

Анализ матриц корреспонденции метро

М.Н. Некраплённая, Д.Е. Намиот

Аннотация— Настоящая статья посвящена анализу транспортных потоков на основе матриц корреспонденции. Такие матрицы описывают количество перемещений между двумя точками за определенный интервал времени. С практической точки зрения, рассматриваются данные, относящиеся к Московскому метрополитену. Соответственно, матрица корреспонденции описывает перемещения между станциями. Теоретически, такие данные описывают все характеристики пассажиропотока. На практике, это зависит, естественно, от выбранной модели обработки данных. Часто, такого рода матрицы используются лишь для простой статистики, типа количества перевезенных пассажиров за какое-то время. При этом теряется интересная для цифровой урбанистики пространственно-временная информация. Например, как были распределены перевозки по времени, насколько стабильны такие распределения и т.д. В работе приводится подробный обзор существующих подходов к анализу данных матриц корреспонденции. В качестве практической задачи рассматривается краткосрочный прогноз пассажиропотока. Отмечается, что краткосрочный прогноз дорожного движения является сложной задачей, которая была предметом многих исследовательских работ в последние несколько десятилетий. Большая часть работ была исторически посвящена анализу транспортных потоков исключительно автомобильного транспорта. Изучение же железнодорожного и, в частности, подземного транспорта с его спецификой, долгое время оставалось без внимания. Соответствующие исследования стали проводиться лишь в последнее время.

Ключевые слова— матрица корреспонденции, транспортный поток, метро.

I. ВВЕДЕНИЕ

В настоящей работе рассматриваются вопросы анализа транспортных потоков. В основание статьи была положена выпускная квалификационная работа, выполненная на факультете ВМК МГУ имени М.В. Ломоносова. Естественно, что транспортная составляющая – это одна из важнейших компонент Умного города. Поэтому исследованиям транспорта в городах уделяется достаточно большое внимание, что отражено во многих наших работах [54, 55].

В данной статье речь идет об анализе матриц корреспонденции (в англоязычной литературе – origin-destination matrix) для метрополитена города Москвы. Телекоммуникационные операторы в Москве работают и в метро. Это позволяет им определять, когда

пользователь мобильной связи (он же – пассажир) переключился на обслуживание базовой станцией, находящейся в метро, а также когда произошло обратное событие – переключение с подземной станции на наземную. Это позволяет мобильным операторам знать начальную и конечную точку маршрута своего абонента в метро. Для сохранения приватности такого рода данные могут быть агрегированы по времени – все поездки между двумя станциями за указанный интервал времени. В таком случае выделить индивидуальный маршрут уже не получится. Такая агрегированная информация и есть матрица корреспонденции. Она показывает количество человек, которые в заданный временной интервал (15 минут, час, день) поехали от станции *A* до станции *B*.

Сейчас, именно для Московского метрополитена такие данные используются для простой статистики (сколько человек пользуется той или иной станцией за месяц, какова суммарная загрузка линий метро и так далее). Естественно, что при этом теряется (просто не используется) вся пространственно-временная информация, характеризующая поездки. А это очень важные характеристики. Транспортная система Москвы постоянно развивается. Это, в частности, выражается в том, что метро (как основной вид транспорта) интегрируется с другими видами транспорта (например, городской железной дорогой). Происходит перераспределение транспортных потоков и пассажиропотоков. Другая причина изменения пассажиропотока - высотное жилищное строительство. Все это вызывает необходимость разбираться в том, как организованы пассажиропотоки в метро, как они распределены по времени, что может произойти в случае роста потока, есть ли какие-то резервы по пропускной способности и так далее. Простой эмпирической оценки того, что обычно в день или в месяц этим транспортом пользуется определённое количество пассажиров, становится недостаточно. Данная работа посвящена именно анализу пассажиропотоков метро. Объектом исследования являются пассажиропотоки московского метрополитена, а предметом исследования выступает моделирование транспортных потоков для выполнения наиболее точных краткосрочных прогнозов состояния транспортной системы.

Существует несколько общепринятых подходов к моделированию и краткосрочному прогнозированию на основе матрицы корреспонденций с использованием методов машинного обучения, однако библиографический анализ показал, что в литературе нет сведений о построении таких моделей именно для метрополитена города Москвы.

Статья получена 30 мая 2019.

М.Н. Некраплённая – МГУ имени М.В. Ломоносова (email: maria240398@mail.ru)

Д.Е. Намиот - МГУ имени М.В. Ломоносова; РУТ (МИИТ) (e-mail: dnamiot@gmail.com).

В мире большое количество учёных уделяет внимание изучению данного вопроса, начиная с 70-х годов прошлого века. В исследованиях применялись различные подходы: сначала аналитическое моделирование (методы классической статистики, анализ временных рядов), а затем моделирование на основе данных (нейронный и эволюционный подходы).

Оставшаяся часть статьи структурирована следующим образом. В разделе II рассматриваются основные определения. В разделе III приведен библиографический анализ. Раздел IV посвящен анализу данных.

II. ОБЩИЕ ОПРЕДЕЛЕНИЯ

В настоящее время, с увеличением населения мегаполисов, очевидно, растёт нагрузка и на их транспортные системы. В любой из них принципиально важной задачей является наблюдение и регулирование транспортных потоков, т.е. перемещений внутри этой системы.

Математически перемещения пассажиров внутри транспортной системы описываются т.н. *матрицей корреспонденций*.

Рассмотрим транспортную сеть как планарный граф $G = (V, E)$, где V – множество вершин, E – множество дуг сети. С каждой вершиной может ассоциироваться некоторый транспортный узел (станция) как место отправления (исток) или прибытия (сток) пассажиров. Пусть $O \subset V$, $D \subset V$ – множества вершин графа, которые можно назвать соответственно, истоками и стоками сети. Матрица корреспонденций в общем виде $p(t) = \{p_{ij}(t), i \in O, j \in D, t \in T\}$

определяет распределение пассажиропотока в сети и может характеризоваться, например, количеством пассажиров, переместившихся из района с номером i в район с номером j за единицу времени t . Матрица корреспонденций рассматривается как укрупненная транспортная модель, описывающая некоторую топологию транспортной сети города или агломерации. Эта матрица служит основой для построения детальной модели распределения транспортных потоков [47, 48]

Предоставленный для исследования набор данных представляет из себя таблицу, содержащую информацию о передвижениях пассажиров Московского метрополитена в период с 1 по 28 февраля 2018 г с получасовыми интервалами. Каждая запись состоит таких полей: дата и время; номер станции отправления; номер станции прибытия; количество поездок, начатых в заданное время по данному маршруту; количество поездок, завершённых в заданное время. Эта форма представления данных о пассажиропотоках не совсем соответствует общему виду матрицы корреспонденций, однако позволяет не упускать из внимания поездки, происходящие в течение более чем одного временного интервала. В дальнейшем для расчётов будем использовать суммарное значение двух последних полей и называть его общим числом поездок. На момент сбора данных транспортная система включала в себя 213 станций. Из рассмотрения исключены поездки, для которых станции входа и выхода совпадают, а также

поездки длительностью более 4 часов. Полный объём полученных данных — 26 790 535 записей (738 432 Кб).

Задача данной работы заключается в том, чтобы по имеющимся матрицам корреспонденций

$$p(t-n), p(t-n+1), \dots, p(t)$$

построить наиболее точный краткосрочный прогноз

$$\hat{p}(t+1)$$

так, чтобы ошибка была минимальной:

$$\hat{p}(t+1) - p(t+1) \rightarrow \min.$$

III. ОБЗОР СУЩЕСТВУЮЩИХ РАБОТ

Краткосрочный прогноз дорожного движения является сложной нелинейной задачей, которая была предметом многих исследовательских работ в последние несколько десятилетий. Большая часть упомянутых далее методов, хоть и подходят, формально, для любых корреспонденций, но реально были применены для анализа транспортных потоков только автомобильного транспорта. Изучение же железнодорожного и, в частности, подземного транспорта с его спецификой, долгое время оставалось без внимания. Соответствующие исследования стали проводиться лишь в последние 15 лет. Они будут рассмотрены в конце раздела.

Термин «краткосрочный» обычно подразумевает, что представляющие интерес переменные прогнозируются на период до 1 часа вперед, хотя точное определение в значительной степени отличается у разных авторов.

