



# 强化学习

## 与环境交互

智能体、环境、状态、动作、策略、奖励(现实中奖励和惩罚的统称)

基于评估、交互性、序列决策过程

## 离散马尔可夫链

$t+1$ 时刻状态仅与 $t$ 时刻状态相关

离散马尔可夫链: 引入奖励机制:  $R: S \times S \rightarrow \mathbb{R}$ ,  $R(S_t, S_{t+1})$  描述  $t$  与  $t+1$  所获得奖励,  
引入回报: 反映该时刻的累加奖励,

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$\gamma \in [0, 1]$ : 折扣因子

$R_{t+k}$ :  $t+k$ 时刻获得奖励

设计的奖励机制对当前时刻及附近时刻能够带来的奖励更为关注

马尔可夫奖励过程: MRP  $(S, P, R, \lambda)$

马尔可夫决策过程: MDP  $(S, A, P, R, \lambda)$

状态集合 $S$ , 动作集合 $A$ , 状态转移概率  $P(S_{t+1} | S_t, A_t)$

奖励函数,  $R(S_t, A_t, S_{t+1})$  在  $S_t$  执行  $A_t$  后到达  $S_{t+1}$  时, Agent 得到的奖励

折扣因子  $\gamma \in [0, 1]$

与环境交互:  $S_0 \xrightarrow[A_0]{} S_1 \xrightarrow[A_1]{} S_2 \rightarrow \dots$

轨迹: 状态序列  $(S_0, S_1, \dots)$

包含终止状态: 分段问题  $\rightarrow$  一个从初态到终态的完整轨迹叫一段  
不包含: 持续问题

策略函数: 刻画了智能体选择动作的机制.  $\pi: S \times A \rightarrow [0, 1]$ ,  $\pi(S, a)$ : 状态 $S$ 下采取动作 $a$ 的概率

价值函数: 衡量某个状态的好坏, 反映智能体从当前状态转移到该状态时能够为目标完成带来多大好处

若只有一个动作 $a$ , 记 $a = \pi(S)$

一个好的策略函数: 最大回报值  $G_t$

$V: S \rightarrow \mathbb{R}$ ,  $V_\pi(S) = E_\pi[G_t | S_t = S]$ ,  $t$ 时刻处于状态 $S$ 时, 按 $\pi$ 采取行动  
所获回报的期望

动作-价值函数  $q: S \times A \rightarrow \mathbb{R}$ ,  $q_\pi(S, a) = E_\pi[G_t | S_t = S, A_t = a]$

强化学习可转化为一个策略学习问题, 给定一个马尔可夫决策过程

MDP, 学习个最优策略 $\pi^*$ , 对任意 $s$ , 使得 $V$

$\pi(s)$ 最大

## ① 贝尔曼方程

价值函数  $V_\pi(S) = E_\pi[R_{t+1} + \gamma R_{t+2} + \dots | S_t = S]$

动作-价值函数:  $q_\pi(S, a) = E_\pi[R_{t+1} + \gamma R_{t+2} + \dots | S_t = S, A_t = a]$

$$V_\pi(S) = E_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = S]$$

$$= E_{a \sim \pi(S)} [E_\pi[\dots | S_t = S, A_t = a]]$$

$$= \sum_{a \in A} \pi(S, a) \times q_\pi(S, a)$$

$$V_\pi(S) = \sum \pi(S, a) \times q_\pi(S, a)$$

$$q_\pi(S, a) = \sum_{S' \in S} P(S' | S, a) [R(S, a, S') + \gamma V_\pi(S')]$$

$$q_\pi(S, a) = \sum_{S' \in S} P(S' | S, a) [R(S, a, S') + \gamma V_\pi(S')]$$

通过策略计算价值函数的过程叫做策略评估, 通过价值函数优化策略的过程叫做策略优化

策略评估和策略优化交替进行的强化学习求解方法叫做通用策略迭代

## 二、基于价值的强化学习

什么是更好的策略?

