

Вопрос №2.16

Архитектура поисковых систем. PageRank. Автоматическая классификация текстов. Семантические сети, Semantic web. Оценка качества результатов поиска.

1 Архитектура поисковых систем

- Поисковый робот (web crawler). Ходит по ссылкам в Интернете и индексирует найденные документы, занося в базу данных.
- База данных. Накапливает информацию о найденных страницах в формате, удобном для использования алгоритмами поиска: прямой (документ — слова) и обратный (слово — документы) индексы.
- Клиент. Обрабатывает поисковые запросы пользователей, извлекая информацию из базы данных.

Альтернативные архитектуры.

- Распределенная поисковая система: база данных распределена по ряду серверов.
- Метапоисковая система: вместо собственного поискового робота и базы данных используются запросы к другим поисковым системам.

2 PageRank

PageRank — алгоритм ранжирования результатов поиска. Ранжирование производится по рейтингу, вычисляемому через количество ссылок на данный документ и рейтинги ссылающихся документов.

PageRank использует модель случайного блуждания.

- Ориентированный граф: документы — вершины, ссылки — дуги.
- Начинаем в случайной вершине.
- С вероятностью ε переходим в случайную вершину.
- С вероятностью $1 - \varepsilon$ переходим по случайной исходящей дуге.

$\text{PR}_k(i)$ — вероятность находиться в вершине i на k -м шаге. $\text{PR}_0(i) = 1/n$, где n — количество вершин.

$\text{PR}(i) = \lim_{k \rightarrow \infty} \text{PR}_k(i)$. На практике ограничиваются $\text{PR}_{50}(i)$.

Важные соотношения:

$$\text{PR}_k(i) = \varepsilon + (1 - \varepsilon) \sum_{j: j \rightarrow i} \frac{\text{PR}_{k-1}(j)}{|\{l: j \rightarrow l\}|},$$

$$\lim_{k \rightarrow \infty} : \quad \text{PR}(i) = \varepsilon + (1 - \varepsilon) \sum_{j: j \rightarrow i} \frac{\text{PR}(j)}{|\{l: j \rightarrow l\}|},$$

где $x \rightarrow y$ означает наличие дуги из x в y .

PageRank в векторной форме.

- $\overline{\text{PR}}_k = (\text{PR}_k(1), \dots, \text{PR}_k(n))^T$
- $L = (l_{ij})$:
 - $l_{ij} = \varepsilon + (1 - \varepsilon) \frac{1}{|\{l: j \rightarrow l\}|}$, если $j \rightarrow i$,
 - $l_{ij} = \varepsilon$ иначе.
- $\overline{\text{PR}}_k = L^k \overline{\text{PR}}_0$
- $\overline{\text{PR}} = L \overline{\text{PR}}$, т.е. $\overline{\text{PR}}$ — собственный вектор матрицы L .

3 Автоматическая классификация текстов.

Постановка задачи.

- Множество документов $D = \{D_1, \dots, D_n\}$.
- Множество категорий $C = \{C_1, \dots, C_m\}$.
- Неизвестная функция $\Phi: D \times C \rightarrow \{0, 1\}$.
- Задача: построить классификатор Φ' , максимально близкий к Φ .
- Иногда достаточно построить $\Phi'': D \times C \rightarrow [0, 1]$, который задает не точную классификацию, а ранжирование категорий для каждого документа.
- От ранжирования легко перейти к точной классификации, введя некоторый порог $0 < t < 1$: $\Phi'(D_i, C_j) = 1 \iff \Phi''(D_i, C_j) \geq t$.

Применение.

- Фильтрация документов, распознавание спама.
- Наполнение интернет-каталогов.
- Классификация новостей.
- Контекстная реклама.
- Персональные новости.

Пусть каждый документ D_i представлен в виде вектора весов термов (t_1, \dots, t_p) . Рассмотрим пример линейного on-line классификатора.

- $\Phi'(D_i, C_j) = \frac{D_i \cdot C_j}{|D_i| |C_j|}$, где векторы коэффициентов C_j вычисляются динамически по мере обработки обучающего множества.
- Начинаем с $C_j = (1, \dots, 1)$.
- Для каждого учебного документа применяем текущее правило.
- При неудаче вносим поправки $+\alpha/-\beta$ в коэффициенты, соответствующие термам неправильно классифицированного документа.

4 Семантические сети, Semantic web.

Семантическая сеть (semantic network) — информационная модель предметной области, имеющая вид ориентированного графа, вершины которого соответствуют объектам предметной области, а дуги задают отношения между ними.

Семантические сети бывают однородными (только один тип отношения) и неоднородными (разные типы отношений).

Semantic Web — часть глобальной концепции развития сети Интернет, целью которой является реализация возможности машинной обработки информации за счет использования семантических сетей для описания документов и их отношений.

Преимущества Semantic Web.

- Семантический поиск.
- Вопросо-ответные системы.
- Агенты в семантическом Вебе.
- Объединение знаний (интеграция баз данных).

Языки описания для Semantic Web (основаны на XML).

- RDF (Resource Description Framework) — описывает конкретные объекты и их отношения (синтаксис).
- OWL (Ontology Web Language) — описывает типы объектов и типы отношений (семантика).

5 Оценка качества результатов поиска.

- Точность (precision) — отношение числа найденных релевантных документов к общему числу найденных документов.
- Полнота (recall) — отношение числа найденных релевантных документов к общему числу релевантных документов в базе.
- F-мера (F-measure, мера Ван Ризбергена) — среднее гармоническое точности и полноты: $F = 2PR/(P + R)$.

Основной источник информации - курс «Алгоритмы для Интернета»:
<http://yury.name/internet.html>