



Title	スケッチソート法による全点間類似度検索
Author(s)	田部井, 靖生
Citation	2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p.117-124.
Issue Date	2011-06
Doc URL	http://hdl.handle.net/2115/48469
Type	conference presentation
Note	ERATO 세미나2010 : No.16. 2010年8月26日
File Information	16_all.pdf



[Instructions for use](#)

ERATO セミナ 2010 - No. 16

スケッチソート法による全点間類似度検索

田部井靖生

JST ERATO 湊離散構造処理系プロジェクト

2010/8/26

概要

近年、画像やテキストなどのウェブデータ、化合物やタンパク質などの生物学的データなど、ベクトルとして表現されるデータが大量に生成され、それらを効率的に処理する手法の開発が求められている。大規模データから近傍検索を行う代表的な手法に Chariker (2002) らが提案した手法がある。これは、ベクトルデータを2点間の距離関係を保ったままバイナリデータへと射影した後、ソートとバイナリサーチにより近傍点を検索する手法であるが、検索速度が遅いという欠点がある。そこで、我々は複合ソート法による全点間類似度検索法を提案する。提案手法は、ソート容量が大きくなる代わりに、近傍候補の距離計算を減少させることができ、高次元データに対して効率的に類似度検索を行うことができる。大規模画像データを用いた実験では、最近の近傍探索法との比較により提案手法の有効性が示す。



田部井研究員

ERATOセミナー@北大 2010/8/26

スケッチソート法による 高速な全点間類似検索

ERATO湊プロジェクト研究員

田部井靖生

共同研究者:宇野毅明(NII), 杉山将(東工大), 津田宏治(産総研)

発表手順


1. 動機
 - 大規模データ解析法の必要性
 - 全点間近傍検索
2. 手法
 - マルチプルソート法
 - Cosine Locality Sensitive Hashing
 - スケッチソート法
3. 実験
 - 大規模画像データを用いた実験
4. 現在進行中のプロジェクト
 - Jaccard係数への拡張
 - 大規模化合物データへの適応

近年、大規模データが大量に生成されている

- ▶ 80 million tiny images
 - 約8千万画像のデータ
 - <http://groups.csail.mit.edu/vision/TinyImages/>
 - ▶ Pubchem
 - 約2千万(2010年8月現在)の化合物からなるデータ
 - <http://pubchem.ncbi.nlm.nih.gov/>
 - ▶ NCBI Sequence Read Archive
 - 次世代シーケンサーから読まれた大規模配列データ
 - <http://www.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?>
- など

データはベクトルとして表現できる

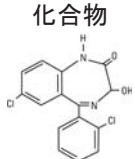
画像



sift特徴量

→ $x_i = (0.2, -0.3, -1.3, 1.2, 2.2, \dots)$

化合物



各次元は部分構造の有る(1)無し(0)を表す (フィンガープリント)

→ $x_i = (1, 0, 1, 0, 1, \dots)$

Locality Sensitive Hashing (Gionis et al. 99)

- ▶ ベクトルデータを文字列に射影する方法
 - **スケッチ**: バイナリの文字列
 - 元の空間の2点間の距離関係を(ハミング距離で)保ったまま射影する
- ▶ 大規模データを効率的に処理することが可能

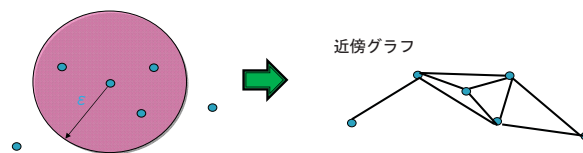
$x = (0.2, -0.3, -1.3, 1.2, 2.2, \dots)$