$\forall s \in S$  若  $q_\pi(S, \pi'(s)) \geq q_\pi(S, \pi(s))$  则  $V_{\pi'}(s) \geq V_\pi(s)$  (策略优化定理)

策略评估: 根据  $\pi$  来计算相应价值函数  $V_\pi$  或  $Q_\pi$

① 动态规划: 初始化  $V_\pi$ .

$$\text{for } s \in S \\ V_\pi(s) \leftarrow \sum_{a \in A} \pi(s,a) \sum_{s' \in S} P(s'|s,a) [R(s,a,s') + \gamma V_\pi(s')]$$

缺点: 智能体需要事先知道状态转移概率, 无法处理状态集合大小无限的情况。

② 蒙特卡洛采样. 选择不同起始状态, 按当前  $\pi$  采样若干轨迹,  $D$ .  
for  $s \in S$ : 计算  $D$  中  $s$  每次出现时对应的反馈  $G_1, \dots, G_k$ .  
 $V_\pi(s) \leftarrow \frac{1}{k} \sum_{i=1}^k G_i$

蒙特卡洛方法的核心思想: 期望可以通过样本均值来估计。  
给定状态  $s$ , 从该状态触发不断采样后续状态, 得到不同的采样序列。通过这些采样序列来分别计算状态  $s$  的回报值, 取均值作为  $s$  价值函数的估计, 避免对状态转移概率的依赖

③ 时序差分. 初始化  $V_\pi$ .

for-  
初始化  $s$  为初状态  
循环  
 $a \sim \pi(s)$ , 执行  $a$ , 观察  $R$  和  $s'$   
 $V_\pi(s) \leftarrow V_\pi(s) + \alpha [R(s,a,s') + \gamma V_\pi(s') - V_\pi(s)]$   
 $s \leftarrow s'$  =  $\frac{(1-\alpha)}{\text{过去}} V_\pi(s) + \frac{\alpha}{\text{学习得到}} [R(s,a,s') + \gamma V_\pi(s')]$   
直至  $s$  终止  
直至  $V_\pi$  收敛

时序差分法可以看作蒙特卡洛方法和动态规划方法的有机结合。时序差分方法从实际经验中获取信息, 无需提前获知环境模型的全部信息(类似于蒙特卡洛树搜索), 与动态规划法相似之处: 利用前序已知信息进行在线实时学习。  
时序差分法根据下一个状态的价值函数来估计, 克服了采样轨迹的稀疏可能带来样本方差较大的不足问题, 同时缩短了反库周期。

④ 价值迭代算法 (每次迭代只对一个状态进行评估和优化)

$$\pi(s) \leftarrow \arg \max_a Q_\pi(s,a) \\ V_\pi(s) \leftarrow Q_\pi(s, \pi(s))$$

⑤ Q 学习算法. 只有  $Q_\pi$  参与

初始化  $Q_\pi$

for 初始  $s$  为初态  
for  $a \leftarrow \arg \max_a Q_\pi(s,a)$ , 执行  $a$ ,  
 $Q_\pi(s,a) \leftarrow Q_\pi(s,a) + \alpha [R + \gamma \max_{a'} Q_\pi(s',a') - Q_\pi(s,a)]$   
 $s \leftarrow s'$   
 $s$  为终  
 $Q_\pi$  收敛

⑥ 探索与利用  $\epsilon$  贪心策略:  $\begin{cases} \arg \max_a Q_\pi(s,a) & 1-\epsilon \\ \text{random } a \in A & \epsilon \end{cases}$

⑦ 参数化与深度强化学习

$$L(\theta) = \frac{1}{2} [R + \gamma \max_{a'} Q_\pi(s',a'; \theta) - Q_\pi(s,a; \theta)]^2 \\ \theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta}$$

## ① DQN. (经验重现)

每探索到一个新的四元组，便将该四元组添加到经验重现表，在对参数进行更新时，则从表里随机选取一批样本计算损失函数。

一方面样本之间相关性显著减弱，另一方面过去的经验可以被重复利用，提高信息利用的效率。

$$L(\theta) = \frac{1}{2} \mathbb{E} [R + \gamma \max_{a'} Q_{\pi}(s, a'; \theta^-) - Q_{\pi}(s, a; \theta)]^2$$

$\theta^-$  相对稳定, 较低频率更新  $\theta^- \leftarrow \theta$ .

