

# ML Lec 11

- Dimensionality Reduction

: learn a mapping  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  from the input feature space to a transformed space, where  $k \ll d$   
(the mapping should preserve certain intrinsic structure of the given data)

## 1. Why reducing feature dimension

⇒ computation time, space requirement ↓

⇒ curse of dimensionality

(feature dimension ↑ ⇒ # of data for nice generalization ↑)

⇒ manifold assumption

(intrinsic dimensionality could be smaller)

⇒ visualization for better understanding and analysis

## 2. Linear Dimensionality Reduction

: Learn a linear transformation  $z = w^T x$ , where  $x \in \mathbb{R}^d$ ,  $z \in \mathbb{R}^k$ ,  $w \in \mathbb{R}^{d \times k}$  and  $k \ll d$

### ★ Principal Component Analysis (PCA)

: orthogonal projection of data onto lower dimensional linear space

(basis transformation that transforms the data to a new coordinate system)

→ maximizes the **variance** of the projected data

→ minimizes the **reconstruction error**

① Standardize the data to have mean 0 and variance 1

② Find the first principal component  $\alpha_1$ .

$$\Rightarrow \max_{\alpha_1} \text{Var}(\alpha_1^T x) = \alpha_1^T \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) \alpha_1 = \alpha_1^T \Sigma \alpha_1$$

s.t.  $\alpha_1^T \alpha_1 = 1$       sample covariance matrix  
 $(= \frac{1}{n} \mathbf{X} \mathbf{X}^T)$

$$\Rightarrow \max_{\alpha_1} \min_{\lambda_1} \alpha_1^T \Sigma \alpha_1 - \lambda_1 (\alpha_1^T \alpha_1 - 1)$$

$$\rightarrow \frac{\partial}{\partial \lambda_1} (\alpha_1^T \Sigma \alpha_1 - \lambda_1 (\alpha_1^T \alpha_1 - 1)) = \alpha_1^T \alpha_1 - 1 = 0$$

$$\rightarrow \frac{\partial}{\partial \alpha_1} (\alpha_1^T \Sigma \alpha_1 - \lambda_1 (\alpha_1^T \alpha_1 - 1)) = 2 (\Sigma \alpha_1 - \lambda_1 \alpha_1) = 0$$

$(\alpha_1 : \text{eigenvector of } \Sigma, \lambda_1 : \text{eigenvalue of } \Sigma)$

$$\Rightarrow \max \lambda_1 (= \text{Var}(\alpha_1^T x))$$

$$\text{s.t. } \sum \alpha_1 = \lambda_1 \alpha_1, \quad \alpha_1^T \alpha_1 = 1$$

$\therefore \alpha_1$ : eigenvector that corresponds to the largest eigenvalue  $\lambda_1$ .

③ Find the second principal component  $\alpha_2$

$$\Rightarrow \max_{\alpha_2} \text{Var}(\alpha_2^\top \chi) = \alpha_2^\top \Sigma \alpha_2$$

$$\text{s.t. } \alpha_2^\top \alpha_2 = 1$$

$$\text{Cov}(\alpha_1^\top \chi, \alpha_2^\top \chi) = \alpha_2^\top \Sigma \alpha_1 = \lambda_1 \alpha_1^\top \alpha_2 = 0$$

$$\Rightarrow \max_{\alpha_2} \min_{\lambda_2, \mu} \alpha_2^\top \Sigma \alpha_2 - \lambda_2 (\alpha_2^\top \alpha_2 - 1) - \mu \alpha_1^\top \alpha_2$$

$$\rightarrow \frac{\partial}{\partial \lambda_2} (\alpha_2^\top \Sigma \alpha_2 - \lambda_2 (\alpha_2^\top \alpha_2 - 1) - \mu \alpha_1^\top \alpha_2) = \alpha_2^\top \alpha_2 - 1 = 0$$

$$\rightarrow \frac{\partial}{\partial \mu} (\alpha_2^\top \Sigma \alpha_2 - \lambda_2 (\alpha_2^\top \alpha_2 - 1) - \mu \alpha_1^\top \alpha_2) = \alpha_1^\top \alpha_2 = 0$$

$$\rightarrow \frac{\partial}{\partial \alpha_2} (\alpha_2^\top \Sigma \alpha_2 - \lambda_2 (\alpha_2^\top \alpha_2 - 1) - \mu \alpha_1^\top \alpha_2) = 2(\Sigma \alpha_2 - \lambda_2 \alpha_2) - \mu \alpha_1 = 0$$

(multiplying  $\alpha_1$ )  $\Rightarrow 2(\alpha_1^\top \Sigma \alpha_2 - \lambda_2 \alpha_1^\top \alpha_2) - \mu \alpha_1^\top \alpha_1 = 2(\lambda_1 \alpha_1^\top \alpha_2 - \lambda_2 \alpha_1^\top \alpha_2) - \mu$

$$\Rightarrow \mu = 0 \text{ so that } \Sigma \alpha_2 = \lambda_2 \alpha_2$$

( $\alpha_2$ : eigenvector of  $\Sigma$ ,  $\lambda_2$ : eigenvalue of  $\Sigma$ )

$$\Rightarrow \max \lambda_2 (= \text{Var}(\alpha_2^\top \chi))$$

$$\text{s.t. } \Sigma \alpha_2 = \lambda_2 \alpha_2, \quad \alpha_2^\top \alpha_2 = 1, \quad \alpha_1^\top \alpha_2 = 0$$

$\therefore \alpha_2$ : eigenvector that corresponds to the largest (possible) eigenvalue  $\lambda_2$

★ Eigenvectors with distinct eigenvalues are orthogonal in symmetric matrix

★ Symmetric matrix is diagonalizable  $\rightarrow$  algebraic multiplicity = geometric multiplicity

★ Gram-Schmidt process enables to find orthonormal basis

★ Magnitude of eigenvalues indicates fraction of variances captured

④ Repeat the process for finding d different principal components

⑤ Choose k components using cross-validation

\* Simple way

⇒ d-dimensional projection with PCA

:  $Z = XA$  where  $\begin{cases} X \in \mathbb{R}^{n \times d} \text{ is a data matrix,} \\ A \in \mathbb{R}^{d \times d} \text{ is a collection of eigenvectors of } \Sigma \end{cases}$

⇒ k-dimensional projection with PCA

:  $Z' = XA'$  where  $A' \in \mathbb{R}^{d \times k}$  so that only k eigenvectors are selected

\* Eigenface

: Principle components of face images derived by PCA

⇒ face image  $x$  can be converted into coordinates of the eigenfaces

$$x \rightarrow ((x - \bar{x}) \cdot v_1, (x - \bar{x}) \cdot v_2, \dots, (x - \bar{x}) \cdot v_k)$$

⇒ a new face image can be generated by linear combination of basis faces

⇒ closest image can be found by Nearest Neighbor search on coefficients.

\* Problems

⇒ assumes normality of the input space distribution.

⇒ fail if data lies on a complicated manifold.

⇒ Euclidian distance may not be a good measure between points in manifold

(⊕ geodesic distance : the distance of two nodes on a graph)

### 3. Non-Linear Dimensionality Reduction

: Learn a low-dimensional non-linear manifolds of given high-dimensional data.

