

Детектирование Data Drift

Цель темы

Научиться пользоваться основными методами для детектирования Data Drift.

Задачи темы

- ✓ Познакомиться с ограничениями в контроле данных
- ✓ Узнать об основных методах детектирования Data Drift

Ограничения очевидных методов контроля

Пример: задача кредитного скоринга.

Методы детектирования

- 1 Статистические методы
- 2 Мониторинг предсказаний модели
- 3 Мониторингов признаков (фичей)

Статистические методы

- 1 Наблюдения за описательными статистиками
- 2 Статистические тесты

Статистические тесты

- 1 Тест Колмогорова — Смирнова для непрерывных признаков
- 2 Кси-квадрат — тест для категориальных признаков
- 3 KPSS-тест для проверки стационарности временных рядов

Пример статистического теста

```
from scipy import stats
```

```
...
```

```
def detect_drift(original, new, alpha=0.05):
```

```
    ks_statistic, p_value = stats.ks_2samp(original, new)
```

```
    print(f"Статистика KS: {ks_statistic:.4f}")
```

```
    print(f"p-значение: {p_value:.4f}")
```

```
    if p_value < alpha:
```

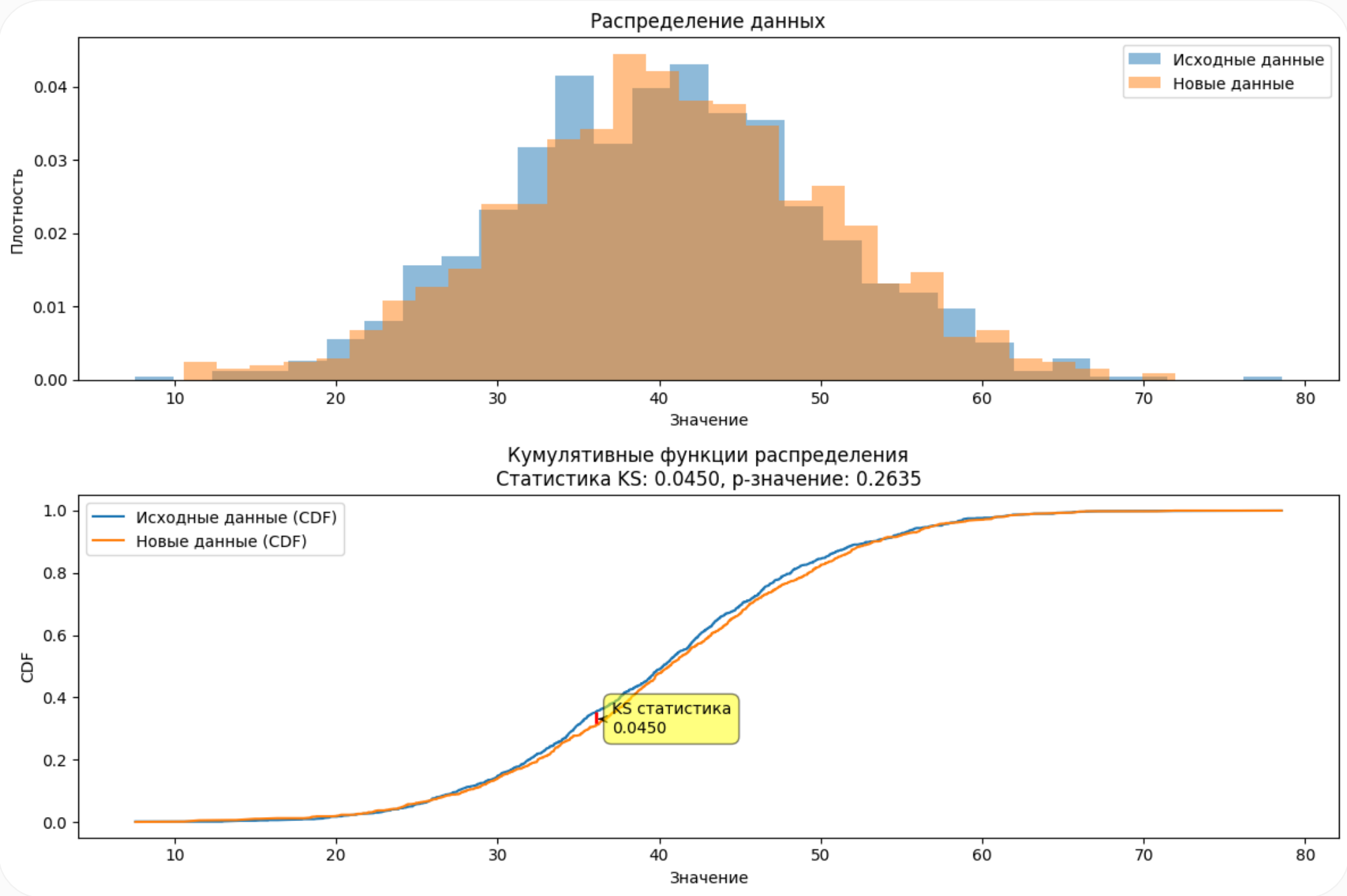
```
        print("Обнаружен значительный дрейф данных.")
```

```
    else:
```

