

机器学习

一. 基本概念与评估指标.

留出法 (hold-out): 直接划分为 train-validation (首次划分, 重复实验求平均)

交叉验证法 (cross-validation): K 互斥子集划分, 每次取 K-1 个子集训练, 行验证, 重复 K 次, 取均值. K=10! 准确但计算开销高

真实 正 反

正 TP FN

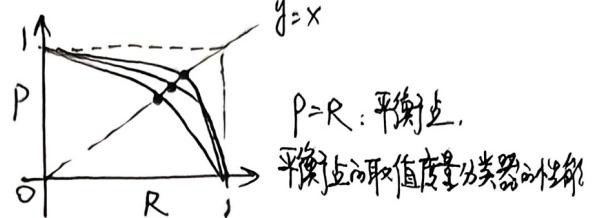
反 FP TN

(假正例)
0.1m

$$\text{Precision} = P = \frac{TP}{TP+FP}$$

$$\text{Recall} = R = \frac{TP}{TP+FN}$$

混淆矩阵.



$$F_1 = \frac{2PR}{P+R} = \frac{2PR}{P+R}, F_\beta = \frac{(1+\beta^2)PR}{\beta^2 P + R}$$

$$\left\{ \begin{array}{l} \beta > 1 \rightarrow R \uparrow (P=R) \\ \beta < 1 \rightarrow P \uparrow (P=R) \end{array} \right.$$

[Definition] ML: 通过经验 (experience) 提升性能表现学习系统.

分类: $\begin{cases} \text{监督: label-classification/regression} & x \rightarrow y \\ \text{无监督: special info-clustering (interesting patterns)} & \end{cases}$

强化: max reward — 强化学习. label exists, credit assignment.

二. 贝叶斯分类:

Naive Bayesian Classifier:

Given $x = (x_1, \dots, x_p)^T$, predict w . ($\max p(w|x) \propto p(x|w)p(w)$).

Eg. $X = (\text{Refund}=\text{No}, \text{Married}, \text{Income}=120k)$

1. compute $p(\text{Refund}=?|\text{No}/\text{Yes})$, $p(\text{Marital Status}=?|\text{No}/\text{Yes})$, sample mean & variance of $\overset{\mu_{ij}}{\text{Income}} |\text{No}/\text{Yes}$.

2. compute $p(X|\text{No}) = \underset{\text{Independence}}{p(\text{Refund}=\text{No}|\text{No}) \times p(\text{Marital}(\text{No}) \times p(\text{Income}=120k|\text{No}))}, p(X|\text{Yes})$

3. compare $p(X|\text{No})p(\text{No})$ and $p(X|\text{Yes})p(\text{Yes})$

$$\rightarrow \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp\left(-\frac{(x_i - \mu_{ij})^2}{2\sigma_{ij}^2}\right)$$

4. the bigger one is the output.

三. MLE

Regression model: $y = f(x, w) + \epsilon$, $\epsilon \sim N(0, \sigma^2)$. $p(y|x, w, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y-f(x, w))^2}$

$$L(D, w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma), w^* = \arg\max L$$

$$\ell(D, w, \sigma) = \log L = \sum_{i=1}^n \log p(y_i|x_i, w, \sigma) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i, w))^2 + C(\sigma).$$

$$\ell'(D, w, \sigma) = 0 \Leftrightarrow \sum_{i=1}^n (y_i - f(x_i, w))^2 = \text{RSS}(f) = 0. J_n = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 = \text{MSE}$$

