

**ПРАВИТЕЛЬСТВО РОССИЙСКОЙ ФЕДЕРАЦИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
«ВЫСШАЯ ШКОЛА ЭКОНОМИКИ»**

Факультет компьютерных наук

Образовательная программа бакалавриата «Программная инженерия»

УДК 004.8

СОГЛАСОВАНО

профессор департамента анализа данных и
искусственного интеллекта факультета
компьютерных наук НИУ ВШЭ, д-р физ.-
мат. наук

_____ В. А. Громов
«__» _____ 2022 г.

УТВЕРЖДАЮ

Академический руководитель
образовательной программы
«Программная инженерия»
профессор департамента программной
инженерии, канд. техн. наук

_____ В. В. Шилов
«__» _____ 2022 г.

**Выпускная квалификационная работа
(академическая)**

на тему: **Прогнозирование хаотических временных рядов: алгоритм self-healing для
прогнозирования с помощью кластеризации**
по направлению подготовки 09.03.04 «Программная инженерия»

Приложения

Выполнила студентка
образовательной программы
09.03.04 «Программная инженерия»
группы БПИ182

_____ А.Т. Антонова
«__» _____ 2022 г.

Оглавление

Приложение А. Условные обозначения алгоритмов	3
Приложение Б. Дополнительные графики	4
Приложение В. Описание классов и функций библиотеки <code>time_series_prediction</code>	9
1. Модуль <code>predictor</code>	9
2. Модуль <code>wishart</code>	13
3. Модуль <code>experiment</code>	14
4. Модуль <code>graph</code>	15
5. Модуль <code>non_pred_model</code>	16

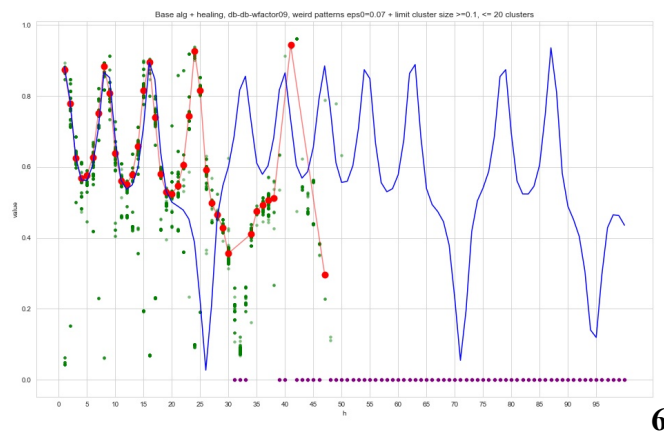
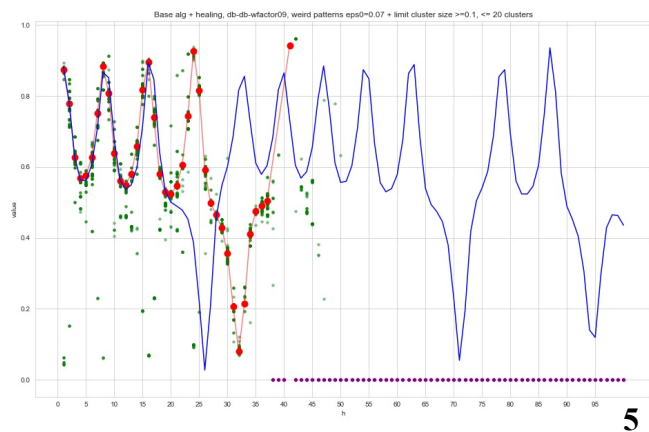
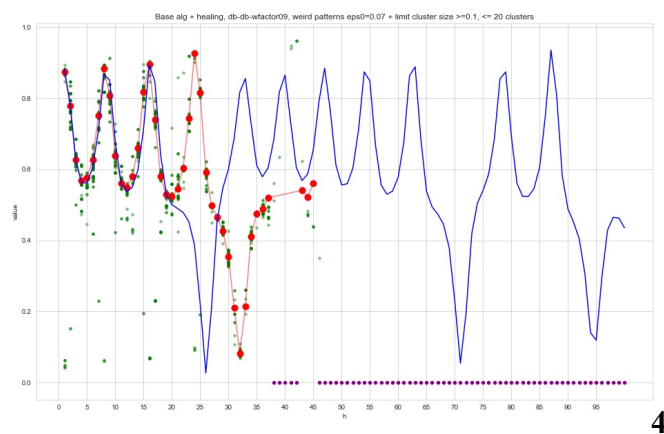
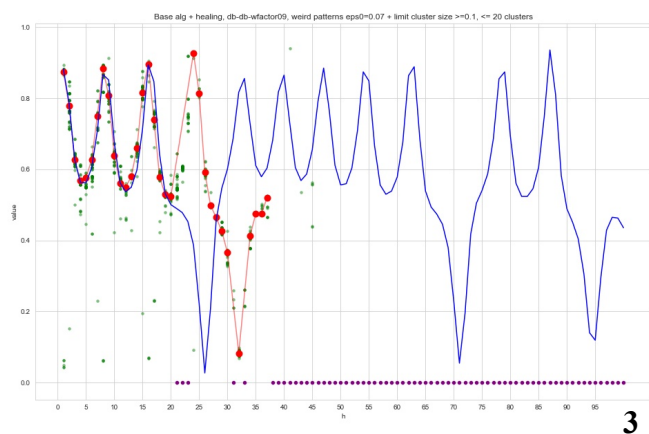
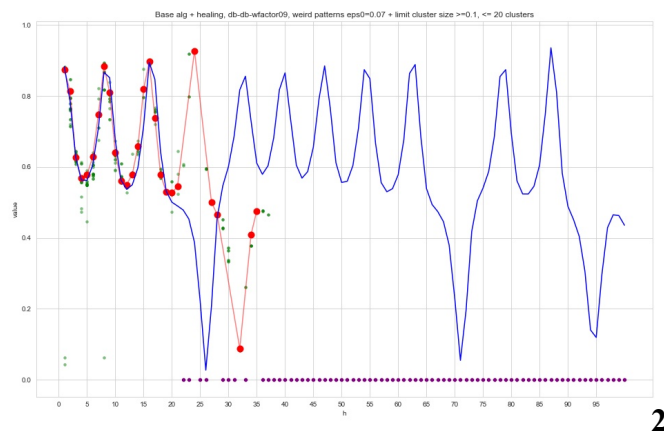
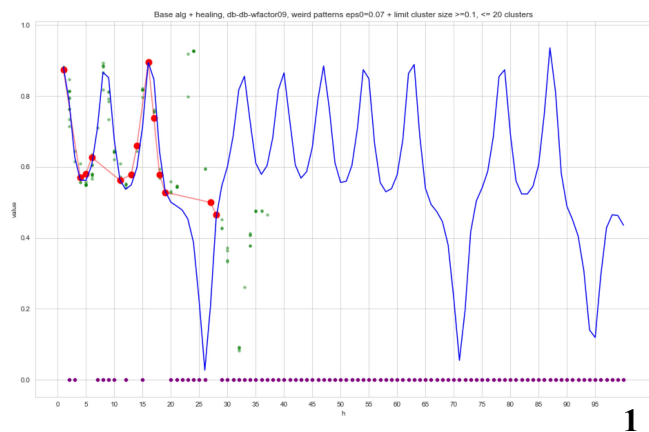
Приложение А. Условные обозначения алгоритмов

