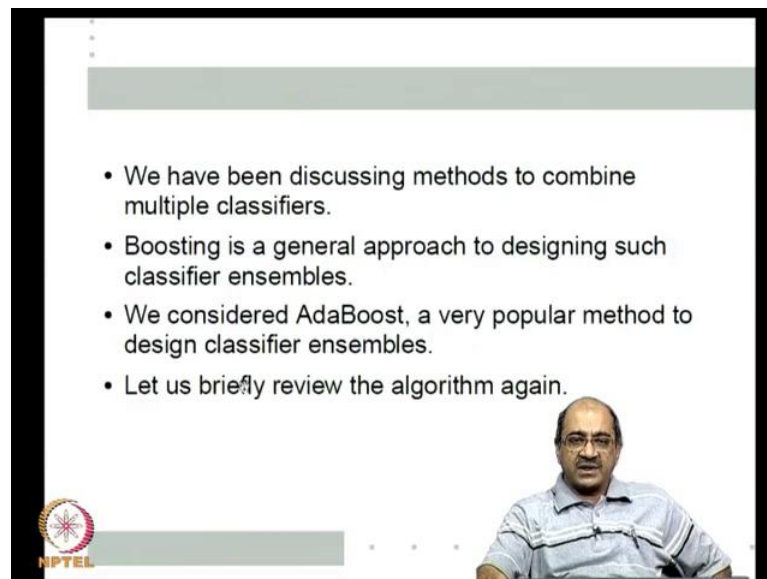


Pattern Recognition
Prof. P. S. Sastry
Department of Electronics and Communication Engineering
Indian Institute of Science, Bangalore

Lecture - 42
Risk minimization view of AdaBoost

Hello and welcome to this final lecture in the course of pattern recognition, we have been last class we been looking at classifier ensembles, that is how can we combined many different classifiers, we learn for the same pattern recognition problem. So, that we can improve the accuracy, we have briefly looked at an algorithmical called AdaBoost today one of the best classifier ensembles algorithms. Today's class will complete our discussion on AdaBoost, we look at AdaBoost also from a risk minimization point of view.

(Refer Slide Time: 00:53)

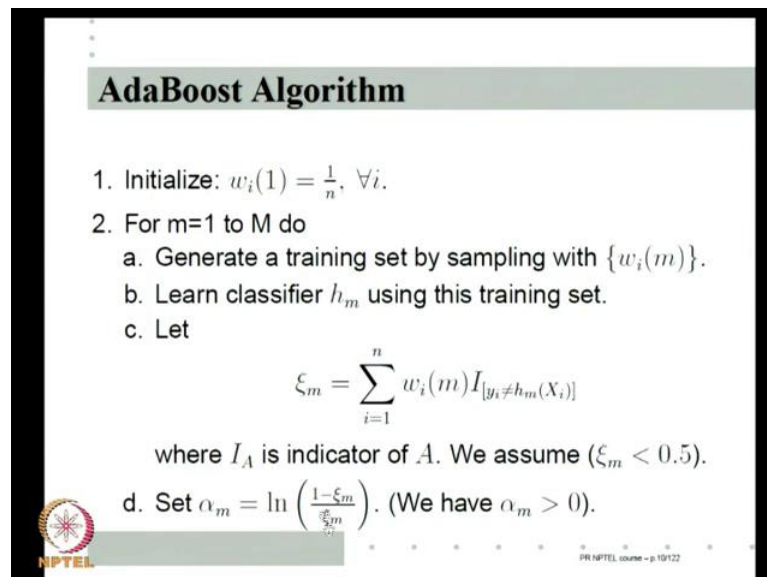


- We have been discussing methods to combine multiple classifiers.
- Boosting is a general approach to designing such classifier ensembles.
- We considered AdaBoost, a very popular method to design classifier ensembles.
- Let us briefly review the algorithm again.

So, to continue we said boosting is a general approach for designing classifier ensembles. So, we looked at no a simple intuitively designed boosting algorithm of learning three classifiers just to recognize about boosting is about. And then in general, there many techniques of boosting, and AdaBoost is a very popular method of designing such ensembles, because it allows to design an ensemble of almost any size, and in practice it works quite well. So, let us go back to the algorithm we have discussed it last class, but

we will look at it again.

(Refer Slide Time: 01:36)



AdaBoost Algorithm

1. Initialize: $w_i(1) = \frac{1}{n}, \forall i$.
2. For $m=1$ to M do
 - a. Generate a training set by sampling with $\{w_i(m)\}$.
 - b. Learn classifier h_m using this training set.
 - c. Let
$$\xi_m = \sum_{i=1}^n w_i(m) I_{[y_i \neq h_m(X_i)]}$$
where I_A is indicator of A . We assume ($\xi_m < 0.5$).
 - d. Set $\alpha_m = \ln\left(\frac{1-\xi_m}{\xi_m}\right)$. (We have $\alpha_m > 0$).

NPTEL
PR NPTEL course - p 10/122

So, the idea of AdaBoost is that I assign weights to each of the examples given, I got an example $x_i y_i$. So, each example I assign a weight changes from iteration to iteration at any given time, the weights are always kept normalize. So, at each iteration I use the weights as a probability distribution, and sample from the examples. That is example $x_i y_i$, I keep and taking n random samples from this example set so at each time is sample the independently chosen. And a fact that a $x_i y_i$ the probability of an $x_i y_i$ coming into my training set is given by w_i , the weight at that time w_i if m if m with y iteration.

So, essentially the examples with higher weights are more likely to come and of course, this is because I am using w_i as a probability distribution and sampling $i i d$, it is like sampling from the example set with replacement. So, multiple times I draw an example I might get the same example. So, this particular true if a particular example is very high weight it may be represented more than once in the training set, but without counting such I mean not bothering about such repetition, I have size n for my training set each training set is of size n and of course, some of the examples may be repeated because they are sample with w_i .

So, the algorithm starts by initializing the weight to $1/n$. That means, initially each example is as likely to be chosen as any other example because we have no preference on examples. Then I keep learning classifiers until M is a total number of classifiers ensembles so that is my iteration count. So, at the m 'th iteration I am learning the m 'th classifier. So, what do I do? I generate a training set by sampling with this weights, then you learn a classifier h_m using this training set, we have not said anything about how you learn this classifier. It can be any method, it can be a neural network, it can be SVM, it can be linear classifier, it can even be decision tree.

We have not consider decision trees except in the first introductory class I talked about we won't say what the classifier is also. As I said on different iteration you may use different classifiers, all that is of that that does not effect the algorithm. So, whatever is our preferred way of learning the base classifier each of the h_m 's are called base classifiers. We learn it and then we calculate the error x_i^m is the weighted error see indicator of y_i not is equal to $h_m(x_i)$ is simply means, on this I example, I met an error right and the error is weight by the weight of that example at this time.

Mind you this is not naturally, the training set because the training set is randomly chosen, but this is an all the examples I have x_i, y_i . For each example is classified by the newly learn classifier and wherever, I misclassified my actual error is weight by the weight of that classifier. So, I want to minimize this weight or some in a sense I can say. I want to learn a classifier h_m which minimizes this weighted error. We assume that each of the classifiers we learn is said that the error less than 0.5, in practice all it means is if I had a classifier whose error is greater than 0.5. I throw it away and learn again then set α_m is $1/n(1 - x_i^m)$ by x_i^m because we assumed x_i^m less than 0.5, $1 - x_i^m$ is greater than x_i^m and hence, α_m is always positive.

(Refer Slide Time: 05:33)

Algorithm contd.

e. Update the weights by

$$w'_i(m+1) = w_i(m) \exp(\alpha_m I_{[y_i \neq h_m(X_i)]})$$
$$w_i(m+1) = \frac{w'_i(m+1)}{\sum_i w'_i(m+1)}$$

3. Output the final classifier:

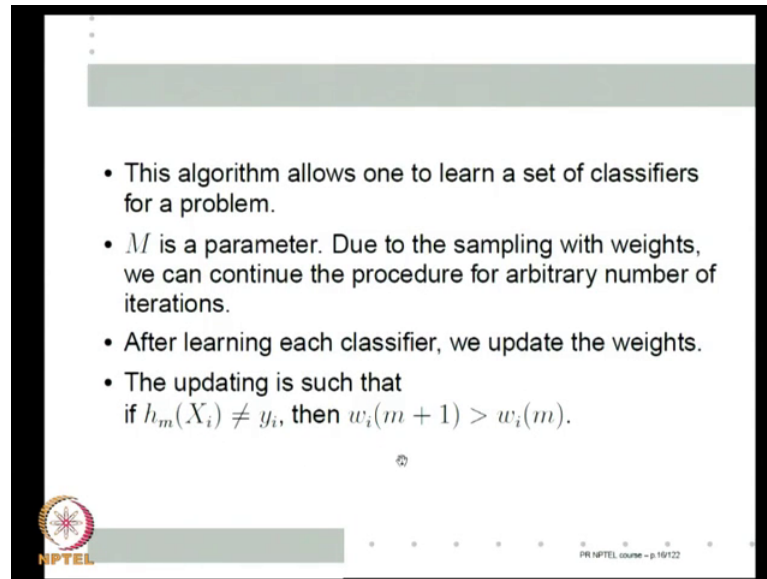
$$h(X) = \text{sgn} \left(\sum_{m=1}^M \alpha_m h_m(X) \right)$$

NPTEL logo and footer text: IPK NPTEL course - p 12/122

Then for the next iteration I learn a new weights. So, actually I want to keep weights normalize. So, the unnormalize once whatever the old weight I increase the weight, if I miss classifier this example, at this stage if y_i is not equal to $h_m(x_i)$. Then I multiply w_m by exponential α_m because α_m is positive, this is greater than 1 exponential α_m is greater than 1. So, I increased the weight and then I finally, normalize all the weights.

And my final learn classifier after learning all the h_m 's my final classifier is simply summation is equal to $\sum_{m=1}^M \alpha_m h_m(x)$. Where this α_m 's are what we calculated, which is $\frac{1}{1 - \sum_{i=1}^n x_i m}$, $x_i m$ where $x_i m$ is the weighted error of the n 'th classifier. So, α_m is some kind of a weight for the accuracy of h_m . So, I take this sum and its sign is what I am using as $h(x)$.

