

Архитектура "трансформер" для анализа рисков группового влияния нутриентов пищи

**Баландина¹, Груздев¹, Савельев², Будамян¹,
Кисиль³, Грачев¹**

1) Физический факультет МГУ имени М. В. Ломоносова

2) ПАО "Вымпелком"

3) Институт русского языка и культуры МГУ имени М. В. Ломоносова

Задача

Актуальность

- Метаболический синдром (МС) - медицинское состояние, которое характеризуется висцеральным ожирением. Основной причиной МС является неправильный образ жизни. МС приводит к сердечно-сосудистым и прочим заболеваниям.
- Раннее выявление и лечение метаболического синдрома может снизить риск развития серьезных заболеваний в будущем.
- Ранняя диагностика нарушений становится одной из ключевых задач в обследовании пациентов с подозрением на МС.

Задача

Данные

- Центром Нутрициологии и Адаптивного Питания были предоставлены анонимизированные данные анамнезов наблюдаемых пациентов.
- Данные были формализованы и приведены в табличном виде, где каждому анамнезу пациента соответствует строка таблицы.
- Данные были предобработаны и очищены от дубликатов. Реальные анамнезы были разбиты на классы: имеет диагноз, связанный с метаболическим синдромом / здоровый.
 - 2885 строк
 - Данные разбиты 70% / 30% (обучающая / тестовая выборки)
 - 85 признаков, 67 – состав диеты, 18 – антропометрия и проч.
 - Баланс классов 1872 / 1013 (~65% / 35%, здоровые / больные)

Возраст	Рост	Профессия	Спорт	ОТ	ОБ
39.0	180.0	Работники преимущественно умственного труда	Не занимаюсь	130.0	120.0
45.0	177.0	Работники преимущественно умственного труда	Регулярный спорт	0.0	0.0
56.0	181.0	Работники среднего по тяжести труда	Не занимаюсь	50.0	150.0
33.0	184.0	Работники среднего по тяжести труда	Регулярный спорт	106.0	101.0
35.0	182.0	Работники преимущественно умственного труда	Легкий спорт	0.0	0.0
31.0	0.0	Работники среднего по тяжести труда	Не занимаюсь	0.0	0.0

Таблица с реальными данными пациентов.

Задача

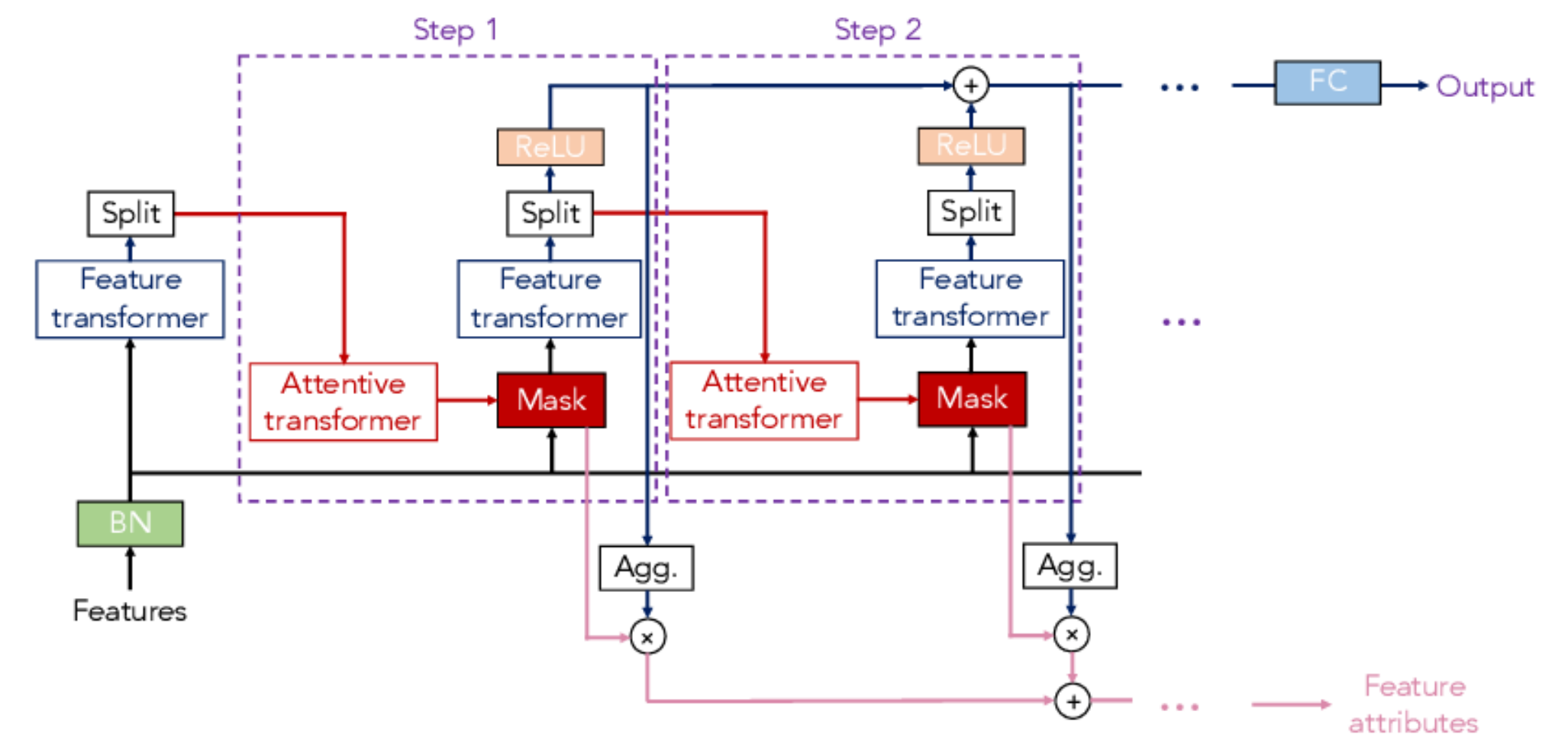
Цель работы и этапы выполнения

- Задача **бинарной** классификации
 - Обучить нейронную сеть оценивать вероятность наличия метаболического синдрома на основе признаков;
 - Оценить точность предсказаний;
 - Использовать обученную модель для того, чтобы найти группы признаков, влияющих на результат.

