

A decorative graphic in the bottom-left corner of the slide, consisting of a grid of colored squares. The grid is 4 squares wide and 4 squares high. The colors of the squares are: Row 1: Teal, Orange, Brown, Tan. Row 2: Orange, Brown, Tan, Brown. Row 3: Orange, Teal, Tan, Brown. Row 4: Tan, Orange, Orange, Brown.

Градиентный бустинг

Повторение

Устойчивость моделей

- $X = (x_i, y_i)_{i=1}^{\ell}$ — обучающая выборка
- Обучаем модель $a(x)$
- Ожидаем, что модель устойчивая
- То есть не сильно меняется при небольших изменениях в X
- \tilde{X} — случайная подвыборка, примерно 90% исходной

Композиция моделей

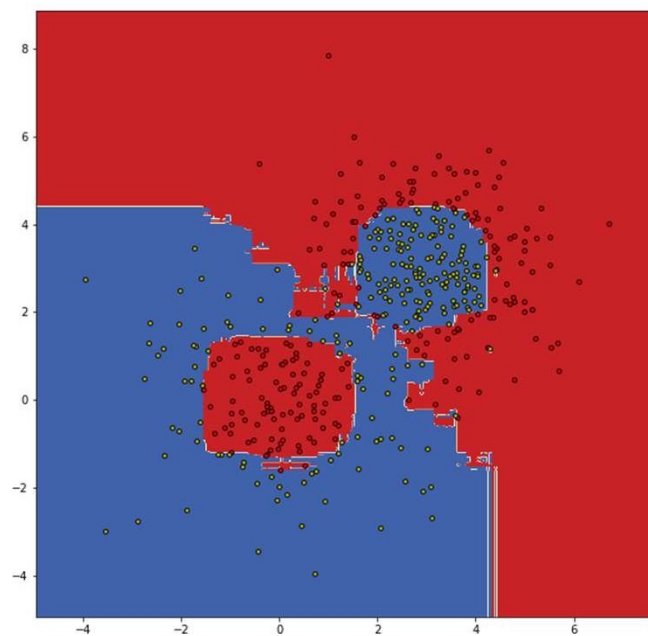
- У нас получилось N деревьев: $b_1(x), \dots, b_N(x)$
- Объединим их через голосование большинством (majority vote):

$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{n=1}^N [b_n(x) = y]$$

Выбираем класс,
который выбрало
большинство
деревьев

Количество деревьев,
выдавших класс y

Композиция моделей



Общий вид: классификация

- $b_1(x), \dots, b_N(x)$ — базовые модели
- Каждая хотя бы немного лучше случайного угадывания
- Композиция: голосование по большинству (majority vote)

$$a_N(x) = \arg \max_{y \in \mathbb{Y}} \sum_{n=1}^N [b_n(x) = y]$$

Общий вид: регрессия

- $b_1(x), \dots, b_N(x)$ — базовые модели
- Каждая хотя бы немного лучше случайного угадывания
- Композиция: усреднение

$$a_N(x) = \frac{1}{N} \sum_{n=1}^N b_n(x)$$

Бэггинг

- Bagging (bootstrap aggregating)
- Базовые модели обучаются независимо
- Каждый обучается на подмножестве обучающей выборки
- Подмножество выбирается с помощью бутстрапа

Бутстрап

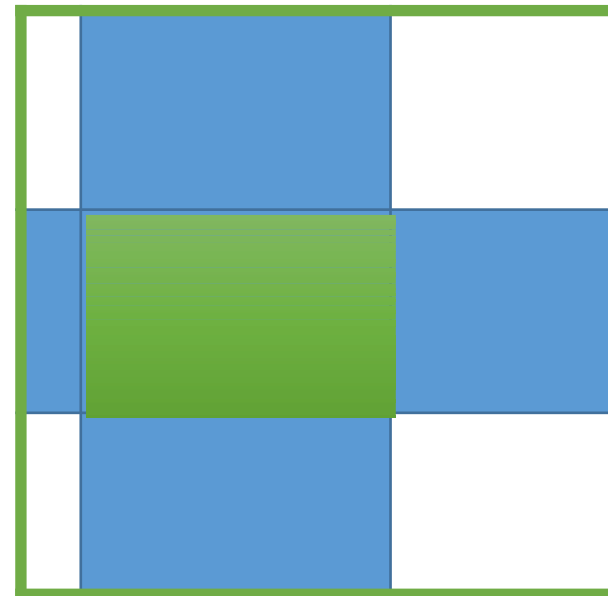
- Выборка с возвращением
- Берём ℓ элементов из X
- Пример: $\{x_1, x_2, x_3, x_4\} \rightarrow \{x_1, x_2, x_2, x_4\}$
- В подвыборке будет ℓ объектов, из них около 63.2% уникальных
- Если объект входит в выборку несколько раз, то мы как бы повышаем его вес

Случайные подпространства

- Выбираем случайное подмножество признаков
- Обучаем модель только на них
- Может быть плохо, если имеются важные признаки, без которых невозможно построить разумную модель

