

1.

N°1

Для того чтобы посчитать смещение ($bias^2$) необходимо выписать следующее мат. ожидание:

$$E_{x,y} [(E_x(\mu(X)(x)) - E[y|x])^2]$$

$$E_x(\mu(X)(x)) = c \quad (\text{т.к. } \mu(X) = c)$$

$$E[y|x] = x^T x \quad (\text{по условию})$$

$$\Rightarrow E_{x,y} [E_x(\mu(X)(x)) - E[y|x]]^2 =$$

$$= E_{x,y} [c - x^T x]^2$$

Так как объекты в выборке независимы,

$$\text{то } E(\zeta_1 + \zeta_2 + \dots + \zeta_n) = E\zeta_1 + E\zeta_2 + \dots + E\zeta_n$$

$$\Rightarrow E_{x,y} [c - x^T x]^2 = \sum_{j=1}^n \frac{E_{x_j} [(c - x_j^T x_j)^2]}{n} = \text{одн. выборки}$$

$$= \frac{1}{n} E_x [(c - x^T x)^2] \quad (\text{т.к. объекты независимы})$$

какой-то
независимый
объект в выборке

$$p(x_1) = \dots = p(x_d) = 1$$

(т.к. равномерное
распределение)

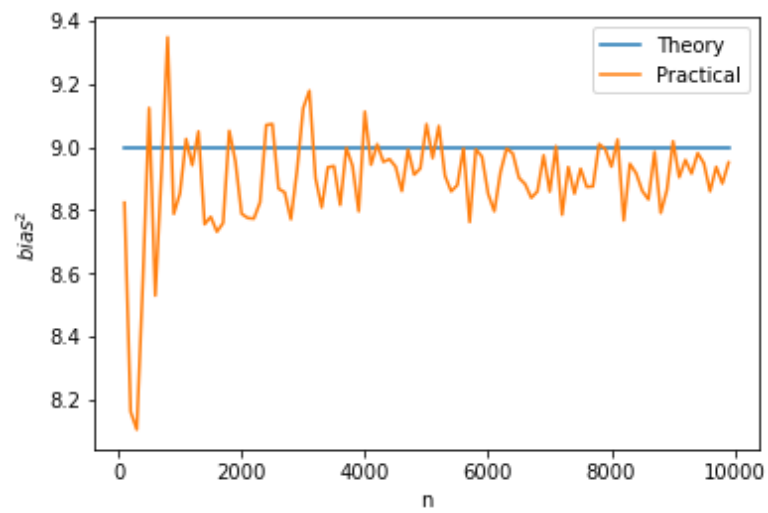
$$E_x [c - x^T x] = \int (c - \sum_{i=1}^d x_i^2) p(x_1) \dots p(x_d) dx_1 \dots dx_d =$$

$$= \int_{x_i \in [0,1]} c p(x_1) \dots p(x_d) dx_1 \dots dx_d - \int_{x_i \in [0,1]} \sum_{i=1}^d x_i^2 p(x_1) \dots p(x_d) dx_1 \dots dx_d =$$

$$= c - \sum_{i=1}^d \int x_i^2 dx_1 \dots dx_d = c - \frac{d}{3}$$

$$\Rightarrow bias^2 = (c - \frac{d}{3})^2$$

Проделаем численный эксперимент и сравним
результаты.



Теоретический расчет и результаты эксперимента не противоречат друг другу, при увеличении размера матрицы видна тенденция стремления к теоретическому значению.

2.

