

INTRODUCTION

In this project variational autoencoders are extended to collaborative filtering with implicit feedback which is a method of making automatic predictions (filtering) about the interests of a user by collecting indirect preferences from many users (collaborating).

MODEL

- $u \in 1, \dots, U$ and $i \in 1, \dots, I$
- $\mathbf{X} \in \mathbb{N}^{U \times I}$ and $\mathbf{x}_u = [x_{u1}, \dots, x_{uI}]^T \in \mathbb{N}^I$
- \mathbf{z}_u : K-dimensional hidden factor for user u
- $f_\theta(\cdot)$: non-linear function $\in \mathbb{R}^I$
- $\pi(\mathbf{z}_u)$: the probability function from which \mathbf{x}_u vector is assumed to be drawn
- $N_u = \sum_i x_{ui}$: Total number of clicks of the given user u

$$\mathbf{z}_u \sim \mathcal{N}(0, \mathbf{I}_k)$$

$$\pi(\mathbf{z}_u) \propto \exp(f_\theta(\mathbf{z}_u))$$

$$\mathbf{x}_u \sim \text{Mult}(N_u, \pi(\mathbf{z}_u))$$

- Multinomial log-likelihood function:

$$\log p_\theta(\mathbf{x}_u | \mathbf{z}_u) = \sum_i x_{ui} \log \pi_i(\mathbf{z}_u)$$

- For comparison, logit log-likelihood with the logit function $\sigma(x)$:

$$\log p_\theta(\mathbf{x}_u | \mathbf{z}_u) = \sum_i [x_{ui} \cdot \log \sigma(f_{ui}) + (1 - x_{ui}) \cdot (1 - \log \sigma(f_{ui}))]$$

METHOD

- **Variational inference:** To learn the generative model we need to estimate θ parameters of $f_\theta(\cdot)$ to calculate the intractable posterior $p(\mathbf{z}_u | \mathbf{x}_u)$ for each data point. Variational inference approximates the true posterior to an instrumental posterior $q(\mathbf{z}_u)$ with Kullback-Leibler divergence $KL(q(\mathbf{z}_u) || p(\mathbf{z}_u | \mathbf{x}_u))$ where,

$$q(\mathbf{z}_u) \sim \mathcal{N}(\boldsymbol{\mu}_u, \text{diag}(\boldsymbol{\sigma}_u^2))$$

- **Amortized inference and variational autoencoders:** The number of variational parameters ($\boldsymbol{\mu}_u, \boldsymbol{\sigma}_u^2$) grows with the number of users and items in the dataset. To cope with this issue, variational autoencoders are used which are in fact data-dependent functions. This is commonly called the inference model:

$$g_\phi(\mathbf{x}_u) = [\boldsymbol{\mu}_\phi(\mathbf{x}_u), \boldsymbol{\sigma}_\phi(\mathbf{x}_u)] \in \mathbb{R}^{2K}$$

Thus the variational distribution becomes also data-dependent:

$$q_\phi(\mathbf{z}_u | \mathbf{x}_u) \sim \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}_u), \text{diag}\{\boldsymbol{\sigma}_\phi^2(\mathbf{x}_u)\})$$

METHOD(CONT'D)

- **Learning VAEs:** The evidence $p(\mathbf{x}_u)$ is a constant term, and KL divergence is a non negative term, we can rearrange the equation as below:

$$\log p_\theta(\mathbf{x}_u) \geq \mathbb{E}_q[\log p_\theta(\mathbf{x}_u | \mathbf{z}_u)] - KL(q_\phi(\mathbf{z}_u | \mathbf{x}_u) || p(\mathbf{z}_u))$$

Right-hand side of the inequality is $\mathcal{L}(\mathbf{x}_u; \theta, \phi)$ which is called evidence lower bound (ELBO).

- **Reparameterization trick:** ELBO is a function of ϕ and θ . We cannot take gradient w.r.t. ϕ when we sample $\mathbf{z} \sim q_\phi$. Instead; we sampled \mathbf{z}_u as follows:

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}_k)$$

$$\mathbf{z}_u = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\epsilon} \odot \boldsymbol{\sigma}_\phi(\mathbf{x})$$

By doing so, the stochasticity in the sampling process is isolated and the gradient w.r.t. ϕ can be backpropagated through the sampled \mathbf{z}_u .

