

3 hrs	10 / Nov / 2023	Venue	EXAMS OFFICE USE ONLY
-------	-----------------	-------	--------------------------

University of the Witwatersrand, Johannesburg

Course or topic No(s)

COMS40XXA/70XXA

Course or topic name(s)
Paper Number & title

Probabilistic Graphical Models

Examination to be held during the month(s) of

August 2023

Year of study

Degrees/Diplomas for which this course is prescribed

BScHons (CS / BDA / CAM), MSc (AI / DS / CS / Robotics/ e-Science)

Faculties presenting candidates

Science

Internal examiner(s)

Prof. Ritesh Ajoodha
x-76188

External examiner(s)

Prof. External Name (Ext Univ)

Special materials

Formula sheet and non-programmable calculator permitted

Time allowance

3 Hours

Instructions to candidates

Please answer all questions in this closed book test. A total of 100 marks are available, which corresponds to 100%. The test comprises of 17 pages.

Question 1 Multiple Choice Questions [10 Marks]

1. For each of the following MCQ questions, circle the correct answer label.
- 1.1 How does representing the joint distribution using the chain rule for probabilities make it intractable? [2]
 - (a) The distribution is computationally expensive to manipulate in memory.
 - (b) Probabilistic inference would take a long time.
 - (c) It is impossible to elicit priors for all the specified parameters from a human expert.
 - (d) A large amount of data is required because of fragmentation.
 - (e) All of the options above.
- 1.2. Suppose that you have a variable X with 2 dependencies (Y and Z). X , Y , and Z can all take one of 3 different values. Which of the following options specifies the length of the Tabular CPD for the variable X ? [2]
 - (a) 8
 - (b) 9
 - (c) 27
 - (d) 30
 - (e) None of the above
- 1.3. Identify the assumption made in dynamic Bayesian networks that pertains to time from the options provided. [2]
 - (a) System is deterministic up to a point, then becomes completely random.
 - (b) Process being modeled remains statistically constant over time.
 - (c) Variable's rate of change determined by a time-based random number generator.
 - (d) Event probability at a time depends on bird-to-temperature ratio.
- 1.4. Which of the following statements is false when selecting a structure for a Bayesian network? [2]
 - (a) The structure should follow a causal ordering.
 - (b) The structure should be sparse
 - (c) The structure should contain as many edges as possible.
 - (d) The structure should be acyclic
 - (e) The structure should contain relevant variables
- 1.5. Which of the following statements is a limitation of variable elimination (VE)? [2]
 - (a) VE produces exact posterior distribution for any query.
 - (b) VE leads to a more compact factorization than other methods.
 - (c) VE is restricted to DAGs and not undirected graphs.
 - (d) VE can produce incorrect results if the network has continuous variables.
 - (e) VE is impractical for dense networks with many variables.

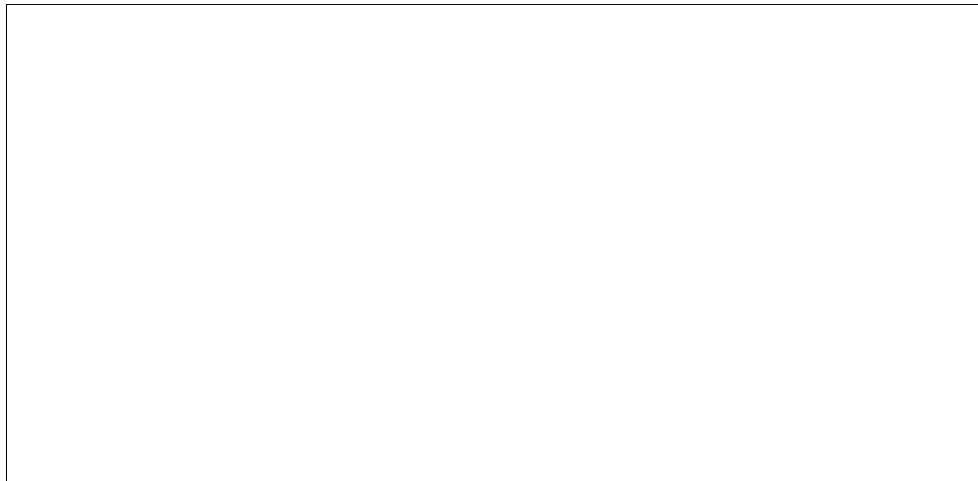
-
-
-
- (b) Does the following set of independences that correspond to d-separation hold true in the context of the graph \mathcal{G} ?

$$\mathcal{I}(\mathcal{G}) = \{(A \perp B \mid C) : \text{d-sep}_{\mathcal{G}}(A : B \mid C)\}$$

Explain your answer.

