

Машинное обучение и анализ данных

Линейные методы

Дьяконов А.Г.

**Московский государственный университет
имени М.В. Ломоносова (Москва, Россия)**



Линейные методы



Линейная регрессия

Гипотеза о линейной зависимости целевой переменной

Ищем решение в виде

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

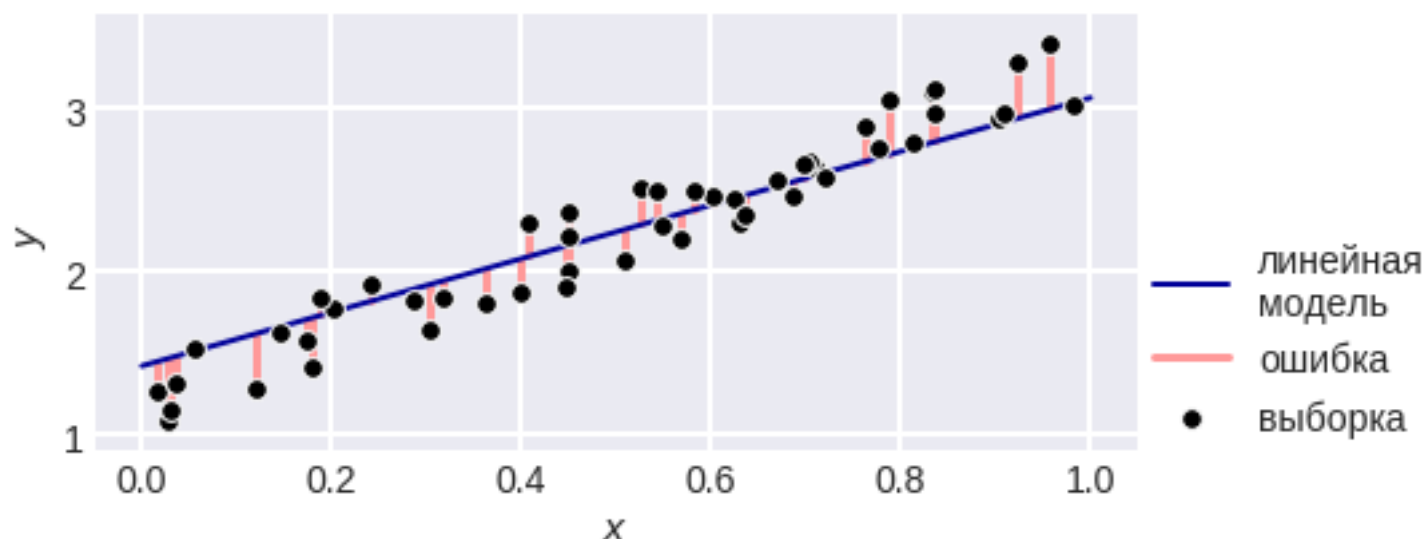
Практика:

- 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 выдачи

...

Линейная регрессия от одной переменной

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



обучение: $\{(x_1, y_1), \dots, (x_m, y_m)\},$

$$x_i \in \mathbb{R},$$

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

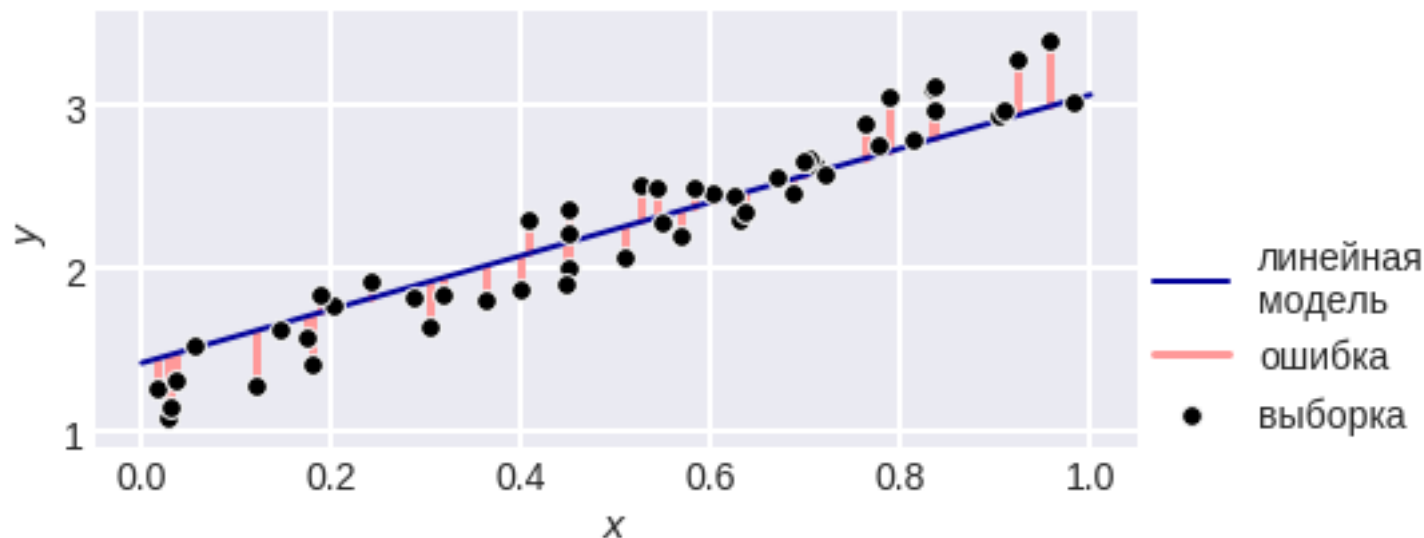
невязки/отклонения (residuals):

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

Линейная регрессия от одной переменной

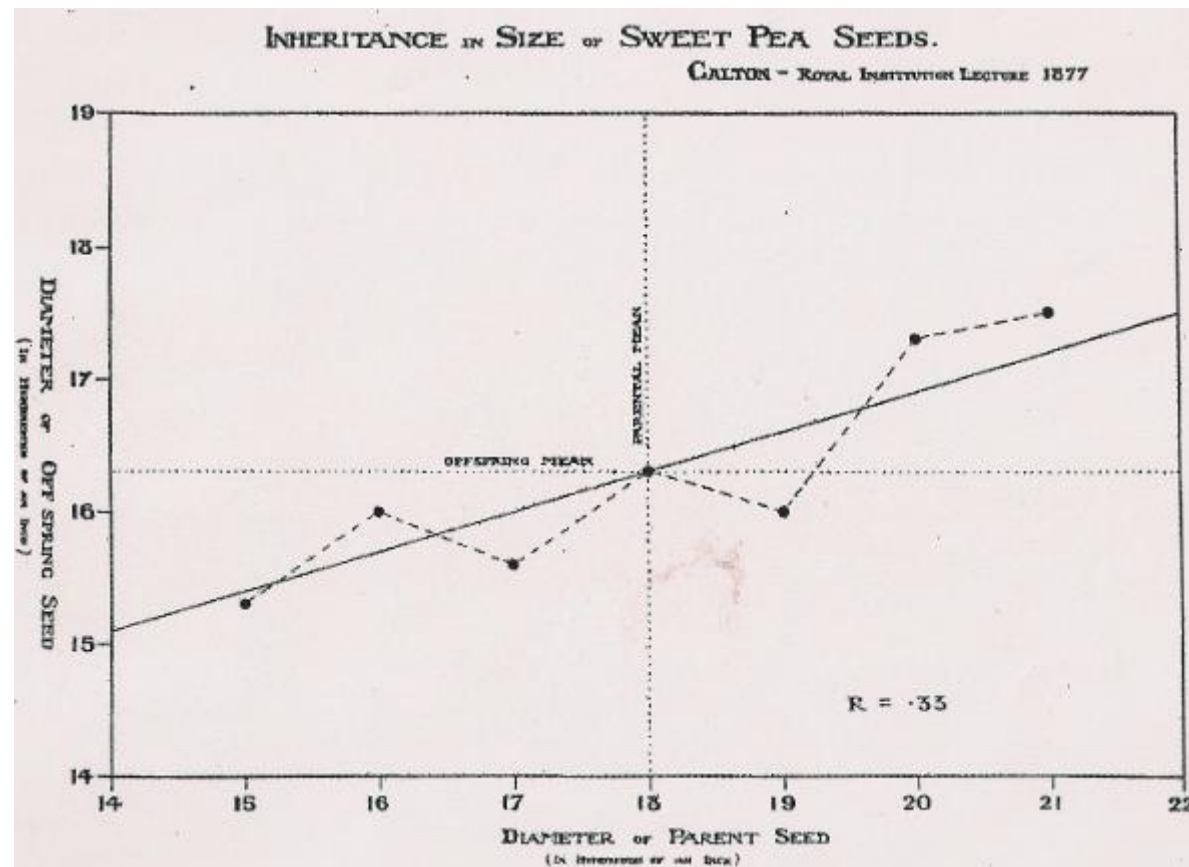
**Задача минимизации суммы квадратов отклонений
(residual sum of squares)**

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



**~ задача описания данных гиперплоскостью (но ф-л качества!)
Есть вероятностное обоснование, но пока... логично**

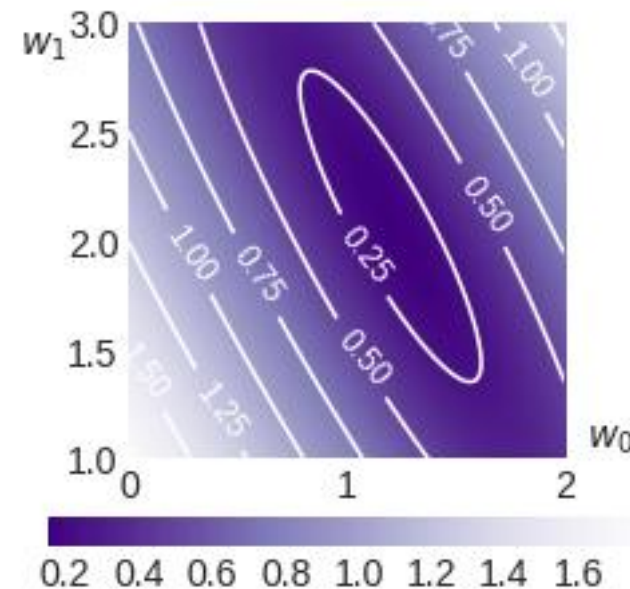
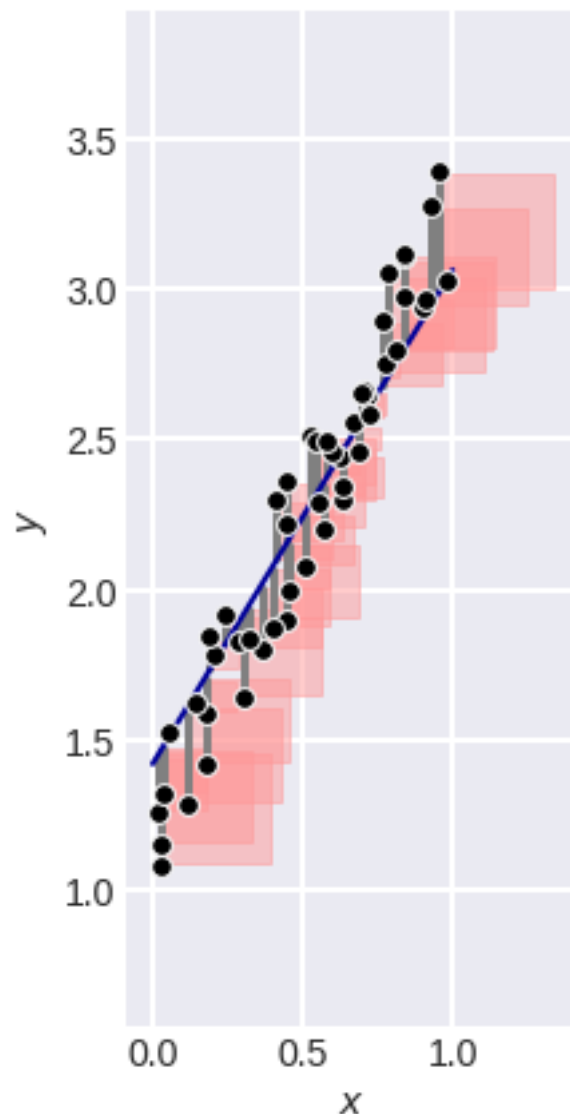
Линейная регрессия от одной переменной



