

# Machine Learning: Part I

UMaine COS 470/570 – Introduction to AI  
Spring 2019

Introduction

Why learn?  
Learning agents  
Kinds of learning?  
Supervised learning:  
Induction  
Unsupervised learning  
Reinforcement learning

“Classical” ML

# Introduction

Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Why learn?

# Why learn?

## Introduction

### Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

## “Classical” ML

- ▶ Why not static agent/knowledge?

# Why learn?

## Introduction

### Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

### “Classical” ML

- ▶ Why not static agent/knowledge?
- ▶ Changing world
- ▶ Different tasks
- ▶ Bad initial design/knowledge
- ▶ Incomplete/uncertain knowledge of world

Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Learning agents

# Agents and learning

## Introduction

Why learn?

Learning agents

Kinds of learning?

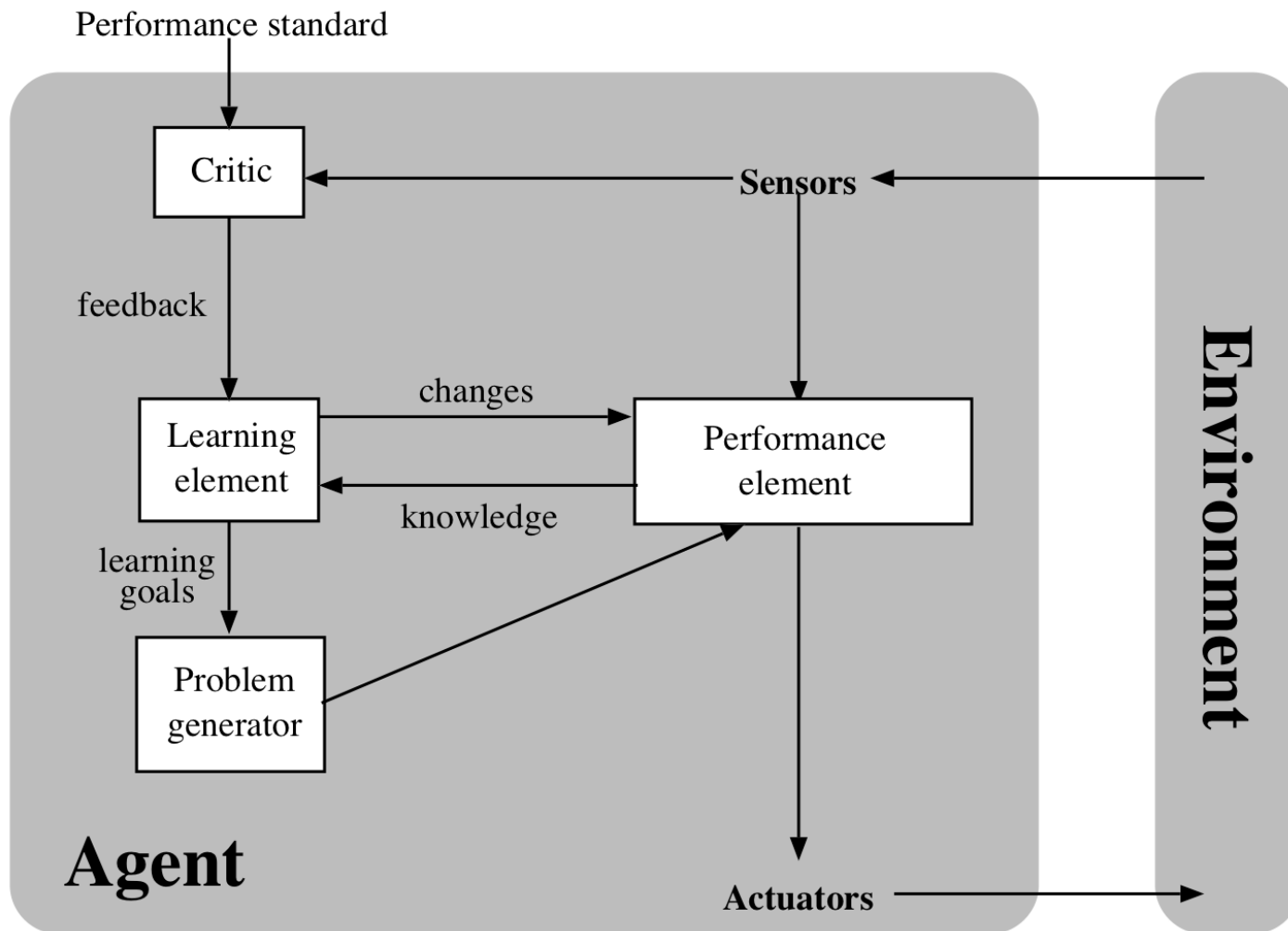
Supervised learning:

Induction

Unsupervised learning

Reinforcement learning

“Classical” ML



Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:

Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Kinds of learning?



# What are the kinds of learning?

## Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:

Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

- ▶ Feedback available: supervised vs unsupervised vs reinforcement
- ▶ What's learned: axioms, rules, plans, weights, . . .
- ▶ Symbolic vs sub-symbolic (neural, GA, LCS)
- ▶ Ensemble learning
- ▶ PAC learning

Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Supervised learning: Induction

# Supervised learning: Induction

Machine Learning:  
Part I

## Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

- ▶ Decision-tree learning
- ▶ Bayesian learning
- ▶ Decision-theoretic
- ▶ Neural network approaches

Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Unsupervised learning

# Unsupervised learning

Machine Learning:  
Part I

## Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:

Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

- ▶ Explanation-based learning
- ▶ Regression
- ▶ Induction: autoencoders, neural networks, support vector machines (SVMs)

Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:  
Induction

Unsupervised learning

Reinforcement learning

“Classical” ML

# Reinforcement learning

# Reinforcement learning

## Introduction

Why learn?

Learning agents

Kinds of learning?

Supervised learning:

Induction

Unsupervised learning

Reinforcement learning

## “Classical” ML

- ▶ Policy learning:  $P(\text{state}) \rightarrow \text{action}$
- ▶ Q-learning:  $Q(\text{action}, \text{state}) \rightarrow \text{value (utility)}$
- ▶ Case-based reasoning
- ▶ GAs, LCSs
- ▶ Deep RL

