

$p_1 \rightarrow$ sketched
 $p_2 \rightarrow$ approx original
 \rightarrow in p_1 & p_2 both. $\sum |y - \hat{y}|$

$$(X^T X)^{-1} \sigma^2 = \text{var}(\hat{w}_1)$$

ssh our name @ 10.192.13.11

Model $\hat{y} = w_0 + w_1 x + \epsilon$ → data derived from this model
(st line with added noise)

Q1 How close are \hat{w}_1 & \hat{w}_0 to true values w_1 & w_0 ?

$$\hat{w}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{w}_0 = \bar{y} - \hat{w}_1 \bar{x}$$

Use variance formulae to compute confidence intervals:

$$\text{var}(\hat{w}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{var}(\hat{w}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

Use var formulae to derive how good are estimates \hat{w}_0 & \hat{w}_1 .

With approx 0.95 prob, the interval

$$[\hat{w}_1 - 2\sqrt{\text{var}(\hat{w}_1)}, \hat{w}_1 + 2\sqrt{\text{var}(\hat{w}_1)}]$$

If we change 0.95, this changes (interval len)

The smaller the variance, the more confident we are that w_1 is in this interval.

$y_i \rightarrow$ uncorrelated
var of $y_i \rightarrow$ constant $= \sigma^2$
 $x_i \rightarrow$ fixed (non-random)

$$\sum x_i^2$$

Page No. _____
Date / / 120

Approximation here it assumes normality.
Hastie & Tibshirani chapter 3.

Std error \rightarrow std deviation.

If std error is small, our estimate is good.

Use similar procedure to obtain confidence set for the parameter vector w [Hastie Tibshirani]

Q) How to check if there exists a relationship between Y and X ?
Ans Use hypothesis testing.

Credit Data set.

get summary of data set (min, 1st Qu, median, etc)

plot density also.

remove outliers (by plotting density & seeing some abnormality)

Bar chart also helps

histogram

Scatter plot of different variables

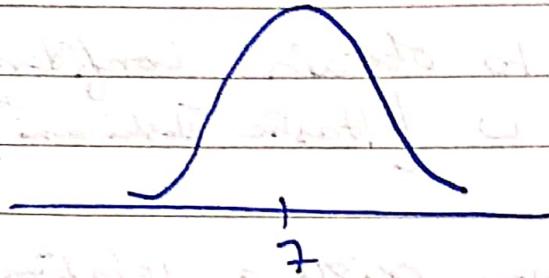
(looking at correlation plots, we get a good idea, what are the variables that are related)

Correlation, R^2 , $R^2 \rightarrow 0.99$

Corr of output y with each individual input variable

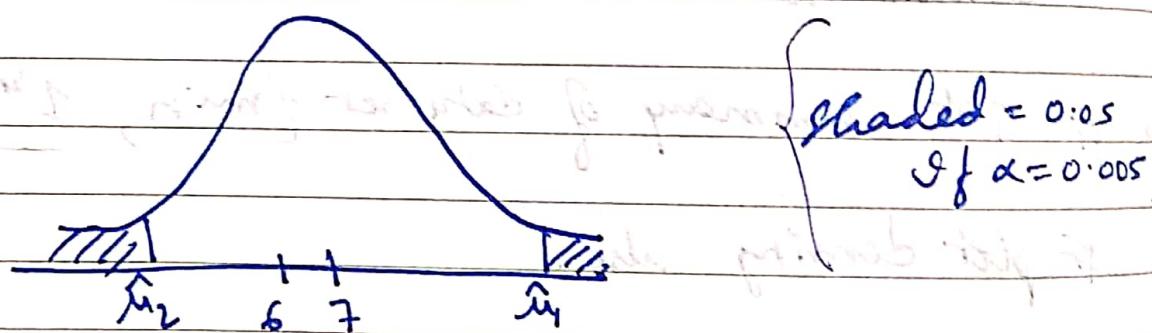
There is other way called hypothesis testing
(rather than drawing plots).

Our hypothesis was that a batch has avg gpa of any class is 7.



Now a new batch comes with avg gpa was 6. ~~was~~ (it rejected our hypothesis)

$\alpha = 0.05$ level of significance.
 $\alpha = 0.01$



If avg is in the shaded region, then our hypothesis cannot be rejected.

Hypothesis assumed to be true \rightarrow null hypothesis.

We need good amt of evidence to prove that null hypothesis is false.

o Null hypothesis, H_0 :- There is no relationship between X and y .

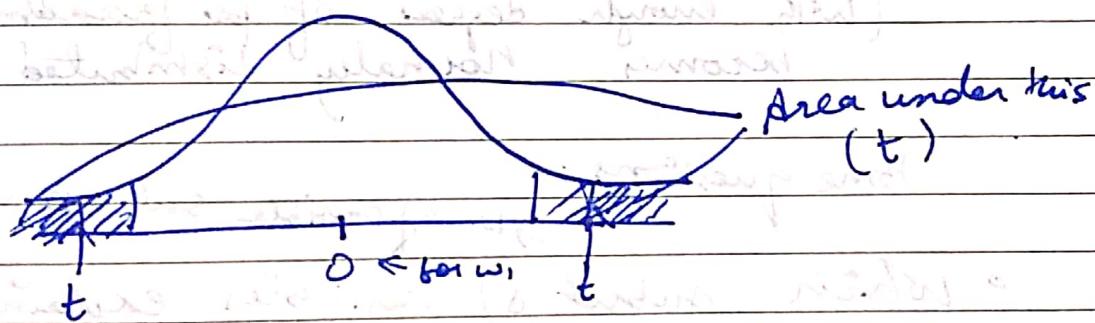
$$H_0: w_i = 0$$

$$\text{So, } y = w_0 + \beta_1 x \quad \text{Should vanish}$$

, Alternate hypothesis, H_a : There is some relationship between X and y .

$$H_a: w_i \neq 0$$

Did we get non zero w_i by chance or
is there ~~heavy~~ a relationship between y & x .



To test null hypothesis, we need to check if \hat{w}_i is sufficiently far from 0.

decided by level of significance.

(Everything true because of normality assumption)

Use t-statistic: $t = \frac{\hat{w}_i - 0}{\sqrt{\text{Var}(\hat{w}_i)}}$

What we are assuming w_i to be.

Either reject the null hypothesis
or do not reject.

Compute p-value using for large sample using
Normal distribution p-value: $P_{\text{H}}(w_1 \geq |t|)$

If p-value is less than the level of significance, then we can reject the null hypothesis. (Then $w_1 \neq 0$, then there exist a relationship between x & y)

A small p-value \Rightarrow reject null hypothesis

\Rightarrow association between predictor & response exists.

Here it is a two sided test.

Here we assumed w_1 is normally distributed.
(With enough degrees of freedom, it becomes normally distributed).

Some questions:

\rightarrow no of possible subsets very very large

° Which subset of variables explains the response?

° How good is the designed model? \rightarrow we can use least squares error for this.

$$w_0 + w_1 x(1) + w_2 x(2) + \dots + w_d x(d)$$

Null hyp \rightarrow all w_1, \dots, w_d are $\neq 0$.

Alternate " \rightarrow at least one of them is non-zero.

$I_{\text{fit}}^{\text{only}} w_0 = \text{non zero}$

$\Rightarrow w_0 \rightarrow \text{mean of target}$

Page No.
Date / /

Some practical approaches:

o Forward stepwise

\rightarrow gives suboptimal soln
(not best optimal soln)
(greedy algo)

Taking all subsets,
not partial.

- \rightarrow Start with only the intercept term in the model.
- \rightarrow Add the feature that improves the fit the most.
- \rightarrow Repeat until the above until some stopping condition satisfied.

keep checking if new feature reduces the least squared error.

If it changes error significantly, we keep it.
 \rightarrow when LSE does not improve.
when we get a flat line, we stop.
Every time we train the model.

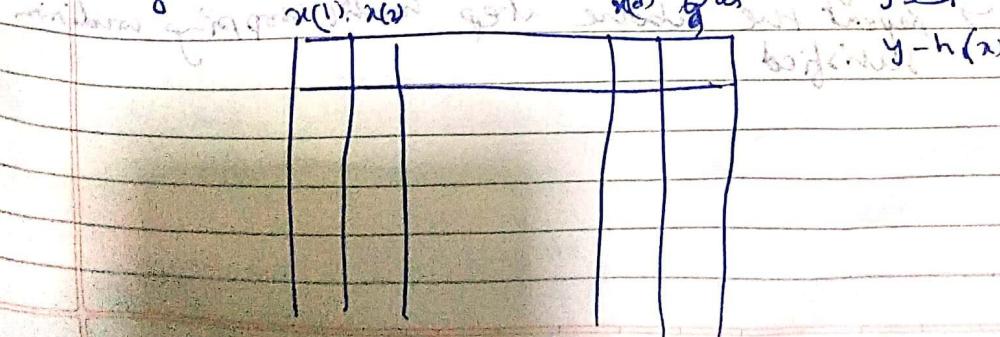
o Backward stepwise

\rightarrow Start with complete model.

\rightarrow Start Remove the feature having the least impact on the model performance.

\rightarrow Repeat the above step until some stopping condition is satisfied.

o Stagewise



1st \rightarrow which var correlated with y .

2nd \rightarrow now " " " with residual.

once $h_1(x)$ is designed, it will be fixed

$$y - h_1(x)$$

$$y - h_1(x) + h_2(x)$$

Since we are retraining, a variable previously selected may drop out
(Not in stagewise)

\rightarrow In forward stepwise \rightarrow everytime we add the feature, parameters change & we find correlation ~~with~~ between y & variables.

In forward stagewise \rightarrow

o Forward stagewise :-

\rightarrow Start with the only intercept term in the model

\rightarrow Identify the variable most correlated with the residual and find its linear regression coeff of the residual

\rightarrow Repeat the above step until stopping condition satisfied

Assume: All variables have 0 mean. (all the
variables have zero means)

Age: 11
Date: 1/12/01

Although computationally expensive gives good result

① First model (using intercept only): $h_1(x) = \bar{y}$, $t = 1$

② while Stopping condition is not satisfied
③ while \rightarrow Find the variable most correlated with
 $y - h_{[t]}(x)$, say k_t
all integers up to t .
 $\rightarrow h_{[t+1]}(x) = h_1(x) + h_2(x) + \dots + h_{k_t}(x)$

\rightarrow Design a regression using the data

$$\{(x_1(k_t), y_1 - h_{[t]}(x_1)), (x_2(k_t), y_2 - h_{[t]}(x_2)), \dots, (x_n(k_t), y_n - h_{[t]}(x_n))\}$$

to get

$$h_{[t+1]}(x) = h_{[t]}(x) + w_{k_t} x(k_t)$$

$$\rightarrow t = t + 1$$

verified that it gives very good result.

