

BAYESIAN LEARNING

CAN BE SUMMARIZED BY THE FOLLOWING:

$$P(X|D) = \sum_{h \in H} P(X|D, h) P(h|D) = \sum_{h \in H} P(X|h) \cdot P(h|D)$$

PREDICTIONS ON
NEW DATA X
(GIVEN WHAT WE LEARNED
FROM THE DATASET)
D

LIKELIHOOD: PROB OF
OBSERVING THE DATASET
UNDER THE
HYPOTHESIS
h

PROBABILITY OF
OBSERVING X UNDER
THE HYPOTHESIS h

PROBABILITY OF
HYPOTHESIS h
ACCORDING TO OUR
DATA

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)}$$

POSTERIOR: PROBABILITY OF
HYPOTHESIS h AFTER
OBSERVING THE DATA

PRIOR: PROBABILITY OF
HYPOTHESIS h BEFORE OBSERVING
THE DATA SET

IT IS CALLED BAYESIAN
BECAUSE THIS TERM
OBTAINED APPLYING
BAYES THEOREM

WITH RESPECT TO THE MAXIMUM LIKELIHOOD HERE WE ARE USING
ALL THE HYPOTHESES IN THE HYPOTHESES SPACE H.

WE CAN APPROXIMATE BAYESIAN INFERENCE USING A SINGLE HYPOTHESIS TO
MAKE PREDICTIONS

MAXIMUM A POSTERIORI

APPROXIMATE $P(X|D)$ WITH $P(X|h_{MAP})$

WHERE

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

MAXIMUM LIKELIHOOD

IF WE FURTHER ASSUME NO PRIOR AMONG THE HYPOTHESES,
THEN WE GET THE FOLLOWING APPROXIMATION

$$P(X|D) \approx P(X|h_{ML}) \quad \text{WHERE } h_{ML} = \operatorname{argmax}_{h \in H} P(D|h)$$

EXAMPLE OF BAYESIAN INFERENCE

Assume I know an industry produces 5 types of candy box (HYPOTHESES) which composition is known and they specify also the number of boxes of each type they produce (prior probability distribution over my hypothesis space).

Types of Candy Boxes (HYPOTHESES)	Prior Probability over the Hypothesis Space
h_1 100% cherry	$P(h_1) = 0.1$
h_2 75% cherry 25% lime	$P(h_2) = 0.2$
h_3 50% cherry 50% lime	$P(h_3) = 0.4$
h_4 25% cherry 75% lime	$P(h_4) = 0.2$
h_5 100% lime	$P(h_5) = 0.1$

Now I pick 3 candies from a box (each time I reinset the candy such that the probability does not change) and I want to use Bayesian inference to predict the next candy given the observed data. Assume I pick 3 lime candies $D = \{l, l, l\}$. First I need to compute the posterior distribution of my hypotheses and then I can use the posteriors to make new predictions.

$$P(h_1 | D = \{l, l, l\}) = d \cdot P(D = \{l, l, l\} | h_1) \cdot P(h_1) = d \cdot 0^3 \cdot 0.1 = 0$$

$$P(h_2 | D = \{l, l, l\}) = d \cdot P(D = \{l, l, l\} | h_2) \cdot P(h_2) = d \cdot 0.25^3 \cdot 0.2 = d \cdot 0.003125 = 0.01$$

$$P(h_3 | D = \{l, l, l\}) = d \cdot P(D = \{l, l, l\} | h_3) \cdot P(h_3) = d \cdot 0.5^3 \cdot 0.4 = d \cdot 0.05 = 0.21$$

$$P(h_4 | D = \{l, l, l\}) = d \cdot P(D = \{l, l, l\} | h_4) \cdot P(h_4) = d \cdot 0.75^3 \cdot 0.2 = d \cdot 0.084375 = 0.36$$

$$P(h_5 | D = \{l, l, l\}) = d \cdot P(D = \{l, l, l\} | h_5) \cdot P(h_5) = d \cdot 1^3 \cdot 0.1 = d \cdot 0.1 = 0.42$$

$$d \left(0.003125 + 0.01 + 0.084375 + 0.1 \right) = 1$$

$$d = \frac{1}{0.2375} = 4.210526$$



I do not divide the denominator and normalize such that the posteriors sum to one.

