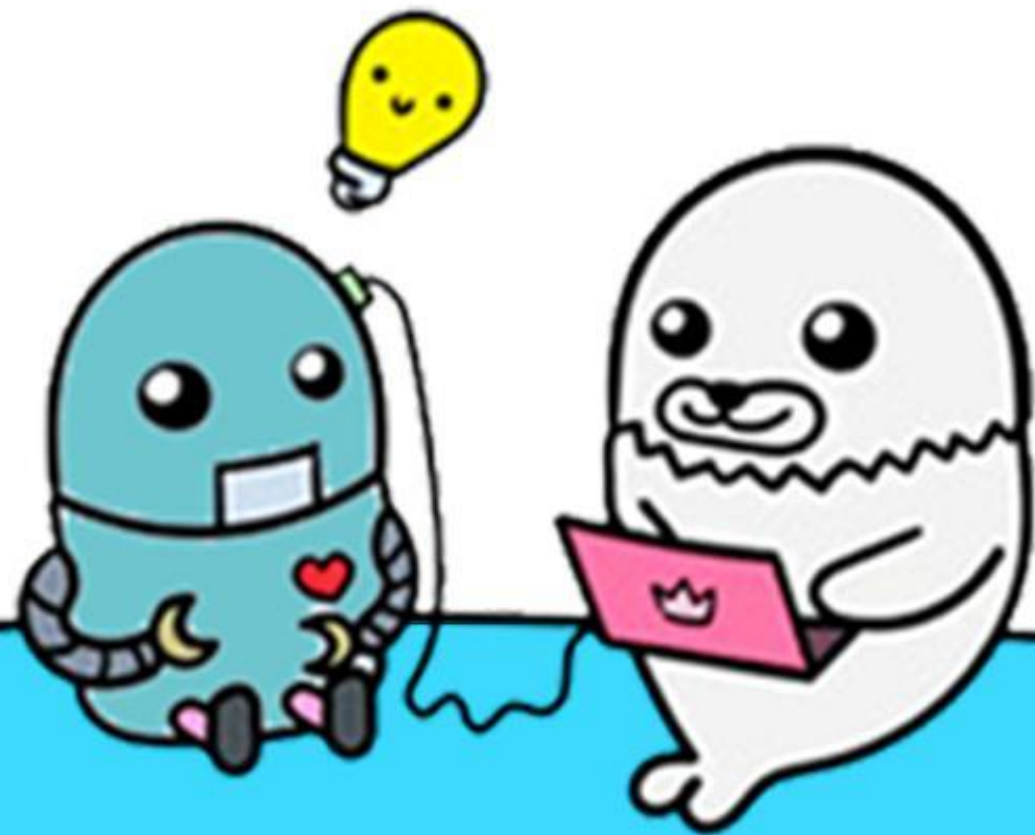


Яндекс Кью



# MLOps и production подход к ML исследованиям



28 марта - 28 мая



# MLOps и production подход к ML исследованиям

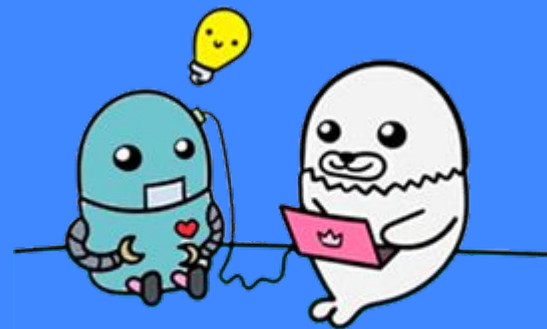
*Концепция воспроизводимых и масштабируемых исследований в ML*

**Павел Кикин**

Газпромнефть ЦР

Руководитель направления NLP

[t.me/pavel\\_kikin](https://t.me/pavel_kikin)





## Для чего этот курс

Всем привет 🙌 Я так понимаю, этот курс не для новичков. Н...

Как уже выше написал Дима, но я все же поясню, этот курс о том как правильно сделать окружение вокруг проекта, что бы перед ровными мльщиками не позориться



8:18

- Работать в команде над одним проектом
- Превратить модель в сервис
- Корректно передать результаты исследований заказчику или другим разработчикам
- Систематизировать исследования
- Отслеживать и сохранять условия и результаты экспериментов
- Автоматизировать эксперименты
- Повысить качество кода
- Выстроить с нуля процессы МЛ разработки в команде
- **Обеспечить воспроизводимость ваших исследований**



## Для чего этот курс

---

- Написать больше умных слов в резюме
- Казаться умнее на собеседованиях





**О курсе**

---



**9 занятий  
теория +  
практика**



**Понедельник  
18:30**



**Рейтинг  
участников**



**Конференции в  
Zoom**



**Онлайн  
трансляции и  
запись в Youtube**



## О курсе

---



### Домашние задания (не оценивается)

*Google Drive*



### Тесты (10 баллов)

*Страничка курса в ODS*



### Дополнительные задания (20 баллов)

- *руководство*



### Вопросы/ответы/договор (5 баллов)



### Итоговый проект (50 баллов)



## О курсе

---



### Примерные темы для докладов:

- Опыт реализации проекта Findmybike.ru
- Исследование сервисов управления ноутбуками (sagemaker, databricks, datalore, ванильный JN)
- Amazon S3 и S3-like хранилища
- Docker
- СУБД
- Особенности ООП в МЛ
- Применение Agile в ML командах
- w&b, neptuna, clearml
- Pachyderm
- Dataflow
- Kubernetes and KubeFlow
- TF serving
- Контейнеризация рабочей среды для DS
- Настройка CUDA для Docker



