



# Машинное обучение в науках о Земле

Михаил Криницкий

К.Т.Н., С.Н.С.

Институт океанологии РАН им. П.П. Ширшова

Лаборатория взаимодействия океана и атмосферы и  
мониторинга климатических изменений (ЛВОАМКИ)



# Задачи классификации

Михаил Криницкий

К.Т.Н., С.Н.С.

Институт океанологии РАН им. П.П. Ширшова

Лаборатория взаимодействия океана и атмосферы и  
мониторинга климатических изменений (ЛВОАМКИ)

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

## 1. формулировка задачи:

- какой тип (классификация, регрессия, другой)? Или переформулировать в легко решаемый тип!
- определиться, что есть объекты (события)
- определиться, что есть целевая переменная
- определить признаковое описание объектов (событий)
- определить критерии качества решения задачи (MSE, MAE, pattern correlation, etc.)

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

## 2. формулировка модели:

- задать вид модели (линейная регрессия, дерево решений, композиционный алгоритм, нейронная сеть, etc.)
- задать сложность модели (задается гиперпараметрами – настройками модели)
- определиться с функцией потерь (MSE, MAE, LogLikelihood, etc.)

## 3. подготовка данных или генератора данных:

- стандартизировать данные (если нужно)
- обработать пропуски, категориальные значения, подготовить кодирование текста, применить понижение размерности данных
- оставить часть данных для проверки качества (train-test split)
- подготовить генератор данных с учетом стратегии скользящего контроля (cross-validation quality estimation)

## 4. оптимизация модели на обучающей выборке:

- $\hat{p} = \underset{\mathbb{P}}{\operatorname{argmin}}(L(\vec{p}, \mathcal{T}))$

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

5. оценка модели:
  - оценить качество по метрикам, определенным на этапе **1.** на тестовой выборке
  - оценить неопределенности параметров модели (если возможно)
  - оценить неопределенности оценок целевой переменной
  - определить наличие недообучения или переобучения
  - оценить соотношение сложности модели и сложности закономерностей в данных; при неадекватной сложности модели вернуться к **п.2**
6. применение модели на вновь получаемых данных:
  - оценка распределения вновь получаемых данных: генерируются ли они из того же распределения, что и обучающая выборка?
  - предобработка новых данных идентично **п.3** с точностью до коэффициентов стандартизации и деталей способов предобработки
  - применение модели к предобработанным новым данным для получения значений целевой переменной
  - построение научных выводов

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

7. публикация результатов:

- описание результатов в виде статей, отчетов об исследованиях
- защита результатов перед лицом научного сообщества
- получение наград, признание успехов, etc.

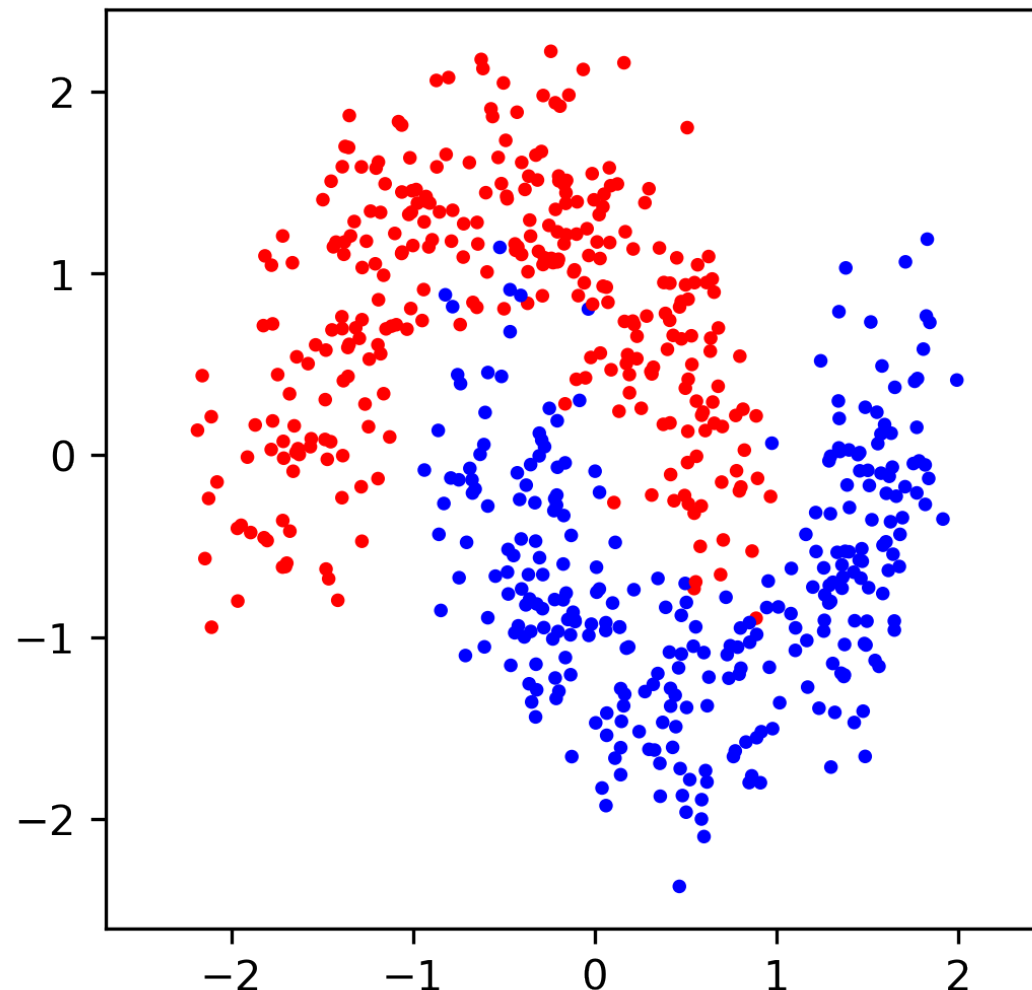
# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

