

# Probabilistic Classification

**Tozammel Hossain**



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

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

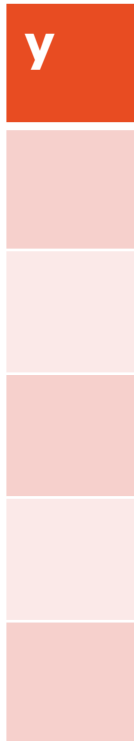
- **attributes**
- **properties**
- **features**
- **dimensions**
- **variables**
- **fields**

# 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)**

# Classification of classifiers

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

# Probabilistic Classification

- **Establishing a probabilistic model for classification**

- Discriminative model

$$P(C | \mathbf{X}) \quad C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_n)$$

- Generative model

$$P(\mathbf{X} | C) \quad C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_n)$$

- **MAP classification rule**

- **MAP**: **M**aximum **A** **P**osterior
- Assign  $x$  to  $c^*$  if

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}) \quad c \neq c^*, c = c_1, \dots, c_L$$



# A Bayesian Classifier

- **Generative classification with the MAP rule**

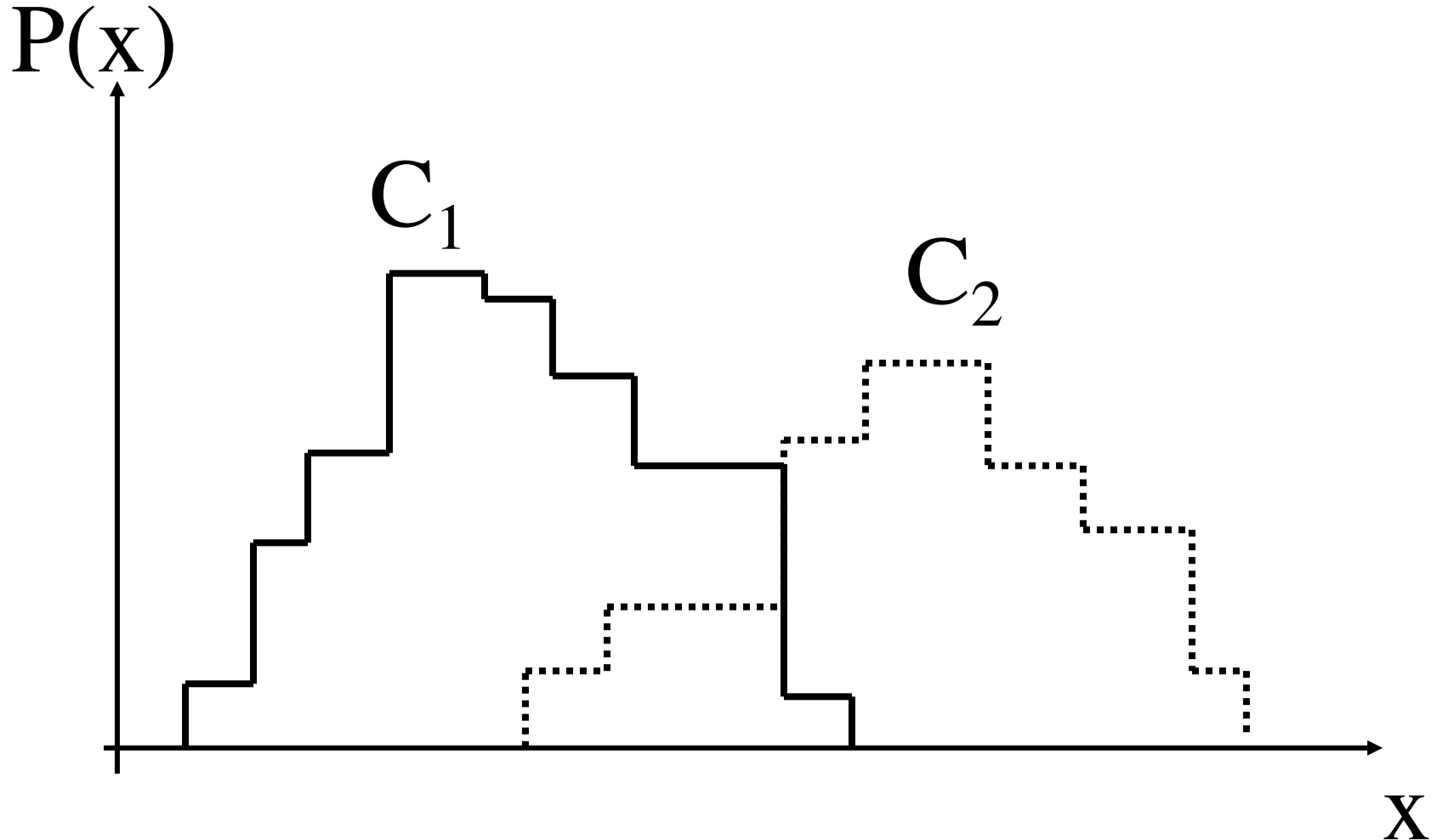
- Apply Bayesian rule to convert

$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})} \propto P(\mathbf{X} | C)P(C)$$

- MAP Classification Rule

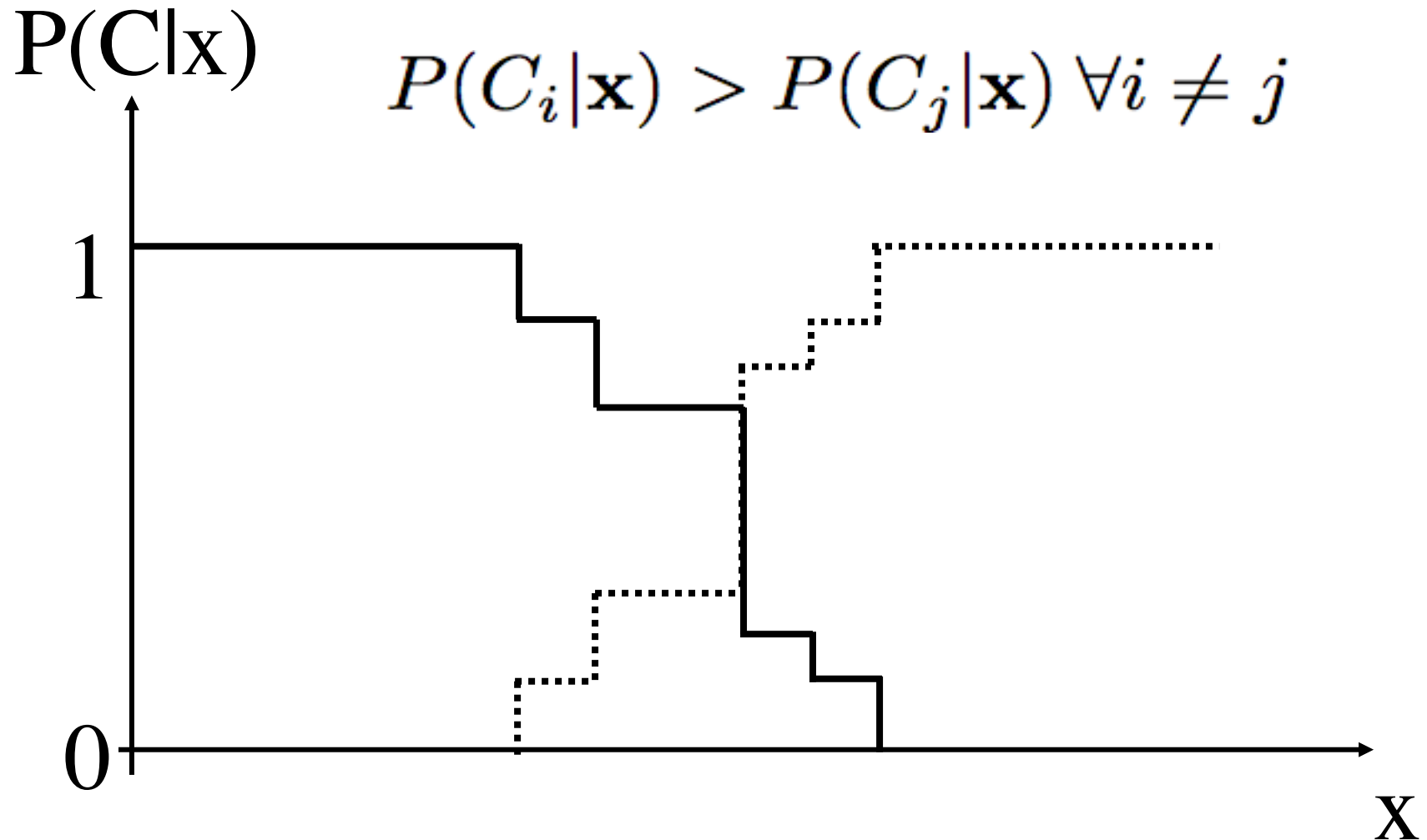
$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}) \quad c \neq c^*, c = c_1, \dots, c_L$$

# Example: A Feature Histograms



Slide by Stephen Marsland

# Example: Posterior Probability



# A Bayesian Classifier (Contd.)

$$P(C | \mathbf{X}) \propto P(\mathbf{X} | C)P(C) = P(X_1, \dots, X_n | C)P(C)$$

- **Why is this 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!

