

CS228 Homework 2

Instructor: Stefano Ermon – ermon@stanford.edu

Available: 01/24/2017; Due: 02/03/2017

1. [8 points] (I-Maps and P-Maps) Consider a probability distribution P over the variables A, B, C and D that satisfies only the following independencies:

1. $A \perp C | B$. ✓

2. $A \perp C | B, D$. ✓

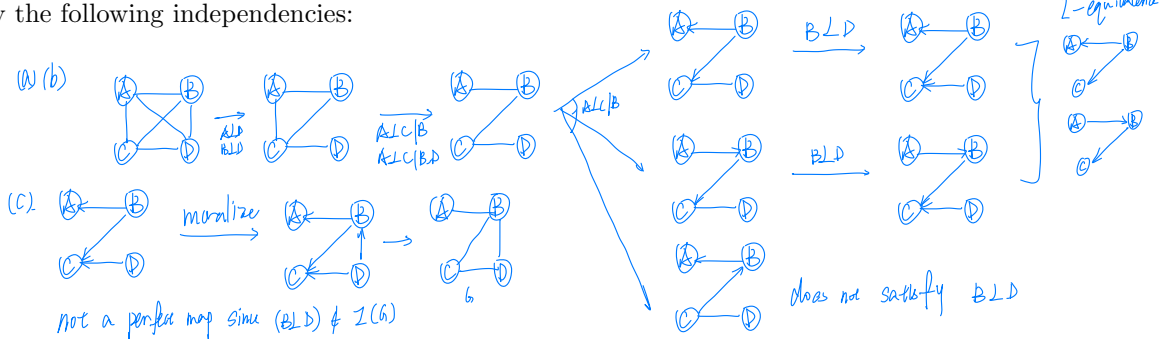
3. $A \perp D | B$.

4. $A \perp D | B, C$.

5. $B \perp D$. ✓

6. $B \perp D | A$.

7. $A \perp D$. ✓



- (a) [3 points] Draw a Bayesian network that is a perfect map for P .

- (b) [2 points] Does this perfect map have any I-equivalent graphs? If so, draw them. If not, explain why not.

- (c) [3 points] Draw a Markov network that is a minimal I-map for P and explain why the Markov network is or is not also a perfect map.

2. [9 points] (I-Maps and P-Maps) Consider a distribution P over four binary random variables (A, B, C, D) , which gives probability $1/8$ to assignments $(0, 0, 0, 0)$, and $(1, 1, 0, 0)$, and gives probability $1/4$ to assignments $(1, 1, 1, 0)$, $(0, 1, 0, 1)$, and $(1, 0, 1, 1)$. A skeleton for the Bayesian network G is also provided; the skeleton contains all the nodes and edges, but the directions of the edges are missing.

(a)

A	B	C	D	P
0	0	0	0	$1/8$
1	1	0	0	$1/8$
1	1	1	0	$1/4$
0	1	0	1	$1/4$
1	0	1	1	$1/4$

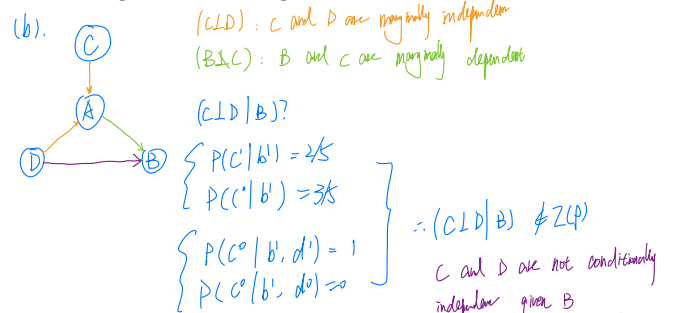
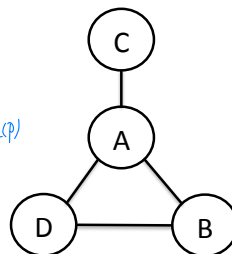
Handwritten calculations for the probabilities:

$$P(A=0) = 1/2, P(A=1) = 1/2$$

$$P(B=0) = 1/2, P(B=1) = 1/2$$

$$P(C=0) = 1/2, P(C=1) = 1/2$$

$$P(D=0) = 1/2, P(D=1) = 1/2$$



(c). Handwritten tables for the conditional probability distributions (CPDs) for each node in the Bayesian network specified in (b).