Francis Galton, 1877

Линейная регрессия от одной переменной

Геометрический смысл ошибки



**Отличается от суммы расстояний
до поверхности!**

Линейная регрессия от одной переменной

Нетрудно показать (Д3):

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

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

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

Общий случай (многих переменных)

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

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

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

для удобства записи вводим фиктивный признак $x_0 \equiv 1$

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

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

$Xw = y$ – **как решать?**

Общий случай (многих переменных)

или 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)

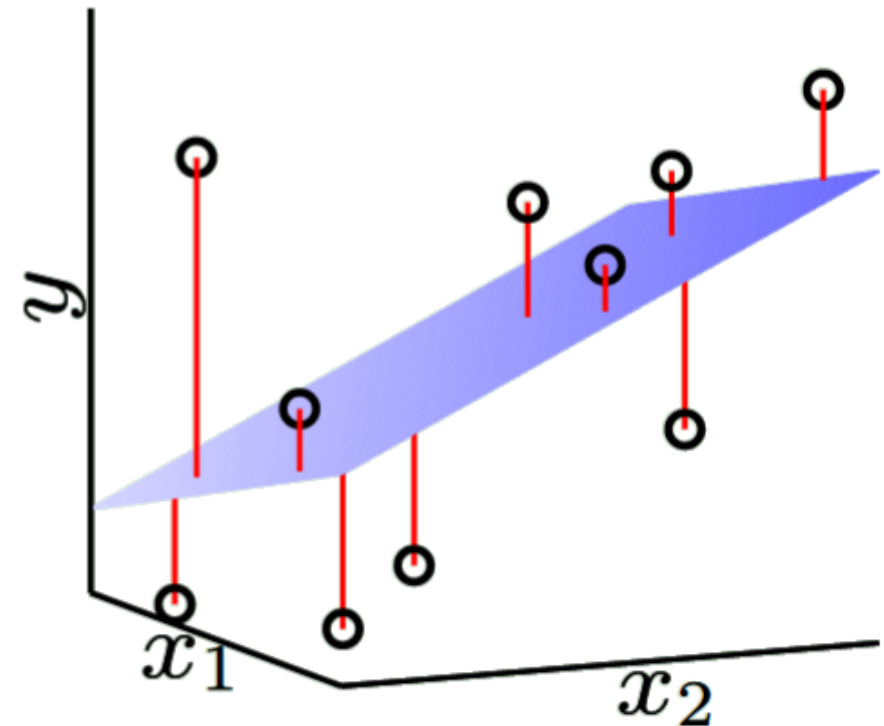
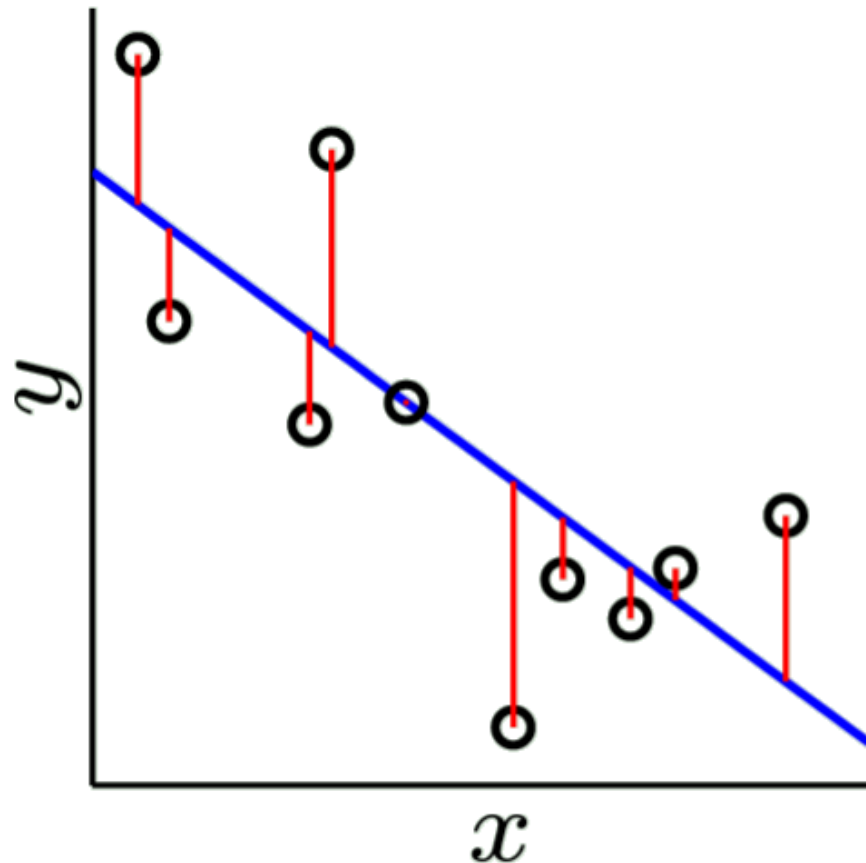
будем решать так:

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

почему?

Общий случай (многих переменных)

геометрический смысл



Решение задачи минимизации

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

$$\begin{aligned}\|Xw - y\|_2^2 &= (Xw - y)^T (Xw - y) = \\ &= w^T X^T Xw - w^T X^T y - y^T Xw + y^T y\end{aligned}$$

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

$$X^T Xw = X^T y$$

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

решение существует, если столбцы линейно независимые

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

Обобщённая линейная регрессия вместо X – что угодно

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

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

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

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

$$\equiv 1$$

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

базисные функции (basis functions)
они фиксированы

Подробности в нелинейных методах...

Проблема вырожденности матрицы

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

Решения:

1. Регуляризация – **здесь и в «сложности»**
2. Селекция (отбор) признаков – **«селекция» / «PZAD»**
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 регуляризации

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

**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 больших весов
(у признаков, по которым могут отличаться похожие объекты)!**

П.С. Плохо, когда модель заточена на один признак!

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

Хотим $\|w\|_2^2 \rightarrow \min$

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**

Всё это справедливо и для общих задач минимизации!

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

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

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

Есть ещё регуляризация Морозова...

Регуляризация

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

ДЗ Доказать!

– гребневая регрессия (Ridge Regression)

Другой смысл – боремся с вырожденностью матрицы!

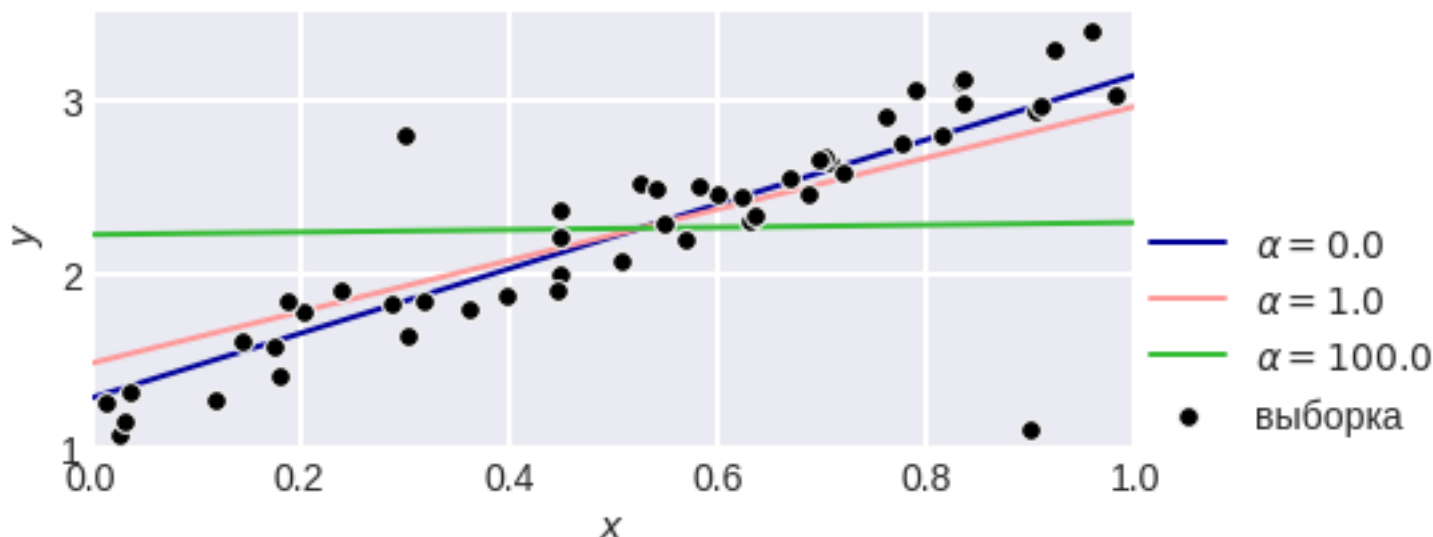
$$\lambda = 0$$
$$\lambda \rightarrow +\infty$$

– получаем классическое решение

**– меньше «затачиваемся на данные» и больше
регуляризуем**

Матрица очевидно становится обратимой!

