

## Supervised learning

### UNIT-IV

#### - Task of Supervised learning

Given a training set of  $N$  example i/p-o/p pairs  
 $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ .

where each  $y_j$  was generated by an unknown  $f$ .

$$y = f(x),$$

discover a  $f^h$  that approximates true  $f^h$   $f$ .

→  $x$  &  $y$  can be any value, need not be numbers.

→  $f^h$  is a hypothesis

→ Learning is a Search through the space of possible hypotheses for one that will perform well, even on new examples beyond training set.

→ To measure accuracy of hypothesis we give it a test set of examples that are distinct from training set.

→ hypothesis generalizes well if it correctly predicts the value of  $y$  for examples.

→ Sometimes function  $f$  is stochastic - it is not strictly a  $f^h$  of  $x$ . ~~so that we~~

→ When o/p  $y$  is one of a finite set of values (such as sunny, cloudy or rainy) the learning problem is called classification.

→ When  $y$  is a no., the learning problem is called regression.  
(like tomorrow's temperature)

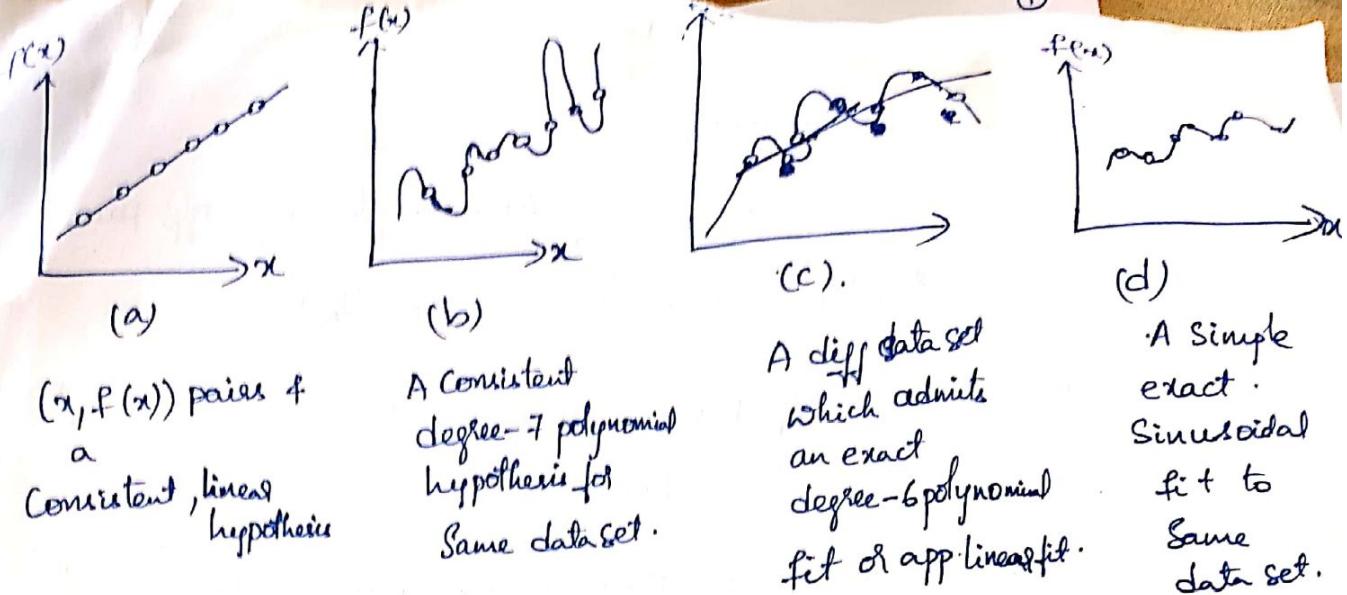


fig shows a: fitting a  $f^n$  of a single variable to some data points.

points in  $(x, y)$  plane where  $y = f(x)$ .

we don't  $f$ , but we will approximate it with  $f^n$ .  
Selected from hypothesis Space  $H$ , we take set of polynomials  $x^5 + 3x^2 + 2$ .

(a) Shows some data with an exact fit by a straight line ( $0.4x + 3$ ).

called a consistent hypothesis b'coz it agrees with all data.

(b) Shows a high degree polynomial that is also consistent with same data.

Illustrates a fundamental problem in inductive learning: how do we choose from among multiple consistent hypotheses?

to prefer simplest hypothesis consistent with data.

Called Ockham's razor principle.

(C) Shows a 2<sup>nd</sup> data set. There is no consistent straight line for this data set, it requires a degree-6 polynomial for an exact fit.

Polynomial with 7 parameters doesn't seem to be finding any pattern in data & we can't expect it to generalize well.

A straight line i.e., not consistent with any of generalize data points, but might fairly well for unseen values 'x'.

- There is a tradeoff b/w complex hypotheses that fit the training data well & simpler hypotheses that may generalize better.

(d) expand the hypotheses<sup>space</sup>  $H$  to allow polynomials over both  $x$  &  $\sin(x)$ . find that data in (C) can be fitted exactly by a simple  $f^n$   $ax+bx\sin(x)$ . Shows the importance of choice of hypotheses space. learning problem is realizable if hypotheses space contains true  $f^n$ .

In some cases, an analyst looking at a problem is willing to make more fine-grained distinctions about hypothesis space, to say - even before seeing any data - not just that a hypothesis is possible or impossible, but rather how probable it is.

Supervised learning can be done by choosing  
probable given data:  
hypotheses  $h^*$  i.e. most

$$h^* = \underset{h \in H}{\operatorname{argmax}} P(h/\text{data})$$

By Bayes' rule this is equivalent to

$$h^* = \underset{h \in H}{\operatorname{argmax}} P(\text{data}/h) P(h)$$

- we can say prior probability  $P(h)$  is high for degree-10<sup>2</sup> polynomial, lower for degree-7 polynomial
- we allow unusual-looking  $f^*$  when data say we really need them, but we discourage them by giving a low prior probability.

Why not let  $H$  be class of all  $\text{Java programs or Turing/m/c}$  ( $\text{TM/c}$ )

- Every computable  $f^*$  can be represented by some TM/c. & i.e., best we can do.
- problem is that it doesn't take into account computational complexity of learning.

There is a tradeoff b/w expressiveness of a hypothesis space & complexity of finding a good hypothesis within that Space.

Ex:- fitting a straight line to data is an easy computation.

Fitting a high degree polynomials is harder

!! TM/c's is undecidable.

Why we prefer simple hypotheses spaces is that we want to use  $h$  after we have learned it & computing  $h(x)$  &  $h$  is linear  $f^*$  is guaranteed to be fast, while TM/c's is undecidable.

