

Машинное обучение в HeadHunter: умный поиск соискателей и работодателей

Александр Сидоров



Александр Сидоров

Руководитель
анализа
данных
HeadHunter

al.sidorov@hh.ru

facebook.com
/asidorov83

Кейс 1: Прескриннинг/модерация резюме

Задача

- Проверять качество заполнения резюме на самом верху воронки подбора

Проблема

- Значительная часть соискателей плохо и неполно заполняют резюме

Решение

- Использовать исторические данные по модерации резюме, чтобы автоматически подтверждать самые качественно заполненные

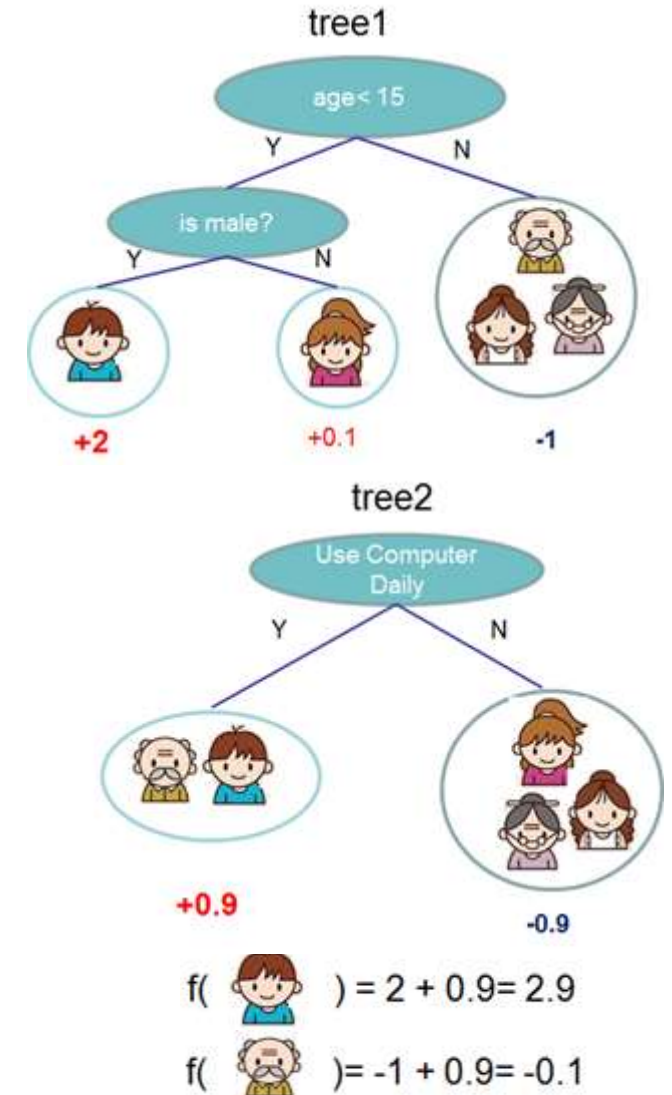
Схема работы автомодерации

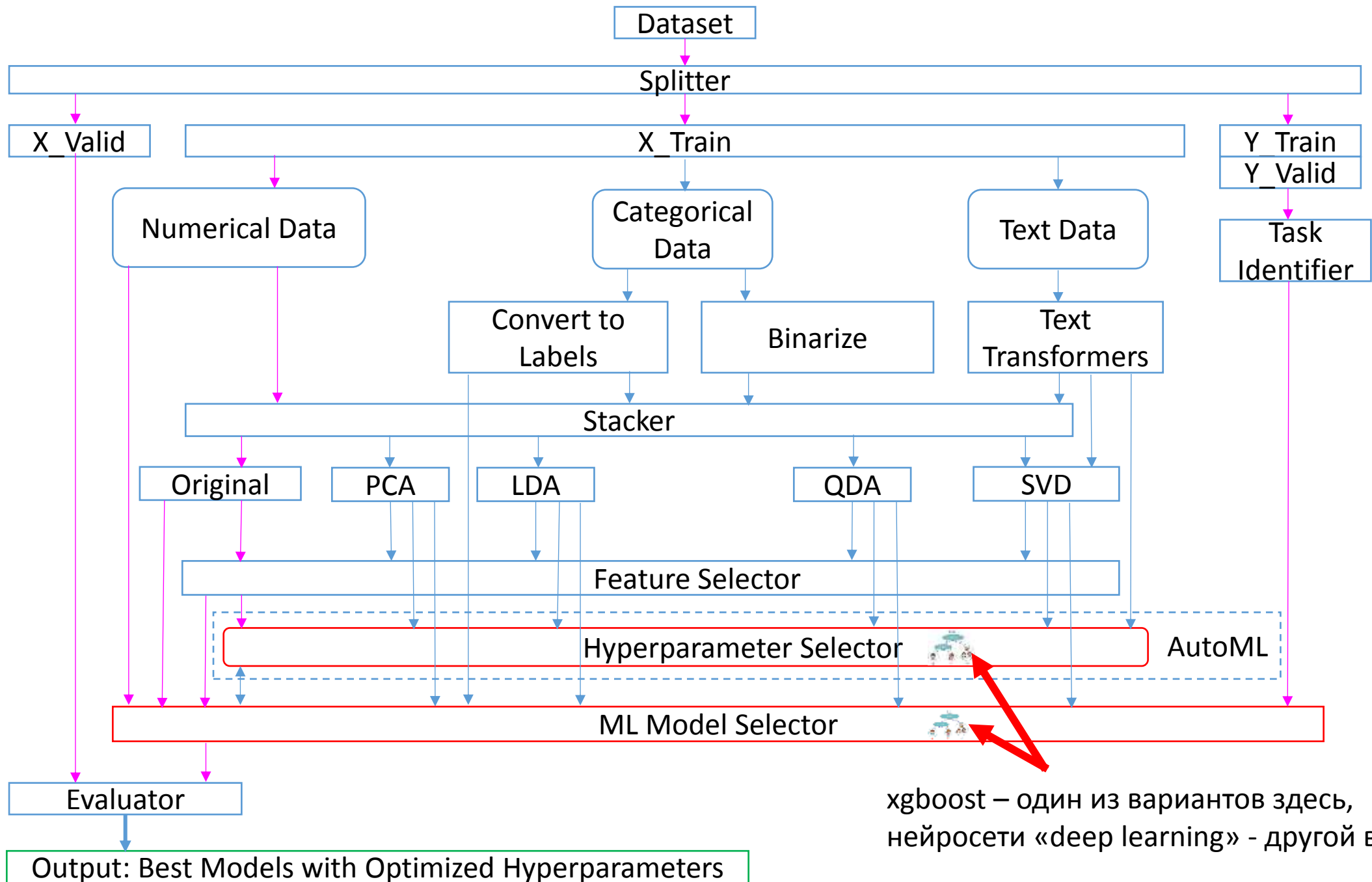
1. Результаты ручной модерации - исторические данные
2. Научили модель автомодерации
3. Новые резюме проверяются автомодератором
4. Качественные резюме автоматически подтверждаются
5. Спорные случаи отправляются людям

