

# Задача классификации

2021

# Что будет ...

- **Задача классификации**
- **Логистическая регрессия**
- **Метрики**



# Классификация

**Обучающая выборка - множество объектов, разделенные на классы**

**Задача - определить класс нового объекта**

**Бинарная классификация - класса задано только 2**

**Многоклассовая классификация - классов более 2**

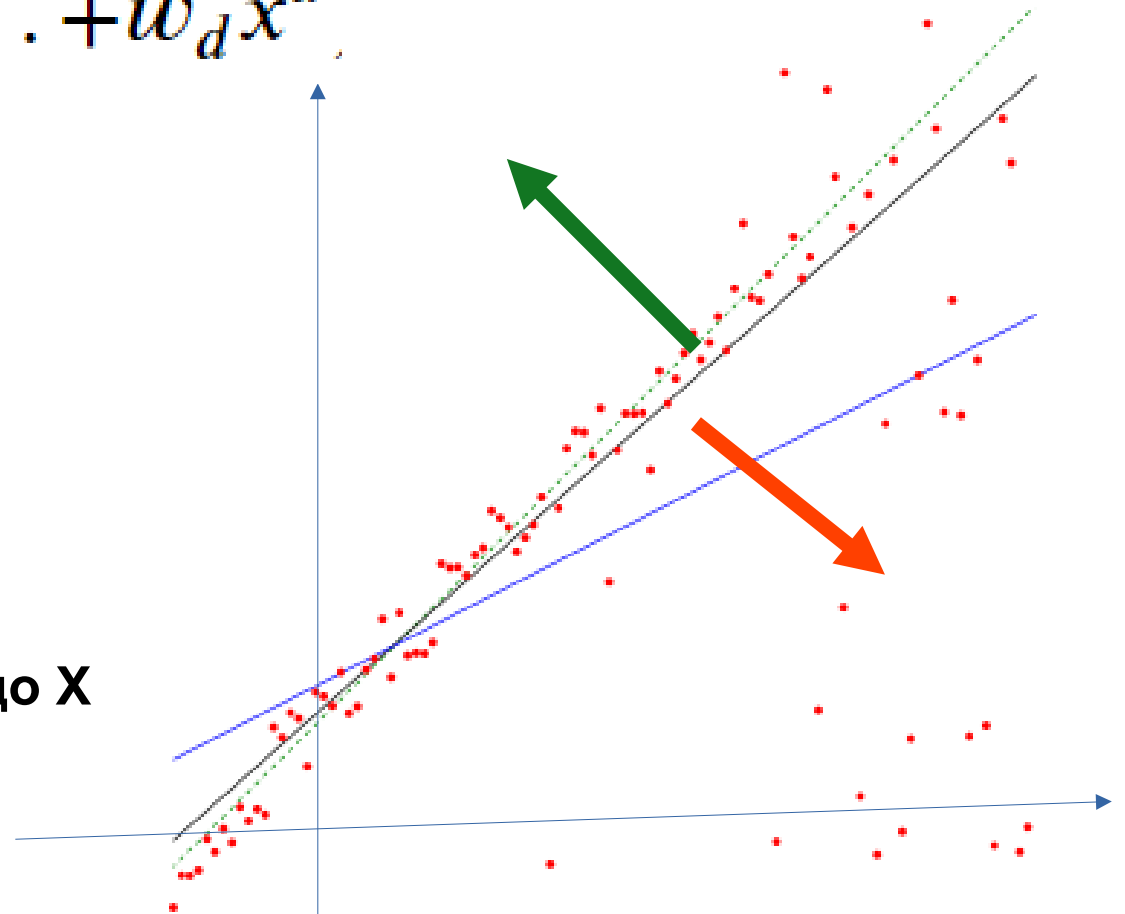
