

Ranking of Genes, SNVs, and Sequence Regions

Ranking elements within various types of biosets for metaanalysis of genetic data.

Introduction

One of the primary applications of the BaseSpace® Correlation Engine is to allow researchers to perform metaanalyses that harness large amounts of genomic, epigenetic, proteomic, and assay data. Such analyses look for potentially novel and interesting results that cannot necessarily be seen by looking at a single existing experiment. These results may be interesting in themselves (eg, associations between different treatment factors, or between a treatment and an existing known pathway or protein family), or they may be used to guide further research and experimentation.

The primary entity in these analyses is the **bioset**. It is a ranked list of elements (genes, probes, proteins, compounds, single-nucleotide variants [SNVs], sequence regions, etc.) that corresponds to a given treatment or condition in an experiment, an assay, or a single patient sample (eg, mutations). For a gene expression experiment, the biosets will consist of gene lists with associated change values and statistical information for each relevant experimental factor. For example, a bioset could consist of a list of Affymetrix probesets, fold-change difference values between treated and control groups, and p-values produced by a t test between treated and control sample probeset intensity values. A sequence-centric bioset from a copy-number analysis experiment, for example, will contain ranked sequence regions with associated gain and loss statistics.

In summary, biosets may contain many of the following columns:

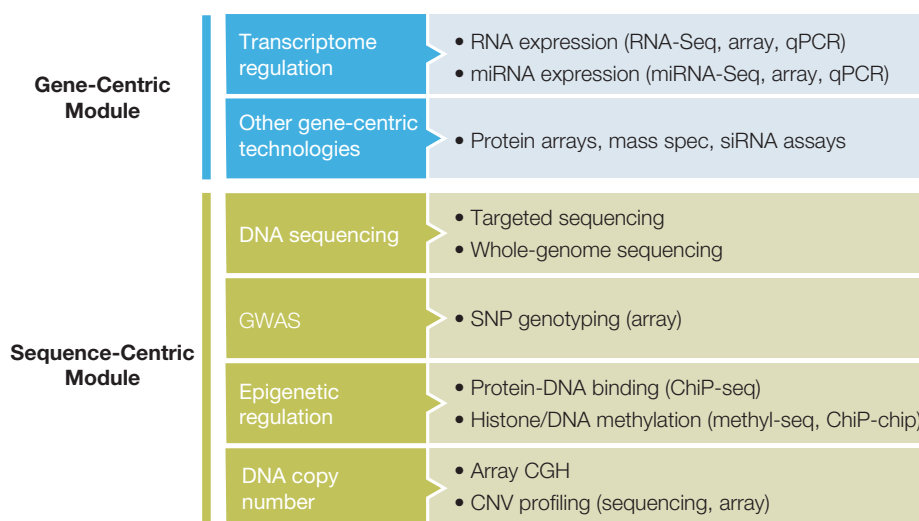
- Identifier of an entity such as a gene, SNV, or sequence region (required)
- Other identifiers of the entity—eg, chromosome, position
- Summary statistics—eg, p-value, fold change, score, rank, odds ratio

Ranking of Elements within a Bioset

Most biosets generated at Illumina include elements that are changing relative to a reference genome (mutations) or due to a treatment (or some other test factor) along with a corresponding **rank** and **directionality**. Typically, the rank will be based on the magnitude of change (eg, fold change); however, other values, including p-values, can be used for this ranking. Directionality is determined from the sign of the statistic: eg, up (+) or down(-) regulation or copy-number gain (+) or loss (-). Some biosets contain signatures from individual samples where the elements may or may not be ranked.