Endwhile

We introduce that variable that reduces the residual or keeps it very low. This stopping condition is met.

$$\text{res}_1 = y - h_1(x)$$

$$\text{res}_2 = y - h_1(x) - h_2(x)$$

a variable correlated to residual that reduces it.

$$\text{res}_t = y - h_{[t]}(x)$$

one dim optimization problem
(easy to solve)

If at any step, error

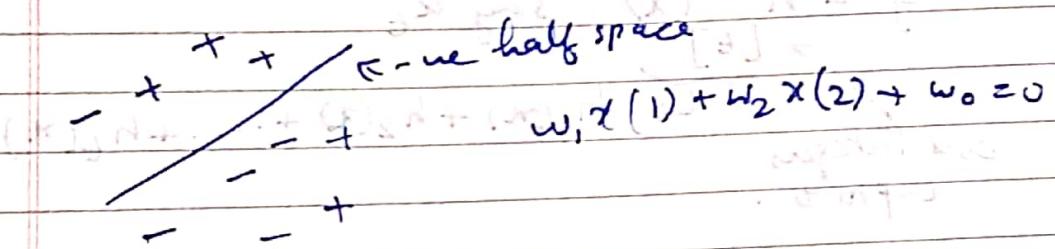
Mean of each col is assumed to be zero, so that finding correlation ~~is~~ between x_i & residuals is easy.

Other ways are ridge regression,

lasso "",

etc, that uses automatically selects the feature.

Decision Trees & Ensemble Methods

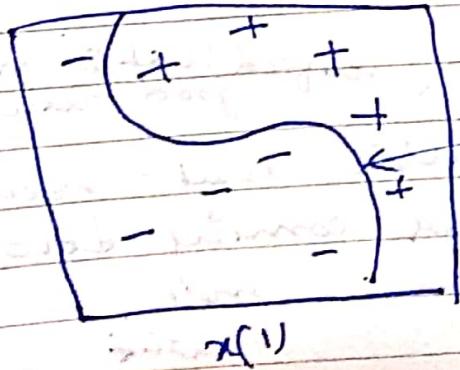


Decision fn: Linear

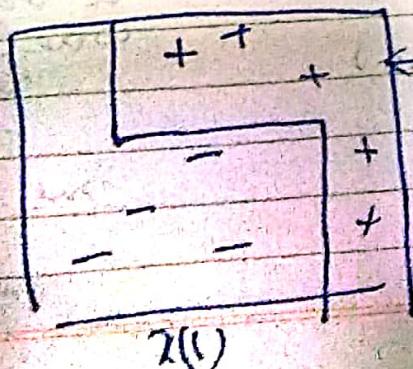
Here least squares can also be used. You will be + or -.

If an example plugged in the eqn, if it is +ve, it is in +ve class, if -ve, in -ve class.

This is not interpretable (\Rightarrow So decision trees used)



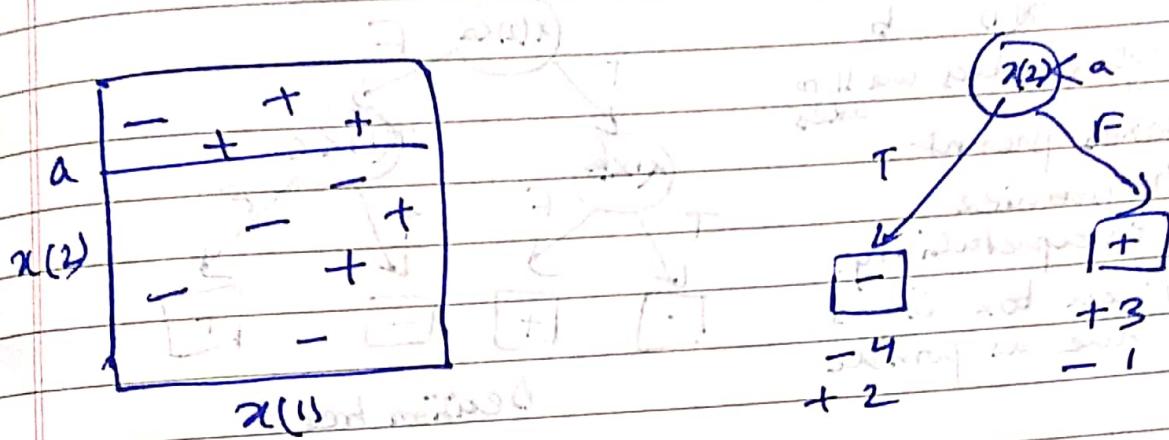
Non linear fn approximated by piecewise fns.



If blood pressure > something then abnormal
e.g. If $x(1)$ is > something
 $x(2)$ is > something
 $x(2)$ is abnormal.

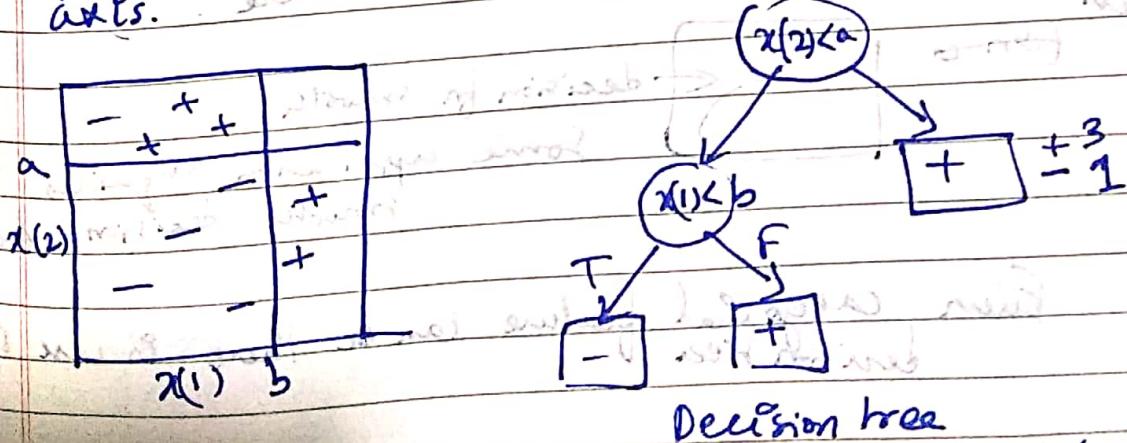
Decision tree other application - bank loan sanction.
Handy for Interpretability.

Decision stump \rightarrow height 1 decision tree



Instead of going downwards here we have made mistakes along the horizontal axis as well. (2 misclassified, 1 one misclassified)

Most Pimp thing is to get the parameter a .
Use any appropriate feature & split the Input space
Split is such that it is parallel to the coordinate axes.



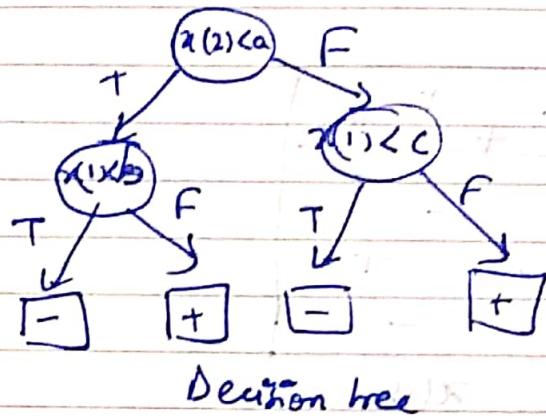
Leaf nodes \rightarrow pure nodes, (no example misclassified)

There is a class annotated with each leaf node.

		c	-	+	+	
	x(2)	a	-	-	-	
			-	-	-	
			-	-	-	
		b	-	-	-	

whose edges are \parallel to axes
 Boxes present to provide interpretability.
 Each box is as pure as possible

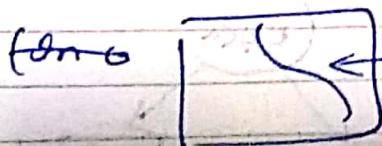
We divide input space into multiple regions such that each region is as pure as possible



Features used at nodes should be pure features
 (decision trees - a way to select imp features)
 "way of doing classification"
 can be used for multiclass classification.

Depends on application whether 2 or more branches should be there.

→ Decision surface is non smooth here.
deterministic process



decision for smooth.

Some application requires smooth decision surface.

Even categorical feature can be there in case of decision tree

We can stop at any level in decision tree.

Possibility to prune some parts of tree.
 (then some features will become irrelevant)

By pruning, we reduce the height of tree.

Page 120
Date: 1/1/20

Sometimes tree is grown from max possible depth & then start pruning \Rightarrow until one feature feels.

Categorical variables - take specific values.

In case previously we ^{were} worked with only real or integer value data.

Here with categorical.

Normalization ^{not} reqd \rightarrow before decision tree is created.

Decision tree can be used for both classification & regression problem.

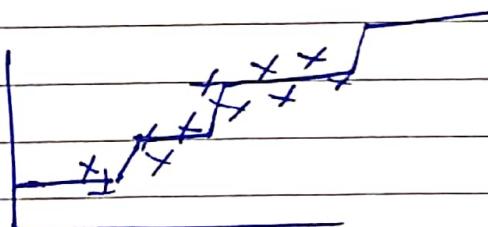
\hookrightarrow divide input region into multiple regions, & treat each region as ^{having} constant fn.

(Discontinuity may be there) At that part of region things may not be continuous.

We could have used non linear fn. that could have been smoother.

\uparrow Disadv of using decision tree?

Decision tree



A decision tree is a hierarchical structure, with each node splitting the data space into disjoint partitions based on feature value.

- With region boundary $\rightarrow \{R_1, R_2, \dots, R_m\}$
- Input space X is partitioned into disjoint regions, e.g.
 - Boundaries of regions parallel to coordinate axes in classical decision trees.

Choose a feature that will reduce the Entropy
(entropy as low as possible).

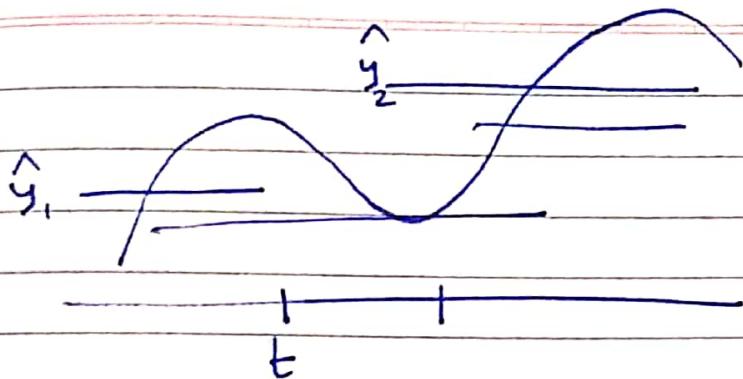
Initially S had some entropy.