Регуляризация – минутка кода



```
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])
```

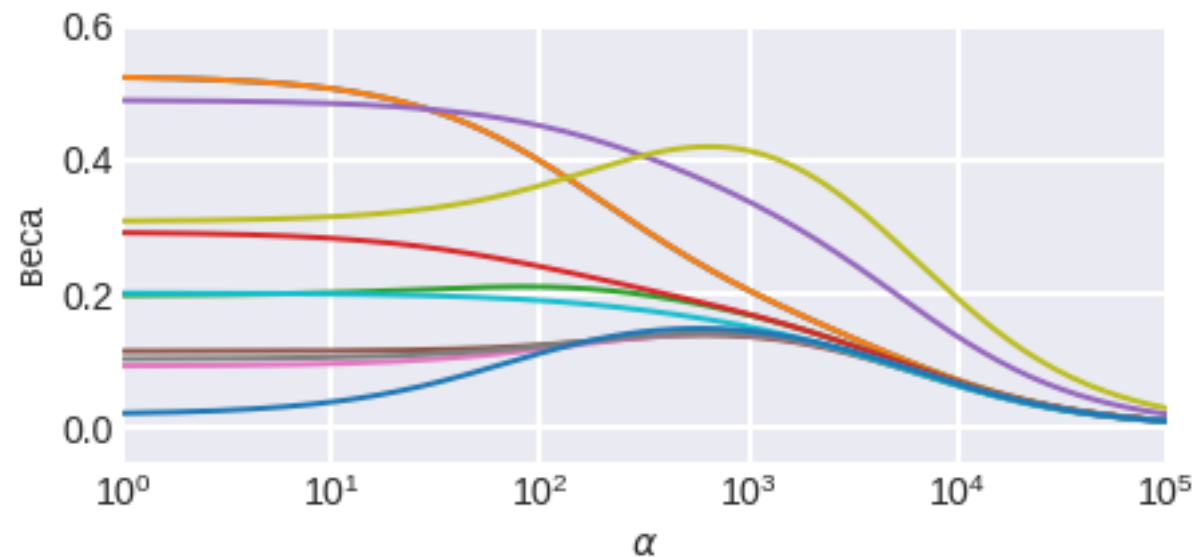
```
a_test = model.predict(x_test[:, np.newaxis])
```

**Интересно, что
рисунок неудачный –
получилась
антиреклама
регуляризации...
почему?**

Ridge-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**

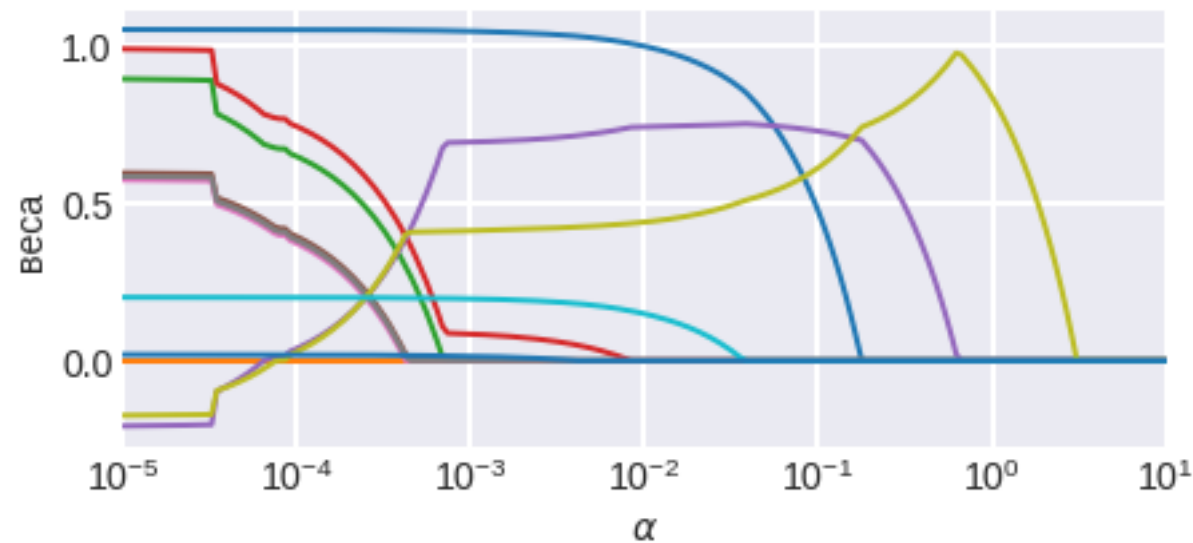
Ridge-регрессия

**Для ridge-регрессии нужна правильная нормировка признаков!
Нет инвариантности (в отличие от линейной) от умножения
признаков на скаляры**

Перед регуляризацией – стандартизация!!!

LASSO

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



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

Не на все коэффициенты нужна регуляризация!

Почему?

Масштаб очень важен! см. дальше

Семейство регуляризированных линейных методов

Ridge

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

LASSO (Least Absolute Selection and Shrinkage Operator)

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

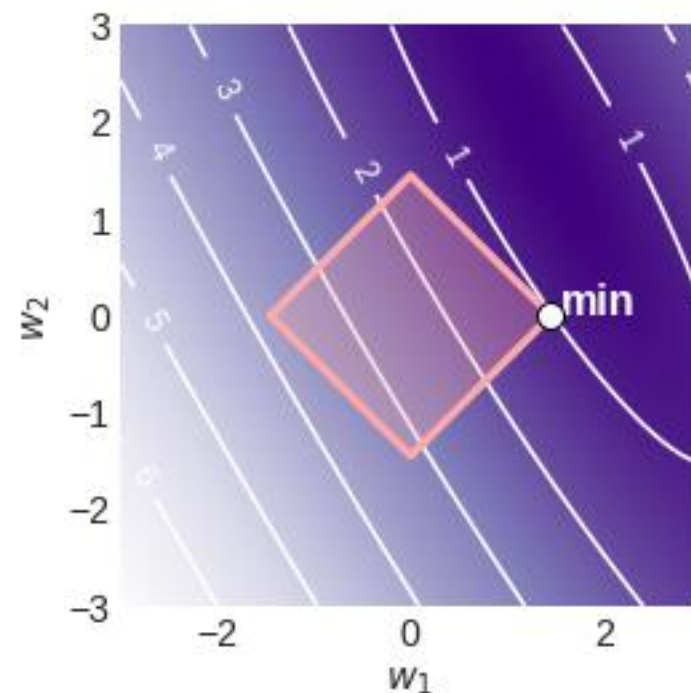
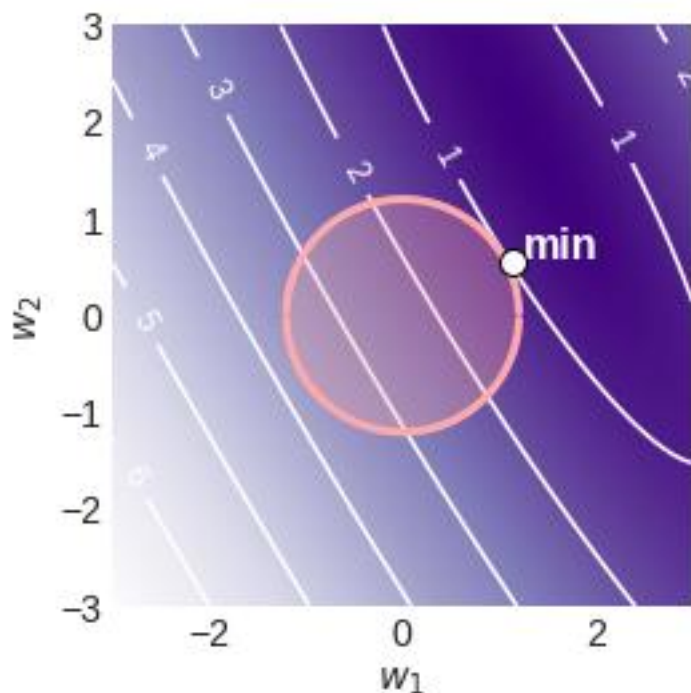
Elastic Net = LASSO + Ridge

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

Геометрический смысл Ridge, LASSO и Elastic Net

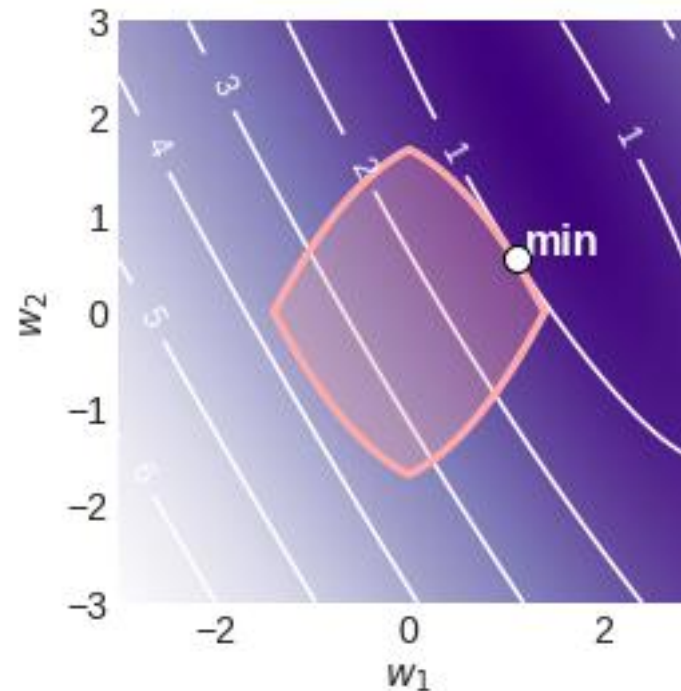
$$\sum_{i=1}^m \left(y_i - w_0 - \sum_{j=1}^n w_j x_{ij} \right)^2 \rightarrow \min_w, \quad \sum_{j=1}^n w_j^2 \leq s$$

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



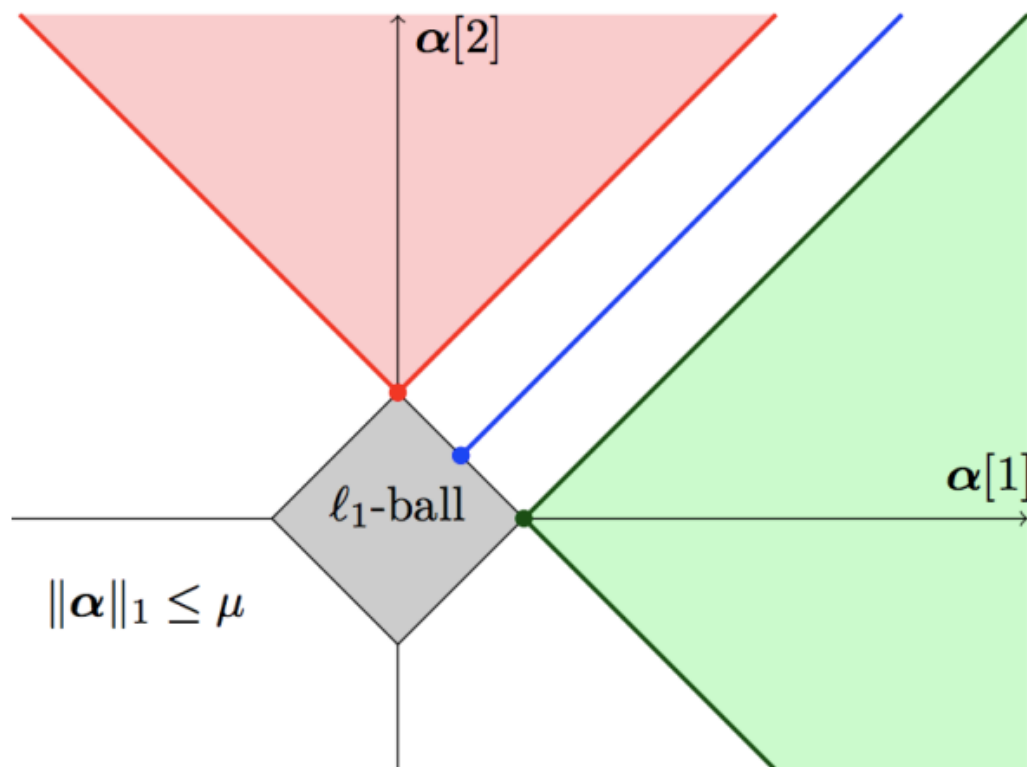
Геометрический смысл Ridge, LASSO и Elastic Net

Эти методы несравнимы... на практике часто модель и не может зависеть от небольшого числа переменных.

