

Representing Uncertainty + Probabilistic Learning

R&N Chapter 13

A bit of 20.2

Uncertainty

- Most real-world problems deal with uncertain information
 - Diagnosis: Likely disease given observed symptoms
 - Equipment repair: Likely component failure given sensor reading
 - Help desk: Likely operation based on past operations
 - Cannot be represented by deterministic rules
Headache => Fever

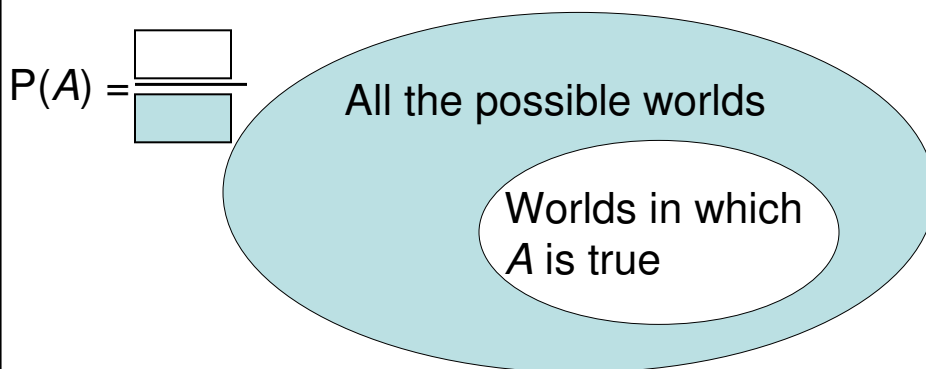
Uncertainty

- Correct framework for representing uncertainty: Probability
- Outline:
 - Review of basic probability tools (much of it well-known, but still important to review)
 - Bayes rule and its use in uncertain reasoning and probabilistic learning

Probability

- $P(A)$ = Probability of event A = percentage of all possible worlds in which A is true.

$$0 \leq P(A) \leq 1$$



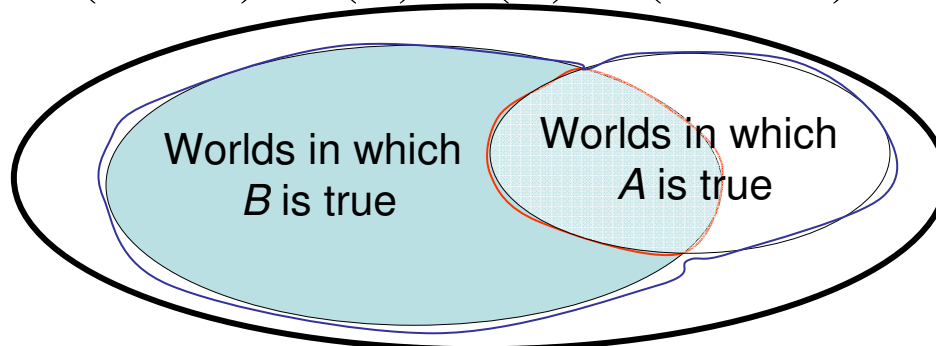
Probability

$$0 \leq P(A) \leq 1$$

$$P(\mathbf{True}) = 1$$

$$P(\mathbf{False}) = 0$$

$$P(A \text{ or } \mathbf{B}) = P(A) + P(\mathbf{B}) - P(A \text{ and } \mathbf{B})$$



Probability

$$0 \leq P(A) \leq 1$$

$$P(\mathbf{True}) = 1$$

$$P(\mathbf{False}) = 0$$

$$P(A \text{ or } \mathbf{B}) = P(A) + P(\mathbf{B}) - P(A \text{ and } \mathbf{B})$$

- Other ideas:
 - Fuzzy logic
 - Non-monotonic logic
 - Multi-valued logic
 - Evidence theory (Dempster-Shafer)
- Probability is the only system that is “consistent”

Probability

- Immediately derived properties

$$P(\neg A) = 1 - P(A)$$

Denotes not- A = All the worlds in which A does not occur

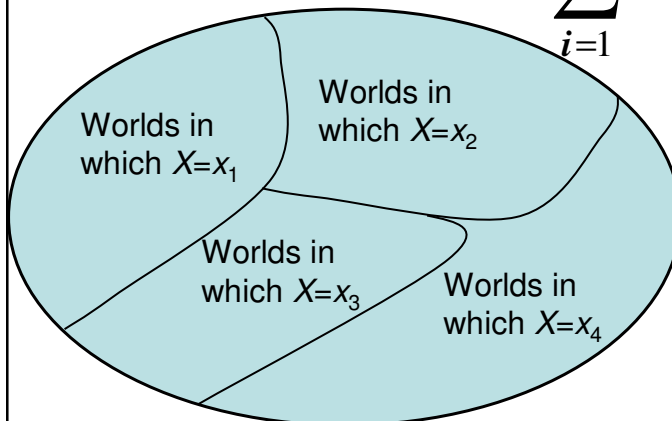
$$P(A) = P(A, B) + P(A, \neg B)$$

Short hand for “ A and B ”

