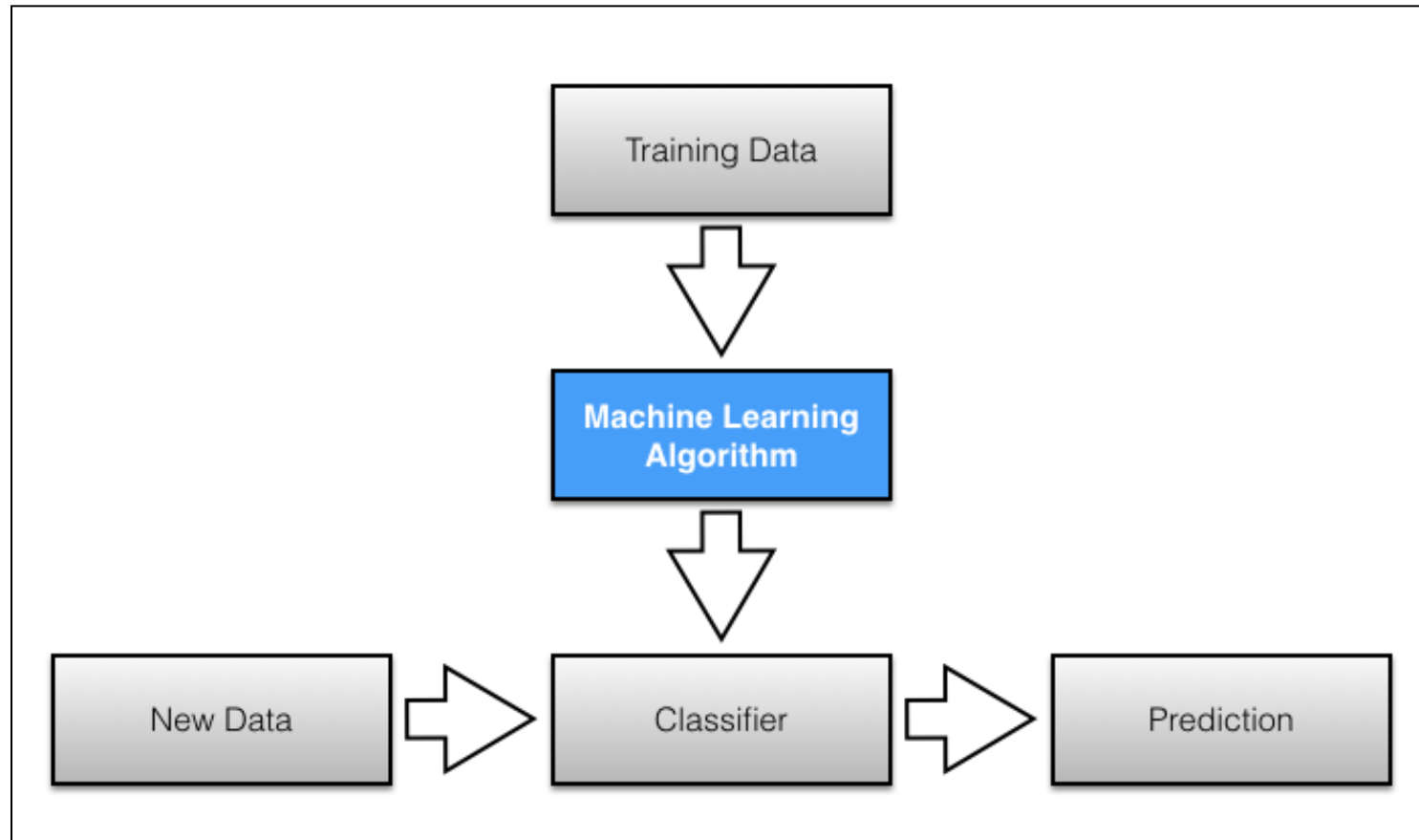




Naïve Bayes Classifier



Machine Learning



Naïve Bayes

- Naive Bayes is a family of probabilistic algorithms that take advantage of probability theory and Bayes' Theorem to predict the probability of being in a class.
 - It has been successfully used for many purposes but it works particularly well with natural language processing (NLP) problems.
-

Multinomial Naïve Bayes:

(for discrete features)

Bayes' Theorem:

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

Example:

Sample Data Set:

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
overcast	cool	normal	true	yes

outlook	temp.	humidity	windy	play
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

Example:

Features:

- Outlook (sunny, overcast, rainy)
- Temperature (hot, mild, cool)
- Humidity (high, normal)
- Windy (true, false)

Class:

- Play
 - Not Play
-

Example:

Frequencies and Probabilities for the Data Set:

outlook			temperature			humidity			windy			play	
	yes	no		yes	no		yes	no		yes	no	yes	no
sunny	2	3	hot	2	2	high	3	4	false	6	2	9	5
overcast	4	0	mild	4	2	normal	6	1	true	3	3		
rainy	3	2	cool	3	1								
	yes	no		yes	no		yes	no		yes	no	yes	no
sunny	2/9	3/5	hot	2/9	2/5	high	3/9	4/5	false	6/9	2/5	9/14	5/14
overcast	4/9	0/5	mild	4/9	2/5	normal	6/9	1/5	true	3/9	3/5		
rainy	3/9	2/5	cool	3/9	1/5								

Example:

■ Classifying an Unseen Example:

outlook	temp.	humidity	windy	play
sunny	cool	high	true	?

$$P(\text{yes}) = 2/9 \cdot 3/9 \cdot 3/9 \cdot 3/9 = 0.0082$$

$$P(\text{no}) = 3/5 \cdot 1/5 \cdot 4/5 \cdot 3/5 = 0.0577$$

Gaussian Naïve Bayes:

(for continuous features)

- If there are n number of classes ($C_1, C_2, \dots, C_k, \dots, C_n$), x is a feature(continuous) and for test data value of feature x is v ($x=v$) then probability of $x=v$ being in class C_k is:

$$p(x = v | C_k) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(v-\mu_k)^2}{2\sigma_k^2}}$$

NLP Example:

Building a classifier that says whether a text is about sports or not.:

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

"A very close game"=? (Sports / Not sports)

Thank You...