We want to find a feature that minimizes the max expected reduction in entropy

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{re Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

expected entropy
based on feature A.

expected reduction.



among all possible t 's we find one that gives least value of obj fn.

This problem is hard, so we divide into regions and then find the mean to find t^* .

We divide into regions by clustering.

In case of least squares error, we keep on reducing error.

But in case of decision tree, there is ^{no} backtracking.

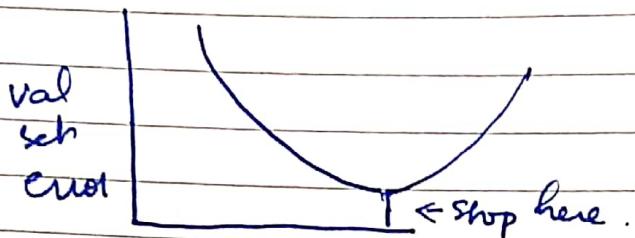
Early stopping possible or pos-i pruning.
(to avoid overfitting)

If in case of validation set performs well at any stage, we can stop.
i.e. we reached a flat state

Validation set used to know when to stop training

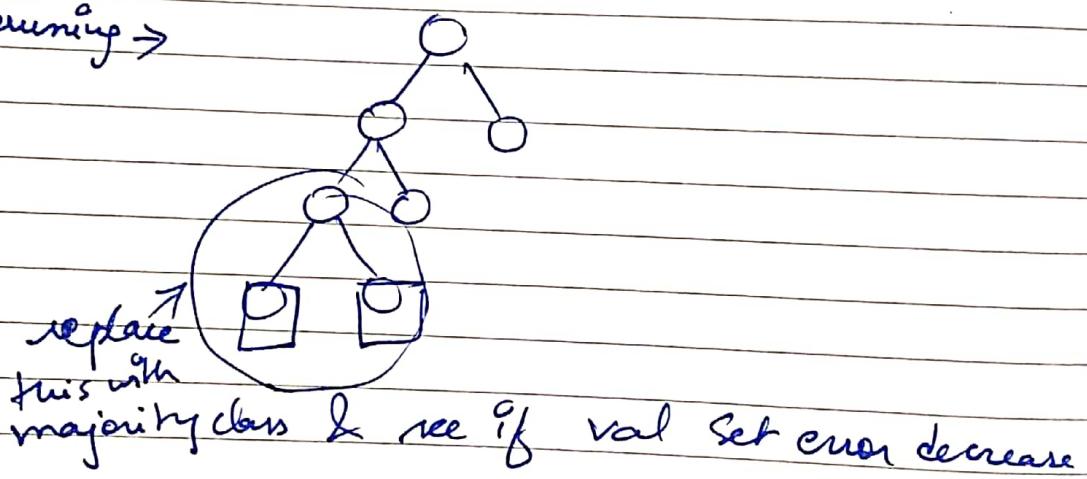
Train	Val	Test
-------	-----	------

Early stopping → As we grow the tree, we test the performance on validation set



Performance on validation set is crucial.
(for early)

Post pruning →



→ Continuous-valued features can be handled by converting them to discrete-valued features.

—	X X X	0 0 0 0	*** —
---	-------	---------	-------

→ Possible to handle training set examples with missing feature values.

→ One can be use mean

→ Other can be put them with some probability. probability of them being present based on some previous experience

→ Use of alternate measures for finding the "best" feature to split on:

→ Gini Index: $2P_0P_1$

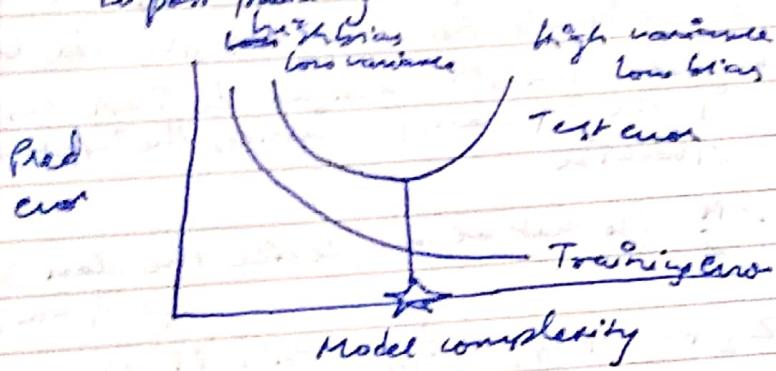
→ Misclassification error: $1 - \max(P_0, P_1)$

→ Squared error (for regression): $\sum_{m=1}^M \sum_{z_i \in R_m} (y_i - g_m)^2$

Banana dataset (2 dim)

Decision tree classifier (overfitting)

→ Do post pruning



Ensemble Methods

- Make an informed decision after consulting experts having diverse domain knowledge
 - Admission to university
 - Making investment in a company's stocks.

- "Large groups of people are smarter than an elite few"
[the wisdom of crowds]

- When are the crowds wise?

→ Diversity of opinion

→ Independence (unbiased)

→ Decentralization (every person has access to some private info not accessible to others)

→ Aggregation

- Are these ideas useful in m/c learning?
 - Combining the decisions of a diverse set of learners
 - Committee m/cs - Boosting, mixture of experts.
- $D = \{(x_i, y_i)\}_{i=1}^n, y_i \in \{-1, +1\}$
- M classifiers trained independently on different views of D .
- $P = \Pr(\text{a classifier correctly classifies an example})$.
- Assumption: For every classifier, the same p holds and $p > 0.5$ (otherwise random classifier will be better)
- $Z_i: R^d$ associated with i th classifier defined as

$$Z_i = \begin{cases} 1 & \text{if the } i\text{th classifier classifies the input correctly} \\ 0 & \text{otherwise.} \end{cases}$$

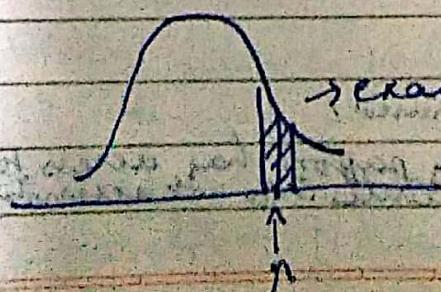
~~At least~~ $Z > \frac{M}{2}$ (so that we can decide the class)
no. of trials (n)

- $S_M = Z_1 + Z_2 + \dots + Z_M$ ($S_M \sim \text{Bin}(M, p)$)
- Prediction: Feed input x to all M classifiers and take a simple majority vote (need at least $\lceil \frac{M+1}{2} \rceil$ votes for correct prediction)
prob of getting class +1
 (all Z s are i.i.d.s)

$$\checkmark P(S_M = n) = \binom{M}{n} p^n (1-p)^{M-n}$$

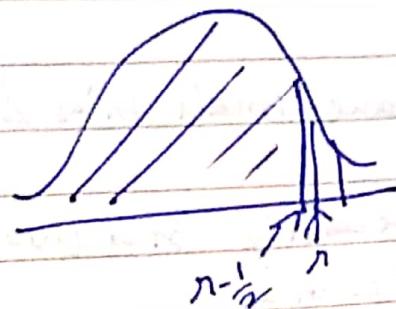
$= \Phi\left(\frac{n + \frac{1}{2} - Mp}{\sqrt{Mp(1-p)}}\right) - \Phi\left(\frac{n - \frac{1}{2} - Mp}{\sqrt{Mp(1-p)}}\right)$

Area under the curve (assuming normal distribution).
 exact value at a pH = 0
 so we take area under the slab around n .



- Probability of getting at least r votes:

$$P_n(S_M \geq r) \approx 1 - \phi\left(\frac{r - \frac{1}{2} - M_p}{\sqrt{M_p(1-p)}}\right)$$



From normal tables
we can
get these
values.

- Find P_n (correct classification of n using majority vote).

Some examples:

- Let $M = 15$, $p = 0.6$

we want at least 8 classified
getting giving
correct

$$P_n(S_M \geq 8) \approx 1 - \phi\left(\frac{6 - \frac{1}{2} - 9}{\sqrt{(15) \cdot (0.24)}}\right) = 0.785$$

- Let $M = 51$, $p = 0.6$

$$P_n(S_M \geq 26) \approx 1 - \phi\left(\frac{26 - \frac{1}{2} - 30 \cdot 0.6}{\sqrt{(51) \cdot (0.24)}}\right)$$

Each one doing just better than random classifier.

Remarks:

- As M increases, $P_n(S_M \geq \lceil \frac{M+1}{2} \rceil)$ increases slightly
(rough rate of increase decreases slowly).

- $p > 0.5 + \delta$ where δ is small & the
could be even 0.51 also (just has to be
better than random classifier)

- Weak classifiers: Perform slightly better than random
classifier
 - e.g. decision stump (depth 1 or height 1)
(just better than random classifier)

- Possible to combine weak classifiers so as to boost their performance.

Decision tree

Ensemble Methods

- Sufficiently Deep trees have small bias & high variance.

- Reduce variance w/o increasing bias (too much).

Decision Stump + high bias

Decision tree to max depth \rightarrow high variance.

Idea: Average reduces variance

Let Z_1, Z_2, \dots, Z_M be i.i.d R.V.s such that

$$\text{Var}(Z_i) = \sigma^2 \ \forall i$$

Define $Z = \frac{1}{M} \sum_{j=1}^M Z_j$. Then

$$\text{Var}(Z) = \frac{\sigma^2}{M} = \frac{\frac{1}{M} (\text{Var} Z_1 + \dots + \text{Var} Z_M)}{M} = \frac{\sigma^2}{M}$$

- Idea: Train multiple models & average their outputs.

Given one train data, only one decision tree possible

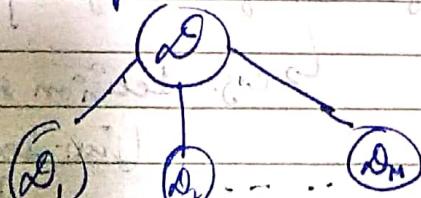
How to train decision trees

(decision trees)

- Given one training set D , how do we train multiple models?

Soln: Generate M bootstrap samples from D and train a different component of a base model using each bootstrap sample

boot samples
draw randomly
done randomly
with replacement



- Bootstrap

\rightarrow draw

D

\rightarrow e.g.

• Boo

• Boo

• Boo

• and

we can

\rightarrow Rema

(1)

a

o F

\rightarrow U

o Bootstrap Sample :

→ draw a sample of size $|D|$ with replacement from D

$$\rightarrow \text{e.g. } D = \{x_1, x_2, x_3, x_4, x_5\}$$

o Bootstrap sample 1: x_3, x_2, x_3, x_1, x_3

$$D_1 = \{x_1, x_2, x_3\}$$

(every sample may
not have
all examples)

• Bootstrap sample 2: $x_1, x_5, x_2, x_4, x_4; D_2 = \{x_1, x_2, x_4, x_5\}$

• Bootstrap sample 3: $x_3, x_4, x_4, x_3, x_4; D_3 = \{x_3, x_4\}$.

