

Тема: Прогноз розничных продаж на Python.

Цель: научиться понимать методы исследовательского анализа данных с помощью языка программирования Python.

Материальное и дидактическое оснащение: Методические рекомендации по выполнению практической работы.

ТЕОРЕТИЧЕСКИЕ СВЕДЕНИЯ

Анализ временных рядов

Временной ряд — это ряд точек данных, проиндексированных (или перечисленных, или нарисованных на графике) во временном порядке. Анализ временных рядов относится к методам извлечения значимой статистики из данных временных рядов. Это обычно используется для прогнозирования и других моделей.[Первоисточник](#)

Цели обучения

Понимание использования анализа временных рядов плюсы и минусы различных методов TSA, включая различие между линейными и нелинейными методами. Проанализируйте результаты на заданных данных магазина rossman.

Описание проблемы

Для этого ноутбука мы будем использовать базу данных продаж магазина rossman. Ниже приводится описание данных с веб-сайта:

«Rossmann управляет более чем 3000 аптеками в 7 европейских странах. В настоящее время перед менеджерами магазинов Rossmann стоит задача прогнозировать свои ежедневные продажи на срок до шести недель вперед. На продажи в магазинах влияет множество факторов, включая рекламные акции, конкуренция, школьные и государственные праздники, сезонность и местоположение. Поскольку тысячи отдельных менеджеров прогнозируют продажи на основе своих уникальных обстоятельств, точность результатов может быть весьма различной».

ПРАКТИЧЕСКАЯ ЧАСТЬ РАБОТЫ

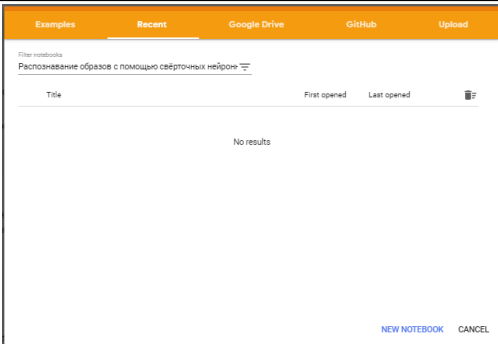
В ходе изучения теоретических материалов теоретических сведений данной работы:

- I. Ознакомьтесь с теоретическими сведениями.
- II. Изучите инструкцию ниже.
- III. Предоставьте ссылку на свою работу в <https://colab.research.google.com/> онлайн, в процессе выполнения преподавателю на uvmranh48@gmail.com, через вкладку Share (поделиться). Для отслеживания Вашей работы преподавателем онлайн и возможности консультирования
- IV. Оформите работу, используя снимки экрана.
- V. Ссылку на работу и отчет по работе выложите в ментальной карте.

- VI. Сделайте алгоритм по Вашей работе в draw.io
 VII. Выкладывайте отчет на странице своего сайта.

Ход работы.

Во время работы не забывайте запускать на каждом этапе код.

	<p>Создайте новый ноутбук по ссылке: https://colab.research.google.com/.</p> <p>Назовите: Исследовательский анализ данных на Python (Ф.И.О.).ipynb. Подпишите в скобках своим Ф.И.О.</p>
<p>Ваш новый ноутбук должен содержать комментарии и иметь законченный вид. Текст для комментариев можем использовать из теоретических сведений.</p>	
<p>Импортируем необходимые библиотеки</p>	<pre># Library Imports import numpy as np import pandas as pd import matplotlib import seaborn as sns import matplotlib.pyplot as plt import matplotlib from scipy.stats import skew from scipy.stats.stats import pearsonr from math import sqrt from sklearn.metrics import mean_squared_error # matplotlib parameters matplotlib.rcParams['axes.labelsize'] = 14 matplotlib.rcParams['xtick.labelsize'] = 12 matplotlib.rcParams['ytick.labelsize'] = 12 matplotlib.rcParams['text.color'] = 'k' %config InlineBackend.figure_format = 'retina' %matplotlib inline</pre>
<p>Чтение данных</p>	<pre># Data Reading train = pd.read_csv('https://raw.githubusercontent.com/RPI-DATA/tutorials-intro/master/rossmann-store-sales/rossmann-store-sales/train.csv', parse_dates = True, low_memory = False, index_col = 'Date')</pre>

