

NEURAL NETWORKS & DEEP LEARNING - COURSE 4 - HANDOUT NO. 1 - PAGES 2

BINARY CLASSIFICATION

1: CAT
0: NOT CAT

LOGISTIC REGRESSION AS A NEURAL NET

FINDING THE MINIMUM WITH GRADIENT DESCENT

- FIND THE DOWNSLOPE DIRECTION (USING DERIVATIVES)
- WALK (UPDATE w & b) AT A α LEARNING RATE
- REPEAT UNTIL YOU REACH BOTTOM (CONVERGE)

PUTTING IT ALL TOGETHER

MINI NEURAL NET

$z(x) = wx + b$
 $\hat{y} = a(z) = \sigma(\text{SIGMOID}(z))$

FORWARD PROPAGATION
BACKWARD PROPAGATION + UPDATE w & b
REPEAT UNTIL IT CONVERGES

LOSS = $\mathcal{L}(\hat{y}, y)$
COST = $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$
COST = LOSS FOR THE ENTIRE DATASET

FIND THE MINIMUM

2 LAYER NEURAL NET

INPUT LAYER
HIDDEN LAYER
OUTPUT LAYER

$z = w_1x_1 + w_2x_2 + w_3x_3 + b$
 $a = \text{ACTIVATION FUNCTION}$

SHALLOW NEURAL NETS

WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION - $a = z$

LAYER 1
 $a^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]}$

LAYER 2
 $a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$

PLUG IN $a^{[1]}$

$a^{[2]} = w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]}$
 $= w^{[2]}w^{[1]}x + w^{[2]}b^{[1]} + b^{[2]}$
 $= w'x + b'$ ← LINEAR FUNCTION

WE COULD JUST AS WELL HAVE SKIPPED THE WHOLE NEURAL NET & USED LIN. REGR.

INITIALIZING w & b

WHAT IF: INIT TO 0

THIS WILL MAKE ALL THE UNITS TO BE THE SAME AND LEARN EXACTLY THE SAME FEATURES

SOLUTION: RANDOM INIT BUT ALSO WANT THEM SMALL SO RAND * 0.01

HYPERPARAM

INTRO TO DEEP LEARNING

SUPERVISED LEARNING

INPUT: X	OUTPUT: Y	NN TYPE
HOME FEATURES	PRICE	STANDARD NN
AD+USER INFO	WILL CLICK ON AD (y/n)	NN
IMAGE	OBJECT (1...1000)	CONV. NN (CNN)
AUDIO	TEXT TRANSCRIPT	RECURRENT NN (RNN)
ENGLISH	CHINESE	NN
IMAGE/RADAR	POS OF OTHER CARS	CUSTOM/HYBRID

AREA # ROOMS
LOCATION
WEALTH
PRICE

STRUCTURED
UNSTRUCTURED

"THE QUICK BROWN FOX"

HUMANS ARE GOOD AT THIS

WHY NOW?

LOTS OF DATA
HARDWARE OPTIMIZED ALGOS

ONE OF THE BIG BREAKTHROUGHS HAS BEEN MOVING FROM SIGMOID TO RELU FOR FASTER GRADIENT DESCENT

PERFORM
AMT. OF DATA "LABELS"

IDEA
EXPERIM.
CODE

FASTER COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS

NETWORK ARCHITECTURES

STANDARD NN
CONVOLUTIONAL NN
RECURRENT NN

Natural Language Processing

в банковской сфере

Семинар МШЭ МГУ, 11.10.2021
Диана Котерева, Сергей Афанасьев

Ренессанс
КРЕДИТ

М Г У
МОСКОВСКАЯ
МШЭ
ШКОЛА ЭКОНОМИКИ

Коротко о нас



Диана Котерева
Head of R&D
КБ «Ренессанс Кредит»



Сергей Афанасьев
Vice-President
Chief Data Scientist
КБ «Ренессанс Кредит»



Банк «Ренессанс Кредит»

75 млн
заявок



13 млн
клиентов

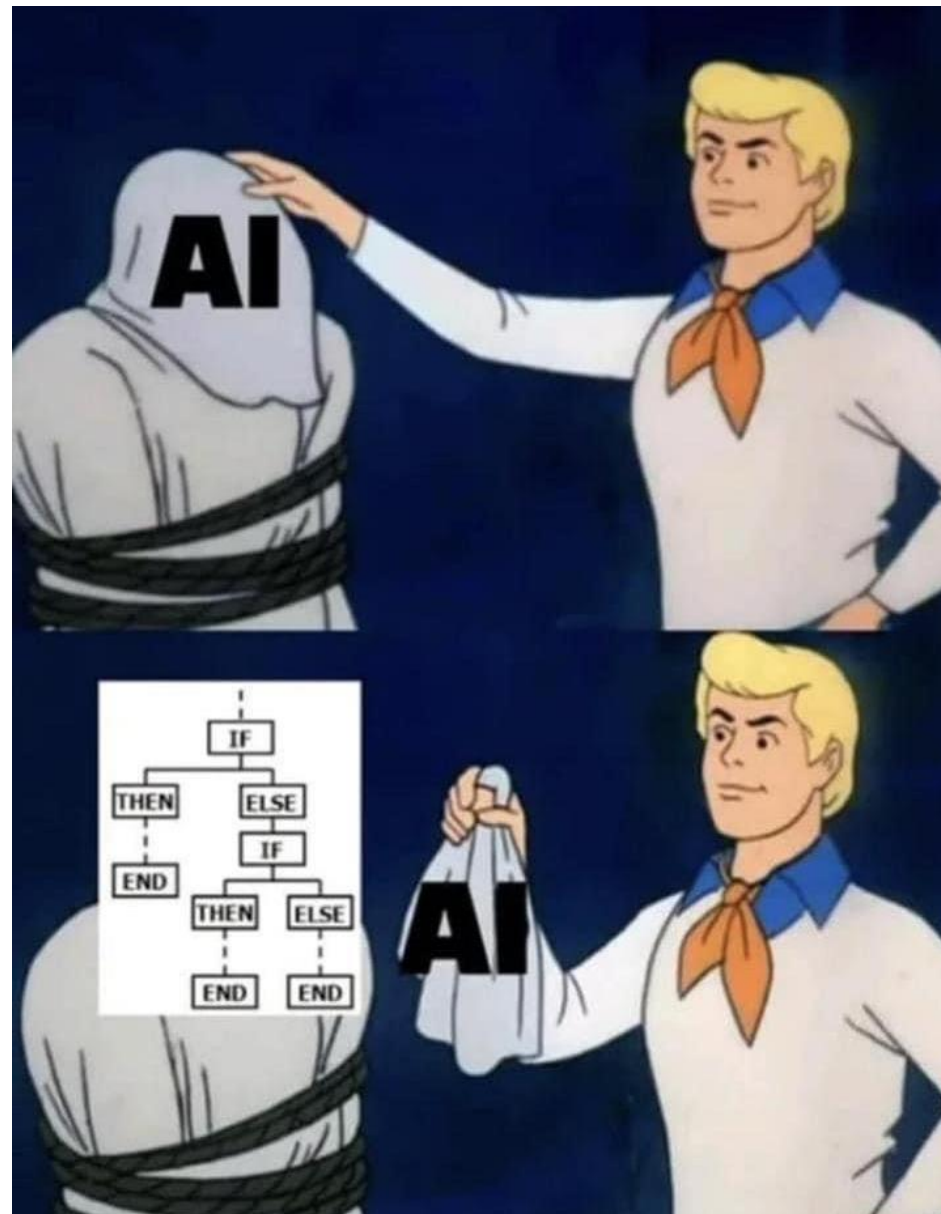


9 млн
кредитов



– Как отличить Machine Learning от **AI**?

