

sscu user manual (2.4.0)

Yu Sun

2016-12-07

Contents

1. Background	1
1.1 S index	2
1.2 Akashi's test	2
1.3 Identify optimal codons	2
1.4 Other functions	2
2. Function overview	2
2.1 akashi_test	2
2.2 genomic_gc3	3
2.3 op_corre_codonW	3
2.4 op_corre_NCprime	3
2.5 op_highly	3
2.6 op_highly_stats	4
2.7 low_frequency_op	4
2.8 s_index	5
3 Workflow	5
3.1 general statistics	5
3.2 optimal codons	6
3.3 selective profile	7
3.4 a function to identify the low frequency optimal codons	8
4 Folders installed in the sscu package	9
4.1 sequences	9
4.2 akashi_test	9
4.3 correlative_test	9

1. Background

The central dogma is the cornerstone of modern molecular biology, as it explains the information processing flow within a biological system is in the direction of DNA to RNA to protein (crick 1970). An important step in this process is to decode the information stored in the cryptic nucleotides sequences into protein molecules. Sixty-four tri-nucleotide units (codons) are correspond to the coding of twenty amino acids, thus almost all the amino acids has more than one matching codons. The codon choice for an amino acid is not random for most species, and it has been shown that mutation and selection are the two major forces shaping the codon usage pattern. For the majority of the gene, mutation bias is the factor determining the codon choice, but the influence from selection is increased with the genes' expression level.

Several programs and packages has been release for the codon usage analysis, such as CodonW. However, most of these programs focus on the overall usage of codons, more specifcly, the mutational influence on the codons usage. Many theorical work has been developed to disentangle the effect of selection on codon usage, such as Sharp's S index, Akashi test, as well as the optimal codons identification by highly expressed genes and correlative method. However, no program or pipeline has been released to do these calculations, thus researchers may either give up on these test, or write their individual code to perform these task. But it is

not only takes a lot of time, but also has the possibility to introduce errors during the coding process. In this R package, we mainly focus on developing functions for identify selection on codon usage, and included functions for the following tests.

1.1 S index

Sharp released a mathematic formula, named as S index, to quantify the strength of selection in the synonymous codon. The method calculates the relative proportion of C- and U-ending codons in the two codon boxes, and also takes into consideration of the genomic mutation bias. The method has been proved to be an effective way to quantify and compare the codon selection among different species.

1.2 Akashi's test

The function calculate Akashi's test for translational accuracy selection on coding sequences. Akashi proposed the translational accuracy theory for codon usage in 1994. The theory suggest that the optimal codons are codons that translated more accurately than the other codons, and these codons are favored in the important sites, such as the evolutionary conserved amino acid sites, whereas the less conserved amino acids sites are more tolerable to the non-optimal codons. For detailed information, see reference (Akashi H. Synonymous codon usage in *Drosophila melanogaster*: natural selection and translational accuracy. *Genetics* 1994 Mar;136(3):927-35.) and Drummond sites.

1.3 Identify optimal codons

Two major method has been used to identify optimal codons: highly expressed gene method (short as highly method) and correlative method. The highly method compare the codon usage in the highly and lowly/all genes. The basic idea is that the highly expressed genes are enriched with more optimal codons, due to the translational efficiency or accuracy reasons. The correlative method identifies the optimal codons based on the method from Hershberg & Petrov (Hershberg R, Petrov DA. 2009. General rules for optimal codon choice. *Plos Genet.* 5:e1001115.). This method take the whole genome into consideration, and predict the optimal codons by making the correlation between the frequency of each codon within each gene and the overall codon bias (N_c or N_c').

1.4 Other functions

Other functions include genomic gc3 calculation, optimal codon statistics and low frequency optimal codon identification.

Overall, the package is useful for the calculation of selective profile in codon usage studies, and it is a good complement to other major codon usage programs, such as CodonW.

2. Function overview

2.1 akashi_test

```
akashi_test(contingency_file=NULL)
```

The function calculate Akashi's test for translational accuracy selection on coding sequences.

The input of the function is a contingency file, you can find an example in the folder `akashi_test` in the `sscu` directory. You can either make the file by yourself, or by the perl script `make_contingency_table.pl` in the

akashi_test folder. You can check the detailed usage of the perl script by reading the first few lines in the perl script. The perl script tabulates and calculate the a,b,c,d entries Drummond sites for the coding sequences, and output the four values for each amino acid for each gene, example as the contingency_file_Gvag. The input of the perl script is a folder contains the codon alignments of the genes that you are interested to calculate.

The function reads and calculates the Z value, p value and odd ratio for the Akashi's test. In addition, it calculates the conserved, variable, optimal and nonoptimal sites for the coding sequences.

2.2 genomic_gc3

```
genomic_gc3(genomic_cds_file)
```

The function calculates the genomic gc3 for an multifasta genomic CDS file. The function first concatenated all the CDS sequences in the file into one long CDS string, than calculated the gc3 from the GC3 function in seqinr package. You can also use the function to calculate the gc3 for a single gene, or a set of genes, depends what CDS sequences you put in the input file. The result can be used as input for the s_index calculation.

2.3 op_corre_codonW

```
op_corre_CodonW(genomic_cds_file=NULL, correspondence_file=NULL)
```

The function identify the optimal codons based on the correlative method from Hershberg & Petrov. This method take the whole genome into consideration, and predict the optimal codons by making the correlation between the frequency of each codon within each gene and the overall codon bias (Nc or Nc'). The input file include the correspondence analysis output file from the program CodonW (to get the Nc value), and the genomic cds file (to get the codon usage information for each gene).

For further details regard how to use CodonW, you can refer to the site <http://codonw.sourceforge.net/>. Note, you must input the same genomic cds file to CodonW and to the op_corre_CodonW funtion, so that the order and number of genes are consistent in the files.

2.4 op_corre_NCprime

```
op_corre_NCprime(genomic_cds_file=NULL, nc_file=NULL)
```

The function identify the optimal codons based on the correlative method from Hershberg & Petrov. This method take the whole genome into consideration, and predict the optimal codons by making the correlation between the frequency of each codon within each gene and the overall codon bias (Nc or Nc'). The input file include the output file from the program ENCPprime (to get the Nc and Nc' value), and the genomic cds file (to get the codon usage information for each gene).

For further details regard how to use ENCPprime, you can refer to the site <https://github.com/jnovembre/ENCPprime>. Note, you must input the same genomic cds file to ENCPprime and to the op_corre_NCprime funtion, so that the order and number of genes are consistent in the two files.

2.5 op_highly

```
op_highly(high_cds_file = NULL, ref_cds_file = NULL, p_cutoff = 0.01)
```

Optimal codons are defined as codons significantly enriched in the highly expressed genes compared to the lowly expressed genes, or other set of appropriate reference genes. In another word, these codons were favored by translational selection, which was strongest among highly expressed genes. This function calculate the

optimal codon list with p-values, thus user could have a general idea of which codons were preferred by selection in the genome.

The argument `high_cds_file` should specific the path for the highly expressed gene dataset. It is up to the users how to define which dataset of highly expressed genes. Some studies use the expression data, or Nc value to divide genes into highly/lowly sets. Other studies use a specific dataset, such as only including the very highly expressed genes (ribosomal genes).

The argument `ref_cds_file` should specific the path for the lowly expressed gene dataset, or any appropriate dataset. In Sharp PM paper (Forces that influence the evolution of codon bias), he used the all gene data set as neutral reference and also get a list of optimal codons.

The argument `p_cutoff` set the cutoff for p values in the chi.square test. Only codons are significantly enriched in the highly expressed genes are marked with + symbol in the ouotput tables. The codons are significantly lower presented in the highly expressed genes are marked with - symbol. The codons are not significantly differently presented compared to the reference dataset are marked with NA symbol.

2.6 op_highly_stats

```
op_highly_stats(high_cds_file=NULL,ref_cds_file=NULL,p_cutoff=0.01)
```

If you want to know the detailed codon usage information for the highly expressed and all genes, you can run the function `op_highly_stats` in the package. The optimal codons are defined as codons significantly enriched in the highly expressed genes compared to the lowly expressed genes, or other set of reference genes. The function calculate the optimal codon list with p-values, thus user could have a general idea of which codons were preferred by selection in the genome.

The argument `high_cds_file` should specific the path for the highly expressed gene dataset. It is up to the users how to define which dataset of highly expressed genes. Some studies use the expression data, or Nc value to divide genes into highly/lowly sets. Other studies use a specific dataset, such as only including the very highly expressed genes (ribosomal genes).

The argument `ref_cds_file` should specific the path for the lowly expressed gene dataset, or any appropriate dataset. In Sharp PM paper (Forces that influence the evolution of codon bias), he used the all gene data set as neutral reference to infer codons used significantly in highly expressed genes.

The argument `p_cutoff` set the cutoff for p values in the chi.square test. Only codons are significantly enriched in the highly expressed genes are marked with + symbol in the ouotput tables. The codons are significantly lower presented in the highly expressed genes are marked with - symbol. The codons are not significantly differently presented compared to the reference dataset are marked with NA symbol.

The function also output the rscu value for the high expressed dataset and reference dataset, actual and expected number of codons, the p value of chi.square test.

2.7 low_frequency_op

```
low_frequency_op(high_cds_file=NULL,genomic_cds_file=NULL,p_cutoff=0.01)
```

The function `low_frequency_op` identify the low frequency optimal codons. Based on the previous study (Sun Y et al. Switches in genomic GC content drives shifts of optimal codons under sustained selection on synonymous sites. Genome Biol Evol. 2016 Aug 18.), in some species, the optimal codons identified by the `op_highly` function has bit strange patterns: the optimal codons do have lower prequency than the non-optimal codons. It occurs most in the mutation-shifting species such as *G. vaginalis*. The function can identify these low frequency optimal codons.

The function first calculated all the optimal codons statistics by the function `op_highly_stats` in the package, then tried to find the low frequency optimal codons with the following settings: 1) it has to be an optimal

codon 2) The RSCU value for the optimal codon is lower than 0.7 (this is the quite arbitrary setting) 3) the RSCU value for the optimal codon is lower than the corresponding non-optimal codons, which has the same first two nucleotide as the optimal codon but the third position experience the point mutation transition. With this detailed filters, we can identify the low frequency optimal codons in the given genome.

The argument `high_cds_file` should specific the path for the highly expressed gene dataset. It is up to the users how to define which dataset of highly expressed genes. Some studies use the expression data, or Nc value to divide genes into highly/lowly sets. Other studies use a specific dataset, such as only including the very highly expressed genes (ribosomal genes). In the example, I used the ribosomal genes as the representative for the highly expressed genes.

The arguments, `genomic_cds_file`, is used to calculate the optimal codons and statistics for highly and all genes.

Same as the `op_highly` and `op_highly_stats`, the p value cutoff was set as 0.01.

2.8 s_index

```
s_index(high_cds_file = NULL, genomic_cds_file = NULL, gc3 = NULL)
```

The function calculate the s index based on Paul Sharp's method. The method take into account of background mutation rate (in the program, two arguments `genomic_cds_file` and `gc3`, are input to calculate the mutation rate), and focus only on codons with universal translational advantages in all bacterial species (in the program, one argument `high_cds_file`, is input to calculate these codons). Thus the s index can be used to quantify the strength of translational selection and is comparable among different species.

The argument `high_cds_file` much be specified with the input filepath for the highly expressed genes. The file should be a multifasta file contains 40 highly, including elongation factor Tu, Ts, G, 50S ribosomal protein L1 to L6, L9 to L20, 30S ribosomal protein S2 to S20. This file can be generated by either directly extract these DNA sequence from genbank file, or parse by blast program. For the four amino acids (Phy, Tyr, Ile and Asn), the C-ending codons are universally preferred than the U-ending codons. Thus, only these four codons were taken into account in the analyses.

The second arguments, `genomic_cds_file` or `gc3`, is used to calculate the genomic mutation rate, and one of them must be specified. The `genomic_cds_file` should be a multifasta file contains all the coding sequences in the genome, and the function use it to calculate the genomic `gc3` and mutation rate. If the `gc3` value for the genome is known already, or calculated by `genomic_gc3` function below, you can specify it in the argument `gc3`. If both the `genomic_cds_file` and `gc3` arguments are specified, the function will use the `genomic_cds_file` to calculate mutation rate, and neglect the `gc3` argument.

3 Workflow

3.1 general statistics

Here is an example of how to use these functions from the package. We will use the genomic CDS file from the *Lactobacillus. kunkeei* and *Gardnerella vaginalis* as exemple for the workflow. These files are included in the sequences folder in the package.

First, we can calculate the genomic `gc3` for the *L. kunkeei* to have a general statistics about the genome. We show how to do it here with `genomic_gc3` function:

```
library(sscu)
genomic_gc3(system.file("sequences/L_kunkeei_genome_cds.ffn",package="sscu"))
```

```
## [1] 0.3083719
```

The genomic gc3 content is 0.308, which is a low gc genome.

Then, we can further check the codon usage pattern in details for *L. kunkeei* by using the `op_highly_stats` function:

```
head(op_highly_stats(high_cds_file=system.file("sequences/L_kunkeei_highly.ffn",package="sscu"),
  ref_cds_file=system.file("sequences/L_kunkeei_genome_cds.ffn",package="sscu")))
```

```
##      codon aa rscu_high rscu_ref high_No_codon high_expect_No_codon
## TTT    TTT F      0.66    1.10           58           88
## TTC    TTC F      1.34    0.90          118           88
## TTA    TTA L      2.58    2.54          198           77
## TTG    TTG L      0.56    1.09           43           77
## TCT    TCT S      1.10    1.04           64           58
## TCC    TCC S      0.12    0.38            7           58
##      ref_No_codon ref_expect_No_codon p_value symbol
## TTT           10540           9539    0.001      -
## TTC            8538           9539    0.005      +
## TTA           16254           6399    0.890     NA
## TTG            6971           6399    0.001      -
## TCT            5263           5062    0.796     NA
## TCC            1937           5062    0.002      -
```

We can see several statistics for each codon, including RSCU values, actual and expected number of codons, the p value of chi.square test. The symbol indicates if the codon is highly (optimal codons), no difference, or lowly presented in the highly expressed genes compared to the reference gene set. Thus, you can get a list of optimal codons for the output. The output has 64 lines with each line for each codon, and we only show the first six lines for simplicity.

3.2 optimal codons

Optimal codons are defined as codons favored by translational selection, by promoting either translational efficiency or accuracy. If we want to get a list of optimal codons by the highly method:

```
op_highly(high_cds_file=system.file("sequences/L_kunkeei_highly.ffn",package="sscu"),
  ref_cds_file=system.file("sequences/L_kunkeei_genome_cds.ffn",package="sscu"))
```

```
## [1] "AAC" "ACT" "ATC" "CGT" "GAC" "GCT" "GGT" "GTT" "TAC" "TCA" "TTC"
```

Now we get a list of optimal codons in *L. kunkeei* “AAC” “ACT” “ATC” “CGT” “CTA” “GAC” “GCT” “GGT” “GTT” “TAC” “TCA” “TTC”. You can input the list to the further analysis `akashi_test`.

If you want to get a list of optimal codons by the correlative method, you need to input a file with Nc or Nc’ information for all the genes you are interested. Now you have two options, either input the Nc file by the correspondence analysis from CodonW, or input the Nc/Nc’ file by ENCprime:

```
op_corre_CodonW(genomic_cds_file=system.file("sequences/Gvag_genome_cds.ffn",package="sscu"),
  correspondence_file=system.file("correlative_test/Gvag.codonw",package="sscu"))
```

```
## [1] "act" "cct" "gct" "ggg" "gtt" "tgt" "ctt" "cgt" "ttg" "att" "aag" "aat"
## [12] "cag" "cat" "gat" "tgc" "ttc"
```

```
op_corre_NCprime(genomic_cds_file=system.file("sequences/LbDelBA1_genome_cds.ffn",package="sscu"),
  nc_file=system.file("correlative_test/LbDelBA1.NCprime",package="sscu"))
```

```
## [1] "acc" "cca" "gcc" "ggc" "gtc" "tcc" "cgc" "ctg" "atc" "aag" "aac"
## [12] "cac" "gaa" "gac" "tac" "tgc" "ttc"
```

Generally, these two test should give similar optimal codon list, but for some cases, especially for genomes with low GC content, the list could be different, for further details, read Hershberg R, Petrov DA. 2009. General rules for optimal codon choice. Plos Genet. 5:e1001115.

3.3 selective profile

Next, if we want to see the selective profile for the species, which is the major function in the package. We can use the `s_index` function to do the calculation. We can use either the genomic gc3 content 0.308 as we got previously:

```
s_index(gc3=0.308,
        high_cds_file=system.file("sequences/L_kunkeei_highly.ffn",package="sscu"))
```

```
## [1] 1.541898
```

or you can directly input the CDS file in the function.

```
s_index(high_cds_file=system.file("sequences/L_kunkeei_highly.ffn",package="sscu"),
        genomic_cds_file=system.file("sequences/L_kunkeei_genome_cds.ffn",package="sscu"))
```

```
## [1] 1.540154
```

We can get the S index as 1.54 from both method. S index = 0 indicates that, in the highly expressed genes, the relative frequency of C/U-ending codons are exactly the same as the genomic gc3 (mutational pattern). S index > 0.3 indicate C-ending codons are used in significantly higher frequency than the U-ending codons, which further indicates that translational selection is affecting codon usage. In this case, S index of 1.54 suggests strong translational selection in *L. kunkeei*, which is even higher than *E. coli* (1.49).

Akashi's test is another way to evaluate the translational selection on codon usage. Akashi's theory suggests that the optimal codons are codons that translated more accurately than the other codons, and these codons are favored in the important sites, such as the evolutionary conserved amino acid sites, whereas the less conserved amino acids sites are more tolerable to the non-optimal codons.

```
akashi_test(system.file("akashi_test/contingency_file_Gvag",package="sscu"))
```

```
## $Z
## [1] 4.443357
##
## $p
## [1] 4.428296e-06
##
## $odd_ratio
## [1] 1.081472
##
## $conserved_optimal_sites
## [1] 17774
##
## $conserved_non_optimal_sites
## [1] 48683
##
## $variable_optimal_sites
## [1] 13235
##
## $variable_non_optimal_sites
## [1] 54958
##
## $con_var_ratio
```

```
## [1] 0.9745428
```

3.4 a function to identify the low frequency optimal codons

The function `low_frequency_op` can be used to identify the low frequency optimal codons. Most of the optimal codons should be quite frequently used, especially in the highly expressed genes. However, Based on the previous study (Sun Y et al. Switches in genomic GC content drives shifts of optimal codons under sustained selection on synonymous sites. *Genome Biol Evol.* 2016 Aug 18.), in some species, the optimal codons do have lower frequency than the corresponding non-optimal codons. These codons were defined as the codons that has the same first two nucleotide as the optimal codon but a different third codon position via point mutation transition. The pattern occurs most in the mutation-shifting species such as *G. vaginalis*. For example:

```
low_frequency_op(high_cds_file=system.file("sequences/Gvag_highly.ffn",package="sscu"),
  genomic_cds_file=system.file("sequences/Gvag_genome_cds.ffn",package="sscu"))
```

```
## $low_frequency_optimal_codons
##      codon aa rscu_high rscu_ref high_No_codon high_expect_No_codon
## GTC   GTC  V      0.67    0.23           122           182
## GCC   GCC  A      0.58    0.22           136           232
##      ref_No_codon ref_expect_No_codon p_value symbol
## GTC             1961             8594      0      +
## GCC             2423             10882     0      +
##
## $corresponding_codons
##      codon aa rscu_high rscu_ref high_No_codon high_expect_No_codon
## GTT   GTT  V      1.87    1.96           341           182
## GCT   GCT  A      2.24    1.83           520           232
##      ref_No_codon ref_expect_No_codon p_value symbol
## GTT             16834             8594    0.604    NA
## GCT             19914             10882    0.012    NA
```

The output is similar to the output of `op_highly_stats`, but it only contains the low frequency optimal codons. The upper dataframe is the low frequency optimal codons, and lower corresponding non-optimal codons. We can clearly see that even in the highly expressed genes, the optimal codons are used less frequently than the non-optimal codons. The pattern suggests that even selection may strive the use the old C-ending codons, the power of mutation in the direction of GC to AT overwhelm this preference. The genomic GC3 value for *G. vaginalis* is 0.34.

```
genomic_gc3(system.file("sequences/Gvag_genome_cds.ffn",package="sscu"))
```

```
## [1] 0.3427021
```

Most of the genomes do not have any low frequency optimal codons, for example, in *L. kunkeei*, no low frequency optimal codons has been identified.

```
low_frequency_op(high_cds_file=system.file("sequences/L_kunkeei_highly.ffn",package="sscu"),
  genomic_cds_file=system.file("sequences/L_kunkeei_genome_cds.ffn",package="sscu"))
```

```
## $low_frequency_optimal_codons
## [1] codon      aa      rscu_high
## [4] rscu_ref    high_No_codon high_expect_No_codon
## [7] ref_No_codon ref_expect_No_codon p_value
## [10] symbol
## <0 rows> (or 0-length row.names)
##
```



```
## $corresponding_codons
## [1] codon          aa          rscu_high
## [4] rscu_ref        high_No_codon high_expect_No_codon
## [7] ref_No_codon    ref_expect_No_codon p_value
## [10] symbol
## <0 rows> (or 0-length row.names)
```

4 Folders installed in the sscu package

There are three additional installed folder in the sscu package directory, named as sequences, akashi_test and correlative_test

4.1 sequences

As the name indicates, the folder has five sequences files: three whole genome coding sequence file and two highly expressed genes coding sequence file

4.2 akashi_test

This folder contains one sub-folder bifidos_alignments, one perl script make_contingency_table.pl and one data file contingency_file_Gvag. All the files are related to the function akashi_test in the package.

The akashi_test function require the input of an contingency file (example as contingency_file_Gvag). The contingency file can be generated by the perl script, you can check in detail about how to use it by reading the first few lines of the perl script. The perl script tabulates and calculate the a,b,c,d entries for the coding sequences, and output the four values for each amino acid for each gene. The input of the perl script is a folder (example as bifidos_alignments) contains the codon alignments of the genes that you are interested to calculate.

4.3 correlative_test

This folder contains two files as input to calculate optimal codons by correlative method. Gvag.codonw is the example input file for op_corre_CodonW function, and LbDelBA1.NCprime is the example input file for op_corre_NCprime function.