## Формулировка задачи (в терминах машинного обучения)

### ○ «Обучение с учителем»

- восстановление регрессии
- классификация

**что я хочу?** – метку класса  
**«красный или синий?»**  
(бинарная классификация)



# ЗАДАЧА КЛАССИФИКАЦИИ

**цель** — метка класса ( $y$  — категориальная переменная)

«спам / не-спам»

«мезоциклон / не-мезоциклон»

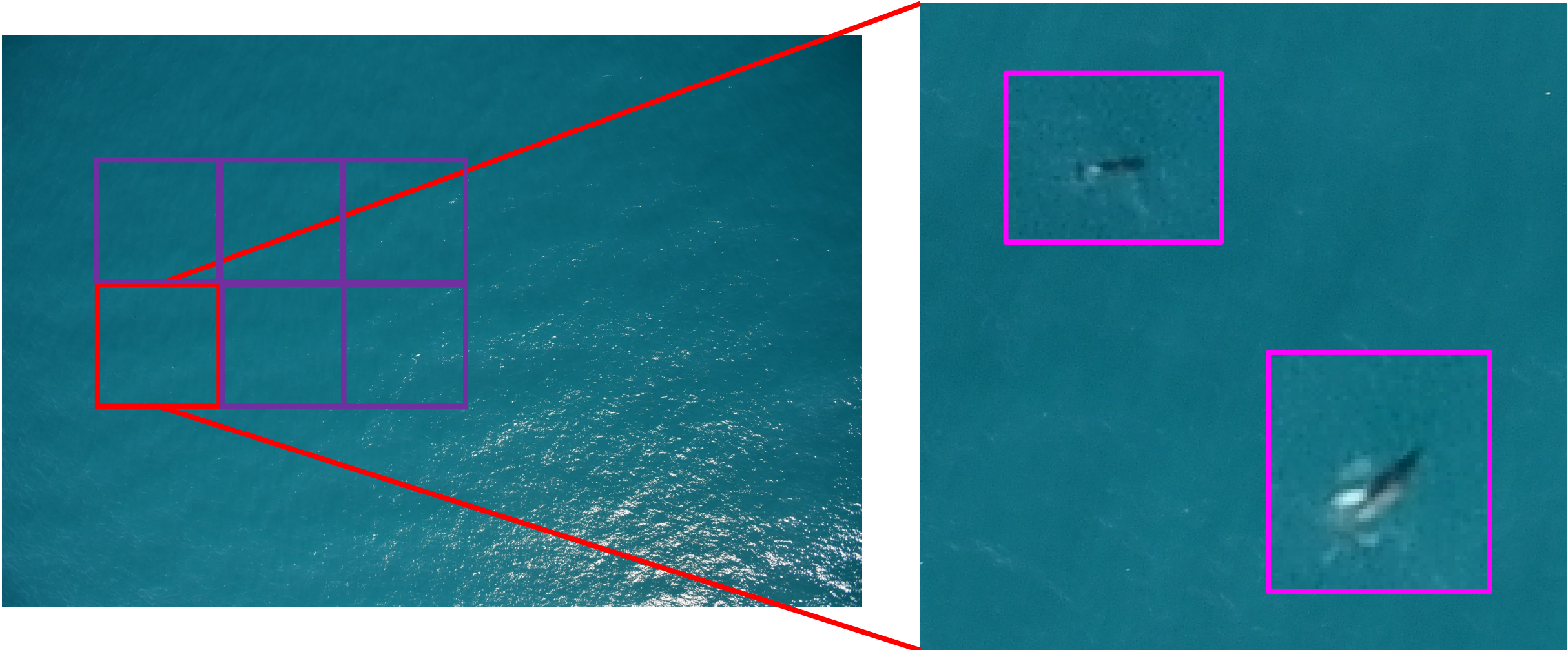
«кот / собака / лошадь»

