

# Artificial Intelligence

Learning from observations  
Chapter 18, AIMA

## Two types of learning in AI

*Deductive:* Deduce rules/facts from already known rules/facts. (We have already dealt with this)

$$(A \Rightarrow B \Rightarrow C) \Rightarrow (A \Rightarrow C)$$

*Inductive:* Learn new rules/facts from a data set  $\mathcal{D}$ .

$$\mathcal{D} = \{\mathbf{x}(n), y(n)\}_{n=1 \dots N} \Rightarrow (A \Rightarrow C)$$

We will be dealing with the latter, *inductive* learning, now

## Two types of inductive learning

*Supervised:* The machine has access to a teacher who corrects it.



*Unsupervised:* No access to teacher. Instead, the machine must search for "order" and "structure" in the environment.



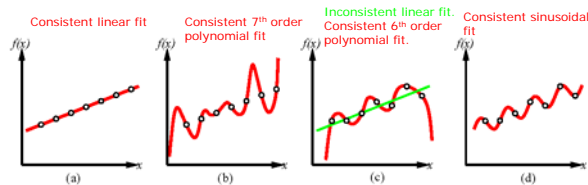
## Inductive learning - example A

<table border="1"><tr><td>○</td><td>○</td><td>×</td></tr><tr><td>×</td><td></td><td></td></tr><tr><td>×</td><td></td><td></td></tr></table>	○	○	×	×			×			<table border="1"><tr><td>○</td><td>○</td><td></td></tr><tr><td>×</td><td>×</td><td>×</td></tr><tr><td>×</td><td></td><td></td></tr></table>	○	○		×	×	×	×			<table border="1"><tr><td>○</td><td>○</td><td>×</td></tr><tr><td>×</td><td></td><td></td></tr><tr><td>×</td><td></td><td></td></tr></table>	○	○	×	×			×			Etc...
○	○	×																												
×																														
×																														
○	○																													
×	×	×																												
×																														
○	○	×																												
×																														
×																														
$\begin{pmatrix} -1 \\ -1 \\ 0 \\ +1 \\ 0 \\ 0 \\ +1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, f(\mathbf{x}) = +1$	$\begin{pmatrix} -1 \\ -1 \\ 0 \\ 0 \\ 0 \\ +1 \\ +1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, f(\mathbf{x}) = -1$	$\begin{pmatrix} -1 \\ -1 \\ 0 \\ 0 \\ 0 \\ +1 \\ +1 \\ 0 \\ +1 \\ 0 \end{pmatrix}, f(\mathbf{x}) = 0$																												

- $f(\mathbf{x})$  is the **target function**
- An **example** is a pair  $[\mathbf{x}, f(\mathbf{x})]$
- Learning task: find a **hypothesis**  $h$  such that  $h(\mathbf{x}) \approx f(\mathbf{x})$  given a training set of examples  $\mathcal{D} = \{[\mathbf{x}_i, f(\mathbf{x}_i)]\}, i = 1, 2, \dots, N$

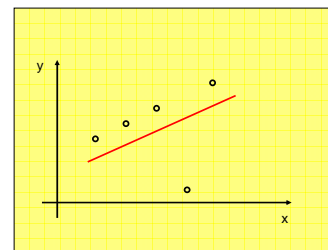
Inspired by a slide from V. Pavlovic

## Inductive learning – example B



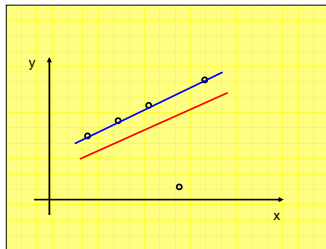
- Construct  $h$  so that it agrees with  $f$ .
- The hypothesis  $h$  is consistent if it agrees with  $f$  on all observations.
- Ockham's razor: Select the simplest consistent hypothesis.
- How achieve good generalization?

