

«Машинное обучение»

Линейные методы (часть 1): линейная регрессии

Александр Дьяконов

28 сентября 2021 года

План

Линейная регрессия

Решение проблемы вырожденности

Регуляризация, гребневая регрессия, LASSO, Elastic Net

Устойчивая регрессия

Градиентный метод обучения

Линейная регрессия

**Гипотеза о линейной зависимости целевой переменной,
ищем решение в виде:**

$$a(X_1, \dots, X_n) = w_0 + w_1 X_1 + \dots + w_n X_n$$

Практика:

- **часто неплохо работает и при монотонных зависимостях**
- **хорошо работает, когда есть много «однородных» зависимостей:**

цель – число продаж

признак 1 – число заходов на страницу продукта

признак 2 – число добавлений в корзину

признак 3 – число появлений продукта в поисковой выдаче

...

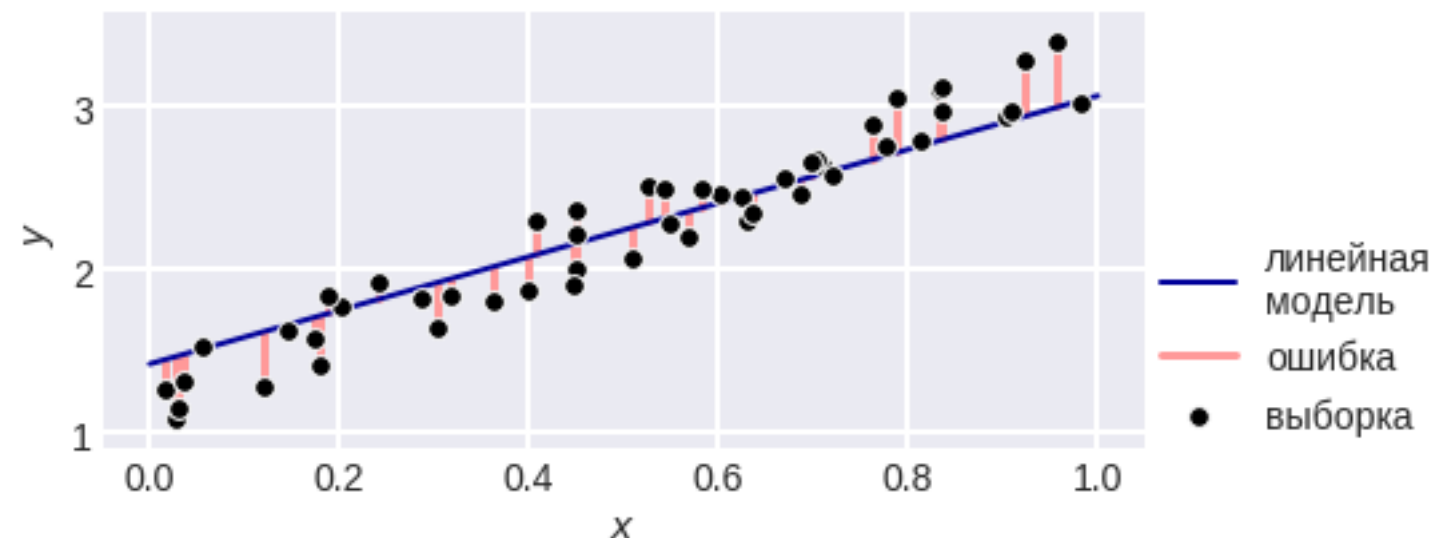
Линейная регрессия от одной переменной

$$a(X_1) = w_0 + w_1 X_1$$

обучение: $\{(x_1, y_1), \dots, (x_m, y_m)\}$, $x_i \in \mathbb{R}$

хотели бы...

$$\begin{cases} w_0 + w_1 x_1 = y_1 \\ \dots \\ w_0 + w_1 x_m = y_m \end{cases}$$



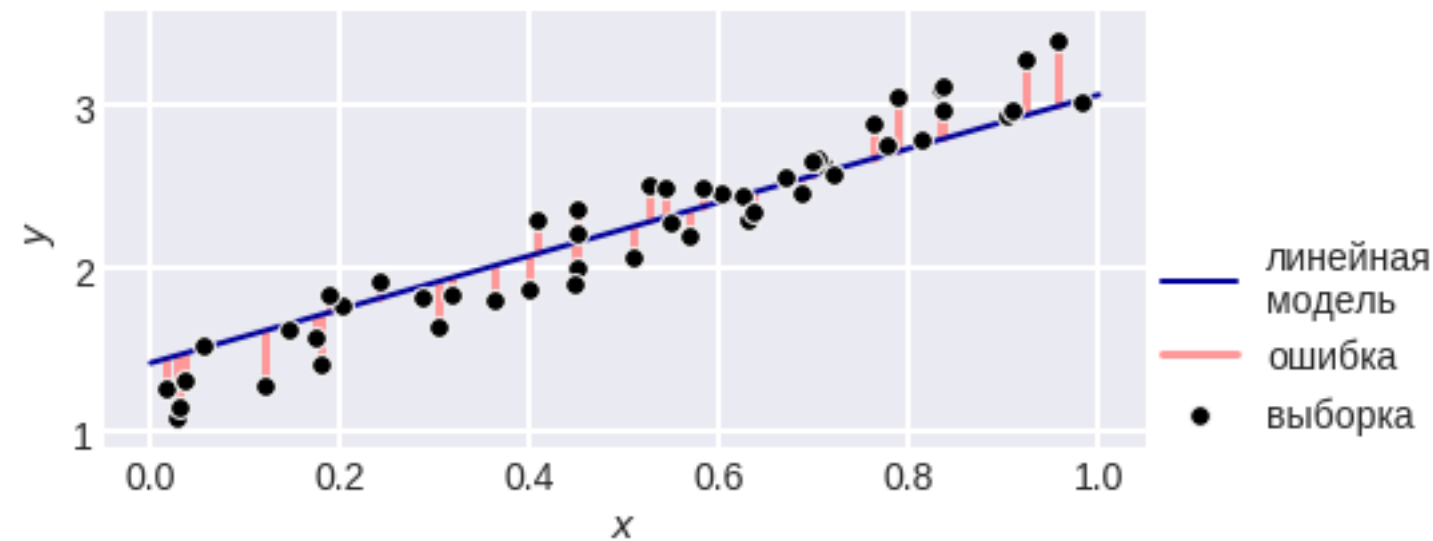
невязки / отклонения (residuals):

$$\begin{cases} e_1 = y_1 - w_0 - w_1 x_1 \\ \dots \\ e_m = y_m - w_0 - w_1 x_m \end{cases}$$

Линейная регрессия от одной переменной

**Задача минимизации суммы квадратов отклонений
(residual sum of squares)**

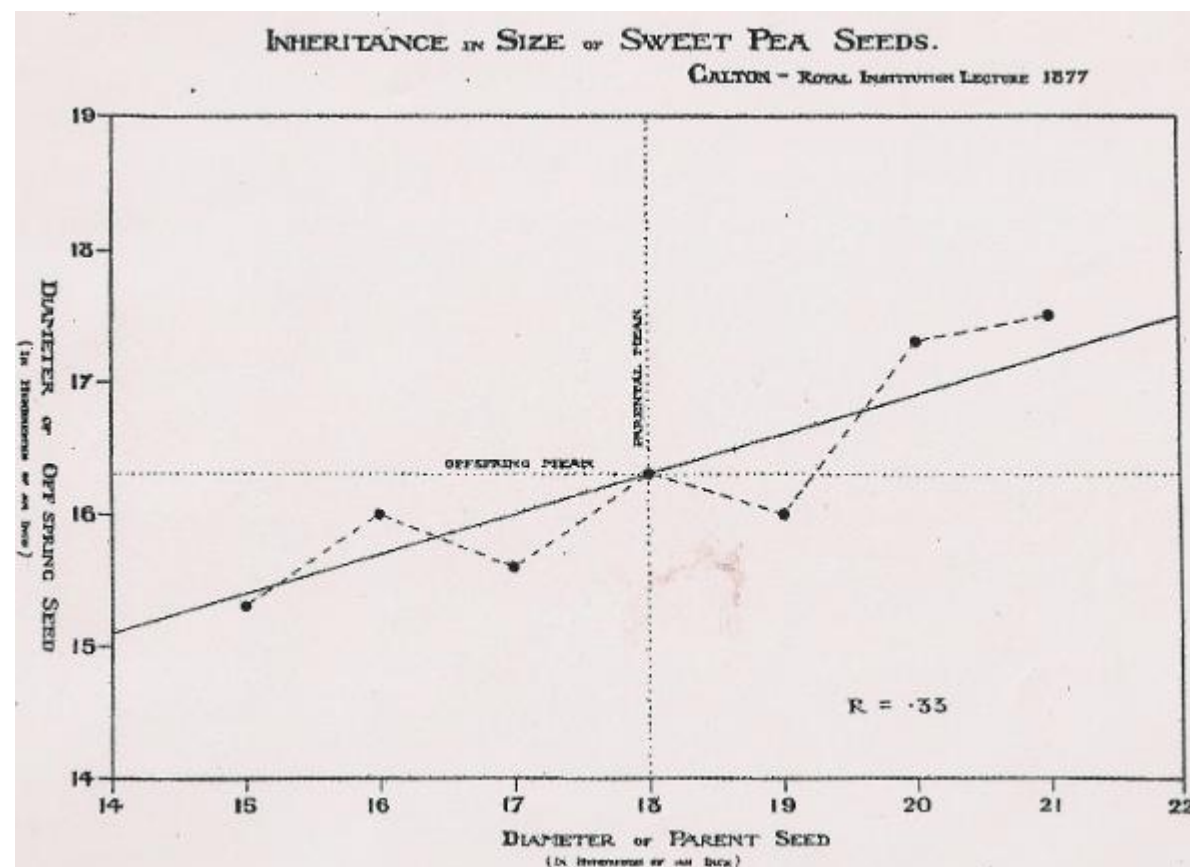
$$\text{RSS} = e_1^2 + \dots + e_m^2 \rightarrow \min$$



~ задача описания данных гиперплоскостью (но тут конкретная ф-я ошибки)

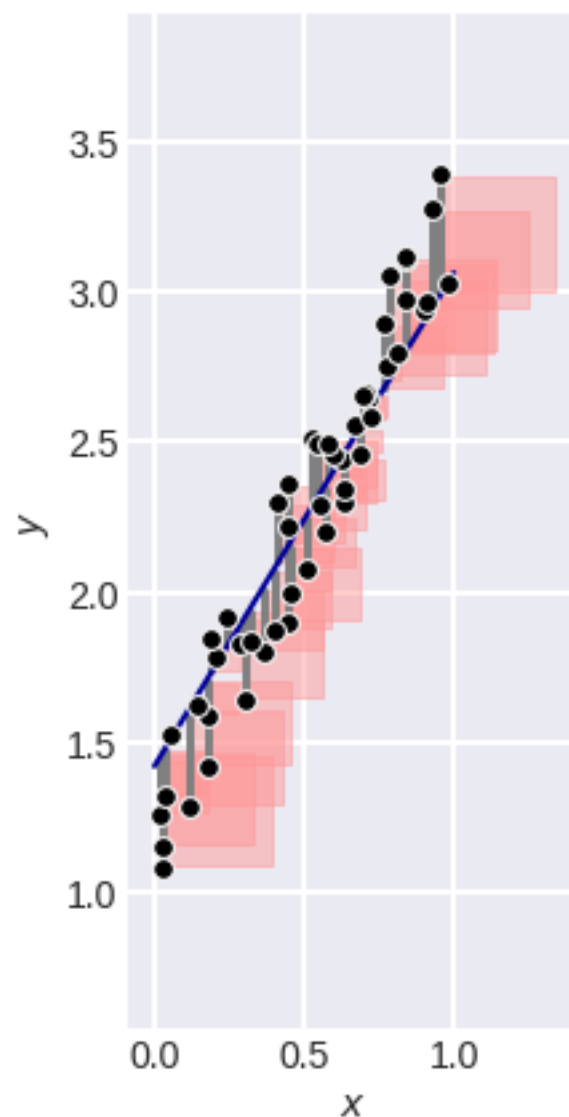
ПОТОМ вероятностное обоснование, пока... довольно логично

Линейная регрессия от одной переменной

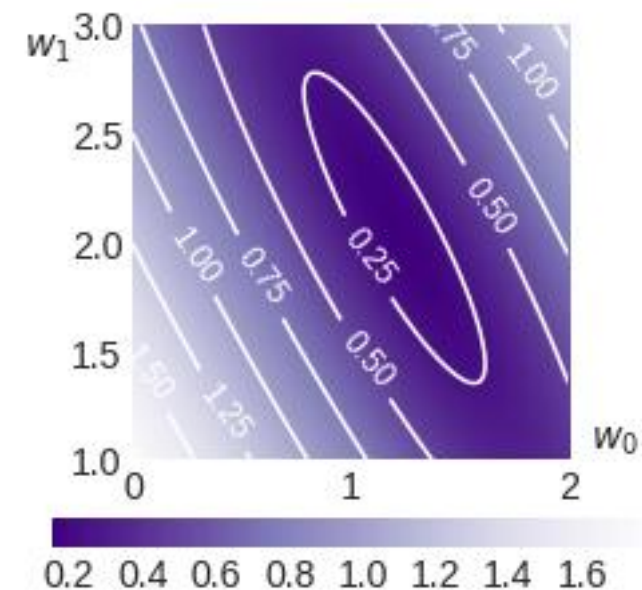


Francis Galton, 1877

Линейная регрессия от одной переменной: геометрический смысл ошибки



$$a(X_1) = w_0 + w_1 X_1$$



$$\sum_{i=1}^m (y_i - w_0 + w_1 x_i)^2$$

Отличается от суммы расстояний до поверхности!

Линейная регрессия от одной переменной

Нетрудно показать (**ДЗ**):

$$w_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2} = \frac{\text{cov}(\{x_i\}, \{y_i\})}{\text{var}(\{x_i\})},$$

$$w_0 = \bar{y} - w_1 \bar{x},$$

где $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$, $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$.

Общий случай (многих переменных)

$$a(X_1, \dots, X_n) = w_0 + w_1 X_1 + \dots + w_n X_n = x^T w$$

веса (параметры) – $w = (w_0, w_1, \dots, w_n)^T$

объект – $x = (X_0, X_1, \dots, X_n)^T$

для удобства записи вводим фиктивный признак $X_0 \equiv 1$

обучение: $\{(x_1, y_1), \dots, (x_m, y_m)\}$, $x_i \in \mathbf{R}^{n+1}$,

опять хотим решить $Xw = y$:

$$\begin{cases} x_1^T w = y_1 \\ \dots \\ x_m^T w = y_m \end{cases}$$

как решать?

