

# Наилучшая модель регрессии. Разложение ошибки на шум, смещение и разброс

Автор материалов: Соколов Евгений Андреевич  
Документ составил: Богачев Владимир Александрович

5 января 2024 г.

# 1 Введение

Рассмотрим задачу регрессии: пусть признаковое описание  $\mathcal{X}$  и целевая переменная  $\mathcal{Y}$  имеют совместное распределение  $(\mathcal{X}, \mathcal{Y}) \sim P_{\mathcal{X}, \mathcal{Y}}$ . Пусть  $(X, Y)$  — выборка мощности  $l$  из совместного распределения  $P_{\mathcal{X}, \mathcal{Y}}$ . Необходимо построить функцию  $f(x) : \text{supp}(\mathcal{X}) \rightarrow \text{supp}(\mathcal{Y})$  из некоторого класса  $\mathcal{F}$ , оптимизирующую функционал (1)

$$\mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathcal{Y})^2 \xrightarrow{f \in \mathcal{F}} \min \quad (1)$$

На практике совместное распределение  $P_{\mathcal{X}, \mathcal{Y}}$  неизвестно, известна лишь реализация выборки  $(x_i, y_i)_{i=1}^l$  из данного распределения. Поэтому вычислить точное значение функционала (1) не представляется возможным. На практике задача оптимизации функционала (1) заменяется приближенной задачей оптимизации эмпирического риска (2)

$$\int_{\mathbb{R}^2} (f(x) - y)^2 \hat{F}_l(dx, dy) = \sum_{i=1}^l (f(x_i) - y_i)^2 \xrightarrow{f \in \mathcal{F}} \min \quad (2)$$

Необходимо построить стратегию  $\mu(x_1, \dots, x_l, y_1, \dots, y_l) : \text{supp } X \times \text{supp } Y \rightarrow \mathcal{F}$ , ставящую в соответствие каждой реализации выборки  $(x_i, y_i)_{i=1}^l$  решение задачи оптимизации (2).

Решение приближенной задачи зависит не только от совместного распределения  $P_{\mathcal{X}, \mathcal{Y}}$ , но и от реализации выборки  $(x_i, y_i)_{i=1}^l$ . Поэтому значение эмпирического риска, достигнутое на данном решении, не позволяет оценить качество стратегии  $\mu$  на всей генеральной совокупности.

Для оценки качества стратегии  $\mu$  на всей генеральной совокупности возможно использовать усредненный по выборкам среднеквадратичный риск (3).

$$L(\mu) := \mathbb{E}_{X, Y} \left[ \mathbb{E}_{x, y} (y - \mu(X, Y)(x))^2 \right] \quad (3)$$

В разделе 2 будет построено решение задачи (1). В разделе 3 будет представлено разложение усредненного среднеквадратичного риска  $L(\mu)$ , позволяющее сравнить качество решения приближенной задачи (2) с решением (1).

## 2 Оптимальная модель регрессии

В данном разделе будет построено решение задачи (1). По заданной совместной плотности  $P_{\mathcal{X}, \mathcal{Y}}$  необходимо выбрать функцию  $f \in \mathcal{F}$ , оптимизирующую функционал  $\mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathcal{Y})^2 \xrightarrow{f \in \mathcal{F}} \min$ .

Далее будем рассматривать класс интегрируемых с квадратом функций  $\mathcal{F} = L_2$ .

$$\begin{aligned} \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathcal{Y})^2 &= \\ &= \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}) + \mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})^2 = \\ &= \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))^2 + \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})^2 + \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y}) \end{aligned}$$

Рассмотрим

$$\begin{aligned} \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y}) &= \\ &= \mathbb{E}_{\mathcal{X}} [\mathbb{E}_{\mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})|\mathcal{X}] = \\ &= \mathbb{E}_{\mathcal{X}} (\mathbb{E}_{\mathcal{Y}} [(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})|\mathcal{X}] \cdot (f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))) = \\ &= \mathbb{E}_{\mathcal{X}} ([\mathbb{E}\mathcal{Y} - \mathbb{E}\mathcal{Y}] \cdot (f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))) = \\ &= 0 \end{aligned}$$

Следовательно, исходный функционал можно представить в виде суммы неотрицательных слагаемых:

$$\mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathcal{Y})^2 = \mathbb{E}_{\mathcal{X}, \mathcal{Y}}(f(\mathcal{X}) - \mathbb{E}(\mathcal{Y}|\mathcal{X}))^2 + 2\mathbb{E}_{\mathcal{X}, \mathcal{Y}}(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})^2 \quad (4)$$

Функционал (4) достигает минимума при выборе  $f(x) = \mathbb{E}(\mathcal{Y}|\mathcal{X} = x)$ .

