

Lecture 14

- Learning
 - Inductive inference
 - Probably approximately correct learning

CS-424 Gregory Dudek

What is learning?

Key point: all learning can be seen as learning the *representation* of a function.

Will become clearer with more examples!

Example representations:

- propositional if-then rules
- first-order if-then rules
- first-order logic theories
- decision trees
- neural networks
- Java programs

CS-424 Gregory Dudek

Learning: formalism

Come up with some function f such that

- $f(x) = y$
for all training examples (x,y) and
- f (somehow) generalizes to yet unseen examples.
 - In practice, we don't always do it perfectly.

CS-424 Gregory Dudek

Inductive bias: intro

- There has to be some structure apparent in the inputs in order to support generalization.
- Consider the following pairs from the phone book.

Inputs	Outputs
Ralph Student	941-2983
Louie Reasoner	456-1935
Harry Coder	247-1993
Fred Flintstone	???-????

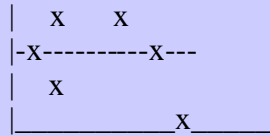
There is not much to go on here.

- Suppose we were to add zip code information.
- Suppose phone numbers were issued based on the spelling of a person's last name.
- Suppose the outputs were user passwords?

CS-424 Gregory Dudek

Example 2

- Consider the problem of fitting a curve to a set of (x,y) pairs.



- Should you fit a linear, quadratic, cubic, piece-wise linear function?
- It would help to have some idea of how smooth the target function is or to know from what family of functions (e.g., polynomials of degree 3) to choose from.
- Does this sound like cheating? What's the alternative?

CS-424 Gregory Dudek

Inductive Learning

Given a collection of examples $(x, f(x))$, return a function h that approximates f .

h is called the **hypothesis** and is chosen from the **hypothesis space**.

- What if f is not in the hypothesis space?

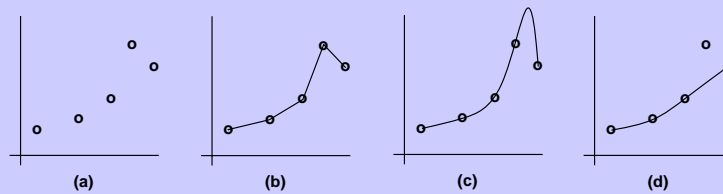
CS-424 Gregory Dudek

Inductive Bias: definition

- This "some idea of what to choose from" is called an inductive bias.
- Terminology
 - H, **hypothesis space** - a set of functions to choose from
 - C, **concept space** - a set of possible functions to learn
- Often in learning we search for a hypothesis f in H that is consistent with the training examples, i.e., $f(x) = y$ for all training examples (x, y) .
- In some cases, any hypothesis consistent with the training examples is likely to generalize to unseen examples. The trick is to find the right bias.

CS-424 Gregory Dudek

Which hypothesis?



CS-424 Gregory Dudek

Bias explanation

How does learning algorithm decide

Bias leads them to prefer one hypothesis over another.

Two types of bias:

- **preference bias** (or search bias) depending on how the hypothesis space is explored, you get different answers
- **restriction bias** (or language bias), the “language” used: Java, FOL, etc. (h is not equal to c).
- e.g. language: piece-wise linear functions: gives (b)/(d).

CS-424 Gregory Dudek

Issues in selecting the bias

Tradeoff (similar in reasoning):
more expressive the language, the harder to find
(compute) a good hypothesis.

Compare: propositional Horn clauses with first-order
logic theories or Java programs.

- Also, often need more examples.

CS-424 Gregory Dudek

Occam's Razor

- Most standard and intuitive preference bias:

Occam's Razor

(aka Ockham's Razor)

*The most likely hypothesis is
the **simplest** one that is
consistent with all of the
observations.*

- Named after Sir William of Ockham.

CS-424 Gregory Dudek

Implications

- The world is simple.
- The chances of an accidentally correct explanation are low for a simple theory.

CS-424 Gregory Dudek

Probably Approximately Correct (PAC) Learning

Two important questions that we have yet to address:

- Where do the training examples come from?
- How do we test performance, i.e., are we doing a good job learning?
- PAC learning is one approach to dealing with these questions.

CS-424 Gregory Dudek

Classifier example

Consider learning the predicate $Flies(Z) = \{true, false\}$.

We are assigning objects to one of two categories: recall we call this a *classifier*.