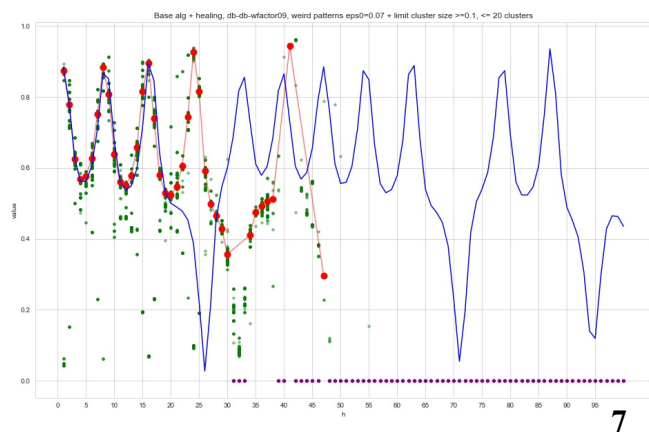
Таблица 1. Условные обозначения алгоритмов, в последнем столбце указан номер страницы подробного описания в основном отчете

Обозначение	Описание	Стр.
Тип алгоритма прогнозирования		
s	Поточенное прогнозирование	16
tp	Траекторное прогнозирование	16
sh	Алгоритм self-healing	20
Алгоритм кластеризации обобщенных z-векторов		
cl_db	Мотивы — это центры кластеров, алгоритм кластеризации — DBSCAN с параметрами $\epsilon = 0.01$, $\text{min_samples} = 5$, если не указано иное.	15
cl_wi	Мотивы — это центры кластеров, алгоритм кластеризации — Wishart с параметрами $\text{significance} = 0.01$, $\text{min_samples} = 5$, если не указано иное.	15
Алгоритмы определения непрогнозируемых точек		
fp	Вынужденное прогнозирование (forced prediction)	17
ls	Большой разброс (large spread)	17
rg	Быстрый рост разброса (rapid growth)	18
rd	Быстрый рост количества кластеров DBSCAN (rapid growth DBSCAN)	18
rw	Быстрый рост количества кластеров Wishart (rapid growth Wishart)	18
lcs	Ограничение на размер максимального кластера и на количество кластеров. Например, lcs_0.1_2 — доля точек в максимальном кластере должна быть хотя бы 0.1, всего кластеров не более 2 (без учета выбросов)	18
bl	Большой скачок (big leap)	23
blbi	Большой скачок между итерациями (big leap between iterations)	23
wp	Странные паттерны (weird patterns)	24
Алгоритмы вычисления единого прогнозного значения		
db	Центр максимального кластера DBSCAN с параметрами $\epsilon = 0.01$, $\text{min_samples} = 5$, если не указано иное.	19
wi	Центр максимального кластера Wishart с параметрами $\text{significance} = 0.01$, $\text{min_samples} = 5$, если не указано иное.	19
wa	Взвешенное среднее	22
dc	Двойная кластеризация	22
factor	Центр кластера DBSCAN, максимального с учетом весов по итерации	22-23
dist	Центр кластера DBSCAN, максимального с учетом весов по расстоянию до мотива	22-23
pl	Центр кластера DBSCAN, максимального с учетом весов по длине паттерна	22-23

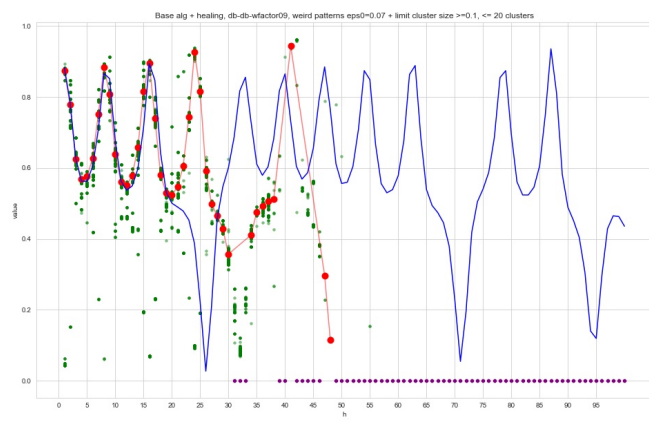
Приложение Б. Дополнительные графики

Рисунки 1-13. Алгоритм self-healing. Кластеризация мотивов — *db*, 20%. Алгоритм вычисления единого прогнозного значения — *db*, $db + factor = 0.9$ для self-healing. Алгоритм определения непрогнозируемых точек — *lcs_0.1_20*, *lcs_0.1_20 + wr* для self-healing. Синим обозначен настоящий временной ряд, зеленые точки — возможные прогнозные значения, красные точки и линии — единые прогнозные значения, фиолетовые точки на прямой $y = 0$ — непрогнозируемые точки. Порядок итераций: слева направо, сверху вниз.