Probability

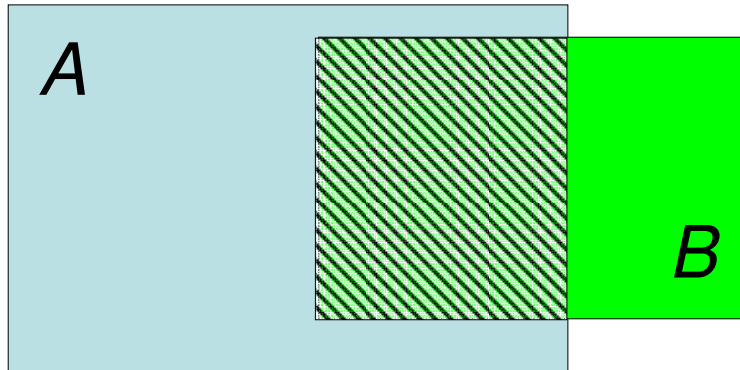
- A random variable is a variable X that can take values x_1, \dots, x_n with a probability $P(X = x_i)$ attached to each $i = 1, \dots, n$

$$\sum_{i=1}^n P(X = x_i) = 1$$



Conditional Probability

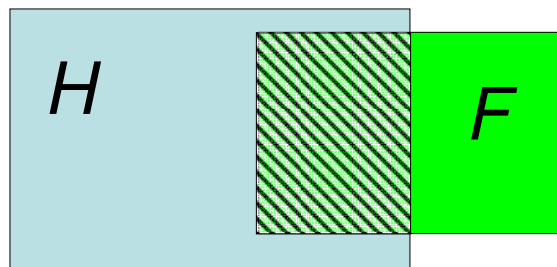
- $P(A|B)$ = Fraction of those worlds in which B is true for which A is also true.



Conditional Probability Example

- H = Headache $P(H) = 1/2$
- F = Flu $P(F) = 1/8$

$$P(H|F) = 1/2$$



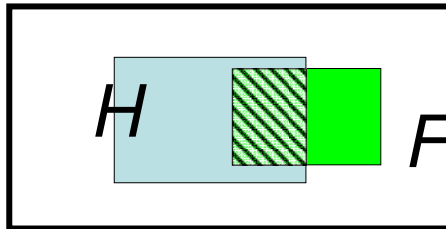
Conditional Probability Example

- $H = \text{Headache}$ $P(H) = 1/2$
- $F = \text{Flu}$ $P(F) = 1/8$

$$P(H|F) = \frac{\text{(Area of "H and F" region)}}{\text{(Area of F region)}}$$

$$P(H|F) = P(H, F)/P(F)$$

$$P(H|F) = 1/2$$



Conditional Probability

- Definition:

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

- Chain rule:

$$P(A, B) = P(A | B) P(B)$$

Conditional Probability

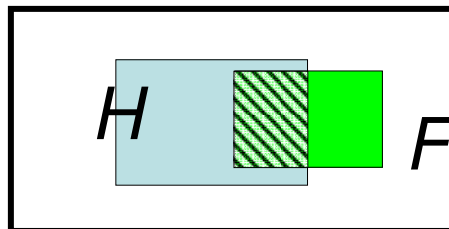
- Other useful relations:

$$P(A | B) + P(\neg A | B) = 1$$

$$\sum_i P(X = x_i | B) = 1$$

Probabilistic Inference

- Suppose H is true
 - Suppose that you know $P(H|F) = 1/2 = 0.5$
 - What is the probability that F is true? 0.5?
- $P(H) = 1/2$
 - $P(F) = 1/8$
 - $P(H|F) = 0.5$



Probabilistic Inference

- Correct reasoning:
- We know $P(H)$, $P(F)$, $P(H|F)$ and the two chain rules:

$$P(H, F) = P(H | F) P(F)$$

$$P(F | H) = \frac{P(H, F)}{P(H)}$$

- Substituting the values:

$$P(H, F) = 0.5 \times 1/8 = 1/16$$

$$P(F | H) = \frac{1/16}{1/2} = 1/8$$

Probabilistic Inference

- Correct reasoning:
- We know $P(H)$, $P(F)$, $P(H|F)$ and the two chain rules:

$$P(H, F) = P(H | F) P(F)$$

$$P(F | H) = \frac{P(H, F)}{P(H)}$$

- Substituting the values:

$$P(H, F) = 0.5 \times 1/8 = 1/16$$

$$P(F | H) = \frac{1/16}{1/2} = 1/8$$

The key difference is that we took into account the fact that catching the flu is unlikely ($P(F)$ is small)