Подходы к анализу матрицы корреспонденций и к прогнозированию трафика отличаются друг от друга во многих аспектах. К ним относятся, например, масштаб прогноза (фиксированное местоположение, маршрут или целая сеть) или тип наблюдаемой дороги или сети (целые автострады или их отдельные участки; перекрёстки, либо другие развязки; контролируемые или нет). Возможно, наиболее важным фактором в оценке и сравнении подходов прогнозирования трафика является тип данных для ввода и вывода, то есть, какие величины, характеризующие трафик, прогнозируются, а какие используются для этого прогноза. В любом случае, проблема краткосрочного прогнозирования трафика сводится к решению следующей задачи регрессии:

$$y_k = G(x_{s|s < k}, e_{s|s < k}, \theta_k) + e_k \quad (1)$$

Где G – выбранная модель прогнозирования; y_k – (вектор из) выходная переменная в момент времени k ; x_k – (вектор из) входная переменная прогнозирования, которая также может включать в себя уже зафиксированные в прошлом y_{k-n} ; θ_k – (вектор из) настраиваемые параметры модели; e_k – шум.

В литературе имеется множество обзоров или сравнительных исследований, например, [1,2]. Авторы данной работы используют классификацию, предложенную в [3], в ней выделяется три больших категории: наивные, параметрические и непараметрические модели.

В параметрических подходах производят анализ, используя методы классической математики (например, исследуют функции времени в пути, или модели массового обслуживания) или имитационные модели движения (микро- или макроскопические) с параметрами (например, пропускная способность дороги, вероятность смены полосы движения в ту или иную сторону). В этом случае модель G в уравнении (1) является фиксированной. Некоторые из этих параметров могут свободно регулироваться, тогда как другие параметры настраиваются в пределах физических границ (например, пропускная способность и свободные скорости должны соответствовать тому, что физически возможно). Хотя эти модели реализуют (правдоподобные) теоретические или физические предположения об изменении характера движения во времени и могут быть откалиброваны для воспроизведения многих дорожных ситуаций, для них необходимо долго «настраивать» внутренние переменные модели. Наиболее популярными подходами являются фильтры Калмана, KF [4]. Можно утверждать, что они могут быть классифицированы как гибридные подходы к прогнозированию трафика, которые сочетают в себе как параметрические, так и непараметрические методы.

Непараметрические подходы, по сути, охватывают все другие методы прогнозирования трафика. Возможно, наибольшее количество моделей прогнозирования трафика и времени в пути относятся к этой категории и в первую очередь, разумеется, это методы машинного обучения и анализа данных. В качестве примеров стоит упомянуть регрессию опорных векторов [7], обобщенную линейную регрессию [8], модель ARIMA [9], нелинейные временные ряды [10], модели пространства состояний и некоторые модификации KF [11,12], нейронные сети [13,14] и их разновидности.

А. Наивные методы

Наивные подходы — это те, в которых ни структура модели (G в уравнении (1)), ни параметры θ_k не зависят от данных. Термин «наивный» довольно субъективен, но его можно свободно интерпретировать как «без дополнительных предположений», кроме измеренных данных и общеизвестных физических законов (например, расстояние = скорость × время). Наивные методы широко применяются на практике из-за их малой вычислительной сложности и простоты реализации, но точность получаемых результатов обычно низкая. Тем не менее, значительное число моделей прогнозирования сочетают в себе наивные методы с более сложными (например, регрессия или кластеризация в [15]).

Самым распространённым предиктором в случае параметрических моделей является текущее время в пути. Основное предположение здесь заключается в том, что преобладающие условия трафика (скорости, плотности, очереди и т. д.) будут оставаться постоянными в течение неопределенного времени. В этом случае хорошо понятные методы восстановления времени в пути могут использоваться в качестве исходных данных для прогнозирования времени в пути по маршруту, состоящему из последовательных

участков i , через $TT_k = \sum_i L_i / u_i$. Такие вычисления обеспечивают точные прогнозы в случаях, когда условия движения являются стационарными и однородными в течение длительных периодов времени, например, в условиях свободного движения или в редком случае очереди с фиксированной длиной, которая рассеивается с постоянной скоростью. Несмотря на то, что этот метод очень быстрый и многие специалисты используют его на практике, прогностическая эффективность метода быстро ухудшается, когда условия движения переходят от свободного потока к перегруженным условиям и обратно.

Во многих случаях текущее время в пути используется просто как базовый предиктор, с которым сравниваются другие подходы [14, 15]. Тем не менее, мгновенное время в пути по-прежнему остается наиболее широко используемым методом для многих систем мониторинга движения.

Историческое среднее время в пути часто используется в качестве базового предиктора, с которым сравниваются другие методы. Для долгосрочного прогнозирования трафика историческое среднее значение во многих случаях является наилучшим доступным подходом. Во многих случаях [14,15], эти средние значения меняются в зависимости от времени суток, дня недели или каких-то других ритмов (например, сезонность). В действительности, на многих участках дороги распределение времени в пути сильно варьируется, то есть, исторический средний показатель обычно является плохим предиктором.

Комбинации мгновенного и исторического среднего [16] дают значительно лучшие результаты по сравнению с каждым из них в отдельности. Вместо простого объединения этих двух подходов, могут, наконец, использоваться методы кластеризации средних значений на кластеры со сходной динамикой трафика. В качестве примера можно привести кластеризацию Уорда [17]. Более сложные методы кластеризации (карты Кохонена) были использованы для предварительной обработки входных данных некоторых непараметрических методов краткосрочного прогнозирования [18].

В. Параметрические модели

Термин «параметрический» в данном контексте означает, что только параметры модели подбираются на основе самих данных. Структура модели при этом строится исходя из теоретических соображений. Большинство параметров при этом характеризуют реальные физические явления, такие как критическая скорость и плотность потока, пропускная способность, свободная скорость, вероятность поворота и т. д.

Самым простым способом прогнозирования времени в пути является использование аналитических формул, таких как функция BPR: $TT_k = \alpha [Q_k / C_k]^\beta$, где TT_k — это время в пути на k -ом участке; Q_k — ожидаемый спрос, C_k — пропускная способность; α и β — регулируемые параметры; или функции очередей вида $TT_k = TT_{своб} + N_k / C_k$ в которых N_k является оценкой количества транспортных средств,

стоящих в очереди в узком месте маршрута. Основная проблема здесь состоит в том, что эти модели эффективны, только если входные переменные (спрос) точно измерены, что в действительности редко имеет место. Более того, все переменные и параметры являются в высокой степени стохастическими (спрос, емкость). Обычно модели очередей неустойчивы, т.е. небольшие входные ошибки приводят к кумулятивным и очень большим ошибкам прогнозирования. Подобные методы используются на практике во многих приложениях управления движением, где прогнозы используются для коротких участков дороги и для коротких горизонтов прогнозирования. Примером является модель TUC [19].

В макроскопических имитационных моделях рассматриваются такие переменные трафика, как плотности, средние скорости и величина потока. В простейшем подходе (первого порядка) используются три уравнения: динамическое уравнение для описания изменения средней плотности транспортного потока ρ , равновесное отношение для расчета потока q (или скорости u) из плотности. В масштабе сети шаблоны выбора маршрутов значительно увеличивают сложность и количество степеней свободы. В контексте краткосрочного прогнозирования обычно рассматривается только совокупный выбор маршрута, оцененный на основе текущих измерений и исторических данных [6].

Чтобы синхронизировать внутренние переменные состояния (скорость и плотность) с фактическими данными, используют KF или методы Монте-Карло. Примеры подходов краткосрочного прогнозирования, основанных на макроскопических моделях потока трафика и расширенных KF, представлены в [4]. KF часто комбинируют с непараметрическим прогнозированием трафика, например, нейронными сетями [5].

Для целей краткосрочного прогнозирования макроскопические модели транспортных потоков являются экономным выбором, однако также были (более) успешно предложены иные подходы. В качестве примера можно назвать Dynasmart [20] и DynaMIT [21], в которых рассматривается выбор маршрута для конкретного пункта назначения — т.н. мезоскопический подход к моделированию трафика.

Другой популярный метод прогнозирования включает в себя клеточные автоматы (КА) [22], — дискретный и экономный метод представления трафика, в котором уделяется внимание алгоритмам выбора маршрута и матрицам корреспонденций [23]. В многоагентной системе [24] каждый водитель моделируется как агент, имеющий индивидуальное психическое состояние. В уме агента есть информация о поездке, счастье о его текущем состоянии, набор планов, например, когда уйти, и так далее. Эти модели были объединены с прогнозами по матрице корреспонденций [25]. Наконец, на основе теории трехфазного трафика Бориса Кернера [26] были разработаны модели ASDA и FOTO. Кернер различает три фазы в трафике: (а) свободный поток, (б) синхронизированный поток, фаза затора, когда движение на некотором участке не происходит вообще, и (в) широко распространенные заторы, фаза затора, в

которой транспорт движется вверх по течению с постоянной скоростью. ASDA и FOTO отслеживают и прогнозируют характеристики этих этапов. Представленные результаты показывают, что время расчетов меньше, чем при микроскопическом моделировании, и что точность является удовлетворительной, хотя используемые точные алгоритмы не раскрываются подробно из-за патентов.

