

Ch.2:

Supervised Learning

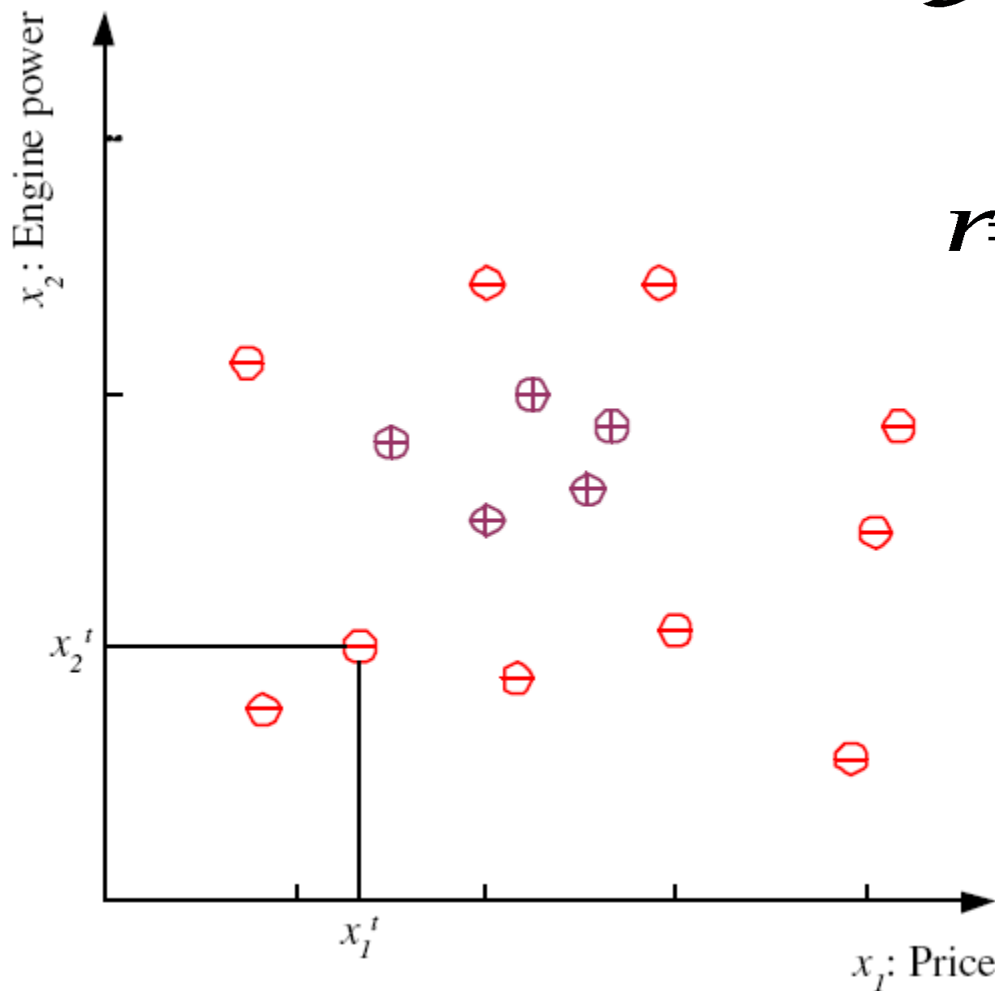
Learning a Class from Examples

- 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 (+) and negative (–) examples
- Input representation:
 - x_1 : price, x_2 : engine power

