

PROBABILITY THEORY

Product $p(x,y) = p(y|x)p(x)$

Bayes $p(x|y)p(y) = p(y|x)p(x)$

Marg. $p(x) = \int_x p(x,y)dy = \int_x p(x|y)p(y)dy$

Entropy $H(x) = -\int p(x) \log p(x) dx$

$H(x,y) = H(x|y) + H(y) = H(y|x) + H(x)$

$H(x|y) = H(x)$ if $x \perp y$

$H(N(\mu, \Sigma)) = \frac{1}{2} \ln |2\pi e \Sigma|$

IG $I(x,y) = H(x) - H(x|y) = I(y;x)$

↳ Monotone $F(A \cup \{x\}) - F(A) \leq$

Submodular $F(B \cup \{x\}) - F(B) \forall A \subseteq B, \forall x$

KL-Div $KL(Q||P) = \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta$

↳ $KL(Q||P) \geq 0$ $KL(Q||P) \neq KL(P||Q)$

Properties

$IE[ax+by] = aIE[x] + bIE[y]$

$Var[X] = IE[(x-IE[x])(x-IE[x])^T]$

$= IE[x^2] - IE[x]^2$

$Cov(x,y) = IE[(x-IE[x])(y-IE[y])^T]$

$= IE[xy] - IE[x]IE[y]$

$V[ax \pm by] = a^2V[x] + b^2V[y] \pm 2abCov(x,y)$

Joint Expectation $IE_{x,y}[.] = IE_x[IE_y[.|x]]$

Tower Rule $IE_x[x] = IE_y[IE_{x|y}[x|y]]$

Law of Unc. Sta. $IE_x[g(x)] = \int g(x)p(x)dx$

Tot. Var. $V[y] = IE_x[V[y|x]] + V[IE_x[y|x]]$

Jensen $g(IE[x]) \leq IE[g(x)]$ $g(\cdot)$ convex

Cauchy-Schwarz $|IE[x,y]|^2 \leq IE[x^2]IE[y^2]$

$IE[X^T] = IE[X]^T$ $V[AX] = AV[X]A^T$

Multi Var Gaussian $MX \sim N(M\mu, M\Sigma M^T)$

Sum: $x \perp y \Rightarrow N(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$

Cond. $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim N(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix})$

↳ $x_1|x_2 \sim N(\mu_{x_1|x_2}, \Sigma_{x_1|x_2})$

↳ $\mu_{x_1|x_2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)$

↳ $\Sigma_{x_1|x_2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$

PDF $\frac{1}{\sqrt{2\pi^n}|\Sigma|} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))$

BAYESIAN LINEAR REGRESSION

Model $y = w^T \phi(x) + \epsilon$ $\epsilon \sim N(0, \sigma^2)$

Likelihood $p(y|x, w) = N(w^T x, \sigma^2)$

Prior $p(w) = N(0, \sigma^2 I)$

Inference $p(w|x, y) = N(\bar{\mu}, \bar{\Sigma})$

↳ $\bar{\mu} = (X^T X + \frac{\sigma^2}{\sigma^2} I)^{-1} X^T y$

↳ $\bar{\Sigma} = (\sigma^2 X^T X + \sigma^2 I)^{-1}$

$w_{MLE} = \arg \max_w p(y|x, w) = (X^T X)^{-1} X^T y$

$w_{MAP} = \arg \max_w p(w|x, y) = \bar{\mu}$

Predict $O(Nd^2)$

$p(y^*|x, y, x^*) = N(\bar{\mu}^T x^*, x^{*T} \bar{\Sigma} x^* + \sigma^2)$

1. Epistemic Uncert. Lack of Data

2. Aleatoric Uncert. Irred Noise

Hyperparams 1. $\hat{\lambda} = \frac{\sigma^2}{\sigma^2}$ by Ridge Regr.

