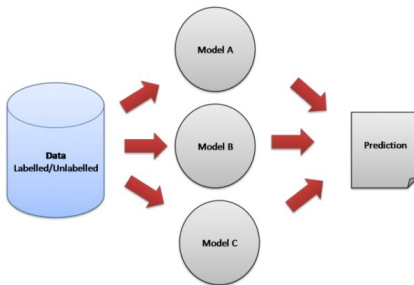


Композиции алгоритмов

Виктор Китов

victorkitov.github.io



Курс поддержан
фондом
'Интеллект'



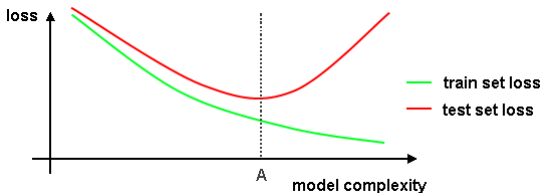
Победитель
конкурса VK среди
курсов по IT



Содержание

- 1 Разложение на смещение и разброс
- 2 Композиции алгоритмов
- 3 Примеры использования ансамблей
- 4 Ансамбли против недообучения
- 5 Ансамбли против переобучения

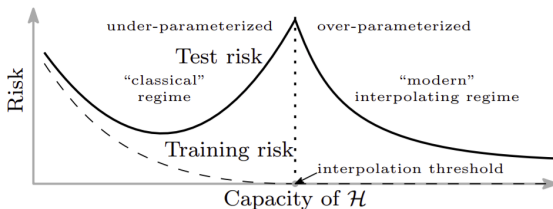
Средние потери в зависимости от сложности модели



Комментарии:

- ожидаемые потери на тестовой выборке выше потерь на обучающей.
- слева от A: модель слишком простая, недообучение.
- справа от A: модель слишком сложная, переобучение

Дальнейшее усложнение модели



- Некоторые эмпирические наблюдения свидетельствуют, что для слишком сложных моделей (многослойные нейросети) качество выше ожиданий.
- Феномен double descent risk curve¹.

¹Belkin, M., Hsu, D., Ma, S., & Mandal, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off.

Разложение на смещение и разброс

- Распределение реальных данных $y = f(x) + \varepsilon$
 - шум не зависит от x и старых наблюдений
 - хотим оценить $f(x)$
- Зависимость оценивается по $(X, Y) = \{(x_n, y_n), n = 1, 2 \dots N\}$.
- Восстановленная зависимость $\hat{f}(x)$.
- x - фикс. объект для прогноза.
- Шум ε не зависит от X, Y , $\mathbb{E}\varepsilon = 0$

Разложение на смещение и разброс

Разложение на смещение и разброс
(bias-variance decomposition):

$$\begin{aligned}\mathbb{E}_{X,Y,\varepsilon}\{[\hat{f}(x) - y(x)]^2\} &= \left(\mathbb{E}_{X,Y}\{\hat{f}(x)\} - f(x)\right)^2 \\ &\quad + \mathbb{E}_{X,Y}\left\{[\hat{f}(x) - \mathbb{E}_{X,Y}\hat{f}(x)]^2\right\} + \mathbb{E}\varepsilon^2\end{aligned}$$

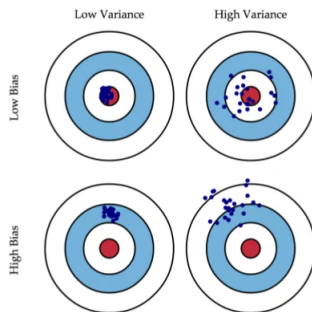
- Интуиция:

MSE = смещение² + дисперсия + неснижаемая ошибка

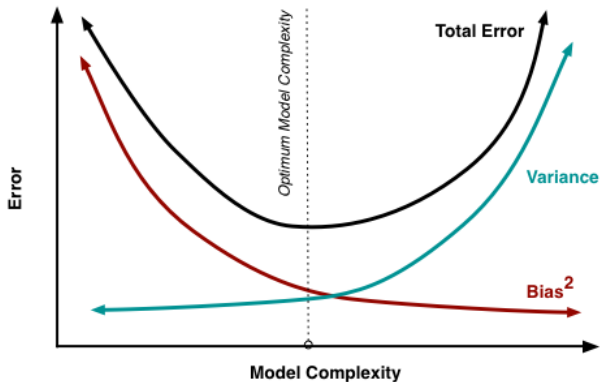
- смещение (bias) - степень недообученности
- дисперсия (variance) - степень переобученности
- неснижаемая ошибка (irreducible error) - определяется шумом в данных

Интуиция разложения

$$\begin{aligned}\mathbb{E}_{X,Y,\varepsilon}\{[\hat{f}(x) - y(x)]^2\} &= \left(\mathbb{E}_{X,Y}\{\hat{f}(x)\} - f(x)\right)^2 \\ &\quad + \mathbb{E}_{X,Y}\left\{[\hat{f}(x) - \mathbb{E}_{X,Y}\hat{f}(x)]^2\right\} + \mathbb{E}\varepsilon^2\end{aligned}$$



Средние потери в зависимости от сложности модели



Доказательство разложения

Обозначим для краткости $f = f(x)$, $\hat{f} = \hat{f}(x)$, $\mathbb{E} = \mathbb{E}_{X,Y,\varepsilon}$.