Виды рандомизации

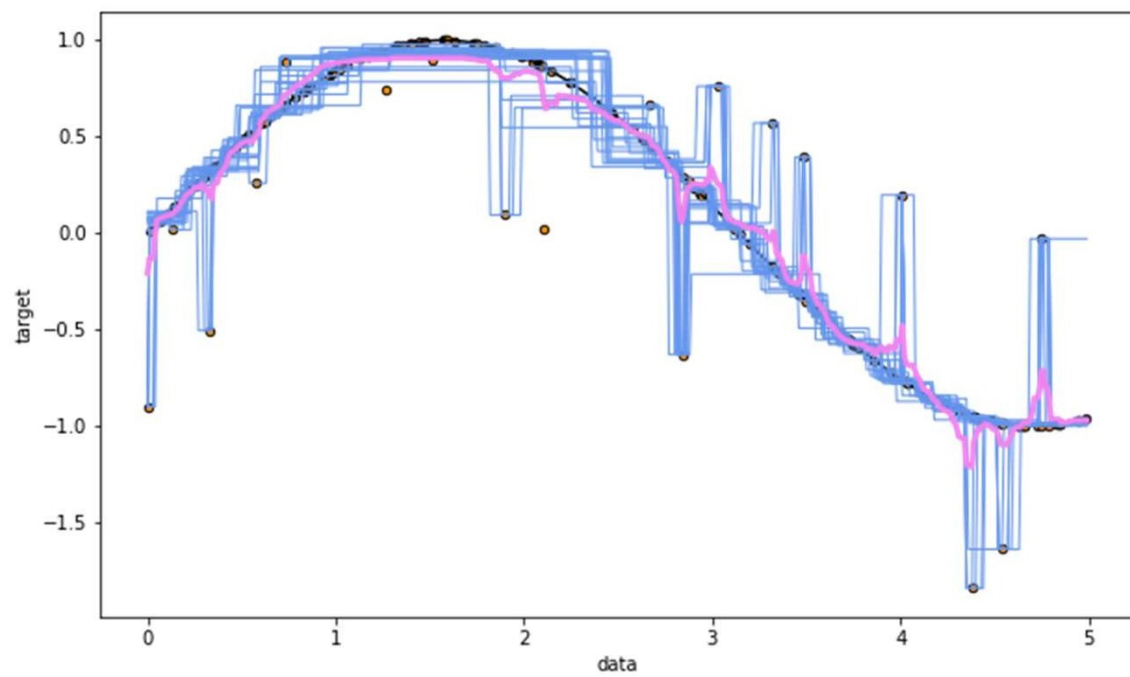
- Бэггинг: случайная подвыборка
- Случайные подпространства:
случайное подмножество признаков



Разложение ошибки на смещение и разброс

- Ошибка модели складывается из трёх компонент
- Шум (noise) — характеристика сложности и противоречивости данных
- Смещение (bias) — способность модели приблизить лучшую среди всех возможных моделей
- Разброс (variance) — устойчивость модели к изменениям в обучающей выборке

Смещение и разброс: деревья



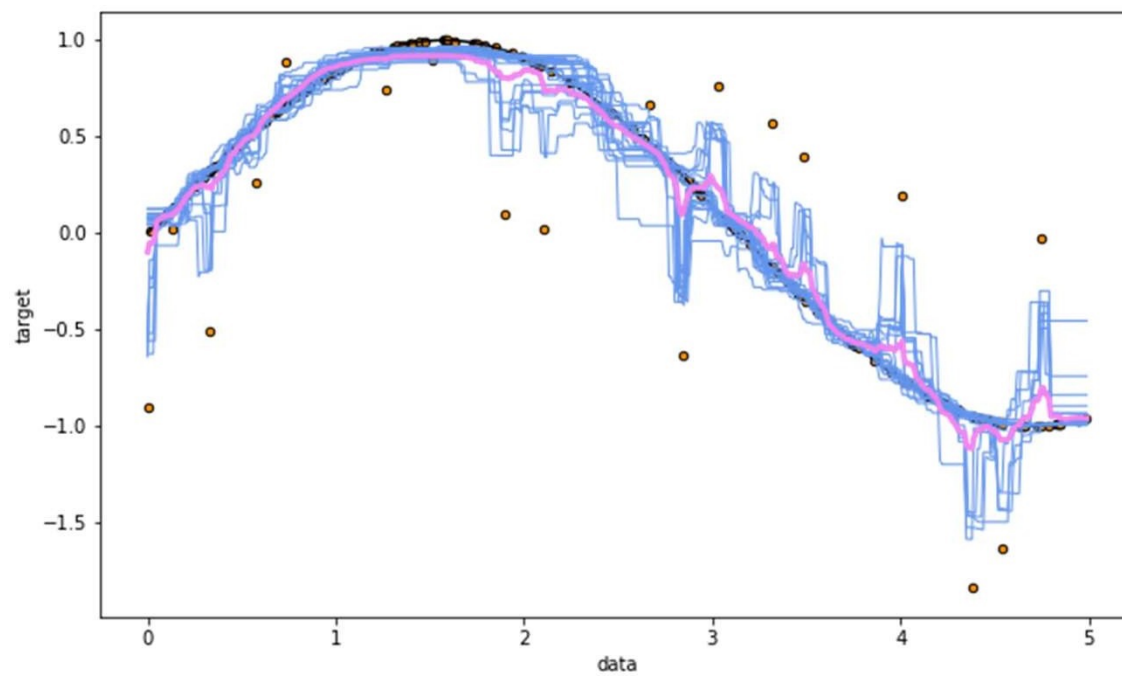
Бэггинг

- Смещение $a_N(x)$ такое же, как у $b_n(x)$
- Разброс $a_N(x)$:

$$\frac{1}{N} (\text{разброс } b_n(x)) + \text{ковариация}(b_n(x), b_m(x))$$

- Если базовые модели независимы, то разброс уменьшается в N раз!
- Чем более похожи выходы базовых моделей, тем меньше эффект от построения композиции

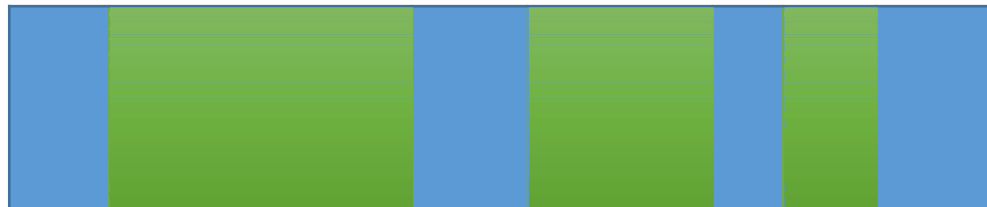
Смещение и разброс: бэггинг



Выбор предиката

$$j, t = \arg \min_{j, t} Q(R_m, j, t)$$

- Будем искать лучший предикат среди случайного подмножества признаков размера q



Случайный лес (Random Forest)

Для $n = 1, \dots, N$:

1. Сгенерировать выборку \tilde{X} с помощью бустрапа
2. Построить решающее дерево $b_n(x)$ по выборке \tilde{X}
3. Дерево строится, пока в каждом листе не окажется не более n_{min} объектов
4. Оптимальное разбиение ищется среди q случайных признаков

Выбираются заново при каждом разбиении!