# Линейная модель

$$a(x) = w_0 + w_1 x^1 + w_2 x^2 + \dots + w_d x^d$$

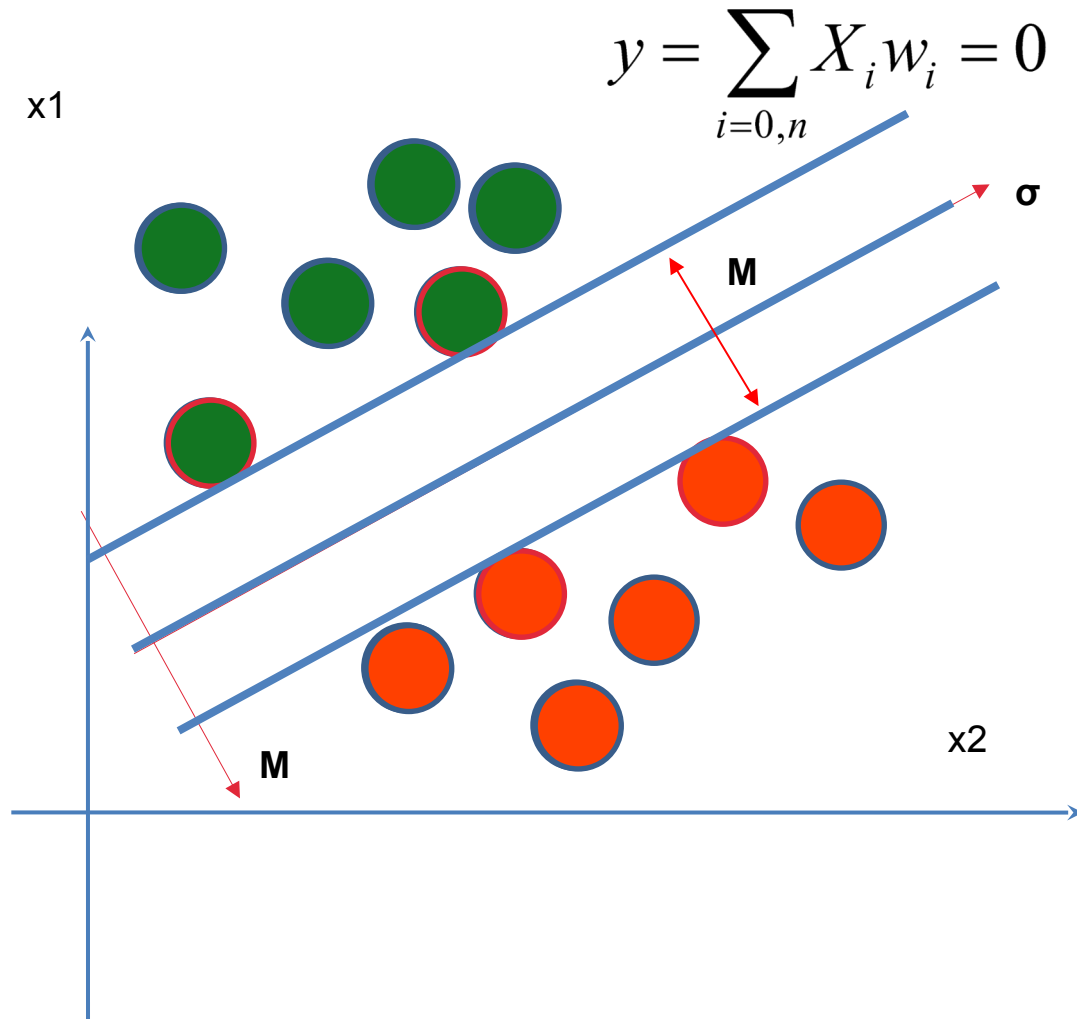
$$Q(a, X) \rightarrow \min$$

$a(x) = 0$  гиперплоскость

$\frac{|a(X)|}{||w||}$  расстояние от гиперплоскости до  $X$



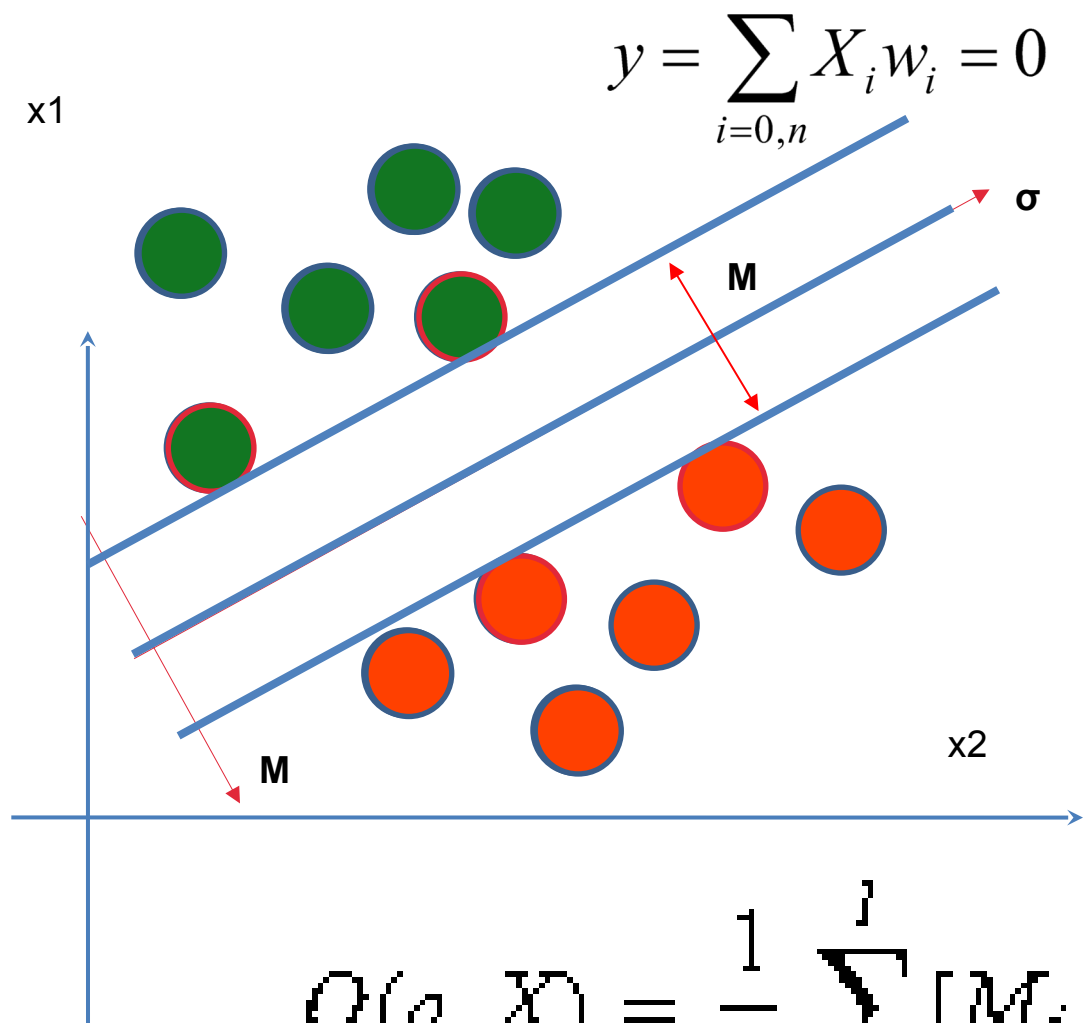
# Отступ (Зазор - Margin)