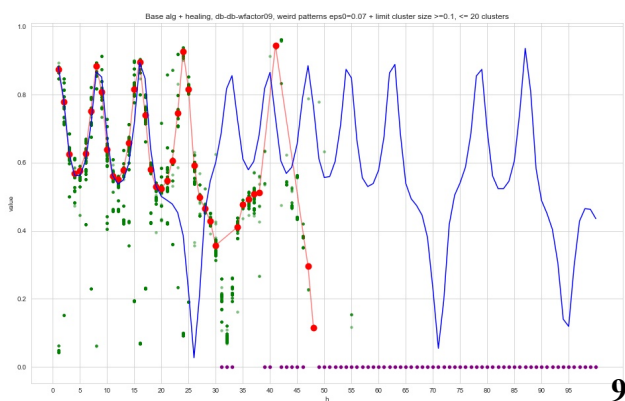




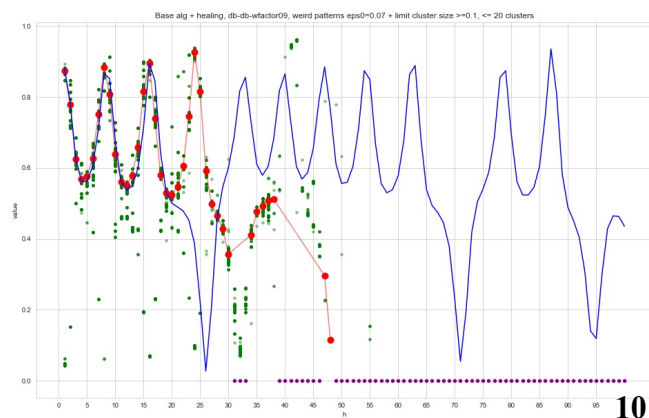
7



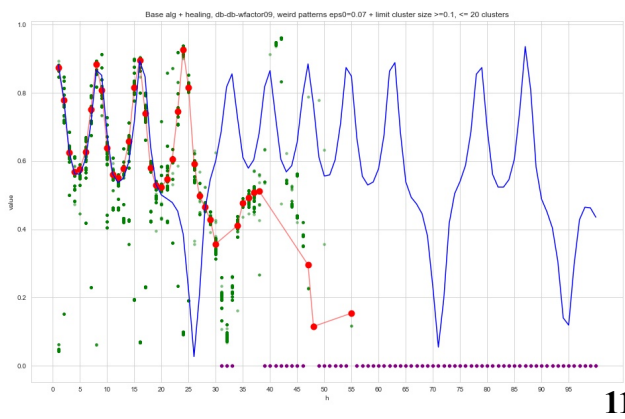
8



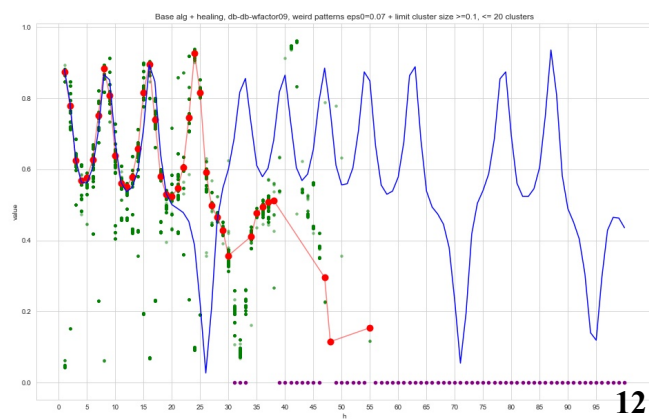
9



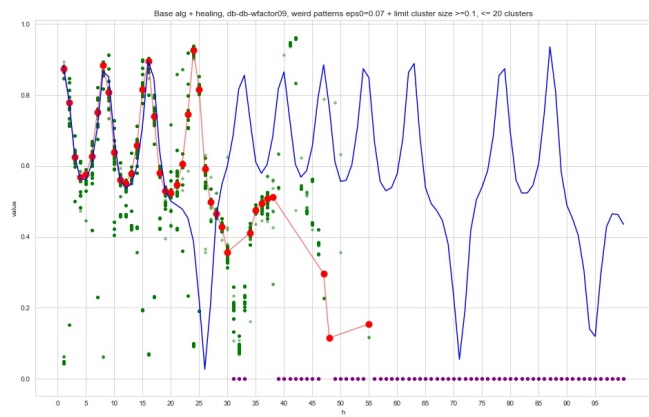
10



11



12



13

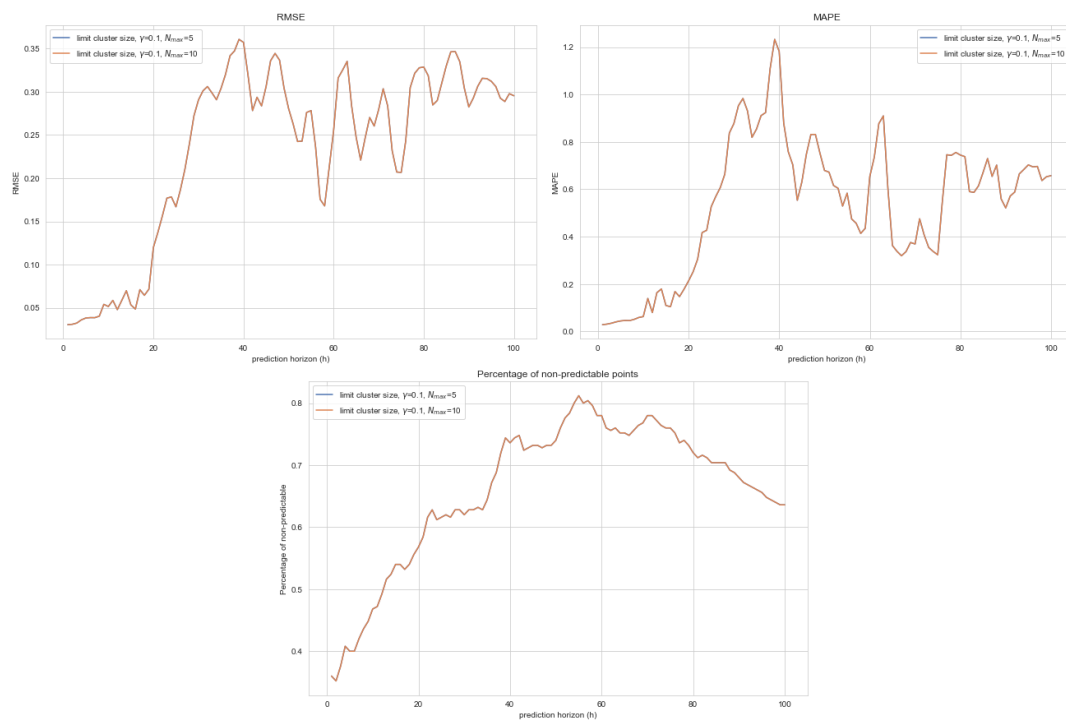


Рисунок 14. Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования. Горизонт прогнозирования $h=100$. Тестовая выборка 250. Кластеризация мотивов — db , 20%. Алгоритм вычисления единого