

Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

➤ PAM matrices, Dayhoff *et al*, 1978

- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

Overall strategy for both PAM and BLOSUM

1. Trusted amino acid alignments
2. Obtain amino acid pair counts (A_{xy}^N) with corrections for
 - Evolutionary divergence
 - Sample biases
3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}^N
4. Log odds substitution matrix: $S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$

Log odds substitution matrices

Two sequences have N PAMs divergence, if, on average, N amino acid replacements per 100 residues occurred since their separation

$$S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$$

Scaling constant

Frequency of x aligned with y in sequences with divergence N

Frequency of x aligned with y in "random" sequences

PAM: A unit of evolutionary divergence

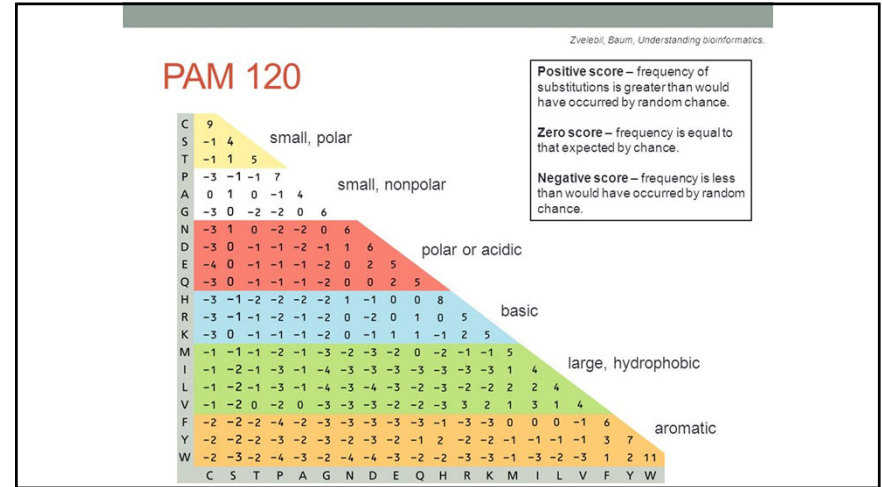
- PAM: Percent Accepted Mutation
 - *Accepted Mutations* are mutations that are retained and passed on to future generations
- We say the divergence between two sequences is N PAMs, if, on average, N amino acid replacements per 100 residues (including multiple substitutions) occurred since their separation.

3. Estimate substitution frequencies, q_{xy}^N , from pair counts, A_{xy}

- Markov model with 20 states (A, C, D, E ... Y)
- Estimate 1 PAM transition matrix P^1 from A_{xy}
- N-PAM transition matrix: $P^N = (P^1)^N$
- $q_{xy}^N = p_x \cdot P_{xy}^N$
- $S^N[x,y] = c \log \frac{q_{xy}^N}{p_x p_y}$

Is P_{xy}^N a symmetric matrix? No. (Check this algebraically).

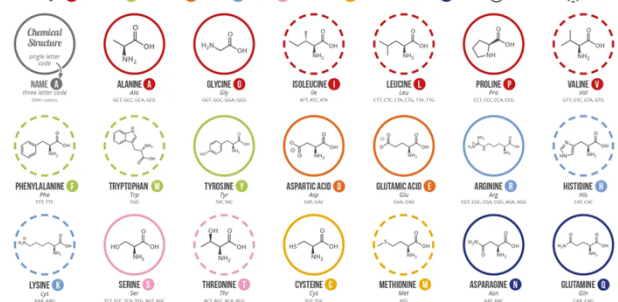
Is $S^N[x,y]$ a symmetric matrix? Yes (Check this algebraically).



A GUIDE TO THE TWENTY COMMON AMINO ACIDS

AMINO ACIDS ARE THE BUILDING BLOCKS OF PROTEINS IN LIVING ORGANISMS. THERE ARE OVER 500 AMINO ACIDS FOUND IN NATURE - HOWEVER, THE HUMAN GENETIC CODE ONLY DIRECTLY ENCODES 20. ESSENTIAL AMINO ACIDS MUST BE OBTAINED FROM THE DIET, WHILST NON-ESSENTIAL AMINO ACIDS CAN BE SYNTHESISED IN THE BODY.

Chart key: ● ALIPHATIC ● AROMATIC ● ACIDIC ● BASIC ● HYDROXYLIC ● SULFUR-CONTAINING ● AMIDIC ● NON-ESSENTIAL ○ ESSENTIAL



Note: This chart only shows those amino acids for which the human genetic code directly codes for. Selenocysteine is often referred to as the 21st amino acid, but is encoded in a special manner. In some cases, distinguishing between asparagine/aspartic acid and glutamine/glutamic acid is difficult. In these cases, the codes aaX (X) and aaY (Y) are respectively used.

© COMPAGND INTEREST 2014. WWW.COMPAGNDINTEREST.COM | Twitter: @compoundchem | Facebook: www.facebook.com/compoundchem
Shared under a Creative Commons Attribution-NonCommercial-NoDerivatives license.

Two widely used families of Amino Acid Substitution Matrices Parameterized for evolutionary divergence (N)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Henikoff & Henikoff, 1991

BLOSUM clustering example

- 1: KKRK
- 2: KKKK
- 3: KNRN
- 4: NRNR
- 5: KNKN
- 6: KRNR

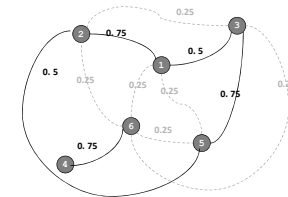
		Percent Sequence Identity				
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	

Unclassified sequences: Every sequence is at least 25% identical

		Percent Sequence Identity				
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	

- 1: KKRK
- 2: KKKK
- 3: KNRN
- 5: KNKN
- 4: NRNR
- 6: KRNR

< 45% identical



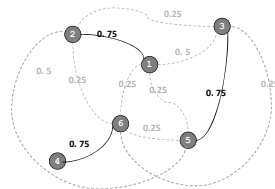
- 1: KKRK
- 2: KKKK
- 3: KNRN
- 5: KNKN
- 4: NRNR
- 6: KRNR

< 65% id

< 65% id

< 65% id

		Percent Sequence Identity				
	[2]	[3]	[4]	[5]	[6]	
[1]	0.75	0.5	0	0.25	0.25	
[2]		0.25	0	0.5	0.25	
[3]			0	0.75	0.25	
[4]				0	0.75	
[5]					0.25	



Two widely used families of Amino Acid Substitution Matrices
Parameterized for evolutionary divergence (*N*)

- PAM matrices, Dayhoff *et al*, 1978
- BLOSUM (Block Sum) matrices, Hennikoff & Hennikoff, 1991

Similarities and differences between PAM and BLOSUM

	PAM	BLOSUM
Evolutionary model	Explicit evolutionary model	None
Data	Full length MSAs of closely related sequences.	Conserved blocks. i.e., ungapped local MSAs
Bias correction	Trees	Clustering
Multiple substitutions	Markov model: $P^n = (P^1)^n$	Implicitly represented in data (clustering)
Evolutionary distance	Markov model: $P^n = (P^1)^n$	Clustering
Matrices	Transition and log odds scoring matrices	Log odds scoring matrix only.
Parameter n	Distance increases with n	Distance decreases with n
Biophysical properties	Derived indirectly from data	Derived indirectly from data

Comparing PAM and BLOSUM matrices

	PAM	Sequence identity	BLOSUM
	20	83%	
	30		
	60	63%	
	70		
	100	43%	90
	120	38%	80
	160	30%	60
	200	25%	50
	250	20%	45

More divergent ↓

↑ Less divergent