

ОБНАРУЖЕНИЕ СОЦИАЛЬНЫХ СОБЫТИЙ В ПРОСТРАНСТВЕННО-ВРЕМЕННЫХ ДАННЫХ С ИСПОЛЬЗОВАНИЕМ CONV LSTM И RESNET50

DETECTING SOCIAL EVENTS IN SPATIO- TEMPORAL DATA USING CONV LSTM AND RESNET50

**Mohammad Hani
V. Pak**

Summary. The article focuses on event detection by analyzing data related to space and time. It examines the use of neural networks to process numerical data from mobile communications companies, with the goal of discovering social activities. The properties of communication data, particularly their connection to time and location, enable the prediction of potential social events in specific places and times.

Recent advancements in deep learning have significantly improved predictive capabilities. Many studies have used deep models, such as LSTM neural networks, to detect anomalies, but these often lack consideration of spatial features or are based on convolutional neural networks (CNNs) alone. However, no previous research has applied a combination of ConvLSTM-based neural networks and ResNet50 to this type of data.

Keywords: spatial-temporal, ConvLSTM, ResNet50, Mahalanobis distance, event detection.

Мохаммад Хани

Аспирант, Санкт-Петербургский Политехнический
Университет Петра Великого
mohammad.h@edu.spbstu.ru

Пак Вадим Геннадьевич

кандидат физико-математических наук, доцент,
Санкт-Петербургский Политехнический
Университет Петра Великого
pak_vg@spbstu.ru

Аннотация. Статья посвящена обнаружению событий путем анализа данных, относящихся к пространству и времени. В нем рассматривается использование нейронных сетей для обработки числовых данных, поступающих от компаний мобильной связи, с целью выявления социальной активности. Свойства коммуникационных данных, особенно их связь со временем и местоположением, позволяют прогнозировать потенциальные социальные события в определенных местах и в определенное время.

Недавние достижения в области глубокого обучения значительно улучшили возможности прогнозирования. Во многих исследованиях использовались глубокие модели, такие как нейронные сети LSTM, для обнаружения аномалий, но в них часто не учитываются пространственные особенности или они основаны только на сверточных нейронных сетях (CNN). Однако ни в одном из предыдущих исследований комбинация нейронных сетей на основе ConvLSTM и ResNet50 не применялась к этому типу данных.

Ключевые слова: пространственно-временные, ConvLSTM, ResNet50, расстояние Махаланобиса, обнаружение событий.

Введение

Власти и советы умных городов полагаются на данные и информацию в качестве основного источника вдохновения при планировании услуг и удобств для лучшего обслуживания и обеспечения безопасности своих граждан в цифровую эпоху.

С. Джеффри разработал модель прогнозирования трафика данных сотовой связи с использованием методов глубокого обучения. Он собирал данные с течением времени, сосредоточившись на единственной характеристике — сеансах Интернета. Затем Джеффри использовал нейронную сеть LSTM (Долговременная кратковременная память) и нейронную сеть прямой связи (FFNN) и сравнил результаты их прогнозирования с базовой моделью ARIMA и FFNN.

Результаты показали, что модель на основе LSTM давала лучшие прогнозы, чем две другие модели, и имела бо-

лее короткое время обучения, чем FFNN. Кроме того, модель ARIMA превзошла FFNN по производительности. [2]

С. ЭлЕлими и С. Мустафа протестировали производительность нейросетевых моделей ARIMA (2, 1, 0) и LSTM в трех разных природных местоположениях. Модели показали неодинаковые результаты в разных местах.

Анализируя данные на еженедельной и почасовой основе с использованием данных только из Интернета, они пришли к выводу, что результаты являются удовлетворительными для обоих подходов к моделированию. [3]

М.С. Парвез использовал алгоритмы кластеризации, включая k-средние, для выявления аномалий в данных. Алгоритмы иерархической кластеризации проанализировали четыре атрибута — входящие вызовы, исходящие вызовы, входящие текстовые сообщения и исходящие текстовые сообщения — в рамках данного набора

данных. После кластеризации группа с наименьшим количеством элементов была идентифицирована как содержащая аномалию. [4]

В литературе TSAD (обнаружение аномалий временных рядов) описаны два основных подхода (как показано на рис. 1): модели, основанные на прогнозировании, и модели, основанные на реконструкции. Модели, основанные на прогнозировании, обучаются предсказывать следующую временную метку, в то время как модели, основанные на реконструкции, предназначены для учета встраивания данных временных рядов. Классификация архитектур глубокого обучения, используемых в TSAD, представлена на рис. 2. [9]

Наш вклад находится в контексте моделей, основанных на прогнозировании, и в разделе RNN на рис. 2, Подходы, основанные на расстояниях, и аномалия точки / подпоследовательности в стабильном наборе данных. [9]

Нейронная сеть со сверточной долговременной кратковременной памятью

Группа исследователей, в том числе К. Сяо, Н. Чен, К. Ху, К. Ван, З. Сюй и Ю. Цай, разработала новую модель глубокого обучения под названием ConvLSTM. Эта архитектура сочетает в себе возможности сверточных нейронных сетей (CNN) и долговременной кратковременной памяти (LSTM), что делает ее хорошо подходящей для обработки многомерных пространственно-временных данных, таких как спутниковые снимки. Исследователи модифицировали предыдущие уравнения LSTM, включив операции свертки в компоненты gate, что позволило ConvLSTM эффективно обрабатывать и моделировать сложные пространственные и временные структуры в данных. [1]

Структура математической модели:

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \sigma_c(c_t)$$
$$i_t = \sigma_g(W_{xi} \otimes x_t + W_{hi} \otimes h_{t-1} + W_{ci} \otimes c_{t-1} + b_i)$$
$$f_t = \sigma_g(W_{xf} \otimes x_t + W_{hf} \otimes h_{t-1} + W_{cf} \otimes c_{t-1} + b_f)$$
$$g_t = \sigma_c(W_{xc} \otimes x_t + W_{hc} \otimes h_{t-1} + b_c)$$
$$o_t = \sigma_c(W_{xo} \otimes x_t + W_{ho} \otimes h_{t-1} + W_{co} \otimes c_t + b_o)$$