• and so on....

we can generate as many bootstrap samples as we want.

→ Remarks:

o Each model trained on less data

$\left(1 - \frac{1}{n}\right)^n \rightarrow$ prob of an example not being picked

as $n \rightarrow \infty$, $\left(1 - \frac{1}{n}\right)^n \rightarrow \frac{1}{e} \approx 33\%$.

Every bootstrap sample will be picked with
65% prob.

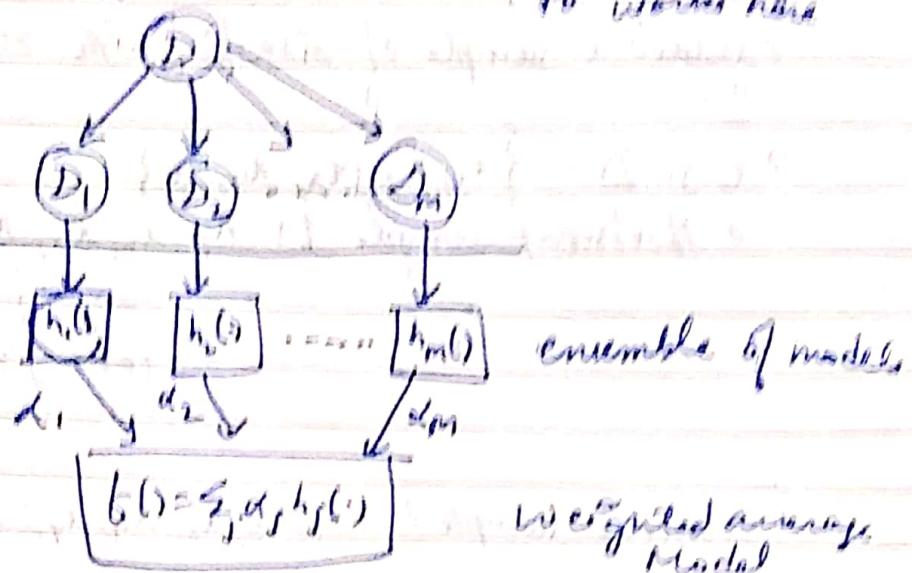
o Every data sample is seen by at least one of
the models. If M is reasonably large.

→ Use weighted average model for prediction.

Bootstrap Aggregation - Bagging

decision trees unstable
models

Bagging only unstable models (high variance)
to work hard



Works well if
base models are
unstable
(e.g. decision trees)

- Training Data: $D = \{(x_i, y_i)\}_{i=1}^n, x_i \in \mathcal{X}, y_i \in \mathcal{Y}$
- Single model learning: Design a $h: \mathcal{X} \rightarrow \mathcal{Y}$ using D
- Ensemble learning: Design an ensemble $\{h_1(\cdot), h_2(\cdot), \dots, h_m(\cdot)\}$ using $\{D_1, D_2, \dots, D_m\}$ respectively where.

$$h_j: \mathcal{X} \rightarrow \mathcal{Y}, j=1, \dots, M$$

- Inference for test input X
 - binary classification

$$\hat{y} = \text{sign}\left(\sum_j \alpha_j h_j(x)\right)$$

$$\rightarrow \text{Regression } \hat{y} = \sum_j \alpha_j h_j(x)$$



tree's unstable
models.
models (high variance)
works here

ensemble of models

Weighted average
Model
used for prediction.

$x \in \mathcal{X}, y \in \mathcal{Y}$

* using

$\{h_1(\cdot), h_2(\cdot), \dots, h_M(\cdot)\}$

where

Bagging

Algorithm

Input : Data $D = \{(x_i, y_i)\}_{i=1}^n$, M: No of models to be built

(1) for $j=1, \dots, M$

(a) Take a bootstrap sample D_j of size n from D

(b) Build a model $h_j(\cdot)$ using D_j .
end for

Output: A set of M models

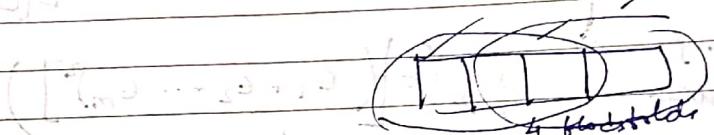
$$H = \{h_1(\cdot), h_2(\cdot), \dots, h_M(\cdot)\}$$

Inference : For an unseen example x , predicted label \hat{y} is calculated as

$$\circ \text{Regression: } \hat{y} = \frac{1}{M} \sum_{j=1}^M h_j(x)$$

$$\circ \text{Classification: } \hat{y} = \arg \max_i |C_i|$$

(where C_i)



Bagging

- o Bagging with k-fold cross-validation
 - train a model on $(k-1)$ folds with early stop on the remaining fold.
 - k models available at the end of k-fold cross-validation

Concern the regression data set derived from the true function $f(x)$ and noise ϵ .

o Suppose we have M bootstrap data sets & train M

of models, $y_M(x)$ be the outputs from

$$\circ y_M(x) = f(x) + \epsilon_m$$

o Average error, made by M models acting individually, is

If we have M independent models, avg expected error

$$\begin{aligned} E_{\text{avg}} &= \frac{1}{M} \sum_{m=1}^M E_x[(y_m(x) - f(x))^2] \\ &= \frac{1}{M} \sum_{m=1}^M E_x[\epsilon_m(x)^2] = \frac{1}{M} (E[\epsilon_1^2] + \dots + E[\epsilon_M^2]) \end{aligned}$$

- The output of the bagged model is

$$y_{\text{bag}} = \frac{1}{M} \sum_{m=1}^M y_m(x)$$

- The expected error from the bagged model is

$$\begin{aligned} E_{\text{bag.}} &= E_x[(y_{\text{bag}} - f(x))^2] \\ &= \frac{1}{M^2} E_x\left[\left(\sum_{m=1}^M \epsilon_m(x)\right)^2\right] \\ &= \frac{1}{M} E_{\text{avg}} \quad (\text{for uncorrelated errors}) \end{aligned}$$

Average expectation
with bagging
reduction

$$\frac{1}{M^2} (E[(\epsilon_1 + \epsilon_2 + \dots + \epsilon_m)^2])$$

No guarantee that we get uncorrelated examples in D_1, D_2, \dots, D_m because bootstrap samples are not totally unrelated.

Cross terms will also be there.
(If errors are uncorrelated, all cross terms will vanish.)

- Models are uncorrelated and so variance reduction is smaller than $\frac{1}{M}$ (difficult to completely uncorrelate training sets)

- Not suitable for stable learning methods (as they typically have low variance) (A model is unstable if small changes in the data cause changes in the predicted values)

- Has little effect on bias.
useful technique for large data sets & unstable models.

Bootstrap samples are actually correlated

Page No. _____
Date / / 120

- Models are correlated \rightarrow variance reduction is smaller.

ρ : correlation between every pair of models.

σ^2 : variance of each model.

$$\begin{aligned}\text{Var}\left(\frac{1}{M} \sum_{m=1}^M Z_m\right) &= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \text{cov}(Z_i, Z_j) \\ &= \frac{1}{M^2} \sum_{i=1}^M \left(\sum_{j=1, j \neq i}^M \text{cov}(Z_i, Z_j) + \text{var}(Z_i) \right) \\ &= \frac{1}{M^2} \sum_{i=1}^M ((M-1)\rho\sigma^2 + \sigma^2) \\ &= M \frac{(M-1)\rho\sigma^2 + M\sigma^2}{M^2} \\ &= \underbrace{\rho\sigma^2}_{M} + \underbrace{(1-\rho)\sigma^2}_{\text{decrease if } \rho \downarrow} \xrightarrow{\text{as } M \uparrow} \text{as } MT\end{aligned}$$

- If we use model averaging for variance reduction, we want

\rightarrow large no. of ~~most~~ models

\rightarrow low correlation between them.

Random forest

\hookrightarrow At every node look at a random subset of features from which we choose the best feature.

If we'll make sure that models are uncorrelated to each other.

- Idea: Use decision tree with different set of features
 \rightarrow reduce correlation between models.

$Q \leftarrow$ set take a set of n features & take best feature

$Q \leftarrow$ at every node, take a random feature subset
among them find the best

Approach: when learning a split, find the best feature
on a reduced random feature subset.

o Random Forest

→ Bootstrapping

→ Random Selection of Features

→ Ensemble of decision trees.

Algorithm:

All undefined

{ Input: d -dimensional data $D = \{(x_i, y_i)\}_{i=1}^n$, M : Number of
models to be built, m : No of features used at each node,
 $m+d$: maximum tree depth

$|D_j| < |D|$

(i) for $j = 1, \dots, M$

(a) Take a bootstrap sample D_j of size n from D

(b) Grow a tree T_j using D_j by recursively repeating
the following steps at each node of the tree,
until the min tree depth $m+d$ is reached.

(i) Draw a random subset Q of size m from

$\{1, 2, \dots, d\}$

(Ten ~~two~~)

(ii) Pick the "best" variable from Q (based on
information gain)

(iii) Split the node into appropriate no of children node

end for.

(Helps in generating better uncorrelated trees than the
simple bagging model)

Output: A random forest of M Decision Trees

$RF = \{T_1(\cdot), T_2(\cdot), \dots, T_M(\cdot)\}$

Combine
Every tree
disadv
here Pkt
if int

In single
axes.

Here
But

To d

→ keep
perf

→ Ca

combine RF just as in bagging.

Every tree generated, we see different set of features.

Disadv

Here Interpretability part of traditional tree is lost. But if Interpretability is not needed this can be used.

In single decision tree all test boundaries are \parallel^l to axes.

