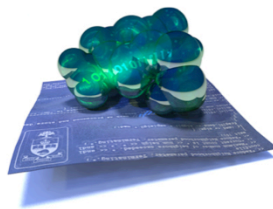


A  
BIOINFORMATICS  
COURSE

# SEQUENCE ANALYSIS: COMPOSITION



---

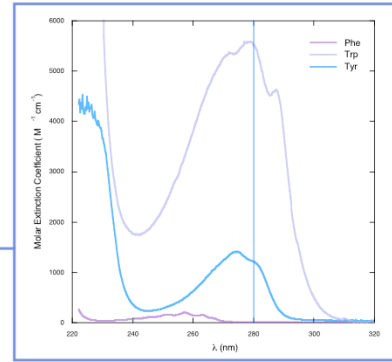
BORIS STEIPE

*DEPARTMENT OF BIOCHEMISTRY – DEPARTMENT OF MOLECULAR GENETICS  
UNIVERSITY OF TORONTO*



Some properties of biomolecules can be simply summed over their components: they depend only on the amino acid (or nucleotide) composition of the molecule.

Molecular weight  
Extinction coefficient  
Isoelectric point



$$e_{280} \text{ (M}^{-1} \text{ cm}^{-1}\text{)} = (5,500)\#\text{Trp} + (1,490)\#\text{Tyr} + (125)(\#\text{Cys-Cys})$$

Gasteiger *et al.* (2005) Protein Identification and Analysis tools on the ExPASy server. *The Proteomics Protocols Handbook*, Human Press 571-607

(# here means “number of”)

A protein's isoelectric point depends on the pK values of the amino acids; the pK values characterize the propensity for an amino acid sidechain to dissociate, which in turn depends on how energetically favourable dissociation is. For example: since a negatively charged amino acid will be stabilized in a positive electrostatic field, such a field will shift a pK value **down**. This means the pH value at which the side chain will be 50% ionized is lower, or in other words, in a positive electrostatic field the concentration of protons must be higher to keep a proton associated to the sidechain. Compositional properties of nucleic acids include hybridization temperature and helix structure.

SEQUENCE COMPOSITION: DATA

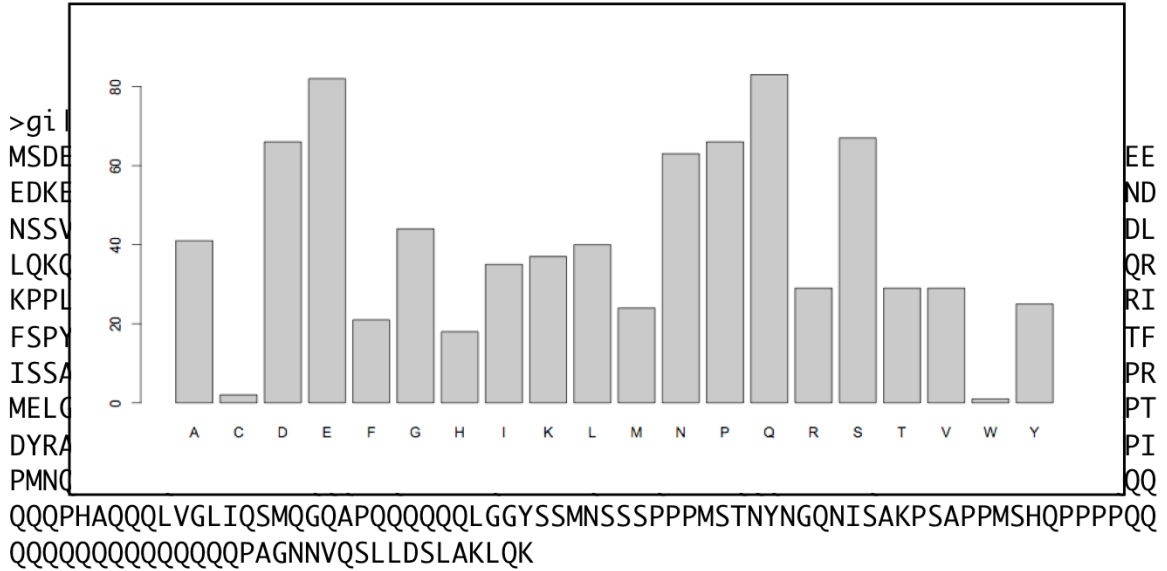
```
>gi16325066|refINP_015134.1| Nab3p [Saccharomyces cerevisiae S288c]
MSDENHNSDVQDIPSELVSDNSNENELMNSSADDGIEFDAPEEEEREAEEREEENEEQHELEDVNDEEE
EDKEEKGEENGEVINTEEEEEEEHQQKGGNDDDDDDNEEEEEEEEDDDDDDDDDDDDEEEEEEEEEGND
NSSVGSDSAEDGEDEEDKKDKTKDKEVELRRETLEKEQKDVEAIKKITREENDNTHFPTNMENVNYDL
LQKQVKYIMDSNMLNLPQFQHLPQEEKMSAILAMLNSNSDTALSVPPHDSTISTTASASATSGARSNDQR
KPPLSDAQRRMRFPADLSKPITEEEHNDRYAAYLHGENKITEMHNIPPKSRLFIGNPLKKNVSKEDLFRI
FSPYGHIMQINIKNAFGFIQFDNPQSVRDAIECESQEMNFGKKLILEVSSSNARPQFDHGDHGTNSSSTF
ISSAKRPFQTESGDMYNDNGAGYKKSRRHTVSCNIFVKRTADRTYAIEVFNRFRDGTGLETDMIFLKPR
MELGKLINDAAYNGVWGVVLVNKTHNVDVQTFYKGSQGETKFDEYISISADDAVAIFNNIKNNRNNRPT
DYRAMSHQQNIYGAPPLPVPNGPAVGPPPQTNYQQGYSMPPPQQQQQPYGNYGMPPPSHDQGYGSQPPI
PMNQSYGRYQTSIPPPPPQQQIPQGYGRYQAGPPPPQPSQTPMDQQQLLSAIQNLPPNVVSNLLSMAQQQ
QQQPHAQQLVGLIQSMQGAPOQQQQQLGGYSSMNSSSPPPMSTNYNGQNISAKPSAPPMSHQPPPPQQ
QQQQQQQQQQQQQPAAGNNVQSLLDSLAKLQK
```

```
> cat(readLines("nab3.fa"), sep = "\n")
```

Let us discuss a simple example of composition analysis for a given protein sequence, the yeast Nab3 protein.

The atypical distribution and clustering of particular amino acids suggests consequences for folding and interactions of the encoded protein.

SEQUENCE COMPOSITION: COUNTS

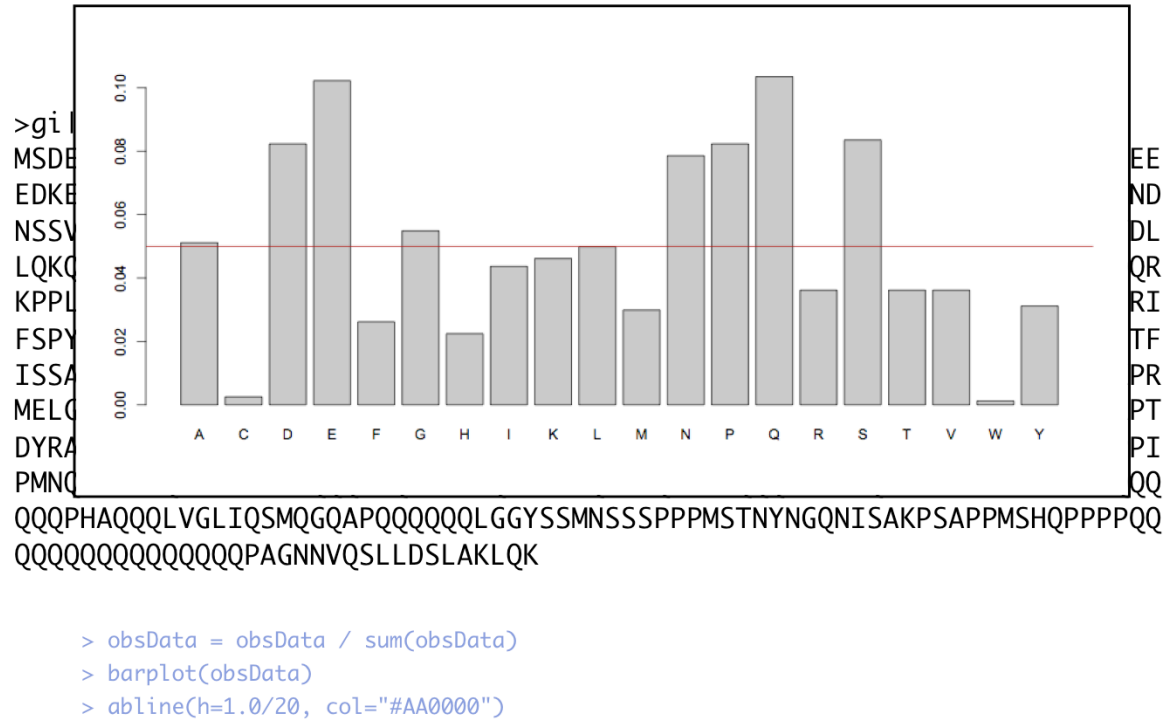


```

> nab3 <- read.fasta("nab3.fa")
> obsData <- table(nab3)
> barplot(obsData)
  
```

Here, **R** was used to tabulate the counts of the different amino acids in the sequence. The values are shown in a barplot, ordered by one-letter code, alphabetically. This ordering makes it hard to evaluate trends quickly.

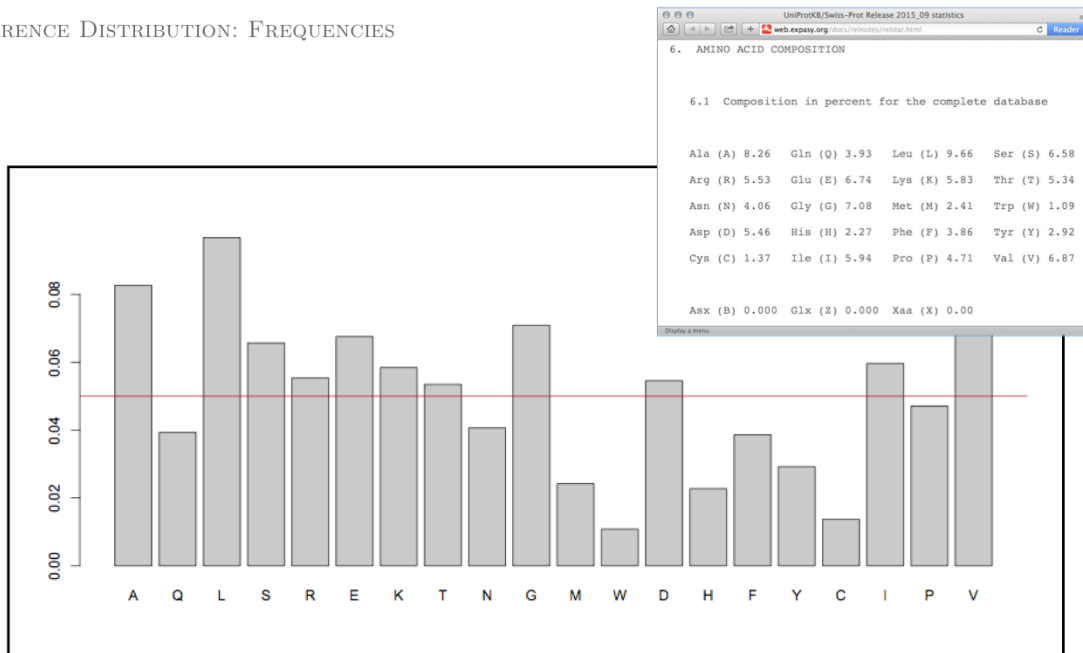
SEQUENCE COMPOSITION: FREQUENCIES



Converting the counts into frequencies, and adding a line to display an expected value, gives us a little more information. We can begin to understand which amino acids are over- and underrepresented.

But what is the “expected value”? The red line in the plot is simply taken as 5% (0.05) – the expected frequency if all 20 amino acids would occur in equal amounts.

## REFERENCE DISTRIBUTION: FREQUENCIES



Data taken from UniProt release notes. Database average over all proteins.

Note: it is often not clear what the best reference dataset is, and the information from comparison depends critically on the choice of reference.

```
> barplot(refData)
> abline(h=1.0/20, col="#AA0000")
```

Equal probability is a poor assumption. But we can access the large sequence databases to evaluate the frequency of amino acids in proteins in general.

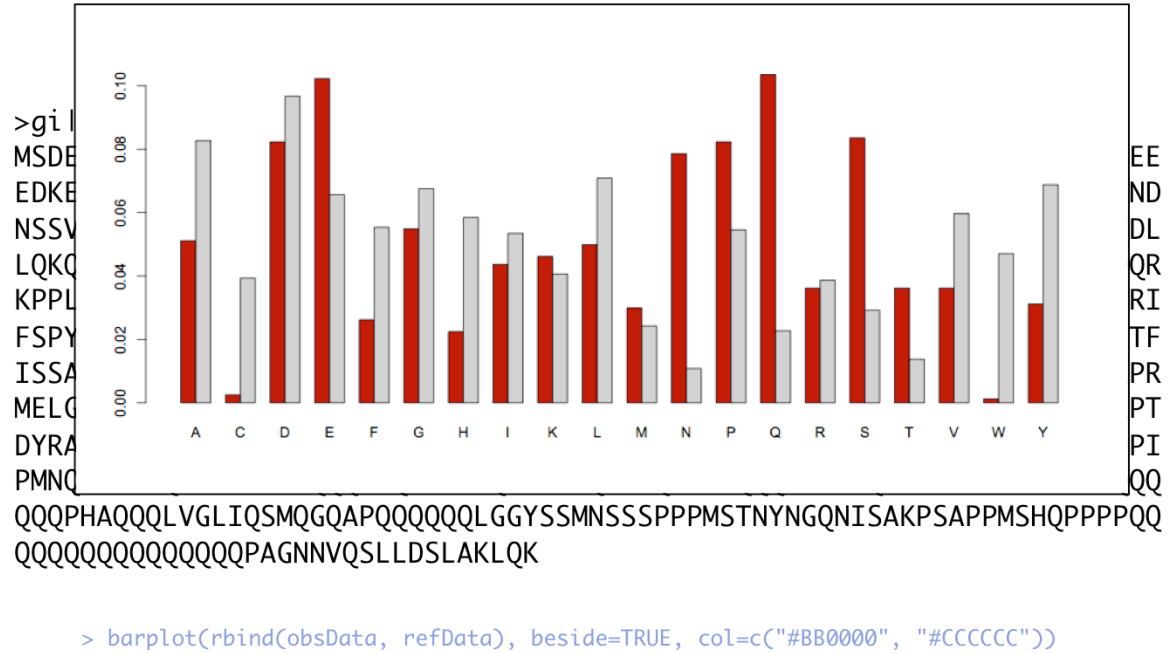
Is this a better assumption? Probably, but it has its own problems. It does not distinguish between highly- and poorly- expressed sequences.

Other considerations may be to look at an organism's total protein, or distinguish between membrane- and cytoplasmic proteins, nuclear proteins, secreted proteins. Or to take the metabolic cost of amino acids into account. Or other biologically motivated distributions we can come up with.

You should note that the definition of an expected distribution is at least as important as to compile the observations.

The inset frequencies are database averages.

SEQUENCE COMPOSITION: FREQUENCIES COMPARED

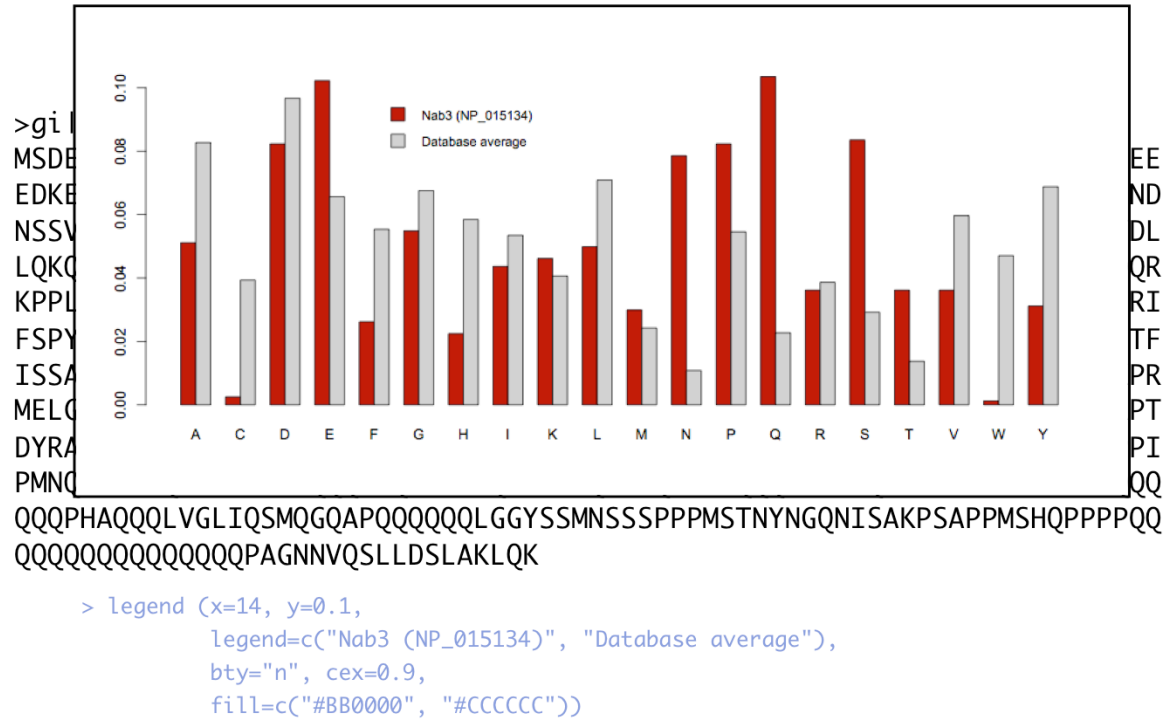


Given this information, we can compare database values with the values in our sequence.

Wait: what do the bars represent? Which is which?



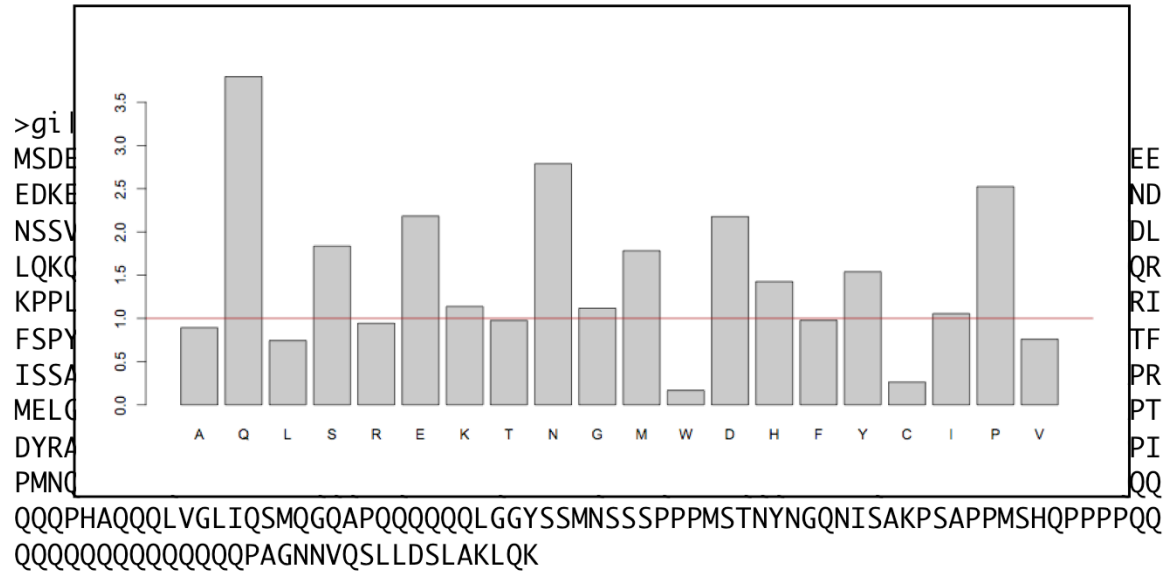
SEQUENCE COMPOSITION: FREQUENCIES COMPARED



We need to add a legend ...

The amino acids that have a larger or smaller than average frequency are becoming apparent. But absolute values are not well suited for this type of comparison. It is much more convenient to express such differences as ratios.

SEQUENCE COMPOSITION: RATIOS



```
> barplot(obsData / refData)
> abline(h=1, col="#AA0000")
```

Now we have a ratio of one if the frequencies are the same, a ratio of 0.5 if the observed frequencies are half that of the database reference, and a ratio of 2.0 if they are double. Therefore the values mean: how much more likely do we observe an amino acid than we expect it.

But we really should make  $\frac{1}{2}$  and 2 times the expected frequency give the same distance on the plot – after all, which is which depends only on our arbitrary choice of which distribution should be in the numerator and denominator of the fraction.

Therefore we express such relationships as  $\log(\text{ratio})$ .

SEQUENCE COMPOSITION: LOG-RATIOS

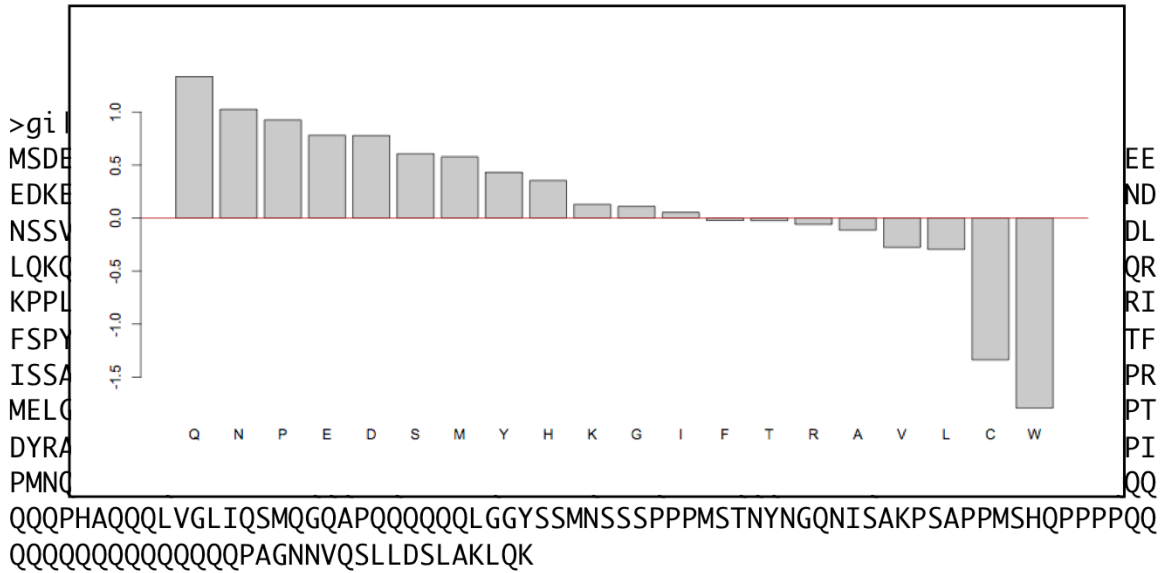


```

> barplot(log(obsData / refData))
> abline(h=0, col="#AA0000")
  
```

In a log ratio, when observed and expected frequencies are the same, the value is zero. Excesses are positive and depletions are negative. The same relative difference (2-fold more, 2-fold less) gives the same distance on the plot, regardless of whether the log-ratio is positive or negative in absolute terms

SEQUENCE COMPOSITION: LOG-RATIOS SORTED

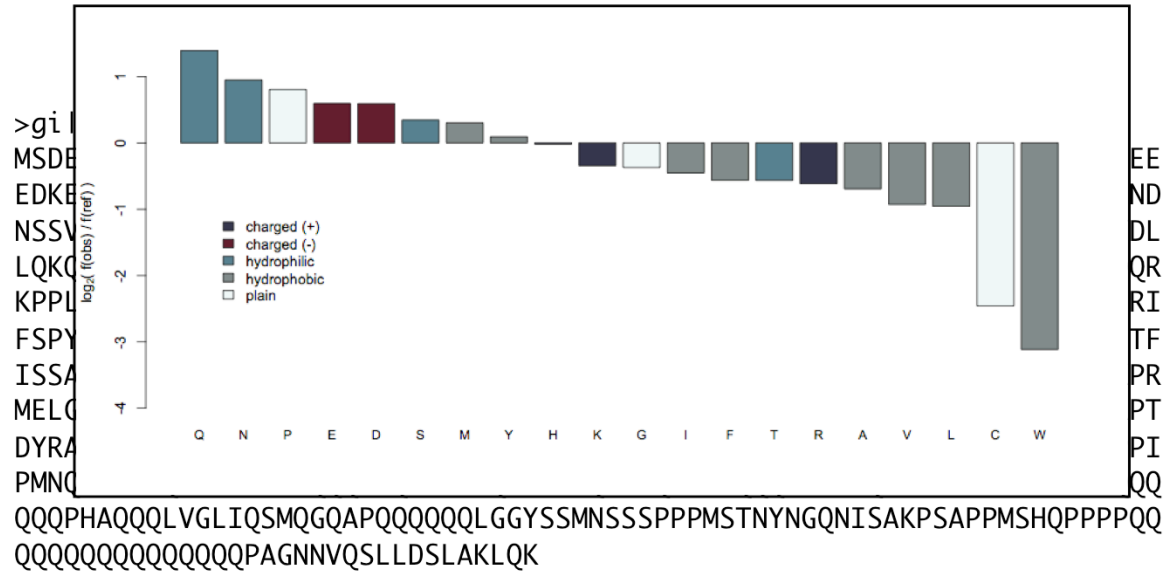


```
>gi |
MSDE |
EDKE |
NSSV |
LQKC |
KPPL |
FSPY |
ISSA |
MELC |
DYRA |
PMNC |
EE
ND
DL
QR
RI
TF
PR
PT
PI
QQ
QQQPHAQQQLVGLIQSMQGQAPQQQQQLGGYSSMNSSSPPPMSTNYNGQNISAKPSAPPMSHQPPPPQQ
QQQQQQQQQQQQQQPAGNNVQSLLDLAKLQK
```

```
> barplot(sort(log(obsData / refData), decreasing=TRUE))
> abline(h=0, col="#AA0000")
```

To further clarify what we are seeing here, we can sort by values.

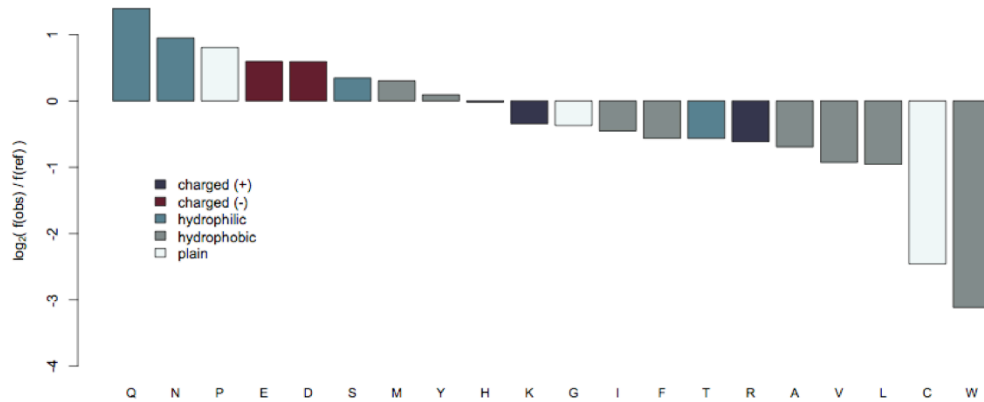
SEQUENCE COMPOSITION: LOG-RATIOS SORTED AND CATEGORIZED



Once we colour the resulig bars by property, we can extract important trends: our sequence has an excess of hydrophilic amino acids and of proline. This would be compatible with unstructured regions that can remain unfolded without aggregating. It also has an excess of negatively charged amino acids, with a depletion of positively charged amino acids. This would be consistent with binding to positively charged molecules. Hydrophobic (aggregation promoting) residues are depleted. That would be consistent with a protein that is natively, functionally unfolded over parts of its sequence.

From this we can build a mental image that this protein might bind to a large, positively charged polymer in a conformation-independent manner.

## SEQUENCE COMPOSITION: INTERPRETATION



```
>gi|6325066|ref|NP_015134.1| Single stranded RNA binding protein; acidic ribonucleoprotein;
required for termination of non-poly(A) transcripts and efficient splicing; interacts with Nrd1p;
Nab3p [Saccharomyces cerevisiae]
MSDENHNSDVQDIPELSVDSNSNENELMNNSSADDGIEFDAPEEEREAREEEENEEOHELEDVNDEEEEDKEEKGEENGEVINTEEEEEEEHQQKG
GNDDDDDDNEEEEEEDDDDDDDDDDEEEEEEGNDNSSVGSDSAEDGEDEEDKDKTKDKEVELRRETLEKEQKDVDEAIKKITREENDN
THFPPTNMENVNYDLLQKQVKYIMDSNMLNLPQFQHLPEEKMSAILAMLNSNSDTALSVPVPHDSTISTTASASATSGARSNDQRKPLSDAQRMRFP
RADLSKPIEEEEHRYAAYLHGENKITEMHNIIPKSRLEIGNLPLKNVSKEDLFRIFSPYGHIMQINIKNAFGFIQFDNPQSVRDAIECESQEMNFGK
KLILEVSSSNARPQFDHGDHGTNSSSTFISSAKRPFQTESGDMYNDNGAGYKKSRRHTVSCNIFVKRTADRITYAIEVFNRFRDGTGLETDMIFLKPR
MELGKGLINDAAYNGVWGVVLVKNKTHNVQVTFYKGSQGETKFDEYISISADDAVAIFNNIKNNRNSRPTDYRAMSHQQNIYGAPPLPVPNGPAVGPP
PQTNYQGYSMPPPQQQQQPPYGNYGMPSPSHDQGYGSQPPIPMNQSYGRYQTSIPPPPPQQQIPQGYGRYQAGPPPQPPSQTMDQQQLLSAIQNL
PNVVSNLLSMAQQQQQPHAQQLVGLIQSMQQAQQQQQQLGGYSSMNSSPPPMSTNYNGQNISAKPSAPPMSHQPQQQQQQQQQQQQQQQQ
PAGNNVQSLDLSLAKLQK
```

These observations are entirely consistent with the annotated function of the protein, Nab3: a single-stranded nucleic acid binding protein. You would expect it to have

- Disordered regions that interact with disordered ligands;
- Negative charge that complements the positive charge of the exposed nucleobases in single-strand nucleic acid molecules;
- A reduced amount of aggregation-promoting residues to keep the disordered structure in solution in the cytoplasm.

Specifically, we can also hypothesize that the protein coats nucleobases and avoids the backbone, keeping the RNA from forming double-stranded secondary structure, perhaps even promoting melting in the first place.

It's remarkable how far we can (sometimes) get with considering composition alone.

<http://steipe.biochemistry.utoronto.ca/abc>

B O R I S . S T E I P E @ U T O R O N T O . C A

DEPARTMENT OF BIOCHEMISTRY & DEPARTMENT OF MOLECULAR GENETICS  
UNIVERSITY OF TORONTO, CANADA