Математические соотношения уравнения 1, представляющие компоненты модуля ConvLSTM

Обратите внимание, что:
Таблица 1.
Описание показателей, входящих в уравнение модели ConvLSTM

c_t	Состояние ячейки в момент времени t	x_t	Цепочка доходности единицы измерения LSTM
h_t	Скрытое состояние на шаге времени t	$w_{xo,xc,xi,xu}$	Ядро свертки применяется к тензору трехмерных входных данных x_t в каждом компоненте
i_t	Выходные данные элемента ввода на шаге времени t	$w_{ho,hc,hi,hu}$	Ядро свертки применяется к тензору трехмерных входных данных в каждом компоненте h_t
u_t	Выходные данные элемента забывания на шаге времени t	$b_{o,j,i,u}$	Коэффициенты смещения для каждого компонента
g_t	ячейка-кандидат на шаге времени t	σ_c	Продолжить активацию портала (сигмовидная)
o_t	Выходной вентиль выводится на шаге времени t	σ_g	Продолжаем активацию портала (tanh)
\odot	изделия Адамара	\otimes	сворачиваем (fold) изделие

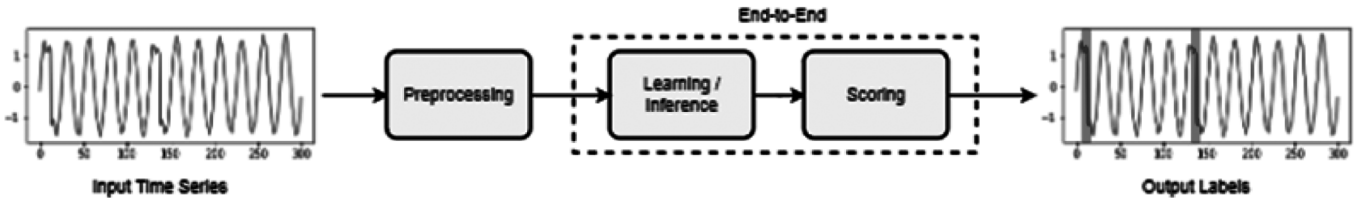


Рис. 1. Общие компоненты моделей глубокого обнаружения аномалий во временных рядах

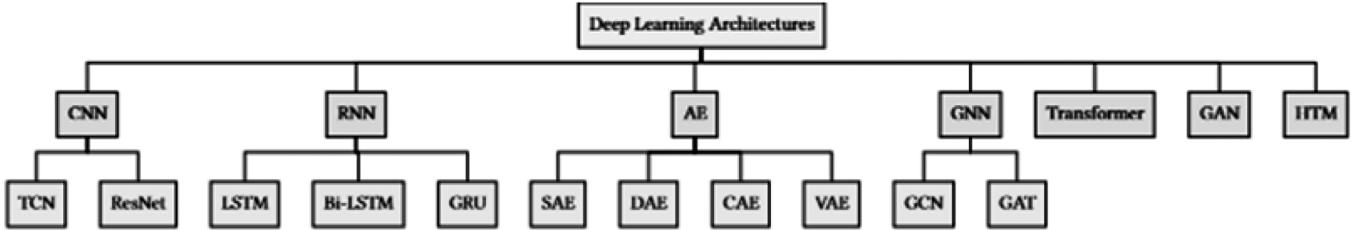


Рис. 2. Архитектуры глубокого обучения, используемые для обнаружения аномалий временных рядов

Остаточная сеть

С увеличением глубины разработанных нейронных сетей стало трудно понять влияние добавления слоев на увеличение сложности и семантики сетей.

Что еще более важно, проектировать сети сложно; увеличение количества уровней является более значительным, чем просто другая структура.

- Функциональные классы:

Рассмотрим F как класс функций, к которым пытаются получить доступ данная сетевая архитектура (с разной скоростью обучения и настройками параметров).

Для любого $f \in F$ существует набор параметров (таких как веса и отклонения), которые могут быть получены путем обучения на соответствующем наборе данных. Давайте предположим, что f^* — это функция «истины», которую мы хотим найти.

Если это в F , то у нас Хороший случай, но обычно нам бы так не повезло. Вместо этого мы попытаемся найти какую-нибудь функцию f_F^* , которая является лучшей функцией внутри F , на которую мы можем сделать ставку. Например, пусть у нас есть набор данных, содержащий атрибуты X и теги y . Мы пытаемся найти их, решая следующие примеры задач:

$$f_F^* \stackrel{\text{def}}{=} \operatorname{argmin}_f (L(X, y, f)) \text{ subject to } f \in F$$

Если вы разрабатываете другую, более мощную структуру, которая пытается достичь F' , чтобы она давала лучшие результаты. Мы ожидаем, что $f_{F'}^*$ лучше, чем f_F^* ; Но это не может быть гарантировано в случае с $F \not\subseteq F'$, наоборот, $f_{F'}^*$ может быть хуже. Мы замечаем на рисунке 1, в случае строк вложенных функций (левая часть),

что самая большая строка не обязательно приближается к функции «истинности» f^* , поскольку F_3 ближе всего к f^* compared to F_1 , в то время как F_6 moves farther away and there is nothing It ensures that increasing complexity can reduce the distance from f^* . Что касается случая строк вложенных функций (правая часть), этой проблемы можно избежать, поскольку $F_1 \subseteq F_2 \subseteq \dots \subseteq F_6$. В результате более крупный функциональный ряд (обозначенный областью) не гарантирует близости к функции. True, за исключением случаев неперекрывающихся функциональных рядов.

В контексте глубоких нейронных сетей, если мы сможем обучить вновь добавленный слой с помощью функции сопоставления $f(x) = x$, новая модель будет такой же эффективной, как и исходная. Поскольку новая модель может получить лучшее решение, соответствующее набору обучающих данных, добавленный уровень может помочь уменьшить количество ошибок при обучении.