>30k новых резюме в день, ↑, кол-во модераторов и
качество модерации const, auc_roc 0.94
200k резюме в обучающем множестве, 1k признаков

Градиентный бустинг с xgboost

- хорошо работает с разнородной информацией
- не очень чувствителен к шуму
- довольно чувствителен к переобучению
- многопоточный, хорошая параллельность
- и бинарная классификация, и регрессия
- подбираем eta, n_estimators, max_depth
- можно полуавтоматически, через hyperopt
- «when in doubt – use xgboost» Owen Zhang, Kaggle

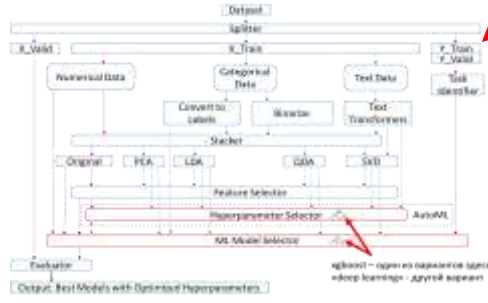




xgboost – один из вариантов здесь,
нейросети «deep learning» - другой вариант

вся Kaggle-подобная работа здесь

Выбор признаков и обучение моделей



Автоматическое тестирование и синхронный выкат

Расчёт байесовской априорной вероятности

Расчёт признаков и работа моделей в runtime

А/В-тесты, баланс «польза/ресурсы»

