



Maximum Likelihood Estimation for Bayesian Network

eg. $(X \rightarrow Y)$



given a graph $D = \{(x[m], y[m])\}$

$$L(\theta; D) = \prod_{m=1}^M p(x[m], y[m]; \theta)$$

$$= \prod_{m=1}^M p(x[m]; \theta) \cdot p(y[m] | x[m]; \theta)$$

$$= \left(\prod_{m=1}^M p(x[m]; \theta) \right) \cdot \left(\prod_{m=1}^M p(y[m] | x[m]; \theta) \right)$$

each of these terms is a "local likelihood" that measures how well the variables is predicted given its parents

$$l(\theta, D) = \sum_{m=1}^M \log p(x[m]; \theta) + \sum_{m=1}^M \log p(y[m] | x[m]; \theta)$$

$$= M[X] \cdot \log \theta_X + M[X] \cdot \log \theta_Y + \sum_{\text{children}(X)} M[X, Y] \cdot \log \theta_{Y|X}$$

$$\min_{\theta} -C^T \log \theta$$

$$f(\theta, v) = -C^T \log \theta + V(L^T \theta - 1)$$

$$\nabla_{\theta} f = -C^{-1} + V L = 0$$

$$L^T \theta = 1$$

$$\theta = \frac{C}{V L}$$

$$g(\theta) = -C^T (\log C - \log V L) + V(L^T \frac{C}{V L} - 1)$$

$$= -C^T \log C + C^T \log V L + V \cdot L^T \frac{C}{V L} - V$$

$$= -C^T \log C + C^T \log V L + L^T C - V$$

$$\nabla_{\theta} g = -\frac{L^T C}{V} - 1 = 0$$

$$V = L^T C$$

$$\theta = C / L^T C$$

$$\theta_X = \frac{M[X]}{M[X] + M[Y]} \quad \theta_Y = \frac{M[Y]}{M[X] + M[Y]}$$

$$\theta_{Y|X} = \frac{M[Y|X]}{\sum_{\text{children}(X)} M[Y|X]} = \frac{M[Y|X]}{M[X]}$$

$$\hat{\theta}_{X|U} = \frac{M[X|U]}{\sum_{\text{children}(U)} M[X|U]} = \frac{M[X|U]}{M[U]} \quad \forall U \in \text{par}(X) \text{ in } \text{par}(X)$$

Gaussian Bayesian Networks

$$p(X|U) = N(\beta_0 + \beta_1 U_1 + \dots + \beta_n U_n; \theta^*)$$

$$l(\theta_{X|U}; D) = \sum_{i=1}^M \log N(X[i]; \beta_0 + \dots + \beta_n U_i, \theta^*)$$

$$= \sum_{i=1}^M -\frac{1}{2} \log(2\pi\theta^*) - \frac{1}{2} \frac{1}{\theta^*} (\beta_0 + \dots + \beta_n U_i - X[i])^2$$

$$\nabla_{\theta} l = \sum_{i=1}^M -\frac{1}{\theta^*} (\beta_0 + \dots + \beta_n U_i - X[i]) = 0$$

$$\sqrt{\frac{M}{M-n}} \beta_0 = \sqrt{\frac{M}{M-n}} \beta_0 + \dots + \sqrt{\frac{M}{M-n}} \beta_n U_i = \beta_0 + \beta_1 \frac{1}{\sqrt{M-n}} \sum_i U_i + \dots + \beta_n \frac{1}{\sqrt{M-n}} \sum_i U_i$$

$$E[\beta_0] = \beta_0 + \beta_1 E[U_1] + \dots + \beta_n E[U_n]$$

$$\nabla_{\beta} \ell = \sum_{m=1}^M -\frac{1}{\theta^T \beta} (\beta_0 + \cdots + \beta_k u_k[m] - X[m]) \cdot u_i[m] := 0$$

$$1/M \sum_{m=1}^M X[m] \cdot u_i[m] = 1/M \sum_{m=1}^M \beta_0 \cdot u_i[m] + \cdots + 1/M \sum_{m=1}^M \beta_k u_k[m] \cdot u_i[m]$$

$$E_D[X \cdot u_i] = \beta_0 E_D[u_i] + \beta_1 E_D[u_i \cdot u_i] + \cdots + \beta_k E_D[u_i \cdot u_i]$$

$$\begin{bmatrix} 1 & E[u_i] & \cdots & E[u_i] \\ E[u_i] & E[u_i u_i] & \cdots & E[u_i u_i] \\ E[u_i] & E[u_i u_i] & \cdots & E[u_i u_i] \\ \vdots & \vdots & \vdots & \vdots \\ E[u_i] & E[u_i u_i] & \cdots & E[u_i u_i] \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = \begin{bmatrix} E[X] \\ E[X u_i] \\ E[X u_i] \\ \vdots \\ E[X u_i] \end{bmatrix}$$