Случайный лес (Random Forest)

- Регрессия:

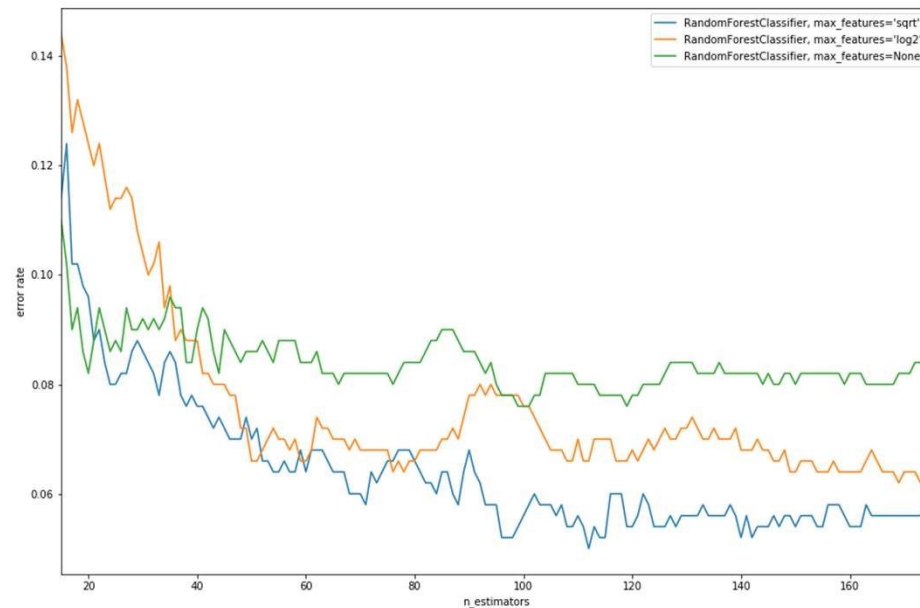
$$a(x) = \frac{1}{N} \sum_{n=1}^N b_n(x)$$

- Классификация:

$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{n=1}^N [b_n(x) = y]$$

Универсальный метод

- Ошибка сначала убывает, а затем выходит на один уровень
- Случайный лес не переобучается при росте N



Out-of-bag

- Каждое дерево обучается примерно на 63% данных
- Остальные объекты — как бы тестовая выборка для дерева
- X_n — обучающая выборка для $b_n(x)$
- Можно оценить ошибку на новых данных:

$$Q_{\text{test}} = \frac{1}{\ell} \sum_{i=1}^{\ell} L \left(y_i, \frac{1}{\sum_{n=1}^N [x_i \notin X_n]} \sum_{n=1}^N [x_i \notin X_n] b_n(x_i) \right)$$

Для каждого
объекта
выборки

Среднее ответов =
предсказание на
объекте

Суммируем ответы для всех
деревьев, в которые не
попал объект

Исправление ошибок моделей и идея бустинга

Бэггинг

- Смещение $a_N(x)$ такое же, как у $b_n(x)$
- Разброс $a_N(x)$:
- $\frac{1}{N} (\text{разброс } b_n(x)) + \text{ковариация}(b_n(x), b_m(x))$
- Если базовые модели независимы, то разброс уменьшается в N раз!
- Чем более похожи выходы базовых моделей, тем меньше эффект от построения композиции

Проблемы бэггинга

- Если базовая модель окажется смещённой, то и композиция не справится с задачей
- Базовые модели долго обучать и применять, дорого хранить

Идея бустинга

- Бустинг (англ. boosting — усиление)
- Возьмём простые базовые модели
- Будем строить композицию последовательно и жадно
- Каждая следующая модель будет строиться так, чтобы максимально корректировать ошибки построенных моделей

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение первой модели (b_1):

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, b_1(x_i)) \rightarrow \min_{b_1(x)}$$

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение первой модели (b_1):

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, b_1(x_i)) \rightarrow \min_{b_1(x)}$$

- Пытаемся подобрать модель b_1 , минимизирующую ошибку

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели (b_N) :

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, \underbrace{a_{N-1}(x_i)}_{\text{фиксировано}} + \underbrace{b_N(x_i)}_{\text{учится}}) \rightarrow \min_{b_N(x)}$$

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели (b_N) :

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Пытаемся подобрать модель b_N , минимизирующую ошибку итоговой композиции $a_N = a_{N-1} + b_N$

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели (b_N) :

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

Идея бустинга

- Обучение N -й модели (b_N) :

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Нет такого алгоритма машинного обучения, который учил бы «добавку»
- Попробуем понять, как можно сформулировать задачу с точки зрения ML

Резюме

- В бустинге базовые модели обучаются последовательно
- Каждая следующая корректирует ошибки уже построенных
- В общем случае получается функционал, на который может быть сложно обучать деревья

Бустинг для среднеквадратичной ошибки

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

Идея бустинга

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- MSE:

$$L(y, \hat{y}) = (y - \hat{y})^2$$

Бустинг для MSE

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (a_{N-1}(x_i) + b_N(x_i) - y_i)^2 \rightarrow \min_{b_N(x)}$$

Бустинг для MSE

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - (y_i - a_{N-1}(x_i)) \right)^2 \rightarrow \min_{b_N(x)}$$

Бустинг для MSE

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - \underbrace{(y_i - a_{N-1}(x_i))}_{s_i^{(N)}} \right)^2 \rightarrow \min_{b_N(x)}$$

$$s_i^{(N)} = y_i - a_{N-1}(x_i) \text{ — остатки}$$

Бустинг для MSE

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

- $s_i^{(N)} = y_i - a_{N-1}(x_i)$ — **остатки**

Бустинг для MSE

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

- $s_i^{(N)} = y_i - a_{N-1}(x_i)$ — остатки
- Если b_N научится выдавать остатки $s_i^{(N)}$, то задача будет решена идеально

Бустинг для MSE

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

- $s_i^{(N)} = y_i - a_{N-1}(x_i)$ — остатки
- Если b_N научится выдавать остатки $s_i^{(N)}$, то задача будет решена идеально

$$y_i = a_{N-1}(x_i) + s_i^{(N)} = a_{N-1}(x_i) + b_N(x_i)$$

Пример

- $y_i = 12$
- $a_{N-1}(x_i) = 10$
- $s_i^{(\overline{N})} = ?$
- $b_N(x_i) = ?$
- $a_N(x_i) = ?$
-

