

検索の仕組み

索引検索



- あらかじめキーワードの索引を作っておく
- 例
 - 本の巻末の索引を調べる
 - OPAC
 - Google などの検索エンジン
- データ量が多くても検索時間は変わらない

本の



Tron 56	ア
tsr-van2 94	アクセス権限 145, 160
U	アクセス制御 142, 143, 145
Ulrich's Periodicals Directory 126	アグリゲータ 53, 111, 112
Uniform Card 5	圧縮方式 59
UNIX 56	イ
URL (Uniform Resource Locator) 65	医学情報 122
USB (Universal Serial Bus) 63	医学中央雑誌 3, 100, 123
W	医学用語ソーラス 14
WAN (Wide Area Network) 64	意匠制度 202
Web 66	一次情報 49
検索エンジン 70	データベース 82, 83
サービス 71	一次資料 8
情報発信 73	医中誌 Web 100
Web 2.0 76	インターネット 64
Web of Science 97	歴史 64
Web 検索技術 72	インタビュー 41
Webcat 108	インバウト・ファクター 36
Webcat Plus 24, 108	インフォトリックス・ビブリアトリックスを参照
Web サイト 131	引用
Wiki 75	Web 情報の 133
Wikimedia 75, 133	著作物の 208

索引検索の例



Index 検索 画像 動画 地図 ニュース ショッピング Gmail もっと見る ログイン

Google 日本

Google 検索 I'm Feeling Lucky

著作権 © 2008 Google Inc. 日本での検索は Google によって提供されています。

広告掲載 ビジネスソリューション サプライズと利用規約

Google Google について Google.com in English

索引検索の例



簡易検索 詳細検索

検索対象
○大妻女子大学
□他大学図書館 (NTL)

資料区分
☐図書
☐一和図書
☐一洋図書
☐雑誌
☐一和雑誌
☐一洋雑誌
☐AI資料
☐オンラインDB

検索条件
AND 夏目漱石 AND 全ての項目から 著者名に在の語を含む 出版社・出版者

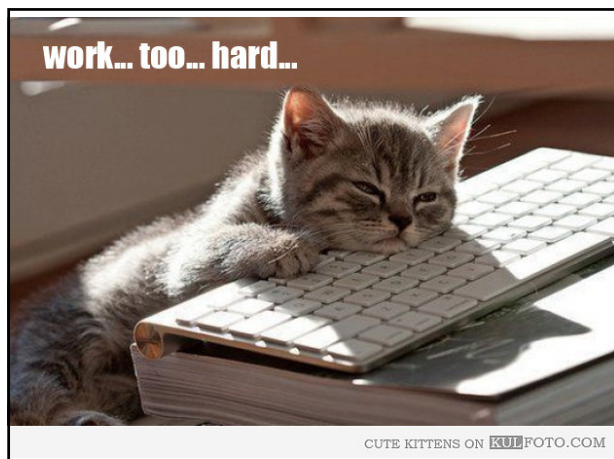
検索 オプション
出版年: 年 月 日
言語: (指定なし)
媒体種別: (指定なし)
配架場所: (指定なし)
並び順: 昇順 降順
一覧表示件数: 20

検索 クリア

索引ファイル



- 索引語 (キーワード) の抽出
 - 対象レコード (文書、ページ) から検索の手がかり (語、句) を取り出す
 - レコード (ページ) 番号をつける
 - 索引語を 50 音順、ABC 順に並べる
- 検索
 - 語句が一致するものを探す
 - 該当レコード (ページ) を取り出す



索引語の抽出



- 昔は人手で抽出
- 今はコンピュータが自動抽出
 - 英語は単語がスペースで区切られているので簡単
 - 日本語は単語の区切りがない
 - 形態素解析法
 - バイグラム法

索引語の抽出 (英語)



- 単語がスペースで区切られている

資料番号: 201833

標題: Information Retrieval in Chemistry and Chemical Patent Law

「Information」「Retrieval」「in」「Chemistry」「and」
「Chemical」「Patent」「Law」

15

索引語の抽出 (日本語)



- 形態素解析法
 - 文法解析と辞書を使って単語を区切る

資料番号: 98170

標題: インターネット時代の化学文献とデータベースの活用法

- まず助詞などで語句を切り出す
「インターネット時代」「化学文献」「データベース」「活用法」
- カタカナの部分は単語として切り出す
 - 「インターネット」「時代」
- 次に語句から辞書を使って単語を取り出す
「化学文献」→「化学」「文献」

16

索引語の抽出 (日本語)



- バイグラム (bigram) 法
 - 1 文字ずつずらしながら、2 文字のペアを機械的に索引

資料番号: 98170

標題: インターネット時代の化学文献とデータベースの活用法

「イン」「ンタ」「ター」「ーネ」「ネッ」「ット」
「ト時」「時代」「代化」「化学」「学文」「文献」
「献と」「とデ」「デー」「ータ」「タベ」「ベー」
「ース」「スの」「の活」「活用」「用法」

17

形態素解析とバイグラム



- 形態素解析
 - 検索結果にノイズが少ない
 - 切り出しのための辞書の整備が必要
 - 検索エンジンなどで用いられる
 - Google, Yahoo!

