

chain rule

$$P(A_1B) = P(A|B)P(B) = P(B|A)P(A)$$

$$P(A_1, A_2 \dots, A_k) = P(A_1)P(A_2|A_1) \dots$$

$$\dots P(A_k|A_1 \dots A_{k-1})$$

if A and B are mutually independent

$$P(A|B) = P(A) \quad P(B|A) = P(B)$$

if A and B are conditionally independent given C

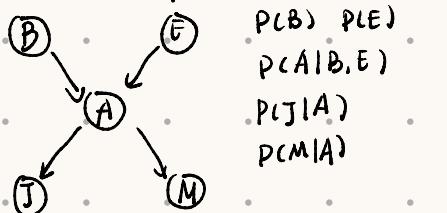
$$P(A, B|C) = P(A|C)P(B|C)$$

$$P(A|B, C) = P(A|C)$$

$$P(B|A, C) = P(B|C)$$

Inference By Enumeration

1. Collect all the rows consistent with the observed evidence variables.
2. Sum out (marginalize) all the hidden variable.
3. Normalize the table so that it is a probability distribution.



$$P(X_1, X_2 \dots, X_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

$$P(C, B, E, A, J, M) = P(C|B)P(B|E)P(A|B, E)$$

$$\dots P(J|A)P(M|A)$$

$$P(\text{Var}|\text{Parent}(\text{Var}), \text{Ancestors}(\text{Var}))$$

$$= P(\text{Var}|\text{Parents}(\text{Var}))$$

Causal Chains

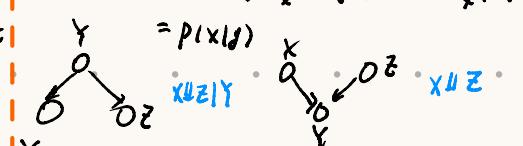
$$O \rightarrow O \rightarrow O \quad p(x, y, z)$$

$$x \quad y \quad z \quad = P(z|y)P(y|x)P(x)$$

(no observation)

$$P(X|Z, Y) = \frac{P(X, Z, Y)}{P(Z, Y)} = \frac{P(Z|Y)P(Y|X)P(X)}{\sum Z P(X, Y, Z)}$$

$$= \frac{P(Z|Y)P(Y|X)P(X)}{P(Z|Y) \sum Y P(Y|X)P(X)} = \frac{P(Y|X)P(X)}{\sum Y P(Y|X)P(X)}$$



D-separation Algorithm

1. Shade all observed values $\{Z_1 \dots Z_k\}$ in the graph
2. Enumerate all undirected paths from X to Y .
3. For each path:
 - a. Decompose the path into triples
 - b. if all triples are active, this path is active and d-connects X to Y .
4. If no path d-separates X and Y , then $X \perp\!\!\!\perp Y | Z_1 \dots Z_k$.

Variable Elimination

$$\text{function ELIMINATION-ASK}(X, e, bn) \text{ returns a distribution over } X$$

inputs: X , the query variable
e, observed values for variables E
bn, a Bayesian network specifying joint distribution $P(X_1, \dots, X_n)$

```

factors ← []
for each var in ORDER(bn.VARS) do
  factors ← [MAKE-FACTOR(var, e)] + factors
  if var is a hidden variable then factors ← SUM-OUT(var, factors)
return NORMALIZE(POINTWISE-PRODUCT(factors))

```

$$\text{by enumeration: } \sum_s \sum_t P(T)p(s|T)p(c|T)p(t|c,s)$$

$$\text{by variable elimination: } 2P(T) \sum_c P(c|T) \sum_s P(s|T)p(t|c,s)$$

prior sampling

simply sampling from given CPT.
disadvantage: require large data samples in order to perform unlikely scenarios. (like $P=0.01\%$).