Residual Sum of Squares 残差平方和

$$(x_i, y_i), f = \sum_{i=1}^m w_i x_i = w^T x, w = [w_1, w_2, \dots, w_m]^T$$

$$\text{MSE} = J_n = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2 \Rightarrow J_n(w) = (y - X^T w)^T (y - X^T w)$$

$$\nabla J_n = -2X(y - X^T w) = 0 \Rightarrow Xy = X X^T w \Rightarrow w = (X X^T)^{-1} X y.$$

$$\hat{y} = X^T w = X^T (X X^T)^{-1} X y.$$

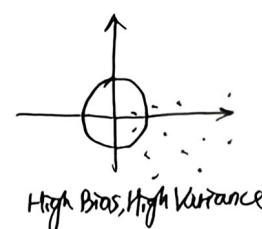
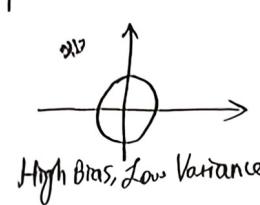
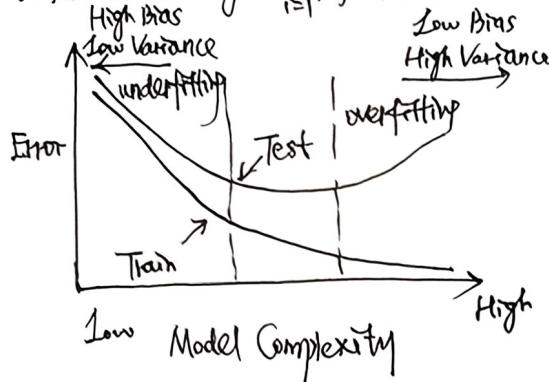


推导①: $w^* = \arg\min \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \sum_{j=1}^p w_j^2$. $\Leftrightarrow a^* = \arg\min \sum_{i=1}^n (y_i - x_i^T w)^2$ subject to $\sum_{j=1}^p w_j^2 \leq t$.

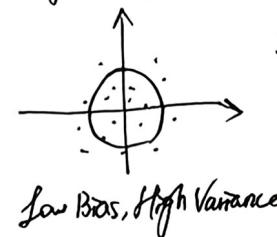
$$J_n = (y - X^T w)^T (y - X^T w) + \lambda w^T w.$$

$$\nabla J_n = -2X(y - X^T w) + 2\lambda w = 0 \Rightarrow (XX^T + \lambda I) w = Xy \Rightarrow w^* = (XX^T + \lambda I)^{-1} Xy.$$

1 LASSO: $\hat{w} = \arg\min \sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda \|w\|_1$, Sparse Model



(large λ): over-regularized
high bias



small λ : under-regularize
(overfitting)
high variance

Classifier:

regression: $f(x) = w^T x \in \mathbb{R}$.

sigmoid / logistic function $\sigma(x) = \frac{1}{1+e^{-x}} \in (0, 1)$

$P(Y_i = \pm 1 | X_i, \alpha) = \sigma(Y_i \alpha^T X_i) = \frac{1}{1+e^{-Y_i \alpha^T X_i}}$

MLE: $L = \prod_{i=1}^n \sigma(Y_i \alpha^T X_i)$

$$l = \log L = \sum_{i=1}^n \log \sigma(Y_i \alpha^T X_i) = \sum_{i=1}^n \log \frac{1}{1+e^{-Y_i \alpha^T X_i}} = - \sum_{i=1}^n \log(1+e^{-Y_i \alpha^T X_i}) = E(\alpha). \text{ Gradient Descent}$$

$$\min J(w). \quad w_{\text{new}} = w_n - \gamma \nabla J(w_n).$$

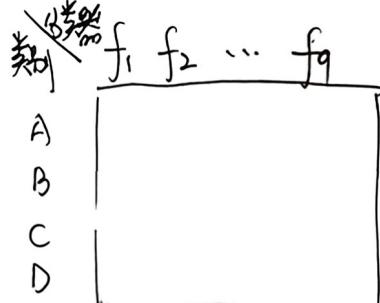
Multi-Class Classifier.