2. $\hat{\sigma}^2 = \frac{1}{n} \sum (y_i - \hat{w}^T x_i)^2$ MAP

3. Solve $\hat{\sigma}^2 = \hat{\sigma}^2 / \hat{\lambda}$

GAUSSIAN PROCESS $O(N^3)$

Model $y = w^T \phi(x) + \epsilon$ $\epsilon \sim N(0, \sigma^2)$

Prior $p(w) = N(0, 1)$

$f \sim GP(\mu, k)$ $f = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \in \mathbb{R}^n$

$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \in N(\mu, k)$ $\mu = \begin{bmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{bmatrix}$

$K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix}$ $K_{x,x} = \begin{bmatrix} k(x, x_1) \\ \vdots \\ k(x, x_n) \end{bmatrix}$

PREDICT $p(f|x_{1:n}, y_{1:n}) = GP(\mu^*, k^*)$

$\mu^*(x) = \mu(x) + K_{x,A}(K_{AA} + \sigma^2 I)^{-1}(y_A - \mu_A)$

$k^*(x, x') = A = \{x_1, \dots, x_n\}$ observed

$= k(x, x') - K_{x,A}(K_{AA} + \sigma^2 I)^{-1} K_{x',A}^T$

Hyperparams $\hat{\theta} = \arg \max_{\theta} p(y_{1:n} | x_{1:n}, \theta)$

$= \arg \max_{\theta} \int p(y_{1:n} | f, x_{1:n}, \theta) p(f | \theta) df$

Kernels $k(x, x') = \phi(x)^T \phi(x')$

$k(x, x') = k(x', x)$ K_{AA} PSD $\forall A \subseteq X$

$k(x, x') = Cov[f(x), f(x')]$ Symmetric

$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ Positive definite $\forall a, d$

$k(x_1, x_2) = Cov(f(x_1), f(x_2))$

$= IE[(f(x_1) - \mu(x_1))(f(x_2) - \mu(x_2))]$

STATIONARY $k(x, x') = k(x-x')$

ISOTROPIC $k(x, x') = k(\|x-x'\|_2)$

KALMAN FILTER $y_t \perp x_{1:t-1} | x_t$

$x_{t+1} \perp x_{1:t-1}, y_{1:t-1} | x_t$

Motion $p(x_{t+1} | x_t) = N(x_{t+1} | Fx_t, \Sigma_x)$

$x_{t+1} = Fx_t + \epsilon_t$ $\epsilon_t \sim N(0, \Sigma_x)$

Sensor $p(y_{t+1} | x_t) = N(y_{t+1} | Hx_t, \Sigma_y)$

$y_t = Hx_t + \eta_t$ $\eta_t \sim N(0, \Sigma_y)$

Update $M_{t+1} = F M_t + K_{tt}(y_{t+1} - H F M_t)$

$Z_{t+1} = (I - K_{tt} H)(F Z_t F^T + \Sigma_x)$

Gain $K_{tt} = (F Z_t F^T + \Sigma_x)^{-1} (H^T (H(F Z_t F^T + \Sigma_x)^{-1} H^T + \Sigma_y)^{-1})$

Bayesian Filtering Keep track x_t . $p(x_1)$

$= N(\mu, \Sigma)$ assume $p(x_{t+1} | y_{1:t-1})$

Predict $p(x_{t+1} | y_{1:t}) = \int p(x_{t+1} | x_t) p(x_t | y_{1:t}) dx_t$

BAYESIAN LOGISTIC REGRESSION

Model $w = \text{sign}(w^T x)$

Likelihood $p(y|x, w) = \text{Bern}(b(y w^T x))$

$= \prod_{i=1}^n p(y_i | x_i; w) = \prod_{i=1}^n b(y_i - w^T x_i)$

Prior $p(w) = N(0, \sigma^2 I)$

Laplace Approx $p(w|x, y) \approx N(\hat{w}, \Lambda^{-1})$

$\hat{w} = \arg \min_w \frac{1}{2\sigma^2} \|w\|_2^2 + \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$

$\Lambda = -\nabla^2 \log p(w|x_{1:n}, y_{1:n})$

$= \sum_{i=1}^n x_i x_i^T \pi_i (1 - \pi_i)$

$= X^T \text{diag}[\pi_i (1 - \pi_i)] X$

$\pi_i = b(\hat{w}^T x_i)$ \wedge does not depend on Z

Predictions $p(y^* | x^*, x_{1:n}, y_{1:n})$

$\approx \int \sigma(y^* w^T x) N(w; \hat{w}, \Lambda^{-1}) dw$

$= \int \sigma(y^* f) N(f; \hat{w}^T x^*, x^{*T} \Lambda^{-1} x^*) df$

VARIATIONAL INFERENCE

Approximates $p(\theta|x, y) \approx q_\lambda(\theta)$

where $q_\lambda(\theta) = N(\mu = \hat{\theta}, \Sigma = \Lambda^{-1})$

↳ $\hat{\theta} = \arg \max_{\theta} p(\theta|x, y)$

$\Lambda = -\nabla^2 \log p(\theta|x, y)$ Hessian

Derivation MAP $\psi(\theta) = \log p(\theta|x_{1:n}, y_{1:n})$

$\hat{\psi}(\theta) = \psi(\theta) + (\theta - \hat{\theta})^T \nabla \psi(\hat{\theta})$