– Если написано на Python – это Machine Learning, а если на PowerPoint, то **AI**.



Data Science в банковской сфере

	Risk, Anti-fraud	Collection	CRM	CS, CC, TM	IT, IT Security	Other
Classical Machine Learning	Credit Scoring <ul style="list-style-type: none"> - Application scoring - Behavioral scoring - Combine scoring Anti-fraud Models <ul style="list-style-type: none"> - Внутр/внеш fraud - Транзакционный fraud IFRS-9/Basel III	Collection Scoring <ul style="list-style-type: none"> - Entry/Self_cure модели - DPD модели - Модели Recovery Rate 	Recommender System <ul style="list-style-type: none"> - РТВ модели Attrition model <ul style="list-style-type: none"> - Модели оттока клиентов (Churn prediction) Targeted Advertising <ul style="list-style-type: none"> - Таргетирование рекламы 		Biometrics <ul style="list-style-type: none"> - Keystroke Dynamics & Mouse/Touch Movements - Device Print для ИБ/МБ Error Detection <ul style="list-style-type: none"> - Выявление сбоев в системах Банка 	AML <ul style="list-style-type: none"> - AML правила и модели Time Normalization <ul style="list-style-type: none"> - Оптимальная загрузка массовых подразделений Staff Recruitment <ul style="list-style-type: none"> - Отбор резюме по ключевым словам
Natural Language Processing	Feature Engineering <ul style="list-style-type: none"> - Предикторы на текстовых данных для Скоринга Voice Recognition <ul style="list-style-type: none"> - Аутентификация клиентов по голосу Fraud Text Mining <ul style="list-style-type: none"> - Анализ корп. переписки 	Feature Engineering <ul style="list-style-type: none"> - Предикторы на текстовых данных для Collection Text-To-Speech <ul style="list-style-type: none"> - Голосовой коллектор 	Feature Engineering <ul style="list-style-type: none"> - Предикторы на текстовых данных для CRM-моделей Recommender System <ul style="list-style-type: none"> - Анализ транзакции с помощью RNN для моделей отклика 	Speech recognition CS Text Mining <ul style="list-style-type: none"> - Извлечение информации - Анализ тональности - Чат-боты (сайт, ИБ/МБ) Text-To-Speech <ul style="list-style-type: none"> - Синтез речи (Роман-2.0) 	Security Text Mining <ul style="list-style-type: none"> - Проверка email-писем на утечку информации и т.п. Topic Model <ul style="list-style-type: none"> - Распределение заявок в Help-Desk по тематикам 	Post Mail Classifier <ul style="list-style-type: none"> - Распределение корреспонденции по тематикам HR Text Mining <ul style="list-style-type: none"> - Анализ переписки и выявление проблем в работе сотрудников
Computer Vision	Photo Biometrics <ul style="list-style-type: none"> - Вход в ИБ/МБ - Вход в системы Банка Object Detection <ul style="list-style-type: none"> - Скоринг по фото клиента - Проверка фото партнеров Spoofing <ul style="list-style-type: none"> - Выявление подделок 			OCR (Text recognition) <ul style="list-style-type: none"> - Перевод скан.копий документов в текстовый формат. 		Photo Biometrics <ul style="list-style-type: none"> - Аутентификация клиентов в ДО/ККО - Пропускные системы для сотрудников

Natural Language Processing (типы задач)

Текстовая аналитика

- Анализ тональности
- Извлечение структурированной информации из текста
- Классификация текстов
- Кластеризация текстов
- Информационный поиск

STT and TTS

- Распознавание речи (Speech-To-Text)
- Синтез речи (Text-To-Speech)

Генерация текста

- Машинный перевод
- Обобщение текста
- Аннотирование текста
- Упрощение текста
- Изменение стилистики текста

Чат-боты

- Разговорные чат-боты
- Вопросно-ответные системы
- Ассистенты (Алиса, Siri)

NLP-моделирование в банковской сфере

Котерева Д.М.
Группа моделирования и оперативного анализа
Управление статистического анализа
11 октября 2021 г.

Основные задачи Data Science в банке

МОДЕЛИ



Risks (scoring)



CRM



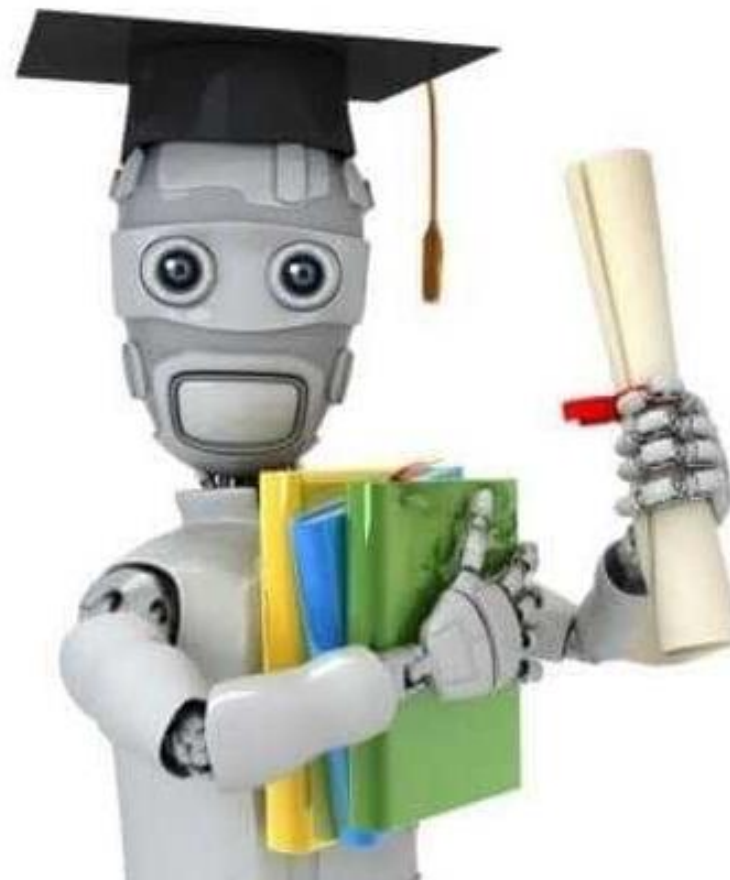
Collection



Antifraud



R&D (NLP, CV, etc)



