

Basic Probability Rules

$$P(A, B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$P(X_1, \dots, X_n) = P(X_1) \cdot P(X_2|X_1) P(X_3|X_1, X_2) \cdots P(X_n|X_{1:n-1})$$

$$P(X_{1:i-1}, X_{i+1:n}) = \sum_x P(X_{1:i-1}, X_i = x, X_{i+1:n})$$

$$P(X_i) = \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} P(X_1, \dots, x_{i-1}, X_i, x_{i+1}, \dots, x_n)$$

$$P(X|Y) = \frac{P(X) \cdot P(Y|X)}{P(Y)} = \frac{P(X) \cdot P(Y|X)}{\sum_x P(X=x) P(Y|X=x)}$$

$$P(X, Y|Z) = P(X|Y, Z) P(Y|Z)$$

$$P(X=x) = \sum_y P(X=x, Y=y) = \sum_y P(X=x|Y=y) P(Y=y)$$

$$X \perp Y|Z \text{ iff } P(X, Y|Z) = P(X|Z) \cdot P(Y|Z)$$

or if $P(Y|Z) > 0$ then $P(X|Z, Y) = P(X|Z)$

$$X \perp Y|Z \Rightarrow Y \perp X|Z \quad X \perp (Y, W)|Z \Rightarrow X \perp Y|Z$$

$$(X \perp Y|Z) \wedge (X \perp W|Y, Z) \Rightarrow X \perp Y, W|Z$$

$$(X \perp Y|W, Z) \wedge (X \perp W|Y, Z) \Rightarrow X \perp Y, W|Z \text{ if } P(Y) > 0$$

Bayesian Network

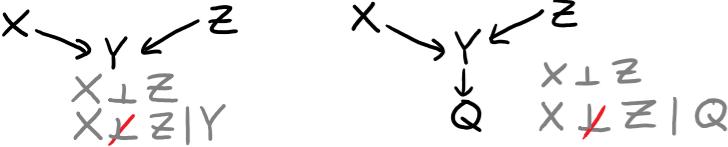
Consists of a structure: directed acyclic graph where each vertex is interpreted as a random variable.

Consists of a set of conditional probability distributions (CPTs) where all nodes S have $P(X_s | P_{\text{parents}}(X_s))$

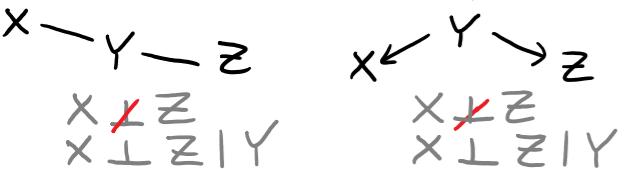
$$P(X_1, \dots, X_n) = \prod P(X_i | P_{\text{parents}}(X_i))$$

d-separation Rules

Colliders (only active trail if observed):



Non-colliders (active trail if unobserved):



$$d\text{-sep}(X; Y|Z) \Rightarrow X \perp Y|Z$$

Inference

tree-structured BN: variable elimination or belief propagation (gives exact marginals)
loopy networks: loopy belief propagation (fast, may not converge), variational inference (fast, converges, may be incorrect), Gibbs sampling (slow, converges to correct marginals)

Variable Elimination Order

For a polytree Bayesian Network (a directed acyclic graph which is a tree when dropping edge directions)
1. Drop edge directions

2. Choose root (e.g. for $P(A)$, pick A)
3. Orient all edges towards root
4. Eliminate from leaves

Always remains bounded: $O(n \cdot 2^{\max_i P_{\text{a.i}}})$

Example Variable Elimination

$$\begin{array}{c} \textcircled{A} \quad \textcircled{C} \\ \textcircled{B} \end{array} \quad \textcircled{D} \quad \textcircled{E} \quad P(A, B, C, D, E) = P(A) \cdot P(C) \cdot P(B|A, C) \cdot P(D|B) \cdot P(E|B)$$