- Остаточный блок

Учитывая входные данные x , наша цель — изучить функцию $f(x)$, которая будет использоваться в качестве входных данных для процесса активации на следующем уровне. Как показано в левой части рисунка 2, часть внутри пунктирной рамки должна непосредственно изучать функцию $f(x)$. С правой стороны часть внутри пунктирной рамки должна изучить остаточную функцию $f(x) - x$. Вот почему эта часть называется остаточной.

Если желаемой функцией соответствия является тождественная функция $f(x) = x$, то легче узнать соответствующий остаток. В этом случае нам нужно только передать нулевые веса и смещения из пунктирной рамки на следующий слой (например, полностью связанный слой или сверточный слой).

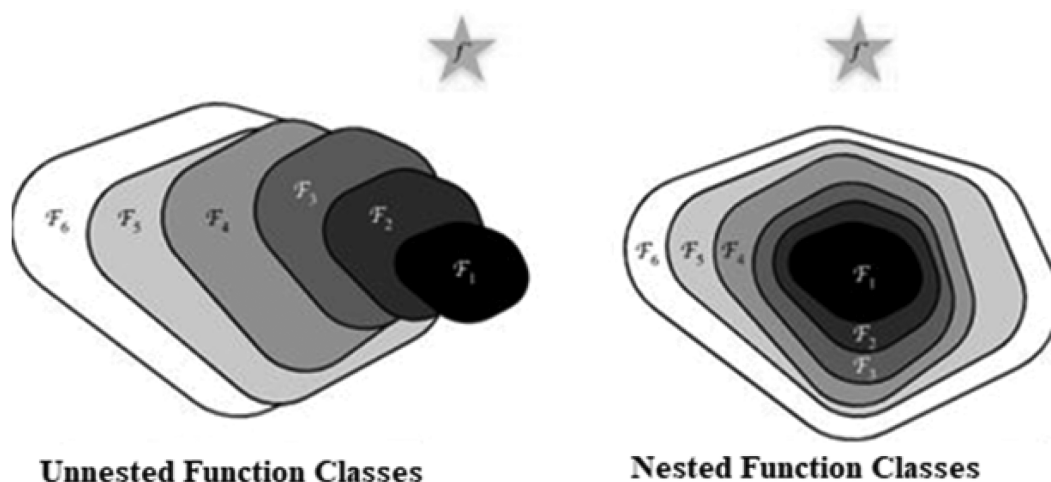


Рис. 3. Строки функций для строк неперекрывающихся функций и строк перекрывающихся функций задана (f^*) функция «истинности»

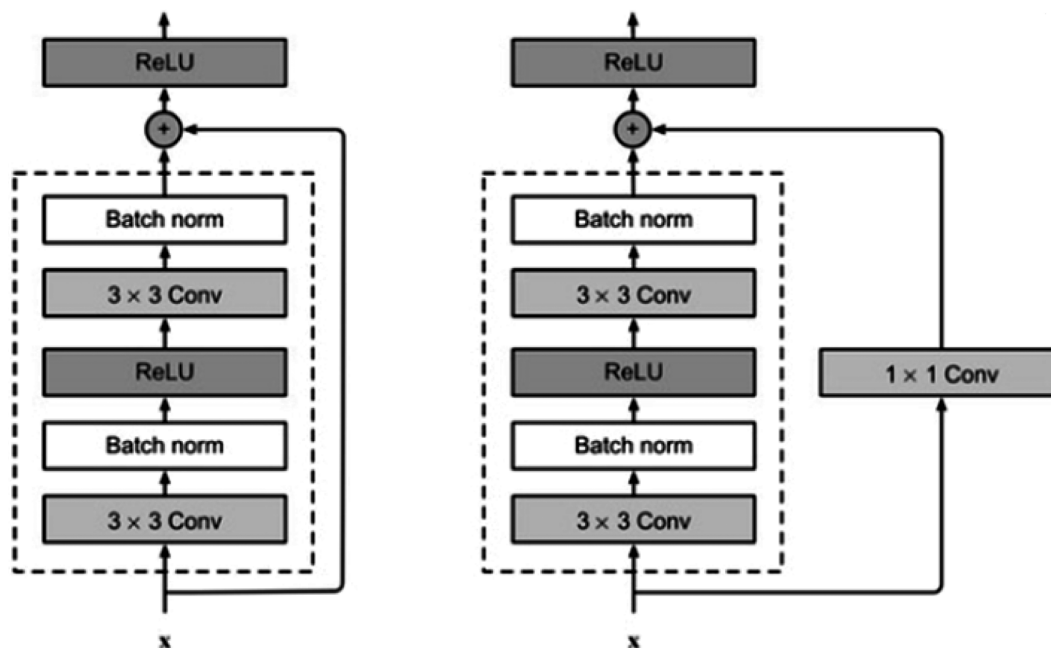


Рис. 4. Остаточный блок с конвблоком 1×1 и без блока идентификатора сверточного слоя

Остаточный блок и сверточный блок считаются фундаментальными строительными блоками любой остаточной сетевой модели, такой как ResNet-50 и ResNet-101. Эти базовые блоки повторяются несколько раз для построения общей сетевой архитектуры. [8]

Конструкция сверточного блока ResNet соответствует схеме сверточного слоя 3×3 в VGG. Этот блок содержит два сверточных слоя 3×3 с одинаковым количеством выходных каналов. За каждым сверточным слоем следует пакетная нормализация и повторная активация. Затем ResNet пропускает эти две операции свертки, напрямую добавляя входные данные к окончательной повторной активации.

Эта конструкция пропускного соединения требует, чтобы выходные данные двух сверточных слоев соответствовали форме ввода, чтобы их можно было сложить вместе. Если количество каналов необходимо изменить, можно использовать сверточные слои 1×1 для преобразования входных данных в желаемую форму или выполнить дополнительные добавления / удаления.

