

NLP 基础技术: 词法分析 (词性标注和词义标注); 句法分析 (判断成分和句法结构, 有完全/浅层句法分析); 语义分析; 语用分析 (具体运用); 篇章分析 (整体理解分析) **NLP 应用技术:** 机器翻译, 信息检索, 情感分析, 自动问答, 自动文摘, 社会计算, 信息抽取, **正则表达式:** [A-Z] (从 A 到 Z), [123] (匹配 1 或 2 或 3), [^Aa] (不是 A 或者 a), [a^] (a 和 ^, 因为只在[后第一个有效), a|bc (替代产生式), ? (上一个字符是可选的), + (出现一次或者多次), * (出现 0 次或者更多次), . (点, 匹配任意单个字符), \$ (在结尾匹配, 加在 RE 的最后), ^ (在开头匹配, 加在 RE 的最前 (注意在[]的外面)), 匹配所有开头/结尾的字符。**有限状态自动机 Finite State Automata 与正则:** 状态 Q, 字母表 Sigma, 初始状态, 终极状态, 转一函数。有限状态机定义了一种形式语言 (Formal Language)。Generator & Acceptor

错误的类别: 假阳性 (1 类, 不该匹配却匹配, 降低假阳性提升精确度), 假阴性 (2 类, 该匹配却没有匹配, 降低假阴性提升覆盖率) **DFA:** 同一个状态/输入只能通向一个状态, **ND-FA:** 同一个状态/输入可能通向多种状态, 需要遍历所有状态, 支持通过空串 epsilon 转向另一等价的字符 (tape 不前进) **分词:** 分词, 规范化, 分句。**词元 (Lemma):** 同一个词干 (stem) 和词性 (part of speech), 大致相同的词义 (cat=cats) **词形 (Wordform):** 词的表面形式 (把单复数啥的变化都加上) (cat=cats) **词型 (Type):** 一个单词

词例 (Token): 词型在文章中的一个实例, **N 词例数>>V 词汇量**, **分词时**的歧义: Finland's, state-of, what're, PhD. 法语, 德语复合词; **方法:** 空格分词, 字符分词, subword-toeknization (BPE, Wordpiece, ...) **Learner:** 从语料库学习词汇报, **Segmenter:** 分词 **中文分词:** Baseline 方法是贪心(最长匹配法)

形态学 (Morphology): 研究单词是如何从语素 (Morphemes) 构造 **语素 (Morpheme):** 词干 (Stem, 核心意思) 和词缀 (Affix, 附加/修改意思) 两种广义的**构造形式:** ①**屈折 (Inflectional):** 不改变词类的词缀 (walk, walking) ②**派生 (Derivational):** 改变意思和词类 (clue, clueless)

词干还原 (Stemming): 只关心词干, 不关心结构, 常用于信息检索应用; 比如 **Porter Stemmer**, 不需要词汇报, 基于规则去词缀, 不保证产生真实词干, 但不影响 IR. (ization->ize->()) **切分 (Segmentation):** 不能单纯按照标点/空格切分句子/单词, 因为没有歧义, "有很多歧义" **断句 (Segmenting Sentences):** 用二分类器 (EOS/NotEOS), 基于规则或 ML 来判断断句", 是否为一句语法的结束

最小编辑距离: 在插入, 删除和替代意义下的最少编辑距离; **应用:** 评估机器翻译和语音识别的效果 (词为单位); **命名实体 (Named Entity) 和 (Entity Conference) 识别; Levenshtein** 插入和删除代价为 1, 替换代价为 2。一般代价都为 1

计算方法: 设 D(i,j) 为 A[1...i] 和 B[1...j] 的最小编辑距离, 目标是让 A 靠近 B. (Levenshtein)

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1 & \text{deletion (X删除Y)} \\ D(i,j-1) + 1 & \text{insertion (X添加Y)} \\ D(i-1,j-1) + \begin{cases} 0 & \text{if } X(i) = Y(j) \\ 1 & \text{if } X(i) \neq Y(j) \end{cases} & \text{substitution (X替换Y)} \end{cases}$$
$$D(0,0) = 0$$
$$D(i,0) = i$$
$$D(0,j) = j$$

