

ch11 为什么特征选择: 降维, 去除不相关特征降低学习难度. 选特征子集: 前向搜索增加相关特征. 子集评价: 属性子集A, 信息增益Gain(A) = Ent(D) - \sum\_{v=1}^V \frac{|D^v|}{|D|} Ent(D^v), 信息熵Ent(D) = -\sum\_{i=1}^{|Y|} p\_k \log\_2 p\_k, Gain(A) \uparrow 则特征子集有助于分类的信息越多. 搜索+评价=特征选择, 前向搜+信息熵则类似决策树. 特征选择方法: 1 过滤式. Relief: 特征子集重要性: 子集中每个特征所对应的相关统计量分量\delta^j = \sum\_i -\text{diff}\left(x\_i^j, x\_{i,\text{nh}}^j\right)^2 + \text{diff}\left(x\_i^j, x\_{i,\text{nm}}^j\right)^2之和, 选阈值\tau以上的. \mathbf{x}\_{i,\text{nm}} 猜错近邻, 属性有益则\mathbf{x}\_i与\mathbf{x}\_{i,\text{nm}}距离大于\mathbf{x}\_i与\mathbf{x}\_{i,\text{nh}}. 2 包裹式. LVM拉斯维加斯方法框架下随机搜索出特征子集, 训练后看误差. 3 嵌入式(特征选择与训练过程融合). LASSO, L\_1范数. 近端梯度下降PGD: \nabla f 满足L-Lipschitz. 迭代 \mathbf{x}\_{k+1} = \arg \min\_{\mathbf{x}} \frac{1}{2} \left\| \mathbf{x} - \left( \mathbf{x}\_k - \frac{1}{L} \nabla f(\mathbf{x}\_k) \right) \right\|\_2^2 + \lambda \|\mathbf{x}\|\_1, 闭式解 \mathbf{x}\_{k+1}^i = z^i - \lambda/L, \lambda/L < z^i; z^i + \lambda/L, z^i < -\lambda/L; 0 其他. 稀疏表示与字典学习: 优点 1 使大多数问题线性可分, 2 存储负担小. 字典学习: 为稠密表达找字典 \rightarrow 稀疏表达: 1 固定B 像LASSO解法: \min\_{\alpha\_i} \|\mathbf{x}\_i - \mathbf{B}\alpha\_i\|\_2^2 + \lambda \|\alpha\_i\|\_1, 2 以\alpha\_i为初值 \min\_{\mathbf{B}} \|\mathbf{X} - \mathbf{B}\mathbf{A}\|\_F^2 = \min\_{\mathbf{B}\_i} \left\| \mathbf{X} - \sum\_{j=1}^k \mathbf{b}\_j \alpha^j \right\|\_F^2 = \min\_{\mathbf{B}\_i} \left\| \left( \mathbf{X} - \sum\_{j \neq i} \mathbf{b}\_j \alpha^j \right) - \mathbf{b}\_i \alpha^i \right\|\_F^2 括号内\mathbf{E}\_i进行奇异值分解, 但直接分解可能破坏\mathbf{A}的稀疏性, \therefore \text{KSVD对}\alpha\_i\text{保留非零元素 } \mathbf{E}\_i\text{仅保留}\mathbf{b}\_i\text{与}\alpha\_i\text{的非零元素乘积项. 压缩感知 } \mathbf{y} = \Phi \mathbf{x} = \Phi \Psi \mathbf{s} = \mathbf{A} \mathbf{s} \text{ 恢复}\mathbf{s} \Rightarrow \text{恢复}\mathbf{x}, k\text{限定等距性k-RIP } (1 - \delta\_k) \|\mathbf{s}\|\_2^2 \leq \|\mathbf{A}\_k \mathbf{s}\|\_2^2 \leq (1 + \delta\_k) \|\mathbf{s}\|\_2^2 \text{ 此时 } \min\_{\mathbf{s}} \|\mathbf{s}\|\_0 \Rightarrow \|\mathbf{s}\|\_1 \text{ s.t. } \mathbf{y} = \mathbf{A} \mathbf{s}. \text{ 矩阵补全 } \min\_{\mathbf{X}} \text{rank}(\mathbf{X}) \text{ s.t. } (\mathbf{X})\_{ij} = (\mathbf{A})\_{ij}, (i, j) \in \Omega, \text{ 转化: } \therefore \text{rank}(\mathbf{X}) \text{ 在集合 } \left\{ \mathbf{X} \in \mathbb{R}^{m \times n} : \|\mathbf{X}\|\_F^2 \leq 1 \right\} \text{ 上的凸包是 } \mathbf{X} \text{ 的“核范数” } \|\mathbf{X}\|\_\* = \sum\_{j=1}^{\min\{m,n\}} \sigma\_j(\mathbf{X}).

ch12 学习算法\mathfrak{L}学得模型对应的假设h尽可能接近目标概念c. 不能完全一致因为D有限D采样偶然. 1PAC辨识: \mathfrak{L}从\mathcal{H}中PAC辨识C: P(E(h) \leq \epsilon) \geq 1 - \delta. C 2PAC可学习: \mathbf{p} = \text{poly}(1/\epsilon, 1/\delta, \text{size}(\mathbf{x}), \text{size}(c)) \forall m \geq \mathbf{p} \text{ 可PAC辨识. 3PAC学习算法: } \mathfrak{L} \text{ PAC可学习且运行时间} = \mathbf{p}. 4样本复杂度: PAC学习算法最小m \geq \mathbf{p}. 有限假设空间可分: \mathfrak{L}以概率1 - \delta找到目标假设的\epsilon近似: \text{All } i, P(h(\mathbf{x}\_i) = y\_i) < (1 - \epsilon)^m. \text{对All } h, \text{上式} < |\mathcal{H}|(1 - \epsilon)^m < |\mathcal{H}|e^{-m\epsilon} \leq \delta, \text{上式可推} m \geq \dots. \text{不可分: } \forall h, P(|E(h) - \hat{E}(h)| \leq \sqrt{(\ln|\mathcal{H}| + \ln(2/\delta))/(2m)}) \geq 1 - \delta, 5\text{不可知PAC可学习: } \forall m \geq \mathbf{p}, P(E(h) - \min\_{h' \in \mathcal{H}} E(h') \leq \epsilon) \geq 1 - \delta. 6\text{增长函数: } \Pi\_{\mathcal{H}}(m) = \max\_{\{\mathbf{x}\_1, \dots, \mathbf{x}\_m\} \subseteq \mathcal{X}} |\{(h(\mathbf{x}\_1), \dots, h(\mathbf{x}\_m)) \mid h \in \mathcal{H}\}| \text{ 集合的数量. 7对分: } \mathcal{H} \text{ 中的} h \text{ 对} D \text{ 中赋予标记的每种可能结果. 8D被}\mathcal{H}\text{打散: 二分类问题}\Pi\_{\mathcal{H}}(m) = 2^m \text{ 实现}\forall\text{对分. 9VC维: } VC(\mathcal{H}) = \max\{m : \Pi\_{\mathcal{H}}(m) = 2^m\}. \text{VC维}d: \exists \text{大小为}d\text{的示例集能被}\mathcal{H}\text{打散, } \forall d \text{大小为}d+1\text{的示例集不能.}

