

Probabilistic Graphical Models

Formula Sheet

Probability Theory

Chain Rule for Probabilities:

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

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:

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

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 :

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

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:

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

Expectation if \mathbf{X} is Continuous:

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

Linearity of Expectation:

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

Conditional Expectation:

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

Variance of \mathbf{X} :

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

Standard Deviation:

$$\sigma_X = \sqrt{\text{Var}_P[X]}.$$

Expectation and Variance of Gaussian distribution $X \sim \mathcal{N}(\mu; \sigma^2)$, then $\mathbb{E}[X] = \mu$ and $\text{Var}[X] = \sigma^2$.

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 $_X \subset \mathbf{X}$. 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 $(X_i \rightarrow X_j) \in \mathcal{E}$, 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:

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

Deterministic CPD: $f: Val(Pa_X) \mapsto Val(X)$ 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:

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

Linear Dynamical Systems:

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

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

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:

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

Inference

Inference:

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

Sum-Product Message Passing:

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

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:

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

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:

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

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:

$$\hat{\theta}_{d^1 \mid c^0} = \frac{\bar{M}_{\theta}[d^1, c^0]}{\bar{M}_{\theta}[c^0]}$$

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:

$$c[m] = \arg\max_c P(c \mid x[m], \theta^t)$$

Hypothesis Testing:

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

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

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

Likelihood:

$$\mathbb{I}_{\hat{P}_{\mathcal{D}}}(\mathbf{X}_i; Pa_{\mathbf{X}_i}^{\mathcal{G}}) = \sum_{\mathbf{u}_i} \sum_{\mathbf{x}_i} \hat{P}(x_i, \mathbf{u}_i) \log \frac{\hat{P}(x_i, \mathbf{u}_i)}{\hat{P}(x_i) \hat{P}(\mathbf{u}_i)}$$

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:

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

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

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

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:

$$P(\mathcal{D} \mid \mathcal{G}) = \prod_i \prod_{\mathbf{u}_i \in \text{Val}(Pa_{\mathbf{X}_i}^{\mathcal{G}})} \frac{\Gamma(\alpha_{\mathbf{X}_i \mid \mathbf{u}_i}^{\mathcal{G}})}{\Gamma(\alpha_{\mathbf{X}_i \mid \mathbf{u}_i}^{\mathcal{G}} + M[\mathbf{u}_i])} \cdot \prod_{\mathbf{x}_i^j \in \text{Val}(\mathbf{X}_i)} \left[\frac{\Gamma(\alpha_{\mathbf{x}_i^j \mid \mathbf{u}_i}^{\mathcal{G}} + M[x_i^j, \mathbf{u}_i])}{\Gamma(\alpha_{\mathbf{x}_i^j \mid \mathbf{u}_i}^{\mathcal{G}})} \right]$$

BIC Score:

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

Decomposibility:

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

Tree weight:

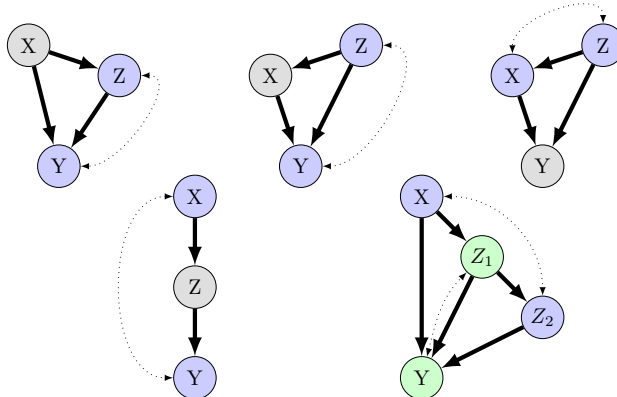
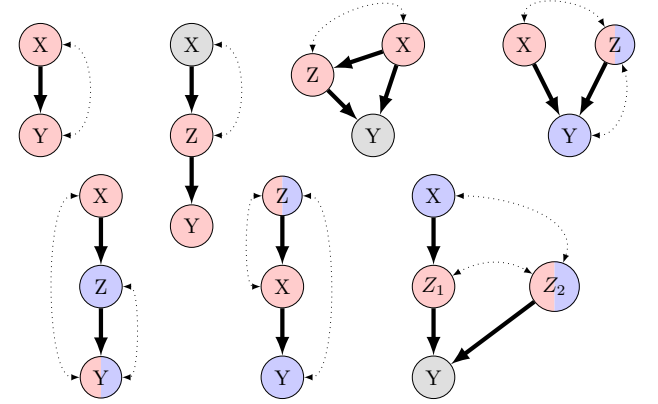
$$w_{i \rightarrow j} = \text{FamScore}(\mathbf{X}_i \mid \mathbf{X}_j : \mathcal{D}) - \text{FamScore}(\mathbf{X}_i : \mathcal{D})$$

Learning Graphs:

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

Causality**Intervention Query:**

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

Identifiable when $P(Y \mid do(X))$:**Not Identifiable when $P(Y \mid do(X))$:****Learning with Intervention Data:**

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

Sufficient Statistics (Intervention Data):

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

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

Likelihood of Data (Intervention):

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

Decision Theory**Expected Utility:**

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

Maximum Expected Utility:

$$a^* = \arg\max_a \text{EU}[D[a]]$$

$$= \arg\max_a \sum_{\mathbf{x}} P(\mathbf{x} \mid a) U(\mathbf{x}, a)$$

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:

$$\arg\max_{\delta_{D_1}, \dots, \delta_{D_k}} \text{EU}[\mathcal{I}[\delta_{D_1}, \dots, \delta_{D_k}]]$$

Value of Information:

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