

ВК

ИИ по следам пользователей

Банда пяти 2.0

the team



Даня
ML Dev



Саша
Analytic



Артём
Social
Dev



Мила
Algorithmist



Антоша
Python Dev

the task

clickstream -> url

the problems



Низкая доля
положительного
класса



Необходимость
экстракции признаков



Пропуски и
шумы



Падающ
VK Clou

data preprocessing



Экстракция признаков: статистика по каждому юзеру и временным интервалам

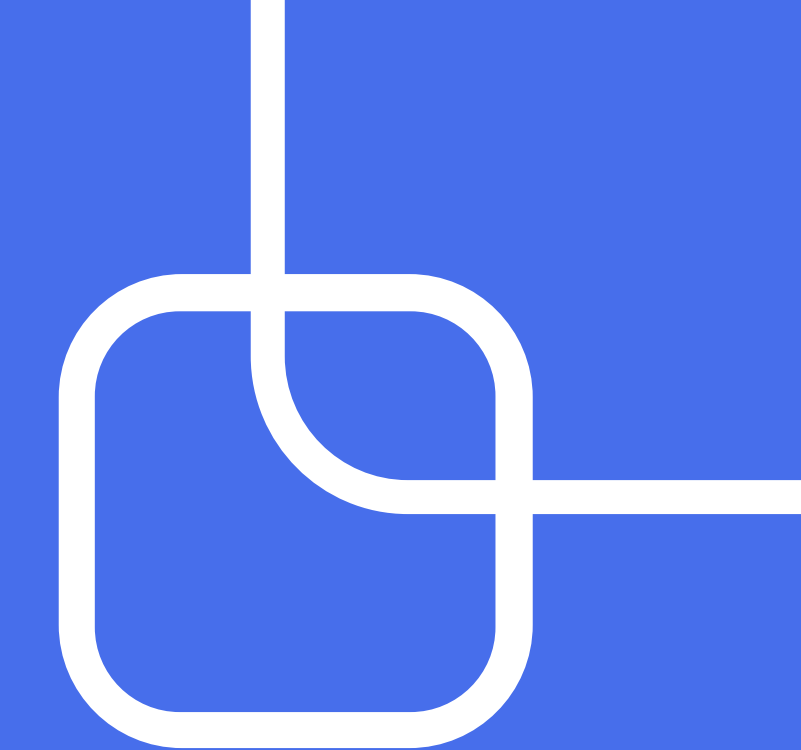


Детекция пропусков и шумов



Приведение токенов к однородному виду
(транслит + приведение к норм. форме)

the facts



> **0.23%**
невалидных
записей

Самые популярные токены:
новости, займы, кредиты,
смотреть, купить

Самые непопулярные токены:
фамилии,
токены с грамматическими
ошибками,
жаргонизмы