#### ★ Locally Linear Embedding (LLE)

① Select neighbors

② Learn linear weight  $w$  that minimizes the reconstruction error by the neighbors given the data points.  $X$

$$\Rightarrow \min_w \sum_{i=1}^n \|x_i - \sum_j w_{ij} x_j\|_2^2$$

③ Learn low-dimensional vectors  $y$  that minimizes the reconstruction error

by the neighbors given the reconstruction weight  $w$

$$\Rightarrow \min_y \sum_{i=1}^n \|y_i - \sum_j w_{ij} y_j\|_2^2$$

# ML Lec 12

## • Clustering

: process of organizing unlabeled data points into clusters where similar data items are collected

### 1. Usage

⇒ similarity / dissimilarity analysis

⇒ dimensionality reduction

⇒ image segmentation

⇒ neural network quantization

### 2. Requirement

① Algorithm to perform clustering

② Number of clusters

⇒ determined by the cross-validation and criterion function

③ Similarity measure ( $S(x_i, x_j)$ )

: large if  $x_i$  and  $x_j$  is similar (e.g. cosine similarity)

④ Dissimilarity measure ( $D(x_i, x_j)$ )

: small if  $x_i$  and  $x_j$  are similar (e.g. euclidian distance)

→ Euclidian distance :  $\sqrt{\sum_{k=1}^d (x_i^k - x_j^k)^2}$

→ Manhatten distance :  $\sum_{k=1}^d |x_i^k - x_j^k|$

→ Minkowski distance :  $\left( \sum_{k=1}^d (x_i^k - x_j^k)^p \right)^{\frac{1}{p}}$

⑤ Criterion function to evaluate the generated clusters

→ intra-cluster cohesion (compactness)

: how near the data points in a cluster are to the cluster centroid.

$$\text{e.g. Sum of squared errors (SSE)} = \sum_{j=1}^k \sum_{x \in D_j} d(x, \mu_j)^2$$

jth cluster      ↪ centroid of cluster j  
                            ↪ Euclidian distance

→ inter-cluster separation (isolation)

: how far the cluster centroids are from each other

e.g. expert judgement

→ purity

: ratio of correctly assigned data points when the label of the most frequent class

is assigned to each data point in clusters

$$\Rightarrow \text{purity}(D, C) = \frac{1}{N} \sum_{j=1}^k \max_i |D_j \cap C_i|$$

jth cluster      ↪ ith class

## ★ K-means clustering

: each data point belongs to one of the k clusters with the nearest mean

$$\Rightarrow \min_s \sum_{j=1}^k \sum_{x \in s_j} \|x - \mu_j\|^2$$

① Assign each data point to the closest cluster center

② Recalculate centers as the mean of the points in a cluster.

(most frequent value is used for categorical data)

③ Repeat until convergence criterion is met

- no (or minimum) reassignment of data points to different clusters
- no (or minimum) change of centroids
- minimum decrease in the sum of squared errors

\* strength / simple : easy to understand and implement

fast :  $O(tkn)$  time complexity

general : work well with multiple distance measure

\* weakness / applicable only when mean is defined

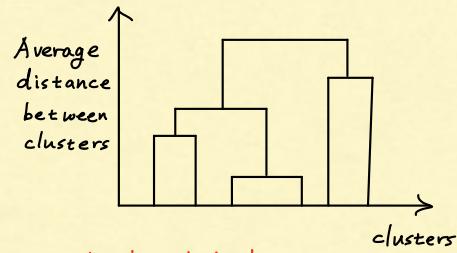
k must be specified in advance

sensitive to outliers, initial centroids

can only create spherical clusters

★ hierarchical clustering

⇒ dendrogram is a visualization tool



① Divisive (top-down) clustering → slower, look global structure

: start with all data points in one cluster, then split into child clusters recursively

until one singleton clusters of individual data points remain.

⇒ can use any flat algorithms that produce a fixed number of clusters (e.g. k-means)

② Agglomerative (bottom-up) clustering → faster, look local structure

: start with each example in singleton cluster, then find 2 nearest clusters and merge

until there is only one cluster

⇒ four common ways to measure cluster distance

1) Minimum distance :  $\min_{x \in D_i, y \in D_j} \|x - y\|$

→ generate minimum spanning tree

→ encourage elongated clusters

→ very sensitive to noise

2) Maximum distance :  $\max_{x \in D_i, y \in D_j} \|x - y\|$

→ encourage compact clusters

→ perform bad for elongated clusters

3) Average distance :  $\frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{y \in D_j} \|x - y\|$

4) Mean distance :  $\|\mu_i - \mu_j\|$

### \* Spectral clustering

: make use of spectrum of the similarity matrix (Graph Laplacian) to reduce dimensionality before clustering (e.g. k-means)

### \* Max-margin clustering

: find maximum margin hyperplanes through data

### \* Gaussian Mixture Model

: assumes that data points are generated from a mixture of finite number of Gaussians with unknown parameters

⇒ EM (expectation and maximization) algorithm can be utilized

# ML Lec 13

- Reinforcement Learning

: learning which actions to take at each state to maximize reward in the long-run

## 1. Notations

⇒ State ( $S_t$ ): representation of the environment

( summarizes the past in a compact way )

⇒ Action ( $A_t$ ): action that an agent takes which leads to a new state

⇒ Reward ( $R_t$ ): numerical value received by performing an action  $A_t$  at state  $S_t$

$$\rightarrow \text{Total reward } (G_t) = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

( discounting is used to determine the present value of the future rewards )

(  $\gamma \rightarrow 0$  : myopic ,  $\gamma \rightarrow 1$  : far-sighted )

⇒ Policy ( $\pi_t$ ): mapping from states to probabilities of selecting each possible action

## 2. Characteristics

⇒ there is no supervisor, but a reward signal

⇒ feedback is delayed and not instantaneous

⇒ time matters (sequential, not i.i.d.)

### 3. Markov Decision Process (MDP)

: learning task for reinforcement learning that satisfies Markov Property , which consists of the following components

- A set of states  $S$
- A set of actions  $A$
- A transition model  $T(s,a,s')$  known : model-based approach.
- A reward function  $R(s,a,s')$  unknown : model-free approach

⇒ when states or actions are finite → "finite MDP"

⇒ agent does not know  $T$  and  $R$  which are estimated by experience

⇒ goal is to find a policy  $\pi(s)$  that maximizes total reward.