## Процесс работы над ML проектом

- Бизнес анализ, инициирование работ
- Обзор и анализ существующих решений – бейзлайнов
- Исследование/собственная разработка
- Подготовка отчета



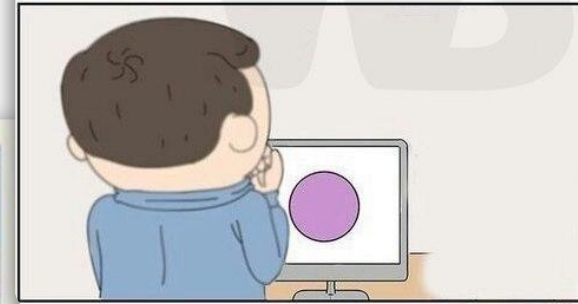


## Бизнес анализ, инициирование работ

- Понять бизнес-цели (помочь заказчику их сформулировать)
- Очертить возможные подходы и решения
- Обозначить технологические пределы возможностей
- Продумать риски
- Повторно продумать и скорректировать бизнес-цели

## Постановка бизнес - задачи

### ТИПИЧНЫЙ ЗАКАЗЧИК



How the customer explained it



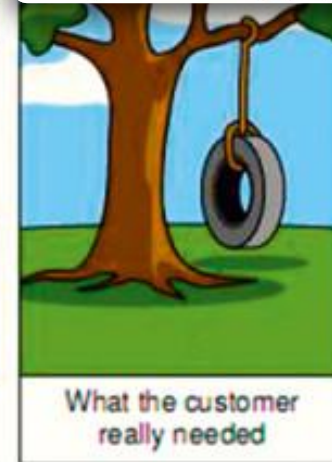
How the engineer designed it



How the project leader understood it



How the programmer wrote it



What the customer really needed



## Бизнес анализ, инициирование работ

### Оценка данных

- Сколько объектов выборки мы имеем?
- Возможно ли ознакомиться со всеми данными до старта проекта?
- Есть ли разметка и какого она качества?
- Есть ли возможность разметить все классы?
- Сколько классов?
- Есть ли дисбаланс?
- Что с выбросами и пропусками?





## Бизнес анализ, инициирование работ

## Оценка заказчика и проекта

- Узнайте заказчика поближе.
- Изучите опыт работы других исполнителей с ним.
- Назначьте встречу со всеми заинтересованными сторонами от заказчика, а не только с одним из представителей.



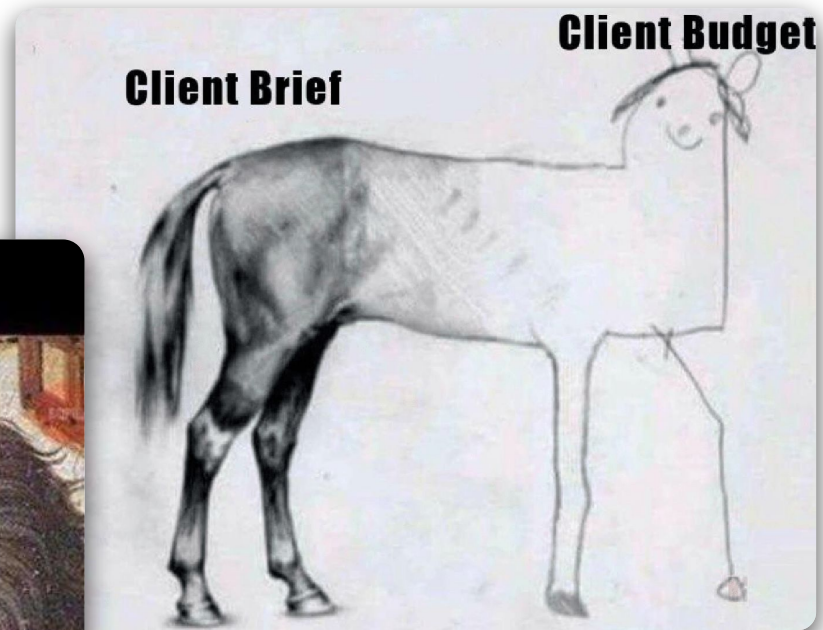




## Бизнес анализ, инициирование работ

## Оценка заказчика и проекта

- Деньги будут?





## Бизнес анализ, инициирование работ

### Подбор метрик

- Выберите метрики, которые будут наилучшим образом согласовываться с бизнес-целями
- Объясните их заказчику
- Убедитесь, что заказчик их понял. Еще раз проверьте, что он их **ТОЧНО** понял
- Если метрик много, нужно выбрать набор, который будет использоваться для валидации решения/ий



## Бизнес анализ, инициирование работ

### Финальный этап

- Четко сформулируйте цели исследования
- Оцените ресурсы с учетом рисков (перезакладывайтесь):
  - Особенности заказчика
  - Неопределенность в требованиях
  - Неопределенность в данных





## Бизнес анализ, инициирование работ

### Финальный этап

- Сформируйте ТЗ
  - Можно ли по этому ТЗ четко сказать, что проект выполнен в полной мере или нет?
  - Точные метрики
  - Критерии приемки







## Обзор, анализ и реализация существующих решений – бейзлайнов

- Формирование критериев оценки:
  - Метрики
  - Производительность
  - Работа под нагрузкой
  - Поддержка
  - Готовность
  - Популярность
- Поиск существующих решений (не беритесь за 1 попавшееся):
  - Популярные сайты для разработчиков
  - Научные источники
  - ODS, TG, Кью, профильные форумы
- Реализация подходящих решений
- Составление валидационного набора данных
- Разработка валидационного дайплайна (башмарка)



*«Не изобрети велосипеда»*





## Подготовка отчета

---

- Постановка задачи
- Описание исходных данных
- Обзор и анализ существующих решений
- Методика, технология и критерии оценки существующих решений
- Результаты оценки существующих решений и выводы

Вариант 1 – мы нашли подходящее решение

- Детальное описание выбранного решения и его адаптации под нашу задачу
- Детали предобработки исходных данных

Вариант 2 – пилим своё

- Гипотезы
- Результаты проверки гипотез
- Заключение:
  - Выводы



# Концепция воспроизводимости



## Воспроизводимость

Мера вероятности того, что, получив один результат эксперимента, вы сможете провести тот же эксперимент с теми же параметрами и получить точно такой же результат. Это способ убедиться, что результаты верны и не являются случайностями.





