

# Naïve Bayes Classifier



Data Science & Analytics  
University of Missouri


# Classification Setting Revisited

- **Given a set of rows/instances/examples**
  - Each row is a set of cols/features
  - Learn a function
    - $\text{class} = f(\text{instance})$
  - And classify
    - New instances (test dataset) into a class
- **Across domains there are different names for these components**

$x_1$	$x_2$	...	$x_n$	$y$
				$y_1$
				$y_2$
				$y_2$
				$y_3$
				$y_1$

# Alternative Terms for Rows

$x_1$	$x_2$	...	$x_n$



- **entities**
- **instances**
- **examples**
- **records**
- **transactions**
- **objects**
- **points**
- **feature-vectors**
- **tuples**

# Alternative Terms for Columns

$x_1$	$x_2$	...	$x_n$

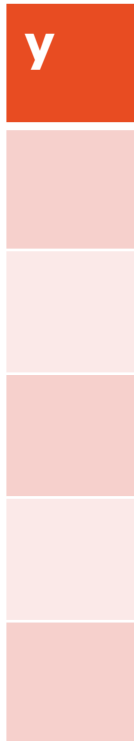
- **attributes**
- **properties**
- **features**
- **dimensions**
- **variables**
- **fields**
- **signal**

# Alternative Terms for Input Variables

$x_1$	$x_2$	...	$x_n$

- independent variable
- predictor
- regressor (regression problem)
- covariate
- manipulated variable
- explanatory variable
- exposure variable (reliability theory)
- risk factor (see medical statistics),
- feature (in machine learning and pattern recognition)
- control variable (econometrics)
- exogenous (economics)

# Alternative Terms for Output variable



- **dependent variable**
- **response variable**
- **regressand (regression)**
- **criterion**
- **predicted variable**
- **measured variable**
- **explained variable**
- **experimental variable**
- **outcome variable**
- **target**
- **class**
- **label**
- **endogenous (economics)**

# Some Types of classifiers

- **Rule-based**
  - Decision tree (M1)
- **Probabilistic**
  - Naïve Bayes classifier (M2)
- **Max-margin classifier**
  - SVM (M5)
- **Neural network (M8)**

# Applying NBC in practice

- **Classical example: text classification**
  - Instances are text samples
    - emails
    - paragraphs
    - sentences
    - documents
    - tweets
    - posts
    - comments
    - reviews
  - Classes are {spam, ~spam}



# Applying NBC in practice

- **Real time prediction**
  - Fast to learn
- **Sentiment analysis**
- **Can be used in multi-class prediction**
  - $\text{Pr}(\text{Class}|\text{instance})$
- **Need less training data**

# Bayes Classifier

- **A probabilistic framework for solving classification problems**
- **Conditional Probability:**  $P(Y | X) = \frac{P(X, Y)}{P(X)}$

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

- **Bayes theorem:**

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$$

# Using Bayes Theorem for Classification

- Consider each attribute and class label as random variables
- Given a record with attributes  $(X_1, X_2, \dots, X_d)$ , the goal is to predict class  $Y$ 
  - Specifically, we want to find the value of  $Y$  that maximizes  $P(Y | X_1, X_2, \dots, X_d)$
- Can we estimate  $P(Y | X_1, X_2, \dots, X_d)$  directly from data?

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Bayesian Classifier: Using Bayes Theorem for Classification

- **Approach:**

- compute posterior probability  $P(Y \mid X_1, X_2, \dots, X_d)$  using the Bayes theorem

$$P(Y \mid X_1 X_2 \dots X_n) = \frac{P(X_1 X_2 \dots X_d \mid Y) P(Y)}{P(X_1 X_2 \dots X_d)}$$

- *Maximum a-posteriori (MAP)*: Choose  $Y$  that maximizes

$$P(Y \mid X_1, X_2, \dots, X_d)$$

- MAP classification rule

- Equivalent to choosing value of  $Y$  that maximizes  $P(X_1, X_2, \dots, X_d \mid Y) P(Y)$

- **How to estimate  $P(X_1, X_2, \dots, X_d \mid Y)$ ?**

# Two Types of Probabilistic Classification

- **Discriminative model**

- Directly estimates the conditional class probability given the features

$$P(Y \mid X_1, X_2, \dots, X_d)$$

- **Generative model**

- Estimates the conditional joint probability of features given the class variable

$$P(X_1, X_2, \dots, X_d \mid Y)$$

# Generative vs Discriminative Classifiers

- **NBC is a generative classifier**
  - because it models  $P(X_1, X_2, \dots, X_n \mid c)$
  - Allows generate new instances by drawing samples the learned joint distribution
    - One can use the Naïve Bayes classifier to “generate” a document for a given class
- **A discriminative classifier, in contrast, will model  $P(\text{class} \mid \text{instance})$ , not  $P(\text{instance} \mid \text{class})$** 
  - Example: logistic regression classifier, coming later

# Why is Bayesian Classifier not Feasible in Practice?

- **Consider 3 classes, 10 features, each with 6 possible discrete instantiations**
  - We need  $3 \cdot 6^{10} + 3$  parameters to be estimated  $\approx 181$  million!