Общий случай (многих переменных): в матричной форме

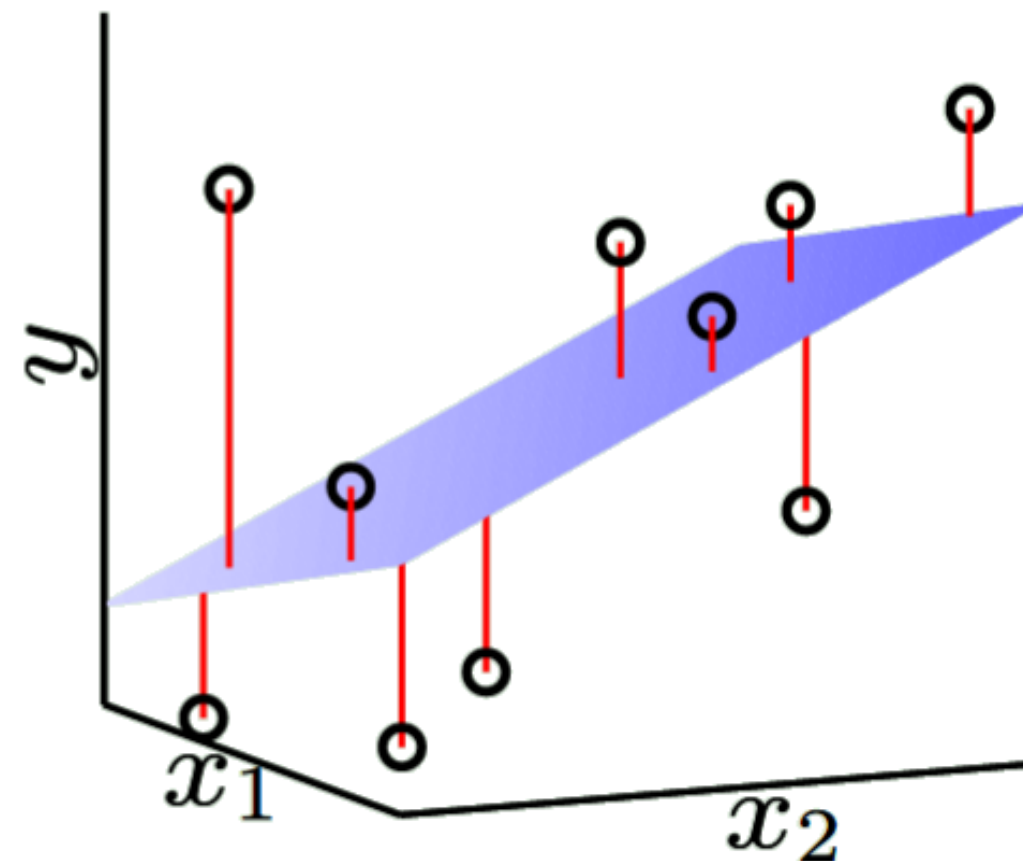
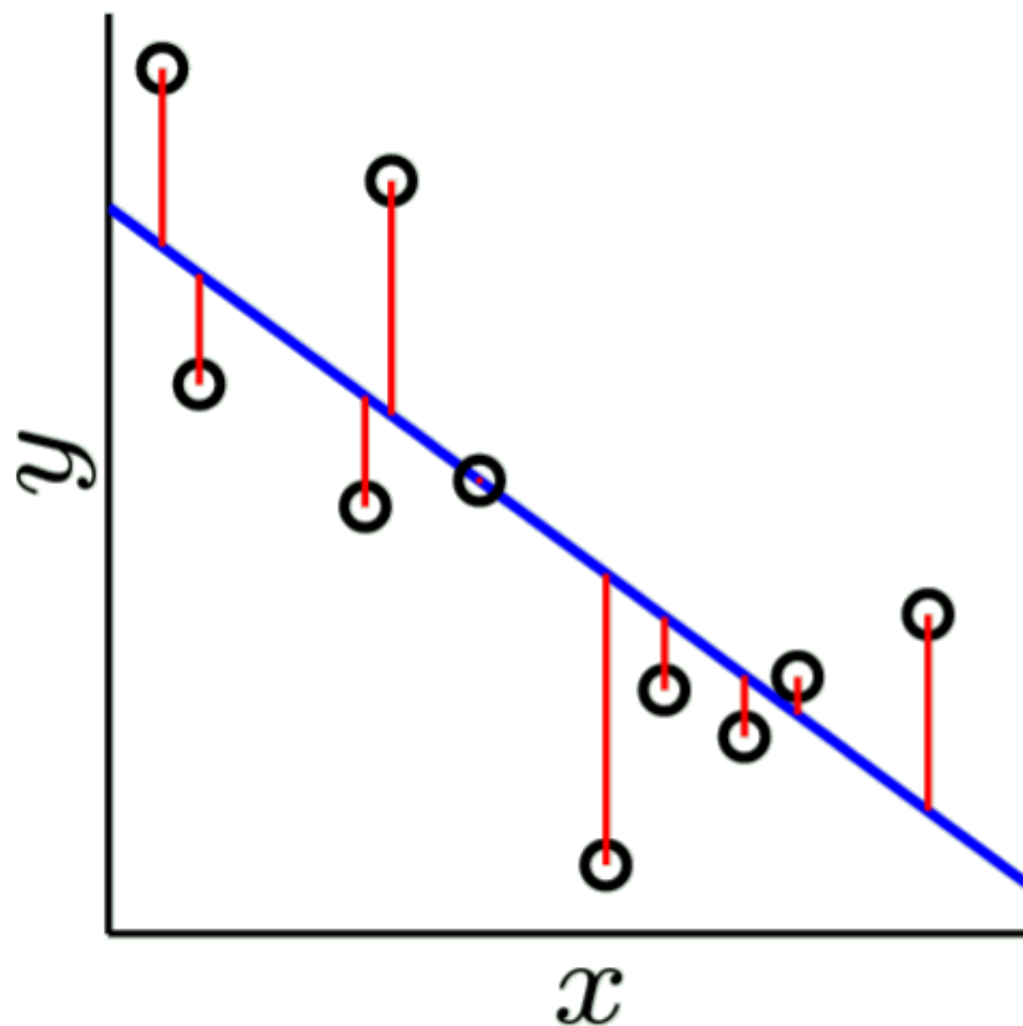
$$Xw = y$$

**в матрице X по строкам записаны описания объектов,
в векторе y значения их целевого признака
(здесь есть коллизия в обозначении y)**

будем решать так:

$$\|Xw - y\|_2^2 \rightarrow \min_w$$

почему?

Общий случай (многих переменных): геометрический смысл

Кстати, полученная задача оптимизации выпукла, единственный глобальный минимум (кроме вырожденного случая)

Решение задачи минимизации: прямой метод

$$\|Xw - y\|_2^2 \rightarrow \min_w$$

$$\|Xw - y\|_2^2 = (Xw - y)^T (Xw - y) = w^T X^T X w - w^T X^T y - y^T X w + y^T y$$

$$\nabla \|Xw - y\|_2^2 = 2X^T X w - 2X^T y = 0$$

$$X^T X w = X^T y$$

**решение существует,
если столбцы л/н**

$$w = (X^T X)^{-1} X^T y$$

помним, что $\text{rg}(X^T X) = \text{rg}(X)$

$(X^T X)^{-1} X^T$ – псевдообратная матрица Мура-Пенроуза
обобщение обратной на неквадратные матрицы

Обобщённая линейная регрессия: вместо X – что угодно

выражаем целевое значение через л/к базисных функций
(они фиксированы)

$$a(X_1, \dots, X_n) = w_0 + w_1 \varphi_1(X_1, \dots, X_n) + \dots + w_k \varphi_k(X_1, \dots, X_n)$$

$$w = (w_0, w_1, \dots, w_k)^T$$

$$x = (X_0, X_1, \dots, X_n)^T$$

$$\varphi(x) = (\varphi_0(x), \varphi_1(x), \dots, \varphi_k(x))^T$$

$\equiv 1$

$$a(x) = \sum_{i=1}^k w_i \varphi_i(x) = \varphi(x)^T w$$

$$\| \varphi(X)w - y \|_2^2 \rightarrow \min_w$$

$$\varphi(X) = \begin{bmatrix} \varphi_0(x_1) & \dots & \varphi_k(x_1) \\ \dots & \dots & \dots \\ \varphi_0(x_m) & \dots & \varphi_k(x_m) \end{bmatrix}$$

Подробности в нелинейных методах...

Проблема вырожденности матрицы

$$w = (X^T X)^{-1} X^T y$$

Только ли вырожденность плоха?

Что делать?

Проблема вырожденности матрицы

$$w = (X^T X)^{-1} X^T y$$

Проблемы, когда матрица X плохо обусловлена...

Решения:

- 1. Регуляризация – **здесь и в «сложности»****
- 2. Селекция (отбор) признаков – **«селекция»****
- 3. Уменьшение размерности (в том числе, PCA) – **USL****
- 4. Увеличение выборки**

если объектов много – то работать с гигантской матрицей невозможно...
но выдели как это делается в оптимизации онлайн-методами

Регуляризация: упрощённое объяснение смысла

$$a(X_1, \dots, X_n) = w_0 + w_1 X_1 + \dots + w_n X_n$$

**если есть два похожих объекта, то должны быть похожи метки,
пусть отличаются в j -м признаке, тогда ответы модели отличаются на**

$$\varepsilon_j w_j$$

Поэтому не должно быть очень больших весов
(у признаков, по которым могут отличаться похожие объекты)

**Поэтому вместе с $\|Xw - y\|_2^2 \rightarrow \min$
хотим $\|w\|_2^2 \rightarrow \min$**

Не на все коэффициенты нужна регуляризация! Почему?

Регуляризация: упрощённое объяснение смысла

Пусть есть какая-то зависимость и лишние признаки, например

$$y = X_1 = X_1 + w'X_2 - w'X_3 \text{ при } X_2 = X_3$$

Если теперь $X_2 \approx X_3$, тогда $\varepsilon = X_2 - X_3$

$$a = X_1 + w'\varepsilon$$

– может быть сколь угодно большим при больших w'

**аналогично при линейных зависимостях!
автоматически, когда объектов мало (сколько?)**

Регуляризация

Иванова

$$\begin{cases} \|Xw - y\|_2^2 \rightarrow \min \\ \|w\|_2^2 \leq \lambda \end{cases}$$

Тихонова

$$\|Xw - y\|_2^2 + \lambda \|w\|_2^2 \rightarrow \min$$

Удобнее: безусловная оптимизация

$$\|w\|_2^2 = w_1^2 + w_2^2 + \dots + w_n^2$$

эти две формы эквивалентны: решение одного можно получить как решение другого

Всё это справедливо и для общих задач минимизации!

$$\begin{cases} L(a) \rightarrow \min \\ \text{complexity}(a) \leq \lambda \end{cases}$$

$$L(a) + \lambda \text{complexity}(a) \rightarrow \min$$

Есть ещё регуляризация Морозова...

Регуляризация и гребневая регрессия

$$\arg \min_w \|Xw - y\|_2^2 + \lambda \|w\|_2^2 = (X^T X + \lambda I)^{-1} X^T y$$
$$\lambda \geq 0$$

ДЗ Доказать!

Такая регрессия называется **гребневой регрессией (Ridge Regression)**

Виден другой смысл регрессии: складываем две матрицы Грама,
неотрицательно определённая + положительно определённая

– боремся с вырожденностью матрицы

потом вернёмся к этому

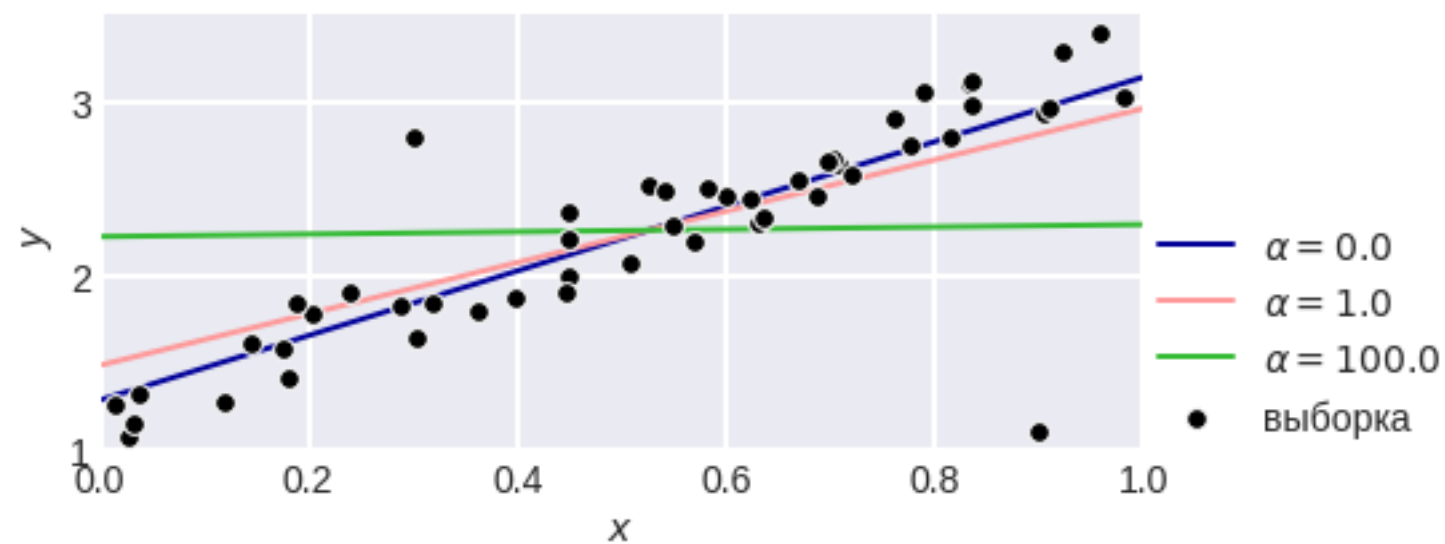
$\lambda = 0$ – получаем классическое решение

$\lambda \rightarrow +\infty$ – меньше «затачиваемся на данные» и больше регуляризуем

Матрица очевидно становится обратимой!

значение параметра регуляризации можно выбрать на скользящем контроле

Минутка кода: регуляризация и гребневая регрессия



```
from sklearn.linear_model import Ridge

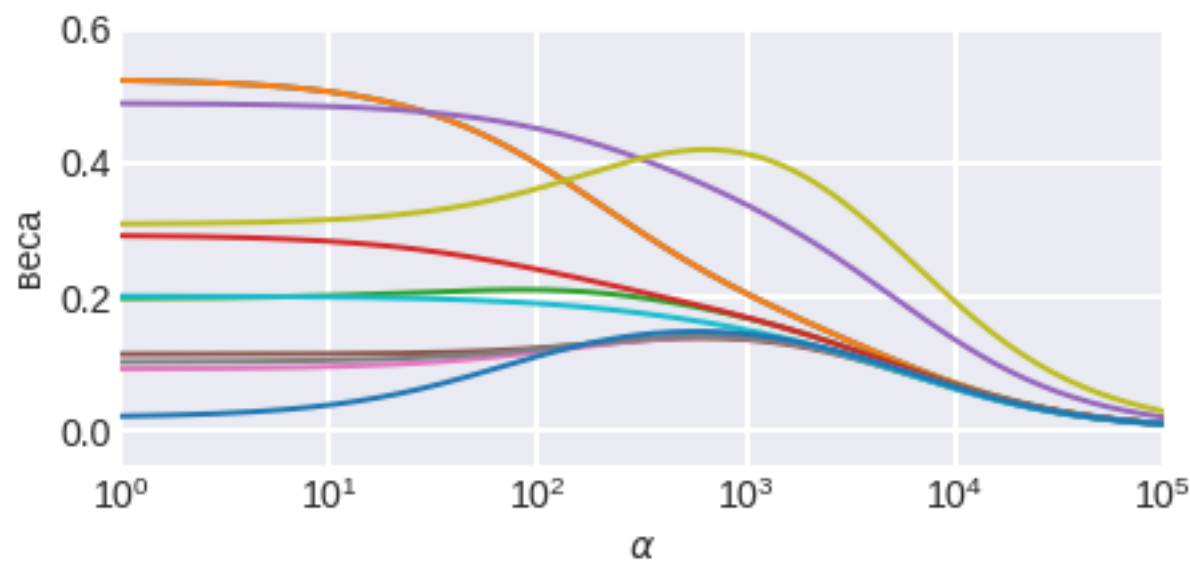
model = Ridge(alpha=0.0) # ридж-регрессия
# обучение
model.fit(x_train[:, np.newaxis], y_train)
# обратите внимание: np.newaxis
# контроль
a_train = model.predict(x_train[:, np.newaxis])
a_test = model.predict(x_test[:, np.newaxis])
```

Кажется, что при регуляризации
отклоняемся к выбросам,
но дело не в этом

Регуляризация и гребневая регрессия

$$\sum_{i=1}^m (y_i - a(x_i))^2 + \lambda \sum_{j=1}^n w_j^2 \rightarrow \min$$
$$\lambda \geq 0$$

добавление **shrinkage penalty** (регуляризатора)



параметр регуляризации может подбираться с помощью скользящего контроля

Регуляризация и гребневая регрессия

Для ridge-регрессии нужна правильная нормировка признаков!

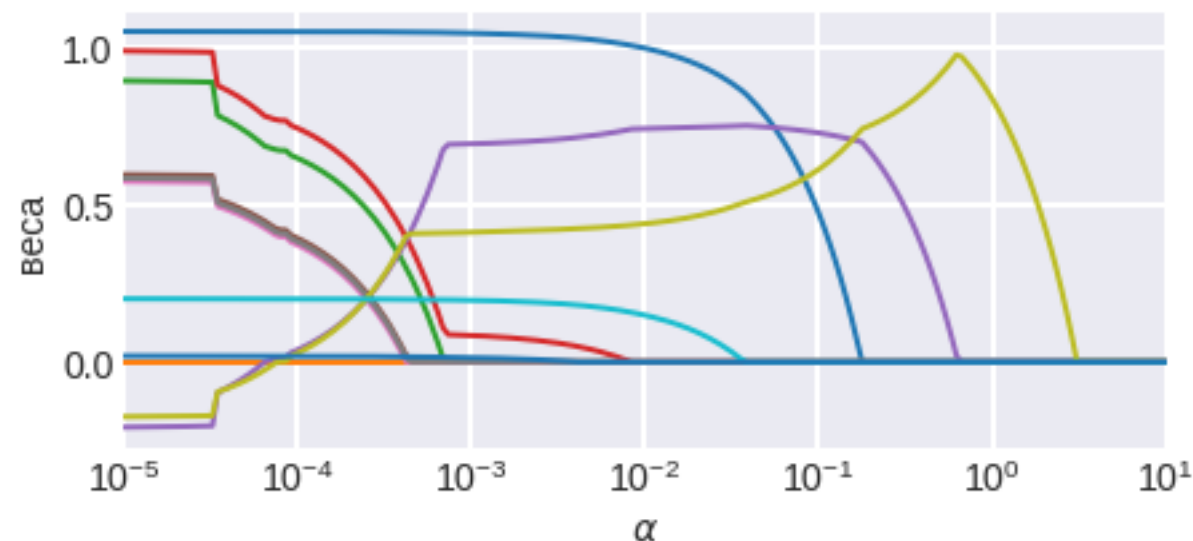
Нет инвариантности (в отличие от линейной) от умножения признаков на скаляры

Перед регуляризацией – стандартизация!!!