Пример

- $y_i = 12$
- $a_{N-1}(x_i) = 10$
- $s_i^{(\overline{N})} = 2$
- $b_N(x_i) = 2$
- $a_N(x_i) = 12$
-

Первая итерация

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (b_1(x_i) - y_i)^2 \rightarrow \min_{b_1(x)}$$

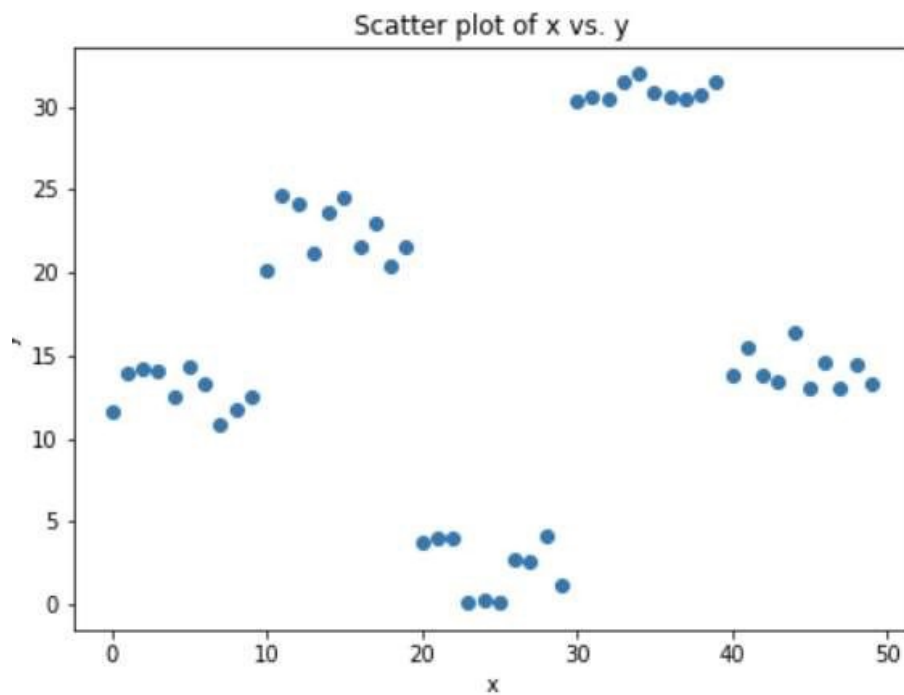
Вторая итерация

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_2(x_i) - (y_i - b_1(x_i)) \right)^2 \rightarrow \min_{b_2(x)}$$

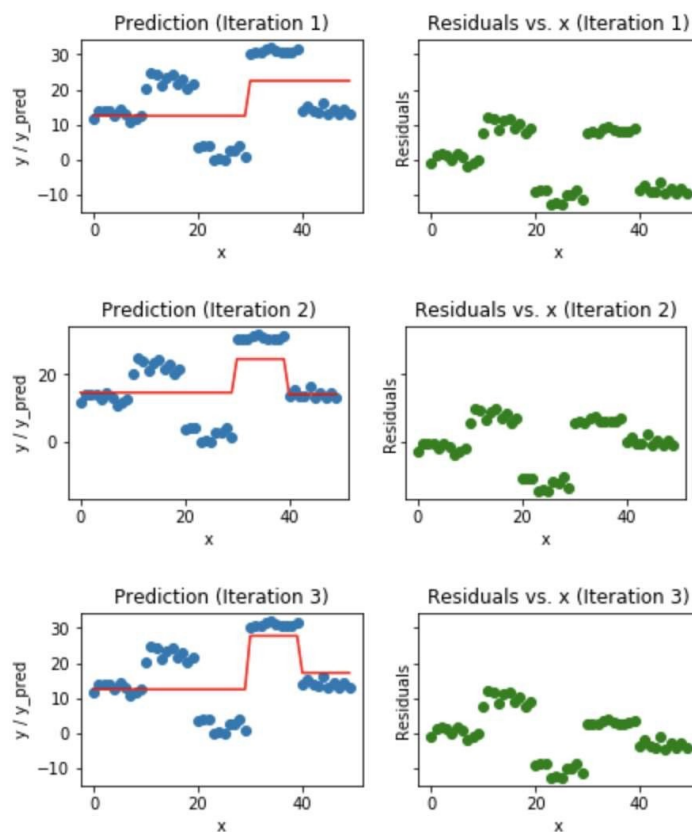
Третья итерация

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_3(x_i) - (y_i - b_1(x_i) - b_2(x_i)) \right)^2 \rightarrow \min_{b_3(x)}$$

