

Разбор алгоритмов:

1. «Наивный алгоритм»

С ним все просто: мы группируем данные по пользователю и считаем кол-во купленных за всю историю продуктов. В результате будем иметь для каждого пользователя такую картину:

	product_id	amount
0	33754	17
1	21903	14
2	24838	14
3	17704	13
4	46667	13

Повторим эту операцию 100 000 раз (для каждого пользователя) и преобразовав к итоговому виду, получим скор 0.27841.

Это будет наш эталон, к нему мы будем стремиться, однако мы не можем назвать такой подход рекомендательной системой, ведь мы ничего нового не предлагаем пользователю, а принимаем в учет только те товары, про которые он сам прекрасно знает и, более того, любит, т к это его топ.

2. Implicit ALS (Alternating Least Squares)

Хорошая удобная библиотека, главная фишка заключается в том, что можно прямо «из коробки» передать топ из скольки товаров мы хотим порекомендовать каждому покупателю. Перед тем, как обучать алгоритм есть 2 главных параметра его настройки: т н факторы (factors) и итерации (iterations). От различных комбинаций этих двух параметров варьируется скор. Однако это изменение не столь существенное, если сравнить его с влиянием, что оказывает подготовка данных: категориальное кодирование. (~0.006 против ~0.025). Было опробовано несколько комбинаций факторов и итераций. Оптимальными показали себя factors=30, iterations=8. Итоговый скор с ними 0.05377. Это существенно НЕ ДОТЯГИВАЕТ до заданных целей.

3. Implicit kNN (k Nearest Neighbours)

Метод из той же библиотеки, но основанный на другом принципе: он сравнивает величину с k ближайшими соседями для выявления сходств. Такой подход сработал лучше всего при 10 соседях, дав 0.12458 точности предсказаний. Однако до заданной планки это также существенно НЕ ДОТЯГИВАЕТ.

4. LightFM

Другая библиотека, также основанна на принципе разреженных матриц, однако при создании модели позволяет добавлять вспомогательные фици для юзеров и айтемов, представляющих строки и столбцы такой матрицы. Кроме этого можно настраивать веса образцов, задавать функцию потерь и выбирать количество эпох обучения, совсем как в нейросетях. Система однако, может работать и по одной лишь матрице, только вот скор такая система по метрике MAP@10 дает наихудший, практически нулевой. Экспериментировать с ней дальше я не стал, не увидел смысла в добавлении фицей, т к вряд ли бы они улучшили скор на несколько порядков. Вердикт: НЕ ДОТЯГИВАЕТ.