$$y_i: \{-1, 1\}$$

$$M_i = y_i \cdot \langle X_i \cdot W \rangle$$

# Классификация - Потери



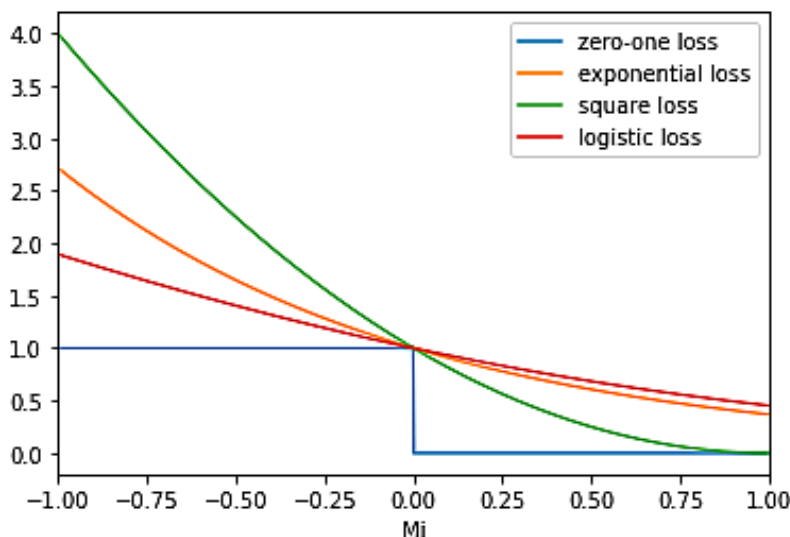
$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) \neq y_i]$$

$$y_i: \{-1, 1\}$$

$$M_i = y_i \cdot \langle X_i \cdot W \rangle$$

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l [M_i < 0] = \frac{1}{l} \sum_{i=1}^l [y_i \langle w, x_i \rangle < 0].$$

# Рабочие Потери



$$[M_i < 0] \leq \tilde{L}(M_i).$$

$$Q(a, X) \leq \tilde{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l \tilde{L}(M_i) \rightarrow \min_{\omega}$$

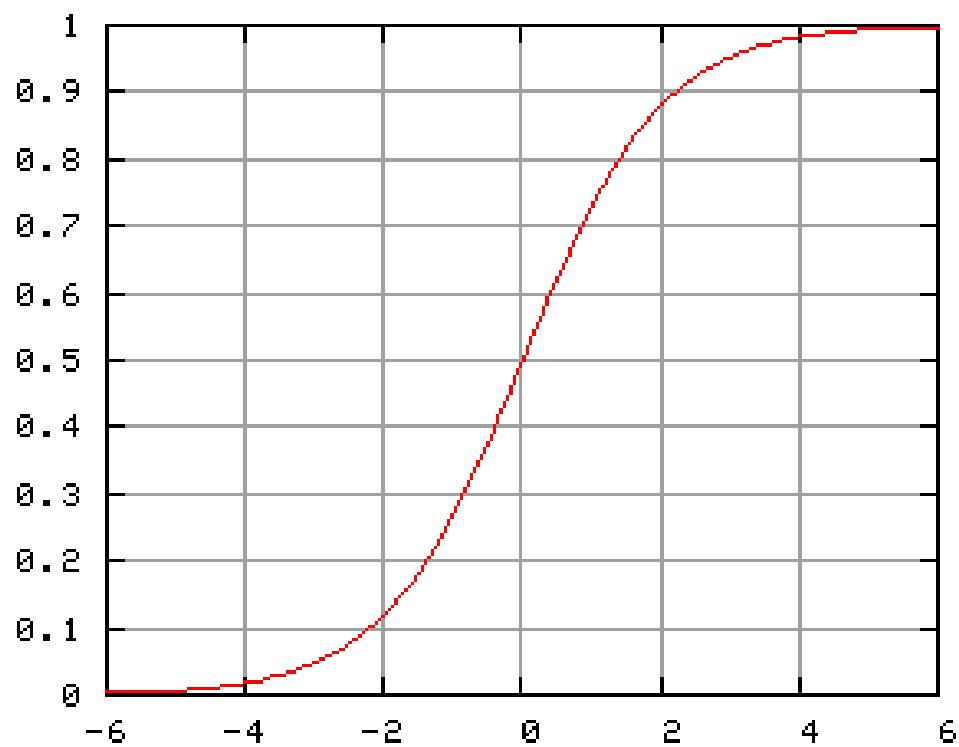
- экспоненциальная функция потерь  $\tilde{L}(M_i) = \exp(-M_i)$
- квадратичная функция потерь  $\tilde{L}(M_i) = (1 - (M_i))^2$
- логистическая функция потерь  $\tilde{L}(M_i) = \log_2(1 + \exp(-M_i))$