LASSO (Least Absolute Selection and Shrinkage Operator)

Попробуем другой «штраф за сложность» (сейчас поймём название)

$$\sum_{i=1}^m (y_i - a(x_i))^2 + \lambda \sum_{j=1}^n |w_j| \rightarrow \min$$
$$\lambda \geq 0$$



Здесь коэффициенты интенсивнее зануляются при увеличении $\lambda \geq 0$

Эксперименты с одинаковыми и зависимыми признаками

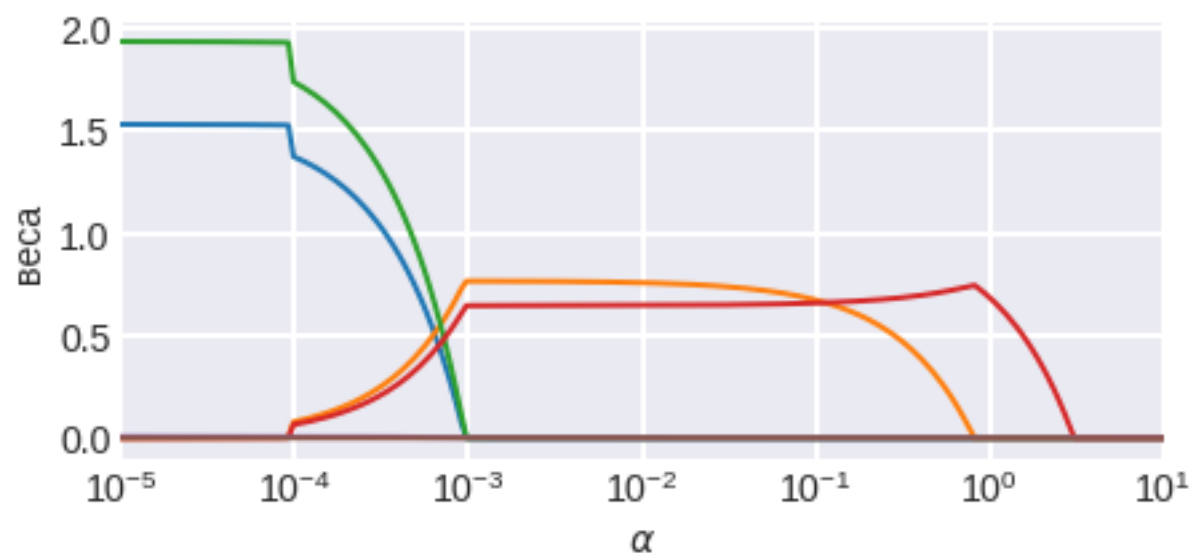
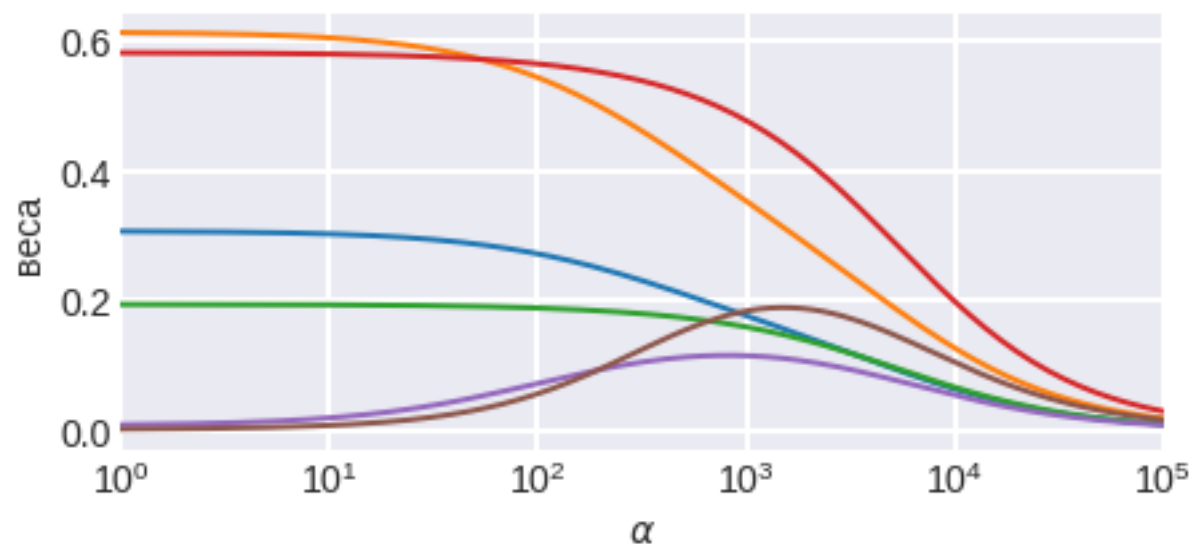
здесь была задача

```
X = np.random.rand(1000, 11)
X[:,1] = X[:,0]
X[:,4] = X[:,2] + X[:,3]
X[:,8] = X[:,5] + X[:,6] + X[:,7]
y = X[:,0] + 0.8*X[:,4] + 0.4*X[:,8] + 0.2*X[:,9] +
    0.5*np.random.randn(1000)
```

зависит от масштаба признаков,
но из-за предварительной нормировки этот эффект не наблюдается

Масштаб очень важен! см. дальше

Эксперименты с одинаковыми и зависимыми признаками



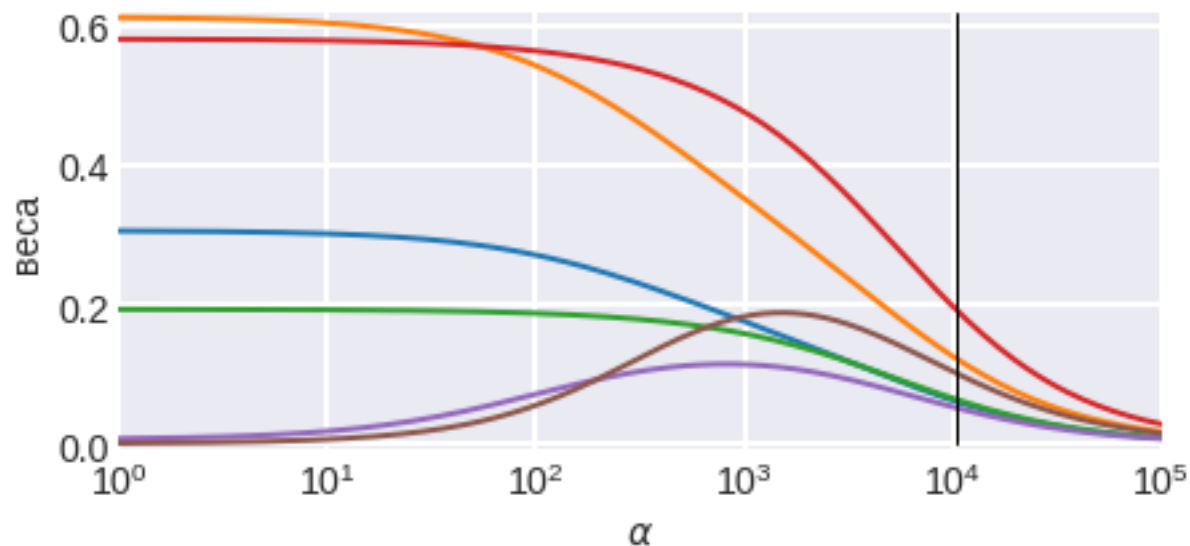
```
np.random.seed(10)
X = np.random.rand(1000, 6)
X[:,1] = X[:,0]
X[:,2] = X[:,3]
```

```
X[:,0] = 1 * X[:,0]
X[:,1] = 2 * X[:,1]
X[:,2] = 1 * X[:,2]
X[:,3] = 3 * X[:,3]
X[:,4] = 1 * X[:,4]
X[:,5] = 2 * X[:,5]
```

```
y = 1.5 * X[:,0] + 2*X[:,2] +
0.5*np.random.randn(1000)
```

Эксперименты с одинаковыми и зависимыми признаками

$$Y = 1.5X_1 + 2X_3 = 0.75X_2 + 0.66X_4$$

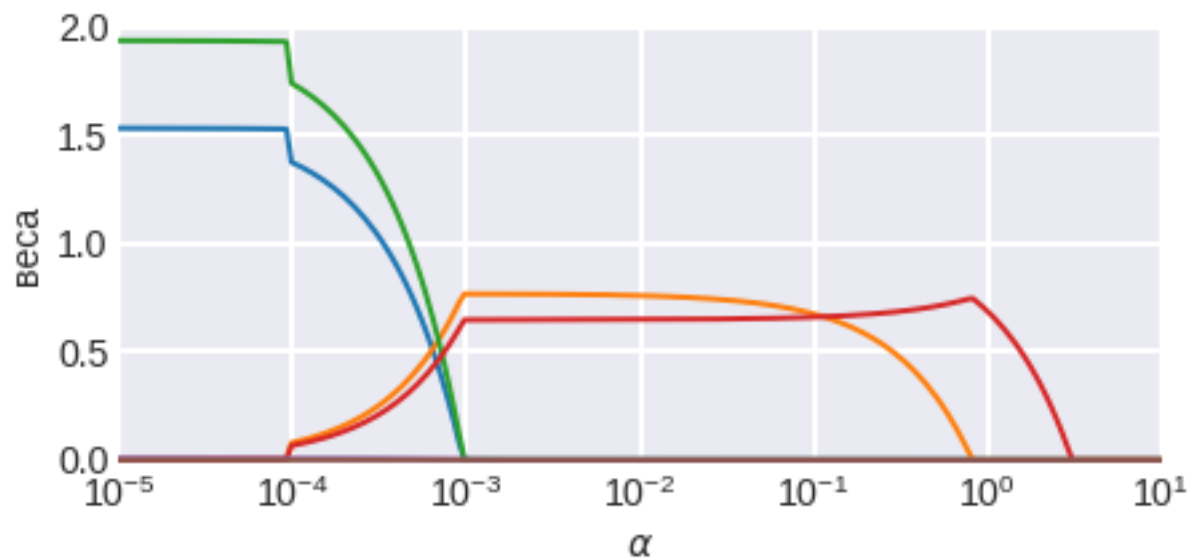


$$\lambda = 1$$

$$Y = 0.31X_1 + 0.61X_2 + 0.19X_3 + 0.58X_4 + 0.01X_5 + 0.0X_6$$

$$\lambda \sim 10500$$

$$Y = 0.06X_1 + 0.12X_2 + 0.06X_3 + 0.19X_4 + 0.05X_5 + 0.1X_6$$



$$\lambda = 10^{-5}$$

$$Y = 1.53X_1 + 1.94X_3$$

$$\lambda \sim 0.01$$

$$Y = 0.76X_2 + 0.65X_4$$

Эксперименты с одинаковыми и зависимыми признаками

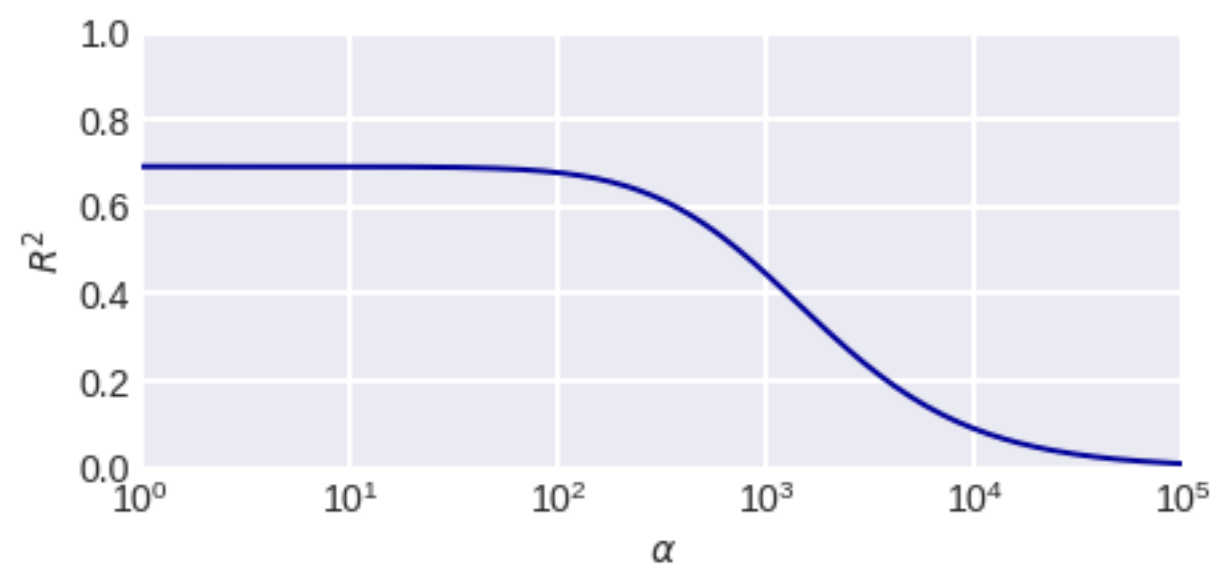
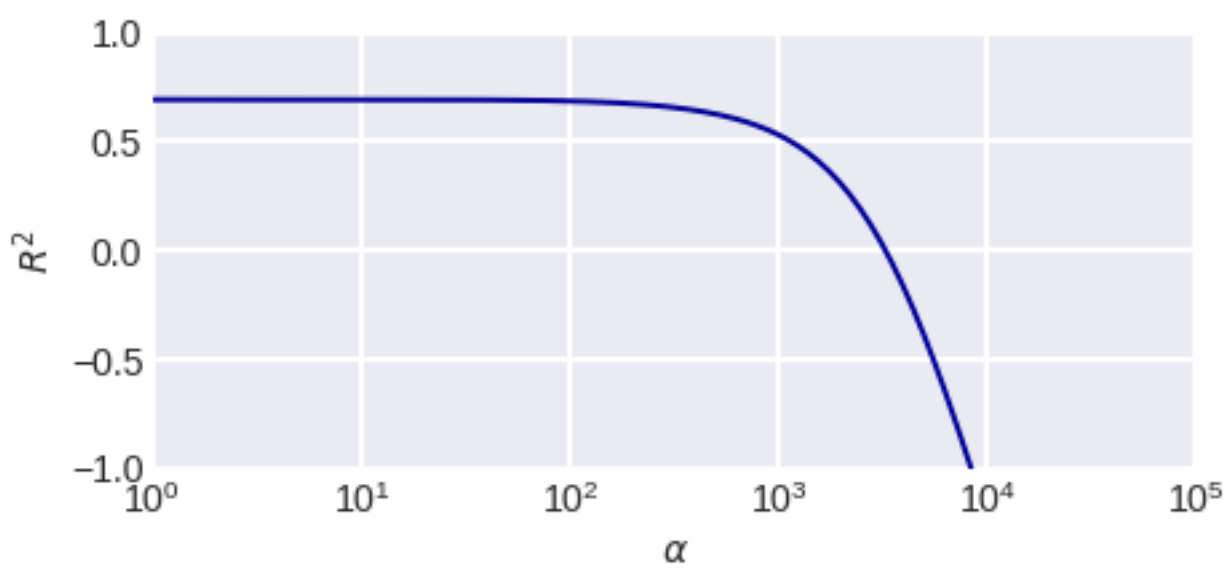
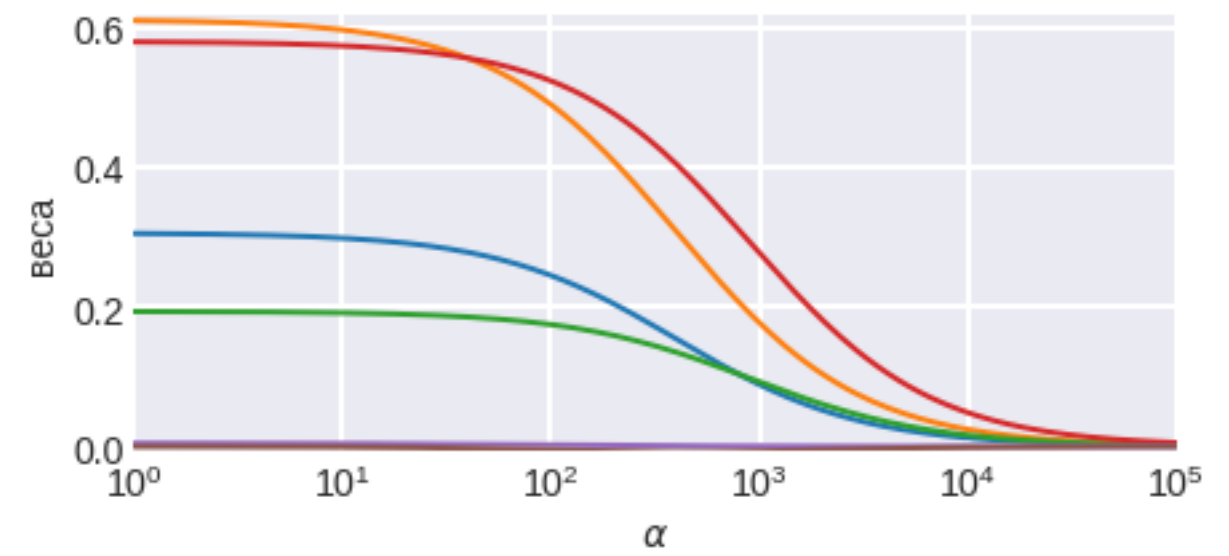
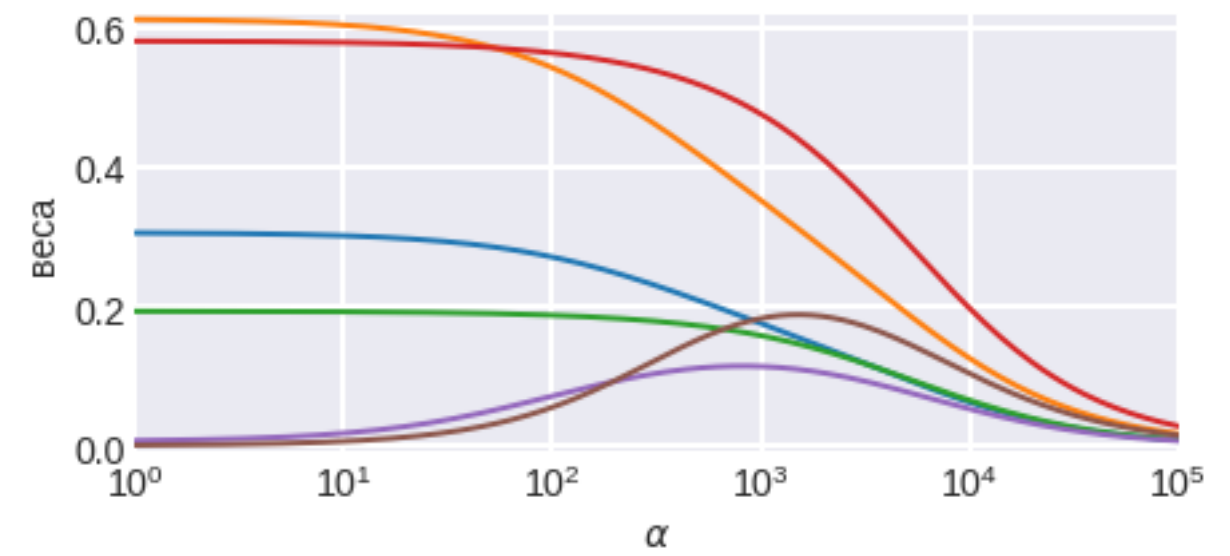
веса зависят от масштаба признаков