Расстояние Махаланобиса

Расстояние Махаланобиса — это эффективная многомерная мера расстояния, которая количественно определяет расстояние точки (наблюдения) от распределения данных. Профессор П.К. Махаланобис ввел эту метрику расстояния в 1936 году.

Что отличает его от евклидова расстояния, которое вычисляет расстояние между двумя точками, так это то, что:

1. Преобразует столбцы (атрибуты) в несвязанные переменные.

2. Масштабирует столбцы так, чтобы их дисперсия была равна 1.
3. Затем вычислите евклидово расстояние.

Горбани обсуждает, насколько расстояние Махаланобиса выгодно, когда объекты в наборе данных коррелированы. В таких случаях значения ковариационной матрицы будут большими. Деление на эту большую ковариацию (или умножение на величину, обратную ковариационной матрице) эффективно уменьшает расстояние между точками. И наоборот, если переменные некоррелированы, ковариация невелика и расстояние Махаланобиса существенно не уменьшается.

Таким образом, расстояние Махаланобиса решает проблемы стандартизации и взаимозависимости переменных, которые не может решить более простое евклидово расстояние. Например, при вычислении евклидовых расстояний расстояние между точкой p_2 и ее ближайшей точкой больше, чем расстояние между точкой p_1 и ее ближайшим соседом. Напротив, расстояние Махаланобиса гарантирует, что эти два расстояния равны.

Таким образом, расстояние Махаланобиса является мощным многомерным показателем, который учитывает корреляции в данных, что приводит к более значимым вычислениям расстояния по сравнению с евклидовым расстоянием. [5]

1. Постановка задачи

1.1. Описание предметной области

Процесс анализа данных включает в себя ряд важных этапов: сбор, хранение, обработку, очистку и анализ

данных. Точность каждого этапа напрямую влияет на надежность последующего этапа, что делает каждый этап важным компонентом общего анализа.

- Сбор данных: Данные, использованные в исследовании, являются телекоммуникационными данными из открытых источников. В рамках проекта 2014 Big Data Challenge, опубликованного Telecom Italia и SpazioDati, для Милана и Торонто в Италии были предоставлены наборы данных SMS, звонков и интернет-коммуникаций с пространственной связью.
- Хранение данных: используемые цифровые данные хранятся в текстовых файлах.
- Обработка данных: На этом этапе текстовые файлы считываются и преобразуются в структуры с использованием языка программирования Python.
- Очистка данных: На этом этапе мы изучили недостающие данные и заменили их средним арифметическим значений для каждого типа информации. Заполните данные, относящиеся к предполагаемой области. Обследование и территория вокруг него также использовались для изучения влияния предлагаемой территории на данные из прилегающих районов.
- Анализ данных: На этом этапе необработанные данные преобразуются в полезную информацию.

2. Системное моделирование

2.1. Обоснование выбора языка моделирования

2.1.1. Обработка исходных данных

В этой статье использовался набор данных Milan city dataset, охватывающий 550 квадратных километров города Милан; он был разделен на группу квадратов, каждый из которых имел идентификатор.

Сбор данных основан на записях данных вызовов (CDR). Эти журналы содержат различные атрибуты активности, которые отражают активность пользователя, а именно:

- SmsIn: представляет значение, пропорциональное количеству SMS, полученных почтовым ящиком за определенный период.
- SmsOut: представляет значение, пропорциональное количеству SMS, полученных почтовым ящиком за определенный период.
- CallIn: представляет значение, пропорциональное количеству вызовов, полученных абонентским ящиком за определенный период.
- ВыноСка: представляет значение, пропорциональное количеству вызовов, отправленных ящиком за определенный период.

- Интернет: количество записей, созданных для начала или завершения подключения к Интернету в пределах квадрата во временной области. Создается, когда сеанс длится 15 минут или каждые потребленные 5 Мбайт.

Если в пространственном блоке не происходит никаких действий, для этого блока не записывается запись. Были собраны данные по каждому квадранту за каждый временной интервал, а также репрезентативные значения для звонков, сообщений и интернет-сессий. Мы отмечаем, что в соответствии с политикой конфиденциальности данных значения атрибутов представляют собой не реальные значения, а скорее значения, пропорциональные реальным значениям. Более высокое значение представляет наибольшую активность для каждого атрибута данных.

В следующей таблице (табл. 2) приведен примерный набор данных по городу Милан. Первый столбец — это квадрат идентификатором после временного индекса.

Если в течение определенного периода внутри квадрата не происходит никаких действий, журнал не создается.

В то время как значение NAN для атрибута представляет отсутствие какой-либо активности во временной области. Короче говоря, эти данные содержат временную информацию, представленную десятиминутным интервалом времени, в дополнение к пространственным характеристикам, представленным идентификатором географического квадрата. Помимо активности коммуникаций, передвижения во времени и пространстве.

Поскольку данные не помечены, мы рассматриваем значения атрибутов активности пользователя в следующий момент времени в качестве целевых.

На этом этапе были собраны данные по территориям, прилегающим к миланскому району Сан-Сиро, протяженностью более 16 квадратов в длину и 16 квадратов в ширину.

Мы собрали значения входящих и исходящих атрибутов предыдущих сообщений за период в 20 минут, значения предыдущих входящих и исходящих сообщений за период в 20 минут и использование Интернета за период в 20 минут.

2.1.2. Статистика, основанная на утвержденных данных

Мы можем определить временной ряд как вектор X такой, что:

Таблица 2.

