

ВЕБИНАРНЫЙ ФОРМАТ
Библиотеки Python для Data Science: продолжение

Урок 3.
Построение модели классификации.

Светлана Медведева

План урока

- Балансировка классов
- Схемы оценки обобщающей способности алгоритма
- Обзор алгоритмов классификации:
 - Логистическая регрессия
 - Метод опорных векторов
 - k ближайших соседей
 - Случайный лес и бустинговые алгоритмы
- Практическая часть

Балансировка классов

Собрать больше данных

Выбрать подходящую метрику качества

Попробовать разные модели, одни модели более устойчивы к несбалансированным данным, чем другие

- Штраф за ошибки при прогнозе меньшего класса
- Undersampling и Oversampling
- Создание синтетических примеров для меньшего класса

Undersampling

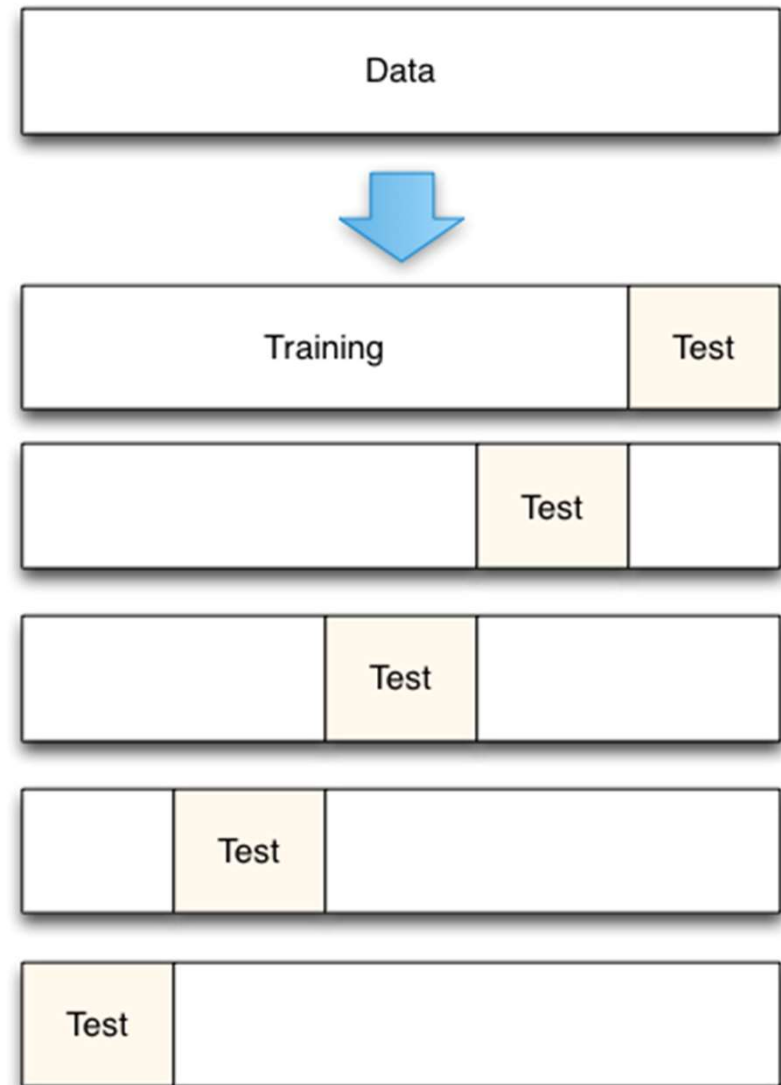


Oversampling



Модель, обладающая хорошей **обобщающей способностью**, способна предсказывать примерно на одном и том же уровне качества, как на обучающем датасете, так и на новых данных.

- Отложенная выборка
- Кросс-валидация



Обзор алгоритмов классификации:

Логистическая регрессия

Логистическая регрессия или логит-модель (англ. *logit model*) -

это статистическая модель, используемая для прогнозирования вероятности возникновения некоторого события путём его сравнения с логистической кривой.

