

数式検索システムの実用化に向けて

- 数式検索システムの基本アルゴリズム
- 論文検索などへの応用の検討

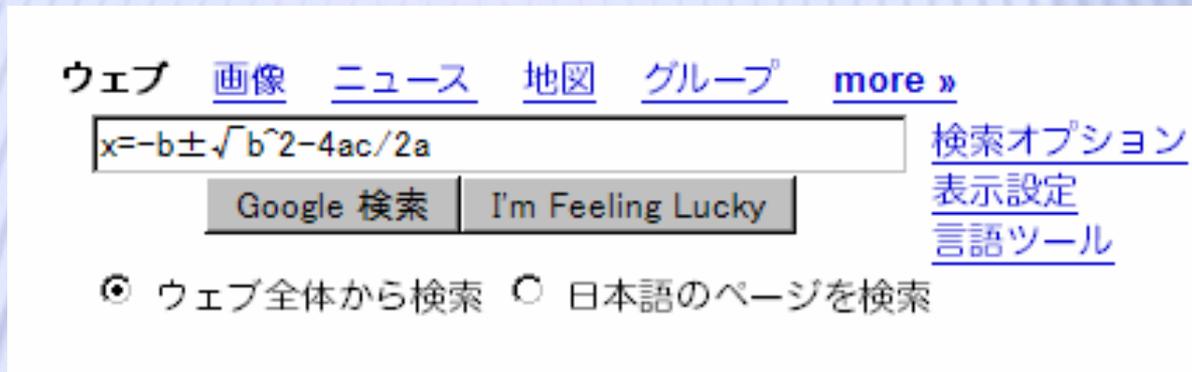
RIMS研究集会「紀要の電子化と周辺の話題」

2008年9月3日

大阪大学基礎工学研究科 橋本 英樹, 土方 嘉徳, 西田 正吾

数式に特化した検索

これまでの検索エンジンでは $\int \frac{1}{\sqrt{1-x^2}} dx$ のような数式をキーワード（クエリ）として数式を含む資料を検索することは困難でした。



数式は科学技術や社会現象などの知識を表現する最良の方法です。

数式検索技術を用いた教材やサービスの提供により、学生、教師、研究者、エンジニアを支援し、知的生産性の向上に寄与することを目標とします。

デモンストレーション

Search Engine DEMO for Mathematical Formulas (Ver.0.5)
by [Research Institute for Mathematical Communications Inc.](#)

Input a query

Formula

General Operators Big Ops. Symbols Matrix Arrows Greek Script and accent

() [] < > √ ∫

[] ||| ∫ √

$\lim_{n \rightarrow \infty} \frac{1-x^n}{1+x^n}$

Search

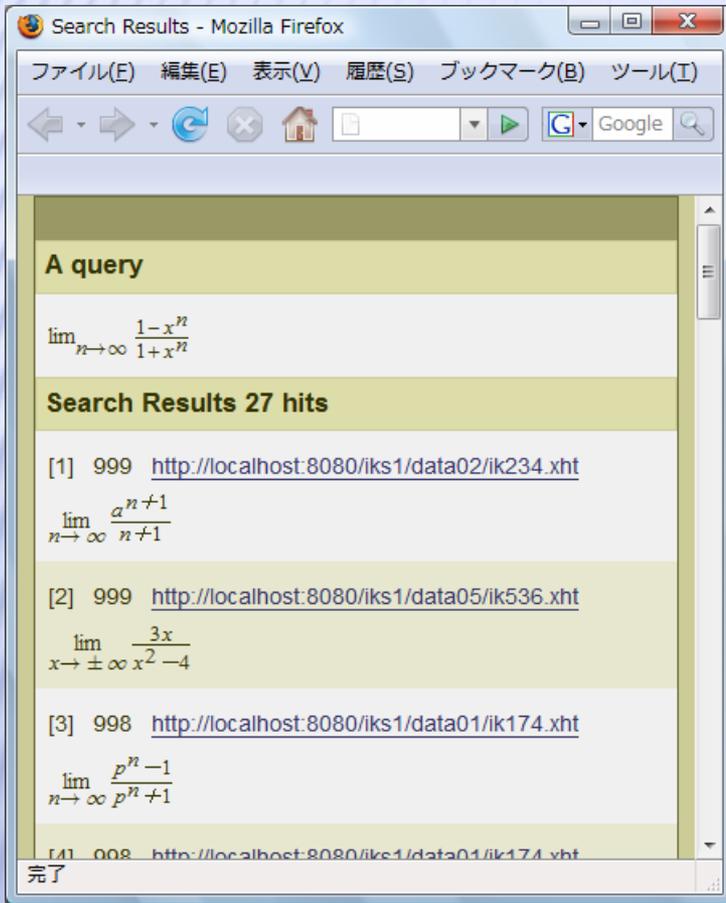
Above cool editor is WIRIS Editor, developed by [Maths for More](#). WIRIS Editor is a powerful tool to create mathematical formulas in HTML browsers.
If you are interested in WIRIS Editor, please contact them: [WIRIS](#)