# Логистическая регрессия

Сигмоида:  $\sigma(z) = \frac{1}{1 + e^{-z}}$

Наша модель

$$a = \text{sigmoid}(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

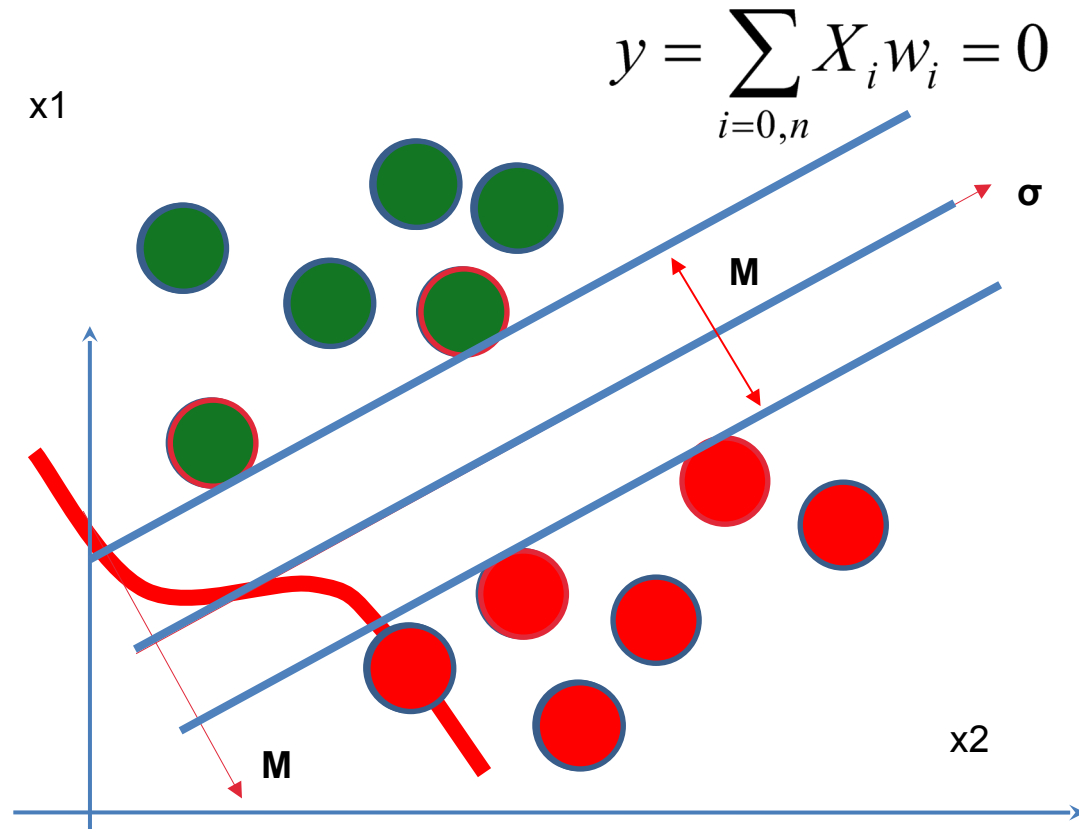




# Логистическая регрессии

$$\sigma = \frac{1}{1 + e^{-M_i}}$$

$$M_i = y_i \cdot \langle X_i \cdot W \rangle$$



$$p_+ = \sigma(\langle w, x_i \rangle) = \frac{1}{1 + \exp(-\langle w, x_i \rangle)}$$

$$\langle w, x_i \rangle = \ln \frac{p_+}{1 - p_+}$$

# Log Loss

## Метод максимального правдоподобия

**Рассматриваем задачу классификации с двумя классами: -1 и 1.**

**Выборка (по Бернулли): каждый объект с вероятностью  $p_i$  относится к классу 1 и с вероятностью  $(1-p_i)$  к классу -1.**

# Log Loss

## Метод максимального правдоподобия

Результат модели:  $a_i = a(x_i|w)$  где  $w$  - параметры модели

Функция правдоподобия:  $P(y = y_i|x_i) = p_+^{[y_i=+1]}(1 - p_+)^{[y_i=-1]}$

$$P(y|X) = L(X) = \prod_{i=1}^I p_+^{[y_i=+1]}(1 - p_+)^{[y_i=-1]}.$$

$$p(y|X, w) = \prod_i p(y_i|x_i, w) = \prod_i a_i^{y_i}(1 - a_i)^{1-y_i}$$

Для конкретной  $a(x_i|w)$   
с классами по  $y_i = \{0,1\}$

# Log Loss

## Метод максимального правдоподобия

Метод максимального правдоподобия:

$$-\ln L(X) = - \sum_{i=1}^I ([y_i = +1] \ln p_+ + [y_i = -1] \ln(1 - p_+)).$$

$$\begin{aligned} -\ln L(X) &= - \sum_{i=1}^I ([y_i = +1] \ln \frac{1}{1 + \exp(-\langle w, x_i \rangle)} + [y_i = -1] \ln(1 - \frac{1}{1 + \exp(-\langle w, x_i \rangle)})) = \\ &= \sum_{i=1}^I \ln(1 + \exp(-y_i \langle w, x_i \rangle)) \end{aligned}$$

# Log Loss

## Метод максимального правдоподобия

**Метод максимального правдоподобия:**

$$\begin{aligned} -\ln L(X) &= - \sum_{i=1}^I ([y_i = +1] \ln \frac{1}{1 + \exp(-\langle w, x_i \rangle)} + [y_i = -1] \ln(1 - \frac{1}{1 + \exp(-\langle w, x_i \rangle)})) = \\ &= \sum_{i=1}^I \ln(1 + \exp(-y_i \langle w, x_i \rangle)) \end{aligned}$$

$$[M < 0] \leq \log_2(1 + e^{-M})$$

Функции равны до  
коэффициента  $1/\ln(2)$

# Log Loss

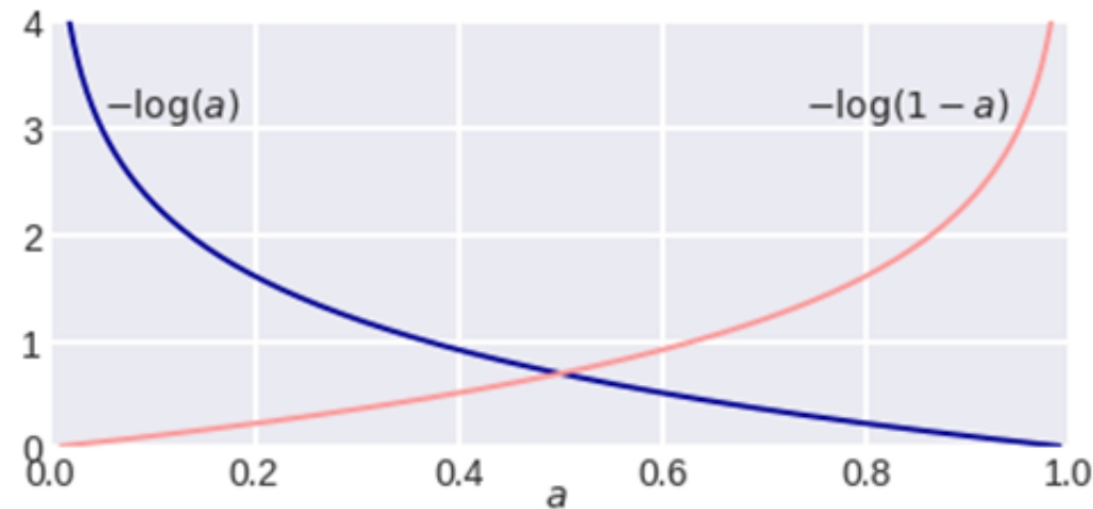
## Метод максимального правдоподобия

$$-\ln L(X) = - \sum_{i=1}^I \left( y_i \ln \frac{1}{1 + \exp(-\langle w, x_i \rangle)} + (1 - y_i) \ln \frac{\exp(-\langle w, x_i \rangle)}{1 + \exp(-\langle w, x_i \rangle)} \right)$$

