

# Введение в искусственный интеллект. Машинное обучение

## Лекция 3. Регрессия и оценка качества

Бабин Д.Н., Иванов И.Е., Петюшко А.А.

MaTIC

1 марта 2019г.

- 1 Регрессия
- 2 Оценка качества

## Постановка задачи и допущения

- $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$

## Постановка задачи и допущения

- $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$
- $a(x) = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ , где  $\theta = (\theta_0, \theta_1, \dots, \theta_n)^T \in \mathbb{R}^{n+1}$  — параметры модели.

## Постановка задачи и допущения

- $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$
- $a(x) = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ , где  $\theta = (\theta_0, \theta_1, \dots, \theta_n)^T \in \mathbb{R}^{n+1}$  — параметры модели.
- Удобно писать в векторном виде

$$a(x) = \theta^T \cdot x,$$

где  $x = (x_0, x_1, \dots, x_n)^T$  и  $x_0 = 1$ .

## Постановка задачи и допущения

- $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$
- $a(x) = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ , где  $\theta = (\theta_0, \theta_1, \dots, \theta_n)^T \in \mathbb{R}^{n+1}$  — параметры модели.
- Удобно писать в векторном виде

$$a(x) = \theta^T \cdot x,$$

где  $x = (x_0, x_1, \dots, x_n)^T$  и  $x_0 = 1$ .

## Метод наименьших квадратов

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2$  — функция потерь

## Постановка задачи и допущения

- $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$
- $a(x) = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ , где  $\theta = (\theta_0, \theta_1, \dots, \theta_n)^T \in \mathbb{R}^{n+1}$  — параметры модели.
- Удобно писать в векторном виде

$$a(x) = \theta^T \cdot x,$$

где  $x = (x_0, x_1, \dots, x_n)^T$  и  $x_0 = 1$ .

## Метод наименьших квадратов

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$

## Теорема

Решением задачи  $\arg \min_{\theta} \left( \sum_{i=1}^{\ell} (\theta^T \cdot x^{(i)} - y_i)^2 \right)$  является  $\hat{\theta} = (X^T X)^{-1} \cdot X^T \cdot y$ , где  $X_{i,j} = x_j^{(i)}$ ,  $y = (y_1, \dots, y_{\ell})$ .

## Доказательство

Запишем задачу в векторном виде  $\|X\theta - y\|^2 \rightarrow \min_{\theta}$ . Необходимое условие минимума в матричном виде имеет вид:

$$\frac{\partial}{\partial \theta} \|X\theta - y\|^2 = \frac{\partial}{\partial \theta} \left( (X\theta - y)^T \cdot (X\theta - y) \right) = 2X^T(X\theta - y) = 0,$$

откуда получаем  $X^T X\theta = X^T y$ , из чего и следует требуемое.



## Определение

Пусть  $\theta = (\theta_1, \dots, \theta_n)$  — вектор столбец, а  $z = z(\theta_1, \dots, \theta_n)$ . Тогда определим

$$\frac{\partial z}{\partial \theta} := \left( \frac{\partial z}{\partial \theta_1}, \dots, \frac{\partial z}{\partial \theta_n} \right)^T$$

## Лемма 1

$$\frac{\partial}{\partial x} x^T a = a$$

## Лемма 2

$$\frac{\partial}{\partial x} x^T A x = (A + A^T)x$$

## Модель шума

$$y(x_i) = f_{\theta}(x_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

## Метод максимума правдоподобия

$$L(\varepsilon_1, \dots, \varepsilon_n | \theta) = \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{\varepsilon_i^2}{2\sigma_i^2}\right) \rightarrow \max_{\theta}$$

$$-L(\varepsilon_1, \dots, \varepsilon_n | \theta) = \text{const}(\theta) + \frac{1}{2} \sum_i \frac{1}{\sigma_i^2} (f_{\theta}(x_i) - y_i)^2 \rightarrow \min_{\theta}$$

# Преимущества и недостатки линейной регрессии

## Преимущества

- Простой алгоритм, вычислительно не сложный
- Линейная регрессия хорошо интерпретируемая модель
- Несмотря на свою простоту может описывать довольно сложные зависимости (например, полиномиальные)

