

# Выбор оптимальных моделей локальной аппроксимации для классификации временных рядов

Сергей Дмитриевич Иванычев

Московский физико-технический институт  
Физтех-школа прикладной математики и информатики  
Факультет управления и прикладной математики  
Кафедра «Интеллектуальные системы»

Научный руководитель: д.ф.-м.н. В.В. Стрижов

Выпускная квалификационная работа бакалавра

Москва 2018

# Классификация временных рядов

## Цель

Предложить способ построения набора моделей локальной аппроксимации для устойчивой классификации сигналов носимых устройств.

## Гипотеза

Суперпозиция моделей локальной аппроксимации доставляет более высокое качество при меньшей сложности чем универсальные модели.

## Прямая задача

Исследование статистических свойств промежуточного параметрического пространства, строящегося моделями локальной аппроксимации.

## Обратная задача

Оптимизировать структурные параметры выбираемых моделей по порождающей выборке с целью получения выборки с оптимальными свойствами.

- Кузнецов М. П., Ивкин Н. П., *Алгоритм классификации временных рядов акселерометра по комбинированному признаковому описанию*, 2015.
- Карасиков М. Е., Стрижов В. В. *Классификация временных рядов в пространстве параметров порождающих моделей*, 2016.
- Артемов А. В., *Математические модели временных рядов с трендом в задачах обнаружения разладки*, 2016.

## Задан временной ряд

$$S : T \rightarrow \mathbb{R}, \text{ где } T = \{t_0, t_0 + d, t_0 + 2d \dots\}.$$

## Определен сегмент временного ряда

$$\mathbf{x}_i = [S(t_i), S(t_i - d), S(t_i - 2d), \dots, S(t_i - (n-1)d)]^T, \quad \mathbf{x}_i \in X \equiv \mathbb{R}^n.$$

Задана выборка  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^I$ ,  $y_i \in \{1, 2, \dots, K\}$ .

$\mathbf{X}$  — набор сегментов данных акселерометра,

$\mathbf{y}$  — метки классов движения (бег, ходьба, подъем и спуск по лестнице).

Задан  $\mathbf{h}$  — конечный набор моделей локальной аппроксимации.

# Постановка задачи классификации

## Модель локальной аппроксимации

$$g_i(\mathbf{w}, \mathbf{x}) \in \mathbf{X}, \text{ где } \mathbf{w} \in \mathbb{R}^{n_g}.$$

Оптимальные параметры определяются как

$$\mathbf{h}_i(\mathbf{x}) = \arg \min_{\mathbf{w} \in \mathbb{R}^{n_g}} \rho(g(\mathbf{w}, \mathbf{x}), \mathbf{x}),$$

$\mathbf{h}_i$  — модель локальной аппроксимации.

Набор функций  $\mathbf{h} = [\mathbf{h}_1 \dots \mathbf{h}_k] : \mathbf{x} \mapsto [w_1^* \dots w_k^*]$  отображает пространство сегментов  $\mathbf{X}$  в промежуточное пространство признаков описаний  $\mathbf{Z}$ .

## Модель классификации

$$T \rightarrow \mathbf{X} \xrightarrow{\mathbf{h}} \mathbf{Z} \xrightarrow{a} Y,$$

$\mathbf{h}$  — набор моделей локальной аппроксимации,  $a(\cdot, \gamma)$  — многоклассовый классификатор.

Минимизация функций ошибки каждой модели локальной аппроксимации

$$\arg \min_{\mathbf{w} \in W} L_g(\mathbf{X}, \mathbf{w}) = \arg \min_{\mathbf{w} \in W} \sum_{i=1}^l \sum_{k=1}^n \|g(\mathbf{w}, \mathbf{x}_i) - \mathbf{x}_i\|_2^2$$

Оптимизация функции ошибки обобщенной линейной модели

$$\arg \min_{\theta \in \Theta} L_a(\mathbf{Z}, \mathbf{y}, \theta) = \arg \min_{\theta \in \Theta} \left[ - \sum_{i=1}^l \sum_{k=1}^K [y_i = k] \log P(y_i = k | \mathbf{z}_i, \theta) \right]$$

## Модели локальной аппроксимации

Модель	Структурные параметры
SEMOR	-
AR-авторегрессия	порядок
Фурье-модель (FFT)	количество главных частот
Сингулярного спектр (SSA)	количество сингулярных чисел