[3]

-
-
-
- (c) Draw the Bayesian network resulting from the minimal I-map constructed using the independence properties observed in [Figure 1](#), which has the variable ordering of D, B, A, and C. [4]



-
-
-
- (d) What does the term ‘I-equivalence’ mean in Bayesian network theory? Is the network in [Figure 1](#) and the network in the answer to (c) above I-equivalent? [3]

- (d) List all of the rules, $\rho = \langle \mathbf{c}; p \rangle$, which together give you the Rule-CPD that hold in \mathcal{T} . [6]

Question 4**Template-based Models****[20 Marks]**

4.1. Use the below parametrization of a Hidden Markov Model (HMM), λ , to answer the following questions.

- i. Number of hidden states: 3.
- ii. Number of observable symbols: 3.
- iii. Initial state probabilities: $\pi_1 = 0.4, \pi_2 = 0.3, \pi_3 = 0.3$
- iv. The transition probability matrix is: $\begin{pmatrix} 0.2 & 0.4 & 0.4 \\ 0.4 & 0.5 & 0.1 \\ 0.3 & 0.3 & 0.4 \end{pmatrix}$
- v. The observation model is provided by the following matrix: $\begin{pmatrix} 0.3 & 0.3 & 0.4 \\ 0.3 & 0.1 & 0.6 \\ 0.7 & 0.1 & 0.2 \end{pmatrix}$

(a) Define the Markov assumption for temporal models. [2]

(b) Draw the structure of \mathcal{H} with the associated probabilities labelled at the correct edges. Remember to label each node. [4]

(c) Calculate the initial state probabilities for the forward variables $\alpha_1(1)$, $\alpha_1(2)$, and $\alpha_1(3)$ using λ for the sequence $O = \{0, 2\}$. [3]

- (d) Compute the forward variables $\alpha_2(1)$, $\alpha_2(2)$, and $\alpha_2(3)$ using λ for the first induction step using the initial state probabilities from the previous question. [6]

- (e) Using the previous question, calculate the $P(O \mid \lambda)$ for the sequence $O = \{0, 2\}$. [2]

- (f) Unroll the plate model illustrated by Figure 3 and specify the variable name and scope as a function in each node. On the unrolled model indicate the shared parameters for all nodes. [3]

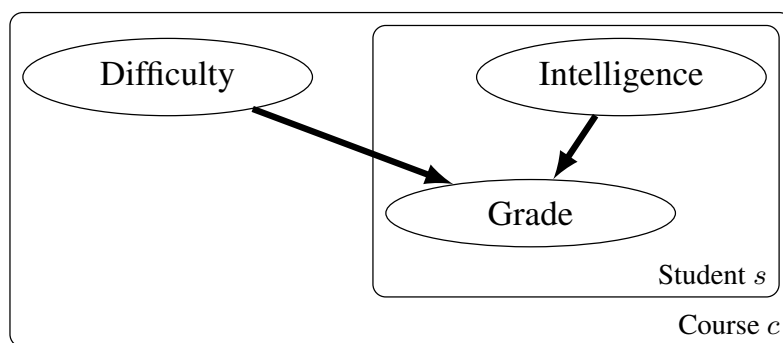
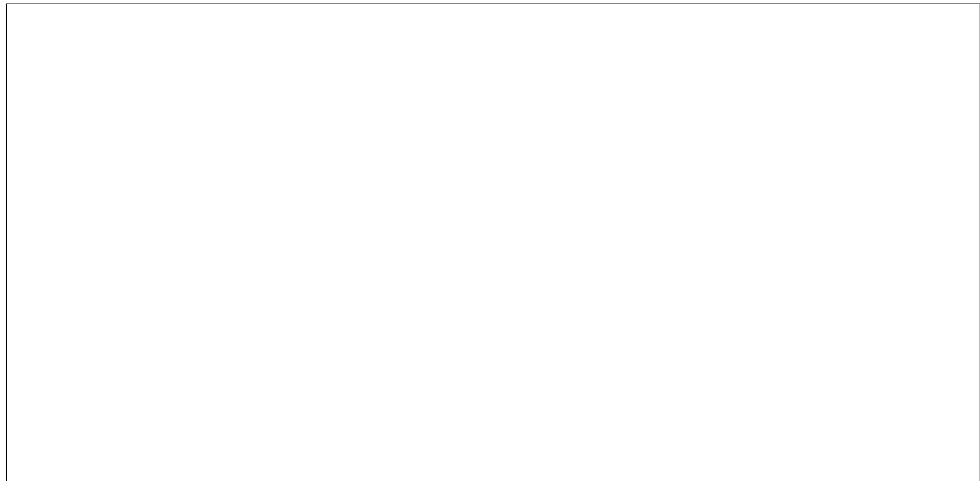


Figure 3: A plate model.