	<pre>store = pd.read_csv('https://raw.githubusercontent.com/RPI-DATA/tutorials-intro/master/rossmann-store-sales/rossmann-store-sales/store.csv', low_memory = False)</pre>																																																																								
<p>Исследовательский анализ данных</p> <p>Начнем с просмотра того, из чего состоят наши данные. Мы хотим увидеть, какие переменные являются непрерывными, а какие — категориальными. Изучив некоторые данные, мы видим, что можем создать функцию. Количество продаж, разделенное на клиентов, может дать нам хороший показатель для измерения средних продаж на одного клиента. Мы также можем сделать предположение, что если в этом столбце отсутствуют значения, то у нас 0 клиентов. Поскольку клиенты стимулируют продажи, мы решили удалить все эти значения.</p>	<pre>train.head()</pre> <table><thead><tr><th></th><th>Store</th><th>DayOfWeek</th><th>Sales</th><th>Customers</th><th>Open</th><th>Promo</th><th>StateHoliday</th><th>SchoolHoliday</th></tr></thead><tbody><tr><td>Date</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td>2015-07-31</td><td>1</td><td>5</td><td>5263</td><td>555</td><td>1</td><td>1</td><td>0</td><td>1</td></tr><tr><td>2015-07-31</td><td>2</td><td>5</td><td>6064</td><td>625</td><td>1</td><td>1</td><td>0</td><td>1</td></tr><tr><td>2015-07-31</td><td>3</td><td>5</td><td>8314</td><td>821</td><td>1</td><td>1</td><td>0</td><td>1</td></tr><tr><td>2015-07-31</td><td>4</td><td>5</td><td>13995</td><td>1498</td><td>1</td><td>1</td><td>0</td><td>1</td></tr><tr><td>2015-07-31</td><td>5</td><td>5</td><td>4822</td><td>559</td><td>1</td><td>1</td><td>0</td><td>1</td></tr></tbody></table>		Store	DayOfWeek	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday	Date									2015-07-31	1	5	5263	555	1	1	0	1	2015-07-31	2	5	6064	625	1	1	0	1	2015-07-31	3	5	8314	821	1	1	0	1	2015-07-31	4	5	13995	1498	1	1	0	1	2015-07-31	5	5	4822	559	1	1	0	1									
	Store	DayOfWeek	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday																																																																	
Date																																																																									
2015-07-31	1	5	5263	555	1	1	0	1																																																																	
2015-07-31	2	5	6064	625	1	1	0	1																																																																	
2015-07-31	3	5	8314	821	1	1	0	1																																																																	
2015-07-31	4	5	13995	1498	1	1	0	1																																																																	
2015-07-31	5	5	4822	559	1	1	0	1																																																																	
	<pre>train.shape</pre> <p>(1017209, 8)</p>																																																																								
<p>Обратите внимание на порядок, в котором перечислены данные. Он упорядочен от самой последней даты до самой старой даты. Это может вызвать проблемы при разработке нашей модели.</p>																																																																									
<p>После этого мы будем использовать удивительную. describe()функцию, которая может предоставить нам большинство статистических элементов.</p>	<pre>train.describe()</pre> <table><thead><tr><th></th><th>Store</th><th>DayOfWeek</th><th>Sales</th><th>Customers</th><th>Open</th><th>Promo</th><th>SchoolHoliday</th></tr></thead><tbody><tr><td>count</td><td>1.017209e+06</td><td>1.017209e+06</td><td>1.017209e+06</td><td>1.017209e+06</td><td>1.017209e+06</td><td>1.017209e+06</td><td>1.017209e+06</td></tr><tr><td>mean</td><td>5.584297e+02</td><td>3.998341e+00</td><td>5.773819e+03</td><td>6.331459e+02</td><td>8.301067e-01</td><td>3.815145e-01</td><td>1.786467e-01</td></tr><tr><td>std</td><td>3.219087e+02</td><td>1.997391e+00</td><td>3.849926e+03</td><td>4.644117e+02</td><td>3.755392e-01</td><td>4.857586e-01</td><td>3.830564e-01</td></tr><tr><td>min</td><td>1.000000e+00</td><td>1.