Доказательство разложения

Обозначим для краткости $f = f(x)$, $\hat{f} = \hat{f}(x)$, $\mathbb{E} = \mathbb{E}_{X,Y,\varepsilon}$.

$$\begin{aligned}\mathbb{E}(\hat{f} - f)^2 &= \mathbb{E}(\hat{f} - \mathbb{E}\hat{f} + \mathbb{E}\hat{f} - f)^2 = \mathbb{E}(\hat{f} - \mathbb{E}\hat{f})^2 + (\mathbb{E}\hat{f} - f)^2 \\ &\quad + 2\mathbb{E}[(\hat{f} - \mathbb{E}\hat{f})(\mathbb{E}\hat{f} - f)] \\ &= \mathbb{E}(\hat{f} - \mathbb{E}\hat{f})^2 + (\mathbb{E}\hat{f} - f)^2\end{aligned}$$

т.к. $(\mathbb{E}\hat{f} - f)$ - константа относительно X, Y ,

$$\mathbb{E}[(\hat{f} - \mathbb{E}\hat{f})(\mathbb{E}\hat{f} - f)] = (\mathbb{E}\hat{f} - f)\mathbb{E}(\hat{f} - \mathbb{E}\hat{f}) = 0.$$

$$\begin{aligned}\mathbb{E}(\hat{f} - y)^2 &= \mathbb{E}(\hat{f} - f - \varepsilon)^2 = \mathbb{E}(\hat{f} - f)^2 + \mathbb{E}\varepsilon^2 - 2\mathbb{E}[(\hat{f} - f)\varepsilon] \\ &= \mathbb{E}(\hat{f} - \mathbb{E}\hat{f})^2 + (\mathbb{E}\hat{f} - f)^2 + \mathbb{E}\varepsilon^2\end{aligned}$$

$$\mathbb{E}[(\hat{f} - f)\varepsilon] = \mathbb{E}[(\hat{f} - f)]\mathbb{E}\varepsilon = 0, \text{ поскольку } \varepsilon \text{ не зависит от } X, Y.$$

Содержание

- 1 Разложение на смещение и разброс
- 2 Композиции алгоритмов**
- 3 Примеры использования ансамблей
- 4 Ансамбли против недообучения
- 5 Ансамбли против переобучения

Композиции алгоритмов

- Композиция алгоритмов (ансамбль моделей, ensemble learning):

$$\hat{y}(x) = G(f_1(x), \dots, f_M(x))$$

- $f_1(x), \dots, f_M(x)$ - базовые модели=признаки для $G(\cdot)$
 - $G(\cdot)$ - агрегирующая модель, мета-модель (meta-model)
- Используется в
 - обучении с учителем (регрессия, классификация)
 - без учителя (кластеризация)

Мотивация композиций

Мотивация:

- борьба с переобучением $f_1(x), \dots, f_M(x)$: простая $G(\cdot)$
- борьба с недообучением $f_1(x), \dots, f_M(x)$: сложная $G(\cdot)$
- каждая $f_1(x), \dots, f_M(x)$ отвечает за свою область признакового пространства (mixture of experts)

$$G(x) = \sum_{m=1}^M w_m(x) f_m(x)$$

$$w_m = \begin{cases} 0, & m \neq i(x) \\ 1 & m = i(x) \end{cases} \quad \text{либо} \quad w = \text{SoftMax}(\dots)$$

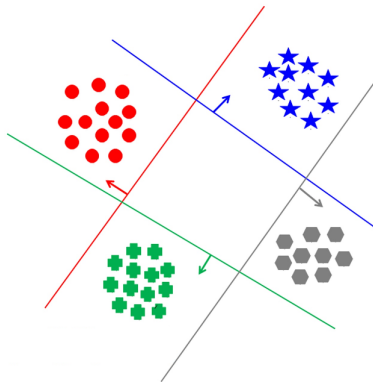
- построение $\hat{y}(x)$ декомпозируется на решение подзадач $f_1(x), \dots, f_M(x)$
- ускорение обучения
 - например, усреднение ядерных SVM на подвыборках

Содержание

- 1 Разложение на смещение и разброс
- 2 Композиции алгоритмов
- 3 Примеры использования ансамблей**
- 4 Ансамбли против недообучения
- 5 Ансамбли против переобучения

Многоклассовая классификация

Многоклассовая классификация бинарными классификаторами (один против всех, один против одного, коды, исправляющие ошибки):



Последовательное решение, признаки разной природы

Последовательное решение

- Разделим классы: 1,2,"3+4"
 - если "3+4", применим модель, разделяющую 3 от 4.
- Прогнозирование стоимости квартир:
 - определяем тип покупки: для жилья/для инвестиций
 - для жилья: комфорт, индивидуальные вкусы и т.д.
 - для инвестиций: обменные курсы, процент по вкладам, рост рынка акций и т.д.
- Определение людей по фото:
 - определяем ракурс: фас/профиль
 - одна модель определяет людей по фото в фас
 - другая определяет людей по фото в профиль

Признаки разной природы

- Идентификация человека по разнородной информации:
по голосу, по лицу, по поведению, и т.д.

Борьба с переобучением и недообучением

Ансамбли моделей используются для борьбы

- с переобучением базовых моделей
 - голосование по большинству, усреднение, бэггинг, метод случайных подпространств, случайный лес, особо случайные деревья.
- с недообучением базовых моделей
 - AND, OR, K-out-of-M, стэкинг, блендинг, бустинг.

