# SNAP: HOW SCORES ARE COMPUTED

AVIAD ROZENHEK

## Contents

## 1. Basic notation

This document describes how SeedSearcher scores motifs, and consists mainly of definitions, which are mathematical in nature, but quite simple.

**Definition: Sequence.** A *sequence* $S$ is a finite vector of characters from an alphabet which we shall keep fixed. The j-th coordinate $S_j$ of the sequence shall be called the j-th position. We shall identify $S$ with the set of all positions of $S$, thus we can (ab)use notation like $S_j \in S$.

Usually we work with gene promotor sequences, where a position is a downstream offset.

**Definition: Database.** An ordered set $G$ of sequences shall be called a *database.* We shall think of $G$ as fixed, with cardinality $N$ = number of sequences. Thus every sequence $S$ shall be denoted by $S^i \in G$ for some $i$.

Usually a database is a collection of named gene promoter sequences. We are interested in scoring different "motifs" against the database. We shall identify a motif with the set of positions in the database, where the motif 'resides' (usually the starting position of the motif). We shall calls these sets Clusters:

**Definition: PCluster.** An *i-Position Cluster* $C^i$ is a set of positions, all belonging to the i-th sequence: $C^i \subseteq S^i$. We shall denote by $\bar{C}^i = \{C^i : C^i \subseteq S^i\}$ the set of all i-Position Clusters.

**Definition: Cluster.** An ordered set (or vector) $C$ of position clusters shall be called a *Cluster.* $C \in \prod_i \bar{C}^i$. that is, a cluster contains a (possibly empty) pcluster for every sequence in the database. We shall denote by $\bar{C} = \{C : C \in \prod_i \bar{C}^i\}$ the set of all Clusters.

We would like to associate Clusters with Scores. A score system is be a mapping of Clusters to classifications. A score function is a mapping of classifications to scores. Formally:

**Definition.**

A *score system* is a mapping $\psi : \bar{C} \rightarrow (\, T_p \ T_n \ F_p \ F_n \,)$
A *score function* is a mapping $\pi : (\, T_p \ T_n \ F_p \ F_n \,) \rightarrow \Re$

Where

$T_p =$ True Positives      $=$ number of correct positive classifications

$T_n =$ True Negatives      $=$ number if correct negative classifications

$F_p =$ False Positives      $=$ number of incorrect positive classifications

$F_n =$ False Negatives      $=$ number of incorrect negative classifications

We can now define the two score functions which SeedSearcher employs:

## 2. Score functions

### 2.1. Hyper-Geometric scoring function.

We define the *Hyper Geometric Pvalue* $P_{hg}$ to be:

$$(1) \qquad P_{hg} (\, T_p \ T_n \ F_p \ F_n \,) = \sum_{x \geq k} D_{hg} (\, x \ n \ K \ N \,)$$

Where

$$(2) \qquad \begin{pmatrix} x \\ n \\ K \\ N \end{pmatrix} = \begin{pmatrix} T_p \\ T_p + F_p \\ T_p + F_n \\ T_p + F_p + T_n + F_n \end{pmatrix}$$

and $D_{hg}$ is the hyper-geometric distribution:

$$(3) \qquad D_{hg}(x, k, K, N) = \frac{\begin{pmatrix} K \\ x \end{pmatrix} \begin{pmatrix} N - K \\ n - x \end{pmatrix}}{\begin{pmatrix} N \\ n \end{pmatrix}}$$

Thus $P_{hg}$ is a pvalue and equals the sum of the tail of the hyper-geometric distribution.

The distribution $D_{hg}$ is more easily understood using the following scenario: Let there be an urn with $N$ balls, $K$ of which are *red* and suppose we draw $n$ balls from the urn, $x$ of those are red. then $H(x, n, K, N)$ is the chance to draw $n$ balls and that *at least* $x$ of these will be red.