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>25%</td><td>2.800000e+02</td><td>2.000000e+00</td><td>3.727000e+03</td><td>4.050000e+02</td><td>1.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>50%</td><td>5.580000e+02</td><td>4.000000e+00</td><td>5.744000e+03</td><td>6.090000e+02</td><td>1.000000e+00</td><td>0.000000e+00</td><td>0.000000e+00</td></tr><tr><td>75%</td><td>8.380000e+02</td><td>6.000000e+00</td><td>7.856000e+03</td><td>8.370000e+02</td><td>1.000000e+00</td><td>1.000000e+00</td><td>0.000000e+00</td></tr><tr><td>max</td><td>1.115000e+03</td><td>7.000000e+00</td><td>4.155100e+04</td><td>7.388000e+03</td><td>1.000000e+00</td><td>1.000000e+00</td><td>1.000000e+00</td></tr></tbody></table>		Store	DayOfWeek	Sales	Customers	Open	Promo	SchoolHoliday	count	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	mean	5.584297e+02	3.998341e+00	5.773819e+03	6.331459e+02	8.301067e-01	3.815145e-01	1.786467e-01	std	3.219087e+02	1.997391e+00	3.849926e+03	4.644117e+02	3.755392e-01	4.857586e-01	3.830564e-01	min	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	25%	2.800000e+02	2.000000e+00	3.727000e+03	4.050000e+02	1.000000e+00	0.000000e+00	0.000000e+00	50%	5.580000e+02	4.000000e+00	5.744000e+03	6.090000e+02	1.000000e+00	0.000000e+00	0.000000e+00	75%	8.380000e+02	6.000000e+00	7.856000e+03	8.370000e+02	1.000000e+00	1.000000e+00	0.000000e+00	max	1.115000e+03	7.000000e+00	4.155100e+04	7.388000e+03	1.000000e+00	1.000000e+00	1.000000e+00
	Store	DayOfWeek	Sales	Customers	Open	Promo	SchoolHoliday																																																																		
count	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06																																																																		
mean	5.584297e+02	3.998341e+00	5.773819e+03	6.331459e+02	8.301067e-01	3.815145e-01	1.786467e-01																																																																		
std	3.219087e+02	1.997391e+00	3.849926e+03	4.644117e+02	3.755392e-01	4.857586e-01	3.830564e-01																																																																		
min	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00																																																																		
25%	2.800000e+02	2.000000e+00	3.727000e+03	4.050000e+02	1.000000e+00	0.000000e+00	0.000000e+00																																																																		
50%	5.580000e+02	4.000000e+00	5.744000e+03	6.090000e+02	1.000000e+00	0.000000e+00	0.000000e+00																																																																		
75%	8.380000e+02	6.000000e+00	7.856000e+03	8.370000e+02	1.000000e+00	1.000000e+00	0.000000e+00																																																																		
max	1.115000e+03	7.000000e+00	4.155100e+04	7.388000e+03	1.000000e+00	1.000000e+00	1.000000e+00																																																																		
<p>Мы проверим все недостающие элементы здесь.</p>	<pre>missing = train.isnull().sum() missing.sort_values(ascending=False)</pre> <pre>Store 0 DayOfWeek 0 Sales 0 Customers 0 Open 0 Promo 0 StateHoliday 0 SchoolHoliday 0 dtype: int64</pre>																																																																								
<p>Затем мы создаем новую метрику, чтобы увидеть средние продажи на одного клиента.</p>	<pre>train['SalesPerCustomer'] = train['Sales'] /train['Customers'] train['SalesPerCustomer'].head()</pre>																																																																								