LogLoss

$$= \sum_i -y_i \log a_i - (1 - y_i) \log(1 - a_i)$$

$$= - \begin{cases} \log a_i, & y_i = 1 \\ \log(1 - a_i), & y_i = 0 \end{cases}$$



# Log Loss

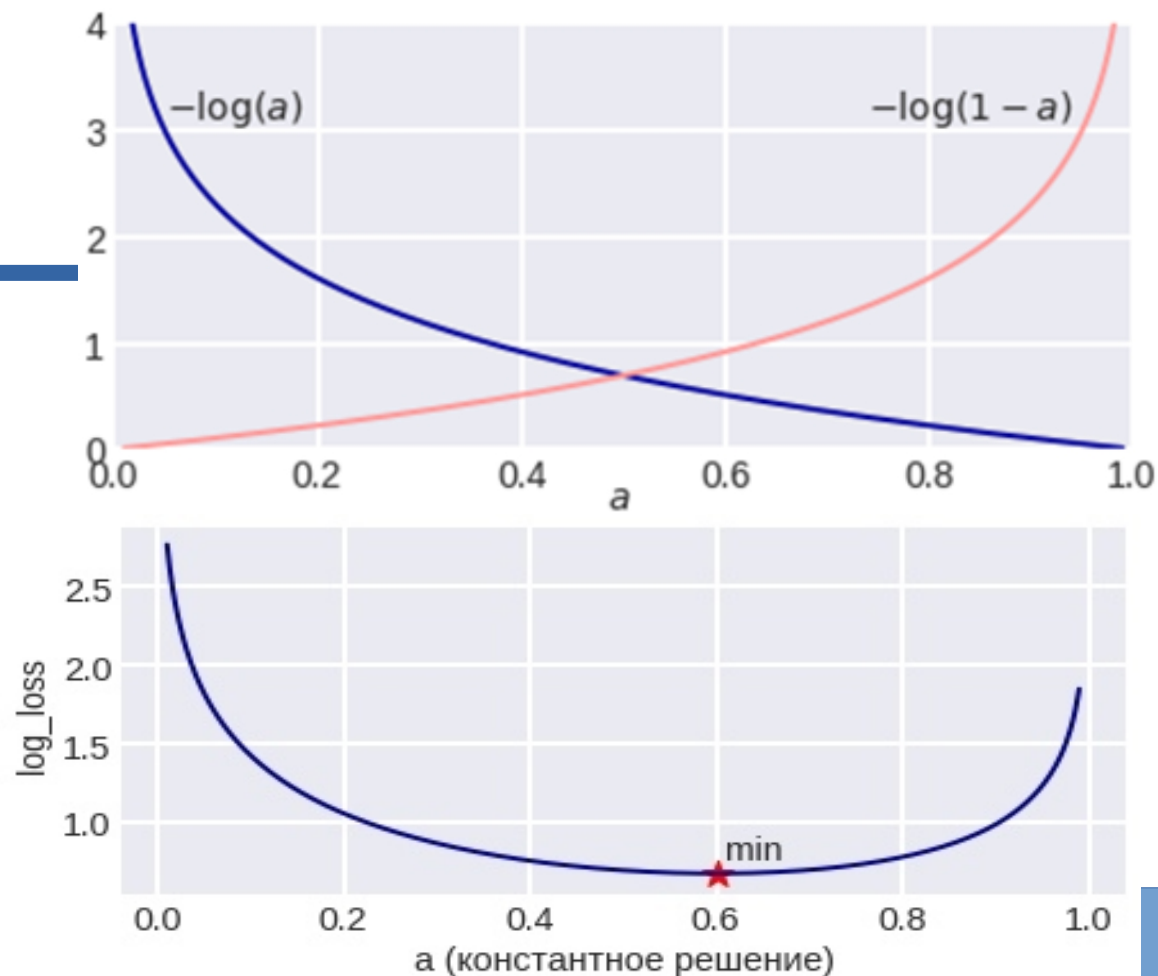
## Метод максимального правдоподобия

$$p = y$$

$$-p \log(a_i) - (1 - p) \log(1 - a_i)$$

$$\frac{p}{a_i} - \frac{1 - p}{1 - a_i} = 0$$