Suppose that $X = \{\text{pigeon}, \text{dodo}, \text{penguin}, 747\}$, $Y = \{true, false\}$, and that

$\text{Pr}(\text{pigeon}) = 0.3$	$\text{Flies}(\text{pigeon}) = \text{true}$
$\text{Pr}(\text{dodo}) = 0.1$	$\text{Flies}(\text{dodo}) = \text{false}$
$\text{Pr}(\text{penguin}) = 0.2$	$\text{Flies}(\text{penguin}) = \text{false}$
$\text{Pr}(747) = 0.4$	$\text{Flies}(747) = \text{true}$

Pr is the distribution governing the presentation of training examples (how often do we see such examples).

We will use the same distribution for evaluation purposes.

CS-424 Gregory Dudek

- Note that if we mis-classified dodos but got everything else right, then we would still be doing pretty well in the sense that 90% of the time we would get the right answer.
- We formalize this as follows.

CS-424 Gregory Dudek

- The approximate error associated with a hypothesis f is
- $$\text{error}(f) = \sum \{x \mid f(x) \neq \text{Flies}(x)\} \Pr(x)$$
- We say that a hypothesis is **approximately correct with error at most ϵ** if $\text{error}(f) \leq \epsilon$

CS-424 Gregory Dudek

- The chances that a theory is correct increases with the number of consistent examples it predicts.
- Or....
- A badly wrong theory will probably be uncovered after only a few tests.

CS-424 Gregory Dudek

PAC: definition

Relax this requirement by not requiring that the learning program necessarily achieve a small error but only that it to keep the error small **with high probability**.

Probably approximately correct (PAC) with probability δ and error at most ϵ if, given any set of training examples drawn according to the fixed distribution, the program outputs a hypothesis f such that

$$\Pr(\text{Error}(f) > \epsilon) < \delta$$

CS-424 Gregory Dudek

PAC

- Idea:
- Consider space of hypotheses.
- Divide these into “good” and “bad” sets.
- Want to assure that we can close in on the set of good hypotheses that are close approximations of the correct theory.

CS-424 Gregory Dudek

PAC Training examples

Theorem:

If the number of hypotheses $|H|$ is finite, then a program that returns an hypothesis that is consistent with

$$\ln(\delta / |H|) / \ln(1 - \epsilon)$$

training examples (drawn according to \Pr) is guaranteed to be PAC with probability δ and error bounded by ϵ .

CS-424 Gregory Dudek

PAC theorem: proof

If f is not approximately correct then $\text{Error}(f) > \epsilon$ so the probability of f being correct on one example is $< 1 - \epsilon$ and the probability of being correct on m examples is $< (1 - \epsilon)^m$.

Suppose that $H = \{f, g\}$. The probability that f correctly classifies all m examples is $< (1 - \epsilon)^m$. The probability that g correctly classifies all m examples is $< (1 - \epsilon)^m$. The probability that one of f or g correctly classifies all m examples is $< 2 * (1 - \epsilon)^m$.

To ensure that any hypothesis consistent with m training examples is correct with an error at most ϵ with probability δ , we choose m so that $2 * (1 - \epsilon)^m < \delta$.

CS-424 Gregory Dudek

Generalizing, there are $|H|$ hypotheses in the restricted hypothesis space and hence the probability that there is some hypothesis in H that correctly classifies all m examples is bounded by

$$|H|(1 - \epsilon)^m.$$

Solving for m in

$$|H|(1 - \epsilon)^m < \delta$$

we obtain

$$m \geq \ln(\delta / |H|) / \ln(1 - \epsilon).$$

QED

CS-424 Gregory Dudek

Stationarity

- Key assumption of PAC learning:

Past examples are drawn randomly from *the same distribution* as future examples: stationarity.

The number m of examples required is called the **sample complexity**.

CS-424 Gregory Dudek

A class of concepts C is said to be PAC learnable for a hypothesis space H if (roughly) there exists an polynomial time algorithm such that:
for any c in C , distribution Pr , ϵ , and δ ,

if the algorithm is given a number of training examples polynomial in $1/\epsilon$ and $1/\delta$ then with probability $1-\delta$ the algorithm will return a hypothesis f from H such that
 $Error(f) \leq \epsilon$.

CS-424 Gregory Dudek

Overfitting

Consider error in hypothesis h over:

- training data: $error_{train}(h)$
- entire distribution D of data: $error_D(h)$

- Hypothesis $h \in H$ **overfits** training data if
 - there is an alternative hypothesis $h' \in H$ such that

- $error_{train}(h) < error_{train}(h')$

but

- $error_D(h) > error_D(h')$