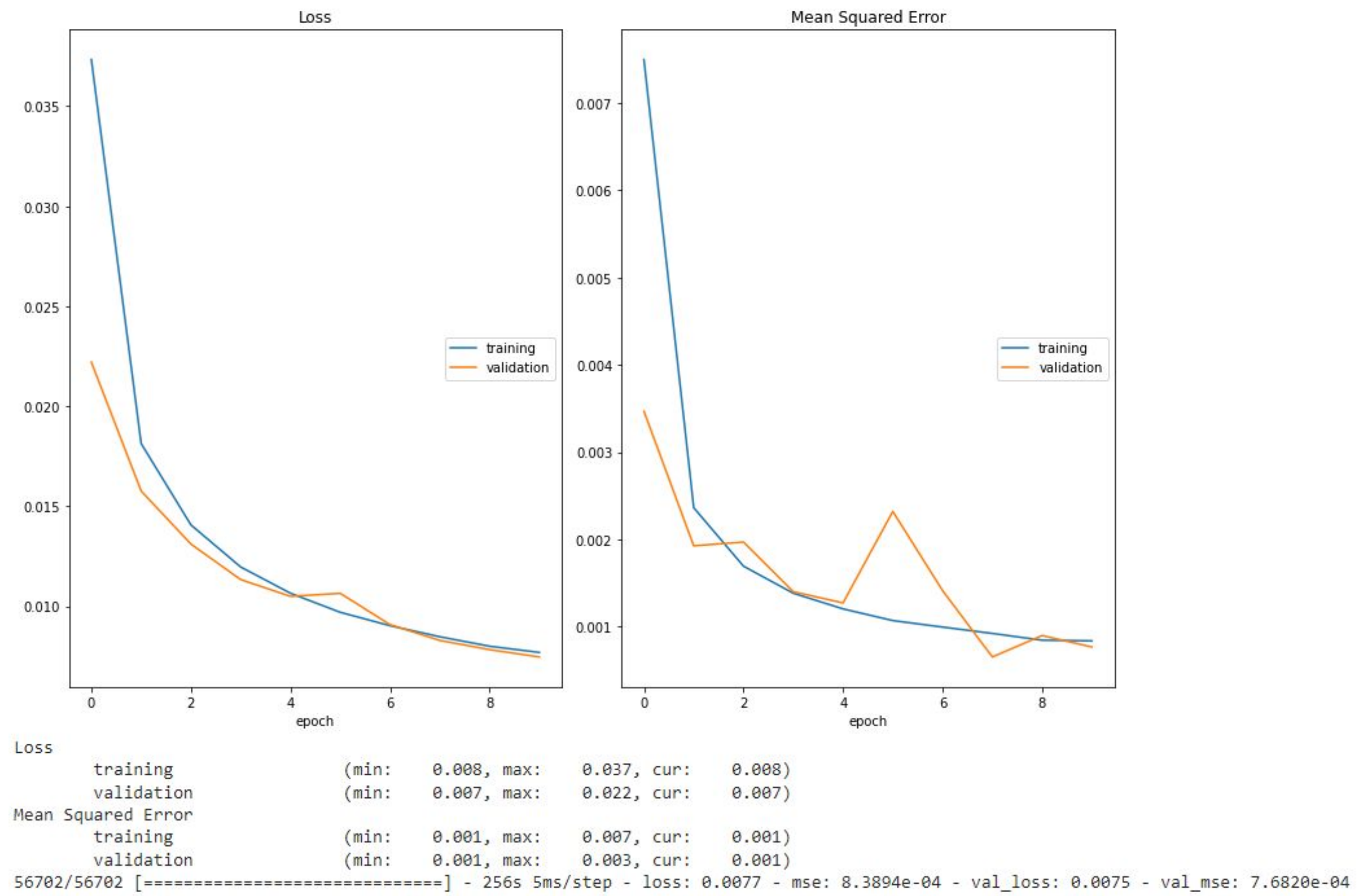
Алгоритм на основе нейронной сети.

- В качестве вишенки на торте и чего-то по-настоящему рабочего было решено воспользоваться нейронной сетью.
- Исходный набор признаков был преобразован и дополнен, и в итоге получилось 3 признака для пользователей и 6 для продуктов. В качестве таргета был взят факт повторного заказа того или иного продукта конкретным юзером. Весь сет был разбит на трэйн и тест выборки по принципу разницы между первым и последним номера заказа. Если разница была больше 1, то такие строки уходили в трэйн, если меньше – то в тест.
- Идея была в том, чтобы по величине сырой сигмоиды после предсказаний судить о «силе» или весе того или иного продукта для конкретного пользователя, отсортировав все эти продукты по убыванию значения сигмоиды, можно было получить топ N. Не всегда это был топ 10, иногда больше, иногда меньше, ведь с некоторыми продуктами некоторые пользователи могли вообще не взаимодействовать.
- Итоговый сет из 10 продуктов добивался случайным образом из исторического топ10 для каждого юзера. Также он добивался и рандомом из всех возможных продуктов, но добивать из топ10 была выгодней для сора приблизительно на 0.015. Повторы естественно исключались.