*Remark.* Naturally, values of $P_{hg}$ are numbers $\in [0, 1]$ and can be very close to 0. So for both numerical stability and ease of use we shall define and use:

$$(4) \qquad h (\, T_p \ T_n \ F_p \ F_n \,) = - \log_{10} Pval_{hg} (\, T_p \ T_n \ F_p \ F_n \,)$$

This h shall be called the *hyper-geometric scoring* function

## 2.2. Exploss scoring function.

This is a much simpler score function, which only "rewards" (or "punishes") according to correct (or incorrect) positive classifications: The score is increased when a correct positive classification is made ($T_p$), and is decreased when an incorrect positive classification is made ($F_p$). negative classifications are ignored ($T_n$, $F_n$). This is still a discriminative function because it relies on classifications of both the positive and negative sets, unlike non-discriminate functions which will rely only on classifications of the positive set, e.g. $T_p$ and $F_n$.

## Definition.

Let $1 < \alpha, \beta$ be any fixed constansts. we define the exploss function $E$ to be:

$$E_{\alpha,\beta}(T_p, F_p) = \frac{\beta^{F_p}}{\alpha^{T_p}} \tag{5}$$

This function receives values $\in (0, \infty)$ and it is easily seen that the better the positive classification, the closer the score is to 0. so we again take:

$$e_{\alpha,\beta}(T_p, T_n, F_p, F_n) = -\log_{10} Exploss_{\alpha,\beta}(T_p, F_p) \tag{6}$$

This $e$ shall be called the *exploss scoring* function

## 3. SCORE SYSTEMS

We shall now recall the definitions in section 1, and define several ways to map Clusters to ($T_p$ $T_n$ $F_p$ $F_n$) classifications. We shall assume to have a constant database of size N, with weights $W^i \in [0, 1]$ on the sequences $s^i$ $i = 1...N$ and also positional weights $W_j^i \in [0, 1]$ on any position in any sequence.

*Remark.* In the simplest case, the prior on the regulation of $s^i$ is discrete, and there is no positional weight prior. thus:

$$\forall i \quad W^i = \begin{cases} 1 & \text{if The sequence } s^i \text{ is positively labeled} \\ 0 & \text{if The sequence } s^i \text{ is negatively labeled} \end{cases}$$

and

$$\forall i \forall j \quad W_j^i = 1$$

To gain some intuition, let us retreat from the abstract and focus our attention on the specific problem which we are trying to solve. We are trying to find a motif (short IUPAC sequence) which will correctly separate a database of gene promoters (sequences) into regulated (positively labeled) and non-regulated (negatively labeled) sets, which are known in advance. Now suppose we have such a motif $m$ of length $l$. $m$ infers a Cluster $c$ which contains all the positions in the promoters where $m$ starts. we are interested in defining several $\rho$ classification functions such that:

$$\rho(c) = \begin{pmatrix} T_p^\rho(c) \\ T_n^\rho(c) \\ F_p^\rho(c) \\ F_n^\rho(c) \end{pmatrix} = \sum_i \rho_i(c^i) = \sum_i \begin{pmatrix} T_p^{\rho_i}(c^i) \\ T_n^{\rho_i}(c^i) \\ F_p^{\rho_i}(c^i) \\ F_n^{\rho_i}(c^i) \end{pmatrix}$$

Thus a classification of a Cluster $c$ is the sum of the classifications on $c^i$ PClusters. How do we classify a PCluster? there are a couple of questions involved, with several good answers:

- $\omega$ *pweight functions:* how much does a position $p_j^i$ (the j-th coordinate of sequence i) weigh? we can define several different weight functions on positions. we call these pweight functions, and they are denoted with the letter $\omega$
    - *Discrete pweight:* given a threshold $t$ (usually $t = 0.5$) we define the $\omega_d^t$ discrete pweight function with threshold $t$ function to be:

$$\omega_d^t(p_j^i) = \begin{cases} 1 & \text{if } t < W^i \text{ (The sequence is positively labeled)} \\ 0 & \text{if } W^i <= t \text{ (The sequence is negatively labeled)} \end{cases}$$

