

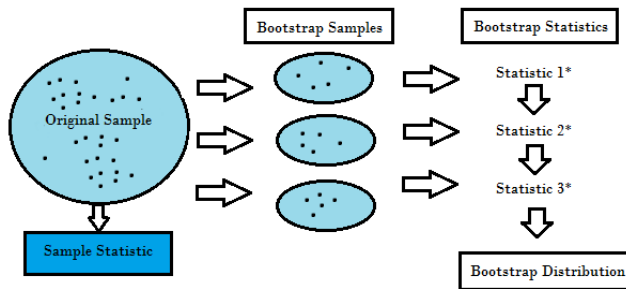
23. Разложение на смещение и разброс

Kazakbaev Rustem

December 2020

1 Идея

Есть обучающая выборка X и есть такая процедура: bootstrap. Из исходной обучающей выборки делаю подвыборку размера l с возвращением.



Делаем n выборок с помощью бутстрэпа: X_1, X_2, \dots, X_n
Сделаю fit модели и получу базовые алгоритмы: b_1, b_2, \dots, b_n
Предположу, что знаю распределение на всех возможных объектах: $p(X)$ распределение на X .

$y(X)$ - правильный ответ на объекте X

Посчитаю мат.ожидание ошибки j -ой модели: $E_x(b_j(X) - y(X)) = E_x(\epsilon_j^2(X))$. Но это ошибка j -ой модели, теперь попробую найти ошибку на всей выборке: $\frac{1}{N} \cdot E_x(\epsilon_j^2(X))$

1. В среднем у модели ошибка 0: $E_x(\epsilon_j(X)) = 0$

2. Независимость между i и j ошибкой.

Построим теперь новую функцию регрессии, которая будет усреднять ответы построенных нами функций:

$$a(x) = \frac{1}{N} \sum_{j=1}^N b_j(x).$$

Найдем ее среднеквадратичную ошибку:

$$E_x = \left(\frac{1}{N} \sum b_j(x) - y(x) \right)^2 = E_x \left(\frac{1}{N} \epsilon_j \right)^2 = \frac{1}{N^2} \cdot E_x \left(\sum \epsilon_j^2 + \underbrace{\sum \epsilon_j \cdot \epsilon_i}_{=0} \right) = \frac{1}{N^2} \cdot E_x \left(\sum \epsilon_j^2 \right)$$
 Будем

считать, что все ошибки одинаково распределены, тогда: $\frac{1}{N} \cdot E_X(\epsilon_1^2)$.

Таким образом, усреднение ответов позволило уменьшить средний квадрат ошибки в N раз!

Следует отметить, что рассмотренный нами пример не очень применим на практике, поскольку мы сделали предположение о некоррелированности ошибок, что редко выполняется, хотя это еще может быть возможно.

Предположил, что некоррелировано, так как это неправда - делал выборки через бутстрэп.

Если это предположение неверно, то уменьшение ошибки оказывается не таким значительным.

Важный вывод из этой идеи: можно объединять алгоритмы для достижения наилучшего результата!

2 Разложение ошибки на смещение и разброс

Дает понимание, какие модели необходимо брать, чтобы сделать правильные композиции.

Ошибка любой модели складывается из трех факторов:

1. Сложность самой выборки
2. Сходства модели с истинной зависимостью ответов от объектов в выборке
3. Богатство семейства, из которого выбирается конкретная модель

Между этими факторами существует некоторый баланс, и уменьшение одного из них приводит к увеличению другого. Такое разложение ошибки носит название разложения на смещение и разброс.

Есть некоторая обучающая выборка: $X = (x_i, y_i)$. Будем считать, что на пространстве всех объектов и ответов \times существует распределение $p(x, y)$, из которого сгенерирована выборка X и ответы на ней.

Рассмотрим квадратичную функцию потерь

$$L(y, a) = (y - a)^2$$

Посмотрим на математическое ожидание по x и y квадрата данной ошибки:

$$R(a) =_{x,y} \left[(y - a(x))^2 \right] = \int \int p(x, y) (y - a(x))^2 dx dy.$$

$R(a)$ называется среднеквадратичным риском. Некоторый способ посчитать ошибку, если знать распределение.