Solve for β by solving linear equation

$$\begin{aligned} Cov_D[X, u_i] &= E_D[X \cdot u_i] - E_D[X] \cdot E_D[u_i] \\ &= \beta_0 E_D[u_i] + \beta_1 E_D[u_i \cdot u_i] + \cdots + \beta_k E_D[u_i \cdot u_i] \\ &\quad - \beta_0 E_D[u_i] - \beta_1 E_D[u_i] \cdot E_D[u_i] - \cdots - \beta_k E_D[u_i] \cdot E_D[u_i] \\ &= \beta_1 Cov_D[u_i, u_i] + \cdots + \beta_k Cov_D[u_i, u_i] \end{aligned}$$

$$\begin{aligned} L(\theta, D) &= \sum_{i=1}^M -\frac{1}{2} \log(2\pi\theta) - \frac{1}{2} \frac{1}{\theta^2} (\beta_0 + \cdots + \beta_k u_k[i] - X[i])^2 \\ \nabla_{\theta} L &= \sum_{m=1}^M -\frac{1}{2\theta^2} 2X + \frac{1}{\theta^3} (\beta_0 + \cdots + \beta_k u_k[m] - X[m]) := 0 \\ \frac{M}{\theta^2} \left(\frac{1}{\theta^2} (\beta_0 + \cdots + \beta_k u_k[m] - X[m]) \right) - \frac{1}{\theta} &= 0 \\ \theta^* &= Cov_D[X, X] - \sum_i \sum_j \beta_j \beta_j Cov_D[u_i, u_j] \end{aligned}$$

MLE as M-Projection

let D be a dataset, $L(\theta, D) = \log L(\theta, D) = M \cdot E_{x \sim p_\theta} [\log p(x; \theta)]$

$$\begin{aligned} L(\theta, D) &= \sum_{m=1}^M \log p(X[m]; \theta) \\ &= \sum_{x \in \text{eval}(X)} \left[\sum_m \mathbb{1}_{\{X[m]=x\}} \right] \cdot \log p(x; \theta) \\ &= \sum_{x \in \text{eval}(X)} M \cdot \hat{p}_\theta(x) \cdot \log p(x; \theta) \\ &= M \cdot E_{x \sim p_\theta} [\log p(x; \theta)] \end{aligned}$$

$$D(P||P') = \int_X p(x) \log \frac{p(x)}{p'(x)} dx = E_{x \sim p} [\log \frac{p(x)}{p'(x)}] = -H(p) - E_{x \sim p} [\log p(x)]$$

$$L(\theta; D) = M \cdot (-H_{p(\theta)}(X) - D(p_\theta || P))$$

The maximum likelihood estimation $\hat{\theta}$ in a parametric family is the param of M-projection of \hat{p}_θ on p_0

Joint Probabilistic Model

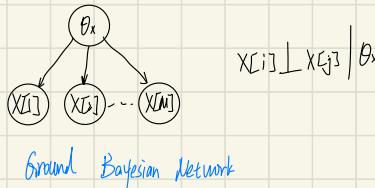
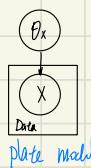
encode the prior knowledge about θ with a probability distribution, which represents how likely different choices of θ are

then create a joint distribution of parameter θ and observation $\{X^{[m]}\}$

e.g. Coin toss

$$P(X^{[m]}=1|\theta) = \theta \quad \theta \sim \text{Uniform}(0,1)$$

$$P(X^{[m]}=0|\theta) = 1-\theta$$



$$\begin{aligned} P(\theta, X^{[1]}, \dots, X^{[m]}) &= P(\theta) \cdot P(X^{[1]}, \dots, X^{[m]} | \theta) \\ &= P(\theta) \cdot \prod_{m=1}^M P(X^{[m]} | \theta) \\ &= P(\theta) \cdot \theta^{N_{[1]}} \cdot (1-\theta)^{N_{[0]}} \end{aligned}$$

