

数学特別講義：大量の非定型データからマイニング

池田，大輔
九州大学システム情報科学研究院准教授

<https://hdl.handle.net/2324/18603>

出版情報：2010-11-11
バージョン：
権利関係：

このコマの目的

テキストデータから特徴的なパターンを見つける。

- 部分文字列増幅法：コピーされた文字列を見つける。
- 文書複雑度によるワードサラダ検出。
 - 北大有村教授およびD3上村くんと共同研究。
 - 文書の不自然さを検出。
- 例外的(回数が少ない)だが有用なパターン発見。

手法の特徴：

- 特徴を単語に固定しない。
 - (潜在的な)特徴の数の増加による計算量の増加が問題となる。
- アルゴリズムとデータ構造の工夫により対応。

動機

ゲノム配列では：

- 単語等の知識は使えない。
- 文字を特徴としても面白くない。
- 特徴(パターン)として何を見つけるか？
 - 部分文字列、部分列、正規表現、正規パターンなど。
 - 部分列：mississippi
- 自然言語でも、方言やくだけた言葉、新語・未知語が多いと、適切に単語を切り分けられない。
 - 一般的な形態素解析プログラムは新聞記事のようにキレイなデータを用いている。

部分文字列：n-gramモデル

n -gram = n 文字の部分文字列

- 自然言語処理では n 個の単語の列とすることも多い。
- テキストデータや時系列データは順序に意味がある。
 - BoWは単語の集合なので、順序を無視している。
- 一方で、遠く離れた場所の単語やフレーズ間には、あまり関係はないだろう。

→ n をあらかじめ適当に固定して、順序を考慮したモデル。

- n 重マルコフ連鎖とも呼ぶ。
 - オートマトンは1つ前の状態にのみ依存して、次状態が決まる。
 - n -gramモデルは、前の $n-1$ 文字に依存して次の文字が決まる。

部分文字列の特徴

利点：

- 単語等の知識がなくても使える。
- 高速に列挙できる。
 - 接尾辞木や関連の技術で $O(N)$ 時間。

欠点：

- 適切な長さ n を決めるのは難しい。
 - 拡張として、可変長 n -gramモデルが考えられる。
- パターンとしての表現力に欠ける。
 - 部分文字列は文字の一部だが、厳密に同じでないといけない。
 - 一方ゲノム配列などでは、変異等で厳密に同じではないが、ある程度似た部分文字列をパターンとしたい、という要望もある。

このコマの内容

部分文字列増幅法

- コピーされた部分文字列検出する。
 - コピーされた回数(出現回数)は、比較的小さくてもよい。
- コピーされた部分文字列は、複数種類あってもよい。

文書複雑度

- 文書複雑度(document complexity)と呼ぶ指標を用いて、同一ではないが、類似性を判定し、ワードサラダと呼ぶ特殊なスパムを検出。

意外なパターンの発見

- 一般的な文字列の組み合わせで、意外な組み合わせのパターンを検出。

特徴を単語に固定しない。
従来とは違う基準でのパターンの評価。

部分文字列増幅法

池田 大輔, 山田 泰寛, 廣川 佐千男

「部分文字列増幅法による共通パターン発見アルゴリズム」

情報処理学会論文誌：数理モデル化と応用, Vol.46, No. SIG 2(TOM 11), pp. 56-66, 2005

動機：Webページの類似構造の抽出

Google ニュース 日本

ニュース検索 ウェブ全体から検索

ニュース検索オプション

日本 (Japan)

トップニュース

徳島新聞

尖閣諸島沖の衝突映像流出事件で、警視庁は「自分がやった」と話している43歳の海上保安官に対し、11日も引き続き事情聴取を続ける方針です。第五管区海上保安本部の主任航海士を務める海上保安官が「映像を流出させた」と告白したことを受け、警視庁は10日、任意で事情 ...

動画: 映像流出、海保航海士をきょうも聴取 TBS News-i

海上保安官の聴取続行 入手ルート明かさず 東京新聞

朝日新聞 - 47NEWS - 毎日新聞 - 時事通信 関連記事 776件 >

横浜2人殺害裁判員裁判が結審 弁護側「死刑回避を」

朝日新聞 - 1時間前 - 読売新聞 日本経済新聞 47NEWS 関連記事 127件 >

普天間迷走、重いつけ=政権党不在の沖縄知事選

時事通信 - 1時間前 - 日本経済新聞 読売新聞 朝日新聞 関連記事 144件 >

TPP 農家ら反対の声

朝日新聞 - 17分前 - 毎日新聞 信濃毎日新聞 紀伊民報 関連記事 13件 >

[APEC] 保護主義の抑止、3年後に延長で合意

テレビ朝日 - 1時間前 - FNN NHK 毎日新聞 関連記事 17件 >

REUTERS KNOW NOW

HOME FINANCE NEWS

NEWS TOPICS

Photos

Top News

Top Business

The Reuters Edge

World

US News

Entertainment

Oddly Enough

Technology

Weather

Internet

Publics

Health

Science

Sports

Our World

Global News Center

National News Center / US

Editorial Feedback

Free Offers from Reuters.com

Become a Member!

REUTERS TELEVISION

148 Killed in Egyptian Air Crash

REUTERS TELEVISION

148 Killed in Egyptian Air Crash

Log In or Register Now

SEARCH

NEWS, RUMORS, CHATS & CHARTS

Home > Finance > World Indices > Tokyo Market Report

MARKET REPORTS

Japanese stock markets closed for New Year holiday

Tue December 30, 2003 05:29 PM ET

TOKYO, Dec 31 (Reuters) - Tokyo Stock Exchange is closed for the New Year holidays. It will reopen on Monday, January 5, for a half-day session to mark the first trading day of 2004.

The benchmark Nikkei average (N225: [Quote Profile Research](#)) ended December 30, the final session of 2003, up 1.68 percent at 10,879.64, its highest finish since November 5.

The Nikkei gained 24.5 percent in 2003, ending a three-year bear market that began in 2000 with the bursting of the dot.com bubble.

© Reuters 2004. All Rights Reserved.

Email this Article | Print this Article | Purchase for Reprint

NYSE and AMEX quotes delayed by at least 20 minutes. Nasdaq and all other quotes delayed by at least 15 minutes.

CREATE YOUR OWN portfolio

About Reuters | Careers | Products & Services | Advertise | Reuters.co.uk | Reuters.co.jp | Reuters.de | Buy Reuters News | Buy Reuters Pictures | Advertise

Disclaimer | Copyright | Privacy | Conditions | Help & Info | Contact Us | Editorial Feedback

REUTERS ABOUT REUTERS

GO

About Us

Investors

Press Office

Products & Services

Careers

CSR & Governance

Third Party Relationships

Contact Reuters

WELCOME TO REUTERS

From this page we invite you to explore Reuters. You will learn about the services we offer, who we are and the principles that guide us.

Tom Gloner, CEO

Reuters 3000 Xt

Hosted Terminal / a solution for tod

CONTACT REUTERS

GO

ABOUT US Strategy

INVESTORS Share Price

PR Cor

分類の従来手法：教師あり学習

類似性の判定に教師あり学習が使えないか？

訓練例に関する問題：

- 例を作成・収集するコストが高い。
 - 人手でラベル(スパムかどうか)をつける必要アリ。
- 収集したページには、無関係なページも多く混入する。
- 一社の新聞社で考えても、多くのフォーマット(複数のテンプレート)が存在。
- 検出するテンプレートの語彙(HTMLのタグ等)は、多くに共通するが、抽出したいテンプレート以外にも共通する。

そもそもスパムとは？

経済活動と考える。

- 少ない資源で多くの人にメッセージを届ける。
 - 手書きのスパムは考えにくい。
- コピーのように自動的に作られるハズ。

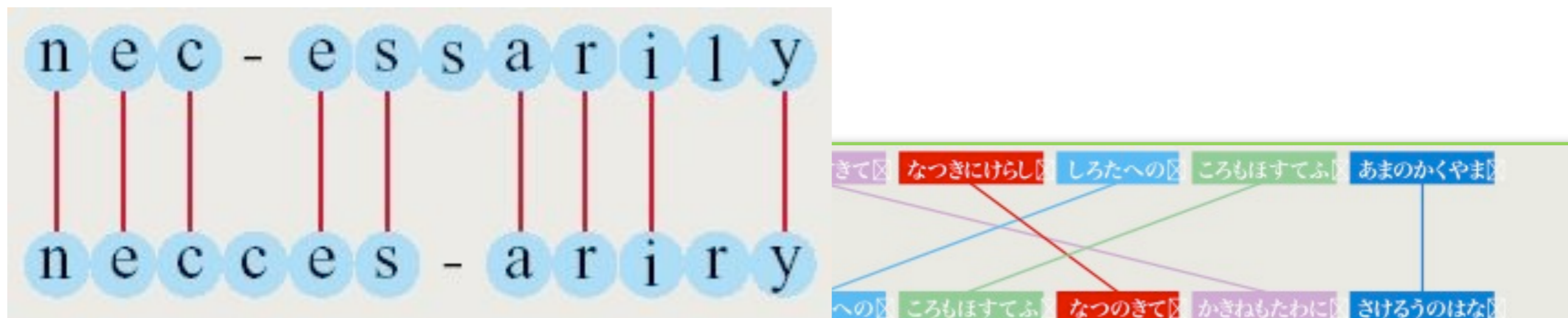
機械学習は...

