Combinatorics and Statistics of Gene Clusters

Laxmi Parida

Computational Biology Center, IBM T J Watson Research, Yorktown Heights & Courant Inst. of Mathematical Sciences, New York University

November 26-30, 2007

A (1) > A (2) > A (2) >

3

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

An everyday fundamental question...

What is the (statistical) significance associated with a *common* gene cluster q?

(1日) (1日) (日) 日

A Naive (straightforward) Approach

Assumption: Let the input be generated by a stationary, *iid* source which emits x_i with probability p_{xi}.

$$q = \{x_1, x_2, \ldots, x_l\}$$

$$\mathbb{P}_q = (p_{x_1})(p_{x_2})\dots(p_{x_l})$$

向下 イヨト イヨト

A more complex model (naive approach)

Assumption: Let the input be generated by a stationary, *iid* source which emits x_i with probability p_{xi}.

$$q = \{x_1(i_1), x_2(i_2), \dots, x_l(i_l)\}$$

$$\mathbb{P}_{q} = \left(\frac{(i_{1}+i_{2}+\ldots+i_{l})!}{i_{1}!\,i_{2}!\,\ldots\,i_{l}!}\right)(p_{x_{1}})^{i_{1}}(p_{x_{2}})^{i_{2}}\ldots(p_{x_{l}})^{i_{l}}$$

(multinomial coefficients)

向下 イヨト イヨト

Some troubling issues....

• What is p_{x_i} ?

- ► A gene x_i is not commonly known to occur too many times
- If $i \neq j$, how do p_{x_i} and p_{x_i} compare?

(日本) (日本) (日本)

臣

Some troubling issues....

▶ What is *p*_{x_i}?

•

- ► A gene x_i is not commonly known to occur too many times
- If $i \neq j$, how do p_{x_i} and p_{x_i} compare?

Can we do without guessing p_{x_i} ?

(日本) (日本) (日本)

臣

- Recall: Genes occur on chromosomes (linearly arranged)
 - Or, on a network (common connected components)

A (1) < A (2) < A (2) </p>

э

- Recall: Genes occur on chromosomes (linearly arranged)
 - Or, on a network (common connected components)
- ► Let the collection of <u>all subclusters</u> within the common cluster *q* that occurs in the species being compared be *S*.

- Recall: Genes occur on chromosomes (linearly arranged)
 - Or, on a network (common connected components)
- ► Let the collection of <u>all subclusters</u> within the common cluster *q* that occurs in the species being compared be *S*.

What is the probability of occurrence of subclusters S, given the occurrence of cluster q?

- lndividual p_{x_i} is irrelevant
- ▶ The occurrence of each x_i is equally likely (or not)
- Well-posed question

Image: A image: A

• 3 3 4

臣

- ▶ Individual p_{x_i} is irrelevant
- ▶ The occurrence of each *x_i* is equally likely (or not)
- Well-posed question, but a combinatorics nightmare!

- E - M

- Individual p_{x_i} is irrelevant
- ▶ The occurrence of each x_i is equally likely (or not)
- Well-posed question, but a combinatorics nightmare!

The 64K \$ question: Is this solvable?

Roadmap

Context

Permutation patterns

$\mathsf{Permutations} \to \mathsf{PQ} \ \mathsf{trees}$

Gene Proximity Analysis Statistics of permutations

3 × 4 3 ×

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

What are the common patterns?

 $s_1 = \dots g_1 g_2 g_3 g_4 g_5 g_6 g_7 \dots$ $s_2 = \dots g_8 g_5 g_2 g_4 g_3 g_9 g_0 \dots$

イロト イヨト イヨト イヨト

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

What are the common patterns?

$$s_1 = \dots g_1 g_2 g_3 g_4 g_5 g_6 g_7 \dots$$

$$s_2 = \dots g_8 g_5 g_2 g_4 g_3 g_9 g_0 \dots$$

Permutation patterns (π patterns)

$$s_1 = \dots g_1 \underbrace{g_2 g_3 g_4 g_5}_{g_2 g_3 g_4 g_5} g_6 g_7 \dots$$

$$s_2 = \dots g_8 \underbrace{g_5 g_2 g_4 g_3}_{g_3 g_9 g_0} g_9 \dots$$

Genes g_i in s_1 and g_i in s_2 are orthologous Block of genes g_2, g_3, g_4, g_5 appear together, albeit in a different order This block is a permutation (pattern)