Algorithm 1: VAE-SGD

Input: Click matrix $\mathbf{X} \in \mathbb{N}^{I \times U}$
Initialize ϕ and θ
while not converged **do**
 Sample a batch of users \mathcal{U}
 for $u \in \mathcal{U}$ **do**
 Sample $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}_k)$
 Compute \mathbf{z}_u
 Compute noisy gradients $\nabla_\phi \mathcal{L}$ and $\nabla_\theta \mathcal{L}$
 end
 Find the average of noisy gradients
 Update ϕ and θ
end

EVALUATION METRIC

Normalized Discounted Cumulative Gain:

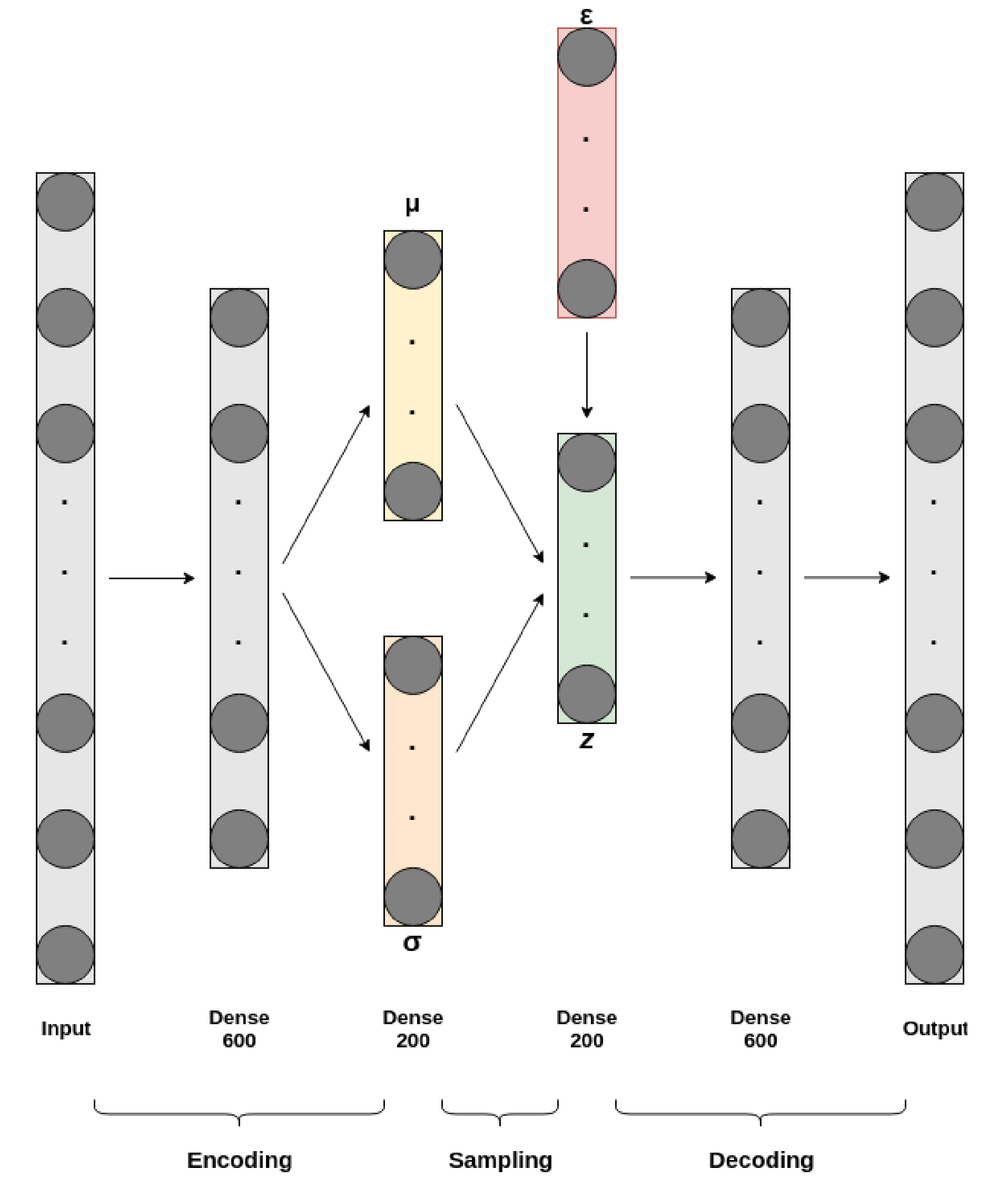
- $w(r)$: Item at rank r
- $I[\cdot]$: Indicator function
- I_u : Set of held-out items user u clicked on

Discounted cumulative gain:

$$DCG@R(u, w) = \sum_{r=1}^R \frac{2^{I[w(r) \in I_u]} - 1}{\log(r + 1)}$$

$DCG@R$ is linearly normalized to $[0, 1]$ by dividing it by the best possible $DCG@R$, where all the held-out items are ranked at the top, and we end up with $NDCG@R$.

