

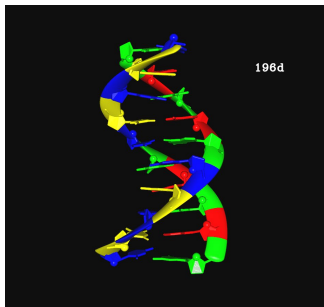
Assessing the significance of Sets of Words

V. Boeva, [J. Clément](#), M. Régner and M. Vandenbergert

Moscow, [Marne-la-Vallée-CNRS](#), INRIA, Biozentrum

CPM 2005 – June 22, 2005

Genome analysis



- Structure of the DNA
- Over-(and under) represented DNA motifs
- Regulation sites in genes

Paradigm: biological/random comparison

Paradigm

Comparing mathematical criteria in biological and random sequences, one can extract biological features.

Example

- If a pattern occurs with different frequencies in a real sequence and a random sequence, then it could have a biological meaning.
- When searching for over-represented or under-represented patterns, we must test that such a pattern is not generated by randomness itself.

Paradigm: biological/random comparison

Paradigm

Comparing mathematical criteria in biological and random sequences, one can extract biological features.

Example

- If a pattern occurs with different frequencies in a real sequence and a random sequence, then it could have a biological meaning.
- When searching for over-represented or under-represented patterns, we must test that such a pattern is not generated by randomness itself.

Over-represented patterns

Biological sequence

TTCATTATCTCCATT**C****GCTGGTGG**GCAAGGACTTGAGCTATCGCCCTTTC . . .
 GCATAAAGTTATTCATAAACTGTCAGGGGTTTCGGTTGCC**GCTGGTGG**AAC . . .
 AG**GCTGGTGG**ACGCCTACGTTATTT**GCTGGTGG**ACTGGAAATCATCTAG . . .
 TCCAACGAAATA**GCTGGTGG**TCTACACTCATATCGTTATTAACAAACGAA . . .
 AGAAACTAATGGGTGTCACAG**GCTGGTGG**GCTCGTATTTTGTAGGAGGTCA . . .

Random sequence

ATATATATATTTATCTTGCAACTCGGAGAATTCTATTAATATATGAACGA . . .
 ACGTAGATGACAACAATTAGCATGTGGATTTGTAAGGTAAGTTTCTTGTG . . .
 CGTTGGTTGGTCATCGATGCAATGAATGAGTCGTTTAAAATAAGACTCGA . . .
 TTGTCTCTCAAGTTTTTTTTTGCATTACCATTCTAAG**GCTGGTGG**ATATAGG . . .
 GTTTACAAGTTTTAACCTTTTGTCACTCGTCACCTTATGTGTGGCTTTAA . . .

→ *Chi* Motif in *E. COLI*.

Over-represented patterns

Biological sequence

TTCATTATCTCCATT**C****GCTGGTGG**GCAAGGACTTGAGCTATCGCCCTTTC . . .
 GCATAAAGTTATTCATAAACTGTCAGGGGTTCGGTTGCC**GCTGGTGG**AAC . . .
 AG**GCTGGTGG**ACGCCTACGTTATTT**GCTGGTGG**ACTGGAAATCATCTAG . . .
 TCCAACGAAATA**GCTGGTGG**TCTACACTCATATCGTTATTAACAAACGAA . . .
 AGAAACTAATGGGTGTCACAG**GCTGGTGG**GCTCGTATTTTGTAGGAGGTCA . . .

Random sequence

ATATATATATTTATCTTGCAACTCGGAGAATTCTATTAATATATGAACGA . . .
 ACGTAGATGACAACAATTAGCATGTGGATTTGTAAGGTAAGTTTCTTGTG . . .
 CGTTGGTTGGTCATCGATGCAATGAATGAGTCGTTTAAAATAAGACTCGA . . .
 TTGTCTCTCAAGTTTTTTTTTGCATTACCATTCTAA**GCTGGTGG**ATATAGG . . .
 GTTTACAAGTTTTAACCTTTTGTCACTCGTCACCTTATGTGTGGCTTTAA . . .

→ *Chi* Motif in *E. COLI*.

Significance of a pattern?

We need to characterize the “probabilistic behaviour” of a pattern.

