

**Многомодальный метод релевантных векторов
в задаче распознавания
вторичной структуры белка**

Сунгуров Д.С.

Московский физико-технический институт
Кафедра «Интеллектуальные системы»

Научный руководитель
Д.т.н., профессор В.В. Моттль

2013

Представление белка

1. Первичная структура $\omega = (\alpha_t, t = 1, \dots, M_\omega)$ - конечная аминокислотная последовательность длины M_ω , где $\alpha_t \in A = (\alpha^1, \dots, \alpha^m)$, $m = 20$ - символы алфавита аминокислот
2. Вторичная структура: $y = (y_t, t = 1, \dots, M_\omega)$, $y_t \in Y = \{h, s, l\}$, где h, s, l - три возможных типа элементов вторичной структуры

Задача предсказания вторичной структуры белка

По обучающей выборке $\{(\omega_l, y_l) = l = 1, \dots, N\}$ для всякой новой аминокислотной последовательности $\omega = (\alpha_t, t = 1, \dots, M_\omega)$ научиться предсказывать вторичную структуру соответствующего белка $\hat{y} = (\hat{y}_t, t = 1, \dots, M_\omega)$.

Существующие методы предсказания вторичной структуры

1. **Общепринятый принцип прогнозирования вторичной структуры белка (Helix, Strand, Coil) заключается в оценивании структуры по локальному фрагменту аминокислотной цепи (симметричному окну).**
2. **Искомая классификация аминокислотного окна неизбежно должна быть результатом его сравнения с множеством «эталонных» окон из обучающей совокупности.**
3. **Известные методы прогнозирования вторичной структуры:**
 - базируются на целом ряде чисто биологических предположений,
 - используют эталонный набор из десятков тысяч фрагментов.

Существующие методы предсказания вторичной структуры

4. Общепринятый принцип прогнозирования вторичной структуры белка (**Helix, Strand, Coil**) заключается в оценивании структуры по локальному фрагменту аминокислотной цепи (симметричному окну).
5. Искомая классификация аминокислотного окна неизбежно должна быть результатом его сравнения с множеством «эталонных» окон из обучающей совокупности.
6. Известные методы прогнозирования вторичной структуры:
 - базируются на целом ряде чисто биологических предположений,
 - используют эталонный набор из десятков тысяч фрагментов.

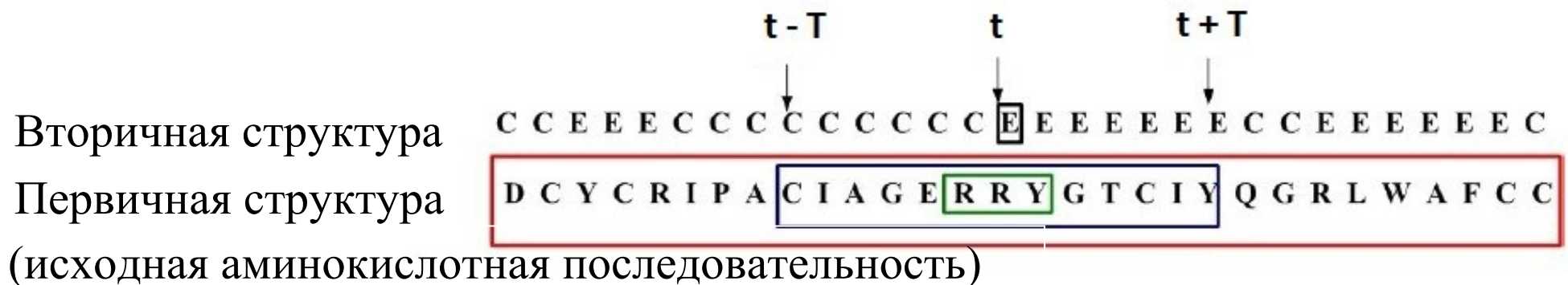
Основная идея предлагаемого подхода

1. Рассматривать задачу классификации аминокислотного фрагмента как задачу распознавания без привлечения биологических предположений.
2. Предлагаемый аппарат обучения опирается на большое множество «пробных» функций парного сравнения аминокислотных окон, из которого автоматически отбирается небольшое число «релевантных» функций сравнения.
3. Разработана новая версия метода релевантных векторов (RVM), названная многомодальным методом релевантных векторов (Multimodal RVM).