Question 5 Undirected Graphical Models [16 Marks]

5.1. Use the below factor graph, \mathcal{H} , to answer the following questions.

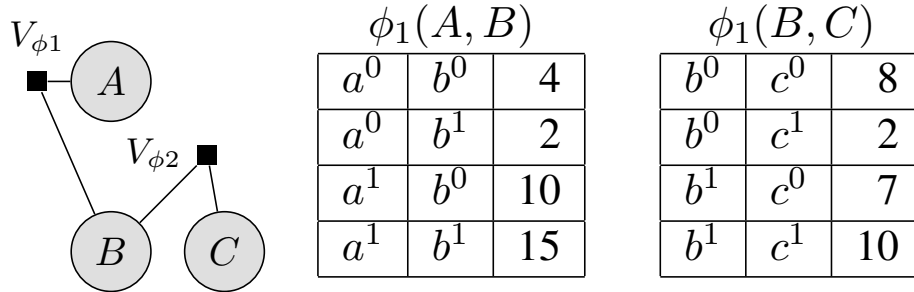


Figure 4: A Markov network with associated factors

- (a) Draw the Markov network structure associated with \mathcal{H} . [2]

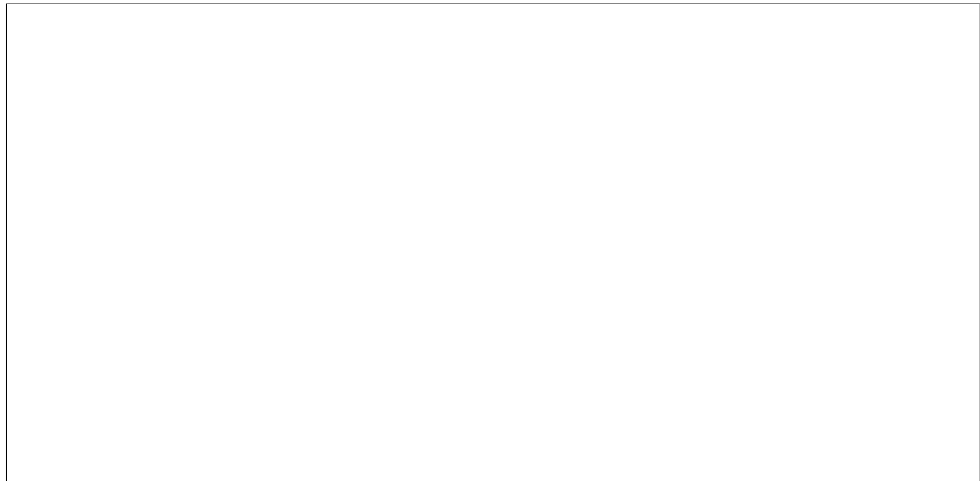
- (b) Does the following set of independences that correspond to variable separation hold true in the context of the graph \mathcal{H} ?

$$\mathcal{I}(\mathcal{H}) = \{(A \perp C \mid MB_{\mathcal{H}}(A))\},$$

where $MB_{\mathcal{H}}(A)$ is the Markov blanket of A. Explain your answer. [2]

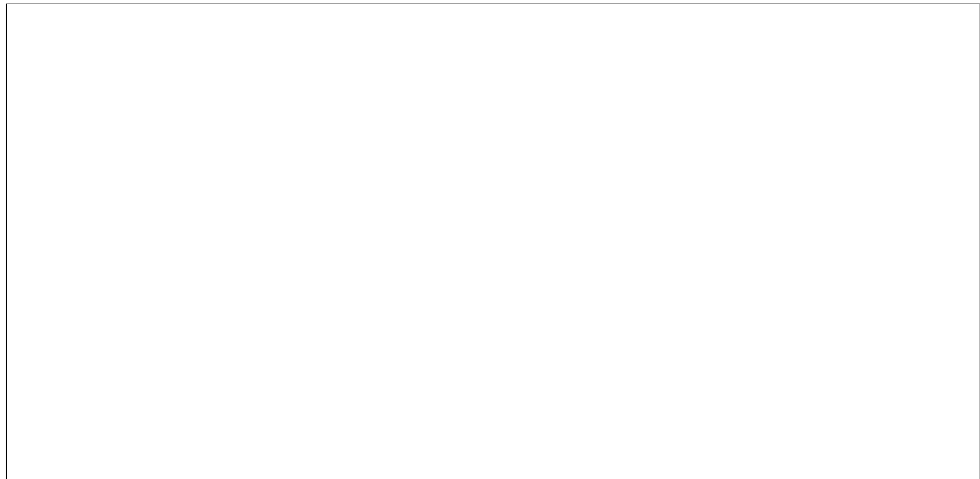
- (c) Write the joint distribution of \mathcal{H} using corresponding maximal clique potentials, Φ . [2]

- (d) Compute the factor product $\psi_3(A, B, C) = \phi_1(A, B) \times \phi_2(B, C)$. [4]



(e) Calculate the value of the partition function. [2]