# Преимущества и недостатки линейной регрессии

## Преимущества

- Простой алгоритм, вычислительно не сложный
- Линейная регрессия хорошо интерпретируемая модель
- Несмотря на свою простоту может описывать довольно сложные зависимости (например, полиномиальные)

## Недостатки

- Алгоритм предполагает, что все признаки числовые
- Алгоритм предполагает, что данные распределены нормально, что не всегда так
- Алгоритм сильно чувствителен к выбросам



## L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь

## L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$



## L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$

## Теорема

Решением задачи  $\arg \min_{\theta} (\sum_{i=1}^{\ell} (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2)$  является

$\hat{\theta} = (X^T X + \alpha I_n)^{-1} \cdot X^T \cdot y$ , где  $X_{i,j} = x_j^{(i)}$ ,  $y = (y_1, \dots, y_{\ell})$ ,  $I_n$  — единичная матрица.

## L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$

## Теорема

Решением задачи  $\arg \min_{\theta} (\sum_{i=1}^{\ell} (\theta^T \cdot x^{(i)} - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n \theta_i^2)$  является

$\hat{\theta} = (X^T X + \alpha I_n)^{-1} \cdot X^T \cdot y$ , где  $X_{i,j} = x_j^{(i)}$ ,  $y = (y_1, \dots, y_{\ell})$ ,  $I_n$  — единичная матрица.

## Доказательство

Аналогично, что и требовалось доказать.

## L1-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \alpha \sum_{i=0}^n |\theta_i| = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \alpha \sum_{i=0}^n |\theta_i|$  — функция потерь

## L1-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + \alpha \sum_{i=0}^n |\theta_i| = \sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + \alpha \sum_{i=0}^n |\theta_i|$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$

## Свойства

- Эта регуляризация обеспечивает отбор признаков

## L1-регуляризация и L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + r\alpha \sum_{i=0}^n |\theta_i| + (1-r)\frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 =$   
 $\sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + r\alpha \sum_{i=0}^n |\theta_i| + (1-r)\frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь

## L1-регуляризация и L2-регуляризация

- $L(\theta, X_{train}) = MSE(\theta, X_{train}) + r\alpha \sum_{i=0}^n |\theta_i| + (1-r)\frac{\alpha}{2} \sum_{i=0}^n \theta_i^2 =$   
 $\sum_i (\theta^T \cdot x^{(i)} - y_i)^2 + r\alpha \sum_{i=0}^n |\theta_i| + (1-r)\frac{\alpha}{2} \sum_{i=0}^n \theta_i^2$  — функция потерь
- Задача найти  $\hat{\theta} = \arg \min_{\theta} (L(\theta, X_{train}))$

## Свойства

- Нет аналитического решения

- Линейная регрессия — простая, хорошо интерпретируемая модель, не устойчивая к выбросам
- Имеет наглядную вероятностную интерпретацию
- Регуляризация — отличный способ борьбы с переобучением и шумом в данных

# Классификация ответов бинарного классификатора

- Задача классификации на 2 класса:  $X \rightarrow Y, Y = \{+1, -1\}$
- Алгоритм классификации  $a(x_i) = y_i$
- Класс с меткой “+1” называется “**positive**”
- Класс с меткой “-1” называется “**negative**”



# Классификация ответов бинарного классификатора

- Задача классификации на 2 класса:  $X \rightarrow Y, Y = \{+1, -1\}$
- Алгоритм классификации  $a(x_i) = y_i$
- Класс с меткой “+1” называется “**positive**”
- Класс с меткой “-1” называется “**negative**”


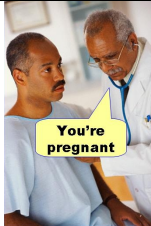


	Выход алгоритма	Правильный ответ
TP (True Positive)	$a(x_i) = +1$	$y_i = +1$
TN (True Negative)	$a(x_i) = -1$	$y_i = -1$
FP (False Positive)	$a(x_i) = +1$	$y_i = -1$
FN (False Negative)	$a(x_i) = -1$	$y_i = +1$

# Матрица ошибок

Более наглядно эти соотношения можно изобразить с помощью **матрицы ошибок (confusion matrix)**