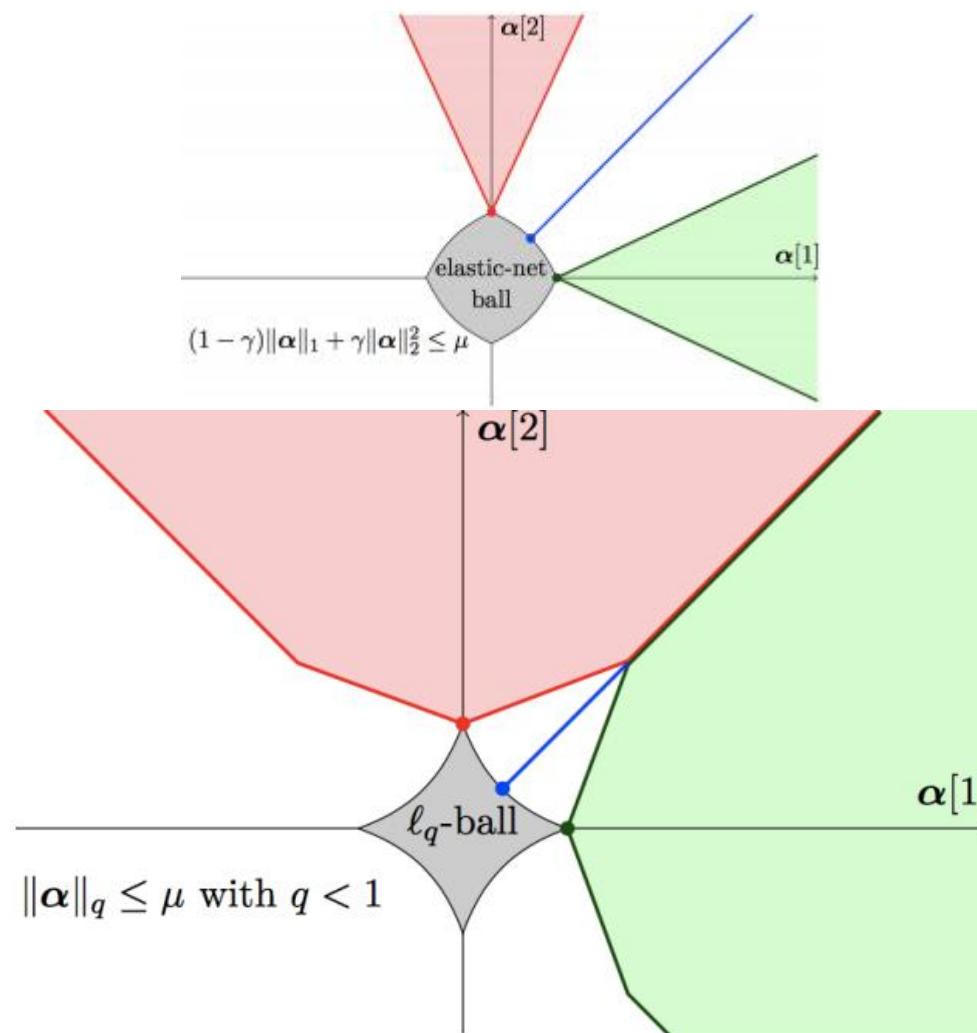


Эффект разреженности

если линии уровня оптимизируемой функции – концентрические окружности...



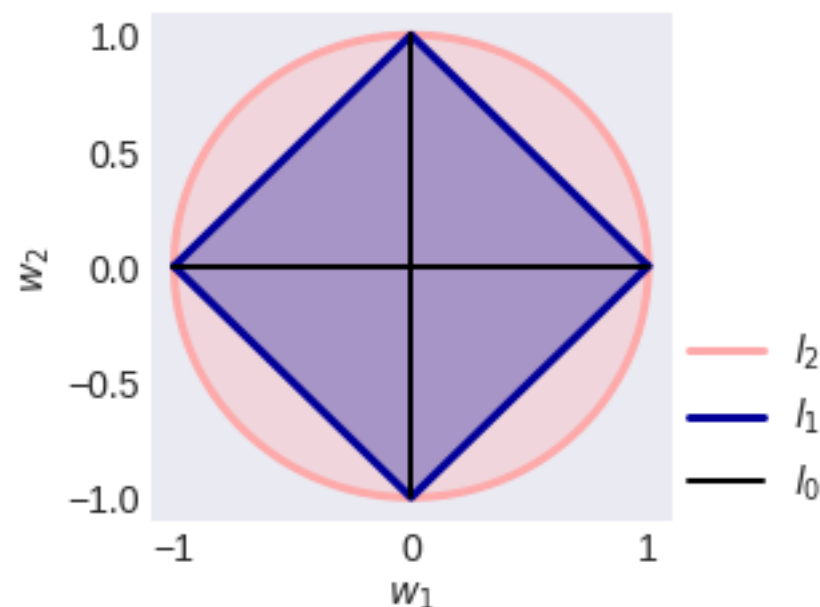
David S. Rosenberg «Foundations of Machine Learning» <https://bloomberg.github.io/foml/>



Почему L1-норма \Rightarrow разреженность

1. См. рис. больше вероятность, что линии уровней функции ошибки касаются области ограничений в точках с нулевыми координатами

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



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..

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

Селекция признаков в линейной регрессии

~ отдельная тема

Какие признаки включить в модель:

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

пока маленький обзор стратегий:

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

Обоснование необходимости селекции

- 1. Проблема вырожденности в линейной регрессии
- 2. Проблема «почти дубликатов»
- 3. Уменьшение модели и интерпретация
- 4. Уменьшение стоимости данных

Проблема вырожденности матрицы

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

Решения:

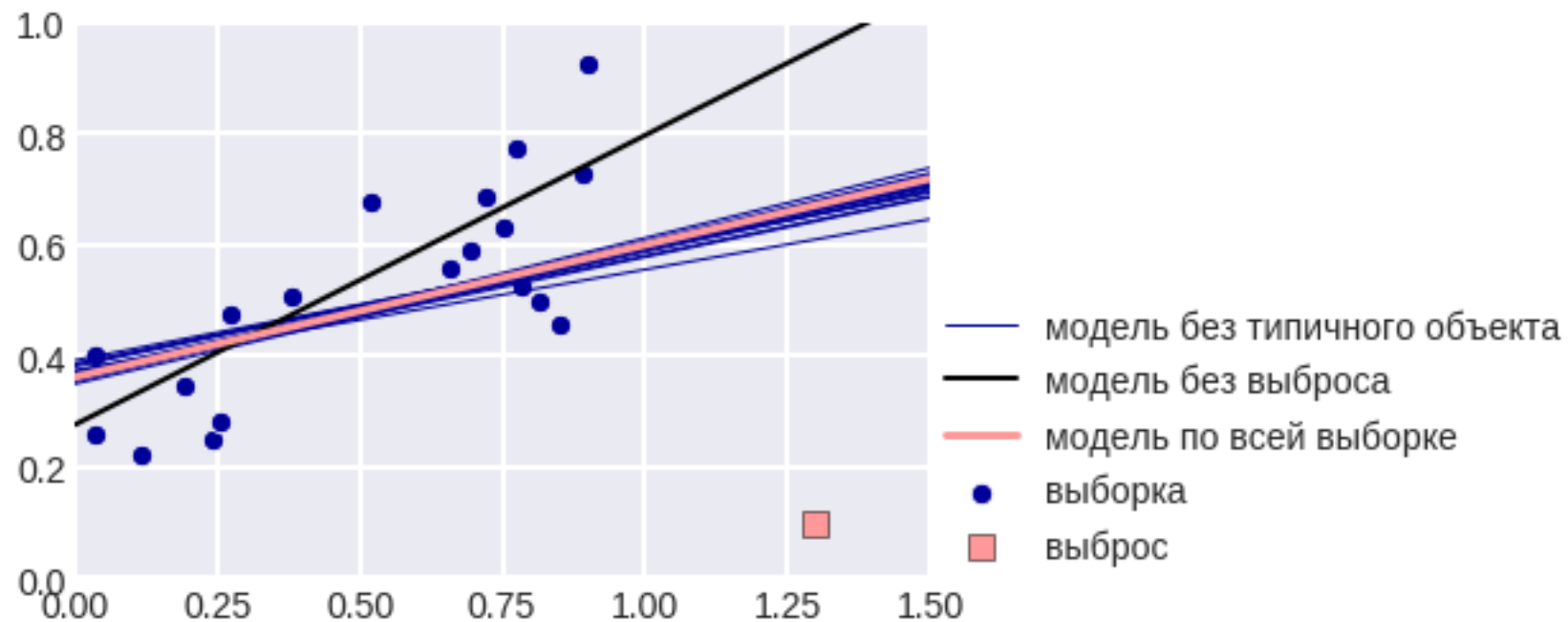
1. Регуляризация
2. Селекция (отбор) признаков
3. Уменьшение размерности (в том числе, PCA)
4. Увеличение выборки

	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 – неустойчивость к выбросам



Ошибка с весами

Если у каждого объекта есть цена ошибки...

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

небольшая переформулировка задачи:

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

если веса целые числа – можно продублировать объекты

**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

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

можно использовать любую регрессионную модель

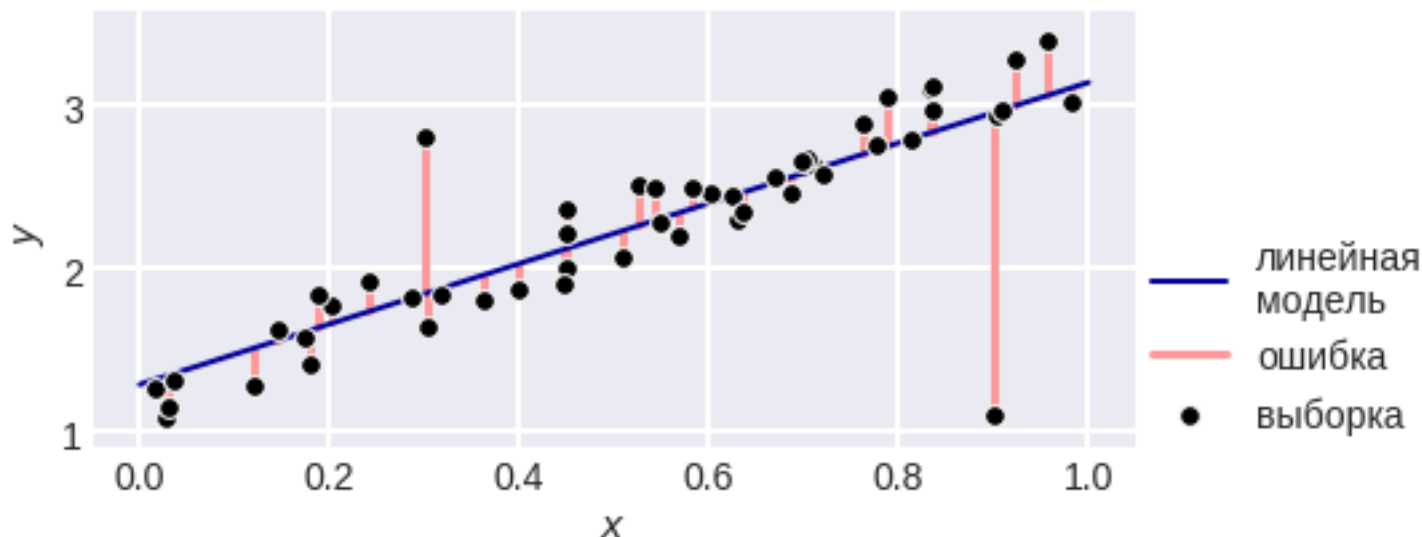
1.2. Вычислить ошибки на обучении

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

1.3. Пересчитать веса объектов

$$v_i = \exp(\varepsilon_i)$$

можно использовать другую невозрастающую функцию; можно (иногда нужно) нормировать



