

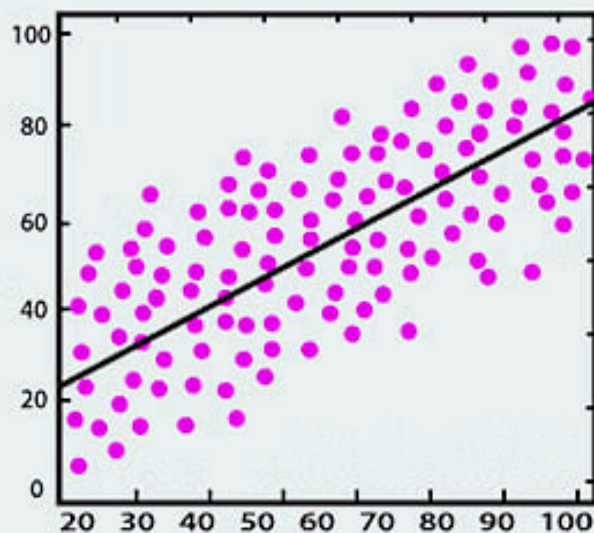
Машинное обучение • Основные алгоритмы

- Основные задачи в машинном обучении
- разберем основные алгоритмы
- ответим на вопрос "какой признак самый важный?"
- визуализируем построенные алгоритмы

Классическое Обучение

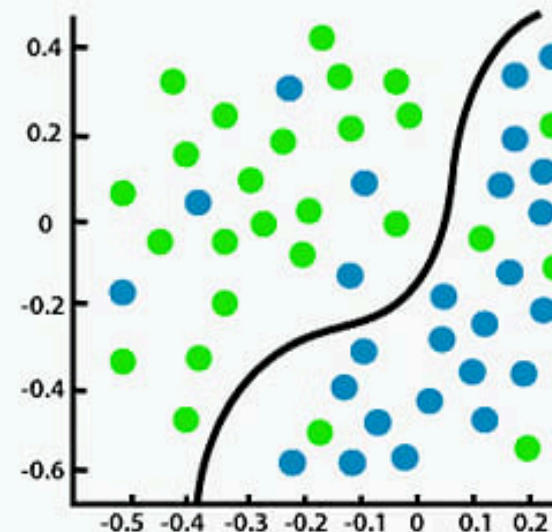


Обучение с Учителем. Регрессия и классификация



Regression

versus



Classification

- Примеры задач **регрессии**:

1. Прогнозирование продаж на следующий месяц.
2. Оценка стоимости недвижимости.
3. Прогнозирование доходов компании.

- Примеры задач **классификации**:

1. Определение, будет ли клиент покупать продукт.
2. Классификация писем как спам или не спам.
3. Определение кредитоспособности клиента.

Задача: оценить стоимость квартиры по одному признаку (например, площади).

Дана таблица с обучающей выборкой:

Площадь, м2	Цена, млн
50	12
33	8
...	...
76	120

Решаем уравнение:

$$y = x_1 w_1 + w_0$$

где x_1 - значение площади, w - "вес" признака, w_0 - свободный параметр.

Если помимо площади есть еще признаки, то уравнение просто увеличивается:

$$y = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + w_0$$

или короче:

$$y = \sum_{i=1}^n x_i w_i + w_0$$

Линейная регрессия: решение

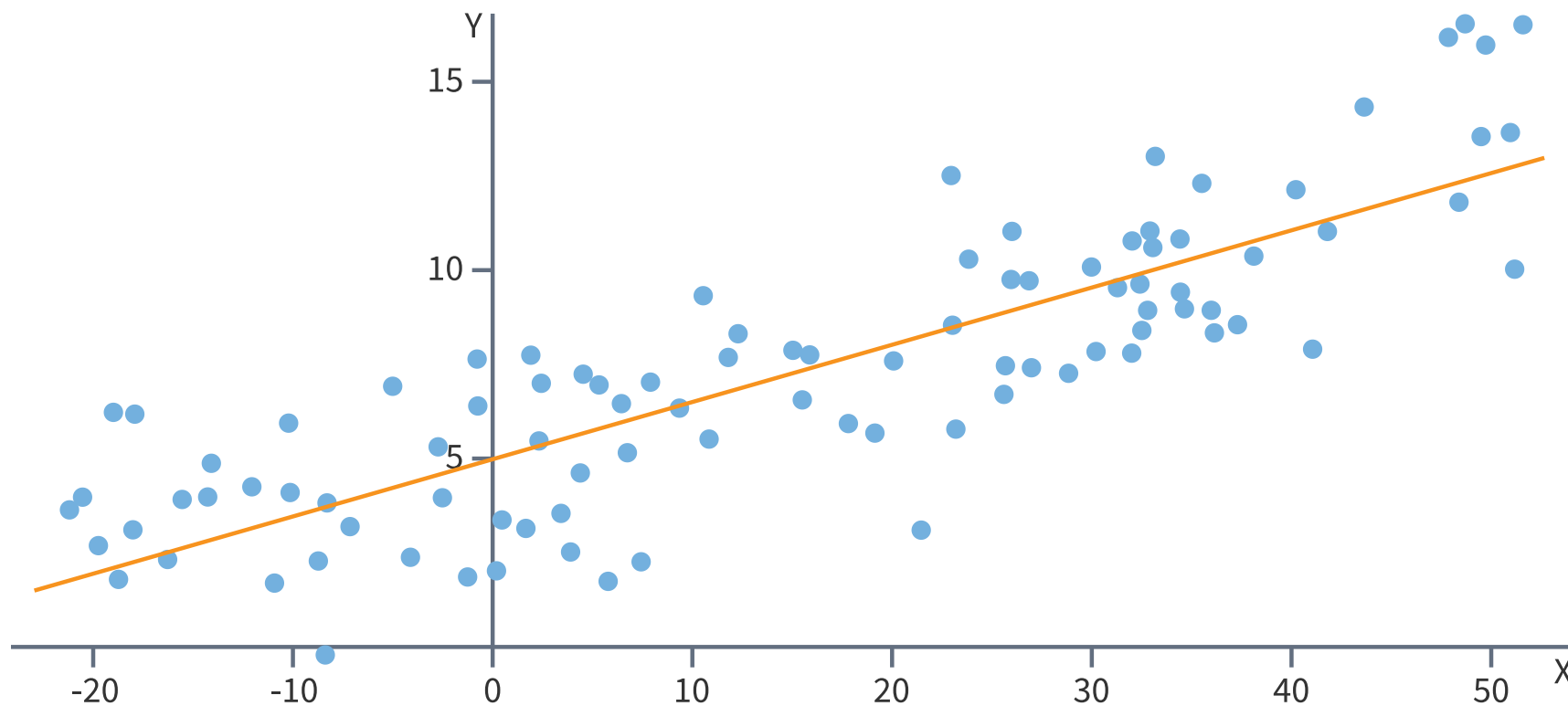
- методом градиентного спуска (так обычно и происходит)
 - будем минимизировать *среднеквадратическую ошибку*:

$$L = \frac{1}{K} \sum_{i=1}^K (y_i - \hat{y}_i)^2 \rightarrow \min$$

$$L = \frac{1}{K} \sum_{i=1}^K (y_i - (x_1 \cdot w_1 + x_2 \cdot w_2 + \dots + x_n \cdot w_n))^2 \rightarrow \min$$

y - настоящее значение, \hat{y} - предсказанное моделью значение, K - число объектов в обучающей выборке, w_1, w_2, \dots, w_n - веса признаков, это и есть наши дифференцируемые параметры, частные производные по ним будут составлять наш градиент.

Линейная регрессия: Визуализация



Логистическая регрессия

- Решает задачу классификации, не смотря на название!