ch13 主动学习: 查询专家后有标签加入再训练, 尽量少的“查询”来获得最优性能. 半监督: 不依赖外界交互/自动利用未标记样本: 基于聚类假设/流形假设相似样本相似输出. 半监督分为: 1纯半监督/2直推学习 1中训练数据中未标记样本不是待测数据而2中是, 1基于开放世界假设希望模型适用于未观察的数据, 2封闭世界假设仅预测未标记数据. 半监督的生成式: 未标记数据看成模型缺失参数. EM: p(\mathbf{x}) = \sum\_{i=1}^N \alpha\_i \cdot p(\mathbf{x} \mid \mu\_i, \Sigma\_i) \therefore \text{MAP: } f(\mathbf{x}) = \arg \max\_{j \in \mathcal{Y}} p(y = j \mid \mathbf{x}) = \arg \max\_{j \in \mathcal{Y}} \sum\_{i=1}^N p(y = j \mid \Theta = i, \mathbf{x}) \cdot p(\Theta = i \mid \mathbf{x}) \text{ 其中 } p(\Theta = i \mid \mathbf{x}) = \alpha\_i \cdot p(\mathbf{x} \mid \mu\_i, \Sigma\_i) / \sum\_{i=1}^N \alpha\_i \cdot p(\mathbf{x} \mid \mu\_i, \Sigma\_i). \Theta \text{ 指向哪个高斯混合成分. } l\_i \text{ 表示第} i \text{ 类的有标记样本数目. } D\_l \cup D\_u \text{ 对数似然: } LL(D\_l \cup D\_u) = \sum\_{(\mathbf{x}\_j, y\_j) \in D\_l} \ln(\sum\_{i=1}^N \alpha\_i \cdot p(\mathbf{x}\_j \mid \mu\_i, \Sigma\_i) \cdot p(y\_j \mid \Theta = i, \mathbf{x}\_j)) + \sum\_{\mathbf{x}\_j \in D\_u} \ln(\sum\_{i=1}^N \alpha\_i \cdot p(\mathbf{x}\_j \mid \mu\_i, \Sigma\_i)), \text{ E步: } \gamma\_{ji} = \frac{\alpha\_i \cdot p(\mathbf{x}\_j \mid \mu\_i, \Sigma\_i)}{\sum\_{i=1}^N \alpha\_i \cdot p(\mathbf{x}\_j \mid \mu\_i, \Sigma\_i)}, \text{ M 步: } \mu\_i = \frac{1}{\sum\_{\mathbf{x}\_j \in D\_u} \gamma\_{ji} + l\_i} \left( \sum\_{\mathbf{x}\_j \in D\_u} \gamma\_{ji} \mathbf{x}\_j + \sum\_{(\mathbf{x}\_j, y\_j) \in D\_l \wedge y\_j = i} \mathbf{x}\_j \right), \Sigma\_i = 1/(\sum\_{\mathbf{x}\_j \in D\_u} \gamma\_{ji} + l\_i) \left( \sum\_{\mathbf{x}\_j \in D\_u} \gamma\_{ji} (\mathbf{x}\_j - \mu\_i)(\mathbf{x}\_j - \mu\_i)^T + \sum\_{(\mathbf{x}\_j, y\_j) \in D\_l \wedge y\_j = i} (\mathbf{x}\_j - \mu\_i)(\mathbf{x}\_j - \mu\_i)^T \right) \alpha\_i = \frac{1}{m} \left( \sum\_{\mathbf{x}\_j \in D\_u} \gamma\_{ji} + l\_i \right). \text{ 此类生成式方法关键: 模型假设必须准确, 假设的生成式模型必须与真实数据分布吻合. 半监督SVM中TSVM: 二分类, 式子① } \min\_{\mathbf{w}, \mathbf{b}, \hat{\mathbf{y}}, \xi} \frac{1}{2} \|\mathbf{w}\|\_2^2 + C\_l \sum\_{i=1}^n \xi\_i + C\_u \sum\_{i=1}^m \xi\_i \text{ s.t. } y\_i(\mathbf{w}^T \mathbf{x}\_i + b) \geq 1 - \xi\_i, \hat{y}\_i(\mathbf{w}^T \mathbf{x}\_i + b) \geq 1 - \xi\_i, \xi\_i \geq 0. \text{ 试图考虑对未标记样本进行尝试各种可能的标记指派: 局部搜索来迭代地寻找近似解 } 1D\_l \text{ 学得SVM, 2给} D\_u \text{ 伪标记, 3代入①求解出新的SVM, 4找预测为异类且很可能发生错误的} D\_u: \exists \{i, j \mid (\hat{y}\_i \hat{y}\_j < 0) \wedge (\xi\_i > 0) \wedge (\xi\_j > 0) \wedge (\xi\_i + \xi\_j > 2)\}, \text{ 再求解①, 这样找一遍} D\_u \text{ 后 } C\_u = \min\{2C\_u, C\_l\}, \text{ 直到 } C\_u = C\_l \text{ 时停止. 