C	P
c^0	$1/2$
c^1	$1/2$

D	P
d^0	$1/2$
d^1	$1/2$

Handwritten table for the joint probability distribution $P(A, B, C, D)$.

C	D	A	P
c^0	d^0	a^0	$1/8$
c^0	d^0	a^1	$1/8$
c^0	d^1	a^0	$1/4$
c^0	d^1	a^1	$1/4$

Handwritten notes: $P(A|C,D) = 1 - P(A^c|C,D)$, $P(B|A,D)$, simple counting.

- (a) [2 points] Decide whether the following two independencies are in $I(P)$: $(C \perp D), (C \perp B)$

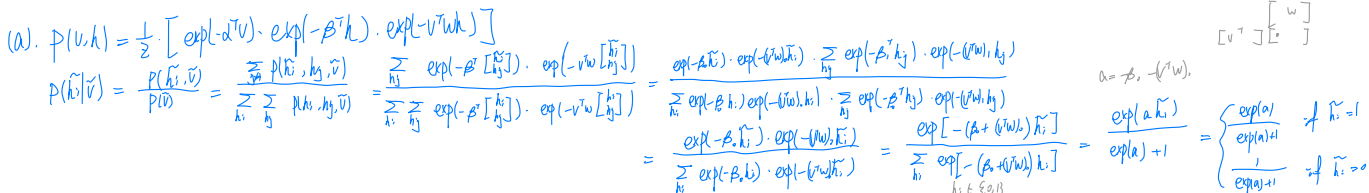
- (b) [5 points] Give a direction to each individual edge in G , such that the resulting Bayesian network is a perfect map of P . Is your solution unique? Briefly state why.

- (c) [2 points] Give the CPDs for each node in your Bayesian network specified in (b).

3. [26 points] (Undirected Graphical Models, Inference) The Restricted Boltzmann Machine

Restricted Boltzmann machines (RBMs) have been widely used as a generative model in many fields of machine learning: image processing, speech recognition and collaborative filtering. Concretely, an RBM

Handwritten equation: $P(A, B, C, D) = P(C) \cdot P(D) \cdot P(A|C, D) \cdot P(B|A, D)$



where

$$\phi(\mathbf{v}, \mathbf{h}) = -\alpha^T \mathbf{v} - \beta^T \mathbf{h} - \mathbf{v}^T W \mathbf{h}$$

(a) **[5 points]** What is the marginal conditional distribution of $P(\mathbf{h}_i \mid \mathbf{v})$ for a single hidden variable \mathbf{h}_i ? Can it be computed tractably?

(c) **[2 points]** Suppose we are given samples for the visible units $\mathcal{D} = \{\mathbf{v}^1, \dots, \mathbf{v}^K\}$ (e.g., the MNIST digits dataset), and we want to learn the parameters of the RBM to maximize the likelihood of these samples. Since we don't know the values of hidden variables \mathbf{h}^k , a natural approach is to look for parameters that maximize the marginal likelihood of each sample \mathbf{v}^k , given by

For a fixed \mathbf{v} , can

be computed efficiently (in time not exponential in n)?

be computed efficiently (in time not exponential in m)?

(e) **[2 points]** Can the normalization constant

be computed efficiently (in time not exponential in m or n)?

4. **[10 points] [Maximum Likelihood Learning of Bayes Nets]** Let G be a valid Bayes Net structure, and \mathcal{D} be some training data. It can be shown that the likelihood of \mathcal{D} for a given choice of the parameters θ (specifying the conditional probability tables in G) is

$$\ell_G(\theta; \mathcal{D}) = \sum_{i=1}^n \sum_{\mathbf{u}_i \in \text{Val}(Pa(X_i))} \sum_{x_i} M[x_i, \mathbf{u}_i] \log \theta_{x_i | \mathbf{u}_i}$$

We have seen in class that the maximum likelihood estimate $\theta^{ML} = \arg \max_{\theta} \ell_G(\theta; \mathcal{D})$ of the parameters of G can be computed as follows

$$\theta_{x_i|\mathbf{u}_i}^{ML} = \frac{M[x_i, \mathbf{u}_i]}{M[\mathbf{u}_i]} \quad (1)$$

where $M[x_i, \mathbf{u}_i]$ is the number of times the tuple $(X_i = x_i, Pa(X_i) = \mathbf{u}_i)$ appears in the data \mathcal{D} .