Problem

There exist measures expressed by expressions and recurrences which can be cumbersome to handle (+ numerical instability)

Our contribution

- A rewriting of exact matricial formula to get tractable formula for the probability of first occurrence of a motif and first co-occurrence of a pair of motifs (here a motif can be a set of words)
- Exhibit a few combinatorial parameters for sets of words
- We consider a positional pattern (\approx affinity matrices) for which efficient computation of these parameters is possible

Significance of a pattern?

We need to characterize the “probabilistic behaviour” of a pattern.

Problem

There exist measures expressed by expressions and recurrences which can be cumbersome to handle (+ numerical instability)

Our contribution

- A rewriting of exact matricial formula to get tractable formula for the probability of first occurrence of a motif and first co-occurrence of a pair of motifs (here a motif can be a set of words)
- Exhibit a few combinatorial parameters for sets of words
- We consider a positional pattern (\approx affinity matrices) for which efficient computation of these parameters is possible

Evaluation of the significance of a pattern H

Let

- $O_n(H)$ = Random variable counting the number of occurrences of the pattern H in a **random** text of length n .
- $\text{Obs}(H)$ = the number of occurrences of the pattern H in the **biological** sequence.

How to estimate the significance?

- **z-score:**

$$Z(H) = \frac{E[O_n(H)] - \text{Obs}(H)}{\sqrt{\text{Var } O_n(H)}}$$

[Meaningful for a normal distribution, not too far from the mean]

- **p-values:** $p(H) = \Pr\{O_n(H) \geq \text{Obs}(H)\}$ [Large deviations techniques]
- **Probability of first occurrence** $\mathcal{F}_n = \Pr\{O_n(H) > 0\}$ [related to waiting time]

Evaluation of the significance of a pattern H

Let

- $O_n(H)$ = Random variable counting the number of occurrences of the pattern H in a **random** text of length n .
- $\text{Obs}(H)$ = the number of occurrences of the pattern H in the **biological** sequence.

How to estimate the significance?

- **z-score:**

$$Z(H) = \frac{E[O_n(H)] - \text{Obs}(H)}{\sqrt{\text{Var } O_n(H)}}$$

[Meaningful for a normal distribution, not too far from the mean]

- **p-values:** $p(H) = \Pr\{O_n(H) \geq \text{Obs}(H)\}$ [Large deviations techniques]
- **Probability of first occurrence** $\mathcal{F}_n = \Pr\{O_n(H) > 0\}$ [related to waiting time]

Evaluation of the significance of a pattern H

Let

- $O_n(H)$ = Random variable counting the number of occurrences of the pattern H in a **random** text of length n .
- $\text{Obs}(H)$ = the number of occurrences of the pattern H in the **biological** sequence.

How to estimate the significance?

- **z-score:**

$$Z(H) = \frac{E[O_n(H)] - \text{Obs}(H)}{\sqrt{\text{Var } O_n(H)}}$$

[Meaningful for a normal distribution, not too far from the mean]

- **p-values:** $p(H) = \Pr\{O_n(H) \geq \text{Obs}(H)\}$ [Large deviations techniques]
- **Probability of first occurrence** $\mathcal{F}_n = \Pr\{O_n(H) > 0\}$ [related to waiting time]

Evaluation of the significance of a pattern H

Let

- $O_n(H)$ = Random variable counting the number of occurrences of the pattern H in a **random** text of length n .
- $\text{Obs}(H)$ = the number of occurrences of the pattern H in the **biological** sequence.

How to estimate the significance?

- **z-score:**

$$Z(H) = \frac{E[O_n(H)] - \text{Obs}(H)}{\sqrt{\text{Var } O_n(H)}}$$

[Meaningful for a normal distribution, not too far from the mean]

- **p-values:** $p(H) = \Pr\{O_n(H) \geq \text{Obs}(H)\}$ [Large deviations techniques]
- **Probability of first occurrence** $\mathcal{F}_n = \Pr\{O_n(H) > 0\}$ [related to waiting time]

Probabilistic models

These criteria suppose an underlying probabilistic model

- **Shuffling** (exact) model: fix a parameter k and keep the same distribution of factors of length k as in a reference sequence [hard to study!]
- **Bernoulli** model: $(p_i)_{i \in \Sigma}$ [memoryless]
- **Markov** model: $\mathbb{P} = (p_{ij})_{i,j \in \Sigma}, (\pi_i)_{i \in \Sigma}$ [finite context]

Our work concerns Bernoulli and Markov model.