Содержание

- 1 Разложение на смещение и разброс
- 2 Композиции алгоритмов
- 3 Примеры использования ансамблей
- 4 Ансамбли против недообучения**
- 5 Ансамбли против переобучения

Борьба с переобучением

- Предположим $f_1(x), \dots, f_M(x)$ - слишком простые модели.
- Можем повысить сложность, применяя сложную мета-модель:

$$\hat{y}(x) = G(f_1(x), \dots, f_M(x))$$

- Примеры:
 - $G(f_1(x), f_2(x)) = f_1(x) \text{ AND } f_2(x);$
 $G(f_1(x), f_2(x)) = f_1(x) \text{ OR } f_2(x);$

Борьба с переобучением

- Предположим $f_1(x), \dots, f_M(x)$ - слишком простые модели.
- Можем повысить сложность, применяя сложную мета-модель:

$$\hat{y}(x) = G(f_1(x), \dots, f_M(x))$$

- Примеры:
 - $G(f_1(x), f_2(x)) = f_1(x) \text{ AND } f_2(x);$
 $G(f_1(x), f_2(x)) = f_1(x) \text{ OR } f_2(x);$
 - Стэкинг, блендинг: $G(\cdot)$ настраиваемая функция.
 - $f_1(x), \dots, f_M(x)$ - сложные, $G(f_1(x), \dots, f_M(x))$ - ещё сложнее!

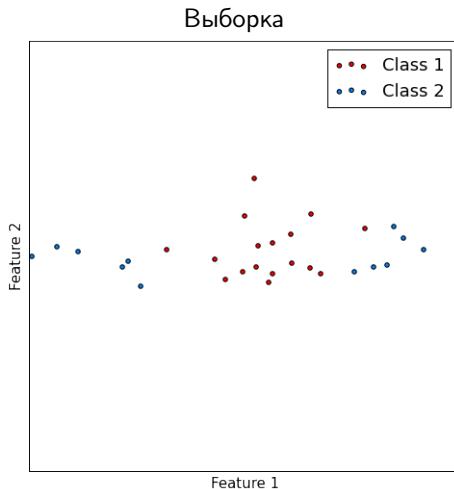
Борьба с переобучением

- Предположим $f_1(x), \dots, f_M(x)$ - слишком простые модели.
- Можем повысить сложность, применяя сложную мета-модель:

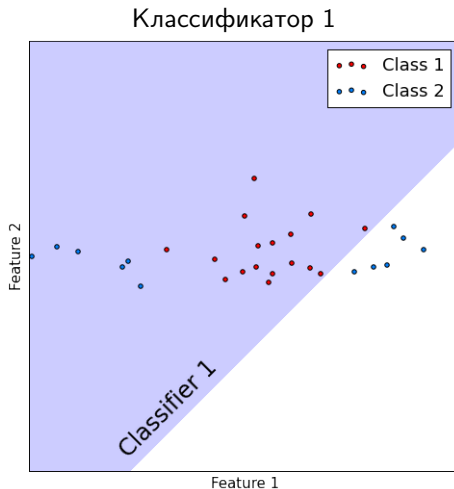
$$\hat{y}(x) = G(f_1(x), \dots, f_M(x))$$

- Примеры:
 - $G(f_1(x), f_2(x)) = f_1(x) \text{ AND } f_2(x)$;
 $G(f_1(x), f_2(x)) = f_1(x) \text{ OR } f_2(x)$;
 - Стэкинг, блендинг: $G(\cdot)$ настраиваемая функция.
 - $f_1(x), \dots, f_M(x)$ - сложные, $G(f_1(x), \dots, f_M(x))$ - ещё сложнее!
 - Бустинг: $G_i(x) = \varepsilon f_1(x) + \varepsilon f_2(x) + \dots + \varepsilon f_i(x)$, $\varepsilon \in (0, 1]$.
 - $f_1(x)$ учится предсказывать y
 - $f_2(x)$ исправляет ошибки $G_1(x)$
 - $f_3(x)$ исправляет ошибки $G_2(x)$
 - ...