Мы собираемся проверить, есть ли какие-либо пропущенные значения с нашей новой метрикой, и удалить их. Либо клиенты, либо продажи должны быть равны нулю, чтобы дать нам нулевое значение SalesPerCustomer.

```
missing = train.isnull().sum()
missing.sort_values(ascending=False)
```

```
SalesPerCustomer    172869
Store                0
DayOfWeek           0
Sales               0
Customers           0
Open                0
Promo               0
StateHoliday        0
SchoolHoliday       0
dtype: int64
```

```
train.dropna().head()
```

	Store	DayOfWeek	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday	SalesPerCustomer
Date									
2015-07-31	1	5	5263	555	1	1	0	1	9.482883
2015-07-31	2	5	6064	625	1	1	0	1	9.702400
2015-07-31	3	5	8314	821	1	1	0	1	10.126675
2015-07-31	4	5	13995	1498	1	1	0	1	9.342457
2015-07-31	5	5	4822	559	1	1	0	1	8.626118

Исследовательский анализ данных (сохранение данных)
Мы делаем то же самое, что и для нашего тренировочного набора. Изучая данные, мы видим, что в CompetitionDistance отсутствуют только 3 значения. Поскольку это такое небольшое количество, мы решили заменить их средним значением столбца. Остальные отсутствующие значения зависят от Promo2. Поскольку эти отсутствующие значения связаны с тем, что Promo2 равен 0, мы можем заменить эти нули на 0.

```
store.head()
```

	Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear	PromoInterval
0	1	c	a	1270.0	9.0	2008.0	0	NaN	NaN	NaN
1	2	a	a	570.0	11.0	2007.0	1	13.0	2010.0	Jan-Apr-Jul-Oct
2	3	a	a	14130.0	12.0	2005.0	1	14.0	2011.0	Jan-Apr-Jul-Oct
3	4	c	c	620.0	9.0	2009.0	0	NaN	NaN	NaN
4	5	a	a	29910.0	4.0	2015.0	0	NaN	NaN	NaN

```
store.shape
```

```
(1115, 10)
```

```
store.isnull().sum()
```

```
Store                0
StoreType            0
Assortment           0
CompetitionDistance    3
CompetitionOpenSinceMonth  354
CompetitionOpenSinceYear  354
Promo2               0
Promo2SinceWeek      544
Promo2SinceYear      544
PromoInterval        544
dtype: int64
```

Так как из этого пропущено только 3 значения, мы заполняем средним значением из столбца

```
store['CompetitionDistance'].fillna(store['CompetitionDistance'].mean(), inplace = True)
```