↓ BAYESIAN INFERENCE

$$\begin{aligned} P(l | D = \{l, l, l\}) &= \sum_i P(l | h_i) \cdot P(h_i | D = \{l, l, l\}) \\ &= 0 \times 0 + 0.25 \times 0.01 + 0.5 \times 0.21 + 0.75 \times 0.36 + 1 \times 0.42 = 0.80 \end{aligned}$$

APPLICATIONS OF BAYESIAN INFERENCE

BAYES OPTIMAL CLASSIFIER

ASSUME I WANT TO CLASSIFY A NEW OBSERVATION X WITH A CLASS

$c_s \in C$ THEN

$$c_{\text{opt}} = \underset{c_s \in C}{\operatorname{argmax}} P(c_s | x) = \underset{c_s \in C}{\operatorname{argmax}} \sum_{h \in H} P(c_s | h) \cdot P(h | x)$$

WITH THE SAME HYPOTHESIS SPACE AND PRIOR KNOWLEDGE NO OTHER CLASSIFICATION METHOD OUTPERFORMS BAYES OPTIMAL CLASSIFIER

HOWEVER IT IS OFTEN IMPRACTICAL BECAUSE

→ I'VE NO ACCESS TO THE ENTIRE HYPOTHESIS SPACE OR IT IS INFINITE, SO THAT THE SUMMATION BECOMES UNRACTABLE

→ DIFFICULT TO MODEL $P(h | d)$

A SIMPLE APPROXIMATION THAT IS OFTEN USED IN PRACTICE IS

NAIVE BAYES CLASSIFIER

ASSUME A POINT X IS REPRESENTED BY A SET OF ATTRIBUTES (OR FEATURES) THAT ARE MUTUALLY INDEPENDENT GIVEN THE CATEGORY $C_S \in C$

$$x = \langle a_1, \dots, a_d \rangle, \text{ AND } C = \{c_1, \dots, c_c\}$$

$$c_{\text{opt}} = \underset{c_s \in C}{\operatorname{argmax}} P(c_s | a_1, \dots, a_d) = \underset{c_s \in C}{\operatorname{argmax}} P(a_1, \dots, a_d | c_s) \cdot P(c_s) = \underset{c_s \in C}{\operatorname{argmax}} P(c_s) \cdot \prod_{i=1}^d P(a_i | c_s)$$

WHERE

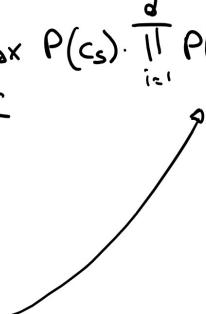
$$P(a_i = l | c_s = h) = \frac{\left(\begin{array}{l} \text{NUMBER OF ELEMENTS WITH} \\ a_i = l \text{ THAT BELONG TO} \\ \text{CLASS } h \end{array} \right) + \lambda}{\left(\begin{array}{l} \text{NUMBER OF ELEMENTS WITH} \\ a_i = l \end{array} \right) + k \lambda}$$

+ λ

+ $k \lambda$

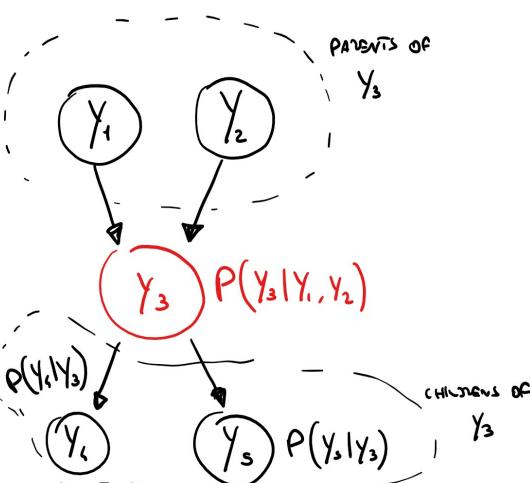
THIS IS A
SMOOTHING TERM,
BENEFIT IF AVG TERM HIGH
IS ZERO THE WHOLE PROBABILITY
BECOMES ZERO

FOR THIS
INSTEAD \rightarrow NO
ASSUMPTION



BAYESIAN NETWORKS (OR BAYES NETWORKS, BAYES NETS, BELIEF NETWORKS, DECISION NETWORK)

DESCRIBES CONDITIONAL INDEPENDENCE BETWEEN SUBSETS OF RANDOM VARIABLES USING A DIRECT ACYCLIC GRAPH



ASSUMPTION: EVERY VARIABLE IS CONDITIONALLY INDEPENDENT WITH RESPECT TO ITS NON DESCENDANT, GIVEN ITS PARENTS

E.G.