$$a_i = p$$



# Log Loss

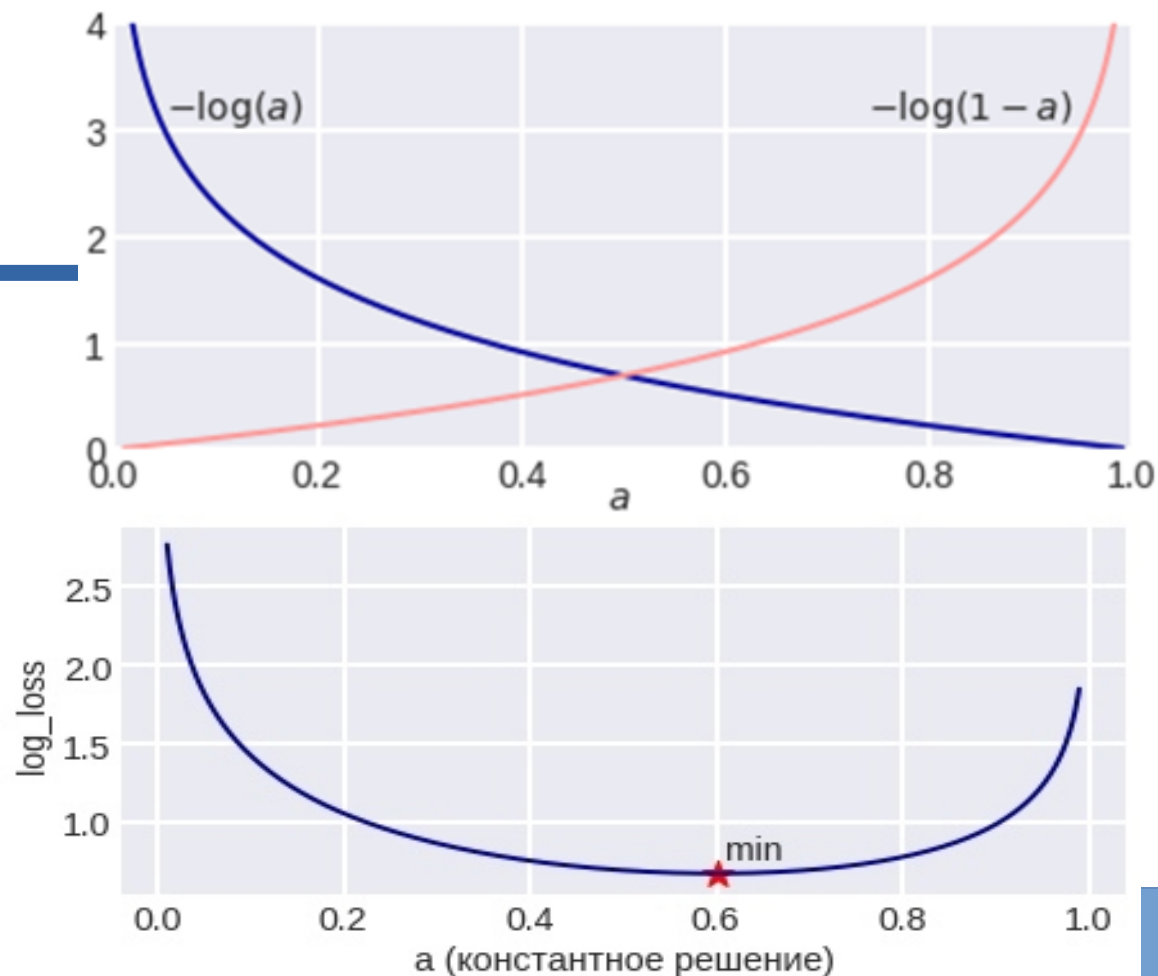
## Метод максимального правдоподобия

$$a = p$$

Энтропия

$$-p \log(p) - (1 - p) \log(1 - p).$$

$$H(y, a) = - \sum_x y(x) \log a(x, \omega)$$

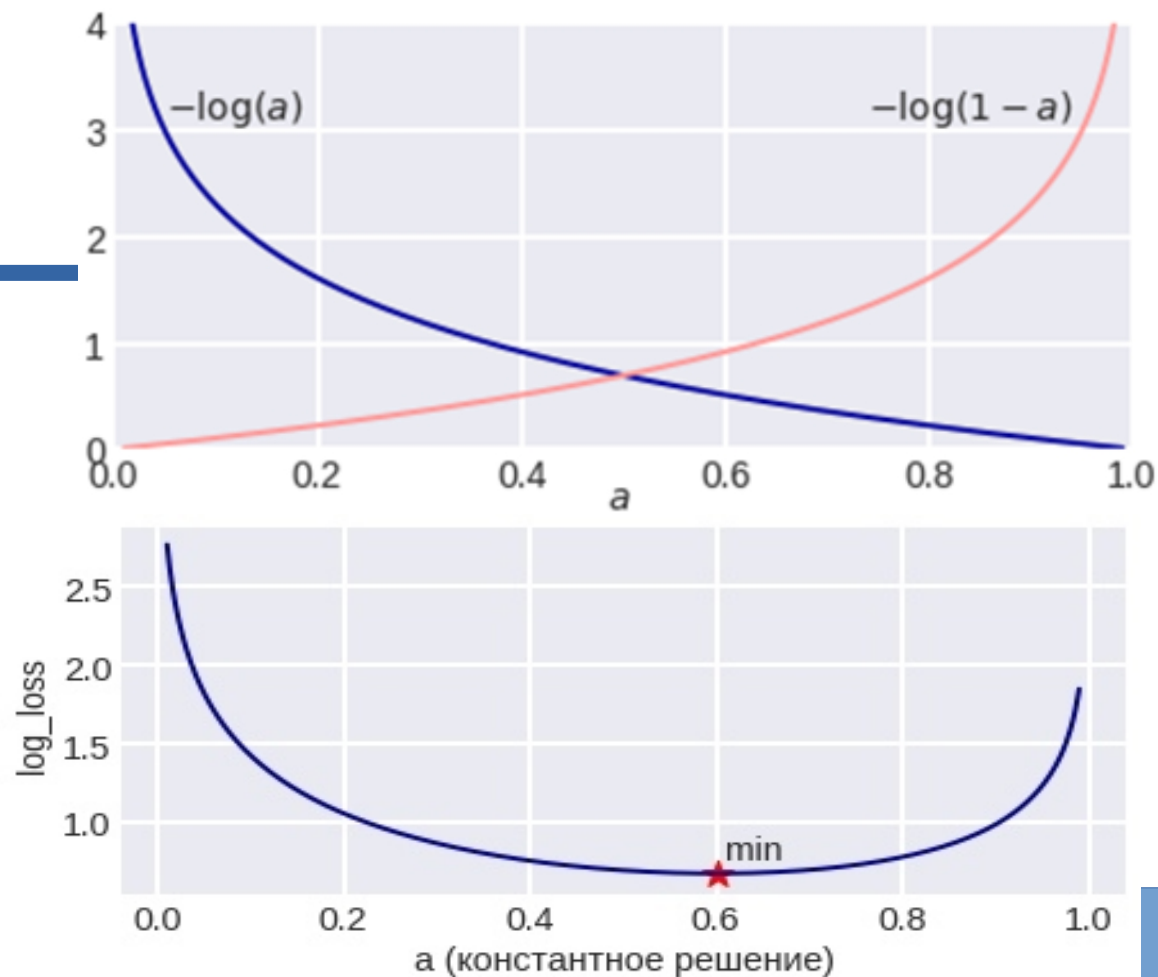




