

## Задание – Дипломный проект по курсу «Python для анализа данных»

### с комментариями Зельберг Ирины

Дан файл HR.csv с данными по опросу уровня удовлетворенности сотрудниками работой.

Файл доступен тут - <https://drive.google.com/file/d/1INgo03nal-vwFJe7Lec5vOUtOwfJdUr1/view?usp=sharing>

Признаки:

1. satisfaction\_level - Уровень удовлетворенности работой
2. Last\_evaluation - Время с момента последней оценки в годах
3. number\_projects - Количество проектов, выполненных за время работы
4. average\_monthly\_hours - Среднее количество часов на рабочем месте в месяц
5. time\_spend\_company - Стаж работы в компании в годах
6. work\_accident - Происходили ли несчастные случаи на рабочем месте с сотрудником
7. left - уволился ли сотрудник
8. promotion\_last\_5years - повышался ли сотрудник за последние пять лет
9. department - отдел в котором работает сотрудник
10. salary - относительный уровень зарплаты

Требуется выполнить следующее задание:

# Задание Баллы	Комментарии (И. Зельберг)
1 Загрузите файл HR.csv в pandas dataframe 5	Выполнено
2 Рассчитайте основные статистики для переменных (среднее, медиана, мода, мин/макс, сред. отклонение). 10	Выведены статистики способами: <ul style="list-style-type: none"><li>▪ Методом describe из pandas</li><li>▪ Рассчитан каждый показатель отдельно методами numpy или pandas, результат статистики по столбам перевернут в сводную таблицу, затем результаты всех статистик добавлены в один датафрейм</li><li>▪ Рассчитан каждый показатель и транспонирован в горизонтальную таблицу, затем добавлены построчно в один датафрейм</li></ul>
3 Рассчитайте и визуализировать корреляционную матрицу для количественных переменных. Определите две самые скоррелированные и две наименее скоррелированные переменные. 10	Выполнено. Величины в представленных данных плохо коррелированы. Коэф-ты корреляции не выше 0,3. Данные были разделены на выборки с поиском более высоких корреляций внутри отдельных выборок (более узких и обобщенных логическим параметром): <ul style="list-style-type: none"><li>▪ работники, с которыми произошел несчастный случай;</li><li>▪ уволившиеся сотрудники;</li></ul>

	<ul style="list-style-type: none"> <li>▪ продвинувшиеся по службе за последние пять лет;</li> <li>▪ с высокими зарплатами</li> <li>▪ со средними зарплатами</li> <li>▪ с низкими зарплатами</li> </ul> <p>Внутри каждой группы существенных корреляций между параметрами не выявлено. Вероятно, это отражает, что на данные по персоналу влияет много субъективных факторов, поэтому прямых зависимостей между переменными нет, на значение показателей и поведение людей влияет много и разных причин в любом сочетании.</p>
4 Рассчитайте сколько сотрудников работает в каждом департаменте. 5	<p>Выполнено:</p> <ul style="list-style-type: none"> <li>▪ Методом <code>value_count</code></li> <li>▪ Группировкой в датафрейме</li> </ul>
5 Показать распределение сотрудников по зарплатам. 5	Выполнено
6 Показать распределение сотрудников по зарплатам в каждом департаменте по отдельности 5	Выполнено
7 Проверить гипотезу, что сотрудники с высоким окладом проводят на работе больше времени, чем сотрудники с низким окладом 10	<p>Применена функция <code>stats.ttest_1samp</code> из <code>scipy</code>. Результаты показаны, выводы сделать затруднительно. Способы проверки гипотез на курсе не раскрыты.</p>
8 Рассчитать следующие показатели среди уволившихся и не уволившихся сотрудников (по отдельности): 10 • Доля сотрудников с повышением за последние 5 лет • Средняя степень удовлетворенности • Среднее количество проектов	Рассчитано, отдельно по всем показателям
9 Разделить данные на тестовую и обучающую выборки Построить модель LDA, предсказывающую уволился ли сотрудник на основе имеющихся факторов (кроме <code>department</code> и <code>salary</code> ) Оценить качество модели на тестовой выборки 20	<p>Выполнено. Точность прогнозной модели – 75% Разделение на два класса (уволился или нет) не очень хорошо делается для этих данных, нечеткость и отсутствие прямых связей между показателями.</p>
10 Загрузить jupyter notebook с решение на github и прислать ссылку 5	

Итого - максимум 85 баллов Для зачета необходимо набрать минимум 55