

似然函数 $P(Y = y) = \pi^y(1 - \pi)^{1-y}$.

则 $f(\mathbf{y} | \mathbf{w}, b) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n (\pi(x_i))^{y_i} (1 - \pi(x_i))^{1-y_i}$

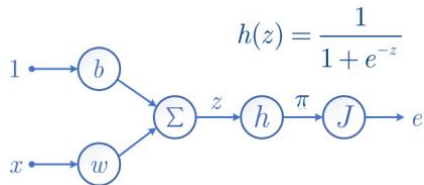
最大似然 $\max \sum_{i=1}^n (y_i \log \pi(x_i) + (1 - y_i) \log (1 - \pi(x_i)))$

信息量 $I = -\log_2 p$. 信息熵 $H = -\sum_{i=1}^n p_i \log_2 p_i$

交叉熵 $H(p, q) = -\sum_{i=1}^n p_i \log_2 q_i$

交叉熵损失: $-\frac{1}{n} \sum_{i=1}^n (y_i \log \pi(x_i) + (1 - y_i) \log (1 - \pi(x_i)))$

其中 $\pi(x_i) = \frac{e^{w x_i + b}}{1 + e^{w x_i + b}}$. 注意最大化对数似然和最小化交叉熵损失是等价的.



从输入算输出 $z = wx + b$, $\pi = h(z) = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}$
 $e = J(\pi) = -y \log \pi - (1 - y) \log (1 - \pi)$

从输出算梯度:

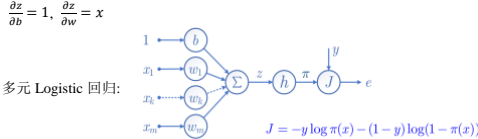
$$\frac{\partial e}{\partial b} = \frac{de}{d\pi} \frac{d\pi}{dz} \frac{dz}{db} = \frac{\pi - y}{\pi(1-\pi)} \pi(1-\pi) = \pi - y;$$

$$\frac{\partial e}{\partial w} = \frac{de}{d\pi} \frac{d\pi}{dz} \frac{dz}{dw} = \frac{\pi - y}{\pi(1-\pi)} \pi(1-\pi)x = (\pi - y)x$$

$$\frac{de}{d\pi} = \frac{d}{d\pi} J(\pi) = -\frac{y}{\pi} + \frac{1-y}{1-\pi} = \frac{\pi - y}{\pi(1-\pi)}$$

$$\frac{d\pi}{dz} = \frac{d}{dz} h(z) = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}+e^{-2z}} = \pi(1-\pi)$$

$$\frac{dz}{db} = \frac{d}{db} z = 1, \frac{dz}{dw} = x$$



$$\frac{\partial e}{\partial b} = \frac{de}{d\pi} \frac{d\pi}{dz} \frac{dz}{db} = \pi - y; \quad \frac{\partial e}{\partial w_k} = \frac{de}{d\pi} \frac{d\pi}{dz} \frac{dz}{dw_k} = (\pi - y)x_k$$

二分类问题的评价

判断为正样本如果 $\pi(x_i) = \frac{e^{a_1 x_i + b}}{1 + e^{a_1 x_i + b}} > c$

TN: 本来就是错的, 预测也是错的. FP: 本来是错的, 预测认为是对的

Sensitivity = True positive rate = TP / (TP+FN)

Specificity = True negative rate = TN / (TN+FP)

1-Sensitivity = False negative rate = FN / (TP+FN)

1-Specificity = False positive rate = FP / (TN+FP)

Recall = TP / (TP+FN); Precision = TP / (TP+FP)

Fall-out = FP / (TP+FP); F1-measure = 2×Precision× Recall / (Precision + Recall)

$$ACC = \frac{TP+TN}{TP+TN+FP+FN}$$

ROC: 横轴: 1-Specificity; 纵轴: Sensitivity

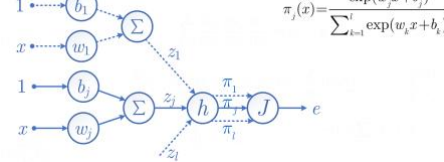
Softmax 回归: $P(Y = j) = \pi_j(x) = \frac{e^{w_j x + b_j}}{\sum_{k=1}^K e^{w_k x + b_k}}$

独热码:如果该观测属于第j类, 则其独热码向量的第 j 维 v_{ij} 为 1, 其他维为 0

softmax $(j, z_1, \dots, z_i) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$ 对数似然函数 $\sum_{i=1}^n \sum_{j=1}^K y_{ij} \log \pi_j(x_i)$

最大化对数似然 $\max \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log \pi_j(x_i)$

最小化交叉熵损失 $\min -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log \pi_j(x_i)$. 两者是一致的.