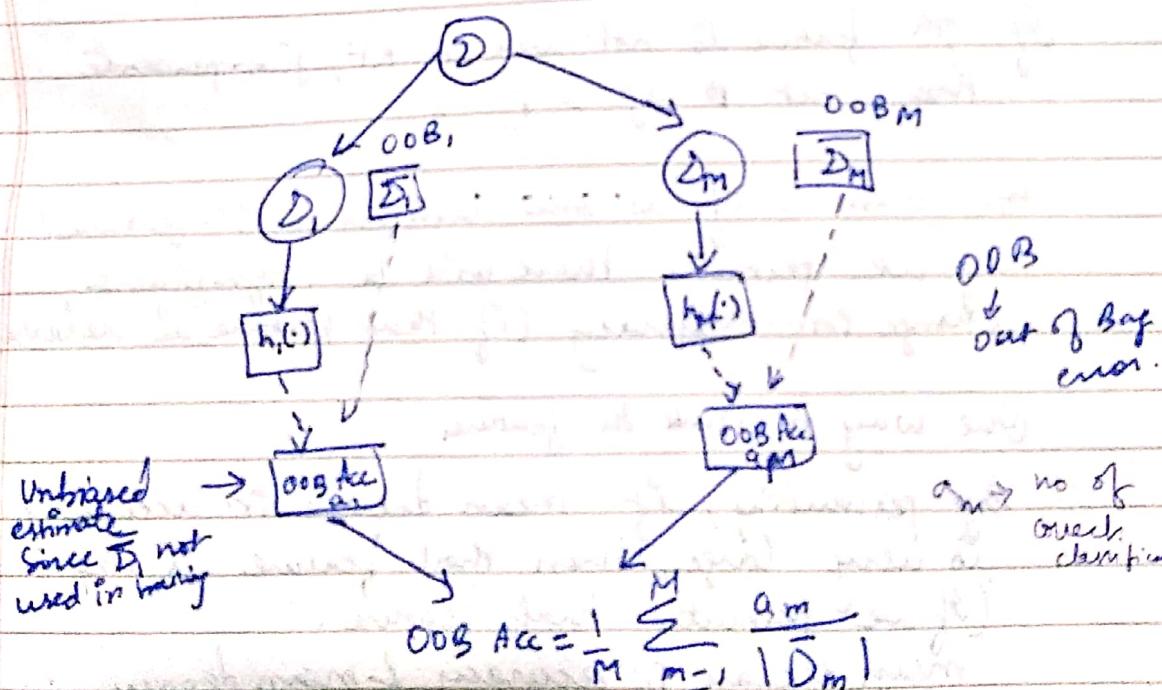
Here need not be \parallel^l . So interpretability part lost.

But it is still a good model.

To decide how many trees in Random Forest?

→ keep a Validation set. keep adding a model until performance deteriorates.

→ Can also work with only training set.



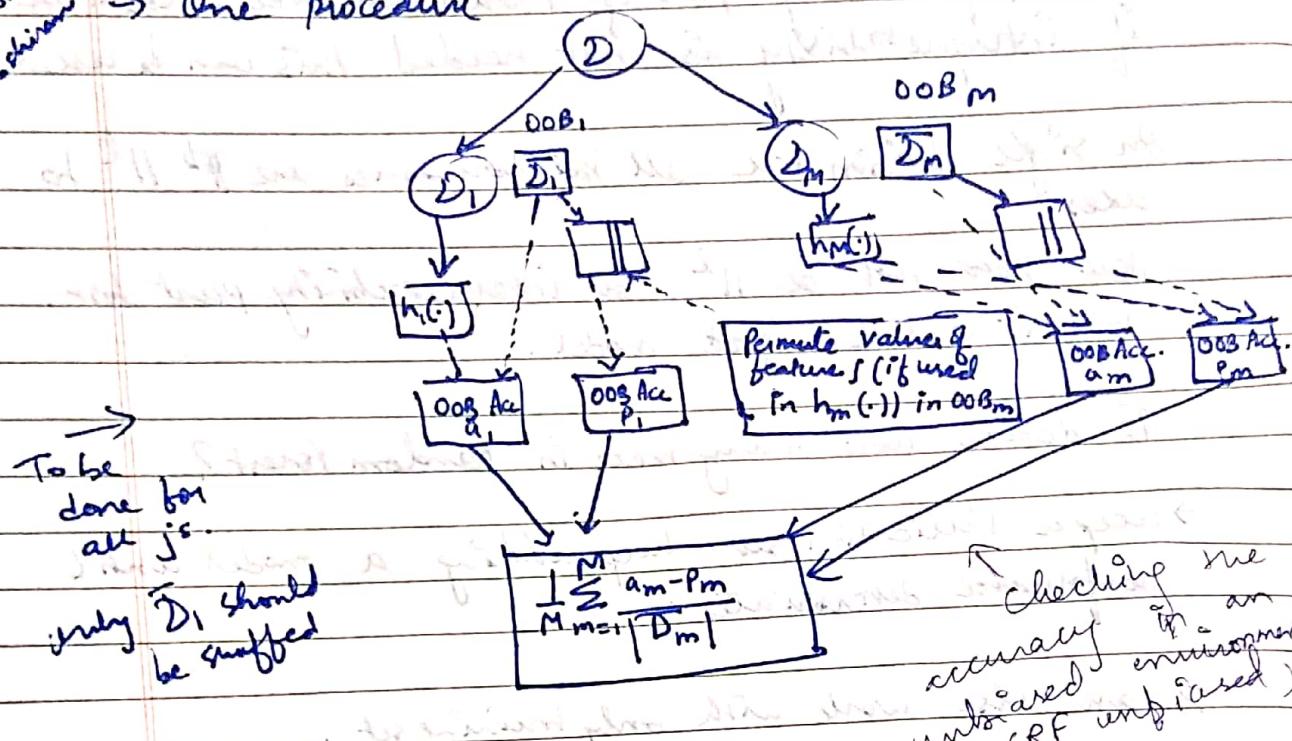
Whenever we add a new model, check OOB Acc when OOB Acc stops increasing, stop.

2nd assignment: build a random forest regression on some dataset.
no of trees, features to use, min tree depth, etc tune them

Random forest can also be used for feature selection.

How to find whether features is good or not?

Hastie
Tibshirani → One procedure



If j th feature is not used in RF, just permute this set θ a_j & $p_j = 0$.

How sensitive is our classifier to j th feature?

If we permute there will be significant change in accuracy (if that feature is relevant).

One way to rank the features.

By permuting, if mean decrease in accuracy is very large, then that feature is imp.
(If we permute that feature.)

mean decrease in accuracy & mean decrease in gini index.

feature selection \rightarrow Lasso, PCA, decision tree, random forest
Q Why RF is used for feature selection?
for feature selection?

Page No.
Date / / 120

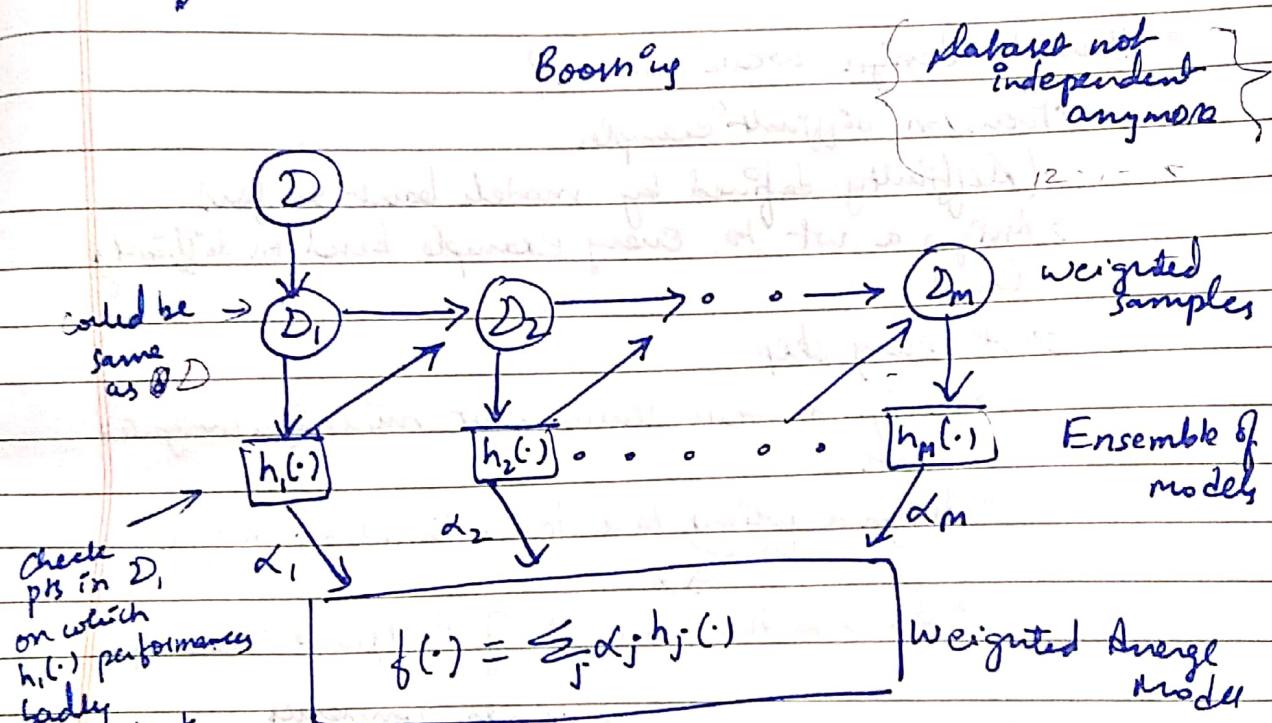
Many times ppl use RF to select features & then give it to linear model.

- Bagging reduces variance by averaging
- Is it possible to average models & reduce bias?

Solution: Boosting

Average weak models & then boost the performance.
High biased

Boosting \rightarrow similar to bagging.



Stage wise approach.
Once α, h , decided, they become fixed.

ie learning

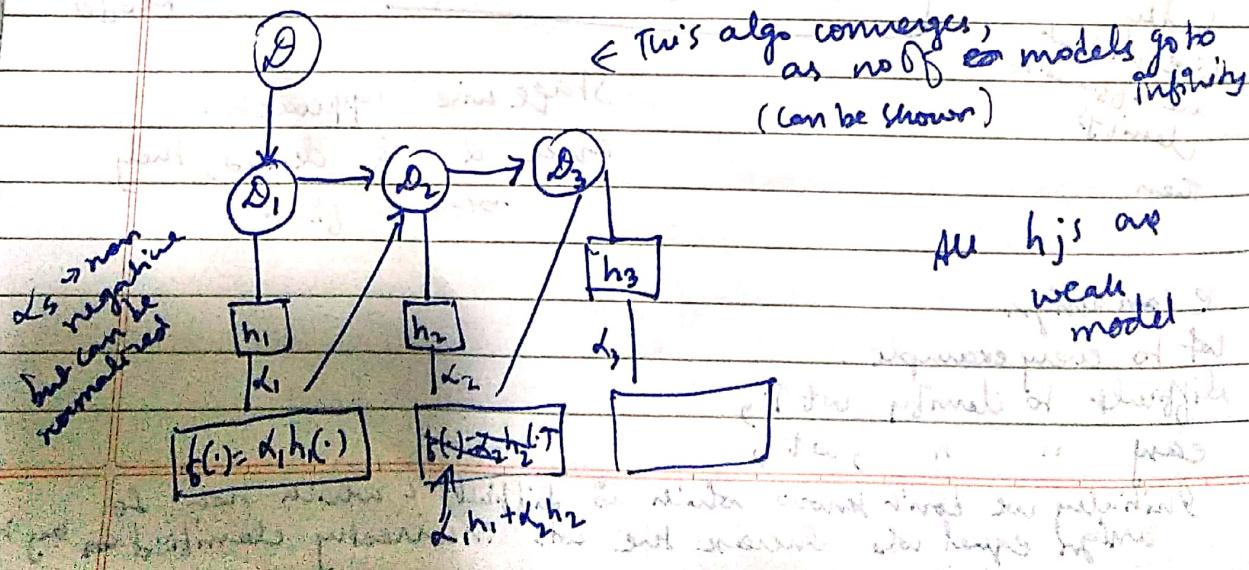
let's see every example.