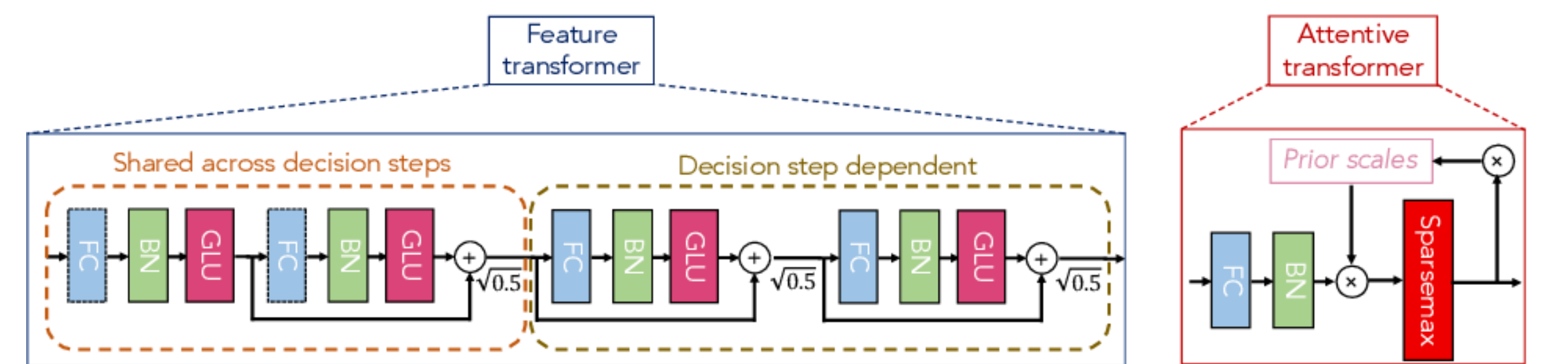
Предыдущие работы

Предыдущие работы по нейросетевым методам обработки табличных данных

TabNet^[1] — рекуррентная нейронная сеть, разработанная для анализа табличных данных.



(a) TabNet architecture



(b) Feature transformer

(c) Attentive transformer

[1]Arik, S. O., & Pfister, T. (2019). TabNet: Attentive Interpretable Tabular Learning. arXiv preprint arXiv:1908.07442.

Предыдущие работы

Предыдущие работы по нейросетевым методам обработки табличных данных. TabNet. Предобучение

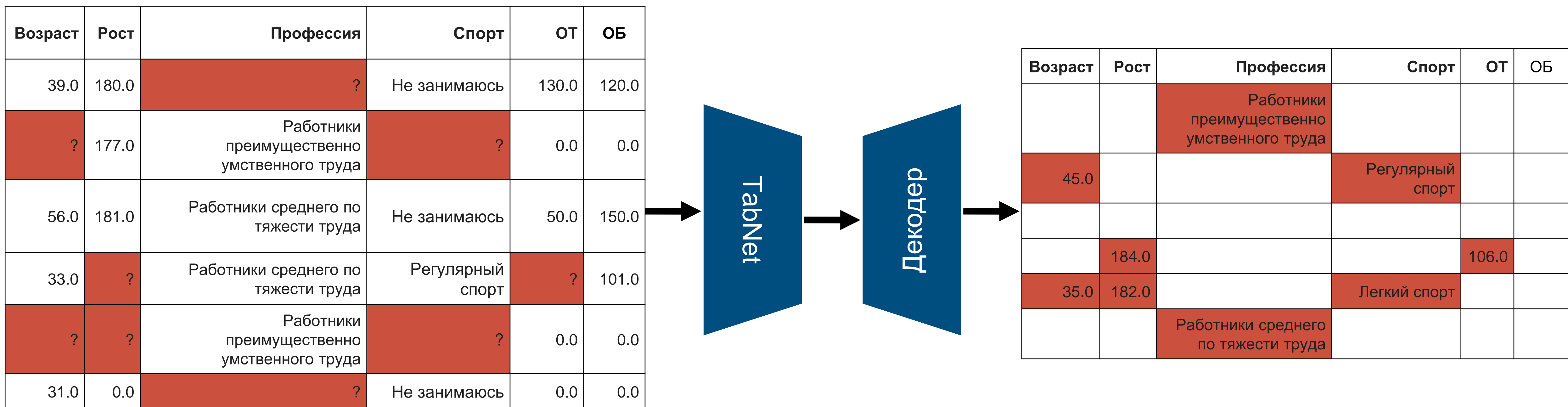


Схема самообучения TabNet

Предыдущие работы

Предыдущие работы по нейросетевым методам обработки табличных данных. TabNet. Дообучение