(Nº2)

$$E_{x,y} [(E_X(\mu(X)(x)) - E[y|x])^2]$$

$$\mu(X): \hat{f}_x = \frac{1}{N} \sum_{i=1}^N I[x_i = x] y_i$$

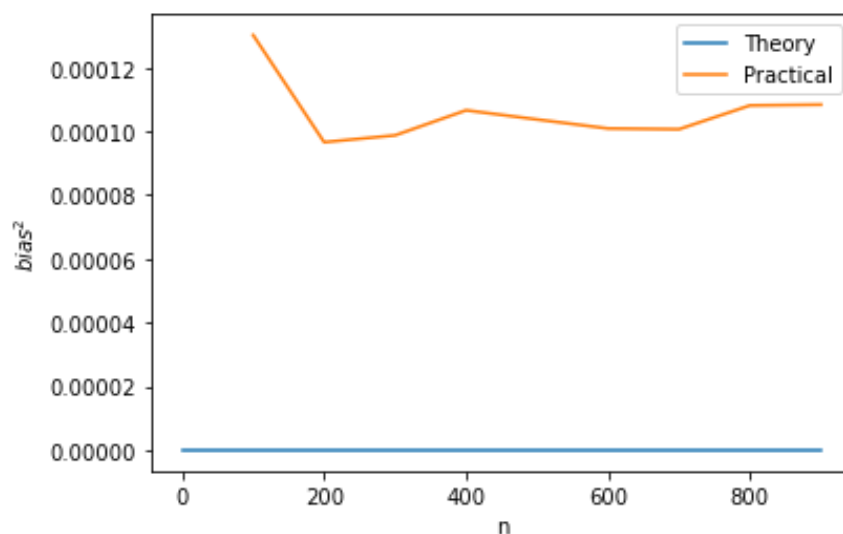
$$E_X(\mu(X)(x)) = E_X \frac{1}{N} \sum_{i=1}^N I[x_i = x] (f_x + \varepsilon) =$$

$$= \left\{ \begin{array}{l} E \varepsilon = 0 \\ \text{по условию} \end{array} \right\} = f_x \text{ т.к. распределение } x \text{ равно-} \\ \text{вероятно, то и } f_x \text{ постоянное, } \Rightarrow \text{это справедливо для } \forall x.$$

$$E[y|x] = E[f_x + \varepsilon | x] = E f_x + \underbrace{E \varepsilon}_0 = f_x$$

$$\Rightarrow \text{bias}^2 = 0$$

Продолает глоссарий экспериментов.



Эксперименты в среднем дают около 0.0001, что близко к теоретически рассчитанному значению 0.

3.

$N=3$

AUC ROC - это пер объектов, которые алгоритм верно упорядочил.

$$\sum_{i=1}^q \sum_{j=1}^q \frac{I[y_i < y_j] I'[a_i < a_j]}{\sum_{i=1}^q \sum_{j=1}^q I[y_i < y_j]}$$

$$I'[a_i < a_j] = \begin{cases} 0, & a_i > a_j \\ 0.5, & a_i = a_j \\ 1, & a_i < a_j \end{cases} \quad I[y_i < y_j] = \begin{cases} 0, & y_i \geq y_j \\ 1, & y_i < y_j \end{cases}$$

Самое важное в данной задаче вспомнить про случай, когда $a_i = a_j$

Таким образом, при использовании функции $\text{ceil}(a; 0.5)$, которая все ответы алгоритма переводит 0.5 в 0.5, можно максимально увеличить/уменьшить ROC AUC на 0.5, т.е. свести к случаю.

Рассмотрим пример увеличения и уменьшения.

уменьшаем на 0.5

```
y_true = [1, 1, 1, 0, 0, 0]
y_pred = [0.9, 0.8, 0.6, 0.5, 0.5, 0.5]

print('ROC_AUC : ', roc_auc_score(y_true, y_pred))

y_true = [1, 1, 1, 0, 0, 0]
y_pred = [0.5, 0.5, 0.5, 0.5, 0.5, 0.5]

print('ROC_AUC : ', roc_auc_score(y_true, y_pred))

ROC_AUC : 1.0
ROC_AUC : 0.5
```

увеличиваем на 0.5

```
y_true = [0, 1, 1, 1, 1, 1]
y_pred = [0.9, 0.8, 0.6, 0.5, 0.5, 0.5]

print('ROC_AUC : ', roc_auc_score(y_true, y_pred))

y_true = [0, 1, 1, 1, 1, 1]
y_pred = [0.5, 0.5, 0.5, 0.5, 0.5, 0.5]

print('ROC_AUC : ', roc_auc_score(y_true, y_pred))

ROC_AUC : 0.0
ROC_AUC : 0.5
```

4.