(Refer Slide Time: 06:28)



- This algorithm allows one to learn a set of classifiers for a problem.
- M is a parameter. Due to the sampling with weights, we can continue the procedure for arbitrary number of iterations.
- After learning each classifier, we update the weights.
- The updating is such that if $h_m(X_i) \neq y_i$, then $w_i(m+1) > w_i(m)$.

NPTEL

PR NPTEL course - p 16/122

So, this is the algorithm that we cancelled last time, the algorithm allows one to learn a set of classifiers for a problem. M is a parameter as we said we can learn arbitrary number of classifiers. So, we can learn ensemble with any number of classifiers in principal. After learning each classifier we update the weights, the whole idea of the weights is weight update is such as a just mentioned is that, if $h_m(x_i) \neq y_i$ then $w_i(m+1)$ is greater than $w_i(m)$.

If $y_i \neq h_m(x_i)$ then this will be exponential α^m this factor and α is greater than 0 α is positive. So, this is greater than 1. So, $w_i(m+1)$ is greater than $w_i(m)$ and hence, of normalization also $1/\alpha$ term. On the other hand if $y_i = h_m(x_i)$ then this factor is 1. So, $w_i(m)$ is $w_i(m)$. So, essentially in $w_i(m)$ I am increasing weight of all examples, whose which is miss classified leaving others where they are and then the renormalizing.

(Refer Slide Time: 07:38)

• This algorithm allows one to learn a set of classifiers for a problem.

• M is a parameter. Due to the sampling with weights, we can continue the procedure for arbitrary number of iterations.

• After learning each classifier, we update the weights.

• The updating is such that if $h_m(X_i) \neq y_i$, then $w_i(m+1) > w_i(m)$.

• If the current classifier misclassifies a pattern, its weight for the next iteration is increased.

NPTEL PR NPTEL course - p.17122

Which means, if I have misclassified error example, then its weight for the next iteration increases.

(Refer Slide Time: 07:52)

• We thus ensure that higher weightage is given to previously misclassified examples while sampling for the training set of next classifier.

• However, this does not mean we choose only such examples next time.

• We also want to maintain proper variability in the successive training sets generated.

• This is also achieved by the weight update scheme of AdaBoost.

NPTEL PR NPTEL course - p.21922

So, we ensure higher weightage is given to previously misclassified example. Essentially, that is the whole idea of generating successive training sets. However, we do

not want only misclassified examples, we already looked at it when we construct general boosting that does not give me enough variation the training set right. So, is not that we can just arbitrary increased the weightage have to know, sample only if in the misclassified examples for the next iteration. We also want to maintain proper variability in the, in the successive training sets such that a generator. As you seen you know we want a good balance between misclassified examples, and properly classified examples. So, this is also achieved by weight update scheme of AdaBoost, this is what we saw last.

(Refer Slide Time: 08:48)

• We have shown earlier that

$$\sum_{i: h_m(X_i) \neq y_i} w_i(m+1) = \frac{1}{2} = \sum_{i: h_m(X_i) = y_i} w_i(m+1)$$

• The weight update is such that at the next iteration, half the total weight is for patterns misclassified by the current classifier and the remaining half is for patterns correctly classified by the current classifier.

• This gives correct level of variability in successive training sets that are generated.

NPTEL

© NPTEL, course - p-24122

A matter of fact last class we proved that the weight of date is such that, if I sum all the weights over all i such that, $x_i \neq y_i$ is misclassified by h_m that sum is equal to some of all weights over i such that $h_m(x_i) = y_i$. So, because total weights is 1 and this two sums are same each is half. What is that means? The weight update is such that at the next iteration half the total weight is for pattern misclassified by the current classifier, and half the total weight is for patterns correctly classified by the current classifier.

This is automatically achieved by my weight update equation, my weight update equations form is such that this is achieve. We proved it last class that this happens to be so. Now, this is a very, very interesting thing about algorithm the algorithm. So, this way the algorithm automatically, can generate as many training sets as you want of required

variability. We do not have to worry about it, the weight update equations are such that they maintain this. This gives us sufficient variability in successive training sets.

(Refer Slide Time: 10:03)

The slide contains the following text:

- The final classification decision on a new X by this classifier ensemble is

$$h(X) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(X) \right)$$

- If α_m are same for all m , this is a simple majority decision. (Recall $h_m(X) \in \{-1, +1\}$).
- Here each component classifier has a different weight for its vote.
- The weight is based on the accuracy of that classifier.

The slide also features the NPTEL logo in the bottom left corner and the text "© NPTEL course - p.20/22" in the bottom right corner.

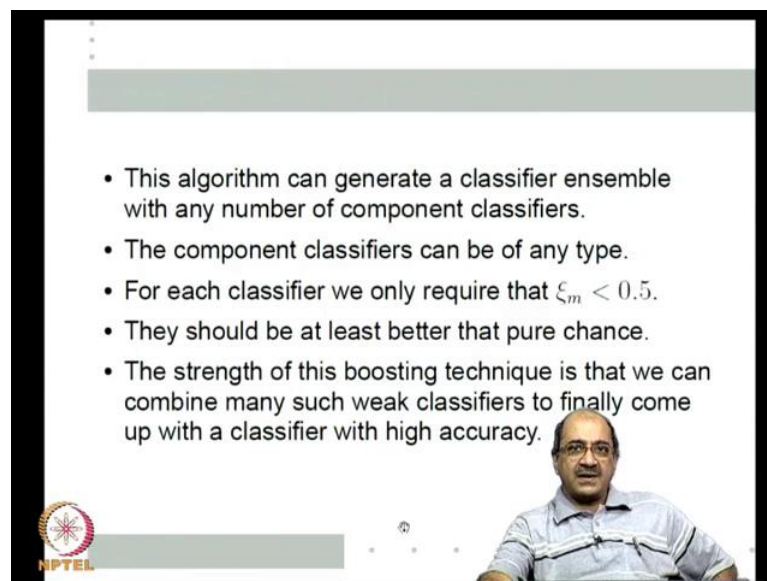
A word about the final classifier, what is the final classifier of this sign ensemble after learning all h_m 's, α_m is as we seen $1/n \sum_{i=1}^n \min(x_i, 1 - x_i)$. Where x_i is the weighted error. So, lower x_i higher is α_m . So, I take $\alpha_m h_m(x)$ and take summation take its sign as my final classifier. Please remember all this algorithm everything, we considering only for the two class case, where the class labels are plus 1 and minus 1.

So, when I take the sign function if the sum is positive than the sign is plus 1 otherwise, it is minus 1. Essentially, can easily see that because each h_m itself is a classifier, we assumed that $h_m(x)$ is either plus 1 or minus 1 for every x . So, because each of the plus 1 if all α_m 's are constant α_m is constant, and positive then it can come out of the sign function. Then if the sum has more plus ones finally, sign will be plus on the sum as more minus ones finally, the sign will be minus 1. So, essentially it would be a simple majority voting if I did not have α_m .

So, if I had α_m what it means is that it is a weighted majority vote. If I did not have α_m each classifier's vote has the same weight as. So, I am asking how many classifieds a plus one and how many classifieds a minus 1 and I am take the majority. Now, I am not doing that because each vote has a weight, I am summing up the weights of all the plus 1 classifieds summing of all the weights in minus 1 classifieds, I am asking which is better.

The weight is α_m weight depends on how good that classifier has been performing, the weight is based on the accuracy of that classifier. So, this a very interesting way to combine and once again a very, very general purpose way to combine classifiers, what is called a weight majority vote, this also very interesting feature of this AdaBoost algorithm.

(Refer Slide Time: 12:02)



- This algorithm can generate a classifier ensemble with any number of component classifiers.
- The component classifiers can be of any type.
- For each classifier we only require that $\xi_m < 0.5$.
- They should be at least better than pure chance.
- The strength of this boosting technique is that we can combine many such weak classifiers to finally come up with a classifier with high accuracy.

The interesting thing about algorithm is it can generate a classifier, ensemble with any number of component classifiers. There are many applications face recognition, image classification, you know many, many areas people have used this and there are in many applications people use 100, 200 classifier ensembles. So, algorithm can generate classifier ensemble with any number of component classifier, as you can see our weight update is such that successively, I can keep generating size n training sets which have the

right properties by updating weights, and sampling with those weights. And hence, the algorithm is you know very general purpose algorithm with lot of power.

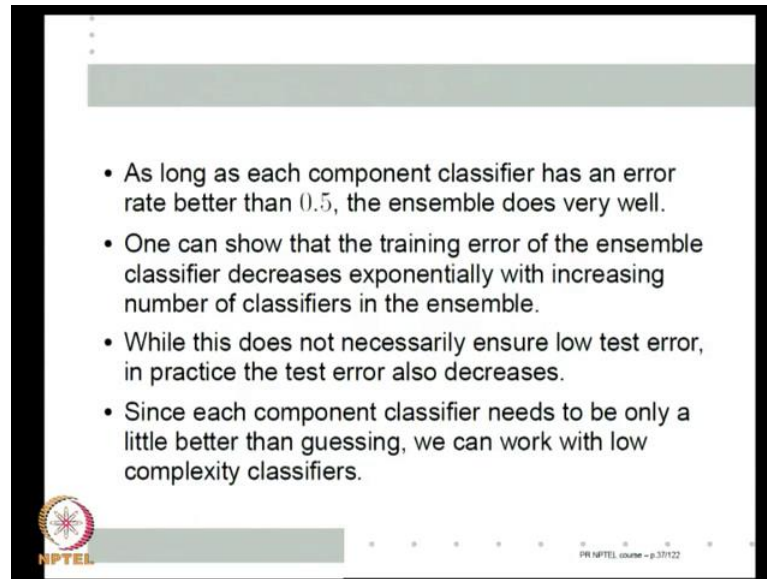
Second thing nice about the algorithm is I have to I do not have to worry about is if not as if applicable only for one type of component classifiers, you can put dish entries you can put neural network. So, you can put things learnt linearly squares you can learn what have you there is no restriction on the type of the component classifiers, that you can use all that we are asking for each classifier is that x_i is less than 0.5. I said in practice, I keep doing it and if any time my algorithm returns a classifier with x_i greater than 0.5 it a way and learning again, but in general what we will asking is that is each of a classifier should at least tell x_i is less than 0.5.