# “Classical” ML

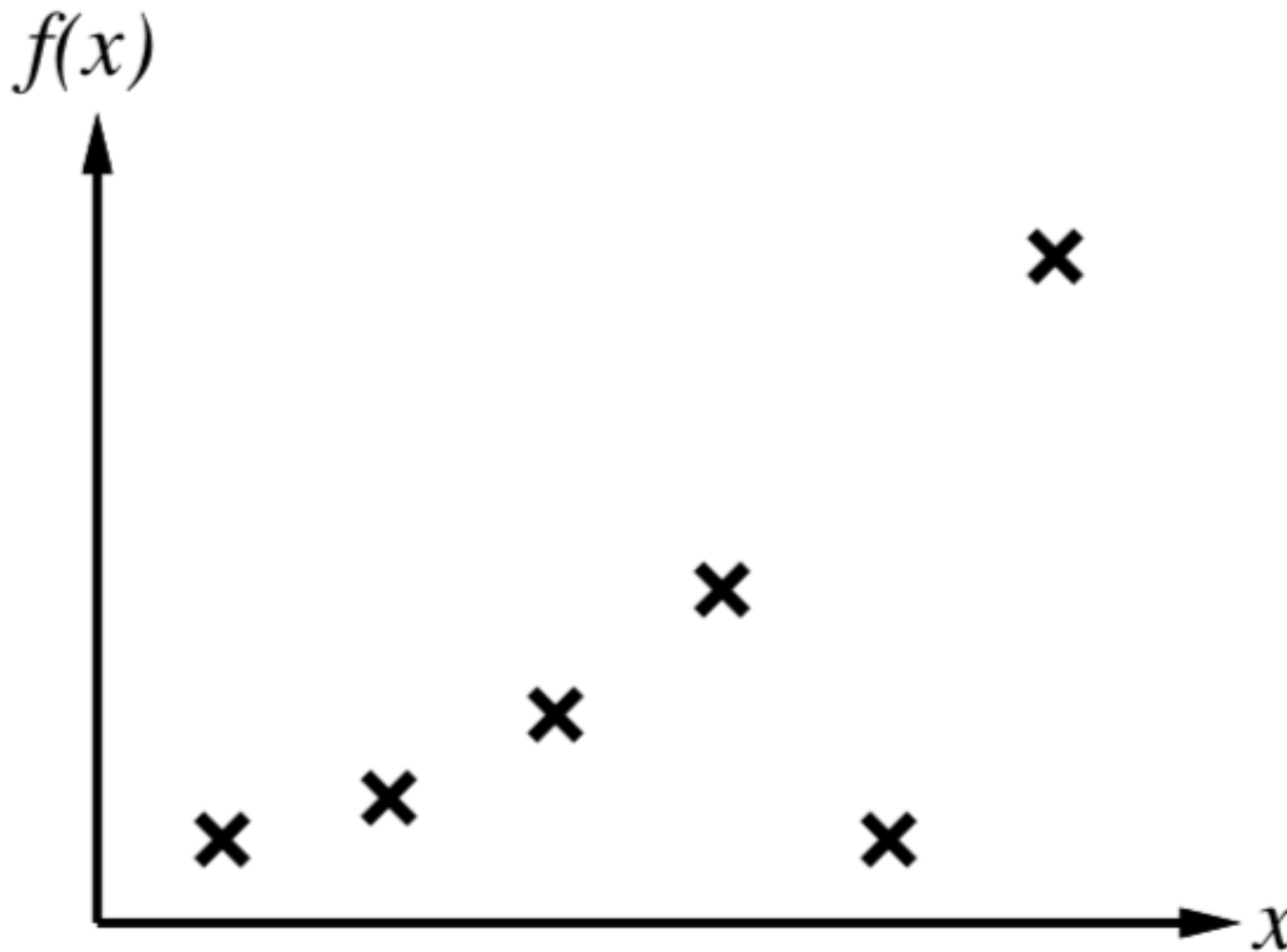


# Induction

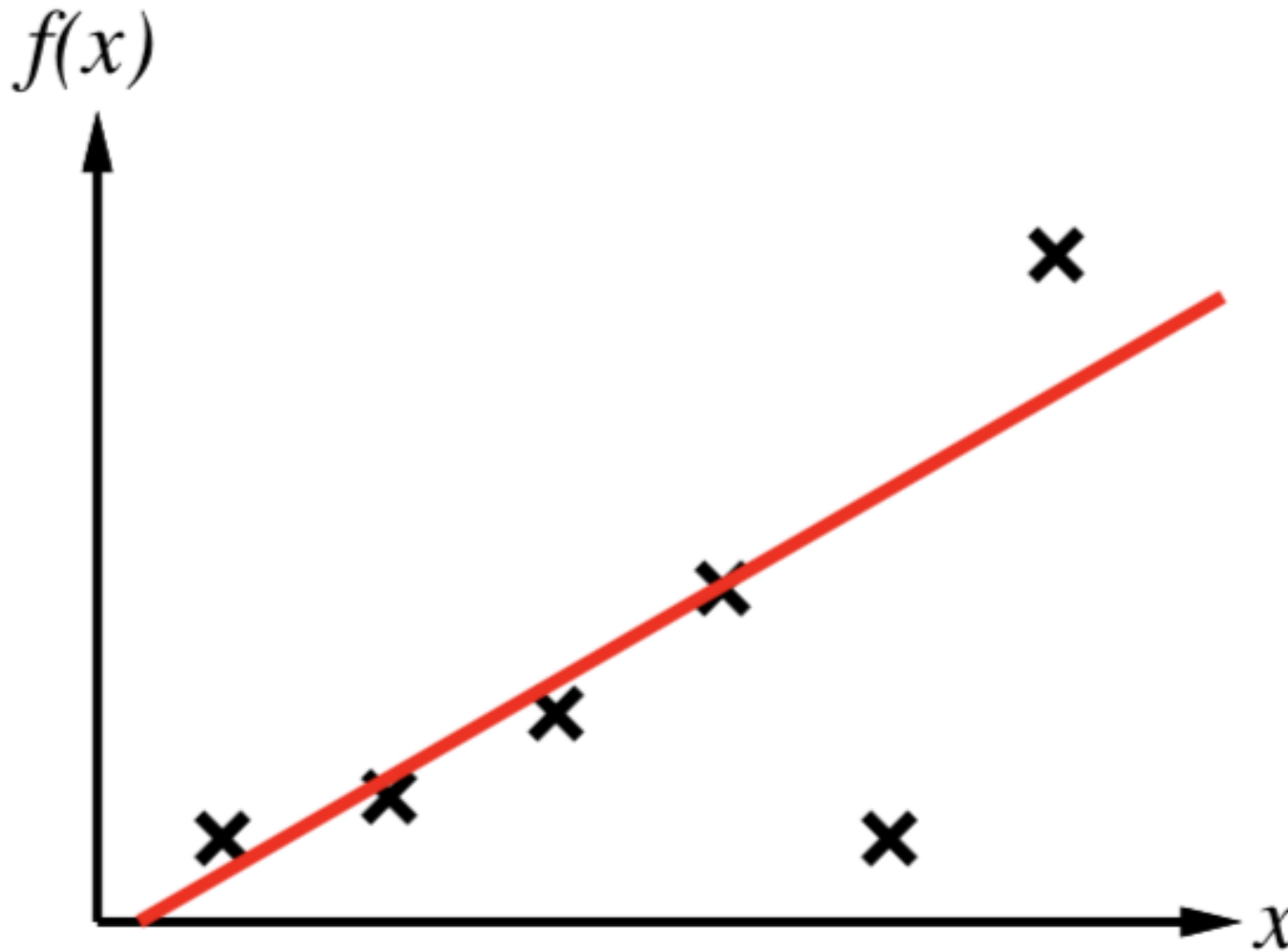
# Inductive learning in general

- ▶ Problem:
  - ▶ Given  $n$  examples of  $(x, f(x))$  pairs
  - ▶ Find  $f(y)$  for some new example  $y$
- ▶ Approach: Find  $h(x)$ , an approximation to  $f(x)$

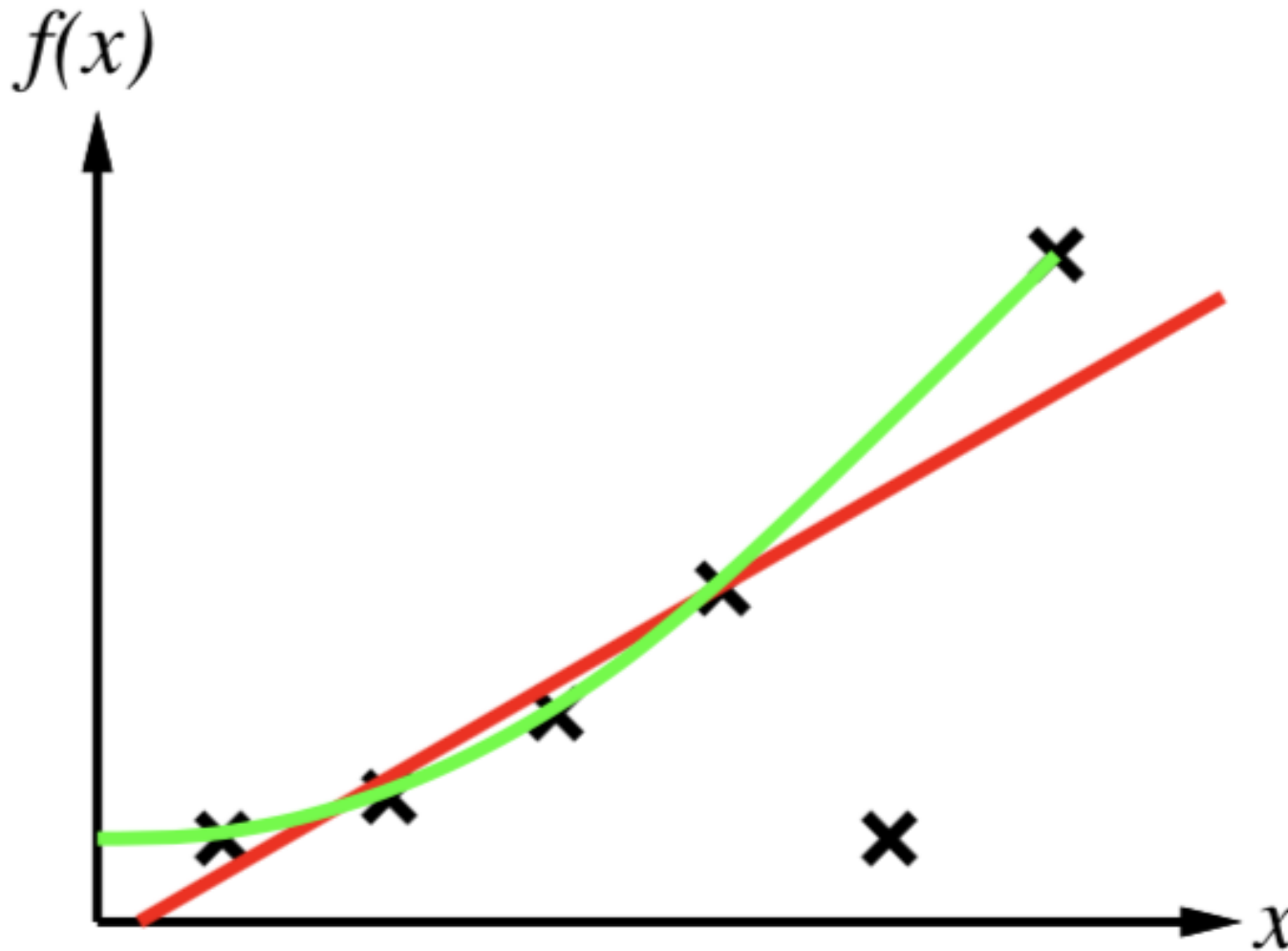
# Example: Regression (curve-fitting)



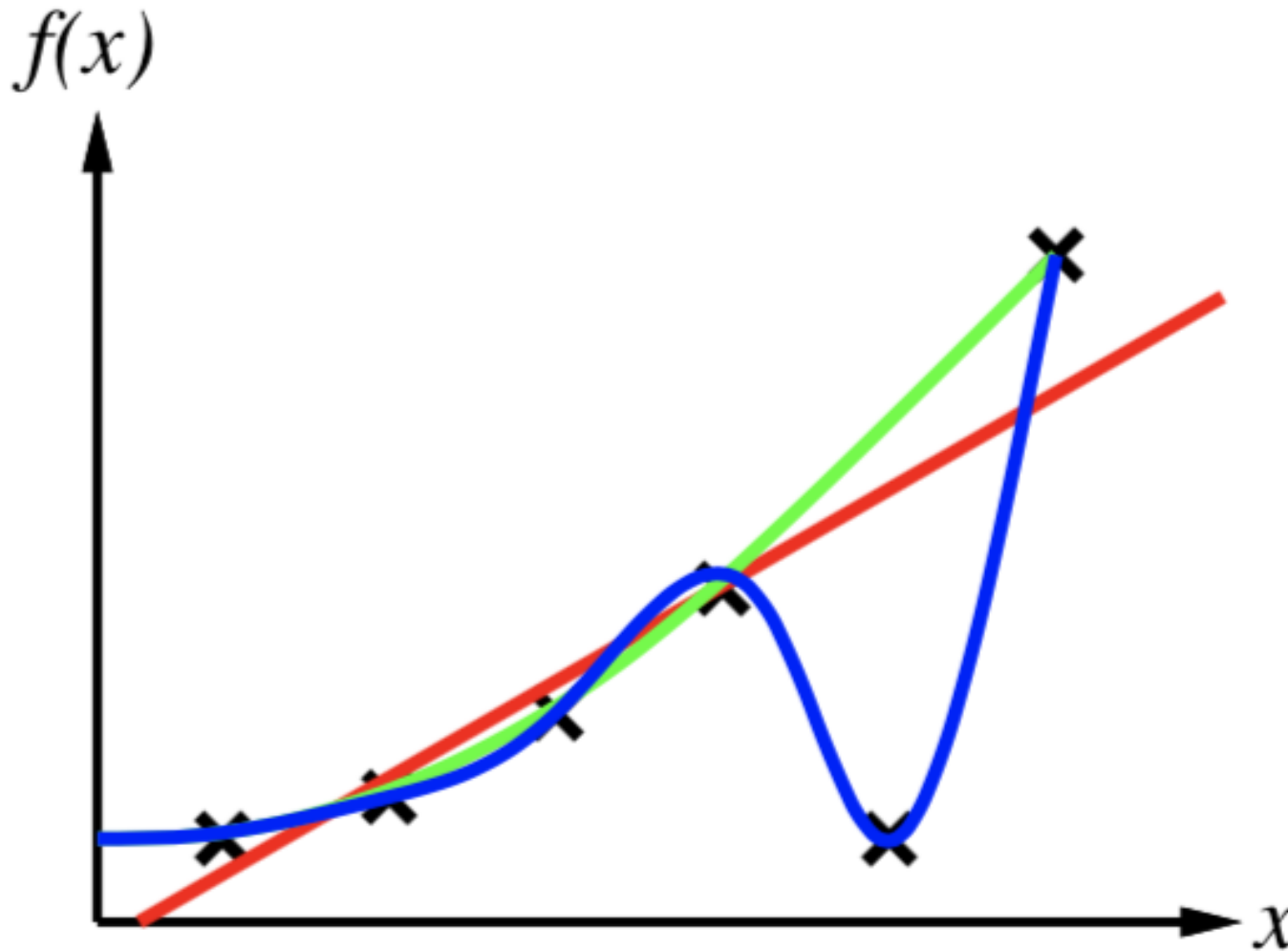
# Example: Regression (curve-fitting)