Training set \mathcal{X}

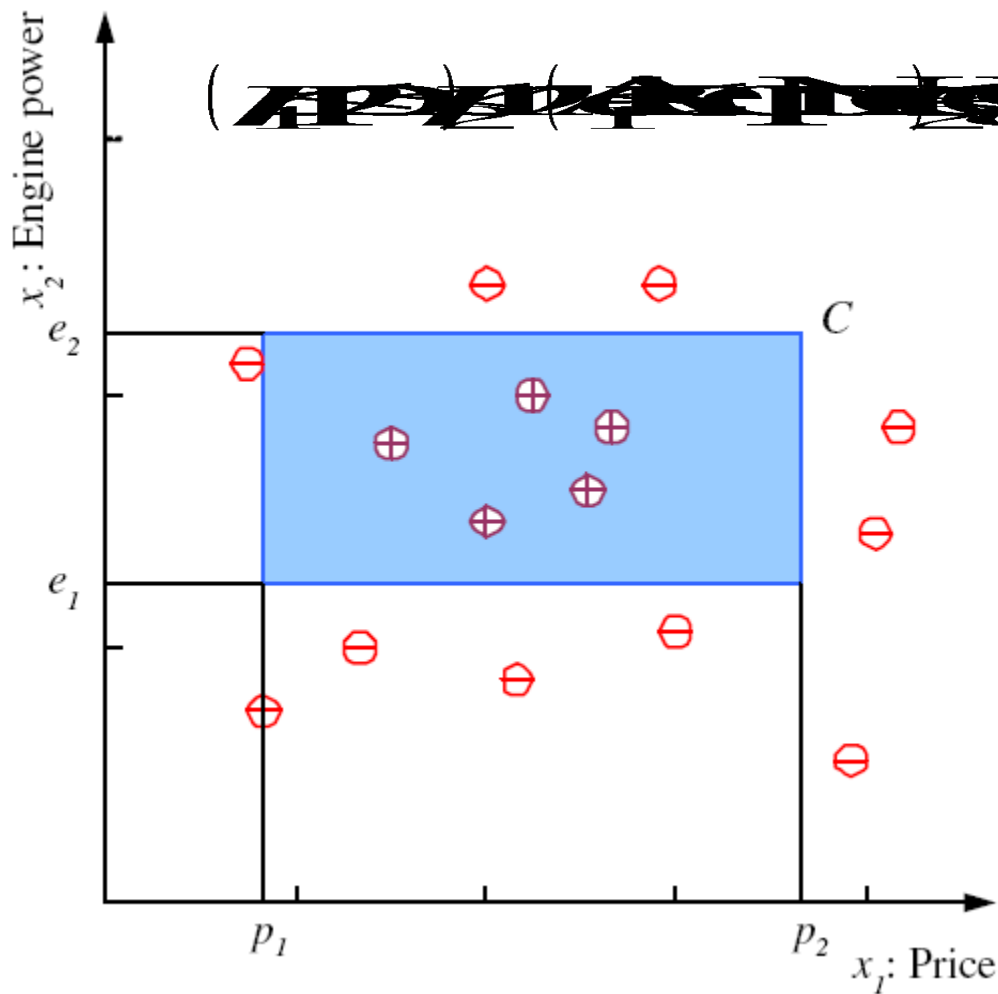
$$\mathcal{X} = \{x^t\}_{t=1}^N$$

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

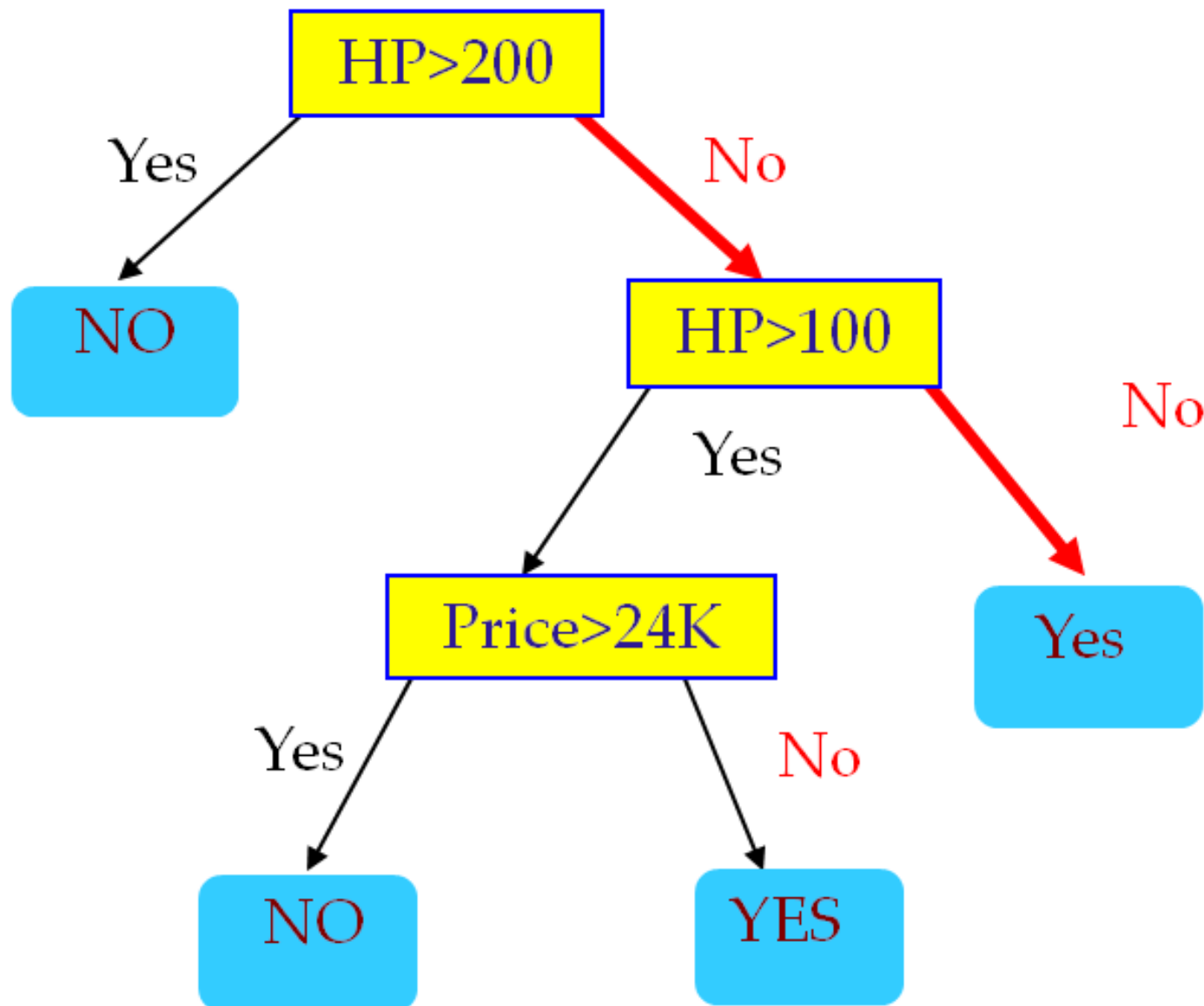


$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

