



МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ
(национальный исследовательский университет)»

ОТЧЕТ ПО КУРСОВОЙ РАБОТЕ

по дисциплине «Статистические методы обработки данных»

Тема: «Методы кодирования категориальных признаков в
задаче бинарной классификации»

Выполнил студент группы М80-101М-21

Фейзуллин К.М.

Руководитель: канд. физ.-мат. наук,
доцент кафедры 804 Платонов Е.Н.

Оценка:

Оглавление

Введение	3
Постановка задачи	4
Описание эксперимента	5
Функционал качества	5
Описание исходных данных	8
Базовая модель	10
Экспериментальная установка	11
Методы кодирования	13
Результаты исследования	15
Заключение	16
Список литературы	17
Приложение	18

Введение

Целью работы является исследование влияния методов кодирования категориальных признаков в задаче бинарной классификации на функционал качества работы прогнозирующей математической модели.

В эру больших данных и расцвета автоматизированных алгоритмов обработки информации и прогнозирования важную роль занимают сами методы обработки. Помимо целочисленных и вещественных признаков нередко приходится сталкиваться и с категориальными данными, которые требуют особый подход для дальнейшей работы с ними.

Актуальность данного исследования состоит в том, что до сих пор точно неизвестно, для какого рода задач и для каких сущностей переменных то или иное кодирование покажет наилучший результат.

Причем, данная задача бинарной классификации нередко используется для решения разного типа проблем, будь то прогноз подорожания или роста акции, решение о предоставлении кредита или ипотеки, определение вредоносности программного обеспечения по показателям системы. Как можно понять, данная задача имеет множество приложений и в том числе выгодных с экономической точки зрения.

Но очень важно учитывать, что качество работы математической модели имеет зависимость с качеством данных, которые обрабатываются алгоритмом прогноза.

Поэтому нахождение оптимальных методов кодирования напрямую скажется на качестве работы прогнозирующей модели.

Постановка задачи

Целью задачи является исследование методов кодирования категориальных признаков в задаче бинарной классификации. Функционалом качества работы алгоритма бинарной классификации является метрика ROC AUC (Receiver Operating Characteristic Area Under Curve), которая представляет из себя интегральную сумму от пороговой функции отношения доли верно определенных положительных меток и доли ложно определенных положительно меток.

Необходимо используя различные способы кодирования категориальных признаков и через создание новых признаков добиться оптимального результата.

Все варианты кодирования (различные наборы признаков) и соответствующую точность решения задачи будут занесены в сводную таблицу.

Описание эксперимента

Функционал качества

Как уже было сказано выше, положительное влияние тех или иных методов кодирования будет определяться с помощью функционала качества ROC AUC[1]. Чтобы определить ее введем карту обозначений, приведенную в таблице 1.

		Прогноз	
		1	0
Правильный ответ	1	TP (True – positives)	FN (False – negatives)
	0	FP (False – positives)	TN (True – negatives)

Табл.1 – карта обозначений.

Определим обозначения, в нашем случае, как TP (True – positives) – число раз, когда объект положительно определен верно, FP (False – positives) – когда объект положительно определен ложно, FN (False – negatives) – когда объект отрицательно определен ложно и TN (True – negatives) – когда объект отрицательно определен верно.

Из этих характеристик определяются TPR (True positive rate) и FPR (False positive rate), доля истинно положительных ответов и доля ложно положительных ответов соответственно:

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

Далее вносит свой вклад природа бинарного классификатора. Простейший линейный классификатор имеет вид:

$$Y = F(\alpha * X + \beta),$$

Где X – входной сигнал, Y – выходной сигнал, $F(*)$ – некая пороговая функций, α – масштабирующий коэффициент, а β – пороговое значение классификатора.

Последовательно варьируя пороговое значение β и находя в каждом

значении β значение TPR и FPR будет определена пороговая функция ROC вида, как на рисунке 1:

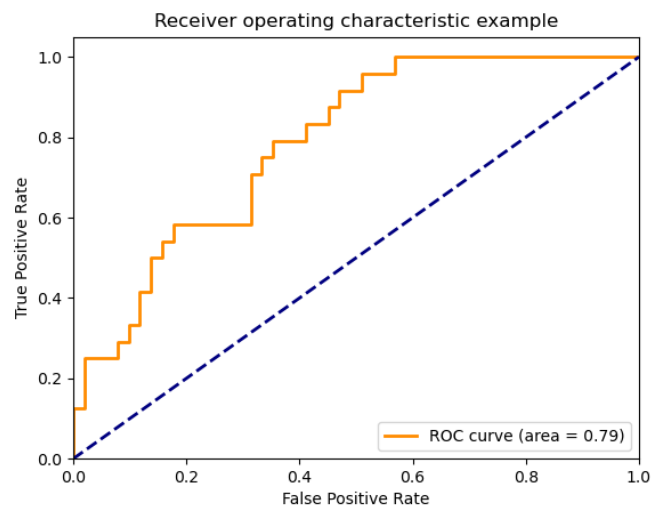


Рисунок 1. График пороговой функции ROC

Как видно, данная функция будет стремиться к левому верхнему углу и ROC AUC (площадь под данной кривой) будет стремиться к единице. То есть, чем больше площадь под кривой AUC, тем больше свободы в выборе порогового значения для наилучшей работы алгоритма классификации, ведь нам желательно максимизировать TPR и минимизировать FPR.

Предельный вид кривой ROC на рисунке 2:

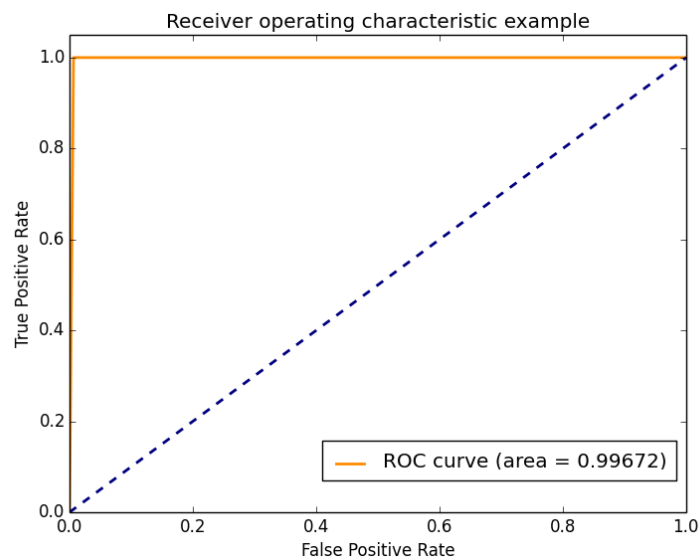


Рисунок 2. Предельный вид пороговой функции ROC

В случае классификатора, имеющий функцию ROC как на рисунке 2, можно выбрать порог классификатора так, чтобы TPR был равен единице, что значит все положительные метки определены верно, а показатель FPR стремится к нулю, что практически значит отсутствие ложно положительных прогнозов.

Добавим, что в сравнительном исследовании будет приниматься значение метрики, полученной на тестовой выборке.

Описание исходных данных

Для исследования были взяты анонимизированные признаки. Набор данных содержит только категориальные функции и включает:

- бинарные признаки
- номинальные признаки с низкой и высокой частотой категорий
- порядковые признаки с низкой и высокой частотой категорий
- (потенциально) циклические признаки

Исходные данные составляют не только численные данные, но и строковые. Пример исходных данных на рисунке 3.

	id	bin_0	bin_1	bin_2	bin_3	bin_4	nom_0	nom_1	nom_2	nom_3	nom_4	nom_5	nom_6	nom_7	nom_8	nom_9	ord_0	ord_1	ord_2	ord_3	ord_4	ord_5	day	month	target
0	0	0.0	0.0	0.0	F	N	Red	Trapezoid	Hamster	Russia	Bassoon	de4c57ee2	a64bc7ddf	598080a91	0256c7a4b	02e7c8990	3.0	Contributor	Hot	c	U	Pw	6.0	3.0	0
1	1	1.0	1.0	0.0	F	Y	Red	Star	Axolotl	NaN	Theremin	2bb3c3e5c	3a3a936e8	1dddb8473	52ead350c	f37df64af	3.0	Grandmaster	Warm	e	X	pE	7.0	7.0	0
2	2	0.0	1.0	0.0	F	N	Red	NaN	Hamster	Canada	Bassoon	b574c9841	708248125	5ddc9a726	745b909d1	NaN	3.0	NaN	Freezing	n	P	eN	5.0	9.0	0
3	3	NaN	0.0	0.0	F	N	Red	Circle	Hamster	Finland	Theremin	673bdf1f6	23edb8da3	3a33ef960	bdaa56dd1	f9d456e57	1.0	Novice	Lava Hot	a	C	NaN	3.0	3.0	0
4	4	0.0	NaN	0.0	T	N	Red	Triangle	Hamster	Costa Rica	NaN	777d1ac2c	3a7975e46	bc9cc2a94	NaN	c5361037c	3.0	Grandmaster	Cold	h	C	OZ	5.0	12.0	0
...
599995	599995	0.0	1.0	0.0	T	N	Red	Polygon	Axolotl	India	Theremin	0147770c0	da5014b01	a7059911d	158183c63	015c63324	3.0	Novice	Freezing	a	R	GZ	5.0	NaN	0
599996	599996	1.0	0.0	0.0	T	Y	Blue	Polygon	Dog	Costa Rica	Oboe	NaN	2023ed4ed	83bdea3a5	e9fde8fa8	a02ae6a63	2.0	Novice	Boiling Hot	n	N	sf	NaN	3.0	0
599997	599997	0.0	0.0	0.0	F	Y	Red	Circle	Axolotl	Russia	Theremin	c7dc5d460	5d7d341ac	114b1dbf3	cccbca824	40f9610c1	2.0	Contributor	Freezing	n	H	MV	7.0	5.0	0
599998	599998	1.0	1.0	0.0	F	Y	NaN	Polygon	Axolotl	NaN	Piano	4d7780407	209e1054e	fba315672	4164322bd	c1a8374a0	1.0	Master	Warm	m	X	Ey	1.0	5.0	0
599999	599999	0.0	0.0	0.0	T	N	Blue	Triangle	Dog	Russia	Theremin	17f06e644	f7706ca16	67a8d4ebb	NaN	e2aea7784	1.0	Contributor	Boiling Hot	b	O	ul	5.0	8.0	0

600000 rows × 25 columns

Рисунок 3. Исходные данные

Помимо признаков переменных данные содержат и ответ (target). В каждой строке определены значения переменных и соответствующий им ответ.

Стоит добавить, что каждая переменная имеет потери данных 3-4% и при исключении всех строк с потерями, количество данных сокращается вдвое.

Так же важным фактором является то, что присутствует сильное преобладание отрицательных ответов над положительными, диспропорция составляет 5:1. Но так как в данной работе функционалом качества определена площадь под кривой AUC, которая не чувствительна к такому роду проблем, никаких действий можно не предпринимать.

Если слепо факторизовать каждую переменную и построить корреляционную матрицу, то получим следующее на рисунке 4.

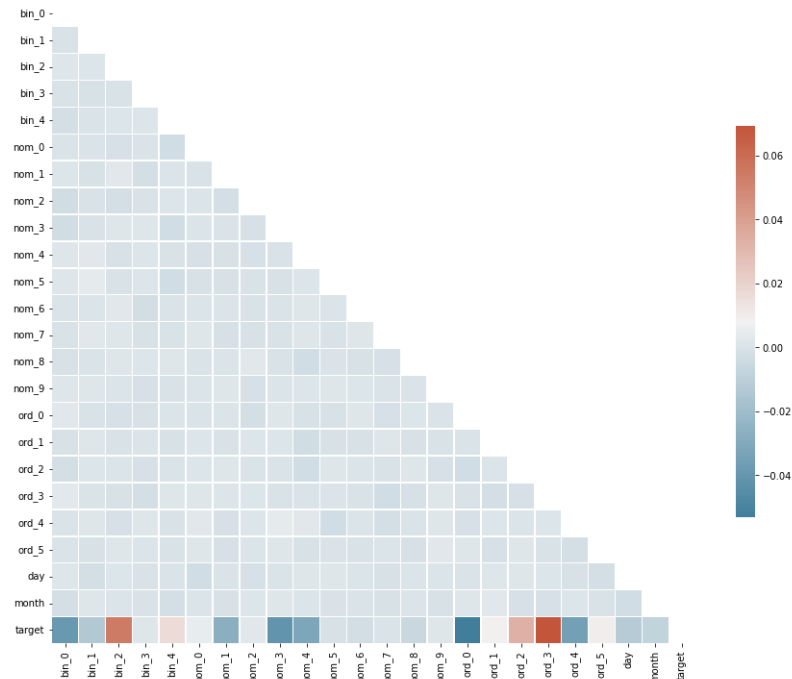


Рисунок 4. Корреляционная треугольная матрица.

Как видно, некоторые переменные выделяются на фоне остальных (bin_2, ord_0, ord_3), имея наибольшие значения корреляции с target переменной, но данная сила стохастической связи все еще очень мала.