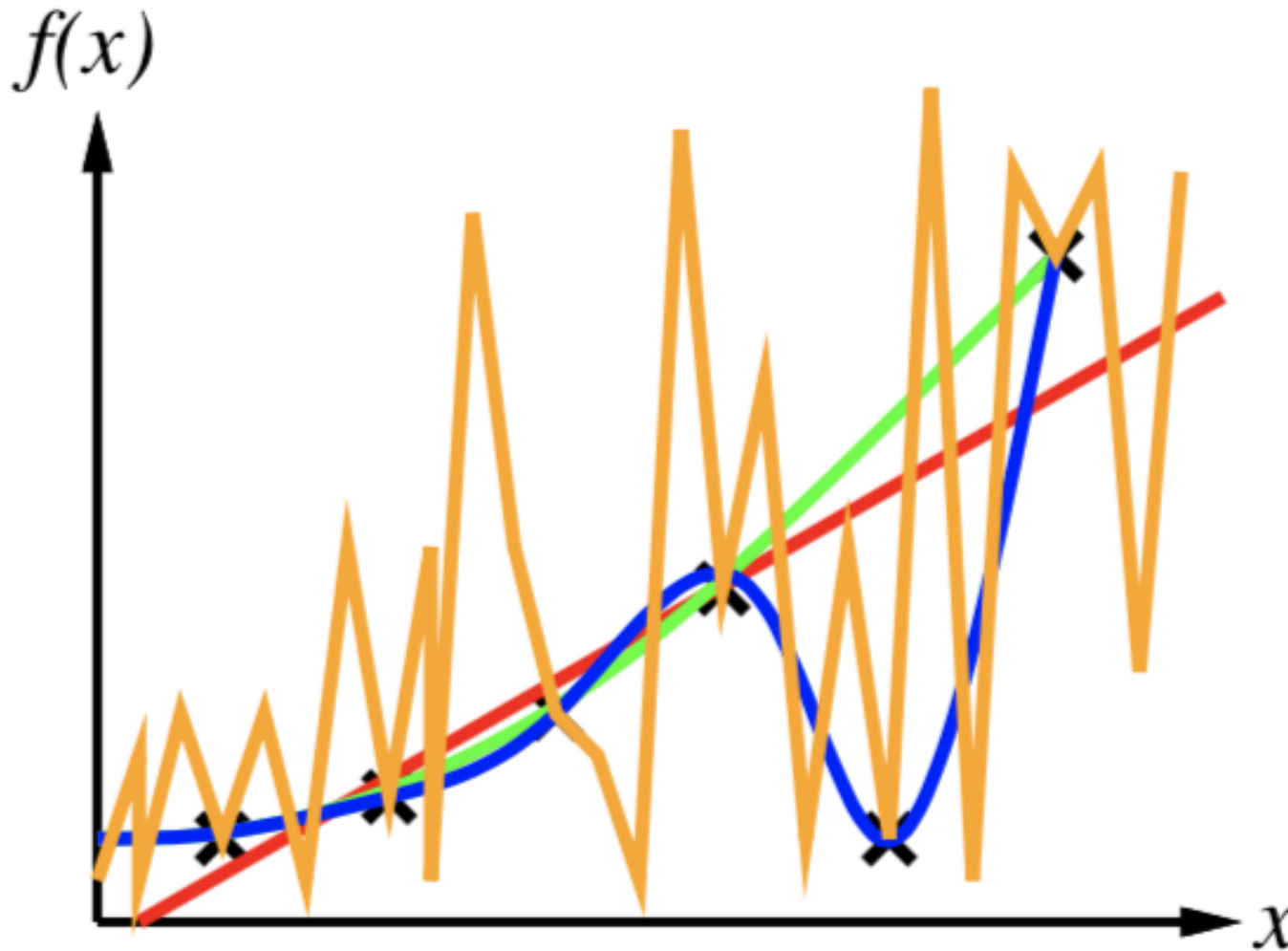
# Example: Regression (curve-fitting)



# Example: Regression (curve-fitting)



# Example: Regression (curve-fitting)

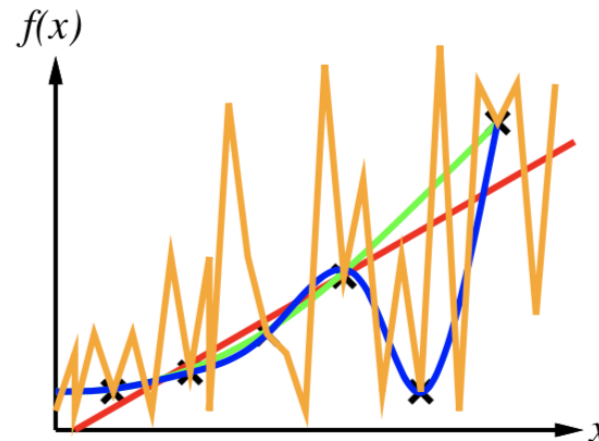


- ▶ *Hypothesis space* contains all the  $h$  you are considering
  - ▶ For regression: lines, polynomials, exponentials, ...
  - ▶ Other example: weight space of a particular neural network architecture
- ▶ Selection critical – needs to contain  $f(x)$  or contain good-enough approximation(s)
- ▶ *Consistent hypothesis*: agrees with all the data (to some  $\epsilon$ , perhaps)
- ▶ Maybe infinite # of consistent hypotheses
- ▶ Use Occam's (Ockham's) razor



# Hypothesis space

- ▶ *Hypothesis space* contains all the  $h$  you are considering
  - ▶ For regression: lines, polynomials, exponentials, ...
  - ▶ Other example: weight space of a particular neural network architecture
- ▶ Selection critical – needs to contain  $f(x)$  or contain good-enough approximation(s)
- ▶ *Consistent hypothesis*: agrees with all the data (to some  $\epsilon$ , perhaps)
- ▶ Maybe infinite # of consistent hypotheses
- ▶ Use Occam's (Ockham's) razor



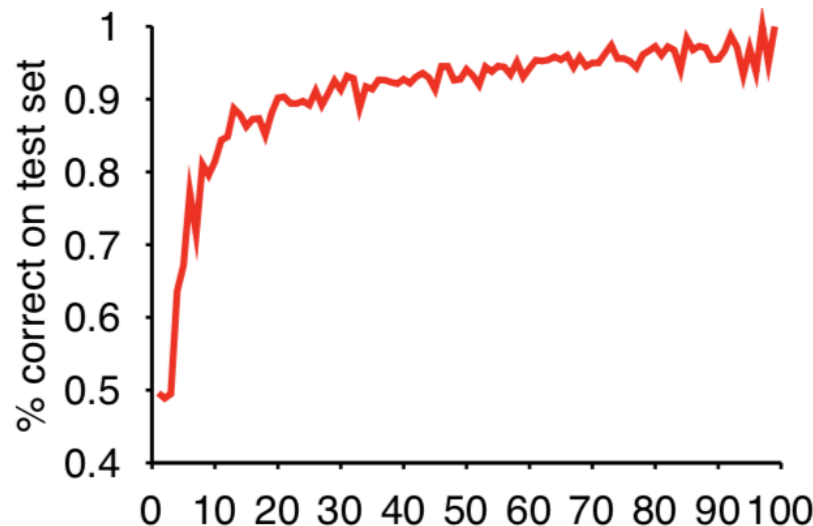
# What if there's no consistent function?

- ▶ Wrong hypothesis space (problem is *unrealizable* in space)
- ▶  $f(x)$  not deterministic
- ▶ Measurement errors for  $(x, f(x))$  pairs

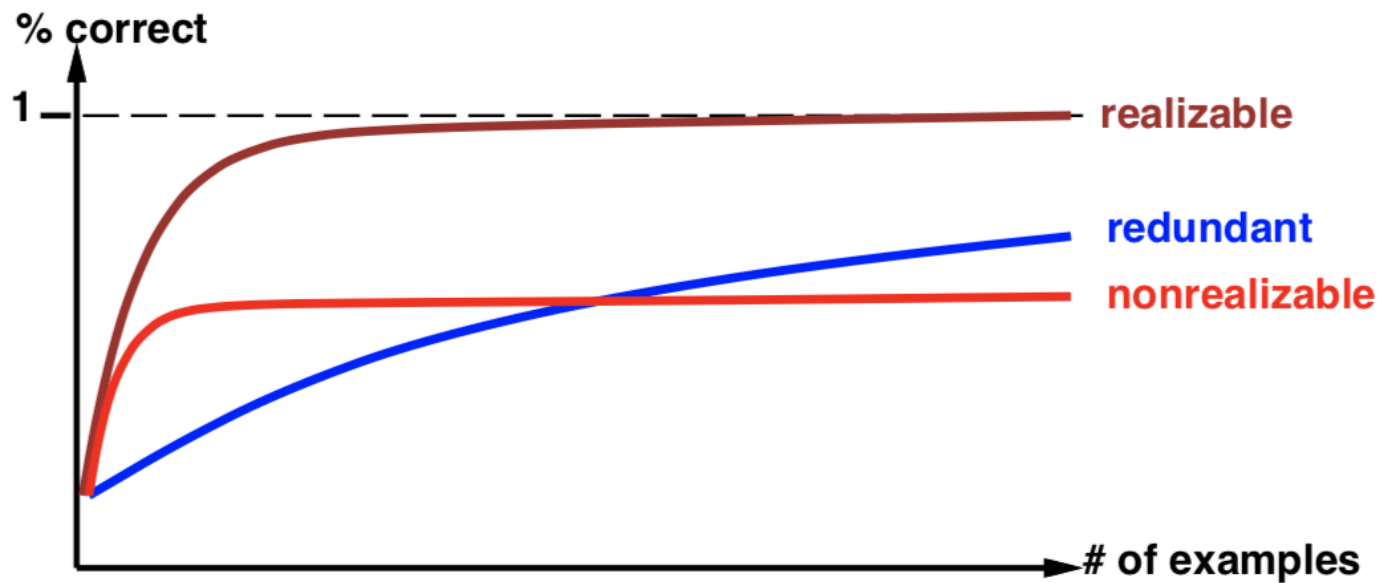
- ▶ How do we know  $h$  sufficiently approximates  $f$  (Hume: *Problem of Induction*)
- ▶ Usual approach:
  - ▶ Try  $h$  on some new test set  $T$  of examples
  - ▶  $T$  should have same distribution over example space as training set
  - ▶ Often: break initial example set into training, test subsets
- ▶ Estimate accuracy directly or create *learning curve*
- ▶ Learning curve = %correct as function of training set size

# Performance measurement

- ▶ How do we know  $h$  sufficiently approximates  $f$  (Hume: *Problem of Induction*)
- ▶ Usual approach:
  - ▶ Try  $h$  on some new test set  $T$  of examples
  - ▶  $T$  should have same distribution over example space as training set
  - ▶ Often: break initial example set into training, test subsets
- ▶ Estimate accuracy directly or create *learning curve*
- ▶ Learning curve = %correct as function of training set size



# Learning curves



# Decision tree learning

# Classification from examples

Machine Learning:  
Part I

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Problem: given a set of properties, find or answer to question
- ▶ Given set of exemplars: sets of features  $\rightarrow$  answers
- ▶ Supervised learning

- ▶ Identify a feature that cleanly divides the exemplars into subsets
- ▶ For each subset, do it again
- ▶ Continue until each subset contain indistinguishable exemplars
- ▶ Each split  $\equiv$  question about new thing to categorize
- ▶ E.g., something like Akinator
- ▶ Result is a *decision tree*
- ▶ Everyday examples: scientific species key, repair manual, etc.