## 4. Value functions

$\Rightarrow$  for finite MDP, policies can be partially ordered.

$$\rightarrow \pi' > \pi \text{ if and only if } v_{\pi'}(s) \geq v_{\pi}(s), \forall s \in S$$

$\Rightarrow$  there is always one or more policies that are better than or equal to others which share the same state-value ( $v_{\pi^*}(s)$ ) and action-value function ( $q_{\pi^*}(s, a)$ )

### ① State-value function for policy $\pi$ ( $v_{\pi}(s)$ )

: expected return starting from the state  $s$ , then following  $\pi$

$$\rightarrow v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

### ② Action-value function for policy $\pi$ ( $q_{\pi}(s, a)$ )

: expected return starting from the state  $s$  taking action  $a$ , then following  $\pi$

$$\rightarrow q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

### ★ Bellman expectation equation

: a relationship between the value of a state and the value of its successor states

$$\Rightarrow v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_{t+1} = s'\right]\right]$$

$$= \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma v_{\pi}(s')\right]$$

$$\Rightarrow q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right]$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_{t+1} = s'\right]\right]$$

$$= \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma v_{\pi}(s')\right]$$

## ★ Bellman optimality equation

$$\begin{aligned}
 \Rightarrow V_{\pi^*}(s) &= \max_a q_{\pi^*}(s, a) \\
 &= \max_a \mathbb{E}_{\pi^*} [G_t | S_t = s, A_t = a] = \max_a \mathbb{E}_{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \\
 &= \max_a \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s, A_t = a \right] \\
 &= \max_a \mathbb{E}_{\pi^*} \left[ R_{t+1} + r \cdot V_{\pi^*}(S_{t+1}) | S_t = s, A_t = a \right] \\
 &= \max_a \underbrace{\sum_{s'} p(s' | s, a) [r(s, a, s') + r \cdot V_{\pi^*}(s')]}_{\text{(any policy that is greedy w.r.t. } V_{\pi^*} \text{ is an optimal policy)}}
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow q_{\pi^*}(s, a) &= \mathbb{E}_{\pi^*} [G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \\
 &= \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s, A_t = a \right] \\
 &= \mathbb{E}_{\pi^*} \left[ R_{t+1} + r \cdot V_{\pi^*}(S_{t+1}) | S_t = s, A_t = a \right] \\
 &= \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma \cdot \max_{a'} q_{\pi^*}(S_{t+1}, a') | S_t = s, A_t = a \right] \\
 &= \sum_{s'} p(s' | s, a) [r(s, a, s') + \gamma \max_{a'} q_{\pi^*}(s', a')]
 \end{aligned}$$

(any action that is greedy w.r.t.  $q_{\pi^*}$  is an optimal action)

(the agent does not have to do one-step ahead search.)

## 5. Approaches in reinforcement learning

### ① Model-based approach

→ build a model of the environment and plan using the model

### ② Model-free approach

→ Value-based approach

→ estimate the optimal value functions  $V_{\pi^*}(s)$ ,  $q_{\pi^*}(s, a)$

→ Policy-based approach

→ search directly for the optimal policy  $\pi^*$

# ML Lec 14

- Model-based Reinforcement Learning

1. Dynamic programming (DP) → use bootstrap  
→ Does not scale well to large problems.

: technique to solve problems with optimal substructure and overlapping subproblems

⎛ optimal substructure : optimal solution can be decomposed into solutions of subproblems.  
 ⎝ overlapping subproblems : subproblems recur many times where solutions are stored and reused

⇒ Assuming full knowledge of MDP and deterministic  $\pi$ , it can be resolved through DP.

⇒ it considers all successor states (full back-up)

★ Policy iteration  $(\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi_k \xrightarrow{E} V_{\pi_k})$

: iteratively update the optimal policy using DP

① Initialize  $\pi$  randomly

② Repeat until convergence

$$\rightarrow V_\pi(s) = V_{\pi'}(s)$$

→ Repeat until convergence (expensive)

$$\begin{aligned}
 & (\text{policy evaluation}) \Rightarrow V_\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma V_\pi(s')] \quad \left( \begin{array}{l} \text{Bellman} \\ \text{expectation} \\ \text{equation} \end{array} \right) \\
 & \qquad \qquad \qquad \text{(full-back-up)} \\
 & \rightarrow \pi'(s) = \operatorname{argmax}_a q_\pi(s,a) \\
 & \qquad \qquad \qquad \text{(given)} \\
 & (\text{policy improvement}) \qquad \qquad \qquad = \operatorname{argmax}_a \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma V_{\pi^*}(s')]
 \end{aligned}$$

(+) Converge in a finite number of iteration (often small in practice)

(-) Each iteration requires a full policy evaluation, which is expensive

## \* Value iteration

: iteratively search for optimal value function using DP

① Initialize  $V(s) = 0, \forall s$

② Repeat until convergence

$$(\text{policy evaluation}) \rightarrow V(s) = \max_a \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma \cdot V(s')] \quad (\text{Bellman optimality equation})$$

③ After obtaining the optimal value function, find the optimal policy

$$(\text{policy improvement}) \rightarrow \pi^*(s) = \arg \max_a \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma V^*(s')]$$

(+) Each iteration is computationally efficient

(-) Convergence is only symptotic

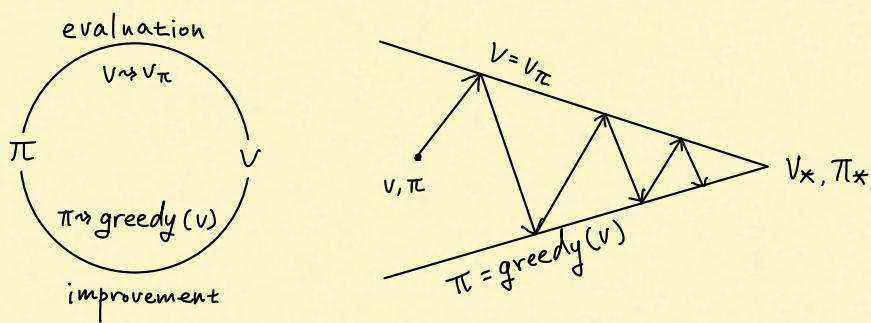
With finite number of states, both converge to optimal solution in polynomial time

## \* Generalized Policy iteration (GPI)

: allows the policy evaluation and the policy improvement to interact

$\Rightarrow$  every reinforcement learning can be described as GPI

$\Rightarrow$  if both the evaluation and improvement stabilize, then the value function and the policy must be optimal.



## • Value-based Reinforcement Learning

: optimize the state-value and action-value functions

(policy is generated directly from these value functions)

1. Monte Carlo method (on-policy: learn  $\pi$  from episodes of  $\pi$ )

⇒ learns from complete sample returns

→ only defined for episodic tasks that have ends

⇒ learns directly from experience

→ online: optimality is attainable even without model

→ simulated: no need for a full model

★ Monte Carlo policy iteration ( $\pi_0 \xrightarrow{E} q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} q_{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi_* \xrightarrow{E} q_{\pi_*}$ )

: estimate  $V_\pi(s)$  for model-based MC and  $q_\pi(s,a)$  for model-free MC of the current behavior policy  $\pi$ , then improve  $\pi$

① Initialize  $\pi(s), V(s)$  (or  $q(s,a)$ ),  $r(s)$  for all  $s$   
 $\xrightarrow{\text{empty list}}$

② Repeat until convergence

→ generate an episode using  $\pi$

→ for each  $s$  appearing in the episode

⇒  $G \leftarrow$  return following the occurrence of  $s$  (followed by  $a$ )

⇒ Append  $G$  to  $r(s)$

⇒  $V(s)$  (or  $q(s,a)$ )  $\leftarrow$  average ( $r(s)$ )

