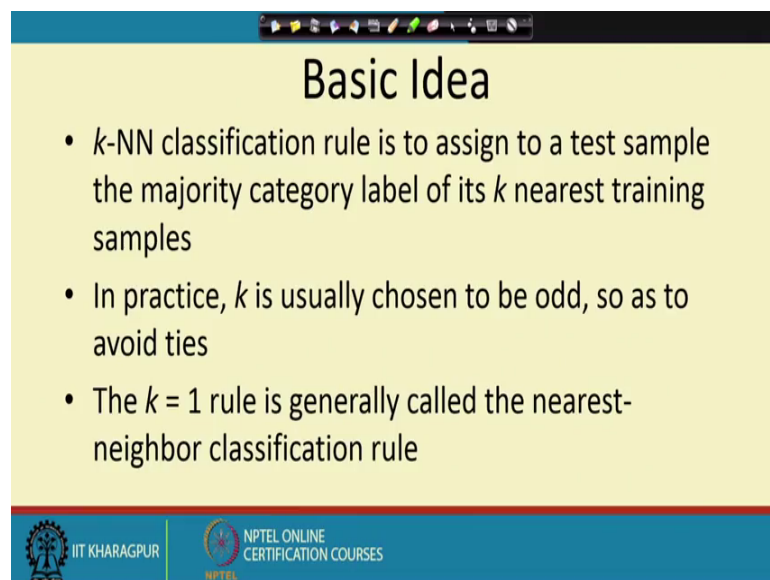**Data Mining**
**Prof. Pabitra Mitra**
**Department of Computer Science & Engineering**
**Indian Institute of Technology, Kharagpur**

**Lecture – 18**
**K-Nearest Neighbor – II**

We continue our discussion on the nearest neighbor rule, to summarize what we did is the following to a test sample, new sample, we find its K nearest neighbor training sample.

(Refer Slide Time: 00:41)



Usually K is chosen odd and then find out the majority class, majority category among the neighbors and put it into that class. As you can understand that if K is odd, if you have 2 classes, one of the class will always be the winner that there cannot be a tie, in more number of classes there can be tie. If there is a tie you can arbitrarily put in any of the class.

The special case where K is one is called the nearest neighbor rule

(Refer Slide Time: 01:28)



So, here is the pictorial description. So, to classify x we take in this case the one neighbor, only one more point, in this case 2 neighbor, in this case 3 neighbor and for example, in 3 neighbor the plus class is the winner. Majority test 2 to 1 and we put x as classify x as plus.

(Refer Slide Time: 02:13)



For the case of one neighbor this has a nice geometric interpretation called the Voronoi diagram or Voronoi cells.

So, each of these plus are my training points, each of this plus as my training point. So, around I tile or it is called a tessellation I split up my entire feature space into small cells, into small cells. What are the cells if you take a point in one of the cell this particular plus. So, note that corresponding to every plus we have a cell, every plus belongs to a cell and there is no other plus one cell has only one point and that is that plus. So, what is the, what is the definition of these cells is that if you take any point in this cell the closest among this training set the plus points is this.

So, if you take any point here suppose I take this point and you take distance to all the training sets the closest would be the representative point of that cell, will be the point of that cell. So, if we look at it the other way round, if we look at it the other way round given every point you define its Voronoi cell to be all points.
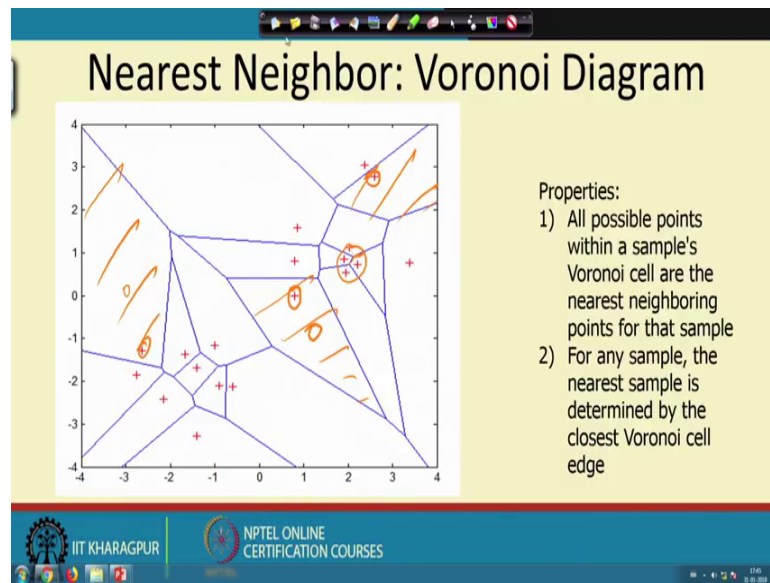
(Refer Slide Time: 04:18)



Whose nearest neighbor is this given point, is this now note that each point in the, is these Voronoi cells they define some kind of a equivalence class, a partition they define a partition.