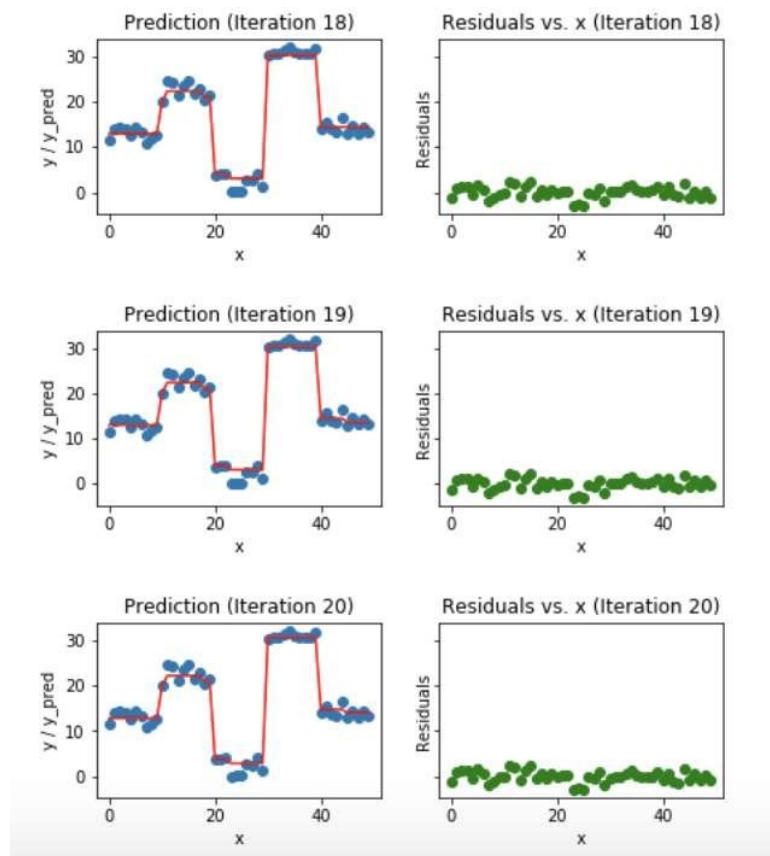
Визуализация



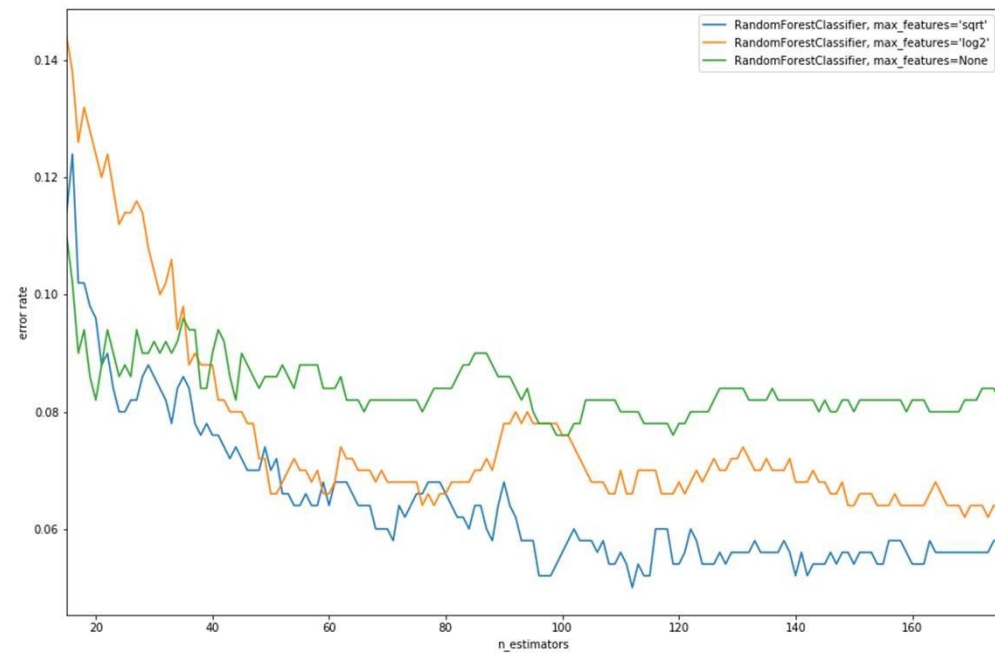
Визуализация



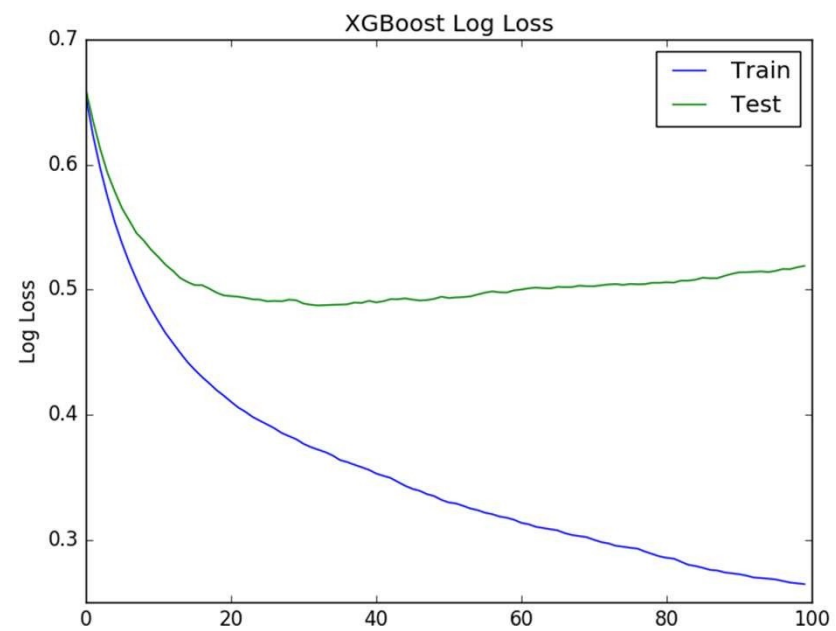
Визуализация



Random Forest



Ошибка бустинга на обучении и тесте



Резюме

- В случае с MSE обучение базовых моделей сводится к обычной процедуре обучения с заменой целевой переменной
- Бустинг может переобучаться, поэтому надо следить за ошибкой на тестовой выборке

Сложности с произвольной функцией потерь

Задача обучения базовой модели

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

Задача обучения базовой модели

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i - a_{N-1}(x_i), b_N(x_i)) \rightarrow \min_{b_N(x)}$$

Задача обучения базовой модели

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i - a_{N-1}(x_i), b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Хотим, чтобы модель b_N выдавала $y_i - a_{N-1}(x_i)$

Логистическая функция потерь

$$a_N(x) = \text{sign} \sum_{n=1}^N b_n(x)$$

$$L(y, z) = \log(1 + \exp(-yz))$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(- (y_i - a_{N-1}(x_i)) b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

Логистическая функция потерь

$$a_N(x) = \text{sign} \sum_{n=1}^N b_n(x)$$

$$L(y, z) = \log(1 + \exp(-yz))$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(- \underbrace{(y_i - a_{N-1}(x_i))}_{\text{выдает } \pm 1} b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

Логистическая функция потерь

$$a_N(x) = \text{sign} \sum_{n=1}^N b_n(x)$$

$$L(y, z) = \log(1 + \exp(-yz))$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(- (y_i - a_{N-1}(x_i)) b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

- Если $y_i = a_{N-1}(x_i)$, то объект не участвует в обучении
- Тогда модель $b_N(x_i)$ может выдавать что угодно и испортить композицию

Логистическая функция потерь

$$a_N(x) = \text{sign} \sum_{n=1}^N b_n(x)$$

$$L(y, z) = \log(1 + \exp(-yz))$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(- (y_i - a_{N-1}(x_i)) b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