Bayes Rule

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

Introduced circa 1763

Probabilistic Inference

We want: *Posterior* probability that B occurs given that A occurs

We know: *Prior* probability that B occurs in the absence of any other information

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

We know: *Likelihood* that A occurs given that B occurs

Bayes Rule

- What if we do not know $P(A)$???
- Use the relation:

$$P(A) = P(A | B)P(B) + P(A | \neg B)P(\neg B)$$

- More general Bayes rule:

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

Bayes Rule

- Same rule for a non-binary random variable, except we need to sum over all the possible events $X = x_i$

$$P(X = x_i | A) = \frac{P(A | X = x_i)P(X = x_i)}{P(A)}$$

$$P(X = x_i | A) = \frac{P(A | X = x_i)P(X = x_i)}{\sum_k P(A | X = x_k)P(X = x_k)}$$

This is actually just $P(A)$

Joint Distribution

- Joint Distribution Table:

- Given a set of variables A,B,C,....
- Generate a table with all the possible combinations of assignments to the variables in the rows
- For each row, list the corresponding joint probability
- For M binary variables \rightarrow size 2^M

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Using the Joint Distribution: Computing Other Probabilities

Compute the probability of event E :

$$P(E) = \sum_{\substack{\text{all rows} \\ \text{containing } E}} P(\text{row})$$

$$P(A, B) = 0.25 + 0.10 = 0.35$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Using the Joint Distribution: Doing Inference

Given that event E_1 occurs, what is the probability that E_2 occurs:

$$P(E_2 | E_1) = \frac{P(E_2, E_1)}{P(E_1)}$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Using the Joint Distribution: Doing Inference

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

$$= \frac{0.10}{0.05+0.05+0.10+0.10} = \frac{0.10}{0.30}$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Inference

- General view: I have some evidence (Headache) how likely is a particular conclusion (Fever)
- Important in many industries: Medical, pharmaceutical, Help Desk, Fault Diagnosis....

Learning the Joint Distribution

- Three possible ways of generating the joint distribution:
 1. Human experts (very difficult!)
 2. Using known conditionally probabilities (e.g., if we know $P(C|A,B)$, $P(B|A)$, and $P(A)$, we know $P(A,B,C) = P(C|A,B)P(B|A)P(A) \rightarrow$ This is the basis for *Bayes Nets*, to be covered later....)
 3. *Learning from data*

Learning the Joint Distribution

Suppose that we have recorded a lot of training data:

(0,1,1)
 (1,0,1)
 (1,1,0)
 (0,0,0)
 (1,1,0).....

The entry for $P(A,B,-C)$ in the table is:

A	B	C	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

$$\frac{\text{\# of data entries with } A=1, B=1, C=0}{\text{Total number of data entries}}$$

Learning the Joint Distribution

Suppose that we have recorded a lot of training data:

(0,1,1)
 (1,0,1)
 (1,1,0)
 (0,0,0)
 (1,1,0).....









More generally, the entry for $P(E)$ in the table is:

A	B	C	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

$$\frac{\text{\# of data entries with } E}{\text{Total number of data entries}}$$









Real-Life Joint Distribution

- UCI Census Database

gender	hours_worked	wealth		
Female	v0:40.5-	poor	0.253122	
		rich	0.0245895	
	v1:40.5+	poor	0.0421768	
		rich	0.0116293	
Male	v0:40.5-	poor	0.331313	
		rich	0.0971295	
	v1:40.5+	poor	0.134106	
		rich	0.105933	

Real-Life Joint Distribution

- UCI Census Database

gender	hours_worked	wealth		
Female	v0:40.5-	poor	0.253122	
		rich	0.0245895	
	v1:40.5+	poor	0.0421768	
		rich	0.0116293	
Male	v0:40.5-	poor	0.331313	
		rich	0.0971295	
	v1:40.5+	poor	0.134106	
		rich	0.105933	

$$P(\text{Male}|\text{Poor}) = 0.4654/0.7604 = 0.612$$

So Far....

- Basic probability concepts
- Bayes rule
- What are joint distributions
- Inference using joint distributions
- Learning joint distributions from data

- Problem: If we have M variables, we need 2^M entries in the joint distribution table → An independence assumption leads to an efficient way to learn and to do inference

Independence

- A and B are independent iff:

$$P(A | B) = P(A)$$

- In words: Knowing B does not affect how likely we think that A is true

Key Properties

- Symmetry:

$$P(\mathbf{A} \mid \mathbf{B}) = P(\mathbf{A}) \Leftrightarrow P(\mathbf{B} \mid \mathbf{A}) = P(\mathbf{B})$$

- Joint distribution:

$$P(\mathbf{A}, \mathbf{B}) = P(\mathbf{A})P(\mathbf{B})$$

- Independence of complements:

$$P(\neg \mathbf{A} \mid \mathbf{B}) = P(\neg \mathbf{A}) \quad P(\mathbf{A} \mid \neg \mathbf{B}) = P(\mathbf{A})$$

Naïve Bayes

- Suppose that A, B, C are *independent*
- Then any value of the joint distribution can be computed easily:

$$P(\mathbf{A}, \mathbf{B}, \mathbf{C}) = P(\mathbf{A})P(\mathbf{B})P(\mathbf{C})$$

$$P(\mathbf{A}, \neg \mathbf{B}, \mathbf{C}) = P(\mathbf{A})P(\neg \mathbf{B})P(\mathbf{C})$$

- In fact, we need only M numbers instead of 2^M for binary variables!!