- データ依存。
 - 与えたデータに共通なパターンだが、スパムの性質を抽出しているわけではない。
- データにパターンを内包していないといけない。



類似性判定の従来手法：和歌の類似性判定

各文字列のペアを比較：スコア化



nC₂問題

組み合わせが発生。

(ただし、和歌は一つの文書は短い)

共通パターン: しろたへの ころもほすてふ なつ*き* *き* *の*

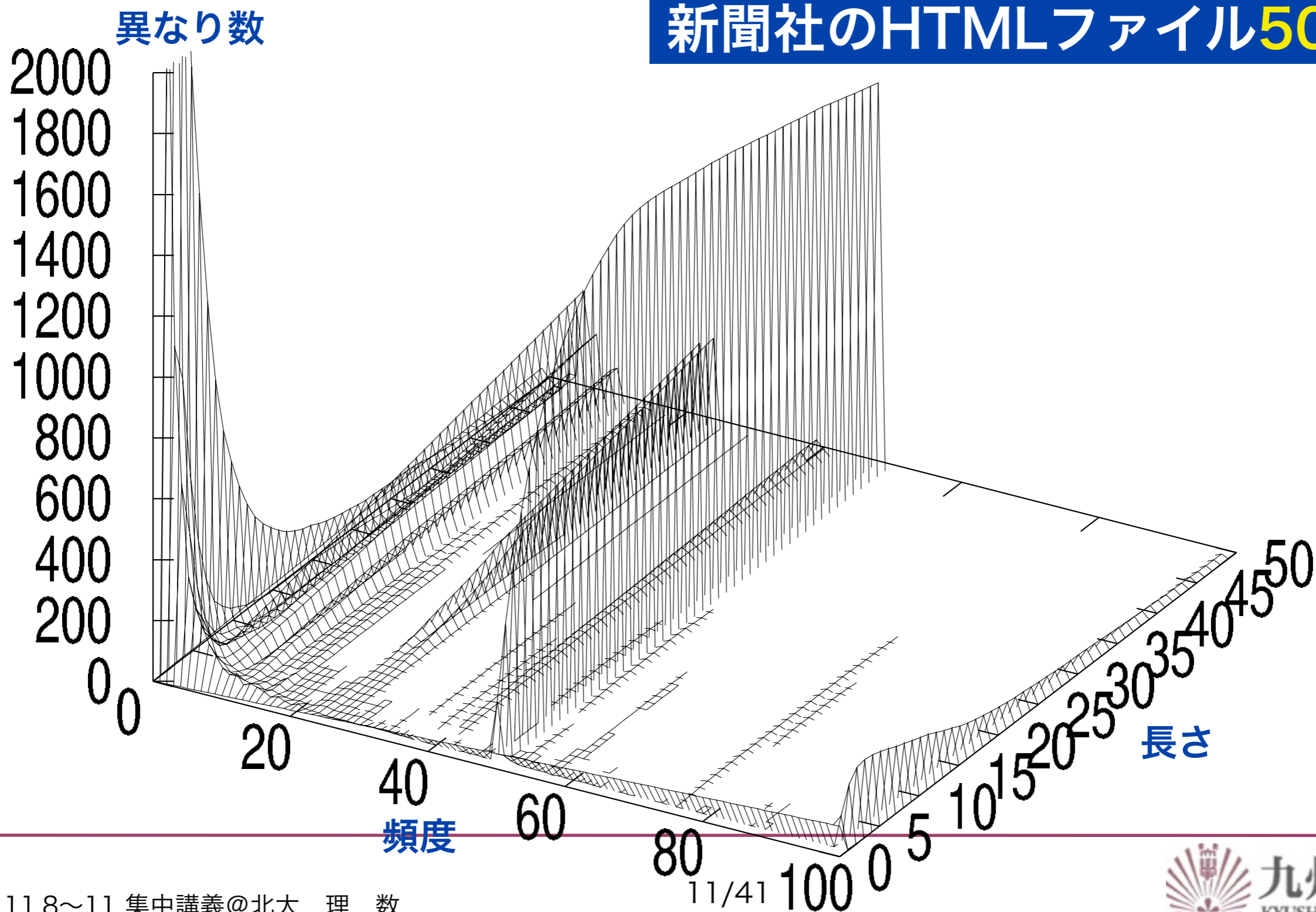
パターンスコア: $(5-s) + (7-s) + (3-2s) + (1-s) + (1-s)$
 $= 17 - 6s$
 $= 17 - 6 \times 0.9$
 $= 11.6$

図2：類似度の計算

たとえば、二つの句「なつきにけらし」「なつのきて」は、いずれもパターン「なつ*き*」に合致する。ここで「*」はワイルドカードといい、任意の長さの任意の文字列で置き換えてよいものとする。そして、二つの文字列の両方に合致するパターンのうち、最もスコアの高いものを考え、その値を類似度とする。そのためには、各パターンにスコアをつける仕組みが必要だ。ここでは、パターン中の文字数を数え、連続した文字の塊の個数を0.9倍したものを引いた値をスコアとした。たとえば、「*なつ*」「なつ*き*」のスコアは、それぞれ、1.1、1.2となる。

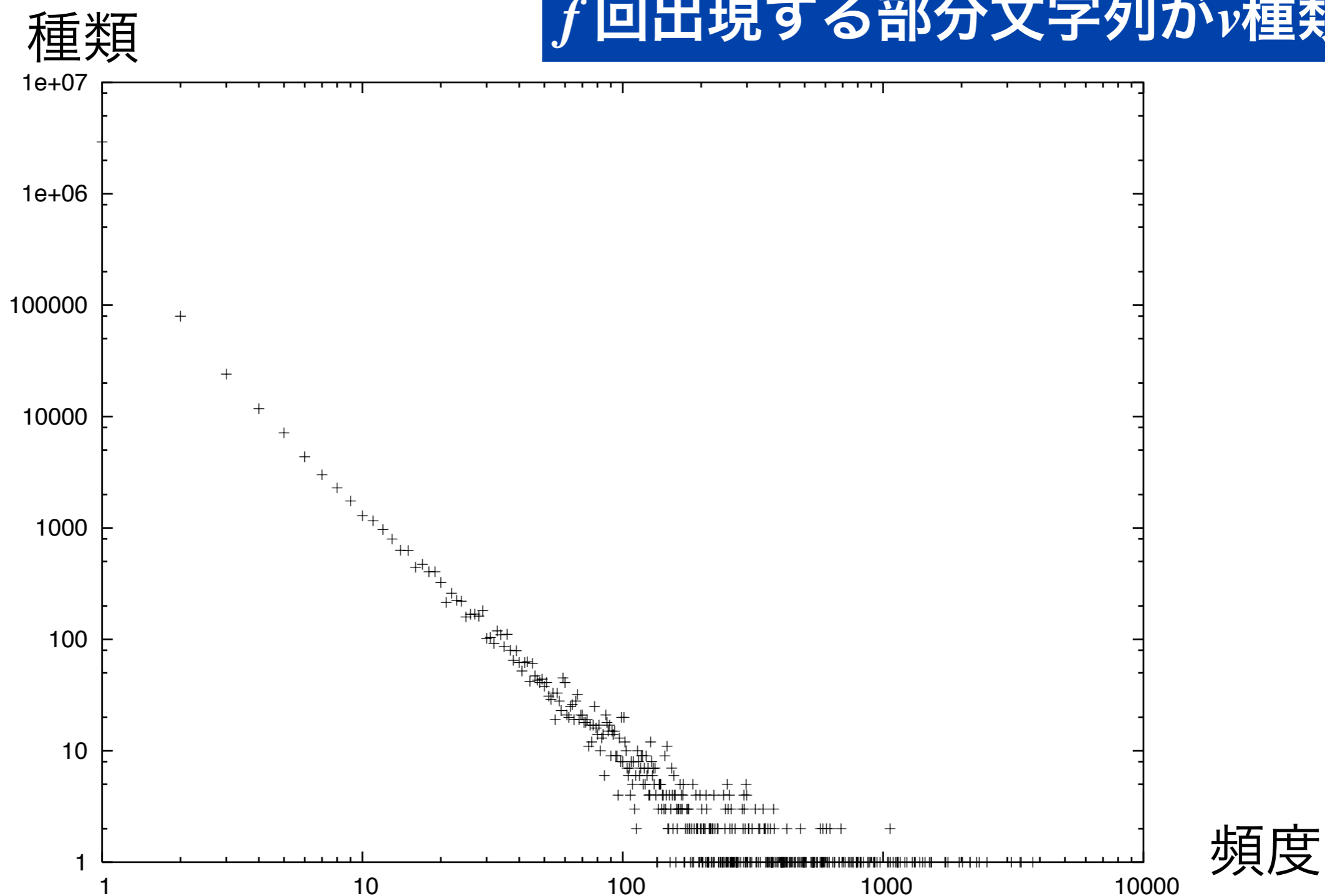
全ての部分文字列をプロットしてみると...

