

Combinatorial Codons

What input does one need to create a set of random sequences with a desired amino acid composition?

Litterature

Arkin and Youvan claim that using 25% of each base in synthesis is inefficient [\[1\]](#). Wolf and Kim made an online tool, which run via a matlab analogue [\[2\]](#) (NB. E. Wolf is not P Wolfe of the [Frank-Wolfe alogorith](#)).

Craig et al prove that the task is impossible to get a near perfect match and use a gene algorimth for the search and use 6 separate codon mixes ("tubes") to meet their requirements [\[3\]](#).

Code

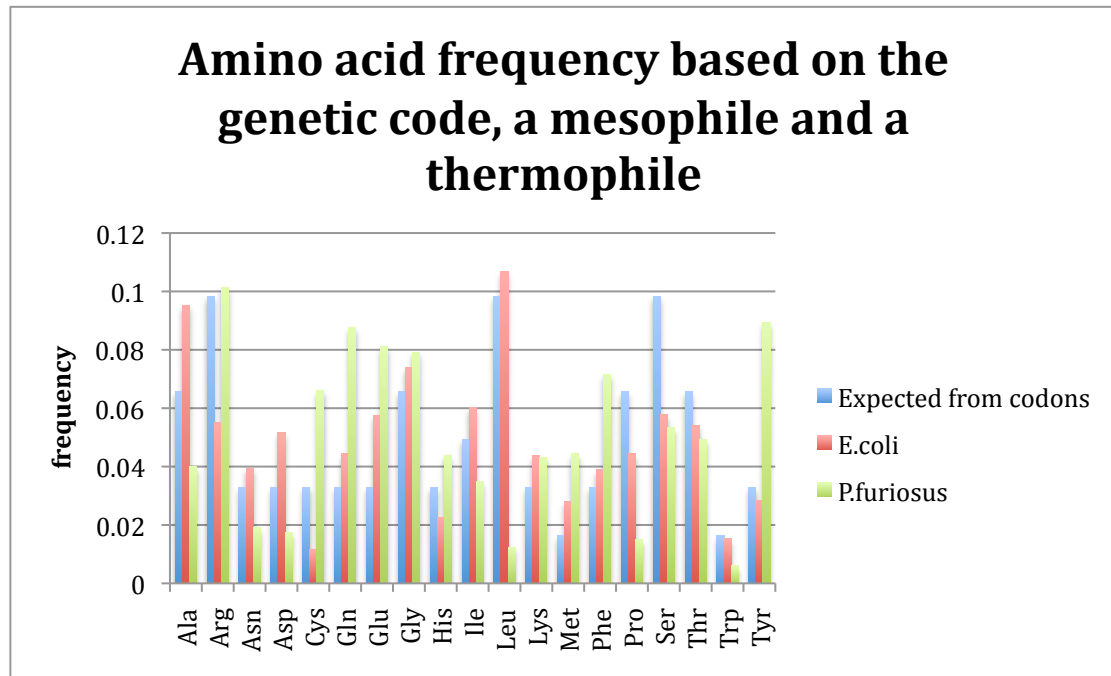
This is an optimization problem, where:

- The matrix to be optimized is the frequency of each base for position 1, 2 & 3 with the constraint that the frequencies are positive and the sum of each base at a position is 1. For simplicity it can be expressed as a 12-dimentional vector.
- The best result is determined by the objective function, for a list of them see [\[4\]](#) or [\[3\]](#). The most common is the least squares approach, where the residues of the fit are minimized. Wolf uses a cosine basin approach.
- Search method.
 - Systematic. for the XYZ situation it takes 30-60 minutes when searching all combinations of 0.1 increments.
 - Heuristic. Several algorithms [exist](#).
 - Wolf uses an algorithm called CFSQP, a very fancy version of Newton method from what I understand.
 - Craig uses a gene algorithm.
 - For simplicity I used a hill-climbing algorithm.
- A translation function to translate the base vector into the amino acid vector.

Stop codons

The genetic code has been described by Crick as a frozen accident.

