

第二节网络爬虫 反爬虫策略 (1):User Agent 利用访问网站时浏览器发出的 Request Headers 信息中的 User-Agent 信息，判断用户是使用何种浏览器。爬虫往往在 U-A 部分空应对：利用 Python 爬虫库 爬取用户自定义请求头信息的手段，在请求头信息中将 User-Agent 的值改为浏览器的请求头标识，从而绕过反爬虫机制。
反爬虫策略 (2):IP/账号访问次数/频率 通过限制特定 IP 地址/账号访问频率和次数进行反爬。本策略判断爬取行为是否人类行为：应对：构造 IP 代理池，然每次访问时随机选择代理池,Github 中有相关资源,通过各种免费 IP 的网站来提供账号,每次爬取行为后间隔一段时间;注册多个账号/IP来保障数据收集(账号本号会出问题,账号会查封的账号)。
反爬虫策略 (3):验证码 通过各类验证码,判断网页是否属于人类还是机器;应对：简单的验证码识别(基于机器学习与模式识别相关技术),复杂的逻辑推理(人工辅助破解);验证码是一把双刃剑,在抵御机器人的同时也伤害用户体验。
反爬虫策略 (4):动态内容加载 从网页的加载时代码之后,会在浏览器页面执行 JavaScript 脚本,网页内容会动态由脚本加载,而直接抓取则得不到空白页面。此类情况在抓去在线播放的音视频文件时尤为常见;应对：核心策略是模拟浏览器请求,使用浏览器元素 ajax 请求,如此便能得到包含数据信息的 json 文文件。
反爬虫策略 (5):蜜罐技术 网页上故意留下一些人类不到或绝对不包含数据的链接,立刻会触发蜜罐代码中获取内容,所以蜜虫可能会访问这些链接。只要网页没有 IP 访问的链接,网页就不会被访问,从而难以继续爬取。应对：依靠思路是干涉爬虫路径,通过工具判断页面上的隐藏元素,使爬虫避开这些元素,可以部分回避蜜罐的诱导。
反爬虫策略 (6):加密-解密技术 对网页中的关键信息进行加密或混淆,页面上显示为正常文本,解密代码为编号,网页加载后解密助手解密,将编号转为文字。字/体/密/钥/可以动态更新。应对：通过加密获取/解密/映射,对于映射进行解密。
反爬虫策略 (7):用户权限限制 不同等级/级别的网页给予不同的内容权限,对于新手：解密,就可以查看。其他（其他的用户限制）:各个相同类型的页面的解锁格式内容不同(双刃剑,增加爬取难度的同时降低用户浏览体验)。多模态的呈现方式:文字转为图像或视频;应对策略:OCR、语音识别、图像/视频标签技术。

第三节网页文本处理 将原始文档转化为词项，建立索引，使文档匹配**精准**而向查询条件。近似**重复**文字的**检测**方法：**指**标表示法：1 对文档进行分词处理，并进行 n-gram 组合，2 挑选部分 n-gram 用于表示这一文档，3 对被选中的 n-gram 进行散列，4 存储散列值作为文档特征。**基于匹配的字符串匹配**（机械故障检测方法）按照一定的策略将待分析的字符串与一个“充分大的”机器结构的词条进行匹配，若在词典中找到某个字符串，则匹配成功。**基于匹配字符串的一般模型**：其形式依表达为 *ASM* (*d*, *a*, *m*)，(Automatic Segmentation Model)。其中 *d* 表示匹配方向，+1 为正向，-1 为反向；*a* 为每次匹配失败后增/减字符数，+1 为增序，-1 为减序；*m* 为最大/最小匹配数，+1 为最大匹配，-1 为最小匹配。如 *ASM* (+, −, +) 即正向最大匹配 (FMM)；对于现代汉语，最大匹配更实用（最小匹配过于精确）**正向最大匹配**词项**Forward Maximum Matching** 从左至右可能查找到的词，直到当前符号与已经处理的字符串不构成词项，输出已经识别的词，并从识别后的字符串接着继续查找下一个词。分词速度较快，但错误率较高（约 1/169）。**反向最大匹配**词项**Reverse/Backward MM** 从右至左可能查找到的词，直到当前符号与已经处理的字符串不构成词项，统计证实 RMM 分词效果更好（约 1/245）。**双向最大匹配**词项**Bi-dir M** 综合比较 FMM 与 RMM 两种方法的切分结果，从而选择最优的切分方法用于识别字符串中的交叉文义。可選擇直接基于与词最接近的方法得到最终结果。**最确切的分词方法** 使句子中切出的词数目最少，等价于在向前句中搜索最短路径的问题。将每个字符串节点，每个形成一条边。边权值都可视为 1，也可根据词频决定（尽量用高频词）结合元音/辅音之后，可视为基于多字分的分方法。**N-最短路径法** *table*(*x*) 表示到 *x* 的路径，保留 *N* 条最短的路径，以提供更多分词方案。每个表项包含路径（路径编号，从 1 到 *n* 开始）、路径长度（路径长度）、前驱。一个表项可能包含多个前驱，即可能有相同长度的多条路径到达 *x*。前驱部分用 *C*(*x*,*y*) 表示，指第 *x* 个节点（亦即 *table*(*x*) 中编号为 *y* 的路径。**基于匹配字符串的优点** 优点是效率高、直接性好；缺点是对于同义性：维护前驱词典成本大、难以应对新生词汇、词汇频率/重量性不影响结果。**基于统计的分词方法** 统计海量文档，字与字相邻共现的概率能够反映词频的可信度。如果某些词语组合在统计上出现的频率非常高，就认为它们是词。**统计分词的形式化表达** *c* = *c*₁*c*₂...*c*_{*n*} 是待分词的句子/字符串，而 *w* = *w*₁*w*₂...*w*_{*m*} 是切分的结果，设 *P*(*w*₁*c*₁) 为 *c* 切分为 *w* 的某种统计概率，*w*₁, *w*₂, ... , *w*_{*k*} 为 *c* 的所有可能的切分结果，基于统计的分词模型就是找到目标词 *w*，使得 *w* 满足 *P*(*w*₁*c*) = max{*P*(*w*₁*c*₁), *P*(*w*₁*c*₂), ..., *P*(*w*₁*c*_{*n*})}**统计分词的一般化过程** 1.建立统计语言模型，2.对于每个子词分方案进行分词，3.计算不同分词方案的概率，选出概率最大的分词结果；理论上，可以不需要词频，但实际操作中第 2 步可以采用机械式的方法进行统计，以获取最优的词组合（配匹分词时切速度快,效率高而第 2 步又可以利用机械式词组合上下文识别生字,自动消除歧义）**N-gram 模型与马尔科夫假设** 马尔科夫假设指 N-gram 由 N 个单词组成的集合，各单词具有先后顺序。模型的概率与假设：当前状态出现的概率只过去有限的历史状态有关，并与其状态无关。具体到分词任务，就是文本中第 N 个词出现的概率仅仅依赖于它前面的 N-1 个词，而与其他词无关。• *N* = 1，一元文法模型（最大概率模型），*P*(*w*₁) = *P*(*w*₁) *P*(*w*₂)...*P*(*w*_{*n*}) • *N* = 2,Bigram 模型,*P*(*w*₁) = *P*(*w*₁) *P*(*w*₂|*w*₁)...*P*(*w*_{*n*}|*w*_{*n*−1}) • *N* = 3,Trigram 模型,*P*(*w*₁) = *P*(*w*₁) *P*(*w*₂|*w*₁) *P*(*w*₃|*w*₁*w*₂)...*P*(*w*_{*n*}|*w*_{*n*−2}*w*_{*n*−1})**N-gram 模型的概率估计** 以 Bigram 模型为例，基于最大似然估计进行推断 *P*(*w*₁*w*_{*n*−1}) = $\frac{C(w_{n-1}w_n)}{C(w_{n-1})}$ 其中，*C*(*c*)：指词类 *c* 在语料库中出现的次数。**N-gram 模型的分词过程** 以 Bi-gram 模型为例，1.首先，构造语料库统计，计算所有 *C*(*w*₁) 与 *C*(*w*_{*n*−1}) 2.其次，对于每一个可能的分词序列 *w*，计算 *P*(*w*) = *P*(*w*₁) *P*(*w*₂|*w*₁)...*P*(*w*_{*n*}|*w*_{*n*−1}) 3.最后，返回最大的 *P*(*w*) 所对应的分词序列作为结果。**特殊形式：一元文法模型** 当 *N*=1 时，N-gram 模型退化为一元文法模型，此时词与词之间独立的。（独立词假设，一元文法模型 *P*(*w*₁) = *P*(*w*₁) *w*₂...*w*_{*n*}，...*w*_{*n*} ≈ *P*(*w*₁) *P*(*w*₂)...*P*(*w*_{*n*}) *P*(*w*₁) *w*₂...*w*_{*n*} 出现频率。基于统计文法模型的优点 在于减轻对词典的依赖，但依赖过多地消除，决策于性能与存储的平衡（如果深度结合机器学习，效率提升并依赖词典；如果太少，需要更多数据，解决于性的问题，解空间巨大）；缺点在于依赖于有监督词频的统计，对于新生词汇或专业词汇不好，冷门词汇的精确词汇往往难以准确划分。易受数据先验偏差（Bias）的影响。**基于序列标注** 对于基于统计模型的进一步简化，标注为开始 B、中间 M、结束 E 与单字词 **S.马萨可夫可转移隐（HMM）** 根据观测序列列到真正隐藏状态序列的推断。核心要素：两个集合：观测值集合、隐藏状态集合；三个矩阵：观测状态概率矩阵（从隐藏状态到观测值的概率）、隐藏状态转移概率矩阵（各种隐藏状态之间的转移概率）、初始状态概率矩阵（第一个字属于和隐含状态的概率）。初始状态概率 *π*，隐含状态转移概率矩阵 *A*，观测概率矩阵 *B*，观测序列为 *S*，隐藏状态值可能为 *S*, $\delta_1(s_t) = \pi_{0,\delta_1} \delta_1(s_t) = \max_{\delta_1} [\max_{s_1} (\alpha_{1,\delta_1}) \delta_1(s_1)]$, $\delta_2(s_2) = \max_{\delta_2} [\max_{s_1} (\alpha_{2,\delta_2}) \delta_2(s_2)]$, $\delta_3(s_3) = \psi_{s_1}(s_3)$ **停用词 stopwords** 指在文本中频繁出现但对实际语义影响不大的词语，如英文中的 The, of, in 中文字“的”、“是”等。数字、副词等与语义关系不大而不适用于词语级搜索。为什么要去除停用词？因为停用词重复率高，会造成对句子的理解偏差很长，影响查询精度，以及对最后结果的排序误差，反而产生干扰。**停用词识别与识别** 停用词的设置与语料库的性质有关（特定学科领域或具有其专用的停用词），如 URL 中的 www.Wikipedia.com 中的 wiki。停用词识别方法：文频数、词频统计、按场景等，更为复杂的算法将结合统计与语法或内容分析。**去除停用词可能导致的隐患** 停用词在特定场景有意义，如“非”、“不”表示否定，“较”、“稍微”表示程度等。停用词的结合有意义，如“的士”、“To be or not to be”。主要依赖分词工具的效果。**未来停用词的使用趋势** 现代搜索引擎的逐渐减少使用停用词，更关注自然语言的特性与特定处理应用场景，如压缩技术（降低停用词的存储冗余）、加入词项权重（将高频词的影响降至最低，吸引去除技术）、**归一化/词根化**处理原语词的特殊形式的问题。往往针对英语等语言，汉语并不需要这一过程。词根化在以前有减少词项而减少冗余、**词干提取 stemming** 去除单词的缀，获得词根的过程。常见的有词干提取有“复数形式”、“过去分词”、“进行时”等。词干还原 **Lemmaisation** 基于词根，将单词的复杂形态变成最基础形态，并不是简单地将其后缀去掉，而是根据词根将单词进行转换，如 am/is/are → be。**词干提取与词干还原的异同** 相同点：目标一致（目标均为将词的不同形态简化