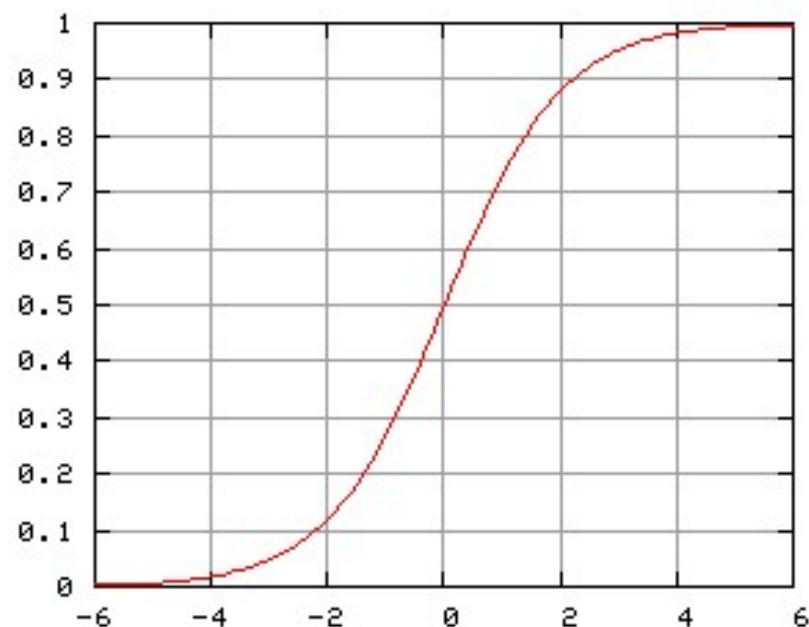
$$P\{y = 1|x\} = f(z),$$

$$\text{где } z = \theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$P\{y = 0|x\} = 1 - f(z) = 1 - f(\theta^T x)$$

$$P\{y|x\} = f(\theta^T x)^y (1 - f(\theta^T x))^{1-y}, y \in \{0,1\}$$



Логистическая регрессия: подбор параметров

Метод максимального правдоподобия

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \prod_{i=1}^m P\{y = y^{(i)} | x = x^{(i)}\}$$

$$\ln L(\theta) = \sum_{i=1}^m \ln \prod_{\theta} P\{y = y^{(i)} | x = x^{(i)}\} = \sum_{i=1}^m y^{(i)} \ln f(\theta^T x^i) + (1 - y^{(i)}) \ln(1 - f(\theta^T x^i))$$

$$\text{где } \theta^T x^{(i)} = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)}$$

Максимизация функции правдоподобия



метод градиентного спуска:

$$\theta := \theta + \alpha \nabla \ln L \theta$$

Обзор алгоритмов классификации:

Метод опорных векторов

Постановка задачи

$$\{(\mathbf{x}_1, c_1), (\mathbf{x}_2, c_2), \dots, (\mathbf{x}_n, c_n)\}$$

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

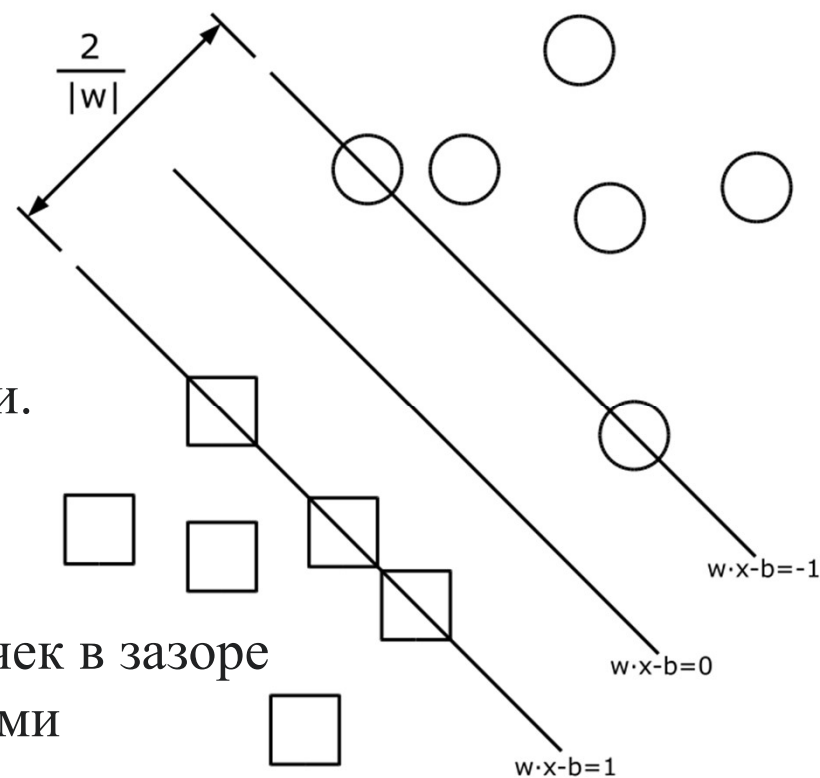
$$\mathbf{w} \cdot \mathbf{x} + b = 1$$

$$\mathbf{w} \cdot \mathbf{x} + b = -1$$

Оптимальная разделяющая гиперплоскость и || ей плоскости.

$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_i - b \geq 1, & c_i = 1 \\ \mathbf{w} \cdot \mathbf{x}_i - b \leq -1, & c_i = -1 \end{cases}$$

Условие отсутствия точек в зазоре между гиперплоскостями



Обзор алгоритмов классификации:

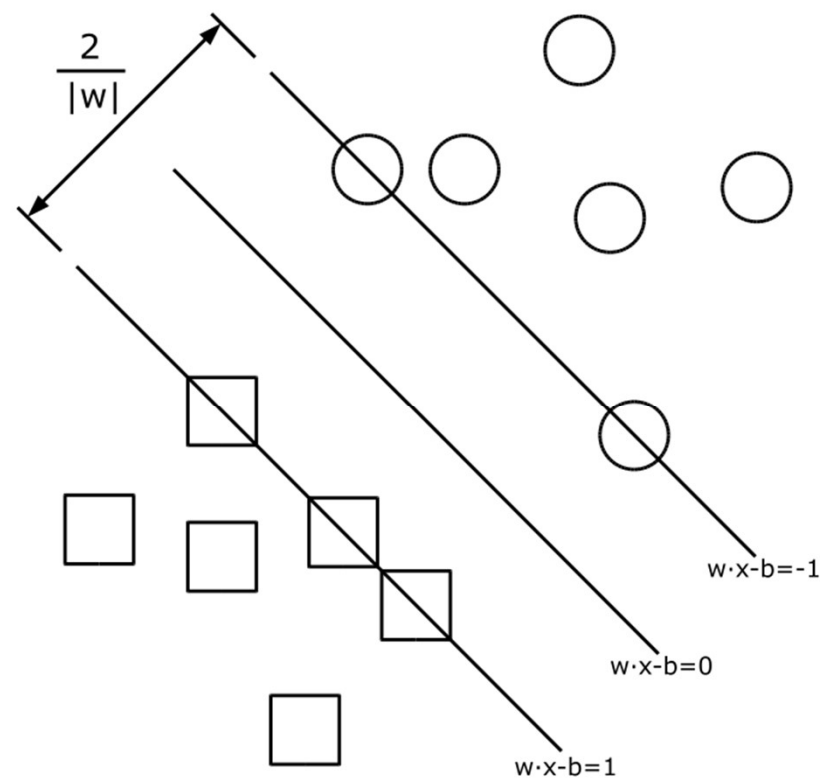
Метод опорных векторов

Случай линейной разделимости

$$\begin{cases} \|\mathbf{w}\|^2 \rightarrow \min \\ c_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1, \quad 1 \leq i \leq n. \end{cases}$$

Случай линейной неразделимости

$$\begin{cases} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \rightarrow \min_{w,b,\xi_i} \\ c_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1 - \xi_i, \quad 1 \leq i \leq n \\ \xi_i \geq 0, \quad 1 \leq i \leq n \end{cases}$$



Обзор алгоритмов классификации: k ближайших соседей

Метод k-ближайших соседей (англ. **k-nearest neighbors algorithm, k-NN**) — метрический алгоритм для автоматической классификации объектов или регрессии.