One vs Rest: $\{f_i\}_{i=1}^K$ 分类器

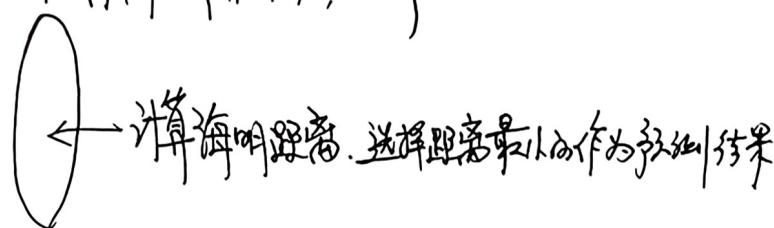
One vs One: $\binom{K}{2}$ 个分类器

ECOC: 多分类问题转化为二分类问题。海明距离最小，任意两个类别之间海明距离最大。

矩阵表示: $A = \{000000000\}, B = \{00011111\}, C = \{111000111\}, D = \{111111000\}$.



测试样本: $\{1, 0, 1, 0, \dots\}$



SVM: $f(x) = w^T x + b = 0$.

$$f(x_p) = w^T (x - \frac{w}{\|w\|}) + b = 0. \quad f(x) = w^T x + b = \|w\|. \Rightarrow r = \frac{|f(x)|}{\|w\|}$$

$$\gamma = y \cdot r = y \frac{\|w\|}{\|w\|} = y. \text{ geometrical margin}$$

Maximum Margin Classifier: $\max_{w,b} \gamma = \max_{w,b} \frac{y_i(w^T x_i + b)}{\|w\|}, \text{ s.t. } y_i > \gamma.$

$$\text{fix } y_i(w^T x_i + b) = 1. \max_{w,b} \gamma \Leftrightarrow \min_{w,b} \frac{1}{2} \|w\|^2 \text{ s.t. } y_i - \frac{y_i(w^T x_i + b)}{\|w\|} \geq \gamma$$

$$y_i - \frac{y_i(w^T x_i + b)}{\|w\|} \geq \gamma \Leftrightarrow y_i(w^T x_i + b) \geq \gamma$$

MMC: $\min_{w,b} \frac{1}{2} \|w\|^2, \text{ s.t. } y_i(w^T x_i + b) \geq \gamma.$

支持向量: margin 上的点, 非支持向量: margin 不满足的支持向量.

Slack variables: $\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i, \text{ s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0.$ Cf. $\xi \downarrow, \text{ overfitting}$. Cf. $\xi \uparrow, \text{ underfitting}$

核方法: 得到线性不可分的问题且避免维数灾难 核函数: 直接计算高维内积.

五. 神经网络.

1. 知识. $f(x) = w^T x + b$. Predict wrong, set $w_{k+1} = w_k + xy$. ($w^T xy > 0$)

OR: $f(x) = x_1 + x_2 + 0.5$, AND: $f(x) = x_1 + x_2 - 1.5$, XOR: x

$$L(w, b) = \sum_{x \in M} -y_i(w^T x_i + b), \quad \nabla L = \sum_{x \in M} -x_i y_i \quad w_{k+1} = w_k + \eta \sum_{x \in M} x_i y_i \quad (\text{gradient descent})$$

2. 反向传播:

forward: $Z = f(\sum_j w_{kj} f(\sum_i u_{ji} x_i + u_{j0}) + a_{k0})$

Activation function

$$\Delta w = -\eta \frac{\partial L}{\partial w}. \quad w_{k+1} = w_k - \eta \frac{\partial L}{\partial w}. \quad \frac{\partial L}{\partial w_{ki}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial Z} \cdot \frac{\partial Z}{\partial a_i} \cdot \frac{\partial a_i}{\partial w_{ki}}$$

例: Input (x_1, x_2) . hidden (h_1, h_2) output \hat{y} , $f = \text{sigmoid}$, $L = \text{MSE} = \frac{1}{2}(y - \hat{y})^2$.