So, in this weighted sense they should be better than pure chance, the idea is that if the weight is high that particular example will be multiply represented in your training set. So, even if I am counting training set error if I misclassify that one example it will actually, be counted as 3 sent, or 4 training set examples are misclassified depending on the weight. So, in that sense essentially because the training set itself comes with the weight.

The weighted error being less than 0.5 is roughly same as saying that I can get I will get right classification more than 50 percent with the time. So, basically what I am asking is that I can't have very bad classifier, classifier should be at least good enough to be better than pure chance cursing. Just that much I am know these are very weak classifier, I am not asking much of I am not saying that accuracy should be very high, all I am saying is that accuracy should be just little better than pure chance.

So, that is the strength of this kind of boosting techniques, right? Such classifiers are called weak classifiers because they can only achieve a little better than pure chance classifier, but there is good enough if you give me a many, many such weak classifiers while nobody alone can get very good accuracy. If you give me 100 or 200 of them combining them properly, I can get very high accuracy. So, this strength figure out an AdaBoost is that it can combine many weak classifiers into coming up with a very good high accuracy classifier.

(Refer Slide Time: 15:07)



The slide contains a bulleted list of four points. At the bottom left is the NPTEL logo, and at the bottom right is the text '© NPTEL, course - p.37122'.

- As long as each component classifier has an error rate better than 0.5, the ensemble does very well.
- One can show that the training error of the ensemble classifier decreases exponentially with increasing number of classifiers in the ensemble.
- While this does not necessarily ensure low test error, in practice the test error also decreases.
- Since each component classifier needs to be only a little better than guessing, we can work with low complexity classifiers.

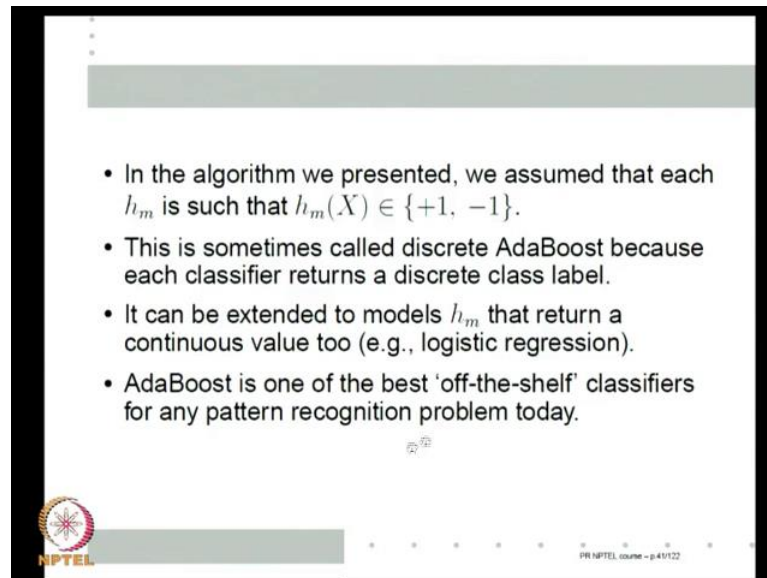
As long as each component classifier has an error rate better than 0.5, the ensemble in practice does very well. A matter of fact just a little algebra, we can show even though in the class we are not showing it, but is not very difficult to show that the training error the ensemble classifier decreases exponentially, with increasing number of classifiers in the ensemble. As long as we maintain this weighted error less than 0.5 in for each component classifier, one can show that the training error of the ensemble decreases exponentially fast with increasing number of classifiers in the ensemble that is a very, very nice property to have.

Of course, we know that low training error does not necessarily, mean low test error we can be doing over fitting. The reason why over fitting may not be occurring AdaBoost is to for firstly, we know it does not in the sense in practice in on many, many interesting applications many, many benchmark problems it does very well. One possibility is that because each component classifier, we ask for need to be only just a little better than pure chance guessing, can work with low complexity classifiers we do not have to take very high, high complexity classifiers.

Because of that right essentially, low training error could mean low test error and because ensemble gives me sufficiently, low training error it also gives me low test error.

Roughly, one way of looking at it is that ensemble allows me to get the higher accuracy, I do not have to go to very, very complex classifier structures. I can use simple classifiers let us say just linear classifier, but I use many of them that is how that is the strength of boosting.

(Refer Slide Time: 16:56)

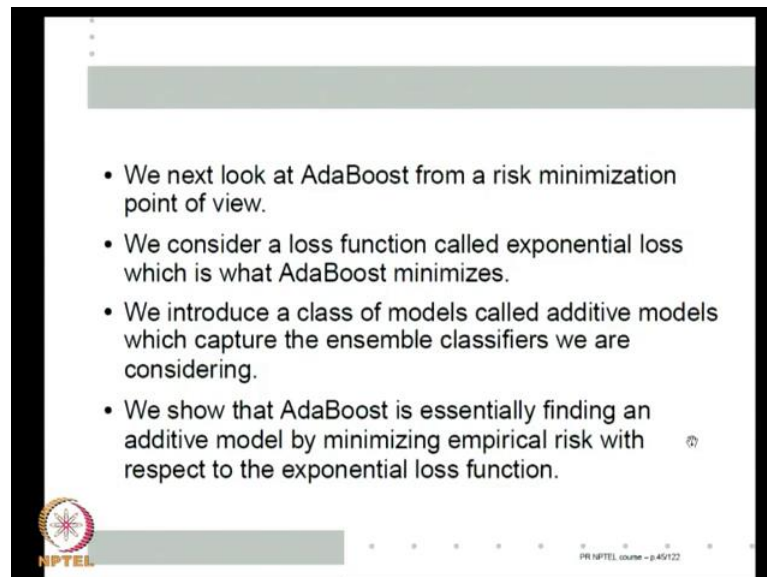


Of course, the particular AdaBoost algorithm we presented is such that we assume that each h_m is a classifier, that is just a binary value output h_m of x takes only values plus 1 and minus 1. So, this AdaBoost is called a discrete AdaBoost because each h_m takes on only discrete values the class labels. Of course, the AdaBoost algorithm itself it can be extended to other cases, where h_m returns a continuous value. Say for regression problems also for classification problems, where h_m could be a neural network. So, h_m is a continuous value or say logistic regression.

So, I can use any such classifiers also the algorithm is slightly different, but still similar conclusion hold. So, one can also extend the algorithm that we presented to the case of h_m that return a continuous value, rather than it discrete class labels as value. So, in this sense AdaBoost is possible one of the best of the shelf classifiers that are available today. Just like kernel base classifiers are good, and as I said the they should be the first choice whenever, looking at a classification problem. Similarly, a boosting AdaBoost is a very

good off the shelf immediately, implementable strategy for classifier learning a classifier, on many problems. And that is that is why it is quite popular and that is the reason why, this is the only classifier ensemble that we considering in this course.

(Refer Slide Time: 18:38)



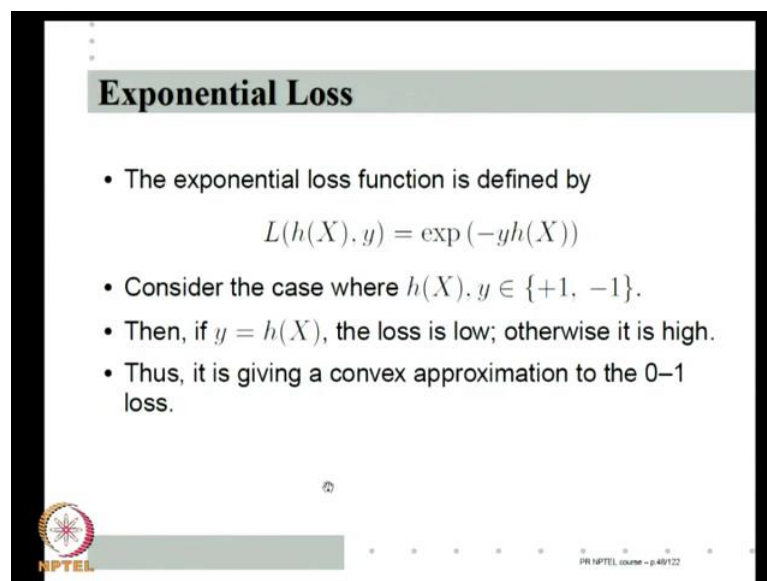
Having said that lets go back and look at AdaBoost a little more closely, all that we done so far is giving you an algorithm, which looks like a black bag box why those kind of weight updates, why that particular α_m should be the weight in the final classifier, nothing is clear have we I mean. As a matter of fact it utilize black magic, because you do not even no where we pulled out all those expressions from. So, what we will do for the rest of this class is we look at AdaBoost from what is called from our standard risk minimization view point.

So, we have been looking at all our classifiers as essentially, you choose a appropriate loss function and then you pass your classifier learning problem, as finding a classifier with minimum risk of course, I can not calculate risk. So, we do empirical risk minimization, all our classification algorithms are essentially empirical risk minimization methods. Except that we choose different loss functions. So, the same thing we did in s v m case, first we present s v m as a interesting optimization problem then shows that it also has a risk minimization view.

So, like that we will now show that AdaBoost has an interesting risk minimization view, the loss function that is appropriate for AdaBoost is what is called an exponential loss function, that is what AdaBoost minimizes. So, we will first briefly defined exponential loss, and the set of classifier the or the hypothesis space over which, this algorithm such is not individual classifiers, but ensemble classifiers. So, we have to introduce a proper you know model class for this so we will introduce a class of models what that a called additive models, which exactly capture the ensemble classifiers.