Copyright 2007 Research Institute for Mathematical Communications Inc. All rights reserved.

分数、√などのアイコンを使って、
検索したい数式をマウスで入力

デモンストレーション

検索結果表示



Search Results - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I)

← → ↻ × 🏠 📄 ▶ 🔍 Google

A query

$$\lim_{n \rightarrow \infty} \frac{1-x^n}{1+x^n}$$

Search Results 27 hits

[1] 999 <http://localhost:8080/iks1/data02/ik234.xht>

$$\lim_{n \rightarrow \infty} \frac{a^n+1}{n+1}$$

[2] 999 <http://localhost:8080/iks1/data05/ik536.xht>

$$\lim_{x \rightarrow \pm \infty} \frac{3x}{x^2-4}$$

[3] 998 <http://localhost:8080/iks1/data01/ik174.xht>

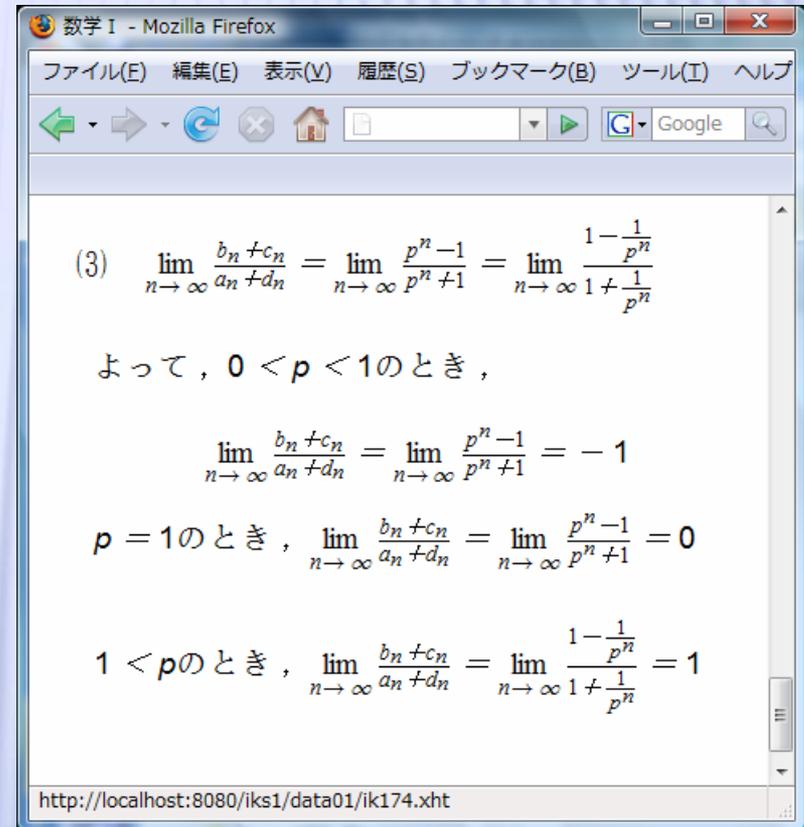
$$\lim_{n \rightarrow \infty} \frac{p^n-1}{p^n+1}$$

[4] 008 <http://localhost:8080/iks1/data01/ik174.xht>

完了



見つかった文書を表示



数学 I - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ

← → ↻ × 🏠 📄 ▶ 🔍 Google

(3)
$$\lim_{n \rightarrow \infty} \frac{b_n+c_n}{a_n+d_n} = \lim_{n \rightarrow \infty} \frac{p^n-1}{p^n+1} = \lim_{n \rightarrow \infty} \frac{1-\frac{1}{p^n}}{1+\frac{1}{p^n}}$$