Let G' be a valid Bayes net structure obtained by adding an edge to G . Show that

$$\max_{\theta'} \ell_{G'}(\theta'; \mathcal{D}) \geq \max_{\theta} \ell_G(\theta; \mathcal{D})$$

Note that G' is strictly more expressive than G (G' makes less conditional independence assumptions). This is confirmed by the fact that to specify a Bayes Net with structure G' we need more parameters, i.e., the vector θ' has more components than θ . The goal of this exercise is to show directly that the more "complex" a Bayesian Network is, the better we can fit it to data.

You may assume that G and G' are identical, except for an extra $x_1 \rightarrow x_n$ edge in G' .

Hint: There are several ways to show this. Jensen's inequality might come handy.

$$\tilde{\ell}(\theta; D) = \prod_{k=1}^n \prod_{x_i \in \mathcal{X}_i} \theta_{x_i | \text{pa}(x_i)}$$

$$\begin{aligned} \ell(\theta; D) &= \log \tilde{\ell}(\theta; D) = \sum_{k=1}^n \sum_{x_i \in \mathcal{X}_i} \log \theta_{x_i | \text{pa}(x_i)} \\ &= \sum_{x_i \in \mathcal{X}_i} \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \sum_{u_i \in \mathcal{U}(x_i)} M[x_i, u_i] \log \theta_{x_i | u_i} \end{aligned}$$

$$\min_{\theta} - \sum_{x_i \in \mathcal{X}_i} \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \sum_{u_i \in \mathcal{U}(x_i)} M[x_i, u_i] \log \theta_{x_i | u_i}$$

$$\text{s.t. } \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \theta_{x_i | \text{pa}(x_i)} = 1 \quad \forall u_i$$

\Downarrow separable

$$\min_{\theta} \sum_{x_i \in \mathcal{X}_i} -M[x_i, u_i] \log(\theta_{x_i | u_i})$$

$$\text{s.t. } \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \theta_{x_i | \text{pa}(x_i)} = 1 \quad \forall u_i$$

$$\begin{aligned} \mathcal{L}(\theta, u) &= \sum_{x_i \in \mathcal{X}_i} -M[x_i, u_i] \log(\theta_{x_i | u_i}) \\ &\quad + V_{u_i} \left(\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \theta_{x_i | \text{pa}(x_i)} - 1 \right) \\ &= \left(\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} -M[x_i, u_i] \log(\theta_{x_i | u_i}) + \theta_{x_i | u_i} V_{u_i} \right) - V_{u_i} \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_{x_i | u_i}} = \frac{-M[x_i, u_i]}{\theta_{x_i | u_i}} + V_{u_i} \stackrel{!}{=} 0$$

$$\theta_{x_i | u_i} = \frac{M[x_i, u_i]}{V_{u_i}}$$

$$g(u) = \inf_{\theta} \mathcal{L}(\theta, u) = \left(\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} -M[x_i, u_i] \cdot \log \frac{M[x_i, u_i]}{V_{u_i}} + M[x_i, u_i] \right) - V_{u_i}$$

$$\forall u_i, g = \left[\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} -M[x_i, u_i] \cdot \frac{V_{u_i}}{M[x_i, u_i]} \cdot M[x_i, u_i] \cdot (-1) + \frac{1}{V_{u_i}} \right] - 1 \stackrel{!}{=} 0$$

$$\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \frac{M[x_i, u_i]}{V_{u_i}} = 1$$

$$V_{u_i}^* = \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} M[x_i, u_i]$$

$$\theta_{x_i | u_i}^* = \frac{M[x_i, u_i]}{V_{u_i}^*} = \frac{M[x_i, u_i]}{\sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} M[x_i, u_i]} = \frac{M[x_i, u_i]}{M[x_i, u_i]}$$

without loss of generality, assume that we add an edge $x_1 \rightarrow x_n$

$$\tilde{\ell}'(\theta'; D) = \prod_{k=1}^n \left(\prod_{x_i \in \mathcal{X}_i} \theta'_{x_i | \text{pa}(x_i)} \right) \cdot \theta'_{x_n | x_1}$$

$$\begin{aligned} \ell'(\theta'; D) &= \log \tilde{\ell}'(\theta'; D) = \sum_{k=1}^n \left(\log \theta_{x_n | x_1} + \sum_{x_i \in \mathcal{X}_i} \log \theta'_{x_i | \text{pa}(x_i)} \right) \\ &= \sum_{x_i \in \mathcal{X}_i} \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \sum_{u_i \in \mathcal{U}(x_i)} M[x_i, u_i] \log \theta'_{x_i | u_i} \\ &\quad + \sum_{\text{pa}(x_n) \in \mathcal{X}_{\text{pa}(x_n)}} \sum_{u_n \in \mathcal{U}(x_n)} M[x_n, u_n] \log \theta'_{x_n | u_n} \end{aligned}$$