    - We define the *real pweight* $\omega_r$ function to be:

$$\omega_r(p_j^i) = W^i$$

      Using this function a position weighs the same as the sequence it belongs to.
    - We define the positional pweight $\omega_p$ function to be:

$$\omega_p(p_j^i) = W^i \cdot W_j^i$$

      Using this function the prior on the coordinates is taken into account.
- $\rho$ *pcluster classification functions:* How do we handle several classifications of same the same sequence? that is, how do we classify a $c^i$ cluster which contains several positions $p_{j_1}^i, ..., p_{j_n}^i$ all belonging to the same sequence $s^i$? given some pweight function $\omega$ we can define several $\rho_\omega$ *pcluster classification* functions:
    - We define the $\rho^1$ *single position classification* function to be:

$$\rho_\omega^1(c^i) = \begin{pmatrix} \max\limits_{p_j^i \in c^i} \{\omega(p_j^i)\} \\ 1 - \max\limits_{p_j^i \in c^i} \{\omega(p_j^i)\} \end{pmatrix}$$

      Using this function only the "best" position in each pcluster is taken into accout.
    - We define the $\rho^\infty$ *total positions classification* function to be:

$$\rho_\omega^\infty(c^i) = \begin{pmatrix} \sum\limits_{p_j^i \in c^i} \omega(p_j^i) \\ \sum\limits_{p_j^i \in c^i} 1 - \omega(p_j^i) \end{pmatrix}$$

      Using this function all the positions of a pcluster are taken into account.

Let $\omega$ be a pweight function $\in \{\omega_d^t, \omega_r, \omega_p\}$
and let $\rho$ be a counting function $\in \{\rho_\omega^1, \rho_\omega^\infty\}$. We shall define the *score system $\rho$* to be

$$\begin{pmatrix} T_p^\rho(c) \\ F_p^\rho(c) \end{pmatrix} = \sum_i \begin{pmatrix} T_p^\rho(c^i) \\ F_p^\rho(c^i) \end{pmatrix} = \sum_i \rho(c^i)$$

$$\begin{pmatrix} F_n^\rho(c) \\ T_n^\rho(c) \end{pmatrix} = \sum_i \begin{pmatrix} F_n^\rho(c^i) \\ T_n^\rho(c^i) \end{pmatrix}$$

$$= \sum_i \left( \rho(s^i) - \rho(c^i) \right)$$

$$= \sum_i \left( \rho(s^i) - \begin{pmatrix} T_p^\rho(c^i) \\ F_p^\rho(c^i) \end{pmatrix} \right)$$

$$= \left( \sum_i \rho(s^i) \right) - \begin{pmatrix} T_p^\rho(c) \\ F_p^\rho(c) \end{pmatrix}$$

We have successfully defined several scoring systems, which map clusters into $\begin{pmatrix} T_p & T_n & F_p & F_n \end{pmatrix}$ classifications. Summing up our method: we classify each pcluster using either the "best" position in the cluster (single position classification) or using all the positions in the cluster (total positions classification). each position is weighed using a discrete/real/positional weight function. this system maps a motif (with its corresponding cluster) to a $\begin{pmatrix} T_p & T_n & F_p & F_n \end{pmatrix}$ classification of the database. We use either an exploss or hyper-geometrical score function to map these classification vectors into real-valued scores.

## 4. REMARKS

Several issues still need further refining:

4.1. **removing overlaps in total position classifications.** When pclusters reflect the starting positions of a motif $m$ with length $l$ it often make sense not to allow overlaping positions. position $p_j^i$ overlaps with $p_k^i$ iff (assuming $w.l.g \quad j < k$) $k \in [j, j+l)$. Thus $\rho(s^i)$ now stands for a classification of the maximum overlap-free set of positions in $s^i$.

4.2. **DNA sequences: using the reverse strand.** When the sequences in concern are DNA, it often makes sense to "look" for motifs in both the positive and reverse strands. This affects scoring in a simple way: The reverse strands are not added as sequences, but rather they have the effect of doubling the length of all sequences. However, when a cluster $c^i$ (or even $s^i$ for that matter) is cleaned off overlaps the overlaps are removed separately from the positive and reverse strands.

4.3. **future changes.** The positional pweight function will change to reflect the positional bias along the entire motif (and not just the starting position)