

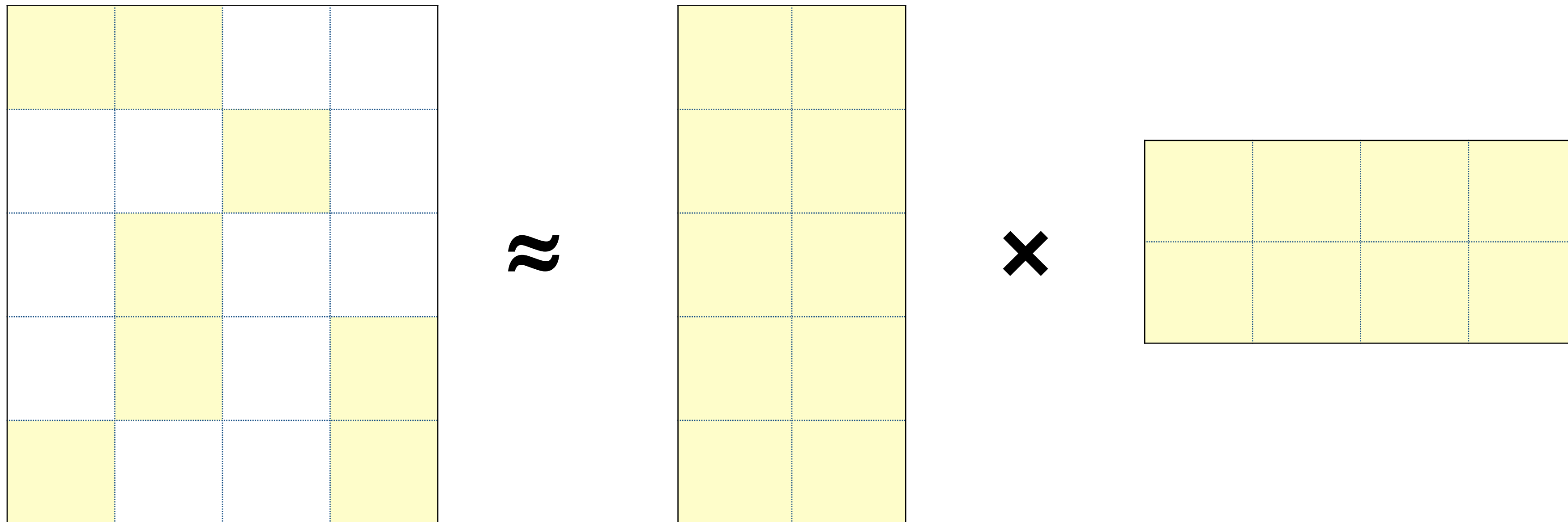
Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla, et al. (2011)

