

Предсказание ухода сотрудников

Канаметов Азамат, Пузач Владислав, Байшев Олег

Московский физико-технический институт

25 декабря 2020 г.



1 Первый этап

2 Второй этап

- Логистическая регрессия
- Нейросетевые модели для классификации
- Нейросетевые модели для временного ряда

Актуальность

Последствия ухода сотрудников:

- ❶ нарушения дедлайнов
- ❷ дополнительные затраты на найм сотрудников
- ❸ временное снижение эффективности

Возможности:

- ❶ выполнение проектов в срок
- ❷ сохранение или повышение эффективности

Данные

- предсказываемые значения — сотрудник в конце года покинет компанию: 1, останется: 0
- соотношение классов примерно 1 к 6
- 4410 записей
- `general_data.csv` — общие данные о сотрудниках: возраст, отдел, уровень образования, доход и т.д
- `employee_survey_data.csv` — результаты первого опроса: уровни удовлетворенности рабочим пространством, удовольствия от выполнения задач, баланса между работой и личной жизнью.
- `manager_survey_data.csv` — результаты второго опроса: оценка сотрудниками менеджеров по их вовлеченности в работу и эффективности.

Датасет

Предобработка

- произведено one-hot кодирование категориальных признаков 'Department', 'EducationField', 'JobRole', 'MaritalStatus', 'Over18'
- бинарные признаки 'Attrition': {'No': 0, 'Yes': 1}, 'Gender': {'Female': 0, 'Male': 1}
- удалены неинформативные признаки 'EmployeeCount', 'Over18_Y', 'StandardHours': у каждого было всего одно значение для всех строк
- заполнены пропуски в столбцах 'NumCompaniesWorked', 'TotalWorkingYears', 'JobSatisfaction', 'EnvironmentSatisfaction', 'WorkLifeBalance': использовалось предсказание `sklearn.ensemble.RandomForestClassifier(n_estimators=30)` по известным признакам отдельно для каждого неизвестного признака
- нормализация: `sklearn.preprocessing.MinMaxScaler`

