

RECITATION 8

GRAPHICAL MODELS

10-601: INTRODUCTION TO MACHINE LEARNING

4/23/2021

1 Warm-Up: Probability Review

| X_1 | X_2 | Probability |
|-------|-------|-------------|
| 0 | 0 | 0.45 |
| 1 | 0 | 0.15 |
| 0 | 1 | 0.2 |
| 1 | 1 | 0.2 |

Table 1: Joint Probability Table

1. What is the joint probability $P(X_1 = 1, X_2 = 0)$?
2. What is the marginal probability $P(X_2 = 0)$?
3. What is the conditional probability $P(X_1 = 1|X_2 = 0)$?

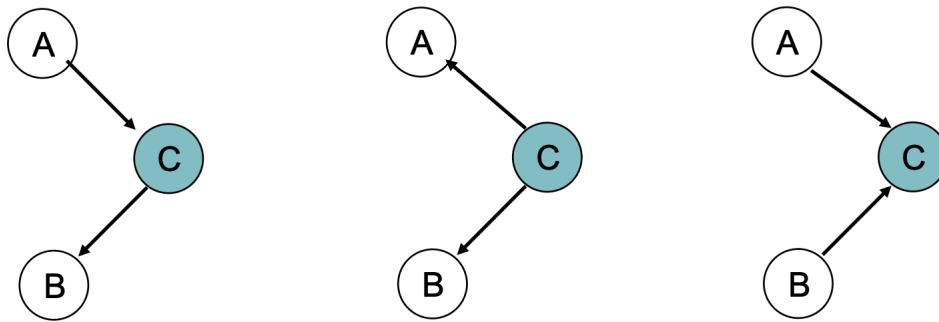


Figure 1: A and B are d-separated by C in the first two graphs, but not in the third.

2 Bayesian Networks

2.1 D-separation

Suppose we have three sets of random variables (X , Y , and Z). X and Y are **d-separated** by Z (and therefore conditionally independent given Z) if and only if every path from every variable in X to every variable in Y is **blocked**. A path is blocked if either:

- Arrows meet head-to-tail or tail-to-tail at a node in Z
- Arrows on the path meet head-to-head at a node, and neither that node, nor any of its descendants, is in Z

2.2 Practice problems

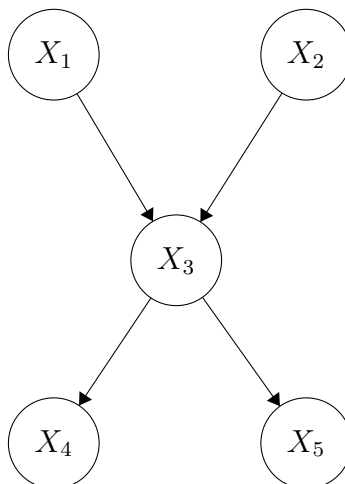


Figure 2: Graphical Model

1. Write down the factorization of the above directed graphical model.

-
2. Given X_3 , what are the relationships (cond. independent or not) between the random variables listed below
 - $(X_1 \text{---} X_4) | X_3$
 - $(X_1 \text{---} X_2) | X_3$
 - $(X_4 \text{---} X_5) | X_3$
 3. Given the graph structure and assuming all variables are boolean valued, how many parameters are required to learn the graphical model?
 4. Without the Bayesian network, how many parameters are required to learn the joint probability model of these five random variables?

3 Expectation Maximization

Consider the following problem set-up. We have two coins, coin A and coin B. Each coin has a probability of coming up heads that we would like to estimate, i.e. θ_A, θ_B . Suppose we now have the following procedure:

- Randomly choose one of the coins.
- Using the chosen coin, perform 6 independent coin flips.
- Repeat the process 3 times.

Now suppose we observe the following data for each of the 3 trials:

| | | | | | | |
|---------|---|---|---|---|---|---|
| Trial 1 | H | H | H | H | H | T |
| Trial 2 | T | T | T | T | T | H |
| Trial 3 | H | T | T | T | T | H |

Table 2: Coin Flip Data

If we knew what coin was selected for each trial, we could easily use MLE estimation to estimate $\hat{\theta}_A, \hat{\theta}_B$. Instead, we will assume that we don't know which coin was selected, and will have to learn the latent factor describing which coin was selected to maximize the likelihood.