## Виды воспроизводимости

---

- **Repeatability:** та же команда, **те же условия эксперимента**. Возможность получения заявленных в результатов на тех же входных данных.
- **Reproducibility:** другая команда, **те же условия эксперимента**. Если наблюдение воспроизводимо, оно должно быть выполнено другой командой, повторяющей эксперимент с использованием тех же экспериментальных данных и методов, в тех же рабочих условиях, в том же или другом месте, в нескольких испытаниях.
- **Replicability:** другая команда, **другие условия эксперимента**. Если наблюдение можно воспроизвести, оно должно быть выполнено другой командой, с использованием другой измерительной системы и набора данных, в другом месте, в нескольких испытаниях. Следовательно, это потребует нового сбора данных.



## Одинаковые условия эксперимента?

---

В целом в науке:

- Место нахождения
- Измерительные инструменты
- Другое оборудование, использованное в эксперименте
- Наблюдатель
- Гипотеза
- Период времени

Для ML:

- Аппаратное обеспечение
- Программное обеспечение и его версии
- Способ/оборудование для получения данных (датчики климата, спутниковые оборудование и т.д.)
- Время получения данных (особенно важно для часто обновляемых источников, новостные ленты, социальные сети и т.д.)



## Кризис воспроизводимости

---

Крупномасштабные усилия по оценке воспроизводимости научных публикаций дали тревожные результаты. Например, в 2015 году группа исследователей психологии, получившая название «Открытое научное сотрудничество», рассмотрела 100 экспериментов, опубликованных в высокорейтинговых рецензируемых журналах. Из этих 100 исследований только 68 репродукций дали статистически значимые результаты, совпадающие с исходными данными.

Согласно данным другого анализа, до 85 % всех проведённых в мире исследовательских работ в области биомедицины не привели к значимым результатам.

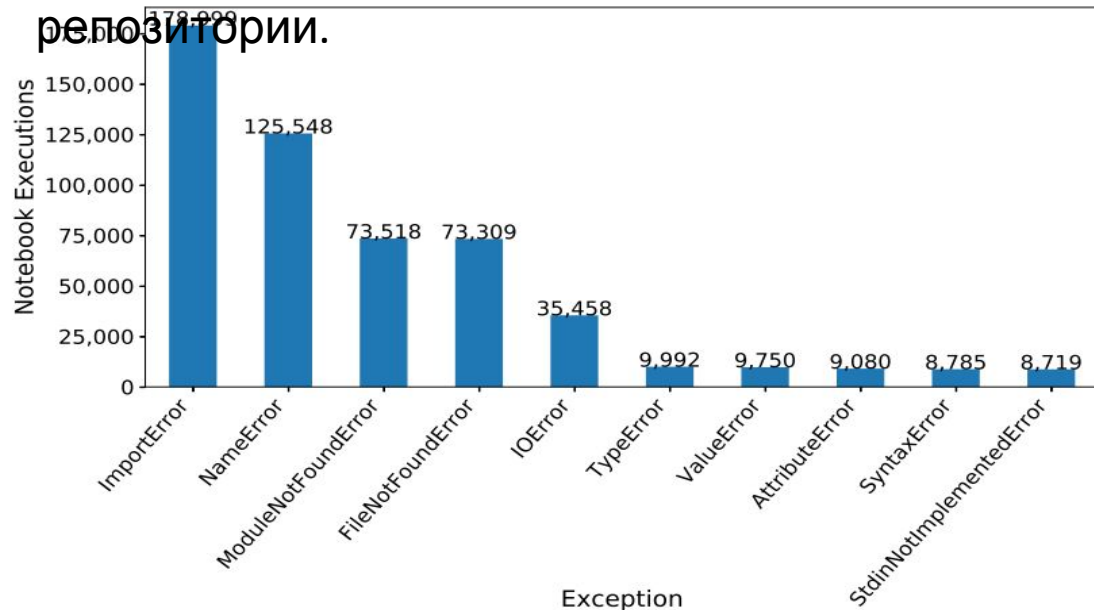
Признаки кризиса воспроизводимости:

- Фактический недостаток формализации воспроизводимости исследований в опубликованной литературе по многим научным направлениям
- широко распространена невозможность воспроизвести результаты опубликованных исследований;
- высокая распространенность «закрытых» методов исследования, которые завышают количество ложноположительных результатов
- отсутствие данных и алгоритмов анализа в научных публикациях.



## Факторы уменьшающие воспроизводимость

- Первое место — проблемы с зависимостями в библиотеках и зависимостями в зависимостях. Часть репозитория использовали requirements.txt, часть setup.py.
- Второе место — порядок исполнения. Тетрадка сохранена без очистки кода, порядок не сохранен и некоторые переменные объявлены или инициализированы после использования.
- Третье место — нет нужных данных, например, указаны абсолютные пути или данных вообще нет в репозитории.