Поиск параметров

model	parameters	accuracy	f1 score	precision	recall
LogRegression	$C = 0.1$	84.4 ± 0.4	13.4 ± 2.4	7.4 ± 1.4	75.6 ± 10.5
LogRegression	$C = 1$	84.2 ± 0.2	22.6 ± 1.3	14.1 ± 1.2	59.2 ± 7.2
LogRegression	$C = 10$	84.5 ± 0.2	25.4 ± 1.2	16.1 ± 1.2	61.0 ± 6.7
LogRegression	$C = 100$	84.5 ± 0.2	25.4 ± 1.3	16.1 ± 1.2	61.0 ± 6.7
SVC	$C = 0.1, \text{kernel} = \text{poly}$	83.7 ± 0.4	1.3 ± 0.1	6.9 ± 0.2	100 ± 0
SVC	$C = 1, \text{kernel} = \text{poly}$	92.5 ± 0.8	71.7 ± 4.5	58.5 ± 6.4	93.6 ± 1.7
SVC	$C = 10, \text{kernel} = \text{poly}$	98.6 ± 0.5	95.7 ± 1.5	94.3 ± 2.0	97.1 ± 1.1
SVC	$C = 100, \text{kernel} = \text{poly}$	98.8 ± 0.4	96.1 ± 1.3	95.2 ± 1.9	97.2 ± 1.0
DecisionTree	$\text{max_depth} = 5$	85.1 ± 0.1	34.1 ± 1.6	23.5 ± 1.9	62.6 ± 6.1
DecisionTree	$\text{max_depth} = 10$	92.4 ± 0.7	74.1 ± 3.1	66.3 ± 4.8	84.2 ± 0.8
DecisionTree	$\text{max_depth} = 20$	98.6 ± 0.3	95.8 ± 0.8	95.9 ± 2.8	96.0 ± 3.2
DecisionTree	$\text{max_depth} = 40$	97.3 ± 0.3	94.4 ± 0.8	93.9 ± 1.8	95.0 ± 1.2
AdaBoost	$n_estimators = 5$	84.0 ± 0.5	20.7 ± 6.6	13.2 ± 5.5	57.6 ± 7.8
AdaBoost	$n_estimators = 10$	84.7 ± 0.2	25.9 ± 1.7	16.4 ± 1.7	63.4 ± 6.9
AdaBoost	$n_estimators = 20$	84.9 ± 0.3	31.0 ± 1.3	20.7 ± 1.3	62.1 ± 1.7
AdaBoost	$n_estimators = 50$	85.0 ± 0.4	35.7 ± 2.4	25.5 ± 3.0	60.1 ± 3.8
AdaBoost	$n_estimators = 100$	85.3 ± 0.4	39.7 ± 2.5	29.5 ± 2.7	61.0 ± 0.5
RandomForest	$n_estimators = 5$	97.7 ± 0.2	92.4 ± 0.7	87.9 ± 2.7	97.7 ± 1.8
RandomForest	$n_estimators = 10$	98.3 ± 0.1	94.6 ± 0.1	89.9 ± 0.3	100 ± 0.0
RandomForest	$n_estimators = 20$	99.0 ± 0.5	96.7 ± 1.6	93.6 ± 3.0	100 ± 0.0
RandomForest	$n_estimators = 50$	99.0 ± 0.3	96.8 ± 0.8	93.8 ± 1.5	100 ± 0.0
RandomForest	$n_estimators = 100$	99.0 ± 0.3	96.8 ± 0.8	93.8 ± 1.5	100 ± 0.0
GradientBoost	$n_estimators = 10$	84.3 ± 0.2	8.7 ± 0.9	4.6 ± 0.5	95.2 ± 6.7
GradientBoost	$n_estimators = 50$	86.7 ± 0.5	35.2 ± 3.1	22.1 ± 2.9	88.5 ± 6.6
GradientBoost	$n_estimators = 100$	87.7 ± 0.1	44.1 ± 1.5	29.5 ± 1.7	87.7 ± 6.6
GradientBoost	$n_estimators = 300$	93.4 ± 0.2	75.7 ± 0.5	62.3 ± 1.3	95.1 ± 1.9
GradientBoost	$n_estimators = 400$	95.4 ± 0.3	84.2 ± 1.1	74.5 ± 1.7	97.0 ± 0.5
GradientBoost	$n_estimators = 20, \text{max_depth} = 10$	97.8 ± 0.4	93.1 ± 1.3	88.9 ± 1.3	97.7 ± 1.8
GradientBoost	$n_estimators = 100, \text{max_depth} = 10$	98.9 ± 0.2	96.5 ± 0.4	93.8 ± 1.5	99.3 ± 1.0
GradientBoost	$n_estimators = 300, \text{max_depth} = 10$	99.0 ± 0.3	96.8 ± 0.8	93.8 ± 1.5	100 ± 0.0

Выбранные параметры и результаты

model	parameters	f1_score
LogRegression	'C': 100	25.4
SVC	'C': 100, 'kernel': poly	96.1
DecisionTree	'max_depth': 20	95.8
AdaBoost	'n_estimators': 100	39.7
RandomForest	'n_estimators': 50	96.8
GradientBoost	'n_estimators': 300, 'max_depth': 10	96.8

Вклад

- 1 Канаметов Азамат: применил метод опорных векторов и композиции алгоритмов
- 2 Пузач Владислав: предобработка и разведочный анализ
- 3 Байшев Олег: нашёл датасет, оформил презентации, добавил комментарии к действиям в тетрадке

Постановка

На 2 этапе необходимо было выполнить задачи:

- 1 построения логистической регрессии как базового классификатора
- 2 определения наиболее важных признаков, исходя из результатов пункта 1
- 3 улучшения классификации при использовании нейросетевых моделей
- 4 прогнозирования временного ряда

Логистическая регрессия

Исходный датасет был предобработан на первом этапе.