Posterior: $P(\theta | X^{[1]}, \dots, X^{[m]}) = \frac{P(\theta, X^{[1]}, \dots, X^{[m]})}{P(X^{[1]}, \dots, X^{[m]})}$

In prediction, instead of selecting a single value for θ , we use it for predicting the probability over the next toss

$$\begin{aligned} P(X^{[M+1]} | X^{[1]}, \dots, X^{[M]}) &= \int P(X^{[M+1]}, \theta | X^{[1]}, \dots, X^{[M]}) d\theta \\ &= \int P(X^{[M+1]} | \theta, X^{[1]}, \dots, X^{[M]}) \cdot P(\theta | X^{[1]}, \dots, X^{[M]}) d\theta \\ &= \int P(X^{[M+1]} | \theta) \cdot P(\theta | X^{[1]}, \dots, X^{[M]}) d\theta \\ &= \int P(X^{[M+1]} | \theta) \cdot \frac{P(\theta) \cdot P(X^{[1]}, \dots, X^{[M]} | \theta)}{P(X^{[1]}, \dots, X^{[M]})} d\theta \\ &= \frac{1}{P(X^{[1]}, \dots, X^{[M]})} \int P(X^{[M+1]} | \theta) \cdot \theta^{N_{[1]}} \cdot (1-\theta)^{N_{[0]}} d\theta \end{aligned}$$

$$P(X^{[M+1]}=1 | X^{[1]}, \dots, X^{[M]}) = \frac{1}{P(X^{[1]}, \dots, X^{[M]})} \int \theta^{N_{[1]+1}} \cdot (1-\theta)^{N_{[0]}} d\theta = \frac{1}{\sum} \frac{(N_{[1]}+1)! \cdot M_{[0]}!}{(N_{[0]}+N_{[1]}+2)!}$$

$$P(X^{[M+1]}=0 | X^{[1]}, \dots, X^{[M]}) = \frac{1}{2} \int_0^1 \theta^{M_{[0]}} \cdot (1-\theta)^{N_{[1]+1}} d\theta = \frac{M_{[0]}! \cdot (N_{[1]}+1)!}{(M_{[0]}+N_{[1]}+2)!}$$

$$\begin{aligned} \sum a+b &= 1 \\ \sum \frac{m}{b} &= \frac{(M_{[0]}+1) \cdot M_{[0]}!}{N_{[1]}! \cdot (M_{[0]}+1)!} = \frac{M_{[0]}+1}{N_{[1]}+1} \Rightarrow P(X^{[M+1]}=1 | X^{[1]}, \dots, X^{[M]}) = \frac{M_{[0]}+1}{N_{[0]}+M_{[0]}+1} \end{aligned}$$

Laplace Smooth

$$\int_0^1 p((k)p)^{pk} dp = \frac{k! \cdot (pk)!}{(pk+1)!}$$

Stony prof:



throw k white balls, m black balls and l red.

$$p(\underline{\text{white red black}})$$

Priors

appropriate distribution is beta distribution

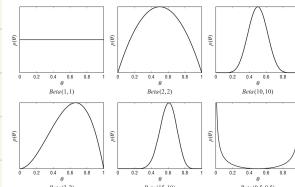
$$\theta \sim \text{Beta}(a_1, a_2) = \frac{\theta^{a_1-1} (1-\theta)^{a_2-1}}{\Gamma(a_1)\Gamma(a_2)}$$

where $\Gamma = \frac{\Gamma(a_1+a_2)}{\Gamma(a_1)\Gamma(a_2)}$ is a normalizing constant, $\Gamma(n) = \int_0^{+\infty} t^{n-1} e^{-t} dt$ is the Gamma function

$$\text{Gamma function } \Gamma(a) = \int_0^{+\infty} t^{a-1} e^{-t} dt \quad (a > 0)$$

$$\Gamma(n) = (n-1)! \quad * \text{ integer } a > 0$$

$$\Gamma(\frac{1}{2}) = \sqrt{\pi}$$



$$\begin{aligned} \Gamma(\frac{1}{2}) &= \int_0^{+\infty} t^{\frac{1}{2}-1} e^{-t} dt \quad u = t^{\frac{1}{2}} \\ &= \int_0^{+\infty} u^{\frac{1}{2}-1} e^{-u^2} du^2 \\ &= 2 \int_0^{+\infty} e^{-u^2} du \\ &= \sqrt{\pi} \\ &= \text{Fn} \end{aligned}$$