$$X = [x_1, x_2] = [0.5, -0.8], \quad y = 1. \quad w^1 = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \end{bmatrix}, \quad b^1 = \begin{bmatrix} 0.1 & -0.1 \end{bmatrix}, \quad w^2 = \begin{bmatrix} 0.2 & -0.5 \end{bmatrix}, \quad b^2 = 0.05.$$

$$\text{forward: } z_1^1 = w_{11}^1 x_1 + w_{12}^1 x_2 + b_1^1 = 0.31. \quad a_1^1 = \sigma(0.31) = 0.58$$

$$z_2^1 = w_{21}^1 x_1 + w_{22}^1 x_2 + b_2^1 = -0.27. \quad a_2^1 = \sigma(-0.27) = 0.43.$$

$$z^2 = w_1^2 a_1^1 + w_2^2 a_2^1 + b^2 = -0.05, \quad a_1^2 = \sigma(-0.05) = \hat{y} = 0.49. \quad (\hat{y} = \frac{1}{1+e^{-z}})$$

$$\text{backward: } \frac{\partial L}{\partial \hat{y}} = \hat{y} - y = -0.51. \quad \frac{\partial \hat{y}}{\partial Z} = \hat{y}(1-\hat{y}) = 0.25. \quad \delta^2 = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} = -0.13$$

$$\frac{\partial L}{\partial w^2} = \delta^2 \cdot a_1^1 = -0.07, \quad \frac{\partial L}{\partial w^2} = \delta^2 \cdot a_2^1 = -0.06. \quad \frac{\partial L}{\partial b^2} = \delta^2 = -0.13.$$

$$\frac{\partial Z}{\partial a_1^1} = w_1^2 = 0.2, \quad \frac{\partial Z}{\partial a_2^1} = w_2^2 = -0.5, \quad \frac{\partial a_1^1}{\partial w_{11}^1} = a_1^1(1-a_1^1) = 0.24, \quad \frac{\partial a_2^1}{\partial w_{21}^1} = 0.25.$$

$$\delta_1^1 = \delta^2 \times w_1^2 \times \frac{\partial a_1^1}{\partial w_{11}^1} = -0.01, \quad \delta_2^1 = \delta^2 \times w_2^2 \times \frac{\partial a_2^1}{\partial w_{21}^1} = 0.02.$$

$$\frac{\partial L}{\partial w_{11}^1} = \delta_1^1 \times x_1, \quad \frac{\partial L}{\partial w_{21}^1} = \delta_2^1 \times x_1, \quad \frac{\partial L}{\partial b_1^1} = \delta_1^1, \quad \frac{\partial L}{\partial w_{12}^1} = \delta_1^1 \times x_2, \quad \frac{\partial L}{\partial w_{22}^1} = \delta_2^1 \times x_2, \quad \frac{\partial L}{\partial b_2^1} = \delta_2^1.$$

$$w_{11}^2_{\text{new}} = w_{11}^2 - \eta \frac{\partial L}{\partial w_{11}^2}, \quad w_{21}^2_{\text{new}} = w_{21}^2 - \eta \frac{\partial L}{\partial w_{21}^2}, \quad b_1^2_{\text{new}} = b_1^2 - \eta \frac{\partial L}{\partial b_1^2}.$$

$w_{11}^1, w_{12}^1, w_{21}^1, w_{22}^1, b_1^1, b_2^1$ 同样更新.



激活函数: non-linear, sigmoid, tanh.

$$\frac{e^x - e^{-x}}{e^x + e^{-x}} \in (-1, 1) = \text{sigmoid}(x) - 1$$

vanishing gradient \Rightarrow ReLU, dropout.
(-overfitting)

六. 决策树:

$$ID3. \text{ Entropy}(S) = -\sum_{i=1}^C p_i \log p_i. \text{ Gain}(S, n) = \text{Entropy}(S) - \sum_{i=1}^{|S_n|} \frac{|S_n|}{|S|} \text{Entropy}(S_n).$$

ID	Age	Gender	Income	Buy?	Age = 20-30, 30-40, 40+
1	27	M	15W	X	
②	47	F	30W	✓	
3	32	M	12W	X	Income = 10-20, 20-40, 40+
4	24	M	45W	✓	
*⑤	45	M	30W	X	
⑥	56	M	32W	✓	$a \in \{\text{Age, Gender, Income}\}$
7	31	M	15W	X	
⑧	23	F	30W	✓	