№4 На картинке: TP - синий
FN - зеленый
FP - красный
TN - белый

$$TPR = \frac{TP}{TP+FN} = TP$$

$$FPR = \frac{FP}{FP+TN} = FP$$

x - значение порога

$$\Rightarrow TPR = \int_1^x (-1,5z^2 + 3z) dz = \frac{x^3}{2} - \frac{3x^2}{2} + 1$$

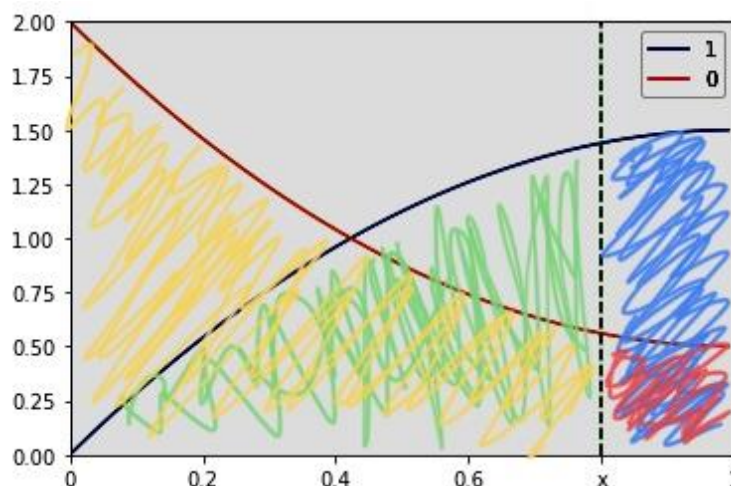
$$FPR = \int_x^1 (2 - (-1,5z^2 + 3z)) dz = \frac{x^3}{2} - \frac{3x^2}{2} + 2x$$

ROC-AUC = $\int_1^0 TPR(t) FPR'(t) dt =$
(площадь под кривой FPR-TPR)

$$= \int_1^0 \left(\frac{t^3}{2} - \frac{3t^2}{2} + 1 \right) \left(\frac{3t^2}{2} - \frac{3 \cdot 2t}{2} + 2 \right) dt = \boxed{0,75}$$

Пункт порог разделения в этот момент:
 $-1,5z^2 + 3z = 2 + 1,5z^2 - 3z$
 $z \approx 0,42$

В этот момент $TPR = TP \approx 0,77$
 $FPR = FP \approx 0,61$



Интересно посчитать accuracy

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$TN = \int_0^x (2 + 1,5z^2 - 3z) dz = \frac{x^3}{2} - \frac{3x^2}{2} + 2x$$

$$\text{При } x = 0,42 \quad TN \approx 0,61, \Rightarrow$$

$$\text{accuracy} \approx 0,69$$

То есть, ROC-AUC зависит максимально
значущую точность.

5.

$N=5$

Мдело заключается в том, что для каждого j -го признака вводится два кластера: зашумленные и незашумленные объекты по признаку j .