Реализация: модель LogisticRegression из библиотеки sklearn

Поиск лучших параметров: GridSearchCV

Название	Значение
C	[1+10*i for i in range(11)]
solver	['lbfgs', 'liblinear']
penalty	['l2', 'none', 'l1']
max_iter	[100, 200, 300]

Лучший набор: {'C': 1, 'max_iter': 100, 'penalty': 'none', 'solver': 'lbfgs'}

f1_score на CV: 0.3026273569309039

f1_score на test: 0.2645502645502646

Длительность поиска параметров: 39.5 s

Характеристики компьютера: Intel(R) Xeon(R) CPU 2.30GHz, RAM 12.72 GB

Интерпретация

Сотрудники склонны к уходу из компании, когда:

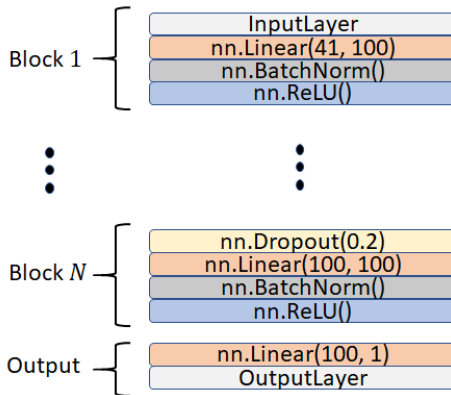
- 1 долго ждут повышения
- 2 часто бывают в командировках
- 3 работали во многих компаниях
- 4 не состоят в браке
- 5 учились на HR

Люди скорее всего останутся в компании на следующий год, если они:

- 1 имеют большой трудовой стаж
- 2 долгое время работают с текущим менеджером
- 3 солидного возраста
- 4 довольны работой
- 5 проходили обучение в течение последнего года

Архитектура

FCNN



Параметры и результаты

model	parameters	accuracy	f1 score	precision	recall
MLP	$max_iter = 500, hidden_layer_sizes = (50, 1)$	88.1 ± 5.4	53.7 ± 22.2	45.2 ± 22.3	67.9 ± 19.5
MLP	$max_iter = 1000, hidden_layer_sizes = (50, 1)$	89.0 ± 6.7	57.1 ± 26.9	49.9 ± 28.9	69.5 ± 21.6
MLP	$max_iter = 500, hidden_layer_sizes = (50, 2)$	92.1 ± 4.9	71.7 ± 18.9	64.2 ± 20.9	82.4 ± 15.3
MLP	$max_iter = 1000, hidden_layer_sizes = (50, 2)$	94.2 ± 6.3	79.0 ± 23.8	75.3 ± 27.5	84.8 ± 17.2
MLP	$max_iter = 500, hidden_layer_sizes = (50, 3)$	97.2 ± 1.5	91.2 ± 4.8	88.5 ± 6.6	94.3 ± 3.0
MLP	$max_iter = 1000, hidden_layer_sizes = (50, 3)$	98.4 ± 0.7	94.9 ± 2.1	94.0 ± 2.7	95.8 ± 1.7
MLP	$max_iter = 500, hidden_layer_sizes = (100, 1)$	85.9 ± 3.5	51.5 ± 11.1	45.4 ± 9.0	59.5 ± 14.5
MLP	$max_iter = 1000, hidden_layer_sizes = (100, 1)$	86.4 ± 4.1	53.0 ± 13.3	46.6 ± 10.6	61.7 ± 17.5
MLP	$max_iter = 500, hidden_layer_sizes = (100, 2)$	90.6 ± 4.8	64.2 ± 20.6	55.8 ± 21.8	77.9 ± 20.4
MLP	$max_iter = 1000, hidden_layer_sizes = (100, 2)$	91.6 ± 5.9	67.5 ± 24.0	59.0 ± 25.7	81.2 ± 22.5
MLP	$max_iter = 500, hidden_layer_sizes = (100, 3)$	96.8 ± 0.3	89.8 ± 1.3	86.1 ± 2.4	94.0 ± 0.4
MLP	$max_iter = 1000, hidden_layer_sizes = (100, 3)$	98.5 ± 0.5	95.2 ± 1.4	94.0 ± 2.2	96.5 ± 1.9