Want to get $P(E) \Rightarrow$

$$P(E) = \sum_b \sum_a \sum_c \sum_d P(a, b, c, d, E)$$

$$= \sum_b P(E|b) \sum_a P(a) \sum_c P(c) P(b|a, c) \sum_d P(d|b)$$

$$= \sum_b P(E|b) \sum_a P(a) \sum_c P(c) P(b|a, c) = 1$$

$$= \sum_b P(E|b) \sum_a P(a) g_1(a, b) = g_2(b)$$

$$= \sum_b P(E|b) g_3(E)$$

Example Variable Elimination for MPE

Idea: Given some evidence J, M we want to find the most probable explanation for the other variables: $J \rightarrow A \rightarrow M$
 $\text{argmax}_{e,b,a} P(e, b, a | J, M) = \text{argmax}_{e,b,a} \cancel{P(e, b, a, J, M)} = \text{argmax}_{e,b,a} P(e) \cdot P(b) \cdot P(a|b, e) \cdot P(J|a) \cdot P(M|a)$
 $\max_a P(J|a) \cdot P(M|a) \max_e P(e) \max_b P(b) \cdot P(a|b)$
 $= \max_a P(J|a) \cdot P(M|a) \cdot \max_e P(e) g_B(a, e)$
 $= \max_a P(J|a) \cdot P(M|a) \cdot g_E(a)$
 $\rightarrow a^* = \text{argmax}_a P(J|a) \cdot P(M|a) \cdot g_E(a)$
 $\rightarrow e^* = \text{argmax}_e P(J|a) \cdot P(M|a) \cdot P(e) \cdot g_B(a, e)$
 $\rightarrow b^* = \text{argmax}_b P(J|a) \cdot P(M|a) \cdot P(e) \cdot P(b) \cdot P(a|b, e)$

Factor Graph

A bipartite graph consisting of variables (the random variables of the Bayesian Network) and factors (each factor is associated with a subset of variables s.t. all conditional probab. dist. of the BN are assigned to one of the factor nodes).



Sum-Product Message Passing

Initialize all messages as uniform distribution (=1)
Until converged do

Pick ordering on the factor graph edges
Update messages according to the ordering
Break if messages change at most ϵ

Message from node v to factor u:

$$M_{v \rightarrow u}(x_v) = \prod_{u' \in N(v) \setminus \{u\}} M_{u' \rightarrow v}(x_v)$$

Message from factor u to node v:

$$M_{u \rightarrow v}(x_v) = \sum_{x_u \sim x_v} f_u(x_u) \prod_{v' \in N(u) \setminus \{v\}} M_{v' \rightarrow u}(x_v)$$

$$P(X_v = x_v) \propto \prod_{u \in N(v)} M_{u \rightarrow v}(x_v)$$

$$P(X_u = x_u) \propto f_u(x_u) \prod_{v \in N(u)} M_{v \rightarrow u}(x_v)$$

Belief Propagation of Polytree BN

- Factor graph is a tree as well
- Choose one node as root, send messages from leaves to root, and from root to leaves
- Converges to correct marginals after two rounds (one such pass)

Max-Product Message Passing

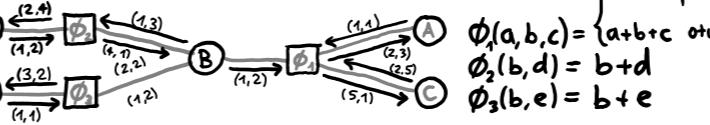
Only difference: $M_{u \rightarrow v}(x_v) = \max_{x_u \sim x_v} f_u(x_u) \prod_{v' \in N(u) \setminus \{v\}} M_{v' \rightarrow u}(x_v)$

$$P_{\max}(X_v = x_v) = \max_{x_u \sim x_v} P(x)$$

For tree BN: $P_{\max}(X_v = x_v) \propto \prod_{u \in N(v)} M_{u \rightarrow v}(x_v)$

Example Belief Propagation

$$P(A, B, C, D, E) \propto \phi_1(A, B, C) \phi_2(B, D) \phi_3(B, E)$$



From these messages, compute the approx. marginals of B.

