

MDP

Markov Model: discrete variables are share the same infinite domain

$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n \rightarrow \dots$$

$$P(x_0) \quad P(x_t|x_{t-1})$$

Transition Model: $P(x_t|x_{t-1}) \rightarrow$ how state evolves

Stationarity assumption: same transition probabilities at all time steps

Joint distribution: $P(x_0, \dots, x_n) = \prod_{t=1}^n P(x_t|x_{t-1})$

Markov assumption: x_{t+1} are independent of x_0, \dots, x_{t-1} given x_t .

\Rightarrow Each time step only depends on previous

k -th order markov model: depends on k earlier steps.

Markov Random walk

State: location on the unbounded integer line

Initial probability: start at 0

Transition model: $P(X_t=k|X_{t-1}=k \pm 1) = 0.5$

Applications: particle motion in crystal
stock price, gambling, ...

Eg. whether state $(\text{rain}, \text{sun})$

$$P(X_0) = \begin{cases} \text{sun: } 0.5 \\ \text{rain: } 0.5 \end{cases} \quad P(X_t|X_{t-1}) = \begin{matrix} \text{sun} & \text{sun} \\ \text{sun} & \begin{matrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{matrix} \\ \text{rain} & \begin{matrix} 0.3 & 0.7 \\ 0.7 & 0.3 \end{matrix} \end{matrix}$$

predict: time 0: 0.5, 0.5

$$\begin{aligned} P(X_1) &= \sum_{X_0} P(X_0, X_1) = \sum_{X_0} P(X_0) P(X_1|X_0) \\ &= 0.5 \cdot 0.9 + 0.1 + 0.5 \cdot 0.3 + 0.7 \\ &= 0.6, 0.4 \end{aligned}$$

$$\begin{aligned} P(X_2) &= \sum_{X_1} P(X_1, X_2) = \sum_{X_1} P(X_1) P(X_2|X_1) \\ &= 0.6 \cdot 0.9 + 0.1 + 0.4 \cdot 0.3 + 0.7 \\ &= 0.66, 0.34 \end{aligned}$$

State at time step t:

$$P(X_t) = \sum_{X_{t-1}} P(X_{t-1}) P(X_t|X_{t-1})$$

Iterate from $t=0$

Stationary distribution

The distribution end up with is called the stationary distribution P_{st} of the chain.

$$\Rightarrow P_{\text{st}}(X) = P_{\text{st}}(X) = \sum_i P(X|t) P_{\text{st}}$$

Eg. Rain and sun

$$P_{\text{st}}(\text{sun}) = 0.9 P_{\text{st}}(\text{sun}) + 0.3 P_{\text{st}}(\text{rain})$$

$$P_{\text{st}}(\text{rain}) = 0.1 P_{\text{st}}(\text{sun}) + 0.7 P_{\text{st}}(\text{rain})$$

$$\text{Also, } P_{\text{st}}(\text{rain}) + P_{\text{st}}(\text{sun}) = 1$$

$$\Rightarrow P_{\text{st}}(\text{sun}) = \frac{2}{3}, P_{\text{st}}(\text{rain}) = \frac{1}{3}$$

$$\text{or by matrix: } P_{\text{st}} = P_{\text{st}} = T P_{\text{st}} \quad \begin{bmatrix} 0.9 & 0.3 \\ 0.1 & 0.7 \end{bmatrix} = \begin{bmatrix} P_{\text{st}} \\ P_{\text{st}} \end{bmatrix} \Rightarrow P_{\text{st}} = 0.75$$

HMM Hidden Markov Model

Joint distribution for MM: $P(x_0, \dots, x_n) = \prod_{t=1}^n P(x_t|x_{t-1})$

Joint distribution for HMM:

$$P(x_0, x_1, \dots, x_n, E_1, \dots, E_n) = P(x_0) \prod_{t=1}^n P(x_t|x_{t-1}) P(E_t|x_t)$$

E_t : observed variables. X_t : hidden variables.

Filtering: $P(X_t|e_{1:t})$ posterior distribution over the most recent state given all evidence.