$s = 10101011101010101\dots$

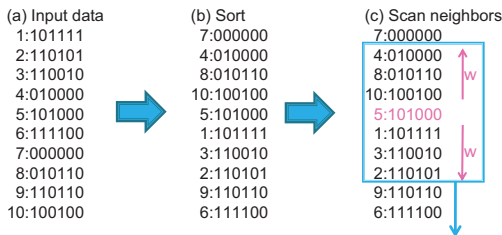
全点間類似度検索

- ▶ データセットからすべての似ているペアを探す
- ▶ データのペア (i, j) , $i < j$, $\Delta(x_i, x_j) \leq \epsilon$,
- ▶ 近傍グラフ(neighborhood graph)を作ることができる
- ▶ 応用
 - 重複データの削除, 大規模クラスタリング,
 - 半教師あり学習, 画像から特徴的な領域発見, など



シングルソート法 (Charikar, 02)

- ▶ スケッチから全点間類似度検索を行う方法
- ▶ ソートしてある上下幅 w のウィンドウでスキャンする
- ▶ 同じウィンドウに入ったペアに関して、ハミング距離を計算する



シングルソート法の欠点

- ▶ 精度を上げるために多くの距離計算が必要
- ▶ 真の近傍ペアを見逃す割合 (false negative) の理論的な見積もりができない



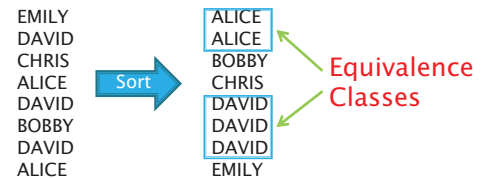
マルチプルソート法の概要 (Uno, 08)

- ▶ 同じ長さ l の文字列集合 $S = \{s_1, \dots, s_n\}$ を入力としてハミング距離 d 以内のペアの全列挙
- ▶ 距離 d 以内のペアの数を m とする
- ▶ 基数ソートを繰り返して、全ペアを $O(n+m)$ で列挙可能
- ▶ ブロックマスクを導入することにより高速化する



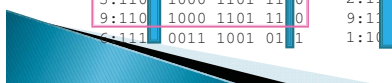
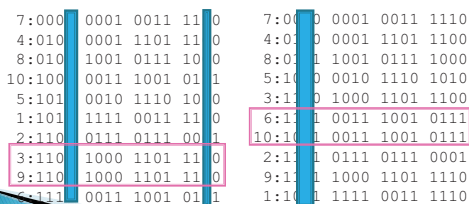
Special case: 同一文字列ペアの発見 ($d=0$)

- ▶ 基数ソートによりソートする、文字列を同一の文字列クラスに分割する: $O(n)$
- ▶ 各同一文字列クラスのすべてのペアに対してエッジを張る: $O(m)$
- ▶ 計算量: $O(n+m)$



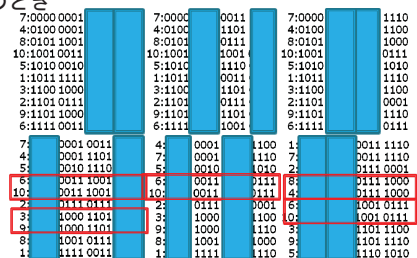
マルチプルソート法 ($d > 0$ のとき)

- ▶ d 個の文字を全通り選んでマスクする
- ▶ 基数ソートを $\binom{l}{d}$ 回繰り返す
- ▶ 計算量は d に対して指数、 l に対して多項式
- ▶ 文字列の数に対しては線形のまま $O(n+m)$



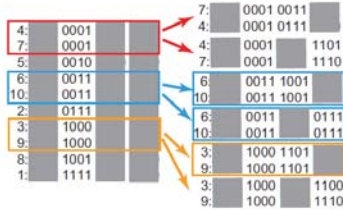
ブロックマスク

- ▶ 文字列を k 個のブロックに分割する
 - ▶ d 個のブロックを、全通りマスクする
 - ▶ **ブロック単位なので、ソート回数が増える**
 - ▶ 近傍でないペアも検出されてしまう
 - 実際に検出されたペアのハミング距離を計算して排除する
- $k=4, d=2$ のとき