## Факторы уменьшающие воспроизводимость

---

- Неуправляемая случайность в данных или алгоритмах (40%) (Random Seed)
- Исходные данные измененные вручную, а не с помощью скриптов
- Зависимость вывода и результатов от функций времени (13%)
- Различия отображения на графиках (некорректное использование matplotlib в том числе) (52%)
- Недоступны внешние данные (3%)
- Различия в выводе чисел с плавающей запятой (3%)
- Непостоянный порядок обхода словарей и др. контейнеров в python (4%)
- Различия в среде исполнения (27%)



# Способы улучшения воспроизводимости

- Пишите код в виде Python скриптов в полноценных IDE

```
File Edit Selection View Go Debug Terminal Help DataScienceCool.ipynb* - connect-petdetector - Visual Studio Code
EXPLORER
OPEN EDITORS
test.py
DataScienceCool.ipynb*
CONNECT-PETDETECTOR
.vscode
images
scripts
.gitignore
classify.ipynb U
DataScienceCool.ipynb U
demo_completed.py
demo.ipynb M
LICENSE
myenv.yml
README.md
score.py
setup.ipynb
test.ipynb U
test.py
OUTLINE
master* Python 3.7.3 64-bit 0 1
```

```
[1] print("hello world")
hello world

[2] import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
p = np.linspace(0,20,100)

[3] # Let's load and review some data
df = pd.read_csv("./data/pima-data.csv") # load Pima data
df.head(5)

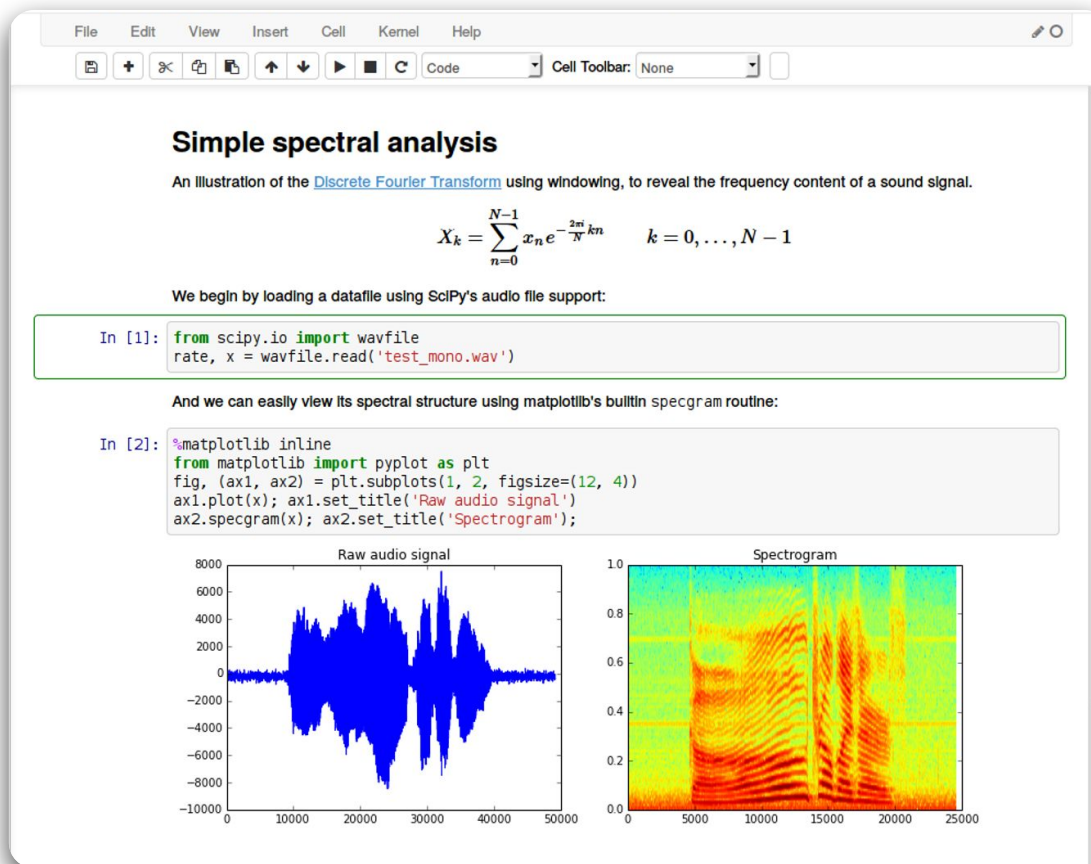
[4] import matplotlib.pyplot as plt # matplotlib.pyplot plots data
def draw_corr(df, size = 11):
    corr = df.corr() # data frame correlation function
    fig, ax = plt.subplots(figsize=(11, 11))
    ax.matshow(corr) # color code the rectangles by correlation value
    plt.xticks(range(len(corr.columns)), corr.columns) # draw x tick marks
    plt.yticks(range(len(corr.columns)), corr.columns) # draw y tick marks
```





# Способы улучшения воспроизводимости

- Используйте ноутбуки для быстрых экспериментов и их описания
  - Не должен являться основным артефактом разработки
  - Максимум абстракции (весь основной код во внешних скриптах)
  - Больше описаний
  - Больше графиков





# Способы улучшения воспроизводимости

- Документируйте эксперимент по ходу исследования, а не после его окончания
- Ведите вики проекта

**Тема** **Задачи** **Актуальность**

**Цель** **Гипотеза** **Проблема**

**Project**

**WIKI**  
Peer Review of Online Learning and Teaching

page discussion edit history move watch

Denlee my talk my preferences my watchlist my contributions log out

## Project overview

Jump to: navigation, search

Currently, as of 22nd July, 2008, the Peer Review of Online Learning and Teaching system is in beta stage of development for the following implemented features. The next major update will incorporate the features listed as "still to be implemented". That update is scheduled for early August, 2008.