		Правильный ответ	
		$y = +1$	$y = -1$
Выход алгоритма	$a(x) = +1$	True Positive	False Positive (Ошибка 1 рода)
	$a(x) = -1$	False Negative (Ошибка 2 рода)	True Negative

# Матрица ошибок

	$y = +1$	$y = -1$
$a(x) = +1$		
$a(x) = -1$		

# Простейшая метрика качества

- Простейшая метрика качества - это доля правильных ответов на тесте (контрольной выборке)
- По-английски - **Accuracy**

## Формула Accuracy

$$Accuracy = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i] = \frac{TP+TN}{TP+FP+TN+FN}$$

# Простейшая метрика качества

- Простейшая метрика качества - это доля правильных ответов на тесте (контрольной выборке)
- По-английски - **Accuracy**

## Формула Accuracy

$$Accuracy = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i] = \frac{TP+TN}{TP+FP+TN+FN}$$

## Недостаток

- Не учитывается дисбаланс классов
- Не учитывается цена ошибки на объектах разных классов

# Метрики по положительному отклику алгоритма

Рассмотрим метрики, которые основаны на подсчёте доли положительных ответов алгоритма.

Доля ложных положительных классификаций

Также известно как False Positive Rate, или **FPR**.

$$FPR(a, X^\ell) = \frac{\sum_{i=1}^{\ell} [y_i = -1][a(x_i) = +1]}{\sum_{i=1}^{\ell} [y_i = -1]}$$

# Метрики по положительному отклику алгоритма

Рассмотрим метрики, которые основаны на подсчёте доли положительных ответов алгоритма.

## Доля ложных положительных классификаций

Также известно как False Positive Rate, или **FPR**.

$$FPR(a, X^\ell) = \frac{\sum_{i=1}^{\ell} [y_i = -1][a(x_i) = +1]}{\sum_{i=1}^{\ell} [y_i = -1]}$$

## Доля верных положительных классификаций

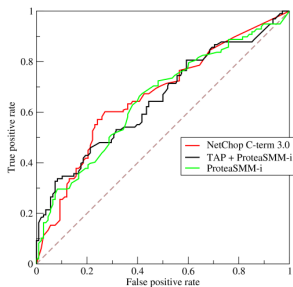
Также известно как True Positive Rate, или **TPR**.

$$TPR(a, X^\ell) = \frac{\sum_{i=1}^{\ell} [y_i = +1][a(x_i) = +1]}{\sum_{i=1}^{\ell} [y_i = +1]}$$

**Замечание.** Обратите внимание на разные знаменатели!

# Кривая ошибок

Наиболее известна как рабочая характеристика приёмника, или Receiver Operating Characteristic (**ROC-кривая**), в который мы смотрим на компромисс между уровнем ложной тревоги и долей верного отклика.



По оси X откладывается FPR, по оси Y - TPR<sup>1</sup>.

**Замечание.** На данной кривой никак не учитываются пропуски.

<sup>1</sup><https://wikipedia.org>



## AUROC

Чем больше для каждого значения ошибки FPR значение правильного предсказания TPR, тем лучше работает классификатор.

Т.о., площадь под кривой (Area Under Curve, AUC / AUROC) необходимо максимизировать.

---

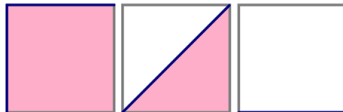
<sup>2</sup><https://dyakonov.org>

## AUROC

Чем больше для каждого значения ошибки FPR значение правильного предсказания TPR, тем лучше работает классификатор.

Т.о., площадь под кривой (Area Under Curve, AUC / AUROC) необходимо максимизировать.

Наглядны ROC-кривые для наилучшего ( $AUC=1$ ), случайного ( $AUC=0.5$ ) и наихудшего ( $AUC=0$ ) алгоритма<sup>2</sup>.



<sup>2</sup><https://dyakonov.org>

# Задача

Предположим, что алгоритм бинарной классификации  $a(x_i)$  принимает решение о присвоении класса на основе некоторого скалярного значения  $g_\theta(x_i) \in \mathbb{R}$ , где  $\theta$  - набор параметров модели, а  $g_\theta(x_i)$  - дискриминантная функция.