Характеристики коммуникативной деятельности данные

time	grid_square	cc	sms_in	sms_out	call_in	call_out	internet	weekday
2013-11-01 00:00:00	5528	39	0.487382	1.127274	0.158860	0.886130	18.682589	Friday
2013-11-01 00:10:00	5528	39	2.734001	2.517733	0.009358	0.136053	18.729035	Friday
2013-11-01 00:20:00	5528	39	0.763770	0.763770	0.481951	2.834657	17.003110	Friday
2013-11-01 00:30:00	5528	39	0.464752	0.634091	0.056240	0.542870	16.878082	Friday
2013-11-01 00:40:00	5528	39	1.054503	0.786900	0.606950	0.872858	18.176881	Friday
...
2014-01-01 23:10:00	7043	39	2.234346	3.138607	1.347595	0.399987	70.288754	Wednesday
2014-01-01 23:20:00	7043	39	0.799974	1.023001	1.060824	0.460843	59.822534	Wednesday
2014-01-01 23:30:00	7043	39	0.460843	1.462171	0.543372	0.740478	80.883076	Wednesday
2014-01-01 23:40:00	7043	39	2.309566	0.343379	1.159250	2.834657	74.720413	Wednesday
2014-01-01 23:50:00	7043	39	1.004216	1.140465	0.199994	2.834657	57.863466	Wednesday

$X = \{x_1, x_2, x_3, x_4 \dots, x_t\}$

Где x_t представляет данные на момент времени $i \in T, T = \{1, 2 \dots t\}$

- Изучив полученные данные, мы пришли к следующим наблюдениям:
- Данные превысили 100 миллионов просмотров. Во-первых, мы хотели проанализировать, как эта аудитория распределяется по дням недели. Мы наблюдали заметное увеличение просмотров по пятницам по сравнению с другими днями.

По средам и четвергам, которые длились 8 дней в течение двухмесячного периода, количество просмотров было ниже, чем в другие дни, когда было 9 дней.

- Хотя в эти среду и четверг было на один день меньше, общее количество просмотров все равно было сопоставимо с другими днями, за исключением меньшего количества просмотров по вторникам и воскресеньям.
- Общее количество просмотров отражает активное участие в разных часовых поясах в течение