### Currently implemented [\[edit\]](#)

- Create categories of learning resources to be reviewed
- Create banks of review criteria
- Ability to create, manage and edit individual review criterions
- The ability to manage response types
- Creating and managing user accounts.

### To still be implemented [\[edit\]](#)

- Improving workflow/interface management
- Implementing end user controls
- Implementing non-administrator controls and management

[Return to Main Page](#)

wiki

- [Main Page](#)
- [Project overview](#)
- [Instructions/Guides](#)
- [Criteria suggestions](#)

links

- [Project homepage](#)
- [Peer review instrument beta](#)
- [Bug reporting/Feature suggestion](#)
- [ALT C Homepage](#)
- [UniSA Homepage](#)

search



## Способы улучшения воспроизводимости

- Соблюдайте Codestyle, пишите комментарии и нотации
- Используйте линтеры и автоформатеры



flake8

```
def indent(string: str) -> int:
    """Count the indentation in whitespace characters.

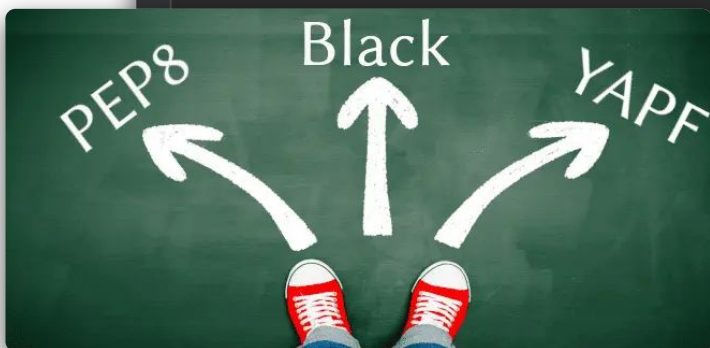
    Args:
        string (str): text with indents

    Returns:
        int: Number of whitespace indentations

    """
    return sum(4 if char == "\t" else 1 for char in string[: -len(string.lstrip())])
```



Bandit



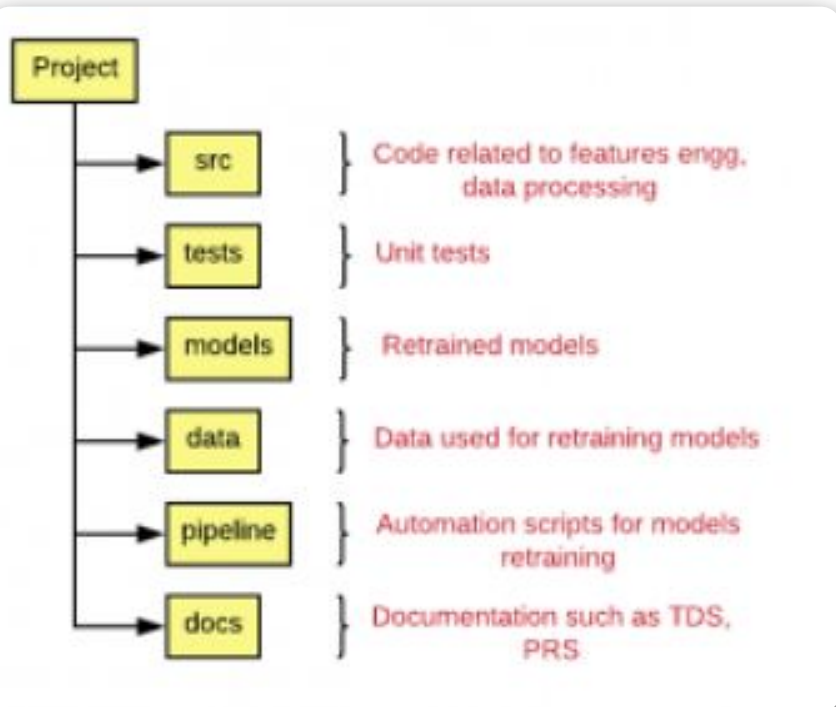
Pylint

Star your Python code!