Naïve Bayes: General Case

- If X_1, \dots, X_M are independent variables:

$$P(X_1 = x_1, X_2 = x_2, \dots, X_M = x_M) =$$

$$P(X_1 = x_1)P(X_2 = x_2) \dots P(X_M = x_M)$$

- Under the “Naïve” assumption, we can compute any value of the joint distribution
- We can answer any inference query
- How do we learn the distributions?

Naïve Bayes: Learning

$$P(X_i = x) = \frac{\text{Number of observations with } X_i = x}{\text{Total Number of observations}}$$

- Learning the distributions from data is simple and efficient
- In practice, the independence assumption may not be met but it is often a very useful approximation (see examples at the end)

So Far....

- Basic probability concepts
- Bayes rule
- What are joint distributions
- Inference using joint distributions
- Learning joint distributions from data
- Independence assumption
- Naïve Bayes

- Problem: We now have the joint distribution.
How can we use it to make decision → Bayes Classifier

Problem Example

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall,short}
 - Country = {Gromland, Polvia}
- Training data: Values of (Eye,Height,Country) collected over population

Joint Distribution Table:

(B,T,G)	(B,T,P)
(D,T,G)	(B,T,P)
(D,T,G)	(B,T,P)
(D,T,G)	(D,T,P)
(B,T,G)	(D,T,P)
(B,S,G)	(D,S,P)
(B,S,G)	(B,S,P)
(D,S,G)	(D,S,P)

Learn Joint Probabilities

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall,short}
 - Country = {Gromland, Polvia}
- Training data: Values of (Eye,Height,Country) collected over population

(B,T,G)	(B,T,P)	$P(B,S,G) = 2/16$
(D,T,G)	(B,T,P)	$P(B,T,G) = 2/16$
(D,T,G)	(B,T,P)	$P(D,S,G) = 1/16$
(D,T,G)	(D,T,P)	$P(D,T,G) = 3/16$
(B,T,G)	(D,T,P)	$P(B,S,P) = 1/16$
(B,S,G)	(D,S,P)	$P(B,T,P) = 3/16$
(B,S,G)	(B,S,P)	$P(D,S,P) = 2/16$
(D,S,G)	(D,S,P)	$P(D,T,P) = 2/16$

Compute other Joint Or Conditional Distributions

$P(B,S,G) = 2/16$	$P(\text{Hair} = B, \text{Height} = S \text{Country} = G) = \frac{P(\text{Hair} = B, \text{Height} = S, \text{Country} = G)}{P(\text{Country} = G)}$ $\frac{2/16}{1/2} = 4/16$
$P(B,T,G) = 2/16$	
$P(D,S,G) = 1/16$	
$P(D,T,G) = 3/16$	
$P(B,S,P) = 1/16$	
$P(B,T,P) = 3/16$	
$P(D,S,P) = 2/16$	
$P(D,T,P) = 2/16$	

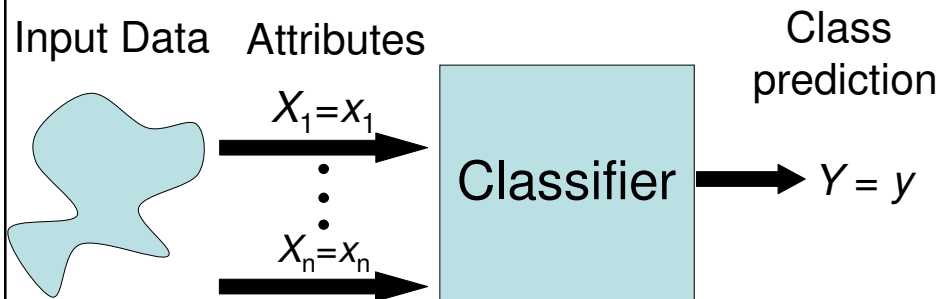
Classifier Example

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall,short}
 - Country = {Gromland, Polvia}
- Training data: Values of (Eye,Height,Country) collected over population