Линейные скоринговые модели в задаче бинарной классификации

Пусть $X = \mathbb{R}^n$, $Y = \{0, 1\}$

**Как решать задачи классификации с помощью линейной модели:
будем получать вероятность принадлежности к классу 1**

$$a(x) \in [0, 1]$$

**Любая линейная функция на \mathbb{R}^n будет получать значения в \mathbb{R} ,
поэтому нужна деформация (transfer function):**

$$\sigma: \mathbb{R} \rightarrow [0, 1]$$

В логистической регрессии

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

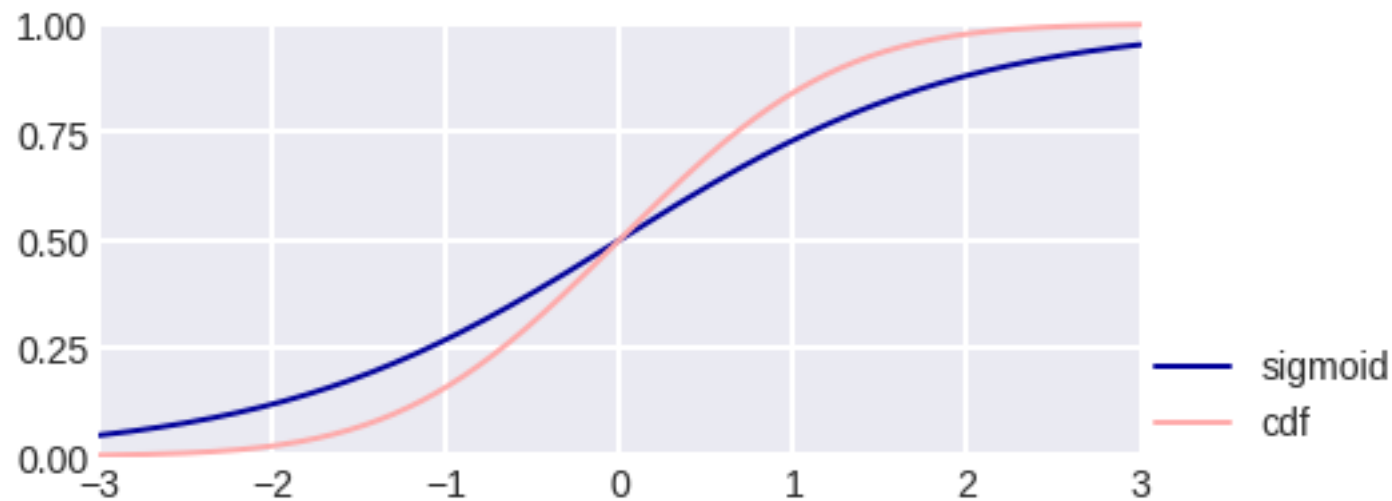
**Логистическая функция
(сигмоида)**

В Probit-регрессии

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp(-t^2 / 2) \partial t$$

**Normal Cumulative distribution
function**

Функции деформации



Логистическая регрессия

$$p(x) \equiv P(Y = 1 | x) = \sigma(z) = \frac{1}{1 + e^{-z}} \in (0, 1),$$

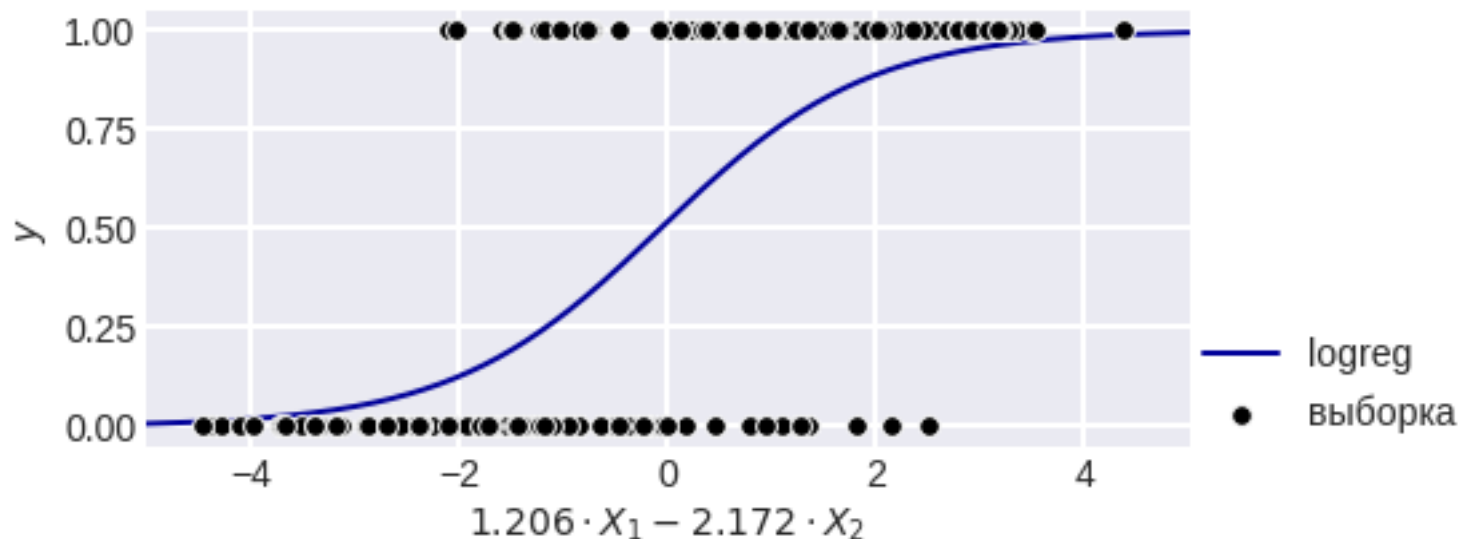
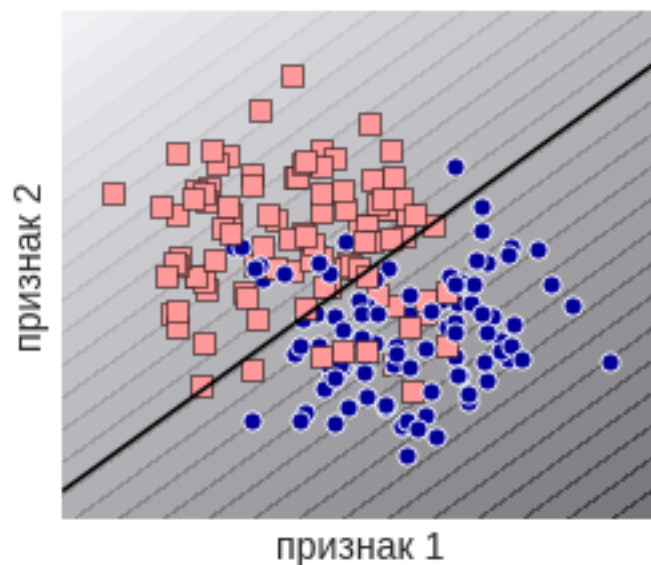
$$z = w_0 + w_1 X_1 + \dots + w_n X_n,$$

$$\log\left(\frac{p(x)}{1 - p(x)}\right) = z$$

– **монотонное преобразование, которое называют logit-transformation**

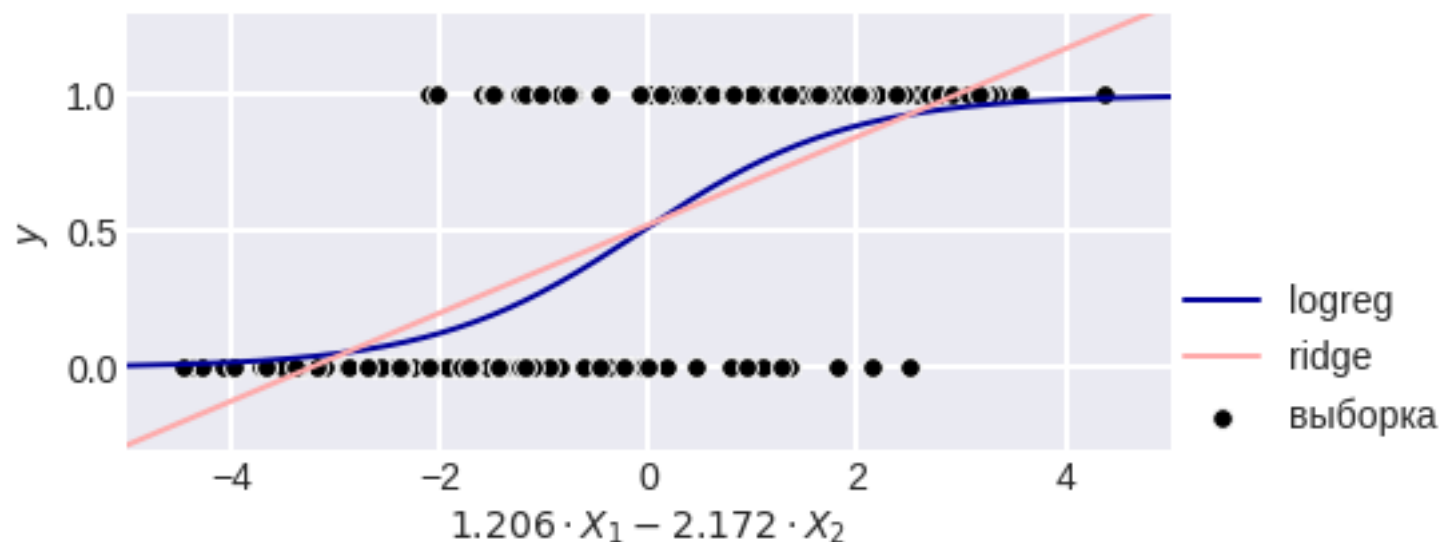
Геометрический смысл логистической регрессии

