

# Probabilistic Artificial Intelligence

## Multivariate Gaussians $\mathcal{N}(y; \Sigma, \mu)$

$\frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(y-\mu)^\top \Sigma^{-1}(y-\mu)\right)$

$\sigma_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)]$ ,  $\sigma_i^2 = \mathbb{V}(X_i)$ ,  $X_i \perp X_j \Leftrightarrow \sigma_{ij} = 0$ , Joint distribution over  $n$  variables requires  $\mathcal{O}(n^2)$  parameters

## Gaussian Conditional Distribution

$p(X_A|X_B=x_b) = \mathcal{N}(\mu_{A|B}, \Sigma_{A|B})$

$\mu_{A|B} = \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_b - \mu_B)$      $\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$

## Mutliples of Gaussians

$X \sim \mathcal{N}(\mu_V, \Sigma_{VV})$ ,  $Y = MX$ ,  $M \in \mathbb{R}^{m \times d}$ , then  $Y \sim \mathcal{N}(M\mu_V, M\Sigma_{VV}M^T)$

## Sums of Gaussians

$X \sim \mathcal{N}(\mu_V, \Sigma_{VV})$ ,  $X' \sim \mathcal{N}(\mu'_V, \Sigma'_{VV})$ ,  $Y = X + X'$  then  $Y \sim \mathcal{N}(\mu_V + \mu'_V, \Sigma_{VV} + \Sigma'_{VV}) \rightarrow$  not the case for product. **Sum of Gaussian** distributed RV is Gaussian distributed. The product of Gaussian distributed RV is NOT Gaussian distributed (but the product of Gaussian PDF is Gaussian).

## Useful Math

### Probabilities

$\mathbb{E}[X] = \int x \cdot p(x) dx$  or  $\sum x \cdot p(x)$

**Tower Law**  $\mathbb{E}[X] = \mathbb{E}_y[\mathbb{E}_x[X|Y]]$

$\mathbb{V}[X] = \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

$\mathbb{V}_x[b + cX] = c^2 \mathbb{V}_x[X]$

$\text{Cov}[X, Y] = \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$

$\mathbb{V}[X+Y] = \mathbb{V}[X] + \mathbb{V}[Y] + 2\text{Cov}[X, Y]$

Change of variables:  $f_y(y) = f_x(x) |dg(x)/dx|^{-1}$

**Jensen's inequality**  $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$  f convex

$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$  f concave

**Robbins Monro (RM) conditions** If sequence  $\epsilon_t$  satisfies  $\sum \epsilon_t = \infty$  and  $\sum \epsilon_t^2 < \infty$  then it will converge to optimum with probability 1. Can be used to check if **learning rate**  $\alpha$  will lead to optimum

## Entropy and Mutual Information

$H(X) = \mathbb{E}_{x \sim p(x)}[-\log P(X)] = -\sum_{i=1}^n P(x_i) \log P(x_i)$

$H(X, Y) = H(Y, X)$      $H(X \cdot Y) = H(X) + H(Y)$

$H(X, Y) \leq H(X) + H(Y)$      $H(X) = \mathbb{E}_X[-\log p(X)]$

$H(X|Y) = \mathbb{E}_{y \sim p(y)}[H(X|Y=y)] = H(X, Y) - H(Y)$

$I(X; Y) = KL(p(X, Y) || p(X)p(Y))$

$I(X; Y) = H(X) - H(X|Y) = I(Y; X) \leq I(X; Y, Z)$

If  $q$  is Gaussian:  $H(q) = \frac{1}{2} \ln |2\pi e \Sigma|$   $q \sim \mathcal{N}(\mu, \Sigma)$

If  $Y = X + \epsilon$  then  $H(Y|X) = H(\epsilon)$

# Bayesian Learning

We can write  $X = aY + b + \epsilon$  if  $X, Y$  are jointly Gaussian

## Ridge Regression

Linear and Ridge Regression fail when multicollinearity, i.e. when more than two explanatory variables are highly linearly related

$\hat{w} = \text{argmin}_w \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$

$\hat{w} = (X^T X + \lambda I)^{-1} X^T y$

Ridge is the same as finding the **MAP parameter estimate** assuming noise is (cond.) iid Gaussian and prior is Gaussian. Ridge can be derived from BLR by only taking the mode of the posterior