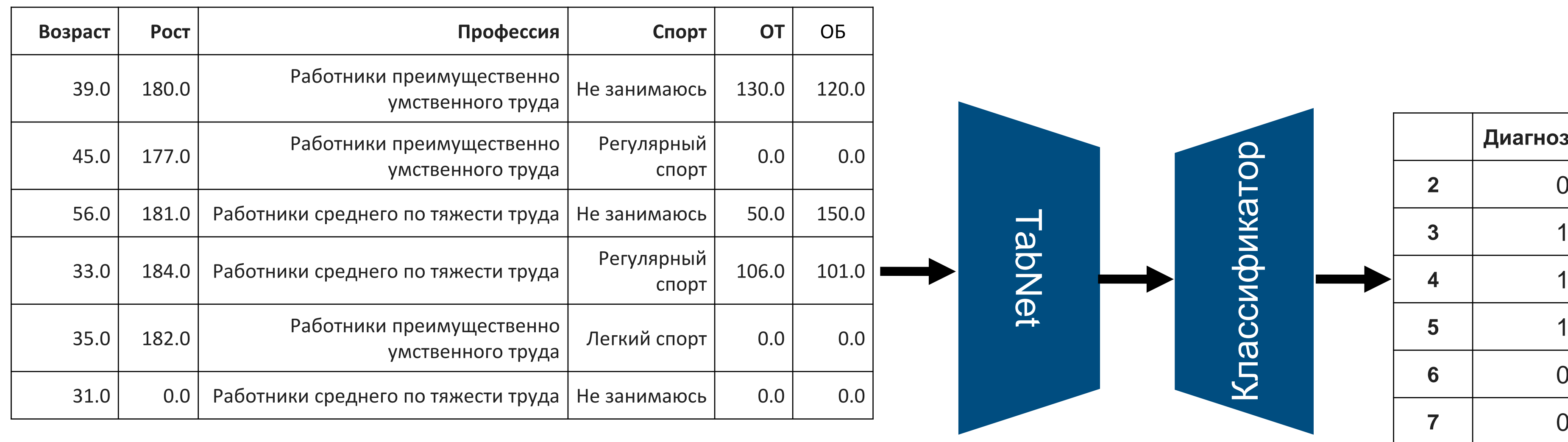
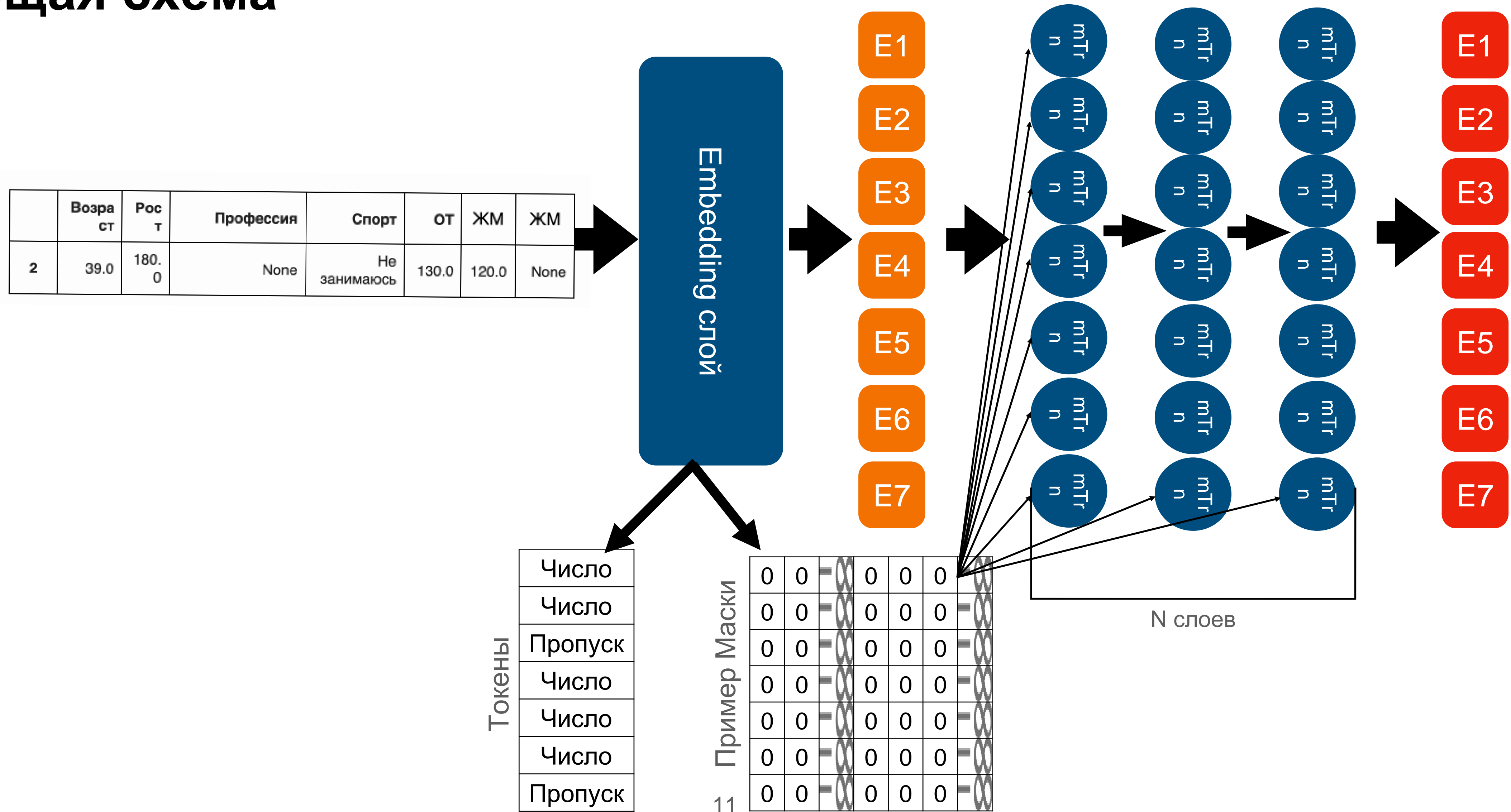


Схема дообучения TabNet

Разработанная архитектура

Общая схема



Разработанная архитектура

Embedding слой

- A. В качестве Embedding слоя для чисел взят слой PAF^[4] (Periodic Activation Functions): $f(x) = Periodic(x) = concat[sin(v), cos(v)]$, где $v = [2\pi c_1 x, \dots, 2\pi c_k x]$
- B. Позиционное кодирование, необходимое, чтобы модель понимала к какому признаку относится embedding.
- C. Производится по следующей формуле:
 $PE_{pos,2i} = sin(pos/1000^{2i/d_{model}})$,
 $PE_{pos,2i+1} = cos(pos/1000^{2i/d_{model}})$ [5]

[4] Yury Gorishniy and Ivan Rubachev and Artem Babenko(2022). On Embeddings for Numerical Features in Tabular Deep Learning. arXiv 2203.05556

[5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems, pages 6000–6010.