Строки, в которых нет Promo2, мы можем заполнить остальные значения 0	<pre>store.fillna(0, inplace = True) store.head()</pre> <table><thead><tr><th></th><th>Store</th><th>StoreType</th><th>Assortment</th><th>CompetitionDistance</th><th>CompetitionOpenSinceMonth</th><th>CompetitionOpenSinceYear</th><th>Promo2</th><th>Promo2SinceWeek</th><th>Promo2SinceYear</th><th>PromoInterval</th></tr></thead><tbody><tr><td>0</td><td>1</td><td>c</td><td>a</td><td>1270.0</td><td>9.0</td><td>2009.0</td><td>0</td><td>0.0</td><td>0.0</td><td>0</td></tr><tr><td>1</td><td>2</td><td>a</td><td>a</td><td>570.0</td><td>11.0</td><td>2007.0</td><td>1</td><td>13.0</td><td>2010.0</td><td>Jan-Apr-Jul-Oct</td></tr><tr><td>2</td><td>3</td><td>a</td><td>a</td><td>14130.0</td><td>12.0</td><td>2006.0</td><td>1</td><td>14.0</td><td>2011.0</td><td>Jan-Apr-Jul-Oct</td></tr><tr><td>3</td><td>4</td><td>c</td><td>c</td><td>620.0</td><td>9.0</td><td>2009.0</td><td>0</td><td>0.0</td><td>0.0</td><td>0</td></tr><tr><td>4</td><td>5</td><td>a</td><td>a</td><td>28910.0</td><td>4.0</td><td>2015.0</td><td>0</td><td>0.0</td><td>0.0</td><td>0</td></tr></tbody></table>		Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear	PromoInterval	0	1	c	a	1270.0	9.0	2009.0	0	0.0	0.0	0	1	2	a	a	570.0	11.0	2007.0	1	13.0	2010.0	Jan-Apr-Jul-Oct	2	3	a	a	14130.0	12.0	2006.0	1	14.0	2011.0	Jan-Apr-Jul-Oct	3	4	c	c	620.0	9.0	2009.0	0	0.0	0.0	0	4	5	a	a	28910.0	4.0	2015.0	0	0.0	0.0	0
	Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear	PromoInterval																																																									
0	1	c	a	1270.0	9.0	2009.0	0	0.0	0.0	0																																																									
1	2	a	a	570.0	11.0	2007.0	1	13.0	2010.0	Jan-Apr-Jul-Oct																																																									
2	3	a	a	14130.0	12.0	2006.0	1	14.0	2011.0	Jan-Apr-Jul-Oct																																																									
3	4	c	c	620.0	9.0	2009.0	0	0.0	0.0	0																																																									
4	5	a	a	28910.0	4.0	2015.0	0	0.0	0.0	0																																																									
Строки, в которых нет Promo2, мы можем заполнить остальные значения 0	<pre>train = train.merge(right=store, on='Store', how='left')</pre>																																																																		
Модель скользящего среднего (Наивная модель)																																																																			
Мы собираемся использовать модель скользящего среднего для прогнозирования акций GM для нашей базовой модели. Модель скользящего среднего будет использовать среднее значение различных «окон» времени, чтобы составить свой прогноз.	<pre>train = pd.read_csv("https://raw.githubusercontent.com/RPI-DATA/tutorials-intro/master/rossmann-store-sales/rossmann-store-sales/train.csv", parse_dates = True, low_memory = False, index_col = 'Date') train = train.sort_index(ascending = True) train['SalesPerCustomer'] = train['Sales']/train['Customers'] train['SalesPerCustomer'].head()</pre> <div><div></div><div><div>Date</div><div>2013-01-01</div><div>3.828049</div></div><div><div>2013-01-01</div><div>4.981328</div></div><div><div>2013-01-01</div><div>4.233600</div></div><div><div>2013-01-01</div><div>5.907021</div></div><div><div>2013-01-01</div><div>5.464286</div></div><div>Name: SalesPerCustomer, dtype: float64</div></div>																																																																		
Мы перезагружаем данные, потому что теперь у нас есть представление о том, как мы хотим манипулировать ими для нашей модели. Прделав те же манипуляции с данными, что и раньше, мы начинаем смотреть на тенденцию наших продаж.	<pre>train = train.dropna()</pre>																																																																		
Здесь мы просто строим график продаж, которые у нас есть. Как видите, количество продаж настолько велико, что наш график просто выглядит как голубая заливка. Однако мы можем понять, как распределяются наши продажи.	<pre>plt.plot(train.index, train['Sales']) plt.title('Rossmann Sales') plt.ylabel('Price (\$)'); plt.show()</pre> <div><div></div><div><div>Rossmann Sales</div><div><div>40000</div><div>30000</div><div>20000</div><div>10000</div><div>0</div></div><div><div>2013-01-01</div><div>2013-07-01</div><div>2014-01-01</div><div>2014-07-01</div><div>2015-01-01</div><div>2015-09-01</div></div></div></div>																																																																		
Чтобы очистить наш график, мы хотим сформировать новый столбец, учитывающий только год продаж.	<pre>train['Year'] = train.index.year # Take Dates from index and move to Date column train.reset_index(level=0, inplace = True) train['sales'] = 0</pre>																																																																		
Разделите данные на набор решений и тестов. Мы используем	<pre>train_store=train[0:675472] test_store=train[675472:]</pre>																																																																		

разделение 80/20. Затем мы смотрим, чтобы начать модели.

test_store означает часть прогнозирования.

```
train_store.Date = pd.to_datetime(train_store.Date, format="%Y-%m-%d")
train_store.index = train_store.Date
test_store.Date = pd.to_datetime(test_store.Date, format="%Y-%m-%d")
test_store.index = test_store.Date
```