Базовая модель

Перед проведением экспериментов следует определить базовую модель, от функционала качества которой нужно будет отталкиваться. Так как в нашем случае признаковое пространство исключительно категориальное, то для базовой модели будет взят вариант реализации градиентного бустинга от компании Яндекс – модель CatBoost.

Инициализация модели происходила с параметрами по умолчанию. Для оценки функционала качества модели была использована кросс-валидация со стратификацией и был получен результат для ROC AUC в 0.7701 ± 0.0013

Экспериментальная установка

Исследование методов категориального кодирования реализовано на высокоуровневом языке программирования Python, с использованием библиотек `scikit-learn`, `keras`, `tensorflow`, `XGBoost`, `CatBoost` для кодирования признаков и реализации алгоритма бинарной классификации. Запуск программы совершен на бесплатной облачной платформе Google Colab, которая позволяет совершать неподъемные для домашней машины вычисления.

Для сравнения методов кодирования используется модель градиентного бустинга, реализованный в библиотеке `XGBoost`.

Чтобы избежать ложных выводов по результатам работы модели на тестовом множестве, в исследовании используется кросс валидация со стратификацией, при разбиении на 5 долей. По итогу кросс валидации будет браться средняя по функционалу качества, на основе которой и будет сравнение.

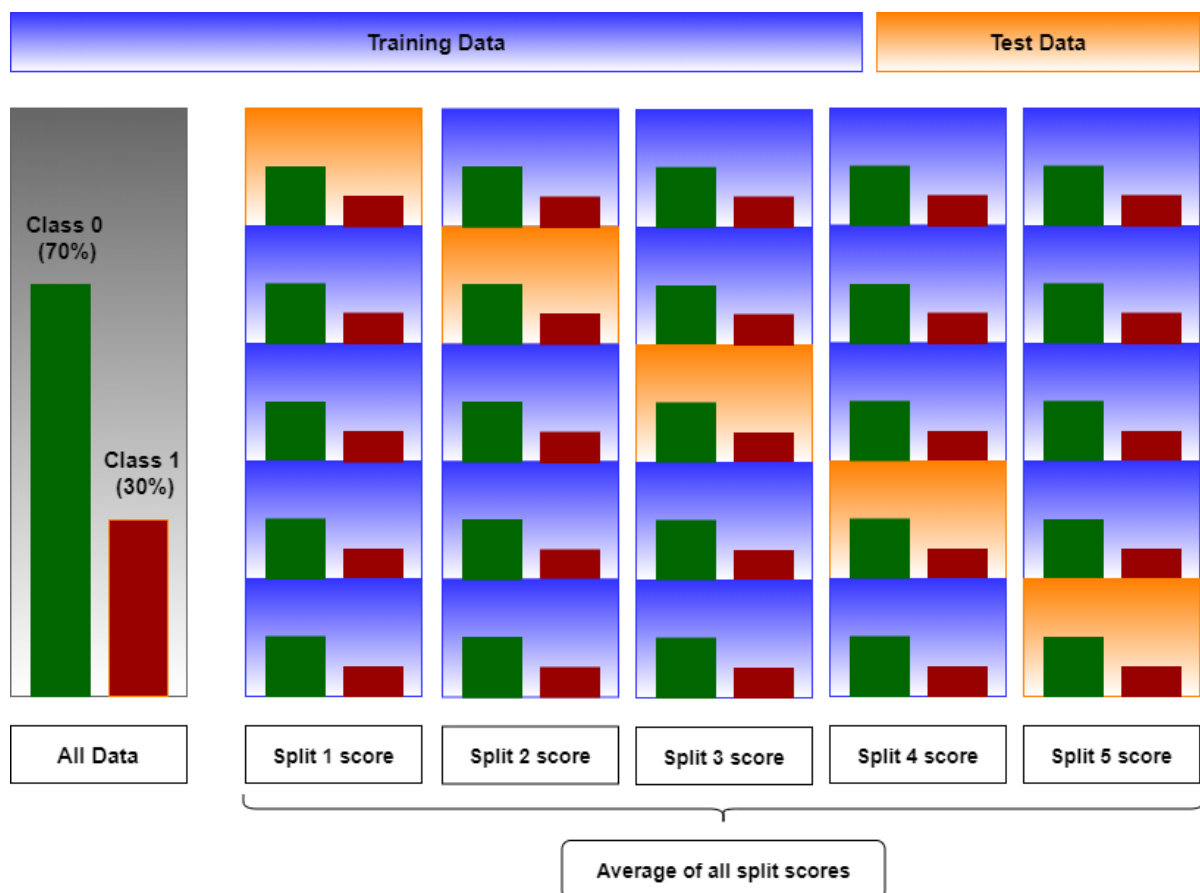


Рисунок 5. Схема стратифицированной кросс валидации

Также в исследовании участвует алгоритм логистической регрессии с параметрами по умолчанию.