(4+, 4-)

$$\text{Entropy}(S) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1.$$

$$\text{Age: } 20-30: 2+, 1-. E = -\frac{1}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.918$$

$$30-40: 0+, 2-. E = 0$$

$$40+: 2+, 1-. E = 0.918$$

$$\text{Gain}(S, \text{Age}) = 1 - \left(\frac{3}{8} \times 0.918 + \frac{2}{8} \times 0 + \frac{3}{8} \times 0.918 \right) = 0.312.$$

$$\text{Income: } 10-20: 0+, 3-. E = 0$$

$$20-40: 3+, 1-. E = 0.811$$

$$40+: 1+, 0-. E = 0$$

$$\text{Age: } 20-30,$$

30-40: 无法划分.

40+:

$$\text{Age: } 20-30: 1+, 0-. E = 0$$

$$30-40: /$$

$$40+: 2+, 1-. E = 0.918$$

$$\text{Gain} = 0.811 - \left(\frac{1}{4} \times 0 + \frac{3}{4} \times 0.918 \right) = 0.$$

$$\text{Gender: } M: 1+, 1-. E = 1$$

$$F: 2+, 0-. E = 0$$

$$\text{Gain} = 0.811 - \frac{2}{4} \times 1 = 0.311$$

CART (Classification and Regression Tree): partition the input space for continuous features (greedy method)

$$\min_{\text{feature } j, \text{ threshold}} \left[\min_{C_1} \sum_{i \in C_1} (y_i - c_1)^2 + \min_{C_2} \sum_{i \in C_2} (y_i - c_2)^2 \right]$$

t. Bagging & Boosting / Ensemble Method

1. Bagging

选择决策树作为 base learner: 1. non-linear, 2. easy to use, 3. easy to overfit.

定义: Bagging 是一种 并行的聚类, 通过重采样投票(分类) / 取平均(回归)来达到预测结果

Random Forest: 1. random sampling

2. randomized feature selection at each split step

no pruning for overfitting

vote for prediction

4.



2. Boosting.

定义：顺序聚合，每步提高被错分样本的权重。

$$\text{AdaBoost: } 1. w_n^{(i)} = \frac{1}{N}$$

2. for $m=1, \dots, M$:

$$3. \text{ fit classifier } y^{(m)}(x), J_m = \sum_{n=1}^N w_n^{(m)} I(y^{(m)}(x_n) \neq t_n)$$

$$4. \Sigma_m = \frac{J_m}{\sum_{n=1}^N w_n^{(m)}}, \alpha_m = \ln \frac{1 - \Sigma_m}{\Sigma_m}$$

5. update w_n^{m+1}

$$6. Y_M(x) = \text{sign}(\sum_{n=1}^M \alpha_n y^{(n)}(x)) \# \text{final model}$$

Gradient Boosting Decision Tree (GBDT): predict remaining data & sum

3. KNN.

$\begin{cases} \text{stored record} \\ \text{distance metrics} \end{cases} \Rightarrow \begin{cases} 1. \text{compute the distance to other record} \\ 2. \text{find the } k \text{ nearest neighbor} \\ 3. \text{vote} \end{cases}$

4. Clustering

K-Means: 1. Randomly pick k data points as μ_k

2. Iterate until converge:

3. $\bar{x}_i = \arg \min_k \|x_i - \mu_k\|^2$, # assign each point to the cluster

4. $\mu_k = \frac{1}{N_k} \sum_{i=1}^{N_k} x_i^{(k)}$ # update seed point.

5. dimensionality reduction.

1. PCA (unsupervised)

1stPC: $\bar{x}_i^{(1)} = \alpha_1^T x_i$, s.t. $\max \text{Var}(\bar{x}^{(1)})$

Solution: $\text{Var}(\bar{x}^{(1)}) = E((\bar{x}^{(1)} - \bar{\bar{x}}^{(1)})^2) = \frac{1}{n} \sum_{i=1}^n (\alpha_1^T x_i - \alpha_1^T \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n \alpha_1^T (x_i - \bar{x})(x_i - \bar{x})^T \alpha_1 = \alpha_1^T S \alpha_1$.