或合并为基础形式）、结果类似（非互斥，结果部分交叉）、方法类似（主流方法均利用语言中的规则或词典实现）。不同点：原理上，词干提取“词根化”，词形还原“复数化”；复杂度上，词干还原需考虑词根化程度、词识别率等，更复杂；实现上，词干提取主要利用规则变化，词形还原更依赖于词典；结果上，词干提取不一定得到完全词根，而词形还原是完整单词。

第四节网络索引 布尔检索 查询简单，易理解，方便用户使用按式控制查询结果；功能弱，不支持部分倒排，不考虑权重和排序。**关联矩阵** 每列对应文档包含某些词，对所有的布尔项行为行向量编码；稀疏矩阵。**倒排索引** 信息检索中最普遍流行的、基于词项的基础文本索引，主要包括倒排表 (dictionary，词典集合) 和倒排表 (posting list，文档 ID 列表) 两部分构成；正排索引引以方便为主，易查找但搜索时间长；倒排索引以项项为项，维护成本低但搜索快。**建立倒排表的流程** 1.检查数据源文档，获得<词项，文档 ID>对，写入临时索引表；2.临时索引表中词项按序 3.遍历临时索引表，合并相同词项的文档 ID **基于倒排表查询** 本上是倒排记录表的“合并”过程，同时扫描每个索引表，所需时间与倒排记录的数量呈线性关系。如果两个倒排表的长度分别为 *x* 和 *y*，则合并需 *O*(*x* + *y*) 次操作。**倒排索引的优化性实现** 使用 AND 连接的查询条件是倒排表的合操作，先处理文档频率高的，再处理低的。一般的优化问题在任意组合的布尔查询，加入 OR：首先将所有的有项的文频数，其保持并设计每个 OR 操作后的结果大小，即考虑 *x* + *y* 的加入顺序，最后按结果从小到大的顺序执行 AND 加入链表后，需要先看链表指针是否小于另一边。链表指针越小，步长越短，可以更多跳但也不需要更多的时间；链表跳表越大，步长越长，跳得少，需要时间多。简单的启发式：如果链表增长为 *L*，则间隔 \sqrt{L} 均方均查表指针，没考虑步长增长因素，未必优化。**倒排表** 索引的动态变化会影响链表指针的设置，如果索引相对固定，建立优化的链表指针相对容易；反之，经常需要索引很难建立合适的链表指针。**动态索引** 主从索引，维护主索引的索引，新文档有储小的辅助索引，检索时会合并；删除索引用无指向位置，定期合并辅助索引。主索引与频繁合并导致大量合并，合并效率低，对于非常大的索引记录表进行分词，对较短的索引记录表进行合并，也可以基于词项常用性分词。**短词查询的例证** 将短词查询中的两个词合并成词项，要求保证语意一致。**第一种解决方案：二元词素表** 将文档中每个连续词对作为一个词组，语意可以构建而用二元词的倒排表，并处理多个词组构成的短语查询。短词短语查询可以分成多个短词查询来处理或使用更多的词素索引来匹配。如果采用二元词素索引链接的方式，对于分布布查询返回的结果，我们并不能确定其中是否真正包含最原始的四个词组。**第二种解决方案：位置信息索引** 多元词素索引最大问题是词表迅速增长，记录词项的位置可以记录它在文档中出现的位位置，达成更广泛的查询。每个词项单独列出出现频率和文档及文档中出现位置。对短词查询，仍采用 AND，查找符合的文档，不同单词单独列两个词是出现在同一文档内，还要检查出现位置情况是否符要求。位置信息索引类型更适用于短词搜索。**倒排表的扩展性增强** 除了文档 ID 和位置，还可以加入词项频率、词项类型、文档来源、文档类型、文档属性等。倒排索引的**存储问题** 可以将词典与倒排表一起存储，便于同时读取，但文档规模大时将导致索引过大，影响性能；也可以分开存储，词典与倒排表分别存储为不同的文件，通过索引链关联（倒排表文件可采用分布存储的方式），性能可以大幅提升，因为词典或者至少一部分内容常驻内存，同时支持并序、分布式存储、顺序存储（字典树、二分查找）、哈希存储（冲突多次效率率）、B/B+树（维护代价低，实现相对简单）、Trie 树（空间浪费问题）。**索引压缩的问题** 前提是快速检索和简单算法；索引提升效率：压缩磁盘空间，提高效率（内存利用率或检索效率）；词典压缩可以直接放入内存中，提升效率；压缩倒排记录表可以减少数据的磁盘空间，可以减少输入数据，减少大量磁盘读写操作所需的时间。**两种索引压缩策略的基本元素分析**：定长存储，易造成空间浪费，但仍存储超增长无法存储。**压缩词项列表：将词典表作为一个字符串** 将所有词项存成字符串，倒排表中记录字符串的存储空间大小；每个词项平均总估计占约 19=8+4+4+1+3 个字节，而并不是原先的 28=20+4+4+4 个字节，每个词项平均长度 8 个字节，词项平均长度为 8 个字节，倒排表平均为 4 个字节不重复，词项指向约 3 字节。**进一步压缩：块状存储 (Blocking)** 一个字符串在词项指上需要占较多存储空间，为节省 *k* 个词项存储一个指针可以减少指针的总数额外使用 1 个字节表示词项长度，节省 3/4 的指针数，但会每个词项加入一个字节的空间的查询。发现问题：未编词表的搜索和节省一半/四的代价，而采用按序存储后，二分查找只能在外部进行，块内只能使用线性时间复杂度。随着 *k* 上升，线性查找部分增多，效率更低。**另一种改进：前端编码 (Front Coding)** 按照词典可排列的连续词项之间，往往具有公共的前缀，使用特殊字节表示前缀使用。**倒排表存储的问题所在** 倒排表所占的空间远大于词典本身，最迫切的需求在于如何紧密地存储每一个倒排表，尤其是其 *z* 值。如果使用 4 字节数据来表式 *z* 值 ID，每个文档 ID 需要 32bit Ziplf 数据；排名第 *i* 个的文档的 *z* 文档频率与 1/ *i* 成正比。**再进一步：可变长度编码** 对于一个词项值 *q*，想用最小的所需字节来表示它，需用类似所需字节来求每个词项值，对小数字使用短码实现这一点。先存储 *C*，并分配 1bit 作为连续位，如果 *C* < 128，则使用第一连续位为 1 的字节，如果 *C* > 128，则使用 128 的连续位，然后取相同算法与 7 位位有二进制编码的格式。如果 *C* > 128，则先取低位的 7 位编码，然后取相同算法对高位进行编码，最后下一个字节 (Sbit) 的连续位为 1，其他字节连续位为 0。相比 4 字节编码，可变字节码在数字上占的码码可以节省更多空间。例如 5 的二进制为 101，加上连续位为 10000101。214577 的二进制为 1101/00011000/0110001，因此拆分为 3 字节，其编码码为 00001101/00001100/0110001。

第五节查询与评估 如何给出 Top *N* 的结果，衡量用户对查询结果的满意度。**相关性反馈 (Relevant feedback)** 用户直接对查询结果进行评价，引导用户表达真实查询意图（查询建议与查询扩展）；间接性反馈：用户在查询后标记相关/不相关，然后迭代更新查询，以获取更好的结果。问题：用户也无法表达想要想的结果，但多少能够判断所看的内容是否；精确查询难以操作无法一键撤回，但可以通过迭代不断趋于精确。**相关性反馈的基本流程** 1.用户输入查询条件 (Query)，2.对于返回的文档，用标出相关与不相关的部分，3.系统根据用户反馈，获得用户所需信息更为准确的反馈。基于相关性反馈，更新新查询条件；基于新查询条件，获取新的结果文档并再次提交用户进行评价。上述过程将根据情况进行一次或多次的迭代，从而不断接近最优查询条件。