## LEARNING DECISION TREES

Decision tree induction is simpler & most successful forms of M/c learning.

We describe representations - the hypothesis Space & then show how to learn a good hypothesis.

### The decision tree representation

- A decision tree represents a function that takes as i/p a vector of attribute values and returns a "decision" - a single o/p value.
- i/p & o/p values can be discrete or continuous.
- Problems with i/p's have discrete values & o/p has exactly two possible values, i.e., Boolean classification, where each example i/p will be classified as true or false.
- Decision tree reaches its decision by performing a sequence of tests.
- Each internal node in tree corresponds to a test of value of one of i/p attributes,  $A_i$  and the branches from node are labeled with possible values of attribute,  $A_i = v_{ik}$ .

Aim here is to learn a definition for goal predicate  
Will Wait.

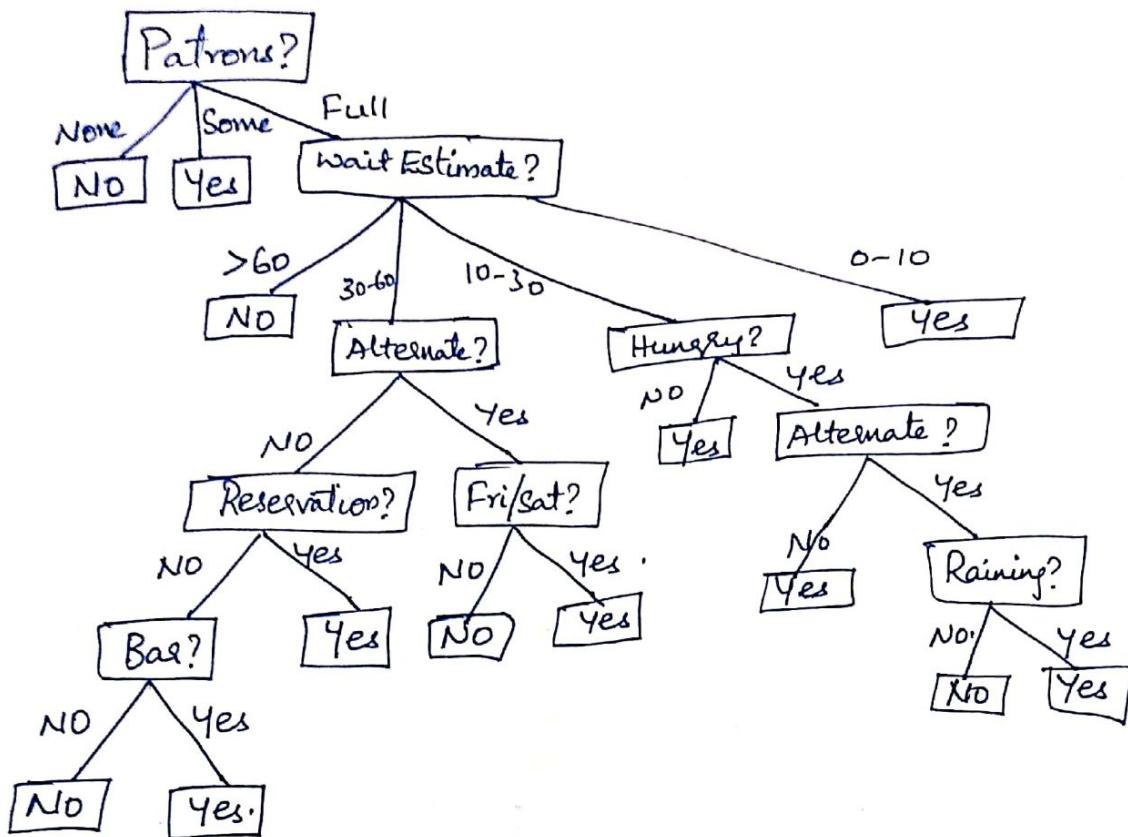
1<sup>st</sup> we list the attributes:

1. Alternate: whether there is a suitable alternative restaurant
2. Bar: whether the restaurant has a comfortable bar area nearby.
3. Fri/Sat: true on Fridays and Saturdays.
4. Hungry: whether we are hungry.
5. Patrons: how many people are in restaurant (values are None, Some & 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. Wait Estimate: the wait estimated by host  
(0-10 minutes, 10-30, 30-60, or >60).

Every variable has a small set of possible values;  
value of WaitEstimate, for ex: is not an integer, rather  
it is one of 4 discrete values 0-10, 10-30, 30-60 or >60.

- Notice that tree ignore Price & type attributes.
- Examples are processed by tree starting at root and following the appropriate branch until a leaf is reached.
- Ex: Patrons = Full & WaitEstimate = 0-10 will be classified as ~~no time yet we will wait for a table~~

A decision tree for deciding whether to wait for a table.



### Expressiveness of decision trees

→ A boolean decision tree is logically equivalent to assertion that the goal attribute is true if & only if i/p attributes satisfy one of the paths leading to a leaf with value true,

In propositional logic, we have

$$\text{Goal} \Leftrightarrow (\text{Path}_1 \vee \text{Path}_2 \vee \dots)$$

Any  $f^n$  in propositional logic can be expressed as a decision tree.

Ex: Path = (Patrons=Full  $\wedge$  WaitEstimate=0-10).

For variety of problems, decision tree format yields a nice, concise result, but some  $f^n$ 's can't be represented concr.

## Inducing decision trees from examples

An example for a Boolean decision tree consists of an  $(x, y)$  pair

where  $x$  is a vector of values for i/p attributes,  
 $y$  is a single Boolean o/p value.

- Positive examples are the ones in which goal will wait in true ( $x_1, x_3 \dots$ )
- Negative examples are the ones in which it is false ( $x_2, x_5 \dots$ )
- We need a tree i.e., consistent with examples & is as small as possible.
- It is an intractable problem to find the smallest consistent tree
- Decision - Tree - LEARNING Algorithm adopts a greedy divide - and - conquer strategy: always test the most important attribute first.
- Type is a poor attribute; b'coz it leaves us with 4 possible outcomes
- Patrons is fairly important, b'coz if value is not
- Four Cases to consider recursive problems:

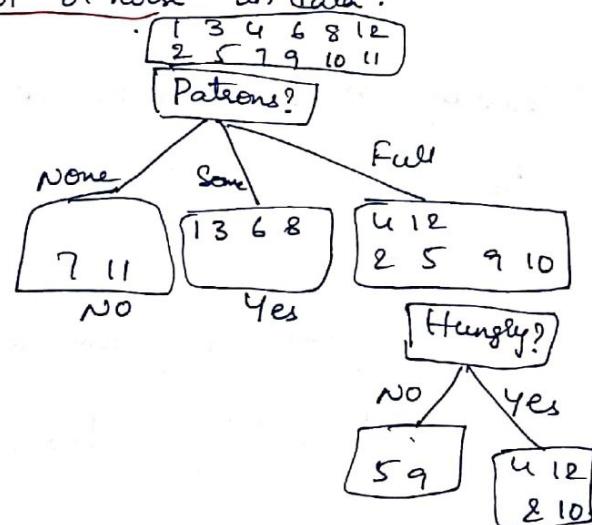
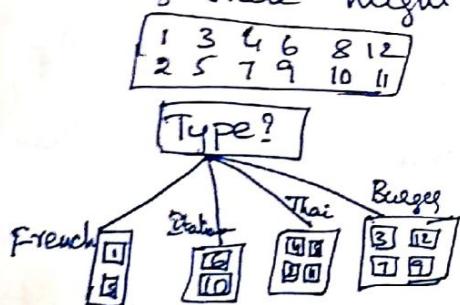
1) If remaining examples are all positive (or all -ve) then we are done: we can answer Yes or NO.  
like None & Some .. in the ex

2) If there are some +ve & some -ve examples, then choose the best attribute to split them.  
Hungry split to remaining examples

3. If there are no examples left, it means that no example has been observed for this combination of attribute values, & we return a default value calculated from the plurality classification of all examples that were used in constructing the node's parent. There are passed along in variable parent-examples.

4. If there are no attributes left, but both +ve & -ve examples, it means that these examples have exactly the same description but diff classifications. This can happen

b'coz there might be error or noise in data.



## Decision-tree learning algorithm

function Decision-Tree-Learning(examples, attributes, parent-example) returns a tree

if examples is empty then return PLURALITY\_VALUE(parent-example)

else if all examples have same classification then return the classification

else if attributes is empty then return PLURALITY\_VALUE(examples)

else

$A \leftarrow \operatorname{argmax}_{a \in \text{attributes}}$

IMPORTANCE(a, examples)

$\downarrow$   
selects most common o/p  
value among a  
set of examples,  
breaking ties.  
randomly.

tree  $\leftarrow$  a new decision tree with root test A.

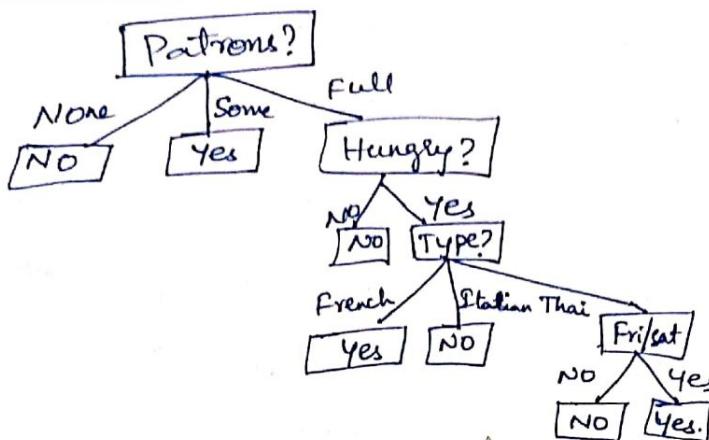
for each value  $v_k$  of A, do

exist e: ee examples and  $e.A = v_k\}$

Subtree  $\leftarrow$  Decision -TREE -LEARNING(exs, attributes,  $-A$ , examples)

add a branch to tree with label ( $A = V_k$ ) and Subtree  
return tree

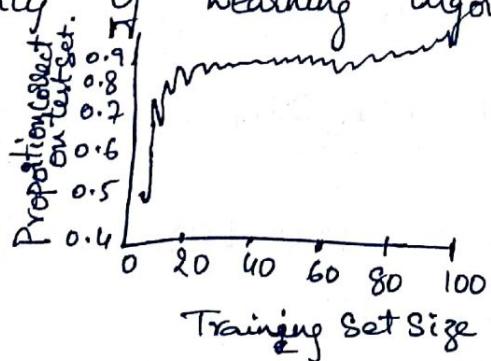
### Decision Tree



→ learning algorithm looks at examples, not at collect  $f^n$  &  
its hypothesis is consistent with all examples & simpler than  
original tree.

→ ~~the~~ the learning algorithm has no reason to include tests  
for Raining & Reservation.

→ Accuracy of learning algorithm is evaluated Learning Curve.



For 100exs, we split into a training set & a test set.  
we learn a hypothesis  $h$ . with training set & measure its accuracy  
with test set, this is done 1<sup>st</sup> for size 1 & increasing up to size 99.  
for each size we repeat process splitting 20 times & avg results 20times  
curve shows training set size grows, accuracy increases  
so learning curves are called happy graphs.

## Choosing Attribute tests

- Greedy Search used in decision tree learning is to minimize the depth of final tree.
  - The Patrons attribute is not perfect but fairly good. useless attribute like Type leaves the example sets with same proportion of +ve & -ve examples.
- So we need a formal measure of "fairly good" & "really useless" & can implement IMPORTANCE f^n.
- We use notion of information gain defined in terms of entropy.
- Entropy is a measure of uncertainty of a random variable
- Acquisition of information corresponds to a reduction in entropy.
- A random variable with only one value - a coin that always comes up heads - has no uncertainty & thus entropy is defined as 0.
- A flip of a fair coin is equally likely to come up heads or tails, 0 or 1 & we count as "1 bit" of entropy.
- Four-sided die has 2 bits of entropy, b'coz it takes 2 bits to describe one of 4 equally probable choices.

- Consider unfair coin that comes up heads 99% of time.  
if we guess heads will be wrong only 1% of time.  
then entropy measure i.e., can be close to zero. but +1
- Entropy of a random variable  $V$  with values  $v_k$ , each with probability  $P(v_k)$ .

$$\text{Entropy: } H(v) = \sum_k P(v_k) \log_2 \frac{1}{P(v_k)} = -\sum_k P(v_k) \log_2 P(v_k)$$

Entropy of fair coinflip is 1 bit.

$$H(\text{Fair}) = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1.$$

if coin is loaded to give 99% heads, we get

$$H(\text{loaded}) = -(0.99 \log_2 0.99 + 0.01 \log_2 0.01) \approx 0.08 \text{ bits.}$$

$B(q)$  as entropy of a Boolean random variable  
i.e.,  $T$  with Probability  $q$ :

$$B(q) = -(q \log_2 q + (1-q) \log_2 (1-q))$$

If a training set contains  $P$  +ve examples &  $n$  -ve examples

then entropy of goal attribute on whole set is

$$H(\text{Goal}) = B\left(\frac{P}{P+n}\right).$$

An attribute  $A$  with distinct values divides training set  $E$  into  
each subset  $E_k$  has  $P_k$  +ve &  $n_k$  -ve examples.  $E_1, \dots, E_d$ .