прогнозного значения — db . Алгоритмы определения непрогнозируемых точек: $lcs_0.1_5$ (синий), $lcs_0.1_10$ (оранжевый).

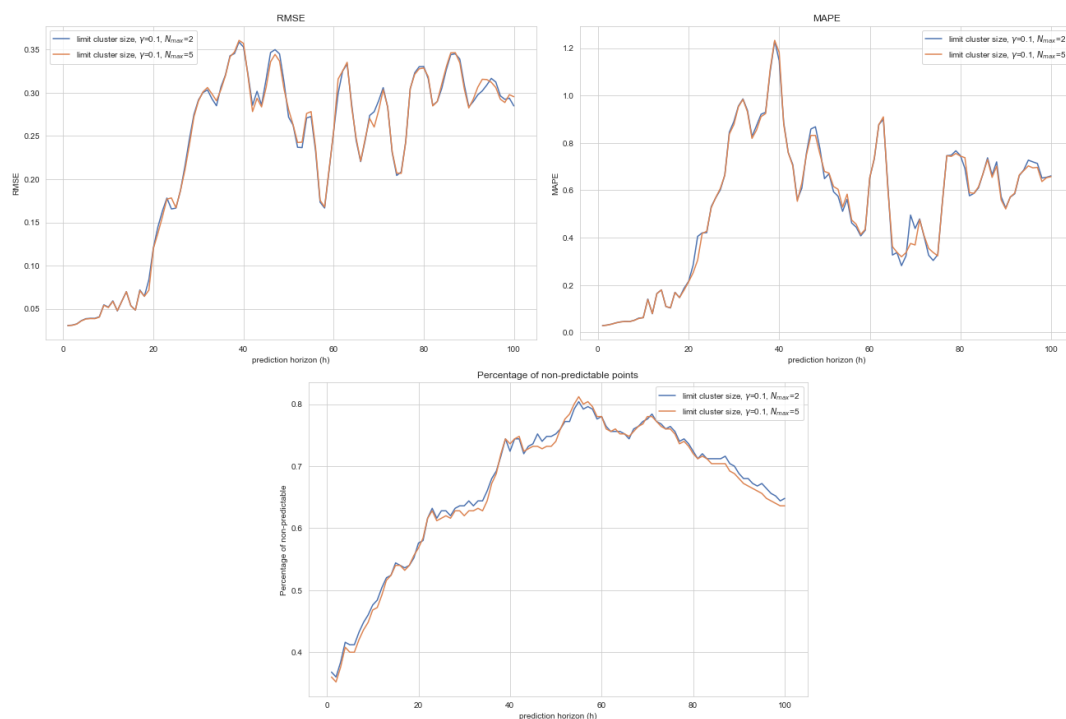


Рисунок 15. Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования. Горизонт прогнозирования $h=100$. Тестовая выборка 250. Кластеризация мотивов — db , 20%. Алгоритм вычисления единого прогнозного значения — db . Алгоритмы определения непрогнозируемых точек: $lcs_0.1_2$ (синий), $lcs_0.1_5$ (оранжевый).