初始化 } C\_u \ll C\_l. \text{ 对未标记样本进行调整过程中可能类别不平衡拆分 } C\_u \text{ 并改进初始化 } C\_u^+ = u\_-/u\_+, C\_u^- = \dots. \text{ 图半监督: 结点 } V \text{ } \mathbf{x}\_i, \text{ 边集 } E \text{ 亲和矩阵 } (\mathbf{W})\_{ij} = \exp(-\|\mathbf{x}\_i - \mathbf{x}\_j\|\_2^2/2\sigma^2), \text{ if } i \neq j \text{ 其他为0, 期望学得 } f: V \rightarrow \mathbb{R}, \text{ 通过 } \text{sign}(f(x\_i)) \text{ 分类, 二分类标记传播: 定义能力函数及其矩阵表达(有点像损失函数): } E(f) = \frac{1}{2} \sum\_{i=1}^m \sum\_{j=1}^n (\mathbf{W})\_{ij} (f(\mathbf{x}\_i) - f(\mathbf{x}\_j))^2 = \frac{1}{2} \left( \sum\_{i=1}^m d\_i f^2(\mathbf{x}\_i) + \sum\_{j=1}^n d\_j f^2(\mathbf{x}\_j) - 2 \sum\_{i=1}^m \sum\_{j=1}^n (\mathbf{W})\_{ij} f(\mathbf{x}\_i) f(\mathbf{x}\_j) \right) = \sum\_{i=1}^m d\_i f^2(\mathbf{x}\_i) - \sum\_{i=1}^m \sum\_{j=1}^n (\mathbf{W})\_{ij} f(\mathbf{x}\_i) f(\mathbf{x}\_j) = \mathbf{f}^T (\mathbf{D} - \mathbf{W}) \mathbf{f}. \text{ 这里 } \mathbf{D} = \text{diag}(d\_1, d\_2, \dots, d\_{l+u}) \text{ } d\_i = \sum\_{j=1}^{l+u} (\mathbf{W})\_{ij}. \text{ 拉普拉斯矩阵 } \Delta = \mathbf{D} - \mathbf{W}. \text{ 按照 } l \text{ 和 } u \text{ 分块矩阵表达: } E(f) = \mathbf{f}\_l^T (D\_{ll} - \mathbf{W}\_{ll}) \mathbf{f}\_l - 2 \mathbf{f}\_l^T \mathbf{W}\_{lu} \mathbf{f}\_l + \mathbf{f}\_u^T (D\_{uu} - \mathbf{W}\_{uu}) \mathbf{f}\_u, \text{ 对 } \mathbf{f}\_u \text{ 求偏导等于0得式②: } \mathbf{f}\_u = (D\_{uu} - \mathbf{W}\_{uu})^{-1} \mathbf{W}\_{ul} \mathbf{f}\_l. \text{ 进一步分块, 令 } \mathbf{P}\_{uu} = D\_{uu}^{-1} \mathbf{W}\_{uu}, \mathbf{P}\_{ul} = D\_{uu}^{-1} \mathbf{W}\_{ul}, \text{ ②化简为: } \mathbf{f}\_u = (D\_{uu} (\mathbf{I} - D\_{uu}^{-1} \mathbf{W}\_{uu}))^{-1} \mathbf{W}\_{ul} \mathbf{f}\_l = (\mathbf{I} - \mathbf{P}\_{uu})^{-1} \mathbf{P}\_{ul} \mathbf{f}\_l. \text{ 多分类标记传播: } \mathbf{F} \text{ 矩阵迭代, } \mathbf{F}(0) = (\mathbf{Y})\_{ij} = \text{当} i \text{ 是有标记样本(第} i \text{ 个)并且类别为} j \text{ 时为1. 标记传播矩阵 } \mathbf{S} = \mathbf{D}^{-\frac{1}{2}} \mathbf{W} \mathbf{D}^{-\frac{1}{2}}, \mathbf{D}^{-\frac{1}{2}} = \text{diag}(1/\sqrt{d\_1}, \dots, 1/\sqrt{d\_{l+u}}), \text{ 于是有迭代式: } \mathbf{F}(t+1) = \alpha \mathbf{S} \mathbf{F}(t) + (1 - \alpha) \mathbf{Y} \text{ 基于此迭代收敛得: } \mathbf{F}^\* = \lim\_{t \rightarrow \infty} \mathbf{F}(t) = (1 - \alpha)(\mathbf{I} - \alpha \mathbf{S})^{-1} \mathbf{Y}. \text{ 对于所有未标记样本} i: y\_i = \arg \max\_{1 \leq j \leq |Y|} (\mathbf{F}^\*)\_{ij}. \text{ 对应于正则化框架: } \min\_{\mathbf{F}} \frac{1}{2} \left( \sum\_{i,j=1}^{l+u} (\mathbf{W})\_{ij} \left\| \frac{1}{\sqrt{d\_i}} \mathbf{F}\_i - \frac{1}{\sqrt{d\_j}} \mathbf{F}\_j \right\|^2 \right) + \mu \sum\_{i=1}^l \|\mathbf{F}\_i - \mathbf{Y}\_i\|^2, \text{ 上}

式中\mu = (1 - \alpha)/\alpha时最优解恰为迭代收敛解\mathbf{F}^\*. 基于分歧的方法: 学习器之间的分歧作用于未标记数据. multi-view中的co-training: 图像属性集是一个视图. 不同视图相容性: 包含的关于输出空间\mathcal{Y}(标签)的信息是一致的. 互补性: 不同视图信息互补. 充分且条件独立视图: 充分每个视图包含足以产生最优学习器的信息, 条件独立在给定类别标记下两个视图独立. 训练: 先在每个视图上训, 让训好的分类器挑选最有把握的未标记数据赋予伪标记, 迭代. 评价基于分歧的方法: 只需采用合适的基学习器具有显著分歧性能尚可的, 就能较少受到模型假设/损失函数非凸性和数据规模问题的影响, 学习方法简单有效理论基础相对坚实适用范围较为广泛. 半监督聚类: 聚类任务中的监督信息 1”必连”(样本属于同一个簇)和”勿连”约束, 2少量的有标记样本. 约束k均值算法利用第1类信息(注意这只有部分样本的约束信息): \mathcal{M}是必连关系集合((x\_i, x\_j)属于它表示同类) \mathcal{C}是勿连关系集合, 训过程: 在k-means聚类过程中保证\mathcal{M}和\mathcal{C}中的约束满足即可. 约束种子k均值算法利用第2类信息: 少量有标记样本: S = \cup\_{j=1}^k S\_j \subset D, S\_j \text{ 表示隶属于第} j \text{ 个聚类簇的样本, 训过程: 直接用有监督信息的样本作为初始化聚类中心, 并在样本簇归属的迭代更新过程中不更新S中的类别隶属.}

ch14 EM算法: \log p(x\_i \mid \theta) = \log(p(x\_i, z\_i)) - \log(p(z\_i \mid x\_i)) = E\_{z\_i}(\log p(x\_i, z\_i)) - H(q\_i) + KL(q\_i(z\_i) \parallel p(z\_i \mid x\_i)), \Rightarrow \text{最大化ELBO, E步把KL设为0; M步最大化ELBO. 概率模型: 一种框架, 将学习任务归结于计算变量的概率分布, “推断”: 利用已知变量推测未知变量的分布; (定义集合: Y 所关心的变量; O 可观测变量; R 其他变量) 则生成式模型考虑联合分布 P(Y, R, O), 判别式模型考虑条件分布 P(Y, R \mid O). 给定O, 需由P(Y, R, O)或P(Y, R \mid O)得到P(Y \mid O). 概率图模型: 一类用图来表达变量相关关系的概率模型, 一个截断表示一个或一组随机变量, 结点之间的边表示变量间的概率相关关系. 