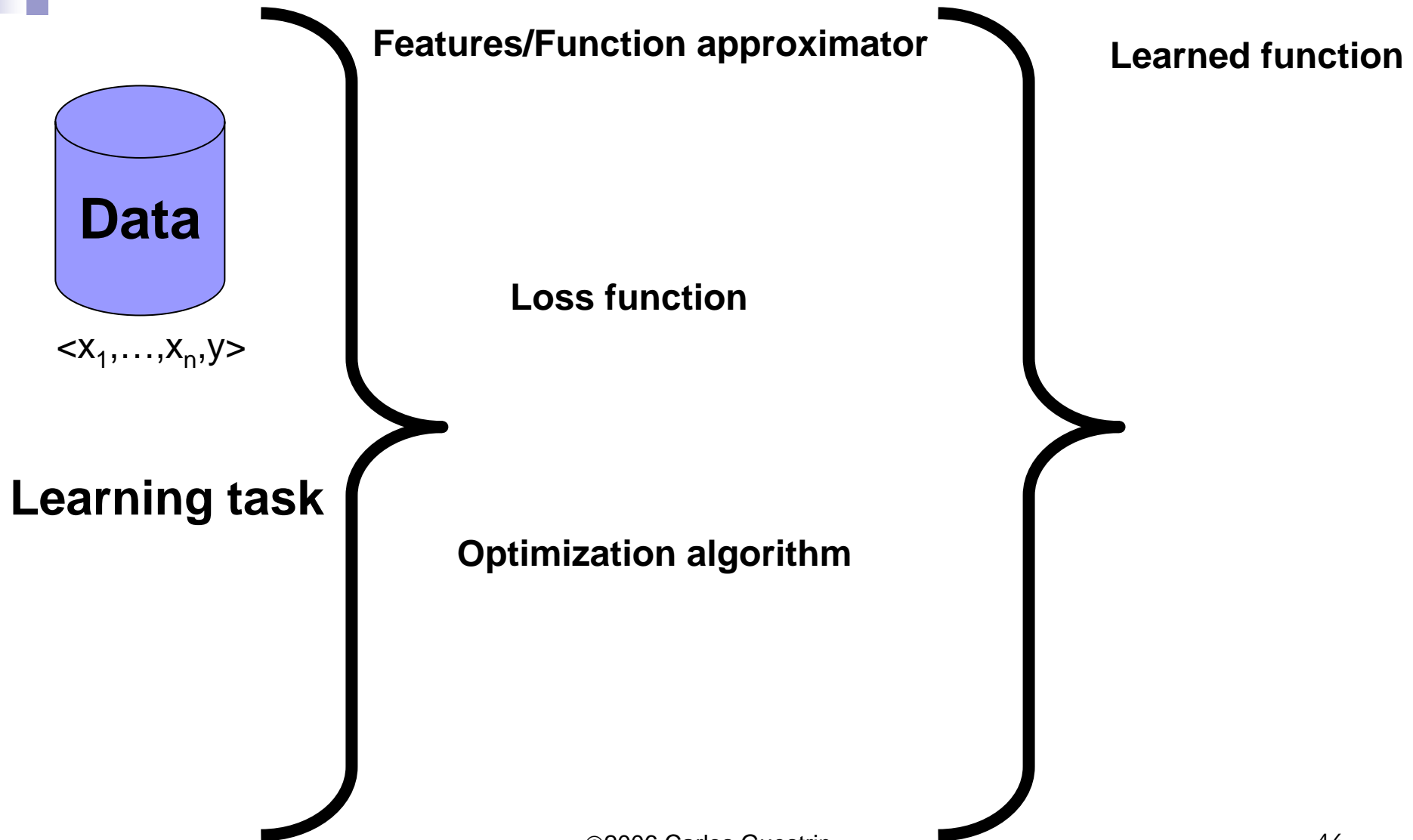
■ Density estimation



Input – Features

Output?

