

# Собеседование на специализацию «Интеллектуальный анализ данных»

---

Киселев Никита Б05-002

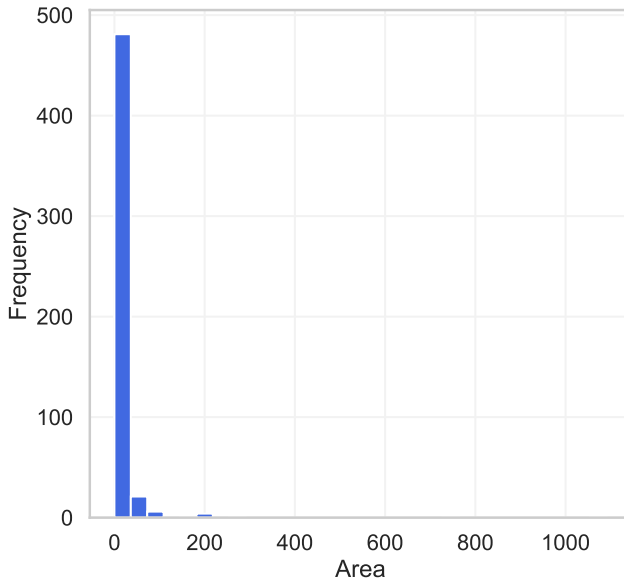
18 апреля 2022 г.

Московский физико-технический институт  
(национальный исследовательский университет)

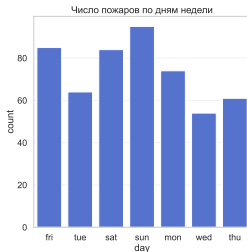
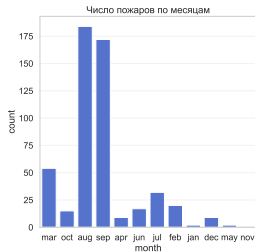
## Задача 21

Предсказание площади лесных пожаров. На основе погодных измерений необходимо предсказать объем выгоревших лесных массивов на севере Португалии. Выборка состоит из 13 признаков и 517 объектов. Для решения задачи предлагается использовать метод наименьших квадратов с регуляризацией. Нарисовать график весов признаков и общей ошибки на кросс-валидации при изменении параметра регуляризации. Какие признаки наиболее важны для нашей задачи? Что изменится, если предварительно все признаки стандартизовать?

# Распределение ответов



# Распределение номинальных признаков



# Корреляция количественных признаков



- Множество объектов  $\mathbb{X} = \mathbb{R}^n$
- Объекту  $\mathbf{x} \in \mathbb{X}$  соответствует признаковое описание  $\mathbf{x} = (f_1(x), \dots, f_n(x))$ , где  $f_j : \mathbb{X} \rightarrow D_j$
- Множество ответов  $\mathbb{Y} = \mathbb{R}$
- Выборка  $\mathbb{D} = \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathbb{X}, y_i \in \mathbb{Y}, i = 1, \dots, m\}$
- Матрица объекты-признаки  $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)^T$ , вектор ответов  $\mathbf{y} \in \mathbb{Y}^m$
- Вектор параметров модели  $\mathbf{w} = (w_1, \dots, w_n)^T$
- Ставится задача минимизации ошибки алгоритма  $Q(\mathbf{w}, X) = \|X\mathbf{w} - \mathbf{y}\|_2^2 \rightarrow \min_{\mathbf{w}}$

## Метод наименьших квадратов

$$Q(\mathbf{w}, X) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \rightarrow \min_{\mathbf{w}}$$

Приравняем к нулю производную по вектору  $\mathbf{w}$ :

$$\begin{aligned}\nabla_{\mathbf{w}} Q(\mathbf{w}, X) &= \nabla_{\mathbf{w}} (-\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y}) = \\ &= -\mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X})\mathbf{w} + 0 - \mathbf{X}^T \mathbf{y} = 0\end{aligned}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\boxed{\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

Могут возникнуть проблемы мультиколлинеарности в случае, если матрица  $X^T X$  плохо обусловлена. Один из способов решения — добавление к этой матрице диагональной:

$$\mathbf{w}^* = (X^T X + \alpha E_n)^{-1} X^T \mathbf{y}$$

При этом значении вектора  $\mathbf{w}$  достигается минимум функционала ошибки

$$Q(\mathbf{w}, X, \alpha) = \|X\mathbf{w} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{w}\|_2^2$$



# Изменение параметра $\alpha$

График зависимости MSE( $\alpha$ )