# Log Loss

## Метод максимального правдоподобия

$$\text{logloss} = -\frac{1}{q} \sum_{i=1}^q \sum_{j=1}^l y_{ij} \log a_{ij}$$



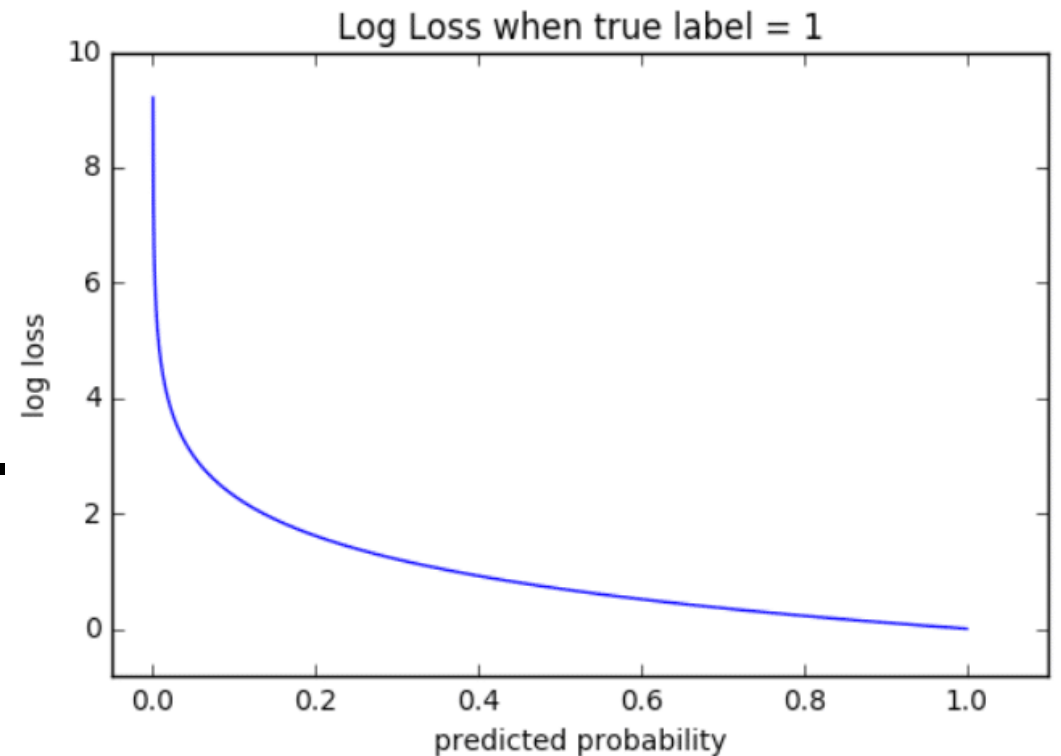
# Градиентный спуск

**LogLoss - функция ошибок**

**Функцию ошибок надо минимизировать**

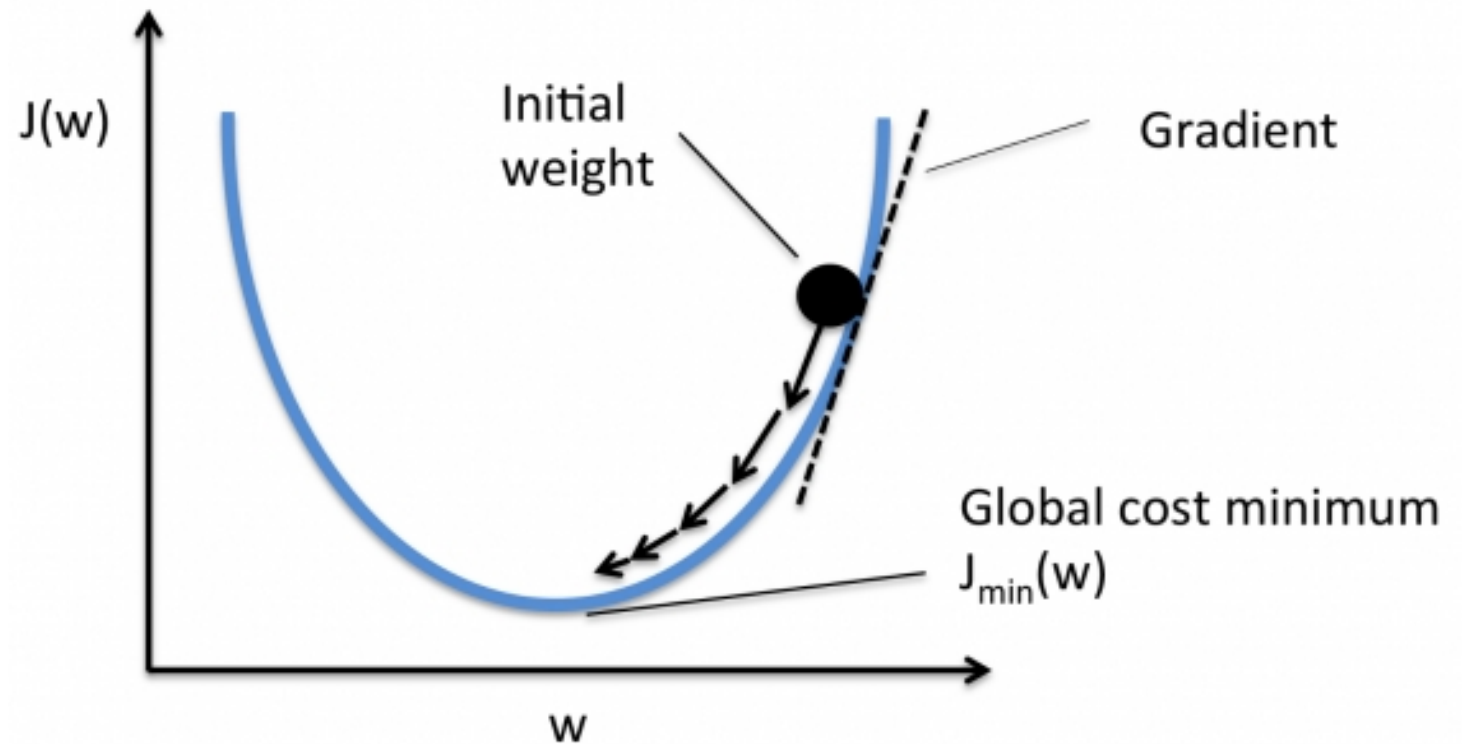
$$\omega_{n+1} = \omega_n - \eta \frac{1}{I} X(A - Y)^T$$