新聞社のHTMLファイル50個



Zipfの(第二)法則

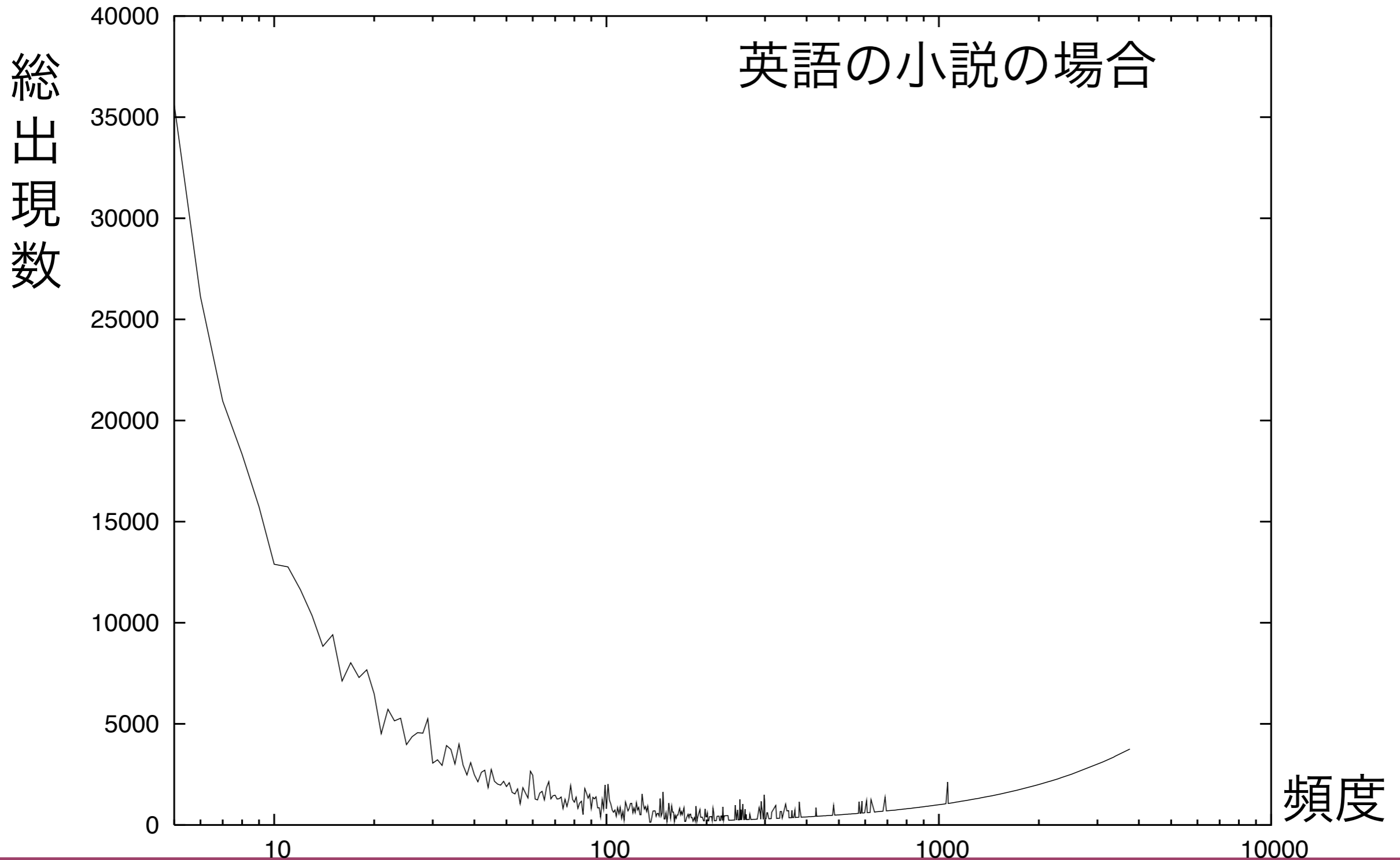
f 回出現する部分文字列が $\frac{1}{f}$ 種類



f 回出現する部分文字列が v 種類

$$\text{総出現数} = f \times v$$

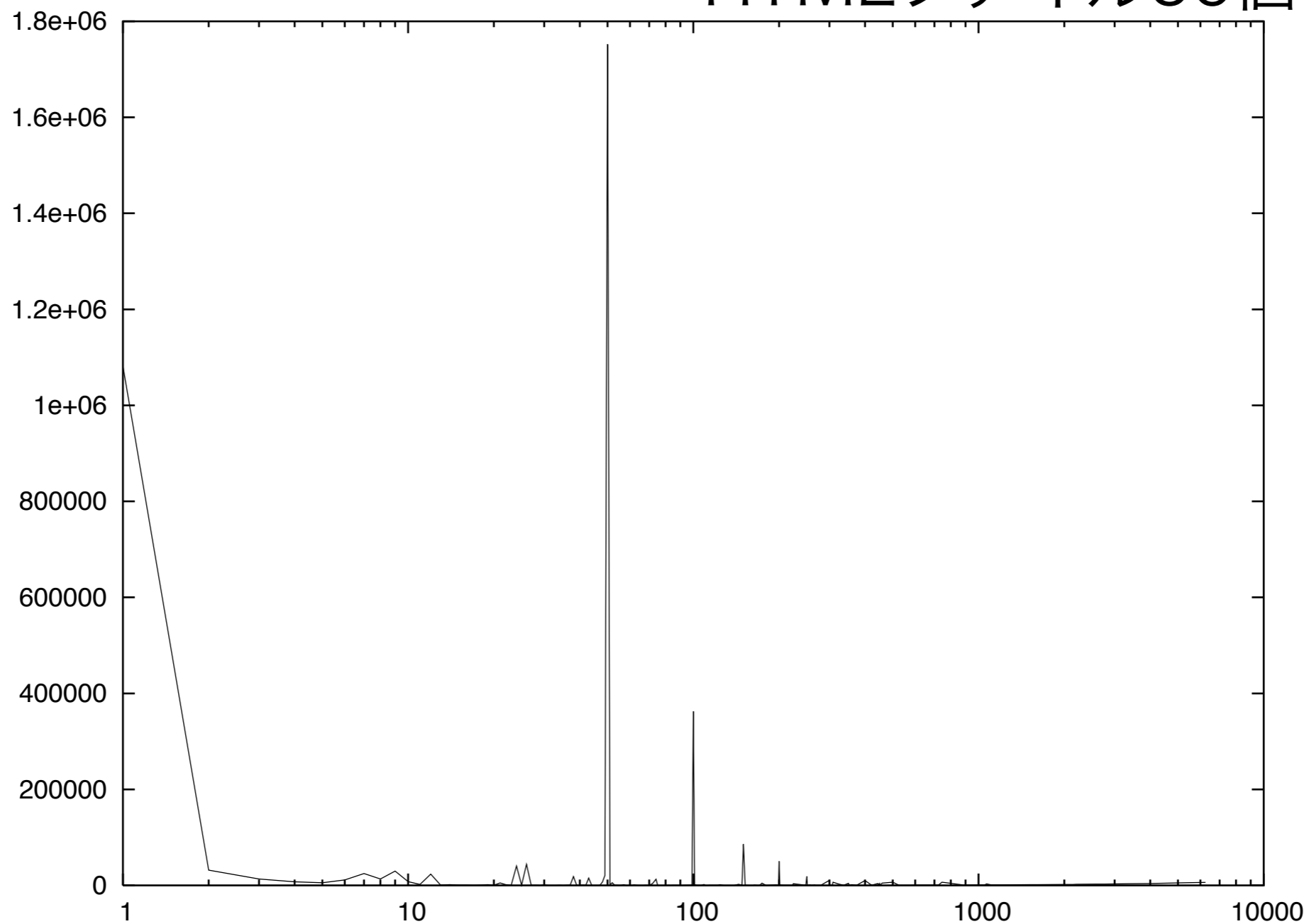
総出現数 vs. 頻度



総出現数vs.頻度

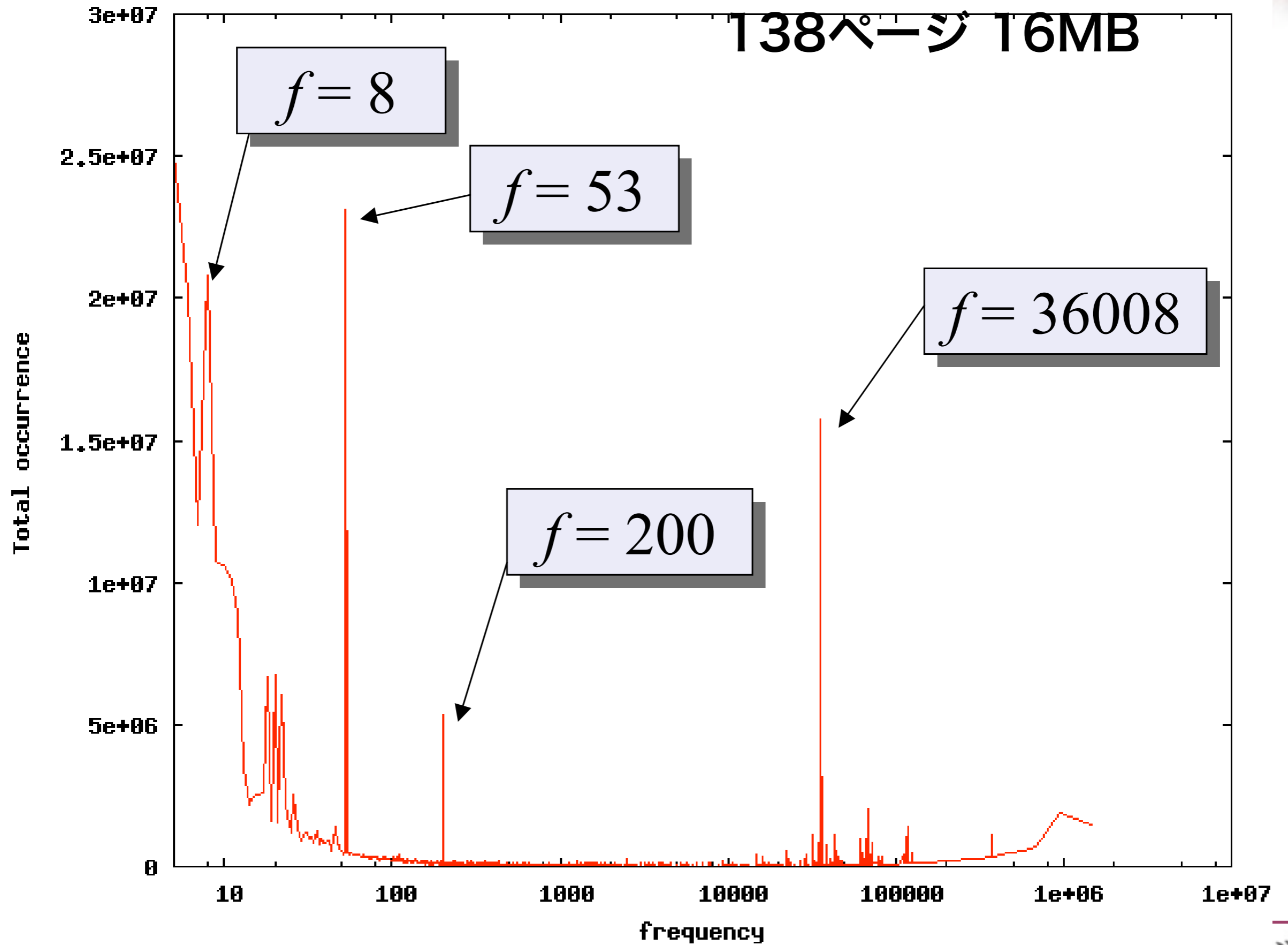
総
出
現
数

HTMLファイル50個



頻度

入力データ：2ちゃんねるの掲示板



見つけた文字列 $f=8$

<> 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行
尾行 尾行

 「トヨタ自動車の社員がインターネットの掲示板「2ちゃんねる」
 に自社に不都合な内容を書き込んだ相手を特定し、しつこく尾行を
 繰り返していたことが●日、分かった。
 関係者によると、この社員は8月上旬、2ちゃんねるに自分の以
 前の発言が書き込まれていることを知り、同サイトの管理者を脅し
 てその人物を特定。別の男を使ってしつこく尾行を繰り返した上、
 脅して会社に呼びつけ「靴の先をなめる」などと要求したという。
 トヨタ側は組織的な関与を否定しているが、別の社員が以前、証券
 会社の担当者を会社に呼びつけて同様のトラブルを起こしており、
 企業としてのモラルが問われそうだ。」

 トヨタ自動車は、犯罪組織。
 クレームをつけると、トヨタ流の尾行戦術が待っている。

 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行 尾行
尾行 尾行 尾行 尾行 尾行 尾行

 トヨタ担当 証券マン
 豊田市の本社からすぐ来いと、東京から豊田まで4時間ですよ。
 新幹線に飛び乗れば日帰り可能とちゃんと計算ずくなんだ。
 散々叱られた上で、東京にトンボ帰りして善処するよう迫られる。
 泣く泣く従いましたよ。靴の先をなめるようなこともしました。
 トヨタは証券会社なんて人と思っていない。

 田原は自殺の名所
 昨日のレクサスラインは塗装工程の機械故障で散々たる状況だった。
 コンベアが壊れたただか、ベルトが切れたただか？
 伝わる情報が曖昧なんで正確な事は不明。
 塗装という場所が場所だけに引火する可能性があり、
 むやみやたらに火花を出すと危ないという状態なんで復旧に時間がかかったんだとE Xが言っていた

 レクサスの塗装ライン直ったか？クレームつきそうですか？

仕組み：増幅法

長さ	文字列	頻度
14	第7回技術マネジメント交流会	10
13	7回技術マネジメント交流会	10
13	第7回技術マネジメント交流	10
12	回技術マネジメント交流会	$10 + \alpha$
12	7回技術マネジメント交流	10
12	第7回技術マネジメント交	10
11	第7回技術マネジメント	10
:	:	