С. Непараметрические модели

Термин «непараметрический» в данной терминологии не подразумевает «без параметров», но указывает на тот факт, что их количество и характер являются гибкими и не фиксированы заранее. Другими словами, как структура модели (например, степень многочлена, количество слоев в нейронной сети, функция расстояния в методах обучения без учителя), так и значения параметров модели определяются на основе, непосредственно, данных.

Эти два популярных семейства подходов прогнозирования трафика возникают в том случае, если функция G в уравнении (1) является линейной по своим параметрам θ_k (но необязательно по своим аргументам x_k , например, в эту категорию попадают многочлены общего вида $y = \sum_m [a_m x^m]$). Успех этих подходов в прогнозировании трафика является то, что выходные временные ряды действительно могут быть аппроксимированы и воспроизведены с помощью взвешенных линейных комбинаций (возможно преобразованных или предварительно обработанных) входных временных рядов.

Прогнозирование временных рядов включает моделирование переменной Y_k как параметризованной (взвешенной) линейной функции прошлых наблюдений

Y_{k-n} этой переменной и прошлых слагаемых ошибок e_k . По сути, они представляют собой особый случай регрессионных моделей, в которых входные данные

$x_{s|s < k}$ состоят из прошлых наблюдений за Y . Так, модели семейства ARMA предполагают, что процесс, генерирующий данные, является стационарным (т.е. во временных рядах нет структурной тенденции). Поскольку движение транспорта, как правило, не является стационарным процессом, обычно используется специальный член для моделирования структурных тенденций в данных. Результирующая модель обычно называется моделью ARIMA(p, d, q), где p , d и q - неотрицательные целые числа. В оригинальной книге Бокса и Дженкинса [27] описана методология для определения этих величин и их весов из данных. Ранний обзор моделей ARIMA, применяемых для прогнозирования трафика, можно найти в [28], и с тех пор в литературе доступно много примеров [19].

Кроме того, было предложено много вариаций или дополнений ARIMA, таких как сезонный ARIMA (SARIMA) [19], в котором добавляются периодические члены (которые обычно относятся ко времени суток, дню недели или другим тенденциям); ARIMA с подмножествами, которые разбивают временные ряды

на подмножества с соответствующими членами и компонентами [29]; вариации экспоненциального сглаживания [30]; модели, которые интегрированы с более сложными методами, например ARIMA с картами Кохонена [44]; модели, в которые включены также экзогенные входы (из других временных рядов) (ARIMAX) [31] и другие: VAR (I) MA, STAR (I) MA [32].

Важно отметить, что не все проблемы прогнозирования трафика, схематизированные уравнением (1), могут быть сведены к проблеме временных рядов. Самое главное, подходы ARIMA не применимы к прогнозированию времени отправления транспорта. Причина в том, что с периодами времени k фиксированного размера Δt (обычно 1, 5 или 10 минут в задачах краткосрочного прогнозирования) нет никакой гарантии, что предыдущее время отправления TT_{k-n} доступно для обработки, с n , обычно изображающим 1 или несколько дискретных периодов времени.

Фактически, TT_{k-n} измеряется в период времени $k^* = k - n + TT_{k-n}$, который в большинстве нетривиальных случаев будет позже, чем k . Эта проблема быстро усугубляется, если проводить прогнозирование на не слишком маленьком участке (более нескольких километров) и перегруженный маршрут со временем прохождения, обычно значительно большим, чем Δt . Подразумевается, что подходы временных рядов к прогнозированию времени в пути работают только в автономном режиме (при наличии данных), но не в реальных данных о движении транспорта или управляющих приложениях. Для прогнозирования времени в пути, а также для прогнозирования переменных трафика в тех местах, где как в пространстве, так и во времени доступны и необходимы дополнительные данные, подходы с линейной регрессией являются более подходящей альтернативой. В регрессионных моделях предполагается, что функция прогнозирования представляет собой линейную комбинацию ее ковариаций, где параметры указывают, какой вклад вносит ковариация в результат [26].

Также большое значение имеют различные модификации этих моделей. Примеры линейных KF-фильтров, применяемых для прогнозирования трафика, включают [33]. Примеры нелинейной (то есть расширенной) фильтрации Калмана, применяемой для прогнозирования трафика, показаны в [11, 24].

Дополнительные примеры моделей прогнозирования в этой широкой категории включают, среди прочего, модель ATHENA. В этой модели трафик моделируется как линейная комбинация исторических и текущих состояний. Для каждого типа трафика применяется нелинейное преобразование. Эта модель превзошла несколько других моделей [34]. В качестве последнего примера приведём так называемую модель AR (SETAR), которая использует линейную комбинацию текущего измерения и одного прошлого измерения для прогнозирования будущего состояния [35]. Но хотя эта модель и быстрая, точность оказалась неудовлетворительной.

Искусственные нейронные сети (ИНС) являются наиболее широко применяемыми моделями для решения проблемы прогнозирования трафика. В некоторых работах прогнозы строятся с использованием ИНС с обратным распространением ошибки. Другой подход — эволюционное обучение (генетические алгоритмы) был рассмотрен в [36].

Модульная нейронная сеть (МНС) основана на стратегии «разделяй и властвуй». Ввод обрабатывается в нескольких подсетях, каждая из которых специализируется на определенной задаче. За счёт этого МНС быстрее обучаются и показывают лучшие результаты [37].

Радиальные базовые частотные сети (RBFNN) используют скрытый уровень так называемых базовых функций для кластеризации входного пространства, причем каждый кластер представлен скрытым нейроном. Применение подходов к поставленной задаче показано в [38].

Так называемая вейвлет-нейронная сеть использует вейвлет-функции вместо стандартной сигмоидальной функции, используемой в BPNN [39], которая, как описано, приводит к повышению точности предсказания.

Спектральная базисная сеть (SNN) использует разложение Фурье входного вектора для получения линейно разделимых входных признаков [40]. Другой альтернативой является использование сети самоорганизующейся карты возможностей (SOFM) Кохонена для кластеризации входных данных перед подачей их в стандартный BPNN. В нейронной сети анализа главных компонент входные векторные данные сжимаются, а не расширяются, что уменьшает количество входов и, следовательно, сложность нейронной сети, что повышает производительность BPNN [37].

Сеть Джордана-Элмана или простая рекуррентная сеть (SRNN) содержит блоки памяти, которые используются для хранения выходных сигналов скрытого уровня на предыдущем временном шаге, обеспечивая механизм для распознавания повторяющихся шаблонов [37].

Нейронная сеть пространства состояний, в свою очередь, является еще одним вариантом сети Элмана [41]. SSNN значительно превосходит наивные методы прогнозирования времени в пути. В работе [41] показано, что SSNN, обученный методом Левенберга-Марквардта, имеет лучшую прогностическую эффективность, чем те же модели, обученные с помощью алгоритмов инкрементного обучения.

Что касается других методов из области машинного обучения, то они также применяются по отношению к изучаемой задаче. Можно назвать метод k -ближайших соседей [42] — быстрый метод, который может превосходить наивные методы прогнозирования, но обычно уступает более сложным моделям машинного обучения.

Нечеткая логика была применена к прогнозированию трафика в нескольких случаях, часто в сочетании с другими методами. Основной принцип заключается в том, что создается т.н. база правил, вручную или

автоматически. В [43] сообщается об обнадеживающих результатах.

Нечеткие модели сравниваются с BPNN и RBFNN в [44], где последние дают лучшие прогнозы.

Байесовская сеть убеждений (BBN), также известная как причинно-следственная модель, представляет собой ориентированную графическую модель, которая представляет условные зависимости между набором случайных величин. Этот метод применяется в реальной жизни одной из крупнейших информационных компаний Inrix [45], дочерней компанией Microsoft. Сравнение с другими методами не проводится по очевидным (коммерческим) причинам.

Наконец, регрессия опорных векторов (SVR) - это (популярный) метод машинного обучения, цель которого состоит в том, чтобы найти функцию, которая минимизирует эмпирическую ошибку, но при этом максимизирует зазор. Применение метода SVR можно найти в [46].

Все приведённые выше исследования касались краткосрочного прогнозирования автомобильного трафика, поскольку долгое время исследованию именно его уделялось значительно больше внимания, чем любому другому виду транспорта. Однако, в последние годы стали появляться и такие исследования тоже.

Например, в 2011 году в работе [47] подход ИНС был применён для прогнозирования состояния одной из линий метро столицы Тайваня, города Тайбэй. Авторы сравнивают модели BPN и SARIMA. Для SARIMA средняя абсолютная ошибка составила 30,39% — результат значительно хуже, чем у ИНС — от 5,04% до 14% в зависимости от конкретной реализации (рассмотрены разные варианты конфигурации).