( $\approx V(s) \leftarrow V(s) + \alpha [G - V(s)]$ )

( $\approx q(s,a) \leftarrow q(s,a) + \alpha [G - q(s,a)]$ )

every-visit MC: consider every time when  $s$  is visited

first-visit MC: consider the first time when  $s$  is visited

→  $\varepsilon$ -greedy policy improvement

: simplest idea for ensuring continual exploration

(all  $m$  actions have non-zero probability)

$$\Rightarrow \pi(a|s) = \begin{cases} 1 - \varepsilon + \frac{\varepsilon}{m} & \text{if } a = \arg\max_a \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma V(s')] \\ & \quad \text{(given)} \\ \frac{\varepsilon}{m} & \text{otherwise} \end{cases}$$

$\Rightarrow \varepsilon \uparrow \Rightarrow$  more exploration ,  $\varepsilon \downarrow \Rightarrow$  more exploitation

\* Theorem. ( $\varepsilon' < \varepsilon$ ) →  $\varepsilon$  would converge to 0 in the end

: For any  $\varepsilon$ -greedy  $\pi$ , the  $\varepsilon'$ -greedy policy  $\pi'$  with respect to  $q_\pi$  is an improvement

$$\begin{aligned} \text{pf)} q_{\pi'}(s, \pi'(s)) &= \sum_a \pi'(a|s) q_{\pi}(s,a) \\ &= \frac{\varepsilon'}{m} \sum_a q_{\pi}(s,a) + (1-\varepsilon') \cdot \max_a q_{\pi}(s,a) \\ &\geq \frac{\varepsilon'}{m} \sum_a q_{\pi}(s,a) + (1-\varepsilon') \cdot \sum_a \boxed{\frac{\pi(a|s) - \frac{\varepsilon'}{m}}{1-\varepsilon'}} q_{\pi}(s,a) \\ &= \sum_a \pi(a|s) q_{\pi}(s,a) \\ &= V_{\pi}(s) \end{aligned}$$

↓  
sum to 1

(+) MC uses return  $G$  that is an unbiased estimate of  $V_{\pi}(s)$

(-) MC uses return  $G$  that has high variance

→ depends on many actions, transitions, rewards

(-) MC have to wait until the episodes end, which could be very long

## 2. Temporal difference

⇒ similar to MC, TD does not need a model

→ learn directly from experience

⇒ similar to DP, TD bootstraps

→ learn a guess from a guess

### ★ Simple Temporal difference (On-policy)

: estimate  $v_\pi(s)$  for the current behavior policy  $\pi$ , then improve policy  $\pi$

① Initialize  $v(s)$  for all  $s$  and  $v(\text{terminal state}) = 0$

② Repeat until convergence

→ initialize  $s$

→ repeat until  $s$  is terminal

→ choose  $a$  from  $s$  using policy  $\pi$

→ take action  $a$  and observe reward  $r$  and next state  $s'$

$$\Rightarrow v(s) \leftarrow v(s) + \alpha (r + \gamma v(s') - v(s))$$

Temporal difference

$$\Rightarrow s \leftarrow s'$$

→  $\epsilon$ -greedy policy improvement

(+) TD can be fully incremental

→ learn before knowing the final outcome

→ work in incomplete, non-terminating environments

→ less memory, less peak computation

→ converges faster than MC

(+) TD uses target  $r + \gamma v(s')$  that has low variance

→ depends on one random action, transition, reward

(-) TD uses target  $r + \gamma v(s')$  that is a biased estimator of  $v_\pi(s)$

## ★ SARSA (On-policy)

: estimate  $q_{\pi}(s, a)$  for the current behavior policy  $\pi$ , then improve policy  $\pi$

① Initialize  $q(s, a)$  for all  $s$  and  $a$  and  $q(\text{terminal state}, \cdot) = 0$

② Repeat until convergence

→ initialize  $s$

→ choose  $a$  from  $s$  using policy derived from  $q(s, a)$  (e.g.  $\epsilon$ -greedy)

→ repeat until  $s$  is terminal

⇒ take action  $a$  and observe reward  $r$  and next state  $s'$

⇒ choose  $a'$  from  $s'$  using policy derived from  $q(s, a)$  (e.g.  $\epsilon$ -greedy)

$$\Rightarrow q(s, a) \leftarrow q(s, a) + \alpha [r + \gamma \underbrace{q(s', a') - q(s, a)}_{\text{on-policy}}]$$

$$\Rightarrow s \leftarrow s', a \leftarrow a'$$

## ★ Q-learning (Off-policy: learn $\pi$ from episodes of $\pi'$ )

: directly approximate  $q_{\pi^*}(s, a)$  independent of the policy being followed

① Initialize  $q(s, a)$  for all  $s$  and  $a$  and  $q(\text{terminal state}, \cdot) = 0$

② Repeat until convergence

→ initialize  $s$

→ repeat until  $s$  is terminal

⇒ choose  $a$  from  $s$  using policy derived from  $q(s, a)$  (e.g.  $\epsilon$ -greedy)

⇒ take action  $a$  and observe reward  $r$  and next state  $s'$

$$\Rightarrow q(s, a) \leftarrow q(s, a) + \alpha [r + \gamma \max_a q(s', a) - q(s, a)]$$

$$\Rightarrow s \leftarrow s' \quad \text{off-policy}$$

✗ SARSA ⇒ learns longer but safe path

✗ Q-Learning ⇒ learns shorter but non-safe path

### 3. Value function approximation

: represent  $v(s)$  or  $q(s,a)$  with a parameterized function instead of a table

⇒ no need to explicitly store or learn every value of  $v(s)$  or  $q(s,a)$

⇒ more compact representations that generalize across  $v(s)$  or  $q(s,a)$

ex) Linear functions, decision trees, nearest neighbors, neural networks  
(differentiable)

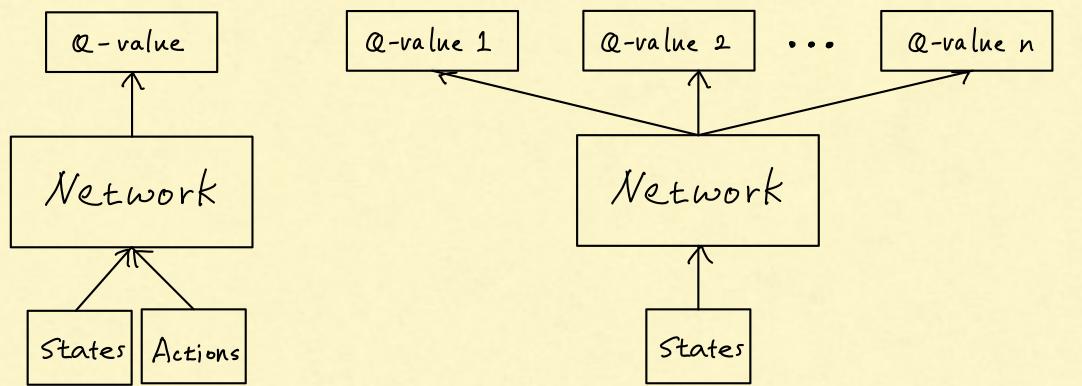
#### ★ Deep Q-learning

: represent value functions by deep Q-network with weight  $w$ .