NLP для подразделений банка

Обогащение банковских моделей



CRM

Обогащение CRM-моделей
текстовыми данными



Collection

Обогащение моделей взыскания
текстами

+2,5% GINI



Risks

Обогащение скоркарт
текстовыми данными

Сервисные NLP-задачи



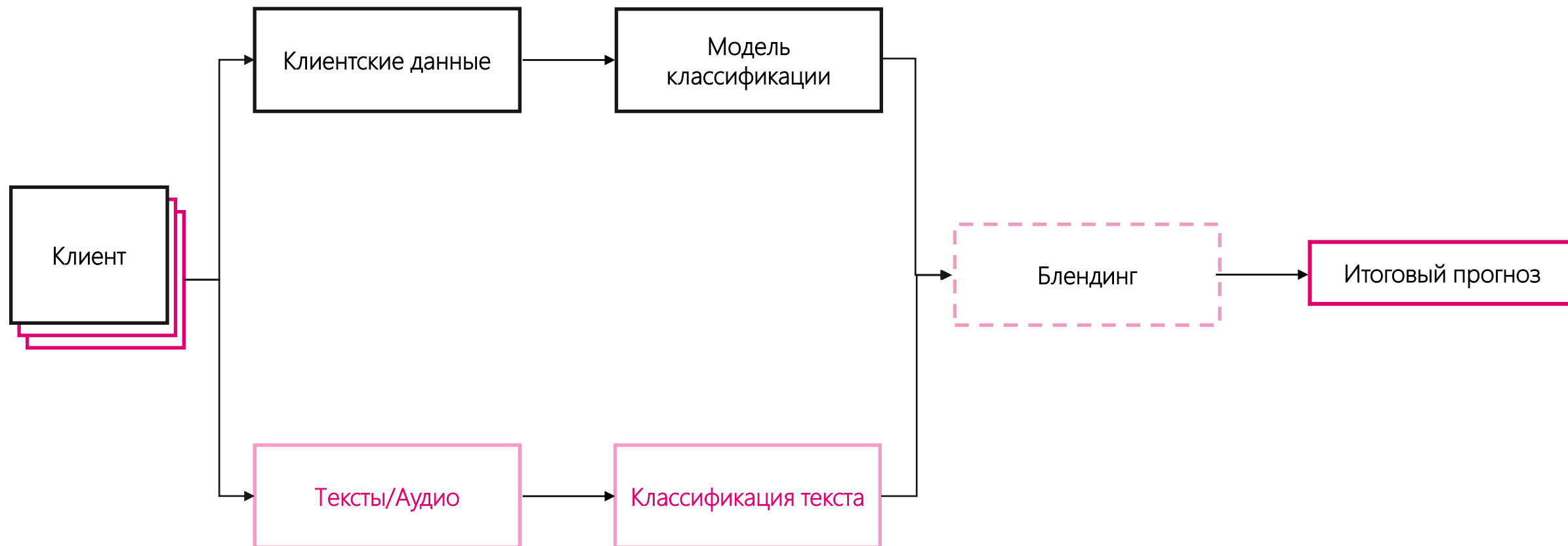
Customer Service



Others

Collection, AML, CRM, Antifraud, IT-
Security и др.

Обогащение банковских моделей

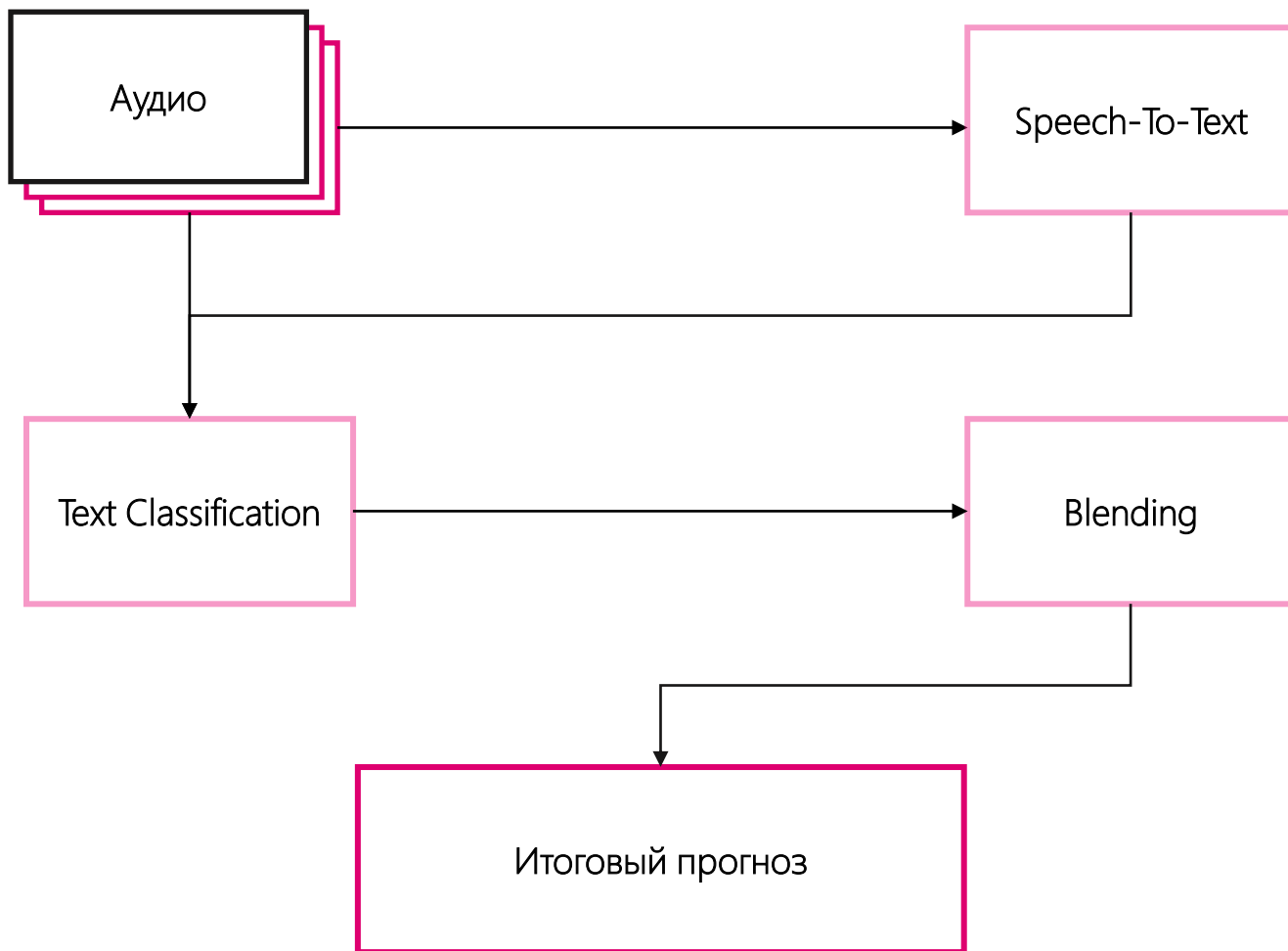


Обработка аудио и текстовых данных

Источники:

- Записи звонков
- Тексты (чаты, обращения)

Обучение классификатора на том же таргете, что и основная модель (Risks, Collection, CRM)



Модели распознавания текста из аудио записи

Объединение прогнозов основной и текстовой моделей с помощью блендинга

Speech-To-Text

1. **Цель модели:** Получить оптимальную последовательность слов для заданной последовательности звуков.