«0 / 1 / 2 / 3 / 4 / 5 / 6 / 7 / 8 / 9»



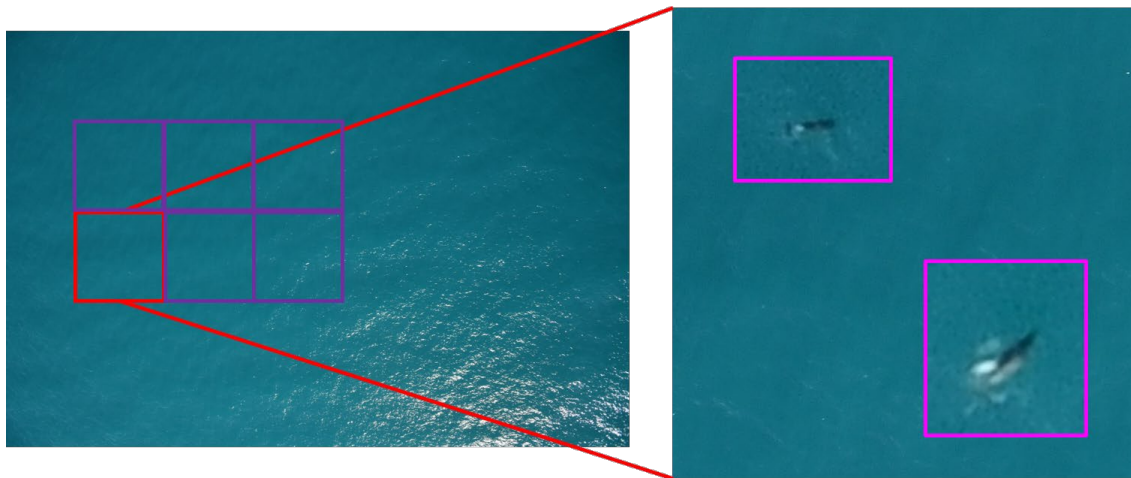
# ЗАДАЧА КЛАССИФИКАЦИИ

**Пример:** классификация участков учетных снимков поверхности океана в отношении признака наличия морских млекопитающих



# ЗАДАЧА КЛАССИФИКАЦИИ

**Пример:** классификация участков учетных снимков поверхности океана в отношении признака наличия морских млекопитающих



- Что есть **объекты** выборки?
- Что есть **признаковое описание**?
- Что есть **целевая переменная**?
- Тип целевой переменной?
- Размерность ц.п.?
- **Мера качества** решения?

# ЗАДАЧА КЛАССИФИКАЦИИ

Простейший пример:

объекты описываются действительным признаком  $x$   
целевая переменная  $y$  – бинарная, классы:  $A$ ,  $B$ ; по 1000 экземпляров каждого класса  
пусть для класса  $y = A$  значения  $x \sim \mathcal{N}(\mu_A, \sigma_A)$ , для класса  $y = B$  значения  $x \sim \mathcal{N}(\mu_B, \sigma_B)$