- Иначе $y_i - a_{N-1}(x_i) = \pm 2$

Логистическая функция потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(-\frac{y_i - a_{N-1}(x_i)}{2} b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

- Если $y_i = a_{N-1}(x_i)$, то объект не участвует в обучении
- Если $y_i \neq a_{N-1}(x_i)$, то базовая модель учится выдавать корректный класс

Логистическая функция потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(-\frac{y_i - a_{N-1}(x_i)}{2} b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

- $y_i = +1, \sum_{n=1}^{N-1} b_n(x_i) = -0.5 \rightarrow \text{надо } b_N(x_i) > 0.5$
- $y_i = +1, \sum_{n=1}^{N-1} b_n(x_i) = -100 \rightarrow \text{надо } b_N(x_i) > 100$

Логистическая функция потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log \left(1 + \exp \left(- \frac{y_i - a_{N-1}(x_i)}{2} b_N(x_i) \right) \right) \rightarrow \min_{b_N(x)}$$

$$y_i = +1, \sum_{n=1}^{N-1} b_n(x_i) = -0.5 \rightarrow \text{надо } b_N(x_i) > 0.5$$

$$y_i = +1, \sum_{n=1}^{N-1} b_n(x_i) = -100 \rightarrow \text{надо } b_N(x_i) > 100$$

- Но на обоих объектах будет одинаково максимизироваться отступ
- На объектах с корректными ответами никак не контролируется выход $b_N(x)$

MSLE

- Mean Squared Logarithmic Error (среднеквадратичная логарифмическая ошибка)

$$L(y, z) = (\log(z + 1) - \log(y + 1))^2$$

MSLE

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

$$L(y, z) = (\log(z + 1) - \log(y + 1))^2$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\log(b_N(x_i) + 1) - \log(y_i - a_{N-1}(x_i) + 1))^2 \rightarrow \min_{b_N(x)}$$

MSLE

$$a_N(x) = \sum_{n=1}^N b_n(x)$$

$$L(y, z) = (\log(z + 1) - \log(y + 1))^2$$

- Может, просто обучаться на остатки, как в MSE?

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\log(b_N(x_i) + 1) - \log(y_i - a_{N-1}(x_i) + 1))^2 \rightarrow \min_{b_N(x)}$$

- Аргумент второго логарифма может оказаться отрицательным

MSLE

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\log(b_N(x_i) + 1) - \log(y_i - a_{N-1}(x_i) + 1))^2 \rightarrow \min_{b_N(x)}$$

y_i	$a_{N-1}(x_i)$	$b_N(x_i)$	Улучшение MSLE композиции	Улучшение функционала базовой модели
1000	100	2	0.09	13.7
2	0	2	1.2	1.2

Резюме

- Нельзя заменить обучение добавки к композиции на обучение базовой модели на отклонение от ответов
- Не учитываются особенности функции потерь

Градиентный бустинг в общем виде

Задача обучения базовой модели

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

Задача обучения базовой модели

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Как посчитать, куда и как сильно сдвигать $a_{N-1}(x_i)$, чтобы уменьшить ошибку?

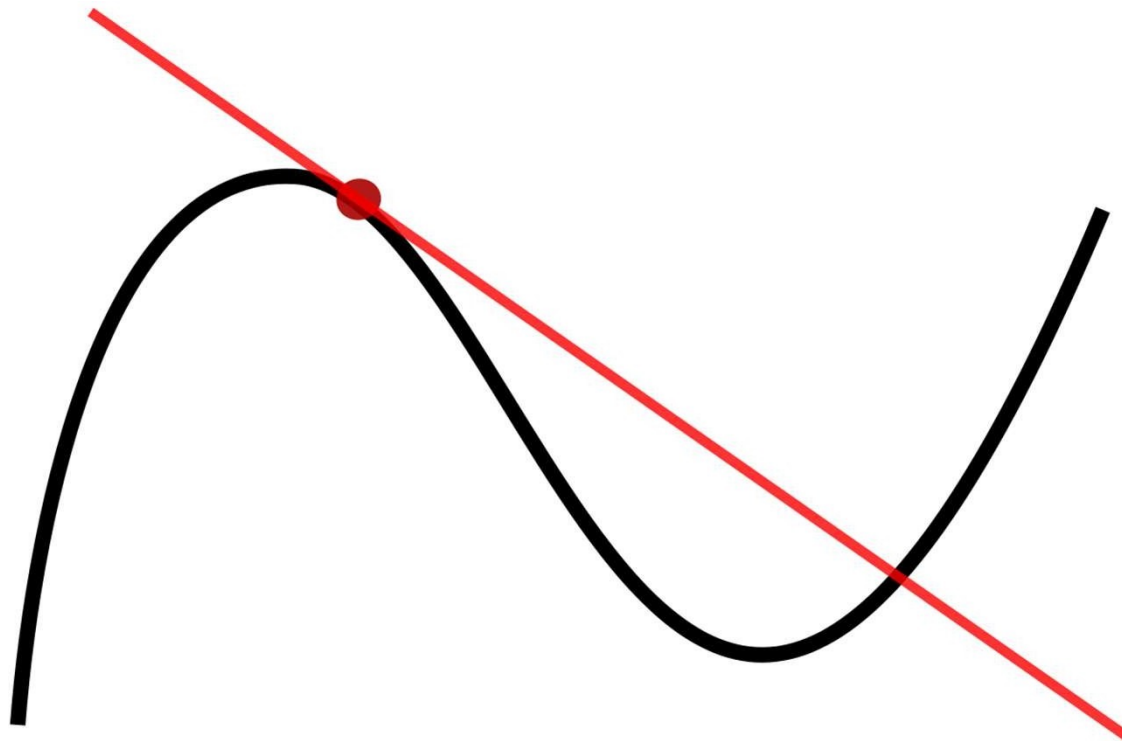
Задача обучения базовой модели

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Как посчитать, куда и как сильно сдвигать $a_{N-1}(x_i)$, чтобы уменьшить ошибку?
- Посчитать производную

Производная



Задача обучения базовой модели

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Посчитаем антипроизводную:

$$s_i^{(N)} = - \frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)}$$

Задача обучения базовой модели

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + b_N(x_i)) \rightarrow \min_{b_N(x)}$$