Rejection Sampling

early reject any sample inconsistent with our evidence.
disadvantage: throw away most of our sample

Likelihood Weighting

manually set all variables equal to the evidence in our query, ensures we never generate a bad sample.

iterate through each variable in the Bayes net as we do for normal sampling, sampling a value if the variable is not an evidence variable, or changing the weight for sample if it is an evidence.

```

function LIKELIHOOD-WEIGHTING(X, e, bn, N) returns an estimate of  $P(X|e)$ 
inputs: X, the query variable
e, observed values for variables E
bn, a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
N, the total number of samples to be generated
local variables: W, a vector of weighted counts for each value of X, initially zero
for j = 1 to N do
  x, w ← WEIGHTED-SAMPLE(bn, e)
  W[x] ← W[x] + w where x is the value of X in x
return NORMALIZE(W)

```

```

function WEIGHTED-SAMPLE(bn, e) returns an event and a weight
w ← 1; x ← an event with n elements initialized from e
foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
  if  $X_i$  is an evidence variable with value  $x_i$  in e
    then  $w \leftarrow w \cdot P(X_i = x_i | \text{parents}(X_i))$ 
    else  $x_i \leftarrow$  a random sample from  $P(X_i | \text{parents}(X_i))$ 
return x, w

```

eg. calc $P(T) + U(e, t)$:

- $w_j = 1 \cdot c = \text{true}, e = \text{true}$
- sampling t_j from $P(T)$
- for C: $w_j = w_j \cdot P(c + t_j)$
- for S: sampling s_j from $P(s|t_j)$
- for E: $w_j = w_j \cdot P(e|t_j, c, s_j)$

Gibbs Sampling

first sampling all variables to some random value, then repeatedly pick one variable at one time, clear its value, and resample it.

```

function GIBBS-ASK(X, e, bn, N) returns an estimate of  $P(X|e)$ 
local variables: N, a vector of counts for each value of X, initially zero
Z, the non-evidence variables in bn
x, the current state of the network, initially copied from e

```

initialize x with random values for the variables in Z

```

for j = 1 to N do
  for each  $Z_i$  in Z do
    set the value of  $Z_i$  in x by sampling from  $P(Z_i | \text{mb}(Z_i))$ 
     $N[x] \leftarrow N[x] + 1$  where x is the value of X in x
return NORMALIZE(N)

```

Axioms of Utility

Orderability: $(A > B) \vee (B > A) \vee (A \sim B)$

Transitivity: $(A > B) \wedge (B > C) \Rightarrow (A > C)$

Continuity: $A > B > C \Rightarrow \exists p \text{ s.t. } P(A; U(p), C) \sim B$

Substitutability: $A \sim B \Rightarrow P(A; U(p), C) \sim P(B; U(p), C)$

Monotonicity: $A > B \Rightarrow (p \geq q) \Leftrightarrow$

$$[P(A; U(p), B) \geq P(A; U(q), B)]$$

$$U(P_1, S_1; \dots; P_n, S_n) = \sum_i P_i U(S_i)$$

Decision Networks

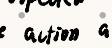
Umbrella



Weather



forecast



Maximum Expected Utility for possible action a and evidence e.

$$\text{EU}(a|e) = \sum_{x_1 \dots x_n} P(x_1 \dots x_n|e) U(a, x_1 \dots x_n)$$

General Formula

maximum expected utility

$$\text{MEU}(a|e) = \max_a \sum_s P(s|e) U(s, a)$$

if observed new evidence e' before acting

$$\text{MEU}(a, e') = \max_a \sum_s P(s|e, e') U(s, a)$$

if we don't know what the new evidence will be like, regard e' as random variable E'