при сильной регуляризации меняется распределение весов зависимых признаков

Пусть $Y = 4X_1, X_1=X_2$

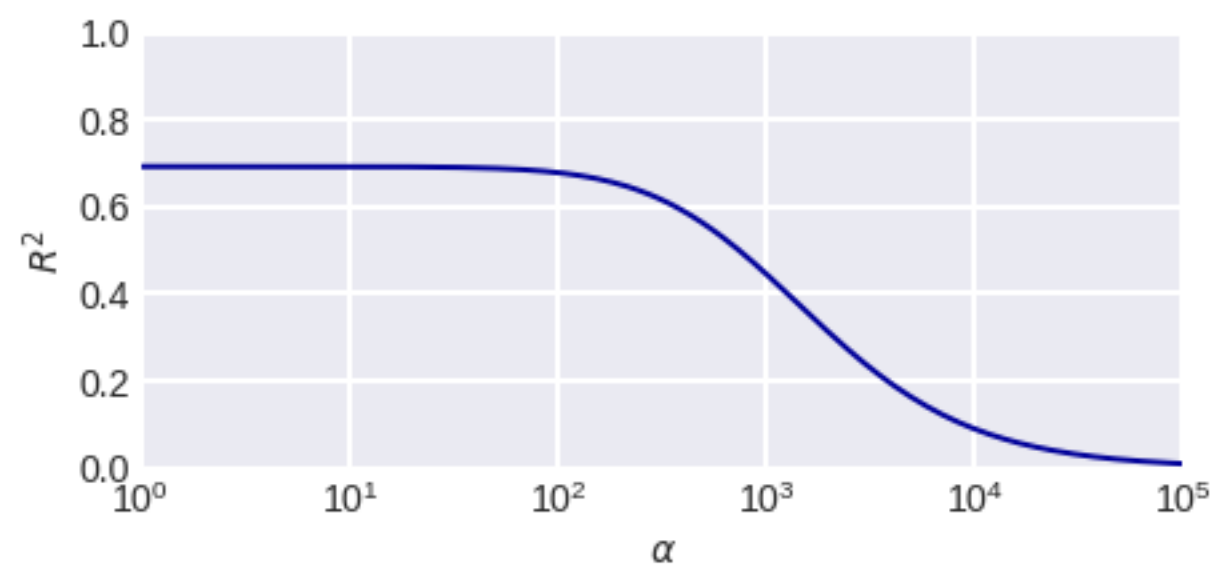
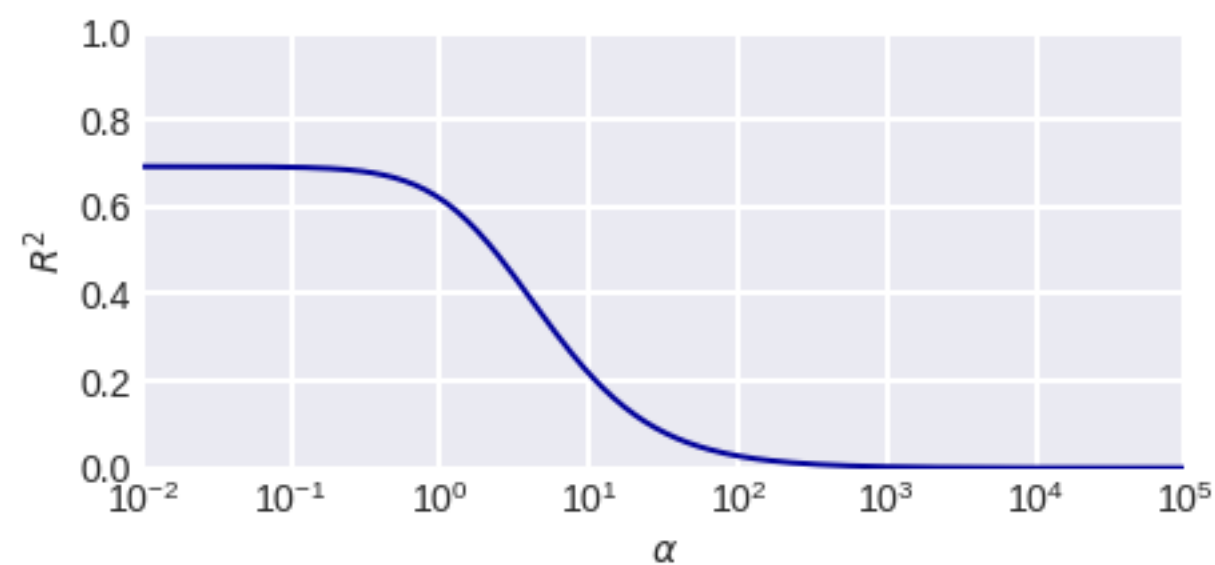
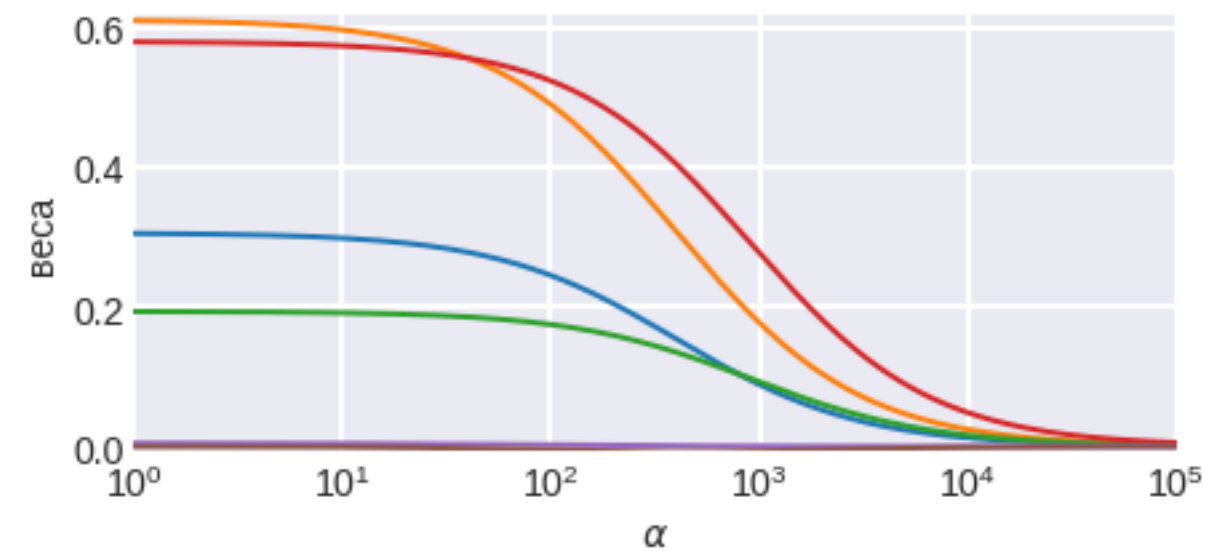
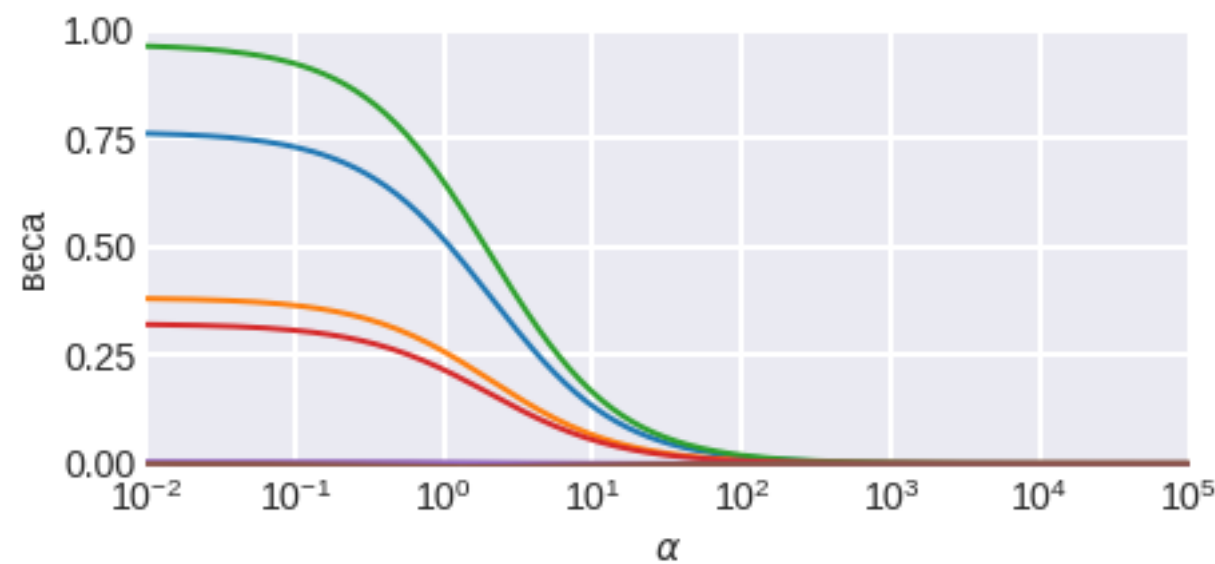
w_1	w_2	$\ w \ _1$	$\ w \ _2^2$
5	- 1	6	26
4	0	4	16
3	1	4	10
2	2	4	8