 $\{g_2, g_3, g_4, g_5\}$

伺下 イヨト イヨト

πPatterns Example

(Pursuit of the Preposterous)

S = abcdefghijabdcefhgij (size 20)

・ 同 ト ・ ヨ ト ・ ヨ ト

3

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Permutation patterns

How bad is the scenario?

Permutation patterns:

$$\mathcal{O}(n^2)$$

Laxmi Parida Combinatorics and Statistics of Gene Clusters

(1日) (1日) (1日)

Maximal π patterns

Let P be the set of all patterns on a given input string s. $(p_1 \in P)$ is *non-maximal* with respect to $(p_2 \in P)$ if both of the following hold.

- (1) Each occurrence of p_1 on s is covered by an occurrence of p_2 on s.
- (2) Each occurrence of p_2 on s covers $l \ge 1$, occurrence(s) of p_1 on s.

A pattern $(p_2 \in P)$ is **maximal**, if there exists no $(p_1 \in P)$ such that p_2 is non-maximal w.r.t. p_1 .

・ロト ・ 同ト ・ ヨト ・ ヨト ……

Maximal π patterns

Let P be the set of all patterns on a given input string s. $(p_1 \in P)$ is *non-maximal* with respect to $(p_2 \in P)$ if both of the following hold.

- (1) Each occurrence of p_1 on s is covered by an occurrence of p_2 on s.
- (2) Each occurrence of p_2 on s covers $l \ge 1$, occurrence(s) of p_1 on s.

A pattern $(p_2 \in P)$ is **maximal**, if there exists no $(p_1 \in P)$ such that p_2 is non-maximal w.r.t. p_1 .

nonmaximal $\pi patterns \Leftrightarrow subclusters$

・ロト ・ 同ト ・ ヨト ・ ヨト ……

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

π patterns (nested & straddling)

$$s_1 = \dots g \boxed{a c d b e f g} e b \dots$$

$$s_2 = \dots b \boxed{g f e d a b c} f b \dots$$

$$p = \{a, b, c, d, e, f, g\},$$
(1)
nonMaximal(p) = {{e, f}, {f, g}, {e, f, g},
{a, b, c, d, }}. (2)
(3)

(本部) (本語) (本語) (二語)

Theorems on π **P**atterns

Theorem:

Let $R = \{Q' | Q' \text{ is non-maximal w.r.t } Q\}$. Then there exists a permutation Q'' of the elements of Q, such that for each Q', a permutation of the elements of Q' is a substring of Q''

Corollary 1: The ordering is not necessarily complete

Corollary 2:

A representation that captures the order of elements of Q along with intervals that captures each of Q' encodes Q

3

What is a PQ Tree?

Is there a sequence where the sets are consecutive?



The answer is YES for this set. All such sequences captured by the PQ tree.

・ロト ・ 同ト ・ ヨト ・ ヨト

PQ Trees Revisited.....

 $\{a,b,c,d\},\ \{b,c\},\ \{g,h\},\ \{h,i\}$

cbadghi dacbihg abcdihg ghidacb ihgdabc



Maximal π **Patterns** (Notation, PQ Tree notation)

- immediate neighbors "-" (Q)
- otherwise "," (P)
- groups "()" levels in the PQ tree

・ 同 ト ・ 三 ト ・ 三 ト

Linear Notation

$$p = \{a, b, c, d, e, f, g\},$$

nonMaximal(p) = {{e, f}, {f, g}, {e, f, g},
{a, b, c, d, }}.



$$p = ((a, b, c, d)-(e-f-g)).$$

▲□ > ▲圖 > ▲ 国 > ▲ 国 > →

π **Pattern (maximal)**

S = abcdefghijabdcefhgij

(input size 20, no of patts 25)