The learning problem



Comparing learning algorithms



- Hypothesis space
- Loss function
- Optimization algorithm

Naïve Bayes versus Logistic regression

Naïve Bayes

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(X|Y) = \prod_i P(X_i|Y)$$

Logistic regression

$$P(Y = 1|x) = \frac{1}{1 + \exp(w_0 + \sum_i w_i x_i)}$$

Naïve Bayes versus Logistic regression – Classification as density estimation

$$P(Y|X)$$

- Choose class with highest probability
- In addition to class, we get certainty measure

Logistic regression versus Boosting

Logistic regression

$$P(Y = y_i | \mathbf{x}) = \frac{1}{1 + \exp(-y_i(\mathbf{w} \cdot \mathbf{x} + b))}$$

Log-loss

$$\sum_{j=1}^m \log [1 + \exp(-y_j(\mathbf{w} \cdot \mathbf{x}_j + b))]$$

Boosting

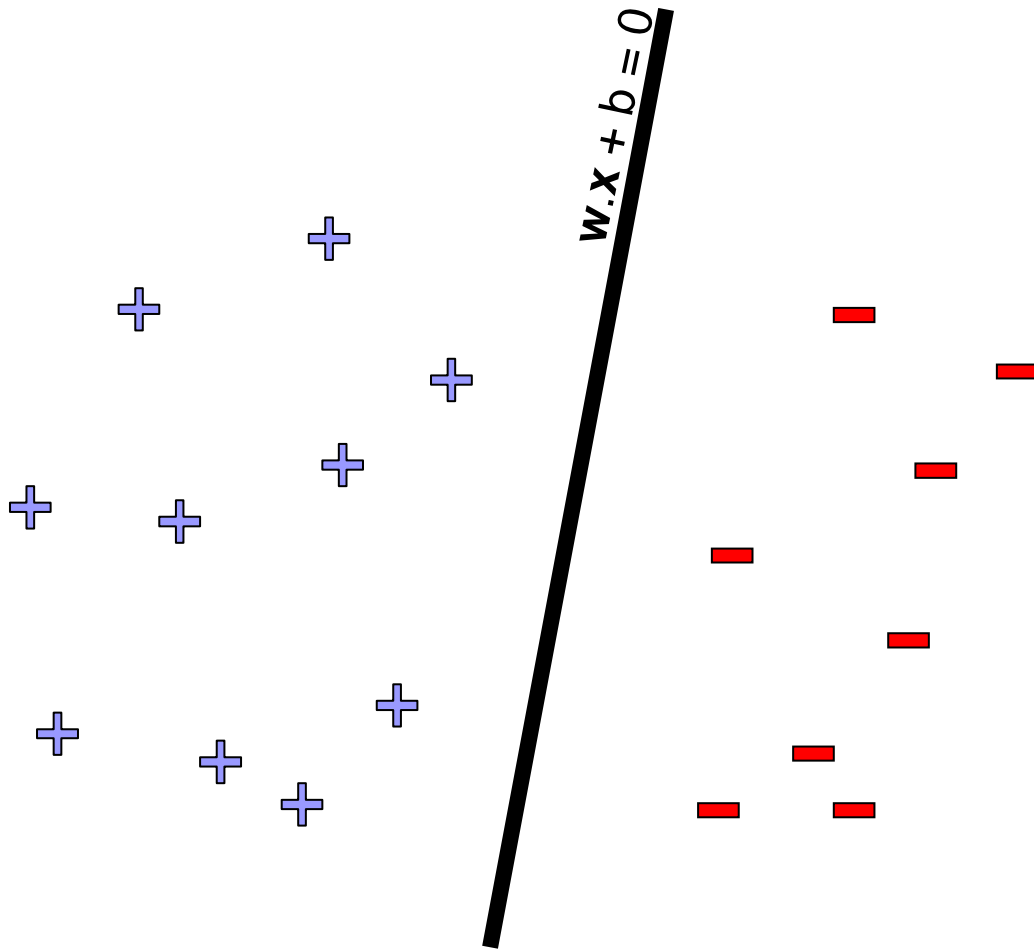
Classifier

$$\text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$$

Exponential-loss

$$\frac{1}{m} \sum_{j=1}^m \exp \left(-y_j \sum_{t=1}^T \alpha_t h_t(\mathbf{x}_j) \right)$$