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
model.fit(X, y)
a = model.predict_proba(X_test)[:,1]
```



Можно и одномерную картинку

Чем логистическая регрессия лучше регрессии



Обучение логистической регрессии

Метод максимального правдоподобия

$$L(w_0, \dots, w_n) = \prod_{i: y_i=1} p(x_i) \prod_{i: y_i=0} (1 - p(x_i)) \rightarrow \max$$

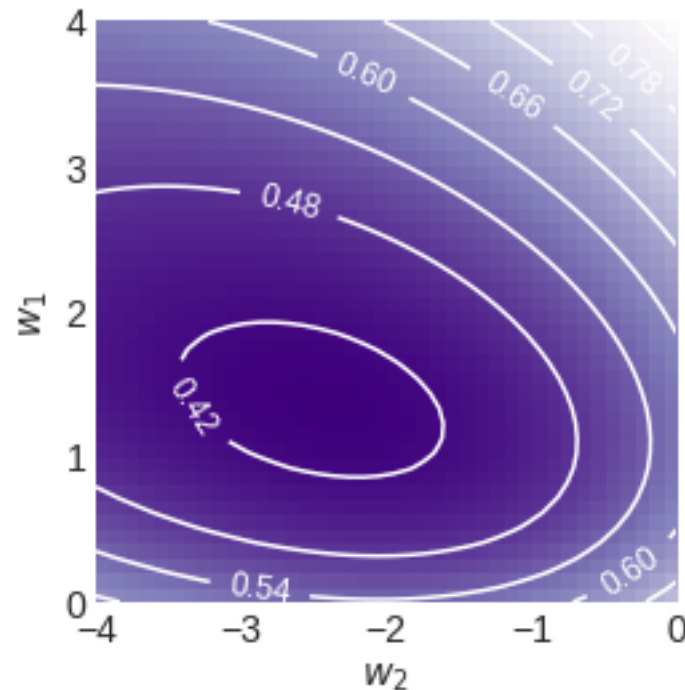
$$\log L = - \sum_{i: y_i=1} \log(1 + e^{-z_i}) - \sum_{i: y_i=0} \log(1 + e^{+z_i}) \equiv - \sum_i \log(1 + e^{-y'_i z_i})$$

$$\begin{aligned} \nabla_w \log L &= \sum_{i: y_i=1} \frac{1}{1 + e^{-w^T x_i}} e^{-w^T x_i} x_i - \sum_{i: y_i=0} \frac{1}{1 + e^{+w^T x_i}} e^{+w^T x_i} x_i = \\ &= \sum_i \frac{y'_i x_i}{1 + e^{+y'_i w^T x_i}} = \sum_i y'_i x_i \sigma(-y'_i w^T x_i) \end{aligned}$$

где (для удобства записи)

$$y'_i = 2y_i - 1$$

Ошибка логистической регрессии



метод SGD

$$w \leftarrow w + \eta \sigma(-y'_i w^T x_i) y'_i x_i$$

Запомним!

Многоклассовая логистическая регрессия Multiclass logistic regression (multinomial regression)

в glmnet такой «симметричный вариант»

$$P(Y = k \mid x) = \frac{e^{w_{0k} + w_{1k}X_1 + \dots + w_{nk}X_n}}{\sum_{j=1}^l e^{w_{0j} + w_{1j}X_1 + \dots + w_{nj}X_n}}$$

Если

$$\text{softmax}(a_1, \dots, a_l) = \frac{1}{Z} [e^{a_1}, \dots, e^{a_l}],$$

где $Z = e^{a_1} + \dots + e^{a_l}$

тогда

$$P(Y = k \mid x) = \text{softmax}(w(1)^T x, \dots, w(l)^T x)$$

Линейные решающие модели в задаче бинарной классификации

Пусть $X = \mathbb{R}^n$, $Y = \{\pm 1\}$

обучающая выборка: $\{(x_i, y'_i)\}_{i=1}^m$

хотим линейную модель:

$$a(x) = \text{sgn}(w^T x + b) = \begin{cases} +1, & w^T x + b > 0 \\ -1, & w^T x + b < 0 \end{cases}$$

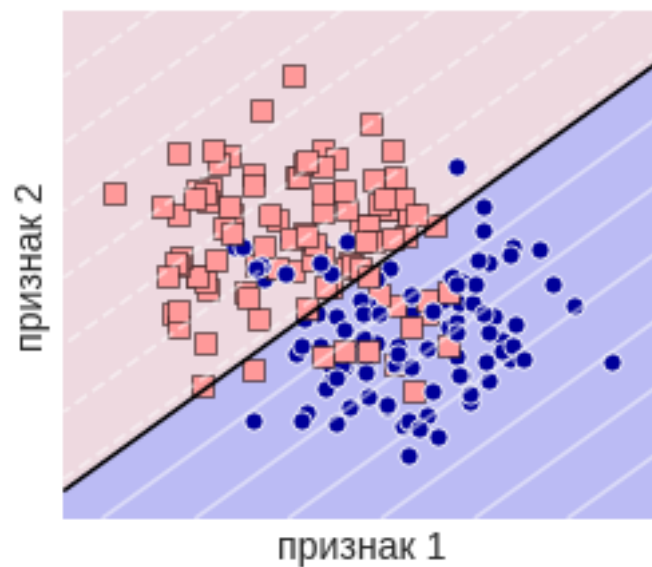
случай $w^T x + b = 0$ нам тут не особо важен

Пополняя признаковое пространство фиктивным признаком,

$$a(x) = \text{sgn}(w^T x)$$

– линейный классификатор

Геометрический смысл линейного классификатора



**Делим пространство
гиперплоскостью на две части**

Линейный классификатор

Общая идея:

$$L(y_t, a(x_t)) = \theta(-y'_t w^T x_t) = \begin{cases} 1, & \text{sgn } w^T x_t = y'_t \\ 0, & \text{sgn } w^T x_t \neq y'_t, \end{cases}$$

$$L(X_{\text{train}}, a) = \sum_{t=1}^m L(y_t, a(x_t)) \rightarrow \min$$

**естественно минимизировать число ошибок,
но**

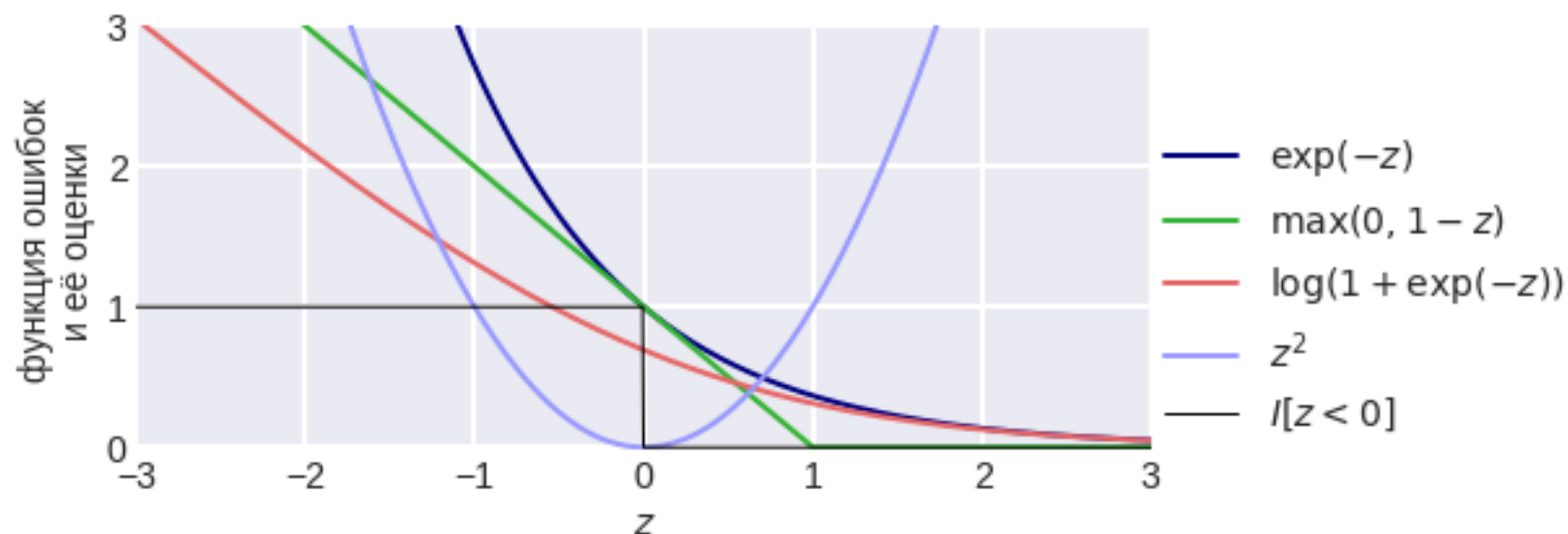
- **ф-ия не дифференцируема**
 - **выдаёт мало информации**
 - **оптимизация здесь – NP-полная задача**
- только число ошибок, а не их «фатальность»**

$y'_t w^T x_t \sim$ **чем меньше, тем хуже (зазор)**

Оценка функции ошибок через гладкую функцию

$$\sum_{t=1}^m L(y_t, a(x_t)) \leq \sum_{t=1}^m L'(y_t, a(x_t)) \rightarrow \min$$

$$\sum_{t=1}^m \theta(-\xi_t) \leq \sum_{t=1}^m f(-\xi_t) \rightarrow \min$$



Оценка функции ошибок через гладкую функцию

Примеры замен:

$$f(-\xi) = \exp(-\xi)$$

$$f(-\xi) = \max(0, 1 - \xi)$$

$$f(-\xi) = \log(1 + \exp(-\xi))$$

**Обучение – минимизация оценки на обучающей выборке
(+ регуляризация)**

Персептрон, SVM, логистическая регрессия минимизируют выпуклые аппроксимации 0-1-loss, сводя NP-трудную задачу к задаче выпуклой оптимизации

Пример персептронного алгоритма

$$f(-\xi) = \max(0, -\xi)$$

Персептрон:

$$\sum_{i=1}^m \max[0, -y'_i(w^T x_i)] \rightarrow \min$$

SGD

$$\frac{\partial \max[0, -y'_i(w^T x_i)]}{\partial w} = \begin{cases} 0, & -y'_i(w^T x_i) < 0, \\ -y'_i x_i & -y'_i(w^T x_i) > 0, \end{cases} = \begin{cases} 0, & \text{sgn}(w^T x_i) = y'_i, \\ -y'_i x_i & \text{sgn}(w^T x_i) \neq y'_i, \end{cases}$$

$$w \leftarrow w + \eta \begin{cases} 0, & \text{sgn}(w^T x_i) = y'_i, \\ +x_i & w^T x_i \leq 0, y'_i = +1, \\ -x_i, & w^T x_i \geq 0, y'_i = -1, \end{cases}$$

Пример персептронного алгоритма

Есть теорема Новикова

Если две выборки линейно разделимы, то разделяющая поверхность находится персептронным алгоритмом за конечное число шагов

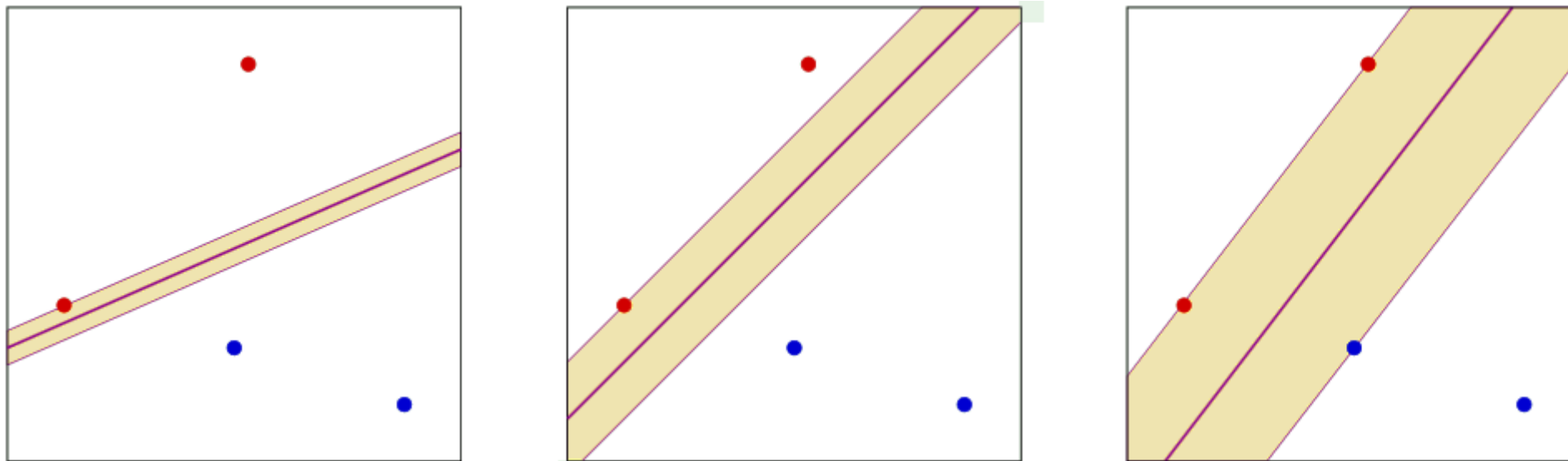
Пример персептронного алгоритма

$$\begin{cases} 2w_1 + w_2 > 0 \\ -w_1 > 0 \\ w_1 - w_2 < 0 \\ -2w_1 - 2w_2 < 0 \end{cases}$$

$$\begin{bmatrix} 2 & 1 \\ -1 & 0 \\ -1 & 1 \\ 2 & 2 \end{bmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} > 0$$

