

CHAPTER 2:

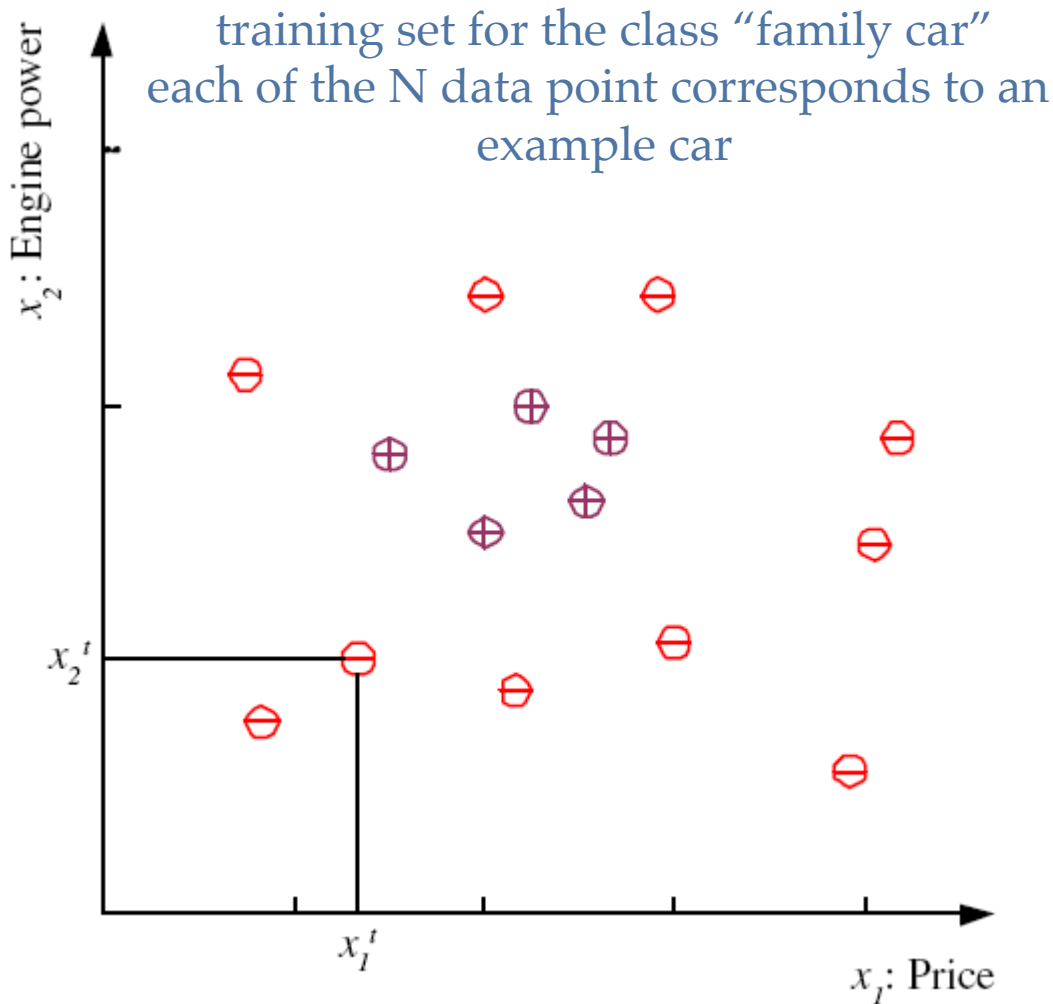
SUPERVISED LEARNING

Learning a Class from Examples

2

- Class C of a “family car”
 - ▣ Prediction: Is car x a family car?
 - ▣ Knowledge extraction: What do people expect from a family car?
- Output:
 - Positive (+) “yes, it’s a family car” and
 - negative (–) “no, it’s not a family car” examples
- Input representation: information we have about each car
 - x_1 : price, x_2 : engine power

Training set X



$$X = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

each car is represented as:

$$\{\mathbf{x}^t, r^t\}$$

where:

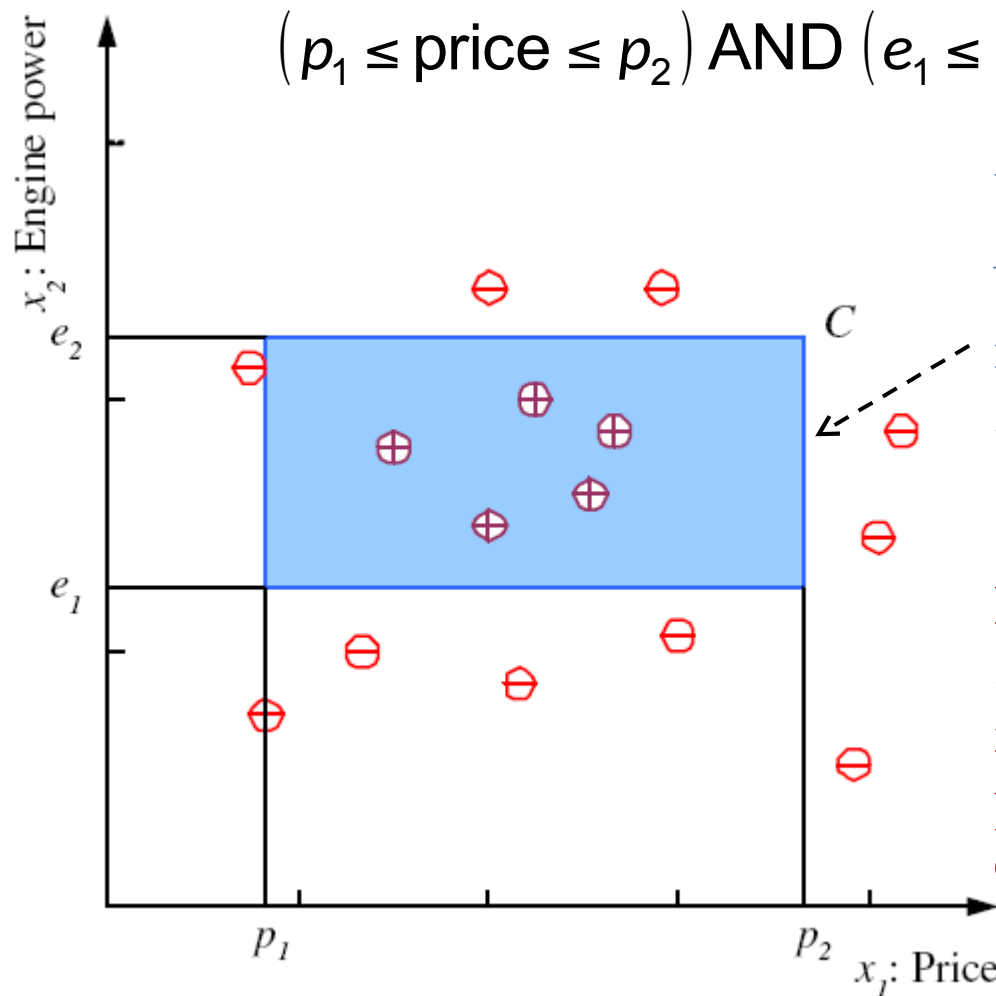
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

r is the example's "class"
also called "label",
"desired class" or "actual class"

$$r = \begin{cases} 1, & \text{if } \mathbf{x} \text{ is a positive example} \\ 0, & \text{if } \mathbf{x} \text{ is a negative example} \end{cases}$$

Class C:

4



$$(p_1 \leq \text{price} \leq p_2) \text{ AND } (e_1 \leq \text{engine power} \leq e_2)$$

Assume that the actual Class ("family car") is characterized by this blue rectangle, defined in the price and engine power space

But that rectangle is not explicitly given to us!

How to learn a good approximation to this rectangle based just on the positive and negative examples given in data set X?

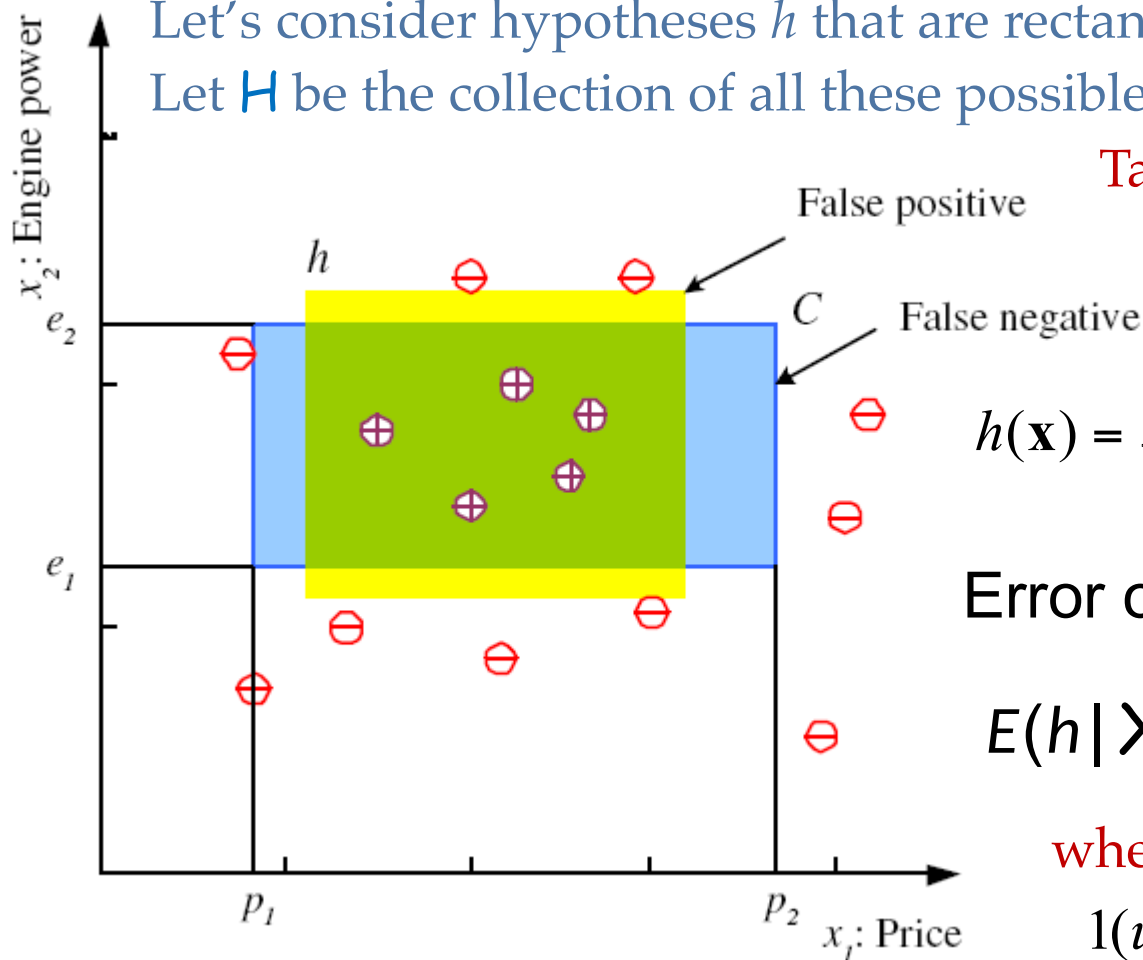
We create hypotheses!