Then we show that AdaBoost is essentially, finding an additive model by minimizing empirical risk with respect to exponential loss function. Though if you think we are using exponential loss function, and want to me minimize empirical risk. While, searching over this class of additive models, the resulting algorithm looks like AdaBoost just to be more precise is not app a full minimization different a global minimum, but let say minimizer which uses a greedy heuristic. I will tell you shortly what a greedy heuristic means, but a slightly approximate way of minimizing empirical risk on the class of additive models is what AdaBoost is that is another way of looking at why AdaBoost gives good performance.

(Refer Slide Time: 21:10)



Exponential Loss

- The exponential loss function is defined by
$$L(h(X), y) = \exp(-yh(X))$$
- Consider the case where $h(X), y \in \{+1, -1\}$.
- Then, if $y = h(X)$, the loss is low; otherwise it is high.
- Thus, it is giving a convex approximation to the 0-1 loss.

NPTEL

PR NPTEL course - p 48/122


So, first let us look at the exponential loss function, exponential loss function simply l of

$h(x) - y$ is exponential minus y to $h(x)$. We have seen that many loss functions can be written as functions of $y - h(x)$. So, the $y - h(x)$ is negative say 0, 1 loss function simply says if $y - h(x)$ is negative that is y and $h(x)$ are opposite sign, then the loss is 1. If they are positive the loss is 0 that is of course, a non convex stuff function. Then we seen the hinge loss of $\max(0, -y(h(x) - 1))$ which tries to give a convex approximation to this.

We also saw that the square loss is another convex approximation to the 0, 1 loss function. So, similarly, exponential also is an approximation in the following sense, if you think of x and y $h(x)$ and y both are plus 1 minus 1. Then it is very easy to see that whenever, there of the same sign the loss will be exponential minus 1 and whenever, there opposite sign loss is exponential plus 1. So, that is a lot of difference. So, say essentially a loss is low if y and $h(x)$ are the same sign high otherwise and of course, the same conclusion holds you even $h(x)$.

For example, takes real values because it will only came so as long as y and $h(x)$ are the same sign it will be exponential minus something. Whenever, a y and $h(x)$ have opposite sign exponential plus something. So, between correct and incorrect classification there is a lot of difference, and in that sense it essentially like a 0, 1 loss function for analyze 0, 1 loss function it gives me a convex loss function is a convex function $y - h(x)$. So, this is a is one more convex approximation to the 0, 1 loss function.

(Refer Slide Time: 22:57)



- Next, we can ask why exponential loss function?
- It can be shown that the minimizer of risk under exponential loss over all real-valued functions h is

$$h^*(X) = \frac{1}{2} \log \left(\frac{P[Y = 1|X]}{P[Y = -1|X]} \right)$$

- Then, using sign of $h^*(X)$ gives us the optimal classifier.
- This is one reason why exponential loss is attractive.

PR NPTEL course - p.52/122

Let us of course, good enough reason to consider exponential loss, but we can still ask is this we have already so many loss functions why one more loss function. Well exponential loss has many, many interesting properties I will just state 1. It can be shown that if you have minimizing risk under exponential loss over all real valued functions h not just bind the valued functions h as we are considering here, but if you minimize the risk over all real valued functions h then the minimizer is half log half odds.

So, y is equal to 1 given x is posterior probability of class one, y is equal to minus 1 given x is posterior probability of class minus 1. So, this is essentially in a true class case this is odds for class 1. So, if the posterior probability class one is greater than posterior probability class minus 1, then this factor will be greater than 1 and log will be positive. And the other hand, if the posterior probability of class one is less than posterior probability of class minus 1. And the factor will be less than 1 and log will be negative.

So, essentially if I use sign of $h^*(x)$ I get the optimal classifier. So, if I can find h^* which is the minimizer of risk, and exponential loss and use its sign that will give me the base optimal classifier. Of course, this is over all real valued functions, but that not the risk any these a very interesting reason why, one should look at exponential value. This is as a matter of fact one of the reason for we because of which, exponential loss function

is attractive and used in many pattern recognition applications. So, with this we will we will take there we want to use exponential loss. So, let us look at what are so this is the loss function that AdaBoost uses and let us look at the model class as I said model class is additive models.

(Refer Slide Time: 24:53)

Additive Models

- Consider a classifier expressed as

$$f_M(X) = \text{sign} \left(\sum_{m=1}^M \beta_m h_m(X) \right)$$

where each h_m is a classifier and $\beta_m \in \mathbb{R}$.

- Since we are considering 2-class classifiers, we used the sign function.
- In general, f_M is a weighted sum of the outputs of h_m .
- Such models are called additive models.

NPTEL

PR NPTEL course - p.16/122

So, this is the kind of classifiers we want to concern f_m of X f subscript m of X is sign of some summation m is equal to 1 to m $\beta_m h_m X$. Where each h_m is a classifier and each β_m is some real number. Now, we take sign essentially, because we considering two class classifiers here because we considering two class classifiers, we take this the sum of all the classifiers. And take it sign for our final classifier, but in general I do not have to take this sign. So, any f_m that can be express as sum of classifiers like this is essentially, a ensemble classifier.

So, in general f_m is a weighted sum of the outputs of h_m of course, f_m anyway have to be classifiers it could be a regressor. So, it could be just a function model. So, in general any of f_m that is a weighted sum of outputs of so many h_m is called a additive model. So, an additive model is simply addition of many other, many other simple models. So, if each h_m is some function some regression function, or classifier. Then summation m is equal to 1 to $\beta_m h_m$. So, this is what an additive model is all about.

So, the set of the classifier additive model are simply obtained as weighted sums of ordinary base classifiers or ordinary sigma functions. Let say if we are learning an ensemble where essentially, learning a negative model. As you already know this is the kind of a, this is an ensemble classifier if each of the h_m are classifiers, this is like a weighted majority we have to somehow combined h_m 's. So, this is one we have combining so we can think of classifier ensembles.

(Refer Slide Time: 26:43)

• When learning a classifier ensemble, we are learning an additive model.

• If we are doing empirical risk minimization for learning an additive model (with fixed M), then we need to solve the optimization problem

$$\min_{\beta_m, h_m} \sum_{i=1}^n L \left(\sum_{m=1}^M \beta_m h_m(X_i), y_i \right)$$

where $L(\cdot, \cdot)$ is the loss function.

NPTEL

PR NPTEL course - p.16/122

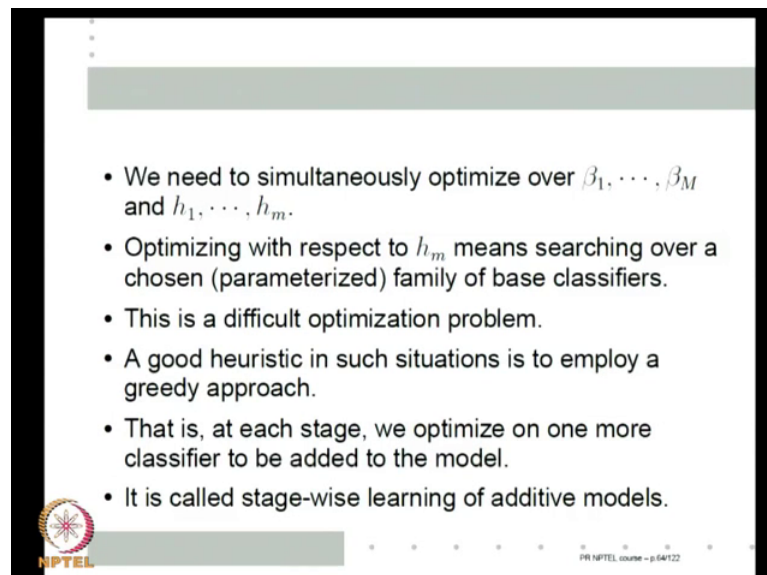
As essentially, learning a negative models. So, let us say we want to learn negative model that means, a model of this form let us assume m is fixed. So, the model of this form any specific member of this model class is specified by specific values of beta m and h_m for m is equal to 1 to n . So, that is what I have to learn actually. So, let us say we want to learn an additive model by doing empirical risk minimization. What is empirical risk?

Summation i is equal to 1 to n of your model say f of x_i if m of x_i coma y_i so of that 1 by n does not really matter in the empirically. So, stop the 1 by n . So, i is equal to 1 to n of what is my model m is equal to 1 to m beta m h_m x_i coma y_i . This is my if l is the loss function, then this is the loss on x_i y_i because this the first term here is the value of my of a model specified by beta m and h_m on x_i this is y_i . So, this is my loss if I sum

over i is equal to 1 to n that is my empirical loss.

So, essentially if you want to do empirical risk minimization, I have to minimize this expression over β_m and h_m . That is what empirical risk minimization would be, in our particular case l is a $1/h$ x comma y is exponential minus y/h x . So, will be exponential minus y_i into this sum, right? That will be the empirical risk with respect to exponential loss and that is what we want to minimize. Of course, this minimization is to over β_m h_m m is equal to one to capital M , we have to learn all of them we have to learn β_1 , β_2 , β_m , h_1 , h_2 , h_m .

(Refer Slide Time: 28:33)



The slide contains a list of six bullet points:

- We need to simultaneously optimize over β_1, \dots, β_M and h_1, \dots, h_m .
- Optimizing with respect to h_m means searching over a chosen (parameterized) family of base classifiers.
- This is a difficult optimization problem.
- A good heuristic in such situations is to employ a greedy approach.
- That is, at each stage, we optimize on one more classifier to be added to the model.
- It is called stage-wise learning of additive models.

The slide also features the NPTEL logo in the bottom left corner and the text "PR NPTEL course - p.64122" in the bottom right corner.

So, we have to simultaneously learn simultaneously optimize over β_1 to β_m and h_1 to h_m . What do you mean by optimize over h_1 to h_m ? That means, see we for each h_m we are searching over a class of base classifiers, may be they will be some parameters parameterized a family of base classifier. So, those parameter vectors is what we have to learn. So, we have to optimize over so many parameter vector at a time. This is a really hard optimization problem, but whenever you have this kind of additive models what standard heuristic is what is called a greedy approach. I want the m classifiers.