概率图模型分为: 1使用有向无环图表示变量间的依赖关系, 有向图模型/贝叶斯网; 2使用无向图表示变量间的相关关系, 无向图模型/马尔可夫网. HMM隐马尔可夫模型结构最简单的动态贝叶斯网(生成式): 变量有 1状态变量\mathcal{Y}是隐藏的不可观测的隐变量, 离散; 2观测变量\mathcal{X}观测值, 离散或连续. 确定HMM所需的参数: 1结构信息联合概率分布: P(指向的结点 \mid 指出来的结点)的相乘; 2状态转移概率: \alpha\_{ij} = P(y\_{t+1} = s\_j \mid y\_t = s\_i); 3输出观测概率 b\_{ij} = P(x\_t = o\_j \mid y\_t = s\_i), 根据当前状态获得各个观测值的概率; 4初始状态概率 \pi\_i = P(y\_1 = s\_i). HMM产生观测序列\{x\_i\}: 1)设置 t = 1, 并根据初始状态概率 \pi 选择初始状态 y\_1; 2)根据状态 y\_t 和输出观测概率 \mathbf{B} 选择观测变量取值 x\_t; 3)根据状态 y\_t 和状态转移矩阵 \mathbf{A} 转移模型状态, 即确定 y\_{t+1}; 4)若 t < n, 设置 t = t + 1, 并转到第(2)步, 否则停止. HMM三个基本问题(基于联合概率分布条件独立性可高效求解): 1)给定模型 \lambda = [\mathbf{A}, \mathbf{B}, \pi], 如何有效计算 P(\mathbf{x} \mid \lambda) (产生观测序列的概率)? 2)给定模型和观测序列 \{x\_1 \dots x\_n\}, 如何找到最匹配的状态序列(观测序列推隐藏模型状态)? (例子: 语音识别) 3) 给定观测序列如何调整模型参数使 P(\mathbf{x} \mid \lambda) 最大(更好地描述观测数据)? MRF马尔可夫随机场无向图模型(生成式): 一组势函数/”因子”: 定义在变量子集上的非负实函数(用于定义概率分布函数). “团”clique: 其中任意两点间有边连接; “极大团”: 加入相邻任一结点不再形成团. 多个变量之间的联合概率分布能基于团分解为多个因子的乘积, 每个因子仅与一个团有关. 例如: 所有团的集合为\mathcal{C}, 联合概率: P(\mathbf{x}) = \frac{1}{Z} \prod\_{Q \in \mathcal{C}} \psi\_Q(\mathbf{x}\_Q), 其中\psi\_Q是团的势函数, Z是规范化(归一化)因子, 因为变量个数\mathbf{x}较多造成团较多, 所以P(\mathbf{x})中\mathbf{C}^\*可基于极大团定义. 一个结点可能出现在多个团中. 全局马尔可夫性: 其中条件独立性借助”分离”概念即结点集A \rightarrow B 都要经过结点集C, 则C是分离集, A和B条件独立: \mathbf{x}\_A \perp \mathbf{x}\_B \mid \mathbf{x}\_C. 此时联合概率 P(x\_A, x\_B, x\_C) = \frac{1}{Z} \psi\_{AC}(x\_A, x\_C) \psi\_{BC}(x\_B, x\_C); 条件独立的数学证明 P(x\_A, x\_B \mid x\_C) = P(x\_A \mid x\_C) P(x\_B \mid x\_C); 等式左边利用此时: P(x\_C) = \sum\_{x'\_A} \sum\_{x'\_B} P(x'\_A, x'\_B, x\_C) = \sum\_{x'\_A} \sum\_{x'\_B} \psi\_{AC}(x'\_A, x\_C) \psi\_{BC}(x'\_B, x\_C) / Z, \text{ 等式右边: } P(x\_A, x\_C) = \sum\_{x'\_B} P(x\_A, x'\_B, x\_C) = \sum\_{x'\_B} \psi\_{AC}(x\_A, x\_C) \psi\_{BC}(x'\_B, x\_C) / Z, P(\cdot \mid \cdot) \text{ 展开约分即可. 局部马尔可夫性: 给定某变量的邻接变量, 则它条件独立于其他变量: } n^\*(v) = n(v) \cup \{v\}, \text{ 有 } \mathbf{x}\_v \perp \mathbf{x}\_{V \setminus n^\*(v)} \mid \mathbf{x}\_{n(v)}. \text{ 成对马尔可夫性: 若 } \langle u, v \rangle \notin E, \text{ 则 } \mathbf{x}\_u \perp \mathbf{x}\_v \mid \mathbf{x}\_{V \setminus \{u, v\}}. \text{ 势函数: 标定模型偏好变量之间的相关性, 非负} \Rightarrow \text{指数函数: } \psi\_Q(\mathbf{x}\_Q) = e^{-H\_Q(\mathbf{x}\_Q)}, H\_Q(\mathbf{x}\_Q) \text{ 是一个定义在变量 } \mathbf{x}\_Q \text{ 上的实值函数常见形式为 } H\_Q(\mathbf{x}\_Q) = \sum\_{u, v \in Q, u \neq v} \alpha\_{uv} x\_u x\_v + \sum\_{v \in Q} \beta\_v x\_v. \text{ CRF条件随机场判别式无向图模型: 目标构建条件概率模型 } P(y \mid \mathbf{x}), y \text{ 可以是结构型变量(分量之间有相关性). } n(v) \text{ 表示} v \text{ 的邻接结点, 每个变量} y\_v \text{ 都满足马尔可夫性: } P(y\_v \mid \mathbf{x}, \mathbf{y}\_{V \setminus \{v\}}) = P(y\_v \mid \mathbf{x}, \mathbf{y}\_{n(v)}), \text{ 则}(\mathbf{y}, \mathbf{x})\text{构成CRF. CRF类似MRF 用势函数和图结构上的团定义条件概率 } P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z} \exp \left( \sum\_j \sum\_{i=1}^{n-1} \lambda\_{ij} t\_j(y\_{i+1}, y\_i, \mathbf{x}, i) + \sum\_k \sum\_{i=1}^n \mu\_k s\_k(y\_i, \mathbf{x}, i) \right) \text{ 其中 } t\_j \text{ 是定义在观测序列的两个相邻标记位置上的转移特征函数: 刻画相邻标记变量间的相关关系和观测序列对它们的影响; } s\_k \text{ 是定义在观测序列的标记位置} i \text{ 上的状态特征函数: 用于刻画观测序列对标记变量的影响. 例子: } t\_j(y\_{i+1}, i, \mathbf{x}, i) = 1, \text{ if } x\_i = \text{”knock”}, \text{ then } y\_{i+1} = [P], y\_i = [V] \text{ 其他为0; } s\_k(y\_i, \mathbf{x}, i) = 1, \text{ if } x\_i = \text{”knock”}, \text{ then } y\_i = [V], \text{ 形式上MRF和CRF没有区别都是势函数, 但MRF是联合概率. 学习与推断: 基于联合概率分布对目标变量的边际分布或以某些可观测变量为条件的条件分布进行推断, 边际分布是对无关变量求和或积分得到的结果. 概率图模型需要确定具体分布的参数} \Rightarrow \text{参数估计用MAP/MLE, 也可以看出待推测的变量用推断做. 