# Enter Naïve Bayes

- **Naïve Bayes assumption: Features are conditionally independent given the class**

$$\begin{aligned}P(X_1, X_2, \dots, X_n | C) &= P(X_1 | X_2, \dots, X_n; C)P(X_2, \dots, X_n | C) \\&= P(X_1 | C)P(X_2, \dots, X_n | C) \\&= P(X_1 | C)P(X_2 | C) \dots P(X_n | C)\end{aligned}$$

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

# Why is this a big deal?

- **Back to: 3 classes, 10 features, each with 6 possible discrete instantiations**
- **Requires 181 million parameters:**

$$\begin{aligned} C_{MAP} &= \operatorname{argmax}_{c \in C} P(C = c | X) \\ &= \operatorname{argmax}_{c \in C} P(X_1, X_2, \dots, X_n | C = c) P(C = c) \end{aligned}$$

- **Requires  $3 \times 60 + 3 = 183$  parameters!**

$$C_{NB} = \operatorname{argmax}_{c \in C} P(C = c) \prod_{i=1}^n P(X_i | C = c)$$

# Example: Play Tennis Data

*PlayTennis: training examples*

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

**How many parameters are required for a naïve Bayes classifier?**

$$= (3 + 3 + 2 + 2) * 2 + 2 = 22$$

# Learning Phase

$$P(\text{Outlook}=o \mid \text{Play}=b)$$

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

$$P(\text{Temperature}=t \mid \text{Play}=b)$$

Temperature	Play=Yes	Play=No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

$$P(\text{Humidity}=h \mid \text{Play}=b)$$

Humidity	Play=Yes	Play=No
High	3/9	4/5
Normal	6/9	1/5

$$P(\text{Wind}=w \mid \text{Play}=b)$$

Wind	Play=Yes	Play=No
Strong	3/9	3/5
Weak	6/9	2/5

$$P(\text{Play}=Yes) = 9/14$$

$$P(\text{Play}=No) = 5/14$$



# Example

- **Test Phase**

- Given a new instance,  $\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$
- Look up tables

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{Yes}) = 2/9$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Play}=\textit{Yes}) = 9/14$$

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{No}) = 1/5$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{No}) = 4/5$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Play}=\textit{No}) = 5/14$$

# Example

- **MAP Rule**

$P(\text{Yes} | \mathbf{x}')$ :  $[P(\text{Sunny} | \text{Yes})P(\text{Cool} | \text{Yes})P(\text{High} | \text{Yes})P(\text{Strong} | \text{Yes})]P(\text{Play}=\text{Yes}) = 0.0053$

$P(\text{No} | \mathbf{x}')$ :  $[P(\text{Sunny} | \text{No}) P(\text{Cool} | \text{No})P(\text{High} | \text{No})P(\text{Strong} | \text{No})]P(\text{Play}=\text{No}) = 0.0206$

Given the fact  $P(\text{Yes} | \mathbf{x}') < P(\text{No} | \mathbf{x}')$ , we label  $\mathbf{x}'$  to be “No”.

# Applying NBC in practice

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

# Sklearn naïve Bayes

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

<code>naive_bayes.BernoulliNB(* [, alpha, ...])</code>	Naive Bayes classifier for multivariate Bernoulli models.
<code>naive_bayes.CategoricalNB(* [, alpha, ...])</code>	Naive Bayes classifier for categorical features
<code>naive_bayes.ComplementNB(* [, alpha, ...])</code>	The Complement Naive Bayes classifier described in Rennie et al.
<code>naive_bayes.GaussianNB(* [, priors, ...])</code>	Gaussian Naive Bayes (GaussianNB)
<code>naive_bayes.MultinomialNB(* [, alpha, ...])</code>	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
  - 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...