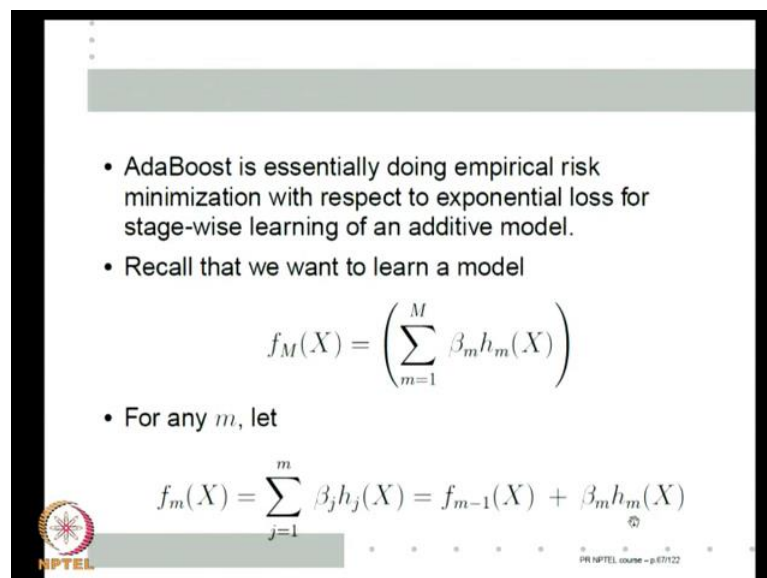
So, really I am asking which particular set of hundred classifiers is best, instead of doing

that what I can say is suppose, I want to use only one, what is the best classifier I take that. Now, if I want to use two classifiers I am saying only given that this is the first classifier, I am using what the second best classifier of course, that may not always be the best. If I am using the only one classifier what is best may not be even a member when I want to use two or three classifiers ensemble, but this called a greedy approach because I accumulate one by one.

So, at each stage we optimize for putting one more classifier in a model. So, if I already learnt m classifiers, when I am learning say first classifier I buck, I act as if I am only one classifier having learnt that one classifier. Now, I want to add one more while adding one more I cannot change, what I already learnt that is the characteristics of greedy approach. So, at each stage we optimize on one more classifier to be added to the model, without changing anything that we have already learnt.

This is called a stage wise learning of additive models. So, of course, I would still be minimizing this with respect to the entire model because I am adding one more to the model so that is the model consider, but while optimizing over that model. I am only optimizing over the new element new component, I am adding to the model that is what the greedy approach means.

(Refer Slide Time: 30:39)




• AdaBoost is essentially doing empirical risk minimization with respect to exponential loss for stage-wise learning of an additive model.

• Recall that we want to learn a model

$$f_M(X) = \left(\sum_{m=1}^M \beta_m h_m(X) \right)$$

• For any m , let

$$f_m(X) = \sum_{j=1}^m \beta_j h_j(X) = f_{m-1}(X) + \beta_m h_m(X)$$

 PR NPTEL course - p.07122

AdaBoost is essentially, doing empirical risk minimization with respect to exponential loss for stage wise learning of an additive model. So, is not learning the optimal additive model, but doing the greedy approximation to the optimal. So, this is the model class right we have to learn (()). So, let say as a as a rotation this capital M can be variables. So, for any little m f m x is sum of this from j is equal to 1 to m x, which also mean f m x is f minus 1 x plus beta m h m x, that is what meant by stage wise so stage m I have already learnt the m minus 1 stage classifier. And now, I want to add one more stage to the classifier. So, this is the new model I consider, I am asking what is the best new model for me, but when I say best I can only tweak beta m h m.

(Refer Slide Time: 31:35)

- At stage m , we have already learnt $m - 1$ component classifiers and are trying to add one more.
- So, we need to find 'best' β_m and h_m to minimize empirical risk of the m -component additive model.
- In the greedy approach, we leave all the previously learned ones unchanged.
- Hence, the solution for the next component of the model can be written as

$$(\beta_m, h_m) = \arg \min_{\beta, h} \sum_{i=1}^n \exp(-y_i [f_{m-1}(X_i) + \beta h(X_i)]) \quad \text{Q}$$

So, at stage m we already learnt m minus 1 component classifiers and I trying to add one more. So, we need to find the best beta m h m to minimize the empirical risk of the m component additive model, that is the whole idea of the approach. In the greedy approach we leave all the previously learned ones unchanged. As I said we are asking what is the best f m at the m th stage, when I say best f m I minimize empirical risk by taking f m minus 1 x plus beta m h m x as my model, but in this model the only thing that is tweak able is beta m h m. I will not tweak f m minus 1 f m minus 1, I consider fixed that is the whole idea of the greedy approach.

So, the greedy approach we leave all the previously learned classifiers unchanged. Hence, the solution for the next component of the model is see this is my model and I am using my exponential loss. So, my empirical risk is sum over the examples i is equal to 1 to n exponential minus y_i into my model my model is $f_{m-1}(x_i) + \beta h(x_i)$ βh is what I want to add at the m th stage β into h . And I am minimizing it all over possible β on h , this is my empirical loss empirical risk under the exponential loss.

So, if I have, if I decide to use the β , and the h as the parameters of the next component I am adding to my model then this is the empirical risk I get. So, I am minimizing it over all β and h and that minimizer is what I am calling β_m and h_m . That is what I am going to add as the n th component in my model.

(Refer Slide Time: 33:30)

• Define

$$w_i(m) = \exp(-y_i f_{m-1}(X_i)), \quad i = 1, \dots, n.$$

• We think of $w_i(m)$ as weight for X_i at the m^{th} stage.

• Since this depends only on f_{m-1} , it is a constant for the optimization problem at the m^{th} stage.

• So, the optimization problem now is

$$(\beta_m, h_m) = \arg \min_{\beta, h} \sum_{i=1}^n w_i(m) \exp(-y_i \beta h(X_i))$$

NPTEL

PR NPTEL course - p.75/22

So, let us define $w_i(m)$ is exponential minus y_i into $f_{m-1}(x_i)$. Why do I want to do that? See if I look at just the empirical risk expression, I have exponential minus y_i into $f_{m-1}(x_i) + \beta h(x_i)$ exponential minus y_i into $f_{m-1}(x_i) + \beta h(x_i)$ does not change by changing β and h so is outside this optimization. So, that is why we can take that factor out. So, $w_i(m)$ is exponential minus y_i into $f_{m-1}(x_i)$ we can think of $w_i(m)$ essentially, at the weight for x_i at stage m this are the we will ultimately show that these are the weights that AdaBoost is actually using.

Since, this depends only on f_{m-1} , it is a constant for the optimization problem at the m th stage, it does not affect the optimization problems at m th stage. Now, I can write the optimization problem effectively, yes minimize over β on h , i is equal to 1 to n with w_i exponential minus $y_i \beta h(x_i)$. So, I find the β that minimizes this and that once is what I called β_m because what are means. So, that will be the solution for my n th stage component classifier. So, this is what I have to do if I am doing stage wise empirical risk minimization, using exponential loss.

Then at the n th stage I want to add n th component, n th component you know model is represented β_m into $h_m \times \beta_m$ to h_m . So, my empirical risk for the model will be some, some unspecified weight at this point because that depends on the f_{m-1} part of the model into exponential minus $y_i \beta h(x_i)$. And which is summed over i is equal to 1 to n to that be the empirical risk, and this is minimized over β and h that the minimizing β and h happened to be my n th component. What we have first going to show is that if I minimize this, this expression over only h by keeping β fixed. For any fixed β I want to minimize it over all possible classifiers h .

(Refer Slide Time: 35:58)

• We first show that for any fixed $\beta > 0$, if we optimize over h , the solution is

$$h_m = \arg \min_h \sum_{i=1}^n w_i(m) I_{\beta}[y_i \neq h(X_i)]$$

• Thus, h_m is essentially the one that has least weighted error.

• This is the error that we have considered in the AdaBoost algorithm.

NPTEL PR NPTEL course - p.79122

Then the h that is that minimizes this is the h that minimizes this i is equal to 1 to n , with w_i indicator y_i not equal to $h(x_i)$ this is nothing but the weighted error. See the if you if

you look at i is equal to 1 to n w_i m indicator y_i not equal to $h(x_i)$, this is the weighted error of a classifier h . Weighted according to this w_i 's, and we were saying my h_m is the classifier that has the least weighted error. If these w_i happened to be the same weight that we use in AdaBoost then this exactly, what AdaBoost do.

H_m is essentially, the one that has the least weighted error and this is the error that we have considered in the AdaBoost algorithm. Of course, we have not at shown that this w_i which we define with respect to f_m minus 1 , has anything to do with our weight, but we will show that later. So, essentially we will first show that if I minimizing this, this is this is my m th stage optimization problem, for stage wise learning of negative models. So, if I did this optimization over only h for fixed beta no matter what value of fixed beta as long as beta is positive, this will be the solution.

There is very interesting because that exactly what the solution AdaBoost also looking because for any given classifier AdaBoost calculate this as the error, and we are looking for classifiers with low error. So, essentially AdaBoost is at the n th iteration AdaBoost is choosing h_m that has least value for this may not be doing this optimization fully, but that is what it is tweaking.

(Refer Slide Time: 37:39)


• Recall that here h is a 2-class classifier and we have $y_i, h(X_i) \in \{-1, +1\}$.

• Hence, $y_i h(X_i) \in \{-1, +1\}, \forall i$.

• Recall

$$(\beta_m, h_m) = \arg \min_{\beta, h} \sum_{i=1}^n w_i(m) \exp(-y_i \beta h(X_i))$$

• In this expression, the exponential term would be either e^β or $e^{-\beta}$ because $y_i h(X_i) \in \{-1, +1\}$.

 © NPTEL course - p. 82/122

Let us show this to show this lets first remember that because both y_i and $h(x_i)$ are minus 1 plus 1. That is the only the binary two class case what we considering because of that a product is also binary because individually, each y_i and $h(x_i)$ minus 1 plus 1 y_i into $h(x_i)$ is also the minus 1 plus 1. Now, this is what I want to optimize. So, inside the exponent I have $y_i h(x_i)$ forget the beta. So, $y_i h(x_i)$ factor takes only minus 1 plus 1 values this means, for each i this exponential fact that is a constant does not depend on i , I mean depends on i , but it can only take one of two value.