# Example: User knowledge

## Introduction

### "Classical" ML

Induction

Decision tree learning

Explanation-based learning

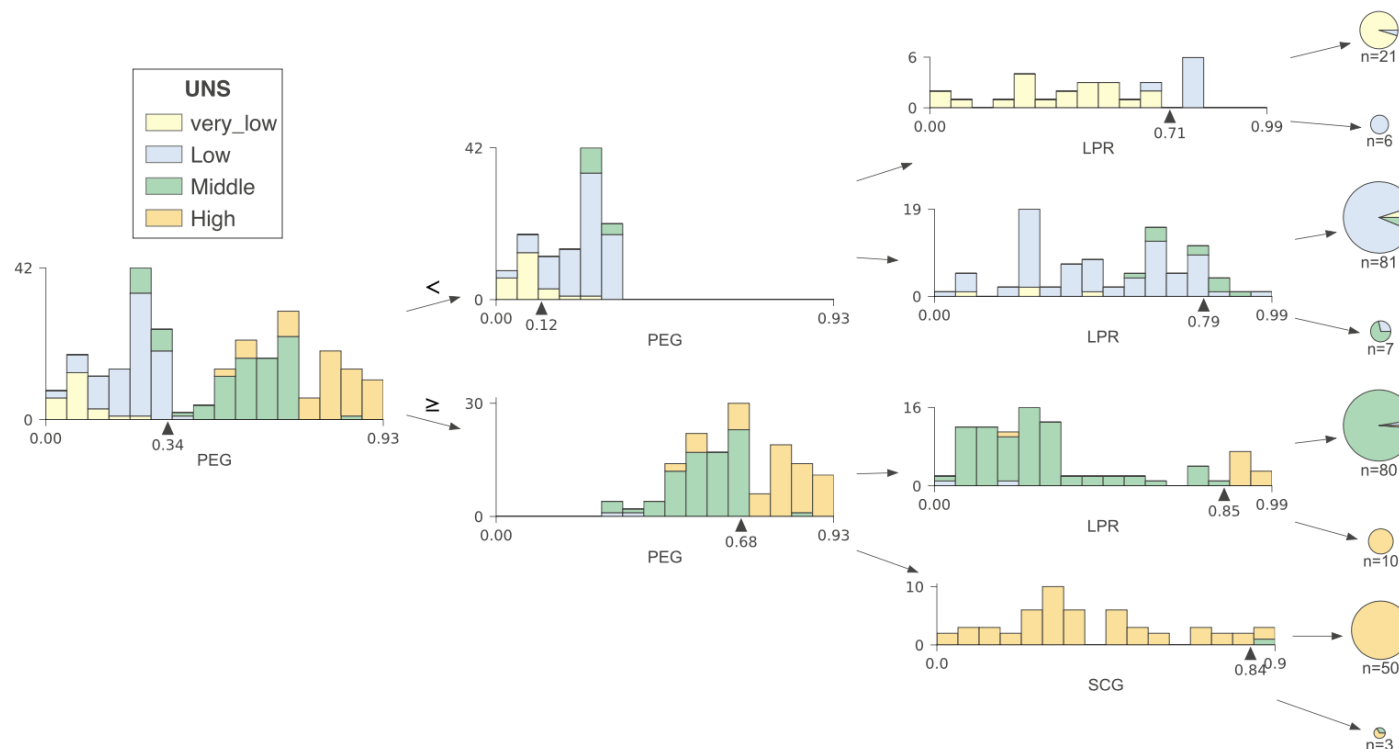
Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

## ► Classify user knowledge in some domain (UNS)



(STG = study time for goal objects, SCG = repetitions, STR = study time for related objects, LPR = exam performance for related objects, PEG = exam performance for goal objects)

(from UCI machine learning repository, explained.ai)

# Example: Breast cancer

## Introduction

### “Classical” ML

Induction

Decision tree learning

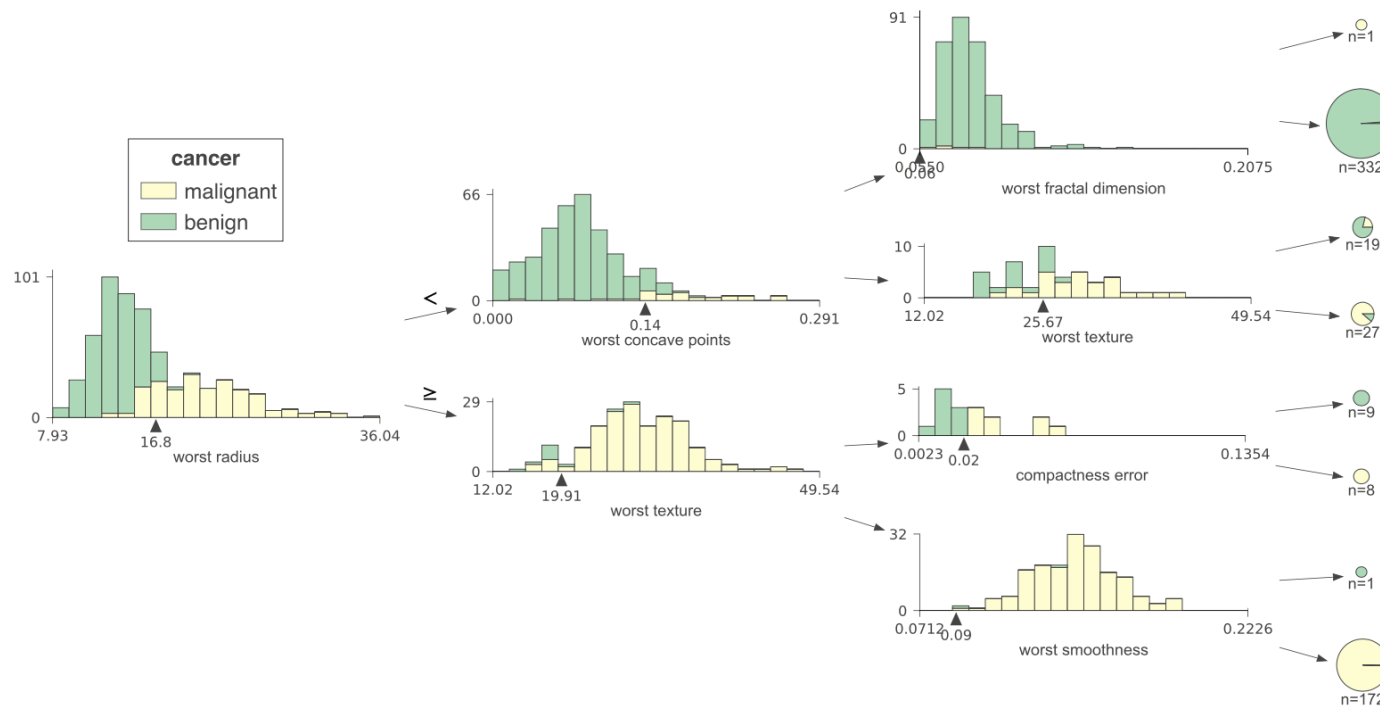
Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning



(from UCI ML Breast Cancer Wisconsin dataset, explained.ai)

# Learning a decision tree

- ▶ Suppose we have a data set described by *attribute values*
- ▶ E.g., to answer: “Should I wait for a table?” (T/F categorization)
- ▶ Attributes:

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry.
5. *Patrons*: how many people are in the restaurant (values are *None*, *Some*, and *Full*).
6. *Price*: the restaurant’s price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (French, Italian, Thai, or burger).
10. *WaitEstimate*: the wait estimated by the host (0–10 minutes, 10–30, 30–60, or >60).

# Learning a decision tree

## ► Have examples:

Example	Attributes										Target
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$X_1$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0–10</i>	<i>T</i>
$X_2$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30–60</i>	<i>F</i>
$X_3$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0–10</i>	<i>T</i>
$X_4$	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10–30</i>	<i>T</i>
$X_5$	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>&gt;60</i>	<i>F</i>
$X_6$	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0–10</i>	<i>T</i>
$X_7$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0–10</i>	<i>F</i>
$X_8$	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0–10</i>	<i>T</i>
$X_9$	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>&gt;60</i>	<i>F</i>
$X_{10}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10–30</i>	<i>F</i>
$X_{11}$	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0–10</i>	<i>F</i>
$X_{12}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30–60</i>	<i>T</i>

## Introduction

## “Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

# Possible tree

## Introduction

### “Classical” ML

Induction

Decision tree learning

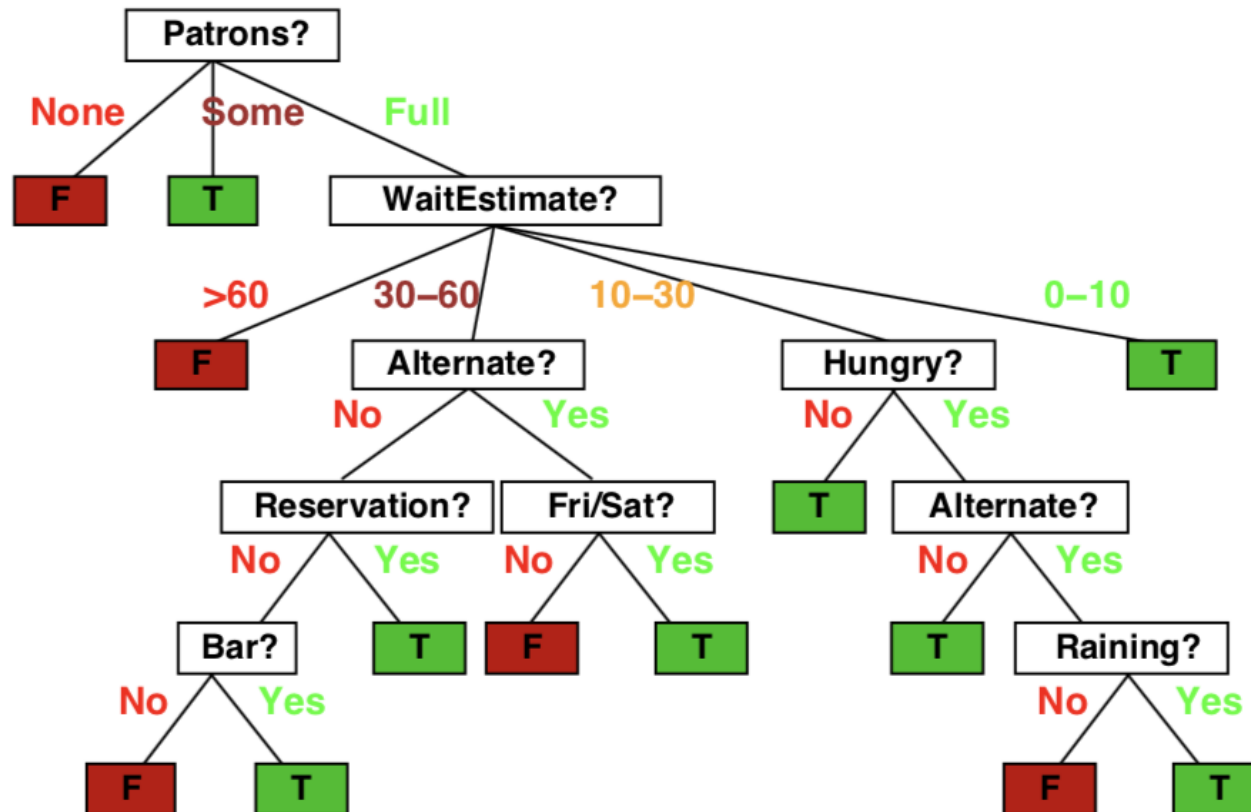
Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning



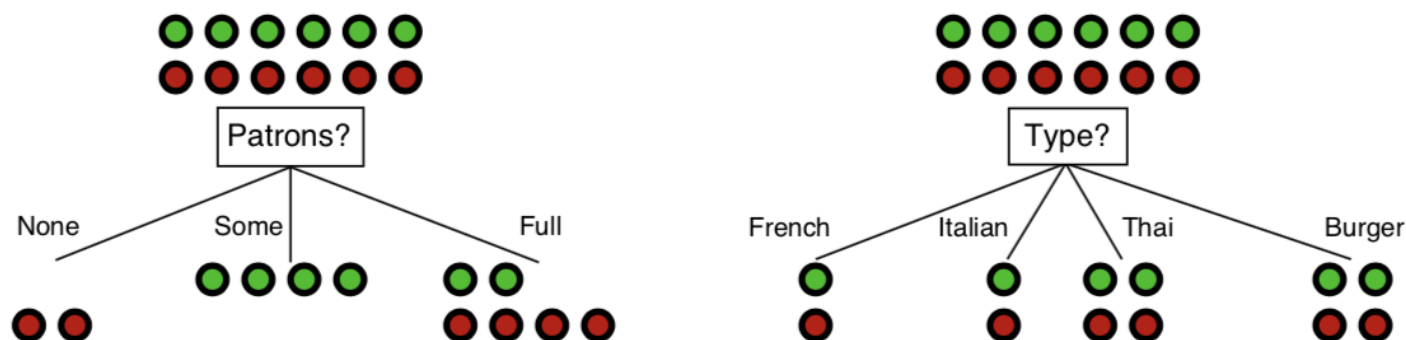
- ▶ Decision trees: can express any function of the attributes
- ▶ E.g., equiv. to binary satisfiability for Boolean functions
  - ▶ Think of a truth table
  - ▶ Internal nodes = input variables
- ▶ Can find a consistent tree for any training set
  - ▶ But trivially, won't generalize –
  - ▶ Want a more compact tree
  - ▶ Decide on attribute to include next based on amount of *information* they provide

# Decision tree learning algorithm

```
function DTL(examples, attributes, default) returns a decision tree
  if examples is empty then return default
  else if all examples have the same classification then return the classification
  else if attributes is empty then return MODE(examples)
  else
    best ← CHOOSE-ATTRIBUTE(attributes, examples)
    tree ← a new decision tree with root test best
    for each value  $v_i$  of best do
       $examples_i \leftarrow \{\text{elements of } examples \text{ with } best = v_i\}$ 
      subtree ← DTL(examplesi, attributes − best, MODE(examples))
      add a branch to tree with label  $v_i$  and subtree subtree
  return tree
```

# Which attribute to choose?

- ▶ Good attribute splits into (ideally) all pos or all neg examples



*Patrons?* is a better choice—gives **information** about the classification



# Which attribute to choose?

- ▶ Consider attributes as questions to be asked
- ▶ Choose one that gives more information when answered
- ▶ Less I initially know, more information answer provides
- ▶ Measure information in *bits*
- ▶ Need 1 bit to answer yes/no question with prior prob. = <0.5, 0.5>
- ▶ If prior = < $P_1, P_2, \dots, P_n$ >, information is the *entropy* of the prior:

$$H(< P_1, P_2, \dots, P_n >) = \sum_{i=1}^n -P_i \log_2 P_i$$

# Which attribute to choose?

- ▶ Suppose we have  $p$  pos and  $n$  neg examples
- ▶  $H(< p/(p + n), n/(p + n) >)$  bits needed to classify new example
- ▶ Each attribute  $i$  splits examples  $E$  into subgroup  $E_i$
- ▶ Hopefully, each  $E_i$  needs less information to complete than initial problem
- ▶ If  $E_i$  has  $p_i$  pos and  $n_i$  neg examples, then need:

$$H(< p_i/(p_i + n_i), n_i/(p_i + n_i) >)$$

bits to classify

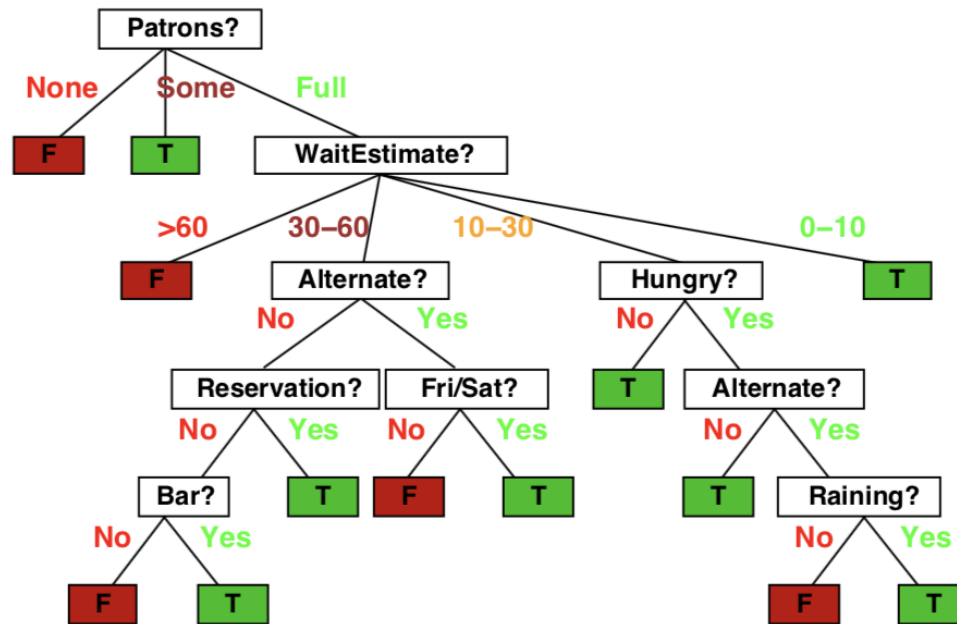
- ▶ *Expected* bits per example over all branches

$$H_E = \sum_i \frac{p_i + n_i}{p + n} H(< p_i/(p_i + n_i), n_i/(p_i + n_i) >)$$

- ▶ Prev example:  $H_E(\text{Patrons?}) = 0.459$  bits,  $H_E(\text{Type}) = 1$  bit
- ▶  $\therefore$  choose *Patrons?*

# Result of using information gain

Russell's



Introduction

"Classical" ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

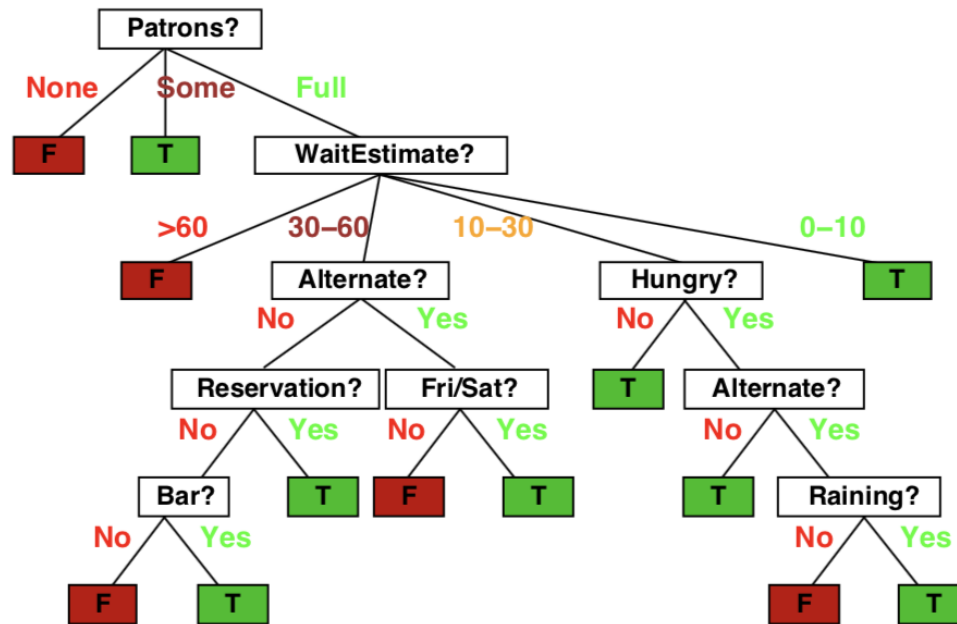
Genetic algorithms

Case-based reasoning

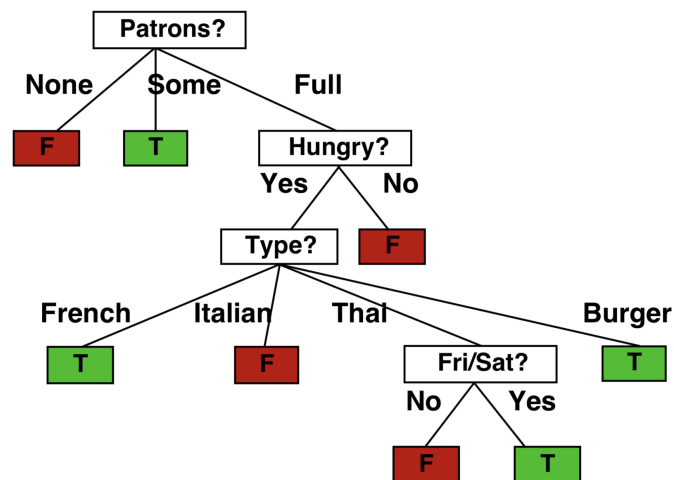
Schema-Based Reasoning

# Result of using information gain

Russell's



Learned tree:



Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

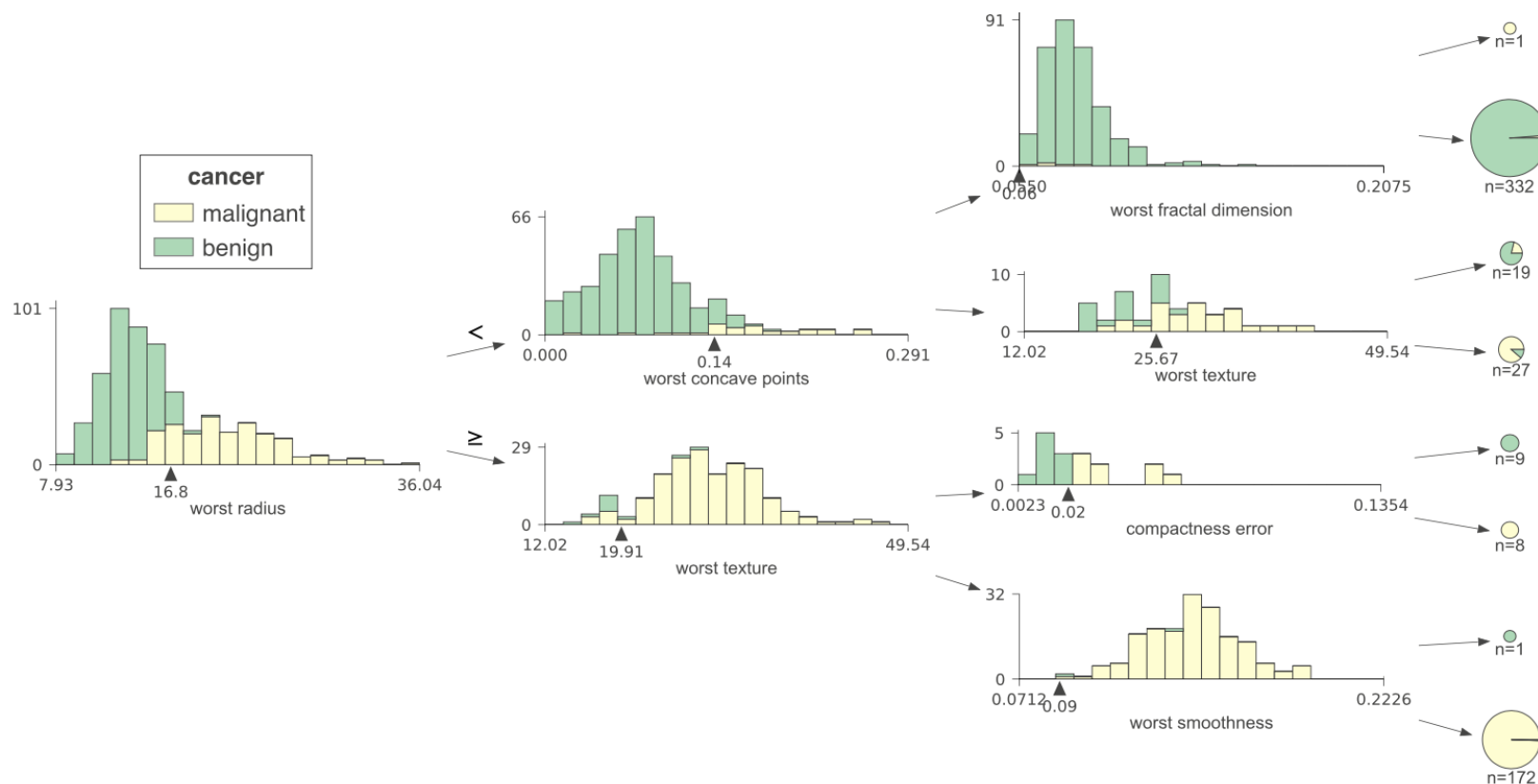
Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

# Drawbacks

- ▶ Leaf nodes with multiple examples may not agree on classification



## Introduction

### "Classical" ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Overfitting
  - ▶ Even if critical info is missing, DT alg. will create *some* DTree
  - ▶ Can make spurious distinctions based on irrelevant attribs
    - ▶ E.g., predict roll of die
    - ▶ Attributes given: hair color, color of die, day of the week, etc.
    - ▶ If no duplicate descriptions for different outcomes  $\Rightarrow$  exact hypothesis
  - ▶ As number of attribute goes up, so does likelihood of overfitting

- ▶ Decision tree learning space
  - ▶ With  $n$  Boolean attributes, tree is equiv. to truth table
  - ▶ So  $2^n$  rows
  - ▶ Space of all possible  $n$ -attribute decision trees = number of  $n$ -variable Boolean functions
  - ▶ There are  $2^n$  rows, same as the answer column
  - ▶ Thus there are  $2^{2^n}$  decision trees
  - ▶ Russell & Norvig: “scary number”

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning



► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$
5	$2^{32} = 4,294,967,296$

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$
5	$2^{32} = 4,294,967,296$
6	$2^{64} = 18,446,744,073,709,551,616$

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$
5	$2^{32} = 4,294,967,296$
6	$2^{64} = 18,446,744,073,709,551,616$
7	$\sim 3.4E38$

► Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$
5	$2^{32} = 4,294,967,296$
6	$2^{64} = 18,446,744,073,709,551,616$
7	$\sim 3.4E38$
8	$\sim 1.2E77$

- ▶ Decision tree learning space (cont'd)

#attributes	# poss. decision trees
1	$2^2 = 4$
2	$2^4 = 16$
3	$2^8 = 256$
4	$2^{16} = 65,536$
5	$2^{32} = 4,294,967,296$
6	$2^{64} = 18,446,744,073,709,551,616$
7	$\sim 3.4E38$
8	$\sim 1.2E77$

- ▶ R&N: “We will need some ingenious algorithms to find consistent hypotheses in such a large space.”



# Recap/Animation

Machine Learning:  
Part I

From: A visual introduction to machine learning

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

# Explanation-based learning

# Explanation-based learning

Machine Learning:  
Part I

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Explanation-based learning (EBL): use prior domain knowledge to learn new facts
- ▶ Traditionally (Mooney, DeJong, others) – no statistical inference, no induction
- ▶ Given a domain theory  $\Delta$  and some example  $X$ :
  - ▶ Use  $\Delta$  to account for  $X$  – e.g., show that  $\Delta \models X$
  - ▶ Then convert reasoning chain to a new rule

# EBL example

- ▶ Suppose we know people want money, people value things they love, people will pay to have things they value, . . .
- ▶ We learn that Johnny took Peter, Peter is Sally's husband, and Sally paid Johnny \$100,000
- ▶ Explain this by such things as: Sally loves Peter, Sally values Peter, Sally pays Johnny to have Peter, etc.
- ▶ Ultimately, compress into rule or schema for kidnapping

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Can be seen more as *operationalizing* existing knowledge
- ▶ Problem: no real new knowledge

# EBL example

*(From R&N)*

*Rewrite( $u, v$ )  $\wedge$  Simplify( $v, w$ )  $\Rightarrow$  Simplify( $u, w$ )*

*Primitive( $u$ )  $\Rightarrow$  Simplify( $u, u$ )*

*ArithmeticUnknown( $u$ )  $\Rightarrow$  Primitive( $u$ )*

*Number( $u$ )  $\Rightarrow$  Primitive( $u$ )*

*Rewrite( $1 \times u, u$ ) Rewrite( $0_u, u$ )*

*...*

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

# EBL example

(From R&N)

$Rewrite(u, v) \wedge Simplify(v, w) \Rightarrow Simplify(u, w)$

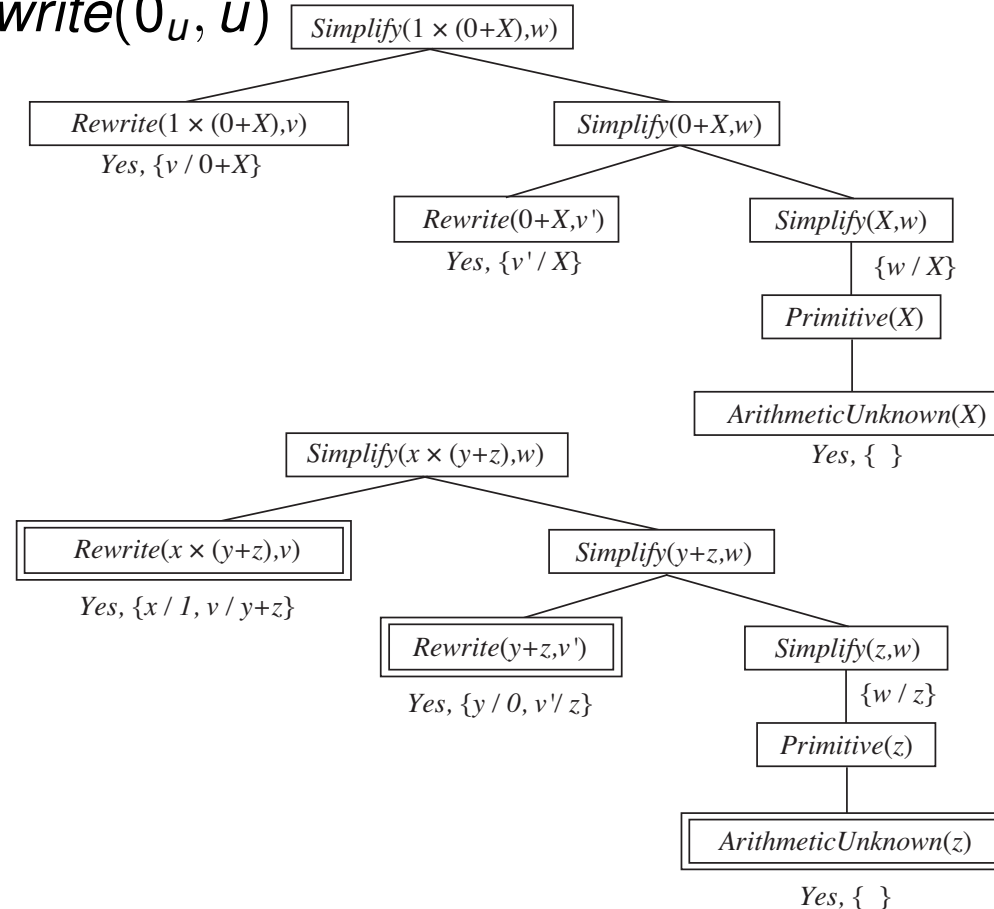
$Primitive(u) \Rightarrow Simplify(u, u)$

$ArithmeticUnknown(u) \Rightarrow Primitive(u)$

$Number(u) \Rightarrow Primitive(u)$

$Rewrite(1 \times u, u)$   $Rewrite(0 + u, u)$

...



Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Can read bottom tree's leaves to extract the (generalized) rule:

$$\text{Rewrite}(1 \times (0 + z), 0 + z) \wedge \text{Rewrite}(0 + z, z) \wedge \\ \text{ArithmeticUnknown}(z) \Rightarrow \text{Simplify}(1 \times (0 + z), z)$$

- ▶ Two rightmost terms: true no matter what  $z$  is
- ▶  $\text{ArithmeticUnknown}(z)$  – not every  $z$  is an arithmetic unknown, so can't drop it
- ▶ Resulting rule:

$$\text{ArithmeticUnknown}(z) \Rightarrow \text{Simplify}(1 \times (0 + z), z)$$



# Learning, expertise, and soundness

*(See DeJong's (2005) explanation-based learning paper, distributed on Slack)*

- ▶ Problem:
  - ▶ Statistical learning (induction) :
    - ▶ unsound
    - ▶ can't readily make use of expert's knowledge
  - ▶ EBL:
    - ▶ sound
    - ▶ can easily incorporate expert's knowledge
    - ▶ but only operationalizes the knowledge
- ▶ DeJong: use EBL to combine expert knowledge with (some) statistical learning

Introduction

"Classical" ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

# Explanation-based learning and expertise

- ▶ Obtain knowledge from expert, encode in representation
- ▶ Problem: Expert knowledge approximate, conflicting
- ▶ Problem: If represented in something like FOPC, doesn't apply to real world:

$$\forall x \text{Bird}(x) \Rightarrow \text{Flies}(x)$$

Statement not so much false as *incomplete*

- ▶ Almost *all* logical statements are incomplete with respect to the real world
  - ▶ Only sound with respect to some axiom set  $\Delta$
  - ▶ In general, subject to the *qualification problem*:
    - ▶ Most universally quantified sentences will have to include an infinite number of qualifications if they are to be interpreted as accurate statements about the world.
- ▶ E.g.:  $\forall x \text{Bird}(x) \wedge \neg \text{Dead}(x) \wedge \dots \Rightarrow \text{Flies}(x)$