再帰的アルゴリズム

- ▶ 初めに第一ブロックをソートし、Equivalence classを発見する
- ▶ それぞれの、Equivalence classに対して、次のブロックを付け加えてソートする
- ▶ k-dのブロックをソートしたら、各Equivalence Classに入っているペアに対して重複除去を行う
- ▶ 生き残ったペアに対して ハミング距離を実際に計算



重複除去

- ▶ ブロックの組み合わせに対して辞書順をつける

1	2	3	4
---	---	---	---

- k=4, d=2のとき
 - (1,2)<(1,3)<(1,4)<(2,3)<(2,4)<(3,4)
 - ▶ あるブロックの組み合わせに対して、完全一致するペアが見つかったら、それが最小の完全一致するブロックの組み合わせかをチェックする
- | | | | |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
|---|---|---|---|
- (2,4)が完全一致したら、(1,2),(1,3),(1,4),(2,3)が完全一致するかをチェックする
 - ▶ 実際には、1と3だけをチェックすればよい

疑似コード

```

1: function MULTIPLESORTINGMETHOD
2:   I ← {1, ..., n}
3:   B ← ∅
4:   RECURSION(I, B)
5:   return
6: end function
7: function RECURSION(I, B)
8:   if |B| = k - d then
9:     for (i, j) ∈ I × I, i < j do
10:      if s_i^b ≠ s_j^b for all b < max(B), b ∉ B then
11:        if HamDist(s_i, s_j) ≤ d then
12:          Report (i, j) to output file
13:        end if
14:      end if
15:    end for
16:    return
17:  end if
18:  for b in (max(B) + 1) .. (k + |B| + 1) do
19:    J ← Sorted indices based on b-th block {s_i^b}
20:    T ← Intervals of equivalence classes in J
21:    for each interval (x, y) ∈ T do
22:      RECURSION(J[x : y], B ∪ b)
23:    end for
24:  end for
25: return
26: end function
  
```

K-dの完全一致するブロックを持つペアに対する重複チェックとハミング距離の計算

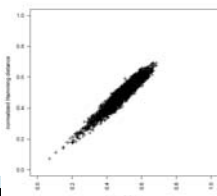
選ばれたブロックを基数ソートして equivalence classに分割する
Equivalence classの各メンバーを再帰処理する

スケッチソートの概要

- ▶ 距離が閾値以内のすべてのペアを列挙する
(i, j), i < j, Δ(x_i, x_j) ≤ ε,
- ▶ 基本的なアイデア
 - ベクトルデータをLSHで文字列データに射影する
 - 文字列データをマルチプルソート法にて近傍ペアを列挙する
- ▶ LSHを変えることで、様々な距離尺度に対応可能
 - ユーグリッド距離 (Raginsky et al. 2010), コサイン距離 (Geomans et al, 1999), Jaccard係数 (Broder, 1998) など
- ▶ ここでは、コサインLSHを使った方法を紹介

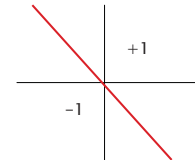
コサインLSH

- ▶ D次元のベクトルを長さlの0/1文字列に射影する
 - コサイン距離 $\Delta(x_i, x_j) = 1 - \frac{x_i^T x_j}{\|x_i\| \|x_j\|}$
- ▶ 射影先の正規化ハミング距離が、元の空間でのコサイン距離に単調増加する
 - Ex) 512bit



コサインLSH

- ▶ $R \in \mathbb{R}^{D \times \ell} : N(0, 1)$ からサンプルされたランダム行列
- ▶ 写像 $s_{ik} := \text{sign}(r_k^T x_i), r_k \in R, k = 1, \dots, \ell$



- ▶ 非衝突確率は2点間の角度に比例

$$\Pr(s_{ik} \neq s_{jk}) = \frac{\theta_{ij}}{\pi}, \forall k,$$

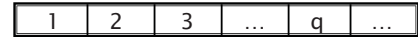
$$\theta_{ij} = \arccos \left(\frac{x_i^T x_j}{\|x_i\| \|x_j\|} \right)$$