相关性反馈的最终目的 通过相关性反馈，通过迭代获得对于表达用户查询意图的最优查询条件，可以减少词项漏掉不相关项，或增加新的词项，而这一过程即称为用户隐藏。**相关性反馈存在的困难** 可以影响用户参与体验（用户不愿意提供显式式反馈，不想避免因反馈显著延迟搜索时间），相关性反馈在查询的查询往往往低相关，降低系统效率，增加计算量（也可以以改变重词项权重而不增加新词项，但效果有限），被相关性反馈链接的词项，未必使用户需要的内容。**显式反馈** 最准确的显式反馈是用点击记录，但存在复杂的问题，只有正样本，而用不了点击，不代表完全不相关！拓展的显式反馈收集与质量评价。更为复杂的显式反馈难以用于评价，可以收集更为完整的相关信息，同时，对于用户问题有更有力量的判断。**隐式反馈** 更适用于用户对当前检索结果采取的作为，而用户检索结果的相关性质量评价，判定为相关的显式反馈难以用于评价，可以收集更为完整的相关信息，同时，对于用户问题有更有力量的判断。但省却了用户对用户的参与行为！包括相关性和相关性用于用户问题预测。鼠标键操作可能揭示用户身份特征，不同用户在击键频率、时延、习惯、错误率等方面存在一定差异，用“眼跟踪操作”，可以揭示用户对文档及关联（视觉感知链接）借助鼠标键操作，还可以了解其他相关的应用，例如，判断文本之间的相关性（甚至检测到答案）。**隐式反馈的优点** 优点在于不需要用户显式参与，减轻用户负担，提升用户体验，用户行为某种程度可以以反应其意图，因此具有可行性；缺点在于对于行为有着较高的要求，准确率低一些，某些情况下需要增加额外设备。**伪相关性反馈** 无需用户参与与反馈过程，而直接根据检索结果自身情况，较为简单，对于用户查询返回的结果，依据前 *k* 篇文档是相关的，在其基础上 *n*，进行相关性反馈。但也有存在伪善的隐患，结果未经用户验证，难以保证准确性，某些查询可能效果很差，甚至出现查询结果（系统 Top *k* 文档符合）含答案**查询扩展：一种用户相关性反馈的特殊方式** 即输入一个词项之后给出一个系列关联选项，用于用来选择，用于用户针对词项的合意程度给出反馈。这些反馈数据用来构建更为完整的查询条件，从而发现和确认的查询扩展能够更好表达其查询意图。**查询扩展的实现** 对于某个查询词项，使用同一词项群中的同义词或相关词进行扩展，相较于原词的查询词项，可以扩展的词项设定权重大小权重。查询扩展有助于提升查询的召回率（找得更多），但可能会影响准确率，尤其在有歧义存在歧义的情况下。

局部和全局的同一词项需要很大的代价。**查询扩展的类型** 利用人工编辑的同义词词库进行扩展，同义词词库的自动生成（基于统计计算之间的共现关系 Co-occurrence，自动构建词库）基本是基于过去进行优化（通过查询日志，挖掘查询的等价类）**相关性反馈的潜在问题** 相关性关联的质量是一个问题，有歧义的同义词可能导致统计上相关，而意思上不相干的词项，同时，由于扩展的词项与原查询词项高度相关，扩展后的查询也未必能够获得更多的相关文档。**结果评价的常见内容** 性能

(effective 返回多少相关文档，是否准确，排序是否靠谱)、效率（响应速度，时空消耗）、其他指标 (efficient 多次查询时，准确性，时效性，效率) **基本评价指标的矩阵化表示** *P*/*N*: Positive or Negative，表示算法判断为正确/正样本 *T*/*F*: True or False，表示算法判断为正确/负样本。TP: True Positive，样本为正例，且被判定为正，真正。FN: False Negative，样本为正例，但错误地被判定为假负。FP: False Positive，样本为负例，但错误地被判定为正，假正。TN: True Negative，样本为负例，且被判定为负，真负。**面向单查询的基本评价性能指标（也称查准率）** (Precision, 检索出的文档中相关文档所占的比例，也称查准率) *TP*/(*TP* + *FP*)；**召回率 (也称查全率/真率/命中率)** (Recall, 所有相关文档中，被检索出来的部分的比例，也称查全率) *TP*/(*TP* + *FN*)。为什么某些网站被搜索？**Accuracy**, (*TP* + *TN*)/(*TP* + *TN* + *FP* + *FN*) 在模式分类中经常使用，它在信息检索中的相关任务并不常见，因为不返回结果就可以保证 Accuracy 接近 100%，不返回结果即认为所有文档都不相关，因此 Accuracy = *TN* / (*TN* + *FN*)，实际中相关文档只占很小一部分，即 FN 相对于 TN 很小，所以正率接近 100%。**准确率与召回率的平衡** 平衡为 *F* (宁愿忘记一些文档数据，也不能错将正常文档)，牺牲对正确文档的召回率，保证较高准确率；智慧医疗（宁愿多断一些疑似患者，不能漏掉一个真实病人），牺牲正确诊断人的准确率，保证较高召回率。**召回率的近似计算** 大规模文档，不可能列出所有相关文档，无法准确计算召回率。用缓冲词池 (Pooling) 方法，针对某一检索问题，各算法分别给出检索结果集中的 Top *N* 个文档，将这些结果收集起来并进入人工标注，从而得到一个相关文档的池。将池放在大多数相关文档存在这个文档 (Doc Pooling) 中，可行性存在，虽然实际上当然无法得到全部相关文档，因此并不能得到召回率的准确值，但它可以比较各个算法的相对优劣。 *N* 通常取 50-200。**P-R 的平衡到 F1 值** *F* 值衡量准确率与召回率的加权平均数 $F_1 = \frac{2 * p * r}{p + r} = \frac{2 * p * r}{p^2 + r^2}$

通常取 *α* = 0.5, *β* = 1（即两者同等重要），得基本 F1 值，*F* = *2PR*/(*P* + *R*)。本质上，目标是 Precision 和 Recall 都比较高，而平均对一个偏值的情况是有惩罚的，必须两者都较高才行。而算数平均几乎任何人在处理机器学习值下效果并不够合理（如果采用平方平均 *F* 值，那么一个返回空文档的搜索索引的情况 F 值就不低于 50%）准确率、召回率与 *F1* 值是信息检索任务中最为重要和常用的三个指标。**P-R 曲线与 ROC 曲线** 在分类问题中，准确率与召回率的平衡是选择定一个阈值 (相关性大于阈值以相关关系) 实现的，通过调控相关性阈值，可以控制检索所得的文档数量。较低的阈值可以使得到更多文档，但也混入大量不相关的文档；较高的阈值可以保证文档的相关性，但会遗漏许多相关的文档。绘制不同阈值的指标化曲线，可以帮助我们做出选择。**P-R 曲线** 以准确率 (*y*) 和召回率 (*x*) 分别作为横纵坐标，通过选定不同的阈值得到不同的 P-R 点并连接成线，可以直观地看出准确率与召回率之间的平衡关系。**ROC (Receiver Operating Characteristic) 曲线** 以真正率 (召回率/命中率) *x* *TP*/(*TP* + *FN*) 和假正率 (误报率) *x* *FP*/(*FP* + *TN*) 作为两条轴轴通过选定不同的阈值得到不同的真正率-假正率点并连接成线。对角线表示区分能力为 0，即随机猜测；在对角线上方越远，效果越好；低于对角线的结果属无效 (无区分度) 如何基于 P-R 曲线或 ROC 曲线选择算法好坏 如果 *A* 曲线 B 完全包含，显然 *A* 对应的算法效果更好；如果两条曲线交叉重合，则计算 AUC (Area under curve, 曲线下面积)，谁更高效果更好。另外，当使用 P-R 曲线时，平衡点 (Precision = Recall) 越低越好。**P-R 曲线与 ROC 曲线的选择** 如果 ROC 曲线难以观察评价，更为全面，而 P-R 曲线则只靠正确例，用户往往更关心正样本，所以面向特定应用效果（如检索）P-R 曲线是更好选择。**单例数据正样本比负例失时**，P-R 曲线更合适（当负样本比重过高时，负例的数目众多导致假正的增长过快，导致 ROC 曲线更直观一个过分散离的效果结论，从而难以体现出假正假负性）P-R 曲线受分布影响较大，多份数据正负比例不同时 ROC 曲线更合适。**ROC 的曲线与两个指标各自对正负样本，而 Precision 只针对正样本，受影响较大。P-R 的曲线与意义** 由于大多数用户只关注第一项或前几项，只针对 P@10, P@20 等对于大规模搜索需求来说是很合适的评估指标。P@*n* 上，即 Precision@*N*，指前 *N* 个搜索结果与大规模搜索（如果相关文档小于 *N*, P@*N* 的理论上限就是 1）由于返回结果有限,Recall@*N* 值，甚至其理论上限往往都小于 1（理论上限为 *N*/相关文档数），即使通过 Pooling 加以控制仍然较小！**特殊例-R-Precision** 检索结果中，在所有相关文档数位置上的正确率。例如，如果相关文档数为 3，那么考察 P@3 作为 P-Recision。而**面向 P-R 的 *q* 值** 的 *P-R* 曲线 在有存储结构的情况下，不再采用不同阈值值为 P-R 值的依据，而是通过依次计算前 *N* 个结果对应的 P-R 值绘制曲线。新的不相关文档被检索时,Recall 变,Precision 下降，因而**假倒查影响**。**平均准确率** *AP* 用于对不同召回率点上的位置进行平均。• 未做值 *AP*：查询 *Q* 共有 6 个相关文档，排序返回了 5 篇相关文档，位置分别是第 1，第 2，第 3，第 10，第 20 位，则 *AP* = (1/1 + 2/2 + 3/3 + 4/10 + 5/20 + 0)/6 (注意补 0) • 值偏 *AP*：事先选定阈值点并平均处理。例如，如果我们计 11 点平均，那么召回率分别为 0（第一条若查不到 1，否则为 0），10%，20%，...，100% 的一个点上的正确率求平均（注意补 0，如查不存在召回率为 100% 的点即没有返回所有相关文档，则需要补 0）。• 简化 *AP*：只对返回的相关文档进行统计，*AP* = (1/1 + 2/2 + 3/3 + 4/10 + 5/20)/5，适合那些快速返回结果的系统，不考虑召回率本身的情况（不补 0）**累计增益 (CG)** 用于衡量位置 *i* 到 *p* 的检索结果的相关度之和，*CG*_{*p*} = $\sum_{i=1}^p rel_i$ *rel* 用于描述文档相关性，可以根据需求设置多个值/级别，级数越高 *CG* 表明文档的整体相关性较高，但因为为文档相关性，并不能体现文档部分文档的质量。**改进的相关度增益 (DCG)** 对不同位置相关不同折扣后折损累加计算 (Discounted Cumulative Gain,DCG) 若搜索系统把相关文档的文档排在后面，则应该给予惩罚。一般用 log 函数来表示这种惩罚，如 *log*(*i* + 1) 其中 *i* 为文档位置，往往有以上两种计算方法（后者使用指数函数，突出相关性）： $DCG_i = rel_i + \sum_{j=2}^i \frac{rel_j}{log_2(i)}$