 $\begin{array}{l} \{a,b\}, \{a,b,c,d\}, \{a,b,c,d,e,f\}, \{a,b,c,d,e,f,g,h\} \\ \{a,b,c,d,e,f,g,h,i\}, \{a,b,c,d,e,f,g,h,i,j\}, \\ \{b,c,d\}, \{b,c,d,e,f\}, \{b,c,d,e,f,g,h\}, \{b,c,d,e,f,g,h,i,j\}, \\ \{c,d\}, \{c,d,e\}, \{c,d,e,f\}, \{c,d,e,f,g,h\}, \{c,d,e,f,g,h,i\}, \\ \{c,d,e,f,g,h,i,j\}, \\ \{c,f\}, \{e,f,g,h\}, \{e,f,g,h,i,j\}, \\ \{f,g,h\}, \{f,g,h,i,j\}, \\ \{g,h\}, \{g,h,i,j\} \\ \{i,j\} \end{array}$

A (1) < A (2) < A (2) </p>

e f

Is the definition any good?

Theorem

Let M be the set of all maximal patterns, i.e.,

$$M = \{p \in P | \text{ there is no } (p' \in P) \text{ maximal w.r.t } p\}$$

Then M is unique.

A (1) > A (2) > A (2) >

Algorithms

Find the patterns (WABI 03, JCB 04)
 For a fixed pattern size, the time taken is

 $\mathcal{O}(|\Sigma| + n(\log t)^2 \log |\Sigma|),$

where

$$t = \mathcal{O}(|\Sigma| + n \log |\Sigma|).$$

 Extract maximal form (CPM 05, JCB 06) The Minimal Consensus PQ Tree Algorithm (linear time)

(software available: http://www.mit.edu/oweimann/BIO/)

A (1) × A (2) × A (2) ×

Gene Proximity Analysis on Whole Genomes (CPM 05, JCB 06)

- Human, rat genomes (http://bio.math.berkeley.edu/slam)
- 25,422 putative orthologous genes
- > 23 human, 21 rat chromosomes

- 4 回 ト 4 ヨ ト 4 ヨ ト

э

Human & Rat

π **Pattern (166)**



Human & Rat

πPattern (303)



Human & Rat

πPattern (250)



Human chromosomes 10,11 Rat chromosome 1

イロト イヨト イヨト イヨト

 $\begin{array}{l} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Gene Proximity Analysis Statistics of permutations

Proximity- Summary

	Number of	Number of
	all patterns	maximal patterns
E Coli K-12 & B Subtilis	15,000	450
human & rat	1,574,312	504

(Joint work with Revital Eres, Oren Weimann, Gadi Landau)

ヘロト 人間 とくほとう ほん

臣

Statistical Significance

Two related but distinct questions:

Naive model: Given n random permutations of k genes,
what is the probability that K of these n contain the cluster q ?

Structured cluster model: Given that a permutation pattern q occurs K times in the input, what is the probability of its maximal form given as a PQ tree T ?

Statistical Significance of Large Gene Clusters, Laxmi Parida, Journal of Computational Biology, 14(9), pp 1145–1159, 2007.

 $\begin{array}{c} \mbox{Context} \\ \mbox{Permutations} \rightarrow \mbox{PQ} \mbox{ trees} \\ \end{array} \begin{array}{c} \mbox{Gene Proximity Analysis} \\ \mbox{Statistics of permutations} \\ \end{array}$

How to answer Question 2 (for K = 2)?

▶ Recall: PQ T in relation with multiple occurrences

- ► a single occurrence has <u>no</u> PQ T
- Consensus PQ of $o_1, o_2 \in Fr(T)$ is not necessarily T
- Use one occurrence as reference WLOG

イロト イポト イヨト イヨト

-

Context Gene Proxin Permutations → PQ trees Statistics of

Gene Proximity Analysis Statistics of permutations

How to answer Question 2 (for K = 2)?

