

ПРОГНОЗИРОВАНИЕ ПРОДАЖ В РИТЕЙЛЕ НА ОСНОВЕ МЕТОДОВ МАШИННОГО ОБУЧЕНИЯ

Е. А. Артемьев, email: artemev.15.01.1998@list.ru, В.В. Мокшин

Казанский национальный исследовательский технический университет

***Аннотация.** В современном мире ритейла недостаточно завозить продукцию в торговые точки, выставлять её на стеллажи. Отчётности от мерчендайзеров или по продажам за месяц стали хорошим инструментом анализа, но это требует времени и больших человеческих ресурсов. Эффективнее использовать машинное обучение для прогнозирования продаж на месяцы вперёд, позволяя заранее принимать решения по кол-ву, типам и ценам товаров в торговых сетях. В данной работе вы узнаете, как реализовать модель машинного обучения, которая может прогнозировать продажи в зависимости от определенных характеристик, таких как день продажи, торговая точка, стоимость, кол-во проданных товаров за предыдущие дни и категория товара.*

***Ключевые слова:** временное прогнозирование, регрессионный анализ, корреляционный анализ, регрессия, продажи, ритейл.*

Введение

Ритейл – это механизм продажи товаров и услуг, запутанная и сложная сеть из покупателей и продавцов. Продавцы закупают товары в определённом кол-ве оптом и продают по определённой цене на различных площадках и в различных розничных сетях.

Покупатели берут выложенный товар исходя из своих мотивов. Собрав датасет из проданных товаров, можно понять мотивы людей, месяцы высокого спроса и малого, интересные продукты и не особо в определённые периоды.

Наша цель в данной статье — изучить прогнозирование будущих продаж основе прошлых продаж. Эти данные поступают с веб-сайта Kaggle, на котором собраны датасеты различных компаний. Важно иметь огромный массив данных для того, чтобы модель могла вычленивать зависимости между данными.

1. Методология

Данные были получены с веб-сайта Kaggle; набор данных был озаглавлен “Predict Future Sales”. Перед началом работы с данными

были выдвинуты несколько различных гипотез относительно того, какие переменные будут иметь большое значение для увеличения продаж. Наиболее интересными для обучения признаками были выделены цена, дата, торговая точка, сам товар. После первых экспериментов стало ясно, что этих параметров недостаточно и их использование в том виде, в котором они представлены в датасете неудачно и требует преобразования.

Данные были разделены на обучающуюся и тестовую выборки в отношении 80:20 и очищены путём удаления дубликатов, выбросов, строк с пустыми значениями, сортировки по дате, замены дат и строковых значений числовыми для более эффективного обучения. Затем был проведён корреляционный анализ, чтобы определить ключевые факторы, которые имеют тесную связь с предсказанием поля, отвечающего за кол-во проданных товаров в определённый промежуток времени. Связь между данными была установлена крайне незначительной, для чего было принято перегруппировать данные, но после проведения экспериментов с линейной регрессией. Алгоритм линейной регрессии сработал неудачно, но дал нам важную информацию насчёт влияния параметров на обучение: id товара, id торговой точки и месяц продаж. Затем заменили алгоритм на метод ансамблевого обучения Extreme Gradient Boosting Regressor для адаптивного изменения распределения обучающих данных, уделяя больше внимания ранее неправильно классифицированным записям для создания базовых учащихся. Алгоритм показал себя лучше, указав нам на то, что всё те же поля оказывают наибольшее влияние на обучение и на то, что можно перегруппировать данные и добавить новые признаки.

Для оценки точности алгоритма были использованы такие метрики как MAE, MSE, RMSE. MAE – В статистике, средняя абсолютная ошибка является мерой ошибок между парными наблюдениями, выражающих то же явление. Вычисляется по формуле:

$$MAE = \sum_{i=1}^n |y_i - x_i| / n = \sum_{i=1}^n |e_i| / n ,$$

- n – число прогнозов;
- y – спрогнозированное значение;
- x – изначальные значения.