```
X = list(np.array([[2,1], [-1, 0],
                  [-1, 1], [2, 2]]))
w = np.array([0, 0])
print ('w=', w)
change = True
while (change):
    change = False
    for x in X:
        if np.dot(x, w) <= 0:
            print ('w=', w, 'x=', x, '-')
            w += x
            change = True
print ('w=', w)
```

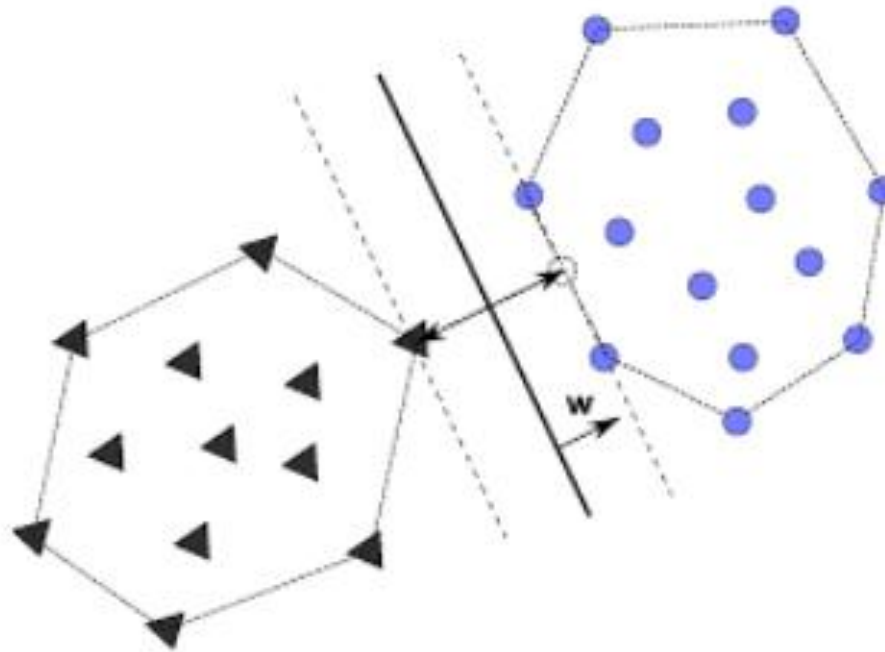
```
w=[ 0 0]
w=[ 0 0] x= [ 2 1] -
w=[ 2 1] x= [-1 0] -
w=[ 1 1] x= [-1 1] -
w=[ 0 2] x= [-1 0] -
w=[-1 2] x= [ 2 1] -
w=[ 1 3] x= [-1 0] -
w=[ 0 3] x= [-1 0] -
w=[-1 3]
```

SVM: идея максимального зазора

**до сих пор хотели разделить точки гиперплоскостью...
а как лучше?**

SVM: идея максимального зазора

Построение SVM эквивалентно нахождению кратчайшего отрезка, соединяющего выпуклые оболочки двух классов



SVM: постановка задачи

Хотим разделить точки двух разных классов гиперплоскостью

$$a(x) = \text{sgn}(w^T x + b)$$

Обучающая выборка:

$$\{(x_i, y_i)\}_{i=1}^m$$

должно быть (здесь пускай нет штрихов)

$$w^T x_i + b \geq 1 \text{ если } y_i = +1$$

$$w^T x_i + b \leq -1 \text{ если } y_i = -1$$

Другая форма записи:

$$y_i(w^T x_i + b) \geq 1$$

можно считать (из-за нормировки), что

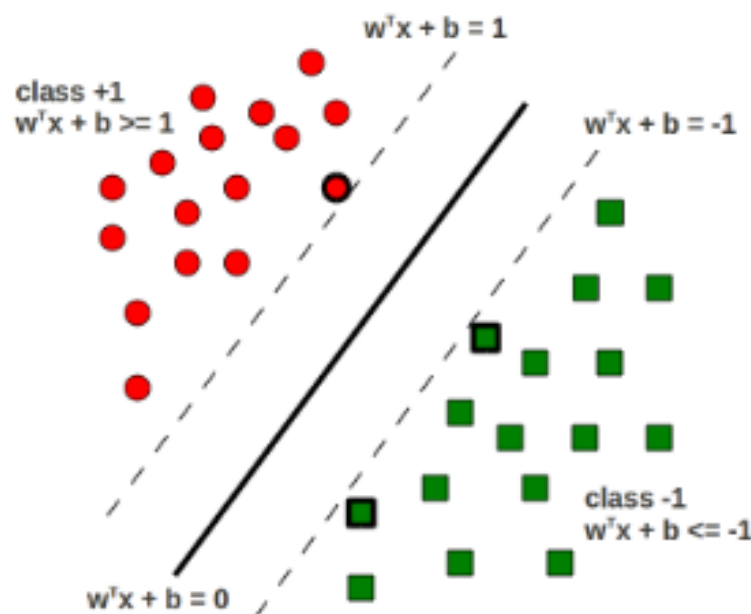
$$\min_i |w^T x_i + b| = 1$$

SVM: постановка задачи**Расстояние от точки до гиперплоскости:**

$$\rho(x_i, w^T x + b) = \frac{|w^T x_i + b|}{\|w\|}$$

хотим, чтобы минимум из этих расстояний был максимален

$$\min_i \frac{|w^T x_i + b|}{\|w\|} = \frac{1}{\|w\|} \rightarrow \max$$

– зазор (margin)

Зазор (margin)

В общем случае, когда

Пусть $X = \mathbb{R}^n$, $Y = \{\pm 1\}$

обучающая выборка: $\{(x_i, y_i)\}_{i=1}^m$

Алгоритм со скоринговой функцией (score function)

$$a(x) = \text{sgn}(b(x)), b(x) \in \mathbb{R}$$

Зазор – $y_i b(x_i)$

~ уверенность в ответе

SVM: постановка задачи

$$\frac{\|w\|^2}{2} \rightarrow \min$$
$$y_i(w^T x_i + b) \geq 1, i \in \{1, 2, \dots, m\}$$

**– задача квадратичного программирования
(QP = Quadratic Program)
с m ограничениями (constraints)**

Заметим, что здесь тоже, как и в регуляризации линейной регрессии, хотим квадрат нормы весов сделать меньше

«квадрат» – для удобства оптимизации

SVM: приближённое решение «в лоб»

$$\frac{\|w\|^2}{2} \rightarrow \min$$
$$1 - y_i(w^T x_i + b) \leq 0, i \in \{1, 2, \dots, m\}$$

не решаем задачу **точно, но стремимся выполнить условия:**

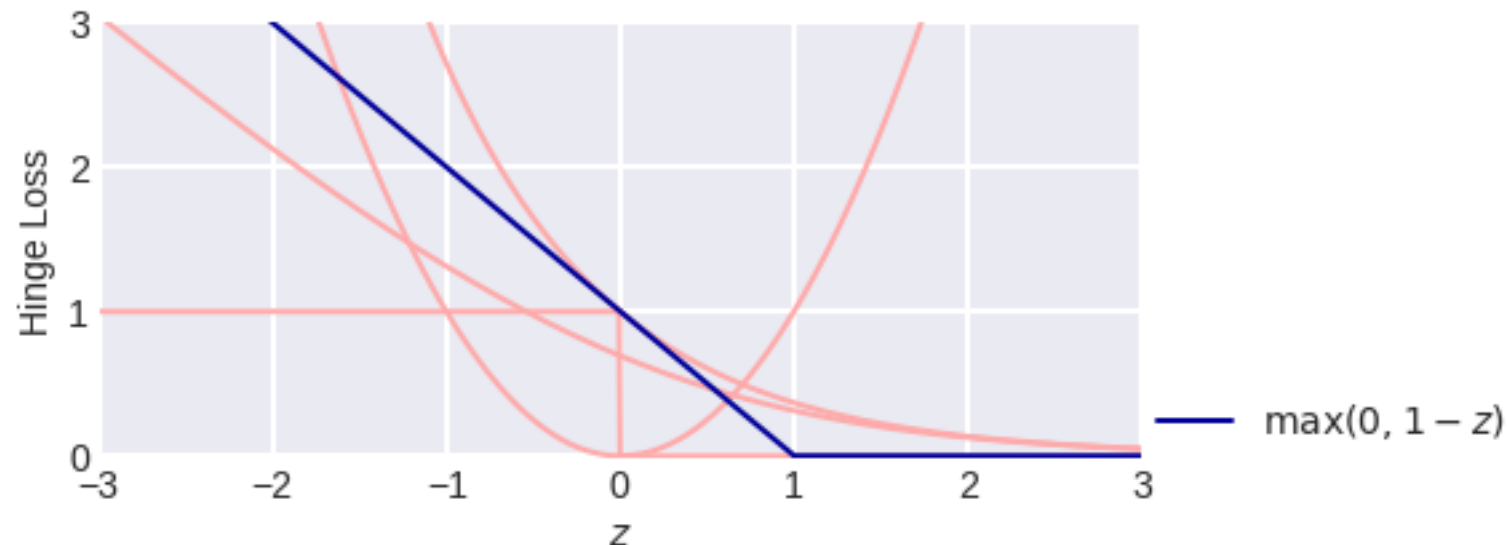
$$\frac{1}{m} \sum_{i=1}^m \max[0, 1 - y_i(w^T x_i + b)] + \lambda \|w\|^2 \rightarrow \min$$

Удивительно:

ошибка $L(y, a) = \max[0, 1 - ya]$ **+ регуляризатор**

но тут нет дифференцируемости из-за max

Hinge loss



**А в логистической регрессии:
логистическая функция ошибки + регуляризатор**

SVM: решение строгое

$$\frac{\|w\|^2}{2} \rightarrow \min$$
$$1 - y_i(w^T x_i + b) \leq 0, i \in \{1, 2, \dots, m\}$$

Вспоминаем оптимизацию с ограничениями:

$$\min_{w, b} \max_{\alpha \geq 0} L(w, b, \alpha) = \frac{w^T w}{2} + \sum_{i=1}^m \alpha_i (1 - y_i(w^T x_i + b))$$

тут будет дифференцируемость, но ограничения

SVM: решение строгое

$$\min_{w,b} \max_{\alpha \geq 0} L(w,b,\alpha) = \frac{w^T w}{2} + \sum_{i=1}^m \alpha_i (1 - y_i (w^T x_i + b))$$

возьмём производные, приравняем к нулю

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

Таким образом, оптимальный вектор весов – взвешенная сумма признаков описаний объектов из обучения
если подумать – аналогично происходит и при настройке персептрона и логистической регрессии...

Переход к двойственной задаче

Задача квадратичного программирования

$$\max_{\alpha \geq 0} \min_{w, b} L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{при условиях } \sum_{i=1}^m \alpha_i y_i = 0$$

Информацию об описаниях объектов мы используем лишь в виде их попарных скалярных произведений!

когда решим задачу...

$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$b = -\frac{1}{2} \left(\min_{i: y_i=+1} w^T x_i + \max_{i: y_i=-1} w^T x_i \right)$$

большинство α_i обратятся в ноль (из-за условий Кунна-Таккера)

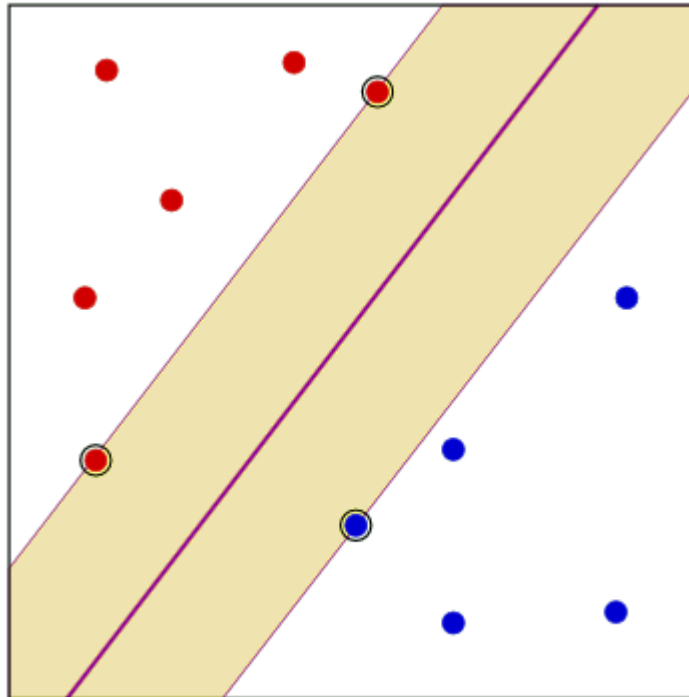
Условия Кунна-Таккера

для решения:

$$\alpha_i(1 - y_i(w^T x_i + b)) = 0$$

если $\alpha_i > 0$, то x_i – **опорный вектор (support vector)**

– лежит на границе $y_i(w^T x_i + b) = 1$



Зачем переходить к двойственной задаче

- размерности

м.б. удобно решать

выгодно, если признаковое пространство большое

- известная задача

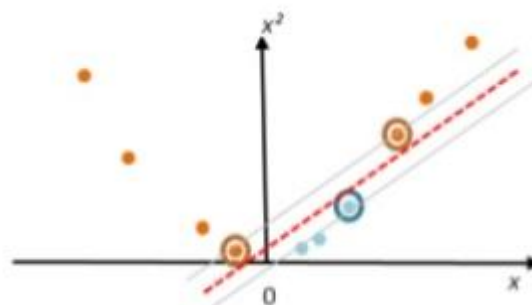
можно использовать солверы из готовых библиотек

- возникли попарные произведения

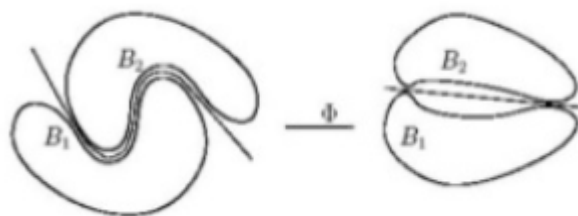
потом используем для **kernel tricks**

Если нет линейной разделимости