Наименьшее значение ошибки:  $\mathbb{E}_{\mathcal{X},\mathcal{Y}}(\mathbb{E}(\mathcal{Y}|\mathcal{X}) - \mathcal{Y})^2$

Как было сказано ранее, на практике совместное распределение  $P_{\mathcal{X},\mathcal{Y}}$  неизвестно, известна реализация выборки  $(x_i, y_i)_{i=1}^l$  из данного распределения, поэтому построить решение  $f(X) = \mathbb{E}(Y|X)$  не представляется возможным. Можно построить решение, основанное на оценке совместного распределения  $P_{\mathcal{X},\mathcal{Y}}$  по реализации выборки (например, KDE). Однако данные методы имеют высокую вычислительную сложность, а построение оценки совместного распределения может оказаться избыточным. Поэтому встает вопрос о сравнении произвольных моделей машинного обучения с наилучшим решением  $f(X) = \mathbb{E}(Y|X)$ .

### 3 Bias-Variance decomposition

В разделе 1 была поставлена приближенная задача (2). Было введено понятие стратегии  $\mu : \text{supp}(X) \times \text{supp}(Y) \rightarrow \mathcal{F}$ . Также был введен усредненный по выборкам среднеквадратичный риск  $L(\mu) := \mathbb{E}_{X,Y} \left[ \mathbb{E}_{x,y} (y - \mu(X, Y)(x))^2 \right]$ .

Далее будет построено разложение  $L(\mu)$  на «шум», «смещение» и «разброс»

$$\begin{aligned} L(\mu) &= \mathbb{E}_{X,Y} \left[ \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y|x])^2 \right]}_{\text{не зависит от } X} + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y|x] - \mu(X, Y))^2 \right] \right] = \\ &= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y|x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mu(X, Y))^2 \right] \right]. \end{aligned}$$

Преобразуем второе слагаемое:

$$\begin{aligned} \mathbb{E}_{x,y} \left[ \mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mu(X, Y))^2 \right] \right] &= \\ &= \mathbb{E}_{x,y} \left[ \mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)] + \mathbb{E}_{X,Y} [\mu(X, Y)] - \mu(X, Y))^2 \right] \right] = \\ &= \mathbb{E}_{x,y} \left[ \underbrace{\mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)])^2 \right]}_{\text{не зависит от } X} + \mathbb{E}_{x,y} \left[ \mathbb{E}_{X,Y} \left[ (\mathbb{E}_{X,Y} [\mu(X, Y)] - \mu(X, Y))^2 \right] \right] + \right. \\ &\quad \left. + 2\mathbb{E}_{x,y} \left[ \mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)]) (\mathbb{E}_{X,Y} [\mu(X, Y)] - \mu(X, Y)) \right] \right] \right]. \end{aligned}$$

Покажем, что последнее слагаемое обращается в нуль:

$$\begin{aligned} \mathbb{E}_{X,Y} \left[ (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)]) (\mathbb{E}_{X,Y} [\mu(X, Y)] - \mu(X, Y)) \right] &= \\ &= (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)]) \mathbb{E}_{X,Y} \left[ \mathbb{E}_{X,Y} [\mu(X, Y)] - \mu(X, Y) \right] = \\ &= (\mathbb{E}[y|x] - \mathbb{E}_{X,Y} [\mu(X, Y)]) \left[ \mathbb{E}_{X,Y} [\mu(X, Y)] - \mathbb{E}_{X,Y} [\mu(X, Y)] \right] = \\ &= 0. \end{aligned}$$