$$\begin{aligned} \Gamma(a+b) &= \int_0^{+\infty} t^{a+b-1} e^{-t} dt = \int_0^{+\infty} t^a e^{-t} dt \\ &= -[t^a e^{-t}]_0^{+\infty} - \int_0^{+\infty} e^{-t} a t^{a-1} dt \\ &= a \int_0^{+\infty} e^{-t} t^{a-1} dt \\ \{ \Gamma(a+1) &= a \cdot \Gamma(a) \\ \Gamma(1) &= 1 \end{aligned}$$

$$\int_0^1 \theta^{a-1} (1-\theta)^{b-1} d\theta = \frac{(a-1)! (b-1)!}{(a-1+b-1+1)!} \quad f = \frac{1}{\int_0^1 d\theta} = \frac{(a+b-1)!}{(a-1)! (b-1)!} = \frac{\Gamma(a+b)}{\Gamma(a) \Gamma(b)} \quad (\text{for discrete } a, b)$$

$$\text{eg. } (X|\theta) \sim \text{Binomial}(M, \theta) \sim \text{Beta}(a, b) \quad \text{from } P(\theta|X)$$

$$P(\theta|X=k) = \frac{P(\theta) \cdot P(X=k)}{P(X)} = \frac{\binom{M}{k} \theta^k (1-\theta)^{M-k}}{\int_0^1 P(\theta) P(X=k|\theta) d\theta} \propto \theta^{a+k} (1-\theta)^{b+n-k}$$

$$(\theta|X) \sim \text{Beta}(a+k, b+n-k)$$

prior: $p(\theta) \sim \text{Beta}$, $X|p \sim \text{Binomial}$ posterior: $p(X|\theta) \sim \text{Beta}$

Beta is "conjugate prior" to Binomial

$$\begin{aligned} P(X[M+1] | X[1] \dots X[M]) &= \int_0^1 P(X[M+1], \theta | X[1] \dots X[M]) d\theta \\ &= \int_0^1 P(X[M+1] | \theta, X[1] \dots X[M]) P(\theta | X[1] \dots X[M]) d\theta \end{aligned}$$

$$P(X[M+1]=1 | X[1] \dots X[M]) = \int_0^1 \theta \cdot 1 \theta^{a+x+1} (1-\theta)^{b+N-x-1} d\theta = \int_0^1 \theta^{a+x} (1-\theta)^{b+N-x} d\theta = \frac{(a+x)! (b+N-x-1)!}{(a+b+N)!} \quad x = \sum_{i=1}^M \mathbb{1}\{X[i]=1\}$$

$$\begin{aligned} P(X[M+1]=0 | X[1] \dots X[M]) &= \int_0^1 (1-\theta) \cdot 1 \theta^{a+x+1} (1-\theta)^{b+N-x} d\theta \\ &= \int_0^1 (1-\theta)^{a+x+1} (1-\theta)^{b+N-x} d\theta = \frac{(a+x+1)! (b+N-x)!}{(a+b+N)!} \\ P(X[M+1]=1 | X[1] \dots X[M]) &= \frac{a+x}{a+b+N} \end{aligned}$$

intuitively, a, b corresponds to the num of imaginary heads and tails we "have seen" before the experiment $[X[1] \dots X[M]]$

• Prior and Posteriors

$$D = \{x_{i1}, \dots, x_{iM}\} \text{ iid } x_{ij} \sim p_\theta$$

$$P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)} = \frac{\underset{\substack{\text{likelihood} \\ \text{prior}}}{\frac{P(D|\theta) \cdot P(\theta)}{\int P(D|\theta') P(\theta') d\theta'}}}{\int P(D|\theta) P(\theta') d\theta'} \quad \text{marginal likelihood}$$

(a priori probability of seeing the particular dataset D given prior beliefs $P(\theta)$)

As Dirichlet distribution is specified by a set of hyperparameters $\alpha_1, \dots, \alpha_k$
 $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \Leftrightarrow p(\theta) \propto \prod \theta_i^{\alpha_i - 1}$

If $P(\theta)$ is Dirichlet ($\alpha_1, \dots, \alpha_k$), then $P(\theta|D)$ is Dirichlet($\alpha_1 + M_{11}, \dots, \alpha_k + M_{k1}$)