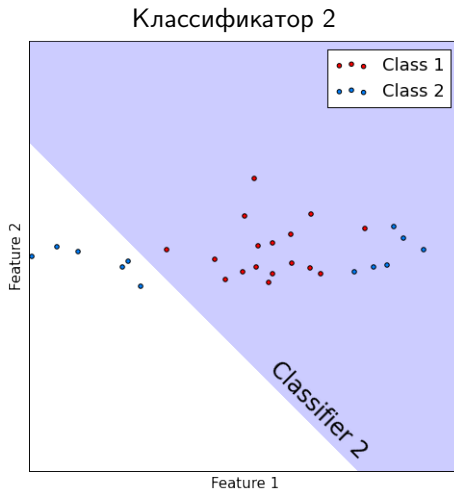
Идея композиции AND



Идея композиции AND

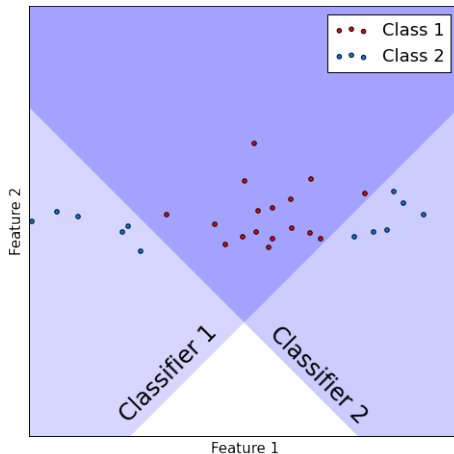


Идея композиции AND



Идея композиции AND

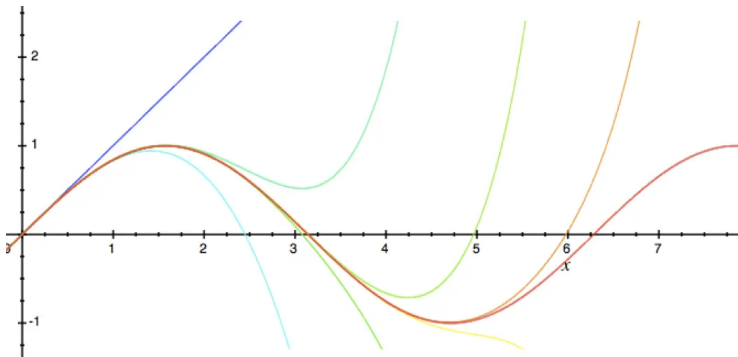
(Классификатор 1) AND (классификатор 2)



Идея стэкинга

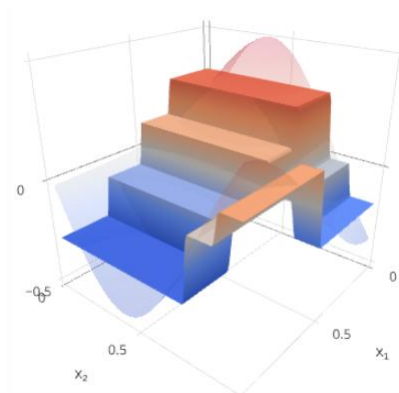
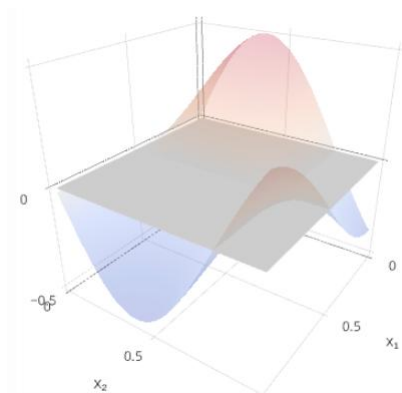
- $f_1(x) = x, f_2(x) = (x)^2, \dots f_M(x) = (x)^M$.
- $G(f_1(x), \dots f(x)_M) = w_0 + w_1 f_1(x) + \dots w_M f_M(x) = w_0 + w_1 x + w_2 (x)^2 + \dots w_M (x)^M$

Последовательное приближение $y = \sin(x)$:



Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

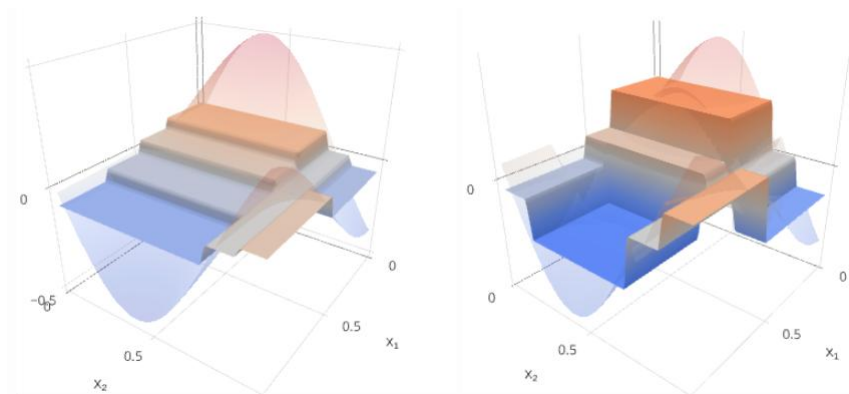
Слева: $y(x)$, $G_0(x)$. Справа: отклонение и $f_1(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

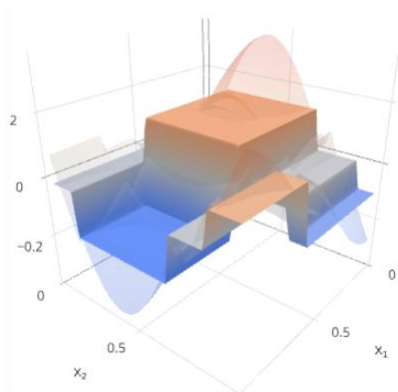
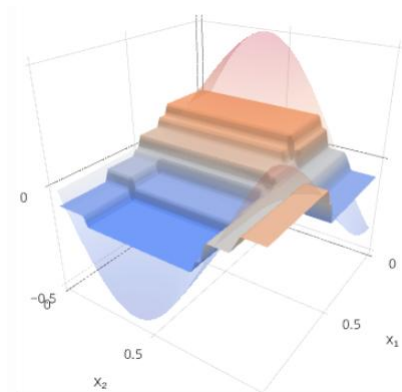
Слева: $y(x)$, $G_1(x)$. Справа: отклонение, $f_2(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