So  $B\left(\frac{P_k}{P_k+n_k}\right)$  we need an additional bits of info to answer question.

(7)

randomly chosen examples from training set has  $k^{th}$  value  
for attribute with probability  $(P_k + n_k)/(P_t + n)$ .

Expected entropy remaining after testing attribute A

is  $\text{Remainder}(A) = \sum_{k=1}^d \frac{P_k + n_k}{P_t + n} B\left(\frac{P_k}{P_k + n_k}\right)$ .

Information gain from attribute 'test' on A is  
expected reduction in entropy:

$$\text{Gain}(A) = B\left(\frac{P}{P_t + n}\right) - \text{Remainder}(A).$$

### Generalization & Overfitting

The Decision-Tree-Learning Algo will generate a large tree when there is actually no pattern to be found.

- Overfitting occurs when target f is not random.
- In overfitting no. of i/p attributes grows & we increase the ~~no~~ no. of training examples.
- Decision tree pruning combats overfitting.
  - Pruning eliminates nodes that are not relevant.
  - we start with a full tree generated by Decision-Tree-learning we look at test node that has only leaf nodes as descendants, if test appears to be irrelevant - detecting only noise in data - then we eliminate the test, replacing it with a leaf node. like this we consider each test with only leaf descendants, until each one has either pruned or accepted.

How do we detect that a node is testing an irrelevant attribute?

Suppose we are at a node consisting of P+ve & n-ve examples.

If attribute is irrelevant, we would expect that it would split the examples into subsets that each have roughly the same proportion of +ve examples as whole set  $P/(P+n)$  & so information gain will be close to zero.

Information gain is a good clue to irrelevance.

How large a gain should we require in order to split on a particular attribute?

- We can answer this question by using a statistical

Significance test. Such a test begins by assuming there is no underlying pattern (null hypothesis).

- Then actual data are analyzed to calculate the extent to which they deviate from a perfect absence of pattern.

- If degree of deviation is statistically unlikely (mean 5% probability or less), then it is considered to be good evidence for presence of a significant pattern in data.

- We need to calculate the probability under null hypothesis, a sample of size  $v=n+p$  which is expected distribution of +ve & -ve examples.

- We can measure the deviation by comparing actual numbers of +ve & -ve examples in each subset

$P_k$  &  $n_k$  with expected nos  $\hat{P}_k$  &  $\hat{n}_k$  assuming true irrelevance.

$$\hat{P}_k = P \times \frac{P_k + n_k}{P+n}$$

$$\hat{n}_k = n \times \frac{P_k + n_k}{P+n}$$

## Evaluating and Choosing the Best Hypothesis

- We want to learn a hypothesis that fits the future dataset.
- To make it precise we need to define "future data" and "best fit".
- We make a Stationarity assumption: that there is a probability distribution over examples that remains stationary over time.
- Each example data point is a random variable  $E_j$  whose observed value  $e_j = (x_j, y_j)$  is sampled from that distribution & is independent of previous examples.

$$P(E_j | E_{j-1}, E_{j-2}, \dots) = P(E_j).$$

& each example has an identical prior probability distribution:

$$P(E_j) = P(E_{j-1}) = P(E_{j-2}) = \dots$$

- Examples that satisfy these assumptions are called independent & identically distributed or i.i.d.
- i.i.d. assumption connects past to future. Future could be anything.

→ Next step is to define "best fit" we define error rate of a hypothesis as the proportion of mistakes it makes - proportion of times that  $h(x) \neq y$  for any  $(x, y)$  example

→ just because a hypothesis  $h$  has a low error rate on training set doesn't mean that it will generalize well.

- + zero but +ve.
- Professor knows that an exam will not accurately evaluate students if they have been already seen the exam questions.
  - In order to get an accurate evaluation of a hypothesis, we need to test it on a set of examples it has not seen yet.
  - Simplest approach is one we have seen already:  
randomly split the available data into a training set from which learning algorithm produces  $h$  & a test set on which accuracy of  $h$  is evaluated.  
This method, sometimes called 'hold-out cross-validation' which has a disadvantage that it fails to use all the available data.
  - We can squeeze more out of data & still get an accurate estimate using a technique called K-fold cross-validation.
  - Each example serves double duty - as training data and test data.
  - 1<sup>st</sup> we split the data into  $K$  equal subsets.  $\Rightarrow$  then perform  $K$  rounds of learning. On each round  $\frac{1}{K}$  of data is held out as a test set & remaining examples are used as training data.  
*data*  
 *$K$  subsets -  $K$  rounds learning*  
 *$\frac{1}{K}$  data as test*
  - The average test set score of  $K$  rounds should then be a better estimate than a single score.
  - Popular values for  $K$  are 5 & 10 - enough to give an estimate i.e., statistically likely to be accurate, at a cost of 5 to 10 times longer computation time.  
Extreme is  $K=n$ , known as leave-one-out cross-validation.

Users frequently invalidate their results by peeking at test data alternatively, we can.

- Peeking: A learning algorithm has various "Knobs" that can be twiddled to tune its behaviour.

Researchers generate hypotheses for various different settings of knobs, measure their error rates on test set and reports the error rate of best hypothesis. Alas, Peeking has occurred!

- The reason is hypothesis was selected on basis of its test set error rate, so information about test set has baked into learning algorithm.

- Peeking is a consequence of using test-set performance to both choose a hypothesis and evaluate it.

- The way to avoid this is to really hold the test set out - lock it away until you are completely done with learning & simply wish to obtain an independent evaluation of final hypothesis.

- If test set is locked away, but you still want to measure performance on unseen data as a way of selecting a good hypothesis, then divide the available data into a training set & a validation set.

### Model Selection: Complexity Versus Goodness of fit

- The higher-degree polynomials can fit the training data better, but if degree is too high they will overfit & perform poorly on validation data.

- choosing degree of polynomial is an instance of problem of model selection.

Task of finding the best hypothesis ~~as~~ may have two tasks:

model Selection defines the hypothesis Space & then optimization finds the best hypothesis within that Space.

- In this section we will see how to select among models that are parameterized by Size.

→ Ex:- with polynomials ~~is~~ Size = 1 for linear  $f^n$   
2 for quadratics.

→ In decision ~~trees~~, Size could be no. of nodes in tree.

→ In all cases we want to find value of

Size parameter that best balances underfitting & overfitting to give the best test set accuracy.

→ Wrapper takes a learning algorithm as an argument (Decision-Tree-learning)

→ Wrapper enumerates models according to a parameter, Size.

→ for each Size, it uses cross validation on learner to compute the average error rate on training & test sets.