Probabilistic models

These criteria suppose an underlying probabilistic model

- **Shuffling** (exact) model: fix a parameter k and keep the same distribution of factors of length k as in a reference sequence [hard to study!]
- **Bernoulli** model: $(p_i)_{i \in \Sigma}$ [memoryless]
- **Markov** model: $\mathbb{P} = (p_{ij})_{i,j \in \Sigma}, (\pi_i)_{i \in \Sigma}$ [finite context]

Our work concerns Bernoulli and Markov model.

Probabilistic models

These criteria suppose an underlying probabilistic model

- **Shuffling** (exact) model: fix a parameter k and keep the same distribution of factors of length k as in a reference sequence [hard to study!]
- **Bernoulli** model: $(p_i)_{i \in \Sigma}$ [memoryless]
- **Markov** model: $\mathbb{P} = (p_{i|j})_{i,j \in \Sigma}, (\pi_i)_{i \in \Sigma}$ [finite context]

Our work concerns Bernoulli and Markov model.

Probabilistic models

These criteria suppose an underlying probabilistic model

- **Shuffling** (exact) model: fix a parameter k and keep the same distribution of factors of length k as in a reference sequence [hard to study!]
- **Bernoulli** model: $(p_i)_{i \in \Sigma}$ [memoryless]
- **Markov** model: $\mathbb{P} = (p_{i|j})_{i,j \in \Sigma}, (\pi_i)_{i \in \Sigma}$ [finite context]

Our work concerns Bernoulli and Markov model.

Over-(or under-)representation of patterns

Input

- model for the sequence
- n , sequence length
- pattern H (or a set of patterns \mathcal{H})

Question

Find the probabilistic law of the pattern in random sequences of size n (expected values, variances, waiting time, ...)

Two different approaches

- Experimental: A. Denise, M.-F. Sagot, L. Marsan
- Analytical approach

Over-(or under-)representation of patterns

Input

- model for the sequence
- n , sequence length
- pattern H (or a set of patterns \mathcal{H})

Question

Find the probabilistic law of the pattern in random sequences of size n (expected values, variances, waiting time, ...)

Two different approaches

- Experimental: A. Denise, M.-F. Sagot, L. Marsan
- **Analytical approach**

Analytical approach

Probabilistic methods

[Prum, Rodolphe, de Turkheim 95], [Schbath 97], [Apostolico, Bock, Xuyan 98], [Reinert, Schbath, Waterman 00], ...

Combinatorial methods

Generating functions of probabilities [Régnier, Szpankowski 98], [Nicodème, Salvy, Flajolet 99], ...

Large deviations

[Denise, Régnier 04]

See also [Lothaire vol.3](#) “Applied Combinatorics on Words” to appear soon with a chapter by Reinert, Schbath, Waterman and another by Jacquet, Szpankowski.

Analytical approach

Probabilistic methods

[Prum, Rodolphe, de Turkheim 95], [Schbath 97], [Apostolico, Bock, Xuyan 98], [Reinert, Schbath, Waterman 00], ...

Combinatorial methods

Generating functions of probabilities [Régnier, Szpankowski 98], [Nicodème, Salvy, Flajolet 99], ...

Large deviations

[Denise, Régnier 04]

See also [Lothaire vol.3](#) “Applied Combinatorics on Words” to appear soon with a chapter by Reinert, Schbath, Waterman and another by Jacquet, Szpankowski.

Analytical approach

Probabilistic methods

[Prum, Rodolphe, de Turkheim 95], [Schbath 97], [Apostolico, Bock, Xuyan 98], [Reinert, Schbath, Waterman 00], ...

Combinatorial methods

Generating functions of probabilities [Régnier, Szpankowski 98], [Nicodème, Salvy, Flajolet 99], ...

Large deviations

[Denise, Régnier 04]

See also [Lothaire vol.3](#) “Applied Combinatorics on Words” to appear soon with a chapter by Reinert, Schbath, Waterman and another by Jacquet, Szpankowski.

Analytical approach

Probabilistic methods

[Prum, Rodolphe, de Turkheim 95], [Schbath 97], [Apostolico, Bock, Xuyan 98], [Reinert, Schbath, Waterman 00], ...

