Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

# Encodings of amino acids and their impact on signal peptide prediction

Michał Burdukiewicz

University of Wrocław, Department of Genomics, Poland

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

# Outline

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

To reduce dimensionality of the problem, signalHsmm aggregates amino acids to four physicochemical groups.

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

| Groups |
| --- |
| H, K, R |
| C, F, I, L, M, U, V, W |
| N, Q, S, T |
| A, D, E, G, P, Y |

Classification of amino acids used by signalHsmm.

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

How optimal is this grouping? How does grouping of amino acids influence detection of signal peptides and cleavage sites?

Motivation        Size
Amino acid properties    Polarity
Methods        PI
Cross-validation    Hydrophobicity
Conclusion and perspectives    α-chain

# Outline

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Size
Polarity
PI
Hydrophobicity
$\alpha$-chain

Motivation
**Amino acid properties**
Methods
Cross-validation
Conclusion and perspectives

Size
**Polarity**
PI
Hydrophobicity
$\alpha$-chain

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Size
Polarity
PI
Hydrophobicity
$\alpha$-chain

Motivation
**Amino acid properties**
Methods
Cross-validation
Conclusion and perspectives

Size
Polarity
PI
**Hydrophobicity**
$\alpha$-chain

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Size
Polarity
PI
Hydrophobicity
α-chain

Motivation
Amino acid properties
**Methods**
Cross-validation
Conclusion and perspectives

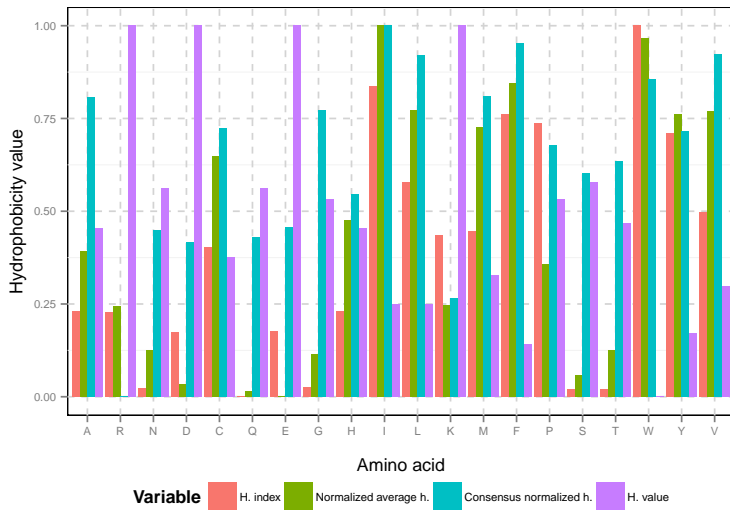Amino acid groupings
Cross-validation

# Outline

1. Motivation

2. Amino acid properties
   - Size
   - Polarity
   - PI
   - Hydrophobicity
   - $\alpha$-chain

3. Methods
   - Amino acid groupings
   - Cross-validation

4. Cross-validation
   - AUC
   - H
   - Specificity
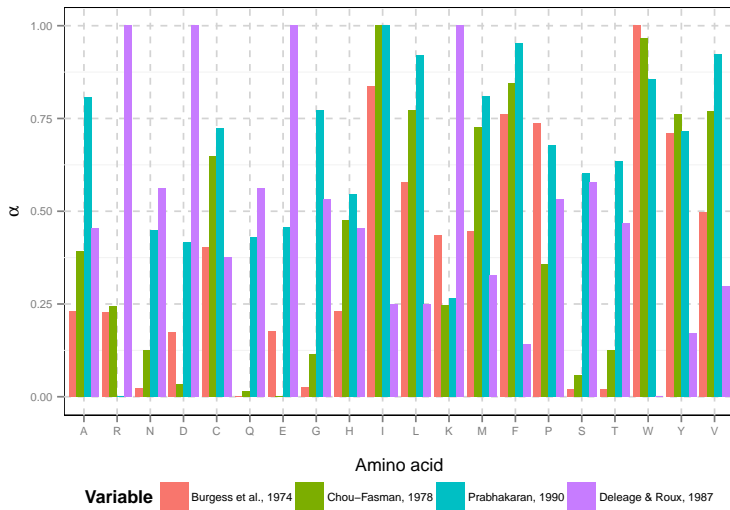   - Mean cleavage site displacement

5. Conclusion and perspectives

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Amino acid groupings
Cross-validation

For each combination of each variant of all properties:

1. calculate euclidean distance between amino acids;
2. cluster amino acids (complete-linkage clustering);
3. extract four highest clusters of amino acids.

67 amino acid groupings in total (signalHsmm, 2 standard, 65 created).

Motivation
Amino acid properties
**Methods**
Cross-validation
Conclusion and perspectives

Amino acid groupings
Cross-validation

For each grouping:

1. sample 3722 from 134044 negative sequences (to have balanced data set);

2. perform 5-fold cross-validation using all groupings (train an instance of signalHsmm using a single grouping);

3. repeat steps 1:2 15 times.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

# Outline

1. **Motivation**

2. Amino acid properties
   - Size
   - Polarity
   - PI
   - Hydrophobicity
   - $\alpha$-chain

3. Methods
   - Amino acid groupings
   - Cross-validation

4. **Cross-validation**
   - AUC
   - H
   - Specificity
   - Mean cleavage site displacement

5. Conclusion and perspectives

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

**AUC**
H
Specificity
Mean cleavage site displacement

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, F, I, L, M, V, W |
| K, R |
| C, G, N, P, Q, S, T, Y |
| D, E, H |

Best grouping of amino acids - AUC.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, L, M, N, Q, S, T, V |
| D, E, H, K, R |
| C, G |
| F, I, P, W, Y |

Worst grouping of amino acids - AUC.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
**H**
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, F, I, L, M, V, W |
| K, R |
| C, G, N, P, Q, S, T, Y |
| D, E, H |

Best grouping of amino acids - H.

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, L, M, N, Q, S, T, V |
| D, E, H, K, R |
| C, G |
| F, I, P, W, Y |

Worst grouping of amino acids - H.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
**Specificity**
Mean cleavage site displacement

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, C, G, S |
| H, K, R |
| F, I, L, M, N, P, Q, T, V, W, Y |
| D, E |

Best grouping of amino acids - specificity.

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, M |
| D, E, H, K, R |
| G, N, Q, S, T |
| C, F, I, L, P, V, W, Y |

Worst grouping of amino acids - specificity.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

Mean cleavage site displacement

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, C, G, N, P, Q, S, T |
| H, K, R |
| D, E |
| F, I, L, M, V, W, Y |

Best grouping of amino acids - mean cleavage site displacement.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, G, I, L, M, P, V |
| H, K, R |
| D, E |
| C, F, N, Q, S, T, W, Y |

Worst grouping of amino acids - mean cleavage site displacement.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, C, G, N, P, Q, S, T, V |
| H, K, R |
| D, E |
| F, I, L, M, W, Y |

Best (ex aequo) grouping of amino acids - median cleavage site displacement.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, C, G, I, N, P, Q, S, T, V |
| H, K, R |
| D, E |
| F, L, M, W, Y |

Best (ex aequo) grouping of amino acids - median cleavage site displacement.

Motivation
Amino acid properties
Methods
**Cross-validation**
Conclusion and perspectives

AUC
H
Specificity
Mean cleavage site displacement

| Groups |
| --- |
| A, G, I, L, M, P, V |
| H, K, R |
| D, E |
| C, F, N, Q, S, T, W, Y |

Worst grouping of amino acids - median cleavage site displacement.

Motivation
Amino acid properties
Methods
Cross-validation
**Conclusion and perspectives**

# Outline

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Grouping on amino acids has strong impact on misclassifications.

Motivation
Amino acid properties
Methods
Cross-validation
Conclusion and perspectives

Detection of signal peptides:

1. charged amino acids should be grouped separately;
2. aromatic amino acids should not be grouped together;

Motivation
Amino acid properties
Methods
Cross-validation
**Conclusion and perspectives**

Detection of cleavage sites:

1. hydrophibic amino acids should be grouped together;