Recall: PQ T in relation with multiple occurrences

- ► a single occurrence has <u>no</u> PQ T
- Consensus PQ of $o_1, o_2 \in Fr(T)$ is not necessarily T
- Use one occurrence as reference WLOG



(1) PQ tree T. (2) Ref PQ Tree T' on 7 consecutive integers.

- 4 周 ト - 4 三 ト - 4 三 ト

Working with T' (on consecutive integers)

- Count the size of the *frontier* set of T'
 - Interpretation: Each $o \in Fr(T')$ is such that the consensus PQ tree of o and the identity permutation is T'
- What is the bottleneck?
 - How many possibilities does a P node (with k children) introduce?
 - Each possibility cannot have any internal structure

A (1) × A (2) × A (2) ×

Working with T' (on consecutive integers)

- Count the size of the *frontier* set of T'
 - Interpretation: Each $o \in Fr(T')$ is such that the consensus PQ tree of o and the identity permutation is T'
- What is the bottleneck?
 - How many possibilities does a P node (with k children) introduce?
 - Each possibility <u>cannot</u> have any internal structure We call this P-Arrangement of size k

(4月) キョン キョン

What is an interval?

_

inte	rval	$\Pi(q_1[k_1k_2])$	size
$[k_1k_2]$			
<u>non-trivial</u> :			
[34]	52431	$\{3, 4\}$	2
[24]	5 243 1	$\{2, 3, 4\}$	3
[14]	5243 1	$\{2,3,4,5\}$	4
<u>trivial</u> :			
[15]	52431	$\{1,2,3,4,5\}$	5

イロト イヨト イヨト イヨト

P-arrangement

An arrangement of size k is a *P*-arrangement if it has no non-trivial intervals.



イロト イボト イヨト イヨト

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Gene Proximity Analysis Statistics of permutations

The central question:

What is the number of P-arrangements of size k?

イロン イヨン イヨン イヨン

The central question:

What is the number of P-arrangements of size k?

Theorem (Par07)

Let q be a P-arrangement of size k + 1. Let q' be obtained by replacing an extreme element (either k + 1 or 1) from its position j in q, with the empty symbol. Then

- 1. q' is a P-arrangement, or,
- 2. every interval $[i_1 \dots i_2]$ in q' is such that $i_1 < j < i_2$.

・ロト ・ 同ト ・ ヨト ・ ヨト

 $\begin{array}{cc} \text{Context} & \text{Gene} \\ \text{Permutations} \rightarrow \text{PQ trees} & \text{Statistical} \end{array}$

Gene Proximity Analysis Statistics of permutations

Theorem illustration (nested arrangements)

$$q = 9 \ 1 \ 5 \ 2 \ \phi \ 3 \ 6 \ 4 \ 7 \ 10 \ 8$$

9 1 5 2
$$\phi$$
 3 6 4 7 10 8 9 1 5 2 ϕ 3 6 4 7 10 8 6(2)

Signature:
$$sig(q) = 2(1) < 5(1) < 6(2) < 10(1)$$

10

Signature Lemma

Let q be an arrangement of size k with symbol ϕ in position j with all the intervals $[i_{r1} \dots i_{r2}]$ satisfying

$$i_{r1} < j < i_{r2},$$

for all $1 \le r \le K$. Let the size of the interval be $i_r = i_{r2} - i_{r1} + 1$.

- 1. (straddling intervals) If two such intervals, where one is not nested in the other, are of size *i* and *i'*, then i = i' and they must overlap in i-1 positions.
- 2. (uniqueness and form)

$$sig(q) = i_1(k_{i_1}) < i_2(k_{i_2}) < \ldots < i_r(k_{i_r}) < \ldots < i_K(k_{i_K}),$$

is unique with $k_{i_1} = 1$, $k_{i_K} = 1$, $i_K = k$, and each k_{i_r} , $1 \leq r \leq K$, is either 1 or 2.

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Gene Proximity Analysis Statistics of permutations

Formula for number of *P*-arrangements

$$Pa(2) = 2,$$

 $Pa(3) = 0,$
 $Pa(4) = 2,$
 $Pa(k) = Nst'(k-1),$ for $k > 4.$

Polynomial time dynamic programming solution.