# Способы улучшения воспроизводимости

- Шаблонизируйте DS проекты

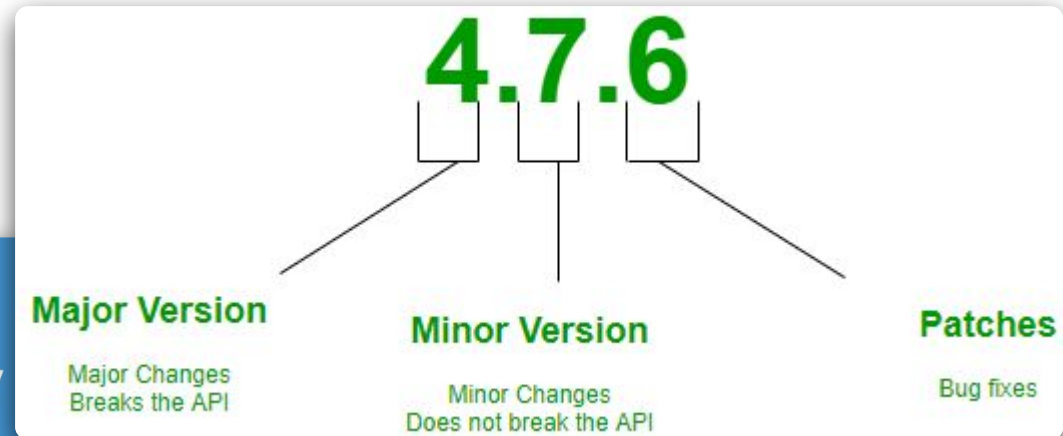
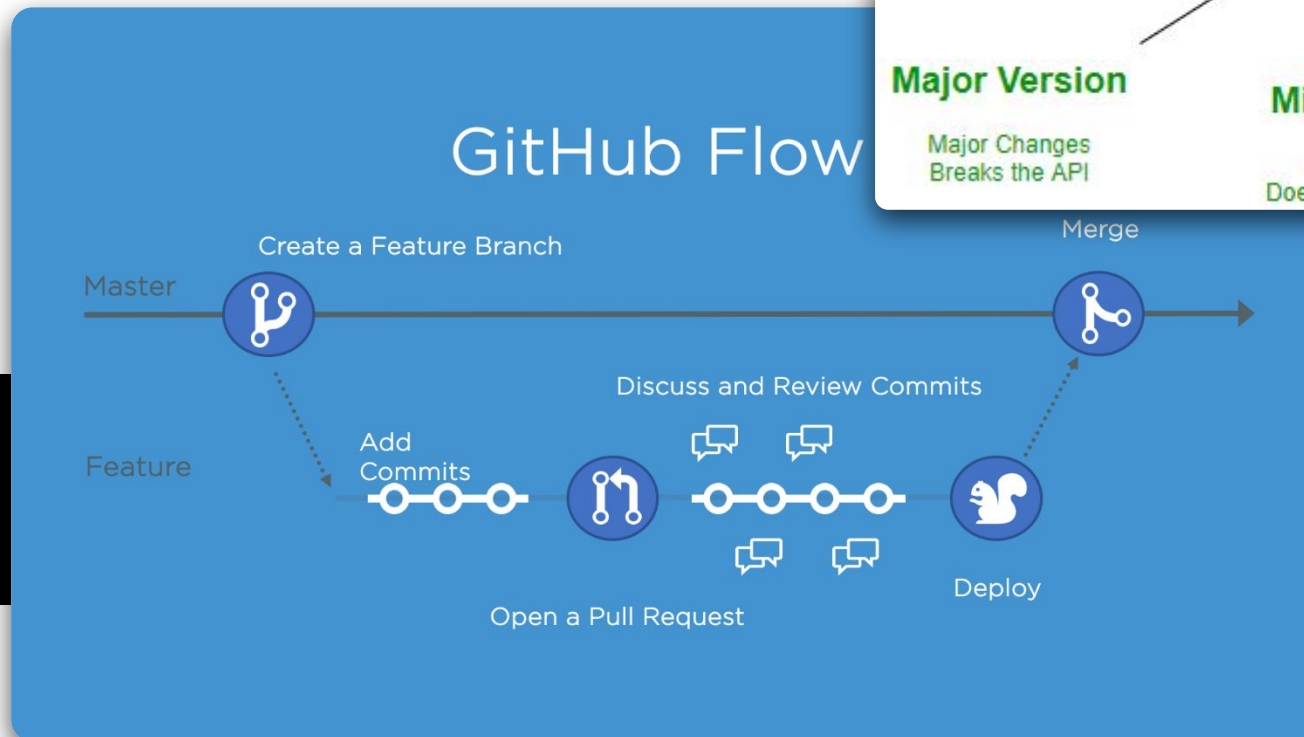


```
— LICENSE
— Makefile      <- Makefile with commands like `make data` or `make train`
— README.md     <- The top-level README for developers using this project.
— data
  — external    <- Data from third party sources.
  — interim     <- Intermediate data that has been transformed.
  — processed   <- The final, canonical data sets for modeling.
  — raw         <- The original, immutable data dump.
— docs          <- A default Sphinx project; see sphinx-doc.org for details
— models        <- Trained and serialized models, model predictions, or model summaries
— notebooks     <- Jupyter notebooks. Naming convention is a number (for ordering),
                  the creator's initials, and a short `-` delimited description, e.g.
                  `1.0-jqp-initial-data-exploration`.
— references    <- Data dictionaries, manuals, and all other explanatory materials.
— reports
  — figures     <- Generated analysis as HTML, PDF, LaTeX, etc.
                  <- Generated graphics and figures to be used in reporting
— requirements.txt <- The requirements file for reproducing the analysis environment, e.g.
                  generated with `pip freeze > requirements.txt`
— setup.py      <- makes project pip installable (pip install -e .) so src can be imported
— src           <- Source code for use in this project.
  — __init__.py <- Makes src a Python module
  — data        <- Scripts to download or generate data
    — make_dataset.py
  — features     <- Scripts to turn raw data into features for modeling
    — build_features.py
  — models       <- Scripts to train models and then use trained models to make
                    predictions
    — predict_model.py
    — train_model.py
  — visualization <- Scripts to create exploratory and results oriented visualizations
    — visualize.py
— tox.ini       <- tox file with settings for running tox; see tox.readthedocs.io
```



## Способы улучшения воспроизводимости

- Используйте системы контроля версий, соблюдайте версионирование кода, проводите качественное code-review