Combinatorial methods

Generating functions of probabilities [Régnier, Szpankowski 98], [Nicodème, Salvy, Flajolet 99], ...

Large deviations

[Denise, Régnier 04]

See also [Lothaire vol.3](#) “Applied Combinatorics on Words” to appear soon with a chapter by Reinert, Schbath, Waterman and another by Jacquet, Szpankowski.

Combinatorial analysis : from generating functions to formulas

Methodology

- Equations on languages translate to equations on generating functions. A language $\mathcal{L} \rightarrow L(z) = \sum_{w \in \mathcal{L}} \alpha(w) P(w) z^{|w|}$
- Computing parameters requires extracting coefficients of generating functions.

Extracting coefficient can be tedious (numerical instability, time efficiency), slow and require the use of formal systems like Maple or Mathematica.

Techniques and tools

- Complex analysis
- Algebraic operations on series

Combinatorial analysis : from generating functions to formulas

Methodology

- Equations on languages translate to equations on generating functions. A language $\mathcal{L} \rightarrow L(z) = \sum_{w \in \mathcal{L}} \alpha(w) P(w) z^{|w|}$
- Computing parameters requires extracting coefficients of generating functions.

Extracting coefficient can be tedious (numerical instability, time efficiency), slow and require the use of formal systems like Maple or Mathematica.

Techniques and tools

- Complex analysis
- Algebraic operations on series

Example

Let $F_n(\mathcal{H})$ be the probability that at least one word in the set \mathcal{H} occurs in a random sequence of size n . One has

$$F_{\mathcal{H}}(z) = \sum_{n \geq 0} F_n(\mathcal{H})z^n = \frac{1}{1-z} H(z) \mathbb{D}(z)^{-1} \mathbf{1}_q^t$$

Here: if $\mathcal{H} = \{H_1, \dots, H_q\}$ with the H_i words of length m ,
 $H(z) = (P(H_1)z^m, \dots, P(H_q)z^m)$,
 $\mathbb{D}(z) = (1 - z)(\mathbb{I} + \mathbb{C}(z)) + \mathbb{H}(z)$.

Applying algebraic and combinatorial identities, we get

$$F_{\mathcal{H}}(z) = \frac{1}{1-z} - \frac{1}{1-z+Q_{\mathcal{H}}(z)} \text{ with } Q_{\mathcal{H}}(z) = \text{Trace}(\mathbb{H}(z)\mathbb{A}(z)^{-1}).$$

The asymptotics comes from dominant singularity in the complex plane $z \in \mathbb{C}$, obtained by bootstrapping

$$P\{O_n > 0\} \approx 1 - (1 + P(\mathcal{H}) - \tilde{C}(\mathcal{H}))^{-n}, \text{ when } nP(\mathcal{H}) < 1$$

Example

Let $F_n(\mathcal{H})$ be the probability that at least one word in the set \mathcal{H} occurs in a random sequence of size n . One has

$$F_{\mathcal{H}}(z) = \sum_{n \geq 0} F_n(\mathcal{H})z^n = \frac{1}{1-z} H(z) \mathbb{D}(z)^{-1} \mathbf{1}_q^t$$

Here: if $\mathcal{H} = \{H_1, \dots, H_q\}$ with the H_i words of length m ,
 $H(z) = (P(H_1)z^m, \dots, P(H_q)z^m)$,
 $\mathbb{D}(z) = (1 - z)(\mathbb{I} + \mathbb{C}(z)) + \mathbb{H}(z)$.

Applying algebraic and combinatorial identities, we get

$$F_{\mathcal{H}}(z) = \frac{1}{1-z} - \frac{1}{1-z+Q_{\mathcal{H}}(z)} \text{ with } Q_{\mathcal{H}}(z) = \text{Trace}(\mathbb{H}(z)\mathbb{A}(z)^{-1}).$$

The asymptotics comes from dominant singularity in the complex plane $z \in \mathbb{C}$, obtained by bootstrapping

$$P\{O_n > 0\} \approx 1 - (1 + P(\mathcal{H}) - \tilde{C}(\mathcal{H}))^{-n}, \text{ when } nP(\mathcal{H}) < 1$$

Example

Let $F_n(\mathcal{H})$ be the probability that at least one word in the set \mathcal{H} occurs in a random sequence of size n . One has

$$F_{\mathcal{H}}(z) = \sum_{n \geq 0} F_n(\mathcal{H})z^n = \frac{1}{1-z} H(z) \mathbb{D}(z)^{-1} \mathbf{1}_q^t$$

Here: if $\mathcal{H} = \{H_1, \dots, H_q\}$ with the H_i words of length m ,
 $H(z) = (P(H_1)z^m, \dots, P(H_q)z^m)$,
 $\mathbb{D}(z) = (1 - z)(\mathbb{I} + \mathbb{C}(z)) + \mathbb{H}(z)$.