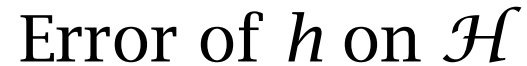
Class C



Family Car Decision Tree

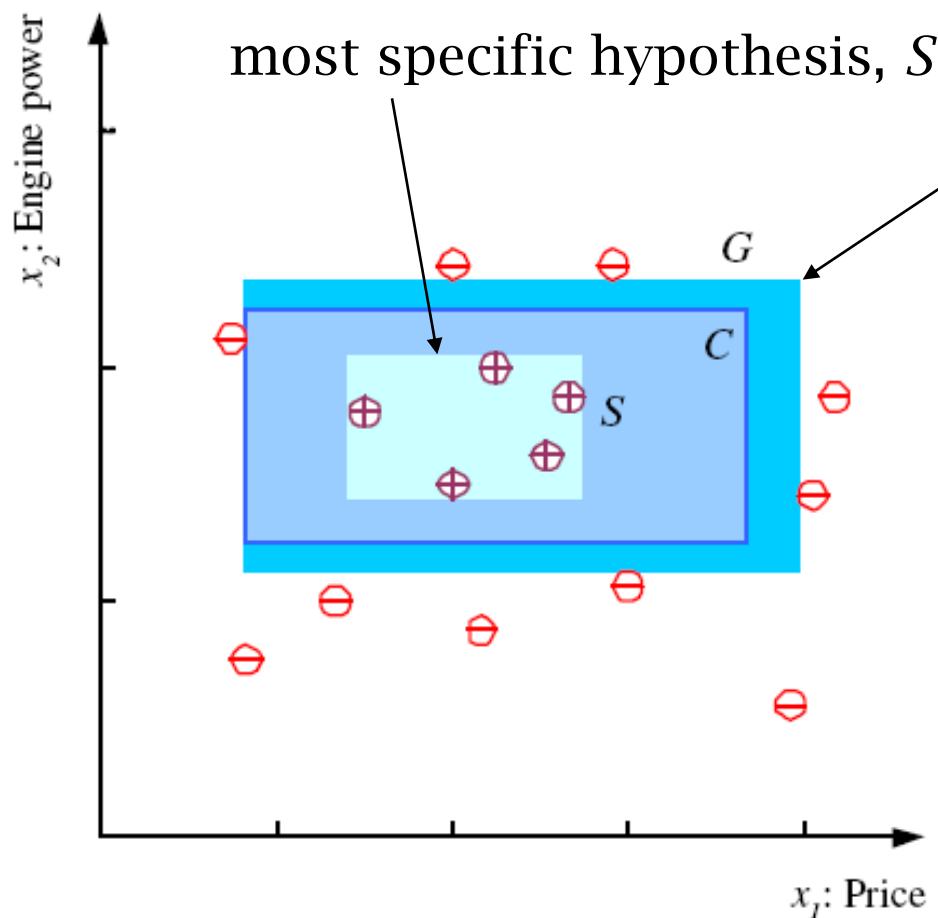


~~A classification~~



$$EhX = \sum_{t=1}^N U(\mathbf{x}^t) = A^t$$

S, G, and the Version Space



most general hypothesis, G

$h \in \mathcal{H}$, between S and G is