Difficult to classify, wt \uparrow ,
easy " " , wt \downarrow

Initially we don't know which is difficult & which not! So assign equal wts. Increase the wts incorrectly classified by $h_i(\cdot)$

We should have ^a classifier, that minimizes the weighted misclassification error.

- ° Weak Learner: Performance on any training set is slightly better than random predictions
e.g. Stump.
- ° Designing a strong learner from many weak learners.
- ° Weak learners relatively easy to design than strong learners.
- ° How to design weak learners?
 - Focus on "difficult" examples.
 - (Difficulty defined by models learnt so far).
 - Assign a wt to every example based on difficulty
 - At every step
 - Design a weak learner that minimizes weighted error.
 - Assign a weight to a learner based on weighted error
 - Increase the weights of "difficult" examples.



Schapire & Freund

is introduction to boosting (paper).

Page No. _____
Date / / 120

- Decrease the lots of "easy" examples
 - Construct a strong learner using weighted combination of weak learners
 - Final model: weighted prediction of weak models
- abs no of examples ↑, training set error asymptotically goes to 0.

Goal: Design a classifier using a linear combination of weak classifiers, $h_j(\cdot)$.

$$f(x) = \sum_{j=1}^m \alpha_j h_j(x) \text{ where } h_j(x) \in \{-1, +1\}.$$

Q: How do we determine α_j 's & h_j 's given $D \in M$?

↑
no of models to
be used.

$$h_{[j-1]}(x) = \sum_{k=1}^{j-1} \alpha_k h_k(x)$$

Given we have already designed $j-1$ classifiers,
how to find α_j & h_j .

$$\begin{matrix} x_1, \dots, x_n \\ w_1 & w_2 & \dots & w_n \end{matrix}$$

Superscript - iteration no
is determined using combined model.

$$h_{[j]}(x) = \underbrace{h_{[j-1]}(x)}_{\text{determined & fixed}} + \alpha_j h_j(x)$$

determined & fixed
after $(j-1)$ iterations.

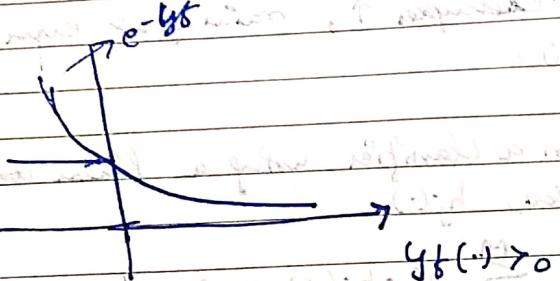
Now how to get D_j , h_j & α_j .

h_2 is designed using D_2 .
But $(\alpha_1 h_1 + \alpha_2 h_2)$ used to test on D_2 .

Determine α_j & $h_j(\cdot)$ by optimizing misclassification loss or its surrogate fn.

An example is correctly classified if $y_f > 0$.

$\begin{cases} 0-1 \\ \text{loss} \end{cases}$ loss = 1
misclassification
loss = 0
for classification



This is discrete loss.

Not continuous so difficult to optimize.

So surrogate loss fn used.

We can use $\approx e^{-yf}$

For example, using exponential loss fn,

$$l(y, h_j(x)) = e^{-y h_j(x)}$$

$$\text{So } w_{j-1}(x_i; h_j) = \underset{\alpha, h}{\operatorname{argmin}} \sum_{i=1}^n e^{-y_i (h_{j-1}(x_i) + \alpha h(x_i))}$$

$$\text{So } w_{j-1} = \underset{\alpha, h}{\operatorname{argmin}} \sum_{i=1}^n e^{-y_i h_{j-1}(x_i)} \cdot e^{-\alpha h(x_i)}$$

w_{j-1} ↑ positive quantity

where w_{j-1} is the wt associated with i th example after $G-1$ iterations.

keep one
the other
Good idea
keeping the
rest the

keep one part fixed & first do optimization w.r.t
the other part.

Good idea to optimize α & $h(\cdot)$ independently while
keeping the other fixed.

Let's let w first fix $\alpha \gg 0$ & solve the following problem:

$$h_j^* = \arg \min_h \sum_{i=1}^n w_i^{j-1} e^{-y_i \alpha h(x_i)}$$

$$= \arg \min_h \sum_{\substack{i: y_i = h(x_i)}} w_i^{j-1} e^{-\alpha} + \sum_{\substack{i: y_i \neq h(x_i)}} w_i^{j-1} e^{\alpha}$$

$$= \arg \min_h \underbrace{\sum_{i=1} w_i^{j-1} e^{-\alpha}}_{\text{Independent of } h} + \sum_{\substack{i: y_i \neq h(x_i)}} (e^\alpha - e^{-\alpha}) w_i^{j-1}$$

\downarrow
 \downarrow

$$= \arg \min_h \sum_{i=1} (e^\alpha - e^{-\alpha}) w_i^{j-1} I(y_i \neq h(x_i))$$

$$= \arg \min_h \sum_{i=1} w_i^{j-1} I(h_i \neq h(x_i))$$

\downarrow
weighted misclassification (h) \rightarrow classifier that

designs a classifier that minimizes weighted misclassification error.

Find $h_j(\cdot)$ which minimizes weighted misclassification error.

Design decision stump such that to reduce weighted misclassification error.

Since every classifier is a weak classifier, they ^{all} make non zero error

- Weighted misclassification error of $h_j(\cdot)$:

$$\epsilon_j = \sum_{i=1}^n w_i^{j-1} I(y_i \neq h_j(x_i))$$

error incurred by h_j on D_j
 (calculated once we design the classifier.)

Now, keep $h_j(\cdot)$ fixed & determine α_j .

$$\alpha_j = \underset{\alpha}{\operatorname{argmin}} \sum_{i=1}^n w_i^{j-1} e^{-\alpha} + \sum_{i:y_i \neq h_j(x_i)} (e^\alpha - e^{-\alpha}) w_i^{j-1}$$

weights missclassified examples
 are normalized. sum = 1.

$$= \underset{\alpha}{\operatorname{argmin}} e^{-\alpha} + (e^\alpha - e^{-\alpha}) \epsilon_j \quad (\text{letting } \sum_i w_i^{j-1} = 1)$$

$$= \frac{1}{2} \ln \frac{1 - \epsilon_j}{\epsilon_j} \quad \text{diff with respect to } \alpha$$

α_j decided based on weighted error incurred by h_j classifier
 if $\epsilon \uparrow$, $\alpha \downarrow$

What if $\epsilon_j > \frac{1}{2}$? $\rightarrow h_j$ is worse than random classifier.

so, we want $0 < \epsilon_j < \frac{1}{2}$

so, $\alpha_j > 0$

α_j just has to be a true quantity.

sum of α_j 's need not be 1,
 because α_j 's are decided by ϵ_j 's

Input:
 Initialize:

(1)

$$w_i^j = \cancel{e^{-y_i h_j(x_i)}} e^{-y_i h_j(x_i)}$$

$$\text{weights of } i\text{th example} = w_i^{j-1} e^{-y_i h_j(x_i)}$$

in j th iteration

if $y_i = h_j(x_i)$

\because if both +1 or both -1

then it will be

$$w_i^{j-1} e^{-\alpha_j}$$

↑ weights of misclassified examples pulled down

$$= \begin{cases} w_i^{j-1} e^{-\alpha_j} & \text{if } y_i \neq h_j(x_i) \text{ (wts if correct predicting)} \\ w_i^{j-1} e^{\alpha_j} & \text{if } y_i \neq h_j(x_i) \text{ (↑ wts if misclassified)} \end{cases}$$

Adaboost algorithm for classification

Input: Data $D = \{(x_i, y_i)\}_{i=1}^n$, M: Max no of iterations

Initialize: $w_i^0 = \frac{1}{n} \forall i$

(1) for $j = 1, \dots, M$

$$(a) h_j = \arg \min_h \sum_{i=1}^n w_i^{j-1} I(y_i \neq h_j(x_i))$$

$$(b) \epsilon_j = \sum_{i=1}^n w_i^{j-1} I(y_i \neq h_j(x_i))$$

$$(c) \alpha_j = \frac{1}{2} \ln \frac{1-\epsilon_j}{\epsilon_j} \quad (\text{assuming } 0 < \epsilon_j < \frac{1}{2})$$

(d) Update weights of examples:

$$w_i^j = \begin{cases} w_i^{j-1} e^{-\alpha_j} & \text{if } y_i = h_j(x_i) \forall i \\ w_i^{j-1} e^{\alpha_j} & \text{if } y_i \neq h_j(x_i) \forall i \end{cases}$$

Paper: Short Introduction to boosting
Stephane &

we upper

Classifier should minimize the weighted misclassification error.

In case of bagging $h_j(x) \rightarrow$ decision tree
 $\& \alpha_j = \frac{1}{n}$.

Thus

so

Will this procedure converge?

→ A correct example can be misclassified after few iterations

does the error get minimized?

$$\text{Average error} = \frac{1}{n} \sum I(y_i \neq f(x_i)) \leq \sum e^{-y_i f(x_i)}$$

\swarrow exp error upper bound of $\overbrace{\text{average error}}$
 \swarrow version $y \cdot f(x)$ has same sign

Also also used to solve regression problems

Adaboost

Learning weak classifier does not take much time, but even with 10 weak classifiers we get reasonably good result.

For regression problem:

o Learn $h_1(x)$ using $\{(x_i, y_i)\}_{i=1}^n$

$$y_i = h_1(x_i) + \text{error}_1 \quad \forall i$$

o Learn $h_2(x)$ using $\{(x_i, \text{error}_1)_i\}_{i=1}^n$

$$\text{error}_2 = h_2(x_i) + \text{error}_2 \quad \forall i$$

We first find the upper bound of error & then show that upper bound converges to 0 as $t \rightarrow \infty$

no. of iterations
Date / / 100

- Learn $h_3(x)$ using $\{(x_i, e_{n2i})\}_{i=1}^n$

$$e_{n2i} = h_3(x_i) + e_{n3i} \forall i$$

Thus, $y_i = h_1(x_i) + h_2(x_i) + h_3(x_i) + e_{n3i} \forall i$

- Learning wts using along with tree will result in

$$y_i = \alpha_1 h_1(x_i) + \alpha_2 h_2(x_i) + \alpha_3 h_3(x_i) + e_{n4i} \forall i$$

We will try to find out $\alpha_1, \alpha_2, \dots, \alpha_3$ so that it is as close as possible to y_i .