Hypothesis class space H

5

Let's consider hypotheses h that are rectangles in the x_1, x_2 space
Let H be the collection of all these possible rectangles

Take for example h to be the yellow rectangle



$$h(\mathbf{x}) = \begin{cases} 1, & \text{if } \mathbf{x} \text{ is inside of } h \\ 0, & \text{if } \mathbf{x} \text{ is outside of } h \end{cases}$$

Error of h on H

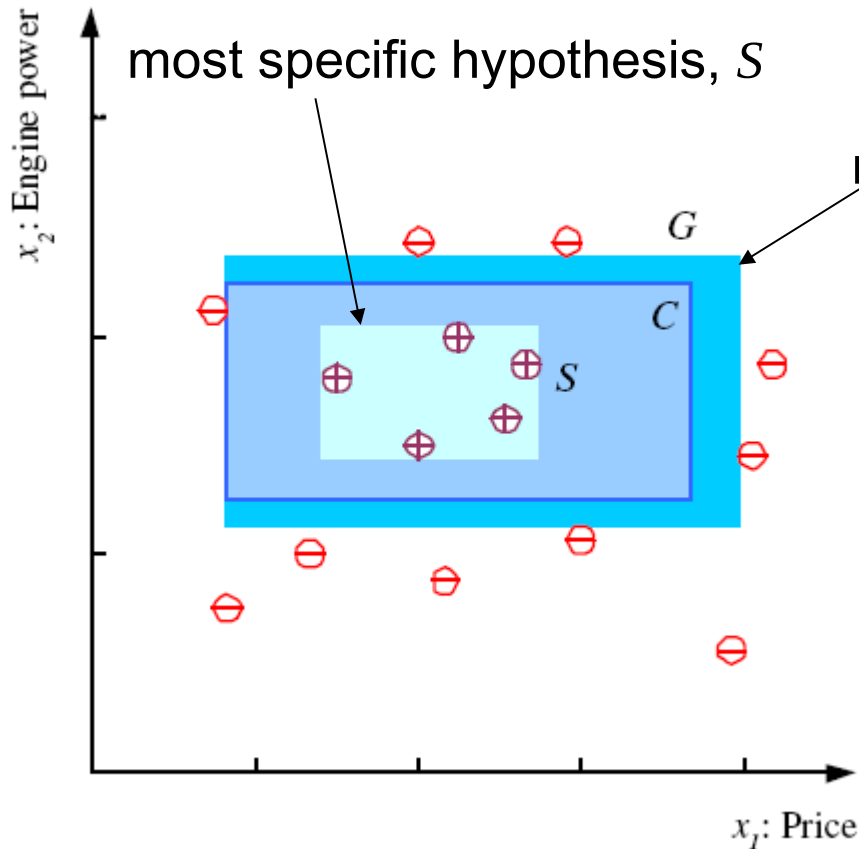
$$E(h | X) = \sum_{t=1}^N 1(h(\mathbf{x}^t) \neq r^t)$$

where

$$1(u) = \begin{cases} 1, & \text{if } u = \text{true} \\ 0, & \text{if } u = \text{false} \end{cases}$$