よって、 $0 < p < 1$ のとき、

$$\lim_{n \rightarrow \infty} \frac{b_n+c_n}{a_n+d_n} = \lim_{n \rightarrow \infty} \frac{p^n-1}{p^n+1} = -1$$

$p = 1$ のとき、
$$\lim_{n \rightarrow \infty} \frac{b_n+c_n}{a_n+d_n} = \lim_{n \rightarrow \infty} \frac{p^n-1}{p^n+1} = 0$$

$1 < p$ のとき、
$$\lim_{n \rightarrow \infty} \frac{b_n+c_n}{a_n+d_n} = \lim_{n \rightarrow \infty} \frac{1-\frac{1}{p^n}}{1+\frac{1}{p^n}} = 1$$

<http://localhost:8080/iks1/data01/ik174.xht>

MathMLの利用

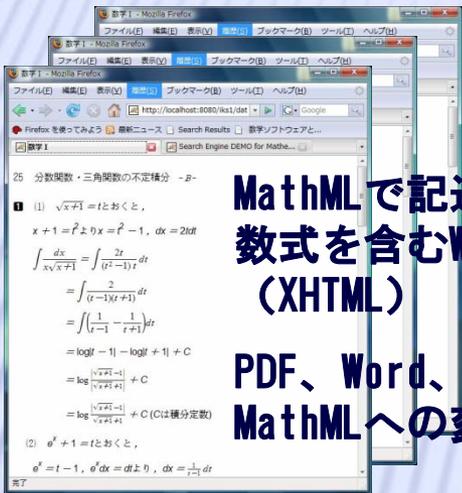
- MathMLで記述された数式を含むWebページをブラウザで表示可能
- 例えばインターネットで見つけたMathMLオブジェクトをMathematicaやMapleなどの数式処理ソフトに入力して計算を行うなど、数式データの再利用が容易
- TeX、WordからMathMLへの変換も可能
- OpenDocumentFormat (ODF) でもMathMLを採用

平成19年3月に各府省情報化統括責任者(CIO)連絡会議において決定された「情報システムに係る政府調達の基本指針」では、ODFファイルフォーマットのよ
うなオープンな国際標準規格の採用を優先するとされています。

- 日本数学会の欧文誌 Journal of Mathematical Society of Japan では
アブストラクトをMathMLで記述
- DLMFプロジェクト (NIST)、Project Euclid、など
- PDF → MathML (InftyProject)

システム構成例

検索対象文書集合

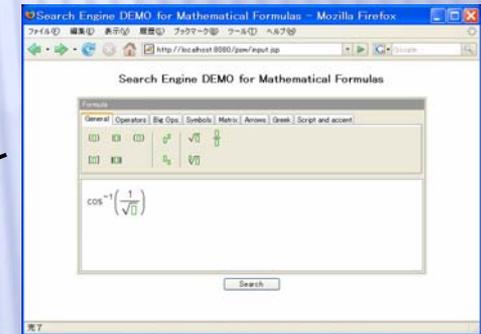


検索サーバー



ユーザーPC

クエリ入力プログラム



検索結果表示プログラム



(全文) 検索エンジンの仕組み

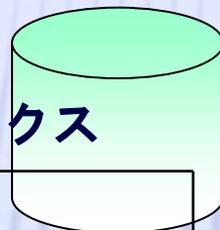


(文書ID=20)
『東京都の輸出産業の、、、』



『東京』 『輸出』 で検索

Nグラムを用いた転置インデックス



Nグラム (N=2)	文書ID
東京	1, 18, 20, 49, , ,
京都	20, 103, , ,
都の	5, 20, 21, 75, , ,
輸出	5, 18, 20, 146,
出産	20, 390,
産業	1, 5, 18, 20, 36, , ,

文書ID:
18, 20, ...