From belief propagation we have

$$P(B) \propto \prod_{\phi_i \in \text{EN}(B)} M_{\phi_i \rightarrow B}^{(+)}(B) = M_{\phi_1 \rightarrow B}^{(+)}(B) M_{\phi_2 \rightarrow B}^{(+)}(B) M_{\phi_3 \rightarrow B}^{(+)}(B)$$

$$\text{for } B=0: \hat{P}(B=0) = 3 \cdot 4 \cdot 1 = 12 \quad \text{for } B=1: \hat{P}(B=1) = 2 \cdot 1 \cdot 2 = 4 \quad \{ P(B) = \frac{12+4}{2} = \frac{16}{2} = 8 \}$$

Compute the message $M_{B \rightarrow \phi_2}^{(+1)}$

$$M_{B \rightarrow \phi_2}^{(+1)} = \prod_{\phi_i \in \text{EN}(B) \setminus \{\phi_2\}} M_{\phi_i \rightarrow B}^{(+)}(B) = M_{\phi_1 \rightarrow B}^{(+)}(B) \cdot M_{\phi_3 \rightarrow B}^{(+)}(B)$$

$$\text{for } B=0: \mu_{B \rightarrow \phi_2}^{(+1)}(0) = 3 \cdot 1 = 3 \quad \text{for } B=1: \mu_{B \rightarrow \phi_2}^{(+1)}(1) = 2 \cdot 2 = 4 \quad \{ M_{B \rightarrow \phi_2}^{(+1)}(B) = (3, 4) \}$$

Compute the message $M_{\phi_1 \rightarrow B}^{(+1)}$

$$M_{\phi_1 \rightarrow B}^{(+1)}(B) = \sum_{V^* \in V_{\phi_1 \setminus \{B\}}} \phi_1(V_{\phi_1}) \prod_{V \in N(\phi_1) \setminus \{B\}} M_{V \rightarrow \phi_1}^{(+)}(V)$$

$$= \sum_a \phi_1(a, B, C) M_{A \rightarrow \phi_1}(a) M_C \rightarrow \phi_1(C)$$

$$\text{for } B=0: M_{\phi_1 \rightarrow B}^{(+1)}(0) = 1 \cdot 1 \cdot 5 + 2 \cdot 1 \cdot 5 = 15 \quad \text{for } B=1: M_{\phi_1 \rightarrow B}^{(+1)}(1) = 2 \cdot 1 \cdot 5 + 3 \cdot 1 \cdot 5 = 25 \quad \{ M_{\phi_1 \rightarrow B}^{(+1)}(B) = (15, 25) \}$$

Example Belief Propagation

Given a factor graph graph LR corresponding to a chain BN. For $0 \leq i \leq N$ we assume $X_i \in \{0, \dots, M\}$. Belief propagation computes $P(X_i = x_i) \forall i \in \{0, \dots, N\} \forall x_i \in \{0, \dots, M\}$.

Suppose we want to compute $P(X_i = x_i | X_0 = x_0, \dots, X_{i-1} = x_{i-1})$. How can this be done? $P(X_0 = x_0, \dots, X_i = x_i) = \prod_{j=0}^{i-1} P(X_j = x_j)$.

We propose the factor $f_i(x_i, x_{i+1}) = \text{if } x_i \leq x_{i+1} \text{ then } f_i(x_i, x_{i+1}) \text{ else } 0$. From belief propagation we know $\hat{P}(X_i = x_i) = \frac{1}{Z} \sum_{x_0, \dots, x_{i-1}} \prod_{k=0}^{i-1} f_k(x_k, x_{k+1})$.

Further $Z = \sum_{x_0, \dots, x_{i-1}} \prod_{k=0}^{i-1} f_k(x_k, x_{k+1})$. We evaluate $P(X_0 = x_0, \dots, X_i = x_i)$ in the original graph: $P(X_0 = x_0, \dots, X_i = x_i) = \frac{1}{Z} \sum_{x_0, \dots, x_{i-1}} \prod_{k=0}^{i-1} f_k(x_k, x_{k+1}) = \frac{Z}{Z} = 1$

$P(X_i = x_i | X_0 = x_0, \dots, X_{i-1} = x_{i-1}) = \frac{1}{Z} \sum_{x_0, \dots, x_{i-1}} \prod_{k=0}^{i-1} f_k(x_k, x_{k+1}) = \frac{1}{Z} \sum_{x_0, \dots, x_{i-1}} \prod_{k=0}^{i-1} \hat{f}_k(x_k, x_{k+1}) = \hat{P}(X_i = x_i)$