$$- \max_{\theta'} \ell'(\theta'; D) = \ln \prod_{x_i \in \mathcal{X}_i} \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \sum_{u_i \in \mathcal{U}(x_i)} -M[x_i, u_i] \log \theta'_{x_i | u_i} + \sum_{\text{pa}(x_n) \in \mathcal{X}_{\text{pa}(x_n)}} \sum_{u_n \in \mathcal{U}(x_n)} -M[x_n, u_n] \log \theta'_{x_n | u_n}$$

$$\left\{ \begin{array}{l} \sum_{x_i} \theta_{x_i | u_i} = 1 \quad \forall u_i \rightarrow x_i, \leq 6 \\ \sum_{u_i} \theta_{x_i | u_i} = 1 \end{array} \right\}$$

$$\text{(separable)} \quad = \inf_{\substack{\sum_{x_i} \theta_{x_i | u_i} = 1 \quad \forall u_i \rightarrow x_i, \leq 6 \\ \sum_{u_i} \theta_{x_i | u_i} = 1}} \sum_{x_i \in \mathcal{X}_i} \sum_{\text{pa}(x_i) \in \mathcal{X}_{\text{pa}(x_i)}} \sum_{u_i \in \mathcal{U}(x_i)} -M[x_i, u_i] \log \theta'_{x_i | u_i} + \inf_{\left\{ \sum_{u_n} \theta_{x_n | u_n} = 1 \right\}} \sum_{\text{pa}(x_n) \in \mathcal{X}_{\text{pa}(x_n)}} \sum_{u_n \in \mathcal{U}(x_n)} -M[x_n, u_n] \log \theta'_{x_n | u_n}$$

$$= -\ell(\theta; D) + \text{hyper number}$$

$$\leq -\ell(\theta; D)$$

$$= - \max_{\theta} \ell(\theta; D)$$

$$\therefore \ell_{\theta'}(\theta'; D) \geq \ell_{\theta}(\theta; D)$$

5. [47 points] (*Programming Assignment*) **Tree Augmented Naive Bayes**

Naive Bayes (NB) classifiers are often competitive classifiers even though their strong independence assumptions may be unrealistic. We will explore an extension of this classification framework, Tree Augmented Naive Bayes (TANB), which will aim to improve the results of NB by modelling correlations between attributes while still being computationally tractable. If C denotes the class variable and (A_1, \dots, A_N) the attributes, then a NB model can be represented as a directed graph with these variables as nodes and edges $\{(C, A_i) : 1 \leq i \leq n\}$. A TANB extends this graph structure by allowing each child A_i to have one additional node A_j as a parent with the restriction that these additional edges (A_i, A_j) must form a tree. Figure 1 illustrates the difference in the graph structure.

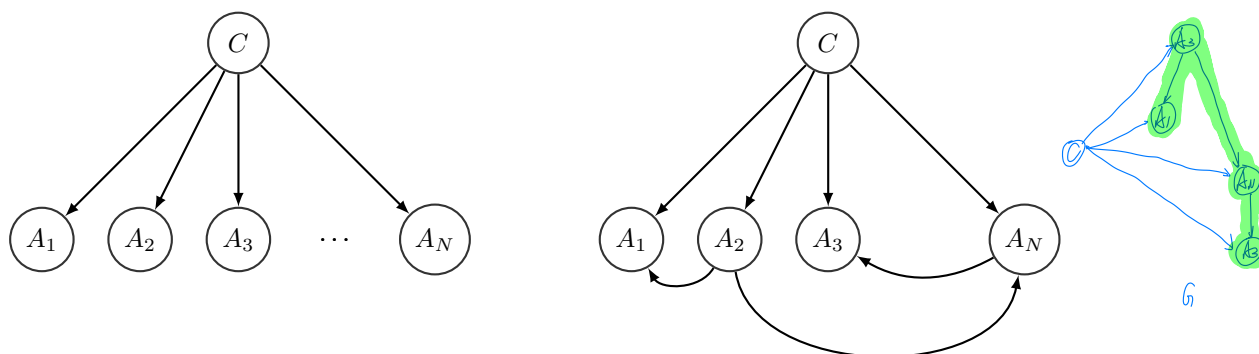


Figure 1: A Naive Bayes model and Tree Augmented Naive Bayes model

In general, learning the structure of this graph can be hard, but we will cover this topic later in the course. For now we have supplied starter code necessary to learn the structure of the graph for the dataset.

