

The slide has a blue background with a faint image of a rainbow and clouds. In the top left corner, there is a dark blue triangle containing the Kent State University logo. In the top right corner, there are four small yellow flower icons. The title "Naive Bayes Classifier" is centered in a white serif font. In the bottom right corner, the website "WWW.KENT.EDU" is written in a small, white, sans-serif font.

Naive Bayes Classifier

In this module, we discuss the Naive Bayes Classifier. The naive Bayes method, and a branch of statistics called Bayesian Statistics, is named after the Reverend Thomas Bayes (1702–1761).

Generative vs. Discriminative Classifiers

- Training classifiers involves estimating $f: X \rightarrow Y$, or $P(Y | X)$
- Discriminative classifiers :
 1. Assume some functional form for $P(Y | X)$
 2. Estimate parameters of $P(Y | X)$ directly from training data
- Generative classifiers:
 1. Assume some functional form for $P(X | Y)$, $P(X)$
 2. Estimate parameters of $P(X | Y)$, $P(X)$ directly from training data
 3. Use Bayes rule to calculate $P(Y | X = x_i)$

Consider an example of classifying an observation, say whether a disease is malignant or benign, based on factors like gender, age, etc. Here, Y , the label, is either benign or malignant, and x , is the set of predictor variables, age, gender, etc. The difference between generative and discriminative classifiers is then as follows:

Generative classifiers learn a model of the joint probability, $p(x, y)$, of the inputs x and the label y , and then make their predictions using the Bayes rule to estimate the probability of y (malignant or benign) given the values of the predictor variables, $p(y | x)$ in notation, and then picking the most likely label. Discriminative classifiers, on the other hand, directly estimate $p(y | x)$. An example of the latter method is logistic regression. An example of the former method is the Naive Bayes approach, which we study in this module.

Naive Bayes Classifier

- Setting: We want to predict the class of Y (our output variable) given a set of predictors (i.e. input variables) X_1, X_2, \dots, X_n
- A good strategy is to choose a class of Y that maximizes the probability of Y given the vectors of predictors:

$$\arg \max_Y P(Y|X_1, \dots, X_n)$$

(no audio)

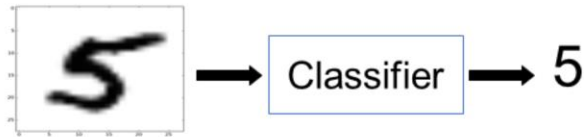
The idea behind NB classifier is simple. We are interested in predicting the class of Y (the label) based on values of a set of predictors. One strategy is to assign the record to the class Y that maximizes the probability of seeing Y given the values seen for the set X in that observation. This approach can be easily implemented as follows for a given observation to be predicted:

1. Find all the other records with the same predictor profile (i.e., where the predictor values are the same).
2. Determine the probability that those records belong to the class of interest.

Let us see a simple example to illustrate this.

Example

- Let assume the input variables are the intensity (or brightness) of each of the pixels of a image. The task is to classify a handwritten digit as 5 or 6.



- We can calculate the probability of the output being 5 given the input, and also calculate to calculate probability of the output being 6 given the input.
- We then compare the probabilities and decide whether the input was a 5 or 6 based on those probabilities.

The objective of the classifier is to determine the probability, or propensity, for an observation to belong to a certain class. In this example, we can calculate the probability that the image is a 5 or a 6 given the specific values of the image intensity. Once we know the probabilities, or the relative values of the probabilities, we can assign the observation to the label with the highest probability.

The probability calculations rely on the concept of conditional probabilities, which we study next.

The Bayes Classifier

- Key idea: the original sample space no longer applies.

$$P(Y|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y)P(Y)}{P(X_1, \dots, X_n)}$$

Diagram labels: Posterior (points to $P(Y|X_1, \dots, X_n)$), Likelihood (points to $P(X_1, \dots, X_n|Y)$), Prior (points to $P(Y)$), Normalization Constant (points to $P(X_1, \dots, X_n)$)

- We use the prior and likelihood probabilities to calculate posterior
- The dominator is not important as it is the same for all class probabilities (i.e. does affect the comparison)

The NB classifier uses the concept of conditional probability, or the probability of event A given that event B has occurred [denoted $P(A|B)$]. In this case, we will be looking at the probability of the record belonging to class Y_i given that its predictor values are x_1, x_2, \dots, x_n . In general, for a response with m classes Y_1, Y_2, \dots, Y_m , and the predictor values x_1, x_2, \dots, x_n , we want to compute $P(Y_i|x_1, \dots, x_n)$. This is called the Posterior Probability, and is shown in the equation above.

To classify a record, we compute its probability of belonging to each of the classes in this way, then classify the record to the class that has the highest probability or use the cutoff probability to decide whether it should be assigned to the class of interest.

From this definition, we see that the Bayesian classifier works only with categorical predictors. If we use a set of numerical predictors, then it is highly unlikely that multiple records will have identical values on these numerical predictors. Therefore, numerical predictors must be converted to categorical predictors.

The Naïve Bayes Model

- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)
- The Naïve Bayes Assumption: Assume that **all features are independent** given the class label Y . That is to say

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

The approach outlined above amounts to finding all the records in the sample that are exactly like the new record to be classified in the sense that all the predictor values are all identical. This may work well for samples with few predictors, but is impractical when there are multiple predictors. The NB model modifies the above approach by making a key assumption, that of conditional independence.

As such, In the naive Bayes solution, we no longer restrict the probability calculation to those records that match the record to be classified. Instead we use the entire dataset.

Returning to our original basic classification procedure outlined before, recall that the procedure for classifying a new record was:

1. Find all the other records with the same predictor profile (i.e., where the predictor values are the same).
2. Determine the probability that those records belong to the class of interest.

The naive Bayes modification (for the basic classification procedure) is as follows:

1. For class Y_1 , estimate the individual conditional probabilities for each predictor

$P(x_j|Y_1)$ —these are the probabilities that the predictor value in the record to be classified occurs in class Y_1 . For example, for X_1 this probability is estimated by the proportion of x_1 values among the Y_1 records in the training set.

2. Multiply these probabilities by each other, then by the proportion of records belonging to class Y_1 . This gives the probability in the numerator for the equation on the previous page.

3. Repeat Steps 1 and 2 for all the classes. This provides the numerator for all classes Y_i in the equation on the previous page. We can further calculate the actual probability by doing the following:

4. Estimate a probability for class Y_i by taking the value calculated in Step 2 for class Y_i and dividing it by the sum of such values for all classes.

The Naïve Bayes Model II

- Naive Bayes is so called because the independence assumptions we have just made are indeed very naive for a practical scenario.
- Despite this assumption, Naive Bayes classifier often does surprisingly well and is widely used because in many cases it outperforms more sophisticated classification methods.
- Naive Bayesian model is easy to build, with no complicated iterative parameter estimation which makes it particularly useful for very large datasets.

A key assumption in the NB model is that of conditional independence. This is unlikely to be true in real-world settings as predictors are often correlated, but surprisingly, the NB model does well from a practical standpoint. While the posterior probabilities calculated using this approach do not usually match the exact probability, the values calculated are nevertheless ranked similarly to the real probability values. Thus the classification rule to assign a record to the highest probability still makes the correct classification, even if the calculated probability is different from the true probability. In other words, the rank order of a record to a class is correct, even if the exact values used are not.

Example Applications of Naïve Bayes Classifiers

- **Real-time Prediction:** As Naive Bayes is super fast, it can be used for making predictions in real time.
- **Multi-class Prediction:** This algorithm can predict the posterior probability of multiple classes of the target variable.
- **Text classification/ Spam Filtering/ Sentiment Analysis:** Naive Bayes classifiers are mostly used in text classification (due to their better results in multi-class problems and independence rule).

The NB algorithm is fast, and does well. But, as mentioned earlier, it is typically used only when all predictor variables are categorical. Further, since the probabilities are not the exact probability, the NB is not used for example in credit scoring, where we require exact values, rather than rank order. Nevertheless, NB is an efficient algorithm to use, and in the next module, we will see how this model can be implemented, and then adapted to handle some of the above limitations.

“Learning a Naïve Bayes classifier is just a matter of counting how many times each attribute co-occurs with each class ”

Pedro Domingos,
Author of “The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World”

WWW.KENT.EDU

This concludes our module on NB.