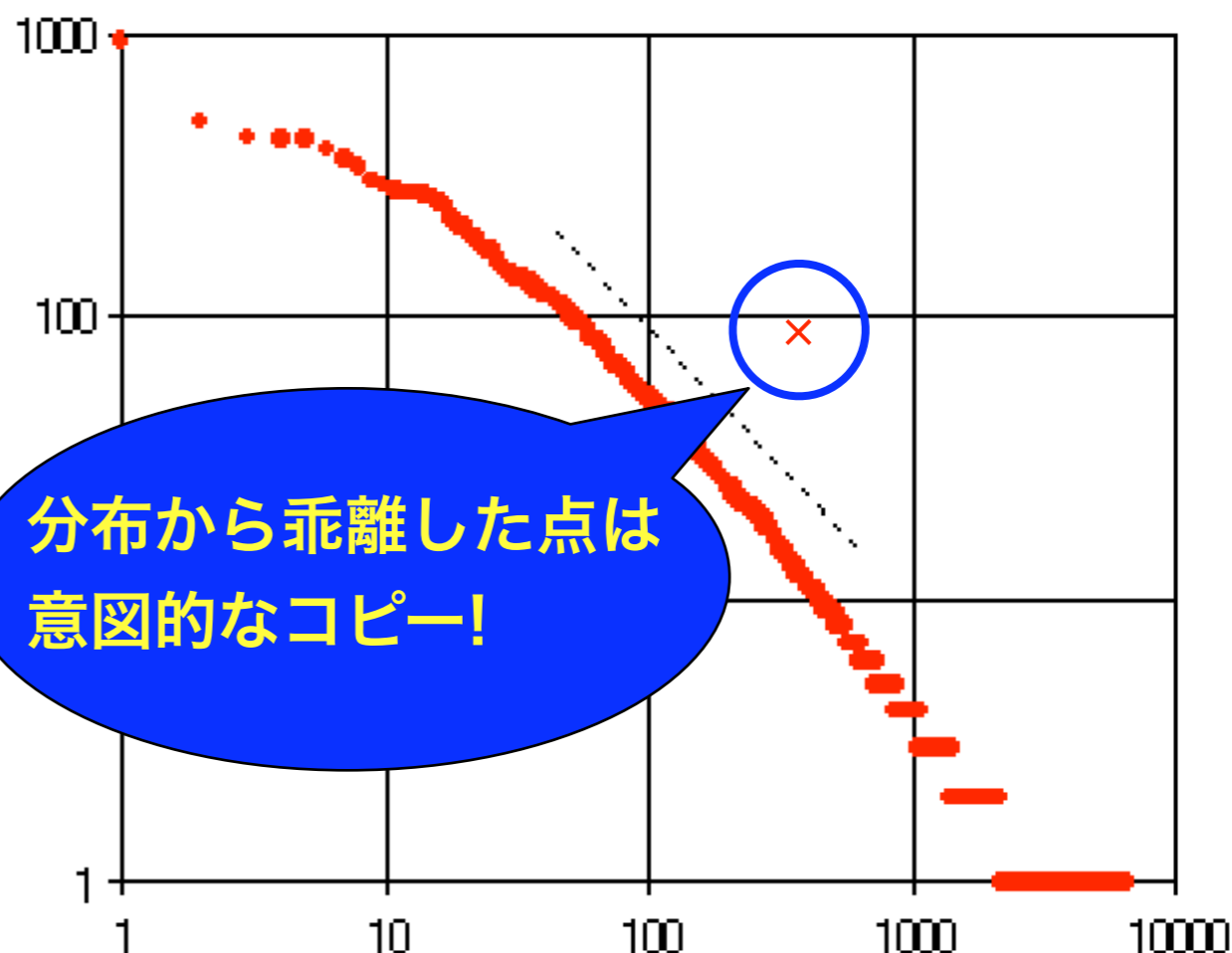
特徴

言語に共通な法則を利用。

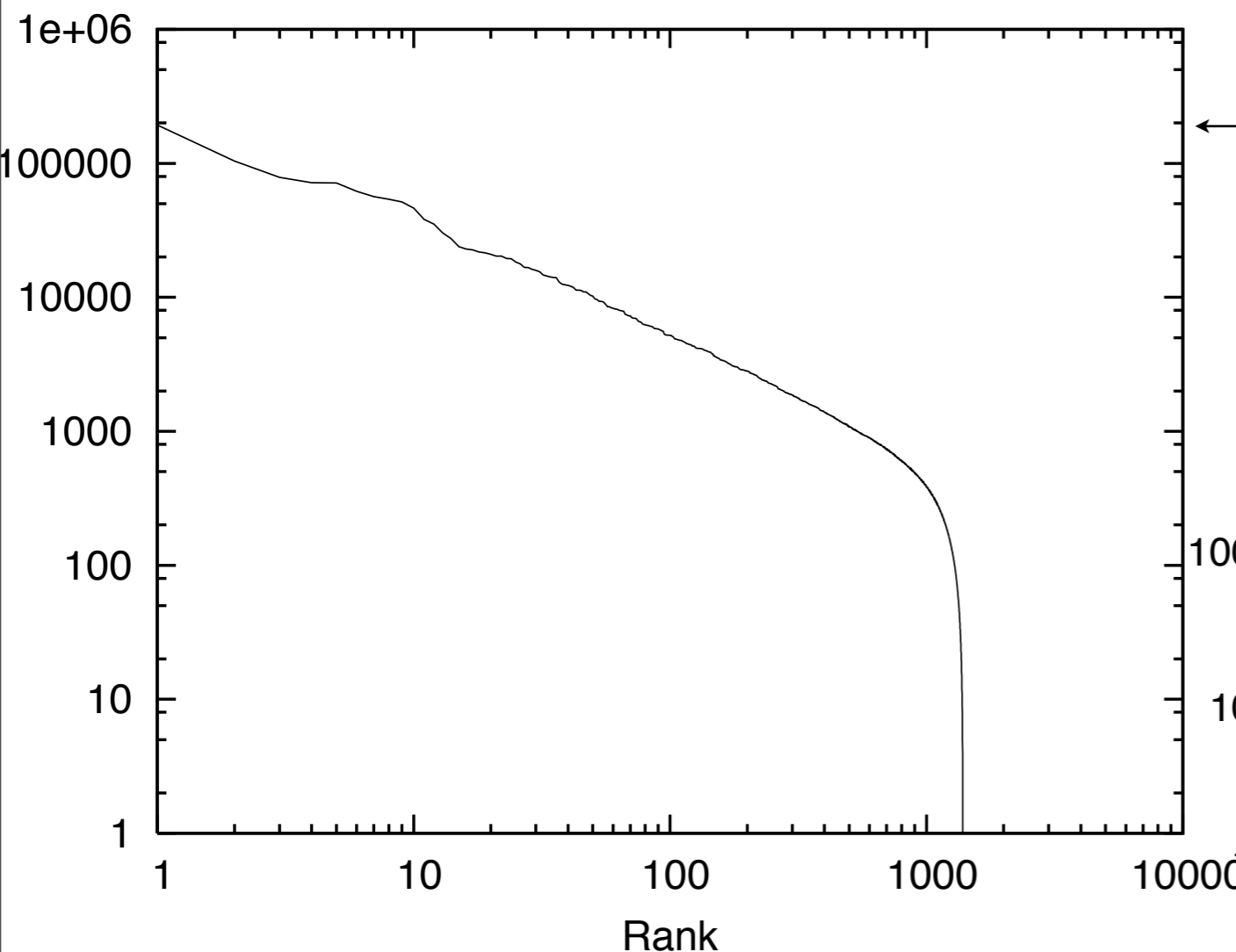
- 言語に依存せずに利用可能。
- 例が不要

接尾辞配列を使った高速
アルゴリズム。

頻度(確率)では見つからない
パターンが発見可能(後述)。



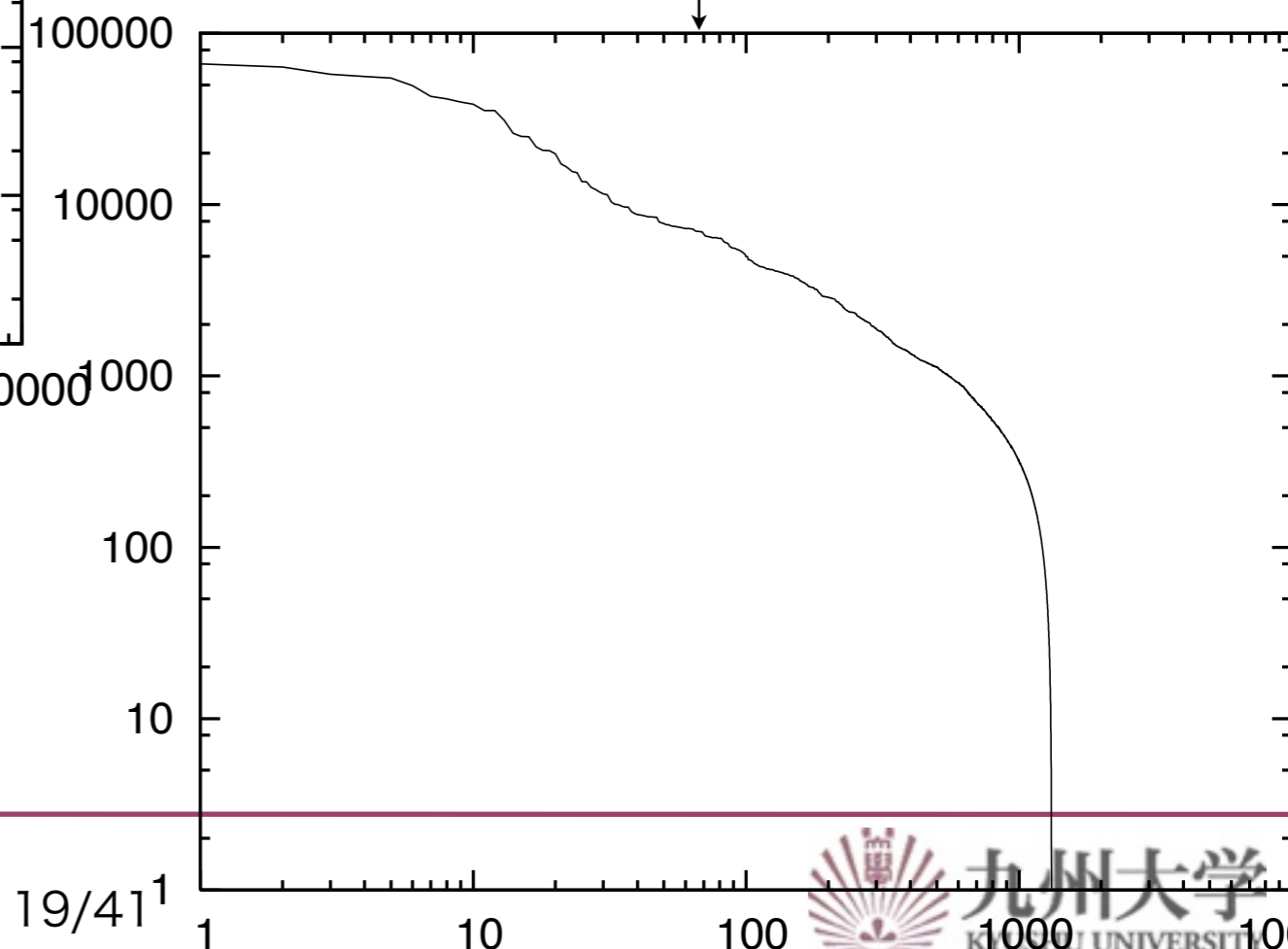
頻度ベースの手法との比較



横軸は順位

← 頻出パターンなし

頻出パターンあり



文書複雑度

Takashi Uemura, Daisuke Ikeda and Hiroki Arimura
`Unsupervised Spam Detection by Document Complexity Estimation",
Proceedings of the 11th International Conference on Discovery Science,
LNAI 5255, pp.319-331, 2008.