## Задача

- Хотим построить ROC-кривую, т.е. найти точки  $\{(FPR_i, TPR_i)\}_{i=1}^\ell$
- Подсчитать площадь под кривой - AUROC

# Задача

Предположим, что алгоритм бинарной классификации  $a(x_i)$  принимает решение о присвоении класса на основе некоторого скалярного значения  $g_\theta(x_i) \in \mathbb{R}$ , где  $\theta$  - набор параметров модели, а  $g_\theta(x_i)$  - дискриминантная функция.

## Задача

- Хотим построить ROC-кривую, т.е. найти точки  $\{(FPR_i, TPR_i)\}_{i=1}^\ell$
- Подсчитать площадь под кривой - AUROC

Подсчитаем количество правильных ответов разного типа:

- $\ell_+ = \sum_{i=1}^\ell [y(x_i) = +1]$
- $\ell_- = \sum_{i=1}^\ell [y(x_i) = -1]$  (понятно, что  $\ell = \ell_+ + \ell_-$ )

Упорядочим обучающую выборку  $X^\ell$  по убыванию значений  $g_\theta(x_i)$ .

Тогда формула для  $AUROC = \frac{1}{\ell_-} \sum_{i=1}^\ell [y_i = -1] TPR_i$ .

## Алгоритм

Первую точку ставим в начало координат:  $(FPR_0, TPR_0) = (0, 0)$ ,  $AUROC = 0$ .

# Решение задачи

## Алгоритм

Первую точку ставим в начало координат:  $(FPR_0, TPR_0) = (0, 0)$ ,  $AUROC = 0$ .

Цикл по упорядоченной выборке  $i = 1 \dots \ell$

Если  $y_i = -1$ :

- $(FPR_i, TPR_i) = (FPR_{i-1} + \frac{1}{\ell_-}, TPR_{i-1})$  (двигаемся по оси X)
- $AUROC = AUROC + \frac{1}{\ell_-} TPR_i$

## Алгоритм

Первую точку ставим в начало координат:  $(FPR_0, TPR_0) = (0, 0)$ ,  $AUROC = 0$ .

Цикл по упорядоченной выборке  $i = 1 \dots \ell$

Если  $y_i = -1$ :

- $(FPR_i, TPR_i) = (FPR_{i-1} + \frac{1}{\ell_-}, TPR_{i-1})$  (двигаемся по оси X)
- $AUROC = AUROC + \frac{1}{\ell_-} TPR_i$

Если  $y_i = +1$ :

- $(FPR_i, TPR_i) = (FPR_{i-1}, TPR_{i-1} + \frac{1}{\ell_+})$  (двигаемся по оси Y)

## В задачах информационного поиска

- Точность, или  $Precision = \frac{TP}{TP+FP}$  (доля релевантных объектов среди найденных)
- Полнота, или  $Recall = \frac{TP}{TP+FN}$  (доля найденных объектов среди релевантных)



## В задачах информационного поиска

- Точность, или  $Precision = \frac{TP}{TP+FP}$  (доля релевантных объектов среди найденных)
- Полнота, или  $Recall = \frac{TP}{TP+FN}$  (доля найденных объектов среди релевантных)

## Как применяются

- **Точность:** позволяет следить, чтобы было мало ложных тревог; но при этом ничего не говорит о пропусках (высока цена ложной тревоги, а цена пропуска - низкая).
- **Полнота:** позволяет следить, чтобы было мало пропусков; но при этом ничего не говорит о ложных тревогах (высока цена пропуска, а цена ложной тревоги - низкая).

**Замечание.** Зачастую задача состоит в оптимизации одной метрики при фиксации другой.

### В задачах медицинской диагностики

- Чувствительность, или  $Sensitivity = \frac{TP}{TP+FN}$  (доля верных положительных диагнозов)
- Специфичность, или  $Specificity = \frac{TN}{TN+FP}$  (доля верных отрицательных диагнозов)

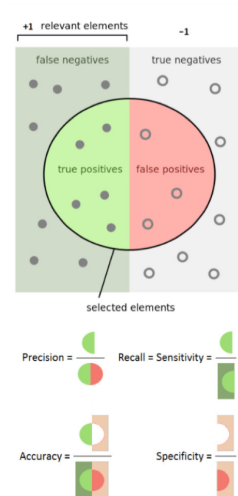
### В задачах медицинской диагностики

- Чувствительность, или  $Sensitivity = \frac{TP}{TP+FN}$  (доля верных положительных диагнозов)
- Специфичность, или  $Specificity = \frac{TN}{TN+FP}$  (доля верных отрицательных диагнозов)