And each i it is either exponential minus beta or exponential plus beta. So, this expression can be either e^{β} or $e^{-\beta}$ because y_i into $h(x_i)$ is always minus 1 plus 1. We use that to rewrite this expression let say A is the expression that is being optimize that is A is this summation, the see we optimizing we said we are optimizing this only over h now keeping beta fix. Let us call this expression some A .

(Refer Slide Time: 38:45)

Denote by A the expression being optimized. Then

$$\begin{aligned}
 A &= \sum_{i=1}^n w_i(m) \exp(-y_i \beta h(X_i)) \\
 &= e^{\beta} \sum_i w_i(m) I_{[y_i \neq h(X_i)]} + e^{-\beta} \sum_i w_i(m) I_{[y_i = h(X_i)]} \\
 &= e^{\beta} \sum_i w_i(m) I_{[y_i \neq h(X_i)]} + e^{-\beta} \sum_i w_i(m) (1 - I_{[y_i \neq h(X_i)]}) \\
 &= (e^{\beta} - e^{-\beta}) \sum_i w_i(m) I_{[y_i \neq h(X_i)]} + e^{-\beta} \sum_i w_i(m)
 \end{aligned}$$

So, A is this expression i is equal to 1 to n $w_i(m)$ exponential minus. Now, we know that the y_i is not equal to $h(x_i)$ this y_i into $h(x_i)$ will be minus 1. So, this factor will be exponential beta y_i is equal to $h(x_i)$ this factor will be exponential minus beta. So, I can write this as sum over i $w_i(m)$ into indicator y_i not equal to $h(x_i)$. So, this will take $w_i(m)$ for only those i for which y_i is not equal to $h(x_i)$. Otherwise this will be 0 so for all of

them it will be e^{β} plus it will be $e^{-\beta}$ for all i 's such that the w_i is equal to $h(x_i)$. So, I can write this summation. Now, like this.

Now, I can further manipulate this indicator y_i is equal to $h(x_i)$ is same as $1 - \text{indicator } y_i \text{ not equal to } h(x_i)$ y_i indicator y_i equal to $h(x_i)$ is one when y_i equal to $h(x_i)$ if y_i equal to $h(x_i)$ y_i equal to indicator y_i not equal to $h(x_i)$ will be 0, and vice versa. So, this indicator is nothing but $1 - \text{indicator } y_i \text{ not equal to } h(x_i)$ now, I can combine this term with this negative term here to get $e^{\beta} - e^{-\beta}$ of w_i $\text{indicator } y_i \text{ not equal to } h(x_i)$ and what is left is the constant term $e^{-\beta}$ w_i , this of course, does not depend on h .

So, minimizing $J(h)$ only need to minimize over this, even here as long as β positive $e^{\beta} - e^{-\beta}$ is a positive quantity. So, it does not effect my minimizing so minimizing this whole expression over h is same as minimizing, only this expression over h that is the result we looking for.

(Refer Slide Time: 40:23)

• The second term in the above expression is a constant.

• Hence for a fixed $\beta > 0$, optimizing the above over h gives us

$$h_m = \arg \min_h \sum_{i=1}^n w_i(m) I_{[y_i \neq h(x_i)]}$$

• Thus, the next classifier that we should add is the one which minimizes this weighted misclassification error.

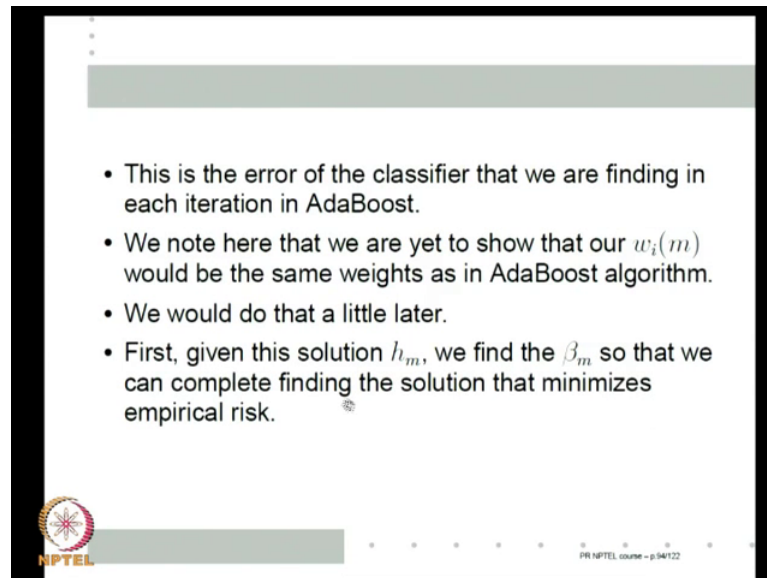
NPTEL

PR NPTEL course - p.10/122

The second term in the above expression is a constant hence, for a fixed β greater than 0 optimizing, the above over h gives us h_m is equal to $\arg \min_h \sum_{i=1}^n w_i$ $\text{indicator } y_i \neq h(x_i)$. That is this term this is what we wanted to show right this is interesting as we

already said the next classifier that we should add is the one which minimizes, the weighted misclassification error. So, the weighted misclassification error that AdaBoost looking at is very, very important of course, provided this w_m trans works to be weights there that will see the end.

(Refer Slide Time: 40:58)



- This is the error of the classifier that we are finding in each iteration in AdaBoost.
- We note here that we are yet to show that our $w_i(m)$ would be the same weights as in AdaBoost algorithm.
- We would do that a little later.
- First, given this solution h_m , we find the β_m so that we can complete finding the solution that minimizes empirical risk.

So, this the error that of the classifier that we are finding in each iteration in AdaBoost also, we note of course, that we have not yet shown that $w_i(m)$ would be the same weights as in AdaBoost algorithm. We will do that before the class is over, but lets complete our minimizes empirical risk to minimize empirical risk, we have to find h_m as well as β_m we show we first showed that no matter, what be what the β value is as long as positive. And we can have only positive weights so for any positive β we know how to find h_m . Now, that we found h_m with this h_m we have to find the best β that will be my β_m .

(Refer Slide Time: 40:35)

• Now, with h_m fixed, we need to find β_m by optimizing (over β),

$$(e^\beta - e^{-\beta}) \sum_i w_i(m) I_{[y_i \neq h_m(X_i)]} + e^{-\beta} \sum_i w_i(m)$$

• The optimization problem will be same if we divide the above by $\sum_i w_i(m)$.

• Define

$$\xi_m = \frac{\sum_i w_i(m) I_{[y_i \neq h_m(X_i)]}}{\sum_i w_i(m)}$$

NPTEL

PR NPTEL course - p 97122

So, with the given h_m we need to find β_m to optimize over β . So, what do I have to optimize. So, you already seen this is what to be optimized originally, and that is same as this. Now, I know this term is optimize by h is equal to h_m . So, I will put h_m there now a putting h_m there, I have to optimize this whole thing with respect to β optimize meaning minimize with this whole thing with respect to β . So, define β_m by optimizing over β minimizing over β this term. So, to make it simpler let us remember that because summation over i $w_i(m)$ will be a constant, see $w_i(m)$'s are exponential something were define it to be exponential y_i of m minus $1 \times i$.

So, the exponential something is always positive. So, this weights are all positive some of weights will be positive. So, if I just multiply this entire objective function by a positive constant it does not change the. So, the optimization problem will be same if we divide it by this we divided we get here summation over i $w_i(m)$ of indicator of y_i not equal to $h_m \times i$ whole divided by summation i $w_i(m)$ will give some name to it. We will call it ξ_m summation over i $w_i(m)$ indicate y_i not equal to $h_m \times i$ whole divided by summation over i $w_i(m)$. And this term this will go away I will get only e power minus β . So, this factor which does not depend on β we give it some name ξ_m .

(Refer Slide Time: 43:18)

• Hence, we can now write β_m as

$$\beta_m = \arg \min_{\beta} (e^{\beta} - e^{-\beta}) \xi_m + e^{-\beta}$$

• Differentiating w.r.t β and equating to zero we get

$$(e^{\beta} + e^{-\beta}) \xi_m - e^{-\beta} = 0$$

which gives us

$$\beta_m = \frac{1}{2} \ln \left(\frac{1 - \xi_m}{\xi_m} \right)$$

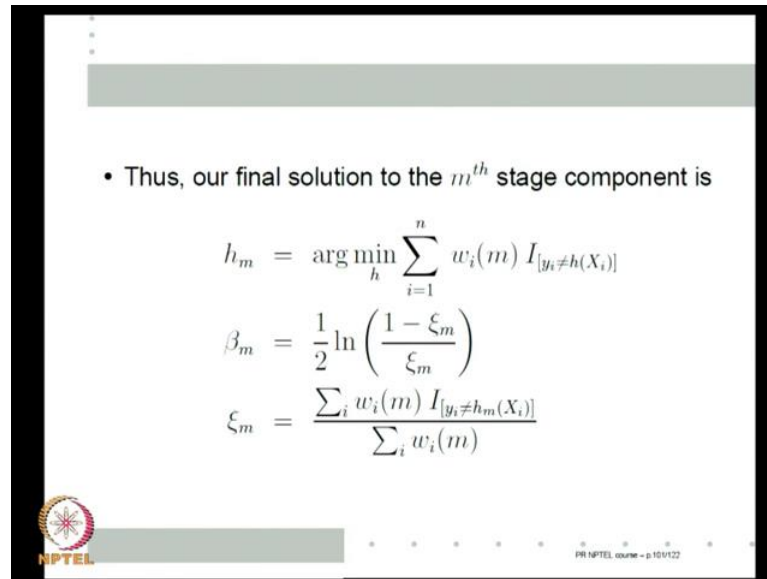
NPTEL

PR NPTEL course - p100122

Now, what we need to minimize over beta is $e^{\beta} - e^{-\beta} \xi_m + e^{-\beta}$. This is minimizing, this over beta is what gives me β_m . Now, this is very simple function of beta let's differentiate it equate to 0. If I differentiate it equate to 0, I get $e^{\beta} + e^{-\beta} \xi_m - e^{-\beta} = 0$. You solve for beta. So, that gives us the optimal beta β_m as $\frac{1}{2} \ln \left(\frac{1 - \xi_m}{\xi_m} \right)$.