$$(Q(s,a,w) = q_{\pi}(s,a))$$

$$\Rightarrow \text{loss function } (\mathcal{L}(w)) = \mathbb{E}_{s,a,r,s' \sim D} \left[ (r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w))^2 \right]$$

$$\rightarrow \frac{\partial \mathcal{L}(w)}{\partial w} = \mathbb{E}_{s,a,r,s' \sim D} \left[ 2(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)) \frac{\partial Q(s, a, w)}{\partial w} \right]$$



$\Rightarrow$  experience replay  $\rightarrow$  break correlations in the states

① store all experiences  $\langle s, a, r, s' \rangle$  in a replay memory

② while training, use random minibatches from the replay memory instead of the recent most transition

$\Rightarrow$  freeze target Q-network  $\rightarrow$  avoid oscillation

① compute q-learning target w.r.t. old, fixed parameters  $w^-$

$$\rightarrow r + \gamma \max_{a'} Q(s', a', w^-)$$

② optimize MSE between Q-network and q-learning targets

$$\rightarrow L(w) = \mathbb{E}_{s, a, r, s' \sim D} [(r + \max_{a'} Q(s', a', w^-) - Q(s, a, w))^2]$$

$\rightarrow$  periodically update the fixed parameter  $w^-$ .

$\Rightarrow$  clip rewards or normalize network adaptively  $\rightarrow$  robust gradients  
[-1, +1]

$\rightarrow$  prevent  $Q(s, a, w)$  from getting too large

$\rightarrow$  ensures the gradient to be well conditioned

$\rightarrow$  can use the same learning rate across multiple games

(cannot tell difference between small and large rewards)

# ML Lec 15

## • Policy-based Reinforcement Learning

: directly learn the policy by approximating it with a parameterized function.

$$(\pi_\theta(a|s) = p(a|s, \theta) \rightarrow \text{find the best parameter } \theta)$$

$\Rightarrow$  advantages

$\rightarrow$  better convergence properties

$\rightarrow$  effective in high-dimensional or continuous action space

$\rightarrow$  can learn stochastic policies

$\Rightarrow$  disadvantages

$\rightarrow$  typically converges to a local rather than global optimum

$\rightarrow$  typically evaluating a policy is inefficient and has high variance

### 1. Policy objective functions

: measure the quality of policy  $\pi_\theta$

$\Rightarrow$  episodic environment

$$\rightarrow \text{start value} : J_1(\theta) = V_{\pi_\theta}(s_1)$$

$\Rightarrow$  continuing environment

$$\rightarrow \text{average value} : J_V(\theta) = \sum_s d_{\pi_\theta}(s) V_{\pi_\theta}(s) := \mathbb{E}[V_{\pi_\theta}(s)]$$

$$\rightarrow \text{average reward per time-step} : J_R(\theta) = \sum_s d_{\pi_\theta}(s) \sum_a \pi_\theta(a|s) r(s, a, \cdot) := \mathbb{E}[r(s, a, \cdot)]$$

( $d_{\pi_\theta}(s)$ : stationary distribution ( $\pi P = \pi$ ) of markov chain for  $\pi_\theta$ )  
 $\pi_\theta \neq \pi$

## 2. Policy optimization

: find  $\theta$  that maximizes  $J(\theta)$

$\Rightarrow$  both gradient-based and non-gradient based methods can be used

(e.g. (stochastic) gradient ascent, conjugate gradient, quasi Newton, genetic algorithm, ...)

## 3. policy gradient algorithm

: searches for a local maximum in  $J(\theta)$  by ascending the gradients w.r.t.  $\theta$

$\Rightarrow$  assume  $\pi_\theta$  is differentiable almost everywhere

(e.g. linear function or neural network)

$$\Rightarrow \Delta\theta = \alpha \cdot \nabla_\theta J(\theta) = \alpha \cdot \nabla \mathbb{E}[V_{\pi_\theta}(s)] \text{ or } \alpha \cdot \nabla \mathbb{E}[r(s, a, \cdot)]$$

$\rightarrow$  Monte Carlo sampling can be used

## 4. score function trick

: estimate the gradient with sampling without consideration of  $d_{\pi_\theta}(s)$

$$\begin{aligned} \Rightarrow \nabla_\theta \mathbb{E}[r(s, a, \cdot)] &= \nabla_\theta \sum_s d(s) \sum_a \pi_\theta(a|s) r(s, a, \cdot) \\ &= \sum_s d(s) \sum_a \nabla_\theta \pi_\theta(a|s) r(s, a, \cdot) \\ &= \sum_s d(s) \sum_a \pi_\theta(a|s) \cdot \underbrace{\frac{\nabla_\theta \log \pi_\theta(a|s)}{(\text{score function})}}_{\text{score function}} r(s, a, \cdot) \\ &= \mathbb{E}[\nabla_\theta \log \pi_\theta(a|s) \cdot r(s, a, \cdot)] \end{aligned}$$

## 5. Policy gradient theorem

: For any differentiable policy  $\pi_\theta(a|s)$  and any policy objective function  $J$ ,  
( $J_L, \frac{1}{1-\gamma} J_R, J_R$ )  
the policy gradient is as follow.

$$\rightarrow \nabla_\theta \mathbb{E}[r(s, a, \cdot)] = \mathbb{E}[\nabla_\theta \log \pi_\theta(a|s) q_{\pi_\theta}(s, a)]$$

## ★ Monte Carlo policy gradient method

$$\Rightarrow \text{assume } J(\theta) = \mathbb{E} [r(s,a,\cdot)]$$

→ cannot simply sample  $R_{t+1}$  and take gradient

( $\because R_{t+1}$  is a scalar and does not depend on  $\theta$ )

⇒ score function trick can be used

$$\rightarrow \theta_{t+1} = \theta_t + \alpha R_{t+1} \nabla_\theta \log \pi_{\theta_t}(A_t | S_t)$$

(stochastic policy gradient update)

→ the probability of selecting  $A_t$  should go up if  $R_{t+1}$  is high

⇒ softmax policy

→ applicable when the action space is discrete and not too large

$$\rightarrow \pi_\theta(a|s) = \frac{\exp(\theta^\top \phi(s,a))}{\sum_{k=1}^m \exp(\theta^\top \phi(s,a_k))}$$

$$\rightarrow \nabla_\theta \log \pi_\theta(a|s) = \phi(s,a) - \mathbb{E}_{\pi_\theta} [\phi(s,\cdot)]$$

⇒ gaussian policy

→ applicable when the action space is continuous

$$\rightarrow \pi(a|s) = N(\mu(s), \sigma^2) \text{ where } \mu(s) = \phi(s)^\top \theta \text{ (linear function)}$$

(variance can be either fixed or further parameterized)

$$\rightarrow \nabla_\theta \log \pi_\theta(a|s) = \frac{(a - \mu(s)) \phi(s)}{\sigma^2}$$

$\Rightarrow$  REINFORCE algorithm

: estimate  $q_{\pi_\theta}(s, a)$  with the return  $G_t$

( $G_t$  is an unbiased estimator of  $v_\pi(s)$ )