### Как применяются

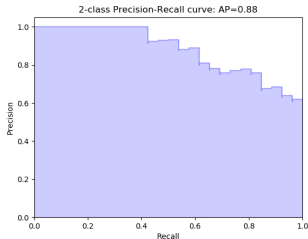
- **Чувствительность:** максимизируем количество верных положительных диагнозов, но не учитываем ложные диагнозы (стоимость лечения низкая, а цена пропуска - высокая).
- **Специфичность:** максимизируем количество верных отрицательных диагнозов, но не учитываем пропуски диагноза (стоимость лечения высокая, а цена пропуска - низкая).

# Иллюстрация метрик



# Агрегированные метрики над Precision-Recall

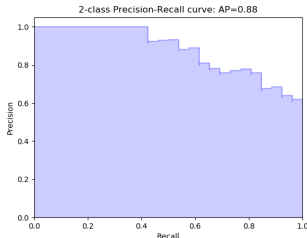
Можно построить кривую Точность-Полнота (PR-кривая) по аналогии с ROC-кривой:



**Замечание.** Обратите внимание, что в данном случае кривая не обязательно монотонна!

# Агрегированные метрики над Precision-Recall

Можно построить кривую Точность-Полнота (PR-кривая) по аналогии с ROC-кривой:



**Замечание.** Обратите внимание, что в данном случае кривая не обязательно монотонна!

## AUPRC

- Аналогично AUROC, можно вычислить площадь под PR-кривой - AUPRC
- Другое название - Average Precision (с некоторым допущениями на способ интегрирования): чем больше, тем лучше

Для каждого класса  $c \in Y$  обозначим через  $TP_c$ ,  $FP_c$  и  $FN_c$  верные положительные, ложные положительные и ложные отрицательные ответы. Тогда:

## Точность и полнота с макроусреднением

- $Precision = \frac{\sum_c TP_c}{\sum_c (TP_c + FP_c)}$
- $Recall = \frac{\sum_c TP_c}{\sum_c (TP_c + FN_c)}$
- Не чувствительно к ошибкам на маленьких классах

# Многоклассовая классификация

Для каждого класса  $c \in Y$  обозначим через  $TP_c$ ,  $FP_c$  и  $FN_c$  верные положительные, ложные положительные и ложные отрицательные ответы. Тогда:

## Точность и полнота с макроусреднением

- $Precision = \frac{\sum_c TP_c}{\sum_c (TP_c + FP_c)}$
- $Recall = \frac{\sum_c TP_c}{\sum_c (TP_c + FN_c)}$
- Не чувствительно к ошибкам на маленьких классах

## Точность и полнота с микроусреднением

- $Precision = \frac{1}{|Y|} \sum_c \frac{TP_c}{TP_c + FP_c}$
- $Recall = \frac{1}{|Y|} \sum_c \frac{TP_c}{TP_c + FN_c}$
- Чувствительно к ошибкам на маленьких классах



- Точность и полнота подходят для задач информационного поиска, когда доля объектов релевантного класса мала

- Точность и полнота подходят для задач информационного поиска, когда доля объектов релевантного класса мала
- Чувствительность и специфичность подходят для задач с несбалансированными классами (как, например, в медицине)

- Точность и полнота подходят для задач информационного поиска, когда доля объектов релевантного класса мала
- Чувствительность и специфичность подходят для задач с несбалансированными классами (как, например, в медицине)
- AUROC подходит для оценки качества при нефиксированном соотношении цены ошибок

# Резюме по оценкам качества классификации

- Точность и полнота подходят для задач информационного поиска, когда доля объектов релевантного класса мала
- Чувствительность и специфичность подходят для задач с несбалансированными классами (как, например, в медицине)
- AUROC подходит для оценки качества при нефиксированном соотношении цены ошибок
- Ещё одна агрегированная оценка качества - F-мера:
$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$
  - Это *гармоническое среднее*, которое стремится к нулю когда хотя бы одно из значений стремится к нулю