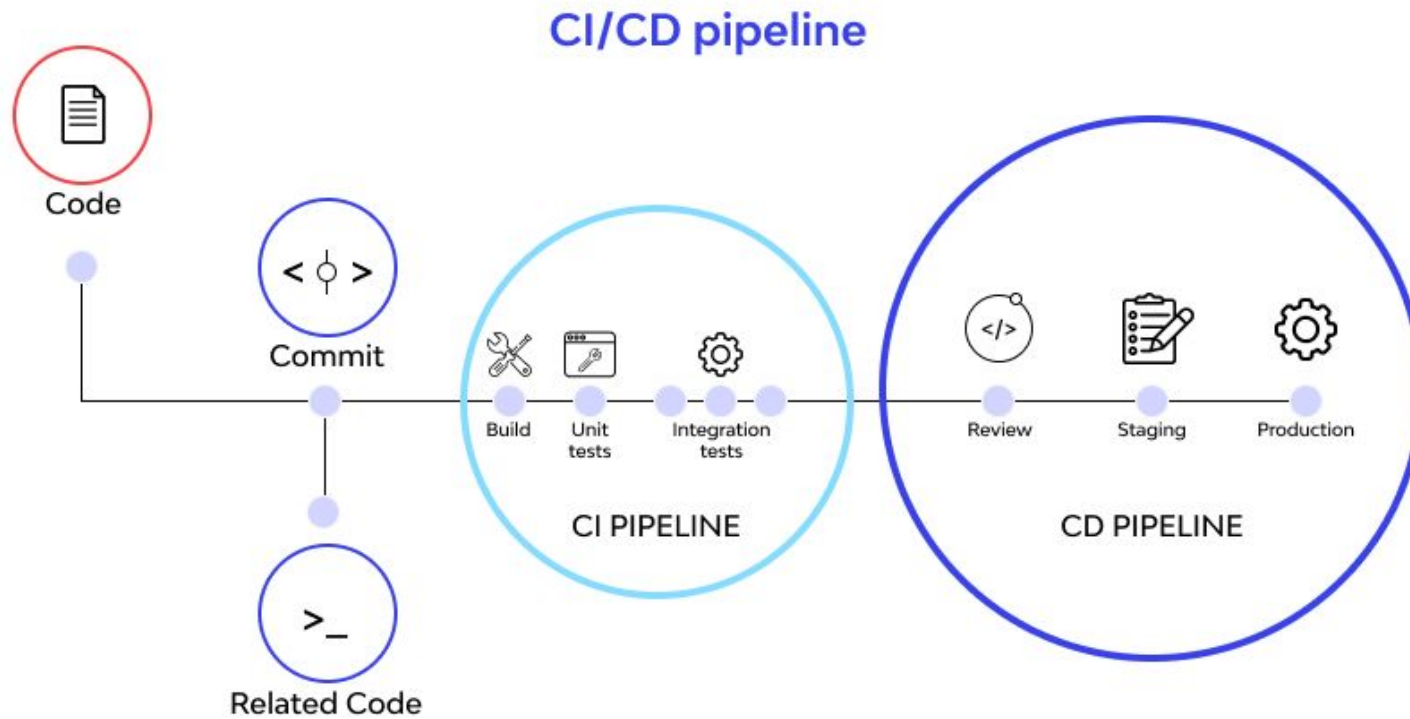






# Способы улучшения воспроизводимости

- Настройте CI/CD (Continuous Integration, Continuous Delivery)





# Способы улучшение воспроизводимости

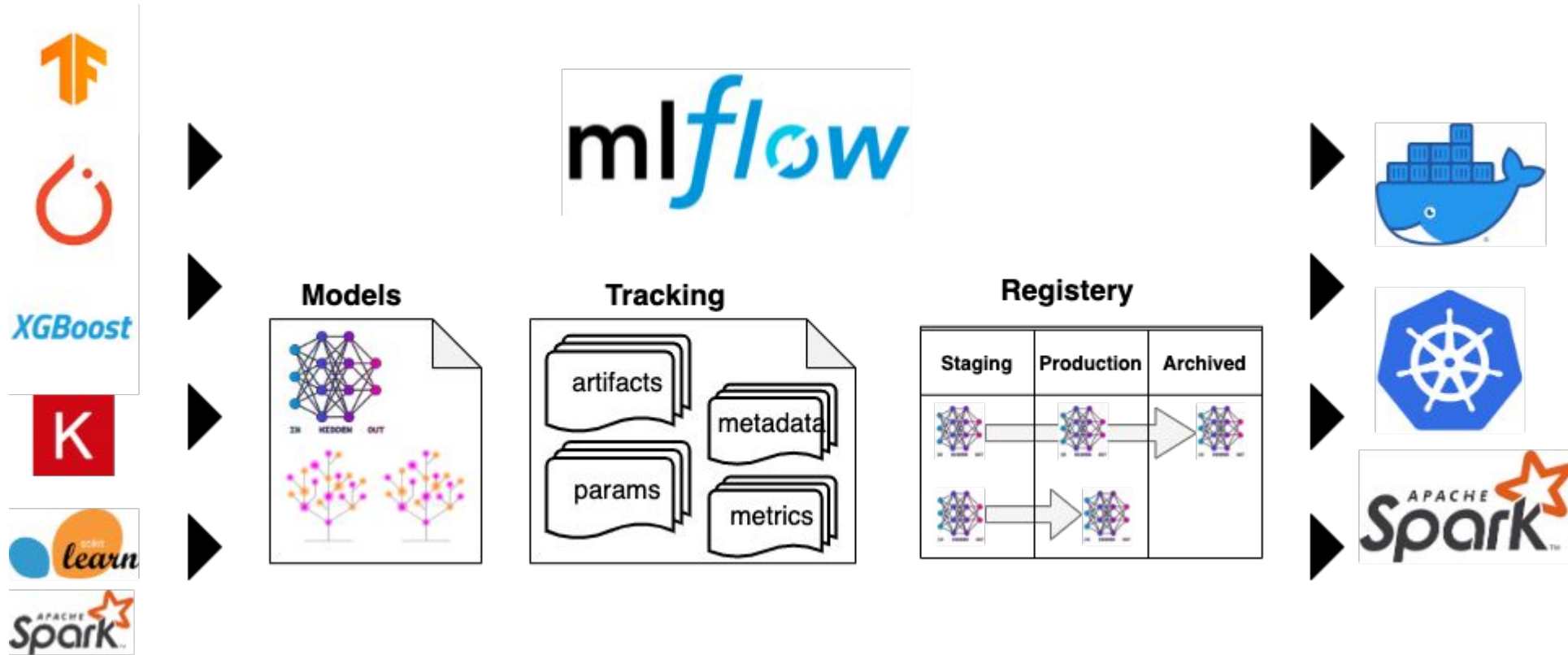
---

- **Покрывайте код, данные и модели тестами**
  - Юнит тесты
  - Интеграционные тесты
  - Тесты данных
  - Тесты моделей
  - Нагрузочные тесты