推断: } P(\mathbf{x}\_F \mid \mathbf{x}\_E) = P(\mathbf{x}\_E, \mathbf{x}\_F) / P(\mathbf{x}\_E) = P(\mathbf{x}\_E, \mathbf{x}\_F) / \sum\_{\mathbf{x}\_F} P(\mathbf{x}\_E, \mathbf{x}\_F), P(\mathbf{x}\_E, \mathbf{x}\_F) \text{ 由概率图模型获得, 关键在于高效地计算边际分布} P(\mathbf{x}\_E). \text{ 精确推断: 计算精确值, 一类动规算法, 利用图模型描述的条件独立性来削减计算目标概率的计算量. 变量消去: 设贝叶斯网络结构: } 1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5. \text{ 推断的目标是计算} P(x\_5): P(x\_5) = \sum\_{x\_4} \sum\_{x\_3} \sum\_{x\_2} \sum\_{x\_1} P(x\_1, x\_2, x\_3, x\_4, x\_5)

= ∑x4 ∑x3 ∑x2 ∑x1 P(x1)P(x2|x1)P(x3|x2)P(x4|x3)P(x5|x3), 一层一层地加, 采用加法顺序1-4: {x1,⋯,x4}: P(x5) = ∑x3 P(x5|x3)∑x4 P(x4|x3)∑x2 P(x3|x2)∑x1 P(x1)P(x2|x1) = ∑x3 P(x5|x3)∑x4 P(x4|x3)∑x2 P(x3|x2)m12(x2) = ∑x3 P(x5|x3)m23(x3)∑x4 P(x4|x3) = m35(x5). 最后仅与x5的取值有关. 对无向图也同样适用. 变量消去缺点: 计算多个边际分布, 重复使用变量消去将会造成大量冗余计算. 信念传播: 将变量消去法中的求和操作看作一个消息传递过程, 较好地解决了求解多个边际分布时的重复计算问题: m*ij*(*xj*) = ∑*xi* ψ(*xi*,*xj*) ∏*k*∈*n*(*i*)\j *mki*(*xi*), 其中 *n*(*i*) 是邻接结点, 一个结点仅在接收到来自其他所有结点的消息后才能向另一个结点发送消息, 且结点的边际分布正比于它所接收的消息的乘积: *P*(*xi*) ∝ ∏*k*∈*n*(*i*) *mki*(*xi*). 如果图结构中没有环: 则信念传播两个步骤可完成所有消息传递, 进而计算所有变量上的边际分布: 1)指定一个根结点, 从所有叶结点开始向根结点传递消息, 直到根结点收到所有邻接结点的消息; 从根结点开始向叶结点传递消息, 直到所有叶结点均收到消息. *x4* → *x3* : *m*<sub>43</sub>(*x3*). 近似推断: 两大类 1采样 2变分推断. MCMC采样: 计算概率分布是为了计算期望, 不如直接计算/逼近期望. 关键在构造“平稳分布为*p*的马尔可夫链”来产生样本: 平稳条件如下 **x**, **x'** 是两个状态则 *p*(**x'**) *T*(**x'<sup>t-1</sup> | x'<sup>t</sup>) = *p*(**x'<sup>t-1</sup>) *T*(**x'<sup>t</sup> | x'<sup>t-1</sup>) 此时*p*(**x**) 是该马尔可夫链的平稳分布且满足条件时已收敛到平稳状态. MH算法: MCMC的重要代表, 基于“拒绝采样”来逼近平稳分布*p*, 算法每次根据上一轮采样结果 **x'<sup>t-1</sup>** 来采样获得候选状态样本 **x\***, 但这个候选样本会以一定的概率被“拒绝”掉. 假定从状态 **x'<sup>t-1</sup>** 到状态 **x\*** 的转移概率为 *Q*(**x\* | x'<sup>t-1</sup>) *A*(**x\* | x'<sup>t-1</sup>), 其中 *Q*(**x\* | x'<sup>t-1</sup>) 是用户给定的先验概率, *A*(**x\* | x'<sup>t-1</sup>) 是 **x\*** 被接受的概率. 若 **x\*** 最终收敛到平稳状态, 则根据上述平稳条件有: *p*(**x'<sup>t-1</sup>) *Q*(**x\* | x'<sup>t-1</sup>) *A*(**x\* | x'<sup>t-1</sup>) = *p*(**x\***) *Q*(**x'<sup>t-1</sup> | x\***) *A*(**x'<sup>t-1</sup> | x\***). MH算法伪代码: 循环{ 根据*Q*(*x\** | *x'<sup>t-1</sup>)*采样出候选样本*x\**, 根据均匀分布采样出阈值*u*, 如果*u* ≤ *A*(**x\* | x'<sup>t-1</sup>), *x<sup>t</sup>* = *x\**, 否则 *x<sup>t</sup>* = *x'<sup>t-1</sup>* }, 最后返回采样的样本序列*x*<sup>1</sup>... 为了达到平稳状态, 设置接受率为: *A*(**x\* | x'<sup>t-1</sup>) = min(1, *p*(**x\***) *Q*(**x'<sup>t-1</sup> | x\***)/(*p*(**x'<sup>t-1</sup>) *Q*(**x\* | x'<sup>t-1</sup>))). Gibbs采样是MH算法的特例, 以下步骤: 1)随机或以某个次序选取某变量*xi*; 2)根据 **x** 中除 *xi* 外的变量的现有取值, 计算条件概率 *p*(*xi* | **x**<sub>*i*</sub>), 其中 **x**<sub>*i*</sub> = {*x*<sub>1</sub>, *x*<sub>2</sub>, ..., *x*<sub>*i*-1</sub>, *x*<sub>*i*+1</sub>, ..., *x*<sub>*N*</sub>}, 3)根据 *p*(*xi* | **x**<sub>*i*</sub>) 对变量 *xi* 采样, 用采样值代替原值. 变分推断: 通过使用已知简单分布来逼近需推断的复杂分布, 并通过限制近似分布的类型, 从而得到一种局部最优但具有确定解的近似后验分布. *p*(**x** | Θ) = ∏<sub>*i*=1</sub><sup>*N*</sup> *z*<sub>*p*</sub>(*x*<sub>*i*</sub>, **z** | Θ). 上式取对数似然, 用EM算法: 在 **E** 步, 根据 *t* 时刻的参数 Θ<sup>*t*</sup> 对 *p*(**z** | **x**, Θ<sup>*t*</sup>) 进行推断, 并计算联合似然函数 *p*(**x**, **z** | Θ); 在 **M** 步, 基于 **E** 步的结果进行最大化寻优, 即对关于变量 Θ 的函数 *Q*(Θ; Θ<sup>*t*</sup>) 进行最大化从而求得 Θ<sup>*t+1*</sup> = arg max<sub>Θ</sub> *Q*(Θ; Θ<sup>*t*</sup>) = arg max<sub>z</sub> ∑**z** *p*(**z** | **x**, Θ<sup>*t*</sup>) ln *p*(**x**, **z** | Θ). ln *p*(**x**) = ∫ *q*(*z*) ln *p*(**x**) *dz* = ∫ *q*(*z*) (ln  $\frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z})}$  - ln  $\frac{p(\mathbf{z}|\mathbf{x})}{q(\mathbf{z})}$ ) = *L*(*q*) + KL(*q*||*p*), 在现实任务中**E**步对*p*(*z* | *x*, Θ<sup>*t*</sup>)的推断很可能因*z*模型复杂而难以进行, 此时可借助变分推断, 假设*z*服从: *q*(**z**) = ∏<sub>*i*=1</sub><sup>*M*</sup> *qi*(**zi**). LDA话题模型: 词袋词频的直方图, Θ<sub>*t*,*k*</sub> 即表示文档*t*中包含话题*k*的比例: 1) 根据参数为 α 的狄利克雷分布随机采样一个话题分布 Θ<sub>*t*</sub>; 2) 生成文档中的 *N* 个词: 1根据 Θ<sub>*t*</sub> 进行话题指派, 得到文档 *t* 中词 *n* 的话题 *z*<sub>*t*,*n*</sub> 2根据指派的话题所对应的词频分布 β<sub>*k*</sub> 随机采样生成词. 