박민주

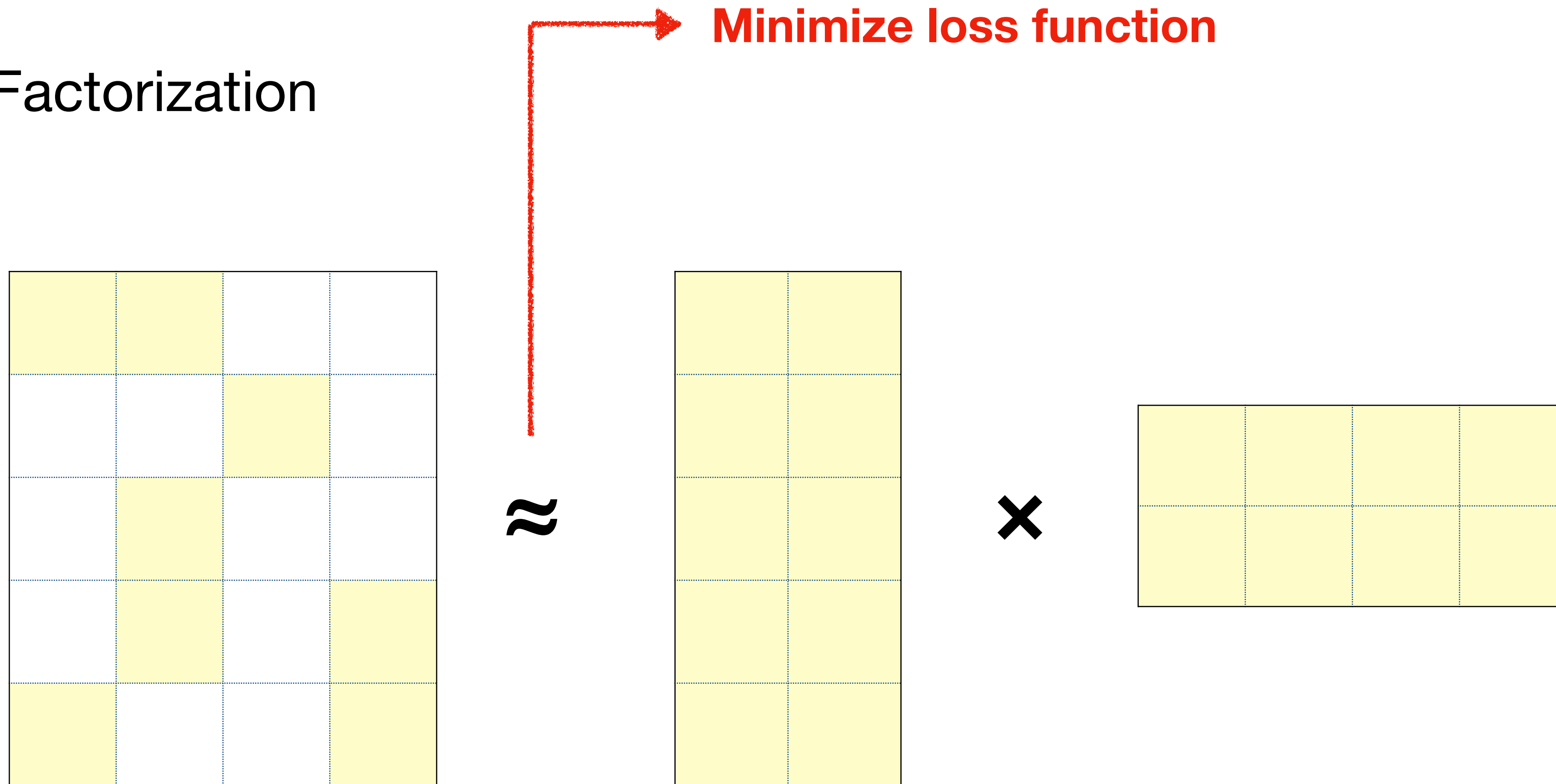
Motivation

- Matrix Factorization



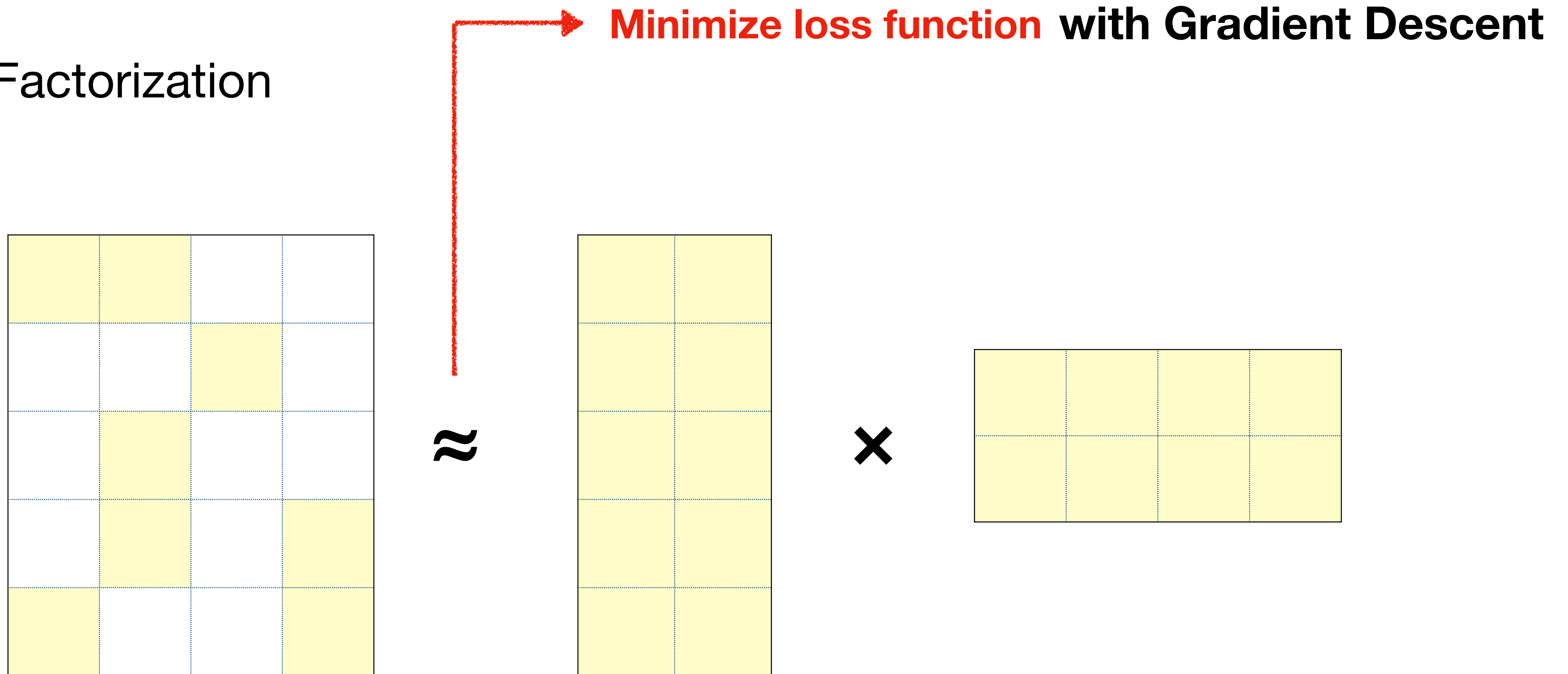
Motivation

- Matrix Factorization



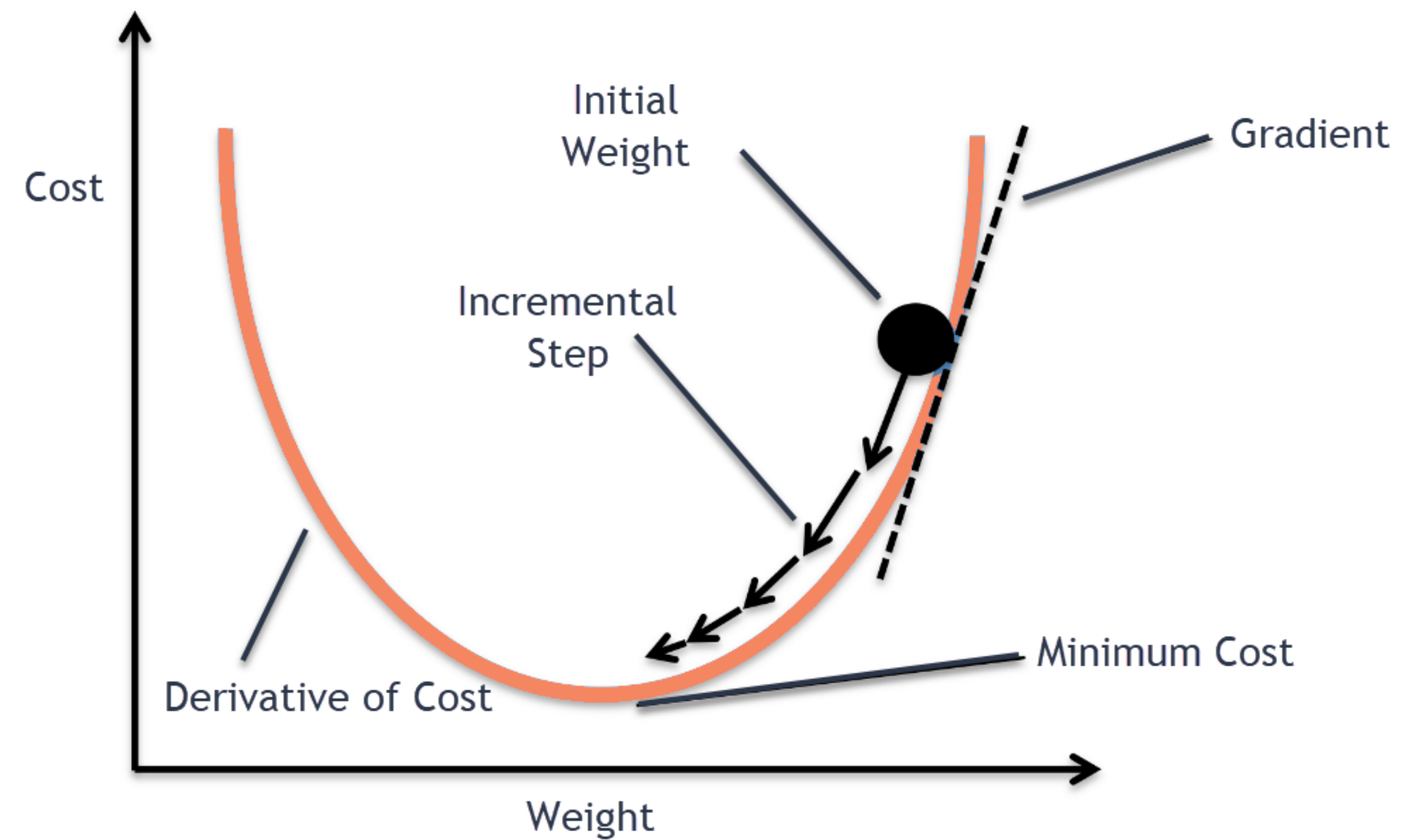
Motivation

- Matrix Factorization



Background

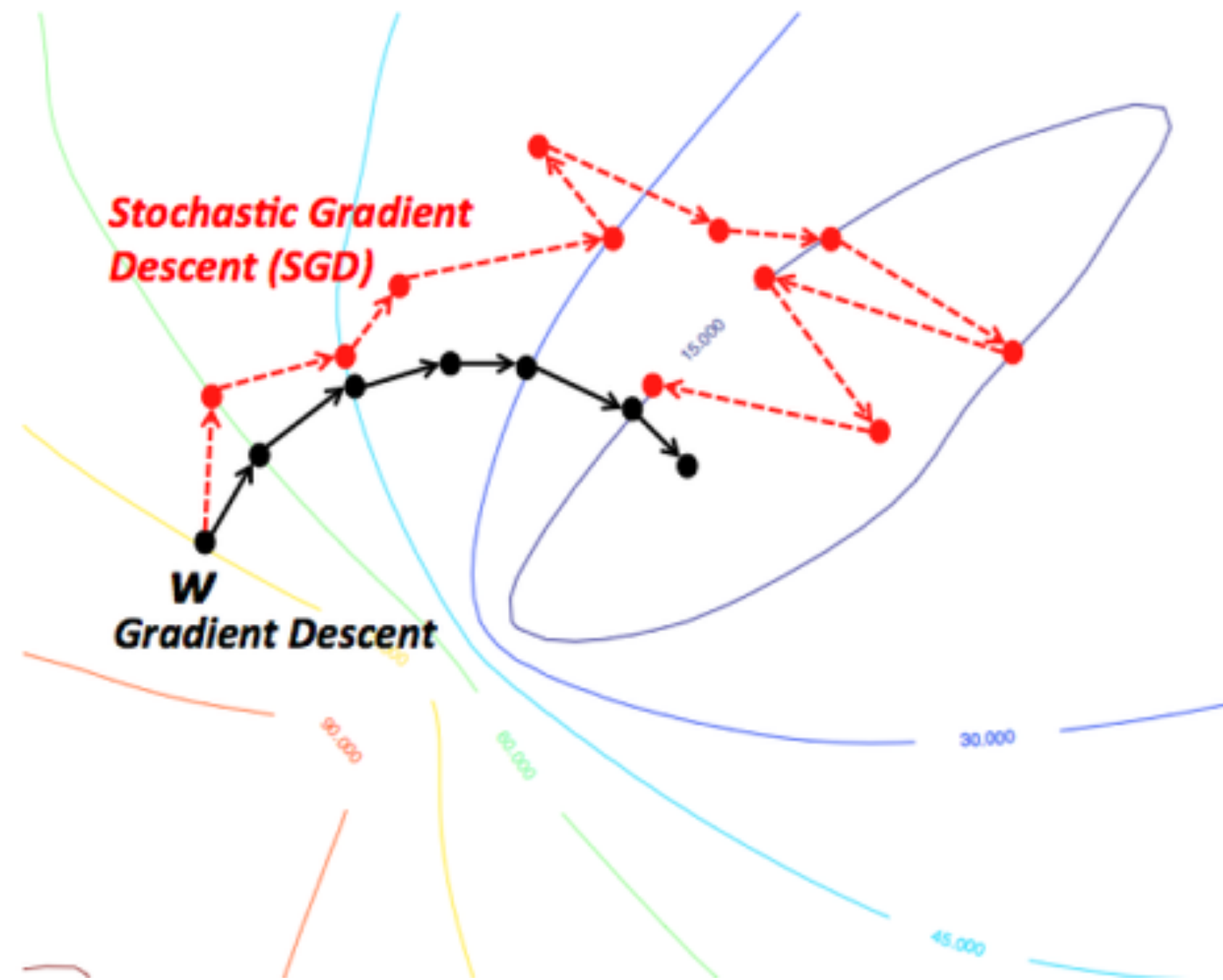
- Gradient Descent



$$\theta_{n+1} = \theta_n - \epsilon \hat{L}'(\theta_n)$$

Background

- Stochastic Gradient Descent



$$\theta_{n+1} = \theta_n - \epsilon \hat{L}'(\theta_n)$$

In practice, with additional projection

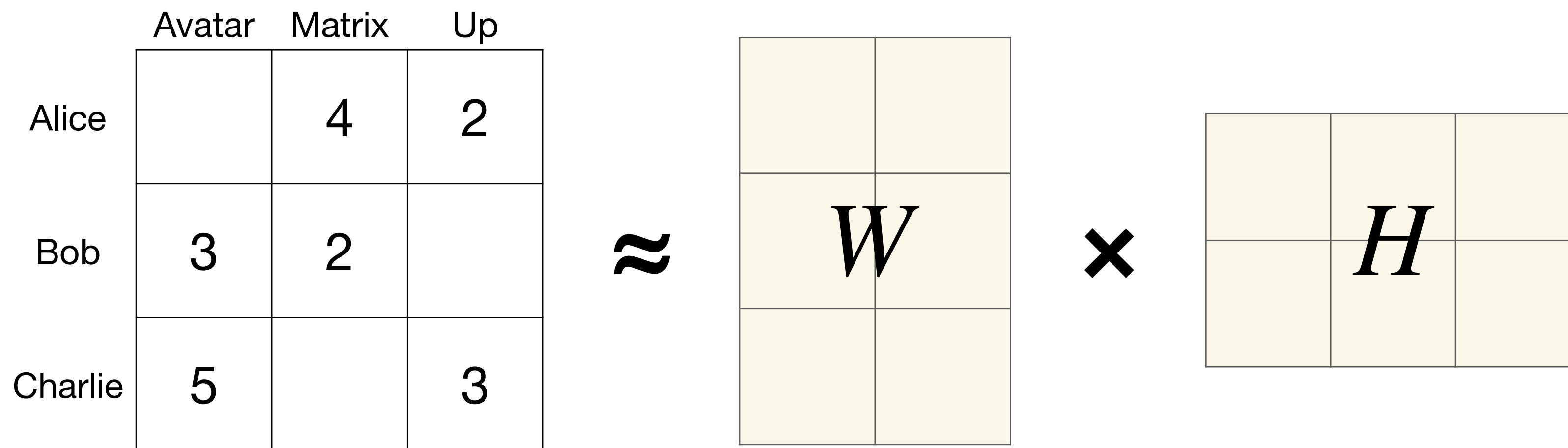