$+ \frac{1}{2} (\theta - \hat{\theta})^T \nabla^2 \psi(\hat{\theta}) (\theta - \hat{\theta})$

$= \psi(\hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})^T H_\psi(\hat{\theta}) (\theta - \hat{\theta})$

Forward KL $KL(p||q)$ covers full pdf

Reverse KL $KL(q||p)$ focuses on Mode_{\max}

VI $q^* \in \arg \min_q KL(q||p)$

$= \arg \max_{q \in \mathcal{Q}} IE_{q \sim p} [\log p(y|\theta)] - KL(q||p)$

$= \arg \max_{q \in \mathcal{Q}} IE_{q \sim q} [\log p(y|\theta)] + H(q)$

ELBO

Reparam. Trick as q dep. on Var. params

$\nabla IE_{q \sim q} [f(\theta)] = \nabla_{\lambda} IE_{p \sim p} [f(g(\lambda, \eta))]$

\uparrow distr. dep. on λ \uparrow new distr. $g(\lambda, \eta) \sim \theta$

e.g. Gaussians $q(\theta|\lambda) = N(\theta; \mu, \Sigma)$

$\lambda = [\mu, C]$ $\Sigma = CC^T$ $\theta = C\epsilon + \mu$

$\phi(\epsilon) = N(0, 1)$ $\epsilon = C^{-1}(\theta - \mu)$ $\phi(\epsilon) = q(\theta|\lambda) |C|$

MARKOV-CHAIN MONTE-CARLO

Hoeffding's Inequality

$P(|IE_p[f(x)] - \frac{1}{N} \sum_{i=1}^N f(x_i)| > \epsilon)$

$\leq 2 \exp(-2N\epsilon^2/C^2)$

Ergodic MC if \exists a finite t s.t.

every state can be reached from

every state in exactly t steps.

Stationary Distributions Ergodic

MC has unique and positive $\pi(x) > 0$

s.t. $\forall x \lim_{N \rightarrow \infty} p(X_N = x) = \pi(x)$

Indep of $p(x_1)$ All MC made ergodic $P_{ij} > 0$

Detailed Balance Eq. $\uparrow \pi = \pi P, \sum \pi_i = 1$

$Q(x)P(x'|x) = Q(x')P(x|x')$

Ergodic Theorem if MC Ergodic

$\Rightarrow \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i) = \sum_{x \in \mathcal{X}} \pi(x) f(x)$

$IE[f(x)]$ via MCMC

$\approx (T-t_0)^{-1} \sum_{\tau=t_0+1}^T f(x^{(\tau)})$

Metropolis Hastings

Proposal $x' \sim R(x'|x)$, Given $x_t \sim \pi$

Acceptance: with Prob.

$\alpha = \min\{1, \frac{Q(x')R(x|x')}{Q(x)R(x'|x)}\}$

\rightarrow Set $x_{t+1} = x'$ $Q(x) = \exp(-f(x))$

Otherwise \rightarrow Prob $1 - \alpha$ set $x_{t+1} = x$

Metropolis Adjusted Langevin Algorithm

Take ∇ info into account to prefer

proposals into Region with high density

$R(x'|x) = N(x'; x - \tau \nabla f(x); 2\tau I)$

+ log-concave \rightarrow MALA Converges to $\pi(x)$

- Requires full energy function $f \downarrow$

S.G. Langevin Dynamics

SGD + Gaussian Noise, guarantee to

converge if $\eta_t \in \Theta(t^{-1/3})$

Gibbs Sampling

$x^{(t)} = [x_1^{(t)}, \dots, x_n^{(t)}]^T$

for $i \in \{1, 2, \dots, n\}$:

$x_i^{(t+1)} \sim p(x_i | x_1^{(t+1)}, \dots, x_{i-1}^{(t+1)}, x_{i+1}^{(t)}, \dots, x_n^{(t)})$