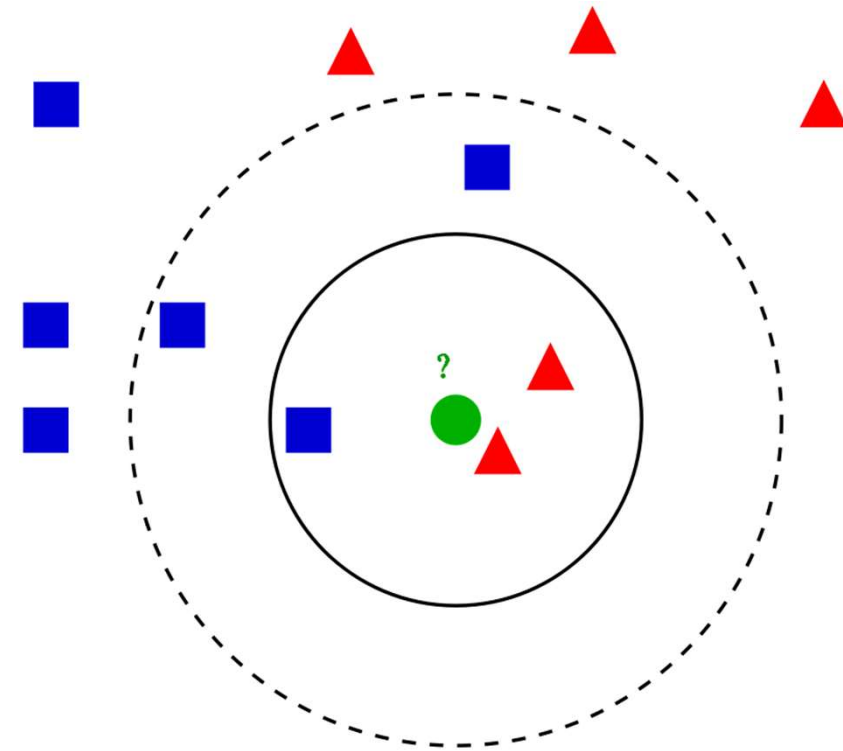
Минимакс-нормализация:

$$x' = (x - \min[X]) / (\max[X] - \min[X]) \rightarrow (0,1)$$

Z-нормализация:

$$x' = (x - M[X]) / \sigma[X] \rightarrow (-3\sigma, 3\sigma)$$

Z-нормализация:



Обзор алгоритмов классификации: k ближайших соседей

Метод k-ближайших соседей (англ. **k-nearest neighbors algorithm, k-NN**) — метрический алгоритм для автоматической классификации объектов или регрессии.

Минимакс-нормализация:

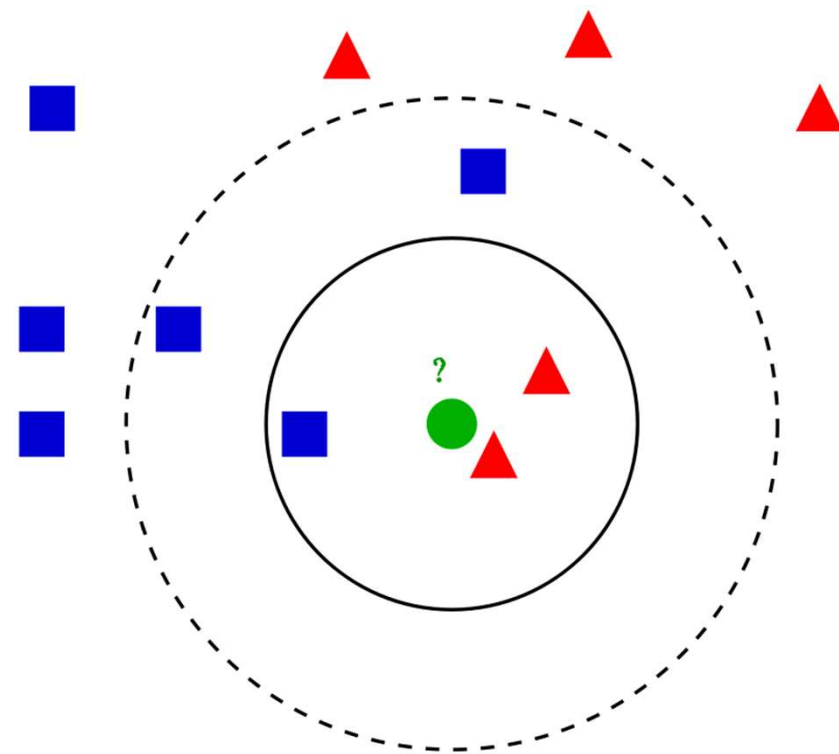
$$x' = (x - \min[X]) / (\max[X] - \min[X]) \rightarrow (0,1)$$