スケッチソート

- ▶ Not good:長い文字列に、マルチプルソート法を適用する
- ▶ 文字列をチャンクへ分割
 - 長さ ℓ の文字列プールを Q 個つくる
 - マルチプルソート法を各プールに適用する
 $E_q = \{(i, j) \mid HamDist(s_i^q, s_j^q) \leq d, i < j\}$.
 - 各プールにおける出力セット E_q の和集合をとる
 $E = E_1 \cup \dots \cup E_Q$.
 - E の中で、 $\Delta(x_i, x_j) \leq \epsilon$ を満たすペアを出力

チャンク単位の重複除去

- ▶ 異なるチャンクで同じペアが見つかったと、重複したペアを出力してしまう
- ▶ 重複除去:チャンク q で、ハミング距離 d 以内のペアが見つかった場合は、チャンク $1, \dots, q-1$ でハミング距離 d 以内のものがない場合だけ出力



それぞれのチャンクのハミング距離をチェックする

各ペアに対して計3回のチェックを行う

1. ブロック単位の重複除去
 2. ハミング距離の計算
 3. チャンク単位の重複除去
- ▶ これらすべてを通過したペアのみコサイン距離を計算する



スケッチソートの疑似コード

```

1: function SKETCHSORT( $x_1, \dots, x_n$ )
2:   Use LSH to obtain sketches  $\{s_1, \dots, s_n\}_{1..Q}$  from
   data  $\{x_i\}_{1..n}$ 
3:    $T \leftarrow \emptyset$ 
4:   for  $q = 1..Q$  do
5:      $B \leftarrow \emptyset$ 
6:     RECURSION( $T, B, q$ )
7:   end for
8:   return  $T$ 
9: end function
10: function RECURSION( $T, B, q$ )
11:   if  $|B| = k - d$  then
12:     for  $(i, j) \in T \times T, i < j$  do
13:       if  $s_i^q = s_j^q$  for all  $h < \max(B), h \notin B$  then
14:         if  $HamDist(s_i^q, s_j^q) \leq d$  then
15:           if  $HamDist(s_i^r, s_j^r) > d$  for all  $r < q$ 
16:             then
17:               if  $\Delta(x_i, x_j) \leq \epsilon$  then
18:                 Report  $(i, j)$  to output file
19:               end if
20:             end if
21:           end if
22:         end for
23:       end if
24:     end if
25:   end if
26:   for  $b \in \max(B) + 1, k + |B| + 1$  do
27:      $J \leftarrow$  Sorted indices based on  $b$ -th block  $\{s_i^q\}_{i \in T}$ 
28:      $T' \leftarrow$  Intervals of equivalence classes in  $\{s_i^q\}_{i \in J}$ 
29:     for each interval  $(x, y) \in T'$  do
30:       RECURSION( $J[x..y], B \cup \{b, q\}$ )
31:     end for
32:   end for
33: end function
    
```

- ▶ 各チャンクに対して呼び出し
- ▶ Equivalence classの各メンバーのペアを発見
- ▶ 3回のチェック
- ▶ Equivalence class に対する再帰呼び出し

2種類のエラー

- ▶ 真のエッジセット E^* , 中間結果 E
- ▶ Type-I error (false positive): 近傍でないペアが、1つ以上のチャンクでハミング距離 d 以内となる事象
 $F_1 = \{(i, j) \mid (i, j) \in E, (i, j) \notin E^*\}$.
- ▶ Type II-error (false negative): 近傍ペアが、全てのチャンクで d より大きいハミング距離となる事象
 $F_2 = \{(i, j) \mid (i, j) \notin E, (i, j) \in E^*\}$.