Prediction: $P(X_{t+k}|e_{1:t})$ for $k > 0$
posterior distribution over a future state given all evidence

Smoothing: $P(X_k|e_{1:t})$ for $0 < k < t$
posterior distribution over a past state given all evidence

Most likely explanation: $\arg \max_{X_{1:t}} P(X_{1:t}|e_{1:t})$

Assumption: 1. $P(X_{1:t}|e_{1:t}) = P(X_{1:t}|x_{1:t})$
2. $P(e_{1:t}|x_{1:t}, e_{1:t}) = P(e_{1:t}|x_{1:t})$

Filtering: recursive filtering

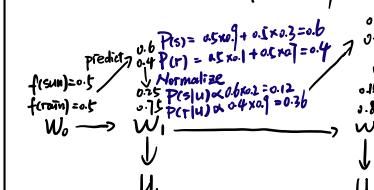
$$\begin{aligned} P(X_{1:t}|e_{1:t}) &= g(x_{1:t}, P(X_{1:t}|e_{1:t})) = P(X_{1:t}|e_{1:t}, x_{1:t}) \\ &= \underset{\text{Normalize}}{\cancel{P(e_{1:t}|X_{1:t})}} \sum_{X_{1:t}} P(X_{1:t}|e_{1:t}) \underset{\text{predict}}{\cancel{P(X_{1:t}|x_{1:t})}} \end{aligned}$$

$$\Rightarrow f_{1:T} = \text{Forward}(f_{1:t}, e_{t+1}) = \sum_{X_{t+1}} P(X_{t+1}|e_{1:t}), f_{1:0} = P(x_0)$$

Time complexity: $O(M^2)$ where M is the number of states.

Example: Weather & Umbrella

	$P(u_{t+1} w_t)$	w_t	$P(u_t w_t)$
w_2	sum rain	sum	sum true false
sum	0.9 0.1	0.2 0.8	
rain	0.3 0.7	0.1 0.9	



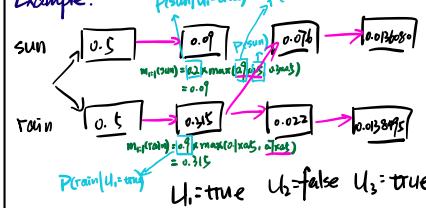
Most likely explanation

Viterbi Algorithm: $\arg \max_{X_{1:t}} P(X_{1:t}|e_{1:t})$

$$m_{t+1} = P(e_{t+1}|X_{t+1}) \max_{x_{t+1}} P(X_{t+1}|x_t) m_{1:t}[x_t]$$

$$m_{1:t}[x_t] = \max_{x_{1:t}} P(X_{1:t}|e_{1:t})$$

Example.



DBN Dynamic Bayes Network

Repeat a fixed Bayes net structure at each time.

Variables from time t can condition on those from t-1

Each HMM is Cartesian product of DBN state variable

\Rightarrow e.g. 2 binary state variables \Rightarrow one state var with 2^2 domain

Advantage DBN vs. HMM:

sparse dependencies \Rightarrow exponentially fewer parameters.

E.g. 2 binary state variables, 2 parent each

$$\Rightarrow \text{DBN has } 2 \times \frac{2^{14}}{2} = 160 \text{ parameters}$$

HMM has $2^2 \times 2^2 \approx 10^4$ parameters.

Partial Filtering

1. Sample N particles by distribution $P(x_0)$

2. transition by $P(X_{t+1}|X_t)$

3. give weight $w_t = P(e_{1:t}|x_{1:t})$ for each particle.

4. resample from weight and particle position.

(将权重之和归一化后，按比例抽样)

Markov Decision Process

Define 1. state $s \in S$ 2. action $a \in A$

3. transition function $T(s, a, s')$

4. reward function $R(s, a, s')$

5. start state s_0 6. Maybe terminal state.

Discounting: γ . values of rewards decay exponentially

Infinite Utilities:

finite horizon: 1. Terminate episodes after a fixed T steps.

2. Given nonstationary policies.

Discounting: use $0 < \gamma < 1$. $U(\Gamma_0, \dots, \Gamma_T) = \sum_{t=0}^T \gamma^t r_t \leq \frac{R_{\max}}{1-\gamma}$

Q-state: value with action

Value Iteration

$$Q^*(s, a) = \max_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a Q^*(s, a) = \max_s \sum_{a, s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

Bellman equation

Value Iteration

$$V_{k+1} \leftarrow \max_s \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

Time complexity of each iteration: $O(SA)$

Theorem: will converge to unique optimal value.