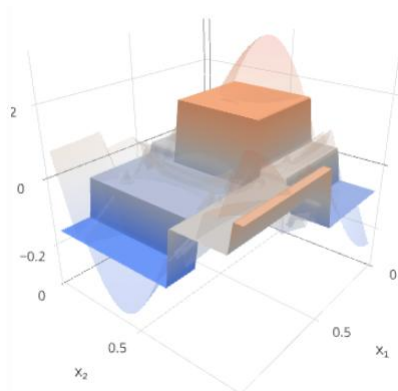
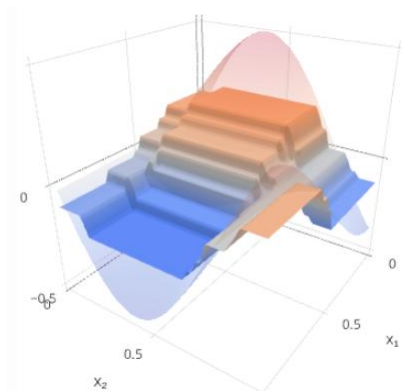
Слева: $y(x)$, $G_2(x)$. Справа: отклонение, $f_3(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

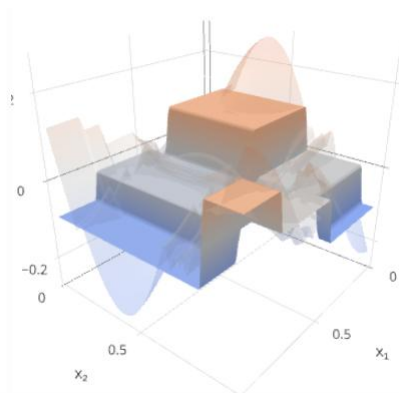
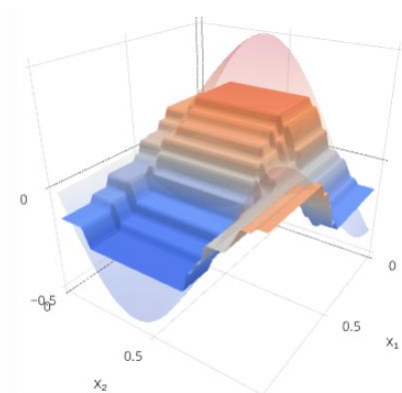
Слева: $y(x)$, $G_3(x)$. Справа: отклонение, $f_4(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

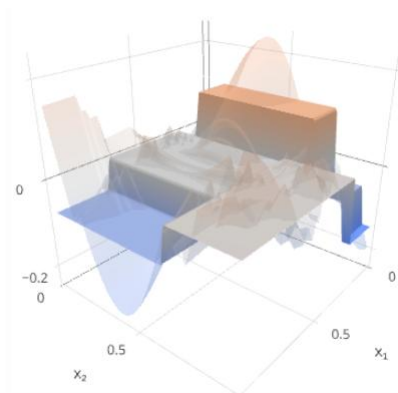
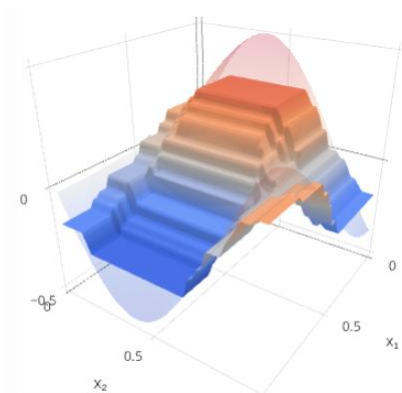
Слева: $y(x)$, $G_4(x)$. Справа: отклонение, $f_5(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

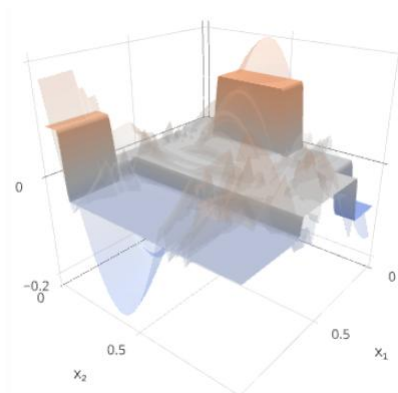
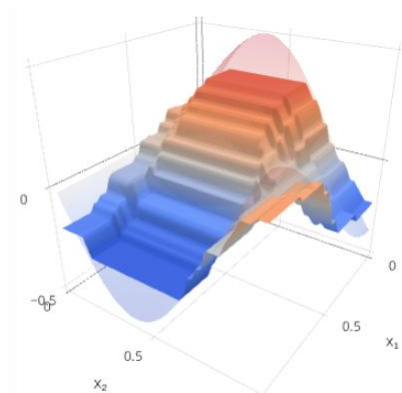
Слева: $y(x)$, $G_5(x)$. Справа: отклонение, $f_6(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