→ The training set error decreases monotonically, while the validation set error decreases at first, then increases when model begins to overfit.

→ cross validation procedure picks the value of Size with lowest validation set error; bottom of U-shaped curve.

## From error rates to loss

- We have been trying to minimize error rate.
- Consider problem of classifying email messages as Spam or non-Spam.
- It is worse to classify non-Spam as Spam than classify Spam as non-Spam.
- So a classifier with a 1% error rate, where almost all errors were classifying Spam as non-Spam, would be better than a classifier with only a 0.5% error rate, if most of those errors were classifying non-Spam as Spam.
- In M/c learning it is traditional to express utilities by means of a loss  $f^n$ .
- loss  $f^n L(x, y, \hat{y})$  is defined as amount of utility lost by predicting  $h(x) = \hat{y}$  when the correct answer is  $f(x) = y$ :

$$L(x, y, \hat{y}) = \text{utility(result of using } y \text{ given an i/p } x) - \text{utility(result of using } \hat{y} \text{ given an i/p } x).$$

$$L(\text{Spam}, \text{nonspam}) = 1, L(\text{nonspam}, \text{Spam}) = 10.$$

let  $\mathcal{E}$  be set of possible i/p-o/p examples.

Then expected generalization loss for a hypothesis  $h$ .  
is

$$\text{Gen Loss}_L(h) = \sum_{(x,y) \in \mathcal{E}} L(y, h(x)) P(x, y),$$

& best hypothesis,  $h^*$  is one with minimum  
expected generalization loss:

$$h^* = \underset{h \in H}{\operatorname{argmin}} \text{Gen Loss}_L(h)$$

$$h^* = \underset{h \in H}{\operatorname{argmin}} \text{Gen Loss}_L(h).$$

$P(x, y)$  is not known, learning agent can only  
estimate generalization loss with empirical loss on  
a set of examples,  $E$ :

$$\text{Emp Loss}_{L,E}(h) = \frac{1}{N} \sum_{(x,y) \in E} L(y, h(x)).$$

estimated best hypothesis  $\hat{h}^*$  is then the one  
with min empirical loss:

$$\hat{h}^* = \underset{h \in H}{\operatorname{argmin}} \text{Emp Loss}_{L,E}(h).$$

There are 4 reasons why  $\hat{h}^*$  may differ from true function,  
 $f$ : unrealizability, variance, noise, computational  
complexity.

- ~~+ max through~~
- 1<sup>st</sup> -  $f$  may not be realizable - may not be in  $H$ - or, may be present in such a way that other hypotheses are preferred.
- 2<sup>nd</sup> - a learning algorithm will return different hypotheses for different sets of examples, even if those sets are drawn from same true  $f$ .  $f$  & those hypothesis will make different predictions on new examples.
- ~~3rd~~ - higher the variance, higher the probability of error.
- 3<sup>rd</sup> -  $f$  may be non-deterministic or noisy - it may return different values for  $f(x)$  each time  $x$  occurs.  
noise can't be predicted.  
in many cases, it arises b'coz observed labels  $y$  are result of attributes of environment not listed in  $x$ .
- 4<sup>th</sup> -  $H$  is complex, it can be computationally intractable to systematically search the whole hypothesis space.  
Best we can do is a local search that explores only part of Space.  
That gives an appr. error.

## Regularization

we saw how to do model Selection with cross-validation  
on model size.

An alternate approach to search for a hypothesis  
that directly minimizes the weighted sum of empirical  
loss & complexity of hypothesis,

total cost.

$$\text{Cost}(h) = \text{EmpLoss}(h) + \lambda \text{Complexity}(h)$$

$$h^* = \underset{h \in H}{\operatorname{argmin}} \text{Cost}(h).$$

$\lambda$  is a parameter, a +ve number that serves  
as a conversion rate between loss & hypothesis  
complexity

This approach combines loss & complexity into one  
metric, to find best hypothesis all at once.

# Theory of Learning

13

- P hy  
— How can we be sure that our learning algorithm has produced a hypothesis that will predict correct value for previously unseen i/p's?
- ~~Qn~~ How do we know that hypothesis h is close to target f? If we don't know what f is?
- How many examples do we need to get a good h?
- What hypothesis Space should we use?
- If hypothesis space is very complex, can we find best h. do we have to <sup>Settle</sup> go for local maximum in Space of hypotheses?
- How Complex should h be?
- How do we avoid overfitting?
- we start with how many examples are needed for learning?
- The learning curve for decision tree on restaurant problem improves with more training data.
- Learning curves are useful, but specific to a particular learning algorithm on a particular problem.
- Computational learning theory will address the principles governing the no. of examples needed.
- Principle is that any hypothesis that is seriously wrong will almost certainly be "found out" with high probability after a small no. of examples, b'coz it will make an incorrect prediction.

That's

- handles 99% of time.

— Thus, any hypothesis that is consistent with a sufficiently large set of training examples is unlikely to be seriously wrong: i.e., it must be probably approximately correct.

→ Any learning algorithm that returns hypotheses that are probably approximately correct is called a PAC learning algorithm.

→ With this we can provide bounds on performance of various learning algorithms.

→ PAC-learning theorems, are logical consequences of axioms.

→ Simplest PAC theorems deal with Boolean functions for which the 0/1 loss is appropriate.

→ Error rate of hypothesis  $h$ , is defined as expected generalization error for examples drawn from stationary distribution:

$$\text{error}(h) = \text{GenLoss}_{\text{0/1}}(h) = \sum_{x,y} h_{0/1}(y, h(x)) P(x, y).$$

→ A hypothesis  $h$  is called approximately correct if  $\text{error}(h) \leq \epsilon$ , where  $\epsilon$  is a small constant.

→ Approximately correct hypothesis can be close to true f in hypothesis Space: it lies inside  $\epsilon$ -ball around the true  $f$ . Hypothesis space outside this ball is called  $H_{\text{bad}}$ .

→ Probability that a seriously wrong hypothesis is consistent with  $N$  examples.

$$\text{error}(h_b) > \epsilon.$$

$$h_b \in H_{\text{bad}}$$

→ probability that it agrees with a given example is at most  $1-\epsilon$ .  
→ Since examples are independent, the bound for  $N$  examples is  
$$P(h_b \text{ agrees with } N \text{ examples}) \leq (1-\epsilon)^N.$$

→ Probability that  $H_{\text{bad}}$  contains at least one consistent hypothesis is bounded by sum of individual probabilities

$$P(H_{\text{bad}} \text{ contains a consistent hypothesis}) \leq |H_{\text{bad}}| (1-\epsilon)^N \leq |H| (1-\epsilon)^N.$$

→ To reduce the probability of this event.

$$|H| (1-\epsilon)^N \leq \delta.$$

Given that  $1-\epsilon \leq e^{-\epsilon}$ , we can achieve this if we allow

$$N \geq \frac{1}{\epsilon} \left( \ln \frac{1}{\delta} + \ln |H| \right).$$

algorithm to see  
PAC learning examples: Learning Decision lists.