$$\text{MEU}(e, E') = \sum_{e'} P(e'|e) \text{MEU}(e, e')$$

VPI is the increment of MEU after a new evidence observed

$$\text{VPI}(E'|e) = \text{MEU}(e, E') - \text{MEU}(e)$$

Property of VPI

Nonnegativity: $\forall E' \text{ s.t. } \text{VPI}(E'|e) \geq 0$

Nonadditivity: $\text{VPI}(E'_j, E'_k|e) \neq$

$$\text{VPI}(E'_j|e) + \text{VPI}(E'_k|e) \text{ in general}$$

Order-Independence: $\text{VPI}(E_j, E_k|e)$

$$= \text{VPI}(E_j|e) + \text{VPI}(E_k|e, E_j)$$

$$= \text{VPI}(E_k|e) + \text{VPI}(E_j|e, E_k)$$

Markov Model

$$W_0 \rightarrow W_1 \rightarrow \dots \rightarrow W_n$$

$$\vdash P(W_0, W_1, \dots, W_n)$$

Initial Distribution Transition Model

Markov Property

W_i only depends on W_{i-1}, independent of others

$P(W_0, W_1, W_2) = P(W_0)P(W_1|W_0)P(W_2|W_1, W_0)$
 $= P(W_0)P(W_1|W_0)P(W_2|W_1)$
 $\vdots P(W_0, W_1, \dots, W_n) = P(W_0) \prod_{i=0}^{n-1} P(W_{i+1}|W_i)$
 if the transition model is stationary, $(P(W_{i+1}|W_i))$ are identical). Then the markov model can be represented by 2 tables: $P(W_0)$, $P(W_{i+1}|W_i)$

The Mini-Forward Algorithm
 $P(W_{i+1}) = \sum_{w_i} P(w_i, W_{i+1})$
 $= \sum_{w_i} P(W_{i+1}|w_i) P(w_i)$
 stationary Distribution
 $\begin{cases} x+y=1 \\ 0.6x+0.1y=x \\ 0.4x+0.9y=y \end{cases} \Rightarrow \begin{cases} x=0.2 \\ y=0.8 \end{cases}$

Hidden Markov Models
 $W_0 \rightarrow W_1 \rightarrow W_2 \rightarrow W_3 \dots$
 $\downarrow \quad \downarrow \quad \downarrow$
 $F_1 \quad F_2 \quad F_3$
 $F_i \perp W_0 \mid W_1$
 $W_i \perp \{W_0, \dots, W_{i-1}, F_1 \dots F_{i-1}\} \mid W_{i+1}$
 $F_i \perp \{W_0, \dots, W_{i-1}, F_1 \dots F_{i-1}\} \mid W_i$
 assumption: $P(W_{i+1}|w_i)$: $P(F_i|w_i)$ are stationary.

Belief Distribution
 $B(W_i) = P(W_i|f_1 \dots f_i)$
 $B'(W_i) = P(W_i|f_1 \dots f_{i-1})$

Forward Algorithm
 $B'(W_{i+1}) = P(W_{i+1}|f_1 \dots i) = \sum_{w_i} P(W_{i+1}, w_i | f_1 \dots i)$
 $= \sum_{w_i} P(W_{i+1}|w_i, f_{i+1}) P(w_i | f_1 \dots i)$
 $= \sum_{w_i} P(W_{i+1}|w_i) P(w_i | f_1 \dots i)$
 $= \sum_{w_i} P(W_{i+1}|w_i) B(w_i)$
 $B(W_{i+1}) = P(W_{i+1}|f_1 \dots f_{i+1}) = \frac{P(W_{i+1}, f_{i+1}|f_1 \dots i)}{P(f_{i+1}|f_1 \dots i)}$
 $\propto P(W_{i+1}, f_{i+1}|f_1 \dots i)$
 $\propto P(f_{i+1}|W_{i+1}, f_1 \dots i) P(W_{i+1}|f_1 \dots i)$
 $\propto P(f_{i+1}|W_{i+1}) B'(W_{i+1})$
 $\propto P(f_{i+1}|W_{i+1}) \sum_{w_i} P(W_{i+1}|w_i) B(w_i)$