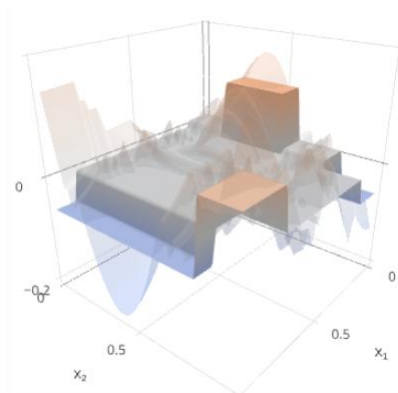
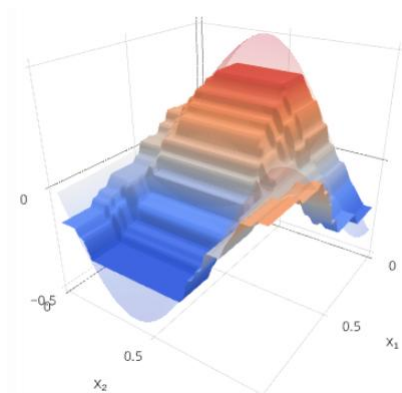
Слева: $y(x)$, $G_6(x)$. Справа: отклонение, $f_7(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

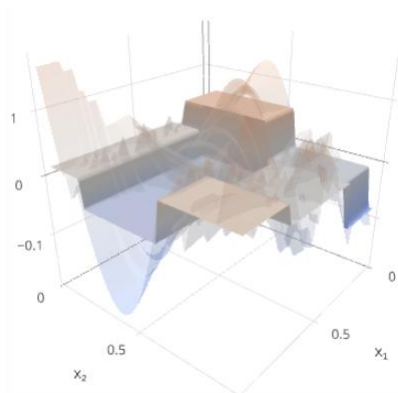
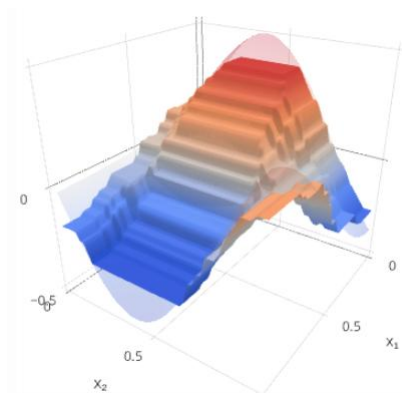
Слева: $y(x)$, $G_7(x)$. Справа: отклонение, $f_8(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

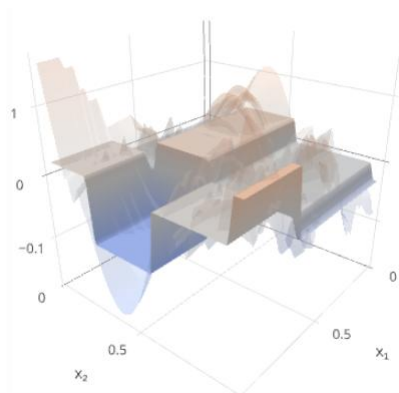
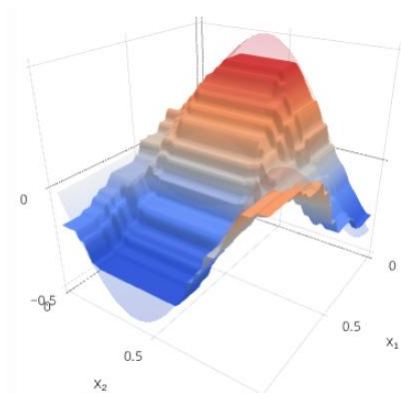
Слева: $y(x)$, $G_8(x)$. Справа: отклонение, $f_9(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

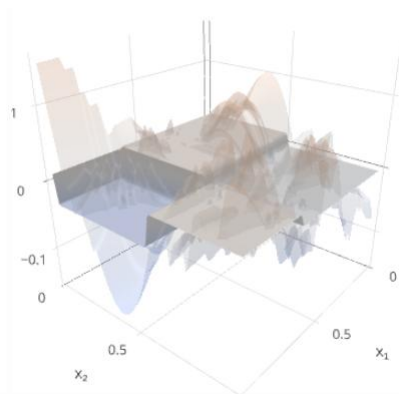
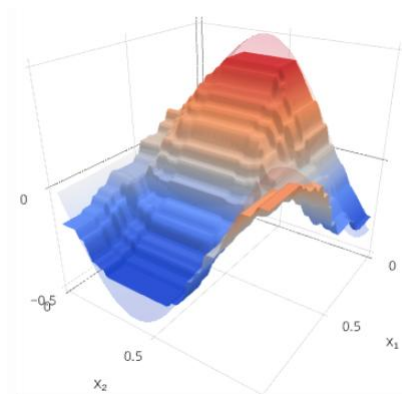
Слева: $y(x)$, $G_9(x)$. Справа: отклонение, $f_{10}(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Идея бустинга (на деревьях глубины 3, $\varepsilon = 0.3$)²

Слева: $y(x)$, $G_{10}(x)$. Справа: отклонение, $f_{11}(x)$



²Интерактивная иллюстрация Алексея Рогожникова.