Тогда  $L(\mu)$  можно представить в виде (5)

$$L(\mu) = \underbrace{\mathbb{E}_{x,y} [y - \mathbb{E}(y|x)]^2}_{\text{шум}} + \underbrace{\mathbb{E}_x [\mathbb{E}_{X,Y} \mu(X, Y) - \mathbb{E}(y|x)]^2}_{\text{смещение (bias)}} + \underbrace{\mathbb{E}_x \mathbb{E}_{X,Y} [\mu(X, Y) - \mathbb{E}_{X,Y} \mu(X, Y)]^2}_{\text{разброс(variance)}} \quad (5)$$

Из вида (5) следует, что  $L(\mu)$  складывается из ошибки наилучшего решения (шума) и неотрицательных компонент **bias** и **variance**

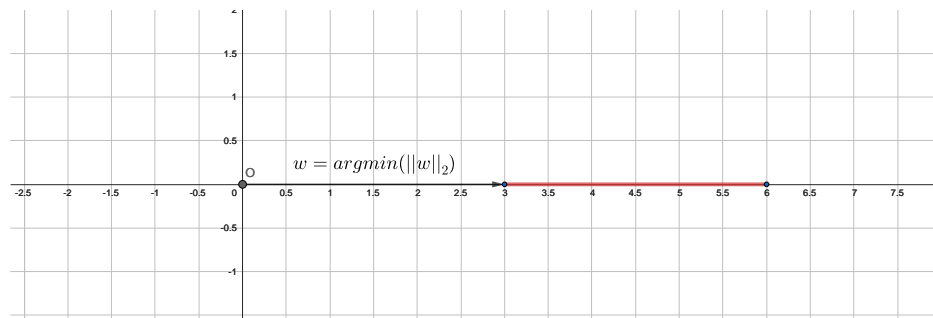
### 4 Bias-Variance tradeoff

В классической теории **bias** убывает с ростом числа параметров модели, **variance**, наоборот, возрастает при увеличении сложности модели. Однако с развитием глубокого обучения стали появляться модели с миллиардами параметров (например, LLaMA2 имеет 70 миллиардов параметров). Существование таких моделей не укладывается в классическую теорию. Из-за этого была разработана новая теория ([Bel+19]).

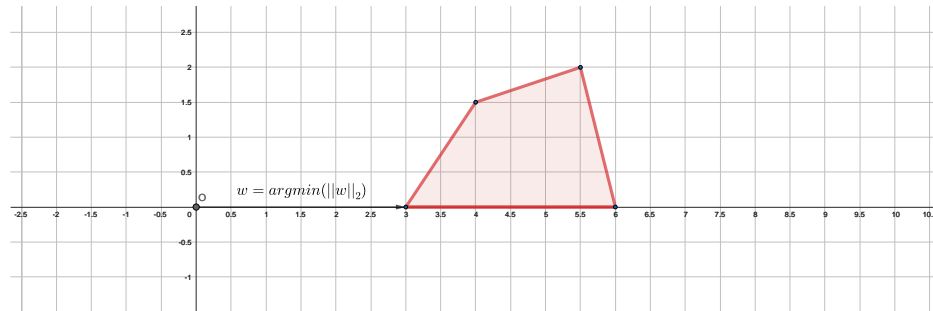
**under-fitting** При малом количестве параметров увеличение сложности приводит к улучшению качества, так как модель начинает выучивать новые зависимости, присутствующие на всей генеральной совокупности (а следовательно, на **train** и на **test**). Происходит значительное уменьшение **bias** при сравнительно небольшом росте **variance**. Это соответствует участку **under-fitted** на графике 2. Данный участок одинаково описывается и «классической» и «новой» теориями.

**over-fitting** При большом количестве параметров увеличение сложности приводит к ухудшению качества, так как модель начинает выучивать закономерности, присущие исключительно обучающей выборке. Происходит значительное увеличение **variance** при малом снижении **bias** (так как усреднение происходит не только по обучающей выборке). Это соответствует участку **over-fitted** на графике 2. В «классической» теории эмпирический риск монотонно возрастает, начиная с некоторого значения сложности. В «новой» теории данный участок ограничен.