We would like to use the EM algorithm to estimate $\hat{\theta}_A, \hat{\theta}_B$. Rather than picking the single most likely completion of the missing coin assignments on each iteration, the expectation maximization algorithm computes probabilities for each possible completion of the missing data, using the current parameters. These probabilities are used to create a weighted training set consisting of all possible completions of the data.

To estimate the unknown values in the E-step at iteration t , for each trial i we will compute weights α_i, β_i representing how likely the data is to come from coin A versus coin B.

As a results,

$$\alpha_i = \frac{P(D_i | \hat{\theta}_A^{(t)})}{P(D_i | \hat{\theta}_A^{(t)}) + P(D_i | \hat{\theta}_B^{(t)})}$$

and,

$$\beta_i = \frac{P(D_i | \hat{\theta}_B^{(t)})}{P(D_i | \hat{\theta}_A^{(t)}) + P(D_i | \hat{\theta}_B^{(t)})} = 1 - \alpha_i$$

where D_i is the data for trial i .

Once we have estimated α_i, β_i for each trial, we will then perform the M-step of maximizing the parameters $\hat{\theta}_A, \hat{\theta}_B$.

To do so, we will compute

$$\hat{\theta}_A^{(t+1)} = \sum_i \frac{\alpha_i * \text{Number of heads in trial } i}{\alpha_i * \text{Number of heads in trial } i + \alpha_i * \text{Number of tails in trial } i}$$

$$\hat{\theta}_B^{(t+1)} = \sum_i \frac{\beta_i * \text{Number of heads in trial } i}{\beta_i * \text{Number of heads in trial } i + \beta_i * \text{Number of tails in trial } i}$$

Essentially, we are performing an MLE estimate weighted by our best guess of which coin was selected for each trial under our current parameters.

Suppose our initial guesses are $\hat{\theta}_A^{(0)} = 0.6$, $\hat{\theta}_B^{(0)} = 0.5$

1. Compute α_i, β_i to 2 decimal places for each trial i .
2. With the found values of α_i, β_i , compute $\hat{\theta}_A^{(1)}, \hat{\theta}_B^{(1)}$ in the M-step

4 Gaussian Naive Bayes

Priors

$$P(Y = 1) = \frac{\#D\{Y = 1\}}{|D|}$$

$$P(Y = 0) = \frac{\#D\{Y = 0\}}{|D|}$$

Class conditional means, variances

$$\hat{\mu}_{m,k} = \frac{1}{\sum_{i=1}^N I(Y^{(i)} = k)} \sum_{i=1}^N \mathbf{x}_m^{(i)} * I(Y^{(i)} = k)$$

$$\hat{\sigma}_{m,k}^2 = \frac{1}{\sum_{i=1}^N I(Y^{(i)} = k)} \sum_{i=1}^N (\mathbf{x}_m^{(i)} - \hat{\mu}_{m,k})^2 * I(Y^{(i)} = k)$$

1. Write down the Naive Bayes assumption.
2. When would we use Gaussian Naive Bayes?
3. Consider a simple Gaussian Naive Bayes example. Let's say we want to use predict if a person will go hiking today based on the temperature outside.

| Temperature (F) | Hiking? |
|-----------------|---------|
| 50.8 | No |
| 67.6 | Yes |
| 39.7 | No |
| 36.5 | No |
| 72.4 | Yes |

- (a) Step 1: Find the class priors.
- (b) Step 2: Find the class conditional mean, variance for X (temperature).
- (c) Calculate the following class-conditional probabilities for Temperature = 65 :
- (d) The red graph represents the class conditional mean for $Y=Yes$, and the green graph represents the class conditional mean for $Y=No$. Based on your results from the last question, interpret the graph below and explain how the output predictions $P(Y = k|X_{temp})$ are affected by different feature values.

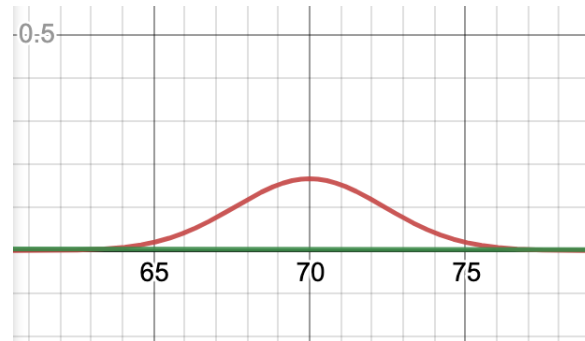


Figure 3: Class conditional probability for temperature, $Y=Yes$