## Inductive learning – example C



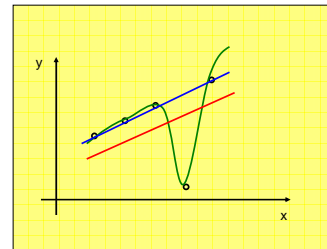
Example from V. Pavlovic @ Rutgers

## Inductive learning – example C



Example from V. Pavlovic @ Rutgers

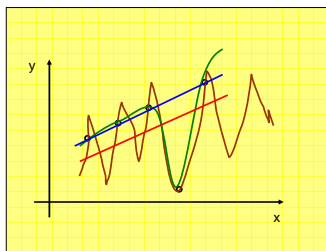
## Inductive learning – example C



Example from V. Pavlovic @ Rutgers

## Inductive learning – example C

Sometimes a consistent hypothesis is worse than an inconsistent



Example from V. Pavlovic @ Rutgers

## The idealized inductive learning problem

Find appropriate hypothesis space  $\mathbf{H}$  and find  $h(\mathbf{x}) \in \mathbf{H}$  with minimum "distance" to  $f(\mathbf{x})$  ("error")

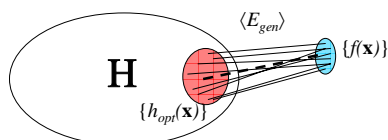


Our hypothesis space

The learning problem is realizable if  $f(\mathbf{x}) \in \mathbf{H}$ .

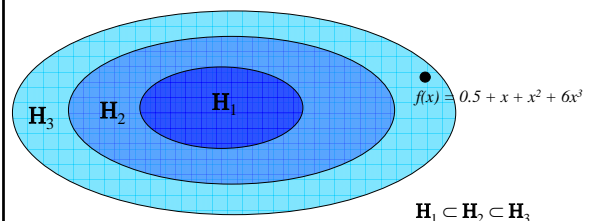
## The real inductive learning problem

Find appropriate hypothesis space  $\mathbf{H}$  and minimize the expected distance to  $f(\mathbf{x})$  ("generalization error")



Data is never noise free and never available in infinite amounts, so we get variation in data and model. The generalization error is a function of both the training data and the hypothesis selection method.

## Hypothesis spaces (examples)



$\mathbf{H}_1 = \{a+bx\}$ ;  $\mathbf{H}_2 = \{a+bx+cx^2\}$ ;  $\mathbf{H}_3 = \{a+bx+cx^2+dx^3\}$ ;  
Linear; Quadratic; Cubic;

## Learning problems

- The hypothesis takes as input a set of attributes  $\mathbf{x}$  and returns a "decision"  $h(\mathbf{x})$  = the predicted (estimated) output value for the input  $\mathbf{x}$ .
- Discrete valued function  $\Rightarrow$  classification
- Continuous valued function  $\Rightarrow$  regression

## Classification

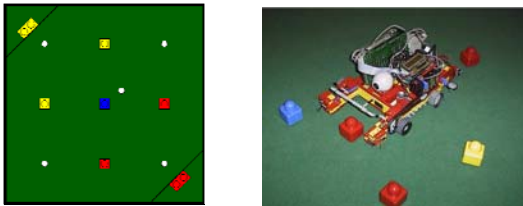
Order into one out of several classes

$$X^D \rightarrow C^K$$

Input space      Output (category) space

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix} \in X^D \quad \mathbf{c} = \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_K \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \in C^K$$

## Example: Robot color vision



Classify the Lego pieces into *red*, *blue*, and *yellow*.  
Classify *white* balls, *black* sideboard, and *green* carpet.  
Input = pixel in image, output = category

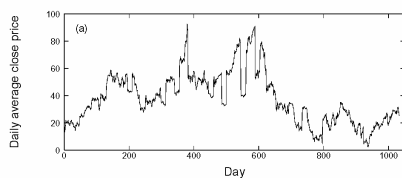
## Regression

The "fixed regressor model"

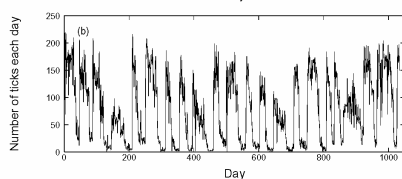
$$f(\mathbf{x}) = g(\mathbf{x}) + \varepsilon$$

$\mathbf{x}$	Observed input
$f(\mathbf{x})$	Observed output
$g(\mathbf{x})$	True underlying function
$\varepsilon$	I.I.D noise process with zero mean

## Example: Predict price for cotton futures

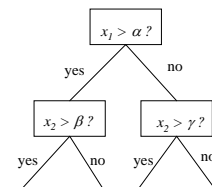


Input: Past history of closing prices, and trading volume  
Output: Predicted closing price