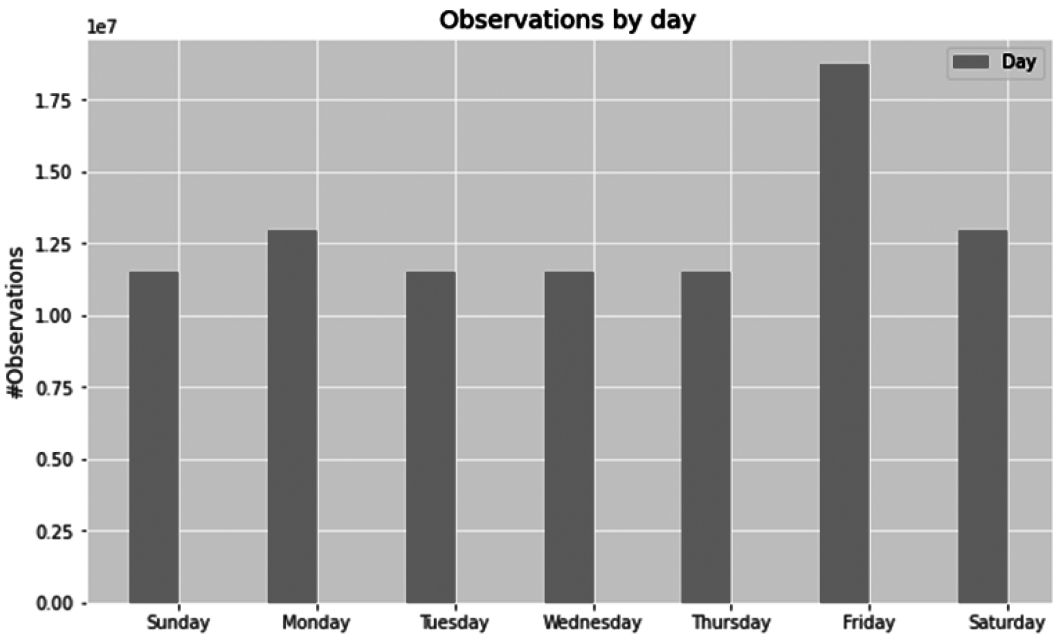


Рис. 5. Общее количество занятий в течение дня

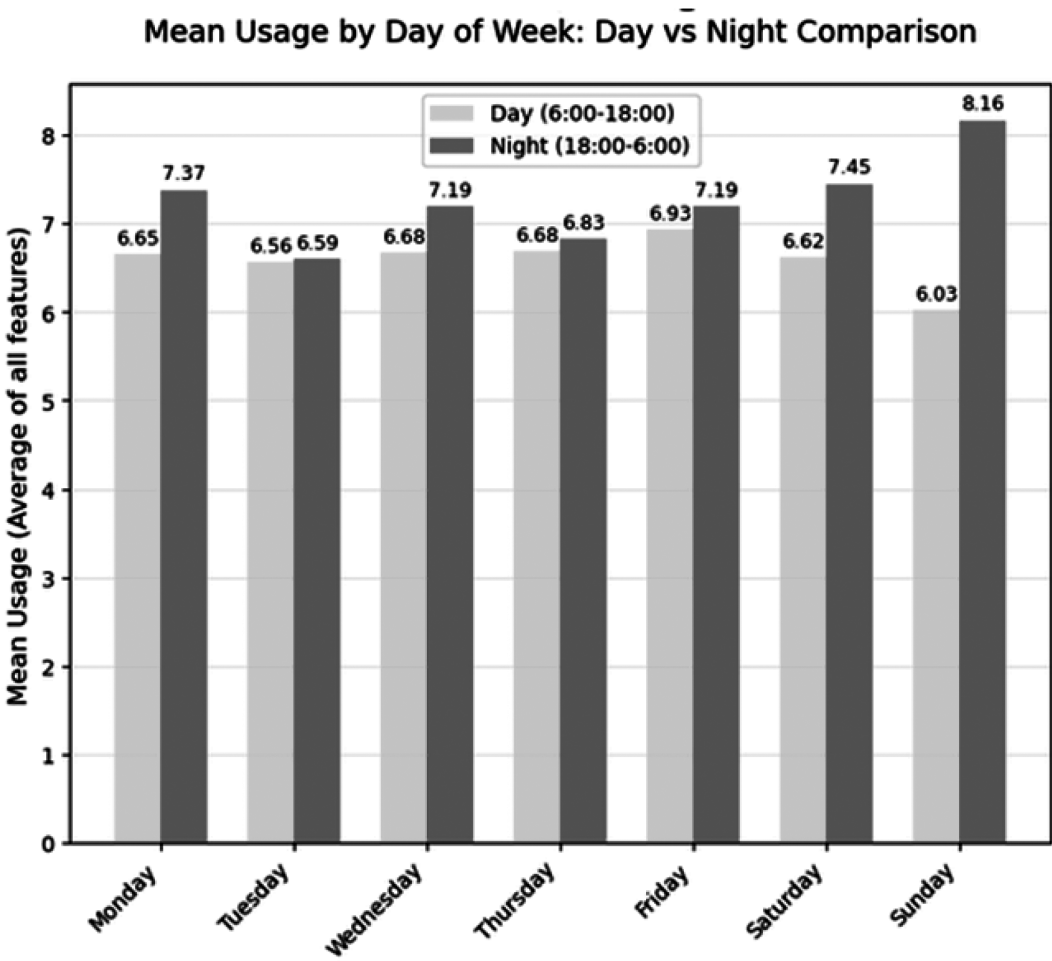


Рис. 6. Общее количество занятий в течение всего дня (утром и вечером)

дня, при этом распределение варьируется в зависимости от дня недели.

- Данные показывают, что по пятницам уровень дневной активности значительно выше по сравнению с другими днями, в то время как ночная активность остается такой же. Мы связываем эту закономерность с общим увеличением активности в последний рабочий день недели.
- Активность в городе меняется утром, в полдень и вечером.
- Активность в городе варьируется в зависимости от будних дней и официальных выходных.
- Существует взаимосвязь между функциями передачи данных (входящие и исходящие сообщения, отправленные и полученные сообщения и использование Интернета).

Данные, подлежащие изучению, должны быть стабильными и предсказуемыми, чтобы исключить любую очевидную корреляцию и коллинеарность с предыдущими данными.

Расширенный тест Дики-Фуллера (ADF) проверяет нулевую гипотезу о единичном корне в выборочном

Таблица 3.

Вставка 5638, Результат теста ADF для пяти атрибутов

Поле данных	Статистики тестирования	p-значение	стационарного
SMS-сообщения	-15.112500	7.660685176721849e-28	Да
SMS-сообщение	-13.772560	9.621714586151542e-26	Да
вызов	-14.049726	3.1934266203633457e-26	Да
вызов	-13.876625	6.321157964243545e-26	Да
Интернет	-12.193763	1.2682583594087165e-22	Да

временном ряду. Это версия теста Дики-Фуллера; но для большего и более сложного набора временных рядов. Статистика ADF, используемая в тесте, представляет собой отрицательное число. Чем более отрицательное значение, тем больше отвергается гипотеза, доказывающая существование авторегрессии.

Процедура тестирования для теста ADF такая же, как и для теста Дики-Фуллера, но она применяется к модели.

$$\Delta y_t = \alpha + \beta t + \gamma y_{t-1} + \delta_1 \Delta y_{t-1} + \dots + \delta_{p-1} \Delta y_{t-p+1} + \varepsilon_t$$

Где α — постоянная, β — коэффициент временной тенденции и p — порядок запаздывания процесса авторегрессии.

2.2. Построение модели

На этапе проектирования модели, поскольку основной целью является использование временной и пространственной информации, имеет смысл установить связи между соседними квадратами в миланской сетке, а также временные связи.

Именно здесь возникла идея использования модели глубокого обучения на основе нейронной сети ConvLSTM, где функция свертки находит пространственные корреляции между квадратами, а LSTM отличается способностью находить временные корреляции.

Идея добавления нейронной сети ResNet-50 возникла из работы [7], которая показала, насколько улучшается прогнозирование при использовании.

Модель, состоящая в основном из слоя ConvLSTM и слоя ResNet50.

- Слои Reshape и Permute — это первые шаги по инициализации соответствующих входных данных для ConvLSTM-слоя фигуры (TimeSteps_Samples, Высота, Ширина, каналы).
- Слой иммерсионного шума добавляет шум к входным данным. Это типичный процесс, который помогает предотвратить перетренированность нейронной сети с обучающими данными.
- Использовался слой ConvLSTM с десятью филь-

трами и ядром свертки 5x5; с ReLU в качестве функции активации; при использовании функции итеративной активации по умолчанию. В дополнение к функции отсева на 20 % для весов; Чтобы избежать переопределения обучающих данных; Он отключает нейроны (или устанавливает некоторые веса равными нулю) случайным образом, чтобы сеть не полагалась на характеристики обучающих данных и могла работать с тестовыми данными, которые значительно отличаются от обучающих данных. Слой Conv2D для получения желаемой выходной формы с размером ядра 5x5.

- Модель состоит из двух последовательностей. В первой последовательности используется ConvLSTM для изучения пространственно-временной структуры данных, за которой следует ResNet50 для изучения пространственной структуры из предыдущего состояния. Вторая последовательность также изучает пространственно-временную структуру данных, используя ConvLSTM, а затем слой Conv2D. Наконец, применяется средний слой для исправления полученной структуры.

Модифицированная модель ResNet50 содержит более 17 миллионов обучающих параметров с изменениями, внесенными на последнем уровне для получения желаемого результата.

2.2.1. Обнаружение событий

Принцип, принятый для проведения различия между выбросами и общественными мероприятиями, заключается в том, что общественные мероприятия длятся более 15 минут.