$DCG_p = \sum_{i=1}^p \frac{2^{rel_i}}{(1+rel_i)}$ **归一化的相关度增益** 和 DCG 局限于随长度单调非减，结果与是否有关因不同与不同算法对比。对 DCG 进行规范化到归一化折损累加函数 (Normalized DCG,NDCG) 将 DCG 除以完美结果得到理想值 *iDCG* (ideal DCG) *NDCG* = *DCG* / *iDCG* 其中 *iDCG* 是根据文档根据相关性从大小到排序到理想化的最优排序计算 DCG 值所得到的。例：假设有 10 个文档，相关性为 0-3 之间的 1,0 个文档的得分依次为 2, 3, 3, 0, 1, 0, 2, 2, 3, 0，由此计算基本的 DCG 为 3, 6.89, 6.89, 6.89, 6.89, 7.28, 7.99, 8.66, 9.61, 9.61。理想的输出排序结果应为 3, 3, 3, 2, 2, 1, 0, 1, 0，由此计算 *iDCG* 依次为 3, 6, 7.89, 8.89, 9.75, 10.52, 10.88, 10.88, 10.88, 10.88。于是可得 *NDCG* 结果如下 1, 0.83, 0.87, 0.76, 0.71, 0.69, 0.73, 0.8, 0.88, 0.88。任何给定相关度位置的 *NDCG* 值都被规范化为 0-1 之间的值。从 **AP 到 MAP** 从单查扩展至多查评价，可以更全面的体现排序算法的综合性。MAP (Mean AP)，对所有查询的 *AP* 求算术平均均可，可以反映全部查询的综合效果，但在查询难度不平衡的条件下有误导。在查询难度不平衡下有误导，如果一个系统在较难系统上有一提升而在简单系统上表现可能更差，引入基于几何平均的 *GMAP* = $\sqrt[p]{\prod_{i=1}^p AP_i}$ ，在查询难度不均或存在有选择性的主题更适合。**注重查询相关性文档的 MRR** 用于带常关心一个相关文档的情况，越靠前越好，该位置的倒数称作 Reciprocal RRR (RRR)，数字越大越好。对于多查询问题的倒排顺序平均,MRR (Mean RR)。一篇文档可能被用户点击的概率为 *PP* = *rel*, (*rel*₁ - 1 *rel*₁)，定义倒排的倒排顺序 *ERR* = $\sum_{i=1}^p \frac{1}{AP_i}$ ，*PP* 表示用户所需数据时停止的位置的倒数的倒数。**为什么要考虑多样性** 用户的单次搜索可能提出多个面的需求，也可能存在歧义，需要更多方面的信息加以确认。同时也要避免信息冗余问题。**多样性的形式化定义** 给定一个查询 *Q*，返回一个多样性的检索结果集合 *R*(*Q*),*R*(*q*) 作为一个整体，应满足 *R*(*Q*) 中所有的检索结果都与查询 *Q* 本身有较大的相关性。总体上要有一定的冗余度，以覆盖 *Q* 的不同方面。目标降低用户无法获知所需信息的风险，尽可能确保排序靠前的结果中至少有一个结果满足用户的需求。**多样性的衡量方式** 衡量不同文档之间的主题差异性，包括隐式多样性（只计算文档之间的差异性，不会考虑用户点击主题与内容）与显式多样性（更加深入地考虑文档所对应的用户意图，会人工文档中抽取主题，并显式地体现主题的多样化）显式 FM-LDA 文档分类器 *F* 的文档相关性 $F = \prod_{j=1}^n (1 - p_{fj}) (1 - p_{fj})^{w_{fj}}$, *p*(*f*_{*j*}) 是

表示文档包含主题 *f_j*，为文档数集中，连续为主题至少一篇文档数。隐式 MMR 最大边界相关性 *MMR* = *arg*max{*P*(*d₁*,*d₂*) - (1 - *λ*) max *P*(*d₁*,*d_j*)}，前半表示文档与查询相似，后半表示文档间多样性，*λ* 调节比例。

第六节第七节网络索引 量化方法 (1):Jaccard 系数 布尔模型要求所有的查询条件都必须满足，因此是二值的。如果 *A* 与 *B* 为两个集合，*JACCARD*(*A*,*B*) = $\frac{|A \cap B|}{|A \cup B|}$ 给出对集合重叠度的描述，但未考虑长度，未考虑频率。**量化方法 (2):词频指数 Term Frequency** 查询词在原文档中出现的次数，该文档相关。词频指数 *t* *f*(*t*,*d*)，词频项 *t* 在文档 *d* 中出现的次数。利用原档的 TF 值，可以粗略计算文档的相关性，但相关性不线性且不线性相关随 TF 基础之上的改善：

引入归一化词频 $wf_{i,d} = \left\{ \begin{array}{l} 1 + \log_{10} tf_{i,d} \\ 0, \end{array} \right.$ 缓和数量级的差异性所造成的影响，可以将文与词项的匹配得分定义为有时可能出现文与文档中的项其对应文档数，即 $\sum_{i=1}^n (1 + wf_{i,d})$ 若没有公共词项，显然得分 0。**量化方法 (3):文档频率 DF** 平列的词信息量更为丰富，而检索词的信息量相对较少，相应的，如果查询词也包含某个平列词，则包含这个词的文档可能相关，引入新的度量指标：文档频率 (Document Frequency) *df_i*，指出包含词 *i* 的文档数，一般采用平方反文档频率 (Inverse DF) 来衡量 *idf_i* = $\log_{10} \frac{N}{df_i}$ 这里同样可以

对数计算方式称 DF 为衡量量级的量化。**量化方法 (4):集成文档的 tf-idf 值** 与 DF 从两个方向对文档为衡量文档相关提供了量化依据 $Wt_{i,d} = (1 + \log_{10} tf_{i,d}) \cdot \log_{10} \frac{N}{df_i}$ (*t*_{*f,d*} = 0 时，*W*_{*t,d*} = 0)，*t*_{*f,d*} 表示词项 *t* 在文档 *d* 中出现的次数，*df_i* 表示含词项 *t* 的文档数，*N* 表示文档总数。随词项频率增大而增大，随词项罕罕程度增大而增大。适合衡量文档相关性的词是在少数文档内出现多次的词。**量化方法 (5):向量空间模型 vector space model** 每个文档和查询操作一个词项权重构成的向量，查询时通过比较向量之间相似性来进行匹配。D 是文档表达，每个文档可视为一个向量，其中每一维对应词项的 tf-idf 值。Q 是查询表达，可视为一个向量，其中每一维对应词项的 tf-idf 值，P 是非完全匹配度，R 是使用每个向量上的相似性来度量文档与查询的相似性。1 首先，将文档与查询表示成向量的 tf-idf 权重向量（也可用其他方法），2 其次，计算两个向量之间的某种相似性（如余弦相似性），3 最后，按相似大小进行排序，将 Top-*k* 的文档返回给用户。优点是简洁直观，可以支持多种不同度量或权重方式，但效果不错；缺点是缺乏足够层面的理解和控制，同时依赖 tf-idf 值也能造成干扰（用户无法控制文档项之间的关联，词项向量的数量性假设在假设上不成立）**查询文档的匹配度** 欧氏距离度量标准，但不是一个好的选择，对于向量长度（文档长度）非常敏感，不能体现文档的相似性。余弦相似性更合适，按照文档向量与查询向量的夹角大小来计算，向量越一致，夹角越小。**基于迭代的查询意图更新**、**罗基奥算法** 用户的查询意图可能无法一蹴而就，而需要通过不断的反馈实现逐步更新。本质上，这一过程就是使查询文档的表达逐步趋近用户目标文档的过程。可以从目标文档的值为出发点，*C* 为文档集合， $\vec{p}(C) = \frac{1}{|C|} \sum_{c \in C} \vec{p}(c)$ 罗基奥算法将以相关文档信息输入到向量空间模型，试图找到一个完美的最查询向量，以满足 $\vec{q}_{opt} = \arg \max [\cos (\vec{q}(\vec{C}_r)) - \cos (\vec{q}(\vec{C}_{nrr}))]$

也就是使查询尽可能离与之相关的文档更近，离为之不相关的文档更远。理想情况，在可知完整的相关/不相关文档的情况下，可以得到 $\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{c \in C_r} \vec{d}_c - \frac{1}{|C_{nrr}|} \sum_{c \in C_{nrr}} \vec{d}_j$ *d_c* 获得完美查询。问题在于事实上并可能使用完整的相/不相关文档集合。实际使用的近似方法如下：

$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nrr}|} \sum_{\vec{d}_j \in D_{nrr}} \vec{d}_j$ 其中，*D_r*

现实场景中，缺点包括主题无关联性，对于恶意链接、广告链接等缺乏区分，旧闻反复出现等，因为新闻网站往往存在冗余，一般不能单独用于排序，需要和上下文信息方法结合，效果更低。
PageRank的计算方法
PageRank的核心公式**P
R
(

p

i

)
=

1
−
d

N

+
d

∑

p

j

∈

M

(

p

i

)

P
R
(

p

j

)

P
(

p

i

)

{\displaystyle PR(p_{i})={\frac {1-d}{N}}+d\sum _{p_{j}\in M(p_{i})}{PR(p_{j}) \over P(p_{i})}}**

PageRank 转移概率矩阵为:

A
=
d
M
+

(
1
−
d

N

)

e

e

T

,

e

e

T

是全1矩阵,
P
R
(

p

i

)为网页

p

i

的PageRank值,
P
R
(

p

j

)为指向网页

p

i

的0-1之间

p

j

的PageRank值,
L
(

p

i

)为网页

p

i

的页链集合,
d为跳转矩阵:邻接矩阵转置后,将各个数除以所在列的非零元素数(即 outgoing),即为跳转矩阵.
(注意这里)
**P
(

p

i

)
=

A

P

i

A

P

i

,**形成马尔可夫过程。
PageRank的实现过程
先给每个网页赋一个初值 1/N, 利用之的方式进行迭代有限次计算，得到近似结果。
迭代收敛后，若

P
(

p

i

)
P
(

p

j

)