Средняя абсолютная ошибка известна как мера точности, зависящая от масштаба и имеет вид:

$$MSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

- n – число прогнозов;
- y – изначальные данные;
- \hat{y} – предсказанное значение.

RMSE (Root Mean Square Error) - это квадратный корень из значения, полученного с помощью функции среднеквадратической ошибки.

$$RMSE = \sqrt{\sum_{i=1}^n (\text{Pr edicted } _i - \text{Actual } _i)^2 / N},$$

- Pr edicted – число прогнозов;
- Actual – изначальные данные;
- N - предсказанное значение.

2. Анализ

Поскольку мы пробуем предсказать кол-во проданных товаров в определённый промежуток времени, мы использовали регрессионный анализ. Регрессионный анализ — набор статистических методов исследования влияния одной или нескольких независимых переменных на зависимую.

В исследовании были некоторые ограничения. Одним из ограничений являются выбросы и определение переменных, необходимых для обучения. Для выявления второго был использован корреляционный анализ. По полученным данным видно, что связь крайне несущественная, поэтому были выбраны все признаки для первых экспериментов (рис. 1).

	date_block_num	shop_id	item_id	item_price	item_cnt_day
date_block_num	1.000000	0.019273	0.009356	0.095010	0.009402
shop_id	0.019273	1.000000	0.029396	-0.024034	-0.005230
item_id	0.009356	0.029396	1.000000	-0.134104	0.016650
item_price	0.095010	-0.024034	-0.134104	1.000000	0.011197
item_cnt_day	0.009402	-0.005230	0.016650	0.011197	1.000000

Рис. 1. Корреляционный анализ

После пропуска датасета по линейной регрессии и методу ансамблевого обучения стало ясно, что сам товар и торговая точка влияют на обучение в большей степени (рис. 2). Но результаты обучения были не оптимистичными, поэтому было принято решение сгруппировать данные по-другому и связь между признаками стала выше (рис. 3) [1].