Также для сравнения была использована нейронная сеть прямого распространения, состоящей из 5 слоев по 512 нейронов в каждом, с использованием нормализации и методов регуляризации[1][2][3]. На выходном слое функция активации – сигмоида. Обучение проходило с использованием бинарной перекрестной энтропии и методом оптимизации RMSProp.

Заполнение пропусков будет одинаковым для любого кодирования. Так как все данные категориальные и пропуски немногочисленные – около 3% на каждый признак, то для заполнения будет браться случайная целочисленная равномерная величина, соответствующая индексу уникального значения в признаке.

Все случайные генераторы инициализировались с начальным числом 42.

Методы кодирования

В данной работе рассмотрены наиболее частые и стандартные методы категорийного кодирования[4]:

- Факторизация (кодирование метки);
- Горячее кодирование (One Hot Encoding);
- Хэширование.

Далее рассмотрим каждый метод подробнее.

Факторизация – каждому значению переменной будет присвоен свой порядковый номер. Данный метод наиболее распространен для перехода от строковых значений к численным, для дальнейшей обработки алгоритмом классификации. В данном методе пропускам в данных будет отведена отдельная категория.

Горячее кодирование – каждая переменная преобразуется в признаковый вектор размерности, равной количеству уникальных признаков. Тогда каждое значение будет кодироваться как вектор, состоящий из единицы в соответствующему значению столбце и нулей.

Хэширование – данный метод преобразует категориальные переменные в пространство целых чисел высших измерений, где расстояние между двумя векторами категориальных переменных приблизительно поддерживается в преобразованном числовом пространстве. При использовании хеширования число измерений будет намного меньше количества измерений с кодировкой, например One Hot Encoding. Однако данный метод позволяет и расширить пространство признаков.

Помимо методов кодирования также были разобраны методы обработки признаков:

- Скрещивание признаков;
- Восстановление семантической значимости.

Также рассмотрим каждый подробнее.

Скрещивание признаков – отличие от прошлого метода в том, что один признак может быть скрещен со многими другими одновременно. Данный метод будет реализован с помощью нейронной сети, реализующая множество скрещиваний благодаря своей природе.

Восстановление семантической значимости – слепое кодирование категориальных строковых признаков не всегда позволит достичь наивысшего показателя функционала качества. Если количество уникальных значений категориального признака не велико, то можно постараться закодировать его самостоятельно. Это поможет с такими признаками, градация уровней мастерства, градация температуры: холодно, тепло, горячо; иначе говоря, при кодировании сохраняется порядок значений. После преобразования данных корреляционная матрица на рисунке значительно изменилась.