初始化 D(i,0) = i, D(0,j) = j
如果 Xi=Yj, 则认为在这里是对齐的; 为了跟踪对齐情况, 用 ptr 维护一个箭头数组跟随着 insert (LEFT), delete (DOWN), subst (DIAG); 时间以及空间复杂度 O(nm), Backtrace 输出 O(n+m)

带权最小编辑距离 (Weighted Edit Dist): 维护单个字母 del 和 ins 权重, 两字母之间 sub 权重; 用于修正一些拼写错误(有些字母对更容易打错; 有些插入/删除更容易发生) **语言模型:** P(w₁ | w₁, ..., w_{n-1}), 用链式法则计算句子 (单词序列) 的概率。**但是由于数据不足, 假设 Markov 性质成立 (P(x₁a, b, c, d, ...) = P(x₁a))** (二元情况, 也可以 N 元); Bigram 即为 Markov 链, 预测用 MLE 选出最大可能的下一个词,

$$PP(w_i | w_{1-i-1}) = \frac{c(w_{1-i-1}w_i)}{c(w_{1-i-1})}$$

句子生成: 根据概率选<s>, 然后根据给定词为条件, 出现下一个词的概率选, 直到选择</s>
Unigram (不用条件概率); Bigram (用上个为条件); Trigram; Quadrigram

封闭词汇任务和开放词汇任务 (没见过的替换为<UNK>)
评价 N-gram 模型-测试集 ①外在评测, Word Error Rate; ②内在评测, 用困惑度 (Perplexity, 多用于先期自测):

$$PP(W) = P(w_1w_2...w_N)^{-1/N}$$

$$PP(W) = \prod_{i=1}^N \frac{1}{P(w_i | w_{1-i-1})}^{-1/N}$$

最小化 PP 就是最大化整个句子在模型中的出现概率 (希望测试集与训练集类似); 开 N 次根号用来做某种关于模型状态空间的归一化补偿; 越小越好

问题: ①过拟合, 测试集和训练集相差很大则效果不好; ②很多概率是 0: 进行平滑: (1)**Laplace 平滑**, 每个 Ci 都加一, 则 P_{Laplace}(w_i) = (c_i + 1) / (N + V); Bigram 如上, c_i' = (c_i + 1)N / (N + V) 加 k 法缺点是对 N 太多的数据集, 非 0 的概率会极大稀释, 适合 0 不多的数据

$P(w_{n-1}w_n) = (c(w_{n-1}w_n) + 1) / (c(w_{n-1}) + V)$
(2)Good-Turing 平滑法, 估计没出现的概率用 P₀ = N₁ / N, c_n' = (c_n + 1)N_{n+1} / N_n, P_n' = c_n' / N, 其中 N_n 为出现 x 次的词的出现次数。对于任何一个出现了 c 次的类别, 都假设

它出现了 c' 次。
回退 (Backoff): 如果更高阶的 Markov 没有出现, 就回退到低阶的 Markov 过程对概率进行估计

$$P(z|xy) \approx P(z|x) \approx P(z)$$

内插 (Interpolation): 将不同阶输出结果线性插值, 权重可以和前面的词相关; 先固定 N-gram 概率, 再找到最优权重使 held-out dataset prob 最大, 可以采用搜索算法找到最优权重 (比如 EM 算法)

$$P(z|xy) = \lambda_1 P(z|xy) + \lambda_2 P(z|x) + \lambda_3 P(z)$$

词汇外词汇(OOV): 不能用 GT, 因为不知道总数, 依旧全转化为一个<UNK>来看

应用: 我们一般用 exp [log p_i] = ∏ p_i 来算对应的乘法: 避免下溢+加快速度

N-gram 优点: 容易构建, 可以使用平滑来适应新数据; 缺陷: 只有在测试集与训练集比较相似的情况下表现较好, 只能捕捉到较短的结果; **神经网络:** 非线性, 前/反向传播, SGD 适应能力强, 但训练消耗相对较大, 词嵌入, 词向量, softmax 需要对所有单词求和