Модель остается линейной:

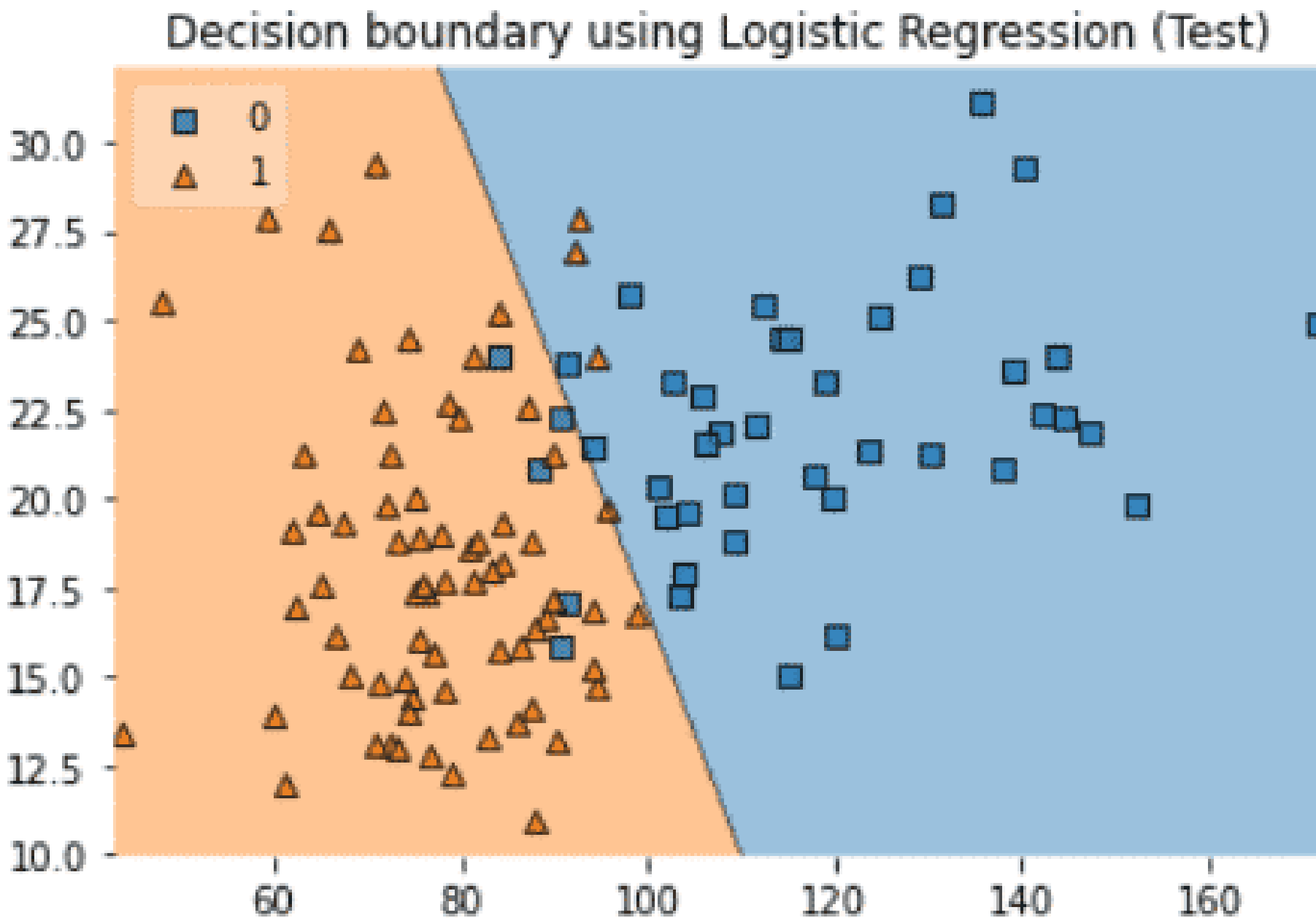
$$z = \sum_{i=1}^n x_i w_i + w_0$$

но полученный z подставляем в логистическую функцию:

$$y = \frac{1}{1 + e^{-z}}$$

выход y будет лежать в диапазоне от 0 до 1: $y \in [0, 1]$

это значение мы будем интерпретировать как *вероятность* того, что объект относится к классу **1**