Monte Carlo Sampling (Forward Sampling)

1. Sort variables in topological order

2. In this order, sample $x_i \sim P(X_i | X_0 = x_0, \dots, X_{i-1} = x_{i-1})$

3. Marginals: $P(X_i = x_i) \approx \frac{1}{N} \sum_{i=1}^N [X_i = x_i] \chi^{(i)}$

Conditionals: $P(X_i = x_i | Y = y) \approx \frac{\text{count}(X_i = x_i, Y = y)}{\text{count}(Y = y)}$ (rejection sampling). If y rare \rightarrow large waste of samples

How many samples? Error decreases exponentially in N

Example Rejection Sampling

d-dimensional target $p(x) = N(x; \mu, \sigma_p^{2/d})$ and the proposal $q(x) = N(x; \mu, \sigma_q^{2/d})$. Optimal acceptance rate can be accomplished with $k = \sigma_q / \sigma_p$. With $d=1000$ and $\sigma_q = 1.01 \sigma_p \rightarrow k = 1/20000$ resulting in a huge waste in samples

Gibbs Sampling: Random Order

Start with initial assignment to all variables
Fix observed variables X_o to their observed value x_o

For $t=1$ to ∞ do

1. Pick a variable i at random (non-observed)

2. Set $v_i = \text{values of all } x \text{ except } x_i$

3. Update x_i by sampling from $P(x_i | v_i)$

Satisfies detailed balance equation

Gibbs Sampling: Practical Variant

Start with initial assignment $x^{(0)}$ for all var.
Fix the fixed values to be conditioned on

For $t=1$ to ∞ do
For each variable X_i (that is not fixed)
Sample value from $\sim P(x_i | v_i)$ other vars from the last iteration

No detailed balance but also has correct stationary distribution.

Computing Expectations via Gibbs Sampling

With GS get samples $X^{(1)}, \dots, X^{(T)}</math$

Ergodic: If every state can be reached in finite steps
↳ have a unique positive stationary distribution.

Ergodic property: An ergodic Markov chain has a unique, positive stationary distribution $\pi(x) > 0$ s.t. $\forall x \lim_{n \rightarrow \infty} P(X_n = x) = \pi(x)$ (π is independent of X_1)

Ergodic theorem: $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i) = \sum_{x \in D} \pi(x) f(x)$
where D is finite state space of x_i , f a function on D

Detailed balance: For $T(x,y) = P(y|x)$ a sufficient condition for ensuring π^\ast is stationary: $\pi(x)T(x,y) = \pi(y)T(y,x)$

Example Detailed Balance

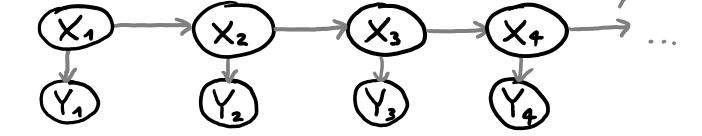
Assume that given a Markov chain with state space S^2 and transition matrix T , stationary distribution π .

Show that if for some $t: X_t$ distributed accord. to π and the chain satisfies the detailed balance equation $\pi(x)T(x,y) = \pi(y)T(y,x)$ for all $x, y \in S^2$, then $\forall k \geq 0$ and $x_0, \dots, x_k \in S^2$ it holds:
 $P(X_t=x_0, \dots, X_{t+k}=x_k) = P(X_t=x_k, \dots, X_{t+k}=x_0)$
 $P(X_t=x_0, \dots, X_{t+k}=x_k) = P(X_t=x_0)P(X_{t+1}=x_1|X_t=x_0) \dots$
 $\dots P(X_{t+k}=x_k|X_{t+k-1}=x_{k-1}) = \pi(x_0)T(x_0, x_1) \dots T(x_{k-1}, x_k)$
 $= T(x_1, x_0)\pi(x_1) \dots T(x_{k-1}, x_k) = \dots = T(x_1, x_0) \dots T(x_k, x_{k-1})T(x_k)$
 $\dots P(X_{t+k}=x_0|X_{t+k-1}=x_1) = P(X_t=x_k)P(X_{t+1}=x_{k-1}|X_t=x_k)$
 $\dots P(X_{t+k}=x_0|X_{t+k-1}=x_1) = P(X_t=x_k, \dots, X_{t+k}=x_0)$