Рекуррентная ИНС была применена для построения прогнозов для транспортной сети города Ренн, Франция в статье [48]. Примечательно, что данные собирались как с наземного, так и с подземного общественного транспорта. При этом среднеквадратичная погрешность составила 4,7% на тестовой выборке. Такая высокая точность объясняется, помимо прочего, тем, что объём данных для исследования достаточно большой — данные собирались в течение 15 месяцев.

В [49] прогноз для метро Лондона был получен при помощи ИНС типа сеть радиально-базисных функций. Величина абсолютной ошибки — 11,8%.

Касательно других методов, то в [50] был предложен гибридный метод для прогнозирования краткосрочных пассажирских потоков в пекинской системе метро, который объединяет вейвлет-преобразование и метод опорных векторов. Ошибка при прогнозировании пассажиропотока между станциями составляет 6,8% - 10,4%, что является довольно хорошим результатом. Там же рассматривается и прогнозирование другой величины — притока пассажиров на станцию. Точность этого прогноза — 9,8% - 18,09%.

В статье [51] сравнивались различные модели: ARMA, ИНС, метод опорных векторов, k ближайших соседей и Adaboost. Прогнозирование производилось отдельно для будних дней и для выходных. Интересно, что вычислительно более сложная Adaboost в точности

проигрывает методу k ближайших соседей. Однако, оба эти метода оказались менее точными, чем ИНС.

Примечательны также работы [52, 53], в которых рассматривается задача кластеризации дней по характеру трафика в транспортной сети. В результате было выделено 13 и 11 кластеров соответственно. При этом характер передвижений обычно специфичен во время больших праздников, выходных дней и каникул.

IV. ПРОГНОЗИРОВАНИЕ ПАССАЖИРОПОТОКА

Для удачного применения алгоритмов машинного обучения, входные данные должны содержать в себе как можно больше информации. Чтобы не ограничиваться одной величиной пассажиропотока от станции А до В, нужно построить более информативное признаковое пространство. Воспользуемся для этого визуализацией данных. На рис. 1 изображена интенсивность притока пассажиров на каждую станцию в течение каждого часа суток (для наглядности были просуммированы значения за целый месяц):

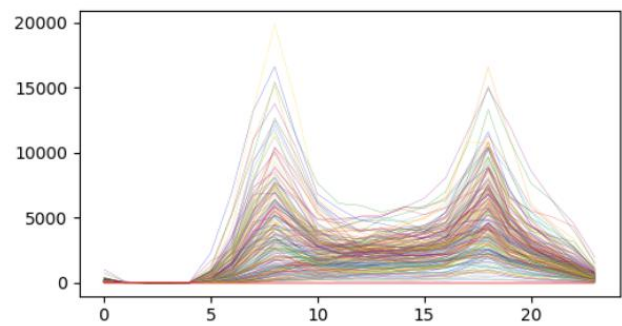


Рис. 1. Интенсивность притока пассажиров на разные станции метро

Хорошо видно, что самая большая нагрузка на транспортную систему приходится в т.н. часы пик. При этом величины этих самых пиков у разных станций сильно отличаются. Станции можно разбить на три большие категории (см. рис. 2). К первой отнесём станции с сильно преобладающим утренним пиком, ко второй — с вечерним, а к третьей — те, у которых величина пиков примерно одинакова. Такое разбиение хорошо показывает причину различия в динамике, а именно — специфику каждой отдельной части города:

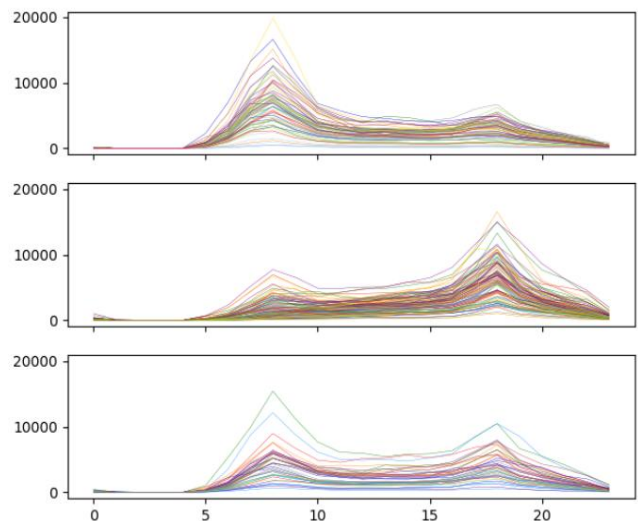


Рис. 2. Разбиение станций по соотношению пиков

ясно, что в спальнях районах утренний пик, когда люди едут на работу или учёбу будет больше; а в тех частях, где находятся деловые центры, образовательные учреждения, предприятия преобладать будет вечерний пик.

Если в соответствии с этим разбиением разметить схему метро в три цвета (рисунок 3), то станет ясно, чем объясняется отличие этих категорий станций: утренний пик преобладает в спальнях районах, а вечерний — там, где находятся деловые центры, образовательные учреждения и т.д. Таким образом, синий цвет, к примеру, — это описание границ рабочей зоны города. Преобладающий поток тут (шаблон использования станций) таков: пассажиры больше выходят на такой станции утром (едут на работу), а больше входят вечером (едут домой). Аналогично выделяется жилищная зона (красный цвет) и зона смешанного типа (зелёный). Однако зонирование, конечно, зависит от коэффициента при разбиении на рисунке 2. В данном случае пассажиропоток на нём превосходит пассажиропоток на другом пике более чем в 1,5 раза.

Отметим, что идея классификации станций по входному потоку (по сути — по соотношениям входы/выходы в пиковые час) сама по себе достаточно интересная. Это может быть темой отдельной работы.

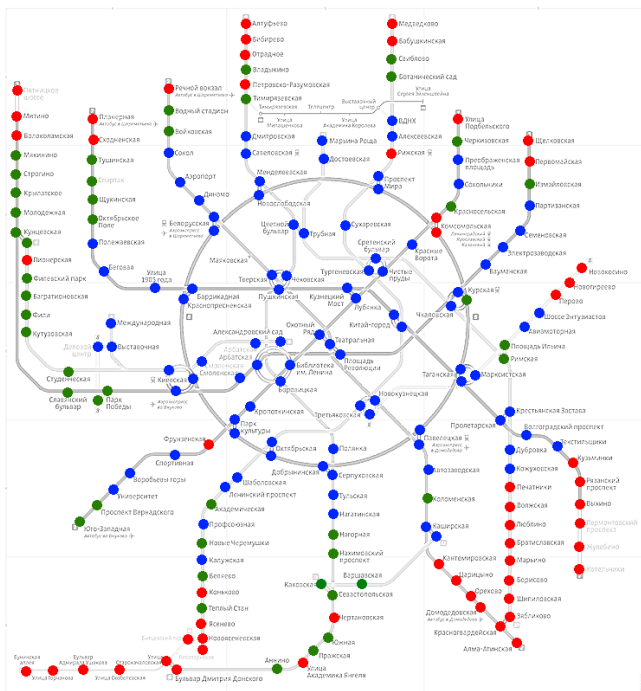


Рис. 3. Зонирование города по величине пиков.

Поскольку соседние станции, очевидно, чаще принадлежат одному району города с его динамикой пассажиропотока, можно выдвинуть предположение о том, что данные на соседних станциях могут быть информативным признаком.

Основная идея в прогнозировании состоит в том, что если на одной станции начинается всплеск активности пассажиров, (например, отток рабочих некоторого предприятия вечером или 'вбросы' пассажиров смежных транспортных систем на крайние станции метро в

утренние часы), то на соседних станциях, вероятно, активность тоже будет увеличиваться.

Для проверки этой гипотезы были рассмотрены корреляции корреспонденций между станциями «Ботанический сад» и «Беляево», а также всех соседних с ними (по 3 ближайших соседа в каждую сторону). Оказалось, что пассажиропотоки первых ближайших соседей, действительно, коррелируют с целевым пассажиропотоком сильнее, чем все остальные. Таким образом, имеет смысл учитывать эту информацию при построении краткосрочного прогноза. Кроме того, очевидно, нужно учитывать исторические данные пассажиропотока.

Так, даже при величине пространственного окна $w=1$ (8 соседних потоков и 1 целевой) и временного окна $t=5$, мы уже получаем 45 признаков для каждого пассажиропотока.

Эти данные исчерпывающе описывают каждый пассажиропоток, но при этом содержат и много избыточной информации, что в последствии приведёт к снижению точности прогноза. Чтобы избежать этого, а также для того, чтобы снизить вычислительную сложность алгоритма, уменьшив количество признаков, используя метод главных компонент — один из самых часто используемых способов понижения размерности с минимальной потерей информации.