## ② 策略梯度法

通过直接参数化策略函数的方法求解强化学习问题。

函数取值表示在状态s下选择动作a的概率。选择一个动作的概率是随着参数的改变而光滑变化的，光滑性对算法收敛有更好保证。 $\pi_{\theta}(s, a)$ ,

最大化目标:  $J(\theta) = V_{\pi_{\theta}}(s_0)$

$$\text{梯度上升: } \nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_s \underbrace{\mu_{\pi_{\theta}}(s)}_{\text{策略分布}} \sum_a Q_{\pi_{\theta}}(s, a) \pi_{\theta}(s, a)$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &\propto \sum_s \mu_{\pi_{\theta}}(s) \sum_a Q_{\pi_{\theta}}(s, a) \nabla_{\theta} \pi_{\theta}(s, a) \\ &= \mathbb{E}_{s, a, j \sim \pi} [G_t \nabla_{\theta} \ln \pi_{\theta}(s, a)] \end{aligned}$$

## ③ Actor-Critic.

# 机器学习

## ① 监督学习, 从假设空间学习得到最优映射 $f$ (决策函数)

训练集, 将训练集中一部分数据作为验证集。

经验风险: 映射函数在训练集上产生的损失

期望风险: 所有数据中计算模型的损失。(真实风险, 真实误差)

机器学习中模型优化目标一般为经验风险最小化。

### ? 模型度量方法?

参数优化  $\left\{ \begin{array}{l} \text{频率学派: 最大似然估计} \\ \text{贝叶斯学派: 最大后验估计} \end{array} \right.$

## ② 回归分析 (刻画不同变量间的关系)

↳ 要从标注数据中学习得到

一元线性回归  $L(a,b) = \sum_{i=1}^n (y_i - ax_i - b)^2$

$$\frac{\partial L(a,b)}{\partial b} = \sum_{i=1}^n 2(y_i - ax_i - b)(-1) = 0 \Rightarrow \sum_{i=1}^n (y_i) - a \sum_{i=1}^n x_i - \sum_{i=1}^n b = 0$$

即  $n\bar{y} - na\bar{x} - nb = 0$

$$\frac{\partial L(a,b)}{\partial a} = \dots \quad \alpha = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}$$

## ③ 线性回归问题: outlier 非常敏感

logistic 回归:  $y = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-(a^T x + b)}}$

## ④ 决策树

建立决策树的过程, 就是不断选择属性值对样本集进行划分, 直至每个子样本为同一个类别。

构建决策树时划分属性的顺序选择是重要的。性能好  $\rightarrow$  决策树的分支结点样本集的纯度越来越高

信息熵: 衡量样本集合纯度的指标, 信息熵越大说明该集合的不确定性越大, 纯度越低。

选择属性划分样本集前后信息熵的减少量称为信息增益, 信息增益被用来衡量样本集合复杂度所减少的程度。

$K$  集合, 集合样本  $D$ :  $E(D) = - \sum_{k=1}^K p_k \log p_k$

信息增益:  $\text{Gain}(D, A) = E_{\text{nt}}(D) - \sum_{i=1}^n \frac{|D_i|}{|D|} E_{\text{nt}}(D_i)$