词类标注 (POS(Parts of speech 词类) Tagging): ①**基于规则的方法** ②**概率方法(HMM):**

基准方法: 无脑选最大类, 查表+无脑选, RE 正则法。HMM: 一些状态 S, 一些观测值 O, 关于状态 S 的转移概率矩阵 A_{ij}, 输出概率矩阵 B_i(k), 即 S=i 时观测值为 k 的概率, 以及初始状态 S 的概率分布。

HMM Tagging: 隐状态是各个 POS, 输出是各个词本身的 HMM; ① **(评估或计算得分问题)** 如何计算给定观察序列出现的概率? 遍历所有可能的状态序列, 乘上对应的状态->输出序列得到概率, 累加。

② **(解码问题)** 给定观察序列, 如何计算最优的隐状态序列? 选出最大概率对应的隐状态序列 t_i = argmax_P(t_i | w_i) = argmax_P(w_i | t_i) P(t_i) ≈ argmax ∏ P(w_i | t_i) P(t_i | t_{i-1}) ③ **(训练问题)** 如何调整模型参数来最大化某特定观察序列的概率? 定义 a 为状态转移概率, b 为观察概率。

对于①, 定义 a_t(i) = P(o₁...o_t, q_t = S_i | λ), 即第 t 个观察值对应的隐状态为 i 时的输出概率, 则有递推

$$a_{t+1}(j) = \sum_{i=1}^N a_t(i) a_{ij} b_j(o_{t+1})$$

成立; P(O | λ) = ∑_{i=1}^N a_T(i); 此方法也可以用来反向计算 (Backward Algorithm): 先初始化 β_T = 1

对于②, 使用 Viterbi 算法, 定义 δ_t(i) = max_P [q₁...q_{t-1}, q_t = i, o₁...o_t | λ], 即给定前 t 个观察序列 q₁...q_{t-1} 个隐状态最优, 第 t 个状态为 i 的概率, 则递推关系如下:

$$\delta_t(j) = \max_{i \in I} [\delta_{t-1}(i) a_{ij} b_j(o_t)]$$

$$\gamma_t(i) = \sum_{j=1}^N \delta_{t-1}(j) \delta_t(i)$$
$$\gamma_t(i) = P(q_t = S_i | O_{1:t})$$

同时, 为了记忆最优状态, 设置 η_t(j) = argmax_i [δ_{t-1}(i) a_{ij}] (每个 t 都要记录 N 次, 最后回溯); 为了加速也可以进行对数化。

$$\delta_t(i) = \pi_i b_i(o_t)$$

对于③, 现在还没有全局最优解, 思路大概有(1)梯度下降(2)EM 或者 Baum-Welch 迭代搜索; 重估计的一种方法如下: 则 ∑_{t=1}^T ξ_t(i,j) 表示从 S_i 转移到 S_j 的期望总转移边数, ∑_{t=1}^T γ_t(i) 表示从 S_i 转移出去的期望总边数; 那么, 可以得到新的估计转移概率和输出概率 π_i' = γ_i(i), a_{ij}' = ∑_{t=1}^T ξ_t(i,j) / ∑_{t=1}^T γ_t(i), b_j(k) = ∑_{t=1}^T γ_t(j) / ∑_{t=1}^T γ_t(j) (v_k 表示在 O_t 观察到符号 k); <UNK> 在处理中可以根据形态学猜测其词性 (比如 -s, -able 等), Forward 就是加起来

$$\xi_t(i,j) = \frac{P[q_t = S_j, q_{t-1} = S_i, O_{1:t}]}{P(O_{1:t})}$$
$$= \frac{a_{ij}(i) b_j(o_t) \beta_{t-1}(j)}{P(O_{1:t})} = \frac{a_{ij}(i) a_{ij} b_j(o_t) \beta_{t-1}(j)}{P(O_{1:t})} = \frac{\sum_{i=1}^N \sum_{j=1}^N a_{ij}(i) a_{ij} b_j(o_t) \beta_{t-1}(j)}{\sum_{i=1}^N \sum_{j=1}^N a_{ij}(i) a_{ij} b_j(o_t) \beta_{t-1}(j)}$$