Для каждого деления каждого квадрата выполните следующие действия:

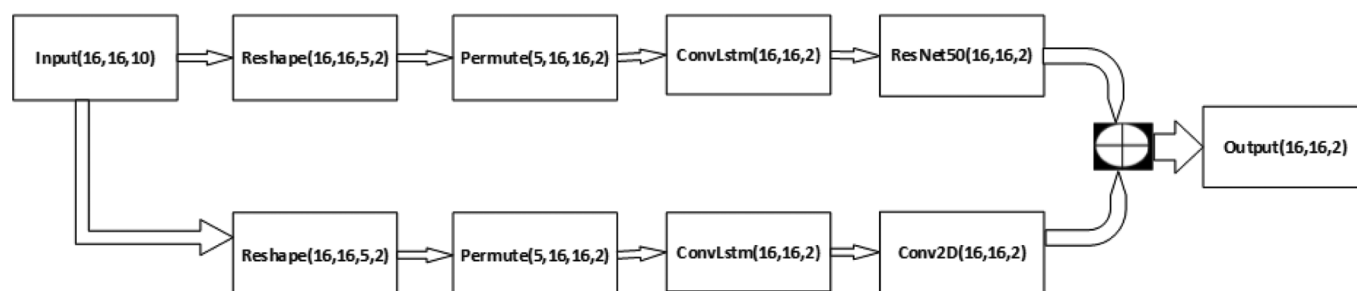


Рис. 7. Архитектура базовой модели на основе ConvLSTM

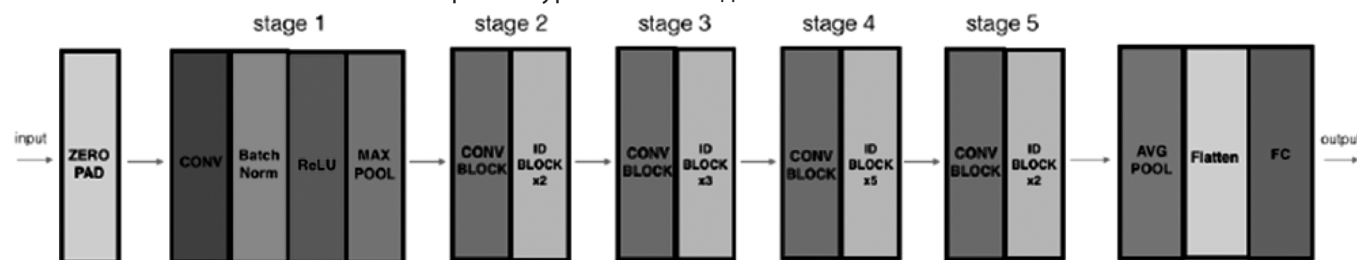


Рис. 8. Архитектура базовой сети ResNet50

Моделирование распределения данных в пределах нормального (ныряющего) распределения; основано на том факте, что оценка максимального правдоподобия может быть рассчитана для нормального распределения данных x_1, \dots, x_N .

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x^{(n)}$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \hat{\mu})(x^{(n)} - \hat{\mu})^T$$

Уравнение 1. Оценка среднего арифметического и ковариационной матрицы при нормальном распределении

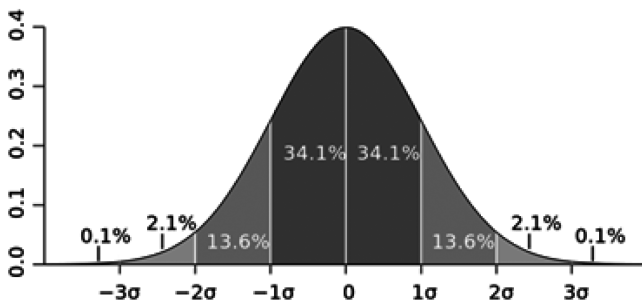


Рис. 9. Нормальное распределение

Однако, поскольку эта оценка очень чувствительна к наличию выбросов (событий) в наборе данных, соответствующие расстояния Махаланобиса вводят в заблуждение и неточны. Поэтому мы использовали робастную оценку ковариации, в которой мы полагались на использование метода MCD для более точной оценки ковариационной матрицы.

В зависимости от вычисления как $\hat{\mu}$, и ковариационной матрицы $\hat{\Sigma}$; Мы вычисляем расстояние Махаланобиса для каждого наблюдения на основе данных.

$$D_M(x) = \sqrt{(x - \hat{\mu})^T \hat{\Sigma}^{-1} (x - \hat{\mu})}$$

На этапе обнаружения аномалии мы полагаемся на ошибку прогноза; мы сравниваем прогнозируемую сеть с реальной сетью и вычисляем разницу между соответствующими квадратами. Таким образом, мы вычисляем аномальные значения в пределах распределения различий в данных между фактическими значениями и ожидаемыми значениями.

Затем мы находим точки с наибольшим расстоянием по распределению различий двух объектов на основе расстояния Махаланобиса и критерия хи-квадрат, которые являются кандидатами на представление общественных событий в городе.

2.2.2. Тестирование и внедрение

Данные были разделены на три группы.

[1] Группа обучения и тестирования, которая работает с 1 ноября 2013 года по 1 декабря 2013 года, где 100 % было посвящено обучению.

[2] Тестовая группа продлится с 1 декабря 2013 года по 1 января 2014 года.

Эксперимент проводился с использованием нескольких оптимизаторов (Adam, Adadelata, SGD, Adagrad, RMSProp, Adamax) и функций потерь (Binary_crossentropy, MSE, MAE).

Таблица 4.