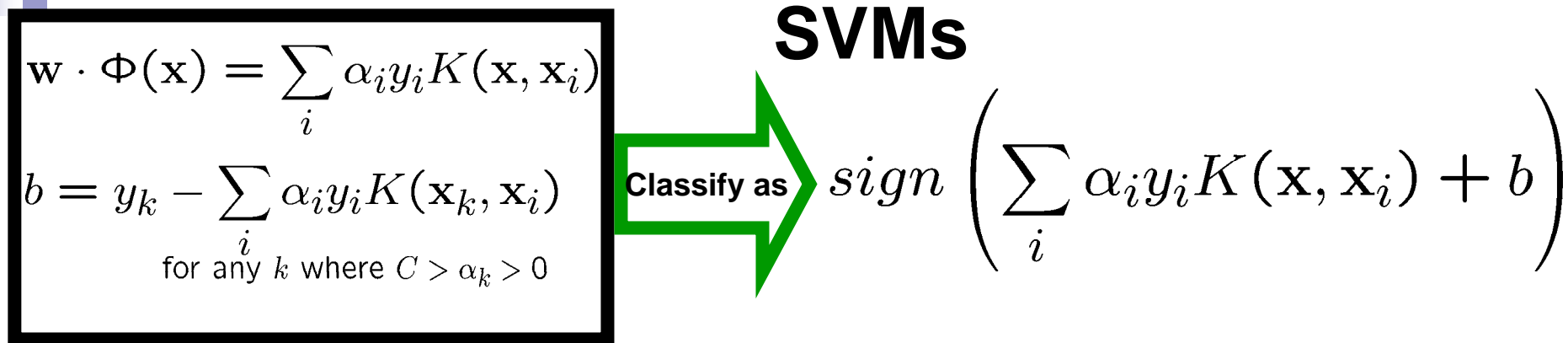
Linear classifiers – Logistic regression versus SVMs



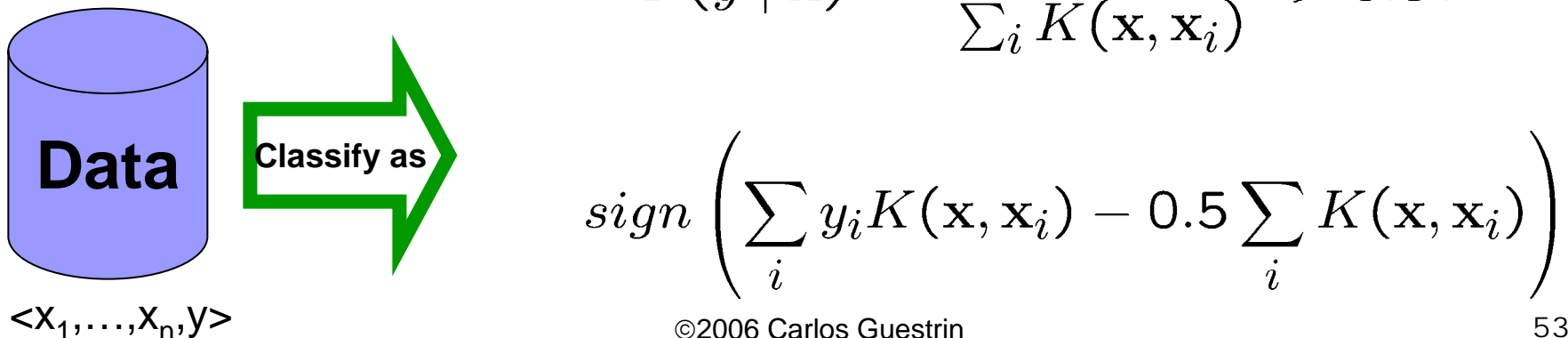
What's the difference between SVMs and Logistic Regression? (Revisited again)