Сбор исторических данных

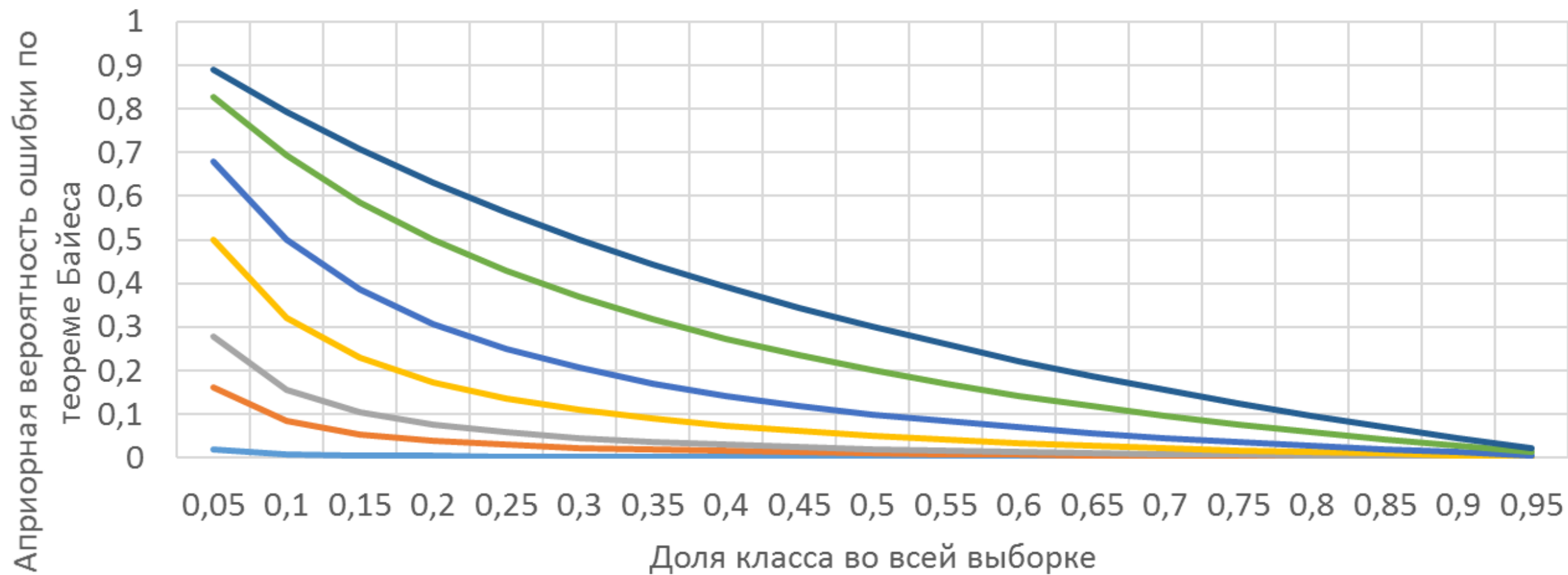
Оптимизация ресурсов

Регулярный пересчёт признаков и переподбор моделей

Мониторинг качества и ресурсоёмкости

Работа с экспериментальным кодом

Априорная вероятность неверной бинарной классификации при разной апостериорной точности модели-классификатора (accuracy)



Апостериорная точность (accuracy): — 0,999 — 0,99 — 0,98 — 0,95 — 0,9 — 0,8 — 0,7

Кейс 2: Подбор вакансий по резюме

Задача

- Подобрать по резюме вакансии, интересные соискателю

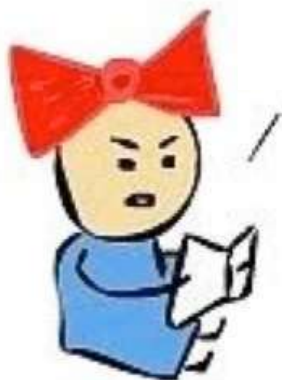
Решение

- Фильтр по параметрам вакансии
- Машинное обучение

Функционал

- Главная страница, рекомендуемые вакансии в списке резюме, рассылки с подходящими вакансиями