文档中的词频是唯一的已观测变量, 它依赖于对这个词进行的话题指派*z*<sub>*t*,*n*</sub>, 以及话题对应的词频β<sub>*k*</sub>, 同时话题指派*z*<sub>*t*,*n*</sub>依赖于话题的分布Θ<sub>*t*</sub>, Θ<sub>*t*</sub>依赖于狄利克雷分布的参数 α, 而话题词频依赖于参数η. 对应的概率分布 *p*(**W**, **z**, **β**, **Θ** | **α**, **η**) = ∏<sub>*t*=1</sub><sup>*T*</sup> *p*(Θ<sub>*t*</sub> | **α**) ∏<sub>*i*=1</sub><sup>*K*</sup> *p*(β<sub>*k*</sub> | **η**) (∏<sub>*n*=1</sub><sup>*N*</sup> *P*(*w*<sub>*t*,*n*</sub> | *z*<sub>*t*,*n*</sub>, β<sub>*k*</sub>) *P*(*z*<sub>*t*,*n*</sub> | Θ<sub>*t*</sub>)) 其中 *p*(Θ<sub>*t*</sub> | **α**) 和 *p*(β<sub>*k*</sub> | **η**) 通常分别设置为以 **α** 和 **η** 为参数的*K*维和*N*维狄利克雷分布, 例如 *p*(Θ<sub>*t*</sub> | **α**) =  $\frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \Theta_{t,k}^{\alpha_k-1}$ . 寻找α和η以最大化对数似然 *LL*(**α**, **η**) = ∑<sub>*t*=1</sub><sup>*T*</sup> ln *p*(*w*<sub>*t*</sub> | **α**, **η**), *p*(*w*<sub>*t*</sub> | **α**, **η**) 不易计算, 采用Gibbs采样或变分法来求取近似解. 参数 **α** 和 **η** 已确定, 则根据词频 *w*<sub>*t*,*n*</sub> 来推断 Θ<sub>*t*</sub>, β<sub>*k*</sub> 和 *z*<sub>*t*,*n*</sub> 求解 *p*(**z**, **β**, **Θ** | **W**, **α**, **η**) =  $\frac{p(\mathbf{W}, \mathbf{z}, \mathbf{\beta}, \mathbf{\Theta} | \mathbf{\alpha}, \mathbf{\eta})}{p(\mathbf{W} | \mathbf{\alpha}, \mathbf{\eta})}$ . ch15 规则学习, 优点: 1更好的可解释性, 2数理逻辑具有极强的表达能力. 冲突消解: 多条规则判别结果不同. 命题规则由“原子命题”和逻辑连接词(与或非蕴含)构成的简单陈述句; 一阶规则加入任意, 存在量词. 序贯覆盖: 逐条归纳, 在训练集上学到一条规则, 就将其覆盖的训练样例去除. 刚开始以第一个样例(青绿蜥蜴...)属性赋值生成规则: “色泽=青绿”如果无法仅覆盖正例, 则继续加“根蒂=蜥蜴”, 加起来判断直到仅覆盖正例. 但常常组合爆炸. 改进: 1自顶向下: 从比较一般的规则开始, 逐渐添加新文字缩小覆盖范围. 2自底向上, 从比较特殊的规则开始, 逐渐删除文字以扩大覆盖范围. 适用于训练样本更少的情形; 一阶规则学习(复杂情况). “生成-测试”法, 适用于产生泛化性能好的规则, 命题规则学习. 覆盖准确率: *m*<sub>+</sub>/*m*, 训过程: 1算出所有属性的所有取值对应的准确率, 2选最高并删掉覆盖的样例, 继续第1步. 如上贪心搜索, 还可以束束搜索: 每层搜索保留*b*个, 在下一层加入新的时候都要评估. 剪枝: 增/删规则逻辑文字前后的性能. 借助统计显著性检验: CN2算法用似然率统计量LRS 衡量规则集覆盖的样例分布与训练集经验分布的差别, LRS越大⇒预测与直接猜(在正/负类比例下的)差别. 后剪枝: REP剪错剪枝复杂度O(*m*<sup>4</sup>), 对规则集穷举所有可能的剪枝操作, 然后评估最好的. RIPPER: 剪枝与后处理优化(将*R*中所有规则放在一起优化)相结合, 缓解贪心的局部性, 循环: **R'** = PostOpt(*R*) *D*<sub>*i*</sub> = NotCovered(**R'**, *D*) *R*<sub>*i*</sub> = IREP<sup>+</sup>(*D*<sub>*i*</sub>) **R** = **R'** ∪ *R*<sub>*i*</sub>. 一阶规则学习(此时数据集也变了, 标签变为更好(*X*, *Y*)): 命题逻辑难以处理对象之间的关系, 这里相互比较: 比... 更好. 关系数据集例子: 根蒂更蜥(*X*, *Y*). 优点容易引入领域知识. FOIL算法: 遵循序贯覆盖框架自顶向下, 最初空规则: 更好(*X*, *Y*) ← . 然后候选文字是所有属性于*X*, *Y*的比较, 使用F.Gain =  $\hat{m}_+ \times \left( \log_2 \frac{\hat{m}_+}{\hat{m}_+ + \hat{m}_-} - \log_2 \frac{\hat{m}_+}{\hat{m}_+ + \hat{m}_-} \right)$  来选择文字, 其中 $\hat{m}_+$ 表示增加候选文字后新规则所覆盖的正样例数, 这里F.Gain仅考虑正例的信息量. ILP归纳****************************