Applying algebraic and combinatorial identities, we get

$$F_{\mathcal{H}}(z) = \frac{1}{1-z} - \frac{1}{1-z+Q_{\mathcal{H}}(z)} \text{ with } Q_{\mathcal{H}}(z) = \text{Trace}(\mathbb{H}(z)\mathbb{A}(z)^{-1}).$$

The asymptotics comes from dominant singularity in the complex plane $z \in \mathbb{C}$, obtained by bootstrapping

$$P\{O_n > 0\} \approx 1 - (1 + P(\mathcal{H}) - \tilde{C}(\mathcal{H}))^{-n}, \text{ when } nP(\mathcal{H}) < 1$$

Combinatorial properties of a set of words

Let $F = \text{ATAA}$ and $G = \text{AACT}$.

ATAA

..AACT

...AACT

- CT and AACT are right complements.
- The set of right complements of F in G is $C_{F,G}$ the *correlation set*.
- When $F = G$, the *autocorrelation set* is $A_F = C_{F,F} + \varepsilon$ with ε the empty word.
- The set $\tilde{C}_{F,\mathcal{H}}$ of *minimal right complements* of F in \mathcal{H} is the set of *minimal* words in $\cup_{G \in \mathcal{H}} C_{F,G}$ for the *prefix order*.

Combinatorial properties of a set of words

Let $F = \text{ATAA}$ and $G = \text{AACT}$.

ATAA
 ..AACT
 ...AACT

- CT and AACT are right complements.
- The set of right complements of F in G is $C_{F,G}$ the *correlation set*.
- When $F = G$, the *autocorrelation set* is $A_F = C_{F,F} + \varepsilon$ with ε the empty word.
- The set $\tilde{C}_{F,\mathcal{H}}$ of *minimal right complements* of F in \mathcal{H} is the set of *minimal* words in $\cup_{G \in \mathcal{H}} C_{F,G}$ for the *prefix order*.

Combinatorial parameters

A few combinatorial parameters

$$P(\mathcal{H}) = \sum_{w \in \mathcal{H}} P(w)$$

$$C(\mathcal{H}) = \sum_{F, G \in \mathcal{H}} \sum_{w \in C_{F, G}} P(Fw)$$

$$\tilde{C}(\mathcal{H}) = \sum_{F \in \mathcal{H}} \sum_{w \in \tilde{C}_{F, \mathcal{H}}} P(Fw)$$

These suffice to express several quantities

$$E[O_n(\mathcal{H})] = (n - m + 1) P(\mathcal{H}),$$

$$\begin{aligned} \text{Var}[O_n(\mathcal{H})] = & (n - m + 1) (P(\mathcal{H}) + (1 - 2m) P(\mathcal{H})^2 + 2 C(\mathcal{H})) \\ & + m(m - 1) P(\mathcal{H})^2 - 2 \tilde{C}(\mathcal{H}) \end{aligned}$$

$$\Pr\{O_n > 0\} \approx 1 - (1 + P(\mathcal{H}) - \tilde{C}(\mathcal{H}))^{-n}, \text{ if } n P(\mathcal{H}) < 1$$

Better than systems of functional equations **only if** we know how to compute these quantities **efficiently!**

Combinatorial parameters

A few combinatorial parameters

$$P(\mathcal{H}) = \sum_{w \in \mathcal{H}} P(w)$$

$$C(\mathcal{H}) = \sum_{F, G \in \mathcal{H}} \sum_{w \in C_{F, G}} P(Fw)$$

$$\tilde{C}(\mathcal{H}) = \sum_{F \in \mathcal{H}} \sum_{w \in \tilde{C}_{F, \mathcal{H}}} P(Fw)$$