I hope all of you still remember $\frac{1}{2} \ln \left(\frac{1 - \xi_m}{\xi_m} \right)$ from the AdaBoost algorithm that is α_m . That is the vote that m th classifier has in the final classifier ensemble, you know model β_m is the vote that the m th classifier had the final ensemble, this β_m happens to be just half, half of α_m .

(Refer Slide Time: 44:01)



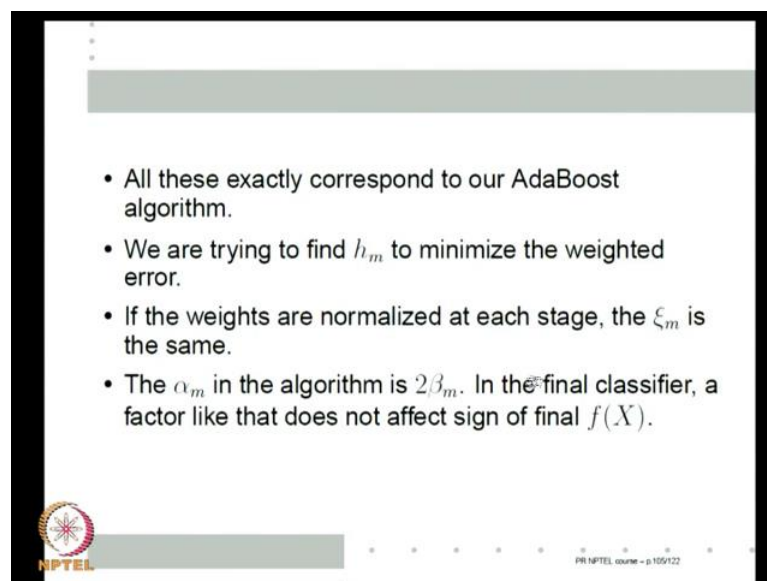
• Thus, our final solution to the m^{th} stage component is

$$h_m = \arg \min_h \sum_{i=1}^n w_i(m) I_{[y_i \neq h(X_i)]}$$
$$\beta_m = \frac{1}{2} \ln \left(\frac{1 - \xi_m}{\xi_m} \right)$$
$$\xi_m = \frac{\sum_i w_i(m) I_{[y_i \neq h_m(X_i)]}}{\sum_i w_i(m)}$$

NPTEL PR NPTEL course - p.105/122

So, this is our final solution at the m^{th} stage, the m^{th} stage component β_m h_m that you want to add. Basically, my final classifier is sign of β_1, h_1 plus β_2, h_2 plus β_3, h_3 so on. So, having length up to $m - 1$ I am learning m^{th} one at the m^{th} stage the best h_m is one that minimizes the weighted error, β_m is half of $\ln \frac{1 - \xi_m}{\xi_m}$ where ξ_m is given by this.

(Refer Slide Time: 44:31)

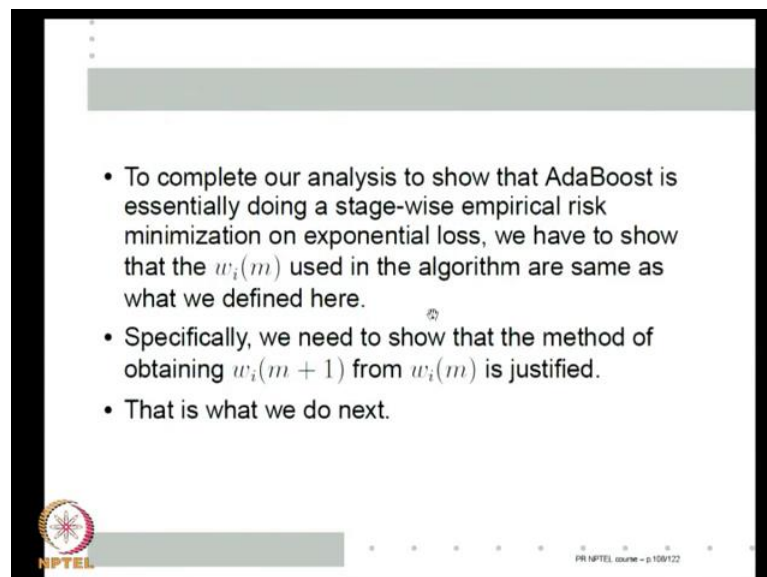


- All these exactly correspond to our AdaBoost algorithm.
- We are trying to find h_m to minimize the weighted error.
- If the weights are normalized at each stage, the ξ_m is the same.
- The α_m in the algorithm is $2\beta_m$. In the final classifier, a factor like that does not affect sign of final $f(X)$.

NPTEL PR NPTEL course - p.105/122

This exactly correspond to AdaBoost algorithm. Let us see one by one by we are trying to find h_m to minimize the weighted error that, that is exactly what AdaBoost also does. Then we will always keeping w_i^m normalize. So, summation w_i^m is 1 so x_i^m is same as summation over i $x_i w_i^m$ indicator y_i not equal to $h(x_i)$. So, that is exactly the x_i^m that we define in the algorithm with respect to the x_i^m , we define in the algorithm in AdaBoost the α_m is $\ln \frac{1}{1 - x_i^m}$ by x_i^m , whereas for us the β is twice α_m . So, α_m met the final vote in the weight of the vote β is the weight of the vote in my model, but essentially because I am taking sign whether I put $2\beta_m$ or whether I put β_m α is only changes a little bit. So, but any case the in the final thing the, what AdaBoost is between all α_m here putting 2β .

(Refer Slide Time: 45:38)



Now, to complete analysis now so what we showing is that essentially, what we learn is what AdaBoost algorithm is doing. So, if essentially doing empirical risk minimization using exponential loss function, the factors will get are same as we get in AdaBoost if this w_i is same as the way weight behaves in the AdaBoost. So, to complete our analysis we have to show that the w_i^m are actually, the weights used in the algorithm there. Specifically, what we have to show because any way at the at the first stage f_0 is arbitrary. So, we can an because we want to keep weight is normalized, it will be 1 by n . So, we only have to show that w_i^{m+1} to w_i^m algorithm is what I get.

(Refer Slide Time: 46:29)

• Recall that we have defined $w_i(m)$ by

$$w_i(m) = \exp(-y_i f_{m-1}(X_i))$$

• Hence we get

$$\begin{aligned} w_i(m+1) &= \exp(-y_i f_m(X_i)) \\ &= \exp(-y_i f_{m-1}(X_i) - y_i \beta_m h_m(X_i)) \\ &= w_i(m) \exp(-y_i \beta_m h_m(X_i)) \end{aligned}$$

NPTEL
PRINTEL course - p114122

This is what we do next, recall our notation f_m is j is equal to 1 to m $\beta_j h_j$ this additive model up to little m . So, which means, $f_m(x)$ is $f_{m-1}(x)$ plus $\beta_m h_m(x)$. Remembering this we know $w_i(m)$ is exponential minus y_i into $f_{m-1}(x_i)$ this what how we defined $w_i(m)$ in our empirical risk minimization. So, $w_i(m+1)$ will be exponential minus y_i into $f_m(x_i)$. Now, we know f_m is this. So, I substitute that will become minus $y_i f_{m-1}(x_i)$, minus $y_i \beta_m h_m(x_i)$.

Now, the first factor is what $w_i(m)$ is so this $w_i(m)$ into exponential minus $y_i \beta_m h_m(x_i)$. So, under our definition $w_i(m+1)$ is $w_i(m)$ into this of course, even in the AdaBoost this next, weights are obtained by multiplying the current weights with the exponential factor. So, it is just show that this is the right exponential factor as we know because ultimately, you are traumatizing this w 's a constant does not make any difference. So, let is see how this factor is same as what we put with the algorithm. So, the factor we have is exponential minus $y_i \beta_m h_m(x_i)$.

(Refer Slide Time: 48:03)

• Since $y_i, h_m(X_i) \in \{-1, +1\}$, we have

$$-y_i h_m(X_i) = 2I_{[y_i \neq h_m(X_i)]} - 1$$

• This gives us

$$w_i(m+1) = w_i(m) \exp(2\beta_m I_{[y_i \neq h_m(X_i)]}) \exp(-\beta_m)$$

• Note that $2\beta_m$ is same as α_m in the algorithm and that $w_i(m)$ are always normalized.

NPTEL PR NPTEL course - p117122

Since, $y_i, h_m(X_i) \in \{-1, +1\}$, we can write $-y_i h_m(X_i)$ as two times indicator $y_i \neq h_m(X_i)$ minus 1. Suppose, y_i and $h_m(X_i)$ are of opposite sign then this indicator will be 1. So, this will be $2 - 1 = 1$ on this side there opposite sign. So, $-y_i h_m(X_i)$ will be plus 1 conversely if y_i and $h_m(X_i)$ are of the same sign then this indicator will be 0. So, it will be minus 1 and this side also it is minus 1 because y_i and $h_m(X_i)$ are of the same side.

So, I can always write $-y_i h_m(X_i)$ as this. Now, I have minus $y_i h_m(X_i)$ here is β_m into minus $y_i h_m(X_i)$. So, I can substitute minus $y_i h_m(X_i)$ there so that gives me $w_i(m+1) = w_i(m) \exp(2\beta_m I_{[y_i \neq h_m(X_i)]}) \exp(-\beta_m)$. So, $2\beta_m$ indicator $y_i \neq h_m(X_i)$ β_m minus is already there right $y_i h_m(X_i)$ into β_m that's why the β_m comes from. So, is exponential there is won't be any minus because this minus $y_i h_m(X_i)$ is what is given by this.