$\Rightarrow \max_{\alpha_1} \alpha_1^T S \alpha_1$, s.t. $\alpha_1^T \alpha_1 = 1$. $L = \alpha_1^T S \alpha_1 - \lambda(\alpha_1^T \alpha_1 - 1) \cdot \frac{\partial L}{\partial \alpha_1} = 0 \Rightarrow \lambda \alpha_1 = S \alpha_1$,

$\Rightarrow \alpha_1$ is the eigenvector of S corresponds to the largest eigenvalue λ_1 , $S = \frac{1}{n} \sum (x_i - \bar{x})(x_i - \bar{x})^T$

2ndPC: $\max_{\alpha_2} \alpha_2^T S \alpha_2$, s.t. $\alpha_2^T \alpha_1 = 0$, $\text{Cov}(\bar{x}^{(1)}, \bar{x}^{(2)}) = 0$.

[例]: Given $(90, 85, 80); (60, 65, 70); (85, 80, 75); (70, 75, 80); (95, 90, 85)$. $\bar{x} = (80, 79, 78)$

$\Rightarrow (10, 6, 2), (-20, -14, -8), (5, 1, -3), (-10, -4, 1), (15, 11, 7)$

$$\Rightarrow S = \frac{1}{5} \left(\begin{pmatrix} 10 \\ 6 \\ 2 \end{pmatrix} (10, 6, 2) + \begin{pmatrix} -20 \\ -14 \\ -8 \end{pmatrix} (-20, -14, -8) + \dots \right) = \begin{pmatrix} 170 & 110 & 50 \\ 110 & 74 & 38 \\ 50 & 38 & 26 \end{pmatrix}$$

$$(S - \lambda I) = \begin{pmatrix} 170 - \lambda & 110 & 50 \\ 110 & 74 - \lambda & 38 \\ 50 & 38 & 26 - \lambda \end{pmatrix} = 0$$

$$(S - \lambda_1 I) v_1 = 0$$

$$v_1 = \begin{pmatrix} 0.805 \\ 0.533 \\ 0.260 \end{pmatrix} = \vec{v}_1$$

$$\vec{v}_1 = (x_i - \bar{x}) \cdot v_1$$



2. LDA (supervised).

$$\text{Rayleigh quotient} = J(\alpha) = \frac{\alpha^T S_B \alpha}{\alpha^T S_W \alpha}, S_W: \text{类内散度}, S_B: \text{类间散度}$$

$$\max J(\alpha) \Leftrightarrow S_W^{-1} S_B$$

例1: JD
类别 故 活

	1	2	3	4	5	6	
1	雄	90	70	}	$\mu_1 = [85 \ 65]$	$\mu = [75 \ 75]$	
2	雌	85	65				
3	老	80	60	}	$\mu_2 = [65 \ 85]$		
4	少	60	85				
5	少	65	90				
6	少	70	80				

$$S_W: \text{对角} = \begin{pmatrix} 5 & 0 \\ 0 & 5 \\ -5 & -5 \end{pmatrix}, S_1 = \sum (X - \mu_1)(X - \mu_1)^T = \begin{pmatrix} 5 \\ 5 \end{pmatrix}(5 \ 5) + \begin{pmatrix} 0 \\ 0 \end{pmatrix}(0 \ 0) + \begin{pmatrix} -5 \\ -5 \end{pmatrix}(-5 \ -5) = \begin{pmatrix} 50 \\ 50 \end{pmatrix}$$