重新估计转移矩阵: 比如 1-1, 1-2, 用有 1-1 转移的序列概率乘转移次数, 1-2 转移的序列概率乘转移次数, 两者归一化(对一个状态出发的所有转移归一化)。**最大熵马尔可夫链:** p(c|x) = (1/z) exp (∑ w_f f_i) 状态是由观察值和上一个状态生成的, 优化可以直接用最大似然估计, 给定观察序列之下哪种隐状态的序列概率最高。**评估:** 用人来分词, 作为金标准 (97%左右), 和机器结果进行比较, 得到 Confusion Matrix, 可以找出最常见错误。**POS 在中文中: 更迷惑, 没见过的字有问题, 神经网络最佳。**

句法和语法: CFG(Context-Free Grammar) 上下文无关语法: 规则或者产生式, 捕捉顺序和组成性
英语核心构成: 句子 (表达完整的思想), 子句 (有一个动词), 短语 (一种词的聚合) S->NP VP, S->VP, S->Aux NP VP, 递归: N->N PP

用 CFG 解析文法的问题: ①**一致关系** (名词词性和动词的单复数形式应该对应, 可以细化语法产生式, 但不够简洁美观, 泛化能力不够强) ②**次范畴化** (谓语和后边的“参数”类型应该满足一些语义上的约束); 对谓语 (基本是动词) 进行次范畴化, 约束框架称为次范畴化框架, 增加语法规则, (find 需要 NP), (补语), V-with-NP-complement, **VP->V-with-NP-comp NP** ③**移位** (倒装等特殊句式使动词和宾语的位置发生变化, 用 CFG 很麻烦, 需要考虑很多基础语法) **依存文法:** 记录词之间的论旨角色 (thematic roles), 一般是二元的; 没有非终结符。适合印度、欧洲。United <= canceled (NSUBJ) 文法库的来源: ①手工构造 ②TreeBanks, 根据 POS 过的句子自动生成文法, 自动解析, 手工修正. 很多冗余。

句法分析: DP 方案: 有效存储二义句子, 避免重复工作, 在多项式时间内完成指数任务。

Earley 自顶向下 O (n^3): 从树库出发, 有可能有和输入无关的树; **CKY 自底向上 O(n^3):** 从输入出发, 有可能有非法树。**CKY:** 文法全部写为 Chomsky 范式 (CNF) (A->BC 或 A->w(两个非终止/一个终止)); 索引 0,1,...,n 分别代表第一个单词前, 第二个, 第 n 个单词后; CKY 表格中填写可能的非终结符, 计算 table[i,j] 时穷举从 i+1 到 j-1 这些可能分隔, 并且检测是否是有效的 Chomsky 范式。

缺陷: 非二元语法需要转成二元的 CNF, 不是 Chomsky 范式的文法可以重写成 Chomsky 的, 可能产生一些无关的成分, S=>ABC, S=>XC, X=>AB, X 不在其他地方出现。**Partial Parsing:** 速度较快, 只处理大的 (某种, 比如 NP) 组块 (chunk)

在哪; **基于规则的 Partial Parsing:** 去掉规则集中的递归元素 (变成正则文法)

POS=>句法短语=>更大的动词/名词组=>句子层面规则
基于 ML 的 Chunking: 一个序列分类任务, 按序列顺序处理, 其把词分为三类: 开始, 内部和外部, 外部表示其不在任何组块中, 内部表示其在某组块中 (比如 LNP 在 NP 中), 开始表示其为某块的开始 (比如 B_NP), 这样就能把所有 NP Chunk 找出来

概率 CFG: 对于特定领域的语言以及难以理解的语言高有效。给文法规则都赋一个概率, 用来解决二义性, 同时模仿人类解析语言的过程; 非二义的情况下就是该文法的概率, 二义的情况是 n 者相加; 解析树的每个文法都放有一个概率值。**树的概率**是树中所有规则概率乘积, 句子的概率是所有可能树的概率之和。求 PCFG 文法概率: 用 Treebank 里面的数据进行统计)

若没有 TB, 假设有确定性语法进行 Parse 后进行重新统计得出概率。对 A->BC 来说概率其实是 P(A->BC) * P(B) * P(C), 只是 DP 的时候 P(B) 和 P(C) 肯定都已经计算出来了; 如果有 B 和 C (来源于不同的分点, 但都合法) 则取最大概率的