In this problem, you will implement the parameter learning for both NB and TANB classifiers and explore the trade-offs of both. You will apply these classifiers to predict the party affiliation of either Democrat or Republican of US Congressmen (the class variable) based on their votes for 16 different measures (the attribute variables) shown in Table 1. Not all Congressmen voted on all 16 measures so sometimes entries in this dataset will have missing attributes, however, we will still be able to utilize our Bayes Network to accurately classify these examples. To keep things simple, the class and attribute variables are all binary with 0, 1 corresponding to a no, yes vote respectively.

When training the models, some of the parameters may not have enough examples for accurate estimation. To mitigate this, we will perform Laplace smoothing with a pseudocount of $\alpha = 0.1$ for every value of an attribute given its parents when learning the parameters.

In order to evaluate the performance of our classifiers on the dataset, we will use 10-fold cross-validation. Under 10-fold cross-validation, the dataset is first partitioned into 10 equally sized partitions. Of these 10 partitions, one of them is used for the test set while the rest of the data are used as the training set to compute test error on this partition. This process is repeated for the other nine partitions, and we can take the average of the resulting test errors to obtain the 10-fold CV test error. We have implemented this procedure for you in the function `evaluate`.

(a) [8 points]

Implement a NB classifier that both learns the parameters from the training data and can use these parameters to score and classify examples in the training data. What is your test error rate using 10-fold cross-validation?

91.81% accuracy, 8.19% error

Note: you can use the `evaluate` function in the starter code, but leave the optional argument `train_subset` to its default value until part d.

(b) [12 points]

We have provided starter code that implements the Chow-Liu algorithm (see Friedman et al. 1997) to learn the tree structure from the training data for a TANB. Using this tree, implement a TANB

Attribute	Name	Incomplete Entry 1
A_1	handicapped infants	1
A_2	water project cost sharing	1
A_3	adoption of the budget resolution	0
A_4	physician fee freeze	0
A_5	El Salvador aid	0
A_6	religious groups in schools	0
A_7	anti satellite test ban	0
A_8	aid to Nicaraguan Contras	?
A_9	mx missile	?
A_{10}	immigration	0
A_{11}	synfuels corporation cutback	0
A_{12}	education spending	?
A_{13}	superfund right to sue	?
A_{14}	crime	?
A_{15}	duty free exports	0
A_{16}	export administration act south Africa	0

Table 1: Attribute names for Congressional Voting Records together with an incomplete example that has some voting records missing for a particular Congressman.

classifier that learns the parameters from the training data and uses these to score examples. What is your test error rate using 10-fold cross-validation? 95.16% accuracy 4.74% error

(c) [15 points]

In general, working with data where the values of attributes and labels are missing is difficult when learning model parameters. However, we can still use our generative model from a fully trained Bayes Network to classify examples in which some of the attributes may be unobserved. Suppose A_i is unobserved. We can still compute $P(C|A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_N)$ by computing $P(C, A_i|A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_N)$ and marginalizing over A_i . Update your TANB implementation to handle the case where some attributes may have missing values and use this new implementation to classify Incomplete Entry 1 in Table 1¹. Given the observed attributes, what is the marginal probability this Congressman is Democrat ($C = 1$) given the votes we did observe? Can you predict how this Congressman voted on education spending (A_{12})?

Note the power of a generative model: it can easily handle missing data and can be used to answer all sorts of probabilistic queries. In contrast, this would *not* be possible with a discriminative model, e.g., if you trained a logistic regression or a random forest classifier to directly predict the affiliation of a Congressman (C), because these models would only provide you with the conditional distribution $P(C|A_1, \dots, A_N)$.

Note: You should train your classifier on the full dataset. You can use the function `evaluate_incomplete_entry`, which both trains on the full dataset and loads Incomplete Entry 1 for classification. `evaluate`.

(d) [12 points]

¹Your implementation does not have to be fully general. It's fine as long as it works on that particular example.

What is the test error when you train both classifiers on a smaller subset of the training data? Explain why the test error for the TANB may not strictly be better than NB.

Note: set the arguments `train_subset=True` when calling `evaluate` so that the classifiers are trained on a much smaller subset of the data.

on small dataset

Naive Bayes : 90.09 %

TANB : 89.66 %

TANB has more parameters, may overfit