$$\theta_{n+1} = \Pi_H[\theta_n - \epsilon \hat{L}'(\theta_n)]$$

* H : Constraint set

SGD for Matrix Factorization

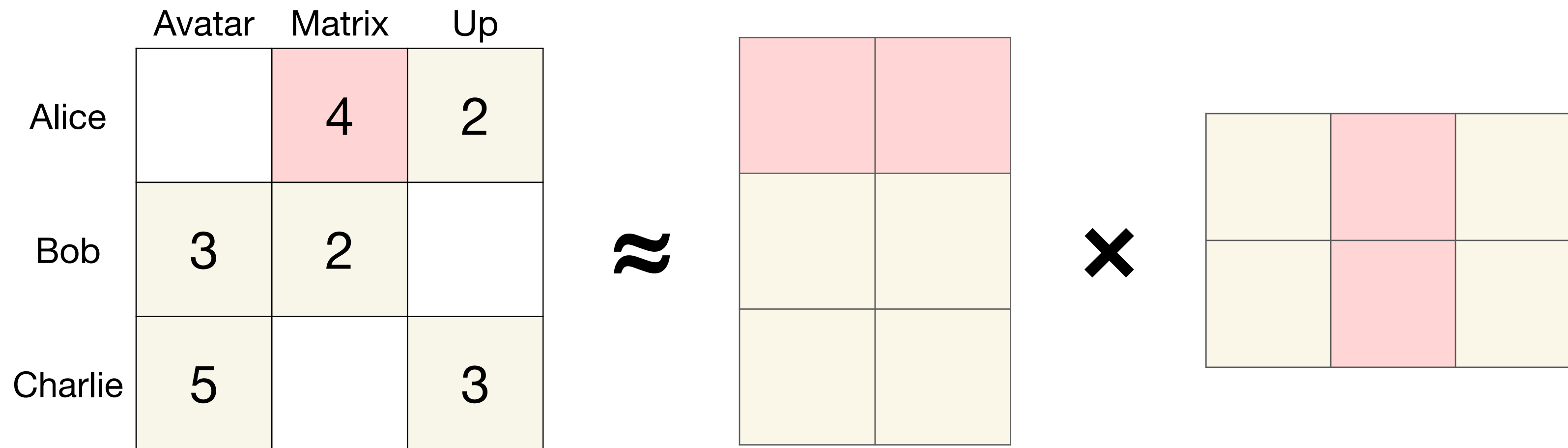
- Stochastic Gradient Descent

Finding the best model $\underset{W,H}{\operatorname{argmin}} L(V, W, H) \rightarrow \theta_{n+1} = \theta_n - \epsilon \hat{L}'(\theta_n)$
 $\theta = (W, H)$



SGD for Matrix Factorization

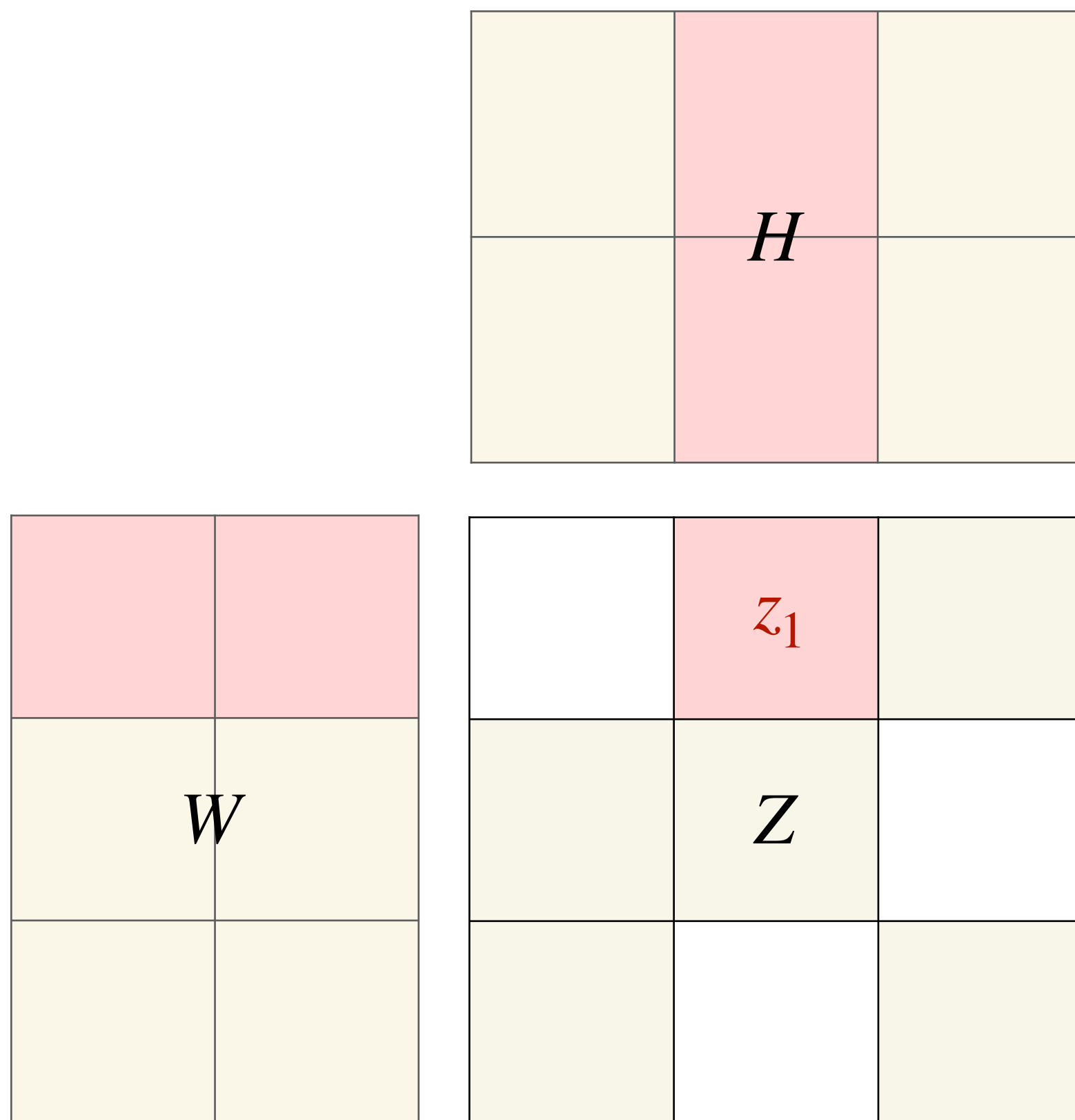
- Stochastic Gradient Descent



- Decompose Loss Function

$$L = \sum_{(i,j) \in Z} l(V_{ij}, W_{i*}, H_{*j}) \quad (Z = \{(i, j) : V_{ij} \neq 0\})$$

SGD for Matrix Factorization



- $L_z(\theta) = L_{ij}(\theta) = l(V_{ij}, W_{i*}, H_{*j})$

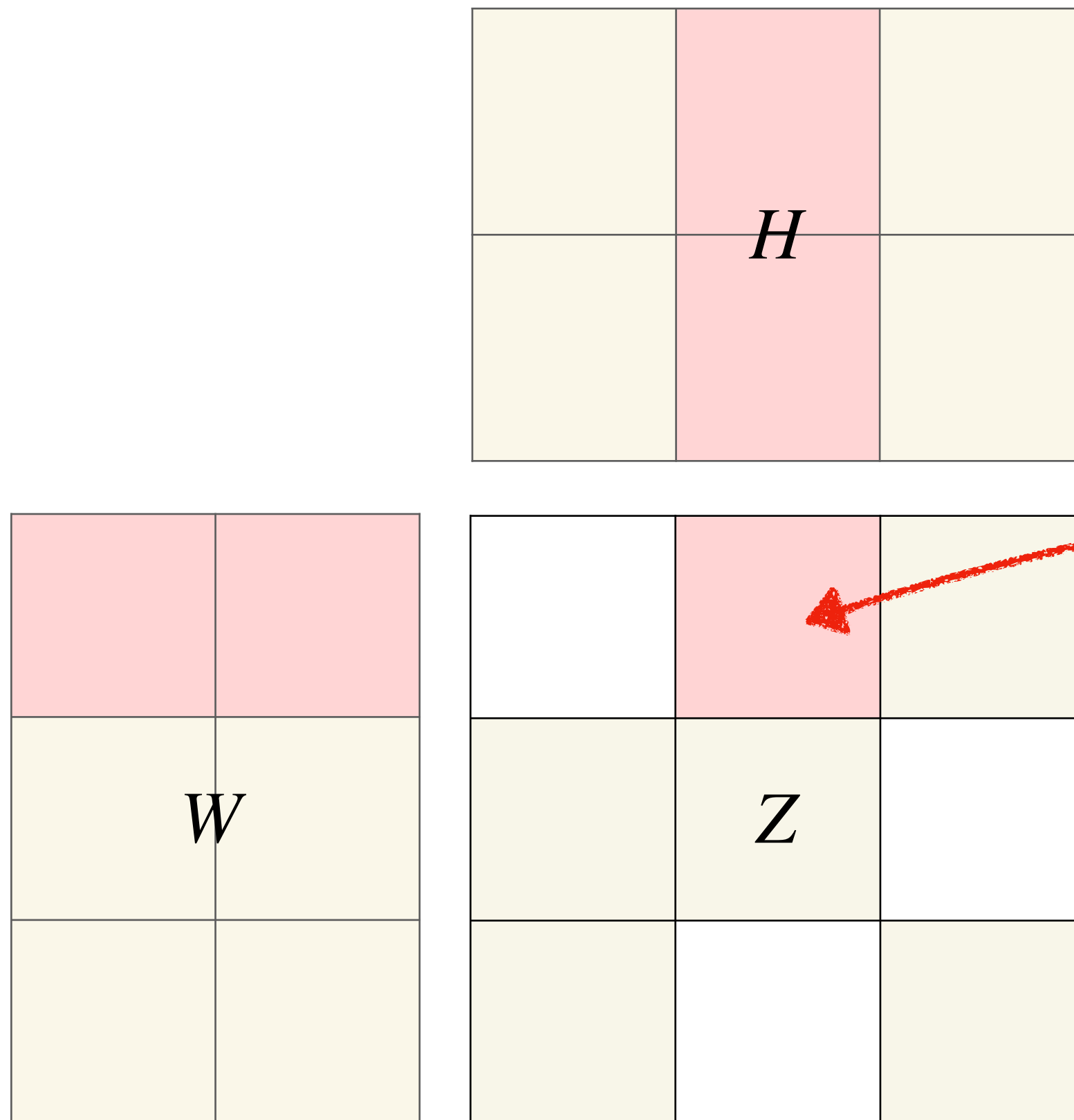
: local loss at position $z = (i, j)$

- $L'(\theta) = \sum_z L'_z(\theta)$

- $\hat{L}'(\theta) = NL'_z(\theta) \quad N = |Z|$

ex) $\hat{L}'(\theta) = 6L'_{z_1}(\theta)$

SGD for Matrix Factorization



Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values W_0 and H_0

while not converged **do** */* step */*

Select a training point $(i, j) \in Z$ uniformly at random.

$$W'_{i*} \leftarrow W_{i*} - \epsilon_n N \frac{\partial}{\partial W_{i*}} l(V_{ij}, W_{i*}, H_{*j})$$

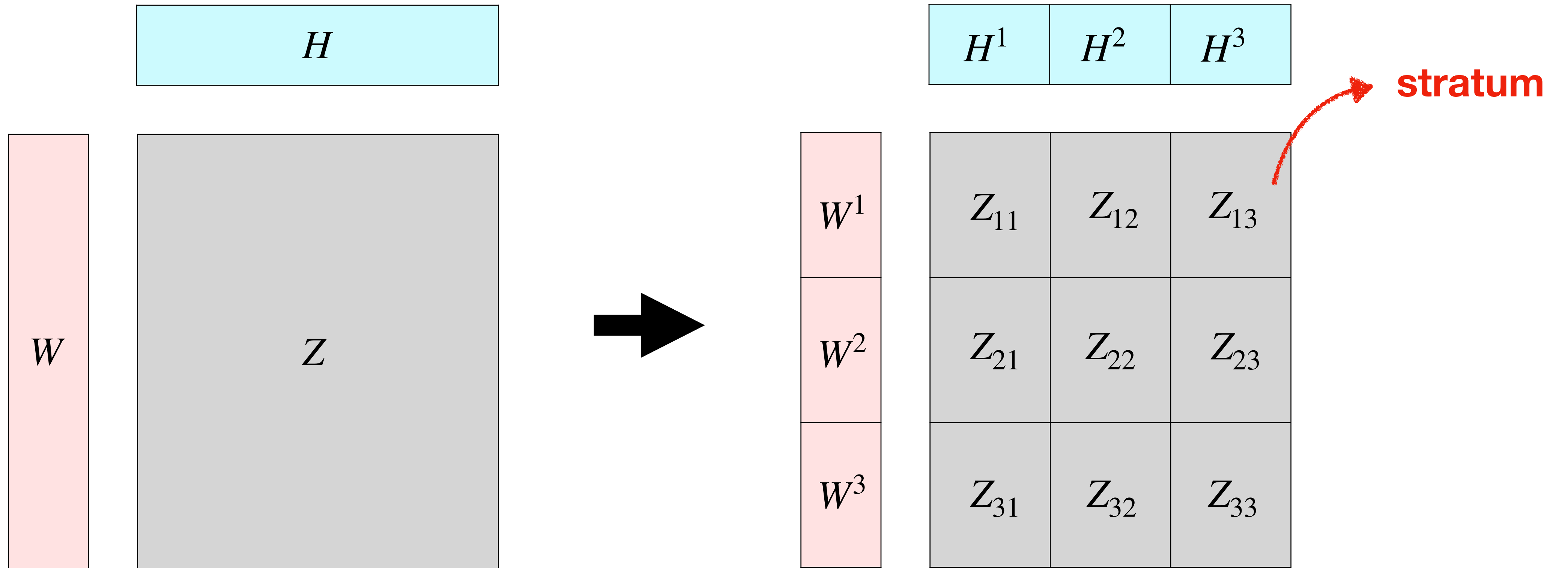
$$H_{*j} \leftarrow H_{*j} - \epsilon_n N \frac{\partial}{\partial H_{*j}} l(V_{ij}, W_{i*}, H_{*j})$$

$$W_{i*} \leftarrow W'_{i*}$$

end while

SSGD

- Stratified SGD



SSGD

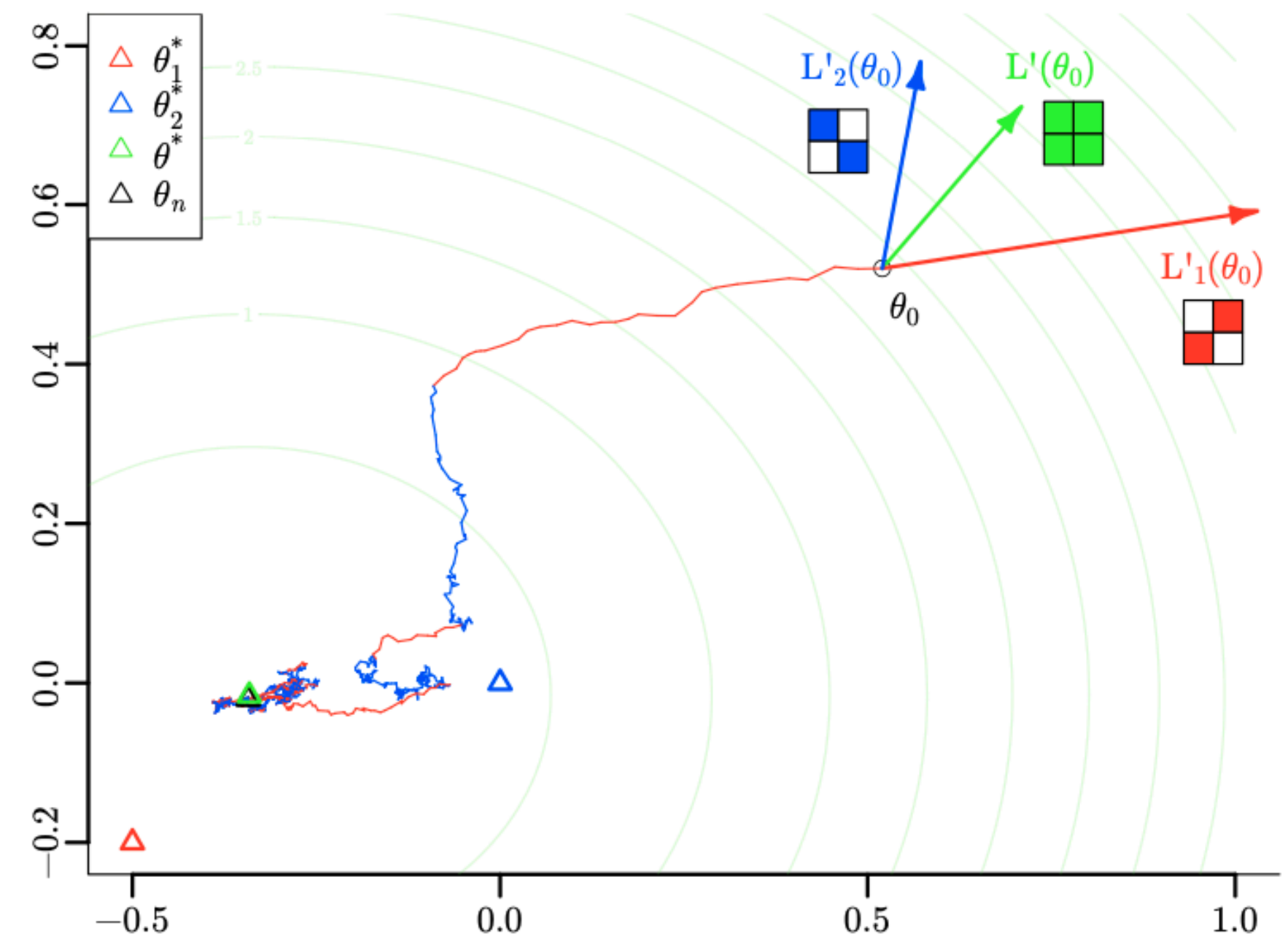
- SSGD Algorithm

- Loss function is decomposed into weighted sum of local loss functions

$$L(\theta) = w_1 L_1(\theta) + w_2 L_2(\theta) + \dots + w_q L_q(\theta)$$

- Stratum sequence $\{\gamma_n\} = \{1, \dots, q\}$

- Update $\theta_{n+1} = \Pi_H[\theta_n - \epsilon_n \hat{L}'_{\gamma_n}(\theta_n)]$



SSGD

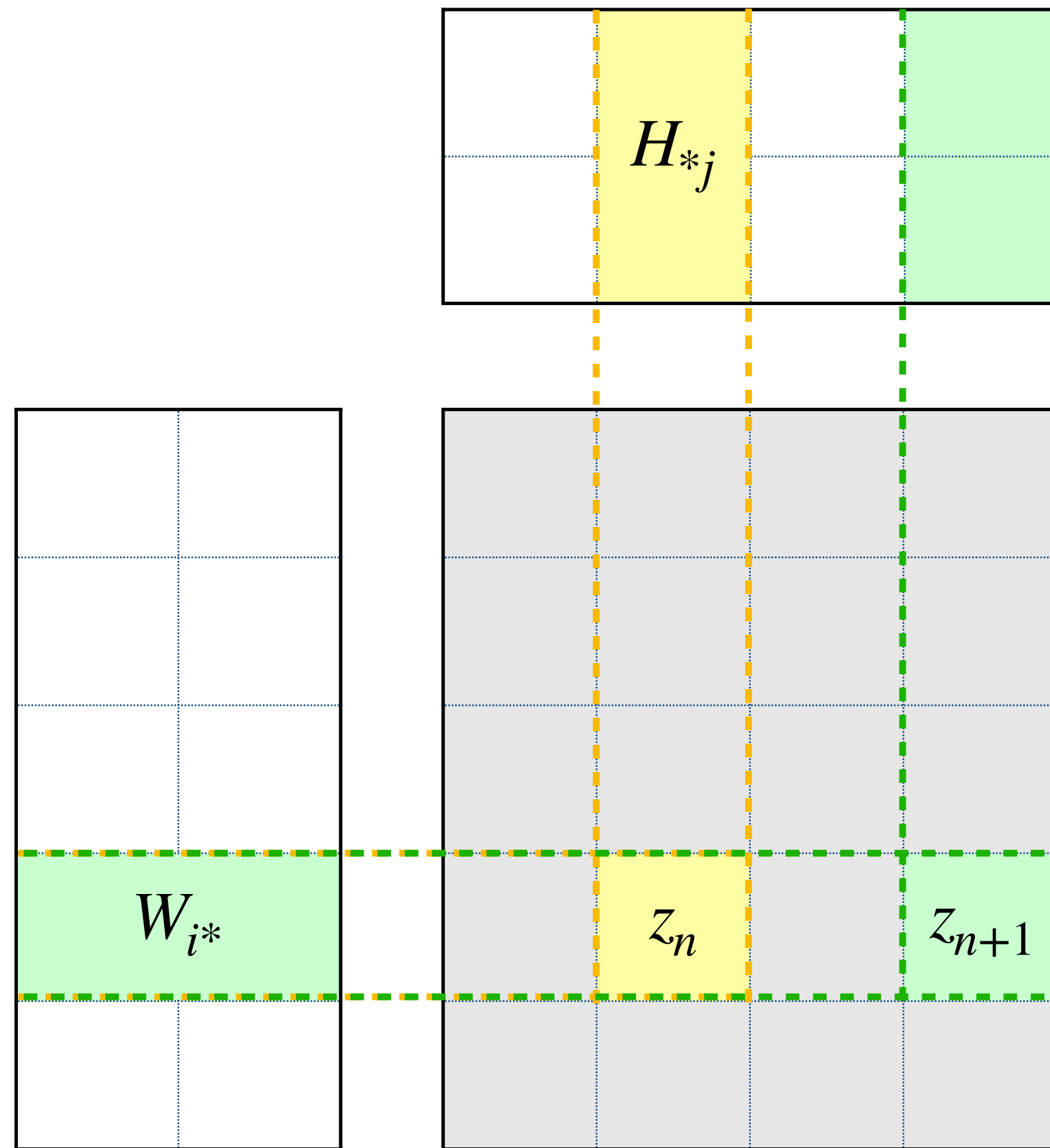
Q) Does SSGD converge?

A) SSGD converges under several conditions

1. Step-size conditions
2. Loss conditions
3. Stratification conditions
4. Stratum-sequence conditions

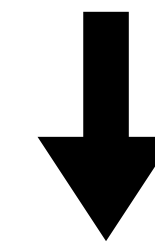
Details in paper

SGD cannot be used directly for rank-r factorization



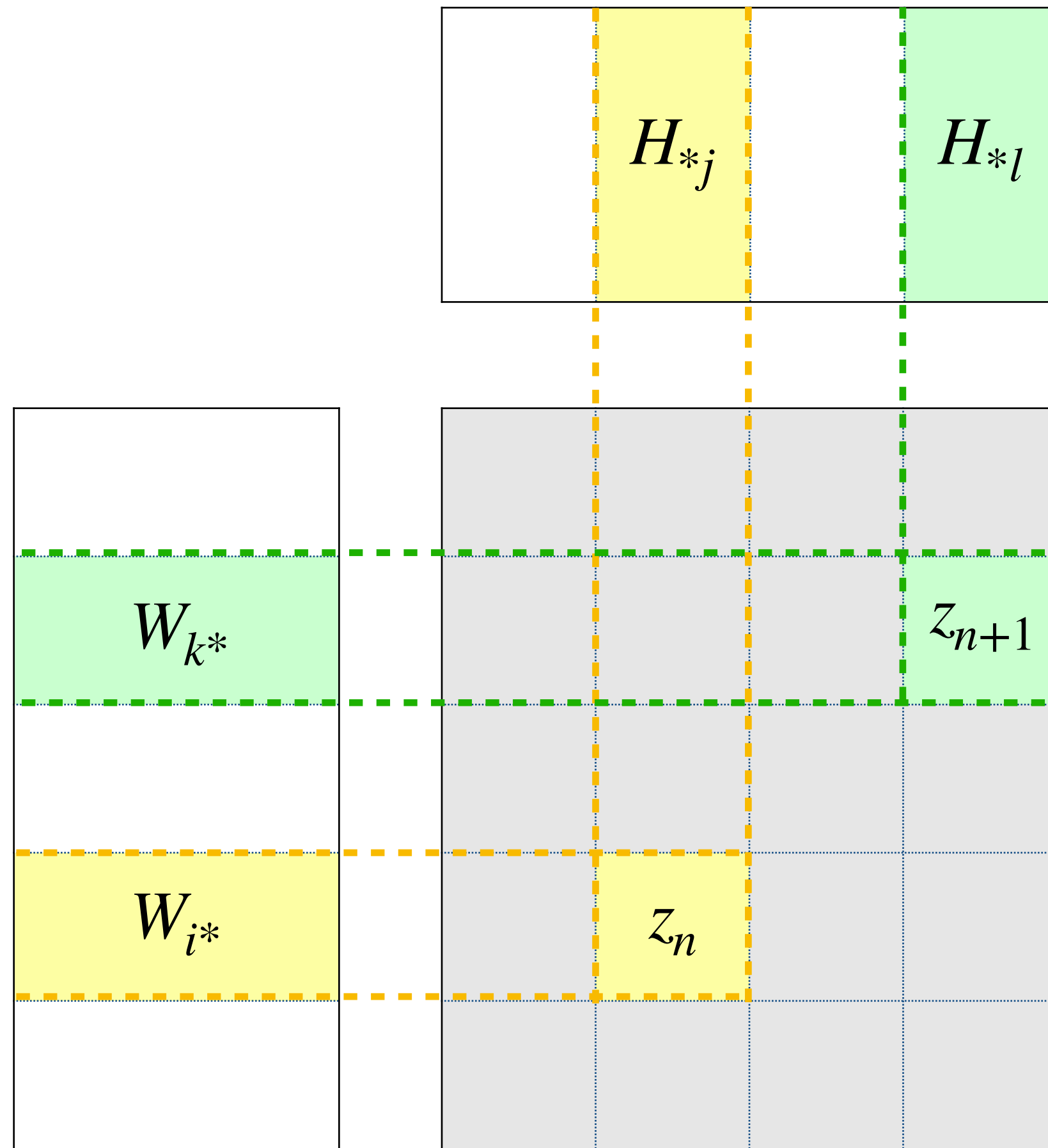
Individual steps depend on each other

$$\theta_{n+1} = \theta_n - \epsilon \hat{L}'(\theta_n)$$



Stratify the training set Z into strata so that each individual stratum are independent

DSGD



- Distributed SGD
 - “Interchangeability”
: share neither row nor column
 - Two machines can be updated at the same time

Read W_{i*} and H_{*j}
Compute gradient of $L'(W_{i*}, H_{*j})$
Update W_{i*} and H_{*j}

Read W_{k*} and H_{*l}
Compute gradient of $L'(W_{k*}, H_{*l})$
Update W_{k*} and H_{*l}

No conflicts!

DSGD

- The idea behind DSGD is to find such blocks that have no data in common and update them simultaneously
- **Goal** : minimize the loss function and update factor vectors (W, H)

→ **Proceeded in stratum**

DSGD

Algorithm 2 DSGD for Matrix Factorization

Require: Z , W_0 , H_0 , cluster size d

$W \leftarrow W_0$ and $H \leftarrow H_0$

Block $Z / W / H$ into $d \times d / d \times 1 / 1 \times d$ blocks

while not converged **do** */* epoch */*

 Pick step size ϵ

for $s = 1, \dots, d$ **do** */* subepoch */*

 Pick d blocks $\{Z^{1j_1}, \dots, Z^{dj_d}\}$ to form a stratum

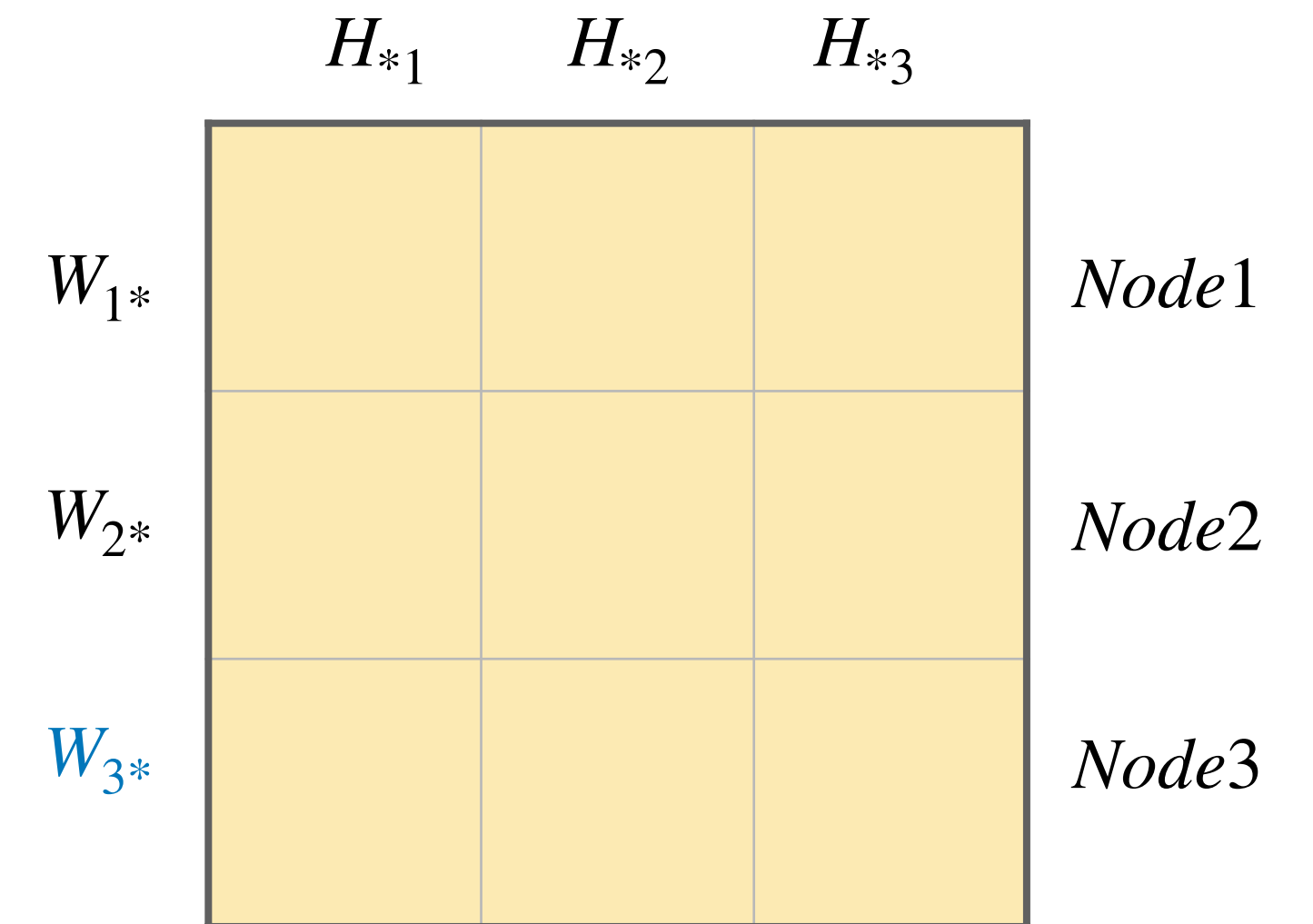
for $b = 1, \dots, d$ **do** */* in parallel */*

 Run SGD on the training points in Z^{bj_b} (step size = ϵ)

end for

end for

end while



DSGD

Algorithm 2 DSGD for Matrix Factorization

Require: Z , W_0 , H_0 , cluster size d

$W \leftarrow W_0$ and $H \leftarrow H_0$

Block $Z / W / H$ into $d \times d / d \times 1 / 1 \times d$ blocks

while not converged **do** */* epoch */*

 Pick step size ϵ

for $s = 1, \dots, d$ **do** */* subepoch */*

 Pick d blocks $\{Z^{1j_1}, \dots, Z^{dj_d}\}$ to form a stratum

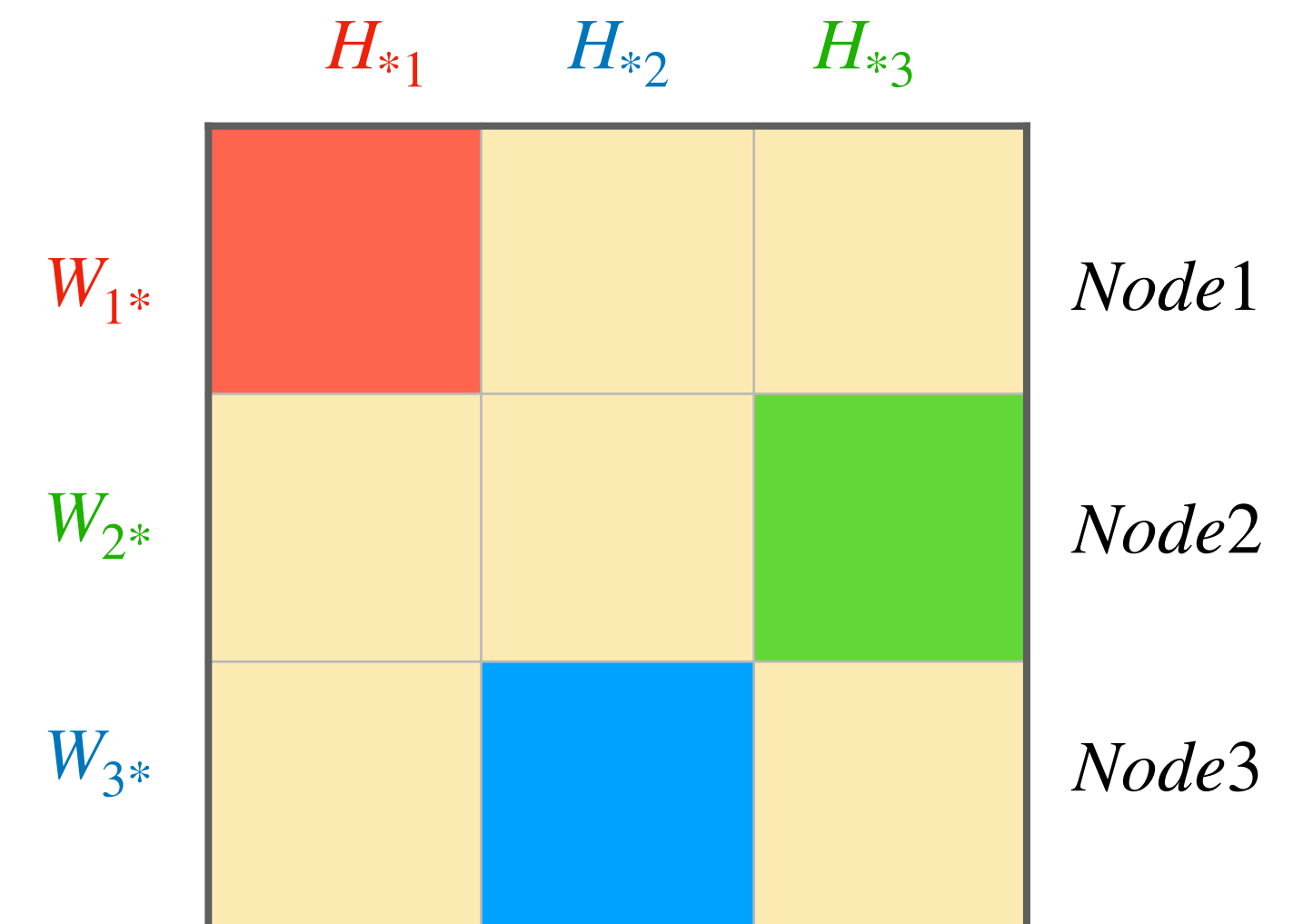
for $b = 1, \dots, d$ **do** */* in parallel */*

 Run SGD on the training points in Z^{bj_b} (step size = ϵ)

end for

end for

end while



Sequence is chosen under SSGD conditions → **DSGD will converge!**

DSGD

Algorithm 2 DSGD for Matrix Factorization

Require: Z , W_0 , H_0 , cluster size d

$W \leftarrow W_0$ and $H \leftarrow H_0$

Block Z / W / H into $d \times d$ / $d \times 1$ / $1 \times d$ blocks

while not converged **do** */* epoch */*

 Pick step size ϵ

for $s = 1, \dots, d$ **do** */* subepoch */*

 Pick d blocks $\{Z^{1j_1}, \dots, Z^{dj_d}\}$ to form a stratum

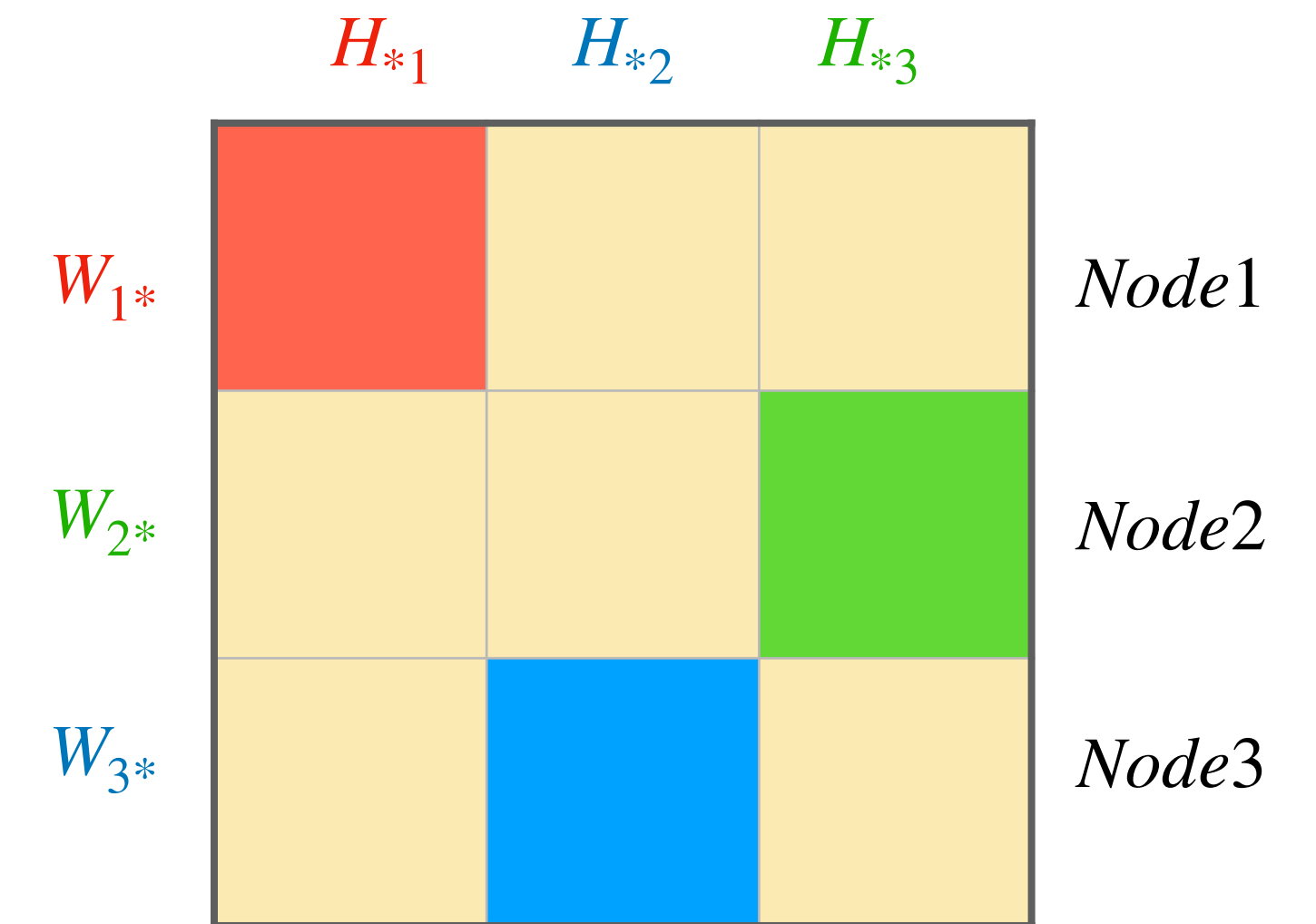
for $b = 1, \dots, d$ **do** */* in parallel */*

 Run SGD on the training points in Z^{bj_b} (step size = ϵ)

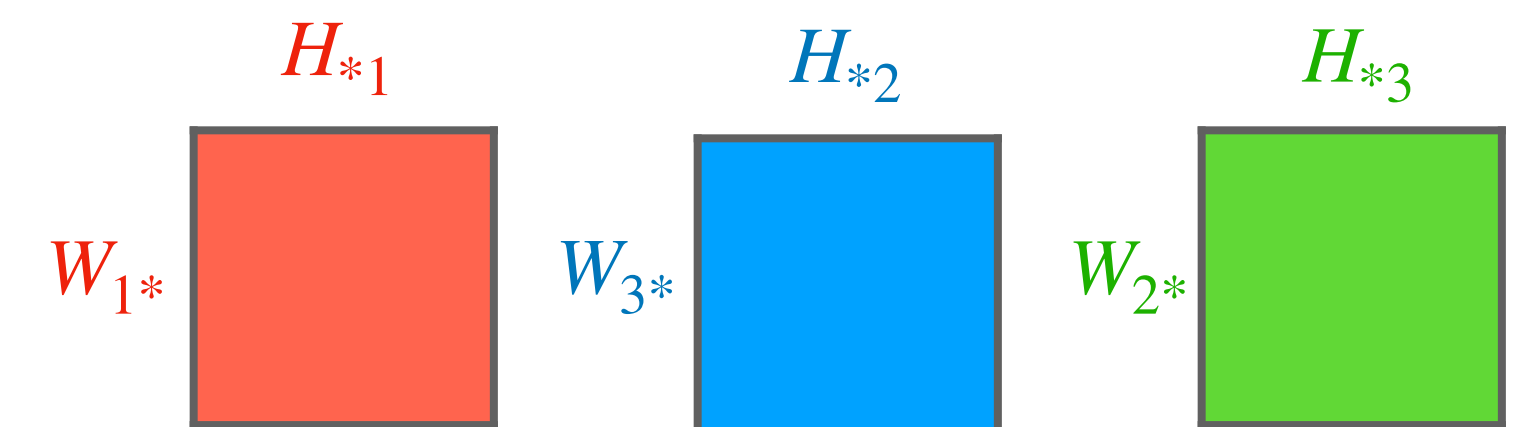
end for

end for

end while



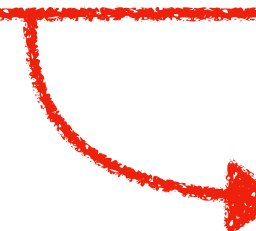
Run SGD on



DSGD

- Sum up losses in each stratum

$$L(W, H) = \sum_{s=1}^q w_s L_s(W, H)$$

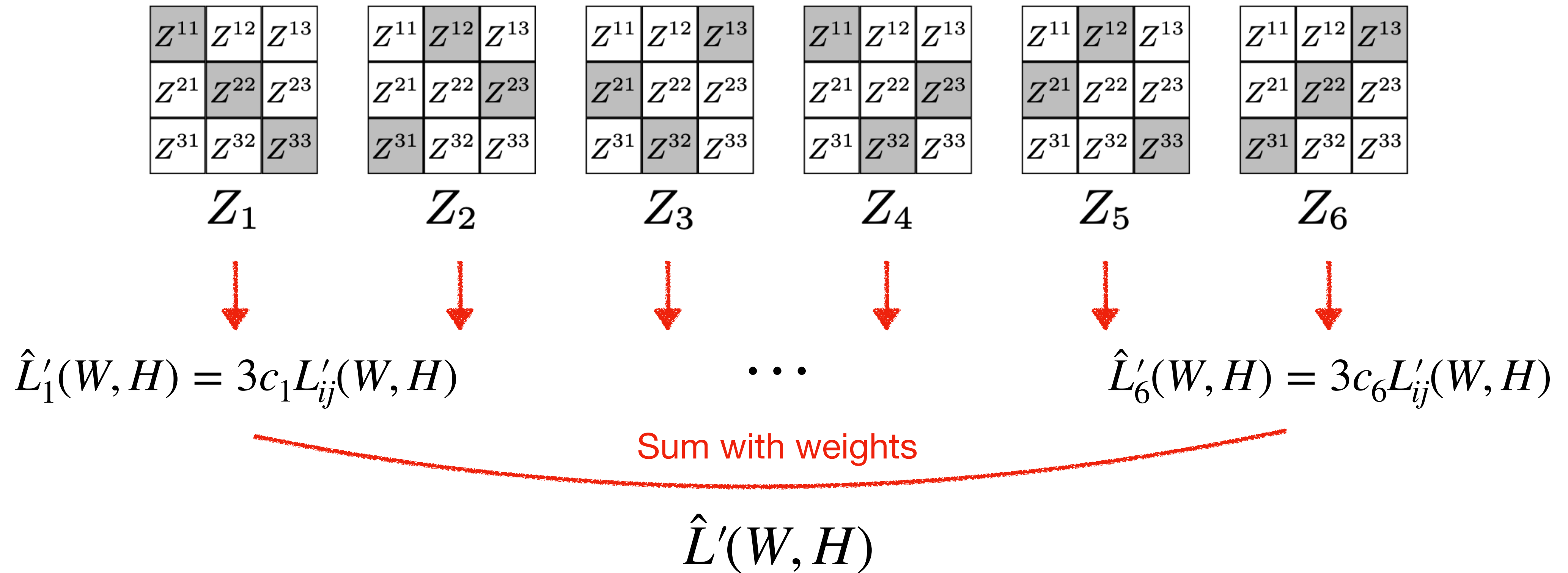

$$L_s(W, H) = c_s \sum_{(i,j) \in Z_s} L_{ij}(W, H) \quad *c_s: \text{stratum-specific constant}$$

- Gradient estimate

$$\hat{L}'_s(W, H) = N_s c_s L'_{ij}(W, H) \quad *N_s = |Z_s|: \text{size of stratum}$$

DSGD

- Distributed SGD

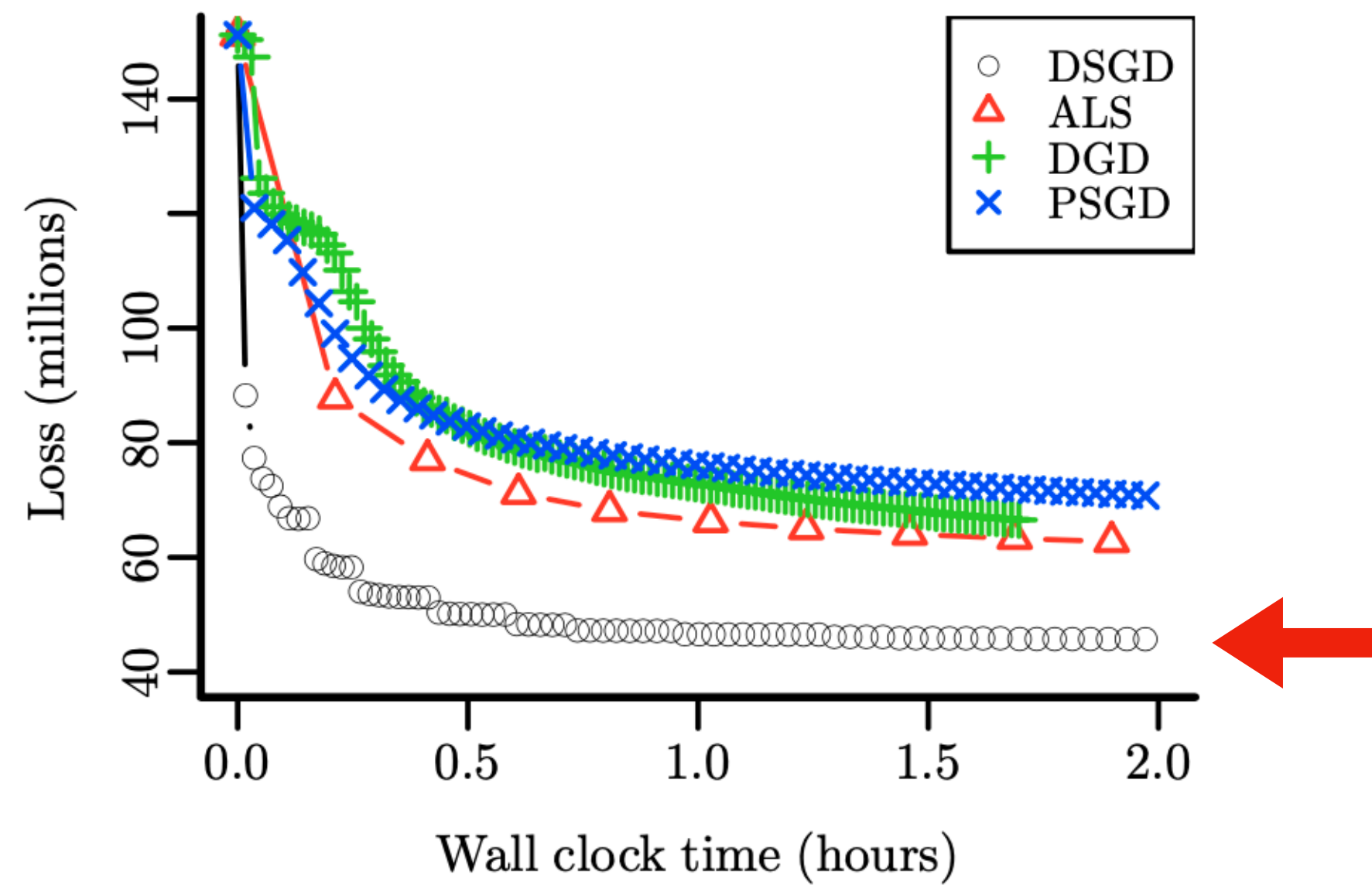


Experiment

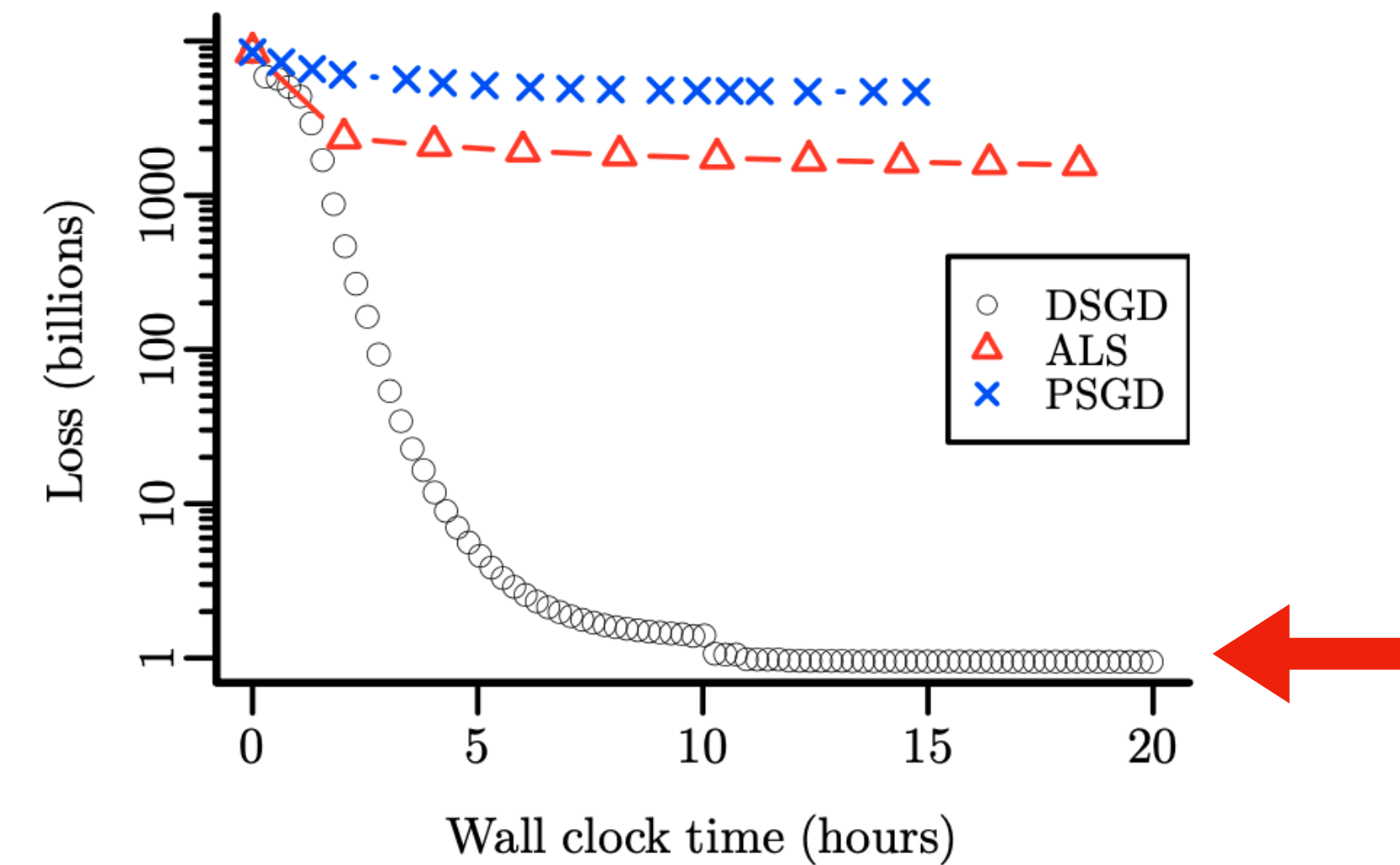
- Compared various factorization algorithms w.r.t
 - Convergence
 - Runtime efficiency
 - Scalability
- Implement on MapReduce
 - In-memory experiment : R, C (Netflix competition dataset)
 - Large scaling experiment : Hadoop (larger synthetic dataset)

Experiment

1. Convergence & 2. Runtime Efficiency