BAYESIAN DEEP LEARNING

Homoscedastic: ϵ Heteroscedastic: $\epsilon(x)$

Prior $p(\theta) = N(\theta; 0, \sigma^2 I)$ on weights

Likelihood

$p(y|x, \theta) = N(y; f_\mu(x, \theta), \exp(f_\sigma(x, \theta)))$

MAP

$\hat{\theta} = \arg \min_{\theta} -\log p(\theta) - \sum_{i=1}^n \log p(y_i | x_i, \theta)$

$= \arg \min_{\theta} \lambda \|\theta\|_2^2 + \frac{1}{2} \sum_{i=1}^n (\sigma(x_i; \theta)^2 \cdot \|y_i - \mu(x_i; \theta)\|^2 + \log \sigma(x_i; \theta)^2)$

Predict $p(y^* | x, y, x^*) = \int p(y^* | x^*, \theta) \cdot p(\theta | x, y) d\theta$ Intractable $\theta \in \mathbb{R}^{n \times n \times n}$

VI Given Variational posterior q Can

approx predictive distributions by sampling

$p(y^* | x^*, x_{1:n}, y_{1:n})$

$\approx IE_{q \sim q(\cdot|\lambda)} [p(y^* | x^*, \theta)]$

$\approx 1/m \sum_{j=1}^m p(y^* | x^*, \theta^{(j)}) = \mu(x^*)$

$$V[y^* | x, y, x^*] = \mathbb{E}[V[y^* | x^*, \theta]] + V[\mathbb{E}_\theta[y^* | x^*, \theta]]$$

Aleatoric Uncert. $\approx \frac{1}{m} \sum_{j=1}^m \sigma^2(x^*, \theta^{(j)}) + \frac{1}{m} \sum_{j=1}^m (\mu(x^*, \theta^{(j)}) - \bar{\mu}(x^*))^2$

Dropout as VI $q(\theta | \lambda) = \Pi_j q_j(\theta_j | \lambda_j)$
 $q_j(\theta_j | \lambda_j) = p \delta_{\theta_j} + (1-p) \delta_{\lambda_j}(\theta_j)$

MCMC for Neural Networks

1. Subsampling $\approx 1/T \sum_{j=1}^T p(y^* | x^*, \theta^{(j)})$

2. Gaussian Approximation

Keep track of a Gaussian Approx of the params $q(\theta | M_{1:d}, \Sigma_{1:d})$
 $M_i = 1/T \sum_{j=1}^T \theta_i^{(j)}$, $\Sigma_i = 1/T \sum_{j=1}^T (\theta_i^{(j)} - M_i)^2$
 $\theta^{(j)} \in \mathbb{R}^d$ sampled after burn-in

3. Probabilistic Ensembles for NNs

Train m NNs on bootstrap datasets $p(y^* | x^*, \theta) \approx 1/m \sum_{j=1}^m p(y^* | x^*, \theta^{(j)})$

ECE $\sum_{m=1}^M \frac{|B_m|}{M} |acc(B_m) - Conf(B_m)|$

acc $(B_m) = |B_m|^{-1} \sum_{i \in B_m} \mathbb{I}[\hat{y}_i = y_i]$

Conf $(B_m) = |B_m|^{-1} \sum_{i \in B_m} \hat{p}_i$

ACTIVE LEARNING

Uncertainty Sampling Pure Explorative

$$x_{t+1} = \arg \max_{x \in X} \Delta_I(x | S_t) = \arg \max_{x \in X} IG(f(x); y_x | y_{S_t}) = \arg \max_{x \in X} \frac{1}{2} \log(1 + \frac{\sigma^2(x)}{\bar{\sigma}^2}) = \arg \max_{x \in X} \sigma^2(x)$$

Const.

$\Delta_I(x | S_t) = F(S_t \cup \{x\}) - F(S_t)$

Greedy Algorithm provides Const. Approx.

$F(S_T) \geq (1 - \frac{1}{e}) \max_{S \subseteq D, |S| \leq T} F(S)$

• Fails to distinguish Epistemic/Aleatoric

• Heteroscedastic Case, **most uncertain outcomes \neq most informative**

Maximizing IG yields

$x_{t+1} \in \arg \max_x \frac{\sigma_f^2(x)}{\bar{\sigma}^2(x)}$

epistemic **Aleatoric**

Bayesian Active Learning By Disagreement

Classif $x_{t+1} = \arg \max_{x \in X} H[y_x | x_{1:t}, y_{1:t}] = \arg \max_{x \in X} H[y_x | x_{1:t}, y_{1:t}] - \mathbb{E}_{\theta | x_{1:t}, y_{1:t}} H[y_x | \theta]$