Разложение ошибки

Bias-Variance Decomposition

- Модель переобучена?
- Плохо предсказывает

$$\text{Error} = \underbrace{\text{Bias}^2(a(x))}_{\text{СМЕЩЕНИЕ}} + \underbrace{\text{Var}(a(x))}_{\text{РАЗБОС}} + \underbrace{\sigma^2}_{\text{ШУМ}}$$

• $\text{Bias}(a(x))$ - средняя ошибка по всем возможным наборам данных - смещение

• $\text{Var}(a(x))$ - дисперсия ошибки, как сильно разнится ошибка при обучении на различных наборах данных - разброс

• σ^2 - неограничиваемая ошибка (ШУМ)

Доказательство:

• Истинная зависимость $\rightarrow y = f(\bar{x}) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$
 Средняя ошибка на всех отборках \rightarrow Модель $\rightarrow \hat{f}(\bar{x})$

$$\text{Err} = E_x[(y - \hat{f}(x))^2] = E(y^2) + E(\hat{f}(x)^2) - 2 \cdot E(y \cdot \hat{f}(x))$$

$$E(y^2) = \text{Var}(y) + (E(y))^2$$

$$\Rightarrow E y^2 = \sigma^2 + f^2$$

$$\Rightarrow E(y \cdot \hat{f}) = E((f + \epsilon) \cdot \hat{f}) = E(f \hat{f}) + E(\epsilon \hat{f})$$

$$= f \cdot E(\hat{f}) + E(\epsilon) \cdot E(\hat{f}) = f \cdot E(\hat{f})$$

т.к. шум в данных не от чего не зависит

$$E(y) = E(f + \epsilon) = E(f)$$

$$\text{Var}(y) = E(y - E(y))^2 = E(y - f)^2 = \sigma^2$$

$$\text{Error} = \sigma^2 + f^2 + \underbrace{E(\hat{f}(x))^2}_{(2)} + 2f \cdot E(\hat{f}) - \underbrace{E(\hat{f}(x))^2}_{(2)} = \text{Var}(\hat{f}) + \text{Bias}^2$$

$$= \sigma^2 + f^2 + \text{Var}(\hat{f}) + (E(\hat{f}))^2 - 2 \cdot f \cdot E(\hat{f}) =$$

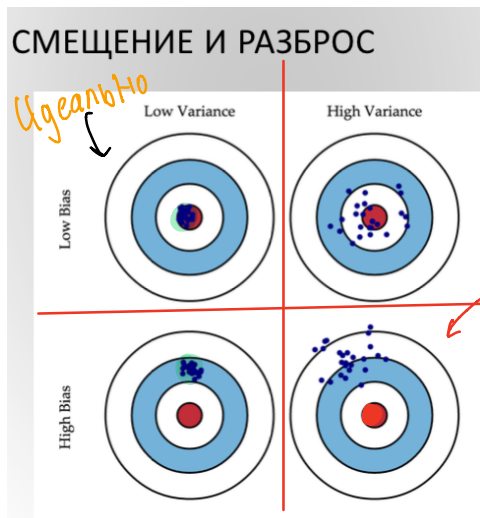
$$= (f - E\hat{f})^2 + \text{Var}\hat{f} + \sigma^2$$

$$\text{Error} = \underbrace{\text{Bias}^2(a(x))}_{\text{СМЕЩЕНИЕ}} + \underbrace{\text{Var}(a(x))}_{\text{РАЗБРОС}} + \underbrace{\sigma^2}_{\text{ШУМ}}$$

• $\text{Bias}(a(x))$ - средняя ошибка по всем возможным наборам данных - смещение

• $\text{Var}(a(x))$ - дисперсия ошибки, как сильно варьируется ошибка при обучении на различных наборах данных - разброс

• σ^2 - неограниченная ошибка (ШУМ)



• От чего зависит смещение и разброс?

• Сложность модели (кол-во параметров)

Совсем плохо

Сложность $\uparrow \Rightarrow \text{Bias} \downarrow \rightarrow \text{Var} \uparrow$
 Сложность $\downarrow \Rightarrow \text{Bias} \uparrow \rightarrow \text{Var} \downarrow$