$\hat{w} = \text{argmax}_w P(w) \Pi_i P(y_i | x_i, w)$

$w_{MAP} = \sigma_y^{-2} (\sigma_w^{-2} I + \sigma_y^{-2} X^T X)^{-1} X^T y$

## Bayesian Linear Regression (BLR) $\mathcal{O}(nd^2)$

Pr:  $p(w) = \mathcal{N}(0, \sigma_p^2 I)$      $\epsilon = \mathcal{N}(0, \sigma_n^2 I)$

Li:  $p(y|x, w, \sigma_n) = \mathcal{N}(y; w^T x, \sigma_n^2) \equiv y = w^T x + \epsilon$

Po:  $p(w|X, y) = \mathcal{N}(w; \bar{\mu}, \bar{\Sigma})$  has **closed form**

$\bar{\mu} = (X^T X + \sigma_n^2 I)^{-1} X^T y$      $\bar{\Sigma} = (\sigma_n^{-2} X^T X + I)^{-1}$

BLR is the same as averaging all  $w$  acc. to posterior

For test point  $x^*$ ,  $f^* = w^T x^*$

$p(f^*|X, y, x^*) = \mathcal{N}(\bar{\mu}^T x^*, x^{*T} \bar{\Sigma} x^*)$

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

(lack of data) and *Aleatoric* (noise) uncertainties captured

Online updating:  $X_{new}^T X_{new} = X^T X + x_{i+1} x_{i+1}^T$

$X_{new} y_{new} = Xy + y_{i+1} x_{i+1}$

Choosing **hyperparameters**  $\hat{\lambda} = \hat{\sigma}_n^2 / \hat{\sigma}_p^2$  via cross-validation. Estimate  $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{w}^T x_i)^2$ , then solve for  $\hat{\sigma}_p^2$ . Otherwise, use marginal likelihood.

## Kalman Filters (KF)

Assume conditional linear Gaussian dependencies between states (X) and observations (Y).  $X_{t+1} \perp y_{1:t} | X_t$

Motion model:  $X_{t+1} = F X_t + \epsilon_t$      $\epsilon_t \sim \mathcal{N}(0, \Sigma_x)$

Sensor model:  $Y_t = H X_t + \eta_t$      $\eta_t \sim \mathcal{N}(0, \Sigma_y)$

## Bayesian Filtering in KFs

Update:  $\mu_{t+1} = \frac{\sigma_y^2 \mu_t + (\sigma_x^2 + \sigma_y^2) y_{t+1}}{\sigma_x^2 + \sigma_y^2}$      $\sigma_{t+1}^2 = \frac{(\sigma_x^2 + \sigma_y^2) \sigma_t^2}{\sigma_x^2 + \sigma_y^2}$

BLR is a form of KF

## Gaussian Processes (GP)

GP is a normal distribution over functions

GP with linear kernel = BLR

**Prediction** gives closed form formulas. Posterior covariance doesn't depend on the observations  $y$ . Exact computation requires  $\mathcal{O}(n^3)$  but speedup with parallelism, local GP methods, kernel function approximations, inducing point methods

# Covariance (kernel) functions

$K_{x, x'} = \phi(x)^T \phi(x')$      $k(x, x') = \text{Cov}(f(x), f(x'))$

Symmetric:  $k(x, x') = k(x', x) \forall x, x'$

Positive semi-definite:  $x^T K x \geq 0 \Leftrightarrow \lambda_K \geq 0 \forall x \in \mathbb{R}^{|A|}$

**Composition rules** 1.  $k(x, x') = k_1(x, x') + k_2(x, x')$

2.  $k(x, x') = k_1(x, x') \cdot k_2(x, x')$  3.  $k(x, x') = c \cdot k_1(x, x')$

for  $c > 0$  4.  $k(x, x') = f(k_1(x, x'))$  with  $f$  a polynomial (with positive coefficient) or exponential

**Stationary** if:  $k(x, x') = k(\tau)$  with  $\tau = x - x'$

**Isotropic** if:  $k(x, x') = k(\tau)$  with  $\tau = \|x - x'\|_2$

RBF kernel is both Stationary and Isotropic

If RBF bandwidth  $h$  large  $\rightarrow$  **smooth** samples

## Optimizing kernel hyperparameters

Max. the marginal likelihood  $\hat{\theta} = \text{argmax}_{\theta} P(y|X, \theta)$

Integrate instead of optimizing prevents overfitting

## Kernel Function Approximation $\mathcal{O}(nm^2 + m^3)$

**Random Fourier Features** interpret kernel as expectation. Requires a stationary kernel and a randomized feature map. Such a kernel has a Fourier transform. Approximates kernel function uniformly well, wasteful.