Прежде чем описывать метод, введём понятие сингулярного числа.

Суть этого метода состоит в том, что для матрицы признаков $R (m \times n)$ может быть выполнено сингулярное разложение следующего вида:

$$R = U \Sigma C^T \quad (2)$$

Рис. 5. Корреляции соседних пассажиропотоков

где m — это количество наблюдений; n — количество признаков, описывающих каждое из них; Σ — прямоугольно-диагональная матрица размерности $m \times n$, на главной диагонали которой лежат неотрицательные числа, названные сингулярными числами; U и C — квадратные матрицы порядка m и n , столбцы которых являются сингулярными векторами для R , левыми и правыми, соответственно. Кроме того, столбцы матрицы C являются собственными значениями матрицы $(1/m)R^T R$, которая представляет собой матрицу ковариации наблюдений [], обозначим их как $\lambda_1, \lambda_2, \dots, \lambda_n$.

Обозначим векторы главных компонент как c_1, c_2, \dots, c_n , при чём c_1 будет обозначать направление, при проекции на ось которого выборка обладает самой большой дисперсией, c_2 — направление с самой большой возможной дисперсией при условии ортогональности к c_1 и т.д. по убыванию. Обозначим матрицу из первых k направлений как

$$C' = [c_1, c_2, c_3, \dots, c_k] \quad (3)$$

Тогда вектор из первых k главных компонент $Z = [z_1, z_2, \dots, z_k]^T$ может быть вычислен по формуле

$$z \approx (C')^T r \quad (4)$$

Где r — некоторая строка матрицы признаков R .

Для определения числа k (количества компонент в построенном пространстве меньшей размерности) используют правило Кайзера: значимы те главные компоненты, для которых

$$\lambda_i > \frac{1}{n} \text{tr}C \quad (5)$$

то есть собственное значение превосходит среднюю выборочную дисперсию.

Альтернативным подходом является правило сломанной трости. В этом случае нормированные собственные числа ($\lambda_i / \text{tr}C$) сравниваются с распределением длин обломков трости (длиной единица), сломанной в $(n-1)$ -ом случайно выбранном месте. Так, если L_i — длины обломков в порядке убывания длины, то математическое ожидание составляет

$$l_i = E(L_i) = \frac{1}{n} \sum_{j=i}^n \frac{1}{j} \quad (6)$$

И i -ю компоненту следует сохранить, если

$$\frac{\lambda_i}{\text{tr}C} > l_i \quad (7)$$

Для построения краткосрочных прогнозов на имеющихся данных будет испытано три разных подхода из области искусственного интеллекта: метод k -ближайших соседей, метод стохастического градиента и нейронная сеть.

Первый из них — метод k -ближайших соседей основывается на предположении о том, что близким (в смысле признакового пространства) объектам соответствуют близкие ответы (прогнозируемая величина в нашем случае). Для нового объекта X нужно найти k ближайших к нему объектов обучающей выборки x_1, x_2, \dots, x_k и, зная их метки, вычислить прогноз по формуле

$$\hat{y} = \frac{\sum_{i=1}^k y_i}{k} \quad (8)$$

Основным плюсом этого алгоритма является его интерпретируемость и простота, лёгкость в программной реализации. К минусам же можно отнести

большую вычислительную сложность — $O(kn)$ как для обучения, так и для прогнозирования, поскольку всю выборку необходимо держать в памяти.

Вторая выбранная модель машинного обучения — регрессор на основе метода стохастического градиента.

Из обучающей выборки $X^l = \{x_i, y_i\}_{i=1}^l$, нужно найти вектор весов W , при котором достигается минимум аппроксимированного эмпирического риска:

$$Q(w, X^l) = \sum_{i=1}^l L((w, x_i) y_i) \rightarrow \min_w \quad (9)$$

Для минимизации Q применяется метод градиентного спуска. В соответствии с ним, сначала выбирается некоторое начальное приближение для вектора W , а затем выполняется итерационный процесс, с каждым шагом которого этот вектор изменяется в направлении самого сильного убывания функционала Q . Это направление противоположно вектору градиента

$$Q'(w) = \left(\frac{\partial Q(w)}{\partial w_j} \right)_{j=1}^n :$$

$$w := w - \eta Q'(w) \quad (10)$$

где $\eta > 0$ — шаг в направлении скорейшего спуска, также он называется скоростью обучения. При этом название 'стохастический' означает, что прецеденты перебираются в случайном порядке..

Основным его преимуществом является то, что он подходит для динамического обучения, т.е. компьютеру необязательно держать весь набор данных в памяти одновременно и обучение может быть организовано итеративно. Данные при этом поступают потоком, и вектор весов обновляется при обработке каждого из них. Третьим испытанным методом стала искусственная нейронная сеть типа многослойный перцептрон, поскольку именно она давала наилучшие результаты прогнозирования в уже проведённых работах. По сути, какая сеть является суперпозицией некоторого количества линейных регрессоров. Значения всех признаков одновременно подаются на вход всех нейронов первого слоя. Затем над ними выполняются некоторые операции и пороговые функции (функции активации) генерируют выходные значения каждого из нейронов, которые впоследствии подаются на вход всех нейронов следующего слоя:

$$a^{(i+1)} = \sigma(Wa^{(i)} + b) \quad (11)$$

где $a^{(i)}$ — i -тый слой нейросети; $\sigma()$ — пороговая функция; W — матрица весовых коэффициентов для каждого нейрона; и b — это сдвиг, определяющий, насколько большой должна быть взвешенная сумма, чтобы активировать нейрон. При этом матрица W итеративно пересчитывается по методу стохастического градиента.

В общем случае сеть может содержать произвольное число слоёв. Все слои, за исключением входного и последнего, называются скрытыми. Сложность описанного алгоритма в значительной степени зависит от конкретной программной реализации, но чаще всего при использовании параллельных вычислений является довольно маленькой, и особенно при небольшом (1-3) числе скрытых слоёв. Функцией активации, была выбрана сигмоидальная логистическая функция:

$$\sigma = \frac{1}{1+e^{-\alpha y}} \quad (12)$$

где α — настраиваемый параметр модели.

Подбор всех параметров моделей, а также сравнение различных моделей проводилось по критерию точности построенного прогноза, а именно по средней абсолютной ошибке и средней квадратичной ошибке:

$$MAE = \frac{1}{N} \sum_{i=0}^{N-1} |\hat{y}_i - y_i|,$$

$$MSE = \frac{1}{N} \sum_{i=0}^{N-1} (\hat{y}_i - y_i)^2. \quad (13)$$

Для реализации рассмотренных алгоритмов был выбран язык программирования Python 3 как самый удобный для области машинного обучения и визуализации данных.

	KNN	SGD	MLP
MAE	4.5	11.4	6.9
MSE	104.9	1078.3	250.3

Таблица 1. Точность обученных моделей на тестовой выборке

Сравнение точности полученных моделей приведены в Таблице 1. Интересно, что самая простая модель KNN (лучшая точность при $k=2$) показывает самые лучшие результаты.

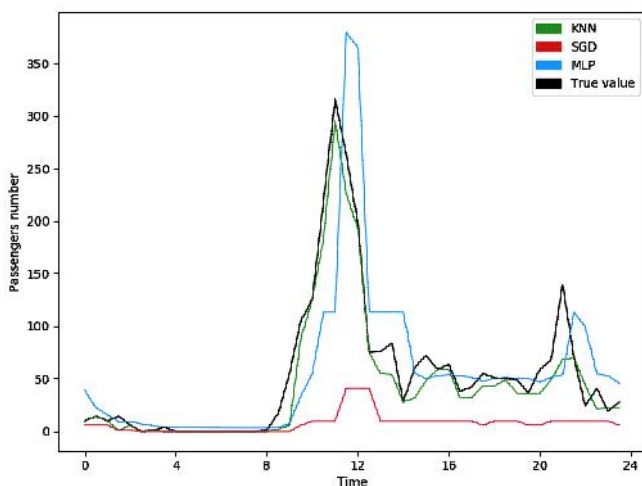


Рис. 4. Сравнение прогнозов с измеренными данными: Тёплый стан — Октябрьская, 13 февраля 2018

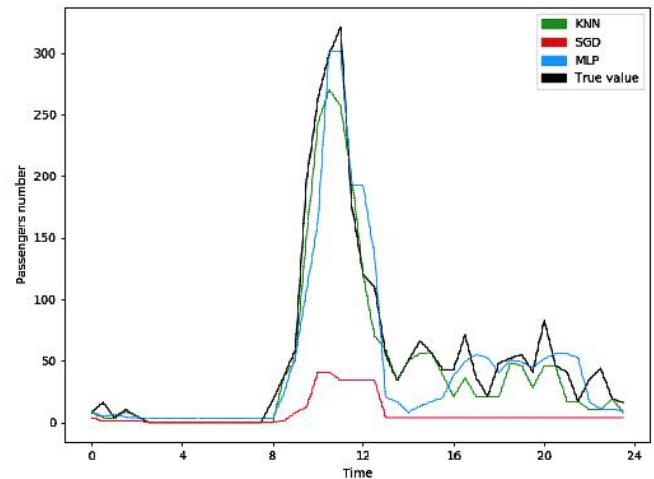


Рис. 5. Сравнение прогнозов с измеренными данными Бабушкинская — Рижская, 15 февраля 2018

V. ЗАКЛЮЧЕНИЕ

В процессе выполнения выпускной квалификационной работы были определены цели и задачи. Проведено библиографическое исследование и анализ литературы по теме. Рассмотрены и изучены понятия матрицы корреспонденций, краткосрочного прогнозирования, а также методов, которые применялись для его построения.