These suffice to express several quantities

$$E[O_n(\mathcal{H})] = (n - m + 1) P(\mathcal{H}),$$

$$\begin{aligned} \text{Var}[O_n(\mathcal{H})] = & (n - m + 1) (P(\mathcal{H}) + (1 - 2m) P(\mathcal{H})^2 + 2 C(\mathcal{H})) \\ & + m(m - 1) P(\mathcal{H})^2 - 2 \tilde{C}(\mathcal{H}) \end{aligned}$$

$$\Pr\{O_n > 0\} \approx 1 - (1 + P(\mathcal{H}) - \tilde{C}(\mathcal{H}))^{-n}, \text{ if } n P(\mathcal{H}) < 1$$

Better than systems of functional equations **only if** we know how to compute these quantities **efficiently!**

A general method to compute parameters

Right complements are related to borders on words \implies They can be computed using a tree-like structure computed thanks to the **Aho-Corasick** algorithm with complexity $O(\sum_{w \in \mathcal{H}} |w|)$.

When \mathcal{H} is the set of words which are at Hamming distance k from a word H of length m , it remains $O(m^k)$, i.e., **exponential** with respect to k .

A general method to compute parameters

Right complements are related to borders on words \implies They can be computed using a tree-like structure computed thanks to the **Aho-Corasick** algorithm with complexity $O(\sum_{w \in \mathcal{H}} |w|)$.

When \mathcal{H} is the set of words which are at Hamming distance k from a word H of length m , it remains $O(m^k)$, i.e., **exponential** with respect to k .

A particular setting: positional pattern with errors

Alphabet $\Sigma (= \{A, C, G, T\})$

- Pattern $\mathcal{H} \in \Sigma^m \equiv (\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_m)$ where $\mathcal{H}_i \subseteq \Sigma$
- Neighborhood $\mathcal{N} = (\mathcal{N}_1, \mathcal{N}_2, \dots, \mathcal{N}_m)$ such that $\mathcal{H}_i \subseteq \mathcal{N}_i \subseteq \Sigma \rightarrow$ states the allowed errors at each position.

Example

$\mathcal{H} = (\{A, T\}, \{A\}, \{G\}, \{A\}, \{C\}),$

$\mathcal{N} = (\{A, C, G, T\}, \{A, T\}, \{G\}, \{A, T\}, \{C, G\})$

(Note that in IUPAC code , $\mathcal{H} = \text{WAGAC}$ and $\mathcal{N} = \text{NWGWS}$)

$\mathcal{H} = \text{AAGAC} + \text{TAGAC}$

$\mathcal{N} = \text{AAGAC} + \text{TAGAC} + \text{AAGAG} + \text{AAGTC} + \text{AAGTG} + \text{ATGAC} +$

$\text{ATGAG} + \text{ATGTC} + \text{ATGTG} + \text{CAGAC} + \text{CAGAG} + \text{CAGTC} + \text{CAGTG} +$

$\text{CTGAC} + \text{CTGAG} + \text{CTGTC} + \text{CTGTG} + \text{GAGAC} + \text{GAGAG} + \text{GAGTC} +$

$\text{GAGTG} + \text{GTGAC} + \text{GTGAG} + \text{GTGTC} + \text{GTGTG} + \text{TAGAC} + \text{TAGTC} +$

$\text{TAGTG} + \text{TTGAC} + \text{TTGAG} + \text{TTGTC} + \text{TTGTG}$

A particular setting: positional pattern with errors

Alphabet $\Sigma (= \{A, C, G, T\})$

- Pattern $\mathcal{H} \in \Sigma^m \equiv (\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_m)$ where $\mathcal{H}_i \subseteq \Sigma$
- Neighborhood $\mathcal{N} = (\mathcal{N}_1, \mathcal{N}_2, \dots, \mathcal{N}_m)$ such that $\mathcal{H}_i \subseteq \mathcal{N}_i \subseteq \Sigma \rightarrow$ states the allowed errors at each position.