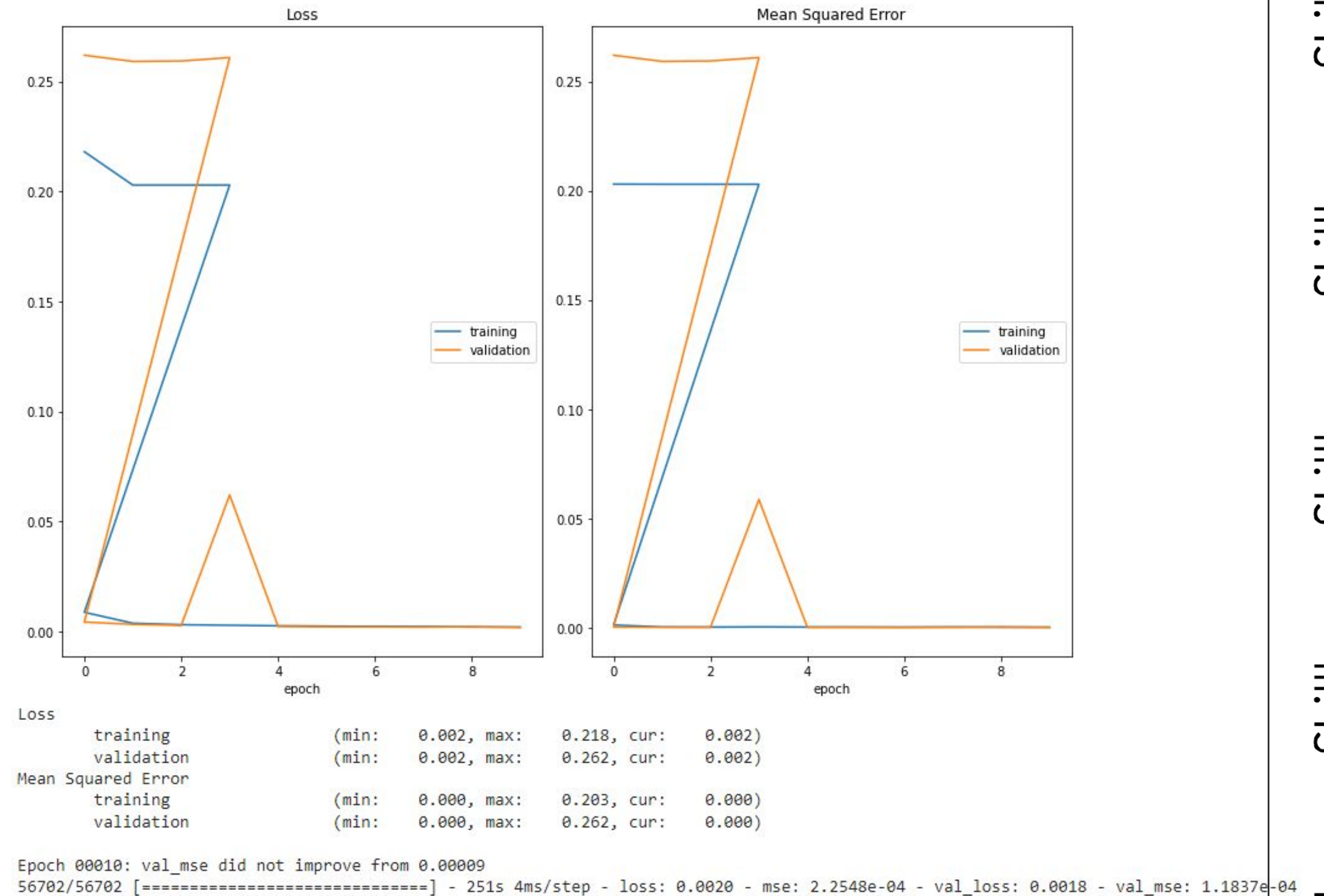
Интересные наблюдения.

Приведу две наиболее удачные конфигурации сетей: одна дала наилучшие финальные рекомендации, другая наивысшую точность предсказаний (accuracy):

Лучшие предсказания (MAP@10 - 0.27080, accuracy – 0.70)



Лучшая точность (MAP@10 - 25571, accuracy - 0.97)



Архитектура показана на пред. слайде, Adam, lr=1e-3, reg. l2=0.01

4 dense(128, 64, 32, 16), Adam, lr=1e-3, reg. l2=0.001

Некоторые технические моменты.

Вид итогового датафрейма, отсортированного по величине выхода сигмоиды:

	user_id	product_id	predictions	
	0	1	196	0.975069
	1098	1	12427	0.974666
	1309	1	10258	0.970747
	1408	1	25133	0.967029
	1590	1	46149	0.929090

	692610	206209	31477	0.695112
	844490	206209	38730	0.683582
	620478	206209	6567	0.670001
	885094	206209	22920	0.654735
	68584	206209	14197	0.507085

1079141 rows × 3 columns

Набор списка из 10 продуктов для пользователя:

```
fetch = rec_df.groupby('user_id')['product_id'].apply(list).reset_index()
```

Как дополнялись предсказания:

```
swap = random.choices(top_10.tolist()[i], k=10)

for i in fetch.index:
    while len(fetch['product_id'][i]) != 10:
        if len(fetch['product_id'][i]) < 10:
            fetch['product_id'][i].append([x for x in swap if x not in fetch[
                'product_id'][i]]) # здесь исключаются повторы
        elif len(fetch['product_id'][i]) > 10:
            fetch['product_id'][i].remove(fetch['product_id'][i][-1]) # удаляем последний (наименее популярный элемент)
```