$$e = -\sum_{k=1}^K y_k \log \pi_k(J(\mathbf{w}, \mathbf{b})) = -\sum_{k=1}^K y_k \log \pi_k$$

$$\frac{\partial e}{\partial b_i} = -\sum_{k=1}^K y_k \frac{\partial \pi_k}{\partial \pi_i} \frac{\partial \pi_i}{\partial z_i} \frac{\partial z_i}{\partial b_i} = -\frac{y_i}{\pi_i} \pi_i + \sum_{k=1}^K \frac{y_k}{\pi_k} \pi_k \pi_i = \pi_i - y_i$$

$$\frac{\partial e}{\partial w_i} = -\sum_{k=1}^K y_k \frac{\partial \pi_k}{\partial \pi_i} \frac{\partial \pi_i}{\partial z_i} \frac{\partial z_i}{\partial w_i} = (\pi_i - y_i)x$$

$$\frac{\partial \pi_k}{\partial z_i} = \frac{de^{\pi_k}}{dz_i} (\sum_{j=1}^K e^{z_j})^{-1} + e^{z_k} \frac{\partial}{\partial z_i} (\sum_{j=1}^K e^{z_j})^{-1} = \begin{cases} \pi_i(1 - \pi_i) & k = i \\ -\pi_i \pi_k & k \neq i \end{cases}$$

$$\frac{\partial e_k}{\partial \pi_k} = \frac{y_k}{\pi_k}; \quad \frac{\partial z_i}{\partial w_i} = 1; \quad \frac{\partial z_i}{\partial w_j} = x.$$

前馈神经网络:输入单元—隐层单元—输出单元

Logistic 单元: $z = \mathbf{w}^T \mathbf{x} + b$. $a = \frac{1}{1 + e^{-z}}$

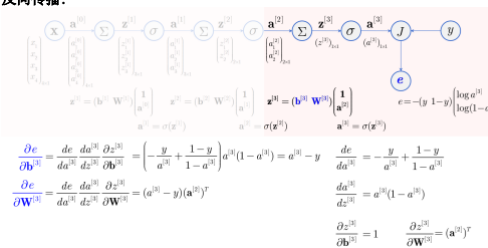
计算梯度: $\frac{da}{dz} = a(1 - a)$; $\frac{\partial z}{\partial b} = 1$; $\frac{\partial w}{\partial w} = \mathbf{x}^T$

$$\frac{\partial a}{\partial b} = \frac{da}{dz} \frac{dz}{db} = a(1 - a); \quad \frac{\partial a}{\partial w} = \frac{da}{dz} \frac{dz}{dw} = a(1 - a)\mathbf{x}^T$$

双曲正切单元: $a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$; $\frac{da}{dz} = 1 - a^2 = (1 + a)(1 - a)$

ReLU 单元: $a = \max(0, z)$; $\frac{da}{dz} = \begin{cases} 1 & z \geq 0 \\ 0 & \text{otherwise} \end{cases}$

反向传播:



卷积神经网络

卷积: $g(x, y) * h(x, y) = \sum_i \sum_j g(i, j) h(i, j)$

填充后卷积: 卷积结果与原图大小一致.

特点一: 稀疏连接; 特点二: 参数共享; 特点三: 等变表示.

多通道卷积: 图像有多个颜色通道,各通道分别进行卷积,再把卷积的结果相加.

一个卷积核提取一个特征.多个卷积核产生多个特征,每个特征对应于一个通道.

原始图 $N \times N$, 卷积核 $m \times m$, 卷积核 $(N - m + 1) \times (N - m + 1)$.

池化: 最大池化, 平均池化, 随机池化.

令 O=输出图像的尺寸, I=输入图像的尺寸, K=卷积层的核尺寸, S=移动步长, P=填充数,则输出图像尺寸

$$O = \frac{I - K + 2P}{S} + 1$$

严格来说这个尺寸是要取整的.

强化学习

马尔可夫过程: $P(S_{t+1} | S_t, S_{t-1}, \dots, S_1) = P(S_{t+1} | S_t)$

考虑两个相邻时刻, 定义状态转移概率: $p_{ss'} = P(S_{t+1} = s' | S_t = s)$

状态转移矩阵: 从行转移到列, 每行之和为 1. 马尔可夫过程由二元组 (S, P) 描述, 其中 $\mathbf{S} = \{s_1, \dots, s_n\}$ 称为状态空间, $\mathbf{P} = (p_{ss'})_{n \times n}$ 称为状态转移矩阵

马尔可夫回报过程

由四元组 (S, P, r, γ) 描述. 状态空间 $\mathbf{S} = \{s_1, \dots, s_n\}$. 状态转移矩阵 $\mathbf{P} = (p_{ss'})_{n \times n}$.

状态期望回报: $\mathbf{r} = (r_1, \dots, r_n)$. $r_s = E[R_{t+1} | S_t = s]$ 折现因子: $\gamma \in [0, 1]$

Return: 从当前时刻开始的折现回报之和

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

状态价值就是状态累积回报的期望值

$$v(s) = E[G_t | S_t = s] = E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

贝尔曼期望方程

所有状态 (矩阵形式): $\mathbf{v} = \mathbf{r} + \gamma \mathbf{Pv}$

求解线性方程组, 可得 $\mathbf{v} = (\mathbf{I} - \gamma \mathbf{P})^{-1} \mathbf{r}$

马尔可夫决策过程:五元组 (S, A, P, R, γ)

多了一个行动空间: $\mathbf{A} = \{a_1, \dots, a_m\}$.