A family of priors $P(\theta; \alpha)$ is conjugate to a model $p(x|\theta)$ if for any $D = \{x_{i1}, \dots, x_{iM}\}$ of iid samples $x_{ij} \sim p(x|\theta)$, $\forall \alpha$, there are hyperparameter α' that describe the posterior

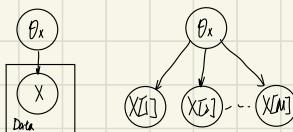
$$P(\theta; \alpha') \propto p(D|\theta) p(\theta; \alpha) \quad (\text{prior and posterior are in the same family})$$

eg. multinomial likelihood and Dirichlet prior

$$\theta \in \mathbb{R}^k \quad \sum \theta_i = 1 \quad p(\theta) \propto \prod \theta_i^{\alpha_i - 1}$$

$$p(x=k|\theta) = \theta_k$$

$$L(\theta; D) = p(D|\theta) = \prod \theta_i^{M_{ik}}$$



$$P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{\int P(D|\theta') P(\theta') d\theta'} \propto \prod \theta_i^{M_{ik}} \cdot \theta_i^{\alpha_i - 1}$$

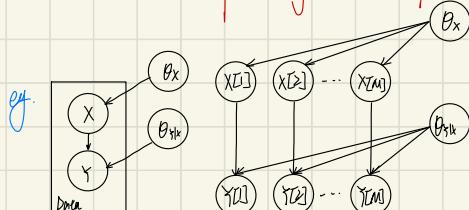
$$\theta | D \sim \text{Dirichlet}(\alpha_1 + M_{11}, \dots, \alpha_k + M_{k1})$$

Dirichlet prior are conjugate of multinomial model

Special case: Beta prior are conjugate of Binomial model

• Parameter Independence and Global Decomposition

the instances are independent given unknown parameters



$$x_{im}, y_{im} \perp x_{im'}, y_{im'} \mid \theta_x, \theta_y$$

$$\theta_x \perp \theta_y \mid (x_{i1}, y_{i1}), \dots, (x_{im}, y_{im})$$



knowing the value of one parameter tells nothing about another

let G be a bayesian network with parameters, $\theta = (\theta_{x|par}) \dots \theta_{y|par}$,

$P(\theta)$ satisfies global parameter independence if it has form $P(\theta) = \prod P(\theta_{x|par})$

$$\theta_x \rightarrow X[m] \rightarrow Y[m]$$

$(X[m], Y[m])$ block the V -structure

complex data α -separates the parameter for different CPTs

$$P(B_x, \theta_{B_x} | D) = P(B_x | D) \cdot P(\theta_{B_x} | D)$$

Once we can store each parameter separately, we can combine the results (like in MLE)

$$\begin{aligned} P(\theta | D) &= \frac{P(D|\theta) \cdot P(\theta)}{P(D)} \\ &= \frac{1}{P(D)} \prod_i^j L_i(\theta_{x_i|par}, D) \cdot \prod_i^j P(\theta_{x_i|par}) \quad \text{global parameter independence} \\ &= \frac{1}{P(D)} \prod_i^j [L_i(\theta_{x_i|par}, D) \cdot P(\theta_{x_i|par})] \end{aligned}$$

let G be a bayesian network over \mathcal{X} , $D = \{X[1], \dots, X[M]\}$ if $P(\theta)$ satisfies global parameter independence

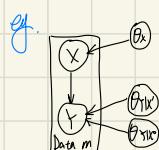


$$\begin{aligned} P(X[M+1], Y[M+1] | D) &= \int_{\theta} P(X[M+1], Y[M+1] | \theta, D) d\theta = \int_{\theta} P(X[M+1], Y[M+1] | \theta) \cdot P(\theta | D) d\theta \\ &= \int_{\theta} \int_{\theta_{x|x}} P(X[M+1] | \theta_x) \cdot P(Y[M+1] | \theta_{y|x}, X[M+1]) P(\theta_x | D) d\theta_x P(\theta_{y|x} | D) d\theta_{y|x} \\ &= \int_{\theta} P(X[M+1] | \theta_x) P(\theta_x | D) d\theta_x \cdot \int_{\theta_{y|x}} P(Y[M+1] | \theta_{y|x}, X[M+1]) P(\theta_{y|x} | D) d\theta_{y|x} \end{aligned}$$

if the prior for params of different CPTs are independent

$$P(X[M+1] \dots X[n+1] | D) = \prod_i^j \int_{\theta} P(X[i+1] | P(x)[i+1], \theta_{x|par}) \cdot P(\theta_{x|par} | D) d\theta_{x|par}$$