Prediction in Markov Chains:

$$P(X_t|X_0=x) = \sum_{x'} P(X_t|X_{t-1}=x'|X_0=x) \\ = \sum_x P(X_t|X_{t-1}=x, X_0=x) \cdot P(X_{t-1}=x'|X_0=x) \\ = \sum_x P(X_t|X_{t-1}=x') \cdot P(X_{t-1}=x'|X_0=x)$$

Hidden Markov Model
 categorical
 X_1, \dots, X_T unobserved (hidden) variables
 Y_1, \dots, Y_T observations → can be categorical or arbitrary



Transition probabilities: $P(X_t|X_{t-1})$

Emission probabilities: $P(Y_t|X_t)$

Tasks: Filtering $P(X_t|Y_{1:T})$
 Smoothing $P(X_t|Y_{1:T})$ $1 \leq t < T$
 Prediction $P(X_{t+\tau}|Y_{1:t})$ $T \geq 1$
 Most probable explanation $\arg\max_{X_{1:T}} P(X_{1:T}|Y_{1:T})$

Example Hidden Markov Models

In each time step a coin is flipped, resulting in heads (h) or tails (t). After each flip, the coin used changes with probability $3/4$. The prior is $P(X_1=b) = 3/5$. $P(X_1=f) = 2/5$.

Derive $\lim_{i \rightarrow \infty} P(X_i=b)$: because i goes to infinity, prior distribution does not matter. Since probability $b \rightarrow f$ and $f \rightarrow b$ is the same, for $i \rightarrow \infty$ $P(X_i=b) = 1/2$

Formal way: $V^i = \begin{bmatrix} \text{prob of fair coin} \\ \text{prob of biased coin} \end{bmatrix} = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$

find transition matrix $T = \begin{bmatrix} P(f|f) & P(f|b) \\ P(b|f) & P(b|b) \end{bmatrix}$
 (this is a stochastic matrix, as rows sum to 1 ⇒ has eigenvalue 1) $= \begin{bmatrix} 1/4 & 3/4 \\ 3/4 & 1/4 \end{bmatrix}$

If we know the initial state V^1 , we can calculate: $V^t = T \cdot V^{t-1} = T \cdot V^1$

We look for $V^\infty = T \cdot V^\infty \sim 1 \cdot V^\infty = T \cdot V^\infty$
 ⇒ find eigenvector for eigenvalue 1:

$$\begin{bmatrix} 1/4 - 1 & 3/4 \\ 3/4 & 1/4 - 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \sim \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} \leftarrow \text{result}$$

Kalman Filter

Special case of HMM where X_t, Y_t are Gaussian and with Linear dynamics and observation model.

Motion model: $P(X_{t+1}|X_t)$ with $X_{t+1} = F X_t + \varepsilon_t$, $\varepsilon_t \sim N(0, \Sigma_x)$

Observation model: $P(Y_t|X_t)$ with $Y_t = H X_t + v_t$, $v_t \sim N(0, \Sigma_y)$

Filtering: $P(X_t|Y_{1:t}) \propto P(X_t|Y_{1:t-1}) \cdot P(Y_t|X_t)$

Prediction: $P(X_{t+1}|Y_{1:t}) = \int P(X_{t+1}|X_t) P(X_t|Y_{1:t}) dX_t$

↳ since Gaussian, can compute integral in closed form

Kalman Update: $\mu_{t+1} = F \mu_t + K_{t+1}(y_{t+1} - HF_{t+1})$ compute on the fly