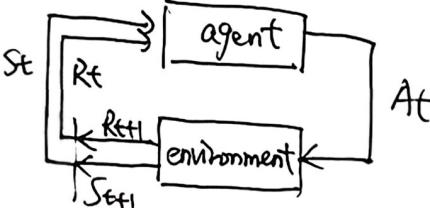
$$2: \begin{pmatrix} -5 & 0 \\ 0 & 5 \\ 5 & -5 \end{pmatrix} S_2 = \begin{pmatrix} 50 & -25 \\ -25 & 50 \end{pmatrix} \quad S_W = \begin{pmatrix} 100 & 25 \\ 25 & 100 \end{pmatrix}$$

$$S_B = S_B = \sum_{i=1}^3 n_i (\mu_i - \mu)(\mu_i - \mu)^T = 3 \begin{pmatrix} 10 \\ -10 \end{pmatrix}(10 \ -10) + 3 \begin{pmatrix} -10 \\ 10 \end{pmatrix}(-10 \ 10) = \begin{pmatrix} 600 & -600 \\ -600 & 600 \end{pmatrix}$$

$$\boxed{[S_B w = \lambda S_W w \Rightarrow \begin{pmatrix} 600 & -600 \\ -600 & 600 \end{pmatrix} w = \lambda \begin{pmatrix} 100 & 25 \\ 25 & 100 \end{pmatrix} w]} \quad \boxed{w = \begin{pmatrix} 1 \\ -1 \end{pmatrix}}$$

+. 强化学习.

agent-environment interaction



$$(State, Action, Reward) = (S, A, R)$$

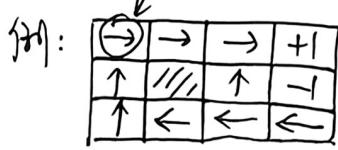
goal: optimal policy to max reward R.

π : policy $\pi(s)$
 reward function
 value function $V_\pi(s)$

Markov Decision Process (MDP): $P_r(S_{t+1}=s' | S_t, A_t, R_t, S_{t+1}, A_{t+1}, \dots, S_0, A_0) = P_r(S_{t+1}=s' | S_t, A_t)$

$$\text{Bellman equation: } q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$$

$$V_\pi(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$$



s	r	$p(s', r s, a)$
↑	1	0.8
↓	-1	0.1
←	-1	0.1

$$\gamma = 0.8 \quad \checkmark$$

$$\gamma = 0.1 \leftarrow, 0.1 \rightarrow$$

$$-0.04 \text{ each step}, \gamma = 1$$

$$V_\pi(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$$

$$= 0.8 \times 1 \times (-0.04 + 1 \times 0.868) + 0.1 \times 1 \times (-0.04 + 1 \times 0.812) + 0.1 \times 1 \times (-0.04 + 0.868)$$

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$$

$$= 1 \times (-0.04 + 1 \times 0.868)$$

6.



$$DP: V_{\pi}(s_t) = \mathbb{E}_{\pi}[R_{t+1} + \gamma V(s_{t+1})]$$

$$\text{policy evaluation: } V_{k+1}(s) = \sum_a \left(\pi(a|s) \sum_{s', r} p(s', r|s, a) (r + \gamma V_k(s')) \right) \quad \{V_k(s)\} \Rightarrow \{V_{k+1}(s)\}.$$

ii improvement: $Q_{\pi}(s, \pi'(s)) > V_{\pi}(s) \Leftrightarrow V_{\pi'}(s) > V_{\pi}(s)$.

$$\pi'(s) = \arg \max_a \sum_{s', r} p(s'|s, a) (r + \gamma V_{\pi'}(s')) : \text{greedy}$$

" $q(s, a)$ "

$$\text{Monte Carlo: } Q_{\pi}(s, a) = \mathbb{E}[G_t]. \quad V(s_t) \leftarrow V(s_t) + \alpha \underbrace{(G_t - V(s_t))}_{\text{return at the end}}$$

exploration starts: first S-A pair

return at the end

$$\text{Temporal-Difference Learning: } V(s_t) \leftarrow V(s_t) + \alpha \underbrace{(R_{t+1} + \gamma V(s_{t+1}) - V(s_t))}_{\text{return in one step}}$$

Sarsa: ϵ -greedy: \sum for exploration (all actions)

$1 - \sum$ for exploitation ($\pi(s)$)

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

Q-learning: fully-greedy exploration. $\pi(s) = \arg \max_a Q(s, a)$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