• Local Decomposition



$$P(Y[m] = y | X[m], \theta_{x|m}, \theta_{y|x}) = \begin{cases} \theta_{y|x} & X[m] = x \\ 0 & X[m] \neq x \end{cases}$$

θ_{yx} and θ_{yz} are dependent given the y (not given x)



θ_{yx} and θ_{yz} are independent given the data



$$P(\theta_{yx} | D) = P(\theta_{yx} | D) \cdot P(\theta_{yz} | D)$$

$$\begin{aligned}
 P(\theta, D) &= P(\theta) \cdot P(D|\theta) \\
 &= P(\theta_x) \cdot P(\theta_{\text{fix}}) \cdot \prod_{m=1}^M P(X_m|x_m) | P(Y_m|X_m; \theta_{\text{fix}}) \\
 &= P(\theta_x) \cdot L_x(\theta_x; D) \cdot \prod_{m=1, m \neq x}^M P(Y_m|X_m; \theta_{\text{fix}}) \\
 &\quad \cdot P(\theta_{\text{fix}}) \cdot \prod_{m=1, m \neq x}^M P(Y_m|X_m; \theta_{\text{fix}})
 \end{aligned}$$

If θ is Dirichlet($\alpha_{\text{fix}}, \dots, \alpha_x$) $P(\theta) = \prod \theta_i^{\alpha_i-1}$ $P(\theta_{\text{fix}}) = \prod \theta_i^{\alpha_{\text{fix}}(k)-1}$

$$\begin{aligned}
 P(\theta_{\text{fix}}|D) &= \int \int P(\theta_{\text{fix}}) \cdot P(D|\theta) \cdot P(\theta_{\text{fix}}|D) d\theta_x d\theta_{\text{fix}} \\
 &= \int \int P(\theta|D) d\theta_x d\theta_{\text{fix}} \\
 &\propto P(\theta_{\text{fix}}) \prod_{m=1, m \neq x}^M P(Y_m|X_m; \theta_{\text{fix}}) \\
 &= \left(\prod \theta_i^{\alpha_{\text{fix}}(k)-1} \right) \cdot \theta_x^{\alpha_x(M-1)+y} \cdot \theta_{\text{fix}}^{M-1-y} \\
 &= \text{Dirichlet}(\alpha_{\text{fix}} + M-1+y, \alpha_x + M-1-y)
 \end{aligned}$$

Let X be a variable with parents U , $P(\theta_{\text{fix}})$ satisfies local parameter independence if
 $P(\theta_{\text{fix}}) = \prod_{u \in \text{parents}(X)} P(\theta_{x|u})$

If the prior $P(\theta)$ satisfies global and local parameter independence, then
 $P(\theta|D) = \prod_{u \in \text{parents}(X)} P(\theta_{x|u}|D)$

$$\begin{aligned}
 P(\theta|D) &= \prod_{u \in \text{parents}(X)} P(\theta_{x|u}|D) \quad // \text{global parameter independence} \\
 &= \prod_{u \in \text{parents}(X)} P(\theta_{x|u}|D) \quad // \text{local parameter independence}
 \end{aligned}$$

If $P(\theta_{\text{fix}})$ is a Dirichlet distribution $P(\theta_{\text{fix}}) \propto \prod_{k=1}^K \theta_{x_k(k)}^{\alpha_{\text{fix}}(k)-1}$
then $P(\theta_{\text{fix}}|D)$ is a Dirichlet distribution $P(\theta_{\text{fix}}|D) \propto \prod_{k=1}^K \theta_{x_k(k)}^{\alpha_{\text{fix}}(k) + M-1+y_k}$

If U is a parent of X $\theta_{\text{fix}} \sim \text{Dirichlet}(\alpha_{\text{fix}}_1, \dots, \alpha_{\text{fix}}_M)$, D is complete observed over \mathcal{X}