## AR-авторегрессия

**Структурный параметр:** порядок  $m$ ,

$$g_{AR}(w, x) = \hat{x}, \text{ где } \hat{x}_i = \begin{cases} x_k & \text{при } k \in [1, m], \\ w_0 + \sum_{i=1}^m w_i x_{k-i} & \text{при } k \in [m+1, n]. \end{cases}$$

## Сингулярный спектр (SSA)

**Структурный параметр:** количество главных собственных значений  $k$ . Сингулярное разложение траекторной матрицы,

$$STS = VHV^T, H = \text{diag}(\lambda_1 \dots \lambda_m),$$

параметры образуют  $k$  главных собственных значения.



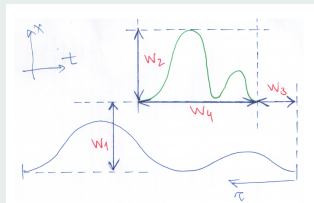
## Фурье-модель (FFT)

**Структурный параметр:**  $k$  частот из прямого преобразования Фурье, соответствующие наибольшим амплитудам

$$w_{2j} = \operatorname{Re} \sum_{k=1}^n x_k \exp\left(-\frac{2\pi i}{n}kj\right), \quad w_{2j+1} = \operatorname{Im} \sum_{k=1}^n x_k \exp\left(-\frac{2\pi i}{n}kj\right)$$

## Self-Modeling Regression

Накладывание шаблона линейными преобразованиями  $x$  и  $t$ .



$$g(\mathbf{x}, \mathbf{w}) = w_1 + w_2 p(w_3 + w_4 t),$$

$$w_{\text{SEMOR}} = [\hat{w}_1, \hat{w}_2, \hat{w}_3, \hat{w}_4, \rho].$$

# Построение промежуточной выборки и оптимизация функции потерь обобщенной линейной модели

- 1 Для каждого  $\mathbf{h}_i \in \mathbf{h}$  вычисляем

$$[\mathbf{z}_i^1 \dots \mathbf{z}_i^k]^T = [\mathbf{h}_i(\mathbf{x}_1) \dots \mathbf{h}_i(\mathbf{x}_k)]$$

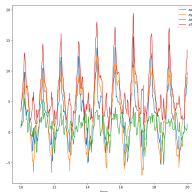
- 2 Конкатенируем вектора параметров  $\mathbf{z}_i = (\mathbf{z}_i^1 \dots \mathbf{z}_i^k)$ , то есть  $\mathbf{z}_i = \mathbf{h}(\mathbf{x}_i)$ . Получили выборку в промежуточном пространстве  $\mathbf{Z}$ .
- 3 Минимизируем функции потерь обобщенной линейной модели

$$\hat{\theta} = \arg \min_{\theta \in \Theta} L(f(\mathbf{Z}), \mathbf{y}).$$

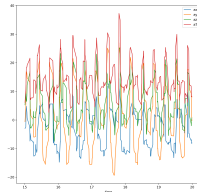
**Данные с акселерометра:** 4 типа движения, частота дискретизации 100 Гц.

**Сегментация:** локальные экстремумы с окном и квантиль по длине сегментов.

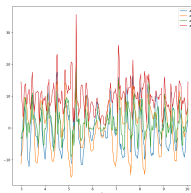
**Нормализация:** приведение к одной размерности с помощью кубических сплайнов.



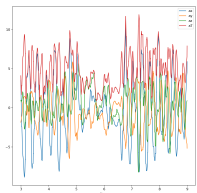
Ходьба



Бег



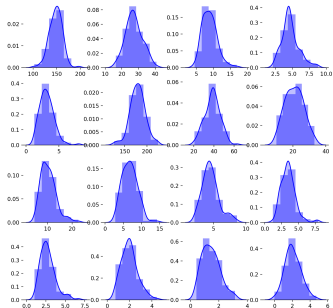
Вверх



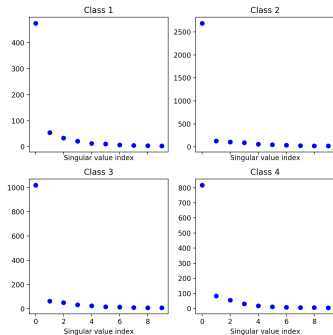
Вниз

Тесты простоты выборки: ( $T$ -тест)  $\mathbb{E}\varepsilon = 0, D\varepsilon = \text{const}$ , а также