So, if I use my nearest neighbor classifier one nearest neighbor classifier.

(Refer Slide Time: 05:05)



So, to say all I need to do to classify a new point is to check which Voronoi cell it belongs to and look at the class of the point corresponding to that cell, note that every training point will belong to one set and will define a set, no other training point will belong to that particular. So, it is as if the region of influence of this point, you might wonder why the design of influence is not circular, it is because it the shape of these cells depend on the shape of the distribution of the other points.

So, here the Voronoi cells are closer where we are spread apart because there are other points here which pull them. So, it is this, Voronoi cell is equivalent to the one nn rule, you just put it into the same class as the point, this is the geometric interpretation you will need it in some future discussion.

(Refer Slide Time: 06:20)



So, all so if we come back to the general K NN rule, what it did was the following take a neighborhood you port among the point.

(Refer Slide Time: 06:28)



So, there is 2 here, 2 here whose ever wins the port is that class, is the class of the new point, I have a I can make a slight modification to this instead of just a equal contribution to the port I can say that points which are closer have higher contribution to the port, have more voting power to the point which are closer to the new point which I want to classify ok.

(Refer Slide Time: 07:04)



So, here this minus has more voting power than these 2 points. So, this rule is known as the distance weighted K nearest neighbor rule.

(Refer Slide Time: 07:29)



So, basically f q is the class, f q is the class, sorry f q is the class, it is some weight function it is the class of the distance where v is the neighborhood, this is the normal K NN and I can have a distance weighted K NN. In fact, what I can do is that not just weighted by the distance, not just inversely proportional to that sorry not just inverse distance I can use that I can use a general weighting function called the kernel function.

(Refer Slide Time: 08:17)



So, I can maybe use something decaying like this. So, a, a function which let me see if I can draw it here. So, as distance increases it goes down its contribution to port goes down. So, this is a example of a kernel function. So, this technique is an is an generalization of the K nearest neighbor technique it is called the Parzen window technique, it is called the Parzen, Parzen window technique. Why do we have a, why do we have a either inverse weight inverse distance or some general function some Gaussian kernel or something in the distance weighted K NN.

In fact, what I can do is to sort of extend this rule to classify problems which are not just classification, but predicting some continuous valued also, contrast values also.

(Refer Slide Time: 09:53)



So, I want to for example, present the temperature of tomorrow, I am sorry I want to predict sorry not present I want to predict the temperature of together depending on other values like the pressure rainfall and so thing and my training point consists of a temperature of say 10 yards in this region. So, what I do I find K other similar days, we have a similar pressure rainfall profile and what that temperature are.

And I predict temperature of tomorrow is the average of all these K neighbor temperatures, may be an weighted average if we consider distance. So, that K NN is a general principle beyond classification also. So, I guess this is clear, what we can what you can possibly do is that you can take some example or say something like this. You assume that there are 3 classes and each of the class sort of easy you make a assumption do not use it actually it is nonparametric.

So, I am not going to assume they are Gaussian, just for generating the data you assume mu 1 sigma 1, mu 2 sigma 2, mu 3 sigma 3 are the 3 Gaussians and you randomly use these Gaussians to generate say thousand points and now they are belong to class 1, 2 and 3, I denote class as integers and for a. So, you have 10000 example, 1000 examples generated from 3 Gaussians, class 1, class 2, class 3 and you take a new, this is a computer experiment. So, you generate this points take a new point x, some new point x and take some value of K say 3 and apply the K NN rule, K nearest neighbor rule that is the non weighted version then the weighted inverse distance weighted version.

And see how much accuracy you get how well how many times it correctly classifies a new example. So, you can do this experiment using a writing a small program in any language and you will see the effect. So, if you are trying to write this program and actually do, you will face the following issues, you have to choose a value of k.

(Refer Slide Time: 13:50)



For 2, 3, 5, 7 what value you have to choose a distance measure because everything is nearest neighbor is defined only when a distance is defined, you have to choose what size training set you want, how many points to keep in the training example and that will depend on how much computational complexity you can afford, how much time you need to classify you can effort to classify.