S, G, and the Version Space

6



most general hypothesis, G

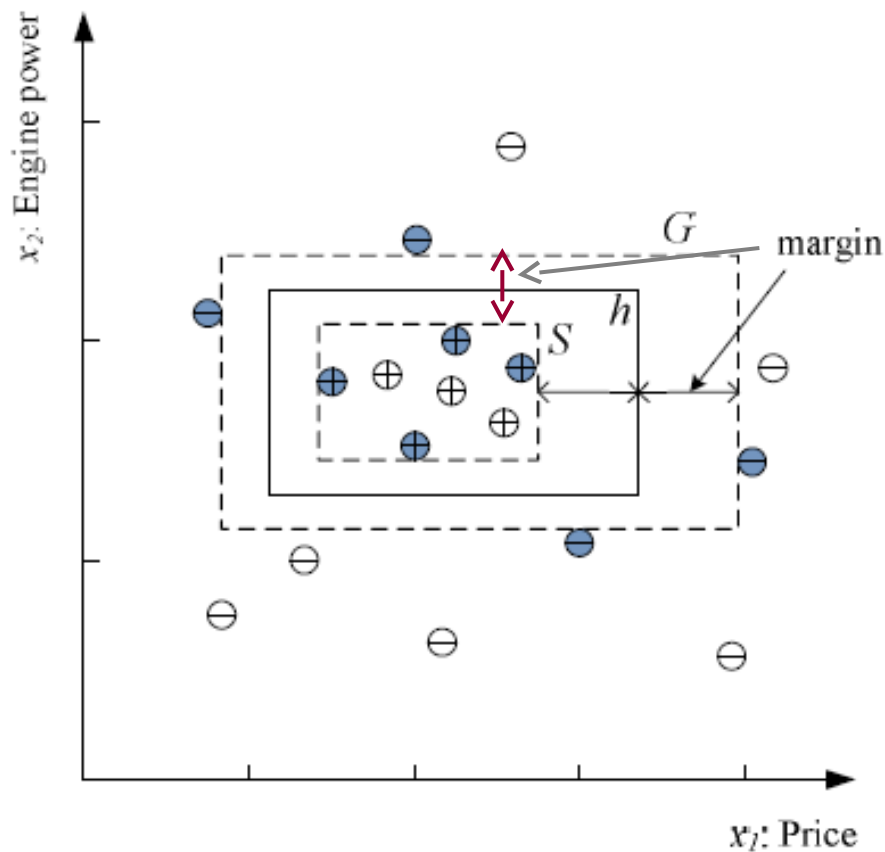
most specific hypothesis, S

$h \in H$, between S and G is
consistent with the training set
 X and make up the version
space
(Mitchell, 1997)

Margin

7

- Choose h with largest margin



Why the largest margin?

because we want to minimize the generalization error.

What's "generalization"?

Generalization: How well a model (i.e., hypothesis) performs on new data

Vapnik-Chernonenkis (VC) Dimension

8

- N points can be labeled in 2^N ways as $+/-$
- \mathcal{H} shatters N if and only if:
 - ▣ there is a set of N points in 2D such that for each of the 2^N possible ways of labelling these N points, there exists $h \in \mathcal{H}$ that is consistent with this labelling (i.e., h correctly separates the $+$ from the $-$ points)
- $VC(\mathcal{H}) =$ maximum N that can be shattered by \mathcal{H}
measures the “capacity” of \mathcal{H}

see example on the next slide

VC Dimension: Example

9

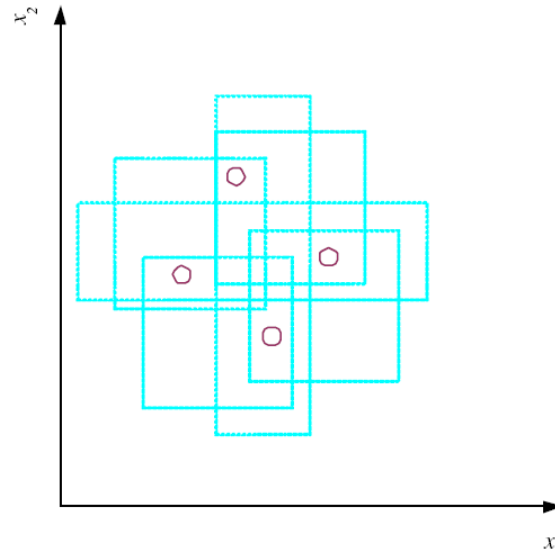
- Let H be the set of 2D axis-aligned rectangles
- $VC(H) = ?$ (value)
- Remember:
 - ▣ $VC(H) =$ maximum N that can be shattered by H
 - ▣ H shatters N if and only if:
 - there is a set of N points in 2D such that for each of the 2^N possible ways of labelling these N points, there exists $h \in H$ that is consistent with this labelling (i.e., correctly separates the + from the – points)

(cont.)

VC Dimension: Example (cont.)

10

- Let H be the set of 2D axis-aligned rectangles
- $VC(H) = 4$



The family of axis-aligned rectangles shatters 4 points only !

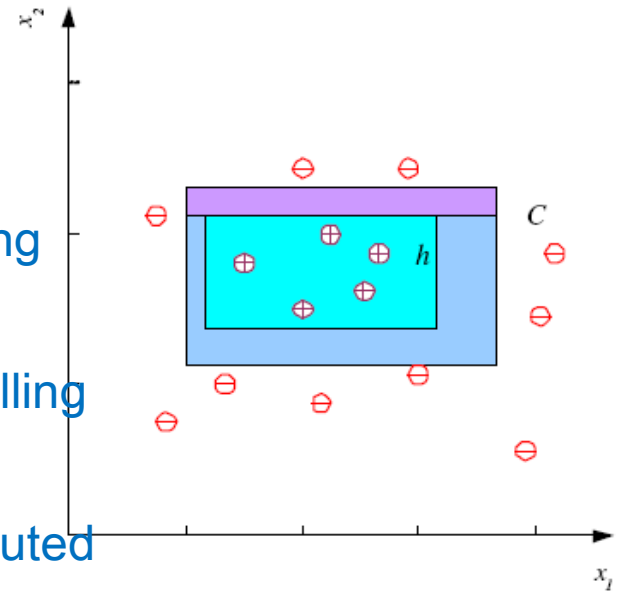
Probably Approximately Correct (PAC) Learning

11

- If we are to use the tightest rectangle h as our chosen hypothesis, how many training examples do we need so that we can guarantee that our hypothesis is approximately correct? In other words:
- Given δ ($1-\delta$ is the desired minimum confidence) and ϵ (desired maximum error): How many training examples, N , should we have such that with probability at least $1-\delta$, h has error at most ϵ ?
that is, $\Pr(E(h|X) \leq \epsilon) \geq 1 - \delta$

(Blumer et al., 1989)