Эксперименты с одинаковыми и зависимыми признаками: L_2 -регуляризация



fit_intercept=True

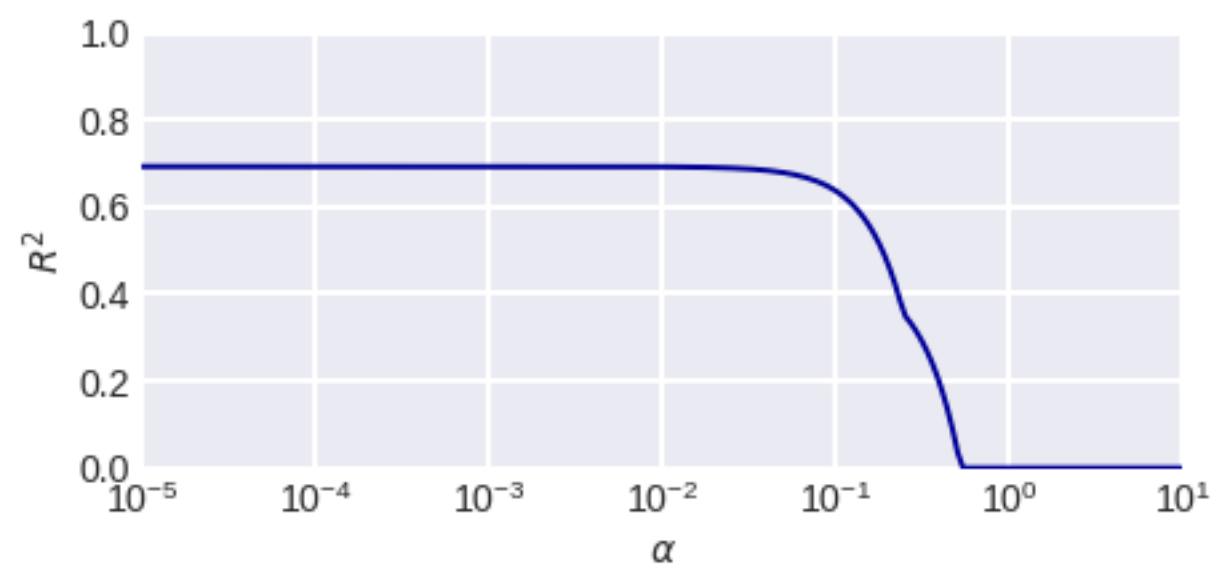
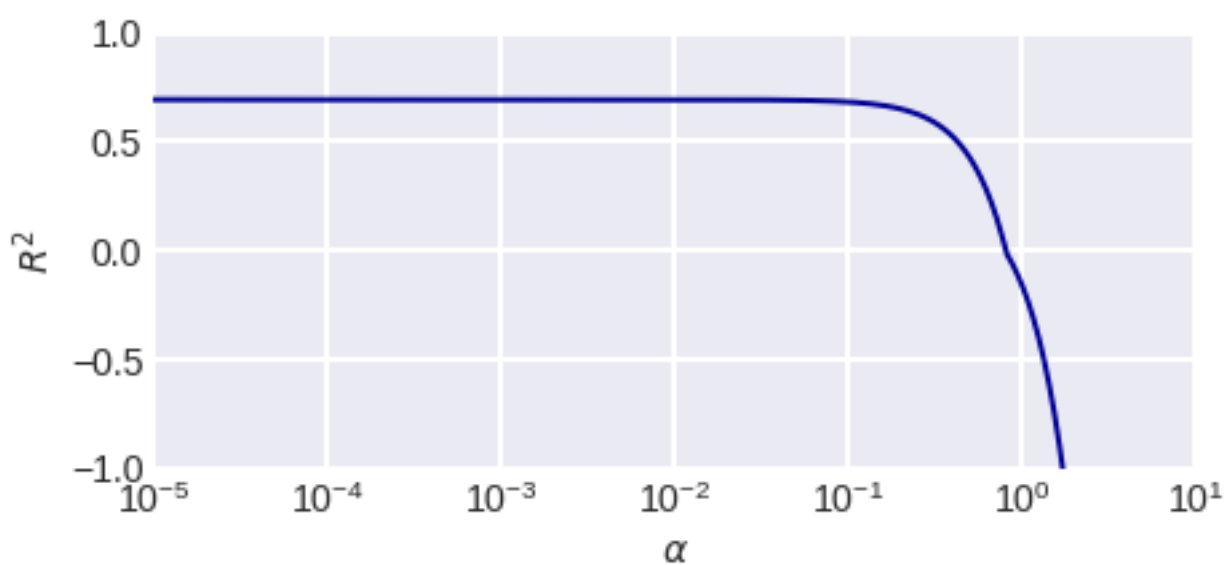
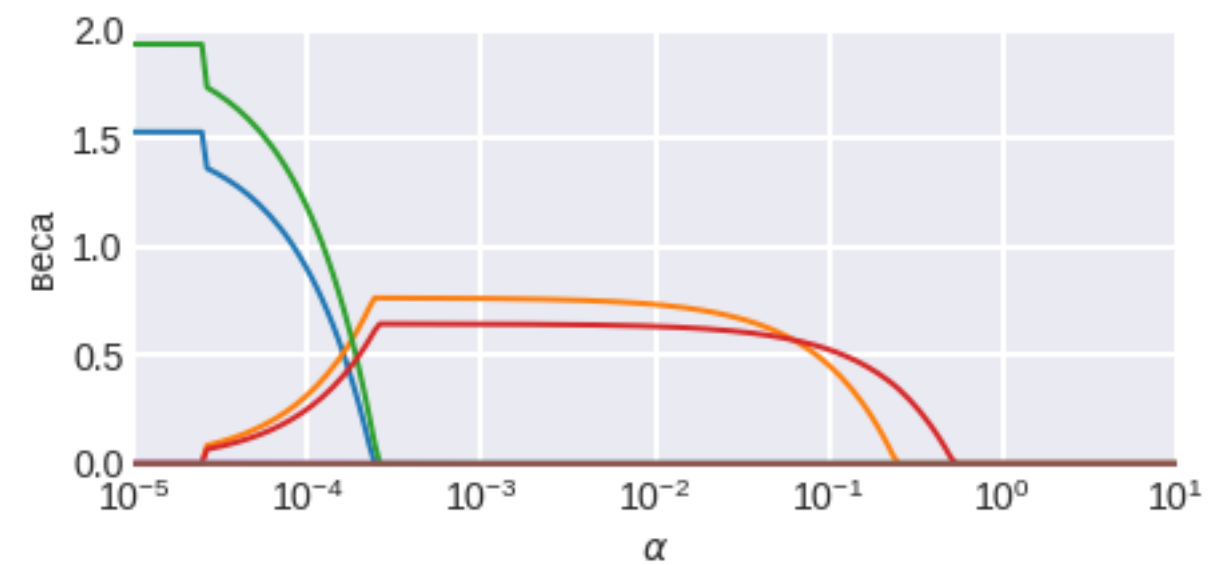
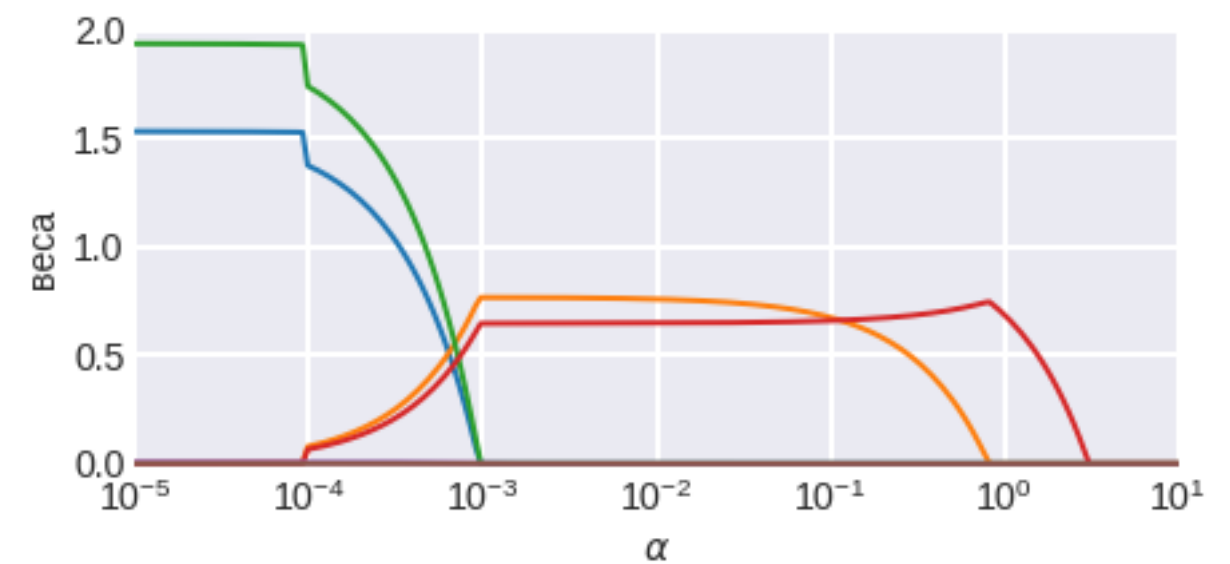
Эксперименты с одинаковыми и зависимыми признаками: L_2 -регуляризация



`fit_intercept=True, normalize=True`

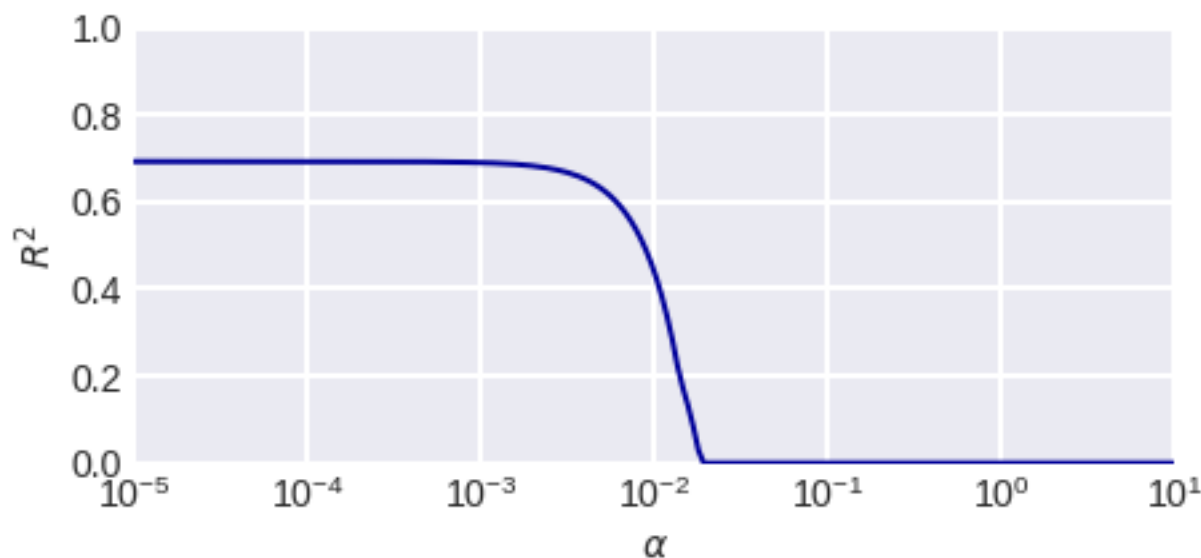
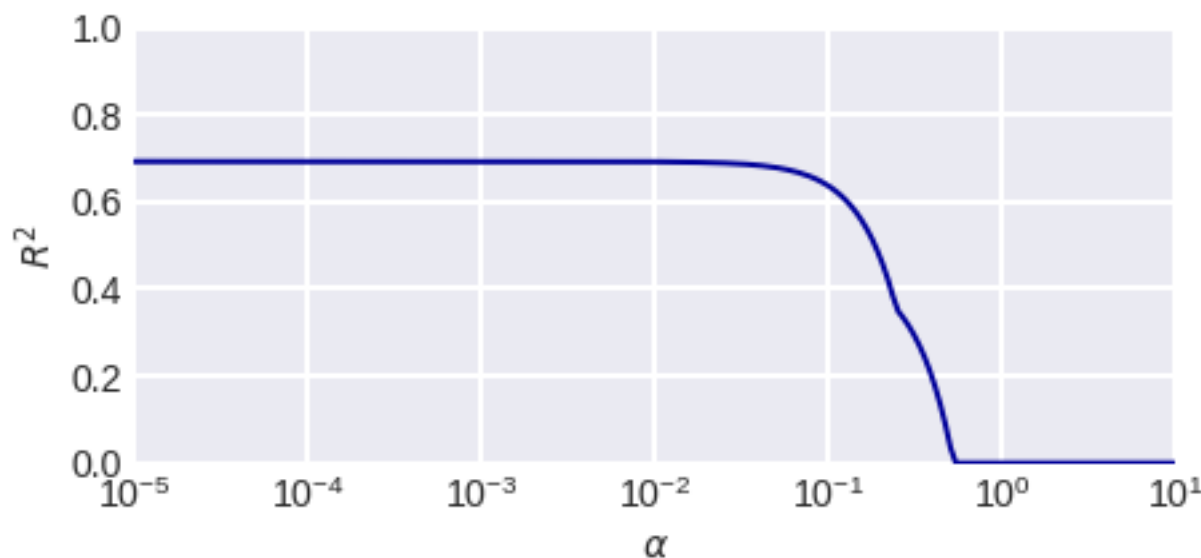
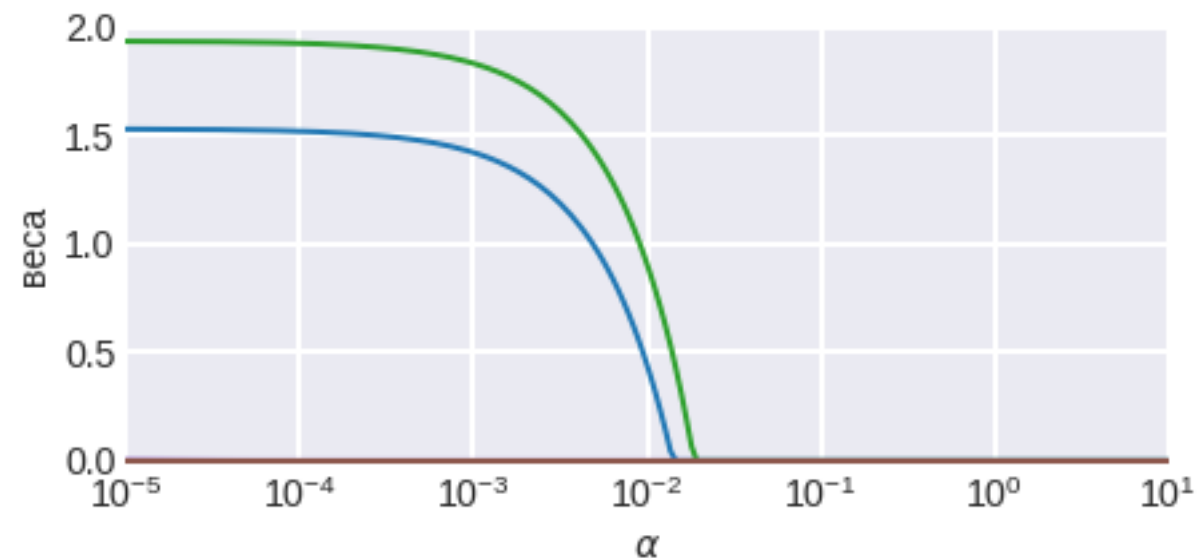
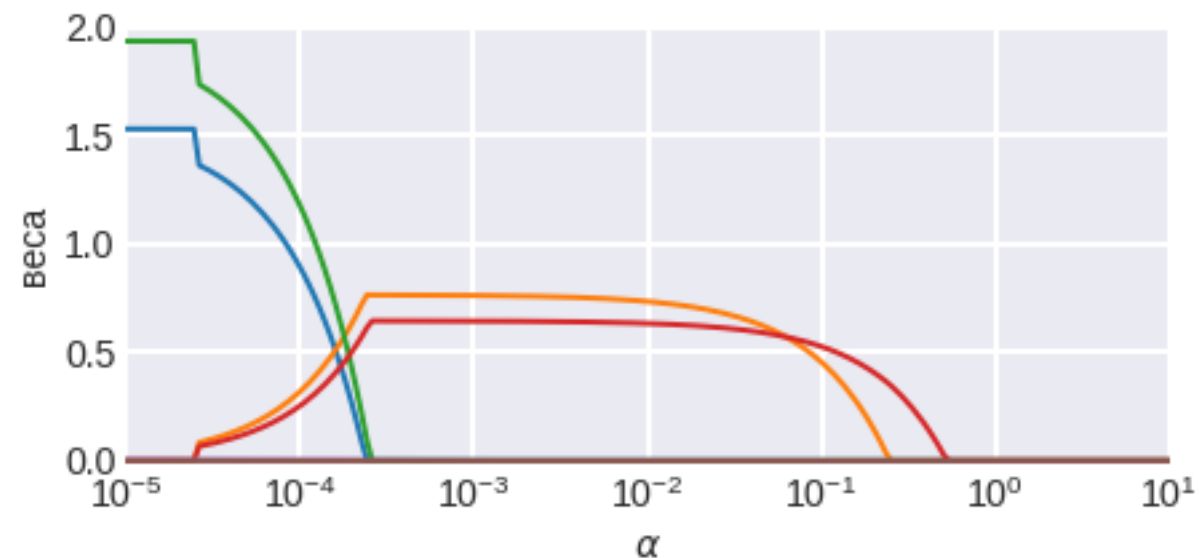
`fit_intercept=True`

Эксперименты с одинаковыми и зависимыми признаками: L_1 -регуляризация



fit_intercept=True

Эксперименты с одинаковыми и зависимыми признаками: L₁-регуляризация



fit_intercept=True

fit_intercept=True, normalize=True

Эксперименты с одинаковыми и зависимыми признаками

Часто важно

- **использовать свободный член**
- **предварительно нормировать данные**

Семейство регуляризированных линейных методов**Ridge**

$$\|y - Xw\|_2^2 + \lambda \|w\|_2^2 \rightarrow \min_w$$

LASSO (Least Absolute Selection and Shrinkage Operator)

$$\|y - Xw\|_2^2 + \lambda \|w\|_1 \rightarrow \min_w$$

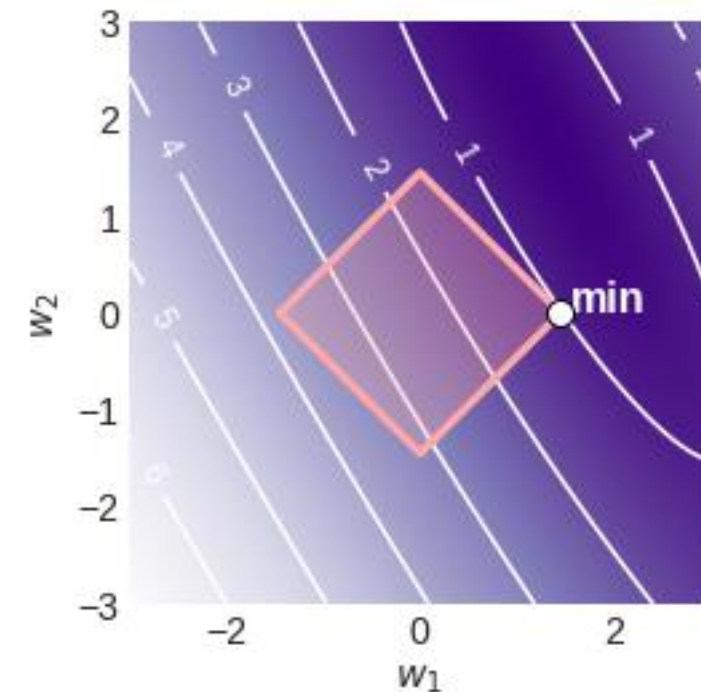
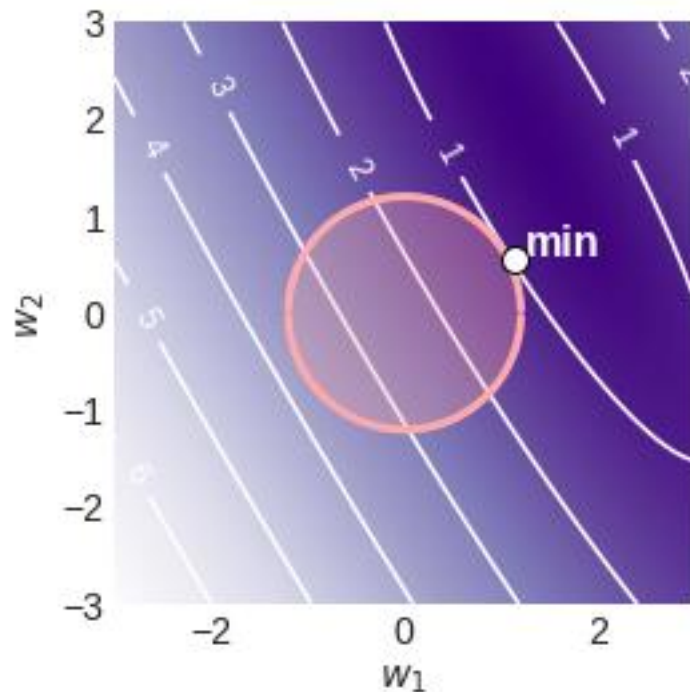
Elastic Net = LASSO + Ridge