→ How to apply PAC learning to new hypothesis space: decision lists.

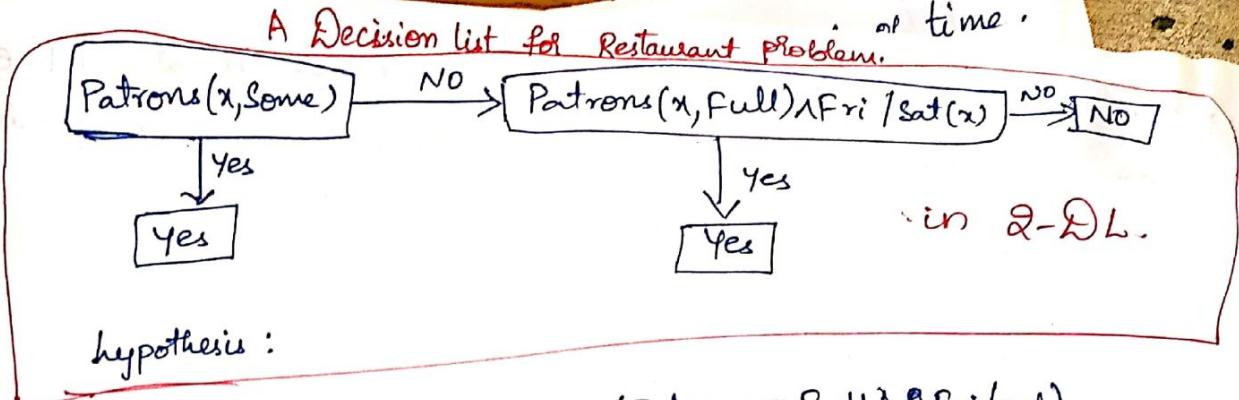
→ A decision list consists of a series of tests, each of which is a conjunction of literals.

→ If a test succeeds when applied to an example description, the decision list specifies the value to be returned.

→ If test fails, processing continues with next test in list.

→ Decision list resembles decision trees, but overall structure is simpler: they branch only in one direction.

→ The individual tests are more complex.



$\text{willWait} \Leftrightarrow (\text{Patrons} = \text{Some}) \vee (\text{Patrons} = \text{Full} \wedge \text{Fri/Sat})$ .

Any

→ If we allow tests of arbitrary size, then decision lists can represent any Boolean f.

→ If we restrict size of each test to at most K literals, then it is possible for learning algorithm to generalize successfully from a small no. of examples. We call this language K-DL.

→ Language referred to by K-DL depends on attributes used to describe the examples.

→ K-DL(n) to denote a K-DL language using n Boolean attributes.

→ 1<sup>st</sup> task is Show that K-DL is learnable. i.e., any f in K-DL can be approximated accurately after training on a reasonable no. of examples.

→ So, we need to calculate no. of hypothesis in lang.

→ Let language of tests - Conjunctions of at most K literals using n attributes - be Conj(n, k).

→ decision list is constructed of tests & b'coz each test can be attached to either Yes or NO or can be absent from decision list, there are at most  $3^{\lceil \text{Conj}(n, k) \rceil}$  distinct sets of component tests.

No. of Conjunctions of  $K$  literals from  $n$  attributes is given by

$$|\text{Conj}(n, k)| = \sum_{i=0}^k \binom{2^n}{i} = O(n^k).$$

$$|K\text{-DL}(n)| = 2^{O(n^k \log_2(n^k))}$$

To show that no. of examples needed for PAC-learning  
a  $K$ -DL  $f^n$  is polynomial in  $n$ :

$$n \geq \frac{1}{\epsilon} \left( \ln \frac{1}{\delta} + O(n^k \log_2(n^k)) \right)$$

→ Any algorithm that returns a consistent decision list  
will PAC-learn a  $K$ -DL function in a reasonable no. of  
examples, for small  $k$ .

→ Next task is to find an efficient algorithm that  
returns a consistent decision list.

→ We use a greedy algorithm called Decision-list-learning  
that repeatedly finds a test that  
agrees exactly with  
some subset of training set.

→ It finds such a test, it adds it to the decision list under construction & removes corresponding examples.

→ This is repeated until there are no examples left.

DECISION-LIST-LEARNING  
function Decision-list learning (examples) returns a decision list or failure

if examples is empty then return the trivial decision list NO  
+  $\leftarrow$  a test that matches a nonempty subset examples<sub>t</sub> of examples.

Such that the members of examples<sub>t</sub> are all positive or all negative

if there is no such t then return failure

if the examples in examples<sub>t</sub> are positive then

OK  $\leftarrow$  Yes else ~~OK~~  $\leftarrow$  NO.

return a decision list with initial test t and outcome  
0 and remaining tests given by

DECISION-LIST-LEARNING(examples - examples<sub>t</sub>).

### Algorithm for Learning Decision lists

- This algorithm doesn't specify method for Selecting next test to add to the decision list.
- Although the formal results given earlier do not depend on Selection method, it would seem reasonable to prefer small tests that match large sets of uniformly classified examples, so that overall decision list will be compact.

# Linear and Classification with Linear Models

- Here See

class of linear-fns of continuous valued fns.

- regression with a univariate linear fn of "fitting a straight line"
- Then Multivariate case
- Then How to turn linear-fns into classifiers by applying hard and soft thresholds.

## Univariate Linear Regression

- A univariate linear fn with i/p x & o/p y has form  
 $y = w_1 x + w_0$ . where  $w_0$  &  $w_1$  are real-valued coefficients.  
 to be learned
- we use  $w$  b'coz we think coefficients as weights.
- y is changed by changing relative weight of one term or another.
- we define  $w$  to be vector  $[w_0, w_1]$ . & define
 

$h_w(x) = w_1 x + w_0$ .
- figure shows an example of a training set of n points in the x, y plane, each point representing the size in Square feet & price of a house offered for sale.
- Task of finding the  $h_w$  that best fits these data is called linear regression.

- To fit a line to the data, ~~we~~ we have to do is find the values of weights  $[w_0, w_1]$  that minimize the empirical loss.