ARCHITECTURE



DATA-RESULTS

- Random 60k users from MovieLens 20M dataset. NDCG is 0.36 for logistic and 0.38 for multinomial likelihood on the test set.

Table 1: Attributes of users and training

# of Training Users	50,000
# of Val and Test Users	5,000
# of Movie Items	17,136
# of Interactions	3.7M
% of Interactions	0.43%
Batch Size	256
# of Epochs	30

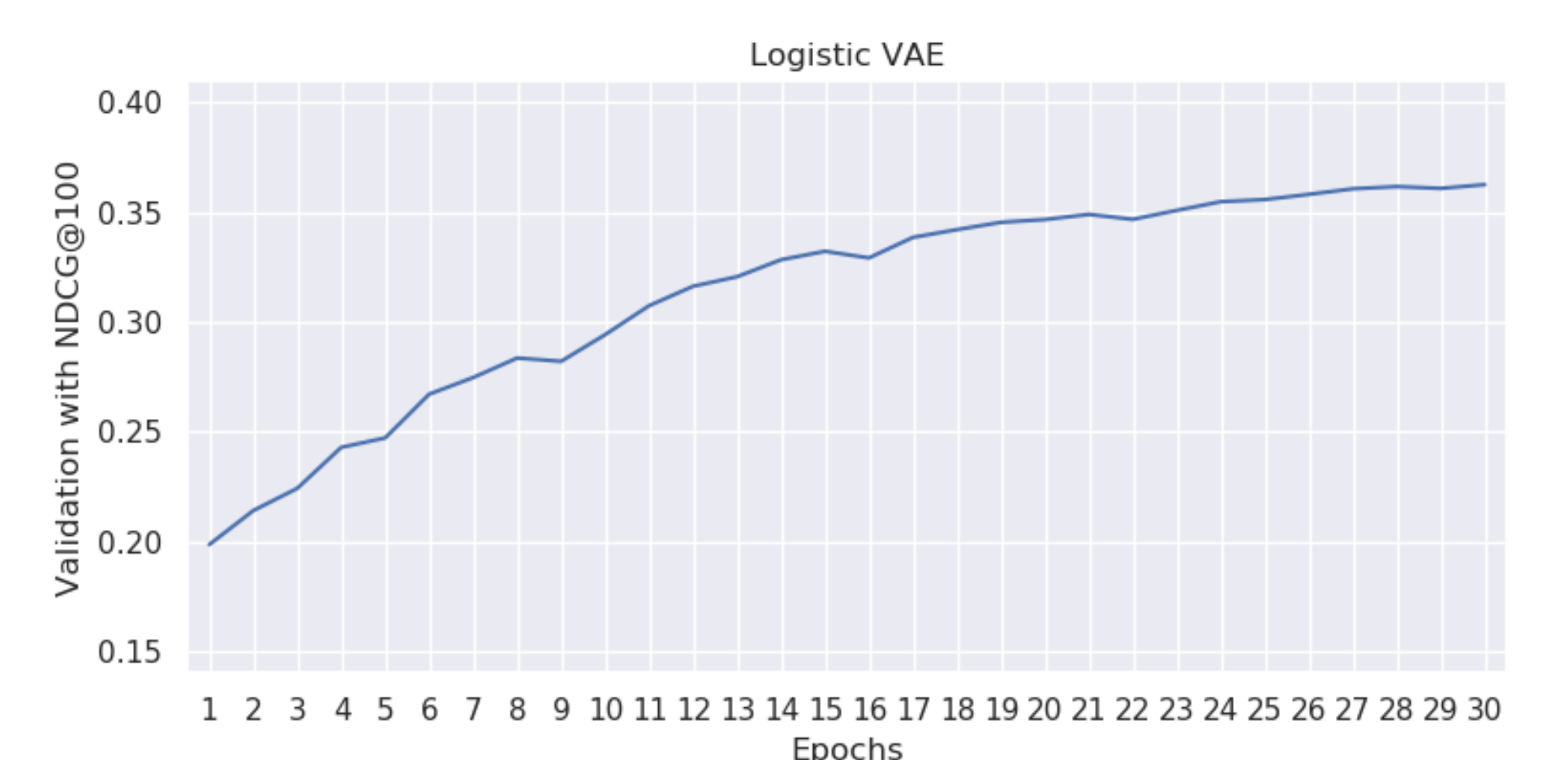


Figure 1: Training results with logistic likelihood