Был проведён анализ полученных данных транспортной системы Московского метрополитена, а также выбраны подходящие методы построения краткосрочных прогнозов по этим данным.

На языке Python3 разработан пакет из 9 программных модулей, каждый из которых решает определённую специфическую подзадачу.

В результате, цели выпускной квалификационной работы были выполнены и задачи решены.

БИБЛИОГРАФИЯ

- [1] Arem, B. v., H. R. Kirby, M. J. M. v. d. Vlist, and J. C. Whittaker. Recent Advances and Applications in the Field of Short-Term Traffic Forecasting. International Journal of Forecasting, Vol. 13, 1997.
- [2] Huisken, G., and M. v. Maarseveen. Congestion Prediction on Motorways: A Comparative Analysis. Proc., 7th World Congress on Intelligent Transport Systems, 2000.
- [3] Van Hinsbergen, C. P. I., J. W. C. Van Lint, and F. M. Sanders. Short-Term Traffic Prediction Models. Proc., 14th World Congress on Intelligent Transport Systems: ITS for a Better Life, Beijing: Research Institute of Highway, Chinese Ministry of Communications, 2007.
- [4] Wang Y., M. Papageorgiou, and A. Messmer. A Real Time Freeway Network Traffic Surveillance Tool. IEEE Transactions on Control Systems Technology, Vol. 14, No. 1, 2006.
- [5] Antoniou, C., M. E. Ben-Akiva, and H. N. Koutsopoulos. Online Calibration of Traffic Prediction Models. In Transportation Research Record: Journal of the Transportation Research Board, No. 1934, Transportation Research Board of the National Academies, Washington, D.C., 2005.
- [6] Calvert, S. C., J. W. C. Van Lint, and S. P. Hoogendoorn. A Hybrid Travel Time Prediction Framework for Planned Motorway Roadworks. Proc., 13th International IEEE Conference on Intelligent Transportation Systems, 2010.
- [7] Wu, C.-H., C.-C. Wei, D.-C. Su, M.-H. Chan, and J.-M. Ho. Travel Time Prediction with Support Vector Regression. Proc., 2003 IEEE

- Conference on Intelligent Transportation Systems, Shanghai, China, 2003.
- [8] Sun, H., H. X. Liu, H. Xiao, R. R. He, and B. Ran. Use of Local Linear Regression Model for Short-Term Traffic Forecasting. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1836, Transportation Research Board of the National Academies, Washington, D.C., 2003.
- [9] Williams, B. M., and L. A. Hoel. Modeling and Forecasting Vehicular Traffic Flow as a Seasonal ARIMA Process: A Theoretical Basis and Empirical Results. *Journal of Transportation Engineering*, Vol. 129, No. 6, 2003.
- [10] Al-Deek, H. M., M. P. D. Angelo, and M. C. Wang. Travel Time Prediction with Nonlinear Time Series. *Proc., Fifth International Conference on Applications of Advanced Technologies in Transportation*, Reston, Va., 1998.
- [11] Al-Deek, H. M., M. P. D. Angelo, and M. C. Wang. Travel Time Prediction with Nonlinear Time Series. *Proc., Fifth International Conference on Applications of Advanced Technologies in Transportation*, Reston, Va., 1998.
- [12] Stathopoulos, A., and M. G. Karlaftis. A Multivariate State Space Approach for Urban Traffic Flow Modeling and Prediction. *Transportation Research Part C: Emerging Technologies*, Vol. 11, No. 2, 2003.
- [13] Park, D., L. Rilett, and G. Han. Spectral Basis Neural Networks for Real-Time Travel Time Forecasting. *Journal of Transportation Engineering*, Vol. 125, No. 6, 1999.
- [14] Rilett, L. R., and D. Park. Direct Forecasting of Freeway Corridor Travel Times Using Spectral Basis Neural Networks. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1752, Transportation Research Board of the National Academies, Washington, D.C., 2001.
- [15] Nikovski, D., N. Nishiuma, Y. Goto, and H. Kumazawa. Univariate Short-Term Prediction of Road Travel Times. *Proc., 8th International IEEE Conference on Intelligent Transportation Systems*, 2005.
- [16] Hobeika, A. G., and C. K. Kim. Traffic-Flow-Prediction Systems Based on Upstream Traffic. *Proc., 1994 Vehicle Navigation and Information Systems Conference, IEEE*, 1994.
- [17] Wild, D. Short-Term Forecasting Based on a Transformation and Classification of Traffic Volume Time Series. *International Journal of Forecasting*, Vol. 13, 1997.
- [18] Voort, M. v. d., M. Dougherty, and S. Watson. Combining Kohonen Maps with ARIMA Time Series Models to Forecast Traffic Flow. *Transportation Research*, Vol. 4, No. 5, 1996.
- [19] Aboudolas, K., M. Papageorgiou, A. Kouvelas, and E. Kosmatopoulos. A Rolling-Horizon Quadratic Programming Approach to the Signal Control Problem in Large-Scale Congested Urban Road Networks. *Transportation Research Part C: Emerging Technologies*, Vol. 18, No. 5, 2010.
- [20] Mahmassani, H. S. DynaSMART-X Home. 2004.
- [21] Ben-Akiva, M. E., N. Bierlaire, D. Burton, H. N. Koutsopoulos, and R. G. Mishalani. Network State Estimation and Prediction for Real-Time Transportation Management Applications. Presented at 81st Annual Meeting of the Transportation Research Board, Washington, D.C., 2002.
- [22] Nagel, K., and M. Schreckenberg. A Cellular Automaton Model for Freeway Traffic. *Journal de Physique I*, Vol. 2, 1992.
- [23] Chopard, B., A. Dupuis, and P. O. Luthi. A Cellular Automata Model for Urban Traffic And Its Application to the City of Geneva. *Proceedings of Traffic and Granular Flow*, 1997.
- [24] Wahle, J., and M. Schreckenberg. A Multi-Agent System of On-Line Simulations Based on Real-World Traffic Data. *Proc., Hawaii International Conference on System Science, IEEE*, 2001.
- [25] Miska, M. P. Real-Time Traffic Management by Microscopic Online Simulation. Delft University of Technology, 2007.
- [26] Kerner, B. S. *The Physics of Traffic*. Springer-Verlag, 2004.
- [27] Box, G. E. P., and G. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-Day, 1976.
- [28] Nihan, N. L. Use of the Box and Jenkins Time Series Technique in Traffic Forecasting. *Transportation*, Vol. 9, 1980.
- [29] Lee, S., and D. B. Fambro. Application of Subset Autoregressive Integrated Moving Average Model for Short-Term Freeway Traffic Volume Forecasting. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1678, TRB, National Research Council, Washington, D.C., 1999.
- [30] Chrobok, R., O. Kaumann, J. Wahle, and M. Schreckenberg. Different Methods of Traffic Forecast Based on Real Data. *European Journal of Operational Research*, Vol. 155, 2004.
- [31] Williams, B. M. Multivariate Vehicular Traffic Flow Prediction. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1776, TRB, National Research Council, Washington, D.C., 1999.
- [32] Kamarianakis, Y., and P. Prastacos. Forecasting Traffic Flow Conditions in an Urban Network: Comparison of Multivariate and Univariate Approaches. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1857, Transportation Research Board of the National Academies, Washington, D.C., 2003.
- [33] Yang, J.-S. Travel Time Prediction Using the GPS Test Vehicle and Kalman Filtering Techniques. *Proc., 2005 American Control Conference*, 2005.
- [34] Kirby, H. R., S. M. Watson, and M. S. Dougherty. Should We Use Neural Networks or Statistical Models for Short-Term Motorway Traffic Forecasting? *International Journal of Forecasting*, Vol. 13, 1997.
- [35] Al-Deek, H., S. Ishak, and M. Wang. A New Short-Term Traffic Prediction and Incident Detection System on I-4. Florida Department of Transportation, University of Central Florida, 2001.
- [36] Lingras, P., S. Sharma, and M. Zhong. Prediction of Recreational Travel Using Genetically Designed Regression and Time-Delay Neural Network Models. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1805, Transportation Research Board of the National Academies, Washington, D.C., 2002.
- [37] Alecsandru, C.-D. A Hybrid Model-Based and Memory-Based Short-Term Traffic Prediction System. Louisiana State University and Agricultural and Mechanical College, 2003.
- [38] Park, B., C. J. Messer, and T. Urbanik II. Short-Term Freeway Traffic Volume Forecasting Using Radial Basis Function Neural Networks. In *Transportation Research Record 1651*, TRB, National Research Council, Washington, D.C., 1998.
- [39] Xie, Y., and Y. Zhang. A Wavelet Network Model for Short-Term Traffic Volume Forecasting. *Journal of Intelligent Transportation Systems*, Vol. 10, No. 3, 2006.
- [40] Rilett, L. R., and D. Park. Direct Forecasting of Freeway Corridor Travel Times Using Spectral Basis Neural Networks. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1752, Transportation Research Board of the National Academies, Washington, D.C., 2001.
- [41] van Lint, J. W. C. Online Learning Solutions for Freeway Travel Time Prediction. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 9, No. 1, 2008.
- [42] Smith, B. L., B. M. Williams, and R. K. Oswald. Comparison of Parametric and Nonparametric Models for Traffic Flow Forecasting. *Transportation Research Part C*, Vol. 10, 2002.
- [43] Coufal, D., and E. Turunen. Short Term Prediction of Highway Travel Time Using Data Mining and Neuro-Fuzzy Methods. *Neural Network World*, Vol. 3–4, 2004.
- [44] Huisken, G. Soft-Computing Techniques Applied to Short-Term Traffic Flow Forecasting. *Systems Analysis Modelling Simulation*, Vol. 43, No. 2, 2003.
- [45] INRIX. Available at <http://www.inrix.com>. Accessed Jan. 3, 2007.
- [46] Wu, C.-H., J.-M. Ho, and D. T. Lee. Travel-Time Prediction with Support Vector Regression. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 5, No. 4, 2004.
- [47] Введение в математическое моделирование транспортных потоков: Учебное пособие / Издание 2-е, испр. и доп. А. В. Гасников и др. Под ред. А. В. Гасникова. — М.: МЦНМО, 2013. — ISBN 978-5-4439-0040-7
- [48] Хабаров В. И. Марковская модель транспортных корреспонденций / В. И. Хабаров, Д. О. Молодцов, С. Г. Хомяков // Доклады ТУСУР. – 2012. – № 1(25). – Ч. 1. – С. 113–117.
- [49] Yang Li, Xudong Wang, Shuo Sun, Xiaolei Ma, Guangquan Lu, Forecasting short-term subway passenger flow under special events scenarios using multiscale radial basis function networks. *Transportation Research Part C: Emerging Technologies*, Vol. 77, 2017.
- [50] Yuxing Sun a, Biao Leng , Wei Guan, A novel wavelet-SVM short-time passenger flow prediction in Beijing subway system. *Neurocomputing*, Vol. 166, 2015.
- [51] Xiaoqing Dai, Lijun Sun, Yanyan Xu. Short-Term Origin-Destination Based Metro Flow Prediction with Probabilistic Model Selection Approach. *Journal of Advanced Transportation* Vol. 2018.
- [52] Chao Yang, Fenfan Yan, Xiangdong Xu. Daily metro origin-destination pattern recognition using dimensionality reduction and