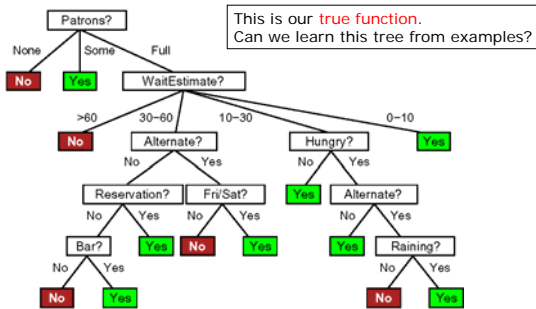


## Decision trees

- "Divide and conquer": Split data into smaller and smaller subsets.
- Splits usually on a single variable



## The wait@restaurant decision tree



## Inductive learning of decision tree

- **Simplest:** Construct a decision tree with one leaf for every example = memory based learning. Not very good generalization.

## Inductive learning of decision tree

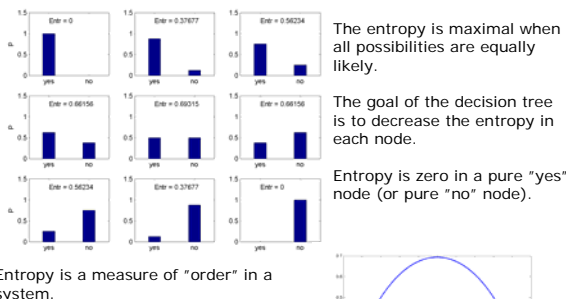
- **Simplest:** Construct a decision tree with one leaf for every example = memory based learning. Not very good generalization.
- **Advanced:** Split on each variable so that the purity of each split increases (i.e. either only yes or only no)

## Inductive learning of decision tree

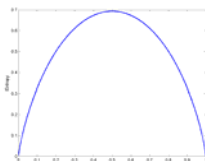
- **Simplest:** Construct a decision tree with one leaf for every example = memory based learning. Not very good generalization.
- **Advanced:** Split on each variable so that the purity of each split increases (i.e. either only yes or only no)
- Purity measured, e.g. with entropy

$$\text{Entropy} = -P(\text{yes}) \ln[P(\text{yes})] - P(\text{no}) \ln[P(\text{no})]$$

General form:  $\text{Entropy} = -\sum_i P(v_i) \ln[P(v_i)]$



The second law of thermodynamics: Elements in a closed system tend to seek their most probable distribution; in a closed system entropy always increases



## Decision tree learning algorithm

- Create pure nodes whenever possible
- If pure nodes are not possible, choose the split that leads to the largest decrease in entropy.

## Decision tree learning example

10 attributes:

- Alternate:** Is there a suitable alternative restaurant nearby? {yes,no}
- Bar:** Is there a bar to wait in? {yes,no}
- Fri/Sat:** Is it Friday or Saturday? {yes,no}
- Hungry:** Are you hungry? {yes,no}
- Patrons:** How many are seated in the restaurant? {none, some, full}
- Price:** Price level {\$, \$\$, \$\$\$}
- Raining:** Is it raining? {yes,no}
- Reservation:** Did you make a reservation? {yes,no}
- Type:** Type of food {French,Italian,Thai,Burger}
- Wait:** {0-10 min, 10-30 min, 30-60 min, >60 min}

## Decision tree learning example

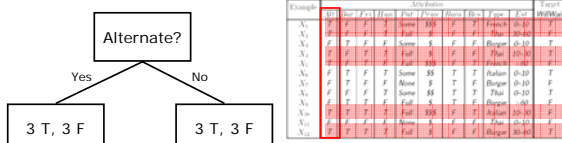
Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X <sub>1</sub>	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X <sub>2</sub>	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X <sub>3</sub>	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X <sub>4</sub>	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X <sub>5</sub>	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X <sub>6</sub>	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X <sub>7</sub>	F	T	F	F	None	\$	T	F	Burger	0-10	F
X <sub>8</sub>	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X <sub>9</sub>	F	T	T	F	Full	\$	T	F	Burger	>60	F
X <sub>10</sub>	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0-10	F
X <sub>12</sub>	T	T	T	T	Full	\$	F	F	Burger	30-60	T