М-а-а-ш-и-и-н-н-о-е
о-б-у-ч-е-е-н-и-и-е



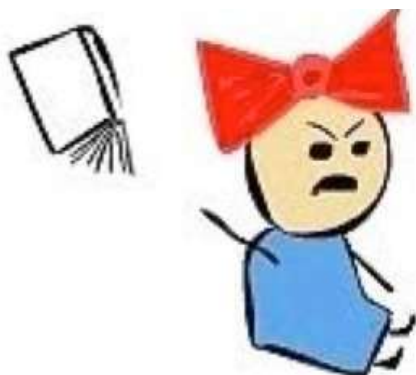
$$p(C | F_1, \dots, F_n) = \frac{p(C) p(F_1, \dots, F_n | C)}{p(F_1, \dots, F_n)}.$$



$$C = \sum C_{ij}, \quad C_{ij} = \log(1 + \exp((-1)^{|I_i > I_j|} (s_i - s_j)))$$
$$\frac{\partial C_{ij}}{\partial s_i} = -\frac{\partial C_{ij}}{\partial s_j} = -\frac{1}{1 + e^{s_i - s_j}} = -\rho_{ij}, \quad \frac{\partial C}{\partial s_i} = \sum_{j: I_i > I_j} (-\rho_{ij}) - \sum_{j: I_i < I_j} (-\rho_{ji})$$
$$\frac{\partial^2 C_{ij}}{\partial s_i^2} = \frac{\partial^2 C}{\partial s_j^2} = \rho_{ij}(1 - \rho_{ij}), \quad \frac{\partial^2 C}{\partial s_i^2} = \sum_j \rho_{ij}(1 - \rho_{ij})$$
$$g_i = \sum_{j: I_i > I_j} |\Delta Z_{ij}| (-\rho_{ij}) - \sum_{j: I_i < I_j} |\Delta Z_{ij}| (-\rho_{ji}), \quad h_i = \sum_j |\Delta Z_{ij}| \rho_{ij}(1 - \rho_{ij})$$



Ну, нафиг...



...буду экспертом
напишу формулу
вручную!