Метод скользящего окна

Решение о классе элемента вторичной структуры на позиции t делается на основе симметричного интервала $\omega = (\alpha_t, t - T \leq t \leq t + T)$ вырезанного из всей аминокислотной последовательности $\omega = (\alpha_t, t = 1, \dots, M_\omega)$.

Мы ограничиваемся здесь распознаванием двух классов структуры $(S, \bar{S}) = (1, -1)$ – Strand, Non-strand



Исходная задача сократилась до последовательности независимых задач определения класса вторичной структуры $(1, -1)$ аминокислоты в очередном окне.

Принцип прогнозирования структуры аминокислотного окна

Аминокислотное окно: $\omega = (\alpha_{t-T} \cdots \cdots$

Вторичная структура: $y_t \in \{-1, 1\}$

Обучающая совокупность аминокислотных окон: $\{(\omega_j, y_j), j = 1, \dots, N\}$

Искомое решающее правило для нового окна:

$$d(\omega | \mathbf{a}, b) = \sum_{j=1}^N a_j \underbrace{\quad}_{\text{функция парного сравнения}} \cdot b \begin{cases} > 0 \rightarrow \hat{y}(\omega) = 1 \\ < 0 \rightarrow \hat{y}(\omega) = -1 \end{cases}$$

неизвестные коэффициенты

Классическая идея релевантных векторов (Relevance Vector Machine – RVM):

Найти a_j , оставив ненулевыми $a_j \neq 0$ только коэффициенты при «нужных» (релевантных) объектах обучающей совокупности.

Многомодальный метод релевантных векторов: множество функций парного сравнения

$$d(\omega | \mathbf{a}, b) = \sum_{i \in I} \sum_{j=1}^N a_{ij} S_i(\omega_j, \omega) + b \begin{cases} > 0 \rightarrow \hat{y}(\omega) = 1 \\ < 0 \rightarrow \hat{y}(\omega) = -1 \end{cases}$$

множество объектов обучающей совокупности и множество функций парного сравнения

Найти a_j , оставив $a_j \neq 0$ только коэффициенты при релевантных объектах и релевантных функциях (модальностях) сравнения

Многомодальный метод релевантных векторов

$$\begin{cases} \sum_{i=1}^n \sum_{l=1}^N [(1-\mu)a_{il}^2 + \mu|a_{il}|] + C \sum_{j=1}^N \delta_j \rightarrow \min(a_{il}, b, \delta_j), \\ y_j \left(\sum_{i=1}^n \sum_{l=1}^N a_{il} S_i(\omega_l, \omega_j) + b \right) \geq 1 - \delta_j, \delta_j \geq 0, j=1, \dots, N. \end{cases}$$

Параметр селективности $\mu > 0$. Чем больше μ , тем больше число нулевых коэффициентов $a_{ij} = 0$

Итоговая дискриминантная гиперплоскость:

$$d(\omega) = \sum_{ij \in \hat{F}} a_{ij} S_i(\omega_j, \omega) + b \geq 0, \quad \underbrace{\hat{F} = \{ \dots \}}_{\substack{\text{Подмножество} \\ \text{релевантных} \\ \text{вторичных признаков}}} = \underbrace{F = \{ \dots \}}_{\substack{\text{полное множество} \\ \text{вторичных признаков}}}.$$

$\hat{J} = \{ j: \exists i (a_{ij} \neq 0) \} \subseteq J = \{1, \dots, N\}$	подмножество релевантных аминокислотных окон в обучающей совокупности
$\hat{I} = \{ i: \exists j (a_{ij} \neq 0) \} \subseteq I = \{1, \dots, n\}$	подмножество релевантных функций парного сравнения окон

В нашей задаче прогнозирования вторичной структуры белка: $n = 6$ шесть функций парного сравнения аминокислотных окон