Proof of convergence of AdaBoost

- Initially, $w_t^0 = \frac{1}{n} \forall i$

- Normalized wts after T rounds:

$$w_i^T = \frac{w_i^{T-1} e^{-y_i \alpha_T h_T(x_i)}}{Z_T} \quad \text{where} \quad Z_T = \sum_{j=1}^n w_j^{T-1} e^{-y_j \alpha_T h_T(x_j)}$$

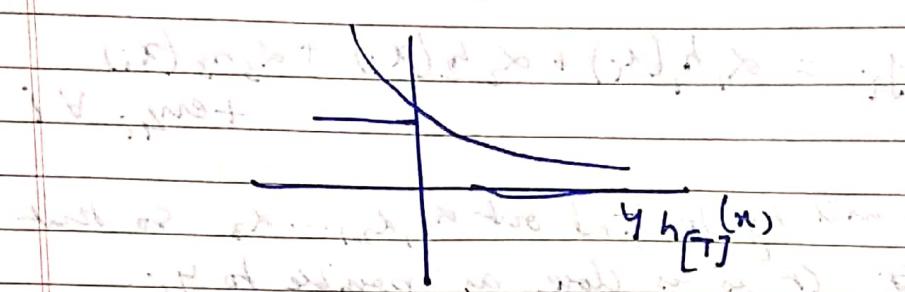
so that $\sum_{i=1}^n w_i^T = 1$. Equivalently,

$$w_i^T = \frac{1}{Z_T} \cdot \frac{e^{-y_i \alpha_T h_T(x_i)}}{\prod_{t=1}^T Z_t} \leftarrow \text{because of normalization.}$$

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \text{sign}(h_{[T]}(x_i)))$$

average 0-1 error incurred after T iterations

$$\frac{1}{n} \sum_i e^{-y_i h_{[T]}(x_i)}$$



Training error due to $h_{[T]}(x)$

$$E_M = \frac{1}{n} \sum_i I(y_i \neq \text{sign}(h_{[T]}(x_i)))$$

$$\leq \frac{1}{n} \sum_i e^{-y_i h_{[T]}(x_i)}$$

$$= \sum_i w_i^\top \prod_{t=1}^T z_t$$

$$= \prod_{t=1}^T z_t$$

$$\text{Now, } z_t = \sum_{i:y_i=1} w_i^{(t-1)} e^{-y_i \alpha_t h_t(x_i)}$$

$$z_t = \sum_{i:y_i=1} w_i^{(t-1)} e^{-\alpha_t} + \sum_{i:y_i \neq h_t(x_i)} w_i^{(t-1)} e^{\alpha_t}$$

$$= (1 - \epsilon_t) e^{-\alpha_t} + \epsilon_t e^{\alpha_t}$$

correct classification
rate



$$\begin{aligned} \therefore \text{Err} &\leq \prod_{t=1}^T z_t \\ &= \prod_{t=1}^T (1 - \epsilon_t) e^{-\alpha_t} + \epsilon_t e^{\alpha_t} \\ &= \prod_{t=1}^T 2 \sqrt{\epsilon_t (1 - \epsilon_t)} \quad (\because \alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}) \\ &\text{Ex: different } \epsilon_t \text{ for every iteration so we remove it} \\ &= \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} \quad (\text{letting } \epsilon_t = \frac{1}{2} - \gamma_t) \\ &\quad \gamma_t \rightarrow \text{how far it is close to random} \end{aligned}$$

Let $\gamma = \min(\gamma_1, \gamma_2, \dots, \gamma_T)$. Note $\gamma > 0$.

$\prod_{t=1}^T$
independent of T

γ always
 > 0
 $\because \epsilon_t$ always
 $< \frac{1}{2}$

\Rightarrow clearly, $\gamma^2 < \frac{1}{4}$ ($\because \gamma_t = \frac{1}{2} - \epsilon_t$)

$$\therefore \text{Err} \leq (1 - 4\gamma^2)^{\frac{T}{2}}$$

Our algo ensures that $\epsilon_t < \frac{1}{2}$

$$= e^{\frac{T}{2} \ln(1 - 4\gamma^2)}$$

$$= e^{\frac{T}{2}(-4\gamma^2)}$$

$$= e^{-2T\gamma^2} \quad \leftarrow \text{error asymptotically goes towards zero in a short time}$$

Thus $E_m \rightarrow 0$ as $T \rightarrow \infty$.

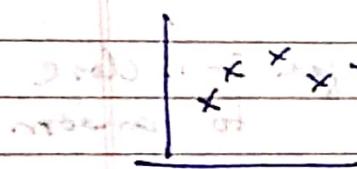
We stop adding models when validation error starts to deteriorate.

How to design a classifier that classifies weighted misclassification error?

Consider now

$$\min_{F} \bar{E} = \sum_i l(y_i, f(x_i))$$

$$F^0$$



$$F' = F^0 - \eta \frac{\partial \bar{E}}{\partial F} \Big|_{F=F^0}$$

optimization in non parametric space (or fn space)

Consider the problem of loss function optimization in fn space:

$$\min_F l(F) = \sum_{i=1}^n l(y_i, F(x_i))$$

$$\text{where } F = (F(x_1), F(x_2), \dots, F(x_n))^T$$

Gradient descent in fn space (assuming sufficient smoothness of $l(\cdot)$)

→ Initialize with F_1

$$\rightarrow F_2 = F_1 - \eta \frac{\partial l}{\partial F} \Big|_{F=F_1}$$

$$\rightarrow F_3 = F_2 - \eta \frac{\partial l}{\partial F} \Big|_{F=F_2}$$

$$\rightarrow F_t = F_{t-1} - \eta_{t-1} \frac{\partial l}{\partial F} \Big|_{F=F_{t-1}}$$

Given $\frac{\partial l}{\partial F} \Big|_{F=F_{t-1}}$, it is easy to determine η_{t-1} .

$$\eta_{t-1} = \operatorname{argmin}_{\eta > 0} \sum_{i=1}^n l(y_i, F_{t-1}(x_i) - \eta \frac{\partial l}{\partial F(x_i)}) \Big|_{F=F_{t-1}}$$

which can be used to update F_{t-1} as

$$F_t = F_{t-1} - \eta_{t-1} \frac{\partial l}{\partial F} \Big|_{F=F_{t-1}}$$

or $F_t = \sum_{j=1}^{t-1} \eta_j \left(-\frac{\partial l}{\partial F} \Big|_{F=f_j} \right)$ (letting $F_i(x_i) = 0 \forall i$)

o Compare with the AdaBoost model after $(t-1)$ iterations

$\xrightarrow{\text{both similar}} h_{[t-1]}(x) = \sum_{j=1}^{t-1} \alpha_j h_j(x)$

$\xrightarrow{\text{similar to } \eta_j}$

Can we align our classifier to the gradient direction of l w.r.t F ?

o $h_j(\cdot)$'s in boosting are analogous to negative gradients in function space.

Paper: Greed for Approximation: A gradient Boosting Machine.
J. H. Friedman
for regression, but can be extended to classification problems.

- Gradient Boosting Machine originally designed for regression problems.
- At t -th iteration, produce a model whose predictions are aligned with the -ve gradient direction of the loss function (design a model based on pseudo response)

$$h_t = \operatorname{argmin}_{h} \sum_{i=1}^n l\left(\frac{-\partial l}{\partial F(x_i)} \Big|_{F=F_{t-1}}, h(x_i)\right)$$

pseudo response

Our t -th model will be designed such that for each x_i it outputs pseudo response.

and determine η_{t-1} using

$$\eta_{t-1} = \operatorname{argmin}_{\eta} \sum_{i=1}^n l(y_i, F_{t-1}(x_i) + \eta h_t(x_i))$$

Helps to design model at t -th iteration in better way.

$$g = \begin{pmatrix} g_1 \\ \vdots \\ g_n \end{pmatrix}$$

$$g_i = -\frac{\partial l}{\partial F(x_i)} \Big|_{F=F_{t-1}}$$

same

$$\min \sum_{i=1}^n (-g(x_i) - h(x_i))^2$$

try to make $h = \text{neg. gradient.}$

$$F_t(x) = F_{t-1}(x) + \eta_{t-1} h_t(x)$$

No class on 27/02
Test #2 06/03

Friday 9:30 - 11:00

Assignment 2: due on 06/03 11:59 pm
extra class: 08/03

Page No. _____
Date / / 120

Review :

$$f(x_i) = w^T x = \sum_{j=1}^d w_j x_j$$

$$\sum_{i=1}^n l(y_i, \hat{y}_i) + \Omega(w)$$

↙ data dependent term

$$f(x_i; w)$$

can be overfit

because of
data dependency

Now

$$\hat{y}_i = \sum_{j=1}^M \alpha_j h_j(x_i) + \sum_{j=1}^M \Omega(h_j)$$

Model is in form of h_j & we want to control
its complexity. Term associated with every classifier
we have designed

$$\hat{y}_i^1 = 0 \quad \hat{y}_i^2 = \hat{y}_i^1 + h_1(x_i)$$

$$\hat{y}_i^3 = \hat{y}_i^2 + h_2(x_i)$$

suppose we have added $t-1$ models & we want to
add t^{th} model :

$$\sum_{j=1}^{t-1} \alpha_j h_j(x_i) + \sum_{j=1}^{t-1} \Omega(h_j)$$

$$\sum_{i=1}^n l(y_i, \hat{y}_i^{t-1}) + \lambda \sum_{j=1}^t \Omega(h_j)$$

Now h_j 's itself are parameters.

$$h^t = \underset{h}{\operatorname{arg\,min}} \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1} + h(x_i)) + \lambda \sum_{j=1}^{t-1} \|h_j\|_2^2 + \lambda \Omega(h)$$

Adding new model such that this obj fn is minimized
we use stagewise strategy. (All other models fixed &

minimizing ① + ③.

may not have a form like least squares loss.

So we get quadratic approximation of this.

Taylor's series approximation.

$$\Psi(x+h) \approx \Psi(x) + h(\Psi'(x) + \frac{h^2}{2!} \Psi''(x))$$

why quadratic?

→ Closed form soln of quadratic possible
for quadratic (may not be possible with
other fns).

F's here are \hat{y}_i .

We are dropping η_t term for simplicity.

o Additive training (models added in stagewise manner keeping previously added in previous round)

Stagewise (keeping all previous models fixed) $\hat{y}_i^1 = 0$ initially no knowledge of test fn.

$$\hat{y}_i^2 = F_2(x_i) = \hat{y}_i^1 + \eta_1 h_2(x_i)$$

η_t can be easily determined by simple optimization problem.

$$\hat{y}_i^t = F_t(x_i) = \hat{y}_i^{(t-1)} + \eta_t h_t(x_i)$$

where \hat{y}_i^t : Output of model on i th example at t th round.

o How do we decide which $h_t(\cdot)$ to add?

could be a decision tree (may not be decision tree always)

o If we use least squares error, $\ell(y, F) = (y - F)^2/2$, then one iteration of the algo is :

$$\rightarrow \text{Pseudo response: } \tilde{y}_i = y_i - F_{t-1}(x_i) = y_i - \hat{y}_i^{(t-1)} \quad (\text{residual error from previous round})$$

$$\rightarrow h_t = \operatorname{argmin}_h \sum_{i=1}^n (\tilde{y}_i - h(x_i))^2$$

$$\rightarrow \eta_{t-1} = \operatorname{argmin}_{\eta > 0} \sum_{i=1}^n (\tilde{y}_i - \eta h_t(x_i))^2$$

$$\rightarrow F_t(x) = F_{t-1}(x) + \underbrace{\eta_{t-1} h_t(x)}_{\text{fixed}} \quad (\text{Friedman's paper Gradient boost}).$$

We use gradient to decide what next model will be.

New model should be aligned with neg gradient of obj

fn. Here it is neg gradient

GradientBoost Algo :

Input : data $D = \{(x_i, y_i)\}_{i=1}^n$, M: max no of iteration

Initialize : $F_0 = \operatorname{argmin}_P \sum_{i=1}^n \ell(y_i, P)$

\leftarrow In Friedman's paper.

(1) for $t = 1, \dots, M$

$$(a) \tilde{y}_i = - \frac{\partial \ell}{\partial F(x_i)} \quad |_{F=F_{t-1}} \quad \forall i = 1, \dots, n$$

$$(b) h_t = \operatorname{argmin}_h \sum_{i=1}^n (\tilde{y}_i - h(x_i))^2$$

$$(c) \eta_{t-1} = \operatorname{argmin}_{\eta > 0} \sum_{i=1}^n \ell$$

Regression tree: partitions input space with axis parallel rectangle

Regression Tree: → will give no output as a model

- Assume each base learner is a regression tree having J leaf nodes

- Partitions input space X into J disjoint regions.

$$h(x; \{w_j, R_j\}_{j=1}^J) = \sum_{j=1}^J w_j I(x \in R_j) = w_r(x)$$

where w_j is the leaf score value associated with region R_j and $r: X \rightarrow \{1, 2, \dots, J\}$.

Every leaf node corresponds to
gives constant output.

w_1	w_2
	w_3

→ 3 leaf nodes.

dec regression tree
can be thought of
as a fn.

Every given pt lies in one of the regions

$$f(x) = \sum_{j=1}^J w_j I(x \in R_j) = w_r(x)$$

$$r: X \rightarrow \{1, 2, \dots, J\}$$

A fn
maps a fn to one of leaf nodes.

o Parameters of

→ No of leaf

Quantities

R_j 's

→ Tree depth

→ norm of w

So far

for the

o Assume

No depth

w_j

R_j

for

(1)

w_j

one

in

- o Parameters associated with regression tree

- No of leaf nodes (J) in the tree
- Quantiles that define the boundaries of the regions R_j 's.
- Tree depth Scores
- norm of the leaf score $\|w\|^2$

So far we have ignored regularization.

for the time being we assume we will use

- o Assume we have designed $(t-1)$ regression trees & want to design a new tree at t -th iteration.

$w_{jt} \rightarrow$ weight associated at j th region at t th instant.

$$F(x) = F_{t-1}(x) + \eta \sum_{j=1}^J h_t(x_j \{ w_{jt}, R_{jt} \})$$

- o R_{jt} : Regions at terminal nodes of the tree at t -th iteration constructed to predict the pseudo-response (Line 1(a) of GB algo)

$$w_{jt} = \text{avg}_{x_i \in R_{jt}} \hat{y}_i$$

one we identify the R_{jt} , we take all examples in that region & take their average to find w_{jt} .

- Obj fn minimized after t-th round:

$$E^{(t)} = \sum_{i=1}^n \underbrace{l(y_i, \hat{y}_i^{(t)})}_{\text{training loss}} + \underbrace{\lambda h(\cdot) - \Omega(h_t)}_{\text{Structural loss}}$$

But, $\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + h_t(x_i)$ — (separately optimize h)

So, the obj fn at t-th iteration is,

$$\min_{h(\cdot)} \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + h(x_i)) + \Omega(h)$$

Optimization prob at t-th iteration (n dropped, we will bring it later.)

~~to be optimized at end of (t-1)th iteration.~~

constant $E^{(t)} = \sum_{i=1}^n l(y_i, F_{t-1}(x_i) + h(x_i)) + \Omega(h).$

\uparrow
loss fn to
be optimized
at t-th
iteration

- Second order Taylor series approximation :-

$$f(x+\epsilon) \approx f(x) + \epsilon f'(x) + \frac{\epsilon^2}{2} f''(x).$$

• Define $p_i = \left. \frac{\partial l}{\partial F} \right|_{F=F_{t-1}(x_i)}, q_i = \left. \frac{\partial^2 l}{\partial F^2} \right|_{F=F_{t-1}(x_i)}$

Instead of F we can write \hat{y} here

- Therefore?

$$E^{(t)} \approx \sum_{i=1}^n \left(\underbrace{l(y_i, F_{t-1}(x_i)) + h(x_i)p_i + \frac{h(x_i)^2}{2} q_i}_{\text{constant}} \right) + \Omega(h)$$

$$l(y_i, F_t) \\ \approx l(y_i, F_t)$$

- Using this

- Using this

$$h_t =$$

- Can define

- $f, \lambda \rightarrow$ const
- Write this

- Recall,

- Obj fn :-

$$E^{(t)} \approx$$

- Define

$$\therefore$$

Step in Taylor's expansion

Page No. _____
Date / / 20

$$\ell(y_i) F_{t-1}(x_i) + h(x_i)$$

$$\approx \ell(y_i) F_{t-1}(x_i) + h(x_i) \frac{\partial \ell}{\partial F} \Big|_{F=F_{t-1}}$$

Using this approx, the $\frac{h^2(x_i)}{2} \frac{\partial^2 \ell}{\partial F^2}$ |
 $F=F_{t-1}$

Using this approx, the opt prob is :

$$h_t = \arg \min_h \sum_{i=1}^n \left(p_i h(x_i) + \frac{1}{2} q_i h(x_i)^2 \right) + \Omega(h)$$

we are
not varying
about this
now.

Can define complexity of a regression line as

$$\Omega(h) = \gamma J + \frac{\lambda}{2} \sum_{j=1}^J w_j^2$$

(different ways
to control
complexity
of free).

γ, λ controls contribution of each term

Write the approx obj fn in terms of w & J

Recall, $h(x) = \sum_{j=1}^J w_j I(x \in R_j) = \sum w_j \pi_j(x)$

Obj fn :

$$E^{(t)} \approx \sum_{i=1}^n p_i h(x_i) + \frac{1}{2} q_i h(x_i)^2 + \Omega(h)$$

$$= \sum_{i=1}^n p_i w_j \pi_j(x_i) + \frac{1}{2} q_i w_j^2 + \gamma J + \frac{\lambda}{2} \sum_{j=1}^J w_j^2$$

Define $I_j = \{ i : \pi_j(x_i) = j \}$

$$\therefore E^{(t)} \approx \sum_{j=1}^J \left[(\sum_{i \in I_j} p_i) w_j + \frac{1}{2} \left(\sum_{i \in I_j} q_i + \lambda \right) w_j^2 \right] + \gamma J$$

Sum of J independent quadratic fns (can be optimized independently)

- Obj fn: $E^{(t)} \approx \sum_{j=1}^J \left[\left(\sum_{i \in I_j} p_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} q_i + \lambda \right) w_j^2 \right] + \gamma J$

- define $G_j = \sum_{i \in I_j} p_i$ & $H_j = \sum_{i \in I_j} q_i$

$$\therefore E^{(t)} \approx \sum_{j=1}^J \left[G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2 \right] + \gamma J$$

- Given a tree with a fixed structure $\pi(x)$ (and also fixed J), the optimal parameters (for leaf j) & optimal obj fn are

$$w_j^*(\pi) = -\frac{G_j}{H_j + \lambda}, \quad E^{(t)*}(\pi) = -\frac{1}{2} \sum_{j=1}^J \frac{G_j^2}{H_j + \lambda} + \gamma J$$

T_1, T_2, \dots, T_m \leftarrow m trees from which we have to choose (their structures fixed)

for each tree find the optimal obj fn value associated with the tree.

$$E^{t^*}(1) \quad E^{t^*}(2) \quad \dots \quad E^{t^*}(m)$$

If we fix a tree, we have a fixed structure $\pi(x)$ & fixed obj fn value.

Choose the one which has min obj fn value.

There could be infinitely many tree structures possible for a given data set.

So heuristic tree used (Greedy approach) to choose from infinitely many tree structures possible.

$E^{(t)*}$ measures "goodness" of a tree structure (smaller the value, better the structure).

• How to structure

• for cu

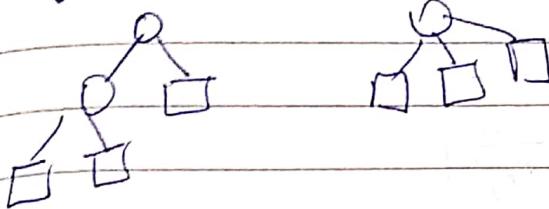
Bes

and

• Unfort

• Soluti

Say tree with 3 leaf nodes = Structures possible



- How to select the best regression tree, given a set T of tree structures at t -th iteration?

- for every tree in $\mathcal{T} \in T$

→ Compute G_j , H_j , w_j^* & j and $E^{(t)*}$

- Best regression tree :

$$T^* = \underset{T \in \mathcal{T}}{\operatorname{argmin}} E^{(t)*} \quad \underset{T \in \mathcal{T}}{\operatorname{argmin}} E_T^{(t)*}$$

and use the corresponding leaf score values w_j^* .

- Unfortunately, getting T is hard! (similar to feature selection (maximize gain))

- Solution: Grow the tree greedily.

(In the paper

Tangjichen

& Carlos Guestrin)
xgboost.

0 0	0 0 0
g_1, h_1 g_4, h_4	$g_2, h_2, g_5, h_5, g_3, h_3$

$$G_L = g_1 + g_4$$

$$G_R = g_2 + g_3 + g_5$$