 对应网页*i*和网页*j*的差值高则继续迭代，直至收敛。若当前

P
i

与

P

j

 差值甚高，则继续迭代。
PageRank中的两类特殊情况
陷阶节点*i*（没有一条指向它的边，没有它 outgoing），终止节点*j*（没有指向它的 outgoing）孤立节点*k*（没有任何入边）只有初始链接，不再更新，也不影响其他节点。
特殊情况的解决:Restart 机制
公式

(
1
−
d

)
 的部分，相当于以一定概率重新选择起点，所有节点以一定概率被选中作为新起点，跳出了陷阶和黑洞的干扰（所有节点全连通）
d一般选择为 0.85 左右（Google 的**选择PageRank中的收敛性讨论** 如前述,PageRank的计算过程, 实际上是一个马尔科夫过程。马尔科夫过程有三个收敛条件：转移矩阵 A 为马尔科夫矩阵 A 矩阵所有元素都大于等于 0，并且每一列之和都为 1；转移矩阵 A 为不可约的，当它是强连通时，A 为不可约，而 Restart 保障了这一条件；转移矩阵 A 为非周期的，这三大条件,PageRank 算数都满足，因此保障了其收敛性。
personalized pagerank 将用户作为起点，跳转概率不一定均等，根据用户偏好改变概率。基本公式如

r
=
(
1
−
d
)
M
r
+

v

 的形式，

v

 体现了用户偏好，可理解为不同用户在用户偏好下的主题偏好。
topic-sensitive pagerank 以主题为中心，体现用户不同的评价才有价值，主题相关的网站指向与推荐行为被削弱，对于主题，重新开始的讨论仅限该主题下网络，最终每个网站得到一个向量，不属于每个主题问题可以获得一确定值。
HITS 算法中的两个基本概念
权威网页 (Authority，某个领域相关高质量网页) 与枢纽网页 (Hub, 中心, 指向高质量权威网页的网页) 假设 1 好的 Authority 会被很多好的 Hub 指向,

V

p
i

,
a
(
p
)
=

∑

h
(
p
i

)
h
(
p
)

 假设 2 好的 Hub 会指向很多好的 Authority,

V

p
i

,
h
(
p
)
=

∑

a
(
p
i

)
a
(
p
)

 在 HITS 算法中，网页间重复计算网页之间。相比 PageRank, HITS 算法可以以不同网页功能，优点是易于更好描述网络特点、主题相关，可以单独用于网页排序；缺点在于需要在线计算，时间代价大，对链接结构变化敏感，可能使链接旁路影响，初始计算选择直接查相邻链接。**HITS 算法的计算过程**
初始化所有网页的

a
(
p
)
 和

h
(
p
)
 为 1。
1. 进行计算直到收敛**注意每一步中的一次归一化过程**。
即按矩阵 **M** (**M** 是邻接矩阵)，Authority 向量

a
,
H
ub 向量

h
，有如下迭代式:

a

k
+
1

=

M

T

h

k

,

h

k
+
1

=
M

a

k

，也可采用

a

k
+
1

=
(

M

T

M

)

a

k

,

h

k
+
1

=
(
M

M

T

)

h

k

。HITS 算法的优缺点：
相比于 PageRank,HITS 算法是一种能够区分网页间功能的排序算法。
优点：更好地描述互联网网络特点；主题相关，因此可以单独用于网页排序；缺点：需要在线计算，时间代价较大；对链接结构变化敏感，且依然可能受到“链接旁路”的影响。

第八章个性化信息检索
信息已经从精确到粗略，因而需要进行有效过滤，个性化检索。

x
 为带长标量

X
 与太矩阵。

X
 (个性检索集),

S
 (项目检索集),

R
 (评分集合)。
定义效用函数

u
:
X
×
S
→
R
。本质上，(个性检索集)是全矩阵

R
 的过程，但将该矩阵严重稀疏化(用用户行为数据，活用用户少，新用户/项目冷启动)。
常用 RMSE 均方根误差

√

(

r

i
,
t

−

r

i
,
t

)

2

/
n
 来评估。也可以视作分类问题，用

P
/
R
/
F
@
N
 表示。
基于内容的推荐: 核心思想
假设用户偏好一般相近对应，给出推荐与用户喜欢的物品类似主题/倾向性的项目，但可能信息检索和推荐**基于内容的推荐: 项目相似**
对于每个查询物品，需要给出相似度的画像 (Profile)，以向量的形式存，抽取关键词或相似度 tf-idf、主题模型等方式。
基于内容的推荐: 用户画像与内容画像
用画像描述可由经过评分后的项目画像描述，一般用加平均方式得到(基于评分进行归化)。
基于用户项目画像，采用相似性度量进行评分。一通用用户之间向量之相似度，也可根据实际选择其他方法。
基于内容的推荐: 优点
每个人的推荐过程相互独立，不需要其他用户的数据，可以不具有独特特好的用户进行有效推荐，不受大众倾向性和热度的影响。
可以推荐新用户和冷门项目
推荐过程有较好的可解释性，可列举推荐理由特征作为推荐的依据。
基于内容的推荐: 缺点
很难找到合适的特征，尤其是对于子结构化信息 (图像、视频等)。
部分特征的提取可能存在误差 (归因错误)。
难以建立新用户的项目画像。
可能出过度度特化 (Over-specialization) 现象，永远只能给用户推荐局限于画像中的内容 (信息茧房)。
难以体现用户的多方面兴趣，难以通过他人评价对推荐结果进行评估。
衍生态
可考虑使用知识图谱进行画像建模，基于图谱进行推荐。
用户反馈可受位置影响。
协同过滤 (Collaborative Filter) 的基本思想
基于内容的推荐方法只基于用户记录向该项目进行推荐，但在实际应用中，其他用户的浏览行为对当前用户有借鉴参考价值，社交关系可带来相似物品。
协同过滤的思想在于基于评分矩阵的行为协助判断本行的数据。
协同过滤的基本类别
基于内存 (Memory-based) 算法：具有数据与数据模型运算进行推荐，包括基于用户与基于项目。
基于模型 (Model-based) 算法：具有有数据训练模型。
通过模型进行推荐，包括矩阵分解与深度学习等
基于用户推荐、User-based
具有相似偏好的用户，往往在评分行为上相似，找到这些相似用户，并基于这些用户的历史行为进行推荐，类似于寻找最近邻的思想，从用户过去行为着手。
定义相似性sim(a,b) =

r

i
,
t

=
p
e
r
o
d
u
c
t
(
P
(

r

a
,
p

)
,
P
(

r

b
,
p

)
)

{\displaystyle r_{i,t}=product(P(r_{a,p}),P(r_{b,p}))}

基于用户推荐的评分预测
在得到用户之间的相似性后，针对待预测的项目，可以根据历史上其他用户行为对于该项目的评分，结合用户之间的相似性作为加权， 预测评分结果。
相似性计算与评分预测中，都通过减去平均来抹去评分平均偏好的影响（不同用户打分范围不同）

p
r
e
d
(
a
,
p
)
=

r

a
+

∑

b
∈
n
e
i
g
h
b
o
r
s
(
a
)

s
i
m
(
a
,
b
)
⋅
(

r

b
,
p

−

r

b
,
p

)

∑

b
∈
n
e
i
g
h
b
o
r
s
(
a
)
s
i
m
(
a
,
b
)

{\displaystyle pred(a,p)={\frac {r_{a,p}+\sum _{b\in neighbors(a)}sim(a,b)\cdot (r_{b,p}-r_{b,p})}{\sum _{b\in neighbors(a)}sim(a,b)}}}

相似性计算与评分预测中，都通过减去平均来抹去评分平均偏好的影响。
有时只考虑 K-最近邻的情况，相应的，前页公式中的 neighbors(N) 即缩小为 K-最近邻的集合。
只考虑对指定 Item 有评分的邻居，跳过无评分的（即使 Similarity 较高），通常忽略小于 0 的相关性。
因此 K-最近邻集合可能实际上小于 K 个邻居。
基于项目推荐、Item-based 类似
与用户行为相似性类似，相似的项目，往往评分也比较接近。
与基于内容的不同，衡量相似度的标准，并不为项目本身的属性，而是不同项目的行为历史。
同一个人给两个项目打分相似度高，说明这两个项目相似；而这样的两个人，两个项目越相似。
计算公式：

r

i
,
x

=

∑

j
∈
I
(
i
,
x
)

s
i
j

r

i
,
j

∑

s
i
j

{\displaystyle r_{i,x}={\frac {\sum _{j\in I(i,x)}s_{ij}r_{i,j}}{\sum _{s_{ij}}}}}

同样可以基于平均分对数进行修正，如果评分不同，则直接设为 0，不同则取最大值。
**最终的

r

i
,
x

 不需要平均修正**
因为只针对该用户自己评分并加权归化。
两方法的区别
往往基于项目的推荐效果更好，因为项目之间的相似性大于不同用户的偏好程度行为丰富多样，某项目受欢迎的原因由相同原因，而用户可能在不同情况下体现出不同的偏好。
在没办法区分这些不同情况时，得到的用户偏好反而不准确。
基于内容的推荐的优势
优点在于于适用用于信息种类不同，效果好坏，不受多数偏好、非结构信息表征与特选选的限制。
缺点在于于不适用于(用户与项目都多样、精确性 (用户历史记录信息稀疏，很难得到评价过同一项目的物品)、热度偏差 (更倾向于推荐热门项目、对具有独特偏好的用户推荐效果差、小众偏好容易被热门项淹没)、**新用户问题** 新用户用户，可以先提供非个性化推荐直至收集足够个性化数据，借助历史信息/记录 (可能有问题隐私问题)、诱导式推荐 (选取代表性和多样特征

目，迭代收集反馈，快速获得近似画像)。
而面向数据项目，可以混合内容与方法，添加一些信息，引入知识图谱。
基于项目的推荐: 基本思路
数据项的精确性、计算最近邻的高复杂度，限制了基于内容推荐的有效性。
用户对项目的评价，本质上是用户偏好与项目属性之间的相似度，相似程度高，评价就越高。
基于项目的推荐: 基本思路
借鉴矩阵分解 (Matrix Factorization) 的思路，揭示潜在因子。
评分矩阵 R 近似视作用户属性矩阵 P 与用户偏好矩阵 Q 的乘积:

R
=

Q

T

P
,

R

为
M
×
N
矩
阵
,
Q
为
M
×
K
矩
阵
,
P
为
N
×
K
矩
阵
,K为潜在因子的数量。
矩阵维度，一方面与用户/项目的数量有关，另一方面潜在因子潜在因子的数量。
用户评分是根据潜在因子的乘积得到，基于子项方法得到的用户评分，应与历史评分记录尽可能接近，即 min_μ_μ_Σ (u_i,c) (r_{i,c} - u_i · v_{c,j})^2。
基于矩阵分解: 过拟合与泛化
上述公式本质是误差平方和 (SSE)，正是通过优化这一 SSE，以获得潜在因子的估计值。
但潜在因子维度 K，也会带来问题。
要训练的参数一共 (M + N) × K，随着 K 的增大而增长，K 足够大，在数据稀疏的情况下，没有那么多训练样本，或导致过拟合问题。
要引入更多正则化项，但很难：要么通过控制“收缩”参数的方法来提升泛化性，但这个本质是对训练样本过分信任，影响模型泛化的能力。
通过控制正则化项，使其不至于“过适”训练样本，引入正则项则会导致太多的参数值，把参数值还原点拉 min_μ_μ_Σ (r_{i,c} - q_i p_{c,j})^2 + λ_1 ∑ ||p_i||_2^2 + λ_2 ∑ ||q_i||_2^2。
正则化的效果：一般而言，用户的历史行为越少，记录越稀疏，模型泛化拟合的风险越大。
相应的，正则化所起到的“回归原点”的作用也越强**MF 的衍生模型**
(1) 非负矩阵分解
基础矩阵分解可能带来负元素，但有时元素要求非负，如用于添加非负参数。
MF 的衍生模型
(2) 概率矩阵分解
属于参数非负高斯分布，在数据稀疏且同时噪声时，引入某些参数可能更好的描述。
假设参数本身与参数均符合高斯分布，两个潜在因子矩阵相互独立，各用户/项目独立分布，两个参数矩阵均服从相同分布 0，预设方差的分布分布。
基于 R 与 U、V 联合分布高斯分布，在参数完成训练后，即可得到得到更完整的概率矩阵分解。
矩阵分解的拓展
给矩阵分解加入，加入更多信息或与假设，往往可以开发更加合理的动机/作用)。
一种非常规化的约束方法: 社交约束
假设好友之间在偏好行为上十分相似，如果好友与友关系作为 Side Info，可以用于模型修正参数。
基于偏好相似性的社交约束效果虽然很好，但未始会终止，社交关系并非非通过数据直接影响效果，而是通过影响参数 (“帮不上忙”的决策行为)。
情感的概念与含义
用户输入的条件特征之精确性，存在歧义/过于精确，借助上下文/环境可优化。
这类信息称作“语境信息” (Context Info)，从计算机科学的视角，语境一词可定义为“所有与人工交互相关，用于区分特定当前场景的信息”。
服务提供者借助语境信息，提供更精确的信息检索和增值服务。
搜索中的基础概念: 上下文
直面上，查询上下文可以帮助更好地理解用户查询意图。
用户的使用行为往往具有一定的连贯性，同一查询会话中的查询和点击的 URL 往往相关。
线下网络 (模型推荐) 将查询词/行为查询意图，避免敏感，将线下模型描述关联关系) 与线上服务 (部分会议，根据已有信息进行了解意图，进行精准推荐)
基于上下文感知方法
查询意图：一组有共同语义的查询词，有助于理解查询词精确性，更规范理解查询上下文。
核心在于从序列抽取查询意图的序列形式。
考虑当前长度上的所有序列数据，保留频率高于阈值的模式，并存储为后缀树 (Suffix Tree)，当查询上下文感知任务时，根据已有序列找到相应节点，从而获得查询意向描述。
进阶的上下文感知方法
引入隐马尔可夫模型 (HMM)，将用户意图

X
 视作变量，查询上下文

x
 与点记录

u
 视作观测值。
因为查询可能不与前一个查询相关，引入变化后的隐马尔可夫模型，打破 HMM 中每一个状态都与前一时间状态相关的约束。
细化的上下文感知不同类别
1 查询意图 (后续查询往往是先前查询的重新表述，目的不变或类似)
先前点击的内容往往不再被点击 (即隐含内容相关)
2 查询词流 (后续查询对先前查询中内容进行了 更为具体、深入的查询)
先前查询或点击为泛指的查询内容被略过
3 查询流化 (后续查询对先前查询中部分内容进行了更广泛的查询)
体现用户对该查询词广泛的意思，而不是局限某单一特定意思。
4 一般关联 (借助历史查询可以完全用于在前者搜索中的特定意图)
更接近于之前所搜索过的“上下文”
语境感知的基本概念
推荐系统中的情境有何意义？
情境是在重新进行同一活动时可能能发生变化的因素，根据应用场景、收集和整理相关数据，可以更加精确地推荐用户所需。
情境感知的应用
描述用户所处状态及意图，但也不是单一下文 (查看或查询数据)，而是一种复杂且相互关联的情境信息。
单一解释某情境要素，可能无法得到用户当前状态的正确描述。
如何获取用户所处的情境？
首先获取，直接利用用户或第三方数据，并让用户参与其中进行选择。
隐式获取，通过使用的行为数据获取用户当前所在位置等信息信息。
推断获取，可以用历史行为 (序列) 进行推断和识别。
情境描述
并非所有的情境信息都会影响用户当下的决策，用户可以忽略用户某些情境信息对于当下是重要的，也可通过特征选择 (主成分分析 PCA 和线性判别分析 LDA)，或者通过统计特征分析 (统计测试，如信息增益、信息熵、Freeman-Halton Test 检验等)。
如何情境信息整合到推荐系统中？
情境描述 (情境信息) 描述信息输入到推荐系统的数据流中，在推荐的各个阶段更好地理解和利用情境信息。
常用的数据采样方法
降噪、离群、缺失、重复等数据问题要求我们对数据进行预处理。
数据采样要求采取大规模样本，起始信息低，降低开支。
最基本的样本为简单随机采样法，对于所有用户，采用同样的数据源和方法进行采样。
无放回采样 (被采中的单元从整体中剔除，仅可选中一次) 和有放回采样 (被抽中的对象会入整体中剔除，可再次被选中) 无差别剔除，但放回有放回采样缺陷 (有误差)。
更复杂的方法包括分层采样 (Stratified Sampling)，对数据进行分层，从预先指定的层中进行采样。
启发式的采样规模确定方法
采样效果会受样本容量影响，因此要确定合适的采样规模。
通过分组采样的方式，可以近似确定适当的样本容量。
组内的数据高度相似，而不同组对象差异较大。
根据不同样本容量每组至少抽取一定个数数据来确保样本容量足够。
渐进式采样
渐进式采样根据预测结果调整样本规模，从小样本开始逐步增大采样规模，在模型准确率达到一定程度时停止采样。
优点是主要需要在一开始就确定好所需的容量，缺点是计算量大 (需要多次迭代)
维度的归并的归并性方法: 主成分分析 (PCA)
描述一个对象涉及许多特征，维度的归并则不具有归并性的特征，保留特征的行为，避免信息失真同时使模型更易理解更可视化。
过程:1. 对所有特征进行中心化，即将原始数据的每一维减去该维的均值。
2. 计算样本协方差矩阵

1
n

X

T

X
，X 为中心化后的矩阵。
3. 对协方差矩阵特征值分解，求出所有特征值及 4. 取最大的 d' 个特征值所对应的特征向量 w_1, w_2, ..., w_d' 对应的特征值的比重反映了主成分的最大信息，一般定义为 0.85)。
5. 计算投影矩阵 W = (w_1, w_2, ..., w_d')。
6. 计算投影式

X
′
=
X
W

,

X
′

即为降维 d' 维后的各样本点的新特征值。
数据离散化 discretization
将连续属性离散化为分类属性，将分类属性多数降低为少数的分类。
对于其数据数据化问题 (尤其无分类问题) 通过合并减少类别数目有 3 类。
非监督离散化
最根本的区别在于其数据化过程是否有监督信息，是否利用类别信息。
不使用类别信息的非监督离散化方法，往往根据数据本身的特性进行离散，常用的方法包括等宽、等频率、k-均值、等深寡数划分等。
监督离散化
更注重面向问题的目的，在于取得更好结果。
基本的方法是：最重要之一，定义区间阈 e_i = - ∑_{j=1}^m p_{ij} log_2 p_{ij}，总定义为加权和 e = ∑_{i=1}^m w_i e_i，然后从小区间阈值高 (标称级) 起，符合要求者，可以进行二分，选择最小的点进行第一次分割，然对于具有较大偏差的部分进行一次进一步分割，以此类推。
第十一节知识图谱概述
查询不再需要提供信息，还须提供信息扩展，一次查询，多重服务，而这也需要信息关联。
信息抽取的 X 类
从语义 (整体文本) 中抽取特定内容、事实等等信息结构化 (预先定义) 的数据 (细粒度、结构化)，是聚合与分解的过程，可以为后续的情报分析、自动挖掘、网络系统等一系列应用提供服务。
与信息检索密切相关但互不鲜明不同，功能检索是找出信息，抽取是获取信息)、处理技术 (检索利用统计关键词、抽取要提取 NLP)、使用领域 (检索与领域无关，抽取与领域相关)。
信息抽取的 X 类
抽取实体，确定关系 (命名) 实体指文本中的基本构成块，如人、机构等；属性类实体的特征任务 (事实/关系) 是实体之间存在的联系；事实类实体的行为为实体参与的活动。
5 属类的信息抽取任务: **命名实体 NE (实体抽取)**
命名实体抽取是信息抽取中最重要任务。
命名实体是文本中基本的信息元素，是正确理解文本的关键。
狭义上指现实世界中具体或抽象的实体，如人、组织、地点等；广义上还包括时间和时间、数量表达式等。
模板元类 TE (属性抽取)
模板元类 (实体属性)，用于更精准、完整描述命名实体、称 (名称、类别、种类等)。“大规模的 TE”，TE 是大规模产出的。
共指关系
CR 不同的命名实体表述相同含义即为其共指关系 (等价概念)，抽取共其共指表达的信息，包括已在命名实体和模板元类任务中标记出的，对某些命名实体所表述。
CR：生成语法和句法树都需要的知识。
模板元类 TR (关系抽取)
通过关系抽取将实体关联起来并为推理奠定基础，职务、婚姻、生产等。
场景模板 ST (事件抽取)
Who,