# Enters Naïve Bayes

- **Conditional Independence Assumption**
  - **X** and **Y** are conditionally independent given **Z** if  $P(\mathbf{X}|\mathbf{Y},\mathbf{Z}) = P(\mathbf{X}|\mathbf{Z})$
  - Example: Arm length and reading skills
    - Young child has shorter arm length and limited reading skills, compared to adults
    - If age is fixed, no apparent relationship between arm length and reading skills
    - Arm length and reading skills are conditionally independent given age

**Disclaimer:** By all accounts, Rev. Thomas Bayes was pretty smart. Naivete is on our part, not Bayes!



# Naïve Bayes Classifier

- **Assume independence among attributes  $X_i$  when class is given:**
  - $P(X_1, X_2, \dots, X_d | Y_j) = P(X_1 | Y_j) P(X_2 | Y_j) \dots P(X_d | Y_j)$
  - Now we can estimate  $P(X_i | Y_j)$  for all  $X_i$  and  $Y_j$  combinations from the training data
  - New point is classified to  $Y_j$  if  $P(Y_j) \prod P(X_i | Y_j)$  is maximal.

# Why is this a big deal?

- **Back to the Bayesian classifier: 3 classes, 10 features, each with 6 possible discrete instantiations**
  - Requires 181 million parameters:
- **Naïve Classifier:**
  - Requires  $3 \times 60 + 3 = 183$  parameters!

# Example Data

**Given a Test Record:**

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- We need to estimate  $P(\text{Evade} = \text{Yes} \mid X)$  and  $P(\text{Evade} = \text{No} \mid X)$

In the following we will replace

**Evade = Yes** by **Yes**, and

**Evade = No** by **No**

# Example Data

**Given a Test Record:**  $X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

## Using Bayes Theorem:

$$\square P(\text{Yes} \mid X) = \frac{P(X \mid \text{Yes})P(\text{Yes})}{P(X)}$$

$$\square P(\text{No} \mid X) = \frac{P(X \mid \text{No})P(\text{No})}{P(X)}$$

$\square$  How to estimate  $P(X \mid \text{Yes})$  and  $P(X \mid \text{No})$ ?

# Naïve Bayes on Example Data

**Given a Test Record:**

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$P(X \mid \text{Yes}) =$

$P(\text{Refund} = \text{No} \mid \text{Yes}) \times$

$P(\text{Divorced} \mid \text{Yes}) \times$

$P(\text{Income} = 120\text{K} \mid \text{Yes})$

$P(X \mid \text{No}) =$

$P(\text{Refund} = \text{No} \mid \text{No}) \times$

$P(\text{Divorced} \mid \text{No}) \times$

$P(\text{Income} = 120\text{K} \mid \text{No})$

# Estimate Probabilities from Data

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- **$P(y)$  = fraction of instances of class  $y$**

– e.g.,  $P(\text{No}) = 7/10$ ,  
 $P(\text{Yes}) = 3/10$

- **For categorical attributes:**

$$P(X_i = c | y) = n_c / n$$

- where  $|X_i = c|$  is number of instances having attribute value  $X_i = c$  and belonging to class  $y$
- Examples:

$$P(\text{Status}=\text{Married}|\text{No}) = 4/7$$
$$P(\text{Refund}=\text{Yes}|\text{Yes})=0$$

# Estimate Probabilities from Data

- **For continuous attributes:**
  - **Discretization:** Partition the range into bins:
    - ◆ Replace continuous value with bin value
      - Attribute changed from continuous to ordinal
  - **Probability density estimation:**
    - ◆ Assume attribute follows a normal distribution
    - ◆ Use data to estimate parameters of distribution (e.g., mean and standard deviation)
    - ◆ Once probability distribution is known, use it to estimate the conditional probability  $P(X_i|Y)$

# Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- **Normal distribution:**

$$P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each  $(X_i, Y_i)$  pair

- **For (Income, Class=No):**

- If Class=No

- ◆ sample mean = 110

- ◆ sample variance = 2975

$$P(\text{Income} = 120 \mid \text{No}) = \frac{1}{\sqrt{2\pi(54.54)}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$



# Example of Naïve Bayes Classifier

## Given a Test Record:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} \mid \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} \mid \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} \mid \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

- $$\begin{aligned} P(X \mid \text{No}) &= P(\text{Refund}=\text{No} \mid \text{No}) \\ &\quad \times P(\text{Divorced} \mid \text{No}) \\ &\quad \times P(\text{Income}=120\text{K} \mid \text{No}) \\ &= 4/7 \times 1/7 \times 0.0072 = 0.0006 \end{aligned}$$