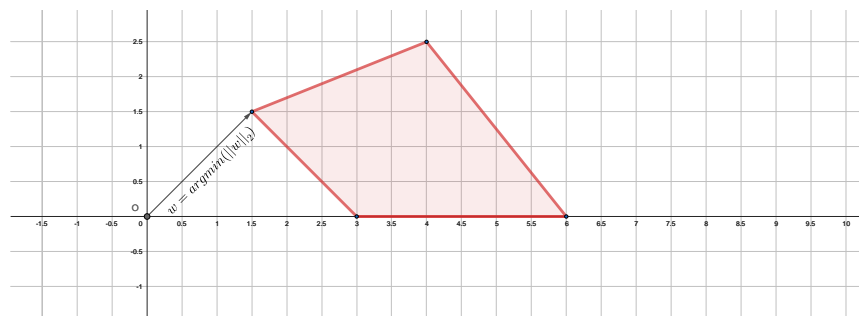
**over-parameterized** При очень большом количестве параметров появляется множество способов описания обучающей выборки. Регуляризация позволяет выбрать решение с минимальной нормой, следовательно, нет переобучения. (см. рисунок 1) Данный участок описывается только «новой» теорией. Начиная с некоторого значения сложности, эмпирический риск на тестовой выборке начинает улучшаться при неизменно низком эмпирическом риске на тренировочной выборке. На данном этапе основной вклад в эмпирический риск вносит шум.



(а) Одномерный случай



(б) Регуляризация не дает увеличению сложности ухудшить качество



(в) Регуляризация позволяет увеличению сложности улучшить качество

Рис. 1: Влияние регуляризации на качество модели при увеличении сложности модели

На рисунке 1 показано влияние регуляризации на модель при увеличении количества параметров. Если увеличение сложности не позволяет улучшить качество модели, то, благодаря регуляризации, будет выбран набор параметров, не учитывающий добавленную степень свободы (рис. 1b).

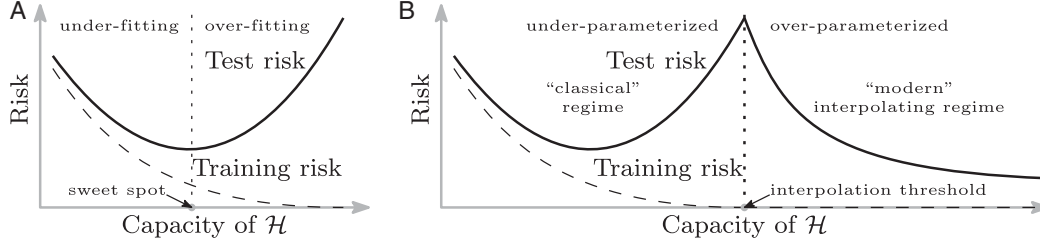


Рис. 2: Bias-variance tradeoff

Для демонстрации данного явления рассмотрим график, опубликованный в [Bel+19]. На рисунке 2А показана зависимость эмпирического риска от сложности модели (*capacity of  $\mathcal{H}$* ), предсказываемая классической теорией. На рисунке 2В показана аналогичная зависимость, предсказываемая «новой» теорией.

## 5 Задачи

Будем считать, что объекты  $x$  генерируются из нормального распределения  $x \sim p(x) = \mathcal{N}(0, \sigma_1^2)$ . Правильный ответ на объекте  $x$  определяется зашумленной функцией  $f(x): y = f(x) + \varepsilon$ ,  $\varepsilon \sim p(\varepsilon) = \mathcal{N}(0, \sigma_2^2)$ . Иными словами,  $y \sim p(y|x) = \mathcal{N}(f(x), \sigma_2^2)$ . Выборка  $X$  составляется из  $\ell$  независимых пар  $(x_i, y_i)$ . Мы будем рассматривать два простых частных случая:  $f(x) = ax$ , когда модель зависимости отвечает искомой зависимости, и  $f(x)$  — четная функция, т. е.  $f(-x) = f(x)$ .