Type II-errorの上限: Missing edge ratio

- ▶ false negativeの方が、致命的
 - false positiveは、コサイン距離の計算により除去することができる
- ▶ Missing edge ratio(false negative rate)は次のようにバウンドすることができる

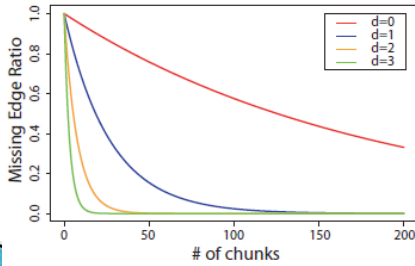
$$E \left[\frac{|F_2|}{|E^*|} \right] \leq \left(1 - \sum_{k=0}^{\lfloor d \rfloor} \binom{\ell}{k} p^k (1-p)^{\ell-k} \right)^Q,$$

- p は、LSHの非衝突確率の上限

$$p = \frac{\arccos(1 - \epsilon)}{\pi}.$$

チャンクの数Qに対する、Missing edge ratioの減少

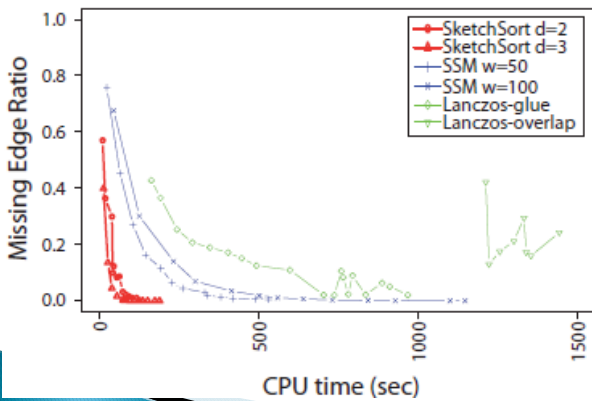
- ▶ チャンクの完全一致を用いる方法($d=0$)では、十分にMissing edge ratioを減らせない



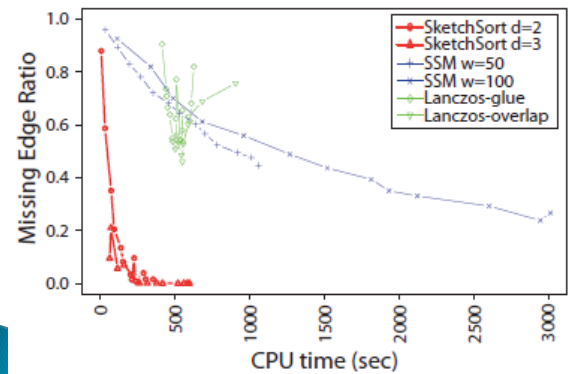
シングルソート法、Lanczos Bisectionとの比較

- ▶ データセット
 - MNIST (60,000 データ 748 次元)
 - TinyImage (100,000 データ, 960 次元)
- ▶ 評価指標はコサイン距離 0.15π でのMissing edge ratio
- ▶ 各チャンクは32ビット
- ▶ マルチプルソート法のハミング距離とブロック数: (2,5)と(3,6)に固定
- ▶ チャンク数, 2, 6, 10, ..., 50
- ▶ 最新の近傍探索法Lanczos Bisection (JMLR,09)との比較
 - データのマーzin最大を最大化する超平面で分割することを再帰的に繰り返し、分割領域内の点でナイーブ探索する