## ⑤ $K$ 均值聚类:

将  $n$  个  $d$  维数据划分为  $K$  个聚簇, 使簇内方差最小化。

$\Rightarrow$  找到一个局部最优, 即没有任何其它聚类结果让簇内方差最小化。

$x$  易受初值影响

1. 初始化聚类质心
2. 对数据进行聚类,  $\text{dist}(x_i, g_j) = \sqrt{\sum_{d=1}^d (x_{i,d} - g_{j,d})^2}$   $\arg \min_{g \in C} \text{mindist}(x_i, g)$
3. 更新聚类质心  $g_j = \frac{1}{|G_j|} \sum x_i$
4. 迭代

另一种解释:  $\arg \min_G \sum_{i=1}^n \sum_{x \in G_i} \|x - G_i\|^2 = \arg \min_G \sum_{i=1}^n |G_i| \text{Var} G_i$

kmeans  $\rightarrow$  最小化聚簇内的数据方差

⑩ 线性判别分析 (LDA)

基于监督学习的降维方法. (Fisher LDA)

对于一组具有标签信息的高维数据样本, LDA利用其类别信息, 将其线性投影到一个低维空间, 在低维空间中同一类别样本尽可能靠近, 不同类别样本尽可能彼此原

理. 样本集  $D = \{(x_i, y_i) | i=1, \dots, n\}$ ,  $y_i = \{c_1, c_2, \dots, c_k\}$   $k$  类样本

第  $i$  类样本协方差阵  $\Sigma_i = \sum_{x \in C_i} (x - m_i)(x - m_i)^T$ ,  $m_i$ : 均值.

对于二分类问题  $y(x) = w^T x$  投影.

$$S_1 = \sum_{x \in C_1} (w^T x - w^T m_1)^2 = w^T \sum_{x \in C_1} [(x - m_1)(x - m_1)^T] w$$

$S_1 = w^T \Sigma_1 w$ ,  $S_2 = w^T \Sigma_2 w$ :  $S_1, S_2$  衡量同一类别样本间分散程度.

$\max_{w^T (S_1 + S_2)} \max (\|m_1 - m_2\|_2^2)$ ,  $m_1 = w^T m_{1L}$ ,  $m_2 = w^T m_{2L}$

$$\Rightarrow \max J(w) = \frac{\|m_2 - m_1\|_2^2}{S_1 + S_2} = \frac{\|w^T (m_2 - m_1)\|_2^2}{w^T \Sigma_1 w + w^T \Sigma_2 w} = \frac{w^T S_b w}{w^T S_w w}$$

$S_b$ : 类间散度矩阵

$S_w$ : 类内散度矩阵

可令  $w^T S_w w = 1$   $L(w) = w^T S_b w - \lambda (w^T S_w w - 1)$

$\Rightarrow$  求导  $S_w^{-1} S_b w = \lambda w$

$S_b = (m_2 - m_1)(m_2 - m_1)^T \triangleq \lambda w = (m_2 - m_1)^T w$

$S_b w = (m_2 - m_1)(m_2 - m_1)^T w = \lambda w (m_2 - m_1)$

$\therefore S_w^{-1} S_b w = S_w^{-1} (m_2 - m_1) \times \lambda w = \lambda w$

$\therefore w = S_w^{-1} (m_2 - m_1)$

最大维度  $\min(k-1, d)$

⑪ PCA. (KL变换, 霍林特变换, 本征正交分解)

样本方差  $\text{Var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$

协方差  $\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - EX)(y_i - EY)$

皮尔逊相关系数  $\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$   
 $|\text{Corr}(X, Y)| \leq 1$ . 刻画线性相关程度

PCA:  $d \gg e$ ,  $d$  维映射到  $e$  维, 去除冗余性.

- 1. 中心化处理
- 2. 协方差矩阵
- 3.  $\Sigma$  特征分解,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_e$
- 4. 对应的  $w_1, w_2, \dots, w_e$  组成映射矩阵.