# Explanation-based learning and expertise

- ▶ Proposal: accept expert's knowledge  $\Delta$  as correct, but not believed
- ▶ Belief only comes from seeing examples that  $\Delta$  can explain
- ▶ Learn new rules based on  $\Delta$  and experience
- ▶ E.g.:
  - ▶  $\Delta$ : birds fly, flying involves volition, dead things have no volition, etc.
  - ▶ New example: cooked chicken – refutes birds fly
  - ▶ Rule still considered correct (from expert)
  - ▶ Explain not flying by proof relying on the chicken being dead
  - ▶ Hypothesize:  $\forall x \text{Dead}(x) \Rightarrow \neg \text{Flies}(x)$
  - ▶ Look for more examples to confirm this
  - ▶ Confirmed to acceptable level of belief: create new rule:

$$\forall x \text{Bird}(x) \wedge \neg \text{Dead}(x) \Rightarrow \text{Flies}(x)$$

# How to derive the new rule?

- ▶ New inference rule – AND introduction

$$\frac{\begin{array}{l} \rho_1 \Rightarrow \psi \\ \rho_2 \Rightarrow \psi \end{array}}{(\rho_1 \wedge \rho_2) \Rightarrow \psi}$$

- ▶ Not logically necessary – but if  $\exists$  weak evidence for  $\rho_1$  and  $\rho_2$ , together might be sufficient to conclude  $\psi$ .
- ▶ Second rule:

$$\frac{\begin{array}{l} \rho_1 \Rightarrow \psi \\ \rho_2 \Rightarrow \neg\psi \end{array}}{(\rho_1 \wedge \neg\rho_2) \Rightarrow \psi}$$

- ▶ Again, logically redundant
- ▶ However, very useful for conditioning knowledge for real world

# Support vector machines

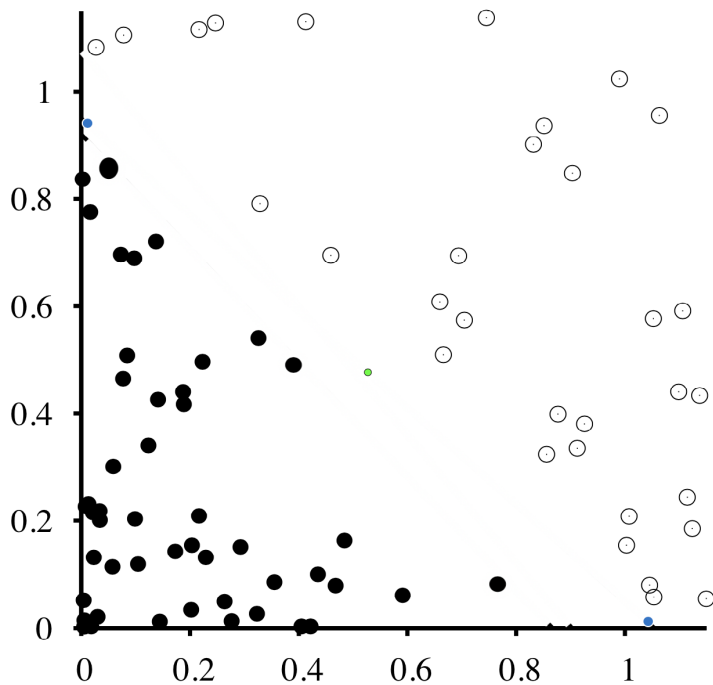
# Support vector machines

- Problem: How to categorize examples of something?

## Introduction

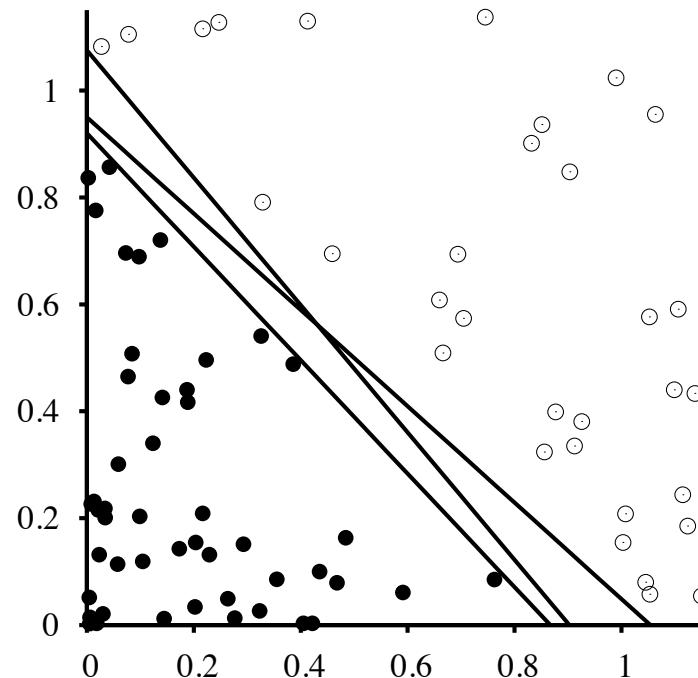
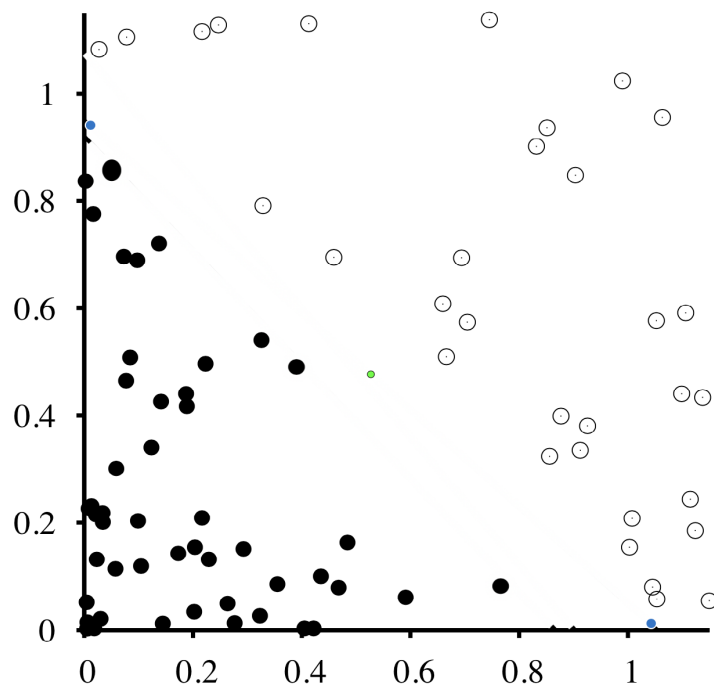
### “Classical” ML

- Induction
- Decision tree learning
- Explanation-based learning
- Support vector machines
- Genetic algorithms
- Case-based reasoning
- Schema-Based Reasoning



# Support vector machines

- Problem: How to categorize examples of something?



## Introduction

### "Classical" ML

- Induction
  - Decision tree learning
  - Explanation-based learning
  - Support vector machines
  - Genetic algorithms
  - Case-based reasoning
  - Schema-Based Reasoning

# Support vector machines

Machine Learning:  
Part I

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

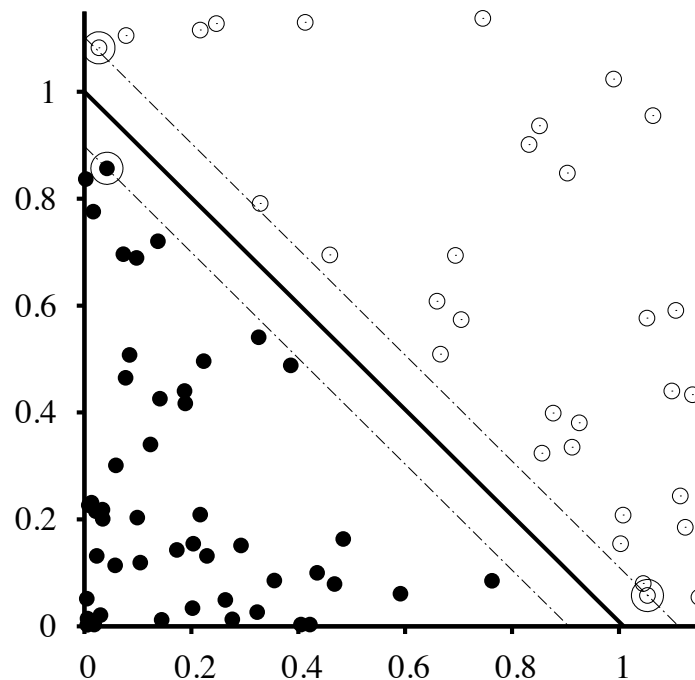
Schema-Based Reasoning

- ▶ Maximize ability to generalize by picking something in the middle
- ▶ *Maximum margin separator*



# Support vector machines

- ▶ Maximize ability to generalize by picking something in the middle
- ▶ *Maximum margin separator*



## Introduction

### “Classical” ML

Induction  
Decision tree learning  
Explanation-based learning  
Support vector machines  
Genetic algorithms  
Case-based reasoning  
Schema-Based Reasoning

# Support vector machines

## Introduction

### “Classical” ML

Induction

Decision tree learning

Explanation-based learning

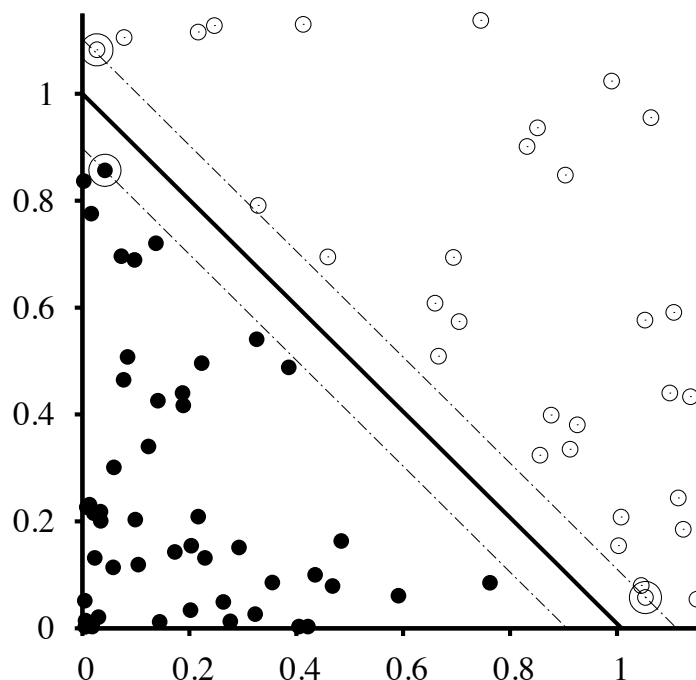
Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Maximize ability to generalize by picking something in the middle
- ▶ *Maximum margin separator*



- ▶ Mathematical formula (in book) to find this

- ▶ What if data no linearly-separable?
- ▶ Change dimensionality of data: map to a different space
- ▶ Look for *hyperplane* that cleanly separates data in that space
- ▶ How?
  - ▶ Formula uses dot product of vectors to data points to find maximal separator
  - ▶ Replace each vector  $\mathbf{x}$  with a higher-dimensional vector  $F(\mathbf{x})$
  - ▶ Dot products of these yield a *kernel function*: e.g.:

$$\vec{F}(\vec{x}_j) \cdot \vec{F}(\vec{x}_k) = (\vec{x}_j \cdot \vec{x}_k)^2$$

- ▶ Different kernel function  $\rightarrow$  different mappings
- ▶ With  $n$  data points, can always separate in  $\geq n - 1$  dimensions

## Example

# SVMs: What are they good for?

- ▶ Categorizing images
- ▶ Detecting possible cancerous lesions (e.g., B. Toner’s work on breast cancer)
- ▶ Basically the same kinds of things that artificial neural networks (ANNs) are good for

# So... SVNs or ANNs?

- ▶ SVMs are non-parametric: meaning, size of model can grow with new data
- ▶ ANNs are parametric: number of weights fixed
- ▶ However, there are parametric variants of SVMs...
- ▶ ....and non-parametric versions of ANNs (if we include neuroevolutionary algorithms)
- ▶ ANNs discover compact representations of input – useful for other things
- ▶ But SVMs don't overfit, like ANNs do
- ▶ Not clear – perhaps a religious war? Need data/theory!