$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{Count}(\alpha \rightarrow \beta)}{\sum_{\gamma} \text{Count}(\alpha \rightarrow \gamma)} = \frac{\text{Count}(\alpha \rightarrow \beta)}{\text{Count}(\alpha)}$$

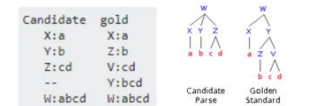
function PROBABILISTIC_CKY(words, grammar) returns most probable parse and its probability
for j ← from 1 to LENGTH(words) do iterate over columns:
 subtable[j-1,j] ← [A | A → words[j]] ∈ grammar
 for i ← from j-2 to 0 do iterate over rows (from bottom up)
 for i ← i+1 to j-1 do iterate over all possible splits
 subtable[i,j] ←
 [A | A → BC ∈ grammar;
 B ∈ subtable[i,i];
 C ∈ subtable[j,j]]

那种。由此可以有改进的 maxprob CKY。**结构依赖:** PCFG 的问题: 没有考虑到推导环境的

上下文 (只基于很小的上下文得到一个概率), 导致最优推导和 TreeBank 中的实际推导存在很大差异

解决方案: ①**父节点标注**, 将非终结符一分为 n, 每个都标上其父节点的信息 (NPNP_{S_{NP}}), 但是增加了语法复杂性, 并且需要更多数据②**中心词** (NP 的中心词是名词, VP 的中心词是动词, PP 的中心词是介词), 现在文法符号变成类似 VP(中心词); 由于句子太少, 规则不变, PCFG 的规则变成类似 VP(dumped)=>V(dumped)NP(sacks)PP(in), P(rule|VP^dumped^sack^in), 以中心词为条件的规则的概率 P(rule|VP^dumped) = Count(这个规则出现次数) / Count(VP^dump^rule 出现次数); 这其实是一种次范畴化捕捉到了中心词和语法的关联度; 同时, 应该加入一种优先级, 描述中心词和上层文法之间的选择关联度, 比较 dumped 作为 head 且有一个 PP 的 head 是 into 的概率和 sacks 作为 head 且有一个 PP 的 head 是 into 的概率即可选择正确的。

Parser 评估: 构成评估 (Constituent-level Evaluation): 句子水平太粗糙了, 分为覆盖率 (Recall, 你的解析结果中正确 [即有相同节点标签和正确划分的节点] 的节点数 / treebank 中的相应 constituent 节点数) 和精准度 (Precision,



你的解析结果中正确的节点数/你的解析结果中的节点数) Recall = 2/5; Precision = 2/4;

交叉括号 (Cross Brackets): treebank has ((X Y) Z) and candidate has (X (Y Z)) 的情形; 这个应该作为客观的函数来最小化, 统计多少对话号穿过 tb 的括号。上述方法的缺点: 会更偏向“安全的, 浅的解析”; 部分错误可能不停向上传播, 一个节点的错误导致很多交叉括号的情况; 将所有节点一视同仁, 而不是更关注核心的语义关系

语义: 单词含义, 三种向量语义模型: ①**稀疏向量表示**(l)以互信息 (Mutual-information) 为权重的单词关联矩阵 ②**稠密向量表示**(l)奇异值分解和潜在语义分析 (Latent Semantic Analysis) (2)各种神经网络模型 (Skip-grams, CBOW) **词项-文档矩阵 (Term-Document Matrix)**, tf_(t,d) 表示词 t 在 d 文档中出现的次数; 两列相近则两个文档相近; 两行相近则两个词相近。**词-上下文 (word-context) 矩阵**, a_{ij} 表示在段落里每次单词 i 的 ±4 (或者别的数) 的范围内出现单词 j 的次数之和; 矩阵会很稀疏, 并且窗口大小和目标有关 (1-3 窗口越小语义信息越多, 4-10 窗口越大语义信息越多) 两个行相近则两个词相近。

两个单词间的两种互相关 (co-occurrence): ①**一阶互相关 (syntagmatic 组合)**, 这些词基本都靠在一起, 比如 wrote 和 book ②**二阶互相关 (paradigmatic 聚合)**, 这些词的邻居相近, 比如 wrote 和 said
直接计数 (raw counts) 的问题: 直接用原始词频考察词-