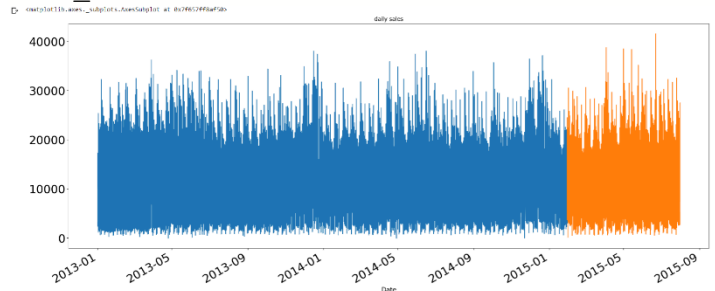
```
train_store = train_store.resample('D').mean()
train_store = train_store.interpolate(method='linear')
```

```
test_store = test_store.resample('D').mean()
test_store = test_store.interpolate(method='linear')
```

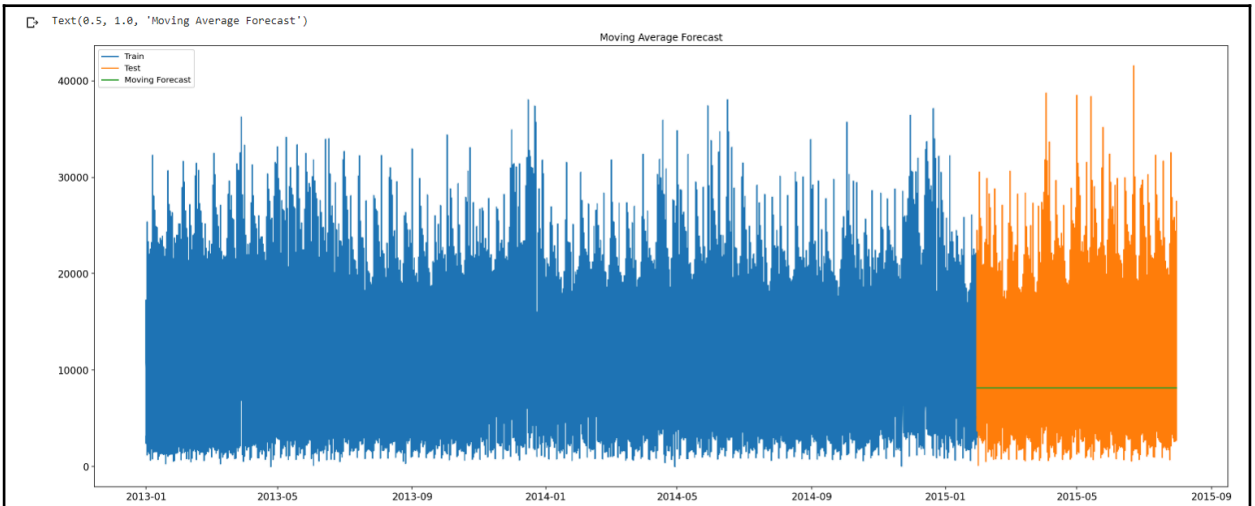
/Users/josefigueroa/anaconda3/lib/python3.7/site-packages/pandas/core/generic.py:4405: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>
self[name] = value

```
train_store.Sales.plot(figsize=(25,10), title='daily sales', fontsize=25)
test_store.Sales.plot()
```



```
y_hat_avg_moving = test_store.copy()
y_hat_avg_moving['moving_avg_forecast'] = train_store.Sales.rolling(90).mean().iloc[-1]
plt.figure(figsize=(25,10))
plt.plot(train_store['Sales'], label='Train')
plt.plot(test_store['Sales'], label='Test')
plt.plot(y_hat_avg_moving['moving_avg_forecast'], label='Moving Forecast')
plt.legend(loc='best')
plt.title('Moving Average Forecast')
```



```
rms_avg_rolling = sqrt(mean_squared_error(
    test_store.Sales, y_hat_avg_moving.moving_
    avg_forecast))
print('ROLLING AVERAGE',rms_avg_rolling)
```

ROLLING AVERAGE 3233.627446795865

Скользящее среднее для нашей модели составляет 1915,88. Этот прогноз кажется очень последовательным в достижении среднего значения будущих продаж. Эта наивная модель определенно выглядит солидной моделью, однако она не самая лучшая.

