

Материалы

- [illegible]

[3]. Постройте картинку, на которой будут отображены попарные зависимости признаков , 'Height', 'Weight' и 'BMI' друг от друга. Используйте метод *pairplot* библиотеки Seaborn.

Часто при первичном анализе данных надо исследовать зависимость какого-то количественного признака от категориального (скажем, зарплаты от пола сотрудника). В этом помогут "ящики с усами" - boxplots библиотеки Seaborn. Box plot - это компактный способ показать статистики вещественного признака (среднее и квантили) по разным значениям категориального признака. Также помогает отслеживать "выбросы" - наблюдения, в которых значение данного вещественного признака сильно отличается от других.

[4]. Создайте в DataFrame *data* новый признак *weight_category*, который будет иметь 3 значения: 1 – если вес меньше 120 фунтов. (~ 54 кг.), 3 - если вес больше или равен 150 фунтов (~68 кг.), 2 – в остальных случаях. Постройте «ящик с усами» (boxplot), демонстрирующий зависимость роста от весовой категории. Используйте метод *boxplot* библиотеки Seaborn и метод *apply* Pandas DataFrame. Подпишите ось y меткой «Рост», ось x – меткой «Весовая категория».

[5]. Постройте scatter plot зависимости роста от веса, используя метод *plot* для Pandas DataFrame с аргументом *kind='scatter'*. Подпишите картинку.

Задание 2. Минимизация квадратичной ошибки

В простейшей постановке задача прогноза значения вещественного признака по прочим признакам (задача восстановления регрессии) решается минимизацией квадратичной функции ошибки.

[6]. Напишите функцию, которая по двум параметрам w_0 и w_1 вычисляет квадратичную ошибку приближения зависимости роста y от веса x прямой линией $y = w_0 + w_1 * x$: $error(w_0, w_1) = \sum_{i=1}^n \{(y_i - (w_0 + w_1 * x_i))\}^2$ Здесь n – число наблюдений в наборе данных, y_i и x_i – рост и вес i -ого человека в наборе данных.

Итак, мы решаем задачу: как через облако точек, соответствующих наблюдениям в нашем наборе данных, в пространстве признаков "Рост" и "Вес" провести прямую линию так, чтобы минимизировать функционал из п. 6. Для начала давайте отобразим хоть какие-то прямые и убедимся, что они плохо передают зависимость роста от веса.

[7]. Проведите на графике из п. 5 Задания 1 две прямые, соответствующие значениям параметров $(w_0, w_1) = (60, 0.05)$ и $(w_0, w_1) = (50, 0.16)$. Используйте метод *plot* из *matplotlib.pyplot*, а также метод *linspace* библиотеки NumPy. Подпишите оси и график.

Минимизация квадратичной функции ошибки - относительно простая задача, поскольку функция выпуклая. Для такой задачи существует много методов оптимизации. Посмотрим, как функция ошибки зависит от одного параметра (наклон прямой), если второй параметр (свободный член) зафиксировать.

[8]. Постройте график зависимости функции ошибки, посчитанной в п. 6, от параметра w_1 при $w_0 = 50$. Подпишите оси и график.

Теперь методом оптимизации найдем "оптимальный" наклон прямой, приближающей зависимость роста от веса, при фиксированном коэффициенте $w_0 = 50$.

[9]. С помощью метода *minimize_scalar* из *scipy.optimize* найдите минимум функции, определенной в п. 6, для значений параметра w_1 в диапазоне [-5,5]. Проведите на графике из п. 5 Задания 1 прямую, соответствующую значениям параметров $(w_0, w_1) = (50, w_{1_opt})$, где w_{1_opt} – найденное в п. 8 оптимальное значение параметра w_1 .

При анализе многомерных данных человек часто хочет получить интуитивное представление о природе данных с помощью визуализации. Увы, при числе признаков больше 3 такие картинки нарисовать невозможно. На практике для визуализации данных в 2D и 3D в данных выделяют 2 или, соответственно, 3 главные компоненты (как именно это делается - мы увидим далее в курсе) и отображают данные на плоскости или в объеме.

Посмотрим, как в Python рисовать 3D картинки, на примере отображения функции $z(x,y) = \sin(\sqrt{x^2+y^2})$ для значений x и y из интервала [-5,5] с шагом 0.25.

```
from mpl_toolkits.mplot3d import Axes3D
```

Создаем объекты типа *matplotlib.figure.Figure* (рисунок) и *matplotlib.axes._subplots.Axes3DSubplot* (ось).

```
fig = plt.figure()  
ax = fig.gca(projection='3d')
```

```
# Создаем массивы NumPy с координатами точек по осям X и Y.
```

```

# Используем метод meshgrid, при котором по векторам координат
# создается матрица координат. Задаем нужную функцию Z(x, y).
X = np.arange(-5, 5, 0.25)
Y = np.arange(-5, 5, 0.25)
X, Y = np.meshgrid(X, Y)
Z = np.sin(np.sqrt(X**2 + Y**2))

# Наконец, используем метод *plot_surface* объекта
# типа Axes3DSubplot. Также подписываем оси.
surf = ax.plot_surface(X, Y, Z)
ax.set_xlabel('X')
ax.set_ylabel('Y')
ax.set_zlabel('Z')
plt.show();

```

[10]. Постройте 3D-график зависимости функции ошибки, посчитанной в п.6 от параметров w_0 и w_1 . Подпишите ось x меткой «Intercept», ось y – меткой «Slope», а ось z – меткой «Error».

[11]. С помощью метода *minimize* из `scipy.optimize` найдите минимум функции, определенной в п. 6, для значений параметра w_0 в диапазоне $[-100, 100]$ и w_1 - в диапазоне $[-5, 5]$. Начальная точка – $(w_0, w_1) = (0, 0)$. Используйте метод оптимизации L-BFGS-B (аргумент `method` метода *minimize*). Проведите на графике из п. 5 Задания 1 прямую, соответствующую найденным оптимальным значениям параметров w_0 и w_1 . Подпишите оси и график.