- Часто модель может переобучиться на какие-либо выбросы в данных.
- Чтобы этого избежать, можно использовать регуляризацию:
регуляризация помогает уменьшить сложность модели и предотвратить переобучение, добавляя штраф за величину коэффициентов в функцию ошибки.
- Три самых распространенных метода регуляризации моделей:
 - i. $L1$ -регуляризация / LASSO
 - ii. $L2$ -регуляризация / Ridge
 - iii. ElasticNET

α - коэффициент регуляризации

Регуляризация линейных моделей

1. $L1$ -регуляризация / LASSO: добавляем в функцию ошибки сумму модулей весов:

$$L(y, \hat{y}) = \frac{1}{K} \sum_{i=1}^K (y - \hat{y})^2 + \alpha \sum_{j=1}^n |w_j|$$

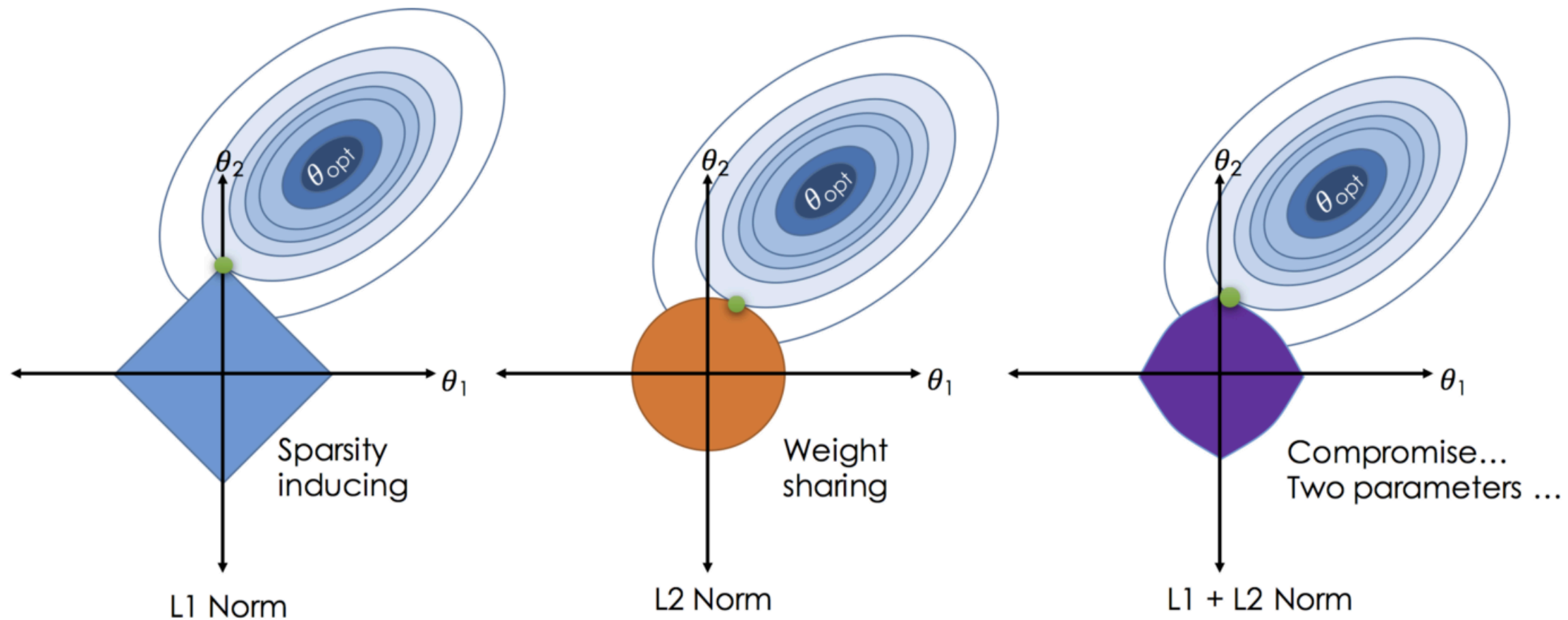
2. $L2$ -регуляризация / Ridge: добавляем в функцию ошибки сумму квадратов весов:

$$L(y, \hat{y}) = \frac{1}{K} \sum_{i=1}^K (y - \hat{y})^2 + \alpha \sum_{j=1}^n w_j^2$$

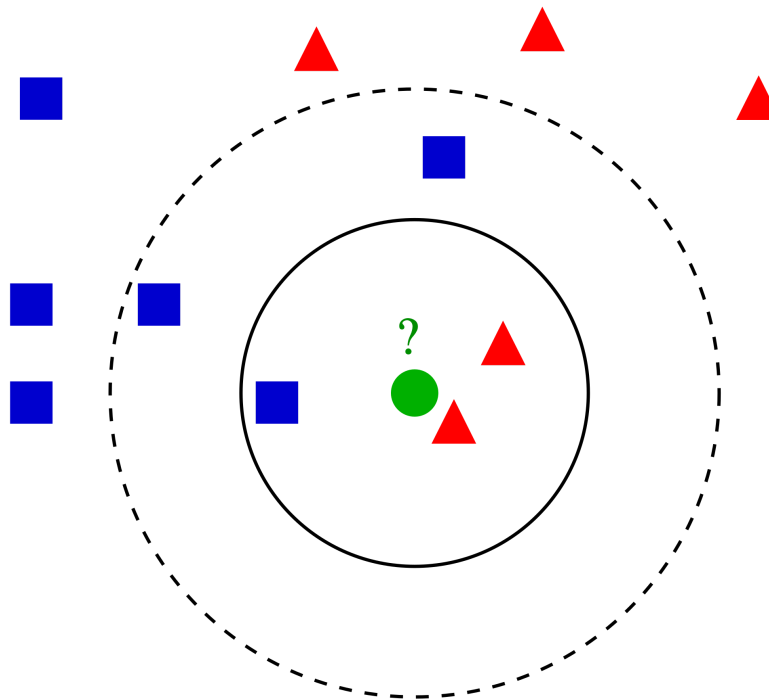
3. ElasticNET: добавляем в функцию ошибки сумму квадратов и сумму модулей весов:

$$L(y, \hat{y}) = \frac{1}{K} \sum_{i=1}^K (y - \hat{y})^2 + \alpha \sum_{j=1}^n w_j^2 + \beta \sum_{j=1}^n |w_j|$$

Регуляризация линейных моделей



Метод ближайших соседей / KNN



- классификация:
 - берем k соседей и смотрим, какой класс встречается чаще
- регрессия:
 - берем и вычисляем среднее (можно средневзвешенное) значение для нового объекта

- Что значит ближайшие?

Близкие по метрике Минковского:

$$\rho(x, y) = \left(\sum_{i=0}^d |x_i - y_i|^p \right)^{1/p}$$

- при $p = 2$ это евклидово расстояние
- при $p = 1$ Манхэттенская метрика
- при $p = \infty$ - метрика Чебышева (наибольшее из всех расстояний)

- число соседей / радиус
- метрика
- способ вычисления весов объектов

Decision Tree

- Классификация: игра состоится?