$S_1(\omega', \omega'')$	классический позиционный метод сравнения аминокислотных окон	$S_4(\omega', \omega'')$	вновь разработанный метод на основе преобразования Фурье
$S_2(\omega', \omega'')$		$S_5(\omega', \omega'')$	
$S_3(\omega', \omega'')$		$S_6(\omega', \omega'')$	

Математическая структура функций парного сравнения

Еще раз шесть функций парного сравнения:

$S_1(\omega', \omega'')$	классический позиционный метод сравнения аминокислотных окон	$S_4(\omega', \omega'')$	вновь разработанный метод на основе преобразования Фурье
$S_2(\omega', \omega'')$		$S_5(\omega', \omega'')$	
$S_3(\omega', \omega'')$		$S_6(\omega', \omega'')$	

Сравниваемые окна ω' и ω''

Векторы признаков $\mathbf{z}(\omega')$ и $\mathbf{z}(\omega'')$ – разные в позиционном и тригономет. сравнениях

Математические способы сравнения векторов признаков

$S_1(\omega', \omega'')$ } $S_4(\omega', \omega'')$ } = $\mathbf{z}^T(\omega')\mathbf{z}(\omega'')$ – классическое скалярное произведение

$S_2(\omega', \omega'')$ } $S_5(\omega', \omega'')$ } = $\exp\left(-\gamma\|\mathbf{z}(\omega') - \mathbf{z}(\omega'')\|_2^2\right)$ – радиальная потенц. функция с нормой l_2

$S_3(\omega', \omega'')$ } $S_6(\omega', \omega'')$ } = $\|\mathbf{z}(\omega') - \mathbf{z}(\omega'')\|_1$ – норма l_1

Функции парного сравнения аминокислотных окон на основе преобразования Фурье

Способ основан на представлении каждой из 20-ти аминокислот

$$A = \{\alpha_1, \dots, \alpha_{20}\}$$

20-мерным вектором ее частотных признаков согласно теории РАМ – Point Accepted Mutation (Margaret Dayhoff, В.В. Сулимова).

$$\alpha^k \rightarrow \mathbf{a}^k \in \mathbb{R}$$

Аминокислотное окно:

$$\begin{array}{c}
 -T \quad \rightarrow \quad s \quad \rightarrow \quad T \\
 1 \left[\begin{array}{ccc} a_{t-T}^1 & \dots & \dots \\ \vdots & \vdots & \vdots \\ a_{t-T}^{20} & \dots & \dots \end{array} \right. \\
 k \downarrow \\
 20
 \end{array}
 \begin{array}{ccc}
 \mathbf{a}_{t-T} & \mathbf{a}_t & \mathbf{a}_{t+T}
 \end{array}$$

Каждая k -я строка рассматривается как скалярный сигнал $a_s^k, s = -T, \dots$...
Этот преобразуется в амплитудный спектр Фурье с пятью гармониками $[z_i^k, i = 0, \dots$

Игнорирование фазовых компонент оставляет спектральное представление окна сильно-чувствительным к последовательности аминокислот в нем ($\alpha_{t-T} \dots$), но делает его малочувствительным к сдвигам окна t .

Общий вектор признаков окна: $\mathbf{z}(\omega) = [z_{il}(\omega), i = 1, \dots, 20, l = 0, \dots, 4] \in \mathbb{R} \quad \mathbb{R}$

Численное решение задачи обучения

Еще раз исходная задача многомодального метода релевантных векторов:

$$\left\{ \begin{array}{l} \sum_{i=1}^n \sum_{l=1}^N \left[(1-\mu)a_{il}^2 + \mu |a_{il}| \right] + C \sum_{j=1}^N \delta_j \rightarrow \min(a_{il}, b, \delta_j), \\ y_j \left(\sum_{i=1}^n \sum_{l=1}^N a_{il} S_i(\omega_l, \omega_j) + b \right) \geq 1 - \delta_j, \delta_j \geq 0, j = 1, \dots, N. \end{array} \right. \quad \begin{array}{l} \text{Задача выпуклого} \\ \text{программирования} \end{array}$$