анализ унимодальности распределений



анализ спектра выборки



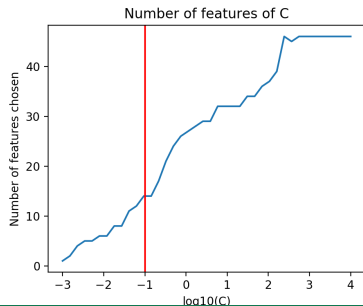
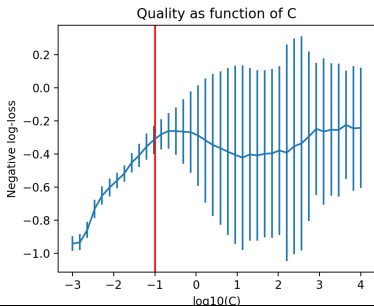
# Обобщенная линейная модель: отбор признаков

Сравниваем обобщающую способность обобщенной линейной модели (GLE) с универсальной моделью при одинаковой сложности.

Определим сложность модели как

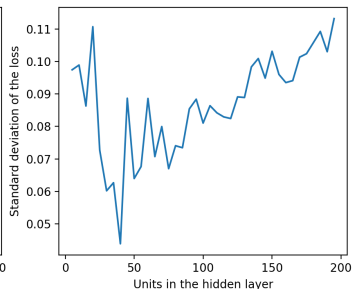
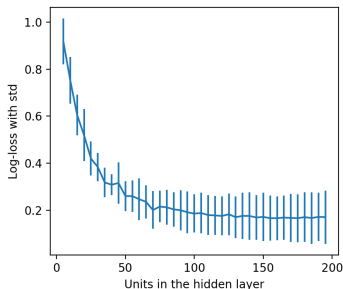
$$\text{Comp}(\mu) = \#|\text{neurons in the hidden layer}|$$

Отбираем признаки в  $(\mathbf{Z}, \mathbf{y})$  для обобщенной линейной модели. Логистическая регрессия с  $L_1$  регуляризацией.



На выборке  $(\mathbf{X}, \mathbf{y})$  оптимизируем параметры двуслойной нейронной сети (NN). Получаем зависимости

$$L(\text{Comp}), D_L(\text{Comp}).$$



# Сравнение ошибки при разных сложностях универсальной модели

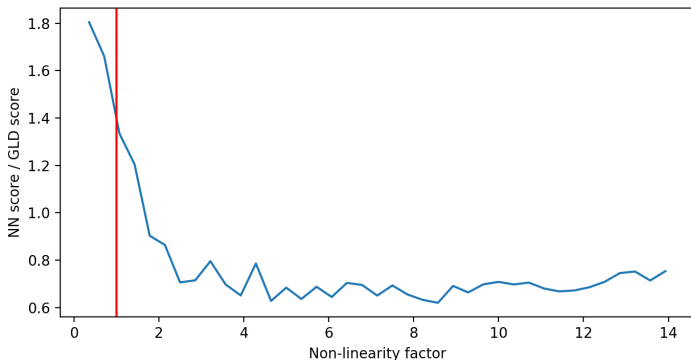


Рис.: Отношение ошибок от отношения сложностей

Результат: при  $\text{Comp}(\text{GLE}) = \text{Comp}(\text{NN})$ , имеем

$$\frac{L(\text{NN})}{L(\text{GLE})} = 1.4, \frac{D_L(\text{NN})}{D_L(\text{GLE})} > 1.$$

- Предложен и реализован способ построения набора моделей локальной аппроксимации для устойчивой классификации сигналов носимых устройств, тест простоты выборки в промежуточном пространстве признаков описаний а также методика оценки ее обобщающей способности по сравнению с универсальными моделями.
- Исследованы статистические свойства промежуточного пространства признаков описаний временных рядов. Выборка в промежуточном пространстве простая, а аппроксимирующие ее линейная модель являются адекватной.
- GLM адекватнее разделяет выборку чем универсальная модель, то есть при одинаковой сложности обеспечивает более высокое качество и меньше переобучается.