# Genetic algorithms

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Biological evolution – can be viewed as:
  - ▶ a search process
  - ▶ a learning process
- ▶ Characteristics:
  - ▶ Unsupervised
  - ▶ Highly parallel
  - ▶ Stochastic
  - ▶ Feedback: *fitness*
- ▶ Use ideas from evolution for AI: *genetic algorithms*

# Example: Match a sentence

- ▶ Initially, output is random
- ▶ Feedback: measure of how close output is to target
- ▶ Use biological-like operations (e.g., mutation, reproduction) to evolve better solutions
- ▶ Many generations to get to correct/good enough answer

Example



# GA is a *local search*

- ▶ Recall from CSP:
  - ▶ Didn't use a complete state representation in search
  - ▶ I.e., didn't start with state  $\{v_1 = val_i, v_2 = val_j, \dots\}$
  - ▶ Reason: Search space was *factorial* time exponential
  - ▶ But what if we can prune search effectively?

- ▶ Start with complete candidate solution
- ▶ If not solution (or within  $\epsilon$  of solution): change some small part of state
- ▶ Different changes  $\Rightarrow$  different *neighborhoods*
  - ▶ Defined by operators sometimes (gen. alg., e.g.)
  - ▶ Defined by problem other times (e.g., BSAT)
- ▶ Choose best neighborhood
- ▶ Hill-climbing search; can use sim. annealing
- ▶ How to choose neighbors?

# Biological evolution (again)

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Evolution: can be viewed as highly-parallel hill-climbing search
  - ▶ “Goal”: optimize for environment/produce most surviving offspring
  - ▶ Species population = state
  - ▶ Operators: mating, mutation, crossover, death
  - ▶ Pruning: death or lack of offspring
  - ▶ So parallel hill-climbing beam search
- ▶ Very successful, very flexible
- ▶ How to mimic?

- ▶ *Genetic algorithm:*
  - ▶ Parallel hill-climbing search
  - ▶ Fixed beam-size (cf. evolution: population size)
- ▶ States: populations of individuals
- ▶ Individual: bit string (usually) – candidate solution
- ▶ *Fitness function:* applied to individual  $\Rightarrow$  *fitness*
- ▶ Best individuals reproduce, get rid of some  $\Rightarrow$  next generation

## Introduction

### “Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Select  $n$  most fit individuals
- ▶ Let pairs reproduce  $\Rightarrow$  new children
  - ▶ No paired chromosomes, so no (re)assortment of chromosomes
  - ▶ *Mutation*
  - ▶ *Cross-over*
- ▶ Replace least fit with children

- ## Introduction

## “Classical” ML

## Induction

## Decision tree learning

## Explanation-based learning

## Support vector machines

## Genetic algorithms

## Case-based reasoning

## Schema-Based Reasoning



- Copyright © 2019 UMaine School of Computing and Information Science

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

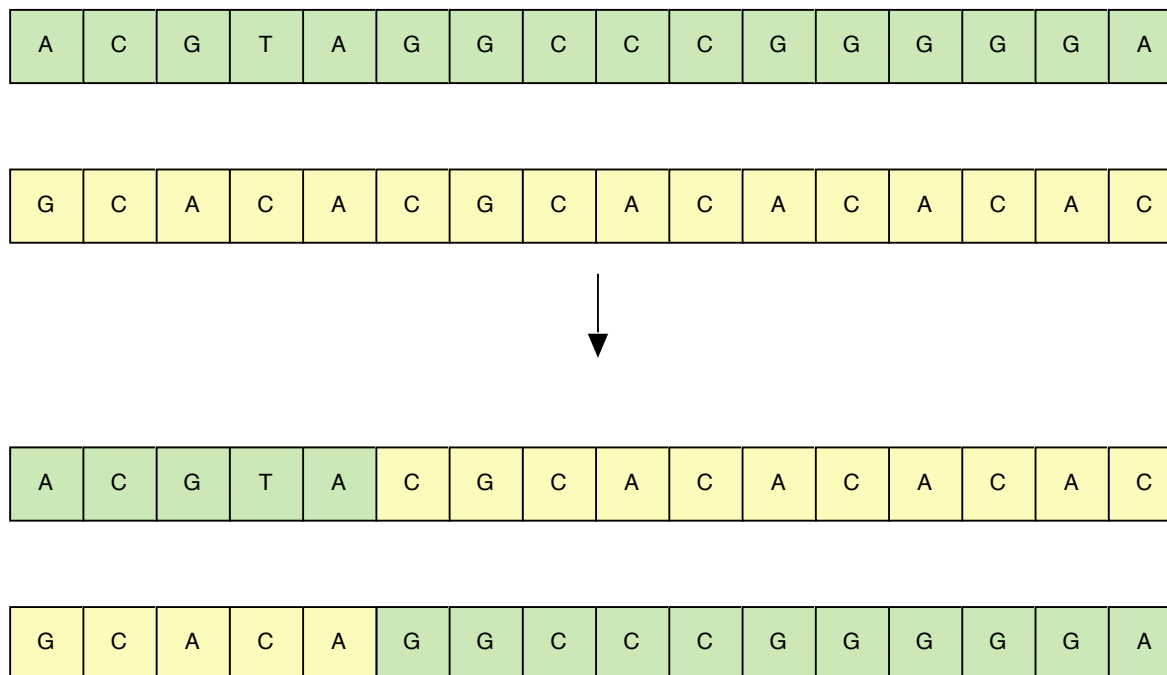
Case-based reasoning

Schema-Based Reasoning

- ▶ Select  $n$  most fit individuals
- ▶ Let pairs reproduce  $\Rightarrow$  new children
  - ▶ No paired chromosomes, so no (re)assortment of chromosomes
  - ▶ *Mutation*
  - ▶ *Cross-over*
- ▶ Replace least fit with children

# Reproduction

- ▶ Select  $n$  most fit individuals
- ▶ Let pairs reproduce  $\Rightarrow$  new children
  - ▶ No paired chromosomes, so no (re)assortment of chromosomes
  - ▶ *Mutation*
  - ▶ *Cross-over*



- ▶ Replace least fit with children



- ▶ Mutation rate
- ▶ Crossover rate

## Introduction

### “Classical” ML

Induction  
Decision tree learning  
Explanation-based learning  
Support vector machines  
Genetic algorithms  
Case-based reasoning  
Schema-Based Reasoning

- ▶ Can use heuristic function
- ▶ Can also simulate solution, measure fitness
- ▶ E.g.:
  - ▶ Learning a string: How many characters are correct?
  - ▶ TSP: combination of penalty for it not being a circuit, not being a Hamiltonian circuit, path cost

## Introduction

### “Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Pros:
  - ▶ Good for scheduling
  - ▶ Good for optimization problems
  - ▶ Design: antennae, solar collection mirrors, robot gait, . . .
- ▶ Cons:
  - ▶ Fitness function:
    - ▶ Applied *many* times, can be costly
    - ▶ Maybe difficult to devise
  - ▶ Not clear what termination conditions should be sometimes
  - ▶ Local optima sometimes

# Case-based reasoning

- ▶ Some problems:
  - ▶ Little domain knowledge
  - ▶ Sparse examples
- ▶ E.g., law, design
- ▶ Even medicine: Rich domain knowledge, but insufficient/too complex
- ▶ Approach for these domains:
  - ▶ Compare current problem to previous, similar ones solved (*cases*)
  - ▶ Transfer information to help solve current one

# CBR examples

- ▶ Labor mediation
- ▶ Meal planning
- ▶ Law
- ▶ Dispute mediation

## Introduction

### “Classical” ML

Induction  
Decision tree learning  
Explanation-based learning  
Support vector machines  
Genetic algorithms  
Case-based reasoning  
Schema-Based Reasoning

# How does it work?

Introduction

“Classical” ML

Induction

Decision tree learning

Explanation-based learning

Support vector machines

Genetic algorithms

Case-based reasoning

Schema-Based Reasoning

- ▶ Solve a problem  $\Rightarrow$  create a case representation  $\rightarrow$  memory
- ▶ New case: get *reminded* of prior case(s)
- ▶ Use old case:
  - ▶ Make predictions
  - ▶ Try same solution
  - ▶ Try similar reasoning

# Finding a case: Dynamic memory

- ▶ Could use anything
- ▶ Often used a *dynamic memory* [Schank, Kolodner]
- ▶ Store cases relative to similarities and differences
- ▶ Similar cases encountered:
  - ▶ create *memory organization packet* (MOP) from similarities
  - ▶ identify *predictive features* to discriminate between cases
  - ▶ *index* each case based on values of predictive features
- ▶ Add new case: traverse MOP structure based on case's features/values
- ▶ Model of human memory: have to *elaborate* indices, can forget
- ▶ Memory organization/reorganization based on use



# Schema-Based Reasoning

- ▶ Problem with CBR:
  - ▶ anecdotal reasoning
  - ▶ doesn't use induction of generalizations
- ▶ But it *does* have generalized cases: MOPs
- ▶ Why not use them and only fall back on cases?
- ▶ MOPs are *schemas*
- ▶ *Schema-based reasoning* [Turner]

# Example: MEDIC

- ▶ Doctors *do* use cases: “I remember one time I saw...”
- ▶ But seems to be the exception rather than rule: “Usually when you see something like this...”
- ▶ MEDIC program: schema-based reasoner in small area of pulmonology
- ▶ Signs and symptoms schemas that represent generalized cases
- ▶ *Interpret* the schemas rather than transfer information

- ▶ Contextual schemas (prev: dxMOPs):
  - ▶ represent problem-solving sessions or parts of sessions
  - ▶ generalized cases
- ▶ Procedural schemas:
  - ▶ generalized action portions of similar cases
  - ▶ can be plans, actions
- ▶ Strategic schemas:
  - ▶ represent problem-solving strategies
  - ▶ e.g., of a novice, expert, . . .
  - ▶ “meta-contextual schemas”
  - ▶ not currently used

## Introduction

### “Classical” ML

- Induction
- Decision tree learning
- Explanation-based learning
- Support vector machines
- Genetic algorithms
- Case-based reasoning
- Schema-Based Reasoning

- ▶ Initial domain: AUV control
- ▶ Uses more mature form of SBR: *context-mediated behavior* (CMB)
- ▶ Identify current context
- ▶ Use c-schema(s) to:
  - ▶ Provide semantic knowledge
  - ▶ Make predictions about situation
  - ▶ Decide goal priorities
  - ▶ Diagnose, assess, and handle unanticipated events
  - ▶ Achieve goals

## Introduction

### "Classical" ML

Induction  
Decision tree learning  
Explanation-based learning  
Support vector machines  
Genetic algorithms  
Case-based reasoning  
Schema-Based Reasoning