So, $2\beta_m$ indicator of this into this minus 1 will simply give me exponential minus β_m . Now, the $2\beta_m$ is exactly equal to α_m . So, this is $w_i(m+1) = w_i(m) \exp(\alpha_m I_{[y_i \neq h_m(X_i)]}) \exp(-\beta_m)$. In the AdaBoost case, we first do $w_i'(m+1) = w_i(m) \exp(2\beta_m I_{[y_i \neq h_m(X_i)]})$ which is same as α_m into indicator $y_i \neq h_m(X_i)$, and then normalizing because I

am normalizing, if I multiplying all weights with the constant that does not depend on i that will go to accumulate normalization. So, this weight update is exactly same as the weight update that is used in AdaBoost, because two beta m is same as alpha m that $w_i(m)$ are always normalized. So, let us go back to as you seen.

(Refer Slide Time: 50:43)

• Thus, our final solution to the m^{th} stage component is

$$h_m = \arg \min_h \sum_{i=1}^n w_i(m) I_{[y_i \neq h(X_i)]}$$

$$\beta_m = \frac{1}{2} \ln \left(\frac{1 - \xi_m}{\xi_m} \right)$$

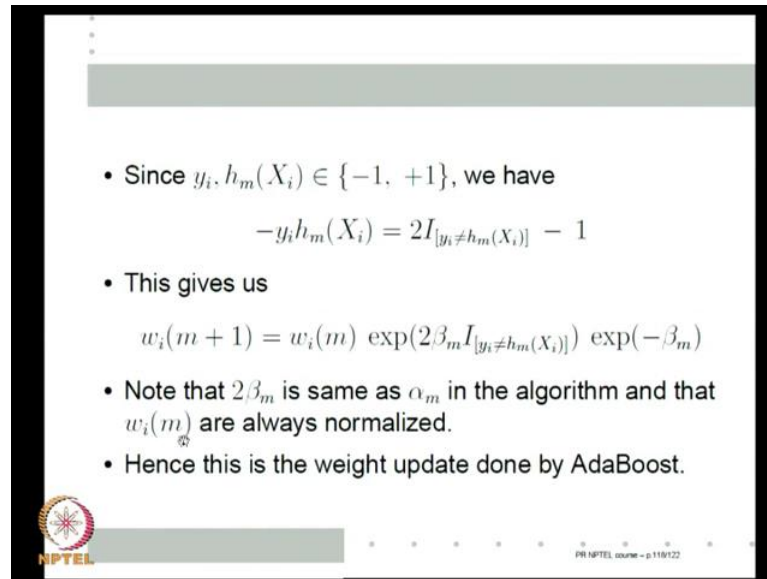
$$\xi_m = \frac{\sum_i w_i(m) I_{[y_i \neq h_m(X_i)]}}{\sum_i w_i(m)}$$

NPTEL

PR NPTEL course - p.10/11/22

If I do stage wise minimization of empirical risk under exponential loss, this is my m th stage solution, you find h_m that minimizes the weighted misclassification rate, calculus misclassification there if $w_i(m)$ are normalized than this is simply $\sum_i w_i(m) I_{[y_i \neq h_m(X_i)]}$. And then the weight itself is given by this right this exactly, what the or algorithm is doing as long as $w_i(m)$'s are normalize weights. And that is what we have shown now, $w_i(m)$ exponential alpha m to y_i not equal to $h_m(X_i)$ into some constant, this constant does not matter because I am normalizing it. So, note that $2\beta_m$ is same as alpha m and the $w_i(m)$ are always normalize.

(Refer Slide Time: 40:58)



• Since $y_i, h_m(X_i) \in \{-1, +1\}$, we have

$$-y_i h_m(X_i) = 2I_{[y_i \neq h_m(X_i)]} - 1$$

• This gives us

$$w_i(m+1) = w_i(m) \exp(2\beta_m I_{[y_i \neq h_m(X_i)]}) \exp(-\beta_m)$$

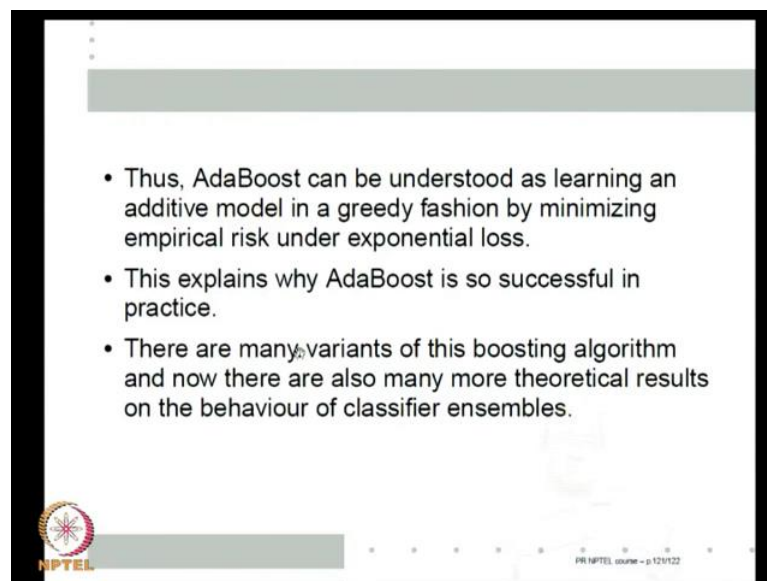
• Note that $2\beta_m$ is same as α_m in the algorithm and that $w_i(m)$ are always normalized.

• Hence this is the weight update done by AdaBoost.

NPTEL PR NPTEL course - p 119/122

Hence, the weight update is done by AdaBoost. So, what we have shown is essentially, the AdaBoost algorithm is finding a minimizer of empirical risk not a true minimize, but in approximate minimizer using a greedy approach, under exponential loss function.

(Refer Slide Time: 51:57)



• Thus, AdaBoost can be understood as learning an additive model in a greedy fashion by minimizing empirical risk under exponential loss.

• This explains why AdaBoost is so successful in practice.

• There are many variants of this boosting algorithm and now there are also many more theoretical results on the behaviour of classifier ensembles.

NPTEL PR NPTEL course - p 121/122

Thus AdaBoost can be understood as learning an additive model, in a greedy fashion by

minimizing empirical risk under exponential loss. So, this is one way in which I can assume AdaBoost successful in practice, AdaBoost as its mentioned earlier is fairly successful even very simple almost simple minded classification algorithms. If I take enough of them, learn through AdaBoost, learn through this proper weighting gives very good classifiers.

A matter of fact to for AdaBoost work well it is often more convenient, if the component classifiers are a little less accurate, in the sense if I am using very high complexity component classifiers, one classifier may have very high accuracy and it may dominate everything. So, it is often better not to do that. All kinds of things have been used including neural networks, (()), and so on. And today, there are many other variants of this boosting algorithm. And as a matter of fact more theoretical also available in terms of information theory and so on, as to avoid this boosting vector.

Boosting is one very major new area in pattern recognition new meaning, it in last ten-fifteen years this is a great idea to come up into pattern recognition. So, there many, many more new results in boosting, but you know this a interactive course. So, we just consult one basic boosting algorithm, but this is one this most often use and very useful in practice. So, this brings us to the close of the course we have completed all our lectures.

So, I will just briefly review back what all we have done. We started with the basic idea pattern recognition feature extraction classification at two step model, we seen that learning functions in general has the same flavor. So, was I am learning any model of f , the basic problem is your you give me examples x_i, y_i , i is equal to 1 to n and we have to learn some model f , such that f of x_i is a good approximation to y_i not just in the training set, but any future things at I may get.

Now, we let this problem in general look many examples and then decided there we look it at only the statistical sense, this entire course is only about statistical pattern recognition. So, we looked at decision statistical decision theory we looked Bayesian decision making. So, we defined loss function risk, what is the optimal base classifier. And shown that essentially, if I can calculate posterior probabilities of classes then I can

implement an optimal classifier, optimal in the sense of minimizing misclassification rate or any other kind of cost for misclassification.

We have also seen that minimizing misclassification may not be the only criterion, then other criteria are such as name, and person where we want to keep one kind of error always low, we see how such classifiers can be derived. Then we spend some time asking, if this is how we want to implement statistical pattern recognition. How can we implement them? So, we need to calculate posterior probabilities which means we need class conditional densities and prior probabilities, and we spend time in asking how we can estimate them.

So, we consider both parametric and non-parametric ways for estimating. In the parametric method, we looked at both maximum likelihood and Bayesian estimation, we looked at nearest neighbor and kernel density estimation, non-parametric estimates. And then we also looked at mixture densities, these are all for single densities, but we also looked at mixture densities and the MLE estimate using the EM algorithm. So, one method of implementing classifiers would be through density estimation, then we looked at other methods. Then we looked at discriminant functions, linear discriminant functions. We looked at perceptron we looked at Adaline which is a linear regression learning algorithm.

We in general looked at least squared method of learning linear models, both for classification, regression looked at the relationship between least squared. And you know as a Bayesian or MLE estimate under some assumed models. We looked at regularization, why regularization needed, why model complexity is an issue. We looked at from the (()) algorithm all the way to logistic regression, and Fisher discriminant. Then we step back and looked at statistical learning theory, we developed the theme of risk minimization a little more, and showed why empirical risk minimization is useful.

When it will work, when it will not work, what is the complexity of model class, we characterize the complexity of model class using (()). We explain what (()) is and essentially, justified empirical risk minimization. Then we moved on to looking at non-linear classifiers, we only consider two classes for non-linear classifiers. One based on

neural networks both feed forward network is sigma activation as well as real bases function networks, and then we looked at SVM's. Unfortunately due to lack of time, we did not concert one other very important topic that is among the first level topic that should be cover in a base pattern recognition course.

One topic that is not covered in this course is dysenteries. That is a very good class of classifiers, but for lack of time we did not cover the dysenteries, but we cover the other two major class of non-linear classifiers neural networks and SVM's. In SVM's we also looked at kernel base methods, which are now the most popular pattern recognition methods and lot of theories the error on them. After that we looked at model assessment, model estimation, bias variance trade off. We looked at how to estimate the final error of the model, how to chose model parameters, we looked at class validation, we looked at bootstrapping, arising out of that we ended the course by looking at classifiers, and some bulls, bagging, and boosting, and AdaBoost. I hope you people benefit it out of this course.

Thank you very much.