Двойственная задача выпуклого программирования:

$$\left\{ \begin{array}{l} W(\lambda_1, \dots, \lambda_N | \mu) = \sum_{j=1}^N \lambda_j - \\ \frac{1}{4(1-\mu)} \sum_{i=1}^n \sum_{l=1}^N \left\{ \min \begin{array}{l} \mu + \sum_{j=1}^N y_j \lambda_j x_{il,j} \\ 0 \\ \mu - \sum_{j=1}^N y_j \lambda_j x_{il,j} \end{array} \right\}^2 \\ \sum_{j=1}^N y_j \lambda_j = 0, \quad 0 \leq \lambda_j \leq C, \quad j = 1, \dots, N. \end{array} \right\} \rightarrow \max, \quad \left\{ \begin{array}{l} \hat{a}_{il} = \frac{1}{2(1-\mu)} \left(\sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{il,j} + \mu \right), \\ \sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{il,j} < -\mu, \\ \hat{a}_{il} = 0, \quad -\mu \leq \sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{il,j} \leq \mu, \\ \hat{a}_{il} = \frac{1}{2(1-\mu)} \left(\sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{il,j} - \mu \right), \\ \sum_{j: \hat{\lambda}_j > 0} y_j \hat{\lambda}_j x_{il,j} > \mu, \end{array} \right.$$

Число переменных равно числу
аминокислотн. окон в обучающей
совокупности

$$\text{Результат обучения} \quad d(\omega) = \underbrace{\nabla_{a, S(\omega, \omega) + b}}_{\text{только релевантные окна и функции сравнения}} \begin{cases} \text{strand} > 0 \\ \text{ne strand} < 0 \end{cases}$$

:
только релевантные окна и функции сравнения

Требования к алгоритму

- Алгоритм поиска должен легко параллелиться
- Годятся алгоритмы делающие небольшое количество итераций
- Каждая итерация должна легко параллелиться
- По этим причинам мы решили использовать метод внутренней точки для поиска максимума исследуемой функции

Требования к алгоритму

- Алгоритм поиска должен легко параллелиться
- Годятся алгоритмы делающие небольшое количество итераций
- Каждая итерация должна легко параллелиться
- По этим причинам мы решили использовать метод внутренней точки для поиска максимума исследуемой функции

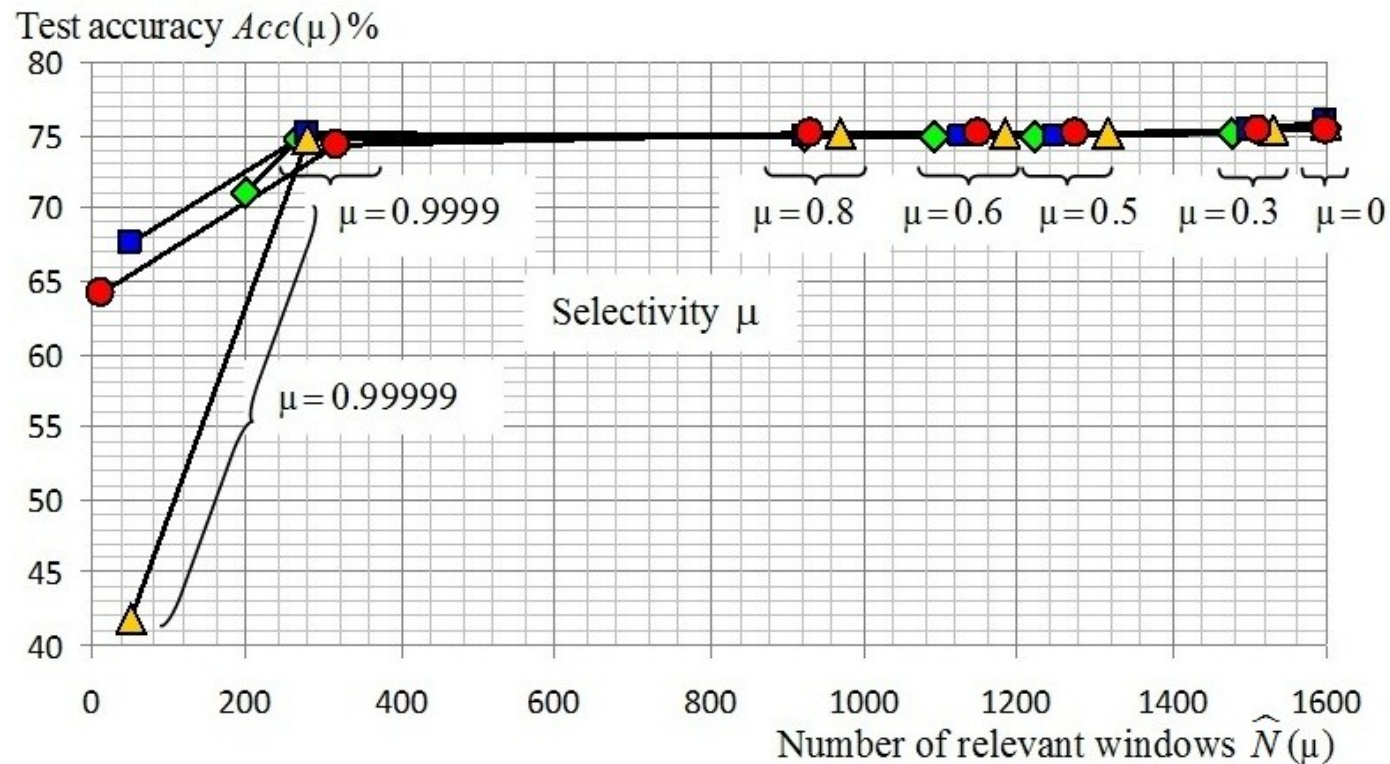
Особенности использования метода внутренней точки