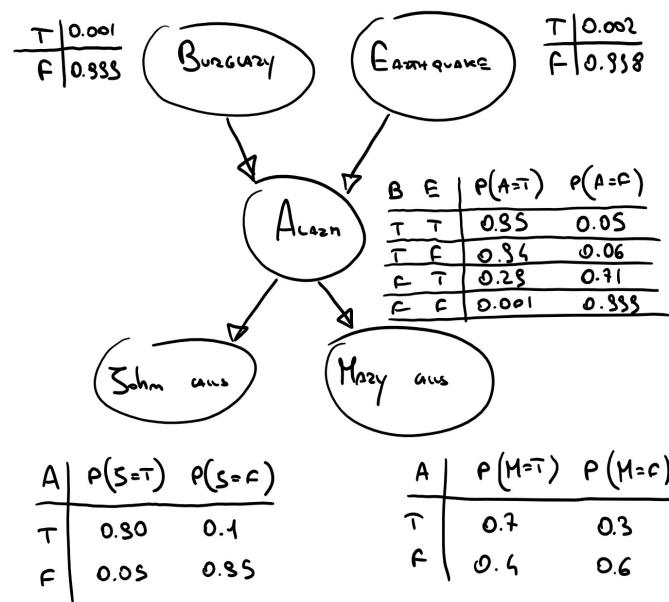
Y_5 IS INDEPENDENT W.R.T. Y_4 GIVEN Y_3

$$P(Y_5, Y_4 | Y_3) = P(Y_5 | Y_3) \cdot P(Y_4 | Y_3)$$

INVERSE: COMPUTE THE POSTERIOR PROBABILITY OF SOME VARIABLES WHEN OTHER VARIABLES ARE OBSERVED

$$P(X_1 = x_1, X_2 = x_2, \dots) = \prod_{i=1}^m P(X_i = x_i | \text{Assignments to Parents}(X_i))$$

EXAMPLE THE CONDITIONAL PROBABILITIES AT EACH NODE ARE SPECIFIED THROUGH CONDITIONAL PROBABILITY TABLES (CPT)



INFERENCE WITH TWO OBSERVED VARIABLES USING VARIABLE ELIMINATION

$$\begin{aligned}
 P(\Sigma) &= \sum_A \sum_B \sum_E P(B) \cdot P(E) \cdot P(A|B,E) \cdot P(\Sigma|A) \\
 &= \sum_A \sum_B P(B) \cdot P(\Sigma|A) \left(\sum_E P(A|B,E) \cdot P(E) \right) \\
 &= \sum_A \sum_B P(B) \cdot P(\Sigma|A) \lambda_1(A,B) \\
 &= \sum_A P(\Sigma|A) \left(\sum_B \lambda_1(A,B) \cdot P(B) \right)
 \end{aligned}$$

$$\begin{aligned}
 \lambda_1(A,B) &= P(A|B, E=T) \cdot P(E=T) + P(A|B, E=F) \cdot P(E=F) \\
 \lambda_1(A,B) &= \begin{array}{c|c|c} A & B & \lambda_1(A,B) \\ \hline T & T & 0,910 \\ T & F & 0,002 \\ F & T & 0,060 \\ F & F & 0,998 \end{array}
 \end{aligned}$$

$$\begin{array}{c|c} A & \lambda_2(A) \\ \hline T & 0,003 \\ F & 0,997 \end{array}$$

$$= \sum_A P(\Sigma|A) \cdot \lambda_2(A) = P(\Sigma|A=T) \cdot \lambda_2(A=T) + P(\Sigma|A=F) \cdot \lambda_2(A=F)$$

Σ	$P(\Sigma)$
T	0,053
F	0,947

INFERENCE WITH EVIDENCES USING VARIABLE ELIMINATION

$$P(M | B=T) = \sum_A \sum_E P(B=T) \cdot P(E) \cdot P(A | B=T, E) \cdot P(M|A)$$

$$P(B=T)$$

$$\lambda_1(A) = P(A | B=T, E=T) \cdot P(E=T) +$$

$$P(A | B=T, E=F) \cdot P(E=F)$$

$$= \frac{1}{0.001} \sum_A 0.001 \cdot P(M|A) \left(\sum_E P(A | B=T, E) \cdot P(E) \right)$$

A	$\lambda_1(A)$
T	0,940
F	0,060

$$= \frac{1}{0.001} \sum_A 0.001 \cdot P(M|A) \cdot \lambda_1(A) =$$

$$= \frac{1}{0.001} \left(\partial_{\lambda_1} \left(P(M|A=T) \lambda_1(A=T) + P(M|A=F) \lambda_1(A=F) \right) \right)$$

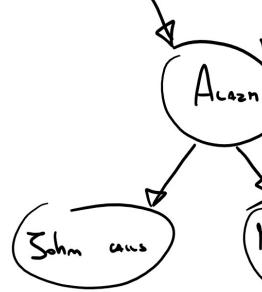
M	$P(M B=T)$
T	0,682
F	0,318

$$\begin{array}{c|cc} T & 0.001 \\ \hline F & 0.999 \end{array}$$

Burglary

Earthquake

$$\begin{array}{c|cc} T & 0.002 \\ \hline F & 0.998 \end{array}$$



John ans

A	$P(S=T)$	$P(S=F)$
T	0.90	0.1
F	0.05	0.95

A	$P(M=T)$	$P(M=F)$
T	0.7	0.3
F	0.4	0.6