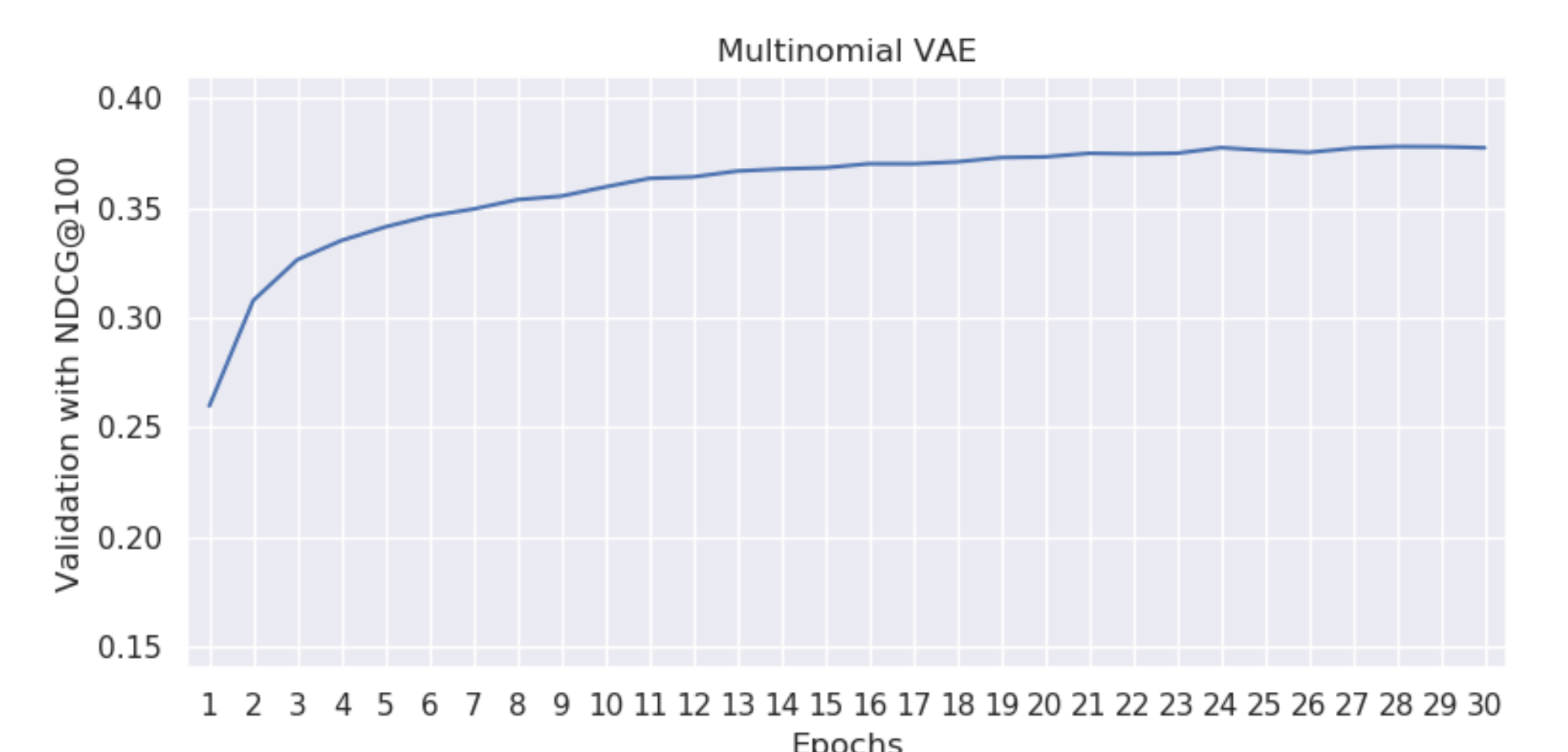


Figure 2: Training results with multinomial likelihood

REFERENCES

- [1] Liang, Dawen, Krishnan, Rahul G, Hoffman, Matthew D, and Jebara, Tony. "Variational autoencoders for collaborative filtering." arXiv preprint arXiv:1802.05814, 2018.
- [2] Variational Autoencoder: Intuition and Implementation, <https://wiseodd.github.io/techblog/2016/12/10/variational-autoencoder/>