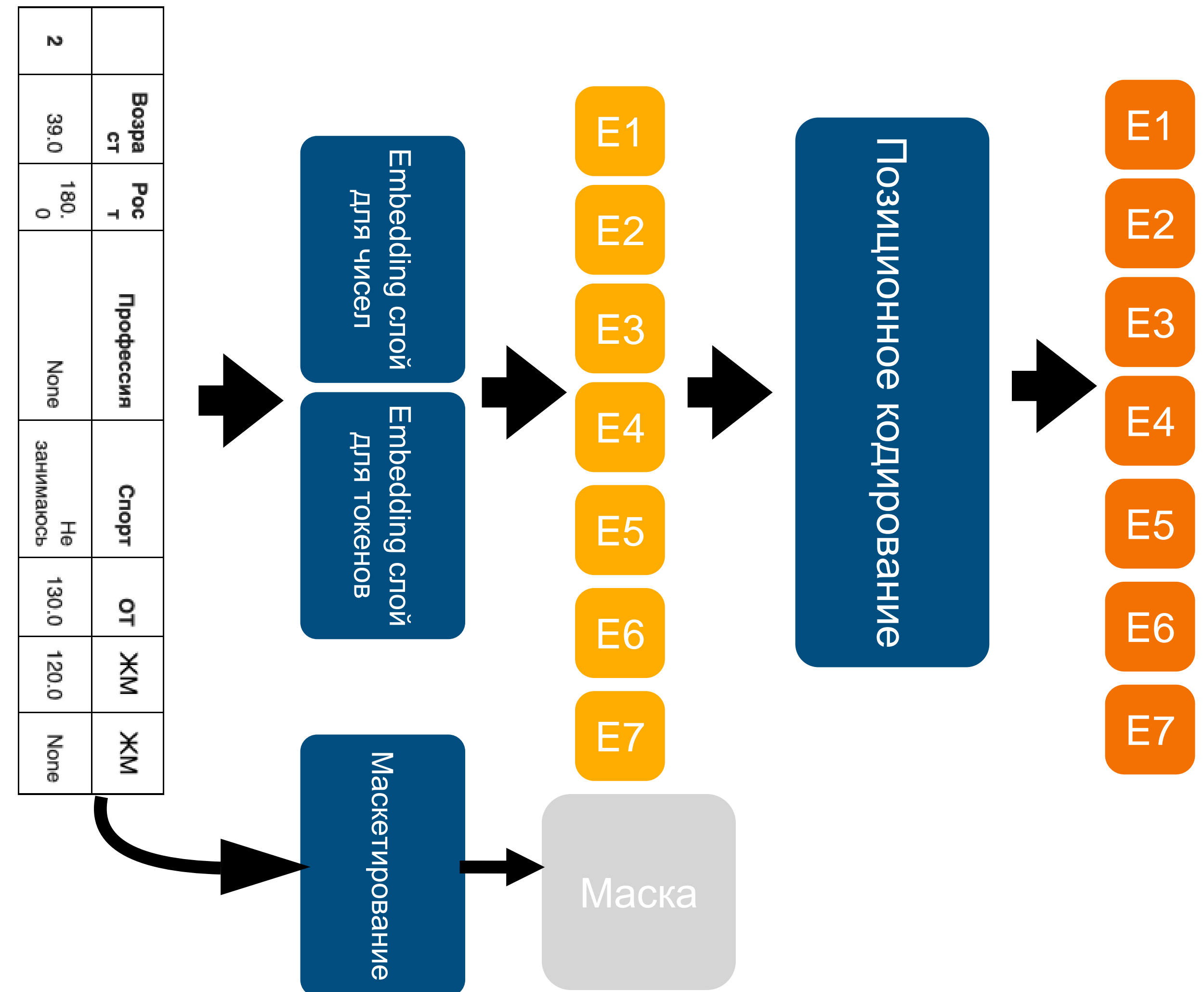


Схема работы Embedding слоя

Разработанная архитектура

Слой трансформера

Главное в трансформере - это механизм внимания (Attention), с помощью которого входные вектора представляются через все прочие вектора (получают контекст).

$$\text{Attention}(Q, K, V, \text{MASK}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_{\text{model}}}} + \text{MASK}\right)V$$

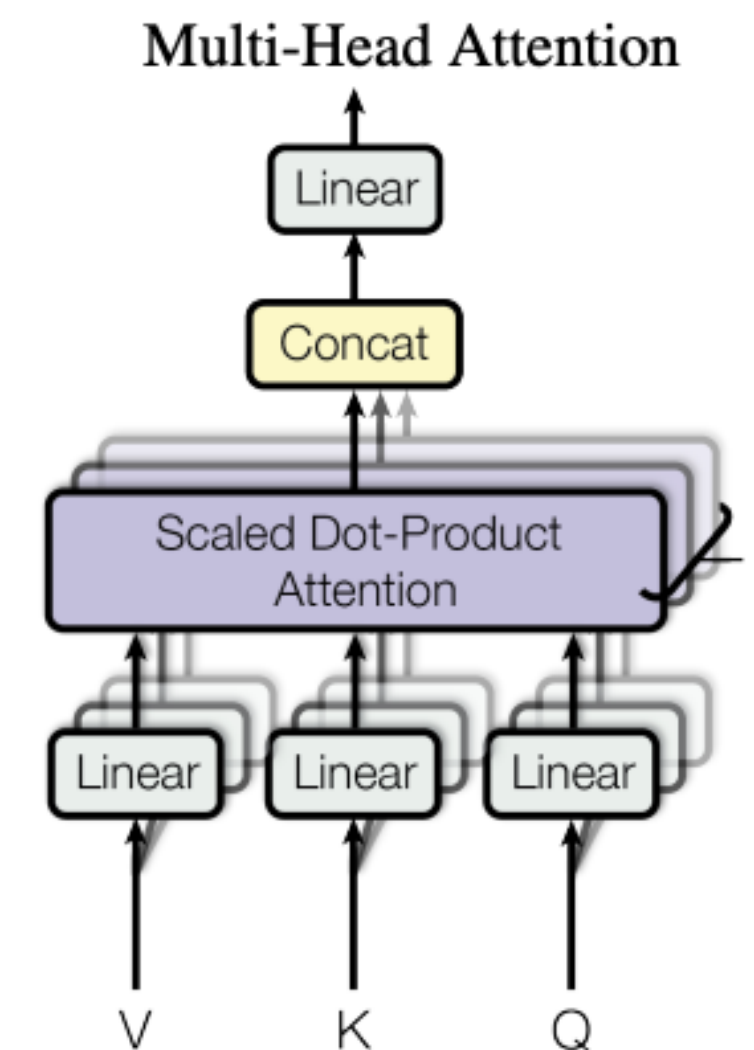
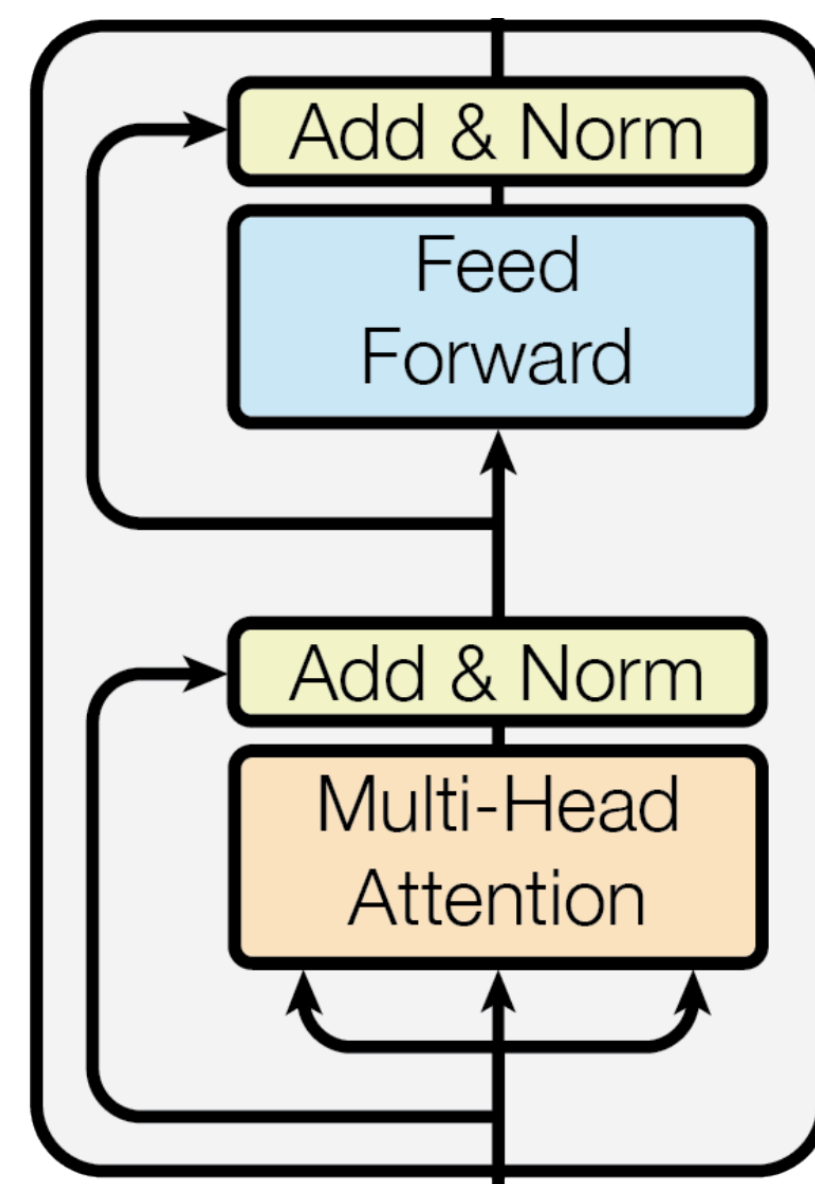
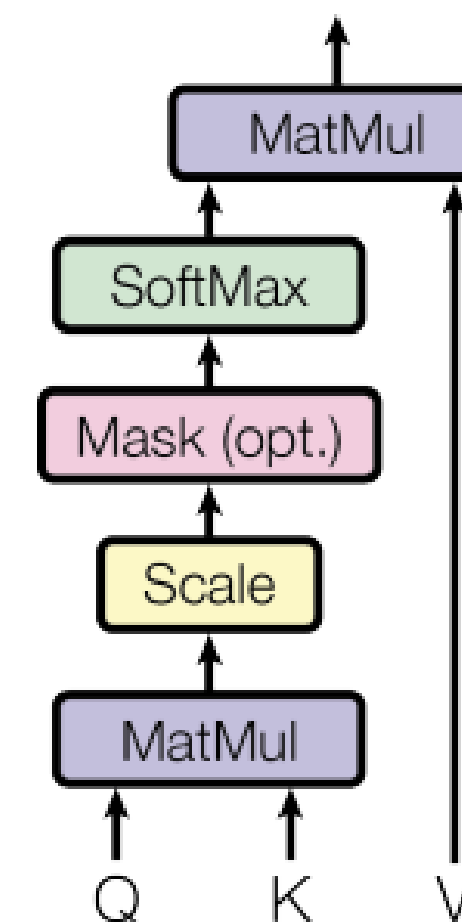