- clustering methods. 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)
- [53] Chao Yang, Fenfan Yan, Xiangdong Xu. Clustering Daily Metro Origin-Destination Matrix in Shenzhen China. *Applied Mechanics and Materials*, Vol. 743, 2015.
- [54] Namiot, Dmitry, Oleg Pokusaev, and Vasily Kupriyanovsky. "On railway stations statistics in Smart Cities." *International Journal of Open Information Technologies* 7.4 (2019): 19-24.
- [55] Namiot, Dmitry, Oleg Pokusaev, and Varvara Lazutkina. "On passenger flow data models for urban railways." *International Journal of Open Information Technologies* 6.3 (2018): 9-14.

Metro correspondence matrix analysis

Mariia Nekraplonna, Dmitry Namiot

Abstract— This article is devoted to the analysis of traffic flows based on correspondence matrices. Such matrices describe the number of displacements between two points for a certain time interval. From a practical point of view, data relating to the Moscow Metro are considered. Accordingly, the correspondence matrix describes the movement between stations. Theoretically, such data describes all the characteristics of passenger traffic. In practice, it depends, of course, on the selected data processing model. Often, such matrices are used only for simple statistics, such as the number of passengers transported over time. At the same time, the space-time information interesting for digital urbanism is lost. For example, how were the trips distributed over time, how stable are these distributions, etc. The paper provides a detailed review of existing approaches to the analysis of data in correspondence matrices. As a practical task, a short-term forecast of passenger traffic is considered. It is noted that the short-term traffic forecast is a challenge that has been the subject of many research papers in the past few decades. Most of the work has been historically devoted to the analysis of traffic flows exclusively by road transport. The study of railways and, in particular, underground transport with its specificity has long been ignored. Relevant studies have been conducted only recently.

Keywords— correspondence matrix, traffic flow, metro.

REFERENCES

- [1] Arem, B. v., H. R. Kirby, M. J. M. v. d. Vlist, and J. C. Whittaker. Recent Advances and Applications in the Field of Short-Term Traffic Forecasting. *International Journal of Forecasting*, Vol. 13, 1997.
- [2] Huisken, G., and M. v. Maarseveen. Congestion Prediction on Motorways: A Comparative Analysis. *Proc., 7th World Congress on Intelligent Transport Systems*, 2000.
- [3] Van Hinsbergen, C. P. I., J. W. C. Van Lint, and F. M. Sanders. Short-Term Traffic Prediction Models. *Proc., 14th World Congress on Intelligent Transport Systems: ITS for a Better Life*, Beijing: Research Institute of Highway, Chinese Ministry of Communications, 2007.
- [4] Wang Y., M. Papageorgiou, and A. Messmer. A Real Time Freeway Network Traffic Surveillance Tool. *IEEE Transactions on Control Systems Technology*, Vol. 14, No. 1, 2006.
- [5] Antoniou, C., M. E. Ben-Akiva, and H. N. Koutsopoulos. Online Calibration of Traffic Prediction Models. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1934, Transportation Research Board of the National Academies, Washington, D.C., 2005.
- [6] Calvert, S. C., J. W. C. Van Lint, and S. P. Hoogendoorn. A Hybrid Travel Time Prediction Framework for Planned Motorway Roadworks. *Proc., 13th International IEEE Conference on Intelligent Transportation Systems*, 2010.
- [7] Wu, C.-H., C.-C. Wei, D.-C. Su, M.-H. Chan, and J.-M. Ho. Travel Time Prediction with Support Vector Regression. *Proc., 2003 IEEE Conference on Intelligent Transportation Systems*, Shanghai, China, 2003.
- [8] Sun, H., H. X. Liu, H. Xiao, R. R. He, and B. Ran. Use of Local Linear Regression Model for Short-Term Traffic Forecasting. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1836, Transportation Research Board of the National Academies, Washington, D.C., 2003.
- [9] Williams, B. M., and L. A. Hoel. Modeling and Forecasting Vehicular Traffic Flow as a Seasonal ARIMA Process: A Theoretical Basis and Empirical Results. *Journal of Transportation Engineering*, Vol. 129, No. 6, 2003.
- [10] Al-Deek, H. M., M. P. D. Angelo, and M. C. Wang. Travel Time Prediction with Nonlinear Time Series. *Proc., Fifth International Conference on Applications of Advanced Technologies in Transportation*, Reston, Va., 1998.
- [11] Al-Deek, H. M., M. P. D. Angelo, and M. C. Wang. Travel Time Prediction with Nonlinear Time Series. *Proc., Fifth International Conference on Applications of Advanced Technologies in Transportation*, Reston, Va., 1998.
- [12] Stathopoulos, A., and M. G. Karlaftis. A Multivariate State Space Approach for Urban Traffic Flow Modeling and Prediction. *Transportation Research Part C: Emerging Technologies*, Vol. 11, No. 2, 2003.
- [13] Park, D., L. Rilett, and G. Han. Spectral Basis Neural Networks for Real-Time Travel Time Forecasting. *Journal of Transportation Engineering*, Vol. 125, No. 6, 1999.
- [14] Rilett, L. R., and D. Park. Direct Forecasting of Freeway Corridor Travel Times Using Spectral Basis Neural Networks. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1752, Transportation Research Board of the National Academies, Washington, D.C., 2001.
- [15] Nikovski, D., N. Nishiuma, Y. Goto, and H. Kumazawa. Univariate Short-Term Prediction of Road Travel Times. *Proc., 8th International IEEE Conference on Intelligent Transportation Systems*, 2005.
- [16] Hobeika, A. G., and C. K. Kim. Traffic-Flow-Prediction Systems Based on Upstream Traffic. *Proc., 1994 Vehicle Navigation and Information Systems Conference*, IEEE, 1994.
- [17] Wild, D. Short-Term Forecasting Based on a Transformation and Classification of Traffic Volume Time Series. *International Journal of Forecasting*, Vol. 13, 1997.
- [18] Voort, M. v. d., M. Dougherty, and S. Watson. Combining Kohonen Maps with ARIMA Time Series Models to Forecast Traffic Flow. *Transportation Research*, Vol. 4, No. 5, 1996.
- [19] Aboudolas, K., M. Papageorgiou, A. Kouvelas, and E. Kosmatopoulos. A Rolling-Horizon Quadratic Programming Approach to the Signal Control Problem in Large-Scale Congested Urban Road Networks. *Transportation Research Part C: Emerging Technologies*, Vol. 18, No. 5, 2010.
- [20] Mahmassani, H. S. DynaSMART-X Home. 2004.
- [21] Ben-Akiva, M. E., N. Bierlaire, D. Burton, H. N. Koutsopoulos, and R. G. Mishalani. Network State Estimation and Prediction for Real-Time Transportation Management Applications. Presented at 81st Annual Meeting of the Transportation Research Board, Washington, D.C., 2002.
- [22] Nagel, K., and M. Schreckenberg. A Cellular Automaton Model for Freeway Traffic. *Journal de Physique I*, Vol. 2, 1992.
- [23] Chopard, B., A. Dupuis, and P. O. Luthi. A Cellular Automata Model for Urban Traffic Andits Application to the City of Geneva. *Proceedings of Traffic and Granular Flow*, 1997.
- [24] Wahle, J., and M. Schreckenberg. A Multi-Agent System of On-Line Simulations Based on RealWorld Traffic Data. *Proc., Hawaii International Conference on System Science*, IEEE, 2001.
- [25] Miska, M. P. Real-Time Traffic Management by Microscopic Online Simulation. Delft University of Technology, 2007.
- [26] Kerner, B. S. *The Physics of Traffic*. Springer-Verlag, 2004.
- [27] Box, G. E. P., and G. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-Day, 1976.
- [28] Nihan, N. L. Use of the Box and Jenkins Time Series Technique in Traffic Forecasting. *Transportation*, Vol. 9, 1980.
- [29] Lee, S., and D. B. Fambro. Application of Subset Autoregressive Integrated Moving Average Model for Short-Term Freeway Traffic Volume Forecasting. In *Transportation Research Record: Journal of the Transportation Research Board*, No. 1678, TRB, National Research Council, Washington, D.C., 1999.
- [30] Chrobok, R., O. Kaumann, J. Wahle, and M. Schreckenberg. Different Methods of Traffic Forecast Based on Real Data. *European Journal of Operational Research*, Vol. 155, 2004.