Time Elapse Update: $B(W_i) \rightarrow B'(W_{i+1})$
Observation Update: $f_{i+1}, B'(W_{i+1}) \rightarrow B(W_{i+1})$
Particle Filter
 gives an approximation of $P(X_N | e_{1:N})$
 particle number $n \ll d$ (d possible states)
 but still enough to generate valid approximation
 ① particles ② passage of time
 $X' = \text{sample}(P(X'|x))$

x_1	x_2	x_3
$1, 2, 2$	0.6	
$1, 2, 3$	0.1	
$2, 2, 3$	0.1	

③ observe
 $w(x) = P(e|x)$
 $B'(X) \propto P(e|X) B(X)$

1. 初始粒子
 粒子数 $n = 10$, 初始状态分布:
 $[15, 12, 12, 10, 18, 14, 12, 11, 11, 10]$
 2. 时间推移更新
 转移模型:
 $P(14|15) = 0.1, P(15|15) = 0.8, P(16|15) = 0.1$
 随机数采样后更新粒子:
 $[15, 13, 13, 11, 17, 15, 13, 12, 12, 10]$
 3. 观测更新
 假设观测 $F = 13$, 传感器模型:
 $P(F|13) = 0.8$, 其他状态 $P(F|T) = 0.02$
 计算权重, 归一化后重新采样:
 粒子更新为 $[13, 13, 13, 13, 13, 13, 13, 13, 15, 13]$
 4. 输出信念分布

T	10	11	12	13	14
$B(T)$	0.0	0.0	0.0	0.9	0.1

Naive Bayes
 $P(Y=\text{spam}|F_1=f_1 \dots f_n) = \frac{\text{count}(Y)}{N}$
 $P(Y=\text{ham}|F_1=f_1 \dots f_n) = \frac{\text{count}(Y)}{N}$
 $P(Y=\text{spam}|F) \propto P(Y=\text{spam}) \prod_{i=1}^n P(F_i|Y=\text{spam})$
 $P(F_1, F_2 \dots F_n|Y) = \prod_{i=1}^n P(F_i|Y)$
 $\text{prediction}(f_1 \dots f_n) = \underset{y}{\operatorname{argmax}} P(Y=y, F_1=f_1 \dots F_n=f_n)$
 $= \underset{y}{\operatorname{argmax}} P(Y=y) \prod_{i=1}^n P(F_i=f_i|Y=y)$

Linear Regression
 $h_w(x) = W_0 + W_1 x^1 + \dots + W_n x^n$
 $\text{Loss } L(h_w) = \frac{1}{2} \sum_{j=1}^n \|y_j - h_w(x_j)\|^2$
 $= \frac{1}{2} \|y - Xw\|_2^2$
 $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_1^1 & \dots & x_1^n \\ 1 & x_2^1 & \dots & x_2^n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^1 & \dots & x_n^n \end{bmatrix}, w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$
 $\nabla_w \frac{1}{2} \|y - Xw\|_2^2 = \nabla_w \frac{1}{2} (y - Xw)^T (y - Xw)$
 $= -X^T y + X^T X w = 0$
 $\therefore \hat{w} = (X^T X)^{-1} X^T y$
 $h_{\hat{w}}(x) = \hat{w}^T x$

Gradient Ascent
 ① Randomly initialize w
 ② while w is not converged:
 $w \leftarrow w + \alpha \nabla_w f(w)$

Gradient Descent
 ① Randomly initialize w
 ② while w is not converged:
 $w \leftarrow w - \alpha \nabla_w f(w)$

Least Squares Gradient Descent
 ① Randomly initialize w
 ② while w is not converged:
 $w \leftarrow w - \alpha (-x^T y + x^T X w)$

Logistic Regression
 $h_w(x) = \frac{1}{1 + e^{-w^T x}}$
 $P(y=+1|f(x); w) = \frac{1}{1 + e^{-w^T f(x)}}$
 $P(y=-1|f(x); w) = 1 - \frac{1}{1 + e^{-w^T f(x)}}$
 $\text{Loss } L(w) = \frac{1}{2} (y - h_w(x))^2$
 $\frac{\partial}{\partial w_i} \frac{1}{2} (y - h_w(x))^2 = -(y - h_w(x)) h_w(x)_i (1 - h_w(x)) x_i$
 $g(z) = \frac{1}{1 + e^{-z}}, g'(z) = g(z)(1 - g(z))$