策略: $\pi(a | s) = P(A_t = a | S_t = s)$

状态转移概率 $P^\pi(S_{t+1} = s' | S_t = s) = \sum_{a \in A} \pi(a | s) p_{ss'}^a$

状态价值函数 $v_\pi(s) = E[G_t | S_t = s]$.

行动价值函数 $q_\pi(s, a) = E[G_t | S_t = s, A_t = a]$

$$q_\pi(s, a) = r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_\pi(s')$$

状态价值的贝尔曼期望方程

$$v_\pi(s) = \sum_{a \in A} \pi(a | s) \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_\pi(s') \right)$$

行动价值的贝尔曼期望方程

$$q_\pi(s, a) = r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a \sum_{a' \in A} \pi(a' | s') q_\pi(s', a')$$

动态规划—策略评价(状态价值计算)、策略改进

策略评价: $v_{k+1}(s) = r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_k(s')$

策略的好坏用状态价值来评价.

最优策略: $\pi_* \geq \pi'$: $v_{\pi_*}(s) \geq v_{\pi'}(s)$, $\forall s, \forall \pi'$

最优状态价值即最优策略下的状态价值 $v_*(s) = \max_{\pi} v_\pi(s)$

最优行动价值即最优策略下的行动价值.

最优状态价值即同时刻最优行动价值 $v_*(s) = \max_{a \in A} q_*(s, a)$

策略改进: 基于原策略, 产生新策略

如果 $v_\pi(s) \leq q_\pi(s, \pi^*(s))$, 那么 $v_\pi(s) \leq v_{\pi^*}(s)$.

策略迭代: 交替进行策略评价和策略改进.

价值迭代

策略改进计算量很大, 故每次策略评价后立即进行策略改进.

状态价值贝尔曼最优方程

$$v_*(s) = \max_{a \in A} \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_*(s') \right)$$

行动价值贝尔曼最优方程

$$q_*(s, a) = r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a \max_{a' \in A} q_*(s', a')$$

价值迭代: $v_{k+1}(s) = \max_{a \in A} \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_k(s') \right)$

则策略: $\pi_*(s) = \arg \max_{a \in A} \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_*(s') \right)$

同步迭代: 算完所有状态后一次更新, 保存两份状态价值, 计算量大收敛较慢.

$$v_{k+1}(s) = \max_{a \in A} \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v_k(s') \right)$$

异步迭代: 算一个更新一个, 只有一份状态价值, 计算量小收敛较快.

$$v(s) = \max_{a \in A} \left(r_s^a + \gamma \sum_{s' \in S} p_{ss'}^a v(s') \right)$$

蒙特卡洛预测

首次访问蒙特卡洛: 根据待评价策略产生观测片段, 若是第一次出现, 计算 $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t} R_{T+1}$, 重复 n 次, 得到 G_1, G_2, \dots, G_n , 则估计状态价值 $V(S_t) = \frac{1}{n} \sum_{i=1}^n G_i$

每次访问蒙特卡洛.

增量式蒙特卡洛预测: $V(S_t) \leftarrow V(S_t) + \frac{1}{k} (G_t - V(S_t))$, $k \leftarrow k + 1$

定步长蒙特卡洛预测: $V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$

行动价值预测: $Q(S_t, A_t) = \frac{1}{n} \sum_{i=1}^n G_i$

增量式预测: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{k} (G_t - Q(S_t, A_t))$, $k \leftarrow k + 1$

定步长预测: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (G_t - Q(S_t, A_t))$

探索与利用:

$$\pi(a | s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{m} & \text{if } a = \arg \max_{a \in A} Q(s, a) \\ \frac{\epsilon}{m} & \text{otherwise} \end{cases}$$

时序差分预测

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

其中 $R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$ 称为 TD Error.

$G_t \approx R_{t+1} + \gamma V(S_{t+1})$ 称之为 TD Target.

使用时序差分代替蒙特卡洛进行策略评价, 有以下两类控制方法:

On-policy: SARSA

Off-policy: Q-learning, Expected SARSA

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

行为策略和目标策略均为 ϵ -贪心. 根据策略取采样值

Q-Learning:

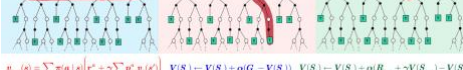
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a) - Q(S_t, A_t) \right)$$

行为策略为 ϵ -贪心, 目标策略为贪心. 根据策略取最大值

Expected SARSA:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \sum_{a \in A} \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t))$$

行为策略为 ϵ -贪心, 目标策略为行动期望. 根据策略取期望值



$$V_{\pi}(s) = \sum_{i=0}^{\infty} \gamma^i \mathbb{E} [r_{t+i} | s_t = s]$$

增量式价值函数近似

状态价值近似:特征提取 $\mathbf{x}(s) = (x_1(s), x_2(s), \dots, x_m(s))^T$, 函数族

$$\hat{v}(s | \mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x} \text{ (线性函数),}$$