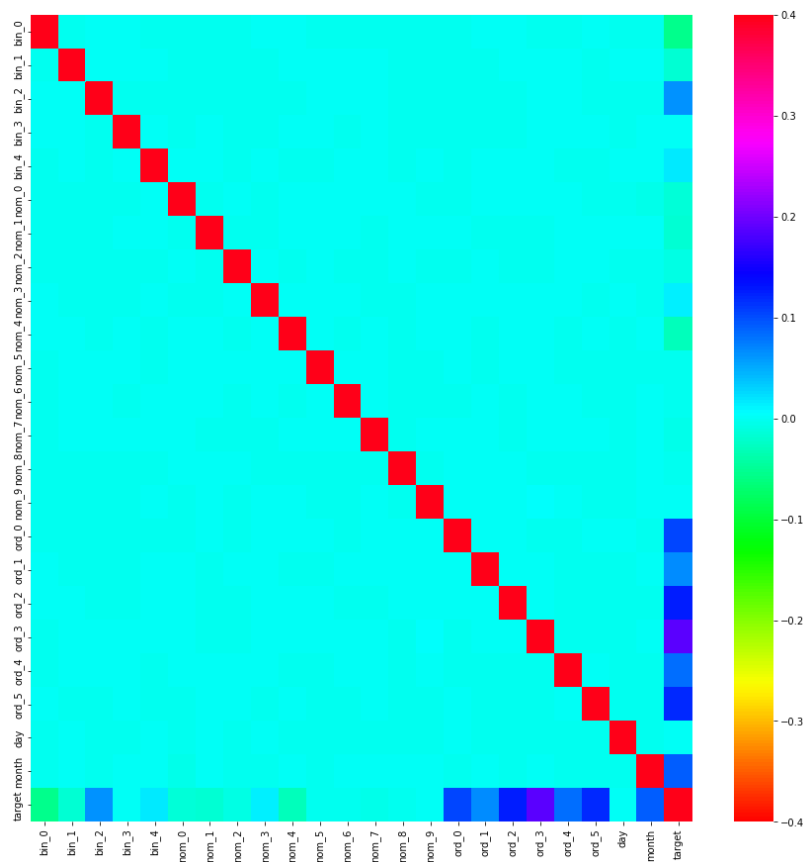


Рисунок 6. Корреляционная матрица признаков.

Результаты исследования

С учетом всего вышесказанного были получены следующие результаты, описанные в таблице 2.

Номер	Метод	ROC AUC
1	CatBoostClassifier + Факторизация	$0,7701 \pm 3*0,0013$
2	NN + восстановление градации + OneHotEncoding	$0,7506 \pm 3*0,0013$
3	XGBClassifier + Хэширование 400 признаков	$0,7454 \pm 3*0,0016$
4	XGBClassifier + восстановление градации + OneHotEncoding	$0,7453 \pm 3*0,0015$
5	XGBClassifier r + Факторизация	$0,7449 \pm 3*0,0017$
6	LogReg r + Факторизация с применением полинома второй степени	$0,7432 \pm 3*0,0015$
7	XGBClassifier + Хэширование 200 признаков	$0,7427 \pm 3*0,0015$
8	XGBClassifier + Хэширование 100 признаков	$0,7419 \pm 3*0,0014$
9	XGBClassifier + Хэширование 300 признаков	$0,7419 \pm 3*0,0014$
10	XGBClassifier + Хэширование 50 признака	$0,7343 \pm 3*0,0016$
11	XGBClassifier + OneHotEncoding	$0,7279 \pm 3*0,0013$
12	XGBClassifier + Хэширование 24 признака	$0,7163 \pm 3*0,0016$

Таблица 2. Результаты экспериментов.

По результатам исследования лучший результат дает модель градиентного бустинга с реализацией Яндекс. Ее функционал качества достигает 0,7701 с СКО в 0,0013. К данному результату не приблизилась ни одна рассматриваемая модель с кодированием или без.

На втором месте оказалась нейронная сеть прямого распространения из пяти слоев по 512 нейронов и с учетом количества входных параметров модель состоит из 1 млн обучаемых коэффициентов модели. Данная модель достигла ROC AUC равному 0,7506 с СКО в 0,0013.

Занимательно, что градиентный бустинг с обычной факторизацией достигает значения ROC AUC в 0,7449 с СКО в 0,0017. Когда как при применении кодирование в превалирующем большинстве случаев значение функционала качества меньше данного.

Заключение

В данной работе были рассмотрены методы кодирования категориальных признаков в задаче бинарной классификации, были рассмотрены такие методы, как факторизация (Label encoding), горячее кодирование (One Hot Encoding) и Хэширование (Hashing).

Также были рассмотрены методы обработки признаков, такие как скрещивание и восстановление цифрового порядка значений.

Все методы были исследованы в одинаковых условиях с заданным стартовым значением генератора случайных чисел и были получены результаты, отраженные в таблице.

На основании результатов исследования, можно сделать вывод, что градиентный бустинг позволяет достичь удовлетворительных значений функционала качества без особых усилий.

Но если исключить из рассмотрения разработку модели градиентного бустинга от Яндекс, то лучший результат приносит нейронная сеть прямого распространения с количеством обучаемых параметров в 1 млн., что не слишком много для современных мощностей ЭВМ.

Список литературы

1. Глубокое обучение / Ян Гудфеллоу, Иошуа Бенджио, Аарон Курвилль // ДМК Пресс, 2018г., второе цветное издание, исправленное.
2. Understanding the difficulty of training deep feedforward neural networks. Xavier GlorotYoshua BengioDIRO, Université de Montréal, Montréal, Québec, Canada
3. Глубокое обучение. / Николенко С., Кадури А., Архангельская Е. // СПб: Питер, 2018. — 480 с.: ил. — (Серия «Библиотека программиста»).
4. Методы кодирования / <https://innovation.alteryx.com/encode-smarter/>

Приложение