Содержание

- 1 Разложение на смещение и разброс
- 2 Композиции алгоритмов
- 3 Примеры использования ансамблей
- 4 Ансамбли против недообучения
- 5 Ансамбли против переобучения**

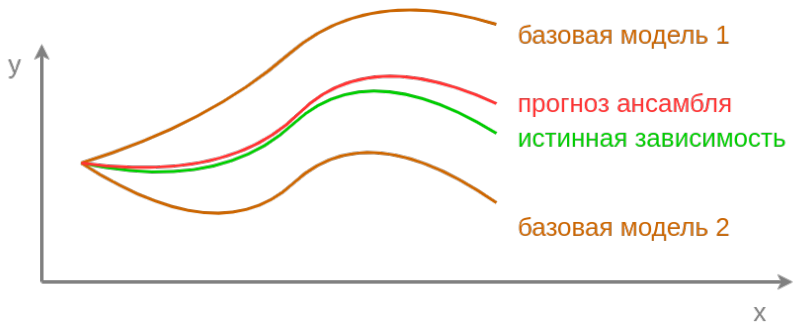
Борьба с переобучением

- $f_1(x), \dots, f_M(x)$ слишком сложные (переобученные модели)
 - решающие деревья большой глубины на разных подвыборках
 - глубокие нейросети
 - обученные из разных начальных приближений
 - разной архитектуры
- Регрессия: сделаем устойчивый прогноз за счет усреднения

$$G(x) = \frac{1}{M} \sum_{m=1}^M f_m(x)$$

- Классификация: прогноз самым частым классом среди $\{f_1(x), \dots, f_M(x)\}$.
 - голосование по большинству, majority voting

Регрессия: переобучение



Усреднение ошибок

Рассмотрим задачу регрессии с усредняющим ансамблем:

$$G(x) = \frac{1}{M} \sum_{m=1}^M f_m(x)$$

Пусть

$$\varepsilon_i = f_i(x) - y(x), \quad \mathbb{E}\varepsilon_i = 0, \quad \mathbb{D}\varepsilon_i = \mathbb{E} \left\{ (\varepsilon_i - \mathbb{E}\varepsilon_i)^2 \right\} = \mathbb{E}\varepsilon_i^2 = \sigma^2$$

$$\text{corr}\{\varepsilon_i, \varepsilon_j\} = \rho = \text{cov}\{\varepsilon_i, \varepsilon_j\} / \sqrt{\mathbb{D}\varepsilon_i \cdot \mathbb{D}\varepsilon_j}$$

$$\text{cov}\{\varepsilon_i, \varepsilon_j\} = \mathbb{E} \{ (\varepsilon_i - \mathbb{E}\varepsilon_i) (\varepsilon_j - \mathbb{E}\varepsilon_j) \} = \mathbb{E} \{ \varepsilon_i \varepsilon_j \} = \rho \sigma^2$$

$$\mathbb{E} \left\{ (f_i(x) - y(x))^2 \right\} = \mathbb{E}\varepsilon_i^2 = \sigma^2, \quad \mathbb{E} \left\{ (G(x) - y(x))^2 \right\} - ?$$

Усреднение ошибок

Ожидаемый квадрат ошибки усредняющего ансамбля:

Усреднение ошибок

Ожидаемый квадрат ошибки усредняющего ансамбля:

$$\begin{aligned}\mathbb{E} \left\{ (G(x) - y(x))^2 \right\} &= \mathbb{E} \left\{ \left(\frac{\sum_{m=1}^M (f_m(x) - y(x))}{M} \right)^2 \right\} \\&= \frac{1}{M^2} \mathbb{E} \left\{ \left(\sum_{m=1}^M \varepsilon_m \right)^2 \right\} = \frac{1}{M^2} \mathbb{E} \left\{ \sum_{m=1}^M \varepsilon_m^2 + \sum_{i \neq j} \varepsilon_i \varepsilon_j \right\} \\&= \frac{M}{M^2} \sigma^2 + \frac{M^2 - M}{M^2} \rho \sigma^2 = \frac{\sigma^2}{M} + \left(1 - \frac{1}{M} \right) \rho \sigma^2\end{aligned}$$

- При $\rho = 1$: σ^2 ; при $\rho = 0$: σ^2/M
 - может быть еще меньше при $\rho < 0$.

Усреднение ошибок

Ожидаемый квадрат ошибки усредняющего ансамбля:

$$\begin{aligned}\mathbb{E} \left\{ (G(x) - y(x))^2 \right\} &= \mathbb{E} \left\{ \left(\frac{\sum_{m=1}^M (f_m(x) - y(x))}{M} \right)^2 \right\} \\&= \frac{1}{M^2} \mathbb{E} \left\{ \left(\sum_{m=1}^M \varepsilon_m \right)^2 \right\} = \frac{1}{M^2} \mathbb{E} \left\{ \sum_{m=1}^M \varepsilon_m^2 + \sum_{i \neq j} \varepsilon_i \varepsilon_j \right\} \\&= \frac{M}{M^2} \sigma^2 + \frac{M^2 - M}{M^2} \rho \sigma^2 = \frac{\sigma^2}{M} + \left(1 - \frac{1}{M} \right) \rho \sigma^2\end{aligned}$$

- При $\rho = 1$: σ^2 ; при $\rho = 0$: σ^2/M
 - может быть еще меньше при $\rho < 0$.
- На практике $\rho > 0$, т.к. модели прогнозируют одинаковый отклик по одинаковой выборке.