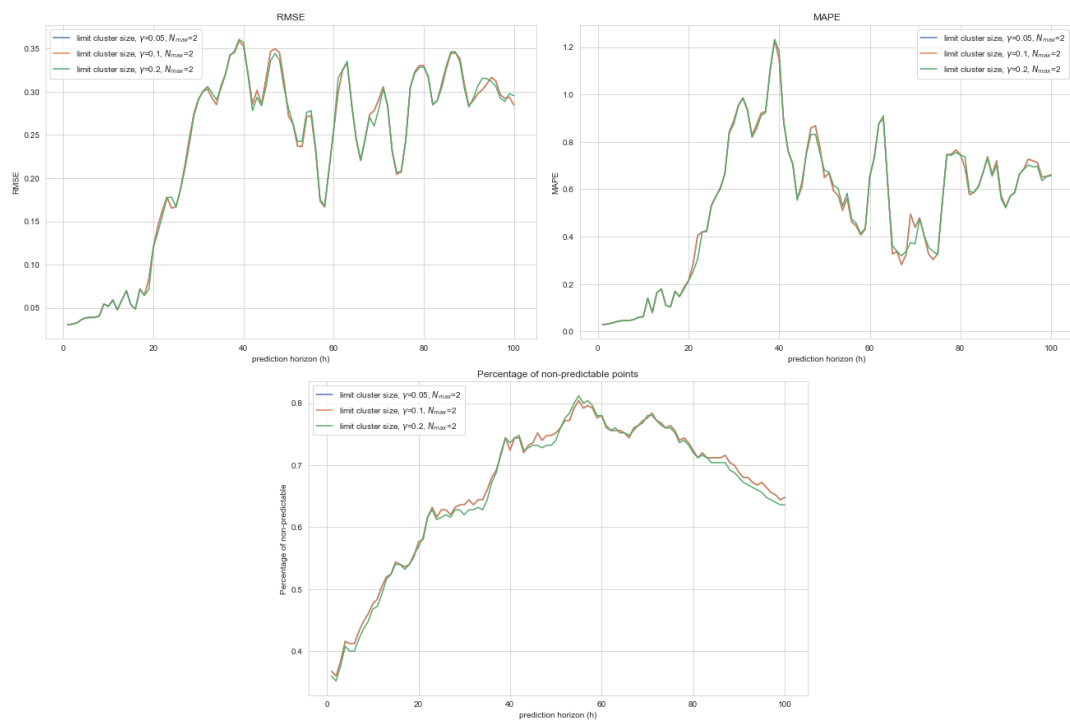


Рисунок 16. Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования. Горизонт прогнозирования $h=100$. Тестовая выборка 250. Кластеризация мотивов — db, 20%. Алгоритм вычисления единого прогнозного значения — db. Алгоритмы определения непрогнозируемых точек: lcs_0.05_2 (синий), lcs_0.1_2 (оранжевый), lcs_0.2_2 (зеленый).

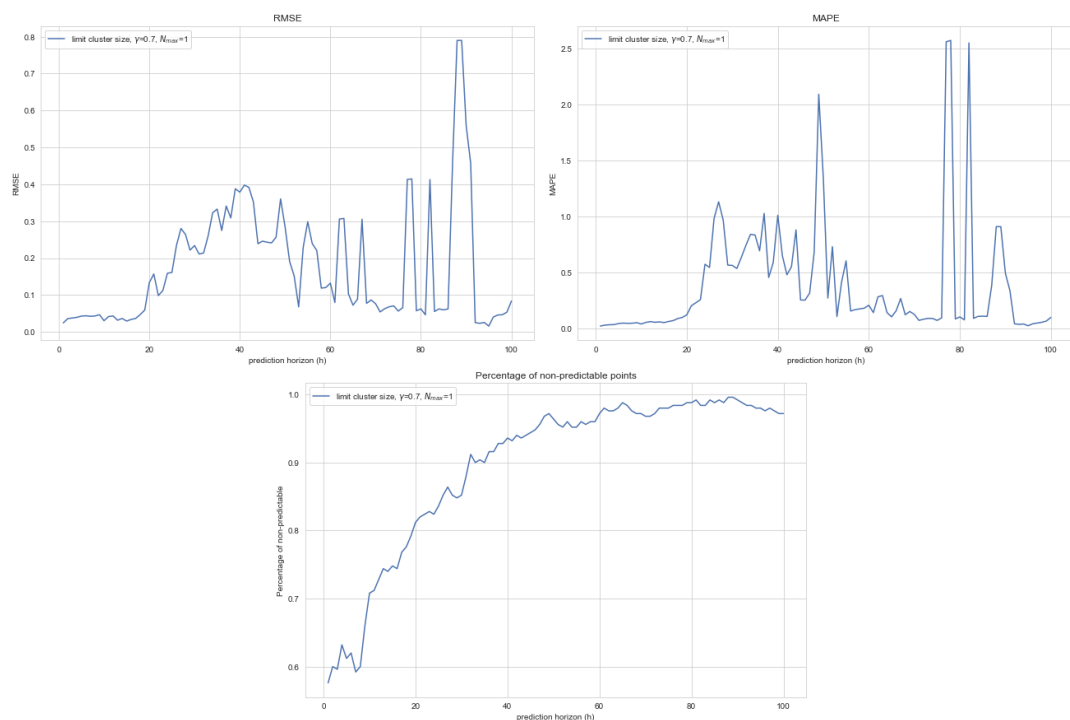
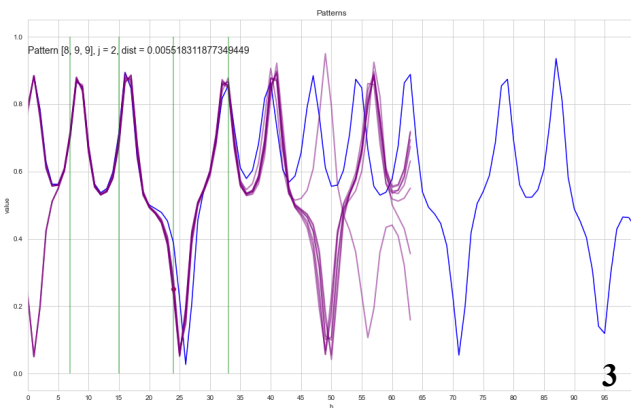
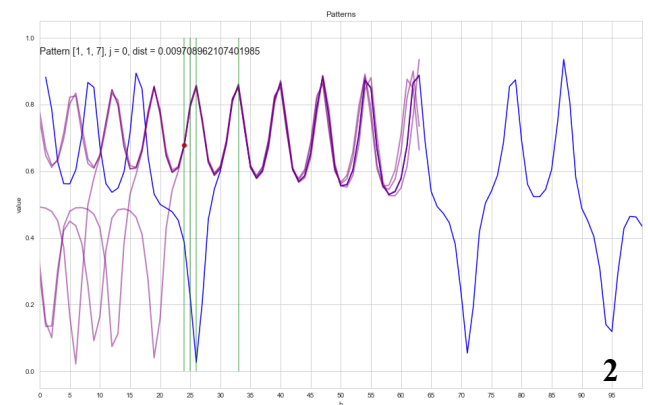
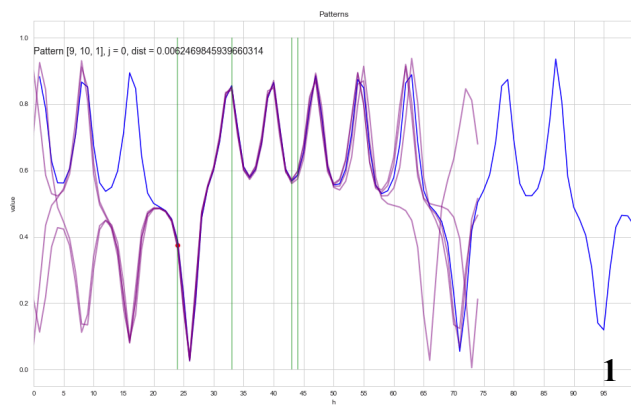