So, we will just quickly discuss some thumb rules, there is no theory some thumb rule to decide on these parameters.

(Refer Slide Time: 14:54)



First the value of K, if K is too small your neighborhood is very small. So, noise and other for example, let me tell this. So, see we will probably expect this x to lie in this plus class or actually it is in the boundary so you can find. So, if we choose K to be only one nearest neighbor, you get minus as the class. If you choose K to be 2 you get a 50 50. So, here it depend, if it is K is small it depends too much on the local property ok.

Similarly, if K is too large it depends on broad range. So, what is a good value of K? Actually let me do it, noise if it is a large K it is too global, can any of you tell you if K is equal to size of the training set, what does this mean? you go back to your previous slide if K equal to size of a, K nearest neighbor is give the probability of nothing, but the prior probabilities of the classes. So, if it is too large this may be the case. So, there is really no theory to decide what is the correct value of K, some rule of thumb at here if one rule of thumb is this, when n is the size of training set.

(Refer Slide Time: 18:03)



Some other people take it this way take it this way. So, this is a option, this is a option it you do not really know which one will work, this is a open problem in data mining nobody really knows. So any of you if you can find the good solution to the problem that will be a quite a contribution to data mining, what about the distance metric.

(Refer Slide Time: 19:02)



So, many possibilities are there, what kind of measure should I take as a distance, I list here a number of distances, a most common is Euclidean, a little bit of generalization is this.

So, this is this plus this is Manhattan or city block, you can have a weighting factor q called the mahalanobis distance this one, others are also possible.
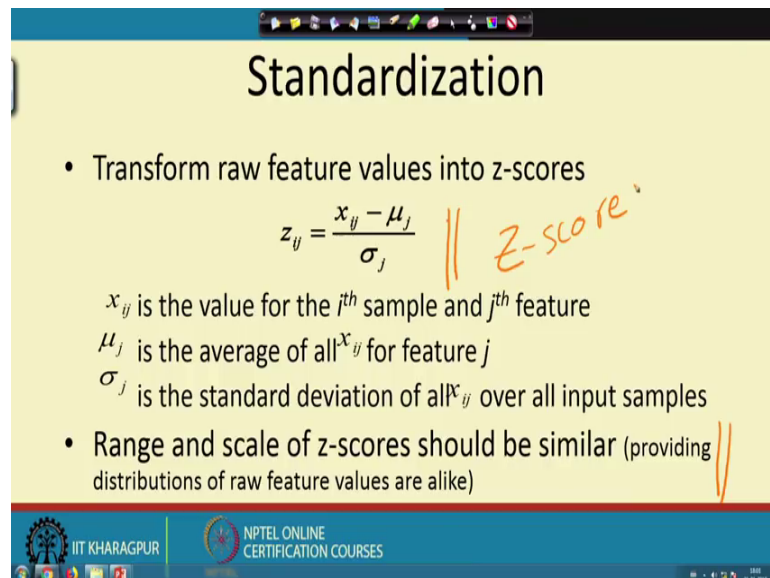
(Refer Slide Time: 20:50)



Some factors to be noted, scale what do you do? You normalize each feature value.

(Refer Slide Time: 21:07)



So, it becomes between 0 to 1 or minus 1 to 1 rather z score is a normalized attribute value.

(Refer Slide Time: 21:50)



Note that many of these distance measures they fail when the number of attributes are high ok.

So, if I have, so sparse if the many of the attributes have 0 value. So, you can see here distance between these two vector is this Euclidean, here it is this whereas, these two are may be similar, these two are very similar this is quite a far apart same d value Euclidean, all right. So, you have to take care of this is a very common problem in applying nearest neighbor to high dimensional data without properly choosing a distance function, many time it depends on the domain.

(Refer Slide Time: 22:58)



How do you define distance in nominal attributes do this, you define how many times at an attribute value is associated with the class and use that all right. So, you just look at the definition I have defined it. So, this is this, is this and c is the number of output class and this is defined this way. So, suppose some attribute a is say low medium or a high like that nominal valued, if you go back to our discussion in data preprocess if we discussed some of this, that is why I am not discussing this in details here you should read that again.

(Refer Slide Time: 24:25)

So, this is for heterogeneous. So, with this I stop this lecture, I complete my discussion of the distance function we will go into other considerations of nearest neighbor which are important in the next lecture so.

Thank you.