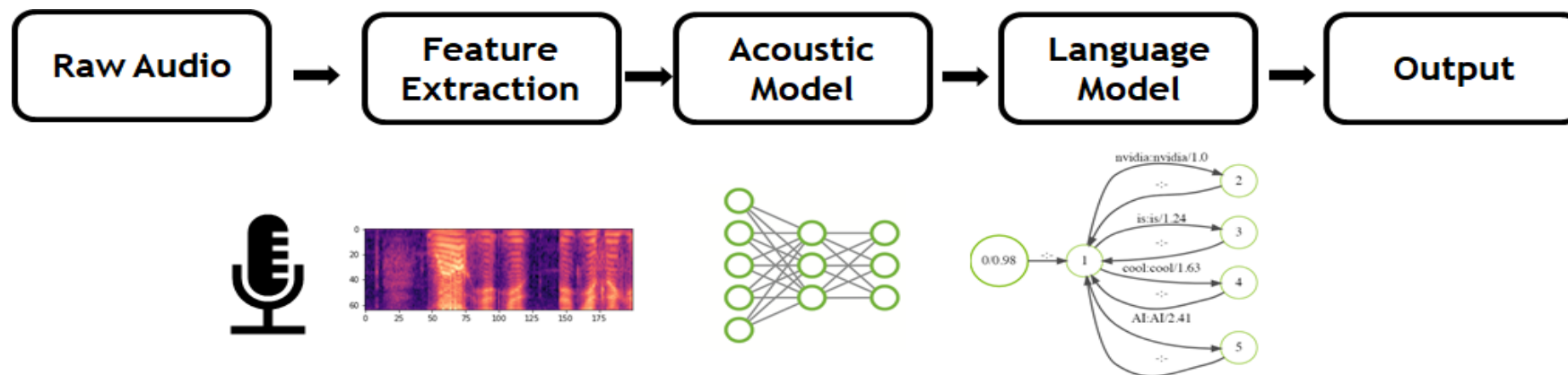
2. **Дано:** Необработанный аудиосигнал.

3. **Процесс обработки:**

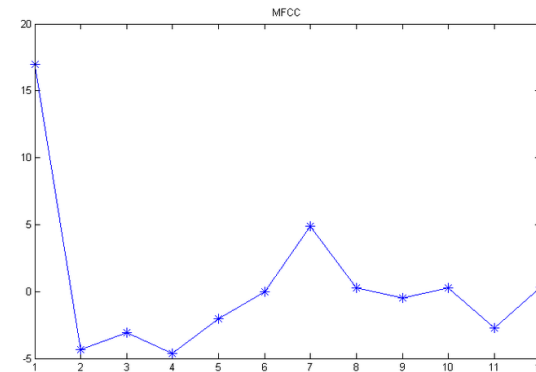
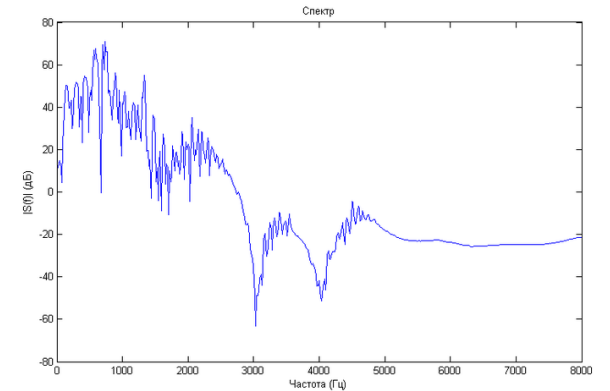
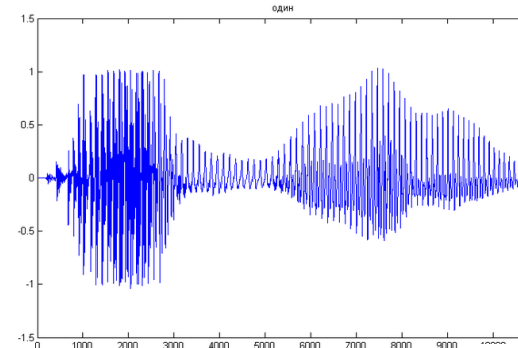
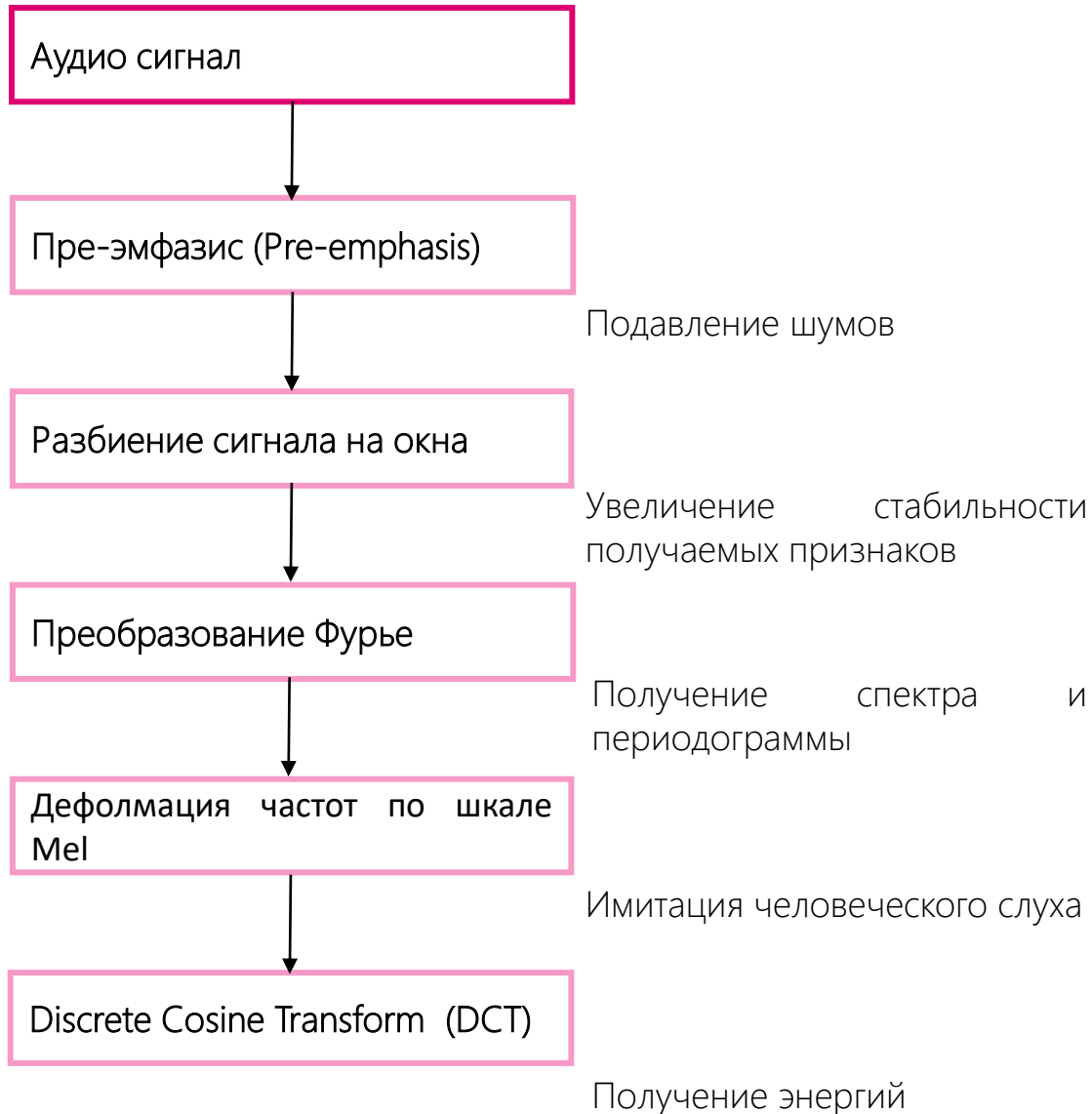
3.1. Получение признаков из аудиосигнала (Извлечение MFCC)