① Initialize policy parameter  $\theta$

② Repeat until convergence

→ generate an episode  $\{s_0, a_0, r_1, \dots, s_T, a_T, r_T\} \sim \pi_\theta(\cdot | \cdot)$

→ repeat  $t=0, \dots, T-1$

$$\Rightarrow G_t = \sum_{k=0}^{T-1} \gamma^k r_{t+k+1}$$

$$\Rightarrow \theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(a_t | s_t) G_t$$

(-) MC policy gradient uses return  $G_t$  that has **high variance**

\* Actor-critic method

$\Rightarrow$  actor adjusts parameter  $\theta$  of  $\pi_\theta(a|s)$  by policy gradient algorithm

$\Rightarrow$  critic  $Q(s, a, w)$  estimates  $q_{\pi_\theta}(s, a)$  using an approximate policy evaluation algorithms (e.g. Monte Carlo method, Temporal difference methods)

→ linear function can be used :  $Q(s, a, w) = w^\top \phi(s, a)$

$\Rightarrow$  actor-critic algorithm

① Initialize parameter  $\theta, w$  and  $s_0$

② Repeat  $t=1, \dots, T$

→ sample  $R_t \sim r(s, a, \cdot)$  and  $s' \sim p(s'|s, a)$

→ sample  $a' \sim \pi_\theta(a'|s')$

(update critic)  $\rightarrow w \leftarrow w + \beta (R_t + \gamma Q(s', a', w) - Q(s, a, w)) \nabla_w Q(s, a, w)$

(update actor)  $\rightarrow \theta \leftarrow \theta + \alpha \cdot \nabla_\theta \log \pi_\theta(a|s) \cdot Q(s, a, w)$

→  $s \leftarrow s'$ ,  $a \leftarrow a'$

\* Reducing variance of policy gradient  $\nabla_{\theta} J(\theta)$

① subtract the baseline function  $b(s)$  that only depends on  $s$

→ state value function can be used :  $b(s) = V_{\pi_{\theta}}(s)$

→ this does not change the expectation of  $\nabla_{\theta} J(\theta)$

$$\begin{aligned}\mathbb{E} [\nabla_{\theta} \log \pi_{\theta}(a|s) b(s)] &= \sum_s d(s) \sum_a \nabla_{\theta} \pi_{\theta}(a|s) b(s) \\ &= \sum_s d(s) b(s) \nabla_{\theta} \underbrace{\sum_a \pi_{\theta}(a|s)}_{(=1)} \\ &= 0\end{aligned}$$

② critic estimates the advantage function  $A(s, a, \omega) = Q(s, a, \omega) - V(s, \bar{\zeta})$

by estimating both  $Q(s, a, \omega)$  and  $V(s, \bar{\zeta})$  through TD-learning

→  $\omega \leftarrow \omega + \beta_1 (R_t + \gamma Q(s', a', \omega) - Q(s, a, \omega)) \nabla_{\omega} Q(s, a, \omega)$

→  $\bar{\zeta} \leftarrow \bar{\zeta} + \beta_2 (R_t + \gamma V(s', \bar{\zeta}) - V(s, \bar{\zeta})) \nabla_{\bar{\zeta}} V(s, \bar{\zeta})$

③ actor updates  $\theta$  with an unbiased estimate of  $A(s, a)$

→  $\delta_{\pi_{\theta}} = R_t + \gamma V_{\pi_{\theta}}(s') - V_{\pi_{\theta}}(s)$  : true TD error → unbiased

(in practice)  $(\delta_v = R_t + \gamma V(s', \bar{\zeta}) - V(s, \bar{\zeta})$  : approximate TD error) → biased

→  $\theta \leftarrow \theta + \alpha \cdot \nabla_{\theta} \log \pi_{\theta}(a|s) \cdot \delta_{\pi_{\theta}}$

(in practice)  $(\theta \leftarrow \theta + \alpha \cdot \nabla_{\theta} \log \pi_{\theta}(a|s) \cdot \delta_v)$

## $\therefore$ Summary

① Stochastic gradient ascent for  $\theta$

$\Rightarrow$  Monte Carlo policy gradient

$$\rightarrow \Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$$

$\Rightarrow$  Temporal Difference (0)-Actor-critic policy gradient

$$\rightarrow \Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) V(s, \bar{s})$$

$\Rightarrow$  Q-Actor-critic policy gradient

$$\rightarrow \Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) Q(s, a, w)$$

$\Rightarrow$  Advantage-Actor-critic policy gradient

$$\rightarrow \Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) A(s, a, w)$$

② Value functions estimation ( $V(s, \bar{s})$ ,  $Q(s, a, w)$  or  $A(s, a, w)$ )

$\Rightarrow$  Monte Carlo method

$$\rightarrow \Delta \bar{s} = \beta (G_t - V(s_t, \bar{s})) \nabla_{\bar{s}} V(s_t, \bar{s})$$

$\Rightarrow$  Temporal Difference (0)

$$\rightarrow \Delta \bar{s} = \beta (R_t + \gamma V(s', \bar{s}) - V_{\bar{s}}(s_t, \bar{s})) \nabla_{\bar{s}} V(s_t, \bar{s})$$

$\Rightarrow$  SARSA

$$\rightarrow \Delta w = \beta (R_t + \gamma Q(s', a', w) - Q(s_t, a_t, w)) \nabla_w Q(s_t, a_t, w)$$

$\Rightarrow$  Q-learning

$$\rightarrow \Delta w = \beta (R_t + \gamma \max_{a'} Q(s', a', w) - Q(s_t, a_t, w)) \nabla_w Q(s_t, a_t, w)$$

# ML Lec 16

## • Introduction to Bayesian Machine Learning

### 1. Uncertainty in Machine Learning

⇒ inherent uncertainty (future forecast)

⇒ insufficient observation

⇒ data ambiguity

$\left( \begin{array}{l} \text{conventionally, point estimate finds a single optimal set of parameters} \\ \text{However, given limited observations, it is difficult to estimate parameters} \end{array} \right)$

### 2. Bayesian Machine Learning

: average over the uncertain variables and parameters rather than optimize. → more expensive

⇒ apply to both discriminative and generative learning

⇒ eliminate the arbitrary loss function and regularizers

⇒ facilitate the incorporation of prior knowledge

⇒ quantifies the uncertainty in hypothesis and prediction

⇒ assume that the hypotheses and predictions follow some distributions → transparent

⇒ find the most likely distributions given the data using the prior beliefs and Bayes' rule

$$\rightarrow \text{Posterior: } p(\theta | D) = \frac{P(D|\theta) \cdot p(\theta)}{P(D)} \rightarrow \text{full Bayesian.}$$

$$\rightarrow \text{Evidence: } P(D) = \begin{cases} \sum_{\theta} P(D|\theta) \cdot p(\theta) & (\text{discrete } \theta) \\ \int_{\theta} P(D|\theta) \cdot p(\theta) d\theta & (\text{continuous } \theta) \end{cases}$$