(f) Calculate $\psi[B = b^1](A, C)$. [2]



(g) Calculate $P(b^1)$. Round off two decimal places. [2]

Question 6 Exact and Approximate Inference [20 Marks]

6.1. Use the Bayesian network, \mathcal{B} , illustrated in Figure 5 to answer the following questions.

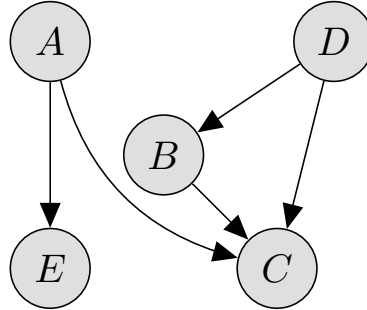


Figure 5: A Bayesian network

- (a) Factorise the joint distribution using the chain rule for Bayesian networks in a way that corresponds to \mathcal{B} . [2]

- (b) Using variable elimination algorithm compute $P(A)$. Show all the required steps of the algorithm using the elimination ordering: $\prec = \{C, B, D, E\}$. [10]

-
-
-
-
- (c) Draw the induced graph $\mathcal{I}_{\mathcal{B}, \prec}$ which results in this ordering of variable elimination. [2]

- (d) Draw the clique tree for $\mathcal{I}_{\mathcal{B}, \prec}$ [2]

- (e) What does it mean for a cluster graph to have the running intersection property? Does the clique tree from the previous question satisfy the running intersection property? [2]
-
-
-
-

- (f) Suppose that you are in the middle of a message-passing procedure using the cluster graph \mathcal{C} as shown in Figure 6. The first message is $\delta_{1 \rightarrow 2}(B)$, what would be the message $\delta_{2 \rightarrow 3}(C)$? [2]

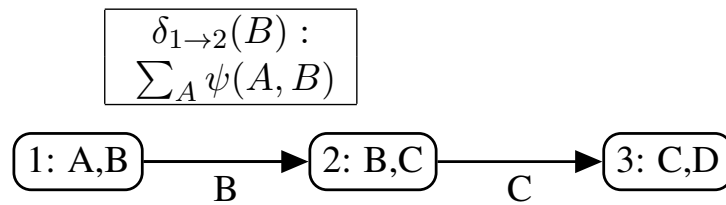


Figure 6: A cluster graph \mathcal{C}

END OF TEST

Working out

Working out

Probabilistic Graphical Models

Formula Sheet

Probability Theory

Chain Rule for Probabilities:

$$P(X_1, \dots, X_n) = P(X_1) \dots P(X_n | X_1, X_{n-1})$$

Bayes Rule:

$$P(\alpha | \beta) = \frac{P(\beta | \alpha)}{P(\alpha)P(\beta)}$$

Probability Density Function: $p : \mathbb{R} \rightarrow \mathbb{R}$ is a

probability density function (PDF) for \mathcal{X} if it is a non-negative integrable function such that:

$$\int_{\mathcal{V}al(\mathcal{X})} p(x) dx = 1.$$

Uniform Distribution: $X \sim \text{Unif}[a, b]$ if it has the PDF:

$$p(x) = \begin{cases} \frac{1}{b-a} & b \geq x \geq a \\ 0 & \text{otherwise.} \end{cases}$$

Gaussian Distribution: X has a Gaussian

distribution: $X \sim \mathcal{N}(\mu; \sigma^2)$ if it has the PDF:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

Joint Density Function: Let P be a joint

distribution over X_1, \dots, X_n . A function

$p(x_1, \dots, x_n)$ is a joint density function of X_1, \dots, X_n if:

1. $p(x_1, \dots, x_n) \geq 0 \forall x_1, \dots, x_n \in X_1, \dots, X_n$.
2. p is integratable.
3. For any choice of a_1, \dots, a_n and b_1, \dots, b_n :

$$P(a_1 \leq X_1 \leq b_1, \dots, a_n \leq X_n \leq b_n) \\ = \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} p(x_1, \dots, x_n) dx_1 \dots dx_n$$

Conditional Density Function: Suppose you would like to condition over the event:

$$x - \epsilon \leq X \leq x + \epsilon.$$

Then

$$P(Y | x) = \lim_{\epsilon \rightarrow 0} P(Y | x - \epsilon \leq X \leq x + \epsilon).$$

If there is a continuous joint density function $p(x, y)$ then

$$P(a \leq Y \leq b | x - \epsilon \leq X \leq x + \epsilon) \\ = \frac{P(a \leq Y \leq b, x - \epsilon \leq X \leq x + \epsilon)}{P(x - \epsilon \leq X \leq x + \epsilon)} = \frac{\int_a^b \int_{x-\epsilon}^{x+\epsilon} p(x', y) dy dx'}{\int_{x-\epsilon}^{x+\epsilon} p(x') dx'}$$