3.2. Получение последовательности текстовых признаков (акустическая модель)

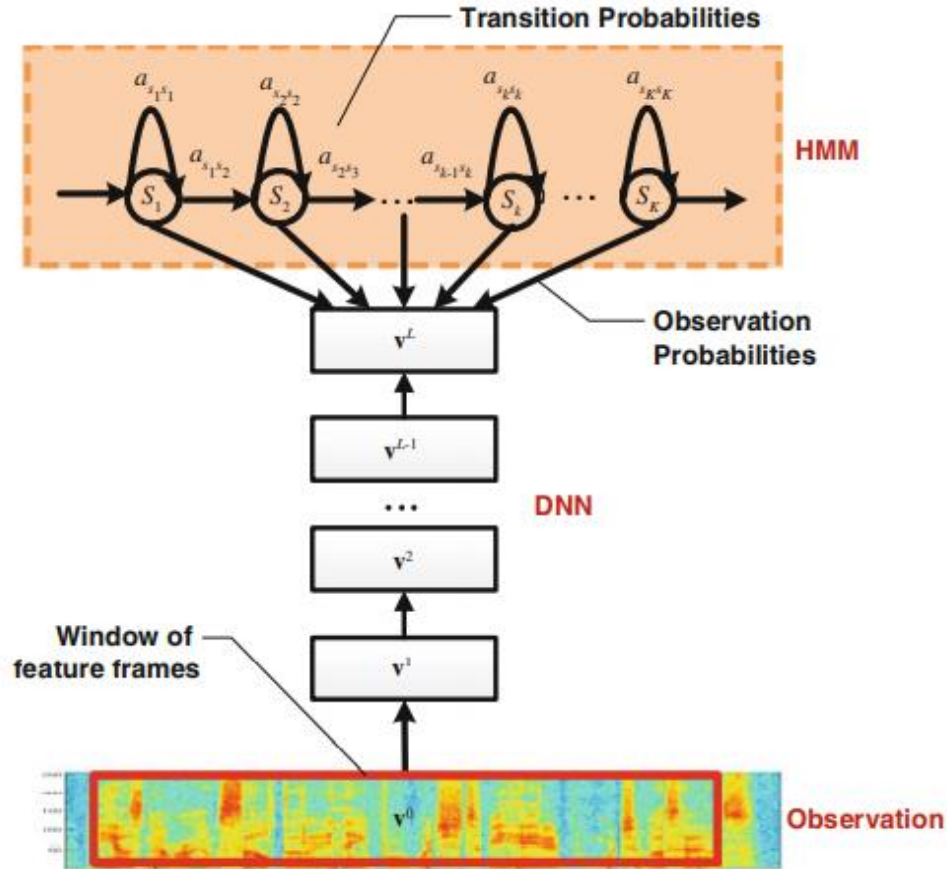
3.3. Получение последовательности слов (лингвистическая модель)



Извлечение признаков аудиосигнала



Лингвистическая и акустическая модели



$$\hat{w} = \operatorname{argmax}_w p(w|x) = \frac{\operatorname{argmax}_w p(x|w)p(w)}{p(x)} = \operatorname{argmax}_w p(x|w)p(w)$$

Акустическая модель

$$p(\text{Acoustic}|\text{HMMState}) = \frac{p(\text{HMMState}|\text{Acoustic})p(\text{Acoustic})}{p(\text{HMMState})}$$

- $p(\text{Acoustic}|\text{HMMState})$ оценивается с помощью HMM
- $p(\text{HMMState}|\text{Acoustic})$ оценивается с помощью глубокой нейронной сети
- $p(\text{HMMState})$ – априорная вероятность получить каждое состояние, оцененное по тренировочным данным

Лингвистическая модель

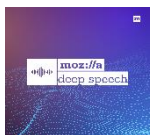
$p(w) = p(w_k|w_{k-1}, \dots, w_1)$ - Оценивается только по текстовому корпусу

Реализации Speech-To-Text моделей



Kaldi

- a. Отдельное использование акустической и лингвистической моделей
- b. В качестве лингвистической модели используется CD-HMM-DNN