Z-нормализация:

$$x' = (x - M[X]) / \sigma[X] \rightarrow (-3\sigma, 3\sigma)$$

Алгоритмическая сложность для тестовой выборки:

$$O(K \cdot N \cdot M)$$

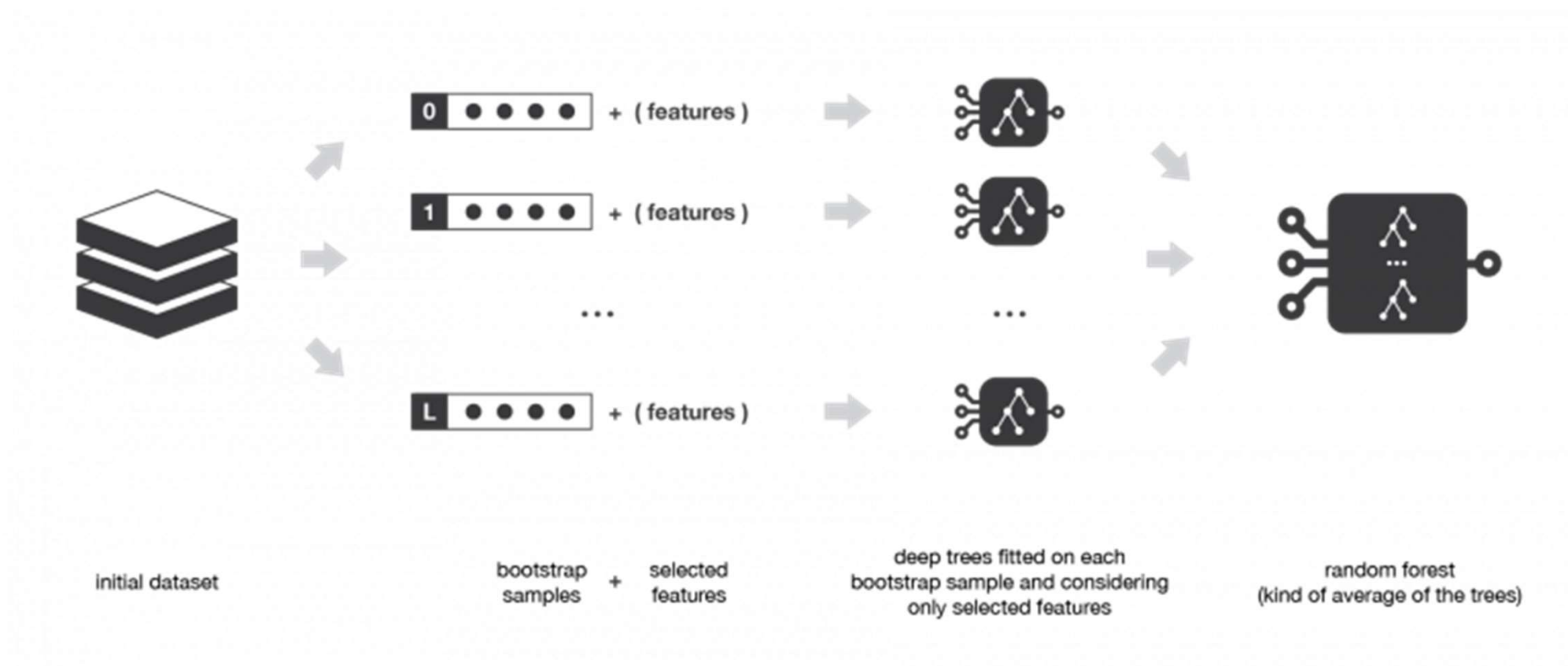


Обзор алгоритмов классификации: Случайный лес (бэггинг) и бустинговые алгоритмы

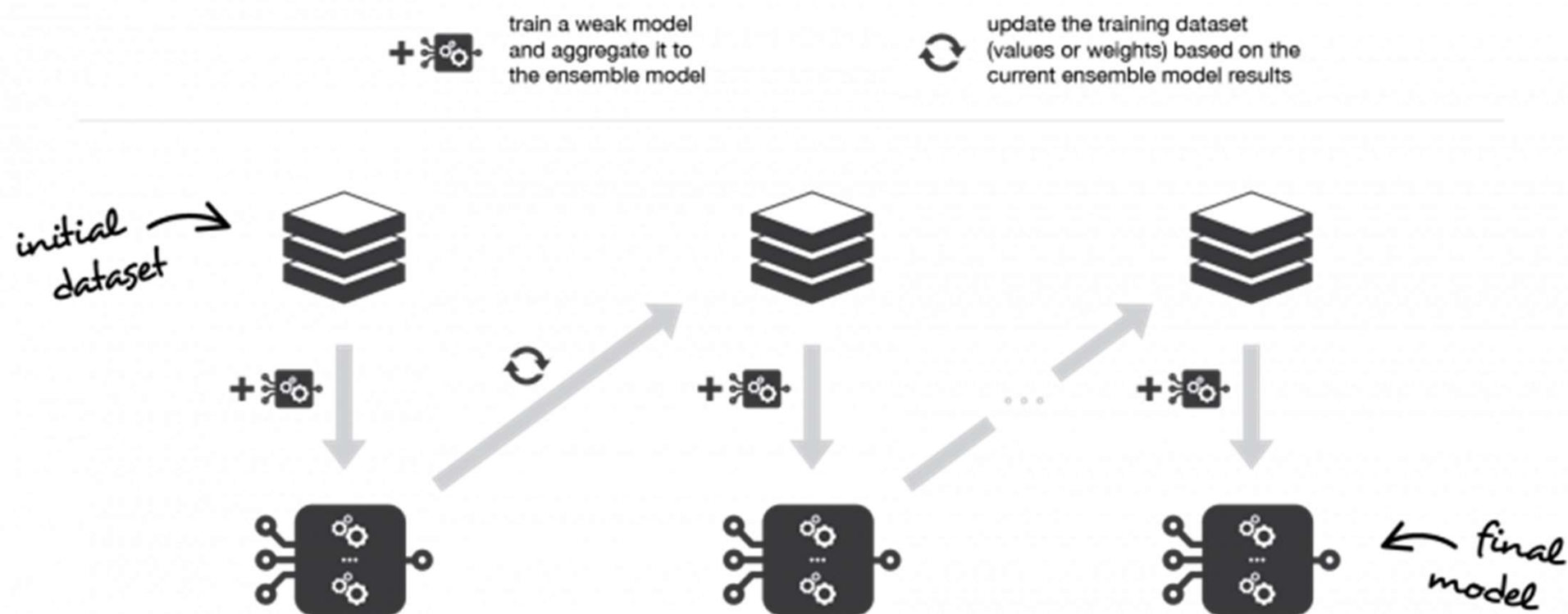
- Бэггинг
- Бустинг
- Стекинг



Обзор алгоритмов классификации: Случайный лес (бэгинг) и бустинговые алгоритмы



Обзор алгоритмов классификации: Случайный лес (бэггинг) и бустинговые алгоритмы



Итоги урока

- Балансировка классов
- Схемы оценки обобщающей способности алгоритма
- Обзор алгоритмов классификации:
 - Логистическая регрессия
 - Метод опорных векторов
 - k ближайших соседей
 - Случайный лес и бустинговые алгоритмы
- Практическая часть

Практическая часть

Спасибо за внимание!