When, Where, What, Why, How.
知识图谱的符号表示
用户可以使用知识图谱直接解决查询问题，不必知道其网络背后信息总结。
知识图谱中节点和节点之间的连接，节点表示实体 (或实体化)，边表示关系 (或属性)。
节点上一般表示有一个指向图，方向表示主动关系、表达相关性等方向如图。
知识图谱的基本元素：
节点一般表示概念或实体
边一般表示关系 (实体间关联) 和属性 (描述实体的特征)。
节点和边知识图谱的基本单元三元组 (如:实体-关系-实体类)
知识图谱的优点
找到最想要的信息，不再需要用户自行阅读、阅读和检索，将信息直接呈现；提供最大的便捷，对搜索对象进行定位，使得用户获得更完整的信息关联；探索更深层次和广度、构建完整知识体系，使用户对信息想不到的新发现。
知识图谱相关应用
该大搜索 (建立事物关联网，更准确直接、快速)、网络系统 (理解用户信息基础上，基于推理能力提升交互体验)、推荐系统 (利用知识图谱提高推荐系统的推荐多样性和可解释性，提升推荐性能)

第十二节 知识图谱抽取与标注 命名实体识别 NER
识别出文本中的人名、地名等专有名词和有意义的日期、日期数量词等命名实体。
信息抽取标注的核心任务，它往往包含多个任务并分别标注边界并识别别类名称。
难点与分词是密切相关，不断有发现存在、存在歧义、构成复杂、多样性的问题同样通过 P/R/F 来衡量。
可以考虑部分正确归化。
命名实体识别的优化方案
按照 MUC-7 的定义，分为实体类 (人名、地名、机构名)、时间类 (日期、时间)、数值类 (货币、百分比)，严格定义为重复出现的普通名词 (飞机、公司)、人的团体名称以及以人名命名的奖项 (共获、诺贝尔奖)、非时间、数值的数字不属于命名实体。
命名实体识别的信息评价/评估方法与检索任务大致相同，采用 Precision / Recall / F-value 加以衡量。
正确率与召回率的计算公式：
方案 1：分子为返回的正确答案数量。
方案 2：分子为返回的正确答案数量 + 1。
得分的计算方式：
部分正确答案的“Severus/Person Sana” (类似正误、边界错误)。
命名实体识别方法 (1): 基于关键词 List Lookup
经常作为 NER 问题的基础解，预先构建命名实体知识库，出现在词典中的词可以识别为命名实体。
词典来自于自领域公开资源，例如黄页、电话簿、公开名单等或有权威的地理信息库。
优点是简单快捷，与具体语料/领域无关，容易部署更新 (无需重新训练)；但 (与基于词典/配词方法存在类似问题) 大部分情况下很难获得所有的命名实体名称。
构建和维护词典的成本代价较大，难以有效处理现实世界数据。
命名实体识别方法 (2): 基于规则
采用手工构造规则模型 (符合规则的词性进行识别标注。
适用特征包括主语信息、标点符号、关联语、指示词和方位词、位置词 (如第一)、中心词)。
以模式和特征与相应配词为主线，可借助语料/半自动化或特定表达达成匹配。
基于规则模型的方法在深度学习方法出现之前被广泛使用。
优点是当抽取的规则能较精确地捕捉语言现象时，性能较好。
缺点是不同表达对应不同规则，导致规则模型庞大，使用不便，规则往往依赖于具体语料、领域和文本风格，不同领域的以往往往差异较大，代价太大，系统维护周期长、修改性差且需要建立不同领域知识库。
命名实体识别方法 (3): 基于统计
下述的三大流方法，分支一、基于字类的 NER，将 NER 作一个多分类问题，通过设计特征提取与分类的方法加以解决，得到维度高而准确，再基于字类进行分类。
分支二、基于序列模型的 NER，与分词中的序列标注方法思路类似，区别在于标注不同，引入更多更细粒度的标注类别，将模型类型变为 HMM、CRF 以及基于序列模型度方法。
1 将词和对应的标注统一分布 (步分布) 向量化
2 基于编码器对文本进行表征。
基于词组 RNN 等
方法描述如下文
3 基于解码器进行分类。
还可以引入词根词源、首等信息。
基本的词性标注方法
与统计特征相关的问题有词性标注问题，词性是词基本面的语法属性，词性标注就是判定每个词语的语法范畴。
对于汉语，因缺乏词形变化信息，无法从词形形态判断词性；常用词可能有多义词，词标注方法需与实体识别方法分开讨论，总体思路，人工或通过大规模语料库训练 (按类类词搭配关系和上下语境语感进行类词判别规则)、HMM 等面向序列标注方法、统计与规则相结合的词性标注方法 (先基于规则排除明显规则，再基于统计模型标注，最后人工校核)
基于统计的命名实体识别方法
特点：对特征选取要求较高，需要从文本中选取对 NER 有帮助的特征来构建特征向量，通常做法是训练语料库所包含的信息语言信息进行统计分析和，从中挖掘出特征，对语料依赖较强、目前缺少通用的大规模语料 (对深度学习产生影响较大，特定专业领域数据更为明显)、大部分特征仍需要进行人工标注和概念标注。
实词消歧 disambiguation 一词多义，对不同含义输入实体仍需要统一标注，抽词并归纳，对关键进行定义和表述，通过相似性分析。
实体对齐 pair alignment/matching 也称实体匹配，对于异构数据源知识库中各个实体，找出属于现实世界中的同一实体。
一般利用实体属性信息判定不同源实体是否可对齐。
基础任务: 基于子集的
知识图谱子集体对齐
针对跨知识图谱的实体对齐任务，有多种基于子集体的模型，不同利用实体的属性和语义信息，还利用实体的关系信息、要求表征 (关系表和实体的表征) 落在同一向量空间。
这些模型更要求子文关系三一致。
跨语言社交网络用户匹配类似，不仅考虑用户画像，也考虑社交关系相似性。
实体消歧和实体对齐的区域标注
该标注旨在消除一词多义的现象，而实体对齐在于在表征同一对象上的多个实体之间构建对齐关系，丰富实体属性。
实体链接
将文本中的提及是链接到知识库中的对象上，应用于文本标注，知识图谱补充，信息检索和**问答系统**
文本中的提及是链接，链接到实体后可能应用数据图例等）等挑战包括数据源的多样性 (同一实体可能有多个不同提及) 和语义性 (同一提及可能对应多个不同实体)
任务框架
类似描述问题 (本质与前述内容/实体识别类似)、候选实体生成 (粗粒度描述)、主体候选子图匹配、候选实体选择 (核心问题)、主要基于子图关系
基于语义的关系抽取
从自由子图关系生成关系模式，然后基于子图模式构建新的关系，得到语义关系后，选择信程度高的关系，作为新种子，再寻找新的模式和新的关系。
不断迭代，直到没有新的关系或新的模式产生。
优点是特定语义关系抽取，但基于字面匹配，没考虑语义/语法/语义，难以保证模式可靠，移植性差，会为每个关系生成模式。
代表方法 1:DIPRE
双重语义关系抽取大思路在于给定一些已知关系类型的种子实体对，找到了这些实体对的 Occurrences，再学习 Occurrences 的模式，进而根据学到的模式，寻找更多符合该模式的数据，并加入到种子集中，不断迭代以实现关系抽取。
DIPRE 的动机：文本中表示三元组时，存在某些重复的“模式”。
DIPRE 的基本元素
元素表示关系实例，如 <Foundation, Isaac Asimov>-<Title, Author>；包含含常量和变量，例如 by 的符号表示为“title” by “author”
DIPRE 的基本流程
元素组往往广泛存在于各个不同源中，各个部分往往在位置上接近，存在这些近似时，存在着某些重复的“模式”。
启发式方法 (类似基于规则)
检查部分网站获取潜在“模式”，并用正则表达式来描述。
缺陷明显，网站/系统有特殊性，某网站模式未必适用于其他网站；人类不可能穷尽所有潜在模式。
而 DIPRE 考虑模式和子组实例之间双向重复关系，既要找到符合模式的内容，也要找到产生生成内容的模式。
DIPRE 的算法流程
1 输入一组种子实例实体 R，如若干 <title, author> 的实例对 (为什么是二元组？因为实体对，只需要包含两个实体)
2 基于种子实例集合 R，找到这些实例在网页中出现的网页 O，注意查找时保留上下文信息
3 基于种子实例实体 O，生成模式 P
4 基于生成模式 P，找到更多符合模式 R 的实例，此时可选择停止，或返回 2 继续基于新实例生成成新模式。
此时生成的新模式可能与之之前的模式有所类似？
Occurrence 可以理解为一组在网页中的呈现形式，只有非常接近才视作 O，否则则因太远可能导致语义互不相关。
将同一关系的不同实例在网页上呈现的形式 Occurrence 中相似内容排除，不同内容用简单取代，即可得到近似模式。
模式往往仅依赖特定网页/体系内，跨网站则模式可能不适用。
URL 前缀可用于判定属于同一网站/体系内，限定范围模式。
基于 DIPRE 算法、生成模式的基本步骤 1
若 O 由网络 (关系属性) 和 Middle (中间部分) 2 定义模式的 Order 和 Middle 即属于 O 集合的 Order 和 Middle，模式为 URLPrefix.Prefix.Suffix，分别 O 集合中最大的公共值。
其他部分用通配符填充。
代表方法 2:snowball
基本思想在于对 DIPRE 算法的提升，仅信任支持度 (Support) 和置信度 (Confidence) 较高模式加以保证模式质量。
支持度满足每个模式的模式的数量，将若干一定数量元素组成的模式予以剔除。
置信度按照如下公式计算：

P
r
o
b
a
b
i
l
i
t
y
=
P
r
o
b
a
b
i
l
i
t
y
+
N
e
g
a
t
i
v
e
.

{\displaystyle Probability=Probability+Negative.}

即考虑将负模式的模式剔除，确保符合相关关系的数据。
例：基本模式 P = <N, ORGANIZATION, <N, <N, >, LOCATION, >, >，可以得到以下三个实例：“Exxon, Irving, >”, “Intel, Santa Clara, cut prices”, “invest in Microsoft, New York-based analyst Jane Smith, >”, 其中前两个符合关系，而最后一个与事实不符，因此为 Negative，置信度为 2/3。
基于机器学习系抽取支持特征
(难点在于找出特征)，有些语言如 NLP 工具难以提取特征，而 NLP 工具输入人工错误信息，人工设计特征又不一定全。
远监督监督的由来和意义
而面向文本的关系抽取方法最大难点在于获取足够数量高质量标注 (人工标注开发具有主观性困难，而标注标注会有累积误差)
远监督监督如果某实体有某类关系，则所有含该实体的句子都用于描述该关系。
可用快速扩充训练数据。
远监督监督的基本思路
将包含指定关系实体的语料打

包，从中训练用于关系识别的模型，并用于判断更多实体之间的关系。
远监督类似于 DIPRE 算法，明显局限性在于字面漂移 (Semantic Drift)，不是所有包含该实体的句子都表该关系，错误链接导致关系判断错误，并通过不断迭代放错错误；
可以在每一轮迭代中观察选出的句子，把不包含该关系的句子删除掉，但过于复杂。
远监督监督的优化方案 (1): 动态转移矩阵
语音数据虽不可避免，仍可对语音数据模式进行统一描述。
引入动态转移矩阵，描述各类之间相互转移的概率。
把算得到关系分布乘转移矩阵，得到更加准确的与初始的关系分布结果。
随机性较差，并不能完全消除可靠性。
远监督监督的优化方案 (2): 远监督学习
远监督监督模型存在可能导致一些句子错误标注，提出新的生成模型，模拟模拟远监督模型的发展式标签过程。
设计相应的否模式类标签 NegPat(r)专门用于检测错误的标注，即某些关系的判断是否错误。
对于单一关系判断，可通过此次学习进行有效收敛，但训练收敛成本较大。
远监督监督的优化方案 (3): 注意力机制
即使是在被打入同一包里的句子，不同句子对训练关系模型的贡献度也不相同，贡献度可采用注意力模型衡量。
采用深度学习技术获取对整句个字的表示。
进而通过注意力机制将最表达该关系的句子挑选出来。
最有为根据依一个高维量的样本集来选。
开放关系抽取
预先定义的关系关系无法涵盖所有关系，可以基于知识图谱结构知识网络抽取，也可以通过识别标注语义/句法关系抽取标注。
事件抽取的概念
是处理事件 (事件推理/推理、推理) 的前提特征是特定人物，在特定时间特定地点进行用户产生的客观事实，一般呈现为句子/事件描述 (关系往往表现为动词/名词)。
事件抽取的基本要素
事件触发式 (事件发生的核心词，多动词/名词，动词与分类基本任务)、事件类型 (对应触发词，往往通过触发词分类问题)
事件元素 (事件要素与、主要由实体/时间/物等组成)
事件元素角色 (事件元素在事件中充当角色)
事件抽取的模板
通过触发词识别事件分类问题及其类型，模板抽取模板抽取。
模板元类 TE 目标在于更完整地描述标注实体。
描述标注实体基本属性信息，内容可包括名称/类别/种类等，不同类型的事件特征描述也不尽相同。
通过事件元素元角色信息的识别，将元素填入模板合适的槽，即完成了事件抽取。

第十四节知识图谱与图计算 图表示学习方法
将图数据进行向量化表征映射到低维向量空间，在这个低维向量空间中，最大程度保留图结构特征语义特征信息。
图的基本表示形式之一：邻接矩阵
每行表示一个节点，0 代表没有边与对应节点为连接。
该行可以看作该节点的代表向量，用于最能精确聚类关系。
局限性在于未能充分关注节点特征信息，无法融入节点属性信息。
几种不同的图表示学习方法
基于随机游走的图表示学习 (基于随机游走的节点序列得到挖掘图结构信息,DeepWalk, Node2vec, Metapath2vec, LINE)
基于图神经网络的网络表示学习 (利用神经网络来学习图网络结构，抽取和挖掘图结构中的特征和模式, GCN, SGCN、VGAE、GraphSAGE、SDNE、GAT)
知识图谱 + 图学习：知识图谱表示学习
为知识图谱中实体 (点) 和关系 (边) 提供向量化表示同时保留语义信息，表征结构信息同时融合语义关系类差别信息，知识图谱推理补全、推荐系统、网络系统。
知识图谱表示学习模型 TransE
: word2vec 翻译不变性发现 (vking)(v+queen)=v(man)(v+woman)，但无法做到一对多/多对一对多/多对多，TransH: 通过向量

u
，将实体投影到关系 R 对应超平面下解决一对多问题 (原关系多维向量变为同一平面内，可以做差)。
TransR: 通过矩阵 M，将实体进行投影坐标到关系 R 对应关系空间内，不同关系对应不同属性。
TransD: 关系 R 对应投影矩阵进行坐标差，描述同一关系不同属性。
DistMult: 每个关系 R 被表示成向量

r
，建模潜在因子构成相互作用，把关系当成线性关系 (限制对角矩阵)。
RotatE: 建模并推断各关系矩阵，包括对称/反转和乘法，将每个关系 R 定义为复向量空间变换。
知识图谱推理补全
推理是从已知事实推出新事实或知识的过程，演绎/归纳/归约/类比。
知识图谱推理主要基于推理引擎将已知/关系推导出未知事实/关系。
现实问题根据预测/因果推理/基于 KG 向图内容都可表为图推理问题。
借助符号逻辑或表征学习，实现基于图谱的推理过程，增强可解释性。
关系推理任务
链接预测 (给定头尾实体，预测是否存在某关系)，实体预测 (给定头尾实体关系，预测未知头尾实体)，事实三元组预测 (给定三元组判断真假)
基本上很多现有 KG 数据集稀疏，需要模型添加全新添加三元组，基于已有事实推理缺失事实。
关系补全类似，推理结果可作为为新知如加入图谱。
包括头/尾实体预测和关系预测：给定三元组 (.,r,.)，预测头实体；给定三元组 (h.,r,.)，预测尾实体；给定三元组 (h.,t,.)，来预测关系。
任务评价方式
Hit@N: 所有预测样本排名在 n 内部比值的。
Mean Recall@MR: 所有预测样本平均排名。
Mean Reciprocal Rank(MRR): 对所有预测样本排名求倒数再平均。
推理/补全方法分类
基于符号推理：显式知识表示方法，一般需人工定义，可解释性较好，但能进行推理相对有限。
知识图谱表示学习：易于拓展表示知识，对数据源友好，但不利于人工理解。
大规模图神经网络
把常规图谱三元组问题全面视作序列分类任务，让模型微调。
头/尾实体类问题中，将三元组转为为文本输入，获得 [CLS] 以征向量向量化分类，判断是否合理。
关系补全中，模型也可完成关系分类，输入两实体，输出多类关系，类似对数似然函数。

HMM 例则
S 是所有可能的状态的集合,S = {T,F}, 观测序列为 {X,Y,Z}, 隐含状态转移概率矩阵为

A
=

[

0.7
0.3

0.4
0.6

]

，初始状态概率分布为

π
=
[
0.6
,
0.4
]
，求最优状态序列。
解:1
确定

δ

1

(

s

1

)
 为时刻 t 状态

s

1

 的概率。
则

t
=
1
 时，

δ

1

(
T
)
=
0.6
×
0.5
=
0.3
,

δ

1

(
F
)
=
0.4
×
0.1
=
0.04
,

t
=
2
 时，

δ

2

(
T
)
=
max
(

δ

1

(
T
)
×
0.7
,

δ

1

(
F
)
×
0.4
)
×
0.4
=
0.3
×
0.7
×
0.4
=
0.084
,

δ

2

(
F
)
=
max
(
delta_1(T) × 0.3,

δ

1

(
F
)
×
0.6
)
=
0.2
,

δ

2

(
F
)
=
T
,

t
=
3
 时，

δ

3

(
T
)
=
max
(

δ

2

(
T
)
× 0.7
,

δ

2

(
F
)
×
0.4
)
×
0.1
=
0.084
×
0.7
×
0.1
=
0.00588
,

δ

3

(
F
)
=
max
(

δ

2

(
T
)
× 0.3,

δ

2

(
F
)
×
0.6
)
=
0.084
×
0.3
×
0.6
=
0.01512
,

δ

3

(
F
)
=
T
,

δ

3

(
F
)
=
T
,

t
=
3
 时，取

T
，因为

P
 的概率更大，所以

t
=
2
 时状态为

δ

3

(
F
)
=
T
，

t
=
1
 时状态为

δ

2

(
T
)
=
T
，所以最优状态序列为

T
→
T
→
F
。
PCA 例则
给定 4 个二维向量

x

1

=
(
4
,
1
)
,

x

2

=
(
2
,
3
)
,

x

3

=
(
5
,
4
)
,

x

4

=
(
1
,
0
)
，使用主成分分析得到特征向量空间降为一维。
解:计算每一维的均值

μ

1

=

4
+
2
+
5
+
1

4

√

p
e
r
o
d
u
c
t
(
P
(

r

a
,
p

)
,
P
(

r

b
,
p

)
)

{\displaystyle {\sqrt {product(P(r_{a,p}),P(r_{b,p}))}}}

Product(P) 为用户 a,b 都评价过的项目集合。
基于用户推荐的评分预测
在得到用户之间的相似性后，针对待预测的项目，可以根据历史上其他用户行为对于该项目的评分，结合用户之间的相似性作为加权， 预测评分结果。
相似性计算与评分预测中，都通过减去平均来抹去评分平均偏好的影响（不同用户打分范围不同）

p
r
e
d
(
a
,
p
)
=

r

a
+

∑

b
∈
n
e
i
g
h
b
o
r
s
(
a
)

s
i
m
(
a
,
b
)
⋅
(

r

b
,
p

−

r

b
,
p

)

∑

b
∈
n
e
i
g
h
b
o
r
s
(
a
)
s
i
m
(
a
,
b
)

{\displaystyle pred(a,p)={\frac {r_{a,p}+\sum _{b\in neighbors(a)}sim(a,b)\cdot (r_{b,p}-r_{b,p})}{\sum _{b\in neighbors(a)}sim(a,b)}}}

相似性计算与评分预测中，都通过