] ⇒ often intractable to compute

$$\rightarrow \text{Likelihood: } P(D|\theta)$$

$$\rightarrow \text{Prior: } p(\theta) \quad \left( \begin{array}{l} \text{objective prior: noninformative prior that attempts to capture ignorance} \\ \text{subjective prior: prior that captures our beliefs} \\ \text{empirical prior: prior that is learned from the data} \end{array} \right)$$

( Wrong prior → worse performance, vague prior would be better )

⇒ conjugate prior yields posterior that is the same family as the prior

(e.g. Beta prior, Bernoulli likelihood → Beta posterior)

$$\rightarrow p(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} \sim \text{Beta}(\beta_H, \beta_T) \text{ where } B(\beta_H, \beta_T) = \frac{\Gamma(\beta_H)\Gamma(\beta_T)}{\Gamma(\beta_H + \beta_T)}$$

$$\rightarrow p(\theta | D) \propto p(D|\theta) \cdot p(\theta)$$

$$\propto \theta^{\alpha_H} (1-\theta)^{\alpha_T} \cdot \theta^{\beta_H-1} (1-\theta)^{\beta_T-1} = \theta^{\alpha_H + \beta_H - 1} \cdot (1-\theta)^{\alpha_T + \beta_T - 1}$$

$$\therefore p(\theta | D) \sim \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T)$$

⇒ prediction forms a distribution, which accounts for uncertainty

$$\rightarrow p(x^* | D) = \mathbb{E}_{p(\theta | D)} [p(x^* | \theta, D)] = \int p(x^* | \theta, D) p(\theta | D) d\theta$$

(Exact posterior inference is difficult due to the intractable evidence term)

→ approximate posterior inference can be used

(This can resolve the scaling issues to large datasets)

① Monte-Carlo-Markov-Chain (MCMC) sampling

: generates samples from Markov chain using algorithms that make it converge to the posterior distribution

② Variational Inference

: find parameters for the parametric distribution that fits the posterior distribution

⇒ Maximum A Posteriori (MAP) estimation

: estimates the most probable parameters  $\theta$  given the data

→ regularization method ⇒ not restricted to Bayesian methods

$$\rightarrow \hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P(\theta | D) = \arg \max_{\theta} \frac{P(D|\theta) \cdot p(\theta)}{P(D)} = \arg \max_{\theta} P(D|\theta) p(\theta)$$

$$(\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} P(D|\theta))$$

⇒ PAC-Bayes

: derive generalization error bound to Bayesian methods for theoretical guarantees  
(convergence, consistency)

\* Discrete variables  $\rightarrow P(X=x|\theta) = \theta^x(1-\theta)^{1-x}, x \in \{0,1\}$

$\Rightarrow$  data points are i.i.d.  $D = \{x_i | i=1,2,\dots,n\}$

### ① Maximum Likelihood Estimation (MLE)

$$\Rightarrow \text{Likelihood} : P(D|\theta) = \prod_i P(x_i|\theta) = \prod_i \theta^{x_i}(1-\theta)^{1-x_i} = \theta^{\alpha_H}(1-\theta)^{\alpha_T} \quad (\alpha_H + \alpha_T = n)$$

$$\Rightarrow \hat{\theta}_{MLE} = \arg \max_{\theta} \log P(D|\theta) = \arg \max_{\theta} \alpha_H \log \theta + \alpha_T \log (1-\theta)$$

$$= \frac{\alpha_H}{\alpha_H + \alpha_T}$$

### ② Maximum A Posteriori (MAP) estimation

$\Rightarrow$  Prior :  $P(\theta) = \text{Beta}(\beta_H, \beta_T)$

$$\Rightarrow \text{Posterior} : P(\theta|D) = \frac{\theta^{\alpha_H + \beta_H - 1}(1-\theta)^{\alpha_T + \beta_T - 1}}{B(\alpha_H + \beta_H, \alpha_T + \beta_T)}$$

$$\Rightarrow \hat{\theta}_{MAP} = \arg \max_{\theta} \log P(\theta|D) = \arg \max_{\theta} (\alpha_H + \beta_H - 1) \log \theta + (\alpha_T + \beta_T - 1) \log (1-\theta)$$

$$= \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H - 1 + \alpha_T + \beta_T - 1}$$

$\rightarrow$  as the number of data points increases, prior is negligible

$$\star \text{ Continuous variables } \rightarrow P(X=x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad x \in \mathbb{R}$$

$\Rightarrow$  affine transformation of gaussian is another gaussian

$$\rightarrow X \sim N(\mu, \sigma^2) \Rightarrow Y = aX + b \sim N(a\mu + b, a^2\sigma^2)$$

$\Rightarrow$  sum of gaussians is another gaussian

$$\rightarrow X \sim N(\mu_X, \sigma_X^2), Y \sim N(\mu_Y, \sigma_Y^2) \Rightarrow X+Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2 + 2\text{Cov}(X, Y))$$

$\Rightarrow$  data points are i.i.d.  $D = \{x_i | i=1, 2, \dots, n\}$

## ① Maximum Likelihood Estimation (MLE)

$$\Rightarrow \text{Likelihood} : P(D|\mu, \sigma) = \prod_i P(x_i|\mu, \sigma) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\Rightarrow \hat{\mu}_{MLE} = \arg \max_{\mu} \log P(D|\mu, \sigma) = \arg \max_{\mu} \left( -\frac{n}{2} \log 2\pi\sigma^2 - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$\rightarrow \sum_{i=1}^n \frac{x_i - \hat{\mu}_{MLE}}{\hat{\sigma}^2} = 0$$

$$\rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Rightarrow \hat{\sigma}_{MLE} = \arg \max_{\sigma} \log P(D|\mu, \sigma) = \arg \max_{\sigma} \left( -\frac{n}{2} \log 2\pi\sigma^2 - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

$$\rightarrow -\frac{n}{\hat{\sigma}_{MLE}^2} + \frac{1}{\hat{\sigma}_{MLE}^3} \sum_{i=1}^n (x_i - \hat{\mu})^2 = 0$$

$$\rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

$$\Rightarrow \mathbb{E}[\hat{\sigma}_{MLE}^2] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (x_i^2 - 2x_i\hat{\mu} + \hat{\mu}^2)\right] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i^2 - \hat{\mu}^2\right]$$

$$= \sigma^2 + \mu^2 - \left(\frac{\sigma^2}{n} + \mu^2\right) = \frac{n-1}{n} \sigma^2 \rightarrow \text{biased}$$

$$\Rightarrow \hat{\sigma}^2_{unbiased} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

## ② Maximum A Posteriori (MAP) estimation

$$\Rightarrow \text{Prior} : \begin{cases} P(\mu) = N(\eta, \lambda) \\ P(\sigma^2) = \text{Wishart distribution} \end{cases}$$

# ML Lec 17

- Naive Bayes

1. Bayesian classification

: given an input with features  $X_1, X_2, \dots, X_D$ , predict its label  $Y$

$$\Rightarrow y = \arg \max_y P(y|X_1, X_2, \dots, X_D) = \frac{P(X_1, X_2, \dots, X_D|y) \cdot P(y)}{P(X_1, X_2, \dots, X_D)}$$