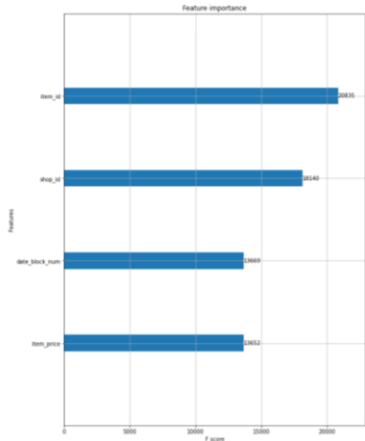


Рис. 2. Параметры, влияющие на обучение

Index	date_block_num	shop_id	item_id	item_cnt_month	shop_category	shop_city	item_category_id
date_block_num	1.0	0.02261958307597148	-0.012621322286894718	-0.00062533539953814	-0.332121547246299	0.0212694803101459403	0.01676265619110706
shop_id	0.02261958307597148	1.0	0.0034466411615666085	0.008252940600719245	0.100518198969574335	0.9861300305824917	-0.0001964694929621935
item_id	-0.012621322286894718	-0.0034466411615666085	1.0	-0.012196804917996203	0.0003326483238240846	-0.00029035849424420363	0.12914878416208933
item_cnt_month	-0.00062533539953814	0.000252940600719245	-0.012196804917996203	1.0	-0.0916701449393382713	0.007829397481211534	-0.2526464791616554
shop_category	0.022121547246299	0.100518198969574335	0.0003326483238240846	-0.012196804917996203	1.0	0.0548622124064911	-0.009932616277374369
shop_city	0.0212694803101459403	0.9861300305824917	-0.00029035849424420363	0.007829397481211534	0.0548622124064911	1.0	0.0001380460216307458
item_category_id	0.01676265619110706	-0.0001964694929621935	0.12914878416208933	-0.2526464791616554	-0.0009358160273374369	-0.0001380460216307458	1.0
name3	0.0201273224181722	0.000856809732270802	0.02771195781094036	0.02771195781094036	0.0005746284720847289	0.000754762816420836	-0.29180512229191184
name3	-0.0798489880103076	-0.08162520111313010	-0.1103222780129966	-0.0177996598489073	0.00263993529297915	-0.016116376123119679	-0.14850143551306
subtype_code	0.10048423034516964	0.002112568279812416	-0.1324324818627913	0.02507496748677888	0.003472492146084114	0.001801801562641998	0.2614970566212584
type_code	0.00644614619343057	-0.0004401424940894664	0.1478840889300076	0.03646398124167445	0.0006206208640943759	-0.000359947064865526	0.9544271077118361
item_cnt_month_lag_1	0.000136088693748072	0.01095337236898129	-0.01599786918970672	0.726736094107479	3.13486801494328666-05	0.009468208436412323	-0.2907494221621217
item_cnt_month_lag_2	2.615611958921048-05	0.01000977326898129	-0.0174858463682894	0.64490571204161	0.010881232329168735	0.011011737634316502	-0.03198922665962825
item_cnt_month_lag_3	0.01259878787080271	0.0113883613480571	-0.01823543894583365	0.026796801874303	0.06103258774209298	0.01089112388411556	-0.03310938422162021
date_avg_item_cnt_lag_1	-0.00726427062110476	0.00846534856321224	-0.0412370616004984	0.00780278330489985	0.002682062804892108	-0.00273903768740848	0.001917018656600112
date_avg_item_cnt_lag_2	4.741952781516818-05	0.00479156295376-05	-0.02492320732746858	0.518149718737881	4.232226595194886-05	9.02415821709949-05	0.04033048300840784
date_avg_item_cnt_lag_3	-0.0001750862588427927	2.16049475287273-05	-0.02608862394310856	0.4688146499633783	0.0015296709520284982	3.21016052709949-05	-0.4807873394883184
item_avg_item_cnt_lag_1	0.0020318802236016727	0.036313673688085-05	-0.028142740768815154	0.4629212932976785	0.05046014828325-05	0.00011246739985893735	-0.0064253488116713
date_shop_avg_item_cnt_lag_1	0.01240888638496223	0.12209618720192943	-0.008724885138118482	0.0707490381656647	-0.0416276227631634	0.11570397902810378	0.00020390714994208
date_shop_avg_item_cnt_lag_2	-0.019334464504894844	0.0378451812842863	0.38129376972816-05	0.000742376749866	-0.0018685881265664	0.12424966018281582	0.00078165161843807
date_shop_avg_item_cnt_lag_3	-0.02210200381443537	0.131442433828213	0.0095928491831545427	0.9504278631988677	0.00199762931911137	0.12467456842213091	0.0006191395285113
date_shop_avg_item_cnt_lag_1	0.000136088693748072	0.01095337236898129	-0.01599786918970672	0.726736094107479	3.13486801494328666-05	0.009468208436412323	-0.2907494221621217
date_shop_avg_item_cnt_lag_2	2.615611958921048-05	0.01000977326898129	-0.0174858463682894	0.64490571204161	0.010881232329168735	0.011011737634316502	-0.03198922665962825
date_shop_avg_item_cnt_lag_3	0.01259878787080271	0.0113883613480571	-0.01823543894583365	0.026796801874303	0.06103258774209298	0.01089112388411556	-0.03310938422162021
date_shop_avg_item_cnt_lag_1	0.000787711172548688	0.010464632932289	-0.01148016895783841	0.733038178732616	6.99944184047463-07	0.00973051301284264	-0.0382129912942117
date_city_avg_item_cnt_lag_1	-0.0190829857399965	0.1392488324232638	-0.01188274668325253	0.9482937378776141	-0.131450059178813728	0.1638789978024001	0.0002770783864716277
date_item_city_avg_item_cnt_lag_1	0.0170501619295376-05	0.003524848899654153	-0.01033008992841172	0.005647292341435	-0.00009951991827072	0.0110041175866665	-0.0336361088925257
delta_price_lag	-0.1461313818613756	-0.0052395971154854	-0.00743519042829805	0.0096779939356483	0.00372729315198479	-0.00505263984823604	0.04132155815103752
delta_revenue_lag_1	-0.013815888668007	0.0081657043795426	-0.00101379348448323	0.0610378362498194	-0.0277868016574268	-0.041994372919969	0.00026236388417817
month	0.02325978058927224	-0.0064650221597759	-0.00333599789586662	0.0048498542214550785	-0.00089358160722846	-0.0005958161824296	0.1466336123244866
days	0.0071712470998804	-0.0030683188986843	0.00739219543714814	0.0028125648498199	-0.0059415297398865	-0.00280413544137487	0.128991334668216-05
item_shop_first_sale	-0.001805062916656	-0.0180944643446689	-0.00863745669777614	-0.02763036396092964	0.009903686298027	-0.01533642975286847	0.00381941978883615
item_first_sale	0.481930682920285	0.01227657748814921	-0.00933461005769093	-0.04131968078102394	-0.04162326354424303	0.01180095907999979	-0.031647762480271