Рисунок 17. Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования. Горизонт прогнозирования $h=100$. Тестовая выборка 250. Кластеризация мотивов — db, 20%. Алгоритм вычисления единого прогнозного значения — db. Алгоритмы определения непрогнозируемых точек: lcs_0.7_1 (синий).



Рисунки 18-20. Иллюстрация к исследованию конкретных отрезков настоящего временного ряда, которые составляют кластеры, центры которых близки к наблюдаемому ряду, для точки 24. На рисунках представлено по одному мотиву из каждого кластера множества возможных прогнозных значений, не являющегося выбросами. В исследовании были проанализированы все мотивы. Синим обозначен настоящий временной ряд, зеленые вертикальные линии — паттерн, по которому взят близкий мотив, фиолетовым — отрезки ряда, отвечающие кластеру обобщенных z -векторов, красным — центр кластера множества возможных прогнозных значений.

Приложение В. Описание классов и функций библиотеки *time_series_prediction*

В ходе работы была использована среда Jupyter Notebook, а также реализована библиотека *time_series_prediction*. Она содержит код алгоритма прогнозирования, в том числе алгоритма self-healing, код вычислительных экспериментов и код для постарения графиков, некоторые из которых представлены в отчете.

1. Модуль *predictor*

Модуль *predictor* содержит класс *TimeSeriesPredictor*, который реализует базовый алгоритм прогнозирования хаотических временных рядов с помощью кластеризации и алгоритм self-healing.

Таблица 2. Атрибуты и методы класса *TimeSeriesPredictor*

Класс TimeSeriesPredictor			
Атрибуты			
Имя	Тип	Описание	
clustered_motifs	list of size (n_patterns)	Список кортежей кластеризованных мотивов: [(pattern, clusters)] Кластеры (clusters) имеют формат: list of length (n_patterns) of np.array of shape (n_clusters, pattern_length + 1)	
non_pred_model	NonPredModel	Модель для определение непрогнозируемых точек	
k_max	int, default=10	Максимальное расстояние внутри паттерна	
pattern_length	int, default=3	Длина паттерна	
Методы			
Имя	Возвращаемое значение	Параметры	Описание
__init__	—	<ul style="list-style-type: none">clustered_motifs : list of size (n_patterns) — список кортежей кластеризованных мотивовnon_pred_model : NonPredModel — модель для определение непрогнозируемых точекk_max : int — максимальное расстояние внутри паттернаpattern_length : int — длина паттерна	Конструктор

cluster_motifs	<ul style="list-style-type: none"> clustered_motifs : list of size (n_patterns) — список кортежей кластеризованных мотивов 	<ul style="list-style-type: none"> Y1 : list or 1d np.ndarray — обучающая часть временного ряда beta : float from 0 to 1, default=0.1 — процент использованных паттернов mc_method : {'wi', 'db'}, default='db' — Метод кластеризации мотивов: 'wi' - Wishart, 'db' - DBSCAN k_max : int, default=10 — максимальное расстояние внутри паттерна pattern_length : int, default=3 — длина паттерна <p>**kwargs : dict</p> <ul style="list-style-type: none"> eps : float from 0 to 1, default=0.01 — максимальное расстояние в одном кластере для кластеризации DBSCAN и Wishart min_samples : int > 1 or float from 0 to 1, default=5 — минимальное число объектов в одном кластере для кластеризации DBSCAN и Wishart 	Кластеризация мотивов, результат также сохраняется в атрибуты класса
set_motifs	—	<ul style="list-style-type: none"> clustered_motifs : list of size (n_patterns) — список кортежей кластеризованных мотивов k_max : int — максимальное расстояние внутри паттерна pattern_length : int — длина паттерна 	Установка кластеризованных мотивов в атрибуты класса
set_non_pred_model	—	<ul style="list-style-type: none"> non_pred_model : NonPredModel — модель для определения непрогнозируемых точек 	Установка модели определения непрогнозируемых точек в атрибуты класса
predict	<ul style="list-style-type: none"> Y_pred : np.ndarray — прогнозные значения, 'N' для непрогнозируемых точек possible_predictions_list : list — список множеств возможных прогнозных значений trajectories : list — прогнозные траектории для alg_type='tp' 	<ul style="list-style-type: none"> Y_preceding : list or 1D np.array — сегмент временного ряда, который предшествует прогнозируемому сегменту h : int — горизонт прогнозирования up_method : str from {'a', 'wi', 'db', 'op'} — алгоритм вычисления единого прогнозного значения alg_type : str from {'s', 'tp'} — тип алгоритма: 's' - поточное, 'tp' - траекторное match_threshold : float, default=0.01 — порог для близкого мотива <p>**kwargs</p> <ul style="list-style-type: none"> n_trajectories : int, default=20 — кол-во траекторий для alg_type='tp' 	Базовый алгоритм прогнозирования на h шагов вперед