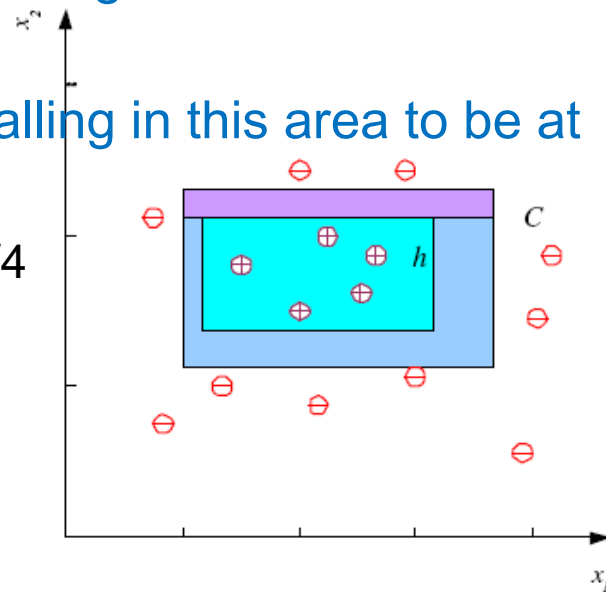
- $E(h|X)$ would come from a positive example falling in the area between C and h
- We need the probability of a positive example falling in this area to be at most ϵ
- Assume that data instances are uniformly distributed in space and are independent from each other



Probably Approximately Correct (PAC) Learning (cont.)

12

- Given δ ($1-\delta$ is the desired minimum confidence) and ϵ (desired maximum error): How many training examples, N , should we have such that with probability at least $1-\delta$, the tightest rectangle h has error at most ϵ ? $\Pr(E(h|X) \leq \epsilon) \geq 1 - \delta$
- $E(h|X)$ would come from a positive example falling in the area between C and h
- We need the Probab. of a positive example falling in this area to be at most ϵ
- Pr of an instance falling in a strip should at most $\epsilon/4$
- Pr that the data instance misses a strip is $1 - \epsilon/4$
- Pr that N instances miss a strip $(1 - \epsilon/4)^N$
- Pr that N instances miss 4 strips $4(1 - \epsilon/4)^N$
- Hence we need: $4(1 - \epsilon/4)^N \leq \delta$



(cont.)

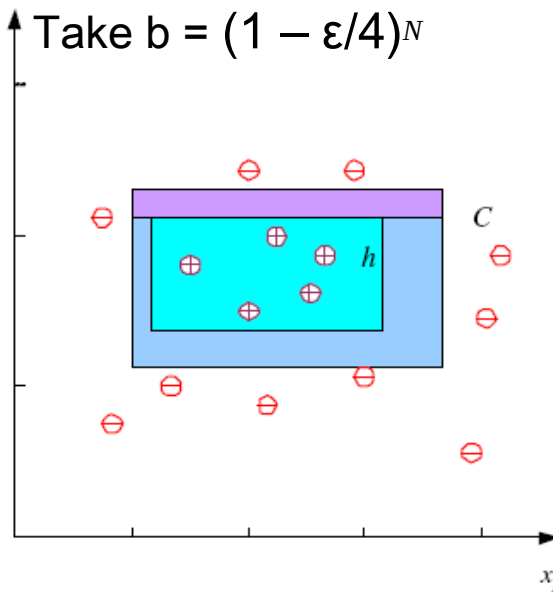
Probably Approximately Correct (PAC) Learning (cont.)

13

- Given δ ($1-\delta$ is the desired minimum confidence) and ϵ (desired maximum error): How many training examples, N , should we have such that with probability at least $1-\delta$, the tightest rectangle h has error at most ϵ ?
 $\Pr(E(h|X) \leq \epsilon) \geq 1 - \delta$
- we need: $4(1 - \epsilon/4)^N \leq \delta$

Note that $b = e^{\ln(b)} = \exp(\ln(b))$ and that $(1 - x) \leq e^{-x}$. Take $b = (1 - \epsilon/4)^N$

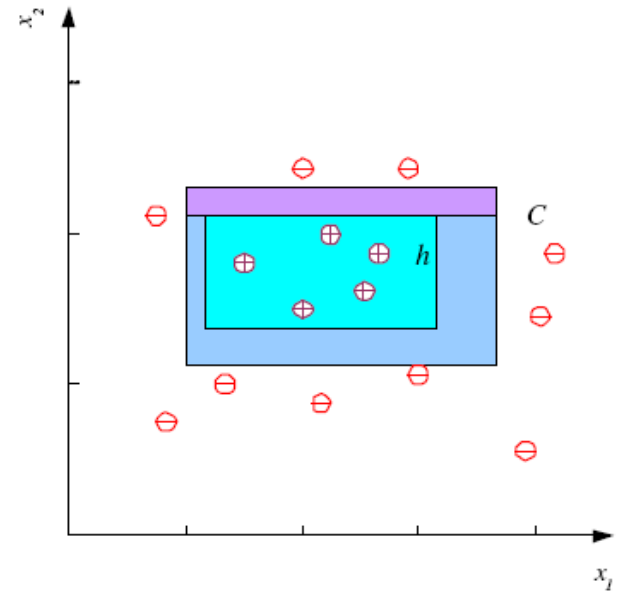
- $4(1 - \epsilon/4)^N \leq \delta$
 $4 \exp(\ln[(1 - \epsilon/4)^N]) \leq \delta$
- $4 \exp(N \ln(1 - \epsilon/4)) \leq \delta$, now using $(1 - x) \leq \exp(-x)$:
- $4 \exp(N \ln(1 - \epsilon/4)) \leq 4 \exp(N \ln(\exp(-\epsilon/4)))$
- Let's make $4 \exp(N \ln(\exp(-\epsilon/4))) = 4 \exp(N (-\epsilon/4)) \leq \delta$
- $4/\delta \leq \exp(\epsilon N/4)$
- $\ln(4/\delta) \leq \ln(\exp(\epsilon N/4)) = \epsilon N/4$
- and so $N \geq (4/\epsilon) \ln(4/\delta)$



Probably Approximately Correct (PAC) Learning (cont.)