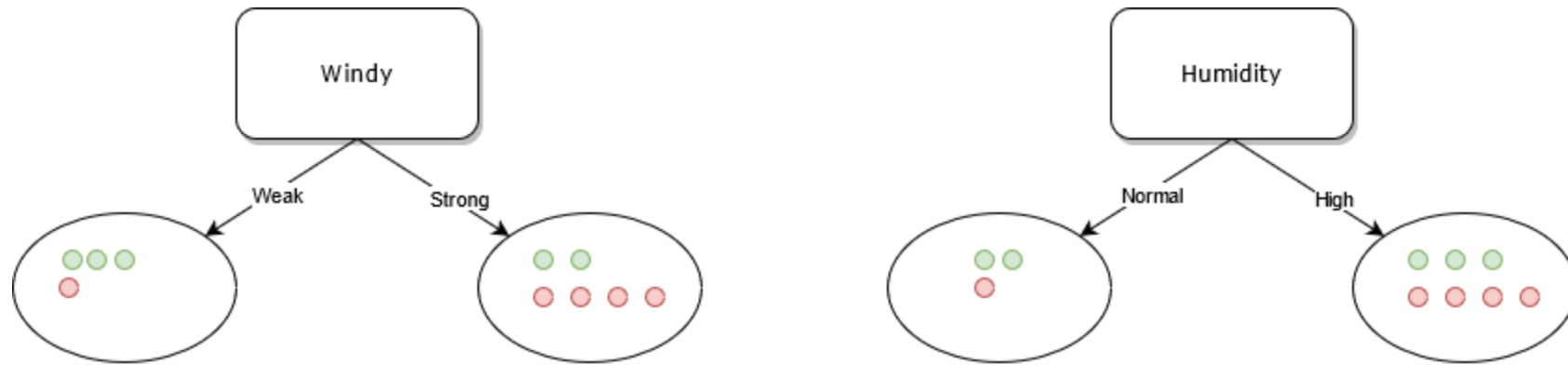
Day	Weather	Temperature	Humidity	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Cloudy	Hot	High	Weak	Yes
3	Sunny	Mild	Normal	Strong	Yes
4	Cloudy	Mild	High	Strong	Yes
5	Rainy	Mild	High	Strong	No
6	Rainy	Cool	Normal	Strong	No
7	Rainy	Mild	High	Weak	Yes
8	Sunny	Hot	High	Strong	No
9	Cloudy	Hot	Normal	Weak	Yes
10	Rainy	Mild	High	Strong	No

Какой признак самый информативный?

Тот, при использовании которого для классификации, мы получаем наиболее "чистые" подмножества объектов выборки:

- Признак `Windy` имеет два значения: `Weak` и `Strong`
 - Для `Windy=Weak`: `Play=Yes` – 3 объекта, `Play=No` – 1 объект
 - Для `Windy=Strong`: `Play=Yes` – 2 объекта, `Play=No` – 4 объекта
- Признак `Humidity` имеет два значения: `High` и `Normal`:
 - Для `Humidity=Normal`: `Play=Yes` – 2 объекта, `Play=No` – 1 объект
 - Для `Humidity=High`: `Play=Yes` – 3 объекта, `Play=No` – 4 объект

Оценки гомогенности для разбиений



Как измерить гомогенность выборки в полученных разбиениях?

- Коэффициент Джини / Gini impurity:

$$G = 1 - \sum_i p_i^2$$

- Энтропия разбиения:

$$H = -\sum p_i \log_2 p_i$$

p_i – частота объектов класса i в разбиении

Разбиение	Entropy
Исходное 	$-(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1$
Windy=Weak 	$-(0.25 \log_2 0.25 + 0.75 \log_2 0.75) = 0.81$
Windy=Strong 	$-(0.33 \log_2 0.33 + 0.66 \log_2 0.66) = 0.92$