predict_one_step	<ul style="list-style-type: none"> possible_predictions : np.ndarray — множество возможных прогнозных значений distances : np.ndarray — расстояния до мотива, используются в некоторых алгоритмах вычисления единых прогнозных значений 	<ul style="list-style-type: none"> Y_preceding : list or 1D np.array — сегмент временного ряда, который предшествует прогнозируемому сегменту Y_pred : list or 1D np.ndarray — спрогнозированные на предыдущих шагах значения match_threshold : float, default=0.01 — порог для близкого мотива 	Один шаг базового алгоритма прогнозирования
unified_prediction	<ul style="list-style-type: none"> avg : float or 'N' — единое прогнозное значение, 'N' для непрогнозируемых точек 	<ul style="list-style-type: none"> possible_predictions : np.ndarray — множество возможных прогнозных значений up_method : str from {'a', 'wi', 'db', 'op'} — алгоритм вычисления единого прогнозного значения (см. Приложение А) 	Алгоритм вычисления единого прогнозного значения
unified_prediction_weighted	<ul style="list-style-type: none"> avg : float or 'N' — единое прогнозное значение, 'N' для непрогнозируемых точек 	<ul style="list-style-type: none"> possible_predictions : np.ndarray — множество возможных прогнозных значений sep_indices : list — список индексов, которые разделяют итерации в списке возможных прогнозных значений up_method : str from {'a', 'wi', 'db', 'op'} — алгоритм вычисления единого прогнозного значения (см. Приложение А) weight_method : str from {'double_clustering', 'weighred_average', 'factor', 'pattern_length', 'pattern_length_dist', 'dist', 'dist_factor'} — метод вычисления единого прогнозного значения с использованием весов (см. Приложение А) 	Алгоритм вычисления единого прогнозного значения с использованием весов

self_healing_one_iteration	<ul style="list-style-type: none"> • new_up : np.ndarray — список новых прогнозных значений • possible_predictions : list of np.ndarrays — список множеств прогнозных значений • sep_indices : list — список индексов, которые разделяют итерации в списке возможных прогнозных значений 	<ul style="list-style-type: none"> • Y_preceding : list or 1D np.array — сегмент временного ряда, который предшествует прогнозируемому сегменту • sep_indices : list — список индексов, которые разделяют итерации в списке возможных прогнозных значений • unified_predictions : list or 1D np.ndarray of length h — список единых прогнозных значений с предыдущей итерации • possible_predictions : list of h np.ndarrays — список множеств возможных прогнозных значений с предыдущей итераций • up_method : str from {'a', 'wi', 'db', 'op'} — алгоритм вычисления единого прогнозного значения (см. Приложение А) • fixed_points_idx : list — список точек, которые не изменяют статус прогнозируемых • weight_method : str from {'double_clustering', 'weighred_average', 'factor', 'pattern_length', 'pattern_length_dist', 'dist', 'dist_factor'} — метод вычисления единого прогнозного значения с использованием весов (см. Приложение А) 	Одна итерация алгоритма self-healing
self_healing	<ul style="list-style-type: none"> • unified_predictions : list or 1D np.ndarray of length h — список единых прогнозных значений • possible_predictions : list of np.ndarrays — список множеств возможных прогнозных значений 	<ul style="list-style-type: none"> • Y_preceding : list or 1D np.array — сегмент временного ряда, который предшествует прогнозируемому сегменту • h : int — горизонт прогнозирования • unified_predictions : list or 1D np.ndarray of length h — список единых прогнозных значений из базового алгоритма • possible_predictions : list of h np.ndarrays — список множеств возможных прогнозных значений из базового алгоритма • healing_up_method : str from {'a', 'wi', 'db', 'op'} — алгоритм вычисления единого прогнозного значения (см. Приложение А) • fixed_points_idx : list — список точек, которые не изменяют статус прогнозируемых • weight_method : str from {'double_clustering', 'weighred_average', 'factor', 'pattern_length', 'pattern_length_dist', 'dist', 'dist_factor'} — метод вычисления единого прогнозного значения с использованием весов (см. Приложение А) 	Алгоритм self-healing

2. Модуль *wishart*

Класс *Wishart* содержит реализации алгоритма кластеризации Wishart

Таблица 3. Атрибуты и методы класса *Wishart*

Класс TimeSeriesPredictor			
Атрибуты			
Имя	Тип	Описание	
wishart_neighbors	int	Минимальное число объектов в одном кластере	
significance_level	int	Максимальный уровень значимости кластера	
Методы			
Имя	Возвращаемое значение	Параметры	Описание
__init__	—	<ul style="list-style-type: none">wishart_neighbors : int — минимальное число объектов в одном кластереsignificance_level : int — максимальный уровень значимости кластера	Конструктор
fit	<ul style="list-style-type: none">result : list — список меток кластеров	<ul style="list-style-type: none">X : list — список объектов для кластеризации	Алгоритм кластеризации Wishart

3. Модуль *experiment*

Таблица 4. Функции модуля *experiment*

Модуль <i>experiment.py</i>			
Функции			
Имя	Возвращаемое значение	Параметры	Описание
<code>experiment_no_pm</code>	—	<ul style="list-style-type: none"> • <code>Y2</code>: list or 1D <code>np.ndarray</code> — тестовая часть ряда • <code>h_max</code>: int — максимальный горизонт прогнозирования • <code>n_iterations</code>: int — размер тестовой выборки эксперимента • <code>iterations_range</code>: tuple — ограничения тестовой выборки • <code>motif_clustering_params</code>: dict — параметры кластеризации мотивов • <code>prediction_params</code>: dict — параметры базового алгоритма • <code>healing_params</code>: dict — параметры алгоритма self-healing • <code>non_pred_model_prediction</code> : <code>NonPredModel</code> — модель определения непрогнозируемых точек для базового алгоритма • <code>non_pred_model_healing</code> : <code>NonPredModel</code> — модель определения непрогнозируемых точек для алгоритма self-healing 	Эксперимент оценки качества прогнозирования, не использует матрицу прогнозов, запись результатов в файл. Есть возможность параллельного выполнения задач.
<code>experiment</code>	—	Те же, что и в функции <i>experiment_no_pm</i>	Эксперимент оценки качества прогнозирования, использует матрицу прогнозов
<code>thrown_points_experiment</code>	—	<ul style="list-style-type: none"> • <code>Y1</code>: list or 1D <code>np.ndarray</code> — обучающая часть ряда • <code>Y2</code>: list or 1D <code>np.ndarray</code> — тестовая часть ряда • <code>h_max</code>: int — максимальный горизонт прогнозирования • <code>healing_params</code>: dict — параметры алгоритма self-healing • <code>logs_filepath</code> : str — путь к файлу записи результатов • <code>step</code> : int — шаг 	Эксперимент для ряда с выкинутыми точками