14

- Given δ ($1-\delta$ is the desired minimum confidence) and ϵ (desired maximum error): How many training examples N should we have, such that with probability at least $1-\delta$, the tightest rectangle h has error at most ϵ ? $\Pr(E(h|X) \leq \epsilon) \geq 1 - \delta$
- Answer: $N \geq (4/\epsilon)\ln(4/\delta)$
- Example: How many training examples, N , do we need so that $\Pr(E(h|X) \leq 0.1) \geq 95\%$?
 - ▣ Here $\epsilon = 0.1$ and $1 - \delta = 0.95$ and so $\delta = 0.05$
 - ▣ $N \geq (4/\epsilon)\ln(4/\delta) = (4/0.1)\ln(4/0.05)$
 - ▣ $N \geq 40 \ln(80) = 175.28$
 - ▣ Hence our training set X should contain at least $N = 176$ data instances.



Noise and Model Complexity

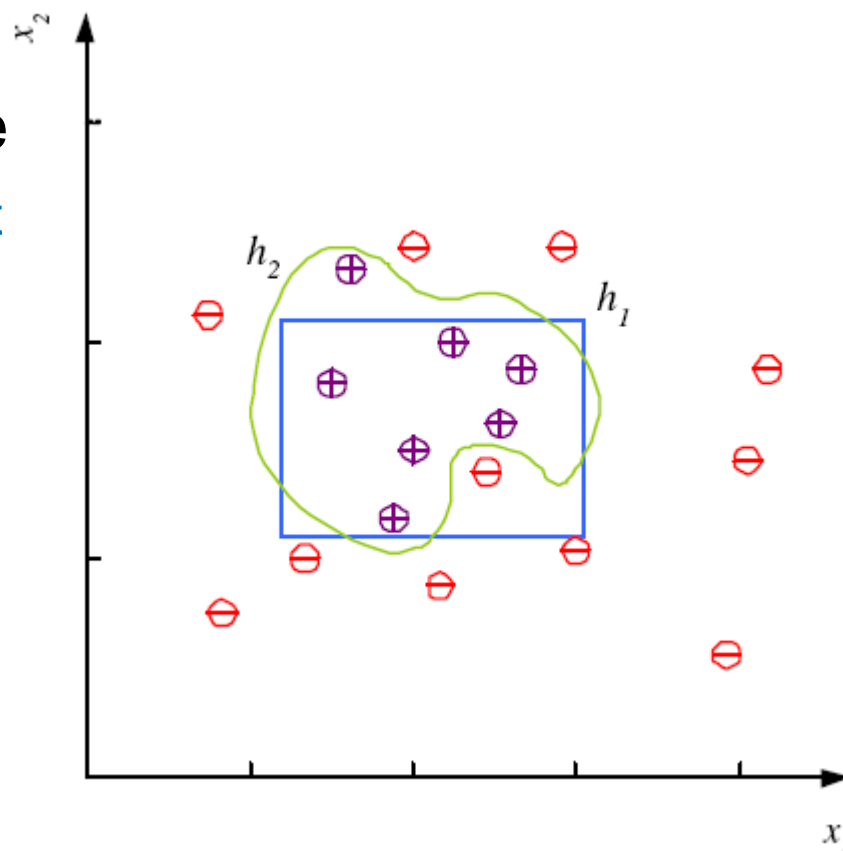
15

Noise:

- Unwanted anomaly in the data

Potential Sources of Noise

- Imprecision in recording input attributes
- Errors in labeling the data points (“teacher noise”)
- “Hidden” or “latent” attributes that affect the labels of the instances, but which are not considered in the data