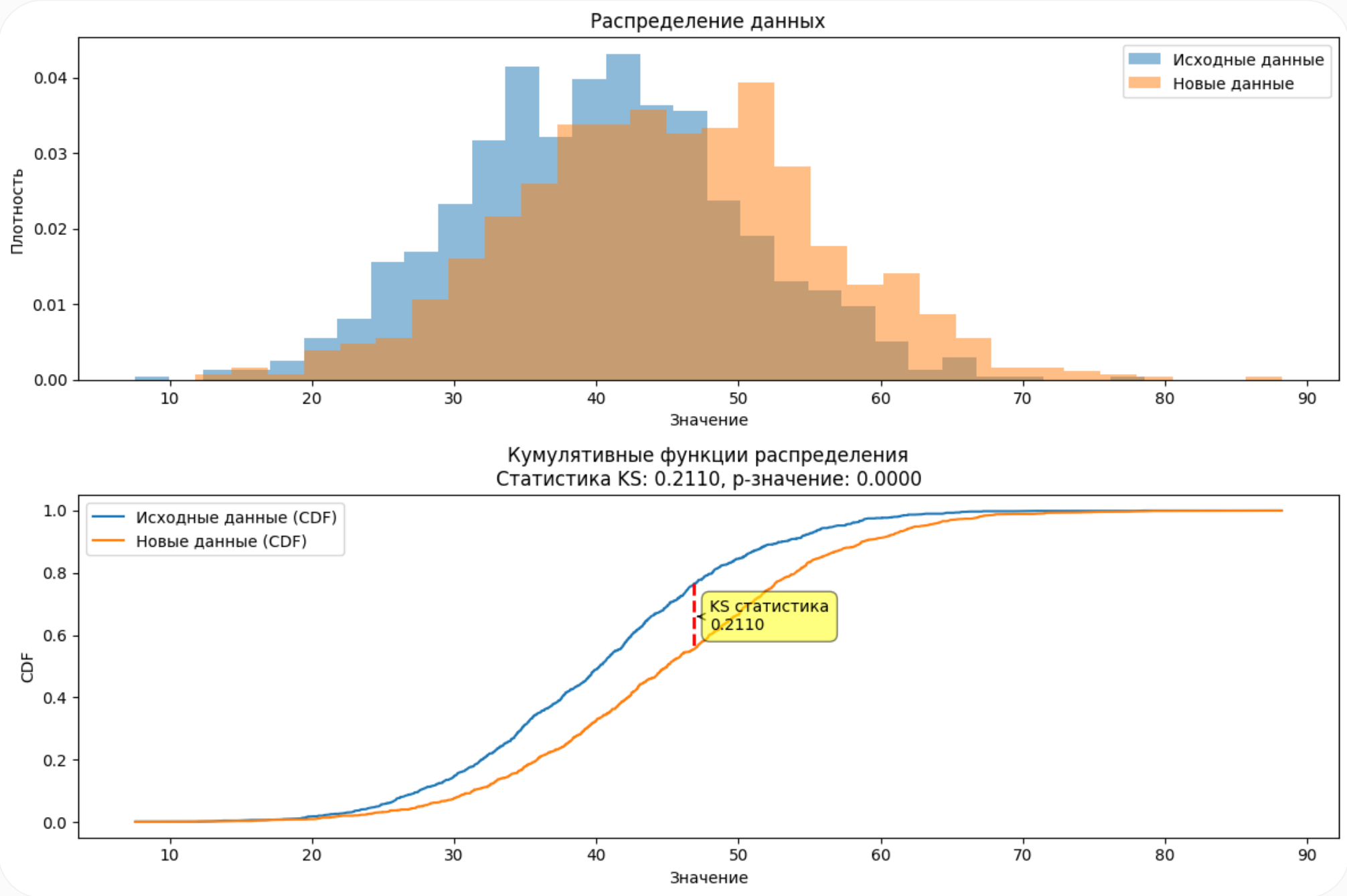
```
        print("Значительного дрейфа данных не обнаружено.")
```

```
...
```


Пример статистического теста



Пример статистического теста



Мониторинг предсказаний модели

Плюсы

- Небольшая размерность
- Легко интерпретировать

Минусы

- Запаздывающие алерты
- Не помогает понять причину

Мониторинг признаков (фичей)

Плюсы

- Отдельные признаки
- Многомерный анализ

Минусы

- Дрифт можно обнаружить заранее
- Требуется тщательного выбора метрик

Вывод темы

Научились пользоваться различными методами для детектирования Data Drift.

Вывод модуля

Научились отслеживать деградацию ML-моделей
после их внедрения.