Рис. 3. Корреляционный анализ преобразованного датасета

3. Прогнозирование

Мы пробуем предсказать кол-во проданных товаров в определённый промежуток времени. Для этого мы используем несколько алгоритмов обучения модели: «Линейная регрессия», «eXtreme Gradient Boosting Regressor».

Применив линейную регрессию, получаем неутешительные результаты. Мы видим такие метрики как MAE, MSE, RMSE и таблицу, по которой можем увидеть, какие столбцы в какой степени влияют на обучение (рис. 4).

```
☞ mse: 18.373, mae: 0.523
Mean Absolute Error: 0.5229426436550997
Mean Squared Error: 18.3727463847258
Root Mean Squared Error: 4.286344174786458
```

Coefficient	
date_block_num	0.004828
shop_id	-0.000826
item_id	0.000008
item_price	0.000021

Рис. 4. Результаты линейной регрессии

Визуализировав данные для ТТ с id=2 в первый месяц 2015 года, видим, насколько плохо справилась модель, если смотреть в разрезе ТТ и периода (рис. 5). Дать советы магазину по такому результату не получится.

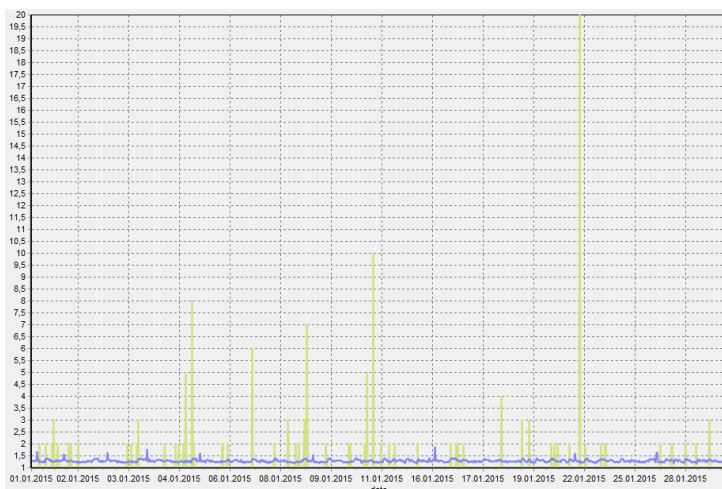


Рис. 5. График полученных и исходных данных

EXtreme Gradient Boosting Regressor - это метод ансамблевого обучения, который позволит создать окончательную модель путем объединения нескольких слабых моделей. Как описано в названии машины для повышения градиента, используется метод градиентного спуска и повышения. Метод Boosting использует итеративную процедуру для адаптивного изменения распределения обучающих данных, уделяя больше внимания ранее неправильно классифицированным записям для создания базовых учащихся. Это можно рассматривать как последовательное добавление моделей до тех пор, пока не прекратятся дальнейшие улучшения. Для минимизации потерь используется алгоритм градиентного спуска. При повышении градиента, когда прогнозы нескольких моделей комбинируются, градиент используется для оптимизации прогнозирования усиленной модели в каждом раунде повышения.

Мы также взяли отсортированный по дате датасет и разделили его на обучающую и тестовую выборки. И подали его на вход модели без ухищрений. Получаем результаты, которые выглядят лучше, чем у линейной регрессии, но недостаточно хороши, чтобы использовать в реальной жизни.