検索結果表示

適合文書
リスト

- ・この例では『京都』『出産』を検索したとき、、、
- ・プラットフォーム ⇔ プラットホーム
- ・パソコン ⇔ PC

数式の構造を反映した検索

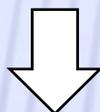
日本語の検索では文書中の文字や単語の出現情報を用いて検索を実行する。これに対して、数式の検索では、、、

xy

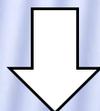


どちらの数式も x と y の2つの記号を含んでいるが、数式の意味と構造は異なる

x^y



従来のシステムでは、 x と y という記号の出現情報だけを用いて検索を行うため、意図した情報を検索することは困難



数式の構造を反映した検索の仕組みが必要

MathMLによる数式表記

$$\sqrt{x^2 + y^2}$$



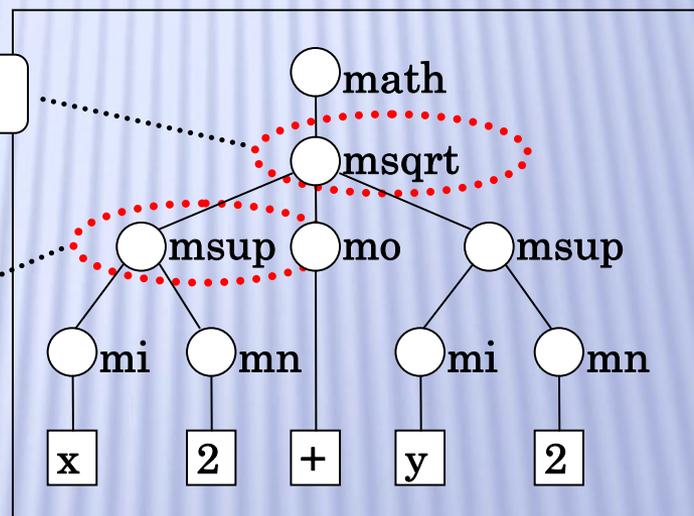
MathMLによる表記

```
<math>  
<msqrt>  
<msup>  
<mi>x</mi>  
<mn>2</mn>  
</msup>  
<mo>+</mo>  
<msup>  
<mi>y</mi>  
<mn>2</mn>  
</msup>  
</msqrt>  
</math>
```

平方根

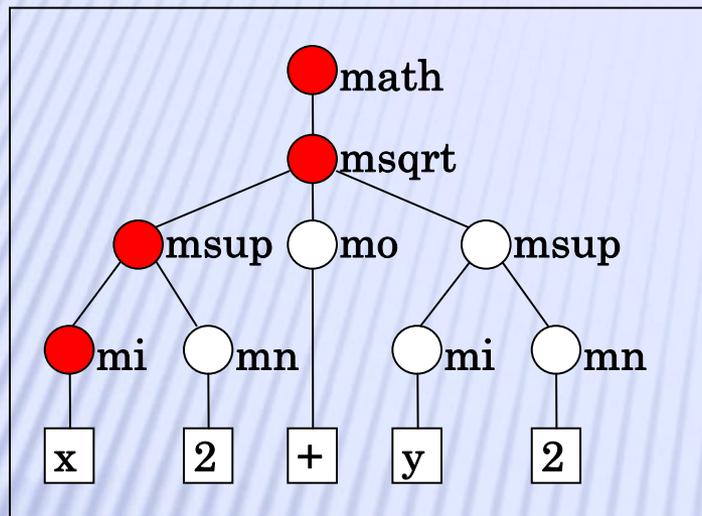
上付き

DOMツリー



MathMLはWeb上で数式を表記するため、World Wide Web Consortium (W3C)により勧告された国際規格です

数式の階層構造を用いたインデックス



MathMLオブジェクトのDOMツリーのルートから葉ノードまでのパスのXPath表記をキーとして転置ファイルを作成

XPath表記	文書IDリスト
<code>/math/msqrt/msup/mi[text()="x"]</code>	1,52,70,271,,,
<code>/math/mfrac/mo[text()=" ∂ "]</code>	2,16,55,102,,,
<code>/math/mfrac/mfrac/mi[text()=" π "]</code>	22,93,181,,,