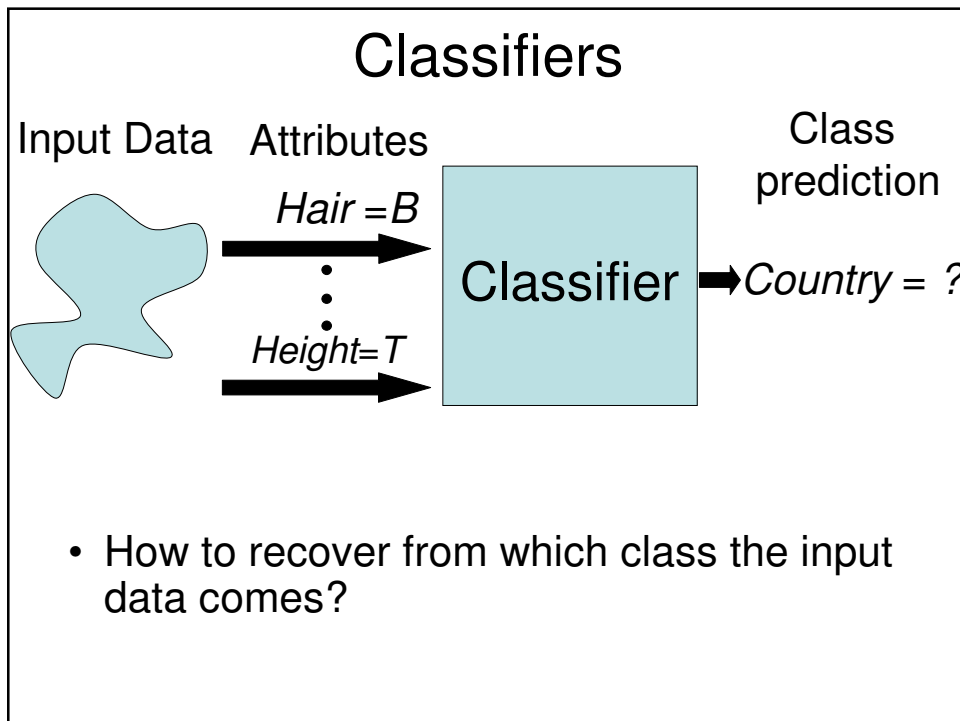
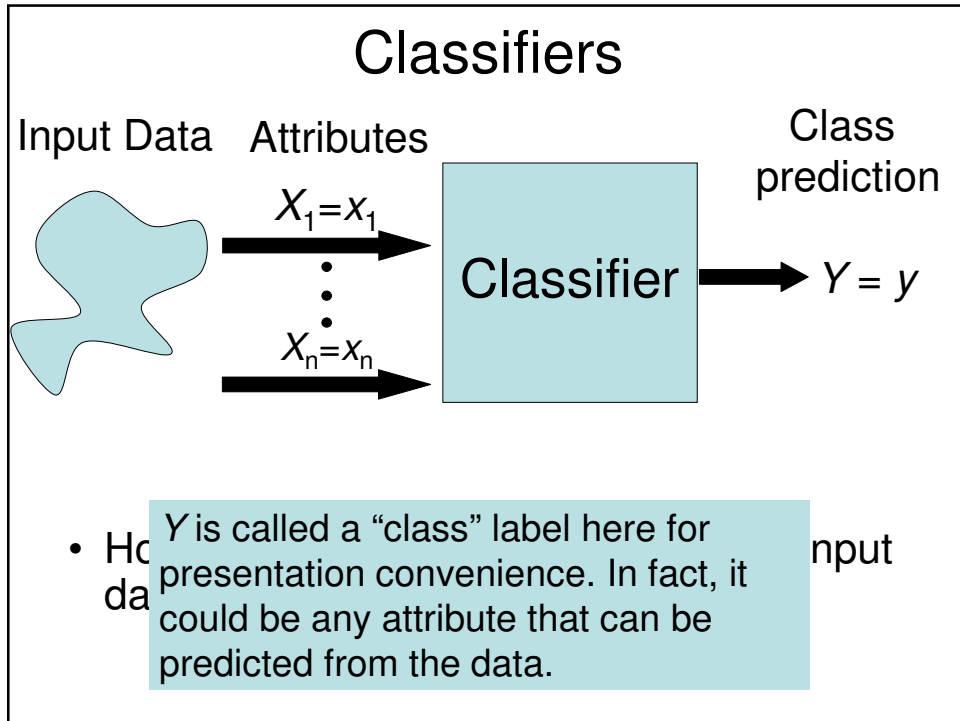
(B,T,G) (B,T,P)
(D,T,G) (B,T,P)
(D,T,G) (B,T,P)
(D,T,G) (D,T,P)
(B,T,G) (D,T,P)
(B,S,G) (D,S,P)
(B,S,G) (B,S,P)
(D,S,G) (D,S,P)

If I observe a new individual tall with blond hair, what is the most likely country of origin?

Classifiers



- How to recover from which class the input data comes?



Classifiers

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum:

$$P(Y = y \mid X_1 = x_1, \dots, X_n = x_n)$$

Classifiers

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum:

$$P(Y = y \mid X_1 = x_1, \dots, X_n = x_n)$$

The maximum is called the *Maximum A Posteriori (MAP)* estimator

Classifiers

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum:

$$P(Y = y \mid X_1 = x_1, \dots, X_n = x_n) = \frac{P(X_1 = x_1, \dots, X_n = x_n \mid Y = y) P(Y = y)}{P(X_1 = x_1, \dots, X_n = x_n)}$$

Classifiers

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum.

Apply Bayes rule

$$P(Y = y \mid X_1 = x_1, \dots, X_n = x_n) = \frac{P(X_1 = x_1, \dots, X_n = x_n \mid Y = y) P(Y = y)}{P(X_1 = x_1, \dots, X_n = x_n)}$$

This denominator does not depend on y . It is a constant (as far as y is concerned) and can be ignored.

Bayes Classifier

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum:

$$P(X_1 = x_1, \dots, X_n = x_n | Y = y) P(Y = y)$$

Bayes Classifier

- We want to find the value of Y that is the most probable, given the observations X_1, \dots, X_n
- Find y such that this is maximum:

$$P(X_1 = x_1, \dots, X_n = x_n | Y = y) P(Y = y)$$

Likelihood of observing (x_1, \dots, x_n) from data of class y . This is learned from training data

Probability of each class, also learned from training data

Bayes Classifier

- Learning:
 - Collect all the observations (x_1, \dots, x_n) for each class y and estimate:

$$P(X_1 = x_1, \dots, X_n = x_n | Y = y) = \frac{\text{\# observations with } (X_1 = x_1, \dots, X_n = x_n) \text{ in class } y}{\text{Total Number of observations in class } y}$$

$$P(Y = y) = \frac{\text{\# observations in class } y}{\text{Total Number of observations}}$$

- Classification:
 - Given a new input (x_1, \dots, x_n) , compute the best class:

$$y^{best} = \arg \max_y P(X_1 = x_1, \dots, X_n = x_n | Y = y) P(Y = y)$$

Classifier Example

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}
- Training data: Values of (E collected over population

If I observe a new individual tall with blond hair, what is the most likely country of origin?

(B,T,G)	(B,T,P)	
(D,T,G)	(B,T,P)	$P(B,T G)P(G) = 2/8 \times 1/2 = 2/16$
(D,T,G)	(B,T,P)	$P(B,T P)P(P) = 3/8 \times 1/2 = 3/16$
(D,T,G)	(D,T,P)	
(B,T,G)	(D,T,P)	Conclusion: Country = P
(B,S,G)	(D,S,P)	
(B,S,G)	(B,S,P)	
(D,S,G)	(D,S,P)	

Classifier Example

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}
- Training data: Values of (E collected over population

If I observe a new individual tall with blond hair, what is the most likely country of origin?

(B,T,G)	(B,T,G)	(B,T,P)	$P(B,T G)P(G) = 2/8 \times 2/3 = 4/24$
(D,T,G)	(D,T,G)	(B,T,P)	$P(B,T P)P(P) = 3/8 \times 1/3 = 3/24$
(D,T,G)	(D,T,G)	(B,T,P)	
(D,T,G)	(D,T,G)	(D,T,P)	
(B,T,G)	(B,T,G)	(D,T,P)	
(B,S,G)	(B,S,G)	(D,S,P)	
(B,S,G)	(B,S,G)	(B,S,P)	
(D,S,G)	(D,S,G)	(D,S,P)	

Conclusion: Country = G

Classifier Example

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}
- Training data: Values of (E collected over population

Note the different conclusion! That's where the "Bayes" part plays a role. We correctly took into account the fact that Gromland is now twice as likely, irrespective of the observation. $P(G) = 2/3$ $P(P) = 1/3$

(B,T,G)	(B,T,G)	(B,T,P)	$P(B,T G)P(G) = 2/8 \times 2/3 = 4/24$
(D,T,G)	(D,T,G)	(B,T,P)	$P(B,T P)P(P) = 3/8 \times 1/3 = 3/24$
(D,T,G)	(D,T,G)	(B,T,P)	
(D,T,G)	(D,T,G)	(D,T,P)	
(B,T,G)	(B,T,G)	(D,T,P)	
(B,S,G)	(B,S,G)	(D,S,P)	
(B,S,G)	(B,S,G)	(B,S,P)	
(D,S,G)	(D,S,G)	(D,S,P)	

Conclusion: Country = G

Naïve Bayes Classifier

- Learning: Collect all the observations (x_1, \dots, x_n) for each class y and estimate:

$$P(X_i = x_i | Y = y) = \frac{\text{Number of observations with } X_i = x_i \text{ in class } y}{\text{Total Number of observations in class } y}$$

$$P(Y = y) = \frac{\text{Number of observations in class } y}{\text{Total Number of observations}}$$

- Classification:

$$y^{best} = \arg \max_y P(X_1 = x_1 | Y = y) \dots P(X_n = x_n | Y = y) P(Y = y)$$

Note that we need only $k \times n$ numbers $(P(X_i = x_i | Y = y))$ to implement this classifier, instead of k^n if we were to use the full model, without independence assumption.