BAYESIAN OPTIMIZATION

Cumulative Regret $R_T = \sum_{t=1}^T (\max_x f(x) - f(x_t))$

Sublinear $R_T/T \rightarrow 0$, implies $\max_t f(x_t) \rightarrow f(x^*)$

Algo $f \sim \mathcal{GP}(M_0, k_0)$

for $t=1$ to T do:

$x_t \leftarrow \arg \max_{x \in X} F(x; M_{t-1}, k_{t-1})$

$y_t \leftarrow f(x_t) + \epsilon_t$

Bayesian update μ_t and k_t **NOT CONVEX**

UCB $(x; M, \epsilon) = \mu(x) + \beta_t \epsilon(x)$

β_t : Controls Exploration $\epsilon(x) = \sqrt{k(x, x)}$

GP-UCB $R_T/T = O(\sqrt{\sigma_T/T})$

$\gamma_T = \max_{S \subseteq X, |S|=T} I(f; y_S)$

Bounds on IG due to submodularity

Linear: $\gamma_T = O(d \log T)$

Gaussian: $\gamma_T = O((\log T)^{d+1})$

Matern $\nu > \frac{1}{2}$: $\gamma_T = O(T^{\frac{d}{2\nu+1}} (\log T)^{\frac{2\nu}{2\nu+1}})$

Guarantees sublinearity / Convergence

PI $(x; M, \epsilon) = P(f(x) > f^*) = \Phi(\frac{\mu(x) - f^*}{\epsilon(x)})$

$\Phi(z) = \frac{1}{2} (1 + \text{erf}(\frac{z}{\sqrt{2}}))$

EI $(x; M, \epsilon) = \mathbb{E}[f(x) - f^* | \mathcal{F}_t] = (\mu(x) - f^*) \Phi(z) + \epsilon(x) \phi(z)$

Thompson Sampling

Sample function $\tilde{f} \sim p(f | x_{1:t}, y_{1:t})$

Select $x_{t+1} \in \arg \max_{x \in D} \tilde{f}(x)$

Randomness of \tilde{f} sufficient to trade exploration / exploitation

RL BASICS

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

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

Action-Value (Q) $Q^\pi(x,a) = r(x,a) + \gamma \sum_{x'} P(x'|x,a) V^\pi(x')$

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

Expected Value $J(\pi) = \sum_{m=0}^{\infty} \gamma^m R_{t+m}$

Multiple Trajectories $J(\pi) = \sum_{i=1}^M P(\tau_i) J(\tau_i)$

V^π(x) closed form

$V^\pi = [V^\pi(1) \dots V^\pi(n)]^T$

$r^\pi = [r(1, \pi(1)) \dots r(n, \pi(n))]^T$

$P^\pi = \begin{bmatrix} p(1|1, \pi(1)) & \dots & p(1|n, \pi(n)) \\ \vdots & \ddots & \vdots \\ p(n|1, \pi(1)) & \dots & p(n|n, \pi(n)) \end{bmatrix}$

Each Row Sums to 1

$V^\pi = r^\pi + \gamma P^\pi V^\pi \Rightarrow V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$

Invertible if $\gamma \in [0, 1)$

Advantage Function $A^\pi(x,a) = Q^\pi(x,a) - V^\pi(x)$

Optimality $V^*(x) = \max_{\pi \in \Pi} V^\pi(x)$

$= \max_{a \in A} Q^*(x,a)$, $Q^*(x,a) = \max_{\pi \in \Pi} Q^\pi(x,a)$

$Q^\pi(x,a)$ one possible (stationary + deterministic) $\pi^*(x) = \arg \max_{a \in A} Q^*(x,a)$

V^* **Unique π^* not always**

\hookrightarrow For infinite Horizon MDP, there exists optimal stationary + deterministic policy.

POMDP can't observe x_t directly, only observation $y_t \rightarrow$ new belief states

$b: \{1, \dots, n\} \rightarrow [0, 1]$, $\sum_x b(x) = 1$

Actions same as original MDP.