Example

$\mathcal{H} = (\{A, T\}, \{A\}, \{G\}, \{A\}, \{C\})$,

$\mathcal{N} = (\{A, C, G, T\}, \{A, T\}, \{G\}, \{A, T\}, \{C, G\})$

(Note that in IUPAC code, $\mathcal{H} = \text{WAGAC}$ and $\mathcal{N} = \text{NWGWS}$)

$\mathcal{H} = \text{AAGAC} + \text{TAGAC}$

$\mathcal{N} = \text{AAGAC} + \text{TAGAC} + \text{AAGAG} + \text{AAGTC} + \text{AAGTG} + \text{ATGAC} +$

$\text{ATGAG} + \text{ATGTC} + \text{ATGTG} + \text{CAGAC} + \text{CAGAG} + \text{CAGTC} + \text{CAGTG} +$

$\text{CTGAC} + \text{CTGAG} + \text{CTGTC} + \text{CTGTG} + \text{GAGAC} + \text{GAGAG} + \text{GAGTC} +$

$\text{GAGTG} + \text{GTGAC} + \text{GTGAG} + \text{GTGTC} + \text{GTGTG} + \text{TAGAC} + \text{TAGTC} +$

$\text{TAGTG} + \text{TTGAC} + \text{TTGAG} + \text{TTGTC} + \text{TTGTG}$

A particular setting: positional pattern with errors

Alphabet $\Sigma (= \{A, C, G, T\})$

- Pattern $\mathcal{H} \in \Sigma^m \equiv (\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_m)$ where $\mathcal{H}_i \subseteq \Sigma$
- Neighborhood $\mathcal{N} = (\mathcal{N}_1, \mathcal{N}_2, \dots, \mathcal{N}_m)$ such that $\mathcal{H}_i \subseteq \mathcal{N}_i \subseteq \Sigma \rightarrow$ states the allowed errors at each position.

Example

$\mathcal{H} = (\{A, T\}, \{A\}, \{G\}, \{A\}, \{C\})$,

$\mathcal{N} = (\{A, C, G, T\}, \{A, T\}, \{G\}, \{A, T\}, \{C, G\})$

(Note that in IUPAC code , $\mathcal{H} = \text{WAGAC}$ and $\mathcal{N} = \text{NWGWS}$)

$\mathcal{H} = \text{AAGAC} + \text{TAGAC}$

$\mathcal{N} = \text{AAGAC} + \text{TAGAC} + \text{AAGAG} + \text{AAGTC} + \text{AAGTG} + \text{ATGAC} +$

$\text{ATGAG} + \text{ATGTC} + \text{ATGTG} + \text{CAGAC} + \text{CAGAG} + \text{CAGTC} + \text{CAGTG} +$

$\text{CTGAC} + \text{CTGAG} + \text{CTGTC} + \text{CTGTG} + \text{GAGAC} + \text{GAGAG} + \text{GAGTC} +$

$\text{GAGTG} + \text{GTGAC} + \text{GTGAG} + \text{GTGTC} + \text{GTGTG} + \text{TAGAG} + \text{TAGTC} +$

$\text{TAGTG} + \text{TTGAC} + \text{TTGAG} + \text{TTGTC} + \text{TTGTG}$

Computing probabilities

For these patterns, we easily compute

$$\sum_{\substack{w \in \mathcal{N} \\ d(w, \mathcal{H}) \leq k}} P(w), \text{ with } d(\cdot, \cdot) \text{ the Hamming distance}$$

- First idea: view sets of words as **formal series**, and **substitute** probabilities to symbols (just as affinity matrices)

$$P : \text{TAAGC} \mapsto p_T p_A p_A p_G p_C \quad (\text{Bernoulli})$$

$$\mapsto \sum_{i \in \Sigma} \pi_i p_{T|i} p_{A|T} p_{A|A} p_{G|A} p_{C|G} \quad (\text{Markov})$$

- **Mark** the errors: introduce a new variable u to count the number of errors with respect to \mathcal{H} .

Mark the errors

The variable u counts the number of errors.

Black: symbols of the patterns marked by $u^0 = 1$, Red: allowed error marked by $u^1 = u$.

Developing

$$(A + uC + uG + T)(A + uT)G(A + uT)(C + uG),$$

gives AAGAC + uAAGAG+ uAAGTC+ u²AAGTG+ uATGAC+ u²ATGAG+
 u²ATGTC+ u³ATGTG+ uCAGAC+ uCAGAG+ u²CAGTC+ u³CAGTG+
 u²CTGAC+ u³CTGAG+ u³CTGTC+ u⁴CTGTG+ uGAGAC+ u²GAGAG+
 u²GAGTC+ u³GAGTG+ u²GTGAC+ u³GTGAG+ u³GTGTC+ u⁴GTGTG+
 TAGAC+ uTAGAG+ uTAGTC+ u²TAGTG+ uTTGAC+ u²TTGAG+
 u²TTGTC +u³TTGTG

Mark the errors

The variable u counts the number of errors.