DeepSpeech&DeepSpeech 2

- a. Основа модели – RNN слои
- b. Единый end-to-end подход



Wav2Letter

- a. Основа модели – CNN слои
- b. Единый end-to-end подход



EspNet

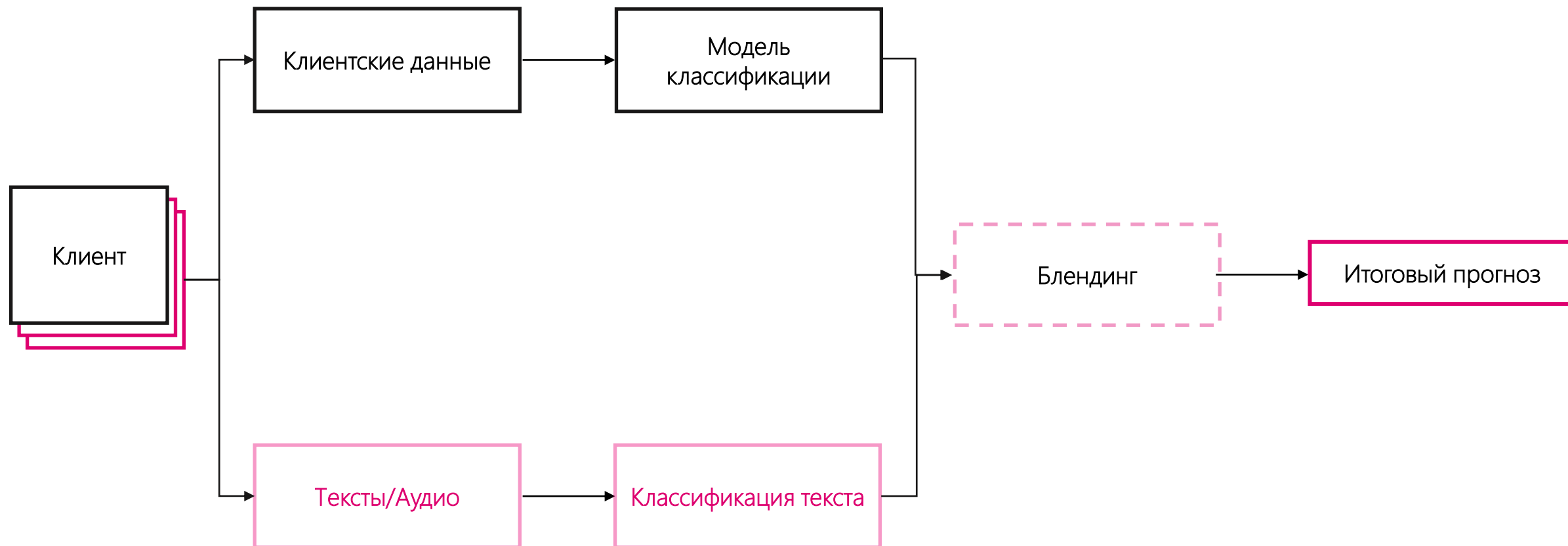
- a. Основа модели – Attention
- b. Единый end-to-end подход
- c. Предобработка сигнала аналогично Kaldi



QuartzNet

- a. Основа модели – 1d Channel Separable Convolutions
- b. Единый end-to-end подход

Обогащение банковских моделей



Обогащение моделей Роботом-Коллектором

Модель: ENTRY

Текстовые данные: пилот с роботом коллектором (аутсорс). Тексты звонков брались до даты скоринга.

	Train	Test	σ
Entry	49,0%	51,2%	
Entry + Robot	53,3%	53,6%	
Δ Gini	4,4%	2,4%	0,7%

Выводы:

По сравнению с Романом диалоги Робота-Коллектора информативнее за счет лучшей адаптации робота под речь оператора

Интерпретатор текста (ELI5)

robot здравствовать вы приветствовать автоматизированный помощник «банка ренессанс кредит». Информировать вас, данный разговор записывается фио я мочь услышать?

customer это я

robot мы звонить вы по повод ваш просрочить задолженность перед банк ренессанс кредит по денежный кредиту. На сегодняшний день у вы имеется задолженность в размер сумма_деньги вы готовый оплатить её в течение три дней?

customer три нет на на следующий неделя оплатить

robot фио возврат в график платёж позволять прекратить звонок и другой вид информирование по повод просрочить задолженность с сторона «банка ренессанс кредит». Вы готовый оплатить просрочить задолженность в течение три дней?

customer в три нет на следующий неделя оплатить в понедельник

robot фио обращать ваш внимание, что весь информация о возникновение просрочить задолженность фиксироваться в бюро кредитный история с указание срок и сумма просрочить задолженности. Вы готовый оплатить просрочить задолженность в течение три дней? К сожаление вы не слышный.

customer в понедельник в понедельник.

Примечание: текст диалога лемматизирован для моделирования

Обогащение моделей Outsource Collection звонками

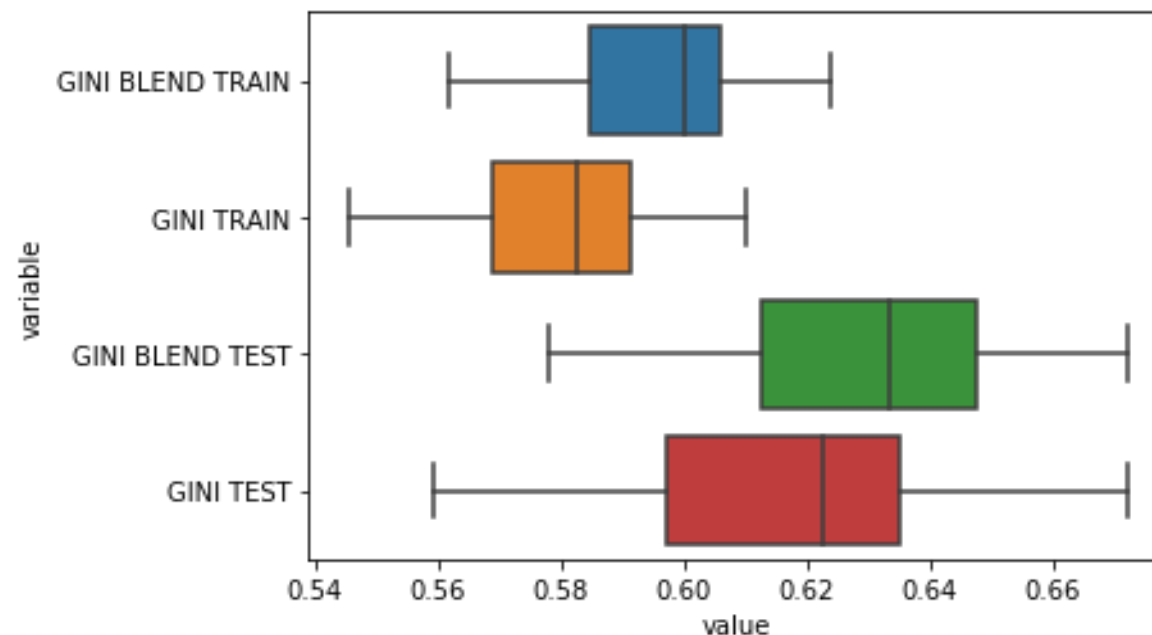
Модель: **SW**

Текстовые данные: записи outsource звонков Collection в течение года до даты скоринга

Blending Test Sample				
Показатель	Среднее	5%	Медиана	95%
GINI_без текстов	61.72%	56.98%	62.25%	66.05%
GINI_с текстами	63.05%	58.64%	63.35%	66.57%
Прирост GINI	1.33%	-0.16%	1.40%	2.96%

Выводы:

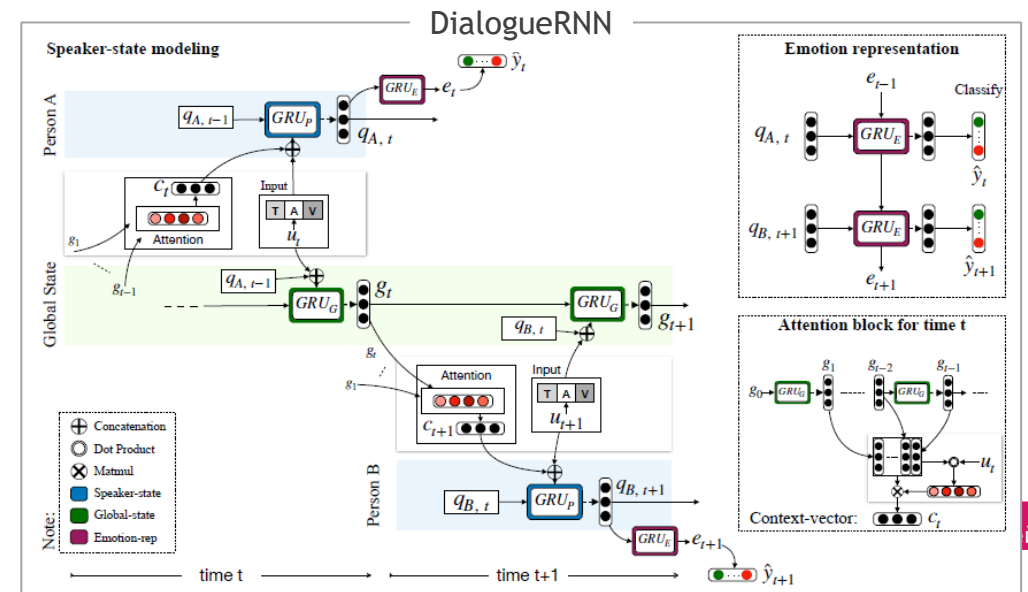
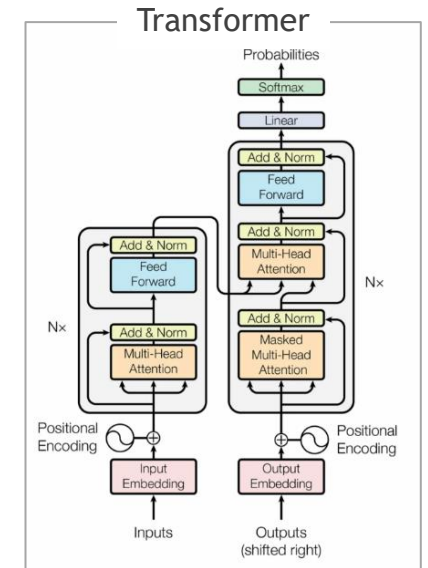
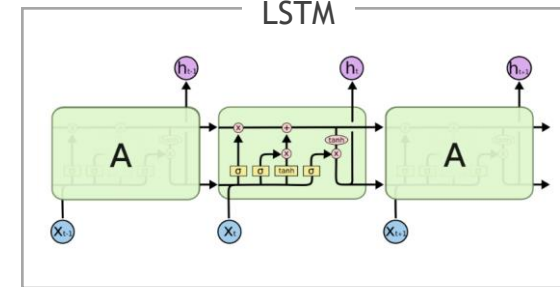
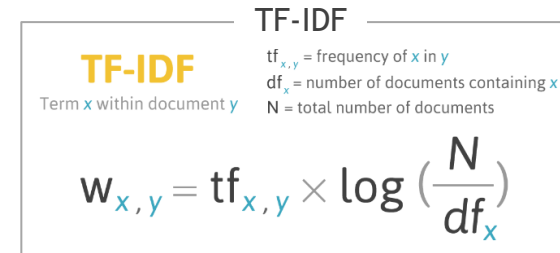
Предыдущие записи контакта клиента с Collection информативны для будущей классификации



Классификация с учетом диалоговой структуры

У оператора и клиента разные паттерны поведения => их надо **разделить!**

- Дополнительные обучающие токены, указывающие на актора
 - Transformers
 - RNN ячейки
- Разделение представлений фраз клиента и оператора
 - DialogueRNN
 - Реализации для Transformers могут быть дороги



Appendix

Примеры моделей банков из топ-1

Риски

1. Анализ тональности новостных текстов для заданного определения тональности
2. Задача выявления нерациональных сделок
3. WWR при расчете цены производных инструментов
4. Разработка моделей эволюции для макроэкономических параметров, релевантных для заданной области промышленности
5. Задача оптимизации процедуры выставления ограничений на торговый портфель банка
6. Модель для оценки надежности сложной системы
7. Задача поиска разладки в условиях ограниченного набора данных
8. Совместное моделирование рыночных риск-факторов при помощи копул.
9. Моделирование рыночных риск-факторов методом Монте-Карло с использованием стохастических процессов со скачками.

Риски ликвидности

1. Кластеризация банков по поведению клиентов и моделям ведения бизнеса
2. Разработка алгоритма определения банков с повышенным спросом на ликвидность
3. Моделирование рынка вкладов физических лиц в валюте
4. Моделирование рынка депозитов юридических лиц в валюте
5. Разработка рекомендательной системы по оптимальным периодам первичного размещения облигаций
6. Определение справедливой ставки привлечения по вкладам физических лиц для различных банков
7. Определение оптимальной структуры розничного портфеля
8. Разработка неценовых инструментов управления клиентским балансом банка на рынке с жестким государственным регулированием
9. Оценка влияния развития sharing есопому на потребительский спрос и вероятность реализации кризиса перепроизводства

Бизнес

1. Извлечение данных из комментариев к объявлениям
2. Построение маршрута с мин. затратами и макс. Эффектом
3. Оценка качества и потенциала бизнеса клиента банка с применением анализа неструктурированных текстов
4. Предсказание совершения покупки физлицом на основе его транзакций
5. Определение поведенческой модели пользователя на основе его трат
6. Определение контекста с предсказанием следующих действий пользователя
7. Оценка качества и потенциала бизнеса клиента банка с применением NLP инструмента
8. Извлечение ценовых параметров из объявлений о продаже/аренде недвижимости
9. Textual Entailment в задаче классификации документов
10. Анализ эффективного использования ценовых участков города

ИТ

1. Оптимизация хранения текстовых данных
2. Надежность, Smart-monitoring (временные ряды, построение регрессий и доверительных интервалов, NLP) – выявление отклонений в работе систем на основе прогнозных значений метрик. Кластеризация, локализация и выявление причин инцидентов, разбор жалоб, анализ логов
3. Автоматизация сервисных запросов внутреннего клиента (NLP, многомерная классификация) – Анализ сервисной заявки (комментарии пользователей и др.), маршрутизация на группы исполнителей, реком. системы, автоисполнение обращений
4. Управление мощностями, запасами, ресурсами (регрессионный анализ, генетические алгоритмы и др)
5. Управление архитектурой, качеством данных и справочной информацией (NLP, CNN и др.)
6. Создание чат-ботов для сервисов внутреннего клиента (NLP)

Collection-модели банков из топ-1

Предконтрактная стадия	Верификация и обогащение данных	Робот-коллектор	Soft Collection	Hard / Legal / Executory Collection
<ol style="list-style-type: none">1. Модели для оценки цессионных портфелей2. Модели для workforce management3. Модели для resource management	<ol style="list-style-type: none">4. Модели для распознавания captcha при автоматизированном сборе данных5. Модели для оценки вероятности найти информацию при помощи skip tracing6. Модели для оценки вероятности принадлежности найденного контакта должнику7. Модели для мэтчинга профилей физлица в различных онлайн-источниках	<ol style="list-style-type: none">8. Модели для определения наличия голоса в звуковом канале9. Модели для распознавания речи10. Модели для выбора направления движения по скрипту общения робота в зависимости от слов должника	<ol style="list-style-type: none">11. Модели для оценки мат.ожидания Recovery Rate12. Модели для оценки вероятности дозвона по телефонному номеру13. Модели для оценки ожидаемого финансового эффекта от звонка14. Модели для прогноза длительности звонка15. Модели для оценки вероятности получить обещание об оплате16. Модели определения наличия автоответчика17. Модели для определения результата переговоров с должником18. Модели для рекомендации оптимальных мотиваторов19. Разговорные модели для текстовых коммуникаций20. Модели для динамического управления процессами взыскания на базе RNN	<ol style="list-style-type: none">21. Модели для оценки мат.ожидания Recovery Rate22. Модели для оценки финансового эффекта от выезда сотрудника выездного взыскания

Примеры внедрений в банках

Post Mail Classifier



В банке Хоум Кредит построили AI, который позволяет разделять почтовую корреспонденцию по тематикам.