→ assume features are generated by the label

→  $2^D - 1$  parameters are needed to specify  $P(X_1, X_2, \dots, X_D|Y)$  for certain  $Y$

2. Naïve Bayes Classifier

: bayesian classifier based on conditional independence assumption

\*  $X$  is conditionally independent on  $Y$  given  $Z$

$$\Leftrightarrow \forall i, j, k, P(X=x_i | Y=y_j, Z=z_k) = P(X=x_i | Z=z_k)$$

ex) thunder is independent of rain given lightening

\* Naïve Bayes assumption

: features are independent given the label

$$\Rightarrow P(X_1, X_2, \dots, X_D | Y) = \prod_d P(X_d | Y)$$

⇒  $D$  parameters are needed to specify  $P(X_1, X_2, \dots, X_D | Y)$  for certain  $Y$

\* Decision rule

$$\Rightarrow y = \arg \max_Y P(X_1, X_2, \dots, X_D | Y) \cdot P(Y) = \arg \max_Y \prod_d P(X_d | Y) P(Y)$$

→ both  $P(X_d | Y)$  and  $P(Y)$  can be obtained from the training set

## ★ Maximum likelihood estimation (MLE)

$$\textcircled{1} X \sim \text{Multinomial}(\theta) \rightarrow p(X=j|\theta) = \theta_j \text{ where } \sum_{j=1}^J \theta_j = 1$$

$$\Rightarrow \text{Likelihood: } p(D|\theta) = \prod_{i=1}^n p(X_i|\theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \cdots \theta_J^{\alpha_J} \quad (\alpha_1 + \alpha_2 + \cdots + \alpha_J = n)$$

$$\Rightarrow \hat{\theta} = \arg \max_{\theta} \log P(D|\theta) = \alpha_1 \log \theta_1 + \cdots + \alpha_J \log \theta_J \text{ s.t. } \sum_{j=1}^J \theta_j = 1$$

$$= \arg \max_{\theta} \min_{\lambda} \alpha_1 \log \theta_1 + \cdots + \alpha_J \log \theta_J - \lambda \left( \sum_{j=1}^J \theta_j - 1 \right)$$

$$\rightarrow \frac{d_j}{\hat{\theta}_j} - \lambda = 0 \quad \forall j, \quad \sum_{j=1}^J \hat{\theta}_j = 1$$

$$\rightarrow \sum_{j=1}^J \alpha_j = \lambda \quad \forall j$$

$$\rightarrow \hat{\theta}_j = \frac{\alpha_j}{\sum_{j=1}^J \alpha_j} = \frac{\# D(X=j)}{\# D}$$

$$\textcircled{2} X_d | y_k \sim \text{Multinomial}(\theta_{dk}) \rightarrow p(X_d=j|y_k, \theta) = \theta_{dkj} \text{ where } \sum_{j=1}^J \theta_{dkj} = 1$$

$$\Rightarrow \text{Likelihood: } p(D(y_k)|\theta) = \prod_{i=1}^n p(X_{di}|y_{ki}, \theta) = \theta_{dk1}^{\alpha_1} \theta_{dk2}^{\alpha_2} \cdots \theta_{dkJ}^{\alpha_J} \quad (\alpha_1 + \alpha_2 + \cdots + \alpha_J = \# D(y_k))$$

$$\Rightarrow \hat{\theta}_{dk} = \arg \max_{\theta} \log P(D(y_k)|\theta) = \alpha_1 \log \theta_{dk1} + \cdots + \alpha_J \log \theta_{dkJ} \text{ s.t. } \sum_{j=1}^J \theta_{dkj} = 1$$

$$= \arg \max_{\theta} \min_{\lambda} \alpha_1 \log \theta_{dk1} + \cdots + \alpha_J \log \theta_{dkJ} - \lambda \left( \sum_{j=1}^J \theta_{dkj} - 1 \right)$$

$$\rightarrow \frac{d_j}{\hat{\theta}_{dkj}} - \lambda = 0 \quad \forall j, \quad \sum_{j=1}^J \hat{\theta}_{dkj} = 1$$

$$\rightarrow \sum_{j=1}^J \alpha_j = \lambda \quad \forall j$$

$$\rightarrow \hat{\theta}_{dkj} = \frac{\alpha_j}{\sum_{j=1}^J \alpha_j} = \frac{\# D(X_d=j, Y=y_k)}{\# D(y_k)}$$

$$\textcircled{3} Y \sim \text{Multinomial}(\pi) \rightarrow p(y_k|\pi) = \pi_k \text{ where } \sum_{k=1}^K \pi_k = 1$$

$$\Rightarrow \hat{\pi}_k = \frac{\# D(y_k)}{\# D}$$

## \* Maximum A Posteriori (MAP) estimation

①  $X_d | y_k \sim \text{Multinomial}(\theta_{dk}) \rightarrow P(X_d=j | y_k, \theta) = \theta_{dk}$  where  $\sum_{j=1}^J \theta_{dj} = 1$

$$\Rightarrow \text{Prior : } P(\theta) = \text{Dirichlet}(\beta_1, \beta_2, \dots, \beta_J) = \frac{\theta_1^{\beta_1-1} \theta_2^{\beta_2-1} \dots \theta_J^{\beta_J-1}}{B(\beta_1, \beta_2, \dots, \beta_J)}$$

$$\Rightarrow \text{Posterior : } p(\theta | D(y_k)) \propto \prod_{i=1}^n p(X_{di} | y_k, \theta) p(\theta) = \theta_{dk}^{\alpha_1 + \beta_1 - 1} \theta_{dk}^{\alpha_2 + \beta_2 - 1} \dots \theta_{dk}^{\alpha_J + \beta_J - 1} (\alpha_1 + \alpha_2 + \dots + \alpha_J = \# D(y_k))$$

$$\Rightarrow \hat{\theta}_{dk} = \frac{\alpha_j + \beta_j - 1}{\sum_{j=1}^J (\alpha_j + \beta_j - 1)} = \frac{\# D(X_d=j, Y=y_k) + \beta_j - 1}{\# D(y_k) + \sum_{j=1}^J (\beta_j - 1)}$$

②  $Y \sim \text{Multinomial}(\pi) \rightarrow P(y_k | \pi) = \pi_k$  where  $\sum_{k=1}^K \pi_k = 1$

$$\Rightarrow \hat{\pi}_k = \frac{\# D(y_k) + r_k - 1}{\# D + \sum_{k=1}^K (r_k - 1)}$$

## \* Naïve Bayes vs Logistic regression

① Naïve Bayes (NB)  $\rightarrow$  generative

$$\Rightarrow \text{maximize data likelihood} \rightarrow \log P(D|\theta) = \sum_{i=1}^n \log P(x_i, y_i | \theta)$$

② Logistic Regression (LR)  $\rightarrow$  discriminative

$$\Rightarrow \text{maximize conditional data likelihood} \rightarrow \log P(Y|X, \theta) = \sum_{i=1}^n \log P(y_i | x_i, \theta)$$

\* When model is correct, both produce identical classifier

\* When model is incorrect, LR is less biased as conditional independence isn't assumed.

\* When # of data points is small, NB is better as it does not overfit