Number of nested arrangements with viable positions to get k + 1-sized *P*-arrangements:

$$Nst'(k) = S(k,2) - S_{cnt}(k,2) + \sum_{l=4}^{k} (l-1)S(k,l) - 2S_{cnt}(k,l).$$

ヘロト 人間 とくほとう ほん

臣

Number of nested arrangements with smallest I and largest u interval sizes:

$$S(u, l) = 4S(u-1, l) + 2S(u-2, l) + \sum_{y=3}^{u-l} \Delta_{u-y} Pa(\Delta_{u-y})S(u-y, l)$$

Number of nested arrangements with the extreme element in the smallest interval:

$$S_{cnt}(u, l) = 2S_{cnt}(u-1, l) + \sum_{y=3}^{u-l} Pa(\Delta_{u-y})S_{cnt}(u-y, l)$$

3

Context Permutations → PQ trees

Gene Proximity Analysis Statistics of permutations

Back to ... Estimating the frontier size

$$\#(A) = \begin{cases} 1 & \text{if } A \text{ is a leaf node,} \\ 2\prod_{j=1}^{c} \#(A_j) & \text{if } A \text{ is a } Q \text{ node,} \\ Pa(c)\prod_{j=1}^{c} \#(A_j) & \text{if } A \text{ is a } P \text{ node.} \end{cases}$$

ヘロト 人間 とくほとう ほん

Estimating the frontier size



A (1) × A (2) × A (2) ×

Context Permutations → PQ trees

Gene Proximity Analysis Statistics of permutations

Estimating Fr(T'): A complete example



(1) The input PQ tree T. (2) Numbering the leaf nodes.

- 4 回 2 4 三 2 4 三 2 4

Example...



Node *A* 3142 2413

Node <i>B</i>	
567	765

Node C		
AB	ВA	

(2) Numbering the leaf nodes& labeling the internal nodes.

(3) The possible arrangmeents.

- 4 回 ト 4 三 ト 4 三 ト

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Gene Proximity Analysis Statistics of permutations

Estimating the frontier size



7654321

(4) The 10 possible arrangements.



gfedcba

(5) Arrangements in input alphabet.

Problems in Combinatorics

- What is the number of simple permutations of degree k?
 What is the number of *P*-arrangements of size k?
 We give the first explicit formula for this number.
- *O*(n²) time to recognize a simple permutation.

 We give *O*(n) time to recognize a *P*-arrangement.

A (10) × (10) × (10) ×

 $\begin{array}{c} \text{Context} \\ \text{Permutations} \rightarrow \text{PQ trees} \end{array}$

Gene Proximity Analysis Statistics of permutations

Combinatorics to probabilities

$$pr(T) = \frac{|Fr(T')|}{n!}$$

Laxmi Parida Combinatorics and Statistics of Gene Clusters

イロト イヨト イヨト イヨト

Back to the clusters







- 4 回 ト 4 三 ト 4 三 ト

166	$0.236 imes 10^{-3}$
303	$0.376 imes10^{-4}$
250	$0.404 imes10^{-4}$

Structured clusters with multiplicities

イロト イヨト イヨト イヨト

Structured clusters with multiplicities

$$p=\{a,b,c(2),d,e,x\},$$

with exactly three occurrences given as

$$o_1 = deabcxc,$$

 $o_2 = cdeabxc,$
 $o_3 = cxcbaed.$



Clusters on networks (common connected components)

ヘロト 人間 とくほとう ほん

Clusters on networks (common connected components)



→

臣

Clusters on networks (common connected components)





< 3 >

Context Gene Proximity Analysis Permutations → PQ trees Statistics of permutations

> Chapman & Hall/CRC Mathematical and Computational Biology Series

Pattern Discovery in Bioinformatics Theory & Algorithms





Laxmi Parida

Chapman & Hall/CRC

Laxmi Parida Combinatorics and Statistics of Gene Clusters

イロト イヨト イヨト イヨト