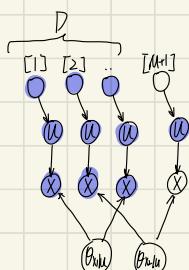
$$P(X_m=x_i | U_m=u, D) = \frac{\alpha_{x_i(u)} + M-1+y_i}{\sum_j (\alpha_{x_j(u)} + M-1+y_j)}$$

$$P(X_m=x_i | U_m=u, D) = \int_{\theta_{\text{fix}} \in \Theta} P(X_m=x_i, \theta_{\text{fix}} | U_m=u, D) d\theta_{\text{fix}}$$

$$= \int_{\theta_{\text{fix}} \in \Theta} P(X_m=x_i | \theta_{\text{fix}}, U_m=u, D) P(\theta_{\text{fix}} | U_m=u, D) d\theta_{\text{fix}}$$

$(X_m \perp D | \theta_{\text{fix}})$
 θ_{fix} blocks the tail

$$= \int_{\theta_{\text{fix}}} P(\theta_{\text{fix}}) P(\theta_{\text{fix}} | D) d\theta_{\text{fix}}$$



$$\begin{aligned}
 &= \frac{P(D)}{P(D)} \int_{\Omega_{X|U}} \prod_{i=1}^n \theta_{x_i u_i}^{k_{x_i u_i}} d\theta_{X|U} \stackrel{\text{defn of } P(D|U)}{=} P(D|U) \rightarrow \text{chain rule} \\
 &= \frac{P(D)}{P(D)} \int_{\Omega_{X|U}} \theta_{x_1 u_1}^{k_{x_1 u_1}} \cdots \theta_{x_n u_n}^{k_{x_n u_n}} d\theta_{X|U} \stackrel{\text{defn of } P(D|U)}{=} P(D|U) \rightarrow \text{chain rule} \\
 &= \frac{P(D)}{P(D)} \cdot \frac{\Gamma(k_{x_1 u_1} + M[x_1, u_1]) \cdots \Gamma(k_{x_n u_n} + M[x_n, u_n])}{\Gamma(k_{x_1 u_1} + M[x_1, u_1] + \sum_{j \neq i} k_{x_j u_j} + M[x_j, u_j])} \\
 &= \frac{1}{\sum_{j \neq i} k_{x_j u_j} + M[x_j, u_j]} \frac{\prod_{j \neq i} (k_{x_j u_j} + M[x_j, u_j])}{\prod_{j \neq i} k_{x_j u_j} + M[x_j, u_j]}
 \end{aligned}$$

// Integrating over a simplex

$$\text{let } S_n = \{(\theta_1, \dots, \theta_n) \mid \theta_i \geq 0, \sum \theta_i = 1\}$$

$$\int_{S_n} \theta_1^{k_1} \cdots \theta_n^{k_n} d\theta$$

$$I(t) = \int_{S_n} \theta_1^{k_1} \cdots \theta_n^{k_n} \int_{t \leq \theta_1 \cdots \theta_n} f(t - \theta_1 \cdots \theta_n) d\theta$$

Laplace transform: $I(s) = \int_0^{+\infty} I(t) e^{-st} dt$

$$= \int_0^{+\infty} \int_0^{t \wedge 1} \theta_1^{k_1} \cdots \theta_n^{k_n} f(t - \theta_1 \cdots \theta_n) d\theta e^{-st} dt$$

$$= \int_{(0, \infty)^n} \theta_1^{k_1} \cdots \theta_n^{k_n} \int_0^{+\infty} f(t - \theta_1 \cdots \theta_n) e^{-st} dt d\theta$$

$$= \int_{(0, \infty)^n} \theta_1^{k_1} \cdots \theta_n^{k_n} \exp(-s(\theta_1 + \cdots + \theta_n)) d\theta$$

$$= \prod_{i=1}^n \int_0^{+\infty} \theta_i^{k_i} \exp(-sk_i) d\theta$$

$$= \prod_{i=1}^n \frac{k_i!}{S^{\sum k_i + n}}$$

$$= \prod_{i=1}^n \frac{k_i!}{S^{\sum k_i + n}}$$

$$= \frac{(\sum k_i + n)!}{(\sum k_i + n)!} \cdot \frac{n!}{(\sum k_i + n)!}$$

$$= \frac{1}{(\sum k_i + n)!} \cdot \int_0^{+\infty} (t^{\sum k_i + n}) (s)$$

$$\int_{\Omega_{X|U}} \theta_1^{k_1} \cdots \theta_n^{k_n} d\theta = I(1) = \frac{\frac{n!}{\sum k_i + n} k_i!}{(\sum k_i + n)!} = \frac{\frac{n!}{\sum k_i + n} \Gamma(k_i + 1)}{\Gamma(\sum k_i + n)}$$

$$\int_{\Omega_{X|U}} \theta_1^{k_1} \cdots \theta_n^{k_n} d\theta = \frac{\frac{n!}{\sum k_i} \Gamma(k_i)}{\Gamma(\sum k_i)}$$