consistent

and make up the
version space

(Mitchell, 1997)

VC Dimension

Vapnik-Chervonenkis

■ N points can be labeled in 2^N ways as $+/-$

■ \mathcal{H} **shatters** N if there exists $h \in \mathcal{H}$ consistent

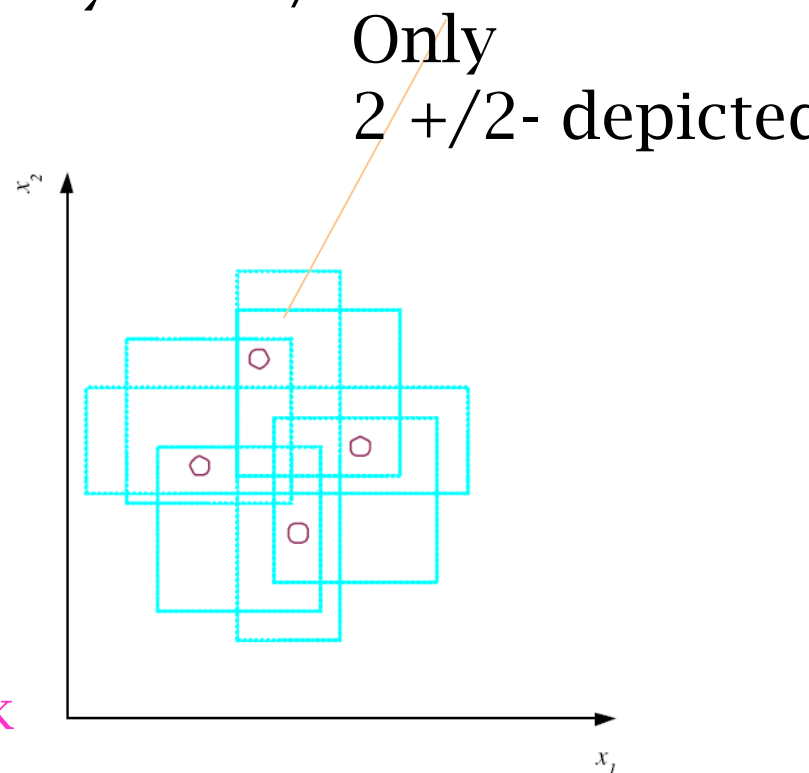
for any of these:

$$VC(\mathcal{H}) = N$$

■ *Does not work for 5 points!*

x
x
x
x
x

try it:
Homework

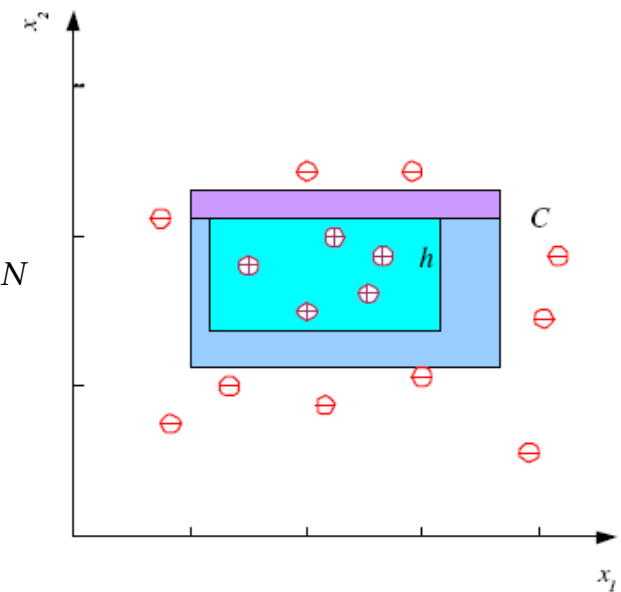


An axis-aligned rectangle shatters 4 points only !

Probably Approximately Correct (PAC) Learning

- How many training examples N should we have, such that with **probability at least** $1 - \delta$, h has **error at most** ϵ ?
(Blumer et al., 1989)

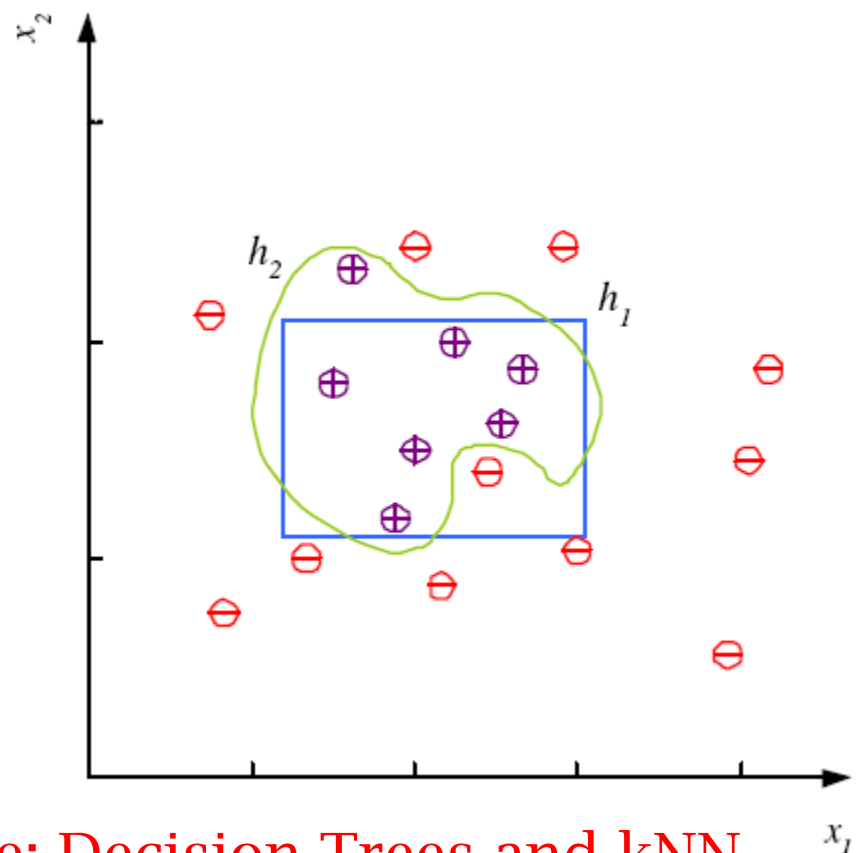
- Each strip is at most $\epsilon/4$
- Pr that we miss a strip $1 - \epsilon/4$
- Pr that N instances miss a strip $(1 - \epsilon/4)^N$
- Pr that N instances miss 4 strips $4(1 - \epsilon/4)^N$
- $4(1 - \epsilon/4)^N \leq \delta$ and $(1 - x) \leq \exp(-x)$
- $4\exp(-\epsilon N/4) \leq \delta$ and $N \geq (4/\epsilon)\log(4/\delta)$



Noise and Model Complexity

Use the simpler one 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)