逻辑程序设计: 在一阶规则学习中引入了函数和逻辑表达式嵌套: 1具备更强表达能力, 2解决基于背景知识的逻辑程序归纳. 难点: 函数和逻辑表达式的嵌套使候选原子公式可能无穷多. ILP一般自底向上直接将一个或多个正例所对应的具体事实作为初始规则, 再对规则逐步进行泛化以增加其对样例的覆盖率. LGG最小一般泛化: 规则对应的特殊关系数据样例是特殊的 → 一般的. 在归纳逻辑程序设计中, 获得LGG之后可将其看做单条规则加入规则集. 逆归结: 演绎是从一般性规律出发来探讨具体事物, 归纳是从个别事物出发概括出一般性规律. ch16 MDP: 环境*E*, 动作空间*A*, 状态空间*X*, 潜在的状态转移函数*P*, 潜在的奖赏函数*R*. 状态转移/奖赏返回不受机器控制, 机器只能选择执行的动作. 学策略 π, ∑<sub>*a*</sub> π(*x*, *a*) = 1. 找长期累积奖赏最大化的策略. *T*步累积奖赏E[ $\frac{1}{T} \sum_{t=1}^T r_t$ ], γ折扣累积奖赏E[ $\sum_{t=0}^{\infty} \gamma^t r_{t+1}$ ], 强化学习是延迟标记信息的监督学习. 探索-利用窘境: 探索平均分配每种动作最后得到期望的近似估计, 利用动作为当前平均奖赏最大的. ε-贪心: 以ε的概率进行探索(从均匀分布选动作), 其他时候用当前最优: *v*<sub>*n*</sub>为当前获得的奖赏则同步更新的平均奖赏为式\*: *Q*<sub>*n*</sub>(*k*) =  $\frac{1}{n} ((n-1) \times Q_{n-1}(k) + v_n)$ . Softmax: 基于当前已知的平均奖赏来对探索和利用进行折中, 由Boltzmann分布: *P*(*k*) =  $\frac{e^{\frac{Q(k)}{\tau}}}{\sum_{i=1}^K e^{\frac{Q(i)}{\tau}}}$ , 其中τ > 0是温度, → 0是仅利用 → +∞是仅探索. model-based: 在已知模型的环境中学习, 对∀状态*x*, *x'*和执行动作*a* 转移概率 *P*<sub>*x*→*x'*</sub><sup>*a*</sup> 以及带来的奖赏 *R*<sub>*x*→*x'*</sub><sup>*a*</sup> 也是已知的. 策略评估: 状态值函数 *V*<sup>π</sup>(*x*) 为从状态*x*出发使用策略π所带来的的累积奖赏, 状态-动作值函数 *Q*<sup>π</sup>(*x*, *a*)表示从状态*x*出发执行动作*a*后再使用策略π带来的累积奖赏, eg: γ折扣累积奖赏: *V*<sub>γ</sub><sup>π</sup>(*x*) = E<sub>π</sub>[ $\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid x_0 = x$ ], *Q*<sub>γ</sub><sup>π</sup>(*x*, *a*) = E<sub>π</sub>[ $\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid x_0 = x, a_0 = a$ ]. MDP具有马尔可夫性, Bellman等式: *V*<sub>*T*</sub><sup>π</sup>(*x*) = E<sub>π</sub>[ $\frac{1}{T} \sum_{t=1}^T r_t \mid x_0 = x$ ] = E<sub>π</sub>[ $\frac{1}{T} r_1 + \frac{T-1}{T-1} \frac{1}{T-1} \sum_{t=2}^T r_t \mid x_0 = x$ ] = ∑<sub>*a*∈*A*</sub> π(*x*, *a*) ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> ( $\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T-1} E_{\pi} \left[ \frac{1}{T-1} \sum_{t=1}^{T-1} r_t \mid x_0 = x' \right]$ ) = ∑<sub>*a*∈*A*</sub> π(*x*, *a*) ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> ( $\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^{\pi}(x')$ ). 对于 γ 折扣累积奖赏: *V*<sub>γ</sub><sup>π</sup>(*x*) = ∑<sub>*a*∈*A*</sub> π(*x*, *a*) ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ*V*<sub>γ</sub><sup>π</sup>(*x'*)). 由上述递归等式计算值函数实际上为动规算法迭代, 停止准则为*V*(*x*) *V'*(*x*)对于max<sub>*x*</sub>相差一个阈值, 由此计算状态-动作值函数(式③): *Q*<sub>*T*</sub><sup>π</sup>(*x*, *a*) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> ( $\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^{\pi}(x')$ ), *Q*<sub>γ</sub><sup>π</sup>(*x*) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ*V*<sub>γ</sub><sup>π</sup>(*x'*)). 策略改进: 理想策略是最大化累积奖赏, 对Bellman等式做改进对动作的求和取最优(式④): *V*<sub>*T*</sub><sup>\*</sup>(*x*) = max<sub>*a*∈*A*</sub> *A* ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> ( $\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} V_{T-1}^*(x')$ ) *V*<sub>γ</sub><sup>\*</sup>(*x*) = max<sub>*a*∈*A*</sub> *A* ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ*V*<sub>γ</sub><sup>\*</sup>(*x'*)). *V*<sup>\*</sup>(*x*) = max<sub>*a*∈*A*</sub> *Q*<sup>π\*</sup>(*x*, *a*), 代入上面式③得到最优Bellman等式: *Q*<sub>*T*</sub><sup>\*</sup>(*x*, *a*) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> ( $\frac{1}{T} R_{x \rightarrow x'}^a + \frac{T-1}{T} \max_{a' \in A} Q_{T-1}^*(x', a')$ ) *Q*<sub>γ</sub><sup>\*</sup>(*x*, *a*) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ max<sub>*a'*∈*A*</sub> *Q*<sub>γ</sub><sup>\*</sup>(*x'*, *a'*)). 改变动作的条件为 *V*<sup>π</sup>(*x*) ≤ *Q*<sup>π</sup>(*x*, π<sup>*i*</sup>(*x*)), *V*<sup>π</sup>(*x*) ≤ *Q*<sup>π</sup>(*x*, π<sup>*i*</sup>(*x*)) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ*V*<sup>π</sup>(*x'*)) ≤ ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ*Q*<sup>π</sup>(*x'*, π<sup>*i*</sup>(*x'*))) = ... = *V*<sup>π</sup>(*x*), 所以值函数对于策略的改进是单调增的, 所以π<sup>*i*</sup>(*x*) = arg max<sub>*a*∈*A*</sub> *Q*<sup>π</sup>(*x*, *a*)的迭代是可行的, 直到π<sup>*i*</sup>和π差不多即满足最优Bellman等式. 