⑫ 特征人脸方法

本质是使用一组特征向量的线性组合来表示人脸  
人脸具有一定的拓扑结构, 可以用低维向量表达  
原始图像大部分信息. 可以用PCA降维, 但是在特征维度较高是, 可以用奇异值分解来实现PCA

$A = UDV^T$   $AA^T = UDV^T V D^T U^T$   
 $\downarrow \downarrow \downarrow$   
 $n \times d \quad n \times n \quad n \times d \quad d \times d$   
 $\Rightarrow (AA^T)U = U D^2$

⑬ 演化算法.

突变重组 + 自然选择, 模拟自然演化过程.

1. 初始化具有若干规模数目的群体. 当前进化代数  $\text{Generation} = 0$ ;
2. 采用评估函数对群体所有的染色体进行评价, 分别计算每个染色体的适应值, 保存适应值最大的染色体  $\text{Best}$
3. 采用轮盘赌选择算法对群体染色体进行选择操作, 产生规模同样的种群.
4. 按照概率从种群中选择染色体进行交配.
5. 按照概率对新种群的基因进行变异操作.
6. 变异后的新群体取代原有群体, 重新计算群体中各个染色体的适应值. 若大于  $\text{best}$  则取代  $\text{best}$
7. 当前进化代数加1

## 知识表达与推理

### 一、命题逻辑

逻辑与推理是人工智能的核心问题

### 二、谓词逻辑

### 三、知识图谱推理

知识图谱由有向图构成，被用来描述显示世界中实体及实体之间的关系。每个节点表示客观世界中的一个实体，两个节点之间的连线表示节点具有某一关系，知识图谱中存在连线的两个实体可以表达为三元组形式。三元组也可以表示为一阶逻辑的形式。

从数据到知识(推导规则): 归纳, 从知识到数据(实例): 演绎

归纳逻辑程序设计是机器学习和逻辑程序设计交叉领域的研究内容。

FOIL通过序贯覆盖学习推理。

FOIL算法:

目标谓词是需要推断规则的结论, 也称为规则头。在给定推理结论后, FOIL算法学习得到使得结论满足的前提条件, 即目标谓词作为结论的推理规则。

正例

反例可以从知识图谱中构造出来, 只有两个实体之间存在的关系且确定相悖时才能用来做反例。

背景知识

算法思路: 从一般到特殊, 逐步添加目标谓词的前提约束谓词, 直到所构成的推理规则不覆盖任何反例。

FOIL信息增益值: 
$$FOIL\_Gain = \hat{m}_+ \cdot \left( \log_2 \frac{\hat{m}_+}{\hat{m}_+ + \hat{m}_-} - \log_2 \frac{m_+}{m_+ + m_-} \right)$$
  
 $\hat{m}_+$ : 添加后,  $m_+$ : 原推理

添加前提约束谓词后所得推理规则的质量好坏由信息增益值来判断

路径排序推理算法:

- 特征提取: 生产并选择路径特征集合

特征计算: 计算每个训练样例的特征值。  $P(S \rightarrow t; \pi_i)$  从S出发通过  $\pi_i$  到达t的概率

分类器训练, 分类器可以用于推理两个实体之间是否存在目标关系

概率图谱推理:  $\begin{cases} \text{贝叶斯网络 (DAG)} \rightarrow \text{局部马尔可夫性, 联合分布} = \prod P(\text{节点} | \text{父节点}) \\ \text{马尔可夫网络 (无向图)} \end{cases}$

因果推理 混淆关联  $\begin{matrix} X \\ \swarrow \searrow \\ T \end{matrix}$   
选择关联  $\begin{matrix} T \\ \swarrow \searrow \\ X \end{matrix}$

① 因果图: DAG,  $P(x_1, \dots, x_n) = \prod_{j=1}^n P(x_j | \text{pa}(x_j))$

② 干预:

改变明确存在关联关系的某变量取值, 研究变量取值改变对结果变量的影响。

do算子:

$$do(X=x), P(Y=y | do(X=x))$$

因果效应差 / 平均因果效应  $P(Y=1 | do(X=1)) - P(Y=1 | do(X=0))$

对用药情况X进行干预并固定其值为x时, 可将所有指向X的边移除

边际概率  $P_m$ :  $\begin{cases} \text{边缘概率 } P(Z=z) \text{ 不变. } P_m(Z=z) = P(Z=z) \\ \text{条件概率 } P(Y=y | X=x, Z=z) \text{ 不变. } P_m(Y=y | X=x, Z=z) = P(Y=y | X=x, Z=z) \end{cases}$

$$\therefore P(Y=y | do(X=x)) = P_m(Y=y | X=x) = \sum_z P(Y=y | X=x, Z=z) \cdot P(Z=z) \quad (\text{调整公式}) \quad Z\text{-调整} / Z\text{-控制}$$

③ 反事实推理



# 人工智能博弈

推动机器学习从数据拟合过程中以求取最优解为核心向博弈对抗过程中以求去均衡解为核心的转变

博弈论主要研究博弈行为中最优的对抗策略及其稳定局势

相关概念：参与者或玩家，策略，某个参与者可采纳策略的全体组合形成了策略集，所有参与者鸽子采取行动后形成的状态被称为局势。如果参与者可以通过一定概率分布来选择若干个不同的策略，称为混合策略，有确定的策略称为纯策略

收益，混合策略下的收益为期望收益

博弈论研究的范式：建模者对参与者规定可采区的策略集和取得的收益，观察当参与者选择若干策略以最大化其收益时会产生什么结果

分类：合作博弈与非合作博弈

静态博弈(所有参与者同时决策，或参与者互相不知道对方的决策)和动态博弈(参与者所采取的行为先后顺序由规则决定，且后行动者知道先行行动者所采取的行为)

完全信息博弈与不完全信息博弈 完全信息：参与者均了解其他参与者的策略集、收益等信息

纳什均衡：博弈的稳定局势：参与者作出的这样一种策略组合，在该策略组合上，任何参与者单独改变策略都不会得到好处。即当所有其他人不改变策略时，没有人会改变自己的策略。

Nash定理：若参与者有限，每位参与者的策略集有限，收益函数为实值函数，则博弈必定存在混合策略意义下的纳什均衡。

策梅洛定理：对于任意一个有限步的双人完全信息零和动态博弈，一定存在先手必胜策略或后手必胜策略或保平策略。

虚拟遗憾最小化算法：N个玩家参与，i策略 $\sigma_i$ ，策略组合 $\sigma = \{\sigma_1, \dots, \sigma_N\}$

除i外， $\sigma_{-i} = \{\sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots, \sigma_N\}$

最优反应策略：玩家i在终局收益 $u_i(\sigma)$

给定 $\sigma_{-i}$ 下，对i的最优 $\sigma_i^*$ ： $u_i(\sigma_i^*, \sigma_{-i}) \geq \max_{\sigma_i' \in \bar{\Sigma}_i} u_i(\sigma_i', \sigma_{-i})$   
 $\hookrightarrow$ 可选的所有策略

在策略组合中，如果每一个玩家的策略相对于其他玩家的策略都是最佳反应策略，则该策略组合就是一个纳什均衡策略。在有限对手，有限策略情况下，纳什均衡一定存在。

$$\sigma^* = \{\sigma_1^*, \dots, \sigma_N^*\}, \text{ 对 } \forall i, u_i(\sigma_i^*) \geq \max_{\sigma_i' \in \bar{\Sigma}_i} u_i(\sigma_i', \sigma_{-i}^*)$$

遗憾最小化算法是一种根据以往博弈过程中所得遗憾程度来选择未来行为的方法。

i在过去T轮采取 $\sigma_i$ 的累加遗憾值

$$\text{Regret}_i^T(\sigma_i) = \sum_{t=1}^T (u_i(\sigma_i^*, \sigma_{-i}^t) - u_i(\sigma_i^t, \sigma_{-i}^t))$$