Результаты обучения выглядят следующим образом (рис. 6):

```
– [1] validation_0-rmse:2.09102 validation_1-
rmse:2.09102
```

- [2] validation_0-rmse:2.05408 validation_1-rmse:2.05408
- [3] validation_0-rmse:2.03131 validation_1-rmse:2.03131
- [4] validation_0-rmse:2.00578 validation_1-rmse:2.00578
- [5] validation_0-rmse:1.98453 validation_1-rmse:1.98453
- [6] validation_0-rmse:1.96968 validation_1-rmse:1.96968
- ...
- [998] validation_0-rmse:1.70407 validation_1-rmse:1.70407
- [999] validation_0-rmse:1.70392 validation_1-rmse:1.70392

Mean Absolute Error: 0.45253064959522105
 Mean Squared Error: 16.83557398677999
 Root Mean Squared Error: 4.103117593584175

	Actual	Predicted
2348679	1.0	1.464342
2348680	2.0	1.464342
2348681	1.0	1.464342
2348682	1.0	1.464342
2348683	1.0	2.632516
...
2935844	1.0	1.029175
2935845	1.0	1.029175
2935846	1.0	1.107367
2935847	1.0	1.029175
2935848	1.0	1.029175

587170 rows x 2 columns

Рис. 6. Результаты EXtreme Gradient Boosting Regressor

Feature_importances позволяет оценить, какие признаки оказывали существенное влияние на обучение (рис. 7).

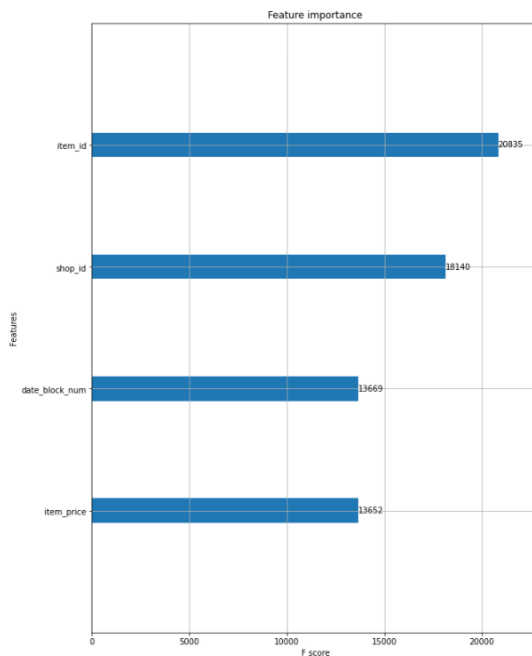


Рис. 7. Feature_importances. EXtreme Gradient Boosting Regressor

Попробуем уменьшить ошибку, доработав алгоритм выше путём группировки данных и добавлением новых признаков. Результаты обучения (рис. 8):

- [0] validation_0-rmse:1.18496 validation_1-rmse:1.11945
- [1] validation_0-rmse:1.12885 validation_1-rmse:1.07960
- [2] validation_0-rmse:1.08653 validation_1-rmse:1.04936
- [3] validation_0-rmse:1.04426 validation_1-rmse:1.02158
- [4] validation_0-rmse:1.01560 validation_1-rmse:0.99966
- [5] validation_0-rmse:0.98959 validation_1-rmse:0.98105
- ...

- [67] validation_0-rmse:0.76339 validation_1-rmse:0.91226
- [68] validation_0-rmse:0.76288 validation_1-rmse:0.91230
- [69] validation_0-rmse:0.76230 validation_1-rmse:0.91229
- [70] validation_0-rmse:0.76176 validation_1-rmse:0.91222