2 Other Classification Technique: Decision Trees and kNN

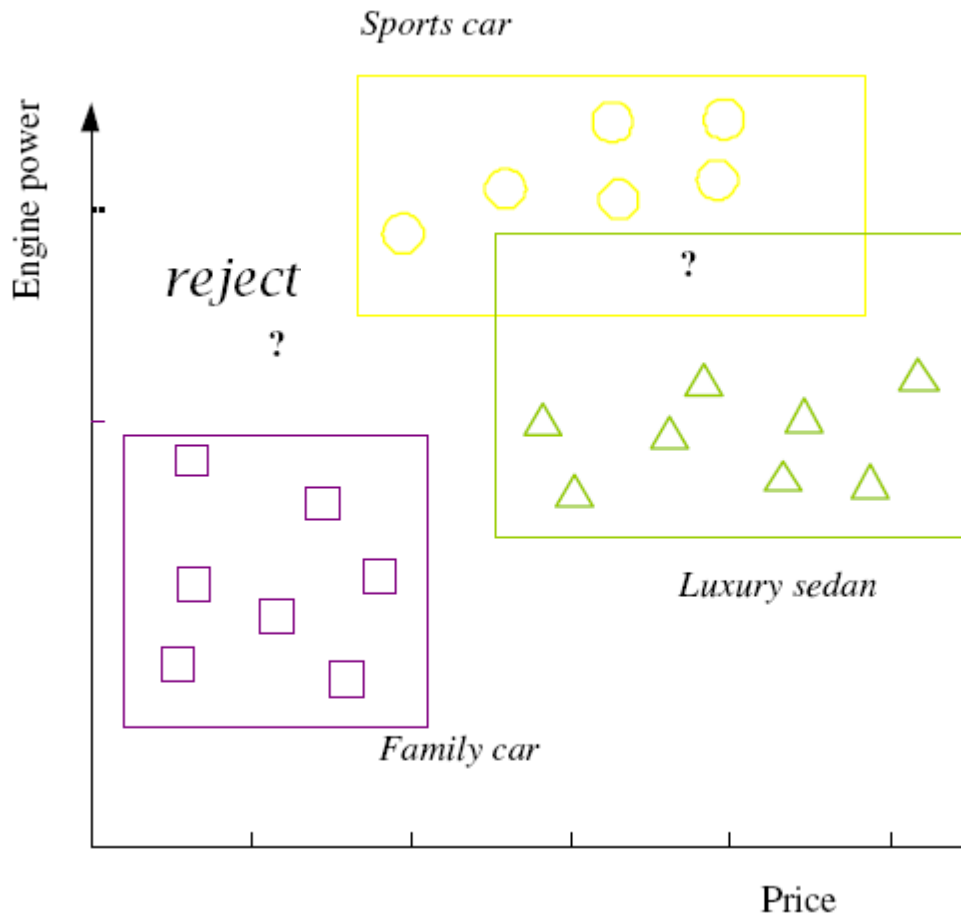
Multiple Classes, C_i $i=1,...,K$

$$\mathcal{X} = \{\mathbf{x}^t, P_{t=1}^N$$

$$r_i^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Train hypotheses
 $h_i(\mathbf{x})$, $i = 1, ..., K$:

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

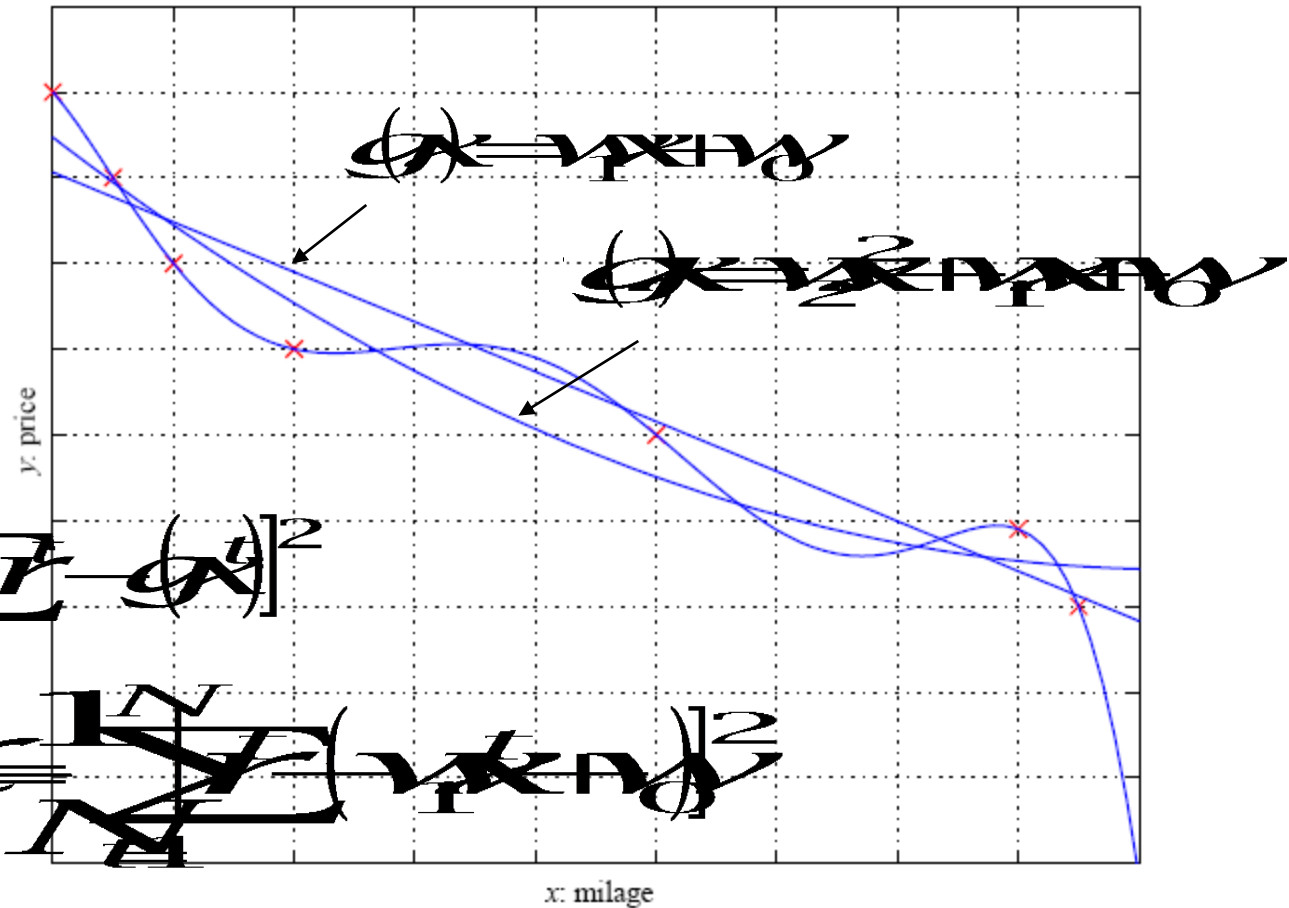


Regression

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

$$r^t \in \mathbb{R}$$

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



$$E(x) = \frac{1}{N} \sum_{t=1}^N \|r^t - \phi(x^t)\|^2$$

$$E_p(x) = \frac{1}{N} \sum_{t=1}^N \|\nabla_p r^t - \nabla_p \phi(x^t)\|^2$$

Finding Regression Coefficients

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

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

$$r^t \in \mathbb{R}$$

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

How to find w_1 and w_0 ?

Solve: $dE/dw_1=0$ and $dE/dw_0=0$

And solve the two obtained equations

Ungraded Homework!

~~$$E(w) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_0 + w_1 x^t)]^2$$~~