Базируясь на этих данных, каково должно быть  
решение (значение  $y$ ) при:

$$x = -10$$

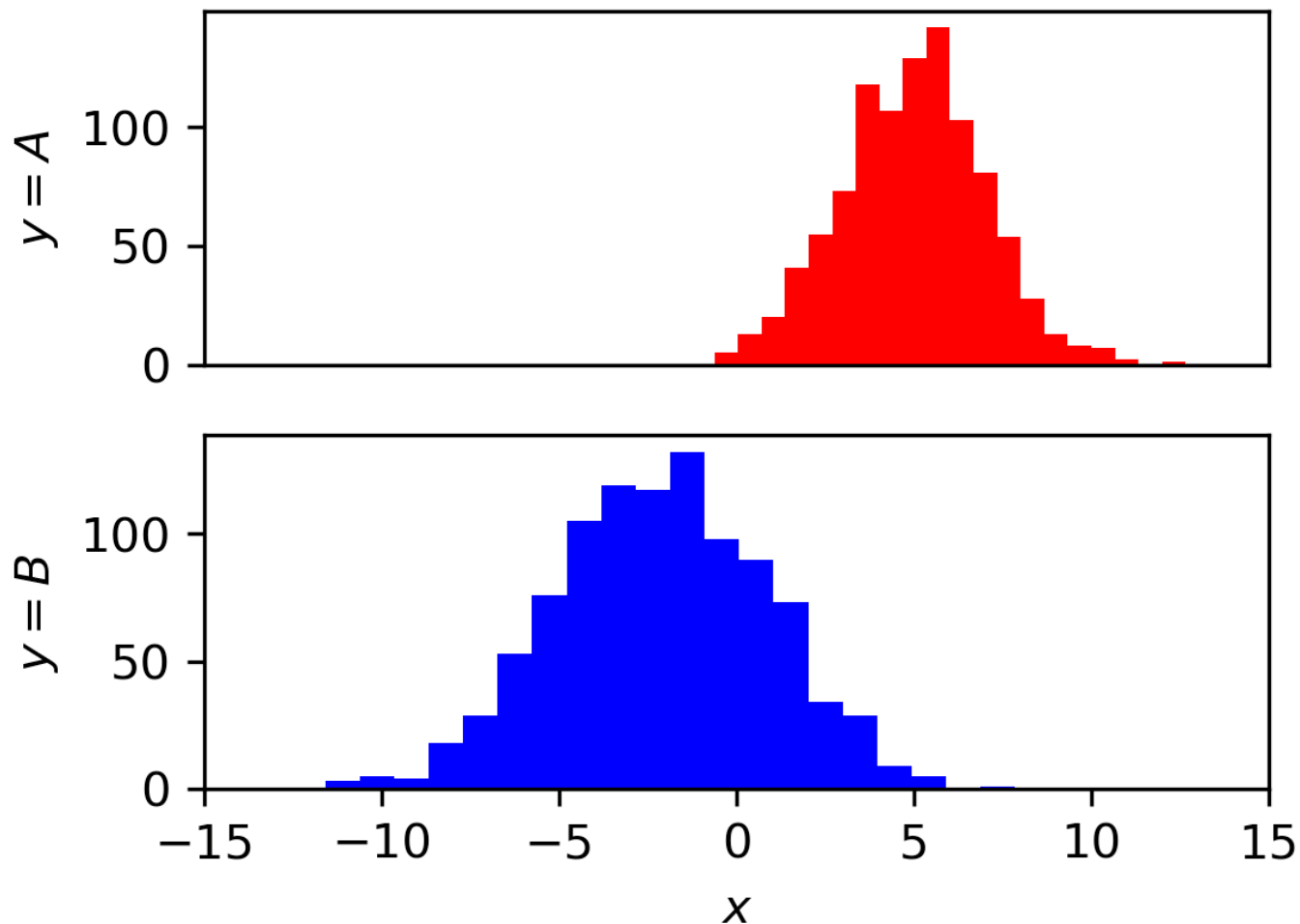
$$x = -5$$

$$x = 2$$

$$x = 5$$

$$x = 10$$

$$x = 15$$

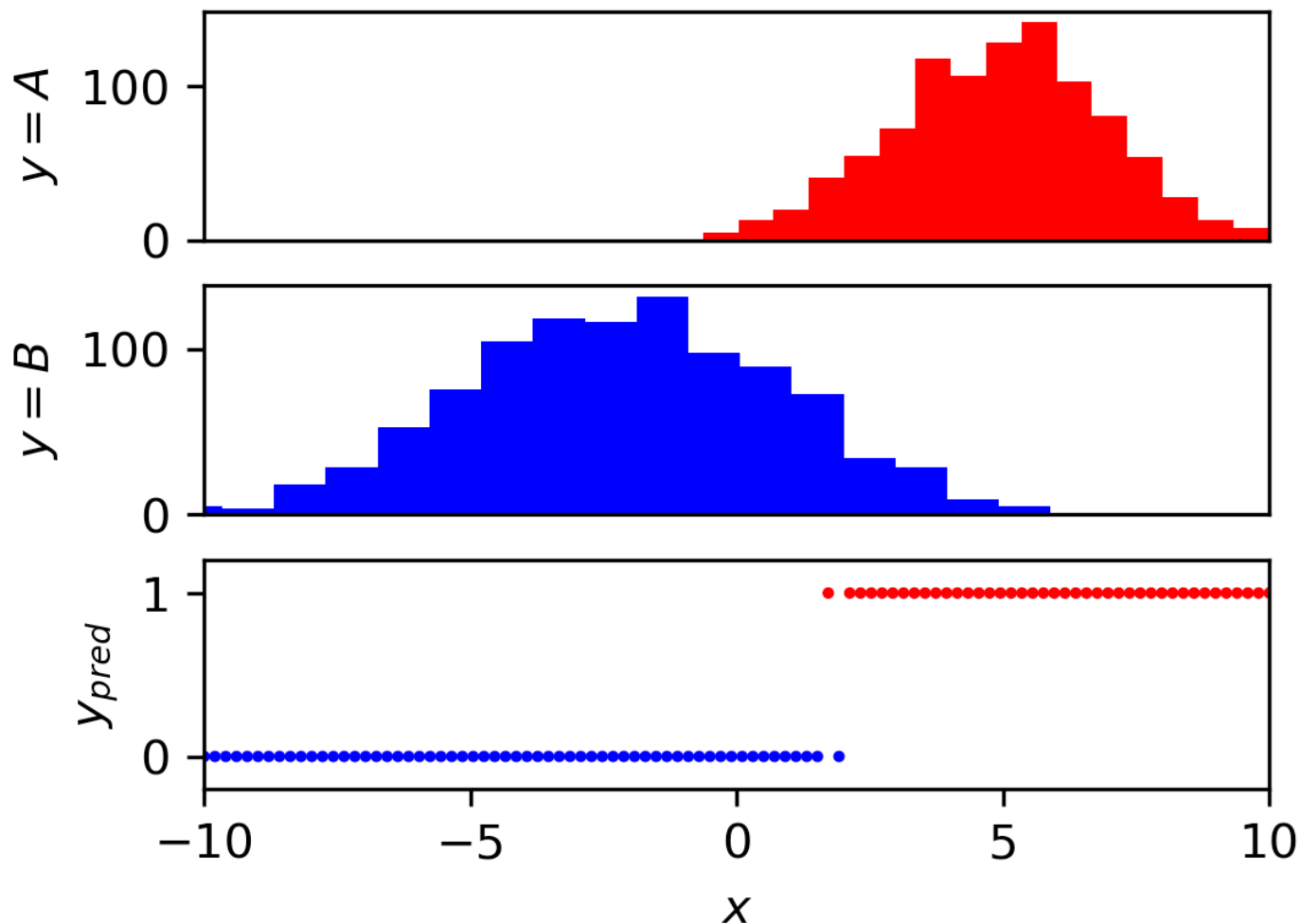


# ЗАДАЧА КЛАССИФИКАЦИИ

Простейший пример: объекты описываются действительным признаком  $x$   
целевая переменная  $y$  – бинарная, классы:  $A$ ,  $B$ ; по 1000 экземпляров каждого класса  
пусть для класса  $y = A$  значения  $x \sim \mathcal{N}(\mu_A, \sigma_A)$ , для класса  $y = B$  значения  $x \sim \mathcal{N}(\mu_B, \sigma_B)$

Подход №1: **KNN** (метод  $K$  ближайших соседей)

1. выбрать  $K$  ближайших соседей для нового объекта (! нужно определить меру близости !)
2. осреднить (можно с разными весами) целевую переменную по этим объектам («простое голосование», «majority vote» или «взвешенное голосование», «weighted vote»)
3. считать полученный результат значением целевой переменной на новом объекте



# ЗАДАЧА КЛАССИФИКАЦИИ

Подход №1: **KNN** (метод  $K$  ближайших соседей)

- простой
- быстрый
- легко настраивается. Гиперпараметр  $K$  регулирует «сложность» модели

## А ЧТО ЕСЛИ ДАННЫХ МАЛО?..

- требуется большое количество обучающих данных
- обучающие данные должны быть распределены достаточно плотно в исследуемой области  $x$
- не обобщает закономерности в данных