類似数式検索

検索キー

$$\pi = 2 \int_0^{\infty} \frac{\sin^2(t)}{t^2} dt$$

ヒットした公式

$$\pi = 2 \int_0^{\infty} \frac{\sin^2(t)}{t^2} dt$$

$$\pi = \frac{8}{3} \int_0^{\infty} \frac{\sin^3(t)}{t^3} dt$$

$$\pi = 3 \int_0^{\infty} \frac{\sin^4(t)}{t^4} dt$$

$$\pi = \frac{384}{115} \int_0^{\infty} \frac{\sin^5(t)}{t^5} dt$$

$$\pi = \left(2^n \int_0^{\infty} \frac{\sin^n(t)}{t^n} dt \right) / \left(n \sum_{k=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{(-1)^k (n-2k)^{n-1}}{k!(n-k)!} \right)$$

検索結果のソート

- ・ Googleライクなランキング
- ・ 論文の閲覧回数、引用回数によってソートして表示
- ・ 出版年月日
- ・ 目次順

数式の類似度に基づくランキング

x^2+y^2 に含まれるパスのXPath表記

math/msup/mi[text()='x']
math/msup/mn[text()='2']
math/mo[text()='+']
math/msup/mi[text()='y']
math/msup/mn[text()='2']

$x^2+2ax+a^2$ に含まれるパスのXPath表記

math/msup/mi[text()='x']
math/msup/mn[text()='2']
math/mo[text()='+']
math/msup/mn[text()='2']
math/msup/mi[text()='a']
math/msup/mi[text()='x']
math/mo[text()='+']
math/msup/mi[text()='a']
math/msup/mn[text()='2']

一致するパスが
多いほど
類似度が高い

数式特有の問題への対応

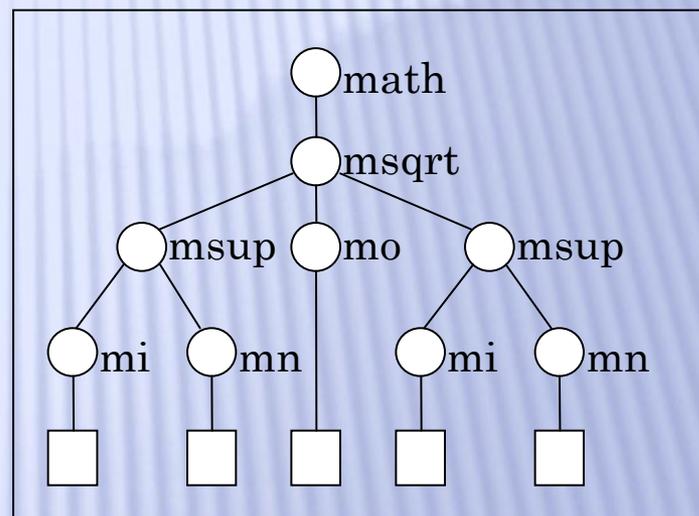
$$\sqrt{x^2 + y^2}$$

```
<math>  
<msqrt>  
  <msup>  
    <mi>x</mi>  
    <mn>2</mn>  
  </msup>  
  <mo>+</mo>  
  <msup>  
    <mi>y</mi>  
    <mn>2</mn>  
  </msup>  
</msqrt>  
</math>
```

$$\sqrt{a^2 + b^2}$$

```
<math>  
<msqrt>  
  <msup>  
    <mi>a</mi>  
    <mn>2</mn>  
  </msup>  
  <mo>+</mo>  
  <msup>  
    <mi>b</mi>  
    <mn>2</mn>  
  </msup>  
</msqrt>  
</math>
```

ツリー構造は同一



ツリー構造だけに着目することにより、異なる文字で表記された数式の検索が可能

複数の表記法

$$\nabla^2 u \Leftrightarrow \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) u \Leftrightarrow \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}$$

$$\frac{1}{x} \Leftrightarrow x^{-1}$$

複数の表記を同じものとして扱うためのルールを用意してシステムに組み込む必要がある

応用事例の構築

数式検索エンジンだけあっても意味がない。
検索対象となる有用なコンテンツと検索エンジンが組み合わさることによって利用価値が生まれる。

- 論文検索、技術情報検索、特許検索
- 教材検索
- 数学辞典検索、数学公式集検索
- Q&Aサイト

誰が使うのか？

- 分野の専門家、異分野の研究者／技術者
- 学生、教師、教材作成者
- 普通の人好奇心を満たすために

対象分野ごとに
コンテンツの
特徴が異なる



カスタマイズが必要

文献検索サービスでの利用イメージ（将来像）

■ 検索条件を指定してください

▼キーワード検索条件

キーワード条件クリア

*複数の語を入力する場合は、スペースで語と語を空けてください。また、語と語の間のスペースは**AND**、**OR**の選択ができます

(*)付フィールドは完全一致検索です

語間のスペースを AND OR とする

キーワード (必須)