our solution

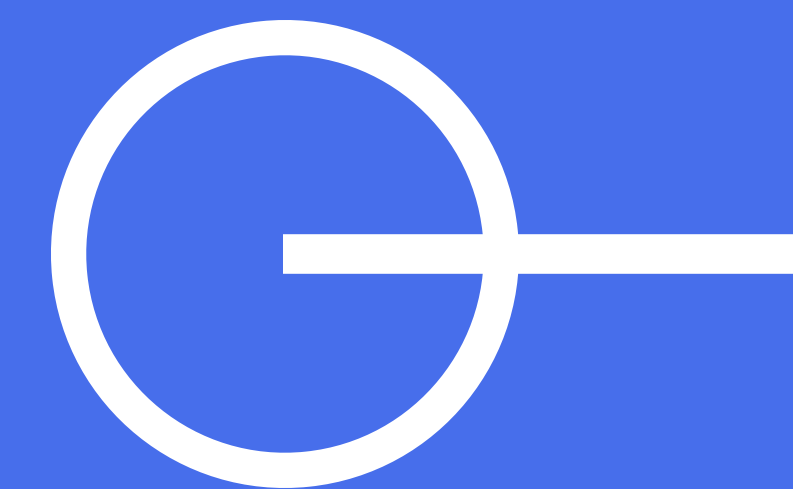


Стек

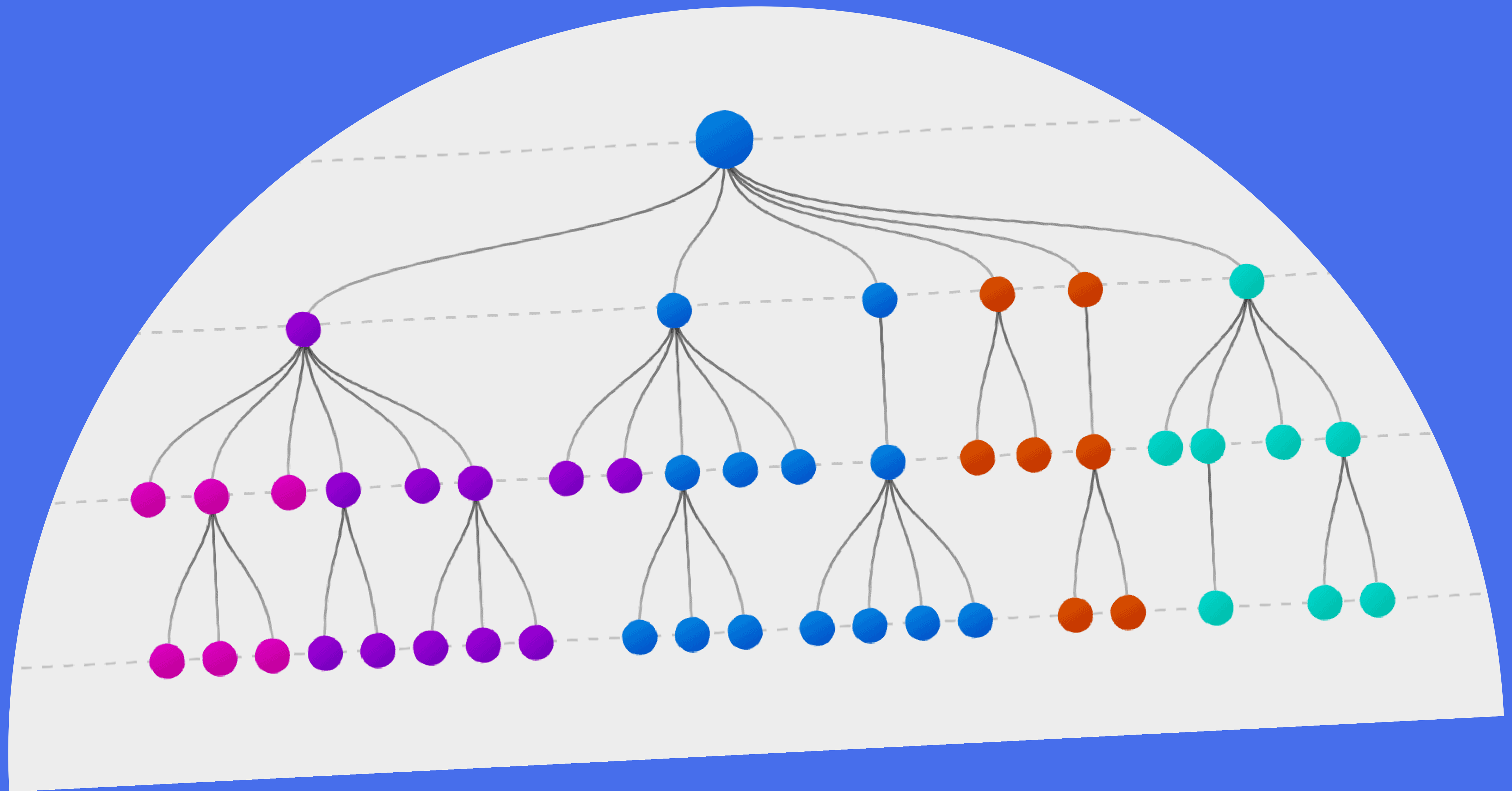
Pandas — для обработки
Numpy, Matplotlib — для
визуализации