- Посчитаем антипроизводную:

$$s_i^{(N)} = - \underbrace{\frac{\partial}{\partial z}}_{\text{производная по } z} \underbrace{L(y_i, z)}_{\substack{\text{прогноз} \\ \text{модели}}} \bigg|_{\underbrace{z = a_{N-1}(x_i)}_{\text{в качестве } z \text{ подставляем} \\ \text{в результат предсказание} \\ \text{композиции}}}$$

правильный ответ

Задача обучения базовой модели

- Посчитаем антипроизводную:

$$s_i^{(N)} = - \frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)}$$

- Знак показывает, в какую сторону сдвигать прогноз на x_i , чтобы уменьшить ошибку композиции на нём
- Величина показывает, как сильно можно уменьшить ошибку, если сдвинуть прогноз
- Если ошибка почти не сдвинется, то нет смысла что-то менять

Градиентный бустинг

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

$$s_i^{(N)} = - \frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} \text{ — СДВИГИ}$$

Градиентный бустинг

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

$$s_i^{(N)} = - \frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} \text{ — сдвиги}$$

- Таким образом, мы обучаем базовую модель b_N так, чтобы на x_i она выдавала $s_i^{(N)}$

Градиентный бустинг

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

$$s_i^{(N)} = -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} \text{ — сдвиги}$$

- Как бы градиентный спуск в пространстве алгоритмов
- Базовая модель будет делать корректировки на объектах так, чтобы как можно сильнее уменьшить ошибку композиции
- Сдвиги учитывают особенности функции потерь

Градиентный бустинг для MSE

$$\begin{aligned} s_i^{(N)} &= -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} = -\frac{\partial}{\partial z} \frac{1}{2} (z - y_i)^2 \Big|_{z=a_{N-1}(x_i)} = \\ &= -(a_{N-1}(x_i) - y_i) = y_i - a_{N-1}(x_i) \end{aligned}$$

Градиентный бустинг для MSE

$$\begin{aligned} s_i^{(N)} &= -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} = -\frac{\partial}{\partial z} \frac{1}{2} (z - y_i)^2 \Big|_{z=a_{N-1}(x_i)} = \\ &= -(a_{N-1}(x_i) - y_i) = y_i - a_{N-1}(x_i) \end{aligned}$$

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - (y_i - a_{N-1}(x_i)) \right)^2 \rightarrow \min_{b_N(x)}$$

Градиентный бустинг для асимметричной функции

$$L(y, z) = \frac{1}{2} ([z < y](z - y)^2 + 5[z \geq y](z - y)^2)$$

$$\begin{aligned} s_i^{(N)} &= -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} = \\ &= [z < y](y - z) + 5[z \geq y](y - z) \end{aligned}$$

Градиентный бустинг для асимметричной функции

$$s_i^{(N)} = [z < y](y - z) + 5[z \geq y](y - z)$$

- $y_i = 10, a_{N-1}(x_i) = 5: s_i = 5$
- $y_i = 10, a_{N-1}(x_i) = 15: s_i = -25$

Градиентный бустинг для логистической функции потерь

$$\begin{aligned} s_i^{(N)} &= -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} = \\ &= -\frac{\partial}{\partial z} \log(1 + \exp(-y_i z)) \Big|_{z=a_{N-1}(x_i)} = \\ &= \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \end{aligned}$$

Градиентный бустинг для логистической функции потерь

$$\begin{aligned} s_i^{(N)} &= -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)} = \\ &= -\frac{\partial}{\partial z} \log(1 + \exp(-y_i z)) \Big|_{z=a_{N-1}(x_i)} = \\ &= \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \end{aligned}$$

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \right)^2 \rightarrow \min_{b_N(x)}$$

Градиентный бустинг для логистической функции потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \right)^2 \rightarrow \min_{b_N(x)}$$

- Отступ большой положительный: $\frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \approx 0$
- Отступ большой отрицательный: $\frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \approx \pm 1$

Градиентный бустинг для логистической функции потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \right)^2 \rightarrow \min_{b_N(x)}$$

- Отступ большой положительный: $\frac{y_i}{\underbrace{1 + \exp(y_i a_{N-1}(x_i))}_{\rightarrow \infty}} \approx 0$
- Отступ большой отрицательный: $\frac{y_i}{\underbrace{1 + \exp(y_i a_{N-1}(x_i))}_{\rightarrow 0}} \approx \pm 1$

Градиентный бустинг для логистической функции потерь

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - \frac{y_i}{1 + \exp(y_i a_{N-1}(x_i))} \right)^2 \rightarrow \min_{b_N(x)}$$

- $y_i = +1, a_{N-1}(x_i) = -0.7: s_i = 0.67$
- $y_i = +1, a_{N-1}(x_i) = 2: s_i = 0.12$

Резюме

- Чтобы учесть особенности функции потерь, можно посчитать её производные в точке текущего прогноза композиции
- Базовую модель будем обучать на эти производные (со знаком минус)

Гиперпараметры и регуляризация в бустинге

Градиентный бустинг

$$a_N(x) = a_{N-1}(x_i) + b_N(x_i)$$

- Обучение N -й модели:

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \left(b_N(x_i) - s_i^{(N)} \right)^2 \rightarrow \min_{b_N(x)}$$

- $s_i^{(N)} = -\frac{\partial}{\partial z} L(y_i, z) \Big|_{z=a_{N-1}(x_i)}$ — сдвиги

Глубина деревьев

- Градиентный бустинг уменьшает смещение базовых моделей
- Разброс может увеличиться

Глубина деревьев

- Градиентный бустинг уменьшает смещение базовых моделей
- Разброс может увеличиться
- Поэтому в качестве базовых моделей стоит брать...