AND

AND

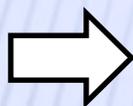
検索条件を追加する

- ・ キーワードと数式を組み合わせた検索によって
正確な絞り込みが可能

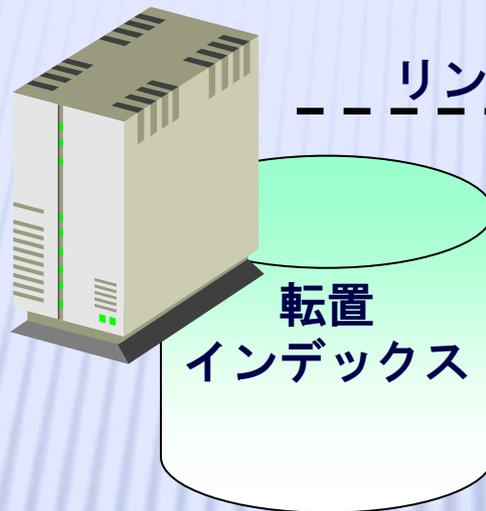
システム構成例



数式をクエリとして
論文を検索

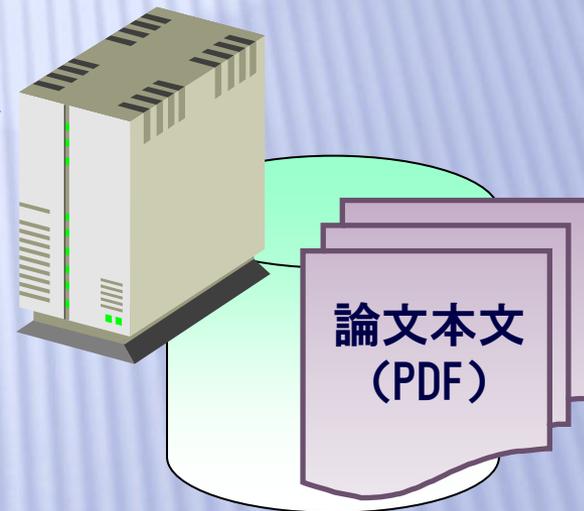


論文検索システム



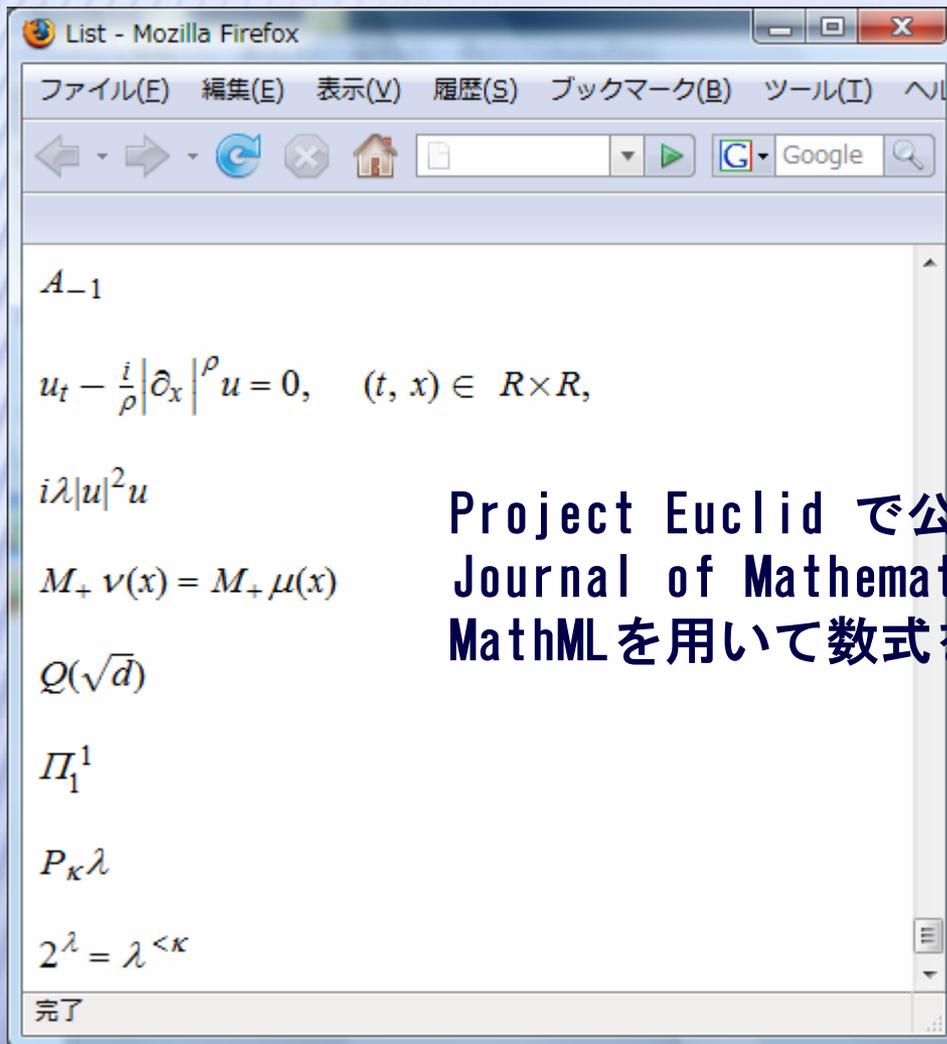
--- リンク --->

論文公開システム



※論文検索システムは論文公開システムとは別の独立したシステム
(既存の論文公開システムのシステム変更は不要)

JMSJの事例



Project Euclid で公開している日本数学会の欧文誌
Journal of Mathematical Society of Japan では、
MathMLを用いて数式を記述したアブストラクトを提供

分野ごとの特色・記号表

流体力学では、

c 比熱

ρ 密度

τ 応力

Re レイノルズ数

別の分野では、

光速度？

時定数？

電磁気学では虚数単位を j で表記

文献サービスに数式検索をプラスすると

- ・ キーワードと数式を組み合わせた検索によって正確な絞り込みが可能

→利便性の向上、省力化、効率化

- ・ 異分野での同一数式の活用事例の発見
- ・ 講師、学生などが学習対象の数式の応用事例を検索

→ユーザー層の拡大