Sklearn — для построения
моделей

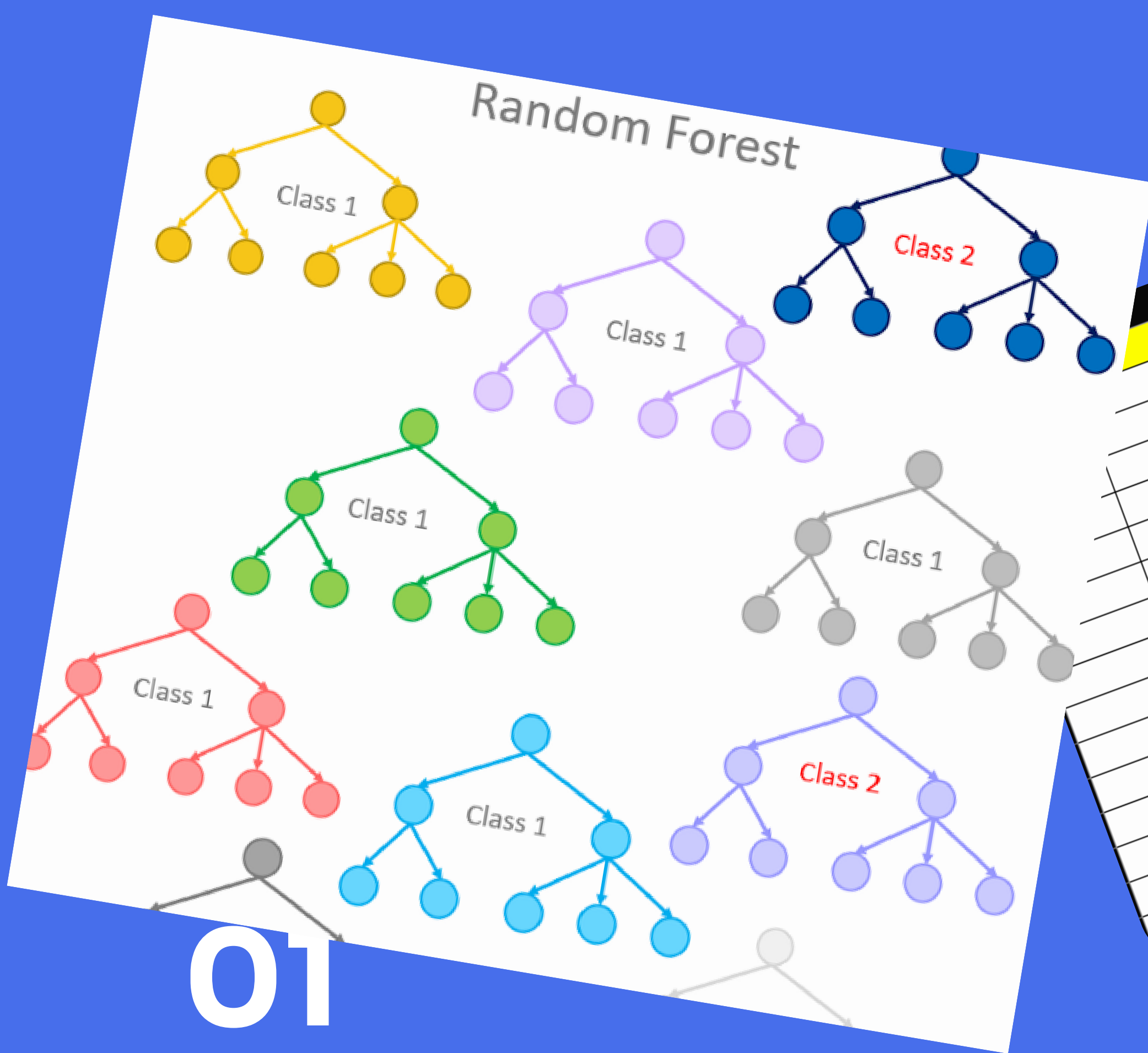
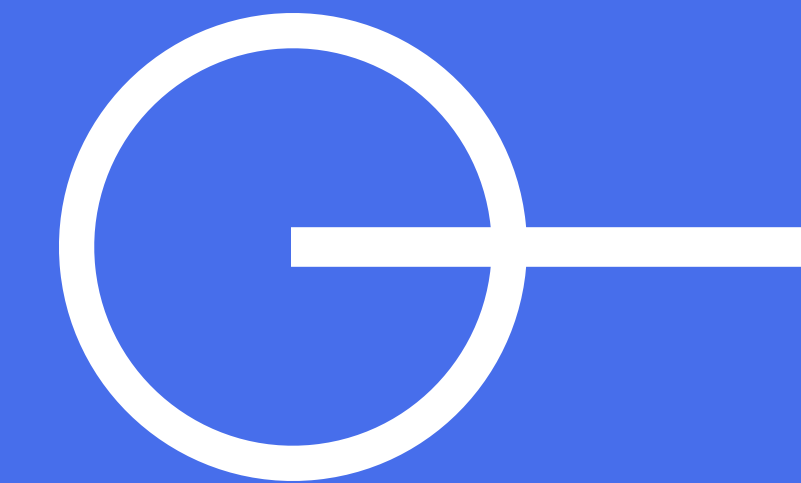
Rumorphy2 — для анализа
тегов