Black: symbols of the patterns marked by $u^0 = 1$, Red: allowed error marked by $u^1 = u$.

Developing

$$(A + uC + uG + T)(A + uT)G(A + uT)(C + uG),$$

gives AAGAC + uAAGAG+ uAAGTC+ u²AAGTG+ uATGAC+ u²ATGAG+
 u²ATGTC+ u³ATGTG+ uCAGAC+ uCAGAG+ u²CAGTC+ u³CAGTG+
 u²CTGAC+ u³CTGAG+ u³CTGTC+ u⁴CTGTG+ uGAGAC+ u²GAGAG+
 u²GAGTC+ u³GAGTG+ u²GTGAC+ u³GTGAG+ u³GTGTC+ u⁴GTGTG+
 TAGAC+ uTAGAG+ uTAGTC+ u²TAGTG+ uTTGAC+ u²TTGAG+
 u²TTGTC +u³TTGTG

Mark the errors

The variable u counts the number of errors.

Black: symbols of the patterns marked by $u^0 = 1$, Red: allowed error marked by $u^1 = u$.

Developing

$$(A + uC + uG + T)(A + uT)G(A + uT)(C + uG),$$

gives AAGAC + uAAGAG+ uAAGTC+ u²AAGTG+ uATGAC+ u²ATGAG+
 u²ATGTC+ uCAGAC+ uCAGAG+ u²CAGTC+
 u²CTGAC+ uGAGAC+ u²GAGAG+
 u²GAGTC+ u²GTGAC+
 TAGAC+ uTAGAG+ uTAGTC+ u²TAGTG+ uTTGAC+ u²TTGAG+
 u²TTGTC

Mark the errors

The variable u counts the number of errors.

Black: symbols of the patterns marked by $u^0 = 1$, Red: allowed error marked by $u^1 = u$.

Developing

$$(A + uC + uG + T)(A + uT)G(A + uT)(C + uG),$$

gives AAGAC + uAAGAG+ uAAGTC+ u²AAGTG+ uATGAC+ u²ATGAG+
 u²ATGTC+ uCAGAC+ uCAGAG+ u²CAGTC+
 u²CTGAC+ uGAGAC+ u²GAGAG+
 u²GAGTC+ u²GTGAC+
 TAGAC+ uTAGAG+ uTAGTC+ u²TAGTG+ uTTGAC+ u²TTGAG+
 u²TTGTC

- k errors max → remove words with more than $k = 2$ red letters
- truncate polynomials in u at order k on the degree

Computing the probability (end)

Of course we don't want to develop! We compute **iteratively** the probabilities by considering successive positions and using **truncated** polynomials in u .

When the number of errors k is **fixed**, we have an algorithm in time $O(mk)$ to compute the probability ($m \times$ cost of multiplying a polynomial of degree k by a monomial of degree 1).

The same principles can be extended to a **Markov** model and the **others parameters** related to the overlapping structure.

Computing the probability (end)

Of course we don't want to develop! We compute **iteratively** the probabilities by considering successive positions and using **truncated** polynomials in u .

When the number of errors k is **fixed**, we have an algorithm in time $O(mk)$ to compute the probability ($m \times$ cost of multiplying a polynomial of degree k by a monomial of degree 1).

The same principles can be extended to a **Markov** model and the **others parameters** related to the overlapping structure.

Computing the probability (end)

Of course we don't want to develop! We compute **iteratively** the probabilities by considering successive positions and using **truncated** polynomials in u .

When the number of errors k is **fixed**, we have an algorithm in time $O(mk)$ to compute the probability ($m \times$ cost of multiplying a polynomial of degree k by a monomial of degree 1).

The same principles can be extended to a **Markov** model and the **others parameters** related to the overlapping structure.

Conclusion

We have rewritten exact matricial expressions in a suitable form depending on a few combinatorial parameters which are easy to compute.

Perspectives

- Extension to dyadic – or structured – patterns (M.-F. Sagot and L. Marsan), palindromic patterns, highly repetitive patterns
- Conditional occurrence problem (artifacts)