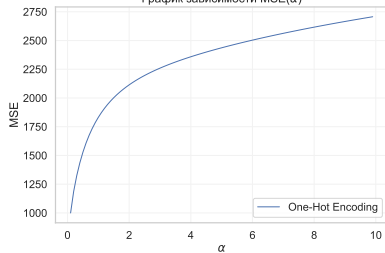


График зависимости весов признаков от  $\alpha$

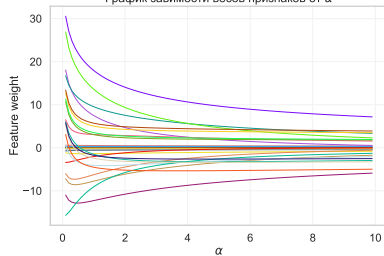


График зависимости MSE( $\alpha$ )

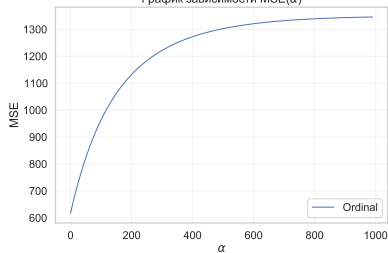
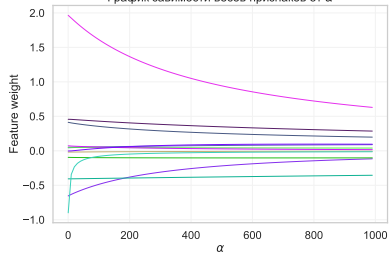


График зависимости весов признаков от  $\alpha$



При стандартизации происходит преобразование признаков:

$$\hat{f}_j(\mathbf{x}_i) = \frac{f_j(\mathbf{x}_i) - \bar{f}_j}{S_j},$$

где

$$\bar{f}_j = \frac{1}{m} \sum_{i=1}^m f_j(\mathbf{x}_i) \text{ — выборочное среднее,}$$

$$S_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (f_j(\mathbf{x}_i) - \bar{f}_j)^2} \text{ — среднеквадратичное отклонение.}$$

# Изменение параметра $\alpha$ при стандартизации

График зависимости MSE( $\alpha$ )

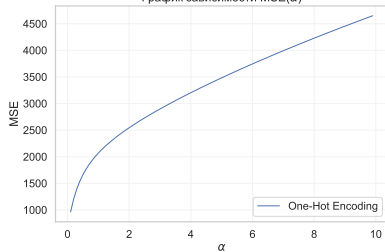


График зависимости весов признаков от  $\alpha$

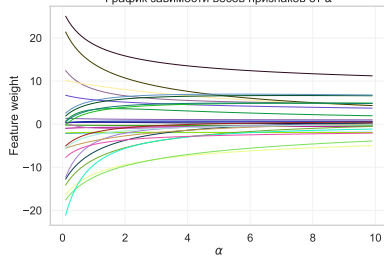


График зависимости MSE( $\alpha$ )

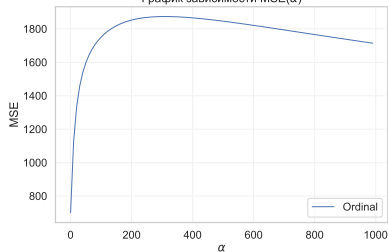
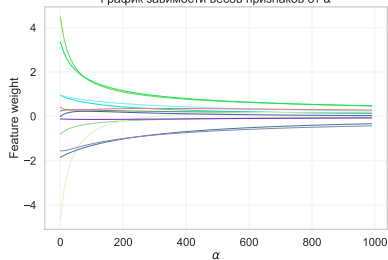
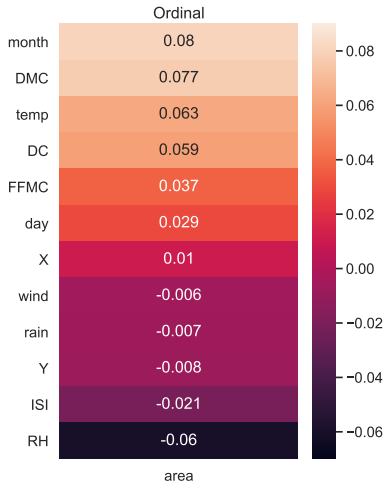
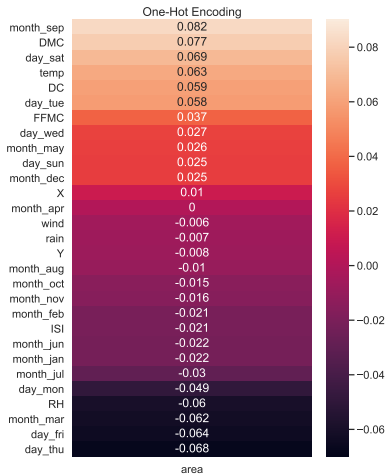