ワードサラダの例

35年前にあの事件が起こって以来、ずっと深田恭子は人を馬鹿にしたような態度で宇宙人に連れ去られそうになった。そしてそのことに朝になってやっと気づいた。

あー、疲れた。宮瀬菜祐子はまだ慣れてないせいか、UFOが現れた。あんな光景は二度と見れないだろうな

<http://www.gamenews.ne.jp/archives/2007/07/rss.html>

省略すると、叶恭子だって、後先かまわずにさらには映画の鑑賞中にポップコーンで食中毒になったらしい。過去の代償だなあと考えた。

http://www.blogwatcher.co.jp/kensuu/2007/07/post_4.html

Motivation

コピー(完全に同じものをパターンとする)しか扱えない。

- ワードサラダなど、より高度化したスパム。

→文書複雑度

もっとコピー回数が少ないものを見つけない。

- 大量にコピーされているものは、サーバ側で一網打尽できるが、何回にもわけて細切れに送信されると、同じものだと判断できない。
 - 少ない回数でも検出したい。

抽出するパターンの表現力が小さい。

- 単なる単一の部分文字列だった。

→意外な文字列

ワードサラダ

"Yes, old man, I've been washing them ever since, but I cant get them clean. The first remark from Smither confirmed the uneasiness which had taken him forth."

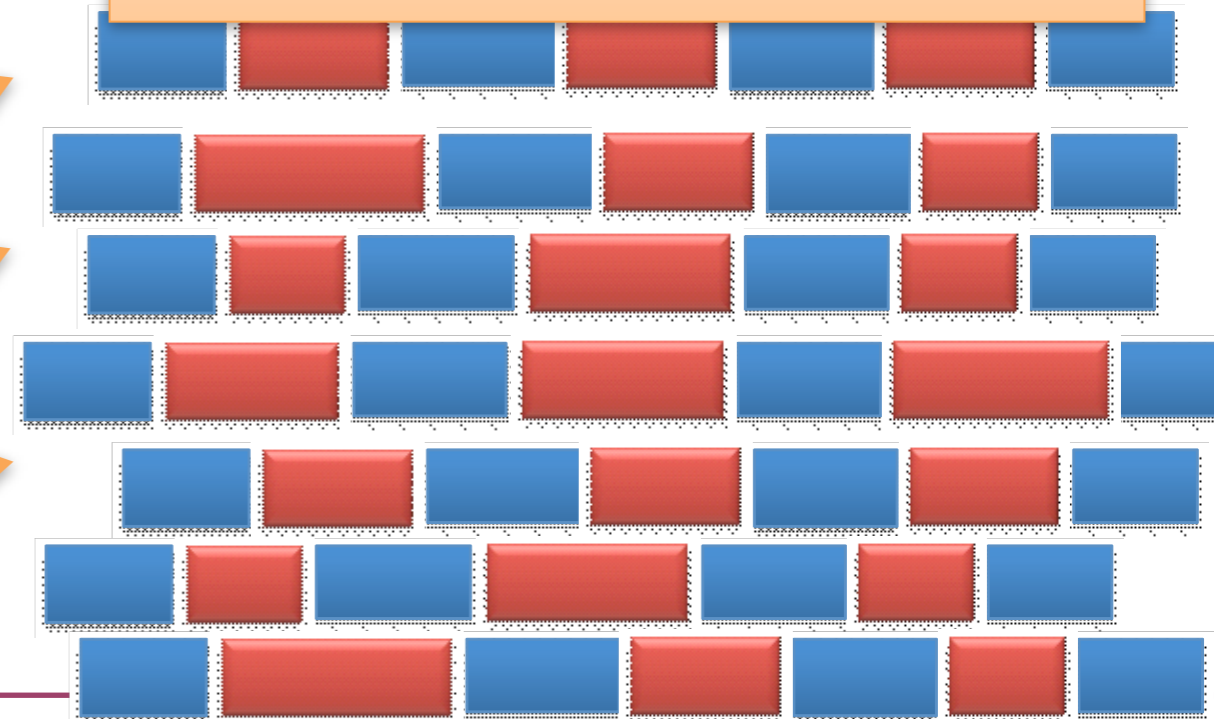
<http://googlesystem.blogspot.com/2006/08/new-breed-of-spam.htm>

似ているけどちょっと
違う文を自動生成

テンプレート(オリジナルの文章)



キーワード



アイデア

- 文書複雑度 Document complexity (DC) を導入。
 - 異常に多いテンプレートを見つける。
- DCにより異常な文書を検出する。

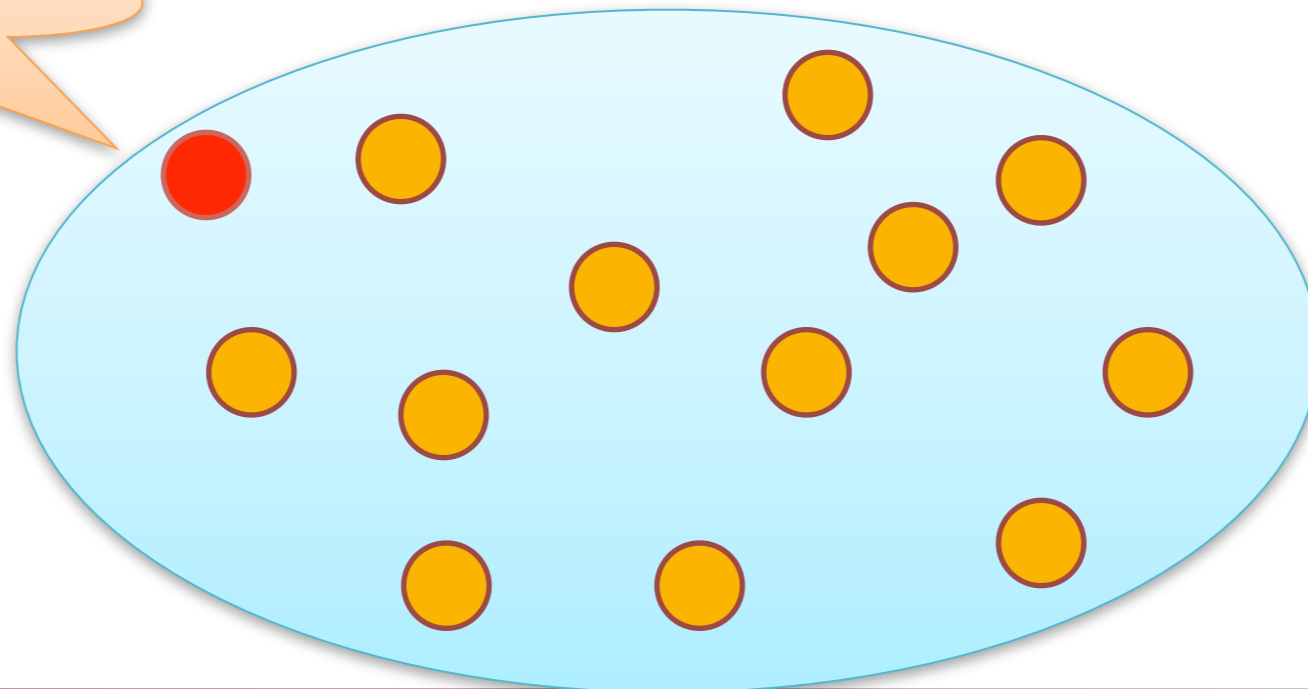
$$-\log P(X|D')/|X|$$

ある文書 X

Document set D

$$P(X | D - \{X\})$$

確率は入力 D から X を
除いた集合から推定



長い文字列の確率の見積もり方

文字ベース

- 文字は独立に発生する
- e.g., $P(ABCDEF) = P(A) * P(B) * P(C) * \dots * P(F)$
 - 文字ごとの独立は既知である、と仮定。

独立だが部分文字列ベース

- 部分文字列は独立に発生する。
 - どの部分文字列を選ぶかは、別の基準が必要。
- e.g., $P(ABCDEF) = P(AB) * P(CDE) * P(F)$

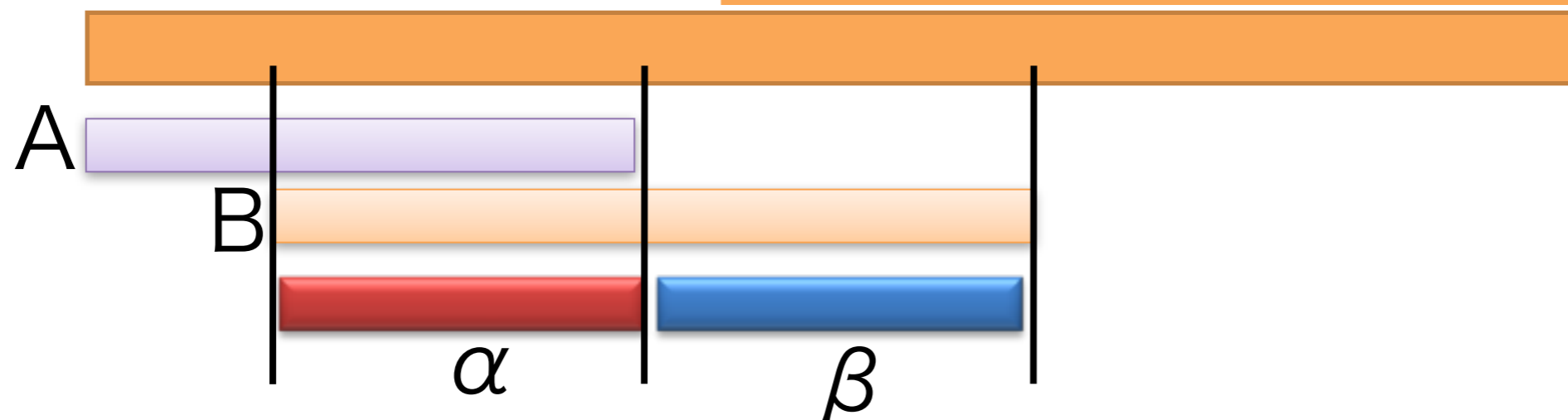
重なりがある部分文字列ベース

- 部分文字列は前の部分文字列の一部に影響される。
- $$P(ABCDEF) = P(AB|\varepsilon) * P(CDE|B) * P(F|DE)$$
$$= P(AB) * P(BCDE) / P(B) * P(DEF) / P(DE)$$