Технология позволяет распознавать печатные сканы писем, разделять их на 250+ тематик и формировать автоматический ответ на распознанные письма.

Для распознавания текста использовали **Open Source** технологию OCR: программу **Tesseract** от компаний Google.

End-to-end время: **10-40 сек**

Качество распознавания печатных текстов: **78%**

В основном тексты «юридического стиля»

HOME CREDIT
MEMBER OF
THE PPF GROUP

Про проект

Задача:
Ускорить процесс анализа бумажной корреспонденции
Понять набор тематик, определить обработчика
Сформировать ответ автоматически

Поставим OCR перед NLU?

Текущий статус:
WIP

18

2018. Frankenstein Forum & Awards. The Retail Finance.

Примеры внедрений в банках

Голосовая биометрия



В банке из топ-26 разработали свою собственную голосовую биометрию, которую внедрили на процессах Колл-центра и Коллекшн (идентификация клиентов).

Рабочая группа: ~10 человек

Срок разработки: ~1,5 года

Запуск: 2018 г.

Боль:

Потребовались дорогие доработки ИТ-инфраструктуры.



Примеры внедрений в банках

Чат-боты



В рамках конференции SDSJ-2017 Сбербанк провел соревнования по разработке чат-бота. Призеры получили предложения работать в лаборатории AI Сбербанка.



В банке Хоум Кредит разрабатывают чат-бота на основе AI, который расширяет возможности «кнопочного» чат-бота. Разработка находится на стадии пилотирования.



В банке Тинькофф разработали чат-бота, который выдает подсказки на вопросы клиентов. В диалоге пока используют живых операторов, которым чат-бот также помогает с подсказками для ответов.

Sberbank
Data Science Journey

A. Определение релевантности вопроса

Бинарная классификация. Можно ли по парам из параграфа текста и заданным по нему вопросам определить, какой из вопросов настоящий и был задан человеком?

Сложность: средняя
Целевая метрика: ROC-AUC
Формат решения: офлайн разметка тестовых данных

В. Построение вопрос-ответной системы

Диалоговая система. Можно ли по паре из параграфа текста и поставленным по нему релевантным вопросам, найти в параграфе точный ответ и вывести его пользователю?

Сложность: высокая
Целевая метрика: (Macro-averaged) F1-score
Формат решения: Docker контейнер с обученной моделью

2017. Sberbank Data Science Journey.

Книги на русском языке



Глубокое обучение: погружение в мир нейронных сетей

Николенко С., Кадурин А., Архангельская Е.

СПб.: Питер, 2018. — 480 с.



Библиотека Keras — инструмент глубокого обучения

Антонио Джулли, Суджит Пал

ДМК-Пресс, 2017. — 296 с.



Глубокое обучение

Ян Гудфеллоу, Иошуа Бенджио, Аарон Курвилль

ДМК-Пресс, 2018. — 652 с.

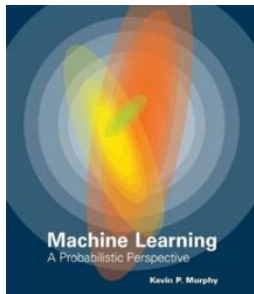


Введение в машинное обучение с помощью Python

Мюллер А., Гидо С.

O'Reilly Media, 2017. — 392 с.

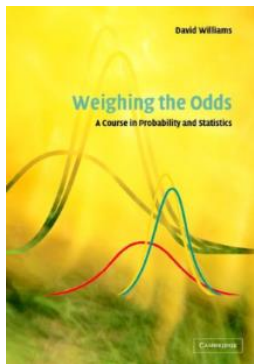
Книги на английском языке



Machine Learning: A Probabilistic Perspective

Murphy K.P.

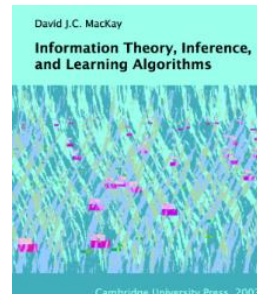
Massachusetts Institute of Technology, 2012. — 1067 p.



Weighing the Odds: A Course in Probability and Statistics

David Williams

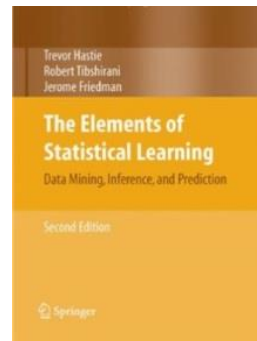
Cambridge University Press, 2001. — 568 p.



Information Theory, Inference, and Learning Algorithms

MacKay D.J.C.

Cambridge University Press, 2003. — 640 p



The Elements of Statistical Learning: Data Mining, Inference, and Prediction

Hastie T., Tibshirani R., Friedman J.

Second Edition (corrected 5th printing). — Springer, 2009. — 763 p.

Лекции и полезные ресурсы

https://www.youtube.com/playlist?list=PLJOzdkh8T5kp99tGTEFjH_b9zqEQiiBtC

Лекции ШАД от **Константина Воронцова** (д.ф.м.н. МФТИ, ВШЭ)

<https://www.lektorium.tv/speaker/2691>

Лекции **Сергея Николенко** (автор книги «Глубокое обучение: погружение в мир нейронных сетей»)

<http://efimov-ml.com/>

Сайт с лекциями **Дмитрия Ефимова**, а также с Python-скриптами (Jupyter Notebooks), презентациями и ссылками

<http://deepbayes.ru/2017/>

Лекции ВШЭ по Байесовским методам в глубинном обучении.
Лекторы: **Дмитрий Ветров**, **Дмитрий Кропотов**, **Евгений Соколов**, **Сергей Бартунов**, **Арсений Ашуха** и др.

<https://arxiv.org/>

Ресурс с самыми свежими препринтами научных статей по ML/AI/DS (см. раздел Computer Science)

<http://scikit-learn.org/>

<http://www.numpy.org/>

<https://pandas.pydata.org/>

<https://keras.io/>

<http://devdocs.io/tensorflow/>

Подробные мануалы популярных библиотек в Python

<https://stackoverflow.com/>

Ресурс для обсуждения практических вопросов по программированию на Python и других языках

<http://ods.ai/>

Русскоязычное сообщество датасайнтистов