T = True, F = False

$$\text{Entropy} = -\left(\frac{6}{12}\right)\ln\left(\frac{6}{12}\right) - \left(\frac{6}{12}\right)\ln\left(\frac{6}{12}\right) = 0.30$$

6 True, 6 False

## Decision tree learning example

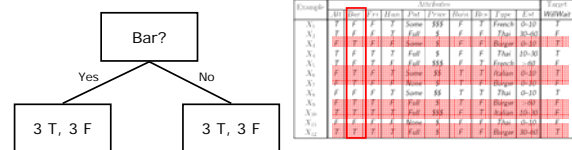


Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X <sub>1</sub>	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X <sub>2</sub>	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X <sub>3</sub>	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X <sub>4</sub>	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X <sub>5</sub>	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X <sub>6</sub>	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X <sub>7</sub>	F	T	F	F	None	\$	T	F	Burger	0-10	F
X <sub>8</sub>	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X <sub>9</sub>	F	T	T	F	Full	\$	T	F	Burger	>60	F
X <sub>10</sub>	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0-10	F
X <sub>12</sub>	T	T	T	T	Full	\$	F	F	Burger	30-60	T

$$\text{Entropy} = \frac{6}{12} \left[ -\left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) \right] + \frac{6}{12} \left[ -\left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) \right] = 0.30$$

$$\text{Entropy decrease} = 0.30 - 0.30 = 0$$

## Decision tree learning example

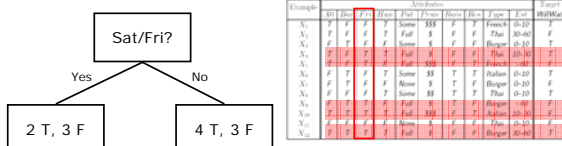


Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X <sub>1</sub>	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X <sub>2</sub>	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X <sub>3</sub>	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X <sub>4</sub>	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X <sub>5</sub>	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X <sub>6</sub>	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X <sub>7</sub>	F	T	F	F	None	\$	T	F	Burger	0-10	F
X <sub>8</sub>	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X <sub>9</sub>	F	T	T	F	Full	\$	T	F	Burger	>60	F
X <sub>10</sub>	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0-10	F
X <sub>12</sub>	T	T	T	T	Full	\$	F	F	Burger	30-60	T

$$\text{Entropy} = \frac{6}{12} \left[ -\left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) \right] + \frac{6}{12} \left[ -\left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right)\ln\left(\frac{3}{6}\right) \right] = 0.30$$

$$\text{Entropy decrease} = 0.30 - 0.30 = 0$$

## Decision tree learning example

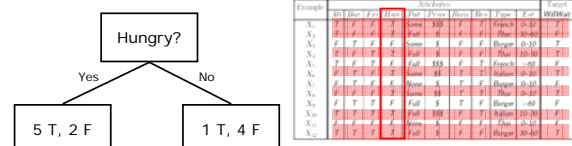


Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X <sub>1</sub>	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X <sub>2</sub>	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X <sub>3</sub>	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X <sub>4</sub>	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X <sub>5</sub>	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X <sub>6</sub>	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X <sub>7</sub>	F	T	F	F	None	\$	T	F	Burger	0-10	F
X <sub>8</sub>	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X <sub>9</sub>	F	T	T	F	Full	\$	T	F	Burger	>60	F
X <sub>10</sub>	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0-10	F
X <sub>12</sub>	T	T	T	T	Full	\$	F	F	Burger	30-60	T

$$\text{Entropy} = \frac{5}{12} \left[ -\left(\frac{2}{5}\right)\ln\left(\frac{2}{5}\right) - \left(\frac{3}{5}\right)\ln\left(\frac{3}{5}\right) \right] + \frac{7}{12} \left[ -\left(\frac{4}{7}\right)\ln\left(\frac{4}{7}\right) - \left(\frac{3}{7}\right)\ln\left(\frac{3}{7}\right) \right] = 0.29$$

$$\text{Entropy decrease} = 0.30 - 0.29 = 0.01$$

## Decision tree learning example



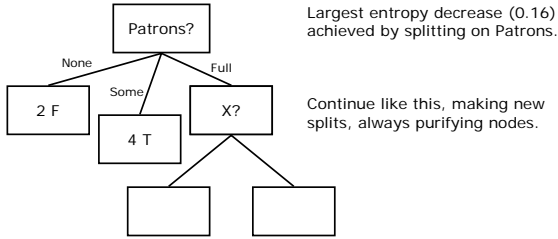
Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
X <sub>1</sub>	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X <sub>2</sub>	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X <sub>3</sub>	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X <sub>4</sub>	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X <sub>5</sub>	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X <sub>6</sub>	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X <sub>7</sub>	F	T	F	F	None	\$	T	F	Burger	0-10	F
X <sub>8</sub>	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X <sub>9</sub>	F	T	T	F	Full	\$	T	F	Burger	>60	F
X <sub>10</sub>	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0-10	F
X <sub>12</sub>	T	T	T	T	Full	\$	F	F	Burger	30-60	T

$$\text{Entropy} = \frac{7}{12} \left[ -\left(\frac{5}{7}\right)\ln\left(\frac{5}{7}\right) - \left(\frac{2}{7}\right)\ln\left(\frac{2}{7}\right) \right] + \frac{5}{12} \left[ -\left(\frac{1}{5}\right)\ln\left(\frac{1}{5}\right) - \left(\frac{4}{5}\right)\ln\left(\frac{4}{5}\right) \right] = 0.24$$

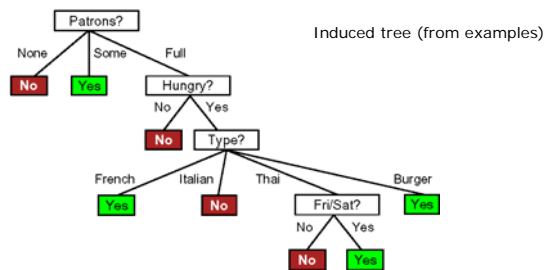
$$\text{Entropy decrease} = 0.30 - 0.24 = 0.06$$



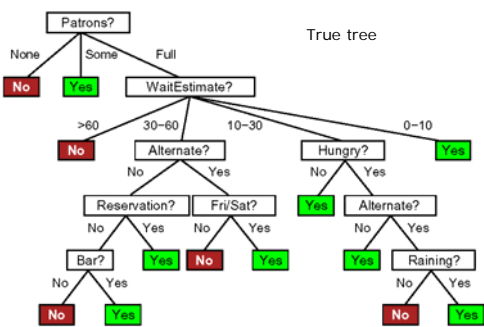
## Decision tree learning example



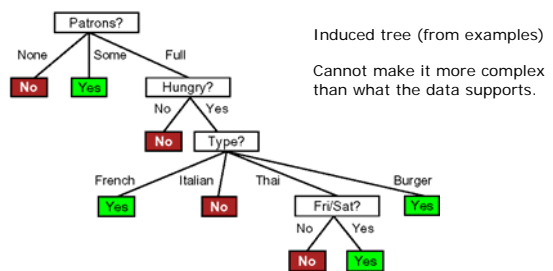
## Decision tree learning example



## Decision tree learning example



## Decision tree learning example



## How do we know it is correct?

How do we know that  $h \approx f$ ?

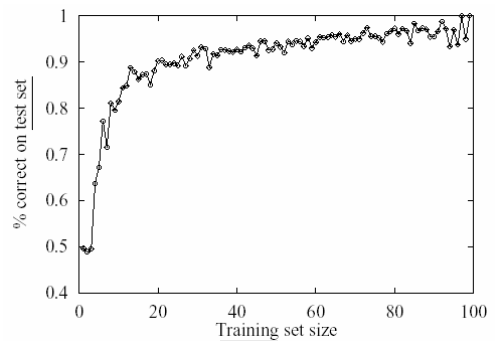
(Hume's Problem of Induction)

- Try  $h$  on a new **test set** of examples (cross validation)

...and assume the "principle of uniformity", i.e. the result we get on this test data should be indicative of results on future data. Causality is constant.

Inspired by a slide by V. Pavlovic

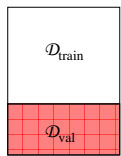
Learning curve for the decision tree algorithm on 100 randomly generated examples in the restaurant domain. The graph summarizes 20 trials.