**Задача 1.** Найти шумовую компоненту для одномерной линейной регрессии

*Решение.* Так как распределение  $p(y|x)$  нормальное, для него легко вычислить мат. ожидание:

$$\mathbb{E}[y|x] = f(x).$$

Тогда

$$\mathbb{E}_{x,y}[(y - \mathbb{E}[y|x])^2] = \mathbb{E}_{x,\varepsilon}[(f(x) + \varepsilon - f(x))^2] = \mathbb{E}_{\varepsilon}\varepsilon^2 = \mathbb{D}\varepsilon + (\mathbb{E}\varepsilon)^2 = \sigma_2^2 + 0 = \sigma_2^2.$$

□

**Задача 2.** Найти смещение алгоритма одномерной линейной регрессии для  $f(x) = ax$  и для произвольной четной  $f(x)$ .

*Решение.* Для начала найдем «средний» по всем выборкам ответ алгоритма на объекте  $x$ :

$$\mathbb{E}_X[\mu(X)(x)] = \mathbb{E}_X[k(X)]x.$$

Итак, нам нужно найти «среднее» по всем выборкам значение коэффициента  $k$ :

$$\mathbb{E}_X[k(X)] = \int \frac{\sum_i x_i (f(x_i) + \varepsilon_i)}{\sum_i x_i^2} \prod_i (p(x_i)p(\varepsilon_i)) dx_1 \dots dx_\ell d\varepsilon_1 \dots d\varepsilon_\ell. \quad (6)$$

Здесь записан несобственный интеграл, в котором каждая переменная принимает значения от  $-\infty$  до  $\infty$ . Значение этого интеграла определяется функцией  $f(x)$  и не всегда вычисляется аналитически. В случае, когда истинная зависимость в данных линейная, мы получим:

$$\begin{aligned} \mathbb{E}_X[k(X)] &= \int \frac{\sum_i x_i (a x_i + \varepsilon_i)}{\sum_i x_i^2} p(\bar{x})p(\bar{\varepsilon}) d\bar{x}d\bar{\varepsilon} = \\ &= a \int \frac{\sum_i x_i^2}{\sum_i x_i^2} p(\bar{x})p(\bar{\varepsilon}) d\bar{x}d\bar{\varepsilon} + \int \frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2} p(\bar{x})p(\bar{\varepsilon}) d\bar{x}d\bar{\varepsilon}. \end{aligned}$$

Мы сократили обозначение для дифференциалов и для плотностей распределений. Первый интеграл равен  $a$  (интеграл по всему пространству от плотности распределения). Второй интеграл берется по симметричным относительно нуля интервалам от нечетной по  $x_i$  и по  $\varepsilon_i$  функции, а значит, равен 0. Итак, «средний» коэффициент равен  $a$ .

Найдем смещение:

$$\mathbb{E}_x[(\mathbb{E}_X[\mu(X)(x)] - \mathbb{E}[y|x])^2] = \mathbb{E}_x[(ax - ax)^2] = 0.$$

Это интуитивно понятный результат: логично, что перебрав все возможные выборки длины  $\ell$  и усреднив по ним значение  $k$ , мы обязательно найдем истинную величину коэффициента.

Теперь найдем «среднее»  $k$  для произвольной четной  $f(x)$ . По аналогии с предыдущим случаем:

$$\mathbb{E}_X[k(X)] = \int \frac{\sum_i (x_i f(x_i))}{\sum_i x_i^2} p(\bar{x}) d\bar{x} + \int \frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2} p(\bar{x}) p(\bar{\varepsilon}) d\bar{x} d\bar{\varepsilon}.$$

Как мы уже выяснили, второй интеграл равен 0. Первый интеграл тоже равен 0, так как подынтегральное выражение — это нечетная по всем  $x_i$  функция. Итак, «среднее» значение коэффициента равно нулю. И это тоже понятный результат, потому что четную функцию логично приближать четной, а единственная четная линейная функция, проходящая через 0 — это  $y = 0$ .