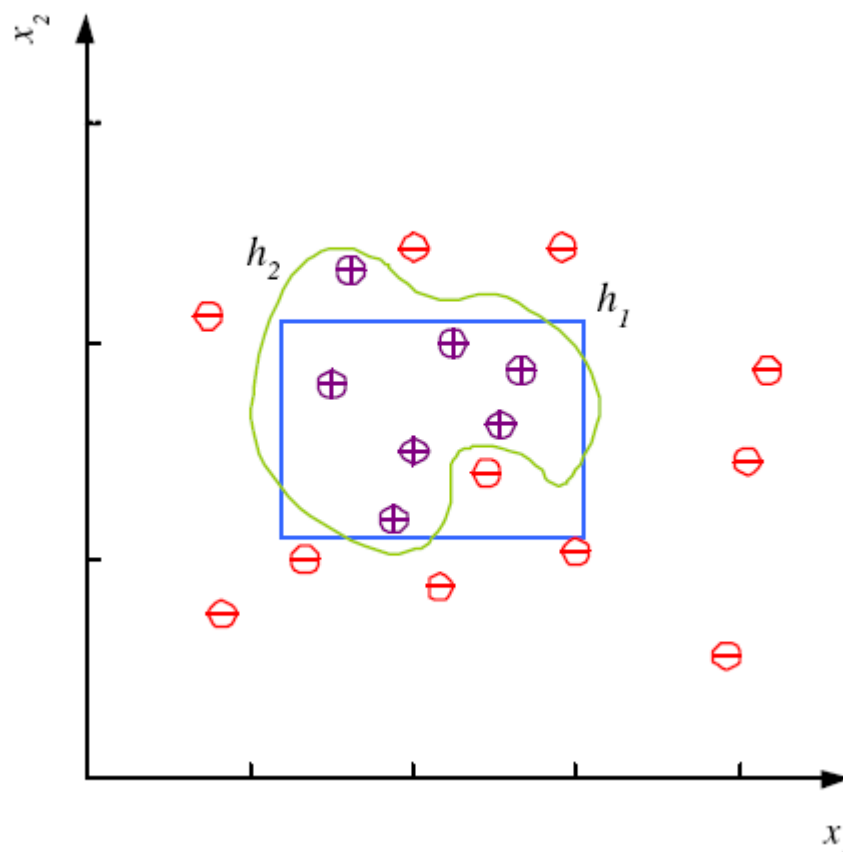
Noise and Model Complexity

(cont.)

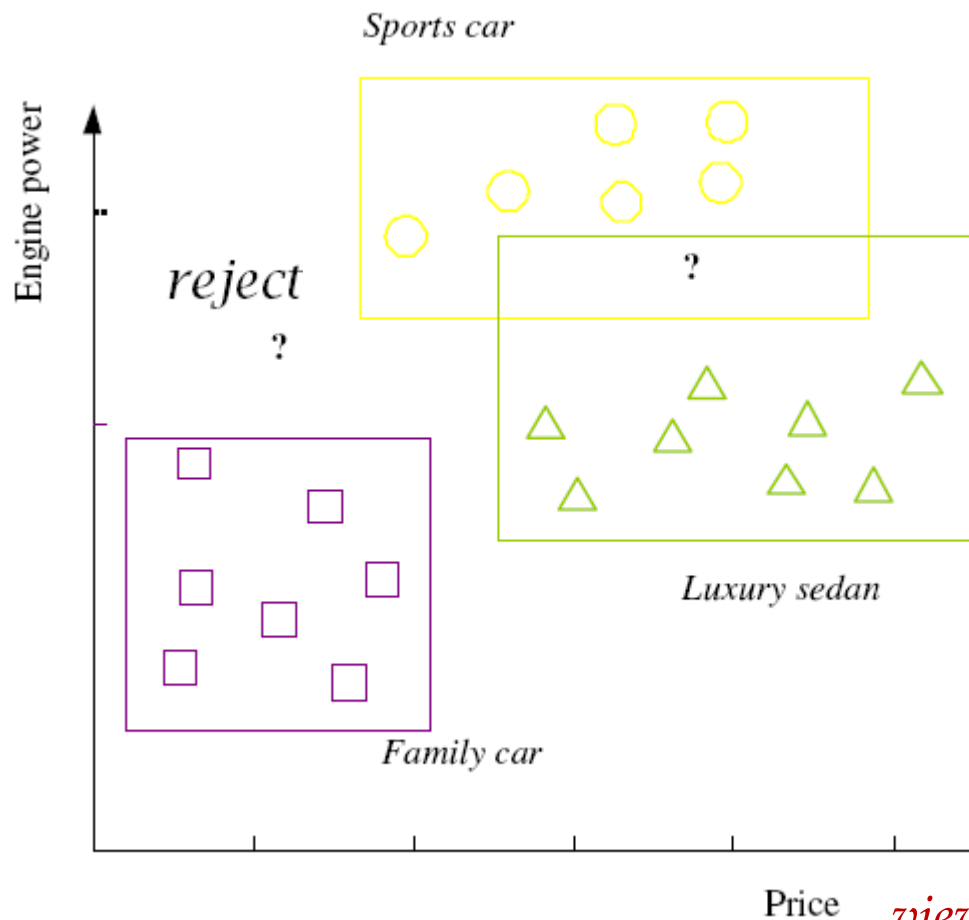
16

Use simplest **model** because

- Simpler to use
(lower computational complexity)
- Easier to train (lower space complexity)
- Easier to explain
(more interpretable)
- Generalizes better (lower variance - Occam's razor)



Multiple Classes, \mathcal{C}_i $i=1,\dots,K$



$$\mathbf{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

each r^t is now a vector

$$r_i^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in \mathcal{C}_i \\ 0 & \text{if } \mathbf{x}^t \in \mathcal{C}_j, j \neq i \end{cases}$$

Train hypotheses

$$h_i(\mathbf{x}), i = 1, \dots, K:$$

one hypothesis for each class

$$h_i(\mathbf{x}^t) = \begin{cases} 1 & \text{if } \mathbf{x}^t \in \mathcal{C}_i \\ 0 & \text{if } \mathbf{x}^t \in \mathcal{C}_j, j \neq i \end{cases}$$