\Rightarrow Definition: Bellman equation $U_{k+1} \leftarrow B U_k$, $|U_{k+1}| = \max_s |U(s)|$

proof: $\|B U_k - B U_{k+1}\| \leq \gamma \|U_k - U_{k+1}\|$

$$B U^* I = V$$

$U_{k+1} = T[U_k]$ converges to V^*

Compute Action from Values.

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

or use Q-value: $\pi^*(s) = \arg \max_a Q^*(s, a)$

Q-Value Iteration

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_a Q^*(s', a)]$$

Policy Iteration

1. Policy Evaluation

$$V_k^*(s) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')]$$

Time complexity: $O(S^2)$. start with $V_0^*(s) = 0$

2. Policy Improvement

$$\pi_{k+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k^*(s')]$$

Reinforcement learning

Don't know Transition or Reward.

Basic Idea: At time step t, agent following a policy $\pi_{\text{old}}(s_t)$

\Rightarrow obtain an observation s_t of the surrounding environment.

\Rightarrow produce action a_t

\Rightarrow then, environment will transmit to r_t , and agent will receive reward r_t

The goal is to maximize the expected cumulative reward as:

$$R = \mathbb{E} \left[\sum_{t=0}^T r_t | s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T \right] \quad \text{Transition}$$

Model Based

1. Learn empirical MDP model

count outcomes s' for each s, a .
directly estimate each entry in $T(s, a, s')$ from counts.
discover each $R(s, a, s')$ when we experience the transition.

2. Solve the learned MDP

use value iteration as before.

Pros: Make efficient use of experiences.

Cons: May not scale to large state space.

RL feed back loop tends to magnify small model errors.
Much harder when the environment is partially observable.

Model Free

use samples instead of probability model.

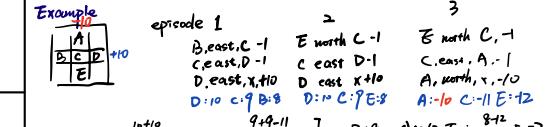
Passive Reinforcement Learning

Input: fixed policy $\pi(s)$ Goal: learn state value $V^*(s)$

根据 policy $\pi(s)$ 进行 -> 观察状态, 然后得到最终的 reward.

根据 reward 回溯到 episode 中每一个 state.

计算 $\pi(s, t)$ 对应 state 的 reward 和 $V^*(s)$.



Pros: 1. easy to understand 2. Doesn't require T or R

3. It converges to right answer in the limit

Cons: 1. Each state need to be learned separately

2. Ignore information about state connection

\Rightarrow take long time to learn.

Temporal Difference learning

Sample Based:

$$\text{Sample} = R(s, \pi(s), s') + \gamma V^*(s')$$

$$V^*(s) = (1-\alpha)V^*(s) + \alpha \cdot \text{sample} = V^*(s) + \alpha (\text{sample} - V^*(s))$$

$$\beta \text{ learning iteration: } Q(s, a) \leftarrow (1-\alpha)Q(s, a) + \alpha \cdot \text{sample} + \beta \cdot \text{sample} (R(s, a) + \gamma V^*(s'))$$

ϵ -greedy

each step, with small probability ϵ , act randomly
with large probability $1-\epsilon$, act on current policy

Does a lot of stupid things \Rightarrow Jumping off a cliff lots of time
keeps doing stupid things forever \Rightarrow decay ϵ towards 0

Optimistic Exploration Functions

Take value estimate u and a visit count n , and returns an optimistic utility, e.g. $f(u, n) = u + \frac{k}{n}$.

\Rightarrow modified α -update:

$$Q(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha [R(s, a, s') + \gamma \max_a f(Q(s', a), n(s'))]$$

Feature Based Representation

$$W(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_m f_m(s)$$

$$Q_w(s, a) = w_1 f_1(s, a) + \dots + w_m f_m(s, a)$$

$$\text{Update: } L = \frac{1}{2} (a - f)^2 \Rightarrow w_i \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$$

$$\Rightarrow w_i \leftarrow w_i + \alpha [R(s, a, s') + \max_a Q(s', a)] - Q(s, a)$$

$$\Rightarrow w_i \leftarrow w_i + \alpha [R(s, a, s') + \max_a Q(s', a)] - Q(s, a)$$

$$\text{total error: } \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \sum w_i f_i(x_i))^2$$

$$\text{minimize error: } w_i \leftarrow w_i + \alpha [\Gamma + \gamma \max_a Q(s', a) - Q(s, a)] f_i(s, a)$$

target prediction

Policy Search

Start with an ok solution (e.g., α -learning), then finetune by hill climbing or gradient ascent on feature weights.