## Способы улучшения воспроизводимости

- Применяйте системы трекинга экспериментов (MLFlow)

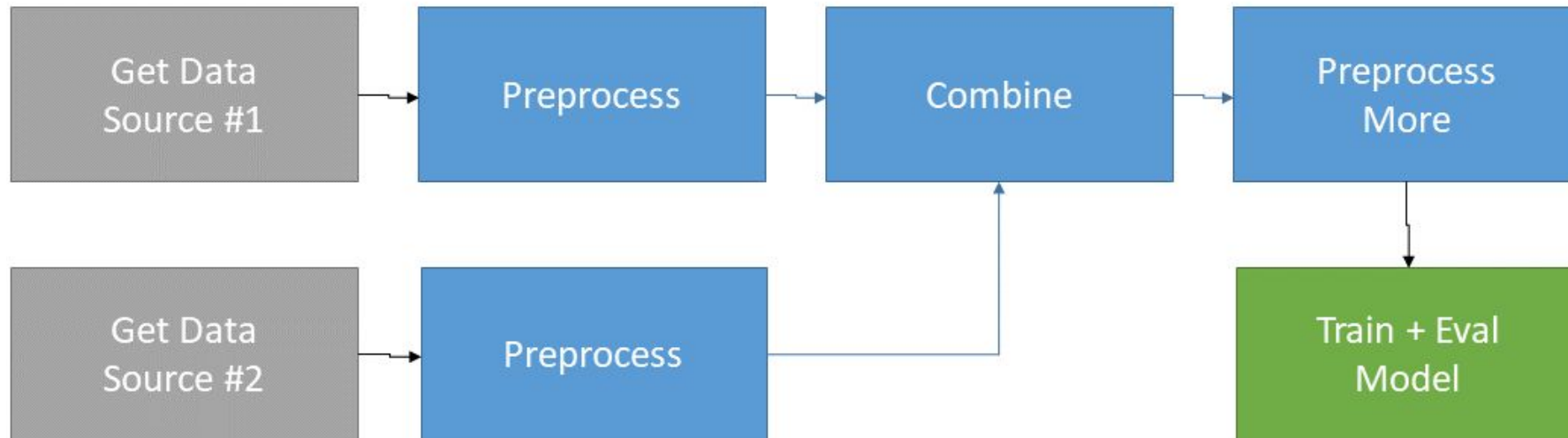






## Способы улучшение воспроизводимости

- Сохраняйте алгоритмы получения для каждого результата (workflow менеджеры с DAG)





## Способы улучшение воспроизводимости

---

- Оформляйте код, как python пакеты
- Создавайте CLI
- Пишите инструкции по запуску/интерфейсам вашего кода
- Избегайте ручного изменения данных



## Способы улучшение воспроизводимости

---

- Управляйте зависимостями и сохраняйте конфигурации с точными версиями (conda, pip)

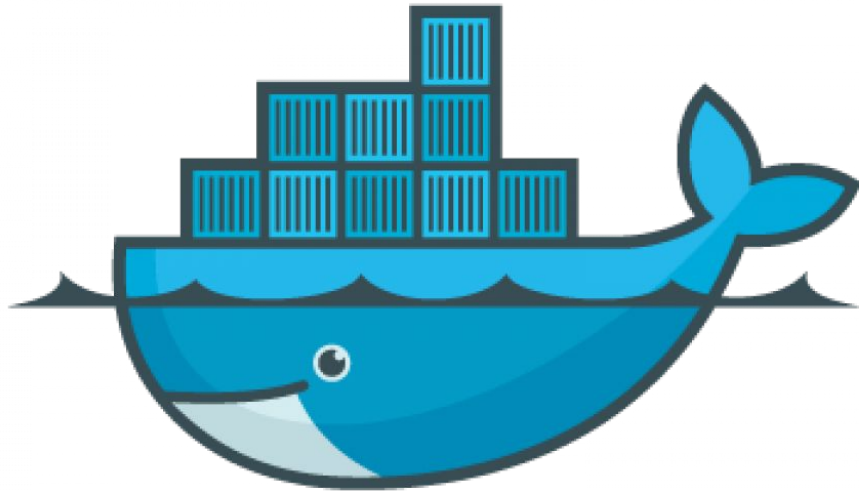




## Способы улучшение воспроизводимости

---

- Контейнеризируйте



docker



## Способы улучшение воспроизводимости

---

- Сохраняйте промежуточные данные (ускорение экспериментов)
- Фиксируйте случайные процессы в алгоритмах оптимизации и изменения данных (Random seed)
- Не работайте в одиночку. Выносите результаты на обсуждение.



# MLOps и production подход к ML исследованиям

*Концепция воспроизводимых и масштабируемых исследований в ML*

**Павел Кикин**

Газпромнефть ЦР

Руководитель направления NLP

[t.me/pavel\\_kikin](https://t.me/pavel_kikin)