有效遗憾值  $\text{Regret}_i^{T+1}(\sigma_i) = \max(\text{Regret}_i^T(\sigma_i), 0)$

$$P(\sigma_i^{T+1}) = \begin{cases} \frac{\text{Regret}_i^T(\sigma_i)}{\sum \text{Regret}_i^T} & \text{if } \sum \text{Regret}_i^T > 0 \\ \frac{1}{|\bar{\Sigma}_i|} & \text{otherwise.} \end{cases}$$

对于任意序贯决策的博弈对抗，可将博弈过程表示成一棵博弈树，博弈树中的每一个中间节点都是一个信息集，信息集中包含了博弈中当前的状态。给定博弈树的每一个节点，玩家都可以从一系列的动作中选择，然后状态发生转化，直到终局。

信息集I，采取 $a$ ： $\sigma_I \rightarrow a$

行动序列h，在策略组合 $\sigma$ 中，h出现的概率 $\pi^\sigma(h)$

在 $\sigma$ 下，信息集出现概率 $\pi^\sigma(I) = \sum_{h \in I} \pi^\sigma(h)$

终局：叶子集合Z，i的收益 $u_i(z)$

给定h，按 $\sigma$ 到达z的概率 $\pi^\sigma(h, z)$

虚拟价值  $v_i(\sigma, h) = \sum_{z \in Z} \pi^\sigma(h, z) \times u_i(h, z)$

遗憾值  $r_i(h, a) = v_i(\sigma_{-a}, h) - v_i(\sigma, h)$  信息集I的遗憾： $r_i(I, a) = \sum_{h \in I} r_i(h, a)$





# 深度学习

## 神经元. 感知机.

在感知机模型中增加若干隐藏层，增强神经网络的非线性表达能力，就会让神经网络具有更强的拟合能力。  
多个隐藏层构成的多层感知机，也被称为前馈神经网络。  
层层递进，逐层抽象；非线性映射；误差反馈调优

$$\frac{1}{1+e^{-x}}$$

激活函数: sigmoid: 概率形式输出，单调递增，非线性变化；梯度消失问题。  
ReLU: 有效缓解梯度消失问题， $x < 0$ ，稀疏性一定程度上克服过拟合现象。但是  
也导致神经元死亡。

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}$$

神经网络参数优化: 监督学习的过程，让神经网络对数据进行拟合。

模型利用反向传播算法将损失函数计算所得误差从输出端触发，由后向前传递给神经网络中每个单元，通过梯度下降对神经网络中参数进行更新。

损失函数（代价函数），计算模型预测值与真实值之间的误差。

均方损失误差

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

交叉熵损失误差

熵用来表示热力学系统所呈现的无序程度，信息熵，通过对数函数来测量信息的不确定性。

交叉熵用来度量两个概率分布间的差异，刻画了两个概率分布之间的距离，旨在描绘通过概率分布来表达的困难程度。

梯度下降: 
$$f(x+\Delta x) - f(x) \approx (\nabla f(x))^T \Delta x < 0$$
  

$$= \|\nabla f(x)\| \|\Delta x\| \cos \theta$$
  

$$\theta = \pi \text{ 时下降最快}$$

## 误差反向传播

利用损失函数来计算模型预测结果与真实结果之间的误差以优化调整模型参数，从输出端向输入端，由后向前递进进行。

批量梯度下降: 在整个训练集上计算损失误差

随机梯度下降: 使用训练集中每个训练样本计算损失函数，梯度方向有很大的波动，收敛慢，但是有助于跳出局部最优解

小批量梯度下降: 选取训练集中小批量样本计算，根据每一批量样本所得到的累加误差来更新参数，保证训练过程更稳定，利用矩阵计算优势。

卷积神经网络:

对于图像这样的数据，不能直接将所构成的像素点向量与前馈神经网络中的每个神经元相连

卷积神经网络的前身: 级联方式（逐层滤波）实现一种满足平移不变性的网络

卷积滤波结果: 特征图

滤波可以视为在给定卷积核权重前提下，记住了领域像素点之间的若干特定空间模式，忽略了某些靠近空间模式

卷积的思想利用了图像中像素点存在空间依赖度的特点，对图像进行了下采样操作

图像平滑操作: 中心位置的权重系数小，且与其它卷积权重系数差吧娇小

填充和步长

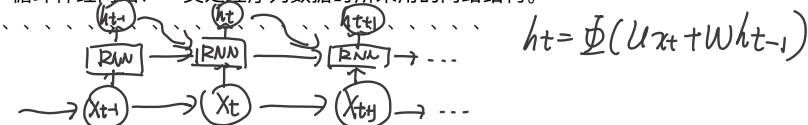
感受野是卷积神经网络每一层输出的特征图上的像素点在输入图像上映射区域的大小。

局部感知、参数共

$$W = \frac{F+2P}{S} + 1$$

池化: 最大/平均/k-max 池化

循环神经网络: 一类处理序列数据时所采用的网络结构。



前序时刻信息影响后续

得到句子的向量编码，最后一个单词: 整个句子向量编码表

每个单词隐式编码加权平均 → 句子向量编码

容易出现梯度消失

# 长短期记忆模型

## ④内部记忆单元[1]

### ④注意力机制 自注意力 $\rightarrow$ 单词之间概率关联

查询向量  $q_i = W^q x w_i$

键  $\dots : k_i = W^k x w_i$

值  $\dots : v_i = W^v x w_i$

$$q_i = W^q x w_i$$

$$q \cdot k \rightarrow a_{ji}'$$

$$a_{ji} = q_i \cdot k_j \xrightarrow{\text{softmax}} a_{ji}'$$

$\sum_j a_{ji}' \times v_j$  为  $w_i$  注意到与其B单词的关联程度

### ④正则化.

\* Dropout. 随机丢掉一部分神经元

批归一化

$L_1$  正则化: 稀疏规则算子

# 搜索

状态, 动作, 状态转移, 路径和代价, 目标测试

搜索过程可视为搜索树的构建

评测标准: 完备性, 最优性, 时间空间复杂度

贪婪最佳优先搜索:

评价函数: 从当前节点n出发, 根据评价函数来选择后续节点: 下一个节点是谁?

启发函数: 计算从节点n到目标节点之间所形成路径的最小代价值: 完成任务还需要多少代价?

贪婪最佳优先搜索: 启发函数等于评价函数, 不一定最优

A\*搜索算法: 在评价函数中考虑从起始节点到当前节点的路径代价 可溶性, 一致性

可溶性:  $\forall n, h(n) \leq \underline{h^*(n)}$  启发函数不会过高估计代价  
最小代价

一致性:  $h(n) \leq c(n, a, n') + h(n')$

一致性必然导致可溶性

minimax搜索: 对抗搜索/博弈搜索

一方想最大化自身的利益, 另一方想最小化对手的利益

Alpha-Beta 剪枝,

对于MAX节点, 若子节点的收益大于 $\alpha$ , 则 $\alpha = \text{收益}$ .

$\alpha = -\infty, \beta = +\infty$  若 $\alpha > \beta$ 则剪枝

① 蒙特卡洛树搜索.

悔值函数  $P_T = T_{\text{root}} - \sum_{t=1}^T \hat{v}_t$

$$C \sqrt{\frac{2 \ln t}{T(i, t-1)}}$$

$\epsilon$ -贪心:  $v_t = \begin{cases} \arg \max \bar{x}_{i, T(i, t-1)} & 1-\epsilon \\ \text{random.} & \epsilon \end{cases}$

上限置信区间  $P(\mu_i - \bar{x}_{i, T(i, t-1)} > \delta) \leq e^{-2T(i, t-1)\delta^2}$

认为  $\bar{x}_{i, T(i, t-1)} + \delta$  是  $\mu_i$  的上界.

$$C \sqrt{\frac{2 \ln t}{T(i, t-1)}}$$

选择, 扩展, 模拟, 反向传播