

# СКОРИНГОВАЯ ФУНКЦИЯ ДЛЯ ОБРАТНОЙ ЗАДАЧИ ВЫЧИСЛИТЕЛЬНОГО ПОСТРОЕНИЯ БЕЛКОВ CPD

---

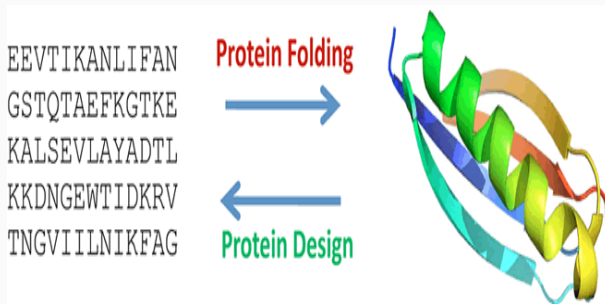
Александр Рубинштейн

8 декабря 2018

Отчет

- Основные идеи
- Составные части проекта
- Выводы

# ОСНОВНЫЕ ИДЕИ



**Figure:** Задача построения белков по последовательности аминокислот

Использование скоринговой функции:

$$\phi_{\vec{b}}(a) := E(\vec{a}, \vec{b}) = \sum_{j=1}^N \sum_{i=1}^N E_{ij}(a_i, a_j) \rightarrow \min_{a_1, a_2, \dots, a_N \in \mathcal{C}}. \quad (1)$$

$$E_{ij}(a_i, a_j) = K \cdot (\vec{w}^T \vec{f}_{ij}(a_i, a_j, \vec{b}_0) + C) \quad (2)$$

Генерация матриц энергии:

$$Q = \begin{bmatrix} [E_{11}] & [E_{12}] & \cdots & [E_{1N}] \\ [E_{21}] & [E_{22}] & \cdots & [E_{2N}] \\ \vdots & \vdots & \ddots & \vdots \\ [E_{N1}] & [E_{N2}] & \cdots & [E_{NN}] \end{bmatrix},$$

$$E_{ij} = \begin{bmatrix} E_{ij}(c_1, c_1) & E_{ij}(c_1, c_2) & \cdots & E_{ij}(c_1, c_{20}) \\ E_{ij}(c_2, c_1) & E_{ij}(c_2, c_2) & \cdots & E_{ij}(c_2, c_{20}) \\ \vdots & \vdots & \ddots & \vdots \\ E_{ij}(c_{20}, c_1) & E_{ij}(c_{20}, c_2) & \cdots & E_{ij}(c_{20}, c_{20}) \end{bmatrix},$$

# ОСНОВНЫЕ ИДЕИ

Решение задачи оптимизации:

$$\begin{aligned} \min_{\vec{x} \in \{0,1\}^{20N}} \quad & \vec{x}^T Q \vec{x} \\ \text{s.t.} \quad & A \vec{x} = \mathbf{1}_N, \end{aligned} \tag{3}$$

где

$$A = \begin{bmatrix} 1 \dots 1 & 0 \dots 0 & \dots & 0 \dots 0 \\ 0 \dots 0 & 1 \dots 1 & \dots & 0 \dots 0 \\ \vdots & \vdots & \ddots & \vdots \\ \underbrace{0 \dots 0}_{20} & \underbrace{0 \dots 0}_{20} & \dots & \underbrace{1 \dots 1}_{20} \end{bmatrix}, \quad A \in \{0,1\}^{N \times 20N}.$$

## Три основные стадии:

- Извлечение признаков из модели белка
- Обучение на CASP и предсказание GDT-TS
- Присвоение энергии (скора) каждому белку

$B(y_i, \hat{y}_j)$  соответствует скору белка при замене  $y_i$  на  $\hat{y}_j$  в  $j$ -ой позиции последовательности.

$$S(\vec{y}_{\text{nat}}, \vec{y}_{\text{pred}}) = \frac{\sum_{k=1}^N B((\vec{y}_{\text{nat}})_k, (\vec{y}_{\text{pred}})_k)}{\sum_{k=1}^N B((\vec{y}_{\text{nat}})_k, (\vec{y}_{\text{nat}})_k)} \quad (4)$$

Хорошее предсказание соответствует  $S > 0$ , а лучшие результаты имеют  $S$  близкую к 1. Отрицательные  $S$  означают предсказания плохого качества.