**Bochner theorem** A stationary kernel is positive-definite  $\Leftrightarrow$  its Fourier transform is non-negative

## Inducing point methods $\mathcal{O}(n)$

Assumes  $f^* \perp f|u$ , training  $f$ , response  $f^*$

$p(f^*, f) \approx q(f^*, f) = \int q(f^*|u) q(f|u) p(u) du$

**SoR approximation**  $q_{SoR}(f|u)$  has the same mean as  $p(f|u)$  but variance is set to 0

## Approximate Inference

**Laplace Approximation** uses a Gaussian approximation to the posterior distribution obtained from a second-order Taylor expansion around the posterior mode. Lead to **poor** approximation if  $p$  has multiple modes.

**Stochastic Gradient Descent** converges to (local) minimum if learning rate satisfies RM conditions

## Variational Inference (VI)

Approx. unnormalized distribution  $p$  by a simple (tractable) distribution  $q$  which depends on some parameters  $\lambda$ :  $p(\theta|y) = \frac{1}{2} p(\theta, y) \approx q(\theta|\lambda)$

Only twice as expensive as MAP inference for diagonal  $q$ . Only need to be able to differentiate the (unnormalized) joint proba. density  $p$  and  $q$ . The quality of inference is hard to analyze.

## Kullback-Leibler Divergence

$KL(q || p) = \mathbb{E}_{\theta \sim q} [\log \frac{q(\theta)}{p(\theta)}] = \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta \geq 0$

$KL(q || p) = 0 \Leftrightarrow p = q$

$KL(q || p) \neq KL(p || q) \Rightarrow KL$  is not a distance

$\text{argmin}_q KL(q || p) = \text{argmax}_q \mathcal{L}(\lambda)$  (c.f. ELBO)

$KL(q || p)$ : *backw.*, **exclusive** (underestimates, if we have two modes, only one of them will be used)

$KL(p || q)$ : *forw.*, **inclusive** (overestimates)

## Evidence Lower Bound (ELBO)

$\mathcal{L}(\lambda) = \mathbb{E}_{\theta \sim q(\cdot|\lambda)} [\log p(y|\theta)] - KL(q_\lambda || p(\cdot))$

$\mathcal{L}(\lambda) \leq \log p(y)$ , where  $p(y)$  is the evidence

Can calculate  $\nabla \mathcal{L}$  via the following trick

**Reparametrization Trick** We want the expectation of the ELBO to depend on an assumed distribution  $\phi$  and not on  $\lambda$  such that:

$\nabla_\lambda \mathbb{E}_{\theta \sim q_\lambda} [f(\theta)] = \mathbb{E}_{\epsilon \sim \phi} [\nabla_\lambda f(g(\epsilon; \lambda))]$

## Makov Chain Monte Carlo (MCMC)

Approx. unnormalized distribution via sampling

**Hoeffding's inequality** prob. of error decreases exp. in  $N$ .

Problem is that normalizing factor  $Z$  is intractable. Use **ergodic** MC (i.e. every states can be reached from every states in a finite number of steps) that has stationary distribution  $\pi(x) = P(X) = \frac{1}{Z} Q(X)$ . Ergodic MC has a **unique**  $\lim_{N \rightarrow \infty} P(X_N = x) = \pi(x) > 0 \forall x$  and it is **independent of the initial state**.

If MC satisfies the **detailed balance equation** (DBE) then MC has  $\pi(x) = P(X)$

**DBE**:  $\frac{1}{2} Q(x) P(x'|x) = \frac{1}{2} Q(x') P(x|x')$

**Markov Chain**  $X_{t+1} \perp X_{t-1} | X_t$

**Metropolis-Hasting** MCMC old state:  $x$

1. Proposal  $x' \sim R(X'|X=x)$  proposed state:  $x'$

2.  $P(X' = x'|X=x) = \alpha$      $P(X' = x|X=x) = 1 - \alpha$

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

Use **Gibbs sampling**, from  $X_i$  given *all* other, to specify the proposal distribution. Random order Gibbs satisfies DBE, practical variant doesn't but still has correct  $\pi(x)$

Variational inference scales better than sampling techniques like Metropolis-Hasting.