Схема работы трансформер-слоя

Scaled Dot-Product Attention



Разработанная архитектура

Предобучение

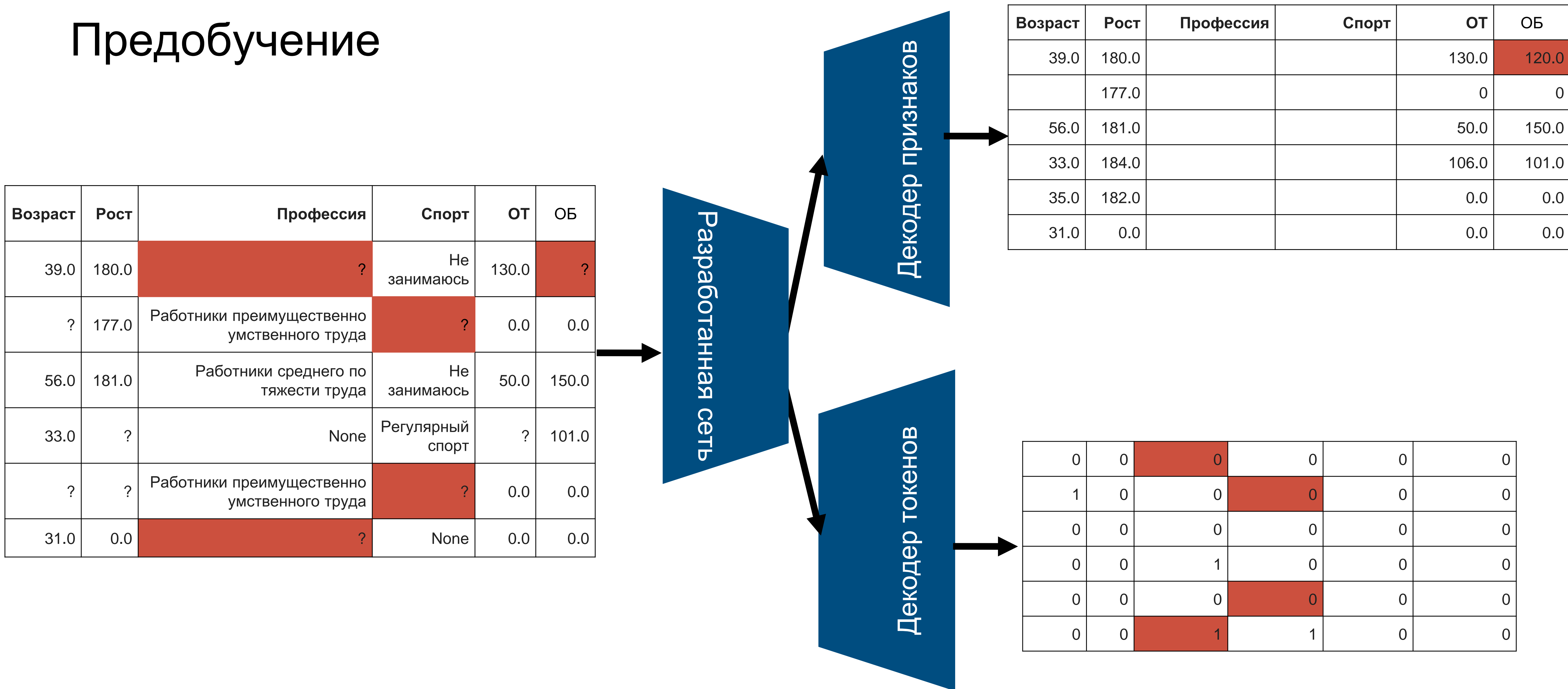


Схема предобучения разработанной нейронной сети.