Два подхода (часто используются вместе)

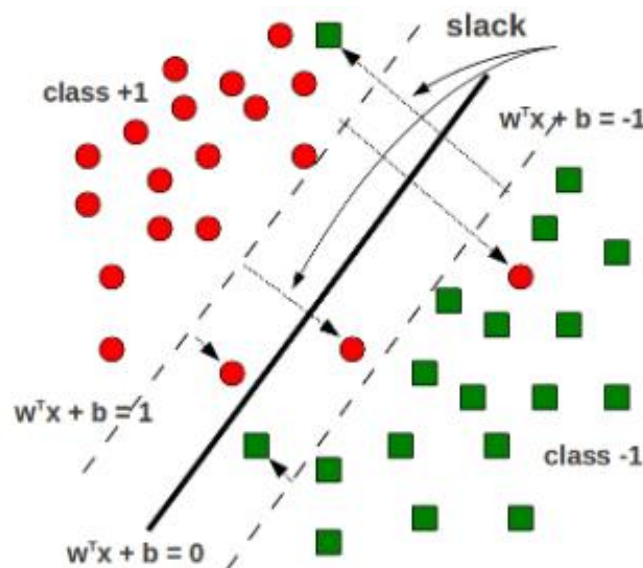


1) разделять так, чтобы ошибок было мало



2) использование нелинейных разделяющих поверхностей
переход в другое признаковое пространство
потом подробно разберём!

Soft-Margin SVM: разделение допуская ошибки



позволить объектам «залезать» на половину другого класса

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$

но не хотим, чтобы было много больших залезаний

Soft-Margin SVM: разделение допуская ошибки**Прямая задача:**

$$\frac{\|w\|^2}{2} + C \sum_{i=1}^m \xi_i \rightarrow \min$$

$$y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i \in \{1, 2, \dots, m\}$$

**тоже задача квадратичного программирования,
но в два раза больше ограничений**

C – баланс между оптимизацией зазора и ошибки на обучении

Soft-Margin SVM решается аналогично... ДЗ вывести!

Двойственная задача:

$$\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \rightarrow \max_{0 \leq \alpha \leq C}$$

появляется лишь ограничение $\alpha \leq C$

одно нетривиальное ограничение $\sum_{i=1}^m \alpha_i y_i = 0$

SVM Regression

хотим решить с ε -точностью

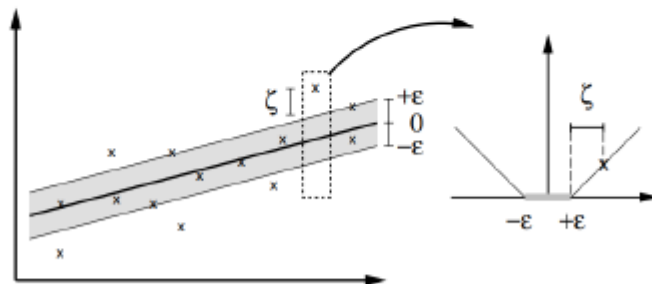
$$\frac{\|w\|^2}{2} \rightarrow \min$$

$$|w^T x_i + b - y_i| \leq \varepsilon, i \in \{1, 2, \dots, m\}$$

Equivalent unconstrained formulation:

$$\frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \mathcal{L}(\langle w, x_n \rangle + w_0 - y_n) \rightarrow \min_w$$

with ε insensitive loss $\mathcal{L}(u) = \begin{cases} 0, & \text{if } |u| \leq \varepsilon \\ |u| - \varepsilon & \text{otherwise} \end{cases}$



Solution will depend only on objects with $|\text{error}| \geq \varepsilon$, called *support vectors*.

SVM

- **должно быть хорошее пространство**
(однородные признаки в одной шкале)
- **тогда работают линейные SVM**
(нелинейные – с ядрами – успешно заменяются другими алгоритмами)
- **не подходят для больших данных**
(особенно нелинейные)
- **требуют хранения опорных векторов**

Проблемы с линейными алгоритмами

- + 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)**
- линейная гипотеза вряд ли верна**
- в теоретическом обосновании ещё предполагается нормальность ошибок**
 - «страдает» из-за выбросов**
 - признаки в одной шкале и однородные**
- статистический вывод регрессии – много предположений**
 - проблема коррелированных признаков**
- \Rightarrow необходимость регуляризации, селекции, PCA, data[↑]**

Проблемы мультиколлинеарности

- **большие коэффициенты**
- **большие изменения коэффициентов при добавлении/удалении признаков**
- **нелогичности**
(чем больше доход, меньше вероятность дать кредит)
- **большое число статистически незначимых оценок коэффициентов**

зависимость от масштабирования**простая модель****нет****с регуляризацией****есть****пайплайн:
нормировка +
линейная****нет**

Зачем нужен дискриминантный анализ?

Когда классы хорошо разделимы оцениваемые параметры, скажем, для логистической регрессии могут быть нестабильны.

Линейный дискриминантный анализ меньше подвержен этой проблеме.

Если размерность малая, распределения нормальные – линейная дискриминантная модель опять лучше.

Также LDA можно приспособить для представления данных в маломерных пространствах.

LDA – для малых размерностей и нормально распределенных данных

Наивный Байес – для больших размерностей

Линейный дискриминант Фишера

рассмотрим случай двух классов:

$$X = \mathbb{R}^n, Y = \{\pm 1\}$$

$$X_\alpha = \{x_i \mid y_i = \alpha\}$$

$$m_\alpha = |X_\alpha|$$

$$m_{+1} + m_{-1} = m$$

$$\mu_\alpha = \frac{1}{m_\alpha} \sum_{x_i \in X_\alpha} x_i$$

$$\sigma_\alpha^2 = \sum_{x_i \in X_\alpha} (x_i - m_\alpha)^2$$

хотим, чтобы проекции на некоторую прямую

$$\frac{(\mu_{+1} - \mu_{-1})|_w|^2}{\sigma_{+1}^2|_w + \sigma_{-1}^2|_w} \rightarrow \max$$

Линейный дискриминант Фишера

$$\begin{aligned}(\mu_{+1} - \mu_{-1})|_w^2 &= (w^T \mu_{+1} - w^T \mu_{-1})^2 = (w^T (\mu_{+1} - \mu_{-1}))^2 = \\ &= w^T (\mu_{+1} - \mu_{-1}) (\mu_{+1} - \mu_{-1})^T w = w^T S w\end{aligned}$$

$$\begin{aligned}\sigma_\alpha^2|_w &= \sum_{x_i \in X_\alpha} (w^T x_i - w^T \mu_\alpha)^2 = w^T \sum_{x_i \in X_\alpha} (x_i - \mu_\alpha)(x_i - \mu_\alpha)^T w = \\ &= w^T S_\alpha w \\ \frac{w^T S w}{w^T (S_{+1} + S_{-1}) w} &\rightarrow \max\end{aligned}$$

междуклассовый разброс
внутриклассовый разброс

Линейный дискриминант Фишера

решение $w \propto (S_{+1} + S_{-1})^{-1}(\mu_2 - \mu_1)$

обосновать!

интересный факт:

можно получить из МНК, выбрав целевые значения

$$\alpha \rightarrow \frac{m_{+1} + m_{-1}}{m_{\alpha}}$$

обосновать!