策略迭代: 评估策略的值函数与策略评估后改进值最优策略结合的迭代. 值迭代算法: 根据式④更新值函数直到更新后的相差一个阈值. model-based方法能归结为基于动规的寻优问题, 估计状态值函数*V*, 最终策略通过状态-动作值函数*Q*获得. model-free: 不依赖环境建模, 策略无法评估因为模型未知导致无法做全概率展开, 蒙特卡罗强化学习: 多次采样替代策略评估, 平均累积奖赏作为期望累积奖赏的近似. 模型未知时*V* → *Q*的转换是困难的. 多次采样就是综合多个轨迹的奖赏后得到状态-动作值函数*Q*的估计, 同策略(on-policy, 被评估与被改进是同一个策略) ε-贪心: 策略评估用式\*, 改进用得到不同的轨迹需要用不同的策略π<sup>ε</sup>: 即以ε概率均匀概率选取动作, 否则采用原始策略 arg max<sub>*a'*</sub> *Q*(*x*, *a'*). 异策略(off-policy)蒙特卡罗强化学习: 使用策略π<sup>*i*</sup>的采样轨迹来评估策略π 则对累积奖赏加权: *Q*(*x*, *a*) =  $\frac{1}{m} \sum_{i=1}^m \frac{P_i^{\pi}}{P_i^{\pi'}} r_i$ , *P*<sup>π</sup> = ∏<sub>*i*=0</sub><sup>*T*-1</sup> π(*x*<sub>*i*</sub>, *a*<sub>*i*</sub>) *P*<sub>*x*<sub>*i*</sub>→*x*<sub>*i*+1</sub></sub><sup>*a*<sub>*i*</sub></sup> 是策略产生第*i*条轨迹的概率, *R*<sub>*i*</sub>是第*i*条轨迹的累积奖赏. 伪代码: 因为ε贪心, *p*<sub>*i*</sub>为1选择原策略的概率: 1 - ε + ε/|*A*|, 或2平均取的概率: ε/|*A*|; 策略评估: *R* =  $\frac{1}{T-1} \sum_{t=1}^T (r_t + \prod_{j=t}^{T-1} \frac{1}{p_j})$ , *Q*(*x*<sub>*t*</sub>, *a*<sub>*t*</sub>) = (*Q*(*x*<sub>*t*</sub>, *a*<sub>*t*</sub>) × count(*x*<sub>*t*</sub>, *a*<sub>*t*</sub>) + *R*)/(count(*x*<sub>*t*</sub>, *a*<sub>*t*</sub>) + 1), 在一条轨迹后进行策略改进 arg max<sub>*a'*</sub> *Q*(*x*, *a'*). TD时序差分学习: (因为蒙特卡罗强化学习没有充分利用MDP结构, 采样一个轨迹后才能更新值估计不够高效), 利用值函数*Q*的增量更新来评估策略 (Sarsa算法): 迭代 1: *a'* = π<sup>ε</sup>(*x'*) (因为 *Q*<sub>*t*+1</sub><sup>*a*</sup>(*x*, *a*) = *Q*<sub>*t*</sub><sup>*a*</sup>(*x*, *a*) +  $\frac{1}{t+1} (r_{t+1} - Q_t^a(x, a))$ ), *Q*<sup>π</sup>(*x*, *a*) = ∑*x'*∈*X* *P*<sub>*x*→*x'*</sub><sup>*a*</sup> (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ ∑*a'*∈*A* π(*x'*, *a'*) *Q*<sup>π</sup>(*x'*, *a'*))) 2: 通过增量求和有 *Q*<sub>*t*+1</sub><sup>*a*</sup>(*x*, *a*) = *Q*<sub>*t*</sub><sup>*a*</sup>(*x*, *a*) + α (*R*<sub>*x*→*x'*</sub><sup>*a*</sup> + γ *Q*<sub>*t*</sub><sup>*a*</sup>(*x'*, *a'*) - *Q*<sub>*t*</sub><sup>*a*</sup>(*x*, *a*)), 其中更新步长α越大则越靠后的累积奖赏越重要 3: 式\*的*Q*更新. *Q*学习算法中仅在策略评估时*a'* = π(*x'*)不同. 值函数近似: 值函数是关于有限状态的表格值函数, 状态空间离散化后, *E*<sub>θ</sub> = E<sub>*w*~π</sub> [(*V*<sup>π</sup>(*x*) - *V*<sub>θ</sub>(*x*))<sup>2</sup>] 梯度下降 -  $\frac{\partial E_{\theta}}{\partial \theta} = \mathbb{E}_{w \sim \pi} \left[ 2 (V^{\pi}(x) - V_{\theta}(x)) \frac{\partial V_{\theta}(x)}{\partial \theta} \right] = \mathbb{E}_{w \sim \pi} [2 (V^{\pi}(x) - V_{\theta}(x)) \mathbf{x}]$ , **θ** = **θ** + α (*V*<sup>π</sup>(*x*) - *V*<sub>θ</sub>(*x*)) **x** 是对于单个样本的更新规则. 借助时序差分学习, 不知道策略的真实函数*V*<sup>π</sup>时用当前估计的值函数代替真实值函数: **θ** = **θ** + α (*r* + γ**θ**<sup>T</sup>*x'* - **θ**<sup>T</sup>*x*) *x*. 模仿学习: 从人类专家的决策过程范例中学习, 直接模仿学习: 模仿人类专家“状态-动作对”数据集, 状态: 特征, 动作: 标记, 学策略模型可作为强化学习初始策略. 逆强化学习: 从人类专家提供的反例数据中反推出奖赏函数. 基本思想: 欲使机器做出与范例一致的行为, 等价于在某个奖赏函数的环境中求解最优策略, 该最优策略产生的轨迹与范例数据一致. 寻找某种奖赏函数使范例数据最优, 使用该奖赏函数训练强化学习策略. ρ<sup>π</sup> = E[ $\sum_{t=0}^{\infty} \gamma^t R(\mathbf{x}_t) \mid \pi$ ] = E[ $\sum_{t=0}^{\infty} \gamma^t \mathbf{w}^T \mathbf{x}_t \mid \pi$ ] = **w**<sup>T</sup> E[ $\sum_{t=0}^{\infty} \gamma^t \mathbf{x}_t \mid \pi$ ], 对于最优奖赏函数 *R*(*x*) = *w*<sup>\*</sup><sup>T</sup>*x* 和任意其他策略产生的 *x*<sup>π</sup>, 有 *w*<sup>\*</sup><sup>T</sup>*x*<sup>π</sup> - *w*<sup>\*</sup><sup>T</sup>*x*<sup>π</sup> = *w*<sup>\*</sup><sup>T</sup>(*x*<sup>π</sup> - *w*<sup>\*</sup><sup>T</sup>*x*<sup>π</sup>) ≥ 0, *w*<sup>\*</sup> = arg max<sub>*w*</sub> min<sub>π</sub> *w*<sup>T</sup>(*x*<sup>π</sup> - *x*<sup>π</sup>) s.t. ||*w*|| ≤ 1. 加油冲冲冲, 越哥罩我!!