$\Sigma_{t+1} = (I - K_{t+1})(F \Sigma_t F^T + \Sigma_x)$ compute offline

Kalman Gain: $K_{t+1} = (F \Sigma_t F^T + \Sigma_x)^{-1} H^T (H(F \Sigma_t F^T + \Sigma_x)^{-1} H^T + \Sigma_y)^{-1}$

$$\text{Multivar. Gaussians: } N(y; \Sigma, \mu) = ((2\pi)^{n/2} |\Sigma|)^{-1} \exp(-\frac{1}{2}(y-\mu)^T \Sigma^{-1} (y-\mu))$$

$$M_{AB} = M_A + \sum_{AB} \sum_{BB} (X_B - \mu_B) \quad \Sigma_{AB} = \sum_{AA} - \sum_{AB} \sum_{BB} \Sigma_{BA}$$

Bayesian Filtering (for HMMs)

We have $P(X_t|Y_{1:t})$ and want to compute $P(X_{t+1}|Y_{1:t})$

1. Conditioning: $P(X_{t+1}|Y_{1:t}) = P(X_t|Y_t, Y_{1:t-1})$

$$= \frac{1}{2} P(X_t|Y_{1:t-1}) P(Y_t|X_t, Y_{1:t-1}) = \frac{1}{2} P(X_t|Y_{1:t-1}) P(Y_t|X_t) O(D)$$

2. Prediction: $P(X_{t+1}|Y_{1:t}) = \sum_{x_t} P(X_{t+1}, x_t|Y_{1:t})$

$$= \sum_{x_t} P(X_{t+1}|X_t, Y_{1:t}) \cdot P(X_t|Y_{1:t}) = \sum_{x_t} P(X_{t+1}|X_t) \cdot P(X_t|Y_{1:t}) O(D)$$

→ recursively compute starting at $P(X_1)$

Particle Filtering

Approximate posterior $P(X_t|Y_{1:t})$ at each time with N i.i.d. particles (samples): $P(x_t|y_{1:t}) \approx \frac{1}{N} \sum_{i=1}^N \delta_{x_i, y_i}$

1. Prediction: $x'_i \sim P(X_{t+1}|x_i, y_{1:t})$ (sample according to before)

2. Conditioning: $w_i \propto P(y_{t+1}|x'_i)$ (compute for each class)

Resample N particles: $x_{i+1} \sim \frac{1}{N} \sum_{i=1}^N w_i \delta_{x_i}$

Why resampling? If we don't, then usually all weight concentrates on a single particle (mode collapse)

Example Particle Filtering / HMM

HMM with $X \in \{N, V, R\}$ and $Y \in \{C, H\}$. Given probabilities: $P(H|V) = 1$, $P(H|R) = 1$, $P(C|N) = 1$.

Also transition probabilities.

If CHC is observed, find max likely states.

$\arg\max P(X_{1:4}|CHHC)$ since we know that $P(C|N)$

$= 1$, it's the only possible scenario. We can simplify:

$\arg\max P(X_{2:3}|X_1=N, X_4=N, CHHC)$. There are four options to compute: RR, RV, VR, VV. → find max

You have 5 particles for N, 3 for V, 2 for R.

The new observation is H, compute the weights.

$w_i \propto P(H|X=x)$. For $x=N$: $w_{1:5} = \frac{1}{2} P(H|N) = 0$

For $x=V/R$: $w_{6:8} = \frac{1}{2} P(H|V) = \frac{1}{2} 1 = \frac{1}{2} P(H|R) = w_{3:10} = 0$

↳ compute $Z = \sum w_i = 5 \sim w_{1:5} = 0$, $w_{6:10} = 1/5$

What is the expected number of particles for N, V, R after resampling?

$$E[\sum_{i=1}^{10} w_i \cdot \mathbb{1}_{X_{i+1}=N}] = 0$$

$$E[\sum_{i=1}^{10} w_i \cdot \mathbb{1}_{X_{i+1}=R}] = \sum_{i=1}^{10} E[w_i \cdot \mathbb{1}_{X_{i+1}=R}] = 10 \cdot \frac{1}{5} \cdot 2 = 4$$

$$E[\sum_{i=1}^{10} w_i \cdot \mathbb{1}_{X_{i+1}=V}] = \sum_{i=1}^{10} E[w_i \cdot \mathbb{1}_{X_{i+1}=V}] = 10 \cdot \frac{1}{5} \cdot 3 = 6$$