Expectation of X under P is:

$$\mathbb{E}_P[X] = \sum_x x.P(x).$$

Expectation if \mathbf{X} is Continuous:

$$\mathbb{E}_P[X] = \int x.p(x) dx.$$

Linearity of Expectation:

$$\mathbb{E}_P[X + Y] = \mathbb{E}_P[X] + \mathbb{E}_P[Y].$$

Conditional Expectation:

$$\mathbb{E}_P[X | \mathbf{y}] = \sum_x x.P(x | \mathbf{y}).$$

Variance of \mathbf{X} :

$$\mathbb{V}ar_P[X] = \mathbb{E}_P[(X - \mathbb{E}_P[X])^2].$$

Standard Deviation:

$$\sigma_X = \sqrt{\mathbb{V}ar_P[X]}.$$

Expectation and Variance of Gaussian

distribution $X \sim \mathcal{N}(\mu; \sigma^2)$, then $\mathbb{E}[X] = \mu$ and

$$\mathbb{V}ar[X] = \sigma^2.$$

Graph Theory

A **Graph** is a data structure $\mathcal{K} = (\mathcal{X}, \mathcal{E})$ consisting

of a set of nodes, denoted $\mathcal{X} = X_1, \dots, X_n$, and

edges, denoted \mathcal{E} .

Induced Subgraph: Let $\mathcal{K} = (\mathcal{X}, \mathcal{E})$, and $\mathbf{X} \in \mathcal{X}$,

then an induced subgraph, denoted $\mathcal{K}[\mathbf{X}]$ is a graph

$(\mathbf{X}, \mathcal{E}')$ where \mathcal{E}' are all the edges $X \rightleftharpoons Y \in \mathcal{E}'$ such

that $X, Y \in \mathbf{X}$.

Complete Graph (Clique): A subgraph over \mathbf{X} is

complete if every two nodes in \mathbf{X} are connected by

some edge. The set \mathbf{X} is called a clique. A clique \mathbf{X}

is maximal if for any superset of nodes $\mathbf{Y} \supset \mathbf{X}$, \mathbf{Y} is

not a clique.

Upward Closure: A subset of nodes $\mathbf{X} \in \mathcal{X}$ is

upwardly closed in \mathcal{K} if, for any $\mathbf{X} \in \mathcal{X}$, we have that

the Boundary $\mathbf{x} \subset \mathbf{X}$. We define upward closure of \mathbf{X}

to be the minimally upward closed subset \mathbf{Y} that

contains \mathbf{X} .

Topological ordering: An ordering of the nodes

X_1, \dots, X_n is a topological ordering if when we have

$$(X_i \rightarrow X_j) \in \mathcal{E}, \text{ then } i < j.$$

Chordal Graph: Let $X_1 - X_2 - \dots - X_k - X_1$ be a

loop in a graph. A chord in a loop is an edge

connecting X_i and X_j for two nonconsecutive nodes

X_i, X_j . An undirected graph \mathcal{H} is said to be chordal

if and loop $X_1 - X_2 - \dots - X_k - X_1$ for $k > 4$ has a

chord. A directed graph \mathcal{K} is said to be chordal if its

underlying undirected graph is chordal.

Bayesian Networks

Naïve Bayes:

$$P(C, X_1, \dots, X_n) = P(C) \prod_{i=1}^n P(X_i | C)$$

Bayesian Network:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Pa}_{X_i}^{\mathcal{G}})$$

Deterministic CPD: $f : \mathcal{V}al(P_{a_X}) \mapsto \mathcal{V}al(X)$ s.t.:

$$P(x | pa_x) = \begin{cases} 1 & \text{if } x = f(pa_x) \\ 0 & \text{if } x \text{ otherwise} \end{cases}$$

Time Granularity Assumption:

$$P(\mathcal{X}^{(0:T)}) = P(\mathcal{X}^{(0)}) \prod_{t=0}^{T-1} P(\mathcal{X}^{(t+1)} | \mathcal{X}^{(0:t)})$$

Markov Assumption:

$$P(\mathcal{X}^{(0:T)}) = P(\mathcal{X}^{(0)}) \prod_{t=0}^{T-1} P(\mathcal{X}^{(t+1)} | \mathcal{X}^{(t)})$$

Time Invariance Assumption:

$$P(\mathcal{X}^{(t+1)} = \xi' | \mathcal{X}^{(t)} = \xi) = P(\mathcal{X}' = \xi' | \mathcal{X} = \xi)$$

Two-TBN:

$$P(\mathcal{X}' | \mathcal{X}) = P(\mathcal{X}' | \mathcal{X}_t) = \prod_{i=1}^n P(X'_i | P_{a_{X'_i}})$$