Схема работы рекомендательной системы



350 000 вакансий

Эвристический фильтр



2 признака

Фильтрующая модель 1



4 признака

Фильтрующая
модель 2



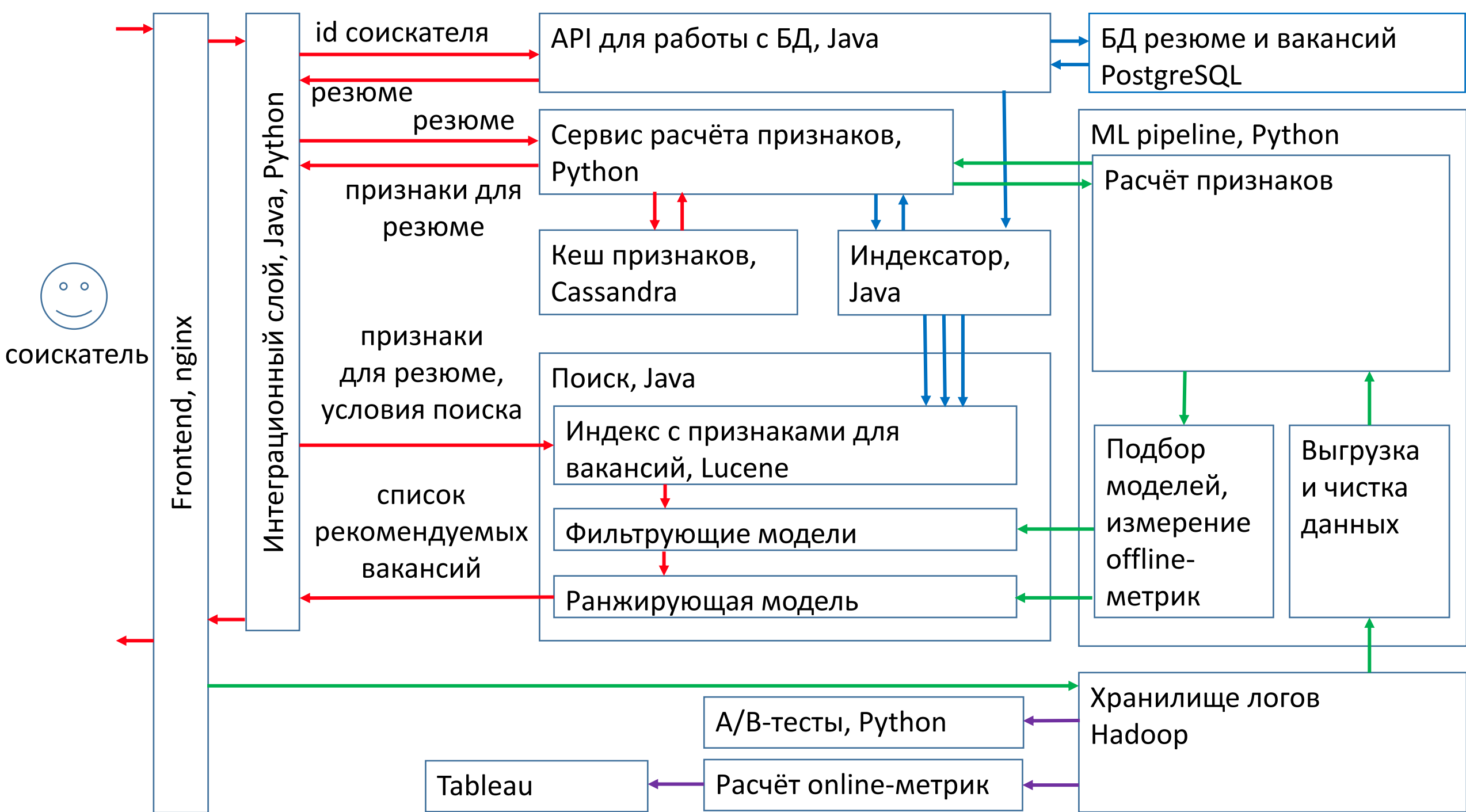
26 признаков

Ранжирующая
модель



85 признаков

500-44m пар в обучающих множествах, переподбор за ночь



Результаты в откликах по результатам А/В-тестов

+ x1000 в сутки
при внедрениях
там, где не было

+ x100 в сутки при
улучшениях
моделей,
факторов, таргетов

Иногда «просто»
уменьшаем
ресурсоёмкость

Кейс 3: Ранжирование откликов

Задача

- На вакансию поступили отклики, необходимо понять кого приглашать на собеседование или телефонное интервью

Проблема

- HR-ры тратят значительное количество времени на разбор резюме и определение списка для собеседования