Example Particle Filter

Given: -1D movement on states $x \in \mathbb{Z}$

- $x_{t+1} = x_t + \varepsilon_t$ where ε_t uniformly distributed with integer values in $[-3, 3]$

- measurements $y_t = (x_t + \eta_t)^2$

where η_t is dist: $P(\eta_t) = \begin{cases} 0.6 & \text{if } \eta_t=0 \\ 0.4 & \text{if } \eta_t=\pm 1 \\ 0 & \text{otherwise} \end{cases}$

You want: particle filter with 6 particles. robot at $t=0$: x_0 . Particles $x_i = 0 \forall i \in \{0, \dots, 5\}$

You get for ε_0 : $(-1, -1, 0, 1, 2, 3)$ (samples)
 ↳ update particles $x'_i = (-1, -1, 0, 1, 2, 3)$

You get measurement $y_1 = 1$. What are the particle weights?

From the measurement model we obtain: $\eta_1 = \pm \sqrt{y_1 - x'_1} \Rightarrow P(\eta_1 = \pm \sqrt{y_1 - x'_1} | x'_1) = P(\eta_1 = \pm \sqrt{y_1 - x'_1} | x'_1)$

The particle weights are computed as:

$$w_i = \frac{1}{2} P(y_1|X'_i) : P(y_1=1|X'_0=-1) = P(\eta_1=0) + P(\eta_1=2) \\ = 0.6 + 0.4 = 0.6 \Rightarrow \text{same for others}$$

$$(0.6, 0.6, 0.4, 0.6, 0.2, 0.0)$$

$$\text{Compute } Z = \sum_{i=0}^5 w_i = 0.6 + 0.6 + 0.4 + 0.6 + 0.2 + 0.0 = \frac{24}{10}$$

⇒ compute weights: $w_0 = \frac{6}{24}, \dots, w_5 = 0$

In reality need more samples to capture probability distributions accurately.

Partially Observable MDPs (POMDP)

Idea: interpret POMDP as an MDP with enlarged state space: $P(x_{t+1}|x_t, A_t)$, new states correspond to beliefs $P(X_t|y_{1:t})$ in the original MDP.

Can calculate optimal action using dynamic programming

Can approximate using policy gradients

Markov Decision Process (MDP)

States $X = \{1, \dots, n\}$, actions $A = \{1, \dots, m\}$, transition probabilities $P(x'|x, a)$, reward function $r(x, a, x')$.

Modes: Finite horizon (T steps) or discounted rewards γx

Value of policy: $V_\pi(x) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(x_t, \pi(x_t)) | X_0 = x]$

$$= \sum_x P(x'|x, \pi(x)) [r(x, \pi(x), x') + \gamma V_\pi(x')] = r(x, \pi(x)) + \gamma \sum_x P(x'|x, \pi(x)) V_\pi(x')$$

Greedy Policy: $\pi_g(x) = \arg\max_a \sum_x P(x'|x, a) (r(x, a, x') + \gamma V_\pi(x'))$

Bellman Equation: Policy is optimal \Leftrightarrow greedy w.r.t. its induced value function. $V(x) = \max_a [r(x, a) + \gamma \sum_x P(x'|x, a) V(x')]$

Policy Iteration: Start with arbitrary (e.g. random) policy Until converged do:

Compute value function $V^\pi(x)$

$$V^\pi(x) = r(x, \pi(x)) + \gamma \sum_{x'} P(x'|x, \pi(x)) V^\pi(x')$$

Compute greedy policy π_π w.r.t. V^π

$$\pi_\pi(x) = \arg\max_a r(x, a) + \gamma \sum_{x'} P(x'|x, a) V^\pi(x')$$

$$\pi \leftarrow \pi_\pi$$

Guaranteed to monotonically improve and converge in $O(n^2 m / (1-\gamma))$ iterations

Cost per iteration: $O(|S|^3 + |S| \cdot |A| \cdot DT)$ can have $\leq |S|$ max paths a action

Value Iteration

Initialize $V_0(x) = \max_a r(x, a)$

For $t=1$