Голосование большинства (против переобучения)

- Рассмотрим M классификаторов $f_1(x), \dots, f_M(x)$.
- Пусть $p(f_m(x) \neq y) = p < 0.5 \forall m$
 - например, $p = 0.49$.
- Пусть модели ошибаются независимо друг от друга.
- Пусть $G(x)$ - выбор самого частого класса.
- Тогда $p(G(x) \neq y) \rightarrow 0$ при $M \rightarrow \infty$

Доказательство

Доказательство

Рассмотрим сл. вел.: $\xi_m = \begin{cases} +1, & f_m(x) = y \\ 0 & f_m(x) \neq y \end{cases} \quad \eta = \frac{\xi_1 + \dots + \xi_M}{M}$

$$G(x) = y \Leftrightarrow \eta > 0.5; \quad P(G(x) = y) = P(\eta > 0.5)$$

$$\frac{\sum_{i=1}^M [\xi_i - \mathbb{E}\xi_i]}{\sqrt{M\mathbb{D}(\xi_1)}} = \frac{\sum_{i=1}^M [\xi_i - p]}{\sqrt{Mp(1-p)}} = \frac{\sqrt{M}}{\sqrt{p(1-p)}} \frac{\sum_{i=1}^M \xi_i - Mp}{M}$$

$$= \frac{\sqrt{M}}{\sqrt{p(1-p)}} (\eta - p) \rightarrow \mathcal{N}(0, 1) \text{ по ЦПТ.}$$

$$\eta \rightarrow \mathcal{N}\left(p, \frac{p(1-p)}{M}\right) \text{ при } M \rightarrow \infty$$

$$\mathbb{E}\eta = p \text{ и } \mathbb{D}\eta = \frac{p(1-p)}{M} \rightarrow 0 \text{ при } M \rightarrow \infty \Rightarrow$$

$$P(\eta > 0.5) \rightarrow 1 \text{ при } M \rightarrow \infty \text{ по неравенству Чебышева.}$$

Взвешенное усреднение (против переобучения)

- **Разложение неоднозначности (ambiguity decomposition):**

пусть (x, y) прогнозируется с помощью регрессии

$G(x) = \sum_{m=1}^M w_m f_m(x)$, $w_m \geq 0$, $\sum_m w_m = 1$. Тогда

$$\underbrace{(G(x) - y)^2}_{\text{ошибка ансамбля}} = \underbrace{\sum_m w_m (f_m(x) - y)^2}_{\text{ошибки базовых моделей}} - \underbrace{\sum_m w_m (f_m(x) - G(x))^2}_{\text{неоднозначность}}$$

- Композиция дает точные прогнозы когда:
 - $f_m(x)$ достаточно точны
 - индивидуальные прогнозы $\{f_m(x)\}_m$ сильно различаются
 - поэтому полезно усреднять по разным моделям

Доказательство разложения неоднозначности

Доказательство разложения неоднозначности

Доказательство:

$$\begin{aligned} & \sum_m w_m (f_m(x) - G(x))^2 = \sum_m w_m (f_m(x) - y + y - G(x))^2 \\ &= \sum_m w_m (f_m(x) - y)^2 + \sum_m w_m (y - G(x))^2 + 2 \sum_m w_m (f_m(x) - y) (y - G(x)) \\ &= \sum_m w_m (f_m(x) - y)^2 + (G(x) - y)^2 + 2 (y - G(x)) \sum_m w_m (f_m(x) - y) \\ &= \sum_m w_m (f_m(x) - y)^2 + (G(x) - y)^2 + 2 (y - G(x)) (G(x) - y) \\ &= \sum_m w_m (f_m(x) - y)^2 + (G(x) - y)^2 - 2 (G(x) - y)^2 \\ &= \sum_m w_m (f_m(x) - y)^2 - (G(x) - y)^2 \end{aligned}$$

Выпуклые потери

Выпуклые потери поощряют использование взвешенных прогнозов вместо индивидуальных.

- Рассмотрим регрессию с выпуклой ф-цией потерь $\mathcal{L}(\hat{y} - y)$.
- Учитываем $f_1(x), \dots, f_M(x)$ с весами $w_1, \dots, w_M \geq 0$,
 $\sum_i w_i = 1$.

Для фикс. x какая стратегия выгоднее?

- 1 усреднять $\hat{y}(x) = \sum_{m=1}^M w_m f_m(x)$
- 2 сэмплировать $m \sim \text{Categorical}(w_1, \dots, w_M)$, $\hat{y}(x) = f_m(x)$.

Выпуклые потери

Выпуклые потери поощряют использование взвешенных прогнозов вместо индивидуальных.

- Рассмотрим регрессию с выпуклой ф-цией потерь $\mathcal{L}(\hat{y} - y)$.
- Учитываем $f_1(x), \dots, f_M(x)$ с весами $w_1, \dots, w_M \geq 0$,
 $\sum_i w_i = 1$.

Для фикс. x какая стратегия выгоднее?

- 1 усреднять $\hat{y}(x) = \sum_{m=1}^M w_m f_m(x)$
- 2 сэмплировать $m \sim \text{Categorical}(w_1, \dots, w_M)$, $\hat{y}(x) = f_m(x)$.

выпуклость функции: $\mathcal{L}(\alpha \varepsilon_1 + (1 - \alpha) \varepsilon_2) \leq \alpha \mathcal{L}(\varepsilon_1) + (1 - \alpha) \mathcal{L}(\varepsilon_2)$

по индукции можно показать, что:

$$\mathcal{L}(w_1 \varepsilon_1 + w_2 \varepsilon_2 + \dots + w_M \varepsilon_M) \leq w_1 \mathcal{L}(\varepsilon_1) + w_2 \mathcal{L}(\varepsilon_2) + \dots + w_M \mathcal{L}(\varepsilon_M)$$