Рис.: MLP

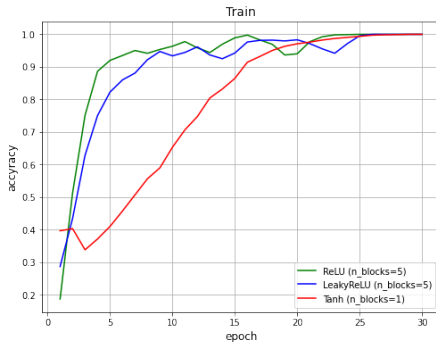
Параметры и результаты 2

model	block structure	N blocks	N parameters	accuracy	f1 score	precision	recall
FCNN	Linear	0	42	84.1 \pm 0.6	8.3 \pm 2.7	4.4 \pm 1.5	80.2 \pm 4.4
FCNN	Linear/Bn/ReLU	1	4501	97.0 \pm 1.1	90.2 \pm 3.9	85.0 \pm 5.3	96.3 \pm 2.2
FCNN	Linear/Bn/ReLU	2	14801	98.9 \pm 0.6	96.8 \pm 1.9	97.9 \pm 1.5	95.7 \pm 2.3
FCNN	Linear/Bn/ReLU	3	25101	99.0 \pm 0.4	96.9 \pm 1.4	96.0 \pm 2.1	97.9 \pm 1.2
FCNN	Linear/Bn/ReLU	4	35401	99.4 \pm 0.5	98.3 \pm 1.5	98.6 \pm 1.5	97.9 \pm 1.5
FCNN	Linear/Bn/ReLU	5	45701	99.7 \pm 0.2	98.9 \pm 0.8	99.1 \pm 0.9	98.8 \pm 0.9
FCNN	Linear/Bn/LeakyReLU	1	4501	94.9 \pm 1.3	82.6 \pm 5.6	75.2 \pm 8.2	92.0 \pm 1.4
FCNN	Linear/Bn/LeakyReLU	2	14801	97.3 \pm 1.2	91.4 \pm 4.3	88.6 \pm 6.0	94.4 \pm 2.4
FCNN	Linear/Bn/LeakyReLU	3	25101	98.3 \pm 0.9	94.6 \pm 3.1	91.4 \pm 4.6	98.2 \pm 1.7
FCNN	Linear/Bn/LeakyReLU	4	35401	98.5 \pm 0.6	95.4 \pm 0.2	95.8 \pm 3.1	94.9 \pm 1.1
FCNN	Linear/Bn/LeakyReLU	5	45701	99.5 \pm 0.3	98.5 \pm 1.0	98.1 \pm 1.2	98.8 \pm 0.9
FCNN	Linear/Bn/Tanh	1	4501	99.4 \pm 0.5	98.1 \pm 1.5	97.4 \pm 2.3	98.8 \pm 0.9
FCNN	Linear/Bn/Tanh	2	14801	98.6 \pm 0.4	95.4 \pm 1.3	91.2 \pm 2.3	100 \pm 0.0
FCNN	Linear/Bn/Tanh	3	25101	99.0 \pm 0.4	97.0 \pm 1.3	97.5 \pm 0.7	96.6 \pm 2.0
FCNN	Linear/Bn/Tanh	4	35401	98.6 \pm 0.7	95.7 \pm 2.3	96.5 \pm 3.2	95.0 \pm 1.6
FCNN	Linear/Bn/Tanh	5	45701	99.5 \pm 0.3	98.5 \pm 1.0	98.1 \pm 1.2	98.8 \pm 0.9

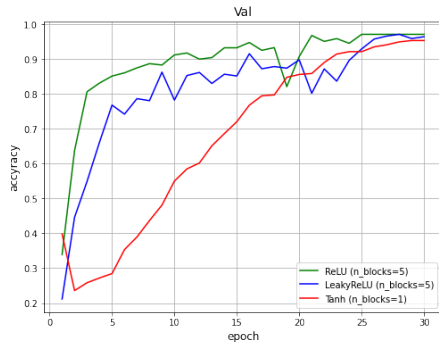
Характеристики компьютера:

PU	name	number of cores
CPU	Intel(R) Xeon(R) CPU 2.30GHz	2
GPU	Tesla T4	2560 (15079MiB)

Кривые обучения



a) accuracy on train



б) accuracy on test

Сравнение с результатами первого этапа

model	parameters	f1_score
SVC	'C': 100, 'kernel': poly	96.1
DecisionTree	'max_depth': 20	95.8
AdaBoost	'n_estimators': 100	39.7
RandomForest	'n_estimators': 50	96.8
GradientBoost	'n_estimators': 300, 'max_depth': 10	96.8
GradientBoost	'n_estimators': 300, 'max_depth': 10	96.8
LogRegression	'C': 1, 'max_iter': 100, 'penalty': 'none', 'solver': lbfgs	26.45
MLP	'max_iter': 1000, 'hidden_layer_sizes': (100, 3)	95.2
FCNN	'block structure': Linear/Bn/ReLU, 'N blocks': 5, 'N parameters': 45701	98.99

Данные

- 2 колонки: месяц; общее количество человек в США, летавших в этот месяц
- с 1949 по 1960 год

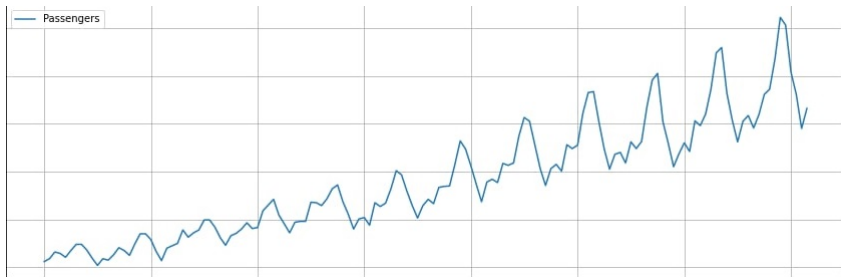


Рис.: Временной ряд

Датасет 2

Параметры, результаты

Результаты для простых моделей:

model	RMSE
mean	122.149
linear	45.7598
seasonal_decompose	9.9855

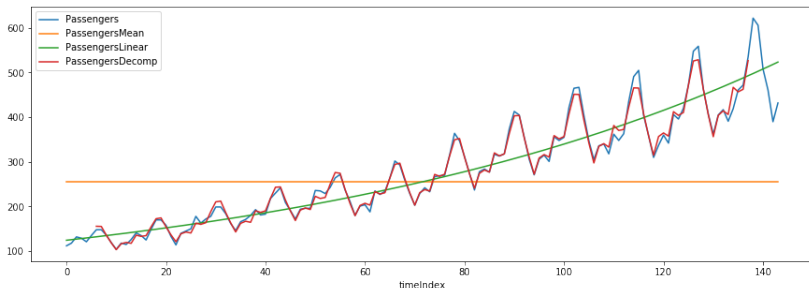


Рис.: Приближение средним, линейное

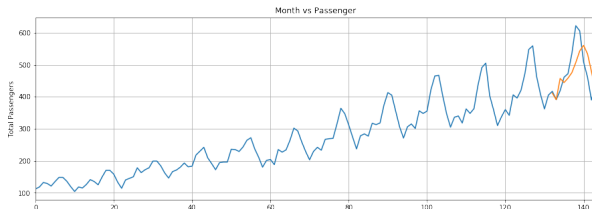
Графики для временного ряда

Параметры LSTM:

- Архитектура: 1 скрытый слой со 100 нейронами
- метод оптимизации Adam
- функция потерь: MSE
- количество эпох: 126
- learning rate: 0.001

Время выполнения для LSTM: train – 44.8 s, test – 0.01 s

Характеристики компьютера: RTX 2080Ti, 11264 Mb, Intel Core i9-9900KF, 3.6GHz



Вклад

Канаметов Азамат: использование нейросетевых моделей для задачи классификации.

Пузач Владислав: построение модели логистической регрессии и оформление результатов.

Байшев Олег: прогнозирование временного ряда.

Репозиторий