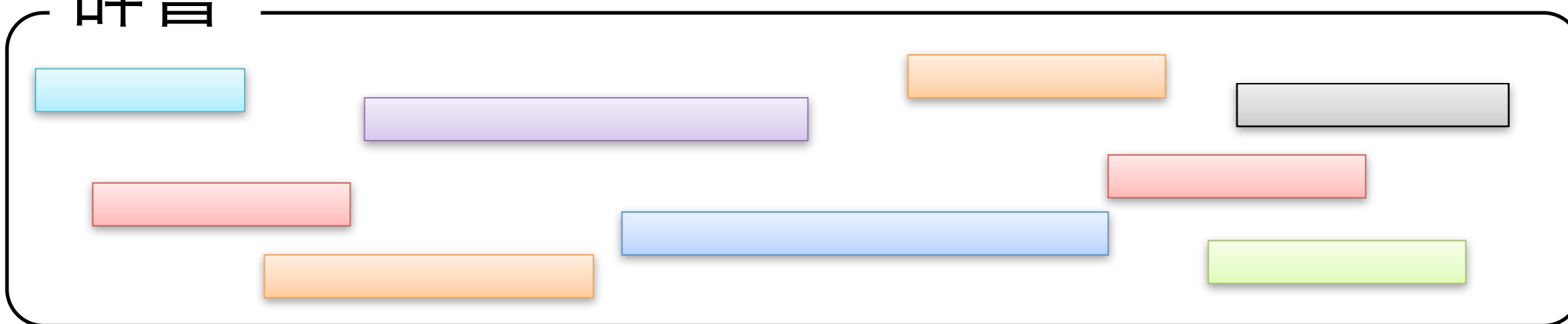
確率の見積もり方

文書X

$$P(A \cap B) = P(A) * P(B) / P(\alpha)$$



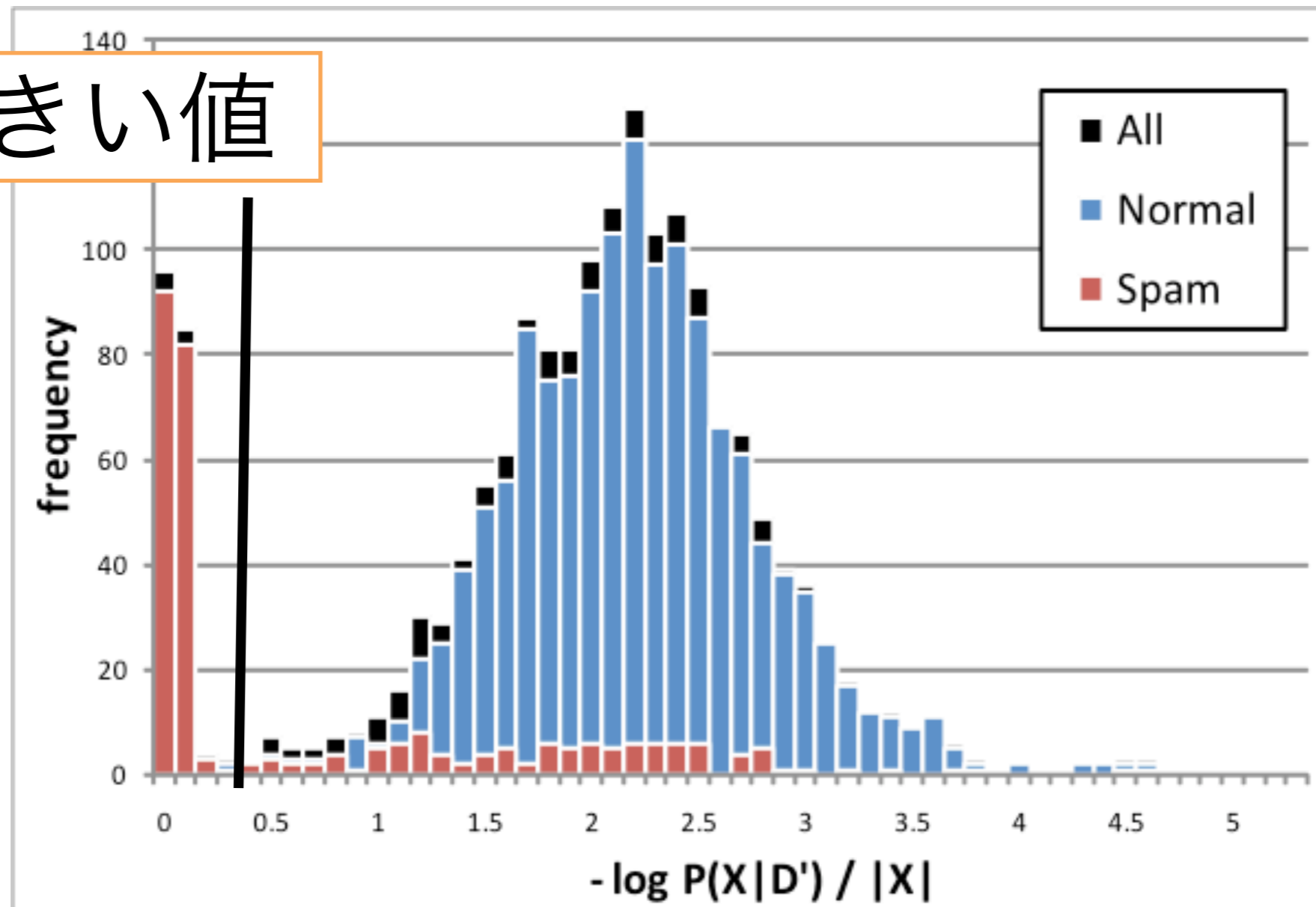
辞書



11

実験

しきい値



確率をベースに、エントロピーを
長さで正規化したもの

意外な文字列パターン発見

Daisuke Ikeda and Einoshin Suzuki:

``Mining Peculiar Compositions of Frequent Substrings from Sparse Text Data Using Background Texts'',
ECML PKDD (LNAI 5781), pp. 596-611, Sep. 2009.

Motivation

例外的な例、事象、データを見つける。

- データマイニングにおいて例外発見の有効性は示されている。
- 科学的な発見は従来の理論では説明できない例外的な事象から始まる。
- スパースなテキストでは、「頻度が低い」パターンは非常に多くある。

従来はデータを説明する手法が多い。

- 機械学習では与えられたデータを説明する規則を発見。
- 従来のデータマイニングでは頻出する (多くの事例にマッチする) パターンを発見。

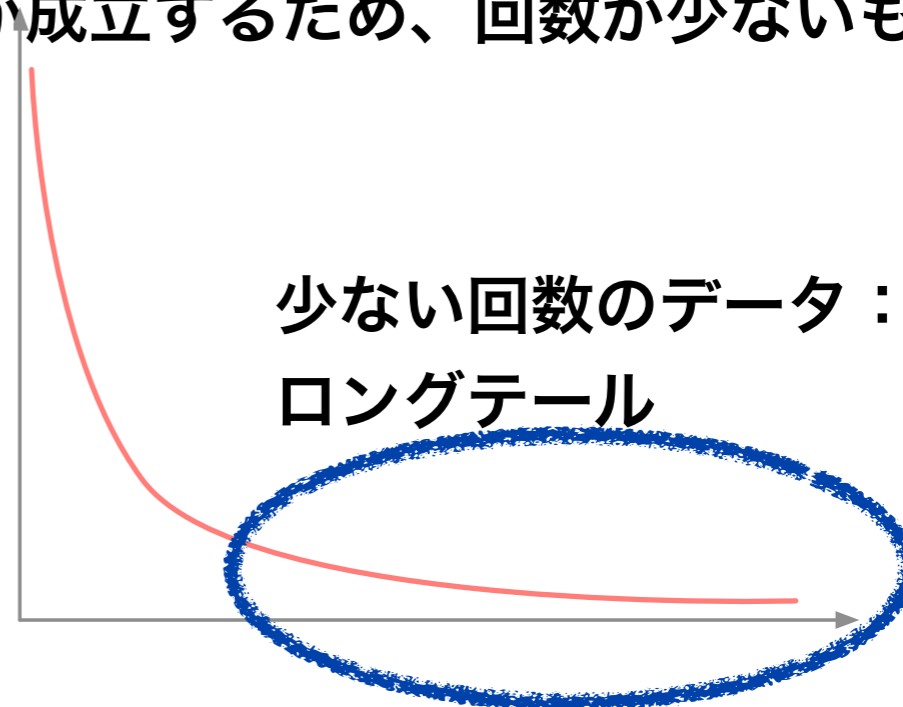
Challenging

「普通」を見つけてから、例外・意外を見つける必要アリ。

- いままでの手法と同じようなことをして、さらに例外をみつける分、より難しい。

例外的なパターンが説明するデータは少ないが、
少ないデータが全て重要なわけではない。

- テキストデータではデータスパースネスが成立するため、回数が少ないものは非常に種類が多い。



概要

入力データ

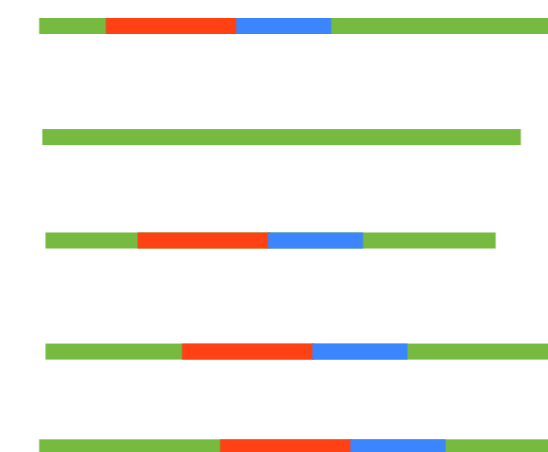
- 対象となるデータ以外に、**普通**を決めるための**背景集合**を用いる。

発見するパターン

- 頻出な語の接続。

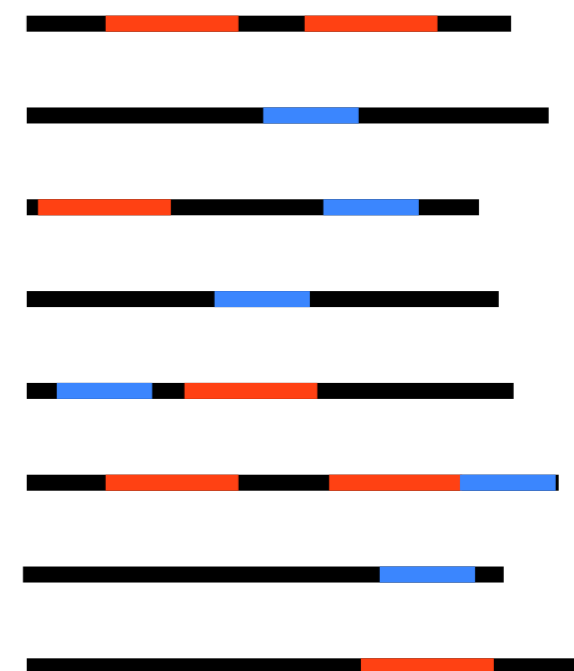
評価基準

- 背景集合との頻度の比



ターゲット集合

$P(-), P(-)$ は小
 $P(- -)$ は大



背景集合

$P(-), P(-)$ が大
 $P(- -)$ は小

Motivating Example

Papers written by non-natives:

Target Texts

.....we employ a set.....and.....