Because joint sample at time  $t$  depends on sample at  $t-1$ , the **law of large numbers doesn't apply**. We use the following theorem to compute expectation with MCMC instead:

**Ergodic Theorem** with  $D$  a finite state space

$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(x_i) = \sum_{x \in D} \pi(x) f(x) = \mathbb{E}_{x \sim \pi} f(x)$

With continuous RVs, use **MALA** for proposal distribution, which converges to  $\pi(x)$  for log-concave distributions, e.g. Bayesian logistic regression

$p(x) = \frac{1}{2} \exp(-f(x))$  is log-concave if the energy function  $f(x)$  is convex. For large data sets, use **SGLD**.

which convergences to  $\pi(x)$  if  $\eta_t \in \mathcal{O}(t^{-1/3})$ .

MALA and SGLD can be improved by adding momentum, resulting in **HMC** (remembering what the gradient was before)

## Bayesian Deep Learning

Consider **nonlinear dependencies** of the likelihood function parameters on the inputs.

**Heteroscedastic** noise depends on input

## Bayesian Neural Network (BNN)

Models are the output of the network, MAP parameter estimate:  $\hat{\theta} = \text{argmin}_{\theta} \mathcal{L}(\theta)$ , Gaussian pr. on the weight are equiv. to applying weight decay.

## Approximate inference for BNNs

**VI for BNNs** aka Bayes-by-backpropagation. Prediction  $p(y^*|x^*, x_{1:n}, y_{1:n}) \approx \mathbb{E}_{\theta \sim q(\cdot|\lambda)} [p(y^*|x^*, \theta)] \approx \frac{1}{m} \sum_{j=1}^m p(y^*|x^*, \theta^{(j)})$  i.e. draw  $m$  sets of weights from posterior and average the neural network predictions. Catches both aleatoric and epistemic uncertainties. Optimize ELBO via SGD.

**MCMC for BNNs** Only requires stochastic gradients of (unnormalized) joint probability, i.e. the same used for MAP estimation. Prediction  $p(y^*|X, y, x^*) \approx \frac{1}{T} \sum_{j=1}^T p(y^*|x^*, \theta^{(j)})$ . Problem is to store all  $T$  samples or models: Can use **Subsampling**, i.e. keep a subset of  $m$  snapshots, or keep track of **Gaussian approximation** of the parameters (e.g. SWAG method).

**Specialized inference techniques for BNNs**  
**Monte-Carlo Dropout** regularization: randomly ignore hidden units during each iteration of SGD with probability  $p$ . Both during training and prediction. Can be viewed as VI. **Probabilistic Ensemble** uses bootstrap sampling, train multiple models on random subsamples of the data

## Active Learning

1. Choose and add the point with the **highest uncertainty** to the training set (often the farthest points from the current evaluation points) 2. Train on the new training set 3. Repeat **Exploration**

We want to pick the points to observe  $f$  which maximizes the **Information Gain** (aka Mutual Information)  $F(S) = H(f) - H(f|y_S)$ . Problem  $F(S)$  is NP-Hard to optimize. We can use **Uncertainty Sampling** (greedy) to get near-optimal result because  $F(S)$  is **monotone submodular**, i.e. satisfies **diminishing returns** property:  $H(X|Y) \geq H(X|Y, Z)$ . It fails for **heteroscedastic** cases because it can't distinguish aleatoric from epistemic uncertainty. Can use BALD.

## Bayesian Optimization

1. Pick next point to evaluate (acc. to the acquisition function by optimizing it) 2. Update the model (Bayes' rule) 3. Repeat **Exploration-Exploitation**

**Bandit** task: multiple options with unknown prob-

ability of producing a reward, the goal is to choose the best option to maximize the overall reward. Idea: pick points that minimize the **Cumulative Regret** Problem is NP-Hard, and we don't know the best possible output. We can use **GP-UCB** acquisition function:  $x_t = \text{argmax}_{x \in D} \mu_{t-1}(x) + \beta_t \sigma_{t-1}(x)$  Picks points for which the **upper confidence bound** is currently the highest. The higher the estimate  $\mu$ , the higher the bound. The higher the uncertainty  $\sigma$ , the higher the bound. **Maximum information gain**  $\gamma_T$  determines the regret. Linear:  $\gamma_T = \mathcal{O}(d \log T)$ , RBF:  $\gamma_T = \mathcal{O}((\log T)^{d+1})$ , Matérn ( $\nu > 2$ ):  $\gamma_T = \mathcal{O}(T^{\frac{d(d+1)}{2\nu+d(d+1)}} \log T)$  **Thompson Sampling**: draw  $\tilde{f} \sim P(f|x_{1:T}, y_{1:T}), x_{t+1} = \text{argmax}_{x \in D} \tilde{f}(x)$ , the randomness of  $\tilde{f}$  is sufficient to trade exploration and exploitation.