Linear Dynamical Systems:

$$P(\mathbf{X}^{(t)} | \mathbf{X}^{(t-1)}) = \mathcal{N}(\mathbf{A}\mathbf{X}^{(t-1)}; Q)$$

$$P(O^{(t)} | \mathbf{X}^{(t)}) = \mathcal{N}(\mathbf{H}\mathbf{X}^{(t)}; R)$$

Gibbs Distribution: A distribution P_{Φ} is a Gibbs

distribution parameterised by a set of factors

$\Phi = \{\phi_1(\mathbf{D}_1), \dots, \phi_K(\mathbf{D}_K)\}$ if it is defined as:

$$P_{\Phi}(X_1, \dots, X_n) = \frac{1}{Z} P_{\Phi}(X_1, \dots, X_n)$$

Inference

Inference:

$$P(\mathbf{Y} | \mathbf{E} = \mathbf{e}) = \frac{P(\mathbf{Y}, \mathbf{e})}{P(\mathbf{e})} = \frac{\sum_w P(\mathbf{y}, \mathbf{e}, \mathbf{w})}{\sum_{y', w} P(\mathbf{e})}$$

Sum-Product Message Passing:

$$\delta_{i \rightarrow j} = \sum \mathbf{C}_i - \mathbf{s}_{i,j} (\psi_i \times \prod_{k \in (N\mathbf{b}_i - \{j\})} \delta_{k \rightarrow i})$$

Tree Calibration:

$$\sum \mathbf{C}_i - \mathbf{s}_{i,j} \beta_i(\mathbf{C}_i) = \sum \mathbf{C}_j - \mathbf{s}_{i,j} \beta_j(\mathbf{C}_j)$$

Graph Calibration:

$$\sum \mathbf{C}_i - \mathbf{s}_{i,j} \beta_i = \sum \mathbf{C}_j - \mathbf{s}_{i,j} \beta_j$$

MAP:

$$\text{MAP}(\mathbf{Y} = \mathbf{y} | \mathbf{E} = \mathbf{e})$$

$$= \text{argmax}_{\mathbf{y}} P(\mathbf{Y} = \mathbf{y} | \mathbf{E} = \mathbf{e})$$

Convergence Bound:

$$\mathbb{E}_{\mathcal{D}}(f) = \frac{1}{M} \sum_{m=1}^M f(\xi[m]).$$

Hoeffding Bound:

$$P_{\mathcal{D}}(\hat{P}(\mathbf{y}) \notin [P(\mathbf{y}) - \epsilon, P(\mathbf{y}) + \epsilon]) \leq 2e^{-2M\epsilon^2}$$

Chernoff Bound:

$$P_{\mathcal{D}}(\hat{P}(\mathbf{y}) \notin [P(\mathbf{y})(\pm\epsilon)]) \leq 2e^{-MP(\mathbf{y})\epsilon^2/3}$$

$$M \geq 3 \frac{\ln(2/\delta)}{P(\mathbf{y})\epsilon^2}.$$

Likelihood Weighting:

$$\hat{P}_{\mathcal{D}}(\mathbf{y} \mid \mathbf{e}) = \frac{\sum_{m=1}^M w[m] \mathbb{1}\{\mathbf{y}[m]=\mathbf{y}\}}{\sum_{m=1}^M w[m]}.$$

MCMC Sampling:

$$P^{(t+1)}(\mathbf{X}^{(t+1)} = \mathbf{x}') =$$

$$\sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} P^{(t)}(\mathbf{X}^{(t)} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$$

Stationary Distribution:

$$\pi(\mathbf{X} = \mathbf{x}') = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} \pi(\mathbf{X} = \mathbf{x}) \mathcal{T}(\mathbf{x} \rightarrow \mathbf{x}')$$

Detailed Balance Equation:

$$\pi(x) \mathcal{T}(x \rightarrow x') = \pi(x') \mathcal{T}(x' \rightarrow x)$$

Acceptance Probability:

$$\mathcal{A}(x \rightarrow x') = \min[1, \frac{\pi(x') \mathcal{T}^Q(x' \rightarrow x)}{\pi(x) \mathcal{T}^Q(x \rightarrow x')}].$$

Metropolis-Hastings Acceptance Probability:

$$\mathcal{A}(x_{-i}, x_i \rightarrow x_{-i}, x'_i) = \min[1, \frac{P_{\theta}(x'_i, x_{-i}) \mathcal{T}_{\theta}^Q(x_{-i}, x'_i \rightarrow x_{-i}, x_i)}{P_{\theta}(x_i, x_{-i}) \mathcal{T}_{\theta}^Q(x_{-i}, x_i \rightarrow x_{-i}, x'_i)}].$$

Learning

Relative Entropy:

$$\mathbb{D}(P^* \parallel \tilde{P}) = \mathbb{E}_{\xi \sim P^*} [\log(\frac{P^*(\xi)}{\tilde{P}(\xi)})],$$