Оценка эффективности предложенной модели на основе используемой функции оптимизации

Оптимизатор	Потерять функцию	Тестовая ошибка	Точность
Adam	MAE	0.0401	0.8285
Adam	MSE	0.0035	0.8319
Adam	binary_crossentropy	0.4578	0.8217

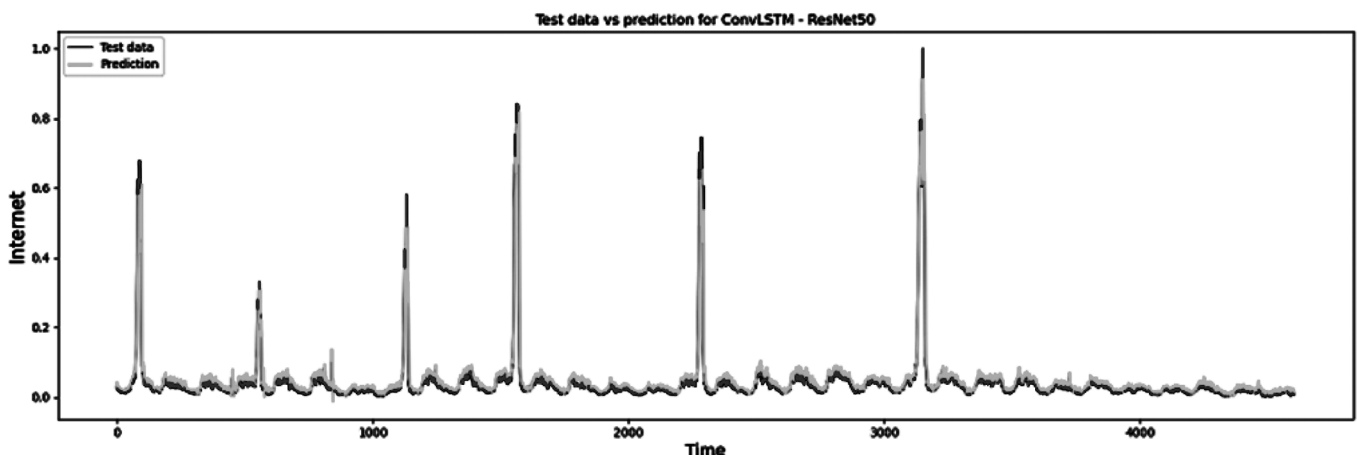


Рис. 10. На рис. 10 тестовые данные выделены черным цветом, а прогнозируемые — серым

Таблица 5.

Совпадения

	Аудитор- ия	Матч	Дата и время	Глава	1
1	43,706	«Интернационале» против Sampdoria 1-1	1 декабря 2013 15:00 pm	Inter Milan season	Рабочий день
2	12,714	Internazionale vs Trapani 3–2	4 декабря 2013 21:00 pm	Coppa Italia	Рабочий день
3	33,732	Internazionale vs. Parma 3-3	8 декабря 2013 20:45	Inter Milan season	Выход-ные
4	61,744	Milan vs. Ajax 0-0	11 декабря 2013 20:45	Inter Milan season	Рабочий день
5	37,987	Milan vs. Roma 2-2	16 декабря 2013 20:45	Inter Milan season	Выход-ные
6	79,311	Internazionale vs. Milan 1-0	22 декабря 2013 20:45	A.C. Milan season	Выход-ные

Использование функции ошибки Binary_crossentropy дает хорошие результаты для функции sms, хотя у нас хорошая точность.

В результате 6 из 6 инцидентов были зафиксированы в районе стадиона «Сан-Сиро» в декабре.

А 16 декабря 2013 года Кристиан Сапата, левый игрок «Милана», празднует вместе со своими товарищами по команде забитый гол в матче чемпионата Италии против «Ромы».

И 01 января 2014 года люди празднуют новый год, найденный с помощью функции Sms.

Заключение и будущая работа

В результате мы обнаружили, что ConvLSTM с ResNet50 эффективен при обнаружении событий и попытках извлечения пространственной и временной информации, но мы также обнаружили, что resnet50 потребляет много ресурсов во время обучения и требует много времени, теперь мы ожидаем обрабатывать другие источники данных для той же области исследования и периода времени в ближайшем будущем и использовать ConvLSTM в генеративной состязательной сети для достижения точного прогнозирования целевых данных.

ЛИТЕРАТУРА

1. Сяо К., Чен Н., Ху К., Ван К., Сюй З., Цай Ю. и др. (2019). «Пространственно-временная модель глубокого обучения для прогнозирования поля температуры поверхности моря с использованием временных рядов спутниковых данных». В разделе Моделирование окружающей среды и программное обеспечение. ELSEVIER.
2. Джеффри С. (2020). «Прогнозирование трафика сотовой связи с помощью рекуррентной нейронной сети». arXiv.org. arXiv: 2003.02807.
3. Элелими С., & Мустафа С. Большие данные в телекоммуникационной отрасли: эффективные методы прогнозирования на CDR. sc., <http://dx.doi.org/10.4108/eai.13-7-2018.164919>, (2020, 6).
4. Парвез М.С., Рават Д.Б. и Гаруба М. «Аналитика больших данных для анализа активности пользователей и обнаружения пользовательских аномалий в мобильной беспроводной сети». IEEE Transactions on Industrial Informatics, стр. 12, (2017).
5. Горбани Х. «Расстояние Махаланобиса и его применение для обнаружения многомерных выбросов». Информационный портал Университета математики. 2019; 34:583–95. doi: 10.22190/FUMI1903583G
6. <https://dandelion.eu/datagems/SpazioDati/telecom-sms-call-internet-mi/description/>, Дата обращения: 1 сентября 2019 г.
7. Хусейн Б., Ду К., Имран А. и Имран М.А. (2019). «Обнаружение аномалий в вычислительных сетях сотовой связи на базе мобильных передовых технологий с использованием искусственного интеллекта». IEEE transactions on industrial informatics.
8. Хэ К., Чжан Х., Рен С. и Сунь Дж. (2016). «Глубокое остаточное обучение для распознавания изображений». Материалы конференции IEEE по компьютерному зрению и распознаванию образов (стр. 770–778).
9. Захра Заманзаде Дарбан, Джеффри И. Уэбб, Шируи Пан, Чару К. Аггарвал и Махса Салехи. 2023 год. Глубокое обучение для обнаружения аномалий временных рядов: Обзор. 1, 1 (май 2023 г.), 42 страницы.

© Мохаммад Хани (mohammad.h@edu.spbstu.ru); Пак Вадим Геннадьевич (pak_vg@spbstu.ru)

Журнал «Современная наука: актуальные проблемы теории и практики»