Найдем смещение:

$$\mathbb{E}_x[(\mathbb{E}_X[\mu(X)(x)] - \mathbb{E}[y|x])^2] = \mathbb{E}_x[(0 - f(x))^2] = \mathbb{E}_x f^2(x).$$

Чем меньше четная функция  $y = f(x)$  похожа на четную линейную  $y = 0$ , тем больше будет смещение. Таким образом, если мы пытаемся приблизить нелинейную функцию  $f(x)$  в классе линейных, мы получаем большое смещение. Обратите внимание, что если бы  $f(x)$  не была четной, мы бы не смогли просто аналитически вычислить интегралы.  $\square$

**Задача 3.** Найти разброс алгоритма одномерной линейной регрессии для  $f(x) = ax$  и для произвольной четной  $f(x)$ .

*Решение.* Для  $f(x) = ax$ :

$$\begin{aligned} \mathbb{E}_x[\mathbb{E}_X[(\mu(X)(x) - \mathbb{E}_X[\mu(X)(x)])^2]] &= \mathbb{E}_x[\mathbb{E}_X[(ax + \frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2} x - ax)^2]] = \\ &= \left(\mathbb{E}_x x^2\right) \left(\mathbb{E}_X\left(\frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2}\right)^2\right) = \sigma_1^2 \mathbb{E}_X\left(\frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2}\right)^2. \end{aligned}$$

Мат. ожидание можно немного упростить, раскрыв квадрат суммы в числителе и внося внутрь суммы мат. ожидание по  $\bar{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_\ell)$ :

$$\mathbb{E}_X\left(\frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2}\right)^2 = \mathbb{E}_{\bar{x}}\left[\frac{\sum_{i \neq j} x_i x_j \mathbb{E}_{\bar{\varepsilon}}[\varepsilon_i \varepsilon_j] + \sum_i x_i^2 \mathbb{E}_{\bar{\varepsilon}}[\varepsilon_i^2]}{(\sum_i x_i^2)^2}\right].$$

Так как  $\varepsilon_i$  и  $\varepsilon_j$  независимы,  $\mathbb{E}[\varepsilon_i \varepsilon_j] = 0$ , а  $\mathbb{E}_{\bar{\varepsilon}}[\varepsilon_i^2] = \sigma_2^2$ . Тогда

$$\mathbb{E}_X\left(\frac{\sum_i x_i \varepsilon_i}{\sum_i x_i^2}\right)^2 = \mathbb{E}_{\bar{x}}\left[\frac{\sum_i x_i^2 \sigma_2^2}{(\sum_i x_i^2)^2}\right] = \sigma_2^2 \mathbb{E}_{\bar{x}}\left[\frac{1}{\sum_i x_i^2}\right],$$

а разброс

$$\mathbb{E}_x[\mathbb{E}_X[(\mu(X)(x) - \mathbb{E}_X[\mu(X)(x)])^2]] = \sigma_1^2 \sigma_2^2 \mathbb{E}_{\bar{x}}\left[\frac{1}{\sum_i x_i^2}\right].$$

Последнее мат. ожидание также можно рассчитать, но мы не будем этого делать. Мы получили, что если шум в ответах небольшой, то и разброс модели будет небольшой.

Для четной  $f(x)$ :

$$\begin{aligned} \mathbb{E}_x[\mathbb{E}_X[(\mu(X)(x) - \mathbb{E}_X[\mu(X)(x)])^2]] &= \mathbb{E}_x[\mathbb{E}_X[(0 - \frac{\sum_i x_i (f(x_i) + \varepsilon_i)}{\sum_i x_i^2} x)^2]] = \\ &= \left(\mathbb{E}_x x^2\right) \left(\mathbb{E}_X\left(\frac{\sum_i x_i (f(x_i) + \varepsilon_i)}{\sum_i x_i^2}\right)^2\right) = \sigma_1^2 \mathbb{E}_X\left[\frac{\sum_i x_i (f(x_i) + \varepsilon_i)}{\sum_i x_i^2}\right]^2. \end{aligned}$$

$\square$

## Список литературы

- [Bel+19] Mikhail Belkin и др. “Reconciling modern machine-learning practice and the classical bias–variance trade-off”. В: *Proceedings of the National Academy of Sciences* 116.32 (2019), с. 15849–15854.