Priors for Bayesian Networks

In a Bayesian network \mathcal{B} over \mathcal{X} , for each var $i \in \mathcal{B}$, X_i has a set of multinomial distribution $\Omega_{X_i|U}$ w.r.t $U = \text{pa}(i)$

Can set the prior of $\theta_{X_i|U}$ as a Dirichlet prior

$$\theta_{X_i|U} \sim \text{Dirichlet}(\#[\mathcal{X}_i|U], \dots, \#[\mathcal{X}_i|U])$$

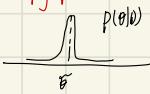
MAP Estimation

the parameters that maximize the posterior probability

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(\theta|D)$$

When we have a large amount of data, the posterior is often sharply peaked around $\hat{\theta}$

$$P(X_{\text{infit}}|D) = \int p(X_{\text{infit}}|\theta) p(\theta|D) d\theta \approx P(X_{\text{infit}}|\hat{\theta})$$



$$\underset{\theta}{\operatorname{argmax}} \log p(\theta|D) = \underset{\theta}{\operatorname{argmax}} \frac{\log p(\theta) p(D|\theta)}{p(D)}$$

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{\log p(\theta)}_{\text{use the prior to provide regularization over log-likelihood.}} + \underbrace{\log p(D|\theta)}_{\text{constant}}$$

$$= \underset{\theta}{\operatorname{argmax}} \log p(\theta) + \underbrace{\log p(D|\theta)}_{\text{when } D \text{ is a big dataset, } \log p(D) \text{ is negligible}}$$

e.g. Beta prior and Bernoulli model

$$p(\theta; \alpha, \beta) = C \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$p(x) = (1-\theta) \exp \left(\sum_i x_i \beta \right), \ln \frac{\theta}{1-\theta} >$$

$$\text{let } \eta = \ln \frac{\theta}{1-\theta} \quad \theta = \frac{1}{1+e^\eta} \quad 1-\theta = \frac{e^\eta}{1+e^\eta}$$

$$\begin{aligned} p(y) &= \frac{d\eta}{d\theta} p(\theta|y) = \frac{d\theta}{dy} p(\theta|y) = -\frac{1}{(1+e^\eta)^2} \cdot e^\eta \cdot \frac{1}{1+e^\eta} \cdot \left(\frac{1}{1+e^\eta} \right)^{\alpha-1} \left(\frac{e^\eta}{1+e^\eta} \right)^{\beta-1} \\ &= \frac{1}{1+e^\eta} \cdot \frac{e^\eta}{1+e^\eta} \cdot c \left(\frac{1}{1+e^\eta} \right)^{\alpha-1} \left(\frac{1}{1+e^\eta} \right)^{\beta-1} \\ &= c \cdot \left(\frac{1}{1+e^\eta} \right)^{\alpha-1} \left(\frac{1}{1+e^\eta} \right)^{\beta-1} \end{aligned}$$

when $\alpha=\beta=1$, $p(\theta)=1$ for $\theta \in [0,1]$, uniform

$$\text{but } p(y) = \frac{1}{1+e^\eta} \cdot \frac{1}{1+e^\eta} = \frac{1}{1+e^\eta+e^\eta} = \frac{1}{2+e^\eta+e^\eta} \text{ is from from uniform}$$

both MLE and Bayesian estimation are "representation-independent". (by change of variable in the integral)

but MAP estimation is not

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(\theta) = \underset{\theta}{\operatorname{argmax}} (\alpha-1) \log \theta + (\beta-1) \log 1-\theta = \frac{\alpha-1}{\alpha+\beta-2}$$

$$\hat{\eta} = \underset{\eta}{\operatorname{argmax}} p(y) = \underset{\eta}{\operatorname{argmax}} \frac{d\eta}{d\theta} \quad b(\hat{\eta}) = \frac{d\eta}{d\theta} \cdot \frac{1}{\hat{\theta}} + \delta(\hat{\eta})$$

MAP estimation is more sensitive to choices in formalizing the likelihood and the prior than MLE or Bayesian inference