4. Модуль *graph*

В таблице перечислены основные функции для построения графиков для анализа результатов исследования. Остальные функции модуля являются служебными.

Таблица 5. Функции модуля *graph.py*

Модуль <i>graph.py</i>			
Функции			
Имя	Возвращаемое значение	Параметры	Описание
<code>plot_stats</code>	—	<ul style="list-style-type: none"> <code>stats_ds</code>: <code>pd.DataFrame</code> columns: <code>points_left</code>, <code>n_iterations</code>, <code>non_predictable</code>, <code>rmse</code>, <code>mape</code> — таблица статистики по исследованию либо <code>stats_filename</code>: <code>str</code> — путь к файлу с таблицей 	Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования из готового файла со статистикой
<code>plot_unified_and_possible_preds</code>	—	<ul style="list-style-type: none"> <code>up</code>: <code>list</code> or <code>1D np.ndarray</code> — список единых прогнозных значений <code>pp</code>: <code>list of lists</code> or <code>list of np.ndarrays</code> — список множеств возможных прогнозных значений <code>Y2</code>: <code>list</code> or <code>1D np.ndarray</code> — настоящие значения временного ряда 	График единых прогнозных значений и множеств возможных прогнозных значений, а также настоящие значения временного ряда
<code>plot_healing_animation</code>	—	<ul style="list-style-type: none"> <code>working_path</code> : <code>str</code> — путь к файлу с логами по алгоритму self-healing <code>Y2</code>: <code>list</code> or <code>1D np.ndarray</code> — настоящие значения временного ряда 	Анимация алгоритма self-healing
<code>plot_experiment_results</code>	—	<ul style="list-style-type: none"> <code>working_directory</code> : <code>str</code> — путь к директории с матрицами прогнозов <code>exp_short_names</code> : <code>list</code> — список коротких наименований алгоритмов, матрицы прогнозов которых используются в исследовании <code>h_max</code> : <code>int</code> — максимальный горизонт прогнозирования <code>n_iterations</code> : <code>int</code> — размер тестовой выборки 	Графики зависимости RMSE, MAPE и количества непрогнозируемых точек от горизонта прогнозирования из матриц прогнозов
<code>plot_thrown_points_exp_results</code>	—	<ul style="list-style-type: none"> <code>working_directory</code> : <code>str</code> — путь к директории с логами алгоритма self-healing <code>exp_short_names</code> : <code>list</code> — список коротких наименований алгоритмов, логи которых используются в исследовании <code>h_max</code> : <code>int</code> — максимальный горизонт прогнозирования <code>n_iterations</code> : <code>int</code> — размер тестовой выборки 	Графики зависимости количества итераций, количества непрогнозируемых точек, RMSE и MAPE от количества выкинутых точек

5. Модуль *non_pred_model*

Таблица 6. Методы корневого класса *NonPredModel*, наследниками которого являются все классы моделей определения непрогнозируемых точек

Класс <i>NonPredModel</i>			
Методы			
Имя	Возвращаемое значение	Параметры	Описание
<code>is_predictable</code>	<code>is_predictable</code> : boolean — прогнозируемая ли точка	<code>possible_predictions</code> : list — множество возможных прогнозных значений <code>**kwargs</code> — другие параметры, необходимые в классах-наследниках	Функция, которая определяет прогнозируемая ли точка по множеству возможных прогнозных значений
<code>reset</code>	—	—	Сброс атрибутов класса
<code>is_predictable_by_up</code>	<code>is_predictable</code> : boolean — прогнозируемая ли точка	<code>unified_predictions</code> : list — единые прогнозных значения <code>**kwargs</code> — другие параметры, необходимые в классах-наследниках	Функция, которая определяет прогнозируемая ли точка по списку единых прогнозных значений
<code>is_predictable_by_up_log</code>	<code>is_pred</code> : list of boolean — прогнозируемы ли точки	<code>up_log</code> : list — лог единых прогнозных значений <code>**kwargs</code> — другие параметры, необходимые в классах-наследниках	Функция, которая определяет прогнозируемая ли точка по логу единых прогнозных значений

Таблица 7. Классы моделей определения непрогнозируемых точек

Модуль <i>non_pred_model.py</i>	
Классы	
Имя	Алгоритм, который реализуется в классе
<code>ForcedPredictionNPM</code>	Принудительное прогнозирование (fp)
<code>LargeSpreadNPM</code>	Большой разброс (ls)
<code>RapidGrowthNPM</code>	Быстрый рост разброса (rg)
<code>RapidGrowthDBSCANNPM</code>	Быстрый рост разброса кластеров DBSCAN (rd)
<code>RapidGrowthWishartNPM</code>	Быстрый рост разброса кластеров Wishart (rw)
<code>LimitClusterSizeNPM</code>	Ограничение на размер максимального кластера и на количество кластеров (lcs)
<code>BigLeapNPM</code>	Большой скачок (big leap)
<code>BigLeapBtwIterationsNPM</code>	Большой скачок между итерациями (big leap between iterations)
<code>WeirdPatternsNPM</code>	Странные паттерны (weird patterns)