

# Машинное обучение 1 2017

Евгений Елтышев

# ОСНОВНЫЕ ПОНЯТИЯ

- Объект (sample)
- Признак (feature)
- Метка (label, ground truth)
- Выборка (dataset)

# ОСНОВНЫЕ ПОНЯТИЯ

- Обучение с учителем (Supervised learning)
- Обучение без учителя (Unsupervised learning)
- Классификация (  $\text{labels.nunique()} < \infty$  )
- Регрессия ( $\text{labels} \subset \mathbb{R}$ )
- Кластеризация

# Признаки

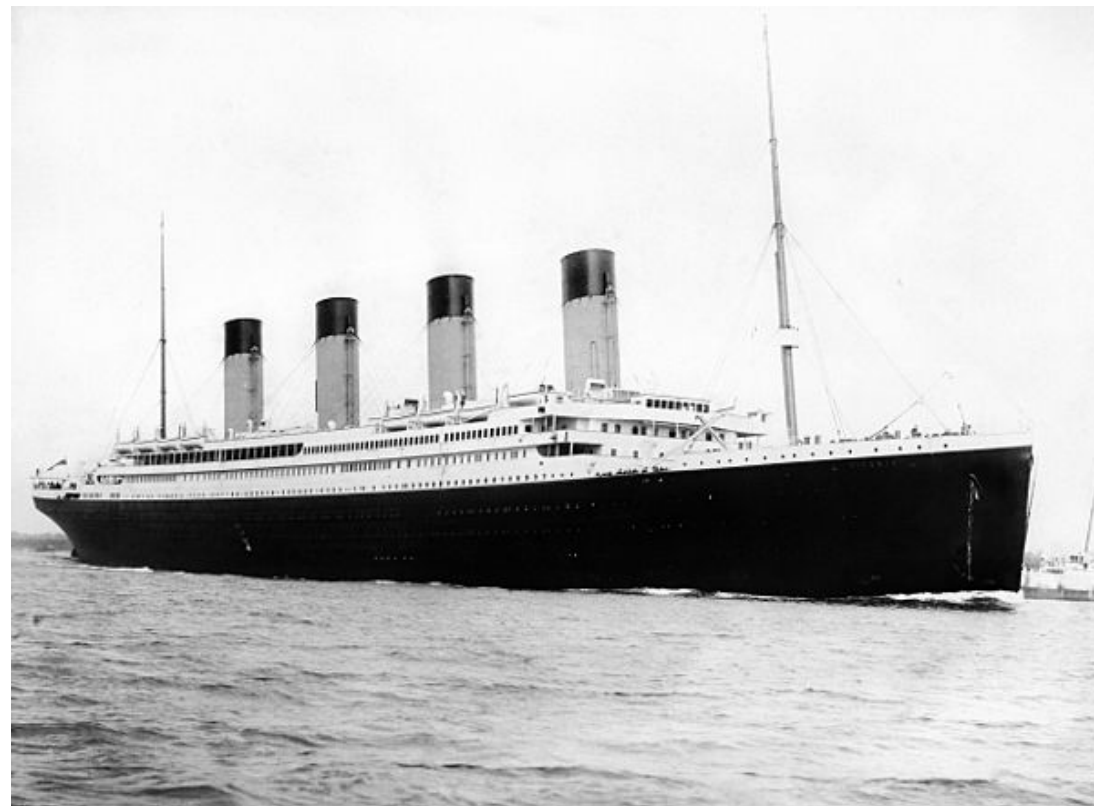
- Категориальные
  - Марка автомобиля
  - Регион пользователя
- Вещественные
  - Пробег автомобиля
  - Возраст пользователя

# Задача: Титаник

Задача: предсказать выжил ли человек?

Погибших	1496 <sup><a href="#">[1]</a></sup>
Пострадавших	712

Какие признаки?



# Оценка качества

# Оценка качества

- Обучение с учителем
  - Сравниваем наши ответы с истинными метками
- Обучение без учителя
  - Менее формальные метрики, например минимальное расстояние между кластерами

# Оценка качества классификации

- Accuracy
- Logloss



# Оценка качества регрессии

- Mean squared error

Окей, а на каких данных измерять  
качество?

# Окей, а на каких данных измерять качество?

- На обучающей выборке
  - Переобучение!
- На отложенной выборке
  - Уже лучше, но можем переобучиться под нее
- K-fold кросс-валидация
  - Совсем хорошо, но долго

# Линейные модели

- Регрессия:  $\hat{y} = \vec{w} \cdot \vec{x}$
- Классификация:  $\hat{y} = \text{sign}(\vec{w} \cdot \vec{x})$
- Легко интерпретировать
- Легко обучать
- Требуется мало данных
- Не могут моделировать сложные зависимости

# Метод ближайших соседей (kNN)

- $\hat{y} = \frac{1}{k} \sum_{i=1}^k y_{(i)}$
- Просто устроен
- Не требует долгого обучения
- При разных k работает «по-разному»
- Плохо работает в больших размерностях