- $$\begin{aligned} P(X \mid \text{Yes}) &= P(\text{Refund}=\text{No} \mid \text{Yes}) \\ &\quad \times P(\text{Divorced} \mid \text{Yes}) \\ &\quad \times P(\text{Income}=120\text{K} \mid \text{Yes}) \\ &= 1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times \end{aligned}$$

$10^{-10}$

Since  $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore  $P(\text{No}|X) > P(\text{Yes}|X)$

$\Rightarrow \text{Class} = \text{No}$

# Make Prediction with Partial Information in the Test Set

**Even in absence of information about any attributes, we can use Apriori Probabilities of Class Variable:**

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} \mid \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} \mid \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} \mid \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

$$P(\text{Yes}) = 3/10$$

$$P(\text{No}) = 7/10$$

**If we only know that marital status is Divorced, then:**

$$P(\text{Yes} \mid \text{Divorced}) = 1/3 \times 3/10 / P(\text{Divorced})$$

$$P(\text{No} \mid \text{Divorced}) = 1/7 \times 7/10 / P(\text{Divorced})$$

**If we also know that Refund = No, then**

$$P(\text{Yes} \mid \text{Refund} = \text{No}, \text{Divorced}) = 1 \times 1/3 \times 3/10 / P(\text{Divorced}, \text{Refund} = \text{No})$$

$$P(\text{No} \mid \text{Refund} = \text{No}, \text{Divorced}) = 4/7 \times 1/7 \times 7/10 / P(\text{Divorced}, \text{Refund} = \text{No})$$

**If we also know that Taxable Income = 120, then**

$$P(\text{Yes} \mid \text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120) = 1.2 \times 10^{-9} \times 1 \times 1/3 \times 3/10 / P(\text{Divorced}, \text{Refund} = \text{No}, \text{Income} = 120)$$

$$P(\text{No} \mid \text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120) = 0.0072 \times 4/7 \times 1/7 \times 7/10 / P(\text{Divorced}, \text{Refund} = \text{No}, \text{Income} = 120)^{26}$$

# Issues with Naïve Bayes Classifier

Consider the table with Tid = 7 deleted

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} \mid \text{No}) = 2/6$$

$$P(\text{Refund} = \text{No} \mid \text{No}) = 4/6$$

$$\rightarrow P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} \mid \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/6$$

$$\rightarrow P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 0$$

$$P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/6$$

$$P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0/3$$

For Taxable Income:

If class = No: sample mean = 91

sample variance = 685

If class = No: sample mean = 90

sample variance = 25

Given  $X = (\text{Refund} = \text{Yes}, \text{Divorced}, 120\text{K})$

$$P(X \mid \text{No}) = 2/6 \times 0 \times 0.0083 = 0$$

$$P(X \mid \text{Yes}) = 0 \times 1/3 \times 1.2 \times 10^{-9} = 0$$

Naïve Bayes will not be able to  
classify  $X$  as Yes or No!

# Issues with Naïve Bayes Classifier

- **If one of the conditional probabilities is zero, then the entire expression becomes zero**
- **Need to use other estimates of conditional probabilities than simple fractions**
- **Probability estimation:**

original:  $P(X_i = c|y) = \frac{n_c}{n}$

Laplace Estimate:  $P(X_i = c|y) = \frac{n_c + 1}{n + v}$

m – estimate:  $P(X_i = c|y) = \frac{n_c + mp}{n + m}$

$n$ : number of training instances belonging to class  $y$

$n_c$ : number of instances with  $X_i = c$  and  $Y = y$

$v$ : total number of attribute values that  $X_i$  can take

$p$ : initial estimate of  $P(X_i = c|y)$  known apriori

$m$ : hyper-parameter for our confidence in  $p$

# Sklearn naïve Bayes

- [https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)

`naive_bayes.BernoulliNB(*  
[, alpha, ...])` Naive Bayes classifier for multivariate Bernoulli models.

`naive_bayes.CategoricalNB(*  
[, alpha, ...])` Naive Bayes classifier for categorical features

`naive_bayes.ComplementNB(*  
[, alpha, ...])` The Complement Naive Bayes classifier described in Rennie et al.

`naive_bayes.GaussianNB(*  
[, priors, ...])` Gaussian Naive Bayes (GaussianNB)

`naive_bayes.MultinomialNB(*  
[, alpha, ...])` Naive Bayes classifier for multinomial models

# Conclusions

- **Naïve Bayes based on the independence assumption**
  - Training is very easy and fast; just requiring considering each attribute in each class separately
    - Works with less data
  - Test is straightforward; just looking up tables or calculating conditional probabilities with normal distributions
- **A popular generative model**
  - Performance competitive to most of state-of-the-art classifiers even in presence of violating independence assumption
  - Many successful applications, e.g., spam mail filtering
  - Apart from classification, naïve Bayes can do more...