зашумленные ← незашумленные

(a) $P(X, Y) = \sum_{j=1}^d \pi_j P(X_j) + (1 - \pi_j) P(Y_j)$

Введем латентные переменные z_{ij} , которые могут принимать значения:

$z_{ij} = 1$, зашумлен по признаку j
 $z_{ij} = 0$, не зашумлен по признаку j

Тогда $\pi_j, \lambda_j^0, j=1, \dots, d$

$$P(X | Y, \theta) = \sum_{i=1}^n \sum_{j=1}^d \sum_{k=0}^1 P(z_{ij}=k) P(x_{ij} | z_{ij}=k, \theta)$$

$$= \sum_{i=1}^n \sum_{j=1}^d (1 - \pi_j) P(Y_j) + \pi_j [P(Y_j) + P_{\text{Pois}}(\lambda_j^0)]$$

$$P_{\text{Pois}}(\lambda_j^0) = \sum_{l=0}^{\infty} \frac{\lambda_j^l e^{-\lambda_j^0}}{l!}, \quad l=0, 1, \dots$$

(b) E-вер:

$$q(z) = \prod_{i=1}^n \prod_{j=1}^d q_{ij}(z_{ij})$$

$$q(z) = \sum_{i=1}^n \sum_{j=1}^d q_{ij}(z_{ij})$$

$$q_{ij}(z_{ij} = k) = \begin{cases} (1 - \pi_j^0) P(y_j^0) & z_{ij} = 0 \\ \pi_j^0 [P(y_j^0) + P_{\text{prior}}(\lambda_j^0)] & z_{ij} = 1 \end{cases}$$

(b) E-вер:

$$\theta = \{\pi_j^0, \lambda_j^0\}_{j=1}^d$$

$$q(z) = \prod_{i=1}^n \prod_{j=1}^d q_{ij}(z_{ij})$$

$$q(z) = \sum_{i=1}^n \sum_{j=1}^d q_{ij}(z_{ij})$$

$$q_{ij}(z_{ij} = k) \sim \begin{cases} (1 - \pi_j^0) P(y_j^0), & z_{ij} = 0 \\ \pi_j^0 [P(y_j^0) + P_{\text{prior}}(\lambda_j^0)], & z_{ij} = 1 \end{cases}$$

// с учетом
доп. информации
(из семантики)

(c) M-вер:

$$E_{q(z)} \log P(x, y, z | \theta) \rightarrow \max_{\theta}$$

$$\sum_{i=1}^n \sum_{j=1}^d (1 - \pi_j^0) P(y_j^0) \log P(x_{ij}^0, z_{ij}^0, y_{ij}^0 | \theta) +$$

$$+ \pi_j^0 [P(y_j^0) + P_{\text{prior}}(\lambda_j^0)] \log P(x_{ij}^1, y_{ij}^1, z_{ij}^1 | \theta) =$$

$$= \sum_{i=1}^n \sum_{j=1}^d f_{ij} \log$$

$$\begin{aligned}
 \text{Eq (17)} \quad \log P(X, Y, Z | \theta) &\rightarrow \max_{\theta} \\
 \sum_{i=1}^n \sum_{j=1}^d \sum_{k=0}^1 d_{ij}^k(z_{ij}=k) \log P(x_{ij}^0, y_{ij}^0, z_{ij}^0 | \theta) &= \\
 = \sum_{i=1}^n \sum_{j=1}^d f_{ij}^0 [\log P(y_{ij}^0) + \log \pi_j^0] + & \\
 + f_{ij}^1 [\log (P(y_{ij}^0) + P_{\text{Pois}}(\lambda_j^0)) + \log(1 - \pi_j^0)] & \\
 &\xrightarrow{\max_{\lambda_j^0, \pi_j^0}}
 \end{aligned}$$

(d) Пусть на вход поступает матрица X' .

Необходимо для каждого признака j определить зашумлен он или нет.

То есть, для $\forall j$ определить значение скрытой переменной z .

① Для этого надо сравнить $d(z=0)$ и $d(z=1)$

признак j зашумлен, если

$$\sum_{i=1}^n \pi_j^0 [P(y_{ij}^0) + P_{\text{Pois}}(\lambda_j^0)] > \sum_{i=1}^n (1 - \pi_j^0) P(y_{ij}^0)$$

② Из зашумленных признаков выбрать Пуассоновский шум с распределением $P_{\text{Pois}}(\lambda_j^0)$