## Cross-validation

Use a "validation set".

$$E_{gen} \approx E_{val}$$



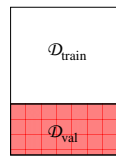
$E_{val}$

Split your data set into two parts, one for training your model and the other for validating your model. The error on the validation data is called "validation error" ( $E_{val}$ )

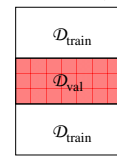
## K-Fold Cross-validation

More accurate than using only one validation set.

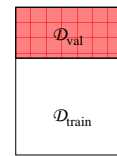
$$E_{gen} \approx \langle E_{val} \rangle = \frac{1}{K} \sum_{k=1}^K E_{val}(k)$$



$E_{val}(1)$



$E_{val}(2)$



$E_{val}(3)$

## PAC

- Any hypothesis that is consistent with a sufficiently large set of training (and test) examples is unlikely to be seriously wrong; it is **probably approximately correct (PAC)**.
- What is the relationship between the generalization error and the number of samples needed to achieve this generalization error?

## The error

$\mathbf{X}$  = the set of all possible examples (instance space).  
 $D$  = the distribution of these examples.  
 $\mathbf{H}$  = the hypothesis space ( $h \in \mathbf{H}$ ).  
 $N$  = the number of training data.

$$\text{error}(h) = P[h(\mathbf{x}) \neq f(\mathbf{x}) \mid \mathbf{x} \text{ drawn from } D]$$

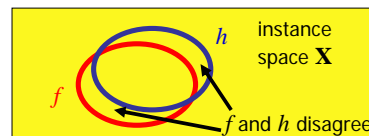


Image adapted from F. Hoffmann @ KTH

## Probability for bad hypothesis

Suppose we have a bad hypothesis  $h$  with  $\text{error}(h) > \epsilon$ . What is the probability that it is consistent with  $N$  samples?

- Probability for being inconsistent with one sample =  $\text{error}(h) > \epsilon$ .
- Probability for being consistent with one sample =  $1 - \text{error}(h) < 1 - \epsilon$ .
- Probability for being consistent with  $N$  independently drawn samples  $< (1 - \epsilon)^N$ .

## Probability for bad hypothesis

What is the probability that the set  $\mathbf{H}_{\text{bad}}$  of bad hypotheses with  $\text{error}(h) > \epsilon$  contains a consistent hypothesis?

$$P(h \text{ consistent} \wedge \text{error}(h) > \epsilon) \leq |\mathbf{H}_{\text{bad}}| (1 - \epsilon)^N \leq |\mathbf{H}| (1 - \epsilon)^N$$



## Probability for bad hypothesis

What is the probability that the set  $\mathbf{H}_{\text{bad}}$  of bad hypotheses with  $\text{error}(h) > \varepsilon$  contains a consistent hypothesis?

$$P(h \text{ consistent} \wedge \text{error}(h) > \varepsilon) \leq |\mathbf{H}_{\text{bad}}|(1-\varepsilon)^N \leq |\mathbf{H}|(1-\varepsilon)^N$$

If we want this to be less than some constant  $\delta$ , then

$$|\mathbf{H}|(1-\varepsilon)^N < \delta \Rightarrow \ln|\mathbf{H}| + N \ln(1-\varepsilon) < \ln \delta$$

## Probability for bad hypothesis

What is the probability that the set  $\mathbf{H}_{\text{bad}}$  of bad hypotheses with  $\text{error}(h) > \varepsilon$  contains a consistent hypothesis?

$$P(h \text{ consistent} \wedge \text{error}(h) > \varepsilon) \leq |\mathbf{H}_{\text{bad}}|(1-\varepsilon)^N \leq |\mathbf{H}|(1-\varepsilon)^N$$

If we want this to be less than some constant  $\delta$ , then

$$N > \frac{\ln(|\mathbf{H}|) - \ln(\delta)}{-\ln(1-\varepsilon)} \approx \frac{\ln(|\mathbf{H}|) - \ln(\delta)}{\varepsilon}$$

Don't expect to learn very well if  $\mathbf{H}$  is large

## How make learning work?

- Use simple hypotheses
  - Always start with the simple ones first
- Constrain  $\mathbf{H}$  with priors
  - Do we know something about the domain?
  - Do we have reasonable a priori beliefs on parameters?
- Use many observations
  - Easy to say...
- Cross-validation...