*view a K -class classification problem
as K two-class classification problems*

Regression

$$\mathbf{X} = \left\{ x^t, r^t \right\}_{t=1}^N$$

$$r^t \in \Re$$

$$r^t = f(x^t) + \varepsilon$$

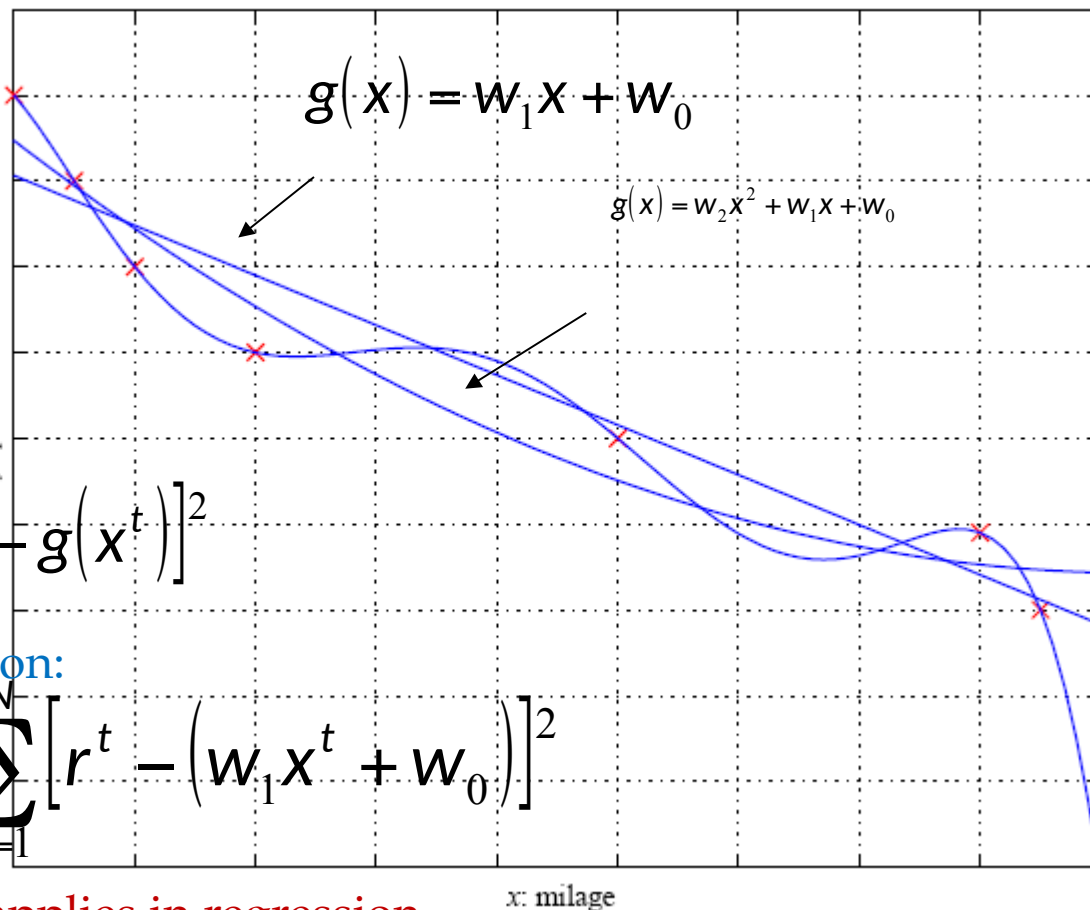
f : underlying function that we want to learn from data

$$E(g | \mathbf{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - g(x^t)]^2$$

Error of g in linear regression:

$$E(w_1, w_0 | \mathbf{X}) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_1 x^t + w_0)]^2$$

linear, second-order, and sixth-order polynomials



Occam's razor also applies in regression

Model Selection & Generalization

19

- Learning is an **ill-posed problem**; data is not sufficient to find a unique solution
- The need for **inductive bias**, assumptions about H
(e.g., our inductive bias in these slides' classification examples was to select rectangles as hypotheses - they could have been circles, squares or something else)
- **Generalization**: How well a model performs on new data
- **Overfitting**: H more complex than C or f
(e.g., fitting a 6th-order polynomial to noisy data sampled from a 2nd – order polynomial)
- **Underfitting**: H less complex than C or f
(e.g., fitting a 2nd-order polynomial to noisy data sampled from a 6th – order polynomial)

Triple Trade-Off

20

There is a trade-off between three factors (Dietterich, 2000):

1. Complexity of H , $c(H)$,
2. Training set size, N ,
3. Generalization error, E , on new data

As $N \uparrow$, $E \downarrow$

As $c(H) \uparrow$, first $E \downarrow$ and then $E \uparrow$

Cross-Validation

21

- To estimate generalization error, we need data unseen during training. We split the data as
 - ▣ Training set (e.g., 50%)
data instances used to construct the model
 - ▣ Validation set (e.g., 25%)
data instances used to test preliminary versions of the model and/or to refine the model
 - ▣ Test (publication) set (e.g., 25%)
data instances used to test the final model after it has been fully constructed
- Resampling when there is few data

Dimensions of a Supervised Learner

1. Model: $g(\mathbf{x} | \theta)$

2. Loss function: $E(\theta | \mathbf{X}) = \sum_t L(r^t, g(\mathbf{x}^t | \theta))$

to compute the difference between the actual and predicted values

3. Optimization procedure: $\theta^* = \arg \min_{\theta} E(\theta | \mathbf{X})$