18

形態素解析とバイグラム

- バイグラム
 - システムが簡単・安価
 - 辞書がいらない
 - 断片的な単語で検索できる
 - ノイズがある
 - 「京都」を検索しようすると「東京都」もヒットする
 - 「タイ」を検索しようすると「スタイル」もヒットする
 - 社内文書など小規模なシステムでよく用いられる
 - OPAC, NDL-OPAC、新聞記事データベース



19



<http://kae.de>

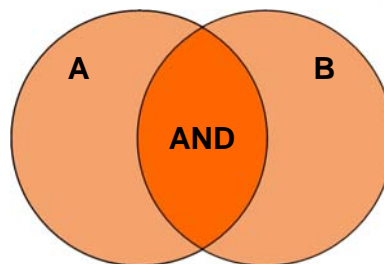
ブール演算

- 論理演算、集合演算ともいう
- 集合間の
 - 積 AND
 - 和 OR
 - 差 NOT
 - を求める
- 情報検索で用いる



21

AND 検索



22

Google の検索オプション

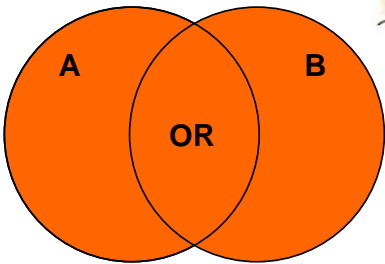


Google の AND



検索の仕組み

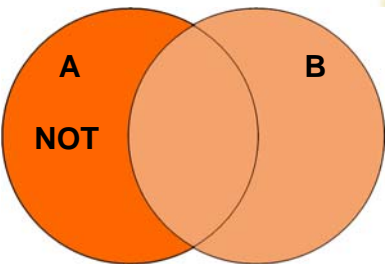
OR 検索



Google の OR

A screenshot of the Google search homepage. The search bar is at the top. Below it, there are several input fields for search options. The "検索オプション" (Search Options) section is visible. The "検索するキーワード" (Search keywords) field is empty. The "すべてのキーワードを含む" (Include all keywords) field is empty. The "記録も含め完全一致" (Include records, exact match) field is empty. The "いずれかのキーワードを含む" (Include any keyword) field is empty. The "含めないキーワード" (Exclude keywords) field is empty. The "検索の範囲" (Search range) field is empty. The "検索" (Search) button is at the bottom right. The text "Google 渡辺 OR 渡辺 OR 渡辺" is visible in the search bar area.

NOT 検索



Google の NOT

A screenshot of the Google search homepage. The search bar is at the top. Below it, there are several input fields for search options. The "検索オプション" (Search Options) section is visible. The "検索するキーワード" (Search keywords) field is empty. The "すべてのキーワードを含む" (Include all keywords) field is empty. The "記録も含め完全一致" (Include records, exact match) field is empty. The "いずれかのキーワードを含む" (Include any keyword) field is empty. The "含めないキーワード" (Exclude keywords) field is empty. The "検索の範囲" (Search range) field is empty. The "検索" (Search) button is at the bottom right. The text "Google 渡辺 - 渡辺 - 渡辺" is visible in the search bar area.

検索の仕組み (AND)

資料番号: 98170
標題: インターネット時代の化学文献とデータベースの活用法

索引語	件数	レコードリスト
インターナル	253	... 50405, 52183, 60031, 64521, ...
インターネット	3862	... 50937, 64742, 68451, 69227, 69228, 85258, 89986, 92454, 94755, 98170, 103220 ...
インターン	28	... 67940, 54697, 54336, ...
データ	3301	... 84161, 84594, 85258, 85850, 90102, 93397, 95151, 99150, 99168, 98170, 98175, 98572, 101090, ...
データベース	485	... 53219, 54083, 77305, 83210, 98170, 100351, 123087, ...
データマイニング	71	... 83943, 88541, 88630, 100455, ...

近接演算

- 語と語のあいだの距離を調べて検索する
 - 距離が近いほど適合性が高い
- 通常、2つの語が続いていることを指定
 - Google では「完全一致」と呼ぶ

検索の仕組み

Google の近接演算 (フリーズ)

検索オプション

検索するキーワード

すべてのキーワードを含む

語句も含め完全一致

いずれかのキーワードを含む

含めないキーワード

検索の範囲

検索




検索の仕組み (近接)

- 単語間の距離をはかって、適合性を見る


資料番号: 98170
標題: インターネット時代の化学文献とデータベースの活用法
1 2 3 4 5 6

索引語	件数	レコードリスト
インターナル	253	... 50405-7, 52183-12, 60031-3, 64521-5, ...
インターネット	3982	... 50637-1, 64742-8, 68451-8, 69227-17, 69228-6, 85258-2, 89866-10, 92454-7, 94755-8, 98170-1, 103220 ...
インターン	28	... 67940-2, 64697-10, 54336-18, ...
データ	3301	... 84181-9, 84594-1, 85258-7, 85850-8, 90102-13, 93397-15, 95151-18, 89150-4, 98168-4, 98170-5, 98175-8, 98572-13, 101090-2, ...
データベース	485	... 53219-8, 54083-7, 77305-11, 83210-5, 98170-5, 100351-8, 123087-9, ...
データマイニング	71	... 83943-17, 88541-3, 88630-2, 100455-9, ...



検索フィールド

- 検索する項目のこと
 - タイトル
 - 著者名
 - 件名 (統制語)
 - フリーワード (タイトル、著者名などすべて)



検索フィールド

詳細検索 ?

キーワード

And

タイトル

And

著者

And

出版者

And

件名

And

ISBN

And

分類

And

選択してください

著者名・件名典拠検索

各種番号・コード

NDLC

NDC

☒ 広範囲に検索(ノイズ多め)(日中韓)

☐ フレーズ検索(欧文) ?