MNIST 閾値 0.15π



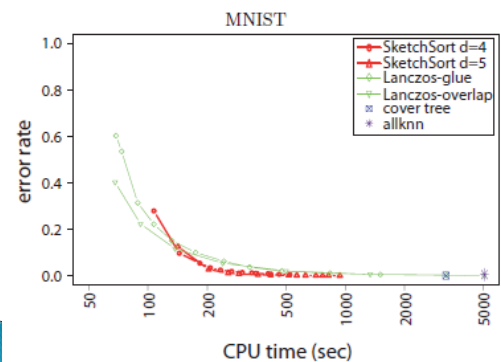
Tiny Image 閾値 0.15π



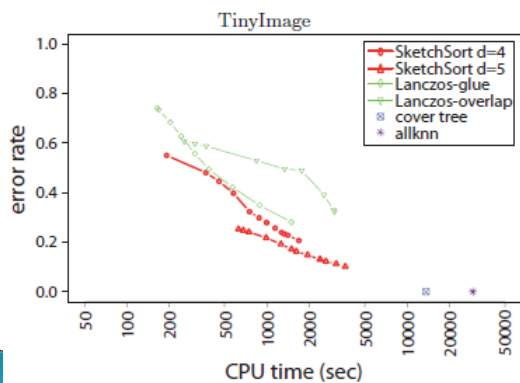
K-nearest neighbor検索実験

- ▶ スケッチソートの出力の各点ペアの距離が近い best-h個をpriority queueで保持
- ▶ 全ての処理が終了した後、近傍グラフ上の各点から距離2以内のノードから、もっとも近い点を選ぶ(リファインメント)
- ▶ 比較手法
 - Cover Tree (Beygelzimer et al., ICML 2006)
 - AllKNN (Ram et al., NIPS 2009),
 - Lanczos-bisection (JMLR, 2009)

5-Nearest Neighbor検索 (MNIST)

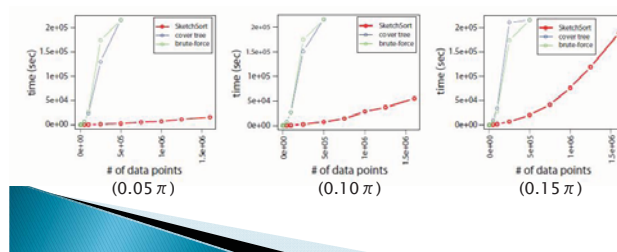


5-Nearest Neighbor検索 (TinyImage)



160万画像での実験

- False negative rate $< 1.0 \times 10^{-6}$
- 160万画像の全ペア類似度検索が、4.3時間 (0.05π)



化合物の大規模解析

- Pubchem
 - 約2千万の化合物のFinger print
- 100万データを使って Kristensen et al. (Algo. for Mol. Biol., 2010) の手法と比較
- 距離尺度として Jaccard 係数を使用
- 2つの集合 X_i, X_j の Jaccard 係数

$$r(X_i, X_j) = \frac{X_i \cap X_j}{X_i \cup X_j}$$

- $1 - r(X_i, X_j)$ で距離として利用

Minwise-independent permutations

- $X \subset [n]$ に対する射影 $\min(\pi_k(X))$
 - $\pi_k: [n]$ のランダム置換
- $n=10$ のとき

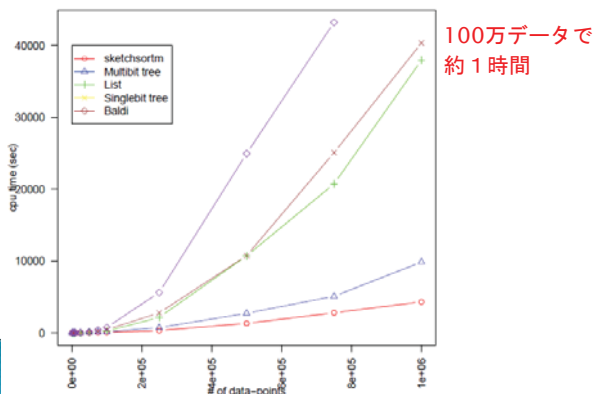
1	2	3	4	5	6	7	8	9	10
---	---	---	---	---	---	---	---	---	----

 $\downarrow \pi$

9	5	10	3	6	1	4	7	8	2
---	---	----	---	---	---	---	---	---	---
- このとき X の各要素は等確率で最小要素に射影される (**Minwise-independent**)
 - $x \in X, \Pr(\min\{\pi_k(X)\} = \pi_k(x)) = 1/|X|$
 - 衝突確率は Jaccard 係数に比例する

$$\Pr(\min\{\pi(X_i)\} = \min\{\pi(X_j)\}) = r(X_i, X_j)$$

100万化合物, 閾値0.1 Missing edge ratio 1×10^{-6} 以下



まとめ

- 高速な全点間類似度検索法を提案
- 大規模なベクトルデータに対して適応可能
- いろいろな応用を考えることができる
- コード
 - <http://code.google.com/p/sketchsort/>