

CS 480 – INTRODUCTION TO ARTIFICIAL INTELLIGENCE

TOPIC: BAYESIAN NETWORKS PARAMETER ESTIMATION



Mustafa Bilgic



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

BAYESIAN NETWORK PARAMETER ESTIMATION

○ Given:

- A set of random variables, V_i
 - E.g., age, gender, cholesterol level, etc.
- A Bayesian network structure over these variables
 - E.g., a doctor can point out the most important correlations and causations
- Data
 - E.g., existing patient records, where ~~some~~ or all V_i are known

CS 583

CS 581

CS 563

○ Goal: Task

- Estimate the parameters needed for the Bayesian network, i.e., $P(V_i \mid \text{parentsOf}(V_i))$

$$P(V_1 \dots V_n) = \prod$$

KNOWN BAYESIAN NETWORK STRUCTURE

- In this class, we assume the structure is given
- How reasonable is this assumption?
 - In some domains, the expert might provide a reasonable structure to start with
- There are many methods that learn the structure of the Bayesian network from data
 - Those topics are covered in the CS583 – Probabilistic Graphical Models course in detail

PARAMETER ESTIMATION FOR BNs

- Assume the network structure is given over variables V_i
- Let d_j be a fully observed instance
 - $d_j = \langle V_1=t, V_2=f, \dots, V_n=t \rangle$
- The data \mathcal{D} consists of fully observed instances
 - $\mathcal{D} = \{d_1, d_2, \dots, d_m\}$
- Estimate the network parameters $P(V_i \mid \text{parents}(V_i))$
- Two approaches
 1. Maximum likelihood estimation
 2. Bayesian estimation

SIMPLEST CASE – ONE VARIABLE

- Imagine we have a thumbtack
- Flip it, and it comes as heads or tails

heads



tails



- $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$
- Assume we flip it 100 times and it comes head 30 times
- What is θ ?

THUMBTRACK TOSSES

- Assume we have a set of thumbtack tosses
 - $\mathcal{D} = \{d_1, d_2, \dots, d_{100}\}$
- Assume we have 30 heads and 70 tails
- $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$
- θ can be any number between 0 and 1
- We have an infinite number of choices
 - $\theta=0, \dots, \theta=0.3, \dots, \theta=0.5, \dots, \theta=1$
- We want to formulate an objective function $f(\theta: \mathcal{D})$, where, given 30 heads and 70 tails, $f(\theta: \mathcal{D})$ achieves its maximum when $\theta=0.3$
 - Any ideas?

Handwritten notes and diagrams:

- A list of tosses: $d_1, d_2, d_3, \dots, d_{100}$ with corresponding outcomes: H, T, T, ..., H.
- A bracket on the right side of the list indicates a total of 30 H and 70 T.
- An arrow points from the text "world is most likely to occur if $\theta = 0.3$ " to the list of tosses.
- At the bottom, two sequences are shown: $\theta = 0 \Rightarrow T, T, \dots, T$ and $\theta = 1 \Rightarrow H, \dots, H$.

LIKELIHOOD

$$P(H) = \theta$$

$$P(T) = 1 - \theta$$

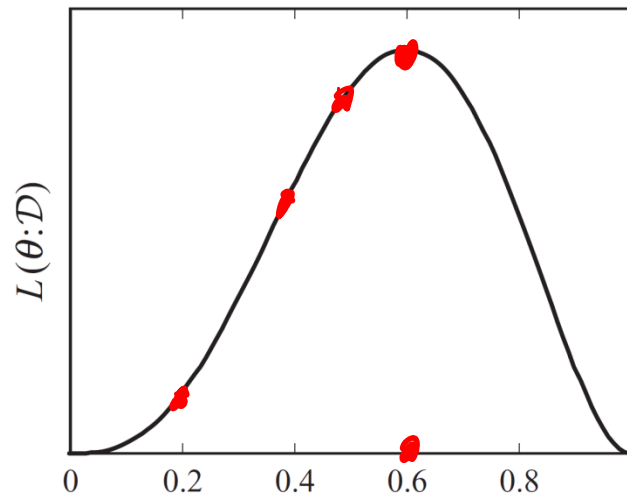
$$\rightarrow P(H, T, T, H, H) = P(H) P(T) P(T) P(H) P(H)$$

- What is the probability, or *likelihood*, of seeing the sequence H, T, T, H, H?

- $\theta * (1 - \theta) * (1 - \theta) * \theta * \theta = \theta^3 (1 - \theta)^2$

$$= f(\theta)$$

$$= \theta^3 (1 - \theta)^2$$



$$\frac{3}{5}$$

When is $L(\theta; D)$ maximum?

LIKELIHOOD/LOG-LIKELIHOOD

- Number of heads = k , number of tails = $m-k$
- Likelihood: $L(\theta: \mathcal{D}) = \theta^k(1-\theta)^{m-k}$
- Log-likelihood: $l(\theta: \mathcal{D}) = k\log\theta + (m-k)\log(1-\theta)$
- Note that $L(\theta: \mathcal{D})$ achieves its maximum for θ that maximizes $l(\theta: \mathcal{D})$
- Find θ that maximizes the log-likelihood
- Take derivate of $l(\theta: \mathcal{D})$ w.r.t. θ and set it to zero

positive find
 f
 $\arg\max_x f(x)$
 $\arg\max_x \log f(x)$

LET'S SEE A FEW EXAMPLES

- Simple structure
 - $X \rightarrow Y$
- General structure
 - The key is that the parameters for each variable can be optimized independently
 - Examples

BAYESIAN ESTIMATION

- Assume we flip a coin 10 times and we get 4 Heads, 6 Tails

- What is $P(C=H)$?

0.5

- What if we repeat the flips 10M times and we get 4M Heads and 6M Tails?

- Bayesian estimation will let us encode our *prior knowledge*

small data \Rightarrow trust your experience more
large data \Rightarrow trust the data more

prior belief distribution



experience + counts

experience + counts

0.4



TO CUT IT SHORT, (I MEAN REALLY SHORT)

- We'll encode our prior knowledge as a set of “imaginary” counts
- For example, we will assume that we have already seen α heads and β tails
- Assume we flip a coin 10 times and we get 4 Heads, 6 Tails
 - $P(C=\text{heads}) = (4 + \alpha) / (10 + \alpha + \beta)$
 - $\alpha = 0, \beta = 0; P(C=h) = 4/10 = 0.4$
 - $\alpha = 1, \beta = 1; P(C=h) = 5/12 = 0.417$
 - $\alpha = 10, \beta = 10; P(C=h) = 14/30 = 0.467$
 - $\alpha = 100, \beta = 100; P(C=h) = 104/210 = 0.495$
- Assume we flip a coin 1000 times and we get 400 Heads, 600 Tails
 - $P(C=\text{heads}) = (400 + \alpha) / (1000 + \alpha + \beta)$
 - $\alpha = 0, \beta = 0; P(C=h) = 400/1000 = 0.4$
 - $\alpha = 1, \beta = 1; P(C=h) = 401/1002 = 0.4002$
 - $\alpha = 10, \beta = 10; P(C=h) = 410/1020 = 0.402$
 - $\alpha = 100, \beta = 100; P(C=h) = 500/1200 = 0.417$

IMAGINARY COUNTS

- Note that imaginary counts can be applied to any categorical variable, not necessarily just binary variables
- Also helps with dealing zero probabilities
- When all imaginary counts are 1, this is called Laplace smoothing
 - E.g, $\alpha = 1$, $\beta = 1$
- Let's see some examples