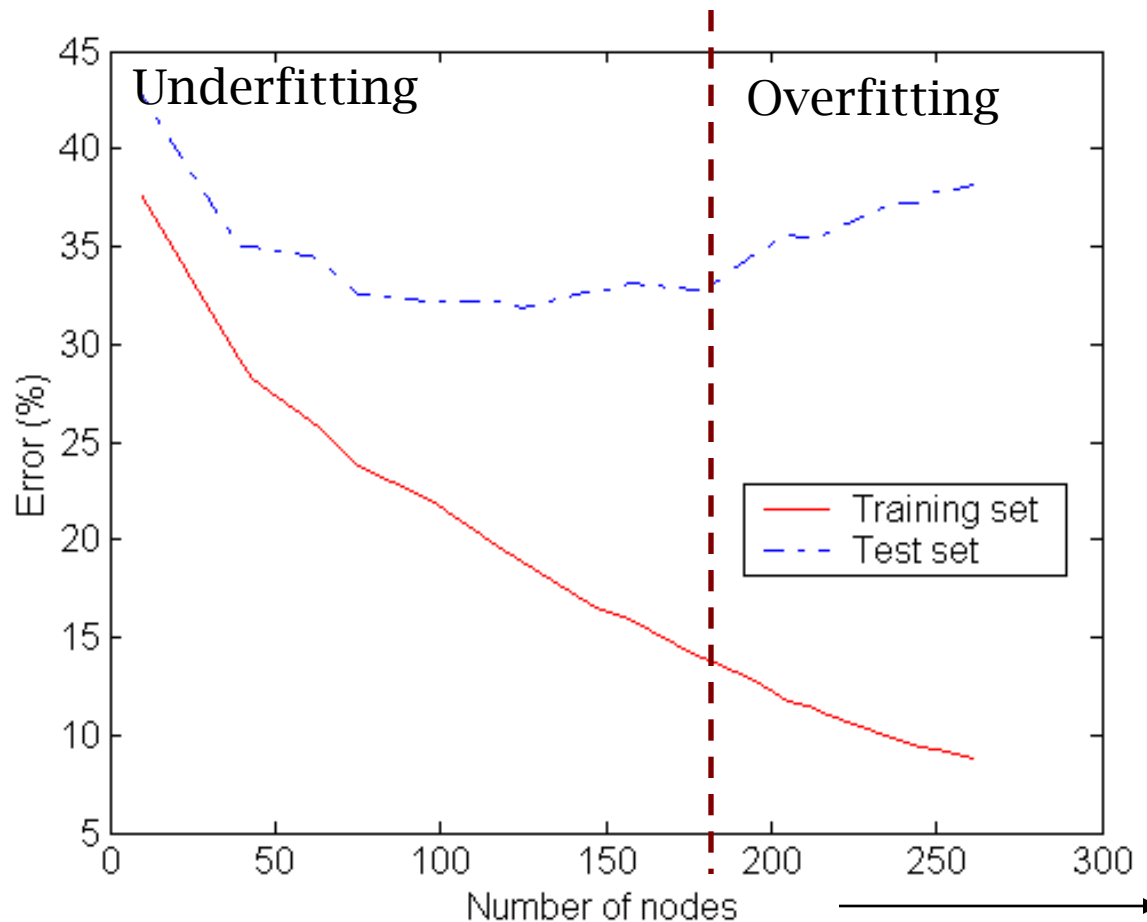
~~$$E(w) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_0 + w_1 x^t)]^2$$~~

http://en.wikipedia.org/wiki/Simple_linear_regression

Model Selection & Generalization

- Learning is an **ill-posed problem**; data is not sufficient to find a unique solution
- The need for **inductive bias**, assumptions about \mathcal{H}
- **Generalization**: How well a model performs on new data
- Overfitting: \mathcal{H} more complex than C or f
- Underfitting: \mathcal{H} less complex than C or f

Underfitting and Overfitting



Complexity of a Decision Tree := number of nodes It uses

Complexity of the classification function

Underfitting: when model is too simple, both training and test errors are large.

Overfitting: when model is too complex and test errors are large although training errors are small.

Cross-Validation

Error on new examples; actually the testing error is used as an estimation of the generalization error!

- Two errors: training error, and testing error usually *called generalization error*. Typically, the training error is smaller than the generalization error.
- To estimate generalization error, we need data unseen during training. We could split the data as
 - Training set (50%)
 - Validation set (25%)→optional, for selecting ML algorithm parameters (e.g. model complexity)
 - Test (publication) set (25%)
- Resampling when there is few data

Triple Trade-Off

overfitting

- There is a trade-off between three factors (Dietterich, 2003):
 1. Complexity of \mathcal{H} , $c(\mathcal{H})$,
 2. Training set size, N ,
 3. Generalization error, E on new data
- As $N \uparrow E \downarrow$
- As $c(\mathcal{H}) \uparrow$ first $E \downarrow$ and then $E \uparrow$
- As $c(\mathcal{H}) \uparrow$ the training error decreases for some time and then stays constant (frequently at 0)

Dimensions of a Supervised Learner

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

2. Loss function: $\ell(\theta; \mathbf{x}) = \sum_t \ell(g(\theta; \mathbf{x}_t), y_t)$

3. Optimization procedure:

$$\theta^* = \arg\min_{\theta} \ell(\theta; \mathbf{X})$$

Remark This procedure is typical for Parametric approaches to supervised learning; Non-parametric approaches work differently!