$$\|y - Xw\|_2^2 + \lambda_1 \|w\|_1 + \lambda_2 \|w\|_2^2 \rightarrow \min_w$$

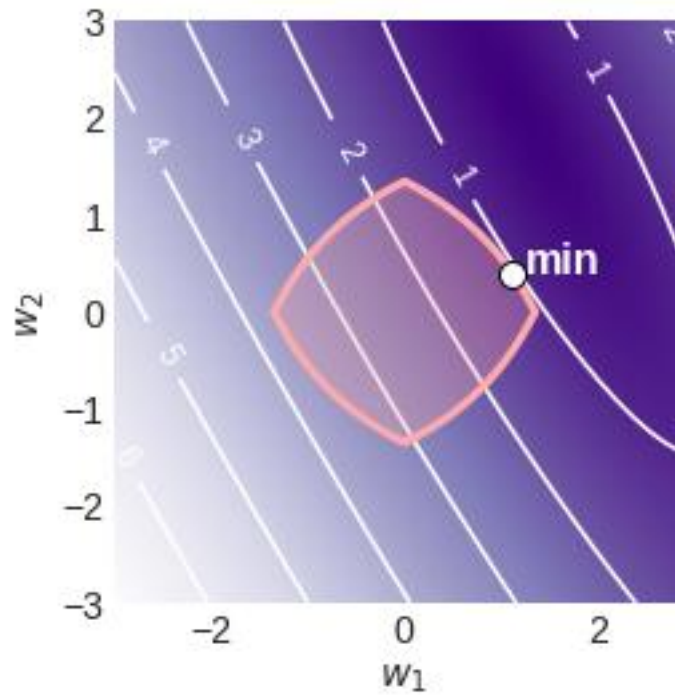
Геометрический смысл Ridge, LASSO и Elastic Net

$$\sum_{i=1}^m \left(y_i - w_0 - \sum_{j=1}^n w_j x_{ij} \right)^2 \rightarrow \min_w, \quad \sum_{j=1}^n w_j^2 \leq s$$

$$\sum_{i=1}^m \left(y_i - w_0 - \sum_{j=1}^n w_j x_{ij} \right)^2 \rightarrow \min_w, \quad \sum_{j=1}^n |w_j| \leq s$$



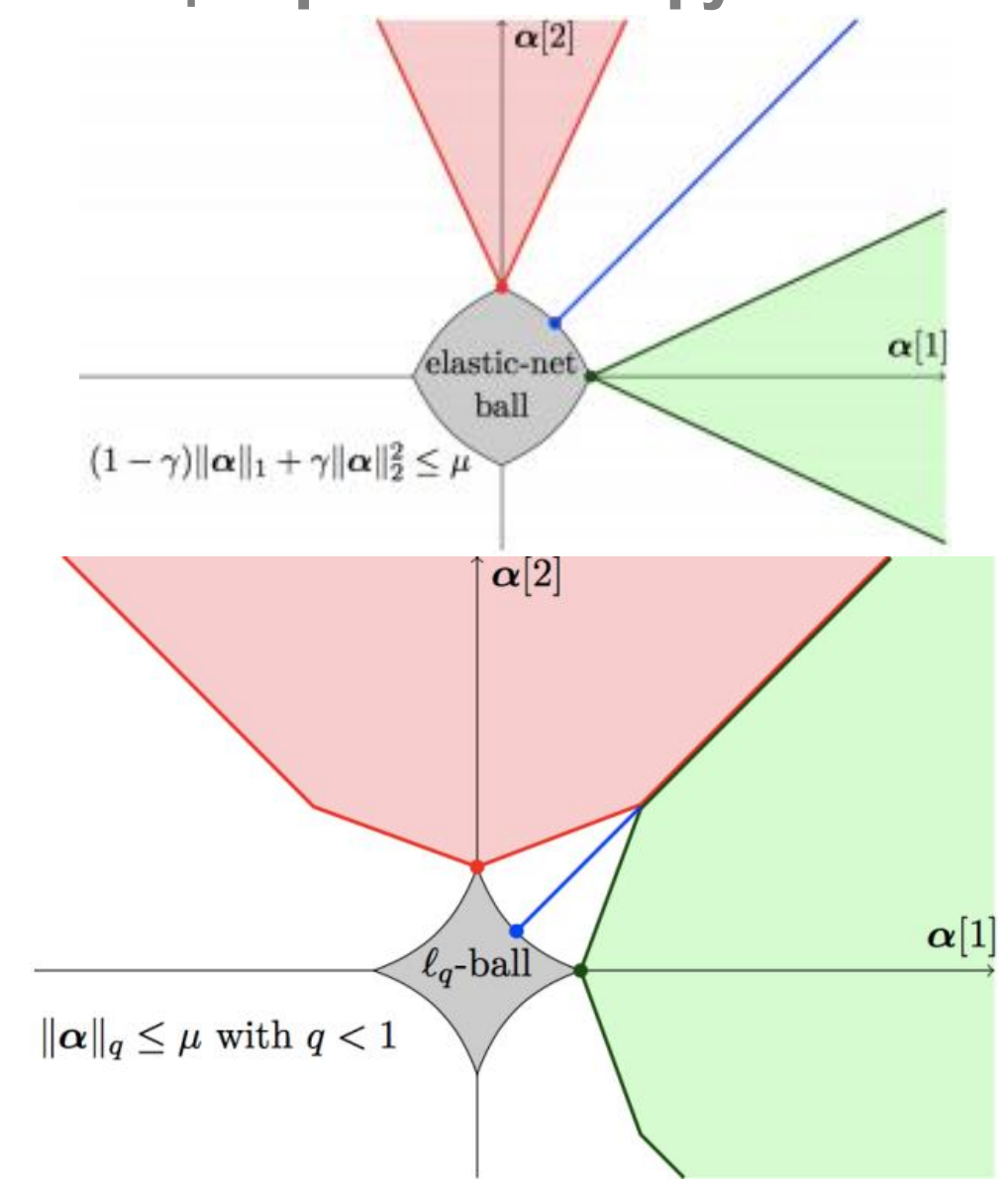
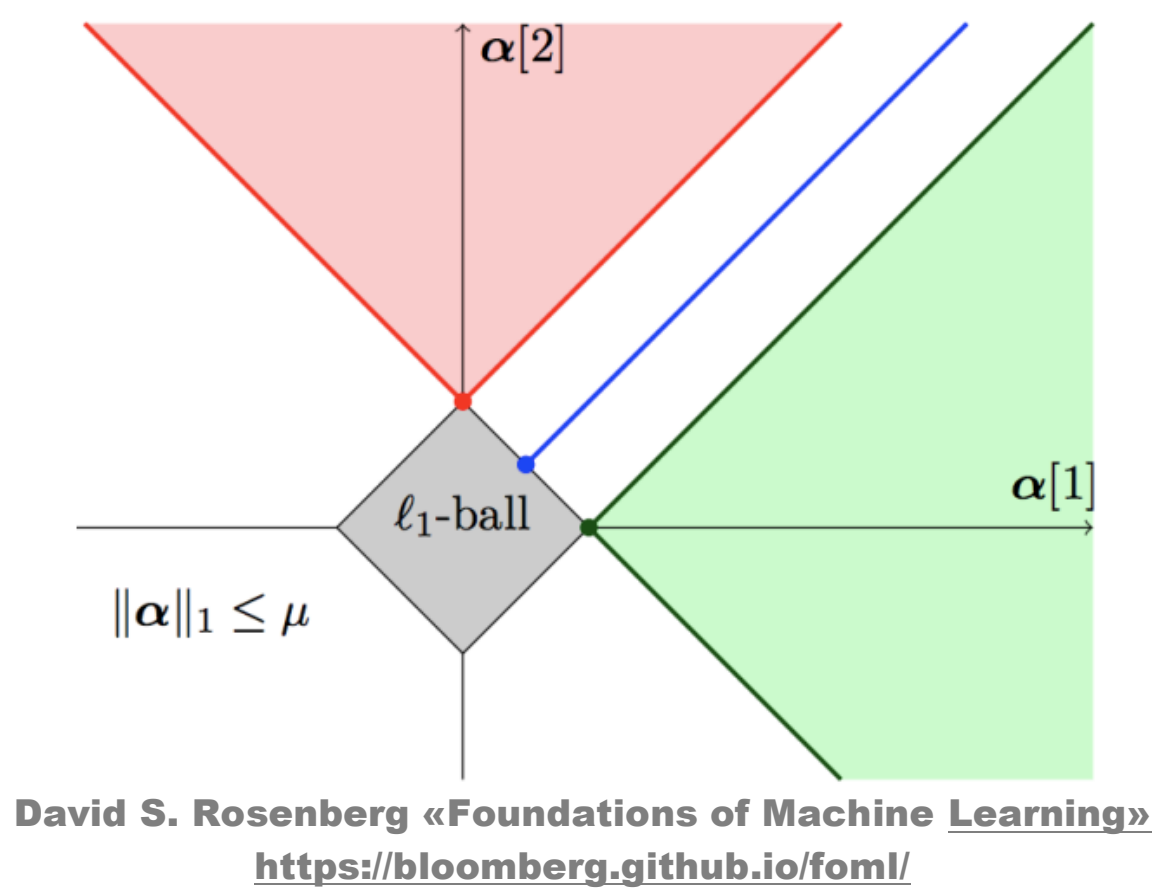
Геометрический смысл Ridge, LASSO и Elastic Net



на практике часто модель и не может зависеть от небольшого числа переменных

Эффект разреженности

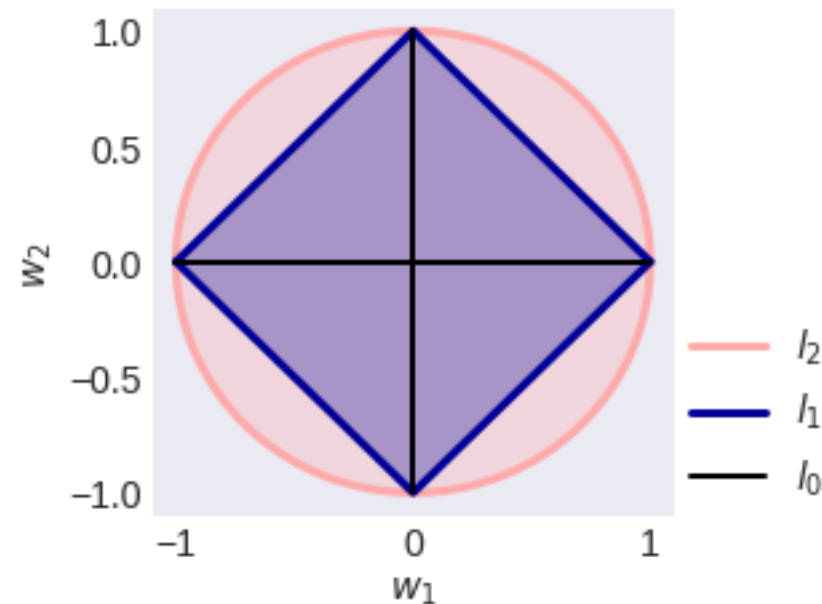
если линии уровня оптимизируемой функции – concentric окружности...



Почему L1-норма \Rightarrow разреженность

1. Больше вероятность, что линии уровней функции ошибки касаются области ограничений в точках с нулевыми координатами, **см. рис.**

2. L1-норма больше похожа на L0, чем L2



$$\|w\|_0 = |\{t \mid w_t \neq 0\}|$$

**При увеличении коэффициента регуляризации веса стремятся к нулю
Обеспечивается автоматическая селекция признаков!**

Регуляризация \Rightarrow упрощение

Соблюдение принципа Оккама

регуляризация \Rightarrow зануление коэффициентов \Rightarrow упрощение модели

**В целом, неверно, что чем меньше коэффициентов, тем проще модель,
но у нас линейная модель..**

потом будет обоснование регуляризации с помощью вероятностных предположений

Проблема вырожденности / плохой обусловленности матрицы

$$w = (X^T X)^{-1} X^T y$$

Решения:

1. Регуляризация
- 2. Селекция (отбор) признаков**
- 3. Уменьшение размерности (в том числе, PCA)**
- 4. Увеличение выборки**

Селекция признаков в линейной регрессии **отдельная тема**

Какие признаки включить в модель

$$a(X_1, \dots, X_n) = w_0 + w_1 X_1 + \dots + w_n X_n$$

Маленький обзор стратегий:

- 1 стратегия – перебор** – умный перебор подмножества признаков
- 2 стратегий – оценка** – оценка качества признаков (фильтры)
- 3 стратегия – автомат** – встроенные методы (ex: LASSO)

Обоснование необходимости селекции

- **Проблема вырожденности в линейной регрессии**
 - **Проблема «почти дубликатов»**
 - **Уменьшение модели и интерпретация**
 - **Уменьшение стоимости данных**

Проблема вырожденности матрицы

$$w = (X^T X)^{-1} X^T y$$

Решения:

1. Регуляризация
2. Селекция (отбор) признаков
3. Уменьшение размерности (в том числе, PCA)
4. Увеличение выборки

	x1	x2	x3	y		x1-x2	y
0	0.44	0.62	0.51	-0.25	0	-0.18	-0.25
1	0.03	0.53	0.07	-0.51	1	-0.50	-0.51
2	0.55	0.13	0.43	0.41	2	0.42	0.41
3	0.44	0.51	0.10	0.04	3	-0.07	0.04
4	0.42	0.18	0.13	0.12	4	0.24	0.12
5	0.33	0.79	0.60	-0.45	5	-0.46	-0.45



обоснование необходимости аналогично селекции

Линейная регрессия: градиентный метод обучения

недостатки прямого...