- Две тяжелые операции: перемножение матриц и решение системы линейных уравнений.
- Мы распараллелили умножение матриц на CUDA.
- Решение системы линейных уравнений распараллелить не удалось.
- Небольшое количество итераций: 40 – 50.
- Метод внутренней точки является методом второго порядка, т.е требует держать в памяти матрицы $(N \times N)$, где N это количество переменных в двойственной задаче.

Эксперимент

- Мы использовали набор данных RS126 который содержит 126 белков представленных первичной и вторичной структурой. Белки в данном наборе имеют схожесть менее 25% для фрагментов более 80 аминокислот.
- Средняя длина белка в наборе 300-400 аминокислот.
- Проблема распознавания стрэндов во вторичной структуре решалась как задача двух классово́й классификации.
- Мы использовали длину окна в $2T + 1 = 35$. При этой длине окна из данных было получено 19075 аминокислотных фрагментов.
- Использовалось 6 различных функций парного сравнения.
- Данные разбивались случайным образом на обучение размера $|\Omega_{tr}| = 1600$ и контроль $|\Omega_{test}| = 17475$.
- Обучение проводилось для 7 различных значений параметра селективности: $\mu = 0$, $\mu = 0.3$, $\mu = 0.5$, $\mu = 0.6$, $\mu = 0.8$, $\mu = 0.9999$ and $\mu = 0.99999$.
- Всякий раз вычислялась доля ошибок прогнозирования вторичной структуры на контрольной совокупности.

Результаты экспериментов



- Точность определения стрендов остается на одном уровне в 75% для всех значений селективности от $\mu \geq 0$ до $\mu = 0.9999$.
- С селективностью $\mu = 0.9999$, только $\hat{n} = 3$ функции сравнения и $\hat{N} = 300$ из 1600 объектов обучающей выборки были признаны релевантными.
- Дальнейшее увеличение селективности $\mu = 0.99999$ влечет сильную потерю в точности.

Заключение

- В этой работе был применен многомодальный метод релевантных векторов, показавший свою эффективность.
- Количество эталонных объектов используемых в решающем правиле в 5 раз меньше обучающей совокупности.
- Только 3 из 6 функций парного сравнения используются в решающем правиле.
- Средняя точность распознавания стрендов оказалась равной примерно 75%, что находится на уровне уже существующих алгоритмов
- Особенно интересно то, что не было замечено переобучения, несмотря на то, что размерность векторов вторичных признаков в несколько раз выше размера обучающей совокупности.