- Learning: Collect all the observations (x_1, \dots, x_n) for each class y and estimate:

$$P(X_i = x_i | Y = y) = \frac{\text{Number of observations with } X_i = x_i \text{ in class } y}{\text{Total Number of observations in class } y}$$

$$P(Y = y) = \frac{\text{Number of observations in class } y}{\text{Total Number of observations}}$$

- Classification:

$$y^{best} = \arg \max_y P(X_1 = x_1 | Y = y) \dots P(X_n = x_n | Y = y) P(Y = y)$$

Naïve Bayes Implementation

- Small (but important) implementation detail: If n is large, we may be taking the product of a large number of small floating-point values \rightarrow underflow \rightarrow avoided by taking the log.
- Take the max of:

$$\log P(X_1 = x_1 | Y = y) + \dots \\ + \log P(X_n = x_n | Y = y) + \log P(Y = y)$$

- Instead of:

$$P(X_1 = x_1 | Y = y) \dots P(X_n = x_n | Y = y) P(Y = y)$$

Same Example, the Naïve Bayes Way

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}
- Training data: Values of (Eye, Height, Country) collected over population

(B,T,G)	(B,T,G)	(B,T,P)	$P(B,T G)P(G) \approx P(B G)P(T G)P(G)$
(D,T,G)	(D,T,G)	(B,T,P)	$8/16 \times 10/16 \times 2/3 \approx 160/768 = 40/192$
(D,T,G)	(D,T,G)	(B,T,P)	
(D,T,G)	(D,T,G)	(D,T,P)	
(B,T,G)	(B,T,G)	(D,T,P)	$P(B,T P)P(P) \approx 4/8 \times 5/8 \times 1/3 = 20/192$
(B,S,G)	(B,S,G)	(D,S,P)	
(B,S,G)	(B,S,G)	(B,S,P)	
(D,S,G)	(D,S,G)	(D,S,P)	Conclusion: Country = G

Same Example with Different Variables

- Three variables:
 - Hair = {blond, dark}
 - Height = {tall, short}
 - Country = {Gromland, Polvia}
- Training data: Values of (Eye Height Country)

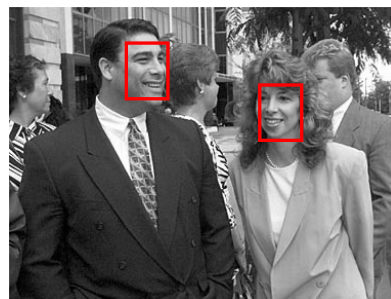
The variables are not independent so it is only an approximation.

(B,T,G)	(B,T,G)	(B,T,P)	$P(B,T,G) \approx P(B G)P(T G)P(G)$
(D,T,G)	(D,T,G)	(B,T,P)	$8/16 \times 10/16 \times 2/3 \approx 16/48 = 40/192$
(D,T,G)	(D,T,G)	(B,T,P)	
(D,T,G)	(D,T,G)	(D,T,P)	
(B,T,G)	(B,T,G)	(D,T,P)	$P(B,T P)P(P) \approx 4/8 \times 5/8 \times 1/3 = 20/192$
(B,S,G)	(B,S,G)	(D,S,P)	
(B,S,G)	(B,S,G)	(B,S,P)	
(D,S,G)	(D,S,G)	(D,S,P)	Conclusion: Country = G

Bayes at Work: Face Detection



Input Image



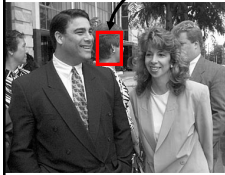
Find the faces (quickly)

Approach:

- Model the likelihood of an image window assuming face/non-face
- Use independence assumption along the way to make computations tractable

Source: Adapted from work by Henry Schneiderman (CMU & Pittsburgh Pattern Recognition)
<http://vasc.ri.cmu.edu/cgi-bin/demos/findface.cgi>

Naïve Bayes at Work



Move a window over an input image

At every position of the window:

1. Compute the values x_1, \dots, x_n of a bunch of features X_1, \dots, X_n from the image content within the window
2. Retrieve the probabilities:

$$P(X_i = x_i | \text{Face}), P(X_i = x_i | \neg \text{Face}) \quad i = 1, \dots, n$$

from tables learned off-line

3. Assuming independence, compute:

$$(1) P(\text{Face}) P(X_1 = x_1 | \text{Face}) \dots P(X_n = x_n | \text{Face})$$

$$(2) P(\neg \text{Face}) P(X_1 = x_1 | \neg \text{Face}) \dots P(X_n = x_n | \neg \text{Face})$$