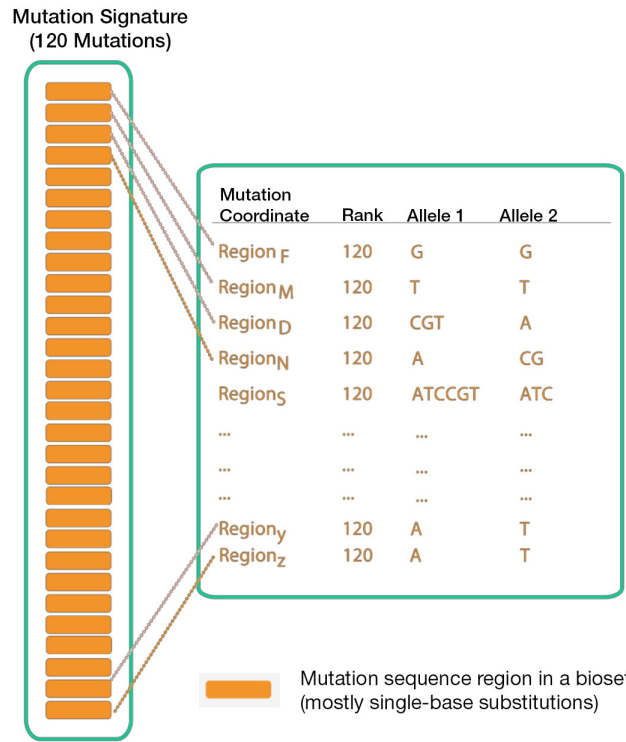
Where applicable, elements within biosets that are generated by Illumina scientists internally are ranked and assigned directionality using predetermined criteria. However, for user-imported biosets, the statistics/values to be used for ranking can be specified during upload.

Figure 1: Overview of Biosets



Biosets within the BaseSpace Correlation Engine comprise gene-centric or sequence-centric signatures from orthogonal types of biological experiments.

Figure 4: SNV Bioset from Mutation Data



An example of an SNV bioset comprising mutation data from an individual sample, with associated statistics.

Mapping

SNVs are mapped to genes and a gene-centric bioset is also created corresponding to each SNV bioset. Genes carry the same rank and normalized rank as the SNVs to which they map.

Illumina is investigating weighted ranking of SNVs by mutation significance, eg, location of SNV (promoter region, intergenic), and impact on protein sequence (nonsense mutation, missense mutation vs. silent mutation).

SNV Biosets from Mutation and Resequencing Data

Mutation and resequencing biosets from individual samples typically contain the mutation identifier and location, and the alleles. The ranking scheme followed is the same as that described for GWAS data. The typical case is when SNVs have no associated p-values or other ranking columns.

Directionality

Elements in an SNV bioset are not assigned a direction.

Biosets from Experiments Investigating Copy Numbers

Figure 5 shows a bioset generated from an analysis of copy numbers in a cell line.

Ranking Criteria

Regions are sorted according to the ranking column (eg, z score) and “raw ranks” are assigned sequentially to each region (Figure 5).

Next, raw ranks are converted to ranks and are assigned to each sequence based on cumulative sequence length and unit region size as follows:

1. For each sequence region, set cumulatedLength to be the sum of sequence lengths of all the regions with better or the same raw rank.
2. Set rank of the region to cumulatedLength ÷ unitRegion, rounded off to the closest integer.
3. Normalized rank is computed as in Equation 4, using the custom platform size and the normalizing constant.

illumina • 1.800.809.4566 toll-free (U.S.) • +1.858.202.4566 tel • techsupport@illumina.com • www.illumina.com

For Research Use Only. Not for use in diagnostic procedures.

© 2016 Illumina, Inc. All rights reserved. Illumina, BaseSpace, Infinium, and the pumpkin orange color are trademarks of Illumina Inc., and its affiliate(s) in the U.S. and/or in other countries. Pub. No. 970-2014-009 Current as of 28 March 2016.



AGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAAC
TCAACGTACCGTAAACGAACGATCATTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTGCAACGAACGAAAAGAAATGATAACAGTAAACACACTTCTGTTAACCTT
CGACGAAAAGAAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAAGATTAATTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTAC
TAACGTACCATTAAAGAGCTACCGTGCAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTGCAACGAACGAAAAGAAATGATA
ACGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAAC
GATTACTTGATCCACTGATTCAACGTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAACGAC
CGTATCAATTGAGACTAAATATAACGTACCATTAAAGATTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTGCAACGAAAAGAAATGATAACAG