

n-gram analysis of biological sequences in R

Michał Burdukiewicz¹, Piotr Sobczyk², Małgorzata Kotulska³, Paweł Mackiewicz¹

¹University of Wrocław, Department of Genomics,

²Wrocław University of Science and Technology, Faculty of Pure and Applied Mathematics,

³Wrocław University of Science and Technology, Department of Biomedical Engineering

Introduction

Biological sequences

Chains of amino acids (proteins) or nucleotides (RNA or DNA).

Sample protein sequence:

MKLLLLLIVSASMLIESLVNADGYIKRRDGCKVACLIGNE

GCDKECKAYGGSYGYCWTWGLACWCEGLPDDKTWKSETNT

CGGKK

Biological sequences



Mesobuthus martensii. Source: <http://www.sciencenews.org>

Toxin produced by *Mesobuthus martensii*:

MKLLLLLIVSASMLIESLVNADGYIKRRDGCKVACLIGNE

GCDKECKAYGGSYGYCWTWGLACWCEGLPDDKTWKSETNT

CGGKK

Biological sequences

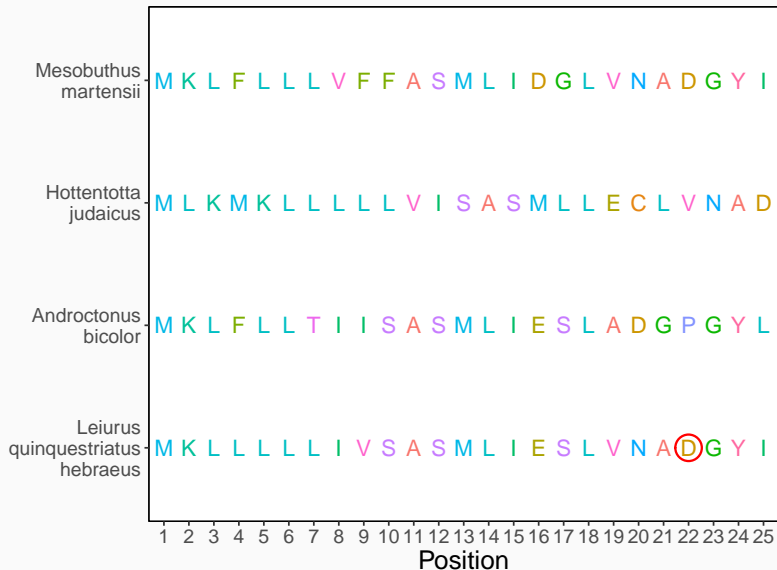
Signal peptide (red): n-terminal amino acid sequence directing proteins to the endomembrane system and next to extracellular localizations.

MKLLLLLIVSASMLIESLVNADGYIKRRDGCKVACLIGNE

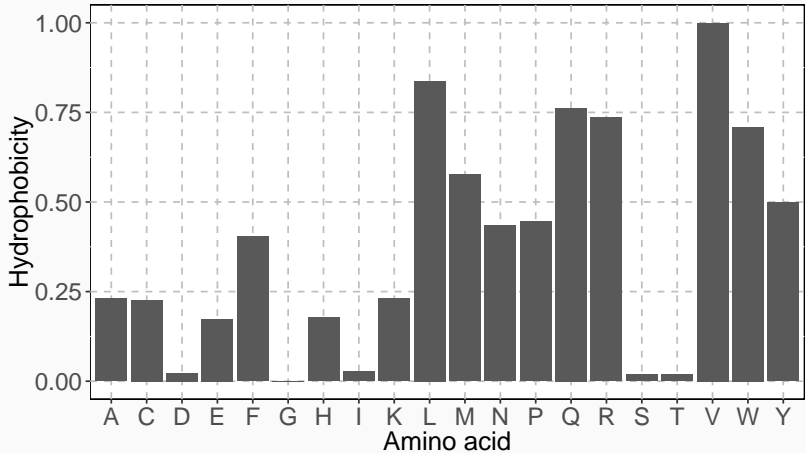
GCDKECKAYGGSYGYCWTWGLACWCEGLPDDKTWKSETNT

CGGKK

Biological sequences



Properties of amino acids



Amino acids may be described using their physicochemical properties.

Biological sequences

n-grams (k-tuples) are vectors of n characters derived from input sequence(s).