

# HDDA Recommendation systems via approximate matrix factorization

Dmitry Beresnev  
Vsevolod Klyushev

Innopolis University — IU — 2024

# Initial problem

We need to solve the following problem:

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{r \times n}} \|W \circ (X - UV)\|_F^2 \quad (1)$$

where:

- $X$  — target matrix ( $m \times n$ );
- $W$  — binary mask;
- $r$  — rank of factorization;

In order to use gradient methods, we need to derive it with respect to each parameter ( $U$  and  $V$ ) (full derivation might be found [here](#)):

$$\frac{\partial \|W \circ (X - UV)\|_F^2}{\partial U} = -2(W \circ X)V^T + 2(W \circ (UV))V^T \quad (2)$$

$$\frac{\partial \|W \circ (X - UV)\|_F^2}{\partial V} = -2U^T(W \circ X) + 2U^T(W \circ (UV)) \quad (3)$$

# Improved problem

As a modified problem, we decided just to add some regularization. This brings us to the following forms:

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{r \times n}} \|W \circ (X - UV)\|_F^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2 \quad (4)$$

$$\frac{\partial(\|W \circ (X - UV)\|_F^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2)}{\partial U} = -2(W \circ X)V^T + 2(W \circ (UV))V^T + 2\lambda U \quad (5)$$

$$\frac{\partial(\|W \circ (X - UV)\|_F^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2)}{\partial V} = -2U^T(W \circ X) + 2U^T(W \circ (UV)) + 2\lambda V \quad (6)$$

where  $\lambda$  — regularization parameter

# Constant + decreasing

We can take constant step, or such step, that would decrease with increase of iterations. For example,  $\gamma = \frac{1}{\sqrt{k}}$ , where  $k$  — number of iterations.

# Estimate $1/L$

---

## Algorithm 1 Estimate $1/L$

---

**Input:**  $\theta$  (point),  $\nabla f(\theta)$  (gradient function),  $p$  (desired direction),  $\beta$  (multiplier),  $t$  (max iterations)

$\alpha \leftarrow 1$

**for**  $i = 1$  **to**  $t$  **do**

$\theta_i \leftarrow \theta + \alpha p$

**if**  $\|\nabla f(\theta) - \nabla f(\theta_i)\| > \frac{1}{\alpha} \|\theta - \theta_i\|$  **then**

$\alpha \leftarrow \beta \cdot \alpha$

**if**  $\alpha < \epsilon$  **then**

**return**  $\alpha$

**end if**

**else**

**return**  $\alpha$

**end if**

**end for**

**return**  $\alpha$

---

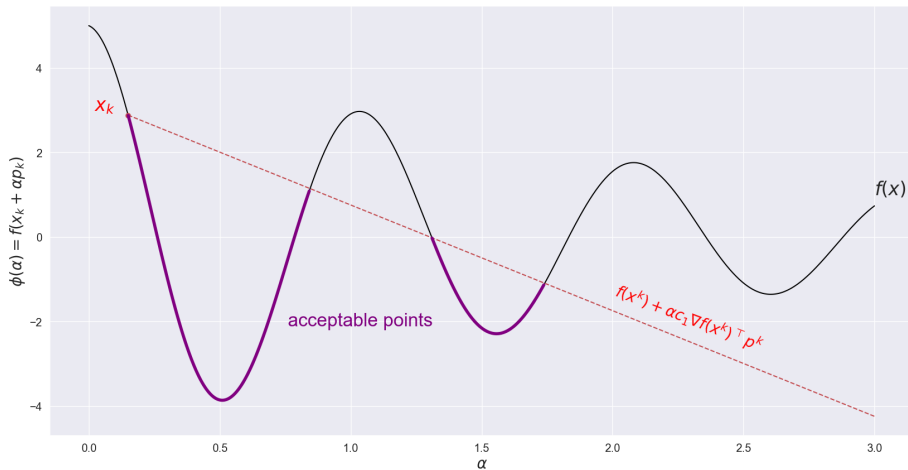


Figure: Armijo rule

---

## Algorithm 2 Armijo Step

---

**Input:**  $\theta$  (point),  $f(\theta)$  (objective function),  $p$  (desired direction),  $\beta$  (multiplier),  $c_1 > 0$ ,  $t$  (max iterations)

$\alpha \leftarrow 1$

**for**  $i = 1$  **to**  $t$  **do**

$\theta_i \leftarrow \theta + \alpha p$

**if**  $f(\theta_i) > f(\theta) + c_1 \alpha \langle \nabla_{\theta} f(\theta), p \rangle$  **then**

$\alpha \leftarrow \beta \cdot \alpha$

**if**  $\alpha < \epsilon$  **then**

**return**  $\alpha$

**end if**

**else**

**return**  $\alpha$

**end if**

**end for**

**return**  $\alpha$

---



# Curvature condition

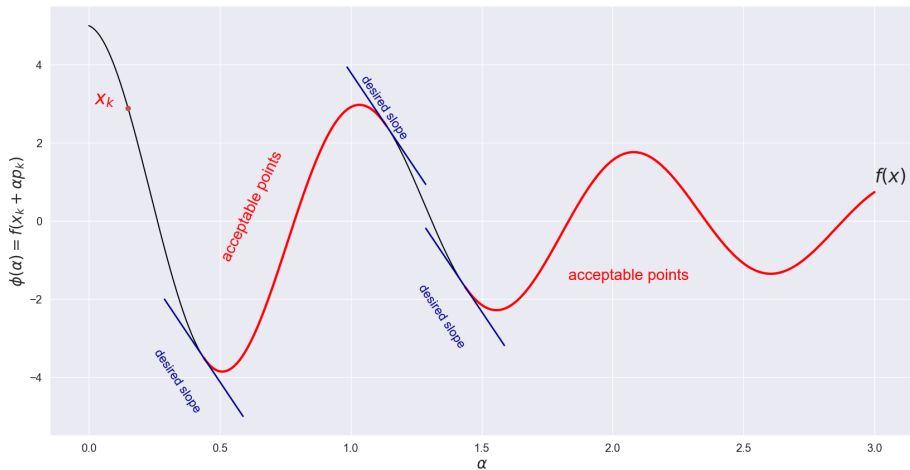


Figure: Curvature condition

# Bisection Weak Wolfe

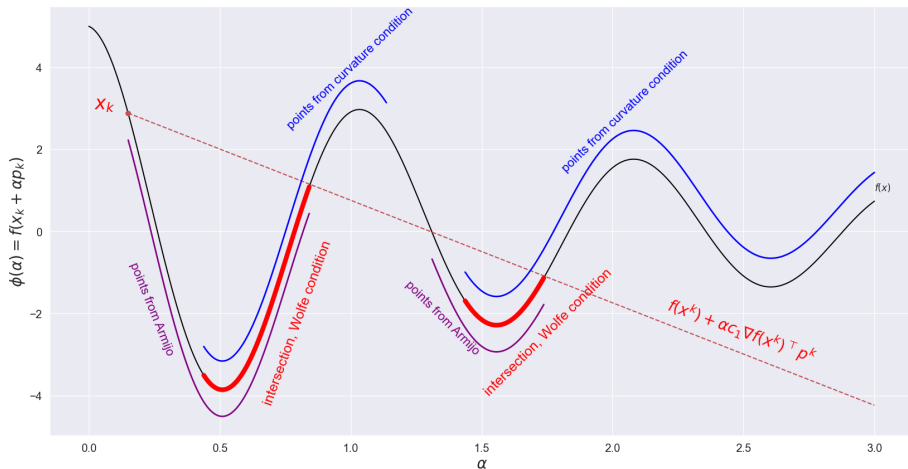


Figure: Weak Wolfe conditions

## Algorithm 3 Bisection Weak Wolfe Step

**Input:**  $\theta$  (point),  $f(\theta)$  (objective function),  $p$  (desired direction),  $\beta$  (multiplier),  $c_1 > 0$ ,  $c_2 > c_1$ ,  $t$  (max iterations)

```
 $\alpha \leftarrow 1$   
 $a \leftarrow 0$   
 $b \leftarrow +\infty$   
for  $i = 1$  to  $t$  do  
   $\theta_i \leftarrow \theta + \alpha p$   
  if  $f(\theta_i) > f(\theta) + c_1 \alpha \langle \nabla_{\theta} f(\theta), p \rangle$  then  
     $b \leftarrow \alpha$   
     $\alpha \leftarrow \frac{1}{2}(a + b)$   
  else if  $\langle p, \nabla_{\theta} f(\theta_i) \rangle < c_2 \langle p, \nabla_{\theta} f(\theta) \rangle$  then  
     $a \leftarrow \alpha$   
    if  $b = +\infty$  then  
       $\alpha \leftarrow 2a$   
    else  
       $\alpha \leftarrow \frac{1}{2}(a + b)$   
    end if  
  else  
    return  $\alpha$   
  end if  
end for  
return  $\alpha$ 
```

▷ Lower bound  
▷ Upper bound  
▷ Armijo condition  
▷ Decrease step  
▷ Curvature condition  
▷ Increase step

# Strong Wolfe

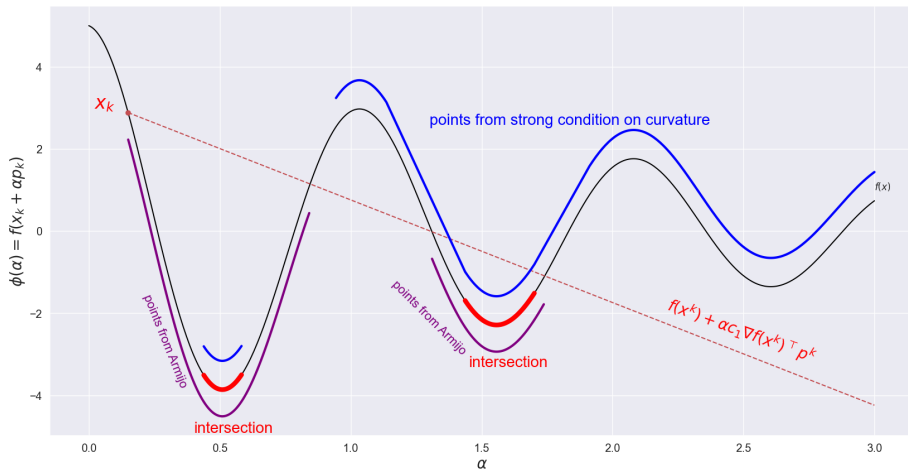


Figure: Strong Wolfe conditions

---

## Algorithm 4 Strong Wolfe Step

---

**Input:**  $\theta$  (point),  $f(\theta)$  (objective function),  $p$  (desired direction),  $\beta$  (multiplier),  $c_1 > 0$ ,  $c_2 > c_1$ ,  $t$  (max iterations)

$\alpha \leftarrow 1$

**for**  $i = 1$  **to**  $t$  **do**

$\theta_i \leftarrow \theta + \alpha p$

**if**  $(f(\theta_i) > f(\theta) + c_1 \alpha \langle \nabla_{\theta} f(\theta), p \rangle)$  **or**  $(\| \langle p, \nabla_{\theta} f(\theta_i) \rangle \| > c_2 \| \langle p, \nabla_{\theta} f(\theta) \rangle \|)$

**then**

$\alpha \leftarrow \beta \cdot \alpha$

**if**  $\alpha < \epsilon$  **then**

**return**  $\alpha$

**end if**

**else**

**return**  $\alpha$

**end if**

**end for**

**return**  $\alpha$

---

## Algorithm 5 Gradient Descent optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\mathcal{L}(p)$  (step size choosing strategy)

**for**  $t = 1$  **to**  $\dots$  **do**

$p_t \leftarrow -\nabla_{\theta} f_t(\theta_{t-1})$  ▷ Step direction

    Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$

**end for**

**return**  $\theta_t$

---

# Adaptive Gradient Descent

---

## Algorithm 6 Adaptive Gradient Descent

---

**Input:**  $x_0$  (parameters to optimize),  $f(x)$  (objective function)

**Initialize:**  $\lambda_0 > 0$  (small start step),  $\theta_0 \leftarrow +\infty$

$x_1 \leftarrow x_0 - \lambda_0 \nabla f(x_0)$

**for**  $t = 1$  **to**  $\dots$  **do**

$\lambda_t = \min \left\{ \sqrt{1 + \theta_{t-1} \lambda_{t-1}}, \frac{\|x_t - x_{t-1}\|}{2 \|\nabla f(x_t) - \nabla f(x_{t-1})\|} \right\}$

$x_{t+1} \leftarrow x_t - \lambda_t \nabla f(x_t)$

$\theta_t \leftarrow \frac{\lambda_t}{\lambda_{t-1}}$

**end for**

**return**  $x_t$

---

# Heavy Ball

---

**Algorithm 7** Heavy Ball optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\beta$  (momentum),  $\mathcal{L}(p)$  (step size choosing strategy)

**for**  $t = 1$  **to**  $\dots$  **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**if**  $\beta \neq 0$  **then**

**if**  $t > 1$  **then**

$b_t \leftarrow \beta b_{t-1} + g_t$

**else**

$b_t \leftarrow g_t$

**end if**

$g_t \leftarrow b_t$

**end if**

$p_t \leftarrow -g_t$

▷ Step direction

    Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$

**end for**

**return**  $\theta_t$



---

**Algorithm 8** Nesterov optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\beta$  (momentum),  $\mathcal{L}(p)$  (step size choosing strategy)

**for**  $t = 1$  **to**  $\dots$  **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**if**  $\beta \neq 0$  **then**

**if**  $t > 1$  **then**

$b_t \leftarrow \beta b_{t-1} + g_t$

**else**

$b_t \leftarrow g_t$

**end if**

$g_t \leftarrow g_t + \beta b_t$

**end if**

$p_t \leftarrow -g_t$

▷ Step direction

    Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$

**end for**

**return**  $\theta_t$

---

**Algorithm 9** AdaGrad optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\mathcal{L}(p)$  (step size choosing strategy)

**Initialize:**  $s_0 \leftarrow 0$  (cumulative square sum)

**for**  $t = 1$  **to**  $\dots$  **do**

$$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$$

$$s_t \leftarrow \alpha s_{t-1} + g_t^2$$

$$p_t \leftarrow -g_t / (\sqrt{s_t} + \epsilon)$$

▷ Step direction

Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$$

**end for**

**return**  $\theta_t$

---

---

**Algorithm 10** RMSProp optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\alpha$  (alpha),  $\mathcal{L}(p)$  (step size choosing strategy)

**Initialize:**  $v_0 \leftarrow 0$  (square average)

**for**  $t = 1$  **to**  $\dots$  **do**

$$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$$

$$v_t \leftarrow \alpha v_{t-1} + (1 - \alpha) g_t^2$$

$$p_t \leftarrow -g_t / (\sqrt{v_t} + \epsilon)$$

▷ Step direction

Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$$

**end for**

**return**  $\theta_t$

---

---

**Algorithm 11** Adam optimizer

---

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\beta_1, \beta_2$  (alpha),  $\mathcal{L}(p)$  (step size choosing strategy)

**Initialize:**  $m_0 \leftarrow 0$  (first moment),  $v_0 \leftarrow 0$  (second moment)

**for**  $t = 1$  **to**  $\dots$  **do**

$$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$$

$$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$$

$$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$$

$$p_t \leftarrow -\hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$$

▷ Step direction

Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$

$$\theta_t \leftarrow \theta_{t-1} + \gamma p_t$$

**end for**

**return**  $\theta_t$

---

**Algorithm 12** BFGS optimizer

**Input:**  $\theta_0$  (parameters to optimize),  $f(\theta)$  (objective function),  $\mathcal{L}(p)$  (step size choosing strategy)

**Initialize:**  $H_0 \leftarrow I$

**for**  $t = 1$  **to**  $\dots$  **do**

$$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$$

$$p_t \leftarrow -H_{t-1}g_t$$

▷ Step direction

Choose step size  $\gamma$  according to  $\mathcal{L}(p_t)$  ▷  $\gamma$  should satisfy Wolfe conditions

$$s_t \leftarrow \gamma p_t$$

$$\theta_t \leftarrow \theta_{t-1} + s_t$$

$$y_t \leftarrow \nabla_{\theta} f_t(\theta_t) - g_t$$

$$H_t \leftarrow H_{t-1} + \frac{(s_t^T y_t + y_t^T H_{t-1} y_t)(s_t s_t^T)}{(s_t^T y_t)^2} - \frac{H_{t-1} y_t s_t^T + s_t y_t^T H_{t-1}}{s_t^T y_t}$$

**end for**

**return**  $\theta_t$

# Experiments. Gradient descent

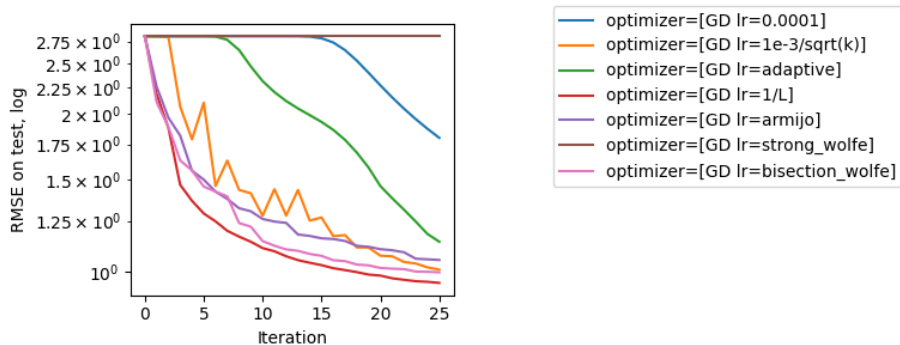


Figure: Gradient descent

# Experiments. Heavy Ball

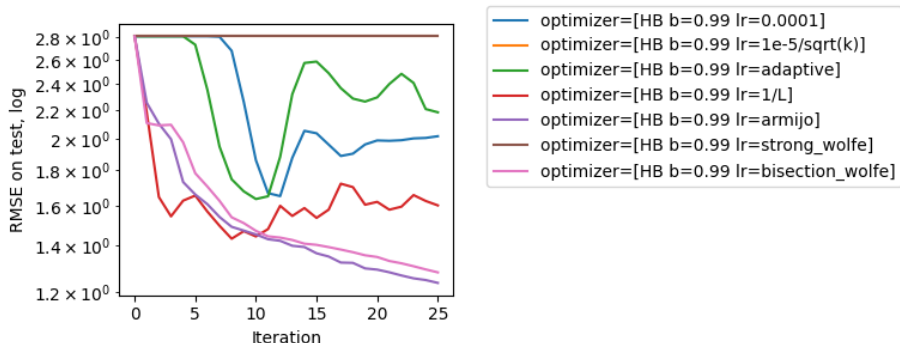


Figure: Heavy Ball

# Experiments. Nesterov

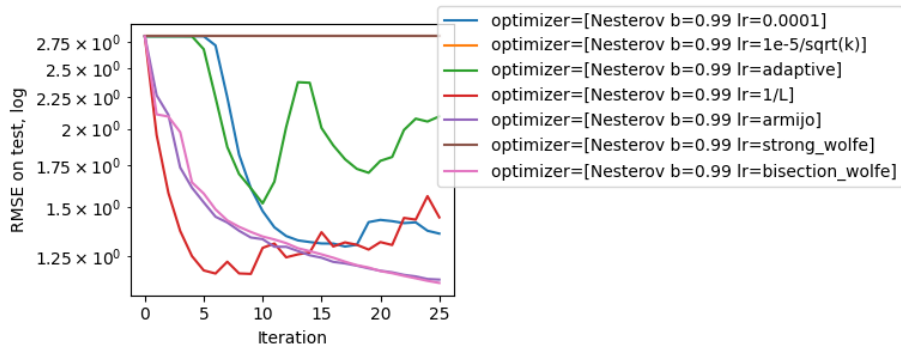


Figure: Nesterov



# Experiments. RMSprop

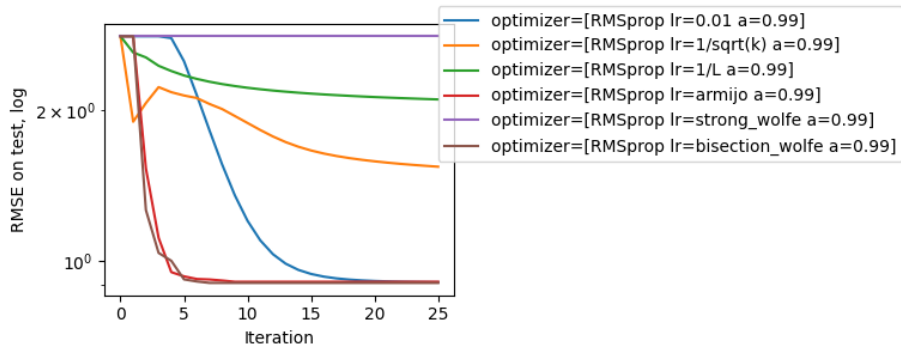


Figure: RMSprop

# Experiments. AdaGrad

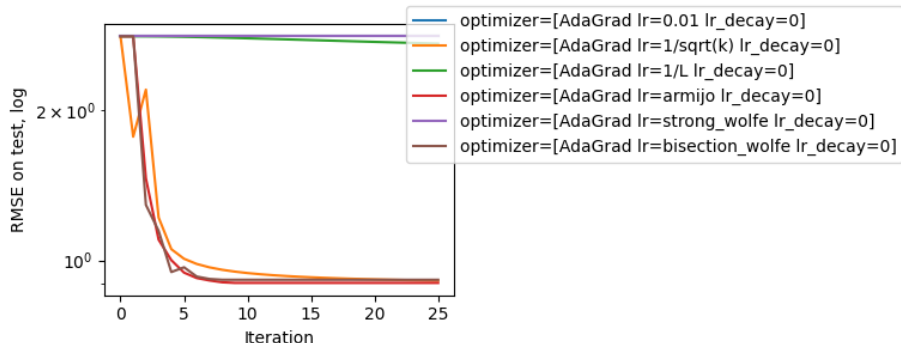


Figure: AdaGrad

# Experiments. Adam

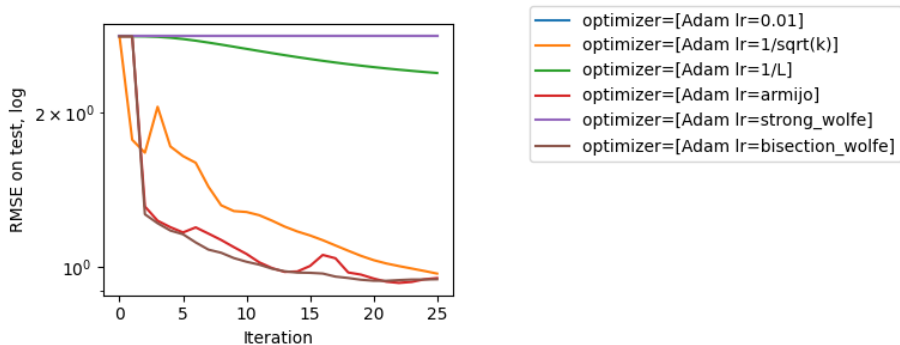


Figure: Adam

# Experiments. Comparison

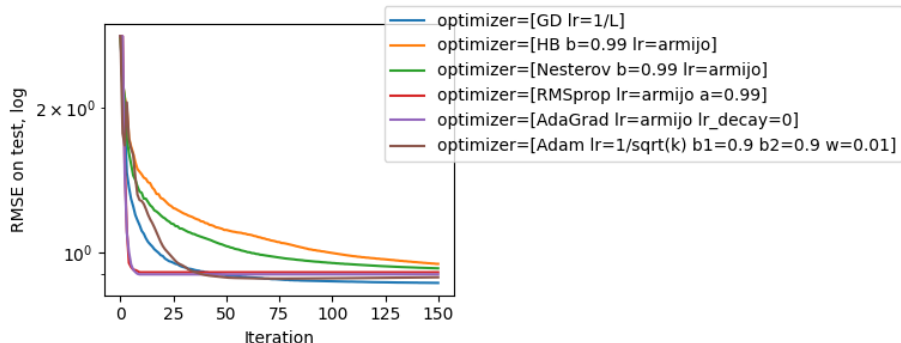


Figure: Comparison

# Experiments. Best model

Our best model with RMSE score 0.86 have the following parameters:  $r = 10$ , GD optimizer with estimate  $1/L$  strategy and  $\lambda = 2$

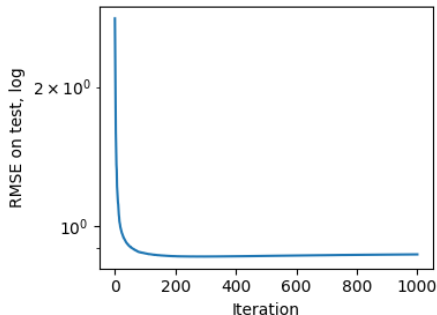


Figure: Best model

# Main idea

Instead of making update for entire  $U$  or  $V$  simultaneously, we can make updates row by row (column by column for  $V$ ). The reasons are the following:

- The objective function becomes  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , so we can apply methods like BFGS
- There will be more updates, and such updates will be more diverse: we will use just updated values for new updates

# Problem formulation

Therefore, the new problem with fixed  $V$  becomes

$$\min_{U_i^T \in \mathbb{R}^r} \|W_i^T \circ (X_i^T - U_i^T V)\|^2 + \lambda \|U_i^T\|^2, \quad \forall i \in \{1, 2, \dots, m\} \quad (7)$$

and the new problem with fixed  $U$ :

$$\min_{V_j \in \mathbb{R}^r} \|W_j \circ (X_j - UV_j)\|^2 + \lambda \|V_j\|^2, \quad \forall j \in \{1, 2, \dots, n\} \quad (8)$$

Gradient of ?? is

$$\begin{aligned} \frac{\partial(\|W_i^\top \circ (X_i^\top - U_i^\top V)\|^2 + \lambda\|U_i^\top\|^2)}{\partial U_i^\top} &= \\ &= -2(W_i^\top \circ X_i^\top)V^T + 2(W_i^\top \circ (U_i^\top V))V^T + 2\lambda U_i^\top \end{aligned} \quad (9)$$

and the gradient of ??:

$$\begin{aligned} \frac{\partial(\|W_j \circ (X_j - UV_j)\|^2 + \lambda\|V_j\|^2)}{\partial V_j} &= \\ &= -2U^T(W_j \circ X_j) + 2U^T(W_j \circ (UV_j)) + 2\lambda V_j \end{aligned} \quad (10)$$



---

## Algorithm 13 Vector Gradient Descent

---

**Input:**  $X, W \in \mathbb{R}^{m \times n}$  — given initial and binary matrices,  $U \in \mathbb{R}^{m \times r}$ ,  $V \in \mathbb{R}^{r \times n}$  — arbitrary matrices,  $k_U, k_V$  — small integers such that  $\frac{m}{n} \approx \frac{k_U}{k_V}$

$d \leftarrow k_U + k_V$

**repeat**

**for**  $t = 0$  **to**  $n + m - 1$  **do**

$r \leftarrow t \bmod d$

**if**  $r < k_U$  **then**

$i \leftarrow k_U \cdot (t \text{ div } d) + r + 1$

**if**  $i > m$  **then**

**continue**

**end if**

      Update  $U_i^\top$  using ??

**else**

---

---

**Algorithm** Vector Gradient Descent. Continue

---

```
...
repeat
  for  $t = 0$  to  $n + m - 1$  do
    ...
    if  $i > m$  then
      ...
    else
       $j \leftarrow k_V \cdot (t \text{ div } d) + r - k_U + 1$ 
      if  $j > n$  then
        continue
      end if
      Update  $V_j$  using ??
    end if
  end for
until convergence

return  $U, V = 0$ 
```

---

# Experiments. Results

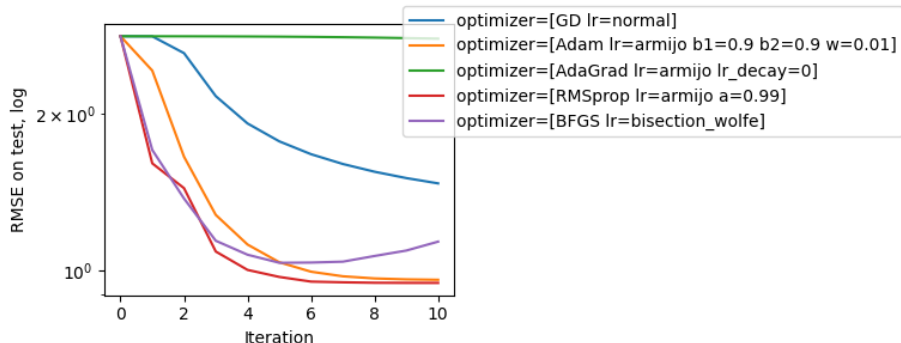


Figure: Vector Gradient Descent

# NNMF. Problem formulation

We want to solve initial problem (??), but with non-negativity constraints:

$$\min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{r \times n}} \|W \circ (X - UV)\|_F^2 \quad (11)$$

s.t.  $U, V \geq 0$

In order to solve such problem more easily, we need to derive multiplicative updates for  $U$  and  $V$ .

With the idea from Lee and Seung paper we defined our MU updates in such way:

$$U \leftarrow U \circ \frac{((W \circ X)V^T)}{((W \circ (UV))V^T + \epsilon)} \quad (12)$$

$$V \leftarrow V \circ \frac{(U^T(W \circ X))}{(U^T(W \circ (UV)) + \epsilon)} \quad (13)$$

# Experiments. Results

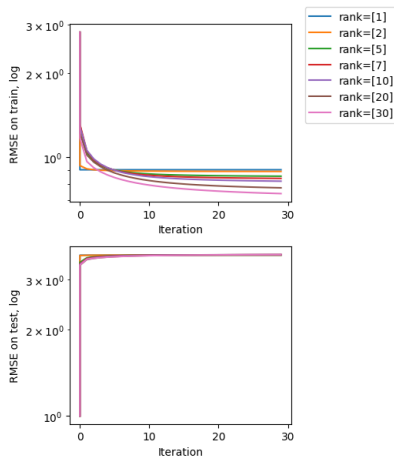


Figure: NNMF

As an alternative approach we decided to train simple neural network with such parameters:

- 4 layers ( $23 \times 64, 64 \times 128, 128 \times 64, 64 \times 1$ )
- ReLU as activation function after each layer
- MSELoss criterion
- Adam optimizer with  $\alpha = 0.001$
- Early stopping (if on test set loss is not decreasing for 3 iterations)
- Maximum number of training epochs = 50

As input to our model we decided to use:

- Min-max scaled 'user\_id'
- Min-max scaled 'movie\_id'
- One-hot encoded 'genres' (we have 18 unique genres)
- Binary encoded 'gender'
- Min-max scaled first two user features ('feature\_1' and 'feature\_2')

In output we have just one number - rating for specific pair of user and film.

As a result we achieved average loss on test set 1.1



# Overall Results

We tested several approaches: Gradient descent, Vector Gradient descent, Non-Negative Matrix Factorization and Neural Network.

- Both Neural Network and Vector Gradient descent showed their applicability for recommendation system task, however, they have been outperformed by other models
- NNMF is inapplicable for recommendation system task mostly because of the mask in the problem formulation. Mask makes NNMF to set almost all unknown values to 1, which is very bad decision
- It happens to be that Gradient descent method showed the best overall performance (both in terms of time and score). We were able to achieved RMSE score 0.86 on test data

**Thanks!**