

The Flow of ML



Probability of observing a dataset

Assume you are flipping a biased coin where p(H) = 0.4. What is the probability that you see this dataset $D = \langle H, H, T, T, H, H \rangle$

- p(H) = 0.4
- p(T) = 1 p(H) = 1 0.4 = 0.6
- If all the trails are independent then $p(D|\theta)$

 $= p(H) \times p(H) \times p(T) \times p(T) \times p(H) \times p(H)$

 $= 0.4^4 \times 0.6^2 = 0.009216$

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 3/16

Note: Order of elements in the data set do not matter in the trial. So $p(\langle H, H, H, H, T, T \rangle)$ is same (in fact any other permutation)

What is θ

It is the parameter. For our case it represents p(H) = 0.4

Hypothesis

Χ	Y	h_1	h_2	
10	0	0	1	
11	0	0	0	
12	0	0	1	
13	1	1	0	
14	0	1	1	
15	1	1	0	
16	0	1	1	
17	1	1	0	
18	1	1	1	

- In this example $h_1, h_2, ...$ are hypothesis.
- Hypothesis is a function that aims to provide value of the Y

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 4/16

- Can you identify h₁ and h₂
- Represent *H* as candidate set of hypothesis, $i.e.h_i \in H$
- Size of *H* is at least 2^{*m*}

Bayesian Learning

It is based on assumption that quantities of interest are governed by probability distribution

Notation

- P(h): initial probability that hypothesis h holds
 P(D): probability that data D will be observed
 P(D|h): probability of observing data D given some world in which hypothesis h holds
- $\overrightarrow{P(h|D)}$: probability of holding hypothesis h when data D is observed

$$P(h|D) = rac{P(D|h)P(h)}{P(D)}$$

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 5/16

Maximum a posteriori (MAP)

• Choose a hypothesis that maximizes P(h|D)

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

=
$$\operatorname{argmax}_{h \in H} \frac{P(D|h)P(h)}{P(D)}$$

=
$$\operatorname{argmax}_{h \in H} P(D|h)P(h)$$
(1)

Because P(D) is independent of h

tificial Intelligence (ZC444)

• If all the hypothesis are equally probable, we may further simplify called maximum likelihood (ML)

$$h_{ML} = \arg\max_{h \in U} P(D|h)$$
(2)

Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 6/16

For our current example



- Let bias for h₁ and h₂ be 2/50 and 6/50
- probability 7/9 and 3/9 respectively
- Posterior is (7/9)*(2/50) and
- Normalized probabilities are 0.4375 and 0.5625 respectively
- So MAP hypothesis corresponds to?

• Which is ML hypothesis? it is h₁

• Brute-force MAP learning algorithm: Evaluates posterior probability for all and returns the one with maximum

 Consistent Learner: learning algorithm is consistent learner if it provides a hypothesis that commits zero error

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 7/16

Bayes Optimal Classifier

Switching the question, from "which is most probable hypothesis?" to at is the most probable classification of the new instance? Is it possible to do better then MAP?

Example: Let posterior probabilities of three hypotheses h_1, h_2, h_3 given the training data are 0.4, 0.3, and 0.3 (obviously h_1 is MAP)

- Let classification of a new instance x is positive by h_1 and negative by h_2 and h_3
- By taking all hypotheses into account, the probability that x is positive is 0.4, and negative is 0.6
 - Most probable classification is negative and it differs from MAP

Bayes optimal classification:

 $\operatorname*{argmax}_{v_j \in V} \sum_{h_i \in H} P(v_j | h_i) P(h_i | D)$

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 8/16

where classification v_i is from V and $P(v_i|D)$ is the correct classification

Bayes Optimal Classifier $\operatorname{argmax}_{v_j \in V} \overline{\sum_{h_i \in H} P(v_j|h_i) P(h_i|D)}$ $V = \{\oplus, \Theta\}$ $P(h_1|D)=0.4$ $P(h_2|D)=0.3$ $P(\ominus | h_1) = 0$ $P(\oplus|h_1)=1$ $P(\ominus | h_2) = 1$ $P(\oplus | h_2) = 0$ $P(h_3|D) = 0.3$ $P(\ominus | h_3) = 1$ $P(\oplus | h_3) = 0$ Therefore, $\sum_{h_i \in H} P(\ominus|h_i) P(h_i|D) = 0.6$ $\sum_{h_i \in H} P(\oplus|h_i) P(h_i|D) = 0.4$ and $\operatorname{argmax}_{v_j \in \{\oplus,\ominus\}} \sum_{h_i \in H} P(v_j|h_i) P(h_i|D) = \ominus$ This type of classifier is called a Bayes optimal classifier, or Bayes optimal learner.

Naive Bayes Classifier

ntelligence (ZC444)

Bayes classifier is a highly practical Bayesian learning method

- In some domains, its performance found to be comparable to neural network and decision tree
- The Bayesian approach to classify a new instance is to assign the most probable target value describing the instance $v_{MAP} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, ..., a_n)$
- · We can use Bayes theorem to rewrite this expression as

$$v_{MAP} = \arg \max_{v_j \in V} \frac{P(a_1, a_2, ..., a_n | v_j) P(v_j)}{P(a_1, a_2, ..., a_n)}$$

=
$$\arg \max_{v_i \in V} P(a_1, a_2, ..., a_n | v_j) P(v_j)$$
(3)

Lecture-14 (Nov 10, 2023) 10/16

Naive Bayes has assumption is that the attribute values are conditionally independent given the target value

Sun (10:30-12:00PM) online@BITS-Pilani

Naive Bayes Classifier

If attribute values are conditionally independent given the target value

Artificial Intelligence (ZC444) Sun (10:30-12:00PM) online@BITS-Pilani Lecture-14 (Nov 10, 2023) 9/16

- Under this assumption.
- Given a target value, the probability of observing the conjunction $< a_1, a_2, ..., a_n >$ is just the product of the probabilities.

 $P(a_1, a_2, ..., a_n | v_i) = \prod_i P(a_i | v_i)$

Naive Bayes classifier

ence (ZC444)

is the one which

 $\operatorname*{argmax}_{v_j \in V} P(v_j) \Pi_i P(a_i | v_j)$

Sun (10:30-12:00PM) online@BITS-Pilani

Lecture-14 (Nov 10, 2023) 11/16

Example: Naive Bayes Classification

Given the data

Day	Outlook	Temperature	numiaity	wina	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rainy	Mild	High	Weak	Yes
D5	Rainy	Cool	Normal	Weak	Yes
D6	Rainy	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rainy	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rainy	Mild	High	Strong	No

Determine classification for < Rainy, Hot, High, Strong >

Example: Naive Bayes Classification



Example: Naive Bayes Classification

P(Yes) = 9/14				P(No) = 5/14				
Outlook				Humidi	ty			
		Yes	No			1		
	Sunny	2/9	3/5	1		res	INO	
-	Overcast	4/9	0/5	1	High	3/9	4/5	
-	Rain	3/9	2/5	1	Low	6/9	1/5	
Wind	Vind Temperature							
	1	Vac	No			Yes	No	
	Strong	3/0	3/5		Hot	2/9	2/5	
	Wook	6/0	2/5		Mild	4/9	2/5	
	Weak	0/9	2/5		Cool	3/9	1/5	
					• • • •	• # • • • •	e ka a	৩৫৫
Artificial Intelligence (ZC444)			Sun (10	30-12:00PM) online@BITS-	Pilani	Lecture-14	(Nov 10, 2023)	14/16

Thank You!

Artificial Intelligence (ZC444)

Thank you very much for your attention!

Queries ?

Example: Naive Bayes Classification



・ロト・イラト・イミト・ミート ラーペー・ Sun (10:30-12:00PM) online@BITS-Pliani Lecture-14 (Nov 10, 2023) 16/16