Pros: Works well for partial observability / stochastic policy

Cons: 1. How do we tell the policy got better?

2. Need to run many sample episodes.

3. If there are a lot of features, this can be impractical

Supervised Machine Learning

Naive Bayes \Rightarrow calculate by sample.

$$\text{Inference: } P(Y, f_1, f_2, \dots) = \frac{P(Y) \prod_i P(f_i | Y)}{\prod_i P(Y) \prod_i P(f_i | Y)}$$

Parameter Estimation

Example

$$P(\text{observation} | \theta) = P(\text{red}, \text{blue} | \theta)$$

$$= P(\text{red} | \theta) P(\text{red} | \theta) P(\text{blue} | \theta) = \theta^2 (1-\theta)$$

$$\theta = \arg \max \theta^2 (1-\theta) \quad (\text{derivative to } \theta)$$

Overfit / Generalize. pass

Smoothing.

$$\text{Laplace Smoothing: } P_{\text{lap}}^{(X)} = \frac{c_{Xj} + 1}{\sum_i c_{Xi}} = \frac{1 + c_{Xj}}{K + \sum_i c_{Xi}}$$

$$\text{Maximum likelihood: } \theta = \arg \max \prod_i P(X_i | \theta)$$

$$\text{Laplace: } \theta_{\text{lap}} = \arg \max_{\theta} P(X | \theta) P(\theta)$$

$P(\theta)$ is Dirichlet distribute parameterized by K

$$\text{Precision: } \frac{t}{k} \quad \text{Recall: } \frac{t}{r}$$

t : current true r : total true k : current sample

Break-even point: Precision = Recall

$$\text{F-measure: } F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Neural Network

$$\begin{array}{l} \text{Positive: output} + 1 \\ \text{Negative: output} - 1 \end{array} \quad \begin{array}{c} f_1 \xrightarrow{w_1} \\ f_2 \xrightarrow{w_2} \\ f_3 \xrightarrow{w_3} \end{array} \boxed{\sum} \rightarrow \boxed{?} \rightarrow \dots$$

$$\text{bias: } y = wf + b$$

Decision Boundary

$$w = w + y^* \cdot f \quad \text{Before: } w \cdot f(y) \quad \text{After: } w \cdot f(y) + y^* \cdot f \cdot f(y)$$

Perception

Separability: true if some parameters get the same training set perfectly correct

Convergence: if the training set is separable, perceptron will eventually converge.

$$\text{Sigmoid: } \sigma(x) = \frac{1}{1 + e^{-x}} \quad (\text{binary case})$$

Logistic Regression

$$\max_w U(w) = \max_w \sum \log P(y^{(i)} | x^{(i)}; w)$$

$$\text{with } P(y=1 | x; w) = \frac{1}{1 + e^{-w \cdot f(x)}}$$

$$P(y=-1 | x; w) = \frac{1}{1 + e^{-w \cdot f(x)}}$$

$$\text{multiclass: } P(y^{(i)} | x^{(i)}; w) = \frac{e^{w^{(i)} \cdot f(x)}}{\sum_j e^{w^{(j)} \cdot f(x)}}$$

Unsupervised Learning

Kmeans pass.

$$\text{optimal: } \phi(s_i x_i), \phi(a_i), \phi(c_i) = \sum \text{dist}(x_i, a_i)$$

points assignments means

GMM

$$P(X=x) = \frac{1}{(2\pi)^{\frac{m}{2}} \| \Sigma \|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

$$(\text{Gaussian}) \quad P(X) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k), \text{ s.t. } \sum_k \pi_k = 1, \sum_k \pi_k \mu_k = \mu$$

$P(Y)$: Distribution over k components (clusters)

$P(X|Y)$: generate data from each μ_i and Σ_i

for m data points, for component i , suppose we have n data points with label i :

$$\mu_i = \frac{1}{n} \sum_{j=1}^n x_j \mid \frac{x}{n} = \frac{1}{n} \sum_{j=1}^n (x_j - \mu_i) (x_j - \mu_i)^T / n = \frac{n}{m}$$

Objective: Maximum the likelihood:

$$\prod_j P(y_j=i | x_j) = \prod_j \pi_i N(x_j | \mu_i, \Sigma_i)$$

$$= \prod_j \pi_i \mathcal{N}(x_j | \mu_i, \Sigma_i)$$

Use EM Algorithm

E-step: Compute probability of each instance having each possible label

$$P(y_j=i | x_j, \theta^t) \propto \pi_i^t N(x_j | \mu_i^t, \Sigma_i^t)$$

M-step: Treating each instance as fractionally having both labels, compute the new parameters.

$$\mu_i^{t+1} = \frac{\sum_j P(y_j=i | x_j, \theta^t) x_j}{\sum_j P(y_j=i | x_j, \theta^t)}$$

$$\Sigma_i^{t+1} = \frac{\sum_j P(y_j=i | x_j, \theta^t) [(x_j - \mu_i^t)(x_j - \mu_i^t)^T]}{\sum_j P(y_j=i | x_j, \theta^t)}$$

$$\pi_i^{t+1} = \frac{\sum_j P(y_j=i | x_j, \theta^t)}{m}$$

Math Behind EM

$$L(\theta; D) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{i=1}^k Q(i | x_j) \log \frac{P(x_i | x_j | \theta)}{Q(i | x_j)}$$

Jensen's Inequality

LLM

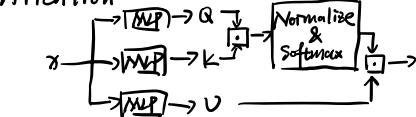
text \rightarrow Tokenize \rightarrow token ids.

word embeddings: $\begin{array}{c} \text{word} \\ \downarrow \\ \mathbb{R}^d \end{array} \quad \begin{array}{c} \text{word} \\ \downarrow \\ \mathbb{R}^d \end{array}$

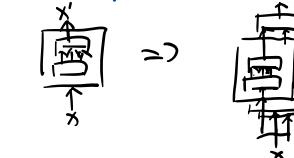
Autoregressive Models



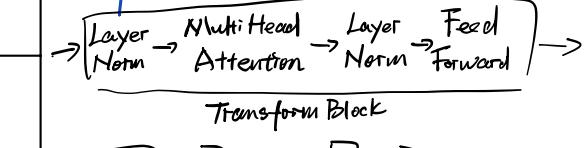
Attention:



Multi Head



Transformer Architecture



pretrain: train a large model with lot of data on self-supervised task

finetune: continue training the same model on task you care about

Policy Search

1. Initialize policy π_θ somehow

2. Estimate policy performance: $J(\theta) = V(s)$

3. Improve policy:

Hill Climbing: Change θ , evaluate new policy. keep if better

Gradient ascent: Estimate $\nabla_\theta J(\theta)$, change θ to ascent gradient: $\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\theta_k)$

$$J(\theta) = \mathbb{E}_{T \sim P(\theta)} \left[\sum_t r(s_t, a_t) \right] \approx \frac{1}{N} \sum_t r(s_t, a_t)$$

$$\theta^* = \arg \max_\theta J(\theta) \quad \text{sum over samples from } T_\theta$$

$$J(\theta) = \mathbb{E}_{z \sim P_\theta(z)} [r(z)] = \int P_\theta(z) r(z) dz$$

$$\nabla_\theta J(\theta) = \int \nabla_\theta P_\theta(z) r(z) dz = \int P_\theta(z) \nabla_\theta \log P_\theta(z) r(z) dz = \mathbb{E}_{z \sim P_\theta(z)} [\nabla_\theta \log P_\theta(z) r(z)]$$

Define Advantage functions: $A^t(s, a) = Q^t(s, a) - V^t(s)$

$$\text{TD error } \mathbb{E}_\pi [s_t] = \mathbb{E}_\pi [Q^t(s_t, a_t) - V^t(s_t)]$$

Let t denote a trajectory from an arbitrary episode:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=0}^{|t|} A^t(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

$$\approx \frac{1}{N} \sum_t A^t(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t | s_t)$$