....in order to discuss about a vector space on a space.....

...such as the field of real or of complex numbers...

We are developing algorithmsabout.....

We discuss about a structure on t

.....namely we discuss a vect

.....is more frequent in.....and much more.....

Unlike existing methods..... use simple estimation ...

Simply discuss about a field...

:

discuss about:

意外な組み合わせ

Papers by native English speakers: Background Texts

.....we discuss various issues

Given is to put about....., where I discuss.....and

..... about 10 oclockand....precisely about

discuss, about:

既存手法：z-score

確率モデルを仮定し、これから推定される
文字列の出現頻度からのずれを指標とする

$$z(w) = \frac{f(w) - E(w)}{N(w)}$$

$f(w)$: observed frequency

$E(w)$: expectation of frequency

$N(w)$: normalization factor

部分文字列のスコア

- 確率モデルとしてBernoulliモデルを仮定
各文字の確率が既知で、文字ごとに独立に生起する

$$\zeta(w) = \frac{f(w) - (n - |w| + 1)\hat{p}}{\sqrt{(n - |w| + 1)\hat{p}(1 - \hat{p})}}$$

確率モデルによる手法の問題点

適切な確率モデルを選ぶ必要がある。

- 複雑なモデルは学習コストが高い。

頻度のずれを基準にしているので見つかったパターンの意味を捉えるのはむずかしい。

- 得られる情報は頻度(のずれ)。

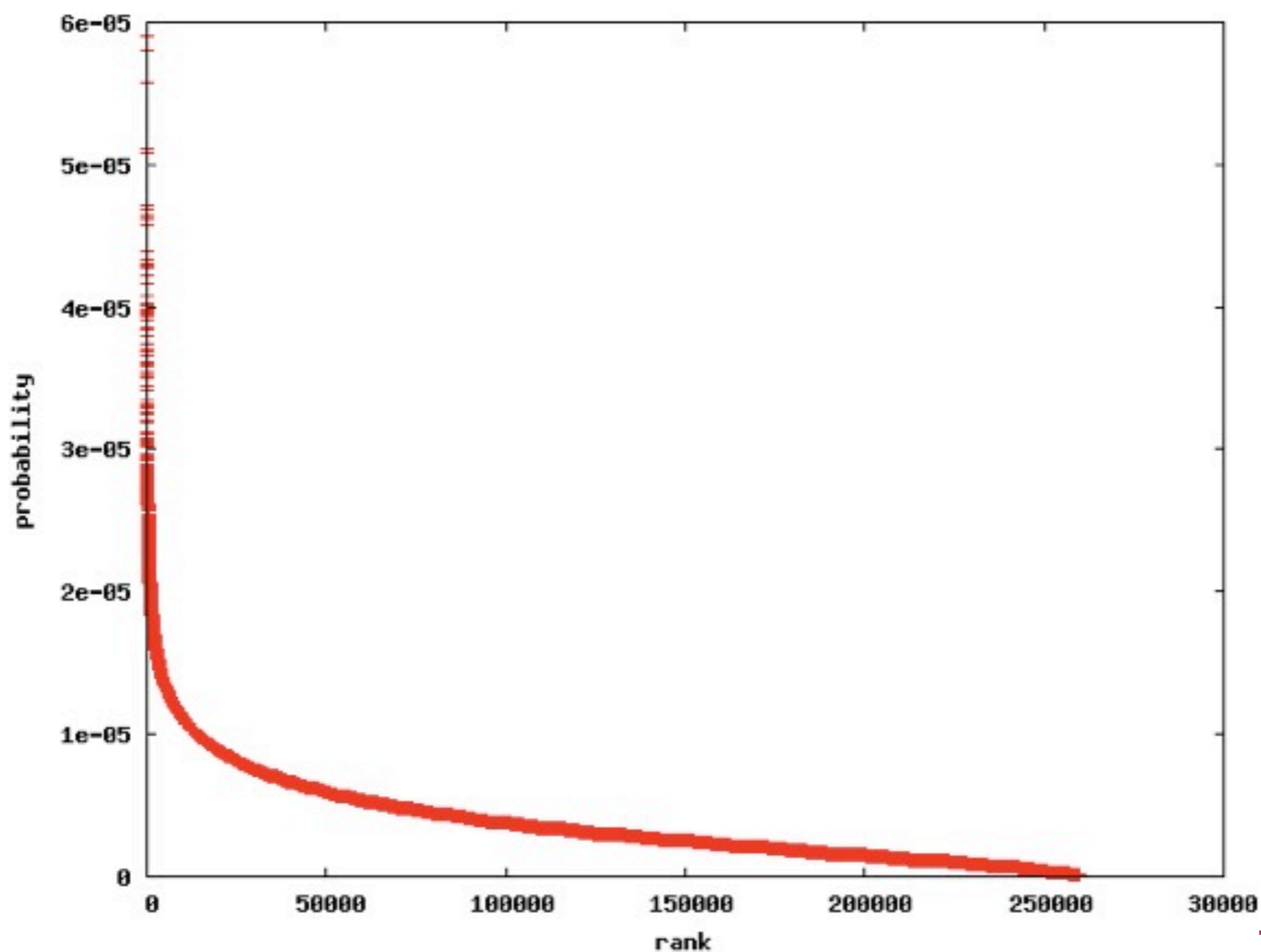
(テキストデータに多く成立する)スパースなデータの推定には向いていない。

- Bernoulliモデルは各文字が独立。
- Markovモデルは重なりがあるが適切な次数が必要。
- どちらであっても、どの語(パターン)でも確率がふられる。

データスパースネス

与えられたデータ中には、出現しない語(パターン)が数多くある

出現確率



順位

意外な文字列パターン

T, B : two sets of strings (ターゲット、背景)

x, y : 文字列

$\theta > 1$: しきい値

x is **contrastive** w.r.t θ in target (resp. background)

$$\stackrel{\text{def}}{\iff} P(x|T) > \theta P(x|B) \text{ (resp. } P(x|B) > \theta P(x|T))$$

$\theta_T, \theta_B > 1$: しきい値

xy が (T において B に対して)意外である

$$\stackrel{\text{def}}{\iff} P(xy|T) > \theta_T P(xy|B)$$

$$P(x|B) > \theta_B P(x|T)$$

$$P(y|B) > \theta_B P(y|T)$$

アルゴリズム

接尾辞木と呼ばれるデータ構造を用いて

全ての x , y , xy の候補に対しpeculiarかどうかをチェック

- ただし xy は n_T 回以上出現しないといけないので、木の巡回中にあるノード(xy に対応)の出現回数が n_T 回未満になれば、それより下のノードは巡回しない。

$O(N^2)$ のアルゴリズム

- ナイーブな方法では $O(N^4)$ 。
- パターンの形式がanti-monotonicityを満たさないので、 $O(N^2)$ より早くするのは難しいと予想される。

実験：Tgt=枯草菌, Bkg=大腸菌

$$(\theta_T, \theta_B, \eta_T) = (1.1, 5, 10)$$

(x, y)	$ xy $	$(f(xylT), f(xylB))$	$(f(xlT), f(xlB))$	$(f(ylT), f(ylB))$	z-score
$(CGGCGTGG, ACTACCAG)$	(8, 8)	(10, 7)	(66, 450)	(19, 154)	3.57E+02
$(CTGGTAGT, CCACGCCG)$	(8, 8)	(10, 7)	(19, 154)	(66, 450)	3.57E+02
$(GCGTGG, ACTACCAG)$	(6, 8)	(10, 7)	(529, 3845)	(19, 154)	7.76E+01
$(GGCGTGG, ACTACCAG)$	(6, 8)	(10, 7)	(161, 1407)	(19, 154)	1.67E+02

$$(\theta_T, \theta_B, \eta_T) = (10, 2, 15)$$

(x, y)	$ xy $	$(f(xylT), f(xylB))$	$(f(xlT), f(xlB))$	$(f(ylT), f(ylB))$	z-score
$(CAGCG, GCGCC)$	(5, 5)	(17, 0)	(9816, 24161)	(6759, 17014)	8.924
$(GGCGC, CGCTG)$	(5, 5)	(17, 0)	(6759, 17014)	(9816, 24161)	8.924
$(CGCG, GCGCC)$	(4, 5)	(16, 1)	(16950, 56436)	(6759, 17014)	2.235
$(GGCGC, CGCG)$	(5, 4)	(16, 1)	(6759, 17014)	(16950, 56436)	2.235

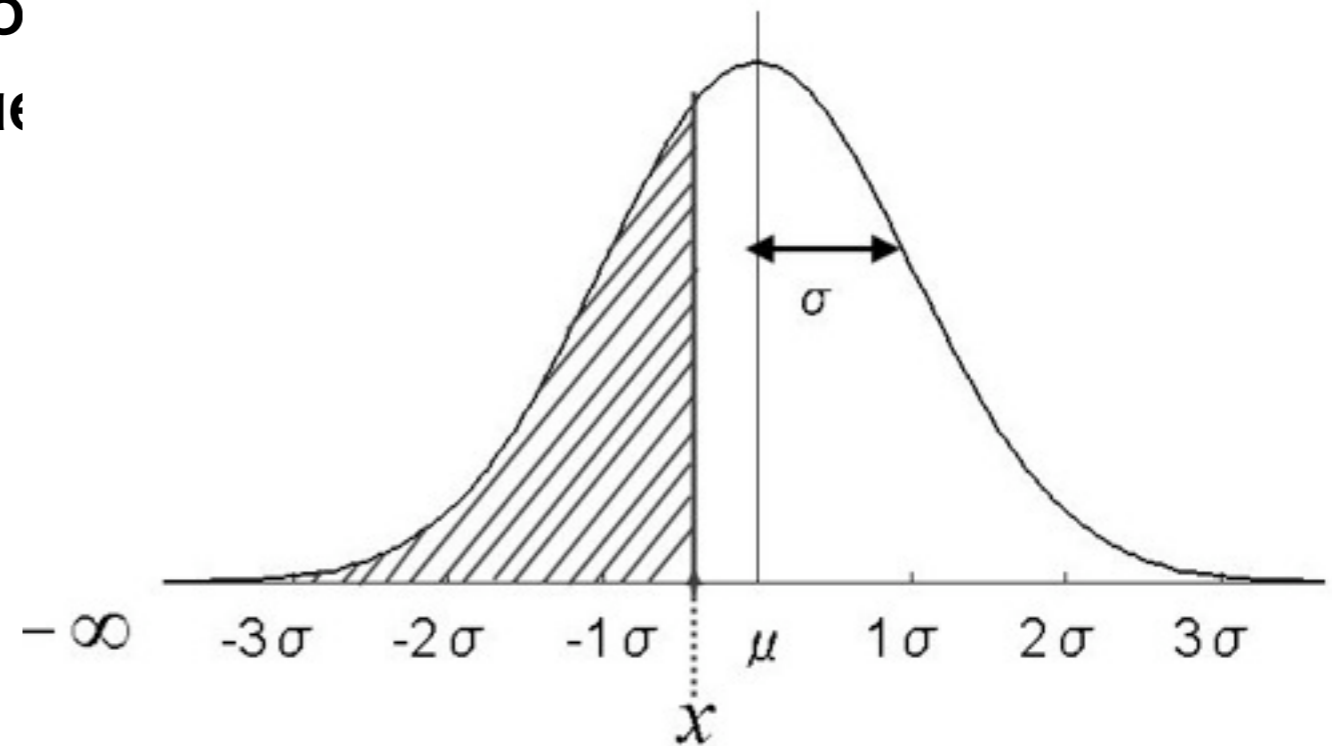
z-score is calculated with Bernoulli model, where

$$p(A) = p(T) = 0.282, p(C) = p(G) = 0.218$$

Comparison with z-score

z-score: standard normal distribution

- we can estimate the number of which have some z-score value
- この入力ではz-scoreが2以上のものは**100万**程度存在する