Negative Empirical Log-loss:

$$\log P(\mathcal{D} : \mathcal{M}) = \sum_{m=1}^M \log P(\xi[m] : \mathcal{M}).$$

Bayesian Parameter Estimation:

$$P(\theta \mid x[1], \dots, x[M]) = \frac{P(x[1], \dots, x[M] \mid \theta) P(\theta)}{P(x[1], \dots, x[M])}$$

Expected Sufficient Statistics:

$$\bar{M}_{\theta}[\mathbf{y}] = \sum_{m=1}^M \sum_{\mathbf{h}[m] \in \text{Val}(\mathbf{H}[m])} Q(\mathbf{h}[m]) \mathbb{1}\{\xi[m]\langle \mathbf{Y} \rangle = \mathbf{y}\}$$

Maximisation of Expected Parameter:

$$\tilde{\theta}_{d^1|c^0} = \frac{M_{\theta}(d^1, c^0)}{M_{\theta}(c^0)}$$

Bayesian Clustering:

$$\bar{M}_{\theta}[c] = \frac{\bar{M}_{\theta}[c]}{M}$$

$$\bar{M}_{\theta}[x_i \mid c] = \frac{\bar{M}_{\theta}[x_i, c]}{\bar{M}_{\theta}[c]}$$

K-means Clustering:

$$c[m] = \text{argmax}_c P(c \mid x[m], \theta^t)$$

Hypothesis Testing:

$$d_{\mathbb{H}}(\mathcal{D}) = \sum_{x,y} \frac{M[x,y]}{M} \log \frac{M[x,y]/M}{M[x]/M \cdot M[y]/M}$$

$$R_{d,t}(\mathcal{D}) \begin{cases} \text{Accept if } d(\mathcal{D}) \leq t \\ \text{Reject if } d(\mathcal{D}) > t \end{cases}$$

$$\text{p-value}(t) = P(\{\mathcal{D} : d(\mathcal{D}) > t\} \mid H_0, M)$$

Likelihood:

$$\mathbb{I}_{\hat{P}_{\mathcal{D}}}(\mathbf{X}_i; Pa_{\mathbf{X}_i}^G)$$

$$= \sum_{\mathbf{u}_i} \sum_{\mathbf{x}_i} \hat{P}(x_i, \mathbf{u}_i) \log \frac{\hat{P}(x_i, \mathbf{u}_i)}{\hat{P}(x_i) \hat{P}(\mathbf{u}_i)}$$

Entropy:

$$\mathbb{H}_{\hat{P}_{\mathcal{D}}}(\mathbf{X}_i) = \sum_{x_i} \hat{P}(x_i) \log \frac{1}{\hat{P}(x_i)}$$

Bayesian Structure Learning:

$$P(\mathcal{G} \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \mathcal{G}) P(\mathcal{G})}{P(\mathcal{D})}$$

$$\text{score}_{\mathcal{B}}(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{D} \mid \mathcal{G}) + \log P(\mathcal{G})$$

$$P(\mathcal{D} \mid \mathcal{G}) = \int_{\Theta_{\mathcal{G}}} P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) P(\theta_{\mathcal{G}} \mid \mathcal{G}) d\theta_{\mathcal{G}}$$

Marginal Likelihood for Binomials:

$$P(x[1], \dots, x[M])$$

$$= P(x[1]) \cdot \dots \cdot P(x[M] \mid x[1], \dots, x[M-1])$$

Marginal Likelihood for Multinomials:

$$P(x[1], \dots, x[M]) = \frac{\Gamma(\alpha)}{\Gamma(\alpha+M)} \cdot \prod_{i=1}^k \frac{\Gamma(\alpha_i + M[x^i])}{\Gamma(\alpha_i)}$$

Bayesian Score:

$$P(\mathcal{D} \mid \mathcal{G}) = \prod_i \prod_{\mathbf{u}_i \in \text{Val}(Pa_{\mathbf{X}_i}^G)} \frac{\Gamma(\alpha_{\mathbf{X}_i \mid \mathbf{u}_i}^G)}{\Gamma(\alpha_{\mathbf{X}_i \mid \mathbf{u}_i}^G + M[\mathbf{u}_i])}.$$

$$\prod_{\mathbf{x}_i^j \in \text{Val}(\mathbf{X}_i)} \left[\frac{\Gamma(\alpha_{\mathbf{X}_i^j \mid \mathbf{u}_i}^G)}{\Gamma(\alpha_{\mathbf{X}_i^j \mid \mathbf{u}_i}^G + M[x_i^j])} \right]$$

BIC Score:

$$\text{score}_{\text{BIC}}(\mathcal{G} : \mathcal{D}) =$$

$$M \sum_{i=1}^n \mathbb{I}_{\hat{P}_{\mathcal{D}}}(\mathbf{X}_i; Pa_{\mathbf{X}_i}^G) - \frac{\log M}{2} \dim[\mathcal{G}]$$