$$IG = S_0 - \sum_{i=1}^q \frac{N_i}{N} S_i,$$

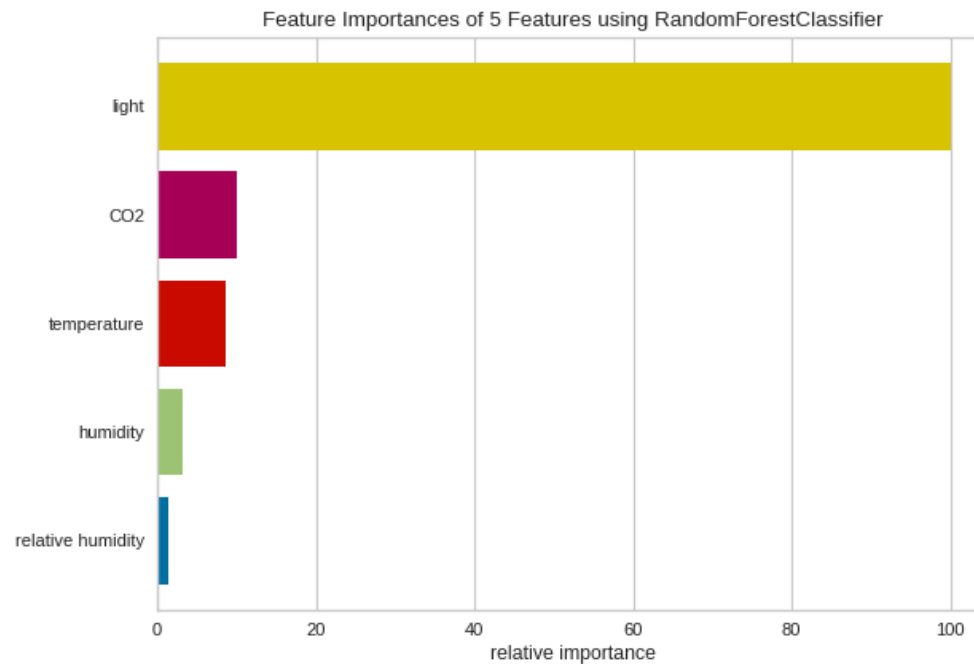
- q – число листьев (обычно 2),
- N_i – число объектов, попавших в i -ое разбиение,
- N – общее число в родительской вершине объектов,
- S_0 – *impurity metric* (*gini* или *entropy*) для исходного разбиения,
- S_i – *impurity metric* для i -го разбиения.

Для разбиения, построенного по признаку Windy :

$$IG = 1 - 0.4 \times 0.81 - 0.6 \times 0.92 = 0.12$$

Важность признаков • Feature importances

Суммарный (а в `sklearn` и нормированный) показатель уменьшения гетерогенности выборки используется для оценки важности признаков: если признак выбирался часто и сильно уменьшал энтропию или коэффициент Джини, то он является **информативным**.



- Критерий ветвления: `criterion: gini, entropy`
- Максимальная глубина: `max_depth`
- Минимальное число объектов в листе: `min_samples_leaf`
- Минимальное значение уменьшения гетерогенности для осуществления деления: `min_impurity_decrease`

`sklearn.tree` : <https://scikit-learn.org/stable/modules/tree.html>

Модель	Регрессия	Классификация
Линейные модели	<code>LinearRegression()</code>	<code>LogisticRegression()</code>
Метод ближайших соседей	<code>KNeighborsRegressor()</code>	<code>KNeighborsClassifier()</code>
Деревья решений	<code>DecisionTreeRegressor()</code>	<code>DecisionTreeClassifier()</code>

Линейные модели

- Простые и интерпретируемые
- Склонны к недообучению на сложных данных
- Примеры: `LinearRegression()`, `LogisticRegression()`

Метод ближайших соседей

- Легко объяснить и интуитивно понятен
- Может быть медленным на больших данных
- Примеры: `KNeighborsRegressor()`, `KNeighborsClassifier()`

Деревья решений

- Мощный алгоритм, склонен к переобучению
- Для уменьшения переобучения можно использовать методы регуляризации
- Может выступать в качестве составного блока для более сложных концепций
- У обученного дерева есть атрибут: `model.feature_importances_`