График зависимости весов признаков от  $\alpha$

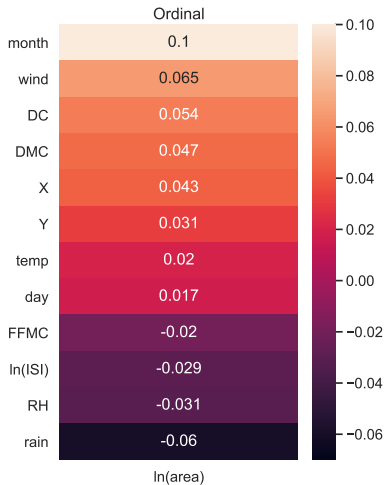
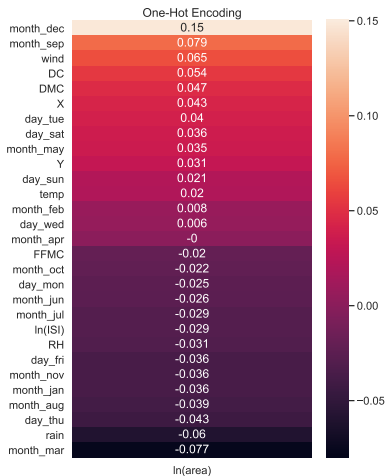


# Отбор признаков

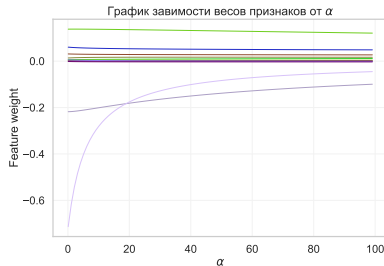
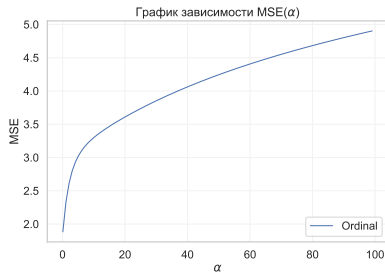
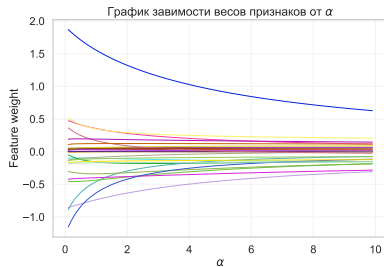
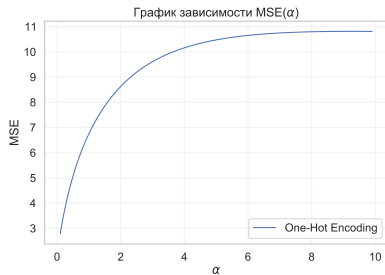


1. rain — номинальный
2.  $FFMC \geq 75$
3.  $ISI \rightarrow \ln(ISI)$
4.  $area \rightarrow \ln(1 + area)$

# Взаимосвязь новых признаков и ответов



# Изменение параметра $\alpha$



# Изменение параметра $\alpha$ при стандартизации

График зависимости MSE( $\alpha$ )

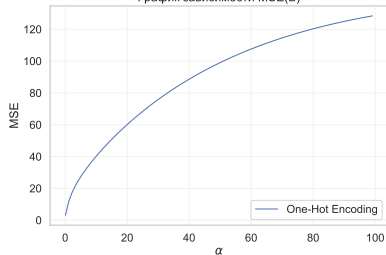


График зависимости весов признаков от  $\alpha$

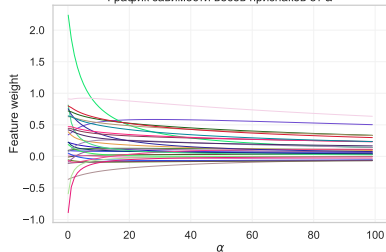


График зависимости MSE( $\alpha$ )

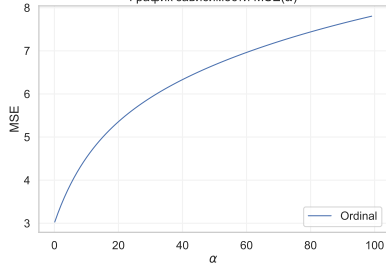
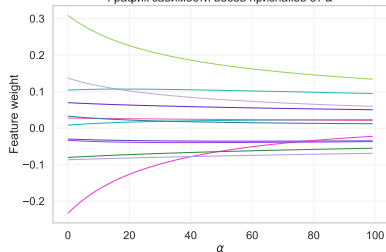


График зависимости весов признаков от  $\alpha$





**Таблица 1:** Лучшее значение MSE на кросс-валидации

Стандартизация \ Преобразование	До	После
	До	После
—	616,50	1,82
+	624,78	1,89

Стоит отметить, что после преобразования ответами являются  $\ln(1 + area)$ .

Значимость признаков при решении задачи лучше всего оценивается на данных после преобразования. Таковыми являются:

- month — месяц года
- wind — скорость ветра
- rain — количество осадков
- DC и DMC — индексы засухи и влажности почвы