Decomposability:

$$\text{score}(\mathcal{G} : \mathcal{D}) = \sum_i \text{FamScore}(\mathbf{X}_i \mid Pa_{\mathbf{X}_i}^G : \mathcal{D})$$

Tree weight:

$$w_{i \rightarrow j} = \text{FamScore}(\mathbf{X}_i \mid \mathbf{X}_j : \mathcal{D}) -$$

$$\text{FamScore}(\mathbf{X}_i : \mathcal{D})$$

Learning Graphs:

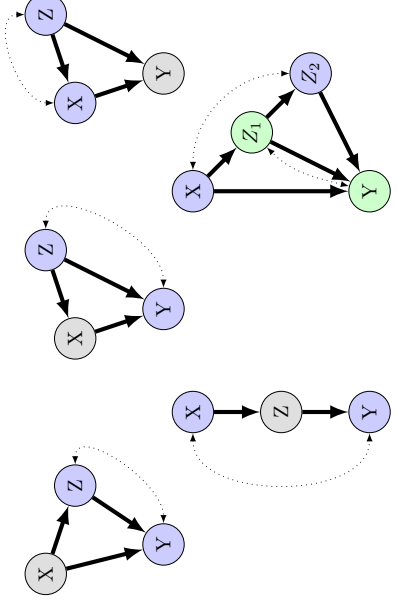
$$\mathcal{G}^* = \text{argmax}_{\mathcal{G} \in \mathcal{G}} \text{score}(\mathcal{G} : \mathcal{D})$$

Causality

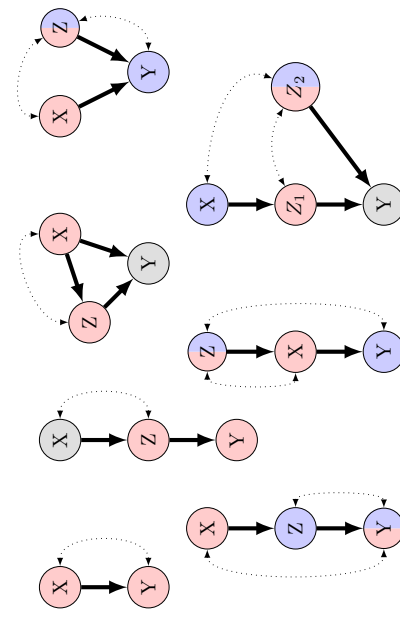
Intervention Query:

$$P_{\mathcal{C}}(\mathbf{Y} \mid do(z), \mathbf{x}) = P_{\mathcal{C}_{z=x}}(\mathbf{Y} \mid \mathbf{x})$$

Identifiable when $P(Y \mid do(X))$:



Not Identifiable when $P(Y \mid do(X))$:



Learning with Intervention Data:

$$P(\xi \mid do(\mathbf{Z} := \mathbf{z}), \mathcal{C}) = \prod_{X_i \notin \mathbf{Z}} P(x_i \mid \mathbf{u}_i)$$

Sufficient Statistics (Intervention Data):

$$M[x_i; \mathbf{u}_i] =$$

$$\sum_{m: X_i \notin \mathbf{Z}[m]} \mathbb{1}\{X_i[m] = x_i, Pa_{X_i}[m] = \mathbf{u}_i\}$$

Likelihood of Data (Intervention):

$$L(\mathcal{C} : \mathcal{D}) = \prod_{i=1}^n \prod_{x_i \in \text{Val}(X_i), \mathbf{u}_i \in \text{Val}(Pa_{X_i})} \theta_{x_i | \mathbf{u}_i}^{M[x_i; \mathbf{u}_i]}$$

Decision Theory

Expected Utility:

$$\text{EU}[D[a]] = \sum_{\mathbf{x}} P(\mathbf{x} \mid a) U(\mathbf{x}, a)$$

Maximum Expected Utility:

$$\begin{aligned} a^* &= \text{argmax}_a \text{EU}[D[a]] \\ &= \text{argmax}_a \sum_{\mathbf{x}} P(\mathbf{x} \mid a) U(\mathbf{x}, a) \end{aligned}$$

Expected Utility with Information:

$$\text{EU}[D[\delta_A]] = \sum_{\mathbf{x}, a} P_{\delta_A}(\mathbf{x}, a) U(\mathbf{x}, a)$$

Maximal Expected Utility (MEU) Strategy:

$$\text{argmax}_{\delta_{D_1}, \dots, \delta_{D_k}} \text{EU}[\mathcal{Z}[\delta_{D_1}, \dots, \delta_{D_k}]]$$

Value of Information:

$$\text{VPI}(A \mid X) := \text{MEU}(D_{X \rightarrow A}) - \text{MEU}(D)$$