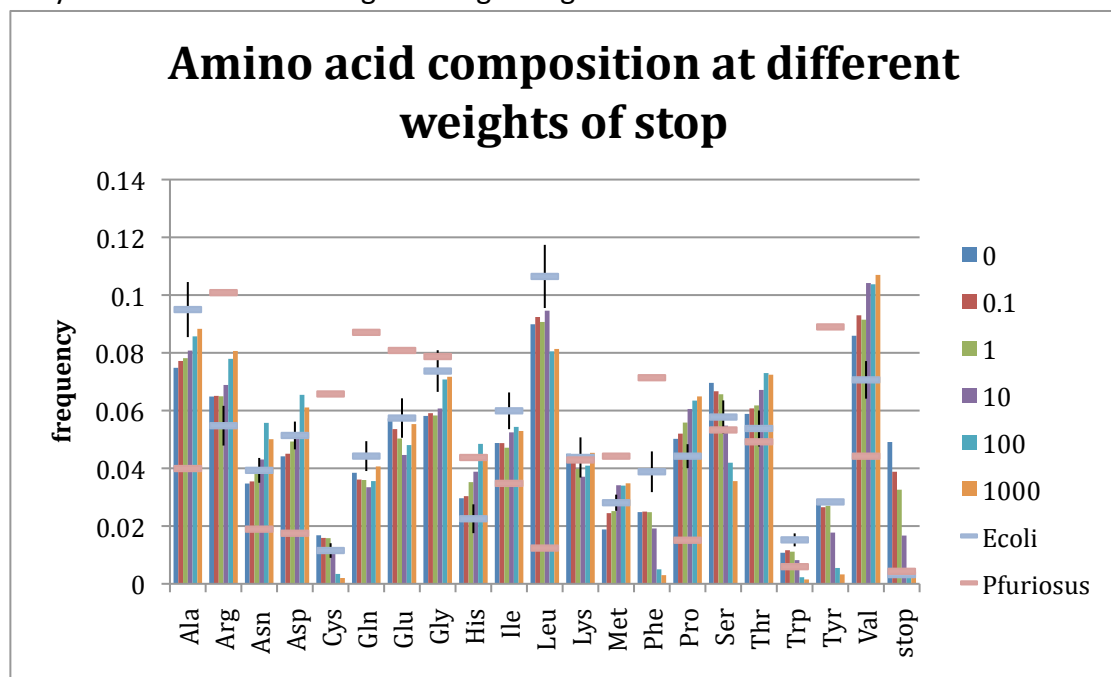
The number of codons and the amino acid frequency unfortunately do not correlate too well (fig 1) otherwise this problem would be non existent



The largest disparity are the 3 stop codons where the average E.coli protein length is 300 amino acids and not 20.

Do stop codons mess stuff up?

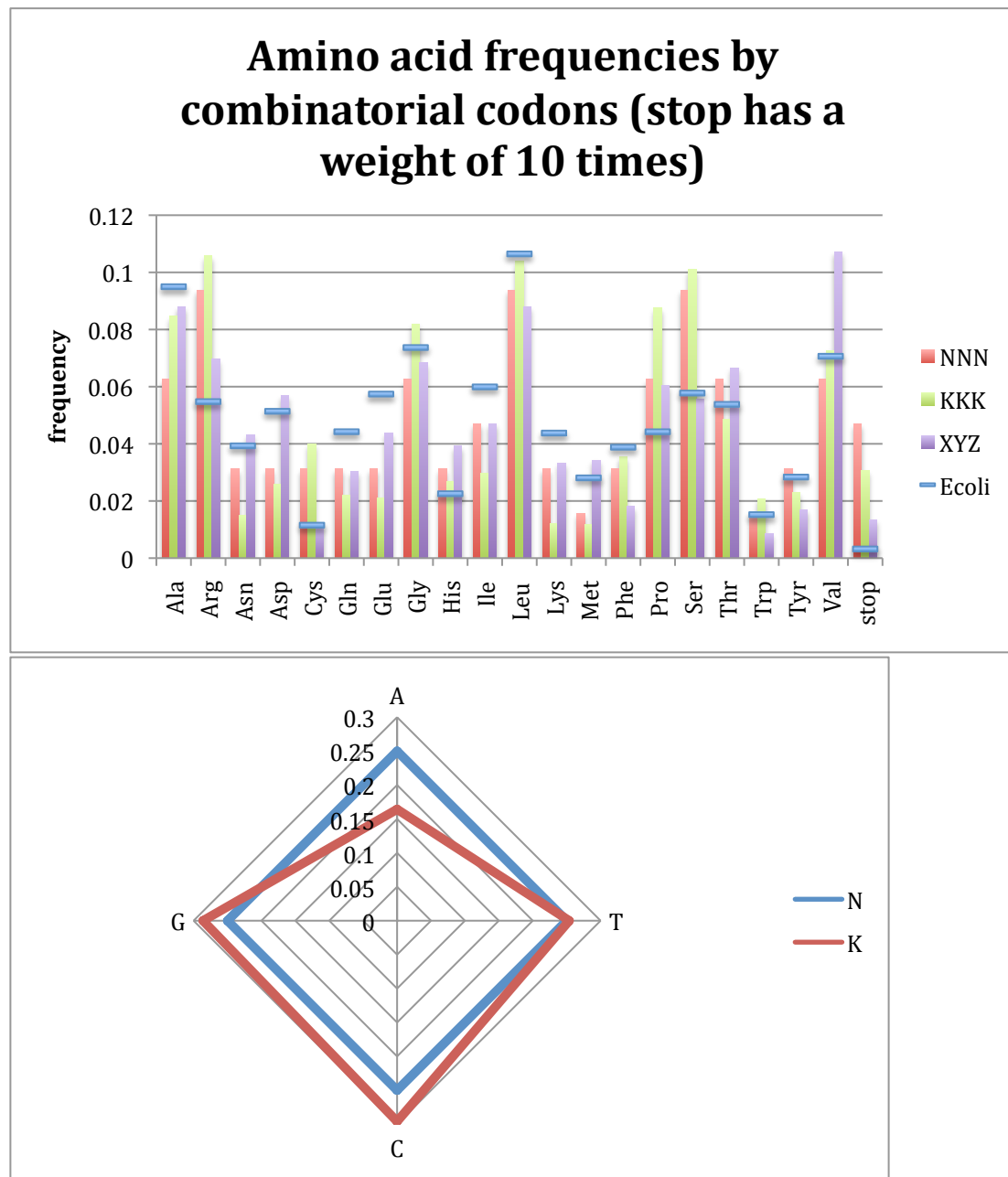
They do but it is not such a good fit ignoring them:



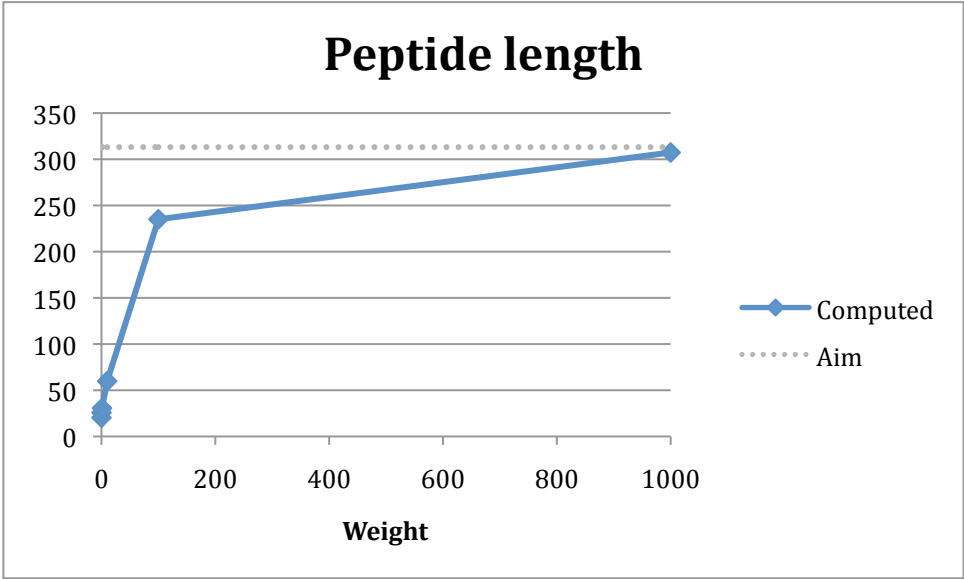
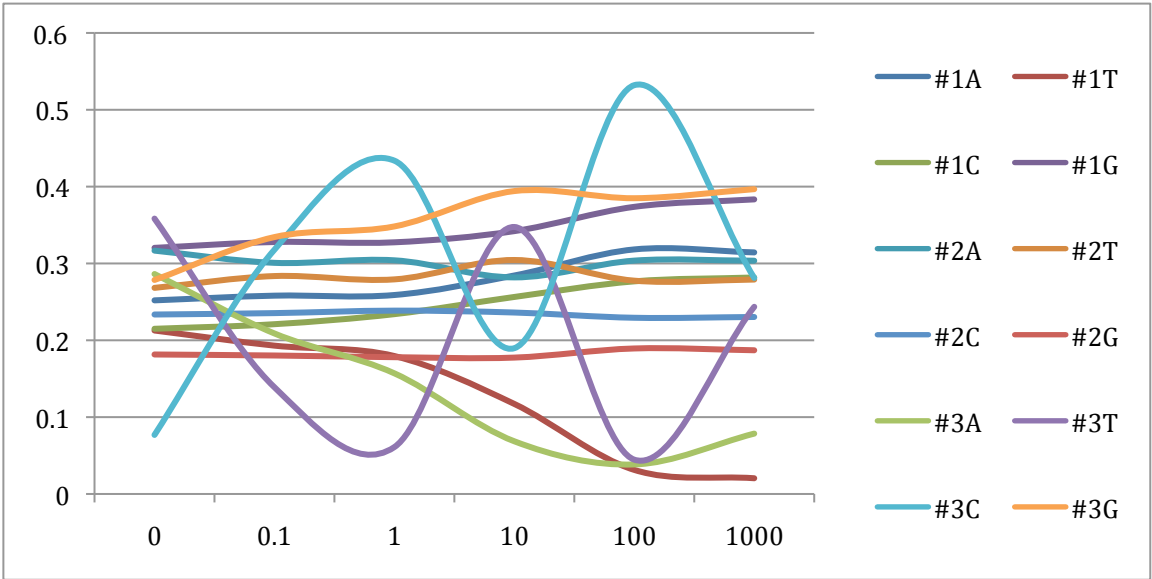
P. furiosus is a thermophile and it is the graph just to show that there is a large diversity in composition. To achieve the correct amount of stop codons, Cys, Phe, Trp and Tyr are lost.

More data beating the bush...

In the KKK model, there is absurdly high amount of proline, as the proportion of A decreases rapidly followed by T (stop = (A:0.5 T:0.3 G:0.2 C:0). Base frequency will be represented as radar plots for ease of visualization (fig 3).



The graph below shows the shift in bases according to several weights on stop (0, 0.1, 1, 10, 100, 1000). At the third position it does not matter if there is a T or C but A is not allowed there as is T in the first codon.



Notes:

matrix to optimize

$$\underline{x} = \begin{bmatrix} x_{a1} & x_{a2} & x_{a3} \\ x_{t1} & x_{t2} & x_{t3} \\ x_{g1} & x_{g2} & x_{g3} \\ x_{c1} & x_{c2} & x_{c3} \end{bmatrix}$$

$$\sum_{i=1}^4 x_{i1} = 1$$

$$\sum_{i=1}^4 x_{i2} = 1$$

$$\sum_{i=1}^4 x_{i3} = 1$$

objective function

$$\min f(x)$$

least squares: $f(\underline{x}) = \sum_{a=1}^{21} w(a) (\text{aim}_a - g_a(\underline{x}))^2$

$$g_a(\underline{x}) = \sum_{\text{codon } a} \left(\prod_{p=1}^3 m_{cp} x_p \right)$$

\underline{m}_c is a masking matrix & p is base 1, 2 or 3 of the codon

eg. $a = \text{tryptophan}$

$$m_{\text{try}} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

↓ ↓ ↓
p=1 2 3

NB. *codon_a depends on a (some amino acids are encoded by a larger number of codons)

1. Arkin, A.P. and D.C. Youvan, *Optimizing nucleotide mixtures to encode specific subsets of amino acids for semi-random mutagenesis*. Biotechnology (N Y), 1992. **10**(3): p. 297-300.
2. Wolf, E. and P.S. Kim, *Combinatorial codons: a computer program to approximate amino acid probabilities with biased nucleotide usage*. Protein Sci, 1999. **8**(3): p. 680-8.
3. Craig, R.A., et al., *Optimizing nucleotide sequence ensembles for combinatorial protein libraries using a genetic algorithm*. Nucleic Acids Research, 2010. **38**(2): p. e10.
4. Wang, W. and J.G. Saven, *Designing gene libraries from protein profiles for combinatorial protein experiments*. Nucleic Acids Research, 2002. **30**(21): p. e120.