работа с большими матрицами (тем более обращение)

В лекции «оптимизация»...

$$\frac{1}{2} \sum_{i=1}^m (a(x_i | w) - y_i)^2 \rightarrow \min$$

$$w^{(t+1)} = w^{(t)} - \eta \sum_{i=1}^m (a(x_i | w^{(t)}) - y_i) \frac{\partial a(x_i | w^{(t)})}{\partial w}$$

Gradient Descent

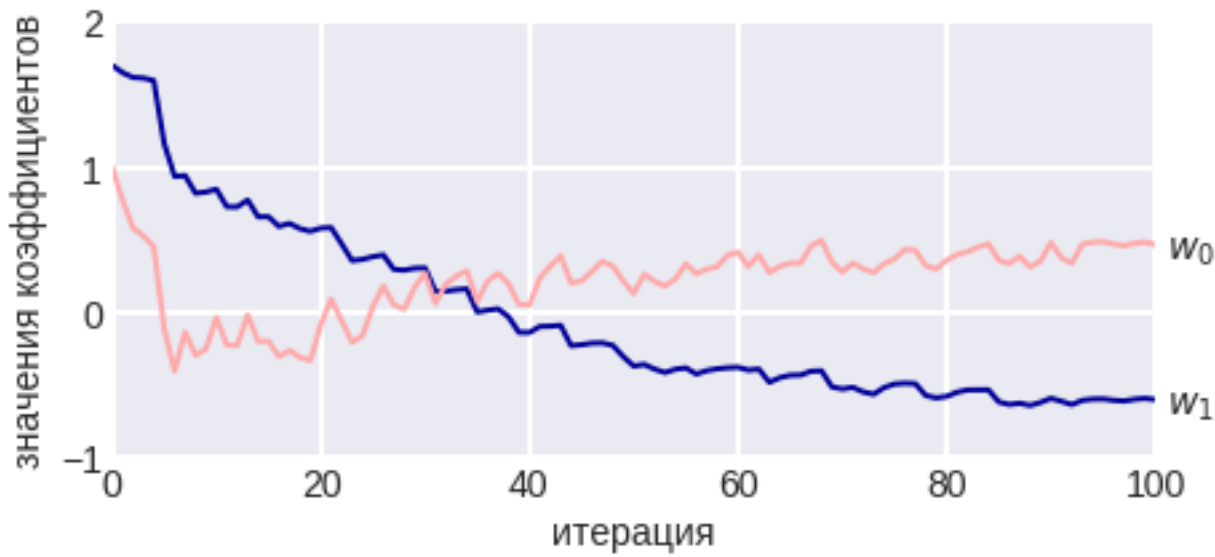
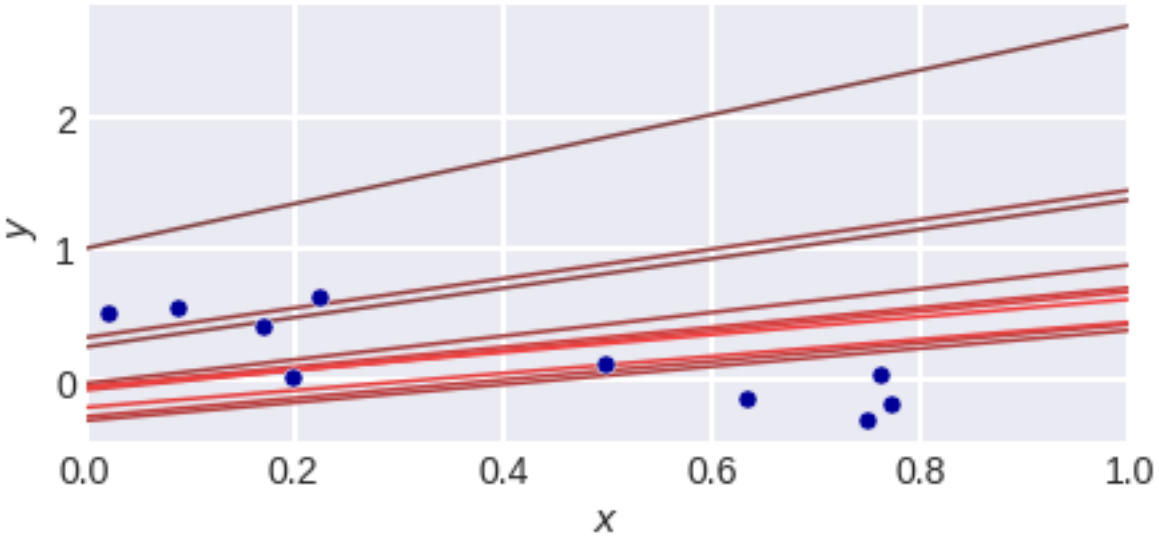
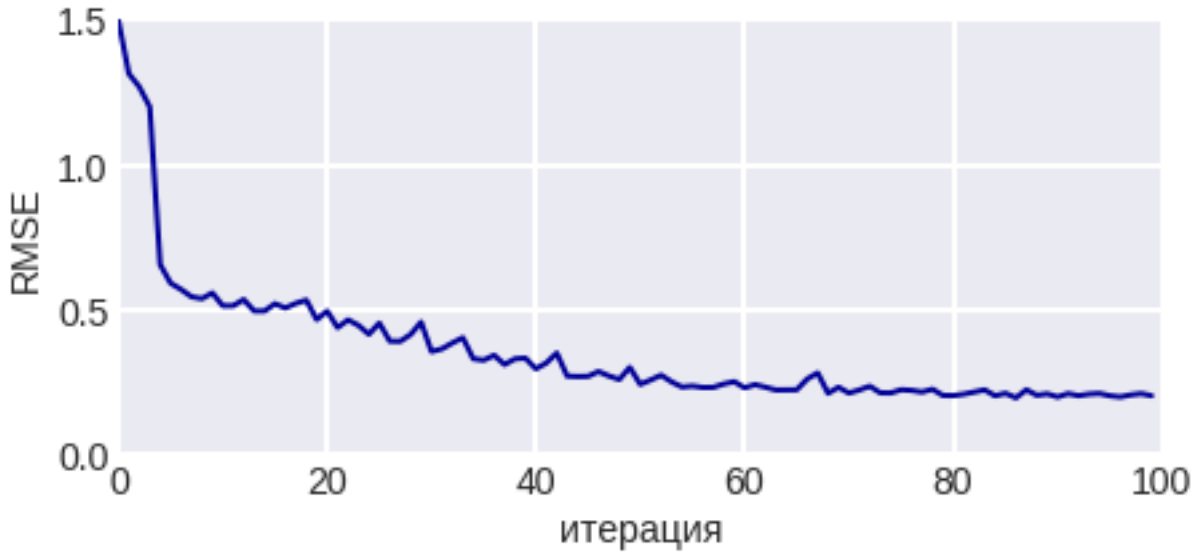
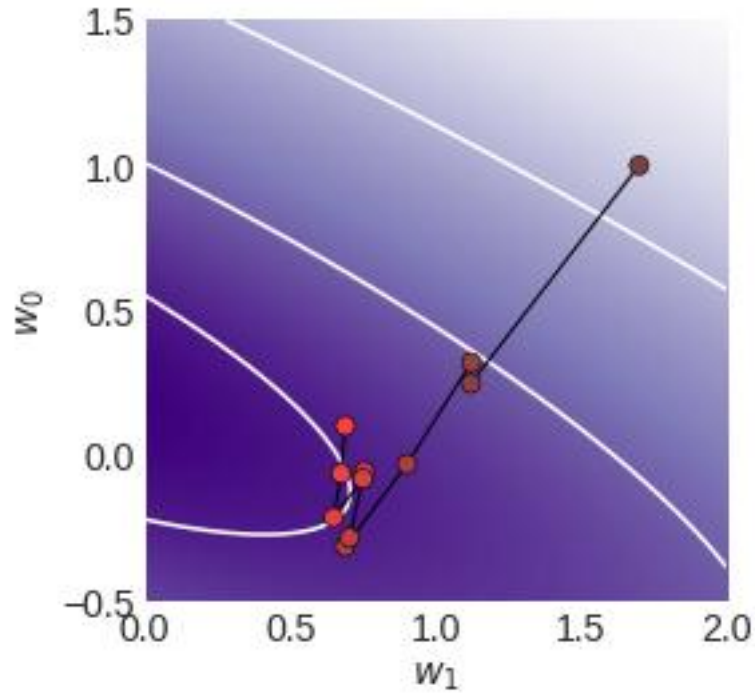
$$a(x | w) = w^T x$$

$$w^{(t+1)} = w^{(t)} - \eta \sum_{i=1}^m (a(x_i | w^{(t)}) - y_i) x_i$$

Stochastic Gradient Descent

$$w^{(t+1)} = w^{(t)} - \eta_t (a(x_i | w^{(t)}) - y_i) x_i$$

Линейная регрессия: градиентный метод обучения



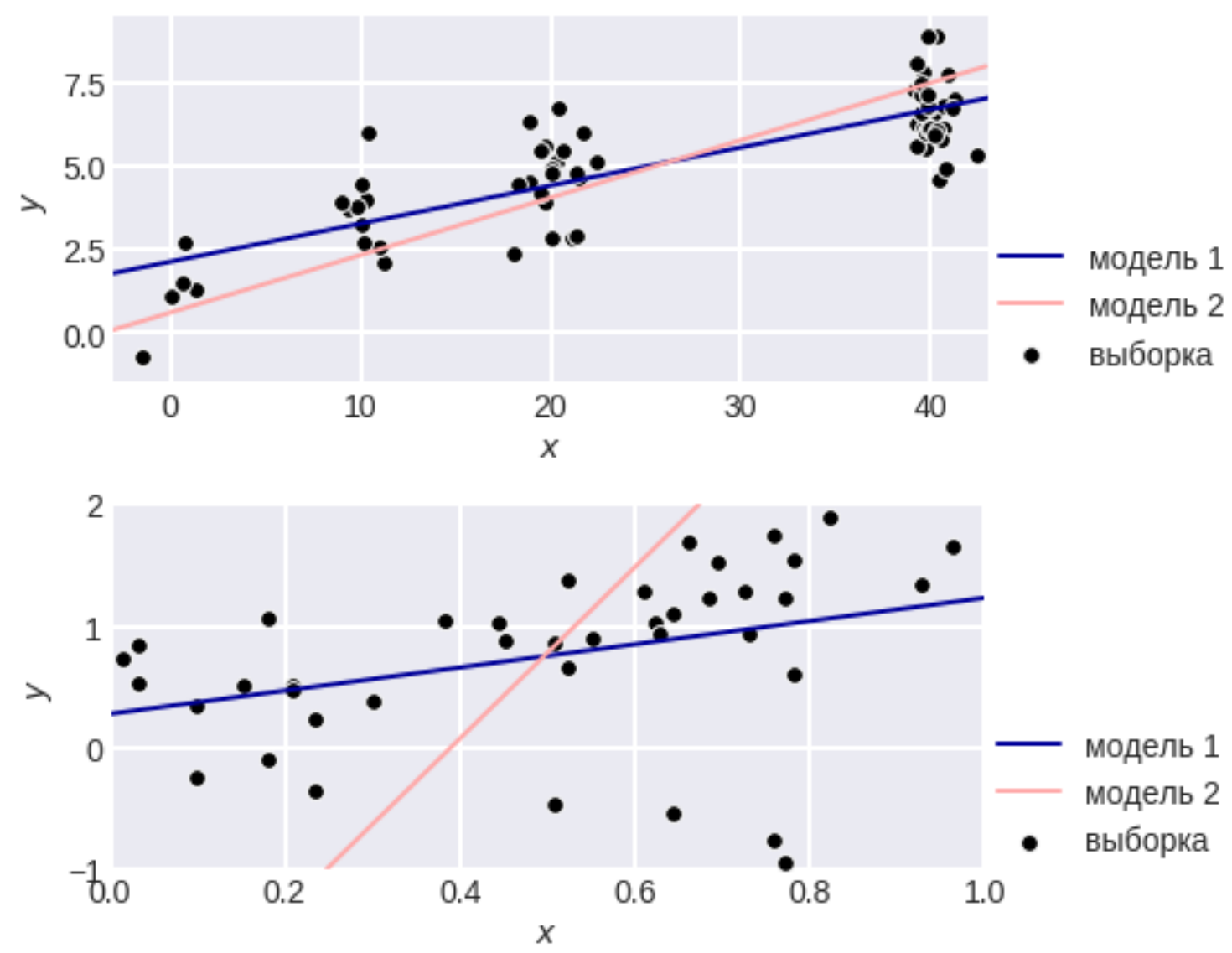
Реализация в **scikit-learn**`sklearn.linear_model.Ridge`

<code>alpha=1.0</code>	Коэффициент регуляризации, больше – сильнее (в отличие от других функций)
<code>fit_intercept=True</code>	Использовать ли свободный член
<code>normalize=False</code>	Нормализация данных Игнорируется без свободного члена
<code>solver="auto"</code>	Метод оптимизации "auto", "svd", "cholesky", "lsqr", "sparse_cg", "sag", "saga"
<code>copy_X=True, max_iter=None, tol=0.001, random_state=None</code>	

```
sklearn.linear_model.Lasso(alpha=1.0, fit_intercept=True, normalize=False,
precompute=False, copy_X=True, max_iter=1000, tol=0.0001, warm_start=False,
positive=False, random_state=None, selection="cyclic")
```

```
sklearn.linear_model.ElasticNet(alpha=1.0, l1_ratio=0.5, fit_intercept=True,
normalize=False, precompute=False, max_iter=1000, copy_X=True, tol=0.0001,
warm_start=False, positive=False, random_state=None, selection="cyclic")
```

Две регрессии



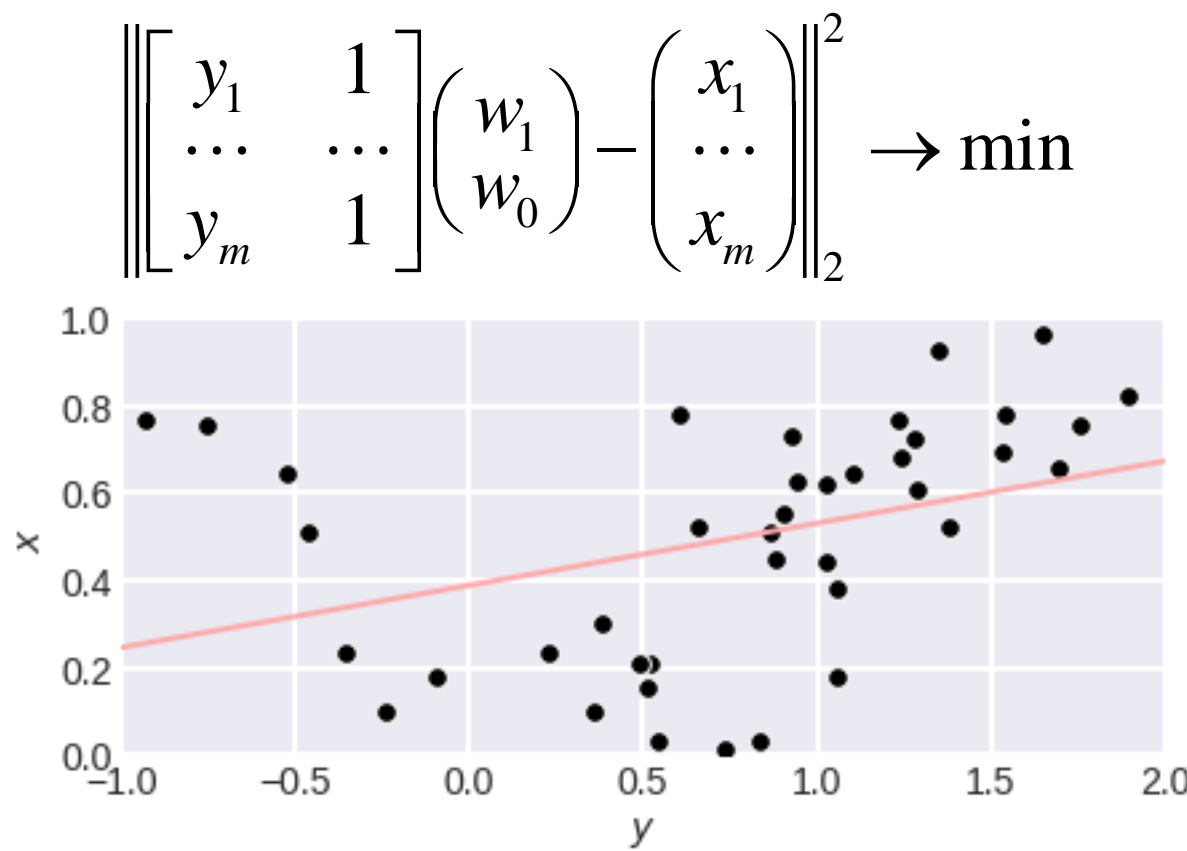
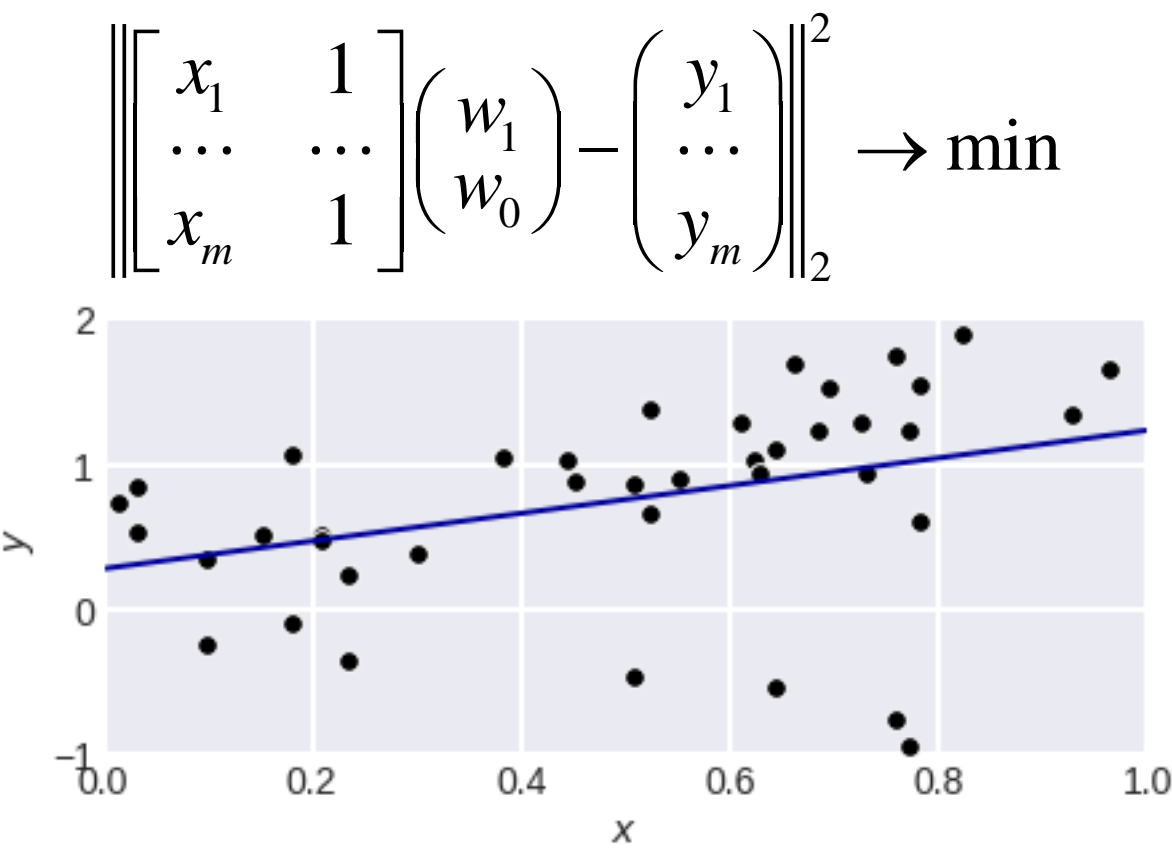
Чем отличаются модели 1 и 2?