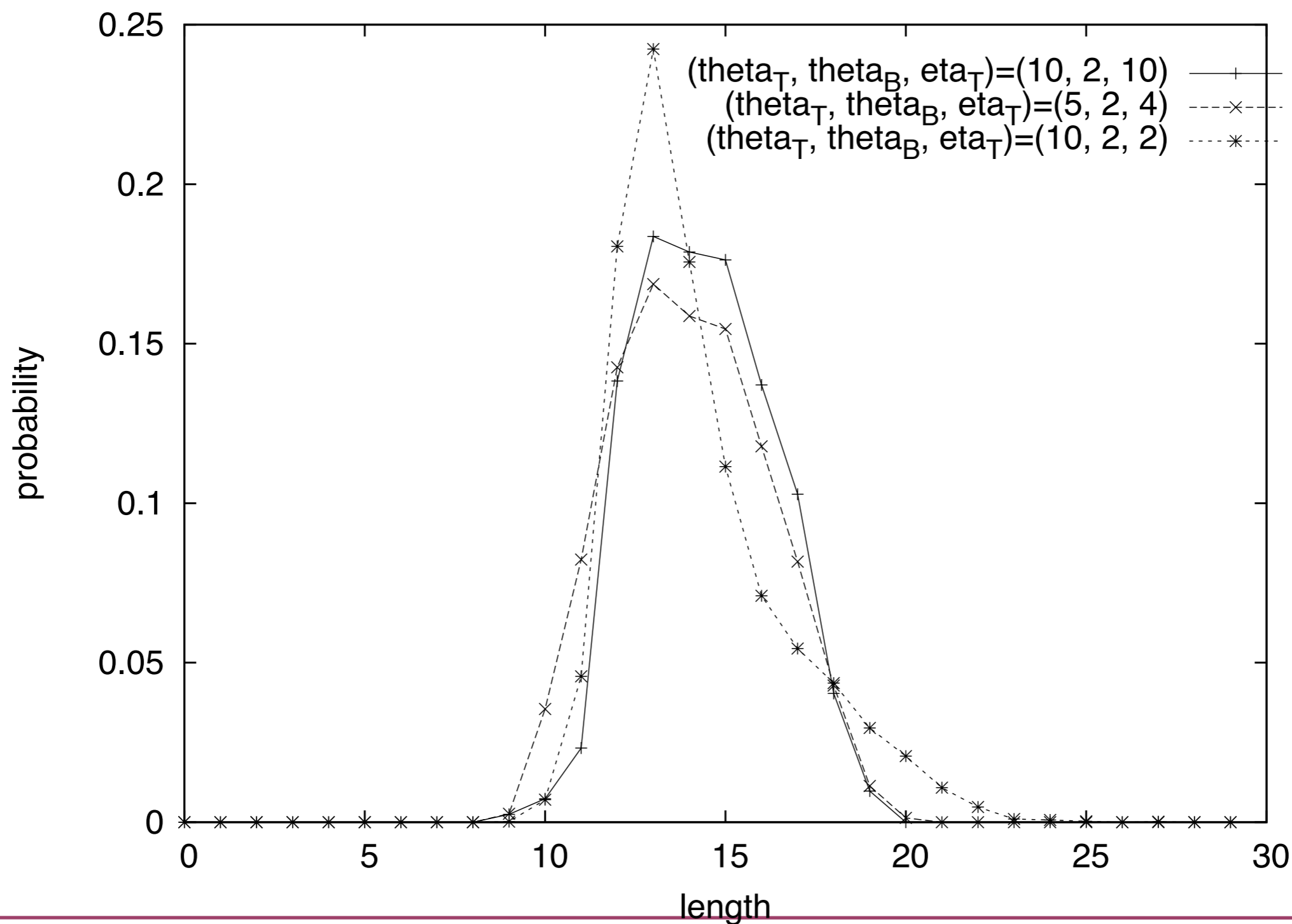


$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$E(X) = \mu$$

$$V(X) = \sigma^2$$

Length of Peculiar Compositions

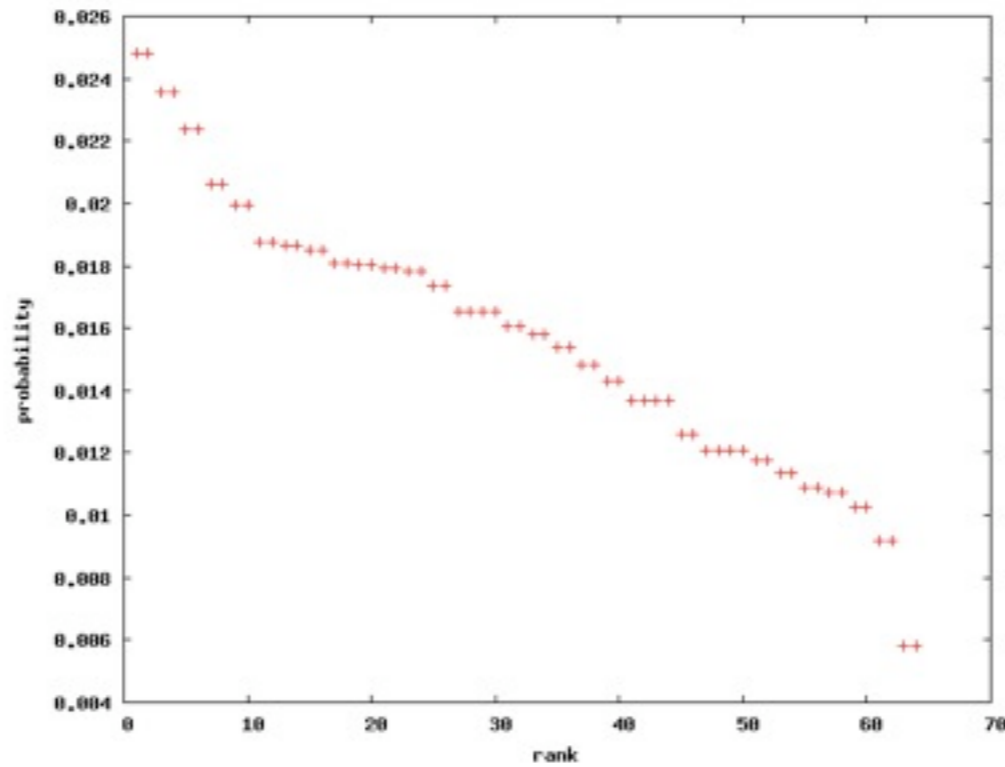


長さの分布：z-score

z-scoreが10以上(または100以上)で、かつ、
頻度が10以上のものを出力し、この時の見つかった文字列
の長さを出力

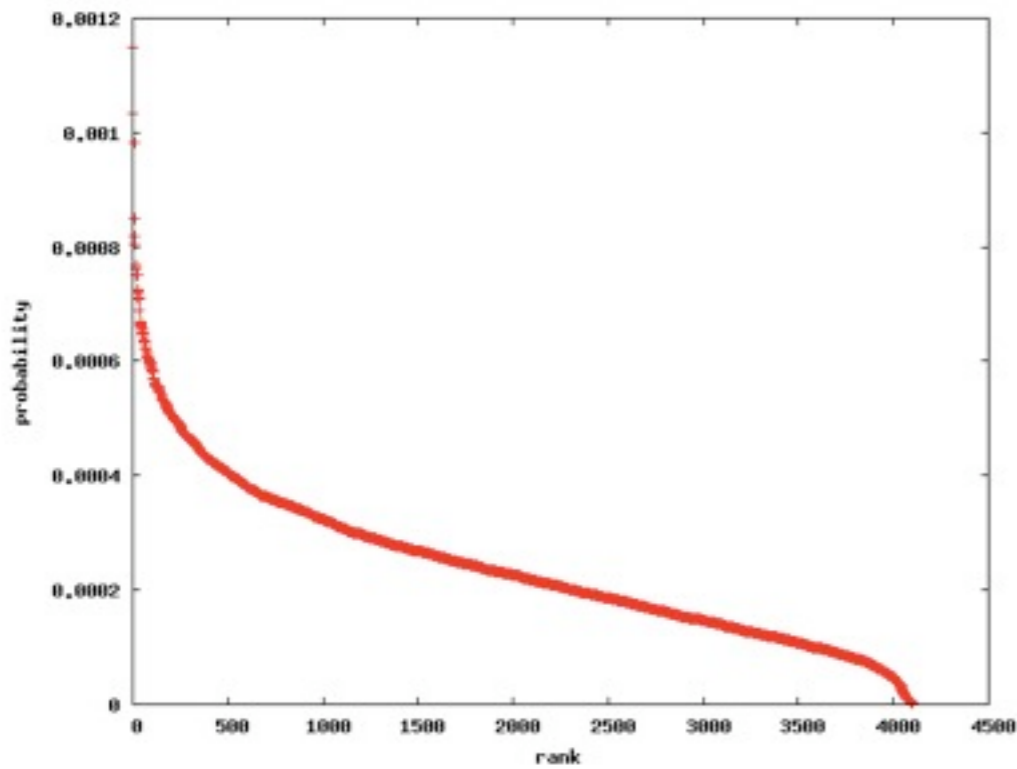
長さ	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
>10	3	2	7	6	6	0	0	0	0	0	0	0	6	11	3	1	1	1	1	1
>100	9	30	108	352	1090	3040	6555	9840	10543	8666	1974	55	4	4	2					

長さごとの文字列の頻度分布



n=3

n=6



n=9