## Markov Decision Processes (MDP)

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

$$\pi_g = \text{argmax}_{\pi} \sum_{x'} P(x'|x, a) [r(x, a, x') + \gamma V_{\pi}(x')]$$

### Bellman's Theorem

Policy is optimal  $\Leftrightarrow$  greedy w.r.t. its induced  $V_{\pi}$

**Policy iteration** converges in  $\mathcal{O}(\frac{1}{1-\gamma})$  iterations

**Value iteration** converges in  $\mathcal{O}(\ln \frac{1}{\epsilon})$  iterations of complexity  $\mathcal{O}(nms)$  per iteration

**Value function for optimal policy**  $V^*(x) = \max_a r(x, a) + \gamma \sum_{x'} P(x'|x, a) V^*(x')$

**POMDP** (aka controlled HMM) Generalization of MDP when partially observable state  $X_i$  via noisy observation  $Y_i$ . The agent has a certain **belief** over the states. Most belief states are never reached. We can use **Policy gradient** methods.

## Reinforcement Learning (RL)

Use observed state transitions and rewards to learn an underlying (unknown) MDP. Data is not iid.

**On-policy** agent has control over which action to pick. Can chose to trade exploration-exploitation **Off-policy** agent only gets observational data. **Online**: after every step, **Offline**: with replay dataset/batches.

**Discount**  $\gamma \ll 1$  focused on instant rewards,  $\approx 1$  less greedy, more on future rewards.

## Model-based RL

**Learn the MDP** and optimize

$$P(X_{t+1}|X_t, A) \approx \frac{\text{Count}(X_{t+1}, X_t, A)}{\text{Count}(X_t, A)} \quad r(x, a) \approx \frac{1}{N_{x,a}} \sum_{r: X_t = x, A_t = a} R_t$$

**$\epsilon$ -greedy** pick random action with proba.  $\epsilon$ , pick best action with proba.  $1 - \epsilon$ . Converges to optimal policy if  $\epsilon$  satisfies RM conditions. Problem: doesn't quickly eliminate suboptimal actions.

**$R_{max}$  Algorithm** (On-policy) *Initialize* 1. add fairy tale state  $x^*$  2. set  $r(x, a) = R_{max}$  and  $P(x^*|x, a) = 1 \forall x, a$  3. chose optimal policy. *Repeat* 4. execute policy 5. update  $r(x, a)$  for each visited state-action pair 6. estimate  $P(x'|x, a)$ . Until  $n \in \mathcal{O}(\frac{R_{max}}{\epsilon^2} \log \frac{1}{\delta})$  observations

## Model-free RL

Value func.  $V^{\pi}(x)$  and Action-Value func.  $Q^{\pi}(x, a)$ .

Advantage func.  $A^{\pi}(x, a) = Q^{\pi}(x, a) - V^{\pi}(x)$

**TD-learning** use Bootstrapping for updating

$$V(x) \leftarrow (1 - \alpha)V(x) + \alpha(r + \gamma V(x'))$$

$V^*(x) = \max_a Q^*(x, a)$ , estimate  $Q^*$  from samples

**Q-learning** pick (greedy)  $a = \text{argmax}_a Q(x, a)$

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

Converges to optimal  $Q^*$  if  $\alpha$  satisfies RM conditions.  $Q^*$  depends on  $\gamma$  used during training, larger  $\gamma$  leads to larger  $Q$ , larger noise leads to smaller  $Q$  value.

TD-learning is On-policy, whereas Q-learning can be Off-policy.

**Large state space**: use **parametric function approximation** of (action) value function to scale up, e.g. using Neural Network leads to DQN. (less biased than Q-learning). Linear function approximation must be linear in the weights (not in the states).

**Large action space**: use **Policy gradient methods**

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [(r(\tau) - b) \nabla_{\theta} \log \pi_{\theta}(\tau)]$$

Note  $\nabla_{\theta} \log \pi_{\theta}(\tau) = \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)}$  For MDP:  $\pi_{\theta}(\tau) = \pi(a|x; \theta)$

Gradient  $\nabla J(\theta)$  is unbiased but can have large variance, using a **baseline**  $b > 0$  reduces variance.