Две регрессии: $y(x)$ vs $x(y)$

разные задачи $y(x)$ и $x(y)$, хотя зависимость линейная

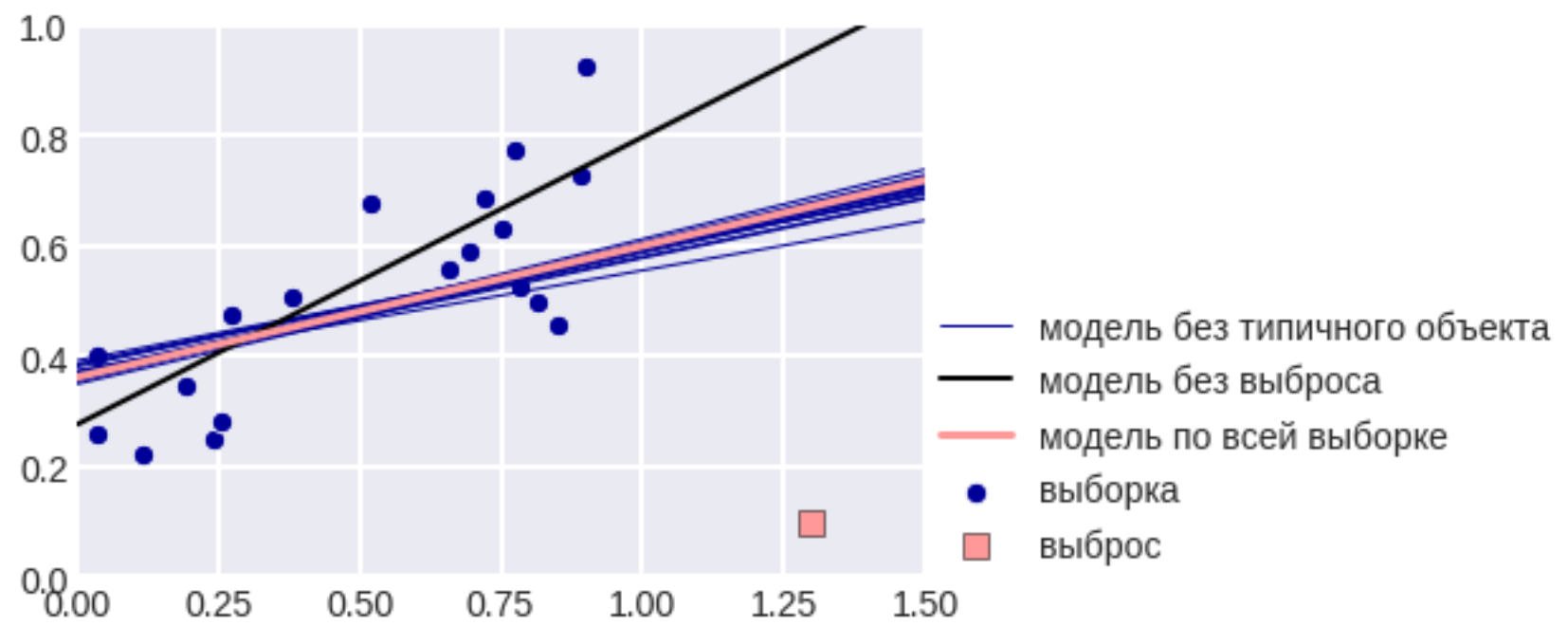
$$Y = w_0 + w_1 X_1$$

$$X_1 = w_0 + w_1 Y$$



есть и «промежуточная стратегия» – **дальше РСА**

Линейная регрессия – неустойчивость к выбросам



Ошибка с весами

Если у каждого объекта есть цена ошибки...

$$\sum_{i=1}^m v_i \left(y_i - w^T x_i \right)^2 + \dots = \sum_{i=1}^m \left(\sqrt{v_i} y_i - w^T (\sqrt{v_i} x_i) \right)^2 + \dots \rightarrow \min_{v_i > 0}$$

небольшая переформулировка задачи:

$$\{(x_1, y_1), \dots, (x_m, y_m)\} \rightarrow \{(\sqrt{v_1} x_1, \sqrt{v_1} y_1), \dots, (\sqrt{v_m} x_m, \sqrt{v_m} y_m)\}$$

$$(y - Xw)^T V (y - Xw) \sim \|V^{1/2} y - V^{1/2} Xw\|_2^2 \rightarrow \min_w$$

$$w = (X^T V X)^{-1} X^T V y$$

Ошибка с весами

$$\{(x_1, y_1), \dots, (x_m, y_m)\} \rightarrow \{(\sqrt{v_1} x_1, \sqrt{v_1} y_1), \dots, (\sqrt{v_m} x_m, \sqrt{v_m} y_m)\}$$

Как реализовать:

- 1) перейти к новым данным («испорченными весам»)**
- 2) если веса целые числа – можно продублировать объекты**
- 3) если веса из отрезка [0, 1] – при численном градиентном решении можно выбирать следующий объект с соответствующей вероятностью**

Устойчивая регрессия (Robust Regression)

0. Инициализация весов объектов

$$v = (v_1, \dots, v_m) = (1/m, \dots, 1/m)$$

1. Цикл

1.1. Настроить алгоритм, учитывая веса объектов

$$a = \text{fit}(\{x_i, y_i, v_i\})$$

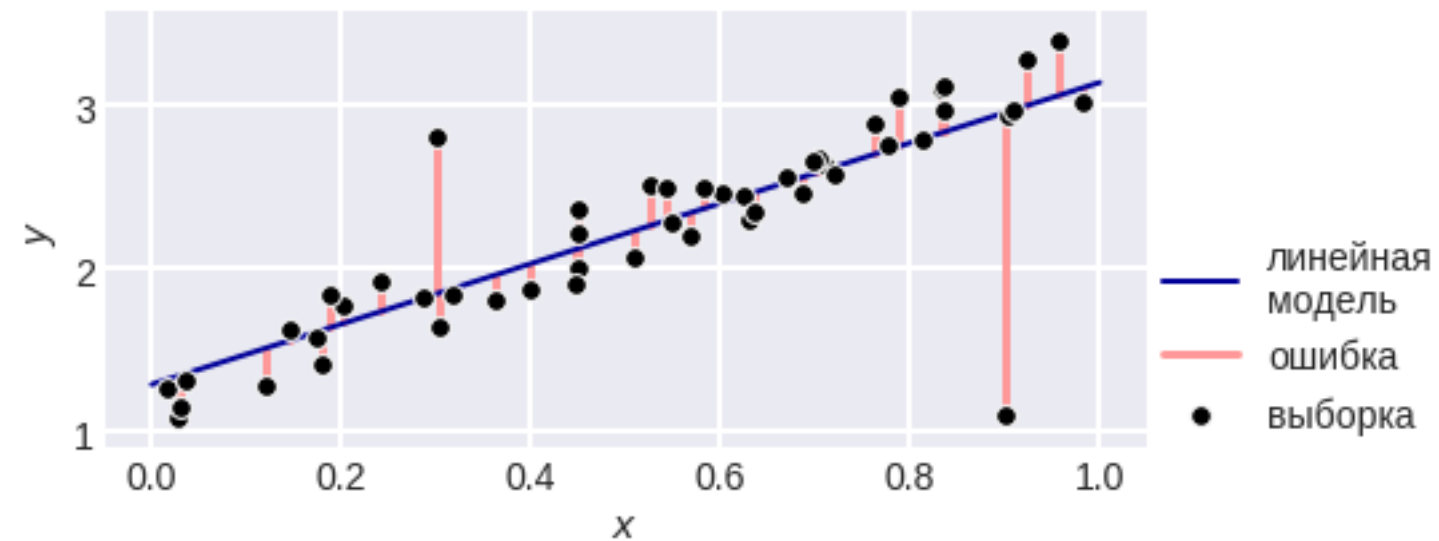
1.2. Вычислить ошибки на обучении

$$\varepsilon_i = a(x_i) - y_i$$

1.3. Пересчитать веса объектов

$$v_i = \exp(-\varepsilon_i^2)$$

нормировать на сумму

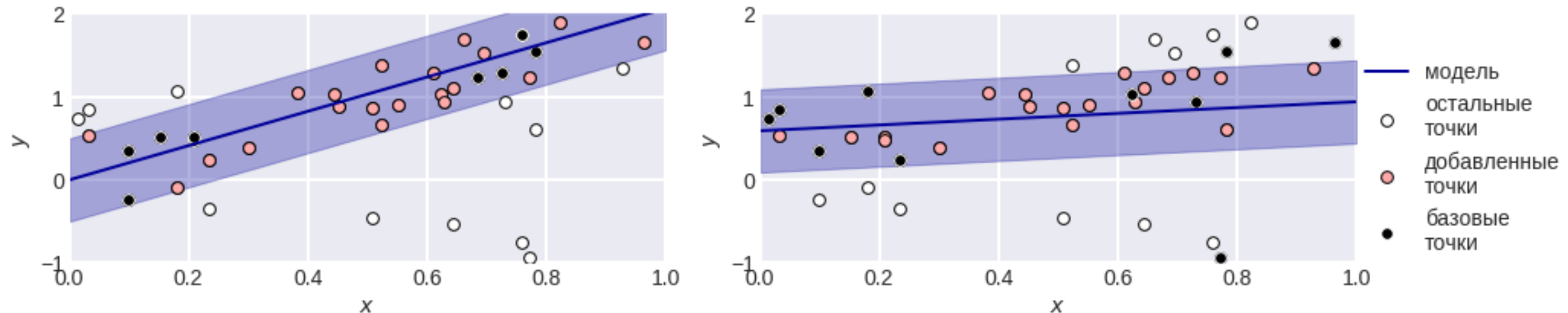


можно использовать любую
регрессионную модель

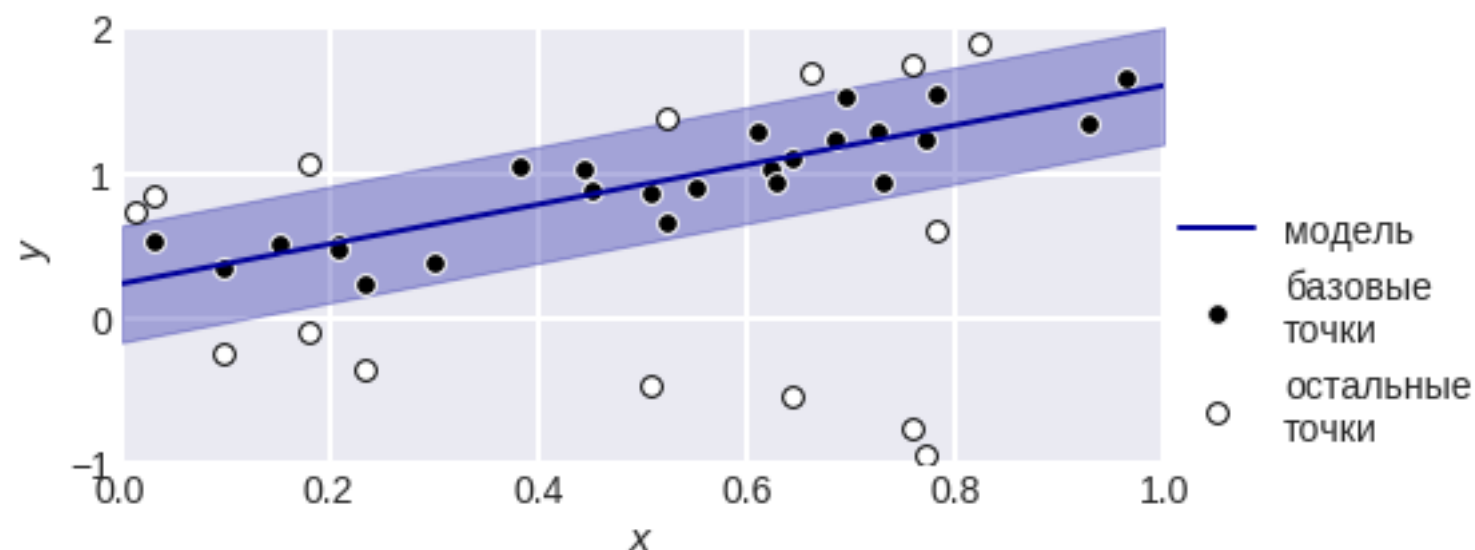
при пересчёте весов можно использовать
другую невозрастающую функцию;
можно (иногда нужно) нормировать

RANdom SAmple Consensus (RANSAC)

- **несколько раз**
 - **выбрать случайное подмножество точек – базовое (inliers)**
 - **обучить модель на базовом подмножестве**
 - **найти все точки, которые хорошо предсказываются моделью**
например, ошибка не больше ε
 - **пополнить ими базовое множество**
 - **(если добавили много) переобучить модель на новом множестве**
- **выбрать модель с наименьшей ошибкой**



Минутка кода: RANSAC в scikit-learn



```
from sklearn.linear_model import RANSACRegressor
# Robustly fit linear model with RANSAC algorithm
ransac = RANSACRegressor()
ransac.fit(x[:, np.newaxis], y)
inlier_mask = ransac.inlier_mask
outlier_mask = np.logical_not(inlier_mask)
```

Минутка кода: RANSAC в scikit-learn

```
sklearn.linear_model import RANSACRegressor
```

<code>base_estimator=None</code>	Базовый алгоритм (по умолчанию – линейная регрессия)
<code>min_samples=None</code>	Число / доля базовых объектов ($n+1$)
<code>residual_threshold=None</code>	Порог для пополнения базового множества ($MAD(y)$)
<code>max_trials=100</code>	Число итераций
<code>stop_n_inliers</code>	Остановить вычисления, если найдено столько базовых точек
<code>loss="absolute_loss"</code>	Как оценивать ошибку

```
is_data_valid=None, is_model_valid=None, max_skips=inf,  
stop_score=inf, stop_probability=0.99, lossrandom_state=None
```

Реализация в scikit-learn немного отличается от некоторых описаний:

Качество = число базовых (inliers) объектов;
лучшая модель выбирается по числу базовых,
если у нескольких моделей число совпадает,
тогда сравнивается ошибка на всей выборке.

Плюсы и минусы линейных алгоритмов

- + простой, надёжный, быстрый, популярный метод**
 - + интерпретируемость (\Rightarrow нахождение закономерностей)**
 - + интерполяция и экстраполяция**
 - + может быть добавлена нелинейность, с помощью генерации новых признаков**
(дальше – это можно автоматизировать)
 - + хороши для теоретических исследований (в Ridge есть явная формула)**
 - + коэффициенты асимптотически нормальны**
(можно тестировать гипотезы о влиянии признаков)
 - + глобальный минимум в оптимизируемом функционале**
 - линейная гипотеза вряд ли верна**
 - в теоретическом обосновании ещё предполагается нормальность ошибок**
(зависит от функции ошибок)
 - «страдает» из-за выбросов**
 - признаки в одной шкале и однородные (см. успешные примеры)**
 - проблема коррелированных признаков**
- \Rightarrow необходимость регуляризации, селекции, PCA, data \uparrow