Sequence Distance Embeddings

Graham Cormode, grahamc@dcs.warwick.ac.uk

(joint work with Mike Paterson, S. Muthukrishnan, Cenk Sahinalp, Uzi Vishkin, Piotr Indyk, Mayur Datar, Nick Koudas)

- Vector distances, streaming models
- Biological distances on permutations
- Editing distances on strings

Vector Distances

Two vectors a[1...n] and b[1...n]. Distance between them:

- Euclidean, $(\Sigma (a[i] b[i])^2)^{1/2}$
- Manhattan, $\sum |a[i] b[i]|$
- Hamming, $\Sigma (a[i] \neq b[i])$

Generalise this to the Lp distance =

$$\| a - b \|_{p} = (\Sigma | a[i] - b[i]|^{p})^{1/p}$$

Euclidean = L_2 , Manhattan = L_1 , Hamming ~ L_0 (as p \rightarrow 0, L_p \rightarrow Hamming. Also if a[i], b[i] = {0,1} then $(L_p)^p$ = Hamming.

Note, L_p may be interesteing for (fractional) p, 0

Stable Distributions

Let X be a random variable distributed with a *stable distribution*. Stable distributions have the property that

$$a_1X_1 + a_2X_2 + a_3X_3 + \dots + a_nX_n \sim ||(a_1, a_2, a_3, \dots, a_n)||_pX$$

if $X_1 ext{ ... } X_n$ are stable with stability paramater p

Gaussian distribution is stable with parameter 2

Stable distributions exist and can be simulated for all parameters 0 .

So, let $X = x_{1,1} \dots x_{m,n}$ be a matrix of values drawn from a stable distribution with parameter p...

Dimensionality Reduction of Vector Problems

$$(a_1 \dots a_n) \quad \bullet \quad (x_{1,1} \dots x_{m,1}) = (s_1, \dots s_m) \ [\ a \ sketch, s]$$

$$(x_{1,n} \dots x_{m,n})$$

$$(b_1 \dots b_n) \quad \bullet \quad (x_{1,1} \dots x_{m,1}) = (t_1, \dots t_m) \ [\ a \ sketch, t]$$

$$(\dots)$$

Then median($|s_1 - t_1|, |s_2 - t_2|, \dots, |s_m - t_m|$)/median(X) is an estimator for $||a - b||_p$

 $(X_{1 n} \dots X_{m n})$

Can guarantee the accuracy of this process: will be within a factor of 1+ ε with probability δ if m = O(1/ ε ² log 1/ δ)

Streaming model

Suppose n is huge (Mb, Gb, Tb, Pb), vectors are too big to store

Data arrives as a stream, eg

$$a[3] = 1027$$
, $a[19] = 76$, $a[61250] = 39$, $a[14235261] = 123$, $a[27] = -4369$, $a[41242] = 2$, $a[2] = 1241$, $a[213] = -16$...

The stream may be ordered, or may be in no particular order. Missing entries are assumed to be 0.

May also be unaggregated eg

$$a[3] \leftarrow +12 \qquad \dots \qquad a[3] \leftarrow +7 \qquad \dots \qquad a[3] \leftarrow -5$$

Overall effect here is a[3] = 14

Sketching the stream

Can compute these sketches on the stream because of linearity of the sketch function, eg

$$a[j] \leftarrow + 12$$

so
$$s_1 = s_1 + 12x_{1,j}$$
 $s_2 = s_2 + 12x_{2,j}$... $s_n = s_n + 12x_{n,j}$

So just add the right multiple of column j when j is updated

But... we need to store all xi, j = O(mn) storage

... which is more than just storing the vector a

Can we get round this?

Reducing space needs

- $x_{i,j}$ must be from stable distribution with fixed parameter
- $x_{i,j}$ must be the same every time it is used

We cheat a little by using pseudo-randomness.

We will use function **stable** to generate values from a stable distribution from a uniform distribution (like Box-Mueller transform for Normal dbn).

We use a random number generator that is good enough so that f(x) appears to be drawn from a uniform dbn.

Then
$$x_{1,j} = stable(f(x))$$
 $x_{2,j} = stable(f(f(x)))$ $x_{3,j} = stable(f(f(x)))$ etc.

Practical Use

So with O(m) space we can create a sketch to allow rapid comparison of huge streaming vectors.

Note m << n, in fact m is independent of n.

Implemented and tested in:

[C, Indyk, Koudas, Muthukrishnan02] - On massive tabular data, looking for clusterings using sketch computations to speed up comparisons for L_1 , L_2 and other Lp distances

[C, Datar, Indyk, Muthukrishnan] - On streaming vectors, to count number of distinct elements, find Hamming norm and Hamming distance.

Other distance measures

What about other more challenging distance measures eg

- string distances based on alignments
- computational biology measures

Outline approach:

