

Jetset: selecting an optimal Affymetrix probe set to represent a gene

Qiyuan Li, Aron C Eklund

August 22, 2011

Contents

1	Introduction	1
2	Contents of the package	2
2.1	jmap	2
2.2	jscores	2
2.3	scores.hgu95av2	2
2.4	eg2ps.hgu95av2	3
3	Simple probe set selection	4
3.1	Selecting probe sets by Entrez Gene ID	4
3.2	Selecting probe sets by ensembl ID	4
3.3	Selecting by gene symbol	4
3.4	Selecting by gene alias	5
4	Probe set quality scores	6
5	Jetset algorithm details	7
6	Reference	8
7	R sessionInfo	8

1 Introduction

On Affymetrix gene expression microarrays, a given gene may be detected by multiple probe sets which may deliver inconsistent or even contradictory measurements. Therefore, obtaining an unambiguous expression estimate of a pre-specified gene can be nontrivial. We developed scoring methods to assess each probe set for specificity, coverage, and degradation resistance. We used these scores to select the optimal probe set for each gene, and thus create a simple one-to-one mapping between gene and probe set.

`jetset` is a package enabling the selection of optimal probe sets from the HG-U95Av2, HG-U133A, HG-U133 Plus 2.0, or U133 X3P microarray platforms.

```
> library(jetset)
```

2 Contents of the package

The `jetset` package contains the following objects:

```
> ls("package:jetset")
[1] "eg2ps.hgu133a"      "eg2ps.hgu133plus2"  "eg2ps.hgu95av2"
[4] "eg2ps.u133x3p"      "jmap"                "jscores"
[7] "scores.hgu133a"     "scores.hgu133plus2" "scores.hgu95av2"
[10] "scores.u133x3p"
```

The functions `jmap` and `jscores` are the intended user-level interface, and `scores.*` and `eg2ps.*` are data sets that support these functions.

In this document, we describe only one (`*.hgu95av2`) of each family of data sets.

2.1 jmap

`jmap` is a function that returns the best probe sets matching a list of Entrez GeneIDs, gene symbols, or gene aliases.

2.2 jscores

`jscores` is a function that returns the `jetset` scores for all probe sets matching a list of Entrez GeneIDs, gene symbols, aliases, or ensembl IDs.

2.3 scores.hgu95av2

`scores.hgu95av2` is a data frame with Entrez IDs and pre-calculated quality control scores for each probe set ID. All scores range from 0 to 1, and a higher score indicates better (predicted) performance.

```
> head(scores.hgu95av2)
```

	nProbes	EntrezID	process	specificity	coverage
1000_at	16	5595	222.0	0.6250	1
1001_at	16	7075	816.0	0.8750	1
1002_f_at	16	<NA>	NA	NA	NA
1003_s_at	16	643	1963.5	0.6875	1
1004_at	16	643	329.0	0.5625	1
1005_at	16	1843	941.0	0.6875	1

- The Entrez GeneID (*EntrezID*) is a unique gene identifier. Note: as in other Bioconductor packages, the GeneID is stored as type *character*.
- The processivity requirement (*process*) is the number of consecutive bases that must be synthesized to generate a target that can be detected by the probe set.
- The *specificity* score is the fraction of the probes in a probe set that are likely to detect the targeted gene and unlikely to detect other genes.
- The *coverage* score is the fraction of the splice isoforms belonging to the targeted gene that are detected by the probe set.

Note that the robustness score and overall score are not stored in this data; instead these scores are calculated on-the-fly when the score data is retrieved using `jscores`.

- The robustness score (*robust*) is intended to quantify robustness against transcript degradation. The robustness score uses the processivity requirement to estimate the signal intensity of a probe set, relative to the ideal case of perfect processivity.
- The *overall* score is the product of the specificity score, coverage score, and robustness score.

2.4 eg2ps.hgu95av2

eg2ps.hgu95av2 is an character vector, of which the values are probe set IDs and the names are Entrez GeneIDs.

```
> head(eg2ps.hgu95av2)
```

10	100	1000
"38912_at"	"41654_at"	"2053_at"
100008588	100008589	10001
"AFFX-HUMRGE/M10098_3_at"	"AFFX-M27830_5_at"	"35432_at"

3 Simple probe set selection

A typical application of the `jetset` packages is to identify probe sets corresponding to a published list of genes.

3.1 Selecting probe sets by Entrez Gene ID

The Entrez GeneID is an unambiguous way to specify a gene; however the numeric ID is not particularly descriptive, which may explain why this is not commonly provided in publications.

```
> jmap("hgu95av2", eg = "2099")  
  
2099  
"1681_at"  
  
> jmap("hgu133a", eg = "2099")  
  
2099  
"205225_at"  
  
> jmap("hgu133plus2", eg = "2099")  
  
2099  
"205225_at"  
  
> jmap("u133x3p", eg = "2099")  
  
2099  
"g4503602_3p_at"
```

Entrez GeneID 2099 corresponds to the estrogen receptor (ESR1) gene, which is an important indicator of breast cancer phenotype.

3.2 Selecting probe sets by ensembl ID

The ensembl ID is another unambiguous way to specify a gene.

```
> jmap("hgu95av2", ensembl = "ENSG00000091831")  
  
ENSG00000091831  
"1681_at"  
  
> jmap("hgu133a", ensembl = "ENSG00000091831")  
  
ENSG00000091831  
"205225_at"
```

3.3 Selecting by gene symbol

Often, we know the official HUGO gene symbols, and want the corresponding probe sets.

```
> jmap("hgu133a", symbol = c("ESR1", "ERBB2", "AURKA"))  
  
ESR1      ERBB2      AURKA  
"205225_at" "216836_s_at" "208079_s_at"
```

Unfortunately, the gene symbol for a given gene can change. Furthermore, in rare cases a HUGO gene symbol can correspond to two distinct genes.

3.4 Selecting by gene alias

If we have gene symbols, but they are not the *official* symbols, the gene aliases might be useful.

```
> jmap("u133x3p", alias = c("P53", "HER-2", "K-RAS"))
```

	P53	HER-2	K-RAS
	"g8400737_3p_at"	"Hs.323910.2.A1_3p_a_at"	"Hs.37003.0.S1_3p_x_at"

4 Probe set quality scores

We might want to compare quality scores for all probe sets corresponding to a gene of interest. For this example, the STAT1 gene is used because it is detected by several probe sets.

```
> jscores("hgu95av2", symbol = "STAT1")
```

	nProbes	EntrezID	process	specificity	coverage
32859_at	16	6772	74.5	0.4375	0.5
32860_g_at	16	6772	405.5	0.7500	0.5
33338_at	16	6772	168.5	0.7500	0.5
33339_g_at	16	6772	1203.0	0.7500	1.0
AFFX-HUMISGF3A/M97935_3_at	20	6772	451.0	0.7500	0.5
AFFX-HUMISGF3A/M97935_5_at	20	6772	2863.5	0.7500	1.0
AFFX-HUMISGF3A/M97935_MA_at	20	6772	1773.5	0.8500	1.0
AFFX-HUMISGF3A/M97935_MB_at	20	6772	1536.0	0.6500	0.5
	robust	overall	symbol		
32859_at	0.883141136	0.193187124	STAT1		
32860_g_at	0.508445547	0.190667080	STAT1		
33338_at	0.754977281	0.283116480	STAT1		
33339_g_at	0.134435239	0.100826429	STAT1		
AFFX-HUMISGF3A/M97935_3_at	0.471284303	0.176731613	STAT1		
AFFX-HUMISGF3A/M97935_5_at	0.008425592	0.006319194	STAT1		
AFFX-HUMISGF3A/M97935_MA_at	0.051907122	0.044121054	STAT1		
AFFX-HUMISGF3A/M97935_MB_at	0.077139816	0.025070440	STAT1		

We can confirm that jmap returns the probe set with the highest overall score:

```
> jmap("hgu95av2", symbol = "STAT1")
```

```
STAT1  
"33338_at"
```

5 Jetset algorithm details

We downloaded probe sequences corresponding to four human gene expression microarrays from Affymetrix: U95Av2, U133A, U133 Plus 2.0, and X3P. We used NCBI BLASTN to search the 25-base probe sequences for matches to the Refseq human RNA database (Pruitt, et al., 2005). The BLASTN search was run with the default parameters, except that filtering was turned off, the word size was set to 8 to increase sensitivity, and the expectation value was set to 1 to reduce output size.

```
blastall -p blastn -d refseq.human.rna -i probe.hgu133a.fa -o hgu133a.refseq.20110817.bls  
-F F -m 8 -e 1 -W 8 -a 8
```

We used the alignment score (bit score) between each probe and cDNA as an indication of probe sensitivity. We defined three levels of alignment: a strong alignment has a score between 48 and 51, indicating that at least 24 bases are identical and that the probe is very likely to detect the target. A moderate alignment has a score between 32 and 47, corresponding to an uninterrupted alignment of length 16 to 23 bases; the probe may or may not respond to the target. A weak alignment has a score less than 32 and is unlikely to respond to the target.

Specificity. A probe was considered to specifically detect a given gene if it aligned strongly to at least one transcript of the gene, but did not have a strong or moderate alignment to a transcript from another gene. The gene specifically detected by the largest number of probes in a probe set was considered the targeted gene of the probe set. We defined the specificity score S_s of a probe set as the fraction of its probes that specifically detect the targeted gene.

Coverage. A transcript of the targeted gene was considered detected by a probe set if the transcript has a strong alignment to the majority of the probes in the probe set. The coverage score S_c of a probe set is defined as the fraction of all transcripts belonging to the targeted gene that are detected by the probe set.

Robustness. The processivity requirement for a probe-transcript alignment is the number of bases between the 5' end of the alignment and the 3' end of the transcript sequence; this corresponds to the length of labeled target that must be synthesized by in vitro transcription to reach the query region. The overall processivity requirement N of a probe set is the median processivity requirement for all strong alignments between probes in the probe set and transcripts in the targeted gene. We define the robustness score S_r of a probe set as the probability that synthesis of the target up to the processivity requirement is achieved without interruption: $S_r = (1 - p)^n$

Here, p is the probability of the IVT synthesis being interrupted at each base, due to either transcript degradation or lack of enzyme processivity. The value of p is likely to be variable in clinical specimens, but for simplicity we use a value corresponding to the manufacturer's design criteria: 1/300 for the X3P array, or 1/600 for the other arrays.

Overall score. We define the overall score S_o as the product of the three scores described above: $S_o = S_s * S_c * S_r$

For a given gene, the probe set targeting this gene with the highest overall score is selected to represent the gene.

6 Reference

Qiyuan Li, Nicolai J. Birkbak, Balazs Gyorffy, Zoltan Szallasi and Aron C. Eklund (2010). Jetset: selecting an optimal microarray probe set to represent a gene. Manuscript in preparation.

7 R sessionInfo

The results in this file were generated using the following packages:

```
> sessionInfo()

R version 2.13.1 (2011-07-08)
Platform: x86_64-apple-darwin9.8.0/x86_64 (64-bit)

locale:
[1] C/en_US.UTF-8/C/C/en_US.UTF-8/en_US.UTF-8

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] jetset_0.99.1 RSQLite_0.9-4 DBI_0.2-5

loaded via a namespace (and not attached):
[1] AnnotationDbi_1.14.1 Biobase_2.12.1      org.Hs.eg.db_2.5.0
[4] tools_2.13.1
```