To use Squared loss f.,  $L_2$  Sum over all training examples.

$$\text{Loss}(h_w) = \sum_{j=1}^N L_2(y_j, h_w(x_j)) = \sum_{j=1}^N (y_j - h_w(x_j))^2 \\ = \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2.$$

$$w^* = \operatorname{arg\!min}_w \text{Loss}(h_w).$$

Sum  $\sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2$  is minimized when  $w_0$  &  $w_1$  are zero.

$$\left. \begin{array}{l} \frac{\partial}{\partial w_0} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0 \\ \frac{\partial}{\partial w_1} \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 = 0. \end{array} \right\} 18.2$$

equations have unique sol<sup>n</sup>.

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum x_j)^2} \quad \left. \right\} 18.3.$$

$$w_0 = (\sum y_j - w_1 (\sum x_j)) / N.$$

Sol<sup>n</sup> for figure  $w_1 = 0.232$ ,  $w_0 = 24.6$ . & line with those weights is shown as a dashed line

- Learning involve adjusting weights to minimize a loss function  
So, it helps to have a picture of what's going on in weight space
- Weight Space is defined by all possible setting of weights.

→ Alternatively, we can ...

1:

- For univariate linear regression, weight space defined by  $w_0 + w_1$  is two-dimensional.
- loss as a  $f^n$  of  $w_0 + w_1$ , in 3D plot. in figure 18.13(b).
- loss  $f^n$  is convex, it is true for every linear regression problem with an log loss  $f^n$  & implies that there are no local minima.
- If we need to fit lines to data, we apply (18.3.).

- Equations defining  $\min \text{loss}$  will often have no closed-form sol<sup>n</sup>.
- we use gradient descent.
  - we choose any starting point in weight space, here a point in  $(w_0, w_1)$  plane and then move to a neighboring point that is down hill, repeating until we converge on min possible loss.

$w \leftarrow$  any pt in parameter space

loop until convergence do

for each  $w_i$  in  $w$  do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(w).$$

- $\alpha$  is step size usually called the learning rate when we are trying to minimize loss in a learning problem.

99% of time.

→ For univariate regression, loss f^n is a quadratic f^n, so partial derivation will be a linear f^n.

on one training example (x, y)

$$\frac{\partial}{\partial w_i} \text{Loss}(w) = \frac{\partial}{\partial w_i} (y - h_w(x))^2$$

$$= 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - h_w(x))$$

$$= 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - (w_0 + w_1 x)) = 18.5$$

apply to both  $w_0$  &  $w_1$ , we get

$$\frac{\partial}{\partial w_0} \text{Loss}(w) = -2(y - h_w(x));$$

$$\frac{\partial}{\partial w_1} \text{Loss}(w) = -2(y - h_w(x)) \times x.$$

In 18.4. folding the 2 into unspecified learning rate, we get.

$$w_0 \leftarrow w_0 + \alpha (y - h_w(x)); w_1 \leftarrow w_1 + \alpha (y - h_w(x)) \times x.$$

if  $h_w(x) > y$  i.e., o/p of hypothesis is too large.

Reduce  $w_0$  a bit & reduce  $w_1$  if x was a +ve i/p.  
but increase  $w_1$  if x was a -ve i/p.

→ Derivative of sum is sum of derivatives.

$$w_0 \leftarrow w_0 + \alpha \sum_j (y_j - h_w(x_j));$$

$$w_1 \leftarrow w_1 + \alpha \sum_j (y_j - h_w(x_j)) \times x_j.$$

- > These update constitute batch gradient descent learning rule for univariate linear regression.
- > Convergence to unique global min is guaranteed, but may be very slow.
- i -> Another possibility, called Stochastic gradient descent, where we consider only a single training pt at a time,
- > It can be used in an online setting.

### Multivariate linear Regression

- each example  $x_j$  is an n-element vector.
  - hypothesis space is set of  $f^n$  of form
- $$h_{sw}(x_j) = w_0 + w_1 x_{j,1} + \dots + w_n x_{j,n} = w_0 + \sum_i w_i x_{j,i}$$
- 
- Matrix product of transpose of weights & i/p vector
- $$h_{sw}(x_j) = w \cdot x_j = w^T x_j = \sum_i w_i x_{j,i}$$
- 
- best vector of weights,  $w^*$ , minimizes Squared-error loss over exam
- $$w^* = \underset{w}{\operatorname{argmin}} \sum_j L_2(y_j, w \cdot x_j)$$
- 
- Multivariate linear regression is actually not much more complicated than univariate case, we covered.
  - Gradient descent will reach <sup>(univ)</sup> minimum of loss  $f^n$ , the update equation for each weight  $w_i$  is

$$w_i \leftarrow w_i + \alpha \sum_j x_{j,i} (y_j - h_w(x_j))$$

→ It is possible to solve for w that minimizes loss.

→ let y be vector of o/p's for training examples, & x be the data matrix, i.e., the matrix of i/p's with one-n-dimensional example per row.

Sol<sup>n</sup>.

$$w^* = (x^T x)^{-1} x^T y \text{ minimizes Squared Err.}$$

→ In multivariate linear regression in high-dimensional spaces, if irrelevant dimension appears by chance to be useful, which results in Overfitting.

→ It is common to use regularization on multivariate linear f<sup>n</sup>'s to avoid overfitting.

→ With regularization we minimize total cost of a hypothesis, counting both empirical loss & Complexity of hypothesis:

$$\text{Cost}(h) = \text{EmpLoss}(h) + \lambda \text{Complexity}(h).$$

$$\text{Complexity}(h_w) = L_2(w) = \sum_i |w_i|^2.$$

→  $\alpha=1$ ,  $L_1$  regularization, which minimizes the sum of absolute values.

→  $\alpha=2$ ,  $L_2$  regularization minimizes the sum of squares.

→ Regularization which we pick depends on specific problem,  $L_1$  has important adv., it tends to produce a Sparse model, it often sets many weights to zero, declaring corresponding attributes to be irrelevant.

- Minimizing  $\text{Loss}(w) + \lambda \text{Complexity}(w)$  is equivalent to minimizing  $\text{Loss}(w)$ . Subject to constraint that  $\text{Complexity}(w) \leq C$ .
- Diamond - Shaped box represents the set of points  $w$  in two dimensional weight space that have  $L_1$  complexity less than  $C$ , our  $\text{Sop}^h$  will have to be somewhere inside this box.
- Concentric ovals represent contours of  $\text{loss } f^n$  with min loss at center.
- we have to find pt in the box i.e., closest to minimum.
- for an arbitrary position of minimum & its contours, it will be common for corner of box to find its way closer to minimum, just b'coz corners are pointy.
- fig 18.16(b) represents a circle rather than a diamond.   
 is there is no reason for intersection to appear on one of the axes, thus  $L_2$  does not tend to produce zero weights.
- Result is that no. of examples required to find a good  $h$  is linear in no. of irrelevant features for  $L_2$ .
- but only logarithmic with  $L_1$  regularization.   
 Empirical evidence on many problems supports this analysis.

## Linear Classifiers with a hard threshold

- Linear fns can be used to do classification as well as regression.
- Example fig 18.15(a) Shows data points of two classes earthquakes & underground explosions.
- Each point is defined by 2 i/p values  $x_1$  &  $x_2$ , that refer to body & surface wave magnitudes Computed from Seismic Signal.
- Given these training data, the task of classification is to learn a hypothesis h that will take new  $(x_1, x_2)$  pts & return either 0 for earthquakes or 1 for explosions.
- Decision boundary is a line that Separates the two classes.
- In fig (a) ~~the~~ decision boundary is a straight line.
- A linear decision boundary is called a linear Separator & data that admit such a Separator are called linearly Separable.
- $x_2 = 1.7x_1 - 4.9$  or  $-4.9 + 1.7x_1 - x_2 = 0$ .
- Explosions, which we want to classify with value 1, are to right of this line with higher values of  $x_1$  & lower values of  $x_2$ ,
- So they are points for which  $-4.9 + 1.7x_1 - x_2 > 0$ , while earthquakes have  $-4.9 + 1.7x_1 - x_2 < 0$ .
- with dummy i/p  $x_0 = 1$ . we can write classification hypothesis as  $h_w(x) = 1$  if  $w \cdot x \geq 0$  and 0 otherwise.

→ Alternatively, we can think of  $h$  as result of passing  
the linear  $w \cdot x$  through a threshold  $f^h$ .

$$h_w(x) = \text{Threshold}(w \cdot x) \text{ where } \text{Threshold}(z) = 1 \text{ if } z \geq 0 \text{ or } 0 \text{ otherwise.}$$

→ <sup>Now</sup> Hypothesis  $h_w(x)$  has a well defined mathematical form,  
we can think about choosing the weights  $w$  to minimize  
the loss.

→ In 18.6.1 & 18.6.2, we did this both by setting  
the gradient to zero & solving for weights.

→ Here, we can't do either of those because the  
gradient is '0' almost everywhere in weight space except  
at those pts where  $w \cdot x = 0$ , & at those pts the  
gradient is undefined.

→ Linear Separator classifies the data perfectly - provided  
the data are linearly separable.

Ex:- Single example  $(x, y)$  we have

$$w_i \leftarrow w_i + \alpha (y - h_w(x)) x_i$$

the update rule for linear regression. This rule  
is called perceptron learning rule.

→ We are considering O/I classification problem.

Both true <sup>value</sup>  $y$  & hypothesis O/p  $h_w(x)$  are either 0 or 1.

So there are 3 possibilities.

- If O/p is correct i.e.,  $y = h_w(x)$ , then the weights are not changed.

- If  $y$  is 1 but  $h_w(x)$  is 0, then  $w_i$  is increased when corresponding i/p  $x_i$  is +ve & decreased when  $x_i$  is -ve.  
This makes sense, b'coz we want to make  $w \cdot x$  bigger so that  $h_w(x)$  o/p is a 1.
- If  $y$  is 0 but  $h_w(x)$  is 1, then  $w_i$  is decreased when corresponding i/p  $x_i$  is +ve & increased when  $x_i$  is -ve.  
b'coz we want to make  $w \cdot x$  smaller so that  $h_w(x)$  o/p is a 0.

→ Fig 18.16(a) shows a training curve for learning rule applied to earthquake/explosion data.

- Training Curve measures the classified performance on a fixed training set as learning process proceeds on that same training set.
- The curve shows update rule converging to a zero error linear separator.
- The "convergence" process is pretty.
- This particular run takes 657 steps to converge, for a data set with 63 examples, so each example is presented roughly 10 times on average.
- The perceptron learning rule converges to a perfect linear separator when data pts are linearly separable, what if they are not?
- This is too common in real world.

## → Linear Classification with Logistic Regression

→ Passing the O/p of a linear f<sup>n</sup> through threshold f<sup>n</sup>

Creates a linear classifier, yet hard nature of threshold  
Causes some problems.

→ The hypothesis  $h_w(x)$  is not differentiable & is in  
fact a discontinuous f<sup>n</sup> of its i/p & its weights.

→ This makes learning with perceptron rule a very  
unpredictable adventure.

→ Linear classifier always announces a completely confident  
prediction of 1 or 0, even for examples that are very  
close to boundary, we need more graduated predictions.

→ All these can be resolved to a large extent by  
softening threshold f<sup>n</sup> - approximating hard threshold with  
a continuous, differentiable f<sup>n</sup>.

logistic f<sup>n</sup> logistic(z) =  $\frac{1}{1+e^{-z}}$ . has convenient  
 mathematical properties.

→ The f<sup>n</sup> (8.17(b)) with logistic f<sup>n</sup> replacing threshold f<sup>n</sup>  
 we now have

$$h_w(x) = \text{logistic}(w \cdot x) = \frac{1}{1+e^{-w \cdot x}}$$

→ example of for two i/p earthquake/explosion problem  
 shown in 18.17(c).

→ Notice that O/p being a no. b/w 0 & 1, can be  
 interpreted as a probability of belonging to  
 class labeled 1.

- The process of fitting the weights of this model to minimize loss on a dataset is called logistic regression.
- There is no easy closed-form Sol<sup>n</sup> to find optimal value of  $w$  with this model. gradient descent is straightforward.
- B'coz our hypotheses no longer O/P just 0 or 1, we will use L2 loss f<sup>"</sup>
- g for logistic f<sup>1</sup>, g<sup>1</sup> its derivative.
- for a single example  $(x, y)$  the derivation of gradient is same as for linear regression up to the point where actual form of  $h$  is inserted.

we need chain rule.

$$\frac{\partial g(f(x))}{\partial x} = g'(f(x)) \frac{\partial f(x)}{\partial x}.$$

$$\begin{aligned}\frac{\partial}{\partial w_i} \text{Loss}(w) &= \frac{\partial}{\partial w_i} (y - h_w(x))^2 \\ &= 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - h_w(x)) \\ &= -2(y - h_w(x)) \times g'(w \cdot x) \times \frac{\partial}{\partial w_i} w \cdot x \\ &= -2(y - h_w(x)) \times g'(w \cdot x) \times x_i\end{aligned}$$

$g'$  of logistic f<sup>"</sup> satisfies  $g'(z) = g(z)(1-g(z))$ .

$$g'(w \cdot x) = g(w \cdot x)(1-g(w \cdot x)) = h_w(x)(1-h_w(x))$$

so, weight update for minimizing loss is

$$w_i \leftarrow w_i + \alpha (y - h_w(x)) \times h_w(x) \cdot (1 - h_w(x)) \times x_i$$