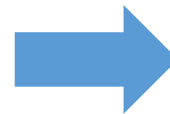
Функционал

- Кнопка «Лучшие» на странице откликов на вакансию

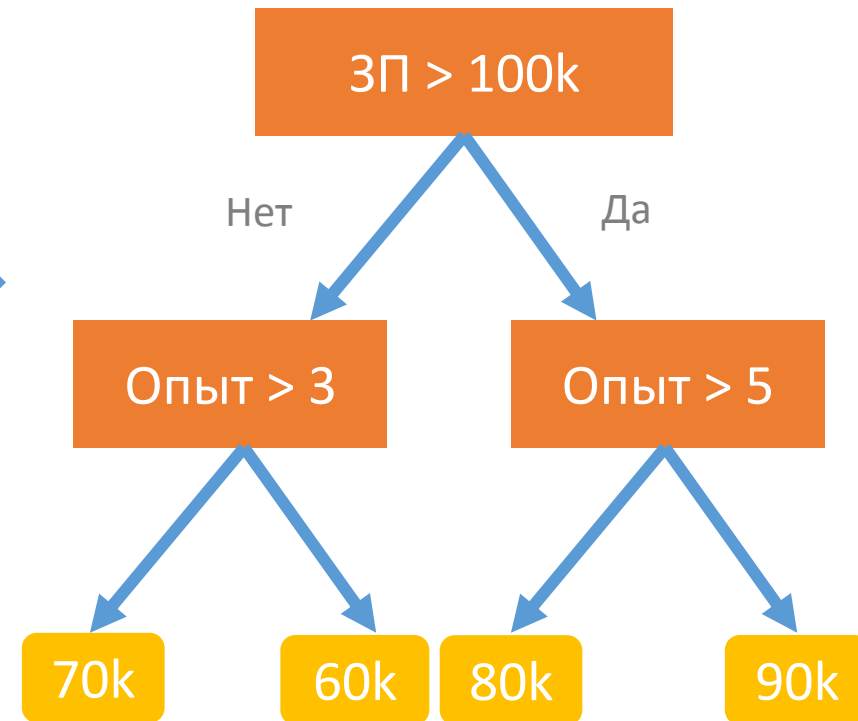
Человек: делает по инструкции, учится на небольшом опыте
ошибается, субъективен, дорогое время

Неформализованный список:

1. Опыт работы
2. Образование
3. Навыки
4. Зарплатные ожидания
5. ...
6. ...
7. ...

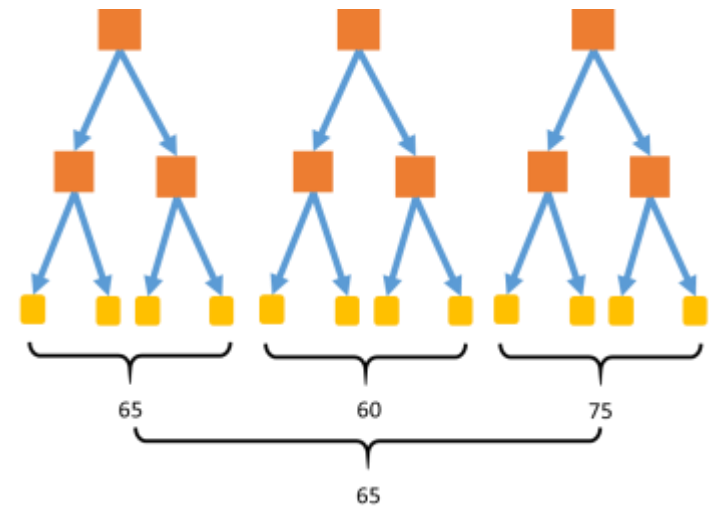


Дерево решений



Машина: учится на большом опыте
тоже ошибается, объективна, дешёвое время

1. Исторические данные: кого пригласили на собеседование
2. Учим модель
3. Применяем модель к кандидатам
4. Улучшаем модель



Результат:

- 3,5m откликов в обучающем множестве
- 500 признаков
- 800k пар в день, 95%: 600ms

Кейс 4: поиск по соцсетям

Задача

- Найти соискателей, у которых нет резюме или у которых они закрыты


Решение

- Собрать все профили людей, сопоставить их между собой, выделить полезное для работы

Функционал

- Профили в поиске по резюме

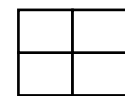
Схема работы сопоставления профилей

 365 млн. записей, 53.856 трлн. потенциальных пар

Нормализация, векторизация

Объединение в блоки

Попарное
сравнение



15 признаков



15 признаков

14k пар в обучающем множестве
точность: 99,9% - 150k пар, 80% - 4m
пересчёт за месяц

Рост использования машинного обучения в рекрутинге

Модерация
резюме

Ранжирование
откликов

Рекомендательные
системы

Умный поиск с
ML



Выводы и прогнозы

Использование машинного обучения дает хороший эффект в рекрутинге: в среднем целевые метрики вырастают на десятки процентов.

В краткосрочной перспективе

- вырастет эффективность подбора за счет экономии времени рекрутеров
- увеличится скорость поиска работы для соискателей

В среднесрочной перспективе

- подбирать сотрудников и искать работу будут ансамбли систем машинного обучения с уменьшающимся участием человека



**17 лет меняем мир рекрутинга и HR к
лучшему!**