4. Classify the window as a face if (1) > (2)

Learning

- Collect the values of the features for training data in tables that approximate the probabilities



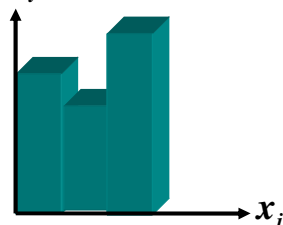
Face examples:
50-2,000 original images

$$P(X_i = x_i | \text{Face})$$

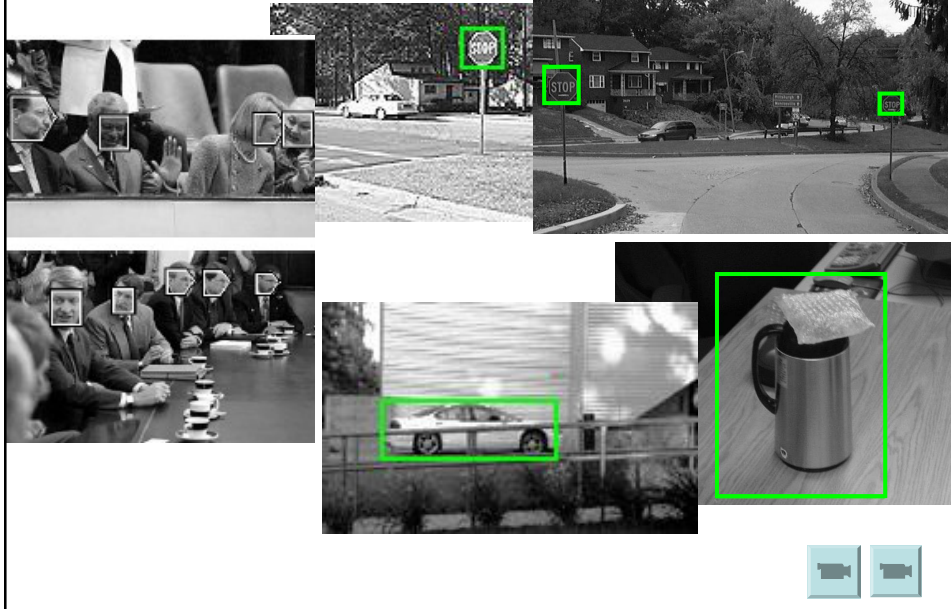


Non-face examples:
~10,000,000 examples

$$P(X_i = x_i | \neg \text{Face})$$



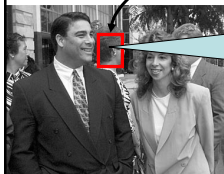
- Yes it works. And in real-time. And with other objects than faces also....



Naïve Bayes

Move a

At every



1. Comp
of featur

within the window

2. Retrieve the probabilities:

$$P(X_i = x_i | \text{Face}), P(X_i = x_i | \neg \text{Face}) \quad i = 1, \dots, n$$

from tables learned off-line

If we did not make the independence assumption, we would have to learn K^n joint probabilities, where K is the number of possible values for each feature!

The independence assumption is clearly violated: All the features X_i are computed from the same data (the image data inside the window). Nonetheless, the approximation is good enough for the classifier to work.

1. Compute the probabilities $P(X_i = x_i)$ from the image content

within the window

2. Retrieve the probabilities:

$$P(X_i = x_i | \text{Face}), P(X_i = x_i | \neg \text{Face}) \quad i = 1, \dots, n$$

from tables learned off-line

3. Using independence, compute:

$$P(X_1 = x_1 | \text{Face}) \dots P(X_n = x_n | \text{Face})$$

$$P(X_1 = x_1 | \neg \text{Face}) \dots P(X_n = x_n | \neg \text{Face})$$

4. Classify the window as a face if (1) > (2)

Naïve Bayes at Work



Move a window over an input image

At every position of the window:

1. Compute the values x_1, \dots, x_n of a bunch of features X_1, \dots, X_n from the image content within the window

2. Retrieve the probabilities:

$$P(X_i = x_i | \text{Face}), P(X_i = x_i | \neg \text{Face}) \quad i = 1, \dots, n$$

from tables learned off-line

3. Assuming independence, compute:

$$(1) P(\text{Face}) P(X_1 = x_1 | \text{Face}) \dots P(X_n = x_n | \text{Face})$$

$$(2) P(\neg \text{Face}) P(X_1 = x_1 | \neg \text{Face}) \dots P(X_n = x_n | \neg \text{Face})$$

4. Classify the window as a face if (1) > (2)

Summary

- Basic probability concepts
- Bayes rule
- What are joint distributions
- Inference using joint distributions
- Learning joint distributions from data
- Independence
- Bayes classifiers
- Naïve Bayes approach