- [31] Williams, B. M. Multivariate Vehicular Traffic Flow Prediction. In Transportation Research Record: Journal of the Transportation Research Board, No. 1776, TRB, National Research Council, Washington, D.C., 1999.
- [32] Kamarianakis, Y., and P. Prastacos. Forecasting Traffic Flow Conditions in an Urban Network: Comparison of Multivariate and Univariate Approaches. In Transportation Research Record: Journal of the Transportation Research Board, No. 1857, Transportation Research Board of the National Academies, Washington, D.C., 2003.
- [33] Yang, J.-S. Travel Time Prediction Using the GPS Test Vehicle and Kalman Filtering Techniques. Proc., 2005 American Control Conference, 2005.
- [34] Kirby, H. R., S. M. Watson, and M. S. Dougherty. Should We Use Neural Networks or Statistical Models for Short-Term Motorway Traffic Forecasting? International Journal of Forecasting, Vol. 13, 1997.
- [35] Al-Deek, H., S. Ishak, and M. Wang. A New Short-Term Traffic Prediction and Incident Detection System on I-4. Florida Department of Transportation, University of Central Florida, 2001.
- [36] Lingras, P., S. Sharma, and M. Zhong. Prediction of Recreational Travel Using Genetically Designed Regression and Time-Delay Neural Network Models. In Transportation Research Record: Journal of the Transportation Research Board, No. 1805, Transportation Research Board of the National Academies, Washington, D.C., 2002.
- [37] Alecsandru, C.-D. A Hybrid Model-Based and Memory-Based Short-Term Traffic Prediction System. Louisiana State University and Agricultural and Mechanical College, 2003.
- [38] Park, B., C. J. Messer, and T. Urbanik II. Short-Term Freeway Traffic Volume Forecasting Using Radial Basis Function Neural Networks. In Transportation Research Record 1651, TRB, National Research Council, Washington, D.C., 1998.
- [39] Xie, Y., and Y. Zhang. A Wavelet Network Model for Short-Term Traffic Volume Forecasting. Journal of Intelligent Transportation Systems, Vol. 10, No. 3, 2006.
- [40] Rilett, L. R., and D. Park. Direct Forecasting of Freeway Corridor Travel Times Using Spectral Basis Neural Networks. In Transportation Research Record: Journal of the Transportation Research Board, No. 1752, Transportation Research Board of the National Academies, Washington, D.C., 2001.
- [41] van Lint, J. W. C. Online Learning Solutions for Freeway Travel Time Prediction. IEEE Transactions on Intelligent Transportation Systems, Vol. 9, No. 1, 2008.
- [42] Smith, B. L., B. M. Williams, and R. K. Oswald. Comparison of Parametric and Nonparametric Models for Traffic Flow Forecasting. Transportation Research Part C, Vol. 10, 2002.
- [43] Coufal, D., and E. Turunen. Short Term Prediction of Highway Travel Time Using Data Mining and Neuro-Fuzzy Methods. Neural Network World, Vol. 3–4, 2004.
- [44] Huisken, G. Soft-Computing Techniques Applied to Short-Term Traffic Flow Forecasting. Systems Analysis Modelling Simulation, Vol. 43, No. 2, 2003.
- [45] INRIX. Available at <http://www.inrix.com>. Accessed Jan. 3, 2007.
- [46] Wu, C.-H., J.-M. Ho, and D. T. Lee. Travel-Time Prediction with Support Vector Regression. IEEE Transactions on Intelligent Transportation Systems, Vol. 5, No. 4, 2004.
- [47] Vvedenie v matematicheskoe modelirovanie transportnyh potokov: Uchebnoe posobie / Izdanie 2-e, ispr. i dop. A. V. Gasnikov i dr. Pod red. A. V. Gasnikova. — M.: MCNMO, 2013. — ISBN 978-5-4439-0040-7
- [48] Habarov V. I. Markovskaja model' transportnyh korrespondencij / V. I. Habarov, D. O. Molodcov, S. G. Homjakov // Doklady TUSUR. — 2012. — # 1(25). — Ch. 1. — S. 113–117.
- [49] Yang Li, Xudong Wang, Shuo Sun, Xiaolei Ma, Guangquan Lu, Forecasting short-term subway passenger flow under special events scenarios using multiscale radial basis function networks. Transportation Research Part C: Emerging Technologies, Vol. 77, 2017.
- [50] Yuxing Sun a, Biao Leng , Wei Guan, A novel wavelet-SVM short-time passenger flow prediction in Beijing subway system. Neurocomputing, Vol. 166, 2015.
- [51] Xiaoqing Dai, Lijun Sun, Yanyan Xu. Short-Term Origin-Destination Based Metro Flow Prediction with Probabilistic Model Selection Approach. Journal of Advanced Transportation Vol. 2018.
- [52] Chao Yang, Fenfan Yan, Xiangdong Xu. Daily metro origin-destination pattern recognition using dimensionality reduction and clustering methods. 2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)
- [53] Chao Yang, Fenfan Yan, Xiangdong Xu. Clustering Daily Metro Origin-Destination Matrix in Shenzhen China. Applied Mechanics and Materials, Vol. 743, 2015.
- [54] Namiot, Dmitry, Oleg Pokusaev, and Vasily Kupriyanovsky. "On railway stations statistics in Smart Cities." International Journal of Open Information Technologies 7.4 (2019): 19-24.
- [55] Namiot, Dmitry, Oleg Pokusaev, and Varvara Lazutkina. "On passenger flow data models for urban railways." International Journal of Open Information Technologies 6.3 (2018): 9-14.