$$A = \frac{1}{1 + \exp(-\langle \omega, x_2 \rangle)}$$



# Градиентный спуск

Градиент (производная) показывает направление роста функции.



# Оптимизация логистической регрессии

Функция потерь:  $\text{LogLoss} = \sum_i -y_i \log a_i - (1 - y_i) \log(1 - a_i)$

Логистическая регрессия:  $a = \frac{1}{1 + e^{-w^T x}}$

Градиент:  $\frac{\partial \text{LogLoss}}{\partial w} = (a - y)x$

# Метрики классификации

	Истина +	Истина -
Предсказано +	True positive	False positive
Предсказано -	False negative	True negative

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = Sensitivity = \frac{TP}{TP + FP}$$

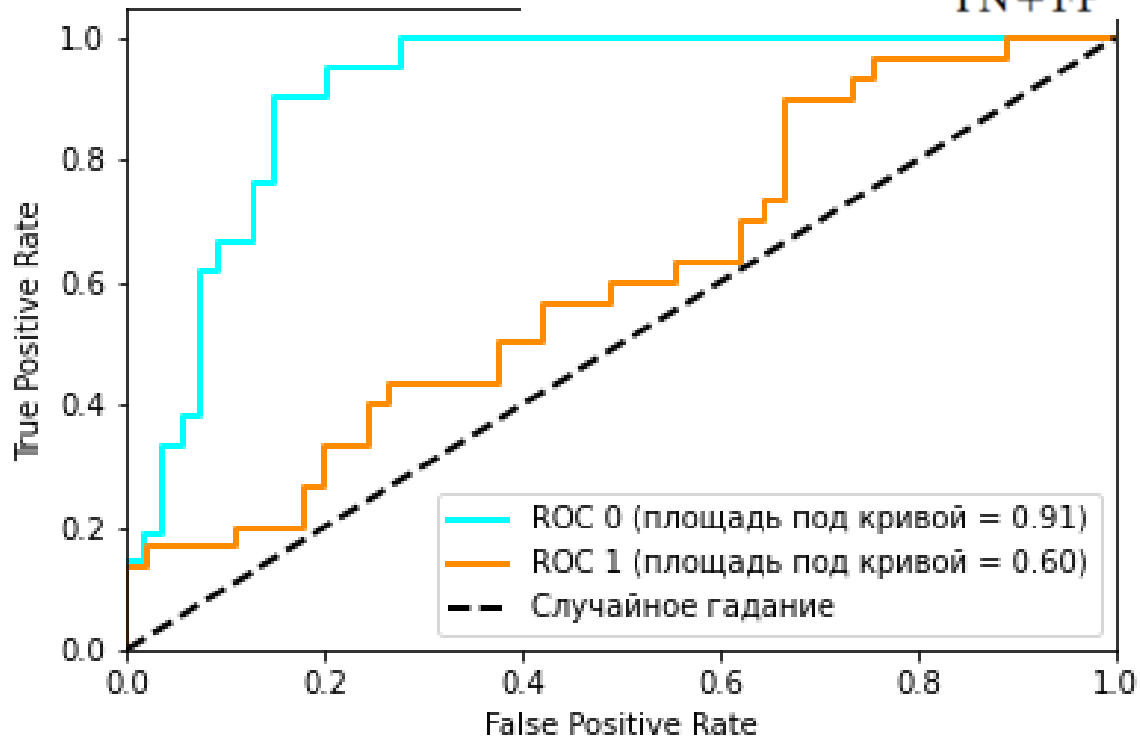
$$Recall = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$

# Метрики классификации

$$\text{TPR (sensitivity)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{FPR (1-specificity)} = \frac{\text{FP}}{\text{TN} + \text{FP}}$$



ROC кривая - зависимость верно классифицируемых объектов положительного класса (Sensitivity) от ложноположительно классифицируемых объектов негативного класса (Specificity)

AUC (Area Under Curve) - площадь под кривой

ROC-AUC - площадь под ROC кривой, численная оценка ROC метрики

# Логистическая регрессия

- Линейный алгоритм классификации оказывается оптимальным байесовским классификатором.
- Однозначно определяется вид функции активации (сигмоидная функция) и функции потерь.
- Делает численные оценки вероятности его принадлежности каждому из классов.
- Является частным случаем обобщённой линейной модели регрессии.

# Логистическая регрессия

- ❑ Оценки вероятностей и рисков могут оказаться неадекватными, если не признаки зависимы
- ❑ Все недостатки метода стохастического градиента.
- ❑ Практичная реализация должна предусматривать :
  - стандартизацию данных,
  - отсев выбросов,
  - регуляризацию,
  - отбор признаков.



# **Логистическая регрессия**

**ВОПРОСЫ?**

**Спасибо!**

**Каждый день  
вы становитесь  
лучше :)**