Разработанная архитектура

Обучение

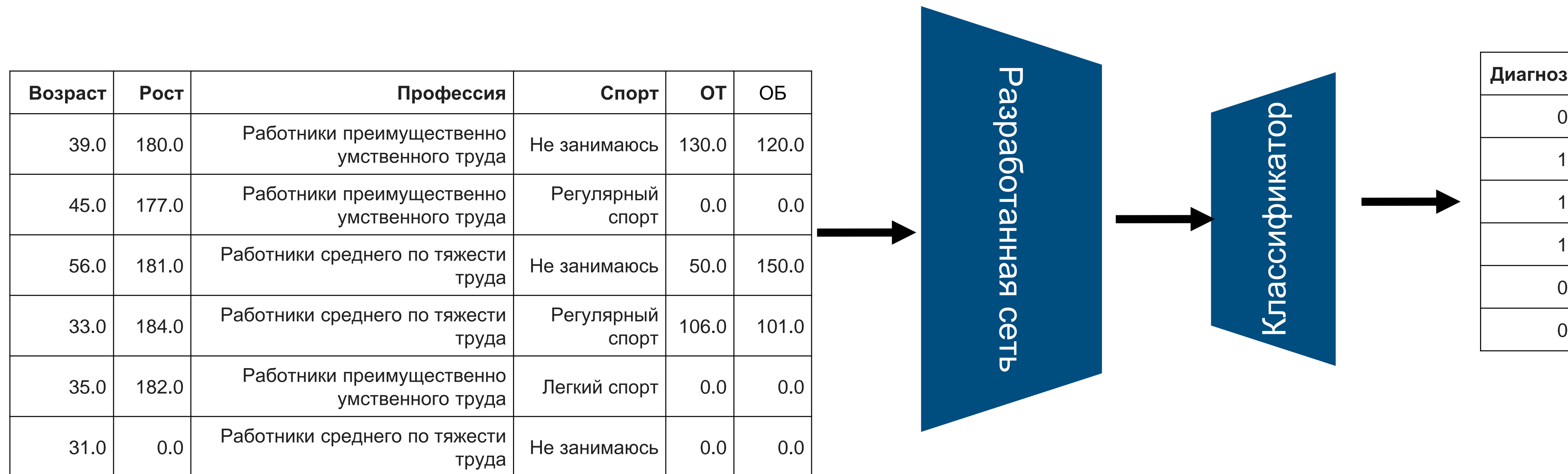


Схема дообучения разработанной нейронной сети.

Результаты

1. При обучении моделей различных архитектур были получены следующие показатели ассигасы на тестовой выборке

- $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$
- TP – количество верно найденных диагнозов с метаболическим нарушением;
- FP – количество раз, когда модель считала, что нарушение есть, однако его не было;
- TN – количество верно найденных диагнозов с отсутствием метаболических нарушений;
- FN – количество раз, когда модель считала, что нарушения нет, однако оно было.

Разработанная сеть	TabNet	CatBoost
72%	68%	74%

Значения Ассигасы на тестовой выборке исследуемых моделей

Марковская кластеризация

- 1) Граф преобразуется в смежную матрицу;
 - 2) Смежная матрица преобразуется в стохастическую матрицу;
 - 3) Пока матрица изменяется повторять:
 - a) Расширение матрицы;
 - b) Инфляция матрицы.
- Хорошо масштабируется с увеличением размера графа.
 - Работает как с взвешенными, так и невзвешенными графами.
 - Выдает хорошие результаты кластеризации. Устойчив к шуму в данных графа.
 - Количество кластеров не задается заранее, но можно настроить гранулярность кластера с помощью параметров.
 - Не подходит для кластеров с большим диаметром.

Результаты

В ходе проведения исследования были кластеры из слоя Attention.

Нездоровые	Здоровые
Масса тела, холестерин, насыщенный жир, магний	Возраст, рост, общая жидкость
Рост, объем талии, галактоза	Возраст, рост, жировая масса
Объем талии, жировая масса, общая жидкость, изолейцин, гамма-линоленовая кислота	Пол, рост, объем талии, холестерин
Объем талии, жировая масса, кремний, изолейцин	Масса тела, насыщенные жиры, сахара, магний
Объем талии, мононенасыщенные жиры, бета-каротин	Насыщенные жиры, сахара, углеводы, фолаты
Объем бедер, жировая масса, мышечная масса, спорт	Насыщенные жиры, сахара олеиновая кислота
Объем талии, жировая масса, йод	Насыщенные кислоты, сахара, лейцин
Жировая масса, мышечная масса, общая жидкость, изолейцин, никель	Натрий, жиры, β -каротин
Жировая масса, мышечная масса, общая жидкость, изолейцин, фтор	Фолаты, омега-3, биотин
Жировая масса, мышечная масса, медь	Фолаты, омега-3, фенилаланин+тирозин
Жировая масса, мышечная масса, молибден	β -каротин, фосфор, метилонин, отсутствие спорта, ИМТ