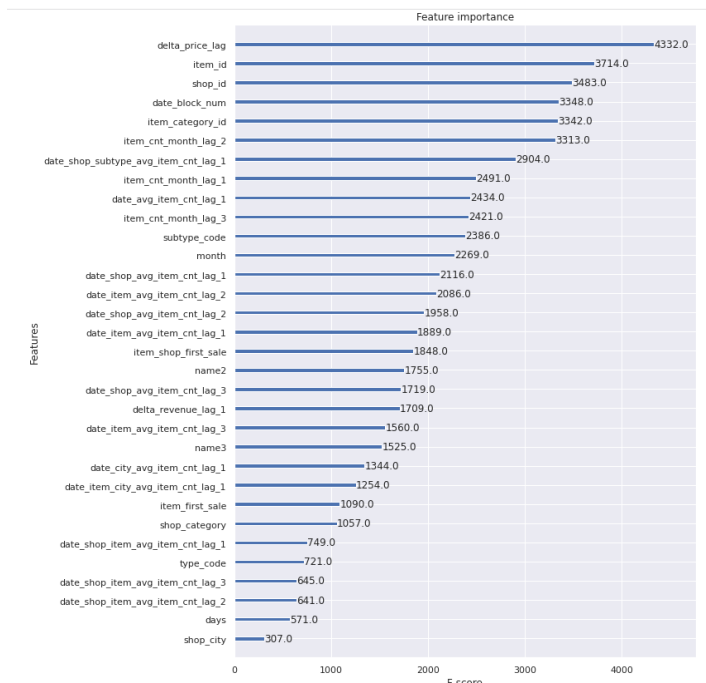


Рис. 8. Feature_importances. EXtreme Gradient Boosting Regressor

Посмотрим, насколько хорошо предсказала модель кол-во проданных товаров за январь 2015 года в ТТ с id=50.

По диаграмме (рис. 9) видно, что результаты несильно отличаются, но опять же стоит отметить, что результаты далеки от совершенных и предсказывать по такой модели опасно.

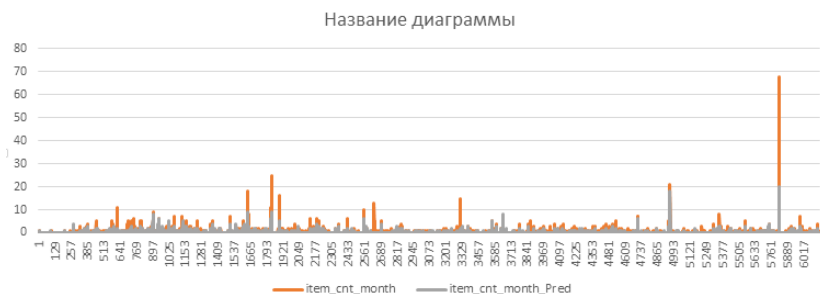


Рис. 9. График исходных и полученных данных EXtreme Gradient Boosting Regressor

Заключение

В данном исследовании представлен подход, основанный на машинном обучении, для оценки продаж на будущие месяцы по данным прошлых месяцев. Полученная модель требует улучшений и, возможно, замены на более сложные модели с обширным функционалом, но положительный результат использования машинного обучения в ритейле виден из графиков выше.

Большой набор данных был обработан, прежде чем передать его предлагаемым классификаторам и спрогнозировать кол-во проданных товаров. Необходимо для каждой задачи правильно сгруппировать данные и выделить необходимые признаки. Оценка основана на сравнении результатов предсказания и реальных значений продаж, производительности модели были оценены с помощью метрик MAE, MSE, RMSE, что показало, насколько модель приближена к реальной. По полученным данным, можно сказать, что разработанный алгоритм имеет место быть для грубого предсказания продаж как подспорье для дальнейшего анализа и выводов. Также результаты указывают на то, что анализ данных по продажам – это непростая задача с множеством подводных камней.

Литература

1. Мокшин, В. В. Рекурсивный алгоритм построения регрессионных моделей сложных вероятностных объектов / В. В. Мокшин, И. Р. Сайфулинов, А. П. Кирпичников // Вестник Технологического университета. – 2017. – № 9. – С. 112-116.
2. Мокшин, В. В. Метод формирования модели анализа сложной системы / В. В. Мокшин, И. М. Якимов // Информационные технологии. – 2011. – № 5. – С. 46-51.