decision tree



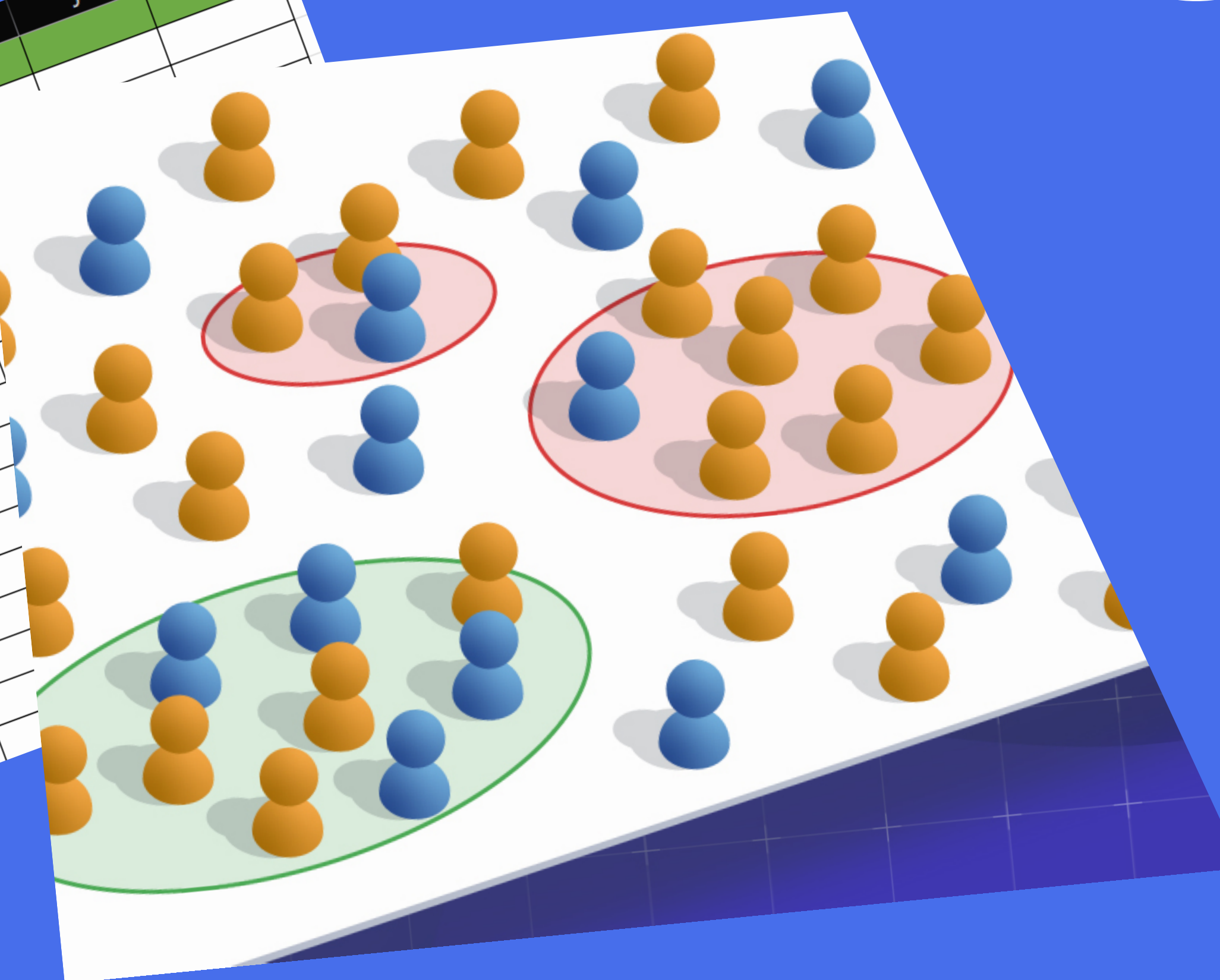
to-do



Расширить модель до леса
(Random Forest)



Взять больше признаков

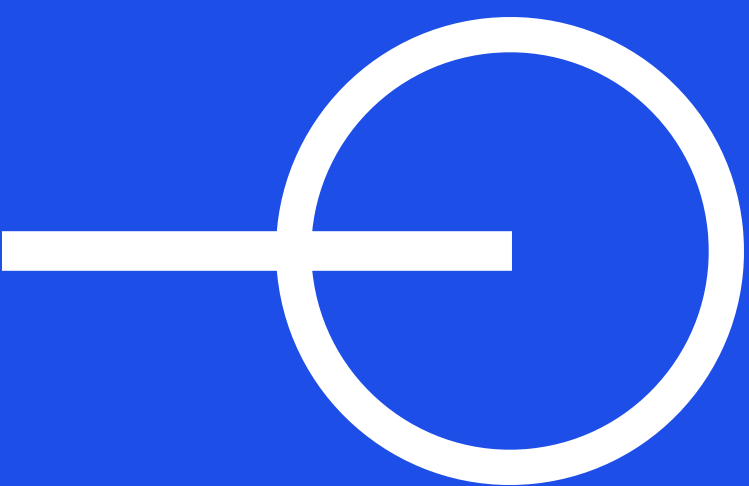


Увеличить количество выборок //
Обучить на другом сете

thank you for your
consideration



band of five 2.0



Если хочешь быть королем джунглей, мало просто вести себя как король, нужно быть королем. И никогда нельзя сомневаться, ведь сомнения порождают хаос и ведут к краху. Этому научила меня моя королева.

— **Mickey Pearson**

БК