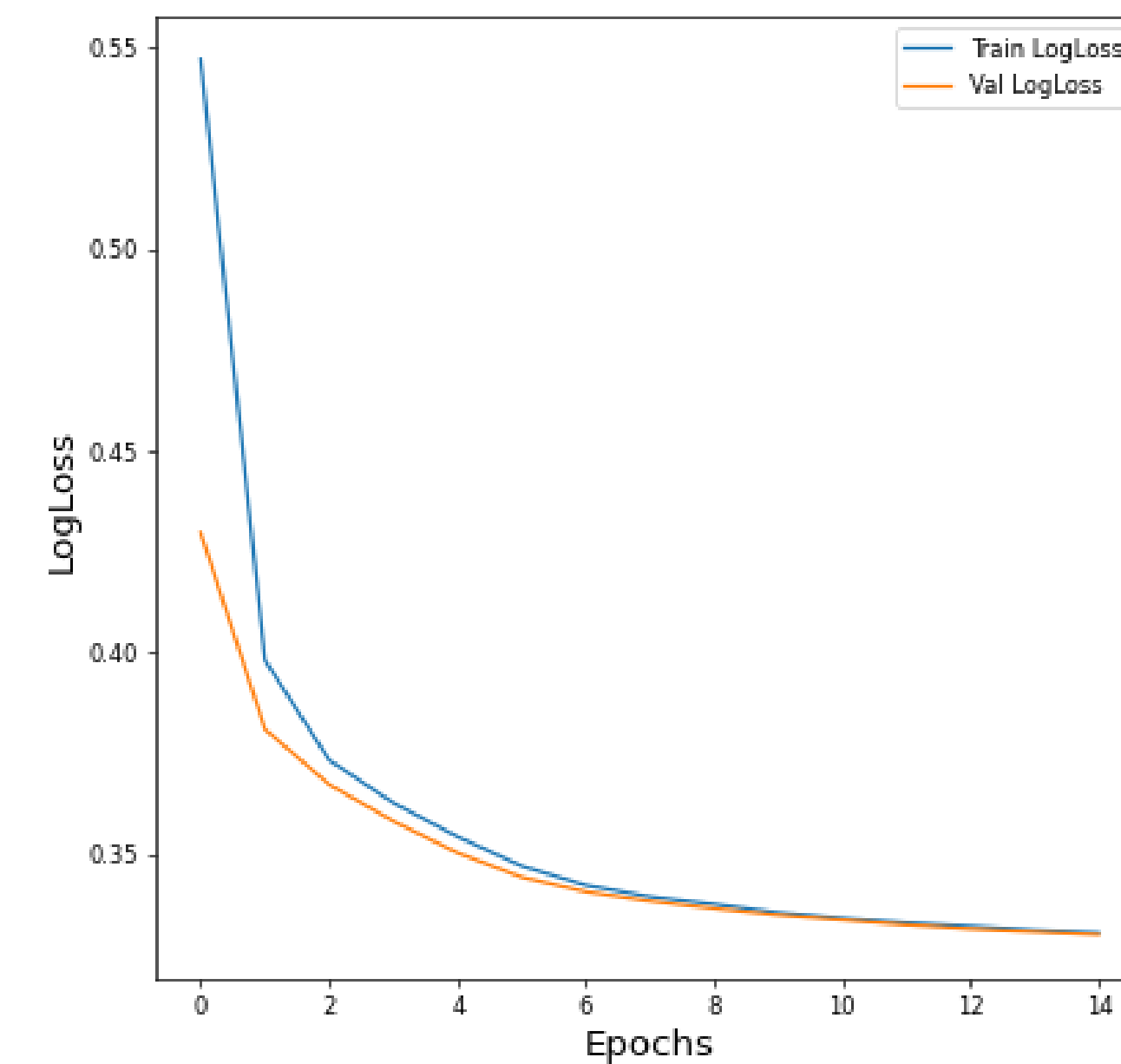
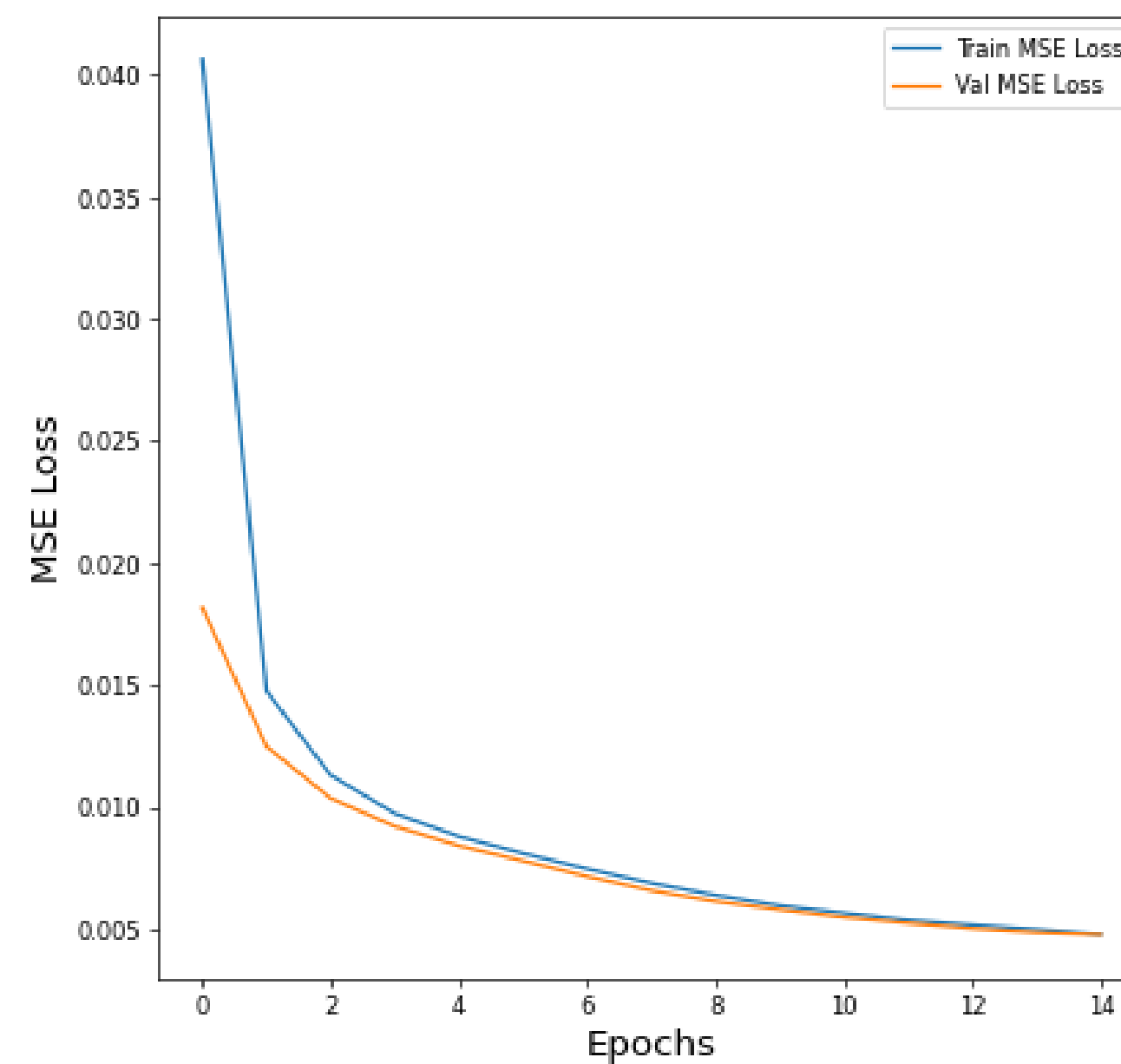
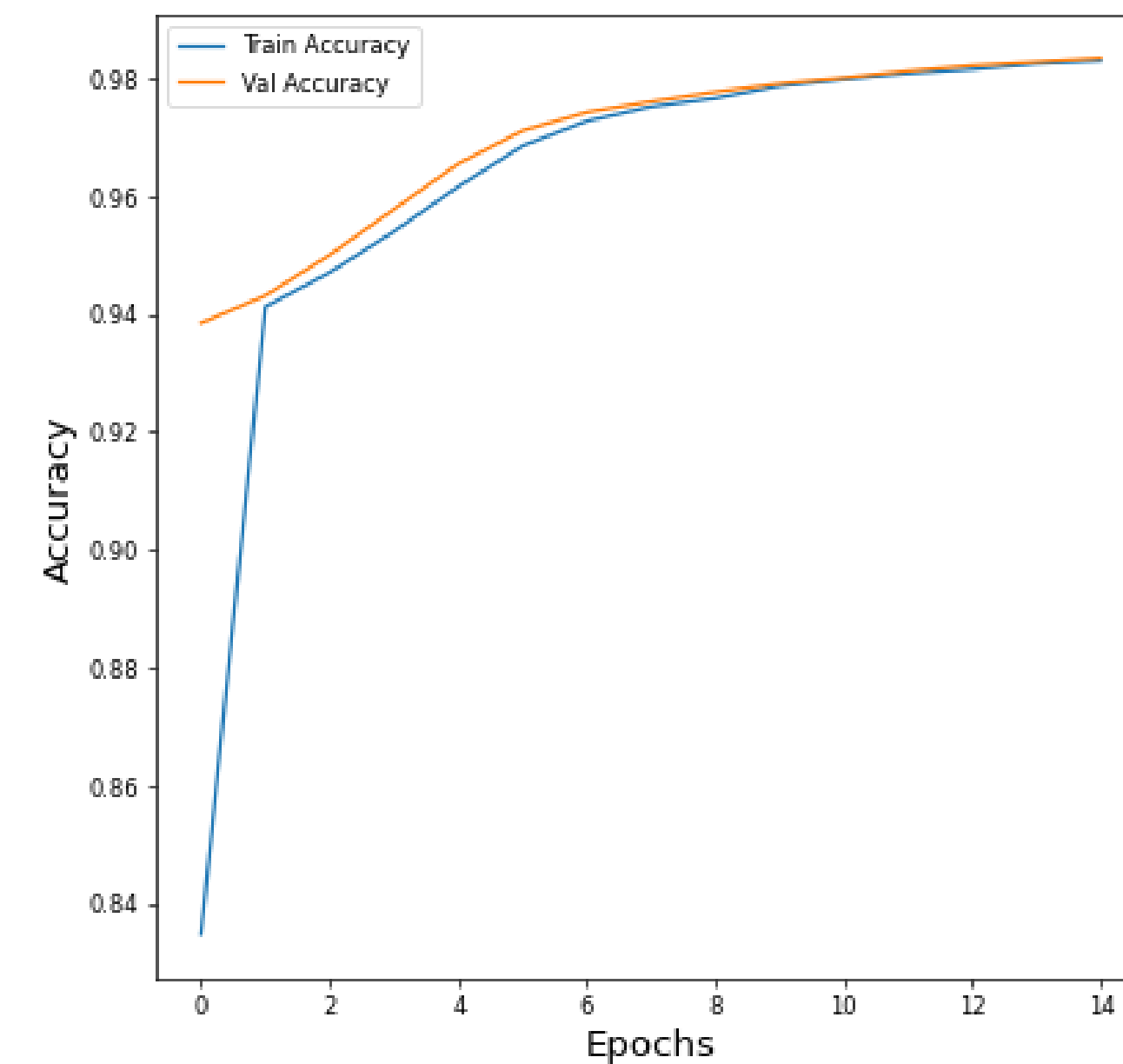
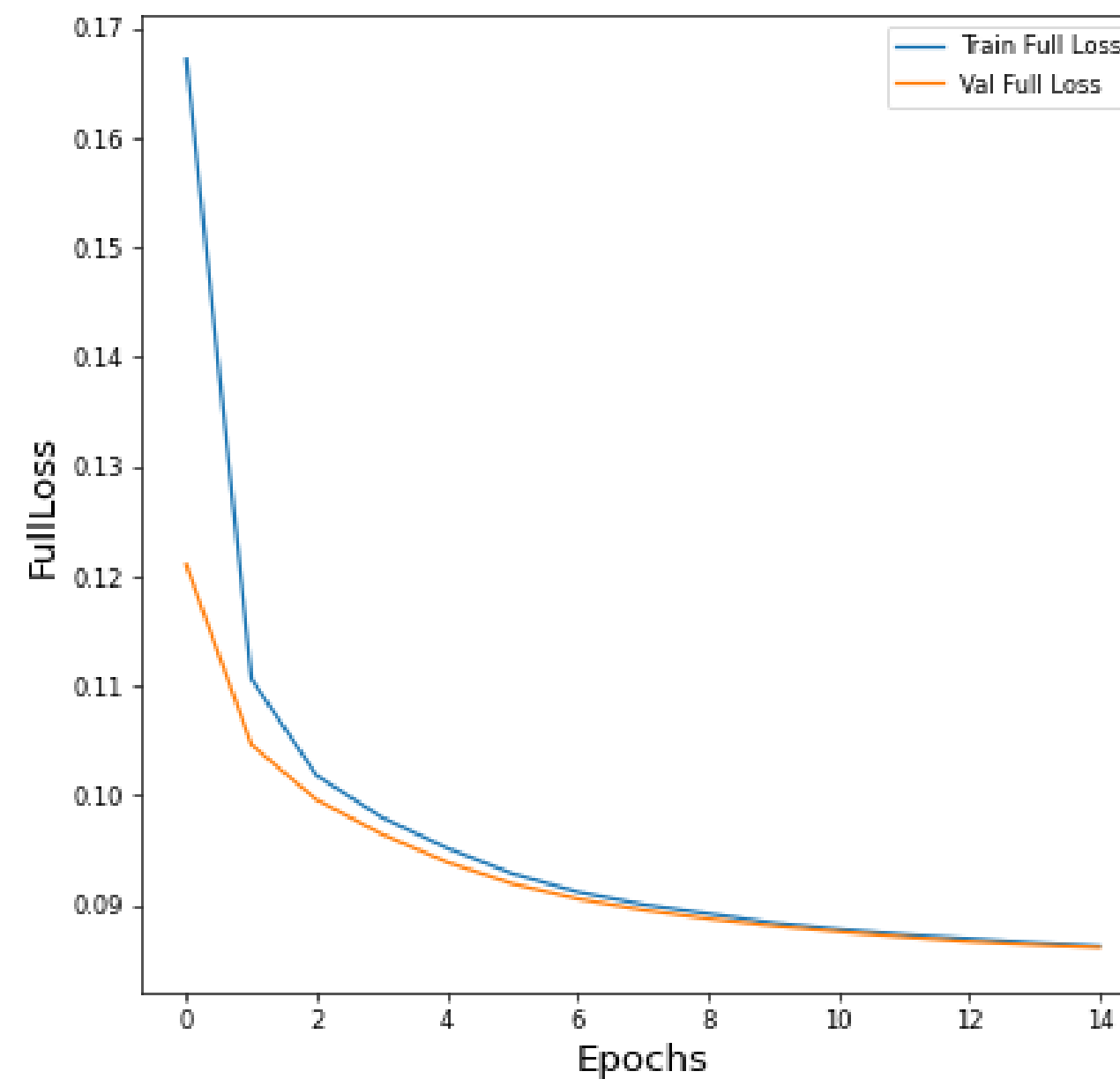
Полученные кластеры для людей с/без метаболического синдрома

Результаты

1. Обучена нейронная сеть архитектуры трансформер на задаче бинарной классификации;
2. Обученная сеть сравнима по качеству с другими, используемыми при задаче бинарной классификации;
3. При помощи марковской кластеризации из слоя внимания сети получены кластеры признаков для людей, классифицированных как имеющих метаболический синдром, так и для людей, классифицированных как не имеющих метаболический синдром.

Графики обучения

Графики самообучения



Графики обучения

Графики дообучения