Глубина деревьев

- Градиентный бустинг уменьшает смещение базовых моделей
- Разброс может увеличиться
- Поэтому в качестве базовых моделей стоит брать **неглубокие** деревья

Гиперпараметры

- Глубина базовых деревьев
- Число деревьев

Проблемы бустинга

- Сдвиги показывают направление, в котором надо сдвинуть композицию на всех объектах обучающей выборки
- Базовые модели, как правило, очень простые
- Могут не справиться с приближением этого направления

Проблемы бустинга

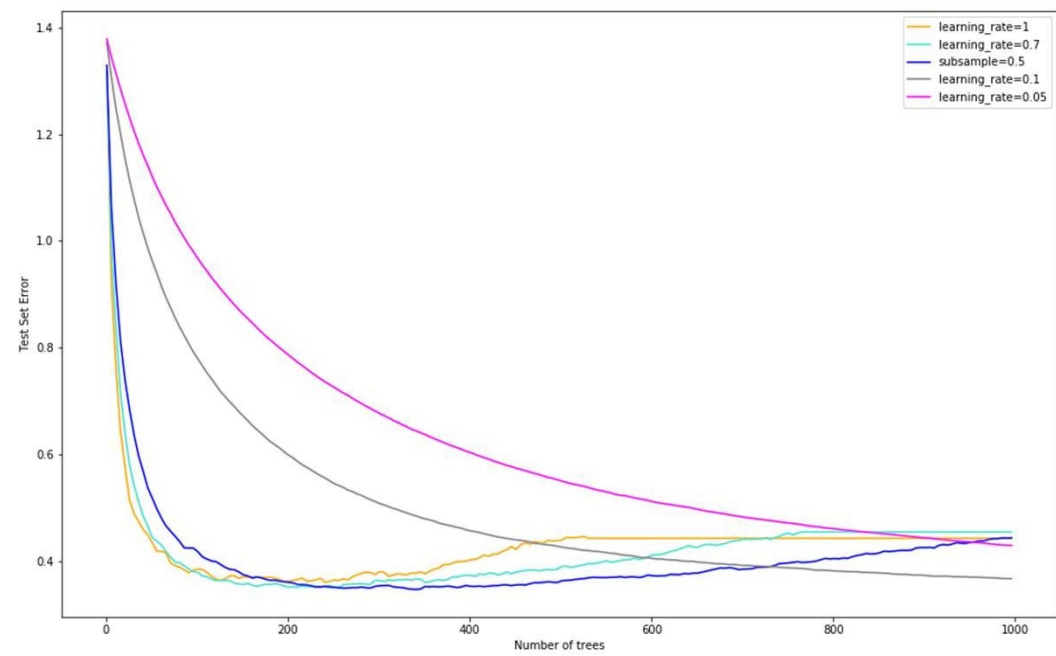
- Сдвиги показывают направление, в котором надо сдвинуть композицию на всех объектах обучающей выборки
- Базовые модели, как правило, очень простые
- Могут не справиться с приближением этого направления
- Выход: добавлять деревья в композицию с небольшим весом

Длина шага

$$a_N(x) = a_{N-1}(x_i) + \eta b_N(x_i)$$

- $\eta \in (0, 1]$ — длина шага (learning rate)
- Можно сказать, что это регуляризация композиции
- Снижает вклад каждой модели в композицию
- Чем меньше η , тем больше надо деревьев

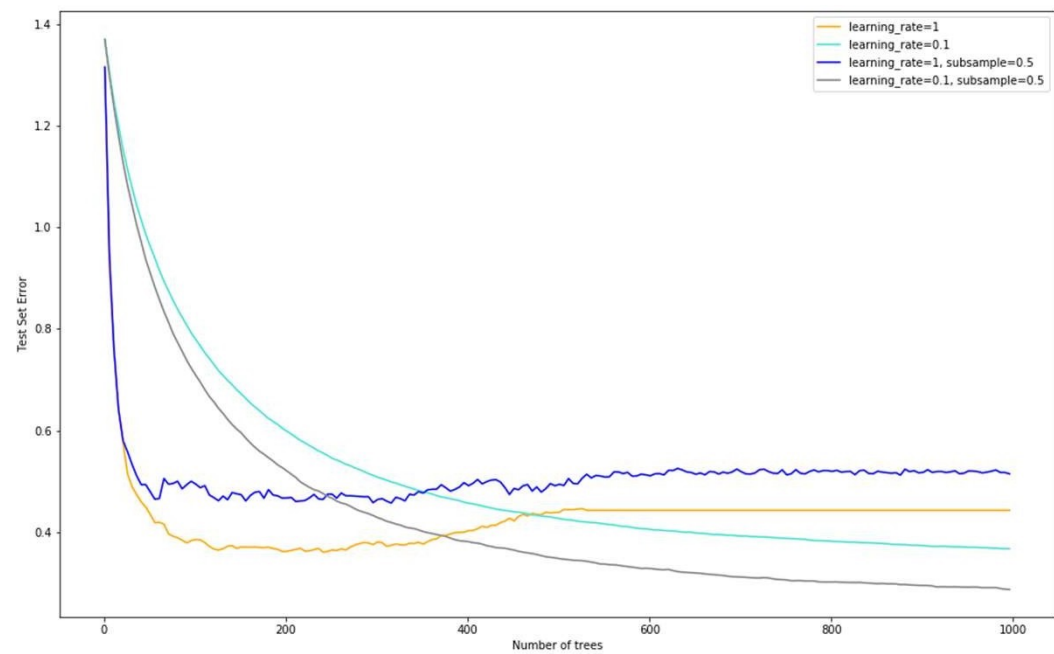
Длина шага



Рандомизация

- Можно обучать деревья на случайных подмножествах признаков
- Бустинг уменьшает смещение, поэтому итоговая композиция всё равно получится качественной
- Может снизить переобучение
- Можно обучать деревья на подмножествах объектов — способ борьбы с шумом в данных

Рандомизация



Гиперпараметры

- Глубина базовых деревьев
- Число деревьев
- Длина шага
- Размер подвыборки для обучения
- и т.д.

Резюме

- Чтобы снизить переобучение, можно добавлять модели в композицию с небольшими весами
- Также может помочь обучение моделей на подвыборках

Полезные ссылки

- <https://www.gormanalysis.com/blog/gradient-boosting-explained/>
- <https://youtu.be/3CC4N4z3GJc>
- https://en.wikipedia.org/wiki/Gradient_boosting
- <https://dyakonov.org/2017/06/09/градиентный-бустинг/>

Спасибо за внимание!



Ildar Safilo

@Ildar_Saf

irsafilo@gmail.com

<https://www.linkedin.com/in/isafilo/>