Transition Model: 1. Stochastic

Observation $P(y_{t+1}=y | b_t, a_t) = \sum_{x, x'} b_t(x) P(x'|x, a_t) P(y|x')$

2. State Update: $b_{t+1}(x') = \frac{1}{z} \sum_x b_t(x) P(x_{t+1}=x' | x_t=x, a_t) P(y_{t+1}|x')$

Reward: $r(b_t, a_t) = \sum_x b_t(x) r(x, a_t)$

Num Det Policies $\Pi_{i=1}^S, A_i$

MODEL-BASED RL

Value Iter. $V_{k+1}(x) = \max_a \{Q_k\}$

E-optimal in polynomial # iterations

Policy Iter. Initialize π_0 $\gamma < 1$

Until Convergence:

Compute $V^\pi(x) \forall x$

Compute Greedy Policy π_G wrt V^π

Set $\pi \leftarrow \pi_G$

Finds **exact solution** in polynomial # iterations! $\pi^* \in O^*(n^2 m / (1-\gamma))$

Guaranteed to monotonically improve

Rmax Implicit Exploration (Agent always tries to exploit)

$r(x,a) = R_{max}$, $P(x^* | x, a) = 1 \forall x, a$

Repeat: update $r(x,a)$, $P(x' | x, a)$

If observed 'enough', Recompute π according to P, r .

E-optimal π in # steps polynomial in $|x|, |A|, 1/\epsilon$, $\log(1/\delta)$ and R_{max} .

MODEL-FREE RL

TD on-policy follow π to obtain (x, a, r, x') , Update via bootstrapping

$\hat{V}^\pi(x) \leftarrow (1-\alpha_t) \hat{V}^\pi(x) + \alpha_t (r + \gamma \hat{V}^\pi(x'))$

α_t satisfies $\sum_t \alpha_t = \infty$ $\sum_t \alpha_t^2 < \infty$

\wedge all x chosen ∞ often $\hat{V}^\pi \rightarrow V^\pi$

Low Variance θ biased

Q-learning off-policy

Estimate $\hat{Q}^*(x,a)$ observe x, a, x', r

$Q(x,a) \leftarrow (1-\alpha_t) Q(x,a) + \alpha_t (r + \gamma \max_{a'} Q(x', a'))$

Same converge

Conditions as TD

Optimistic Q initialize $\hat{Q}^*(x,a)$

$= R_{max} / (1-\gamma) \prod_{t=1}^T (1-\alpha_t)^{-1}$

At time t , pick $a_t \in \arg \max_a \hat{Q}^*(x_t, a)$

Convergence same as Rmax

Q-learning: Mem $O(nm)$, RT: $O(m)/\epsilon$

FUNCTION APPROX RL

$Q(x, a; \theta) = \theta^T \phi(x, a)$

$L(\theta) = \frac{1}{2} (\hat{Q}^*(x, a; \theta) - r - \max_{a'} \hat{Q}^*(x', a'; \theta_{old}))^2$

Until Converged

In state x , pick a

Observe x' , reward r

$\theta \leftarrow \theta - \alpha_t \delta \nabla_\theta Q(x, a; \theta)$

$\delta := Q(x, a; \theta) - r - \gamma \max_{a'} Q(x', a'; \theta)$

DEEP Q-LEARNING (DQN)

$\ell_{DQN}(\theta) = \frac{1}{2} \sum_{(x, a, r, x') \in D} (r + \gamma \max_{a'} Q(x', a'; \theta_{old}) - Q(x, a; \theta))^2$

Suffers "maximization bias" \downarrow

DDQN Current network for $\arg \max$

$\ell_{DDQN}(\theta) = \frac{1}{2} \sum_{(x, a, r, x') \in D} (r + \gamma \max_{a'} Q(x', a'; \theta_{old}) - Q(x, a; \theta))^2$

where $a^*(x; \theta) = \arg \max_{a \in A} Q^*(x, a; \theta)$

POLICY GRADIENT **Unbiased but high Variance**

REINFORCE

$\theta \leftarrow \theta + \eta \gamma^t G_t \nabla_\theta \log \pi(a_t | x_t; \theta)$

Reward-To-Go $G_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$

Online Actor Critic

At time t after observing (x, a, r, x')

$\theta_\pi \leftarrow \theta_\pi + \eta_t Q(x, a; \theta_Q) \nabla \log \pi(a | x; \theta_\pi)$

$\theta_Q \leftarrow \theta_Q - \eta_t (Q(x, a; \theta_Q) - r - \gamma Q(x', \pi(x', \theta_\pi); \theta_Q)) \nabla Q(x, a; \theta_Q)$

Under Compatibility Conditions Guaranteed to improve. Baselines to reduce Var