**REINFORCE** reduces variance by using **reward to go**  $G$  instead of  $r(\tau)$ .

**Actor-critic methods** use value function estimate and policy gradient methods

$$\nabla J(\theta) = \mathbb{E}_{(x, a) \sim \pi_{\theta}} [Q(x, a; \theta_Q) \nabla \log \pi(a|x; \theta)]$$

derived by using the **discounted state occupancy** measure  $\rho(x) = \sum_{t=0}^{\infty} \gamma^t P(X_t = x)$

Allows application in **online** setting

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

$$\theta_Q \leftarrow \theta_Q - \eta_Q (Q(x, a; \theta_Q) - r - \gamma Q(x', \pi(x', \theta_{\pi}); \theta_Q)) \nabla \log \pi(a|x; \theta_{\pi})$$

Updates improvement guaranteed under **compatibility conditions**

Can reduce variance via baselines (using advantage function)  $A(x, a) = Q(x, a; \theta_Q) - V(x; \theta_V)$

$$\theta_{\pi} \leftarrow \theta_{\pi} + \eta_{\pi} A(x, a) \nabla \log \pi(a|x; \theta_{\pi})$$

Actor-critic can be parallelized

**TRPO** and **PPO** are variant of policy gradient/actor-critic methods. PPO is a heuristic variant of TRPO widely used in practice.

### Advantage function properties

$$\max_a A^{\pi}(x, a) \geq 0 \forall \pi, x$$

$$\pi^* \text{ is optimal } \Leftrightarrow A^*(x, a) \leq \sigma \forall x, a$$

$$\pi_g(x) = \text{argmax}_a Q^{\pi}(x, a) = \text{argmax}_a A^{\pi}(x, a)$$

**Off-policy policy gradient methods** If we use differentiable approximation of the maximum  $Q$  and differentiable **deterministic** policy, then we can use chain rule (i.e. backprop.) to obtain stochastic gradient. Problem, policy gradient methods rely on randomized policies for exploration, we have to **inject additional action noise** to encourage exploration in

the off-policy variant, e.g. **DDPG**. Can encourage exploration via **entropy regularized** MDPs, e.g. **SAC**

## Homeworks

$$\text{If } y = Xw + \epsilon \Rightarrow w_{MLE} = (X^T X)^{-1} X^T y$$

Kalman filters is like HMM but for continuous variables. High gain means follow the most recent measurement. Low gain means follow the model predictions.

The support of  $q$  must be a subset of the support of  $p$ , otherwise  $KL(q||p) \rightarrow \infty$

Expected improvement

$$EI(X) = \mathbb{E}[(t-y)_+] = \int_{-\infty}^{\infty} (t-y)_+ P(Y|X) dy \text{ where } (t-y)_+ = \max(0, t-y)$$

REINFORCE update:

$$\theta \leftarrow \theta + \alpha G \cdot E \text{ where } E = \frac{\nabla_{\theta} \pi(a|x, \theta)}{\pi(a|x, \theta)}$$
 is the eligibility vector

Linear MDP with feature map  $\phi(x, a)$  means (1)

$\exists \theta$  s.t.  $r(x, a) = \langle \phi(x, a), \theta \rangle$  (2)  $P(x'|x, a) = \langle \phi(x, a), P(x') dx' \rangle$  Note:  $P(x')$  doesn't have to be a distribution, but  $P(x'|x, a)$  should.

Sum rule  $P(X) + P(\bar{X}) = 1$  Product rule  $P(X, Y, Z) = P(X|Y, Z) \cdot P(Y, Z) = P(Y|X, Z) \cdot P(X, Z) = P(Z|X, Y) \cdot P(X, Y)$

Policy gradient:  $\mathcal{A} = \{0, 1\}$ ,  $\pi(1|s_t) = \theta$  and  $\pi(0|s_t) = 1 - \theta$ .

trajectory  $\tau = (s_1, a_1, s_2, a_2, s_3, a_3)$ , future reward  $R_t^{future} = \sum_{t'=t}^H r_{t'}$  i.e.  $R_1^{future} = r_1 + r_2 + r_3$  but  $R_3^{future} = r_3$ , we can calculate  $\nabla_{\theta} \pi(1|s_t) = 1$  and  $\nabla_{\theta} \pi(0|s_t) = -1$  then  $\nabla_{\theta} \log \pi(1|s_t) = \frac{\nabla_{\theta} \pi(1|s_t)}{\pi(1|s_t)}$  and put together  $Gradient = \sum_t R_t^{future} \nabla_{\theta} \log \pi(a_t|s_t)$