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	Yes!
Solution sparse	Often yes!	Almost always no!
Type of learning		

SVMs and instance-based learning



Instance based learning

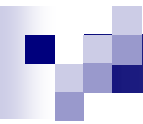


Instance-based learning versus Decision trees

1-Nearest neighbor

Decision trees

Logistic regression versus Neural nets


$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{-(w_0 + \sum_i w_i x_i)}}$$

Logistic regression

Neural Nets

Linear regression versus Kernel regression

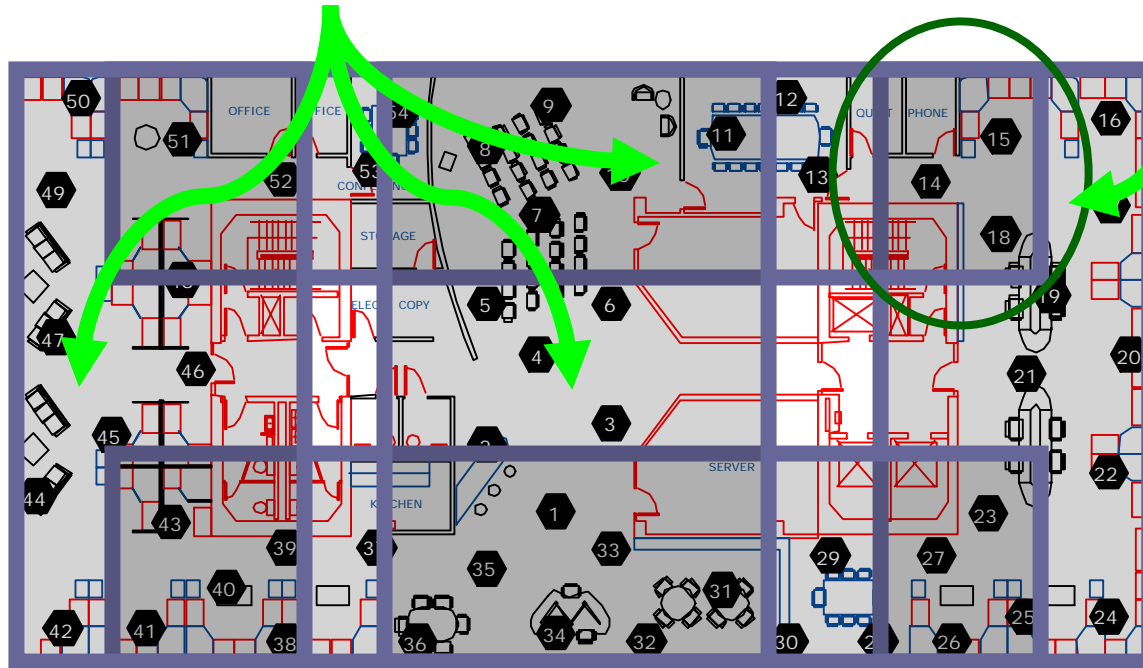
**Linear
Regression**

**Kernel
regression**

**Kernel-weighted
linear regression**

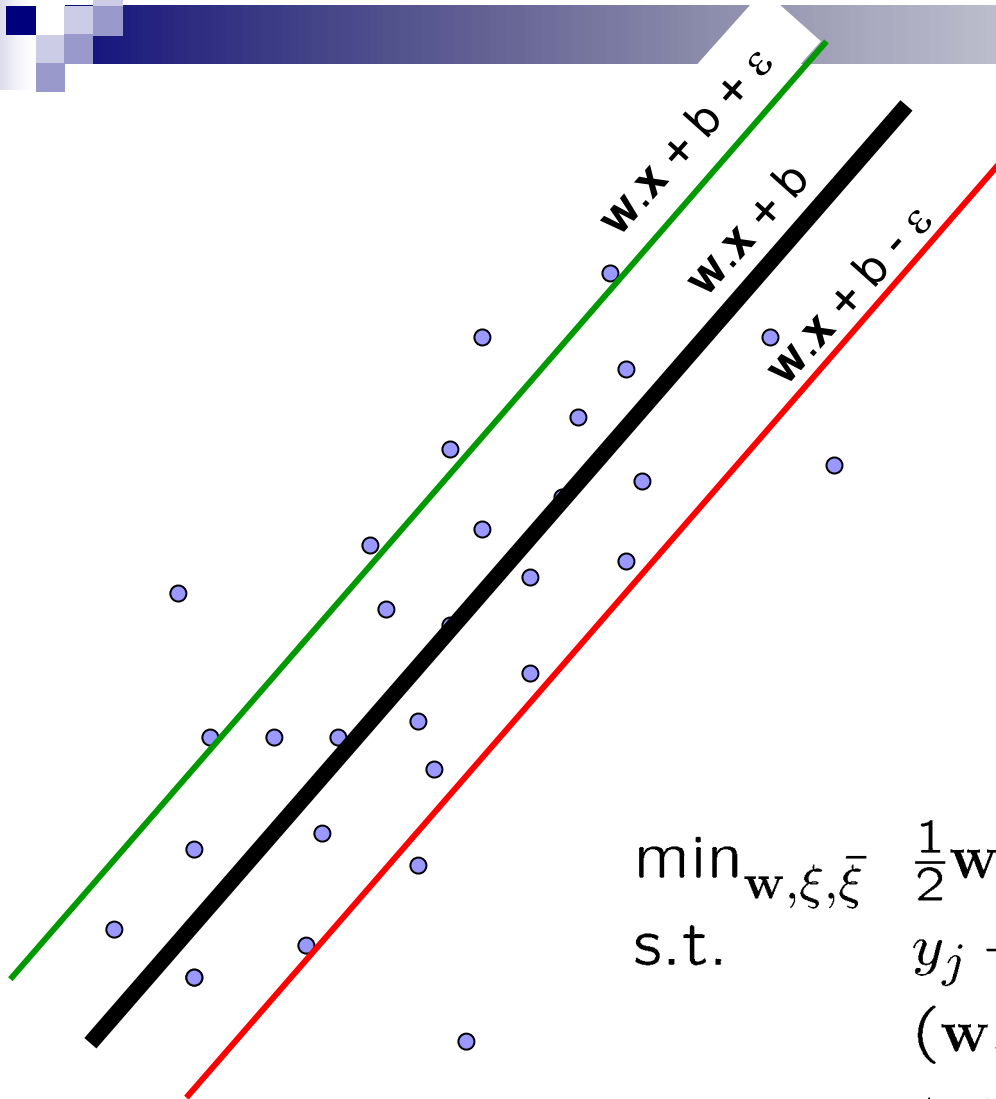
Kernel-weighted linear regression

Local basis functions for each region



Kernels
average
between
regions

SVM regression



$$\begin{aligned} \min_{w, \xi, \bar{\xi}} \quad & \frac{1}{2} w \cdot w + C \sum_{j=1}^m (\xi_j + \bar{\xi}_j) \\ \text{s.t.} \quad & y_j - (w \cdot x_j + b) \leq \epsilon + \xi_j \\ & (w \cdot x_j + b) - y_j \leq \epsilon + \bar{\xi}_j \\ & \xi_j \geq 0, \quad \bar{\xi}_j \geq 0, \quad \forall j \end{aligned}$$

BIG PICTURE

(a few points of comparison)

learning
task

loss
function

DE	density estimation
CI	Classification
Reg	Regression
LL	Log-loss/MLE
Mrg	Margin-based
RMS	Squared error

Naïve
Bayes
DE, LL

Boosting
CI, exp-loss

Logistic
regression
DE, LL

SVMs
CI, Mrg

SVM
regression
Reg, Mrg

kernel
regression
Reg, RMS

Instance-based
Learning
DE, CI, Reg

Neural
Nets
DE, CI, Reg, RMS

Decision
trees
DE, CI, Reg

linear
regression
Reg, RMS

This is a very incomplete view!!!