$$\frac{\partial^2 T}{\partial x^2} + \frac{\partial^2 T}{\partial y^2} + \frac{\partial^2 T}{\partial z^2} = -\frac{\nu}{\kappa} \frac{\partial T}{\partial x} \quad \Leftrightarrow$$

この熱伝導のシミュレーションソフトを開発するとき、熱伝導以外の分野で同じような微分方程式の例を調べると良いヒントがあるかも

文献検索サービスに数式検索は必要なのか？

- ・ その分野の専門家以外の人にも利用してもらうために
(キーワードや著者名の知識が少なくても)

$$\frac{1}{1+e^{-x}} \Leftrightarrow \text{『シグモイド関数』}$$

- ・ より多くの人に論文を見つけてもらい、
知識を活用してもらうために
(社会への還元)

→ 学術コンテンツ活用、技術情報利用の入り口を広げる

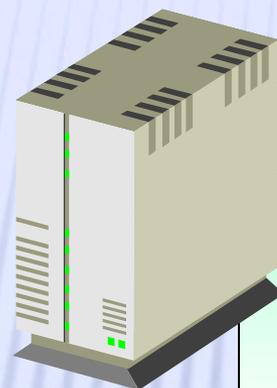
数式をメタデータとして利用

文献情報（2次情報）データベースに数式データを付加

文献検索システム



数式をキーワードとして
論文を検索



文献情報
データベース

抄録、著者名、
キーワードなど



文献（本文PDF）に含まれる重要な数式を
従来のキーワードと同様にデータベースに追加

発展

- ・ 査読支援に活用

 - 査読時に同じような数式を多く含む既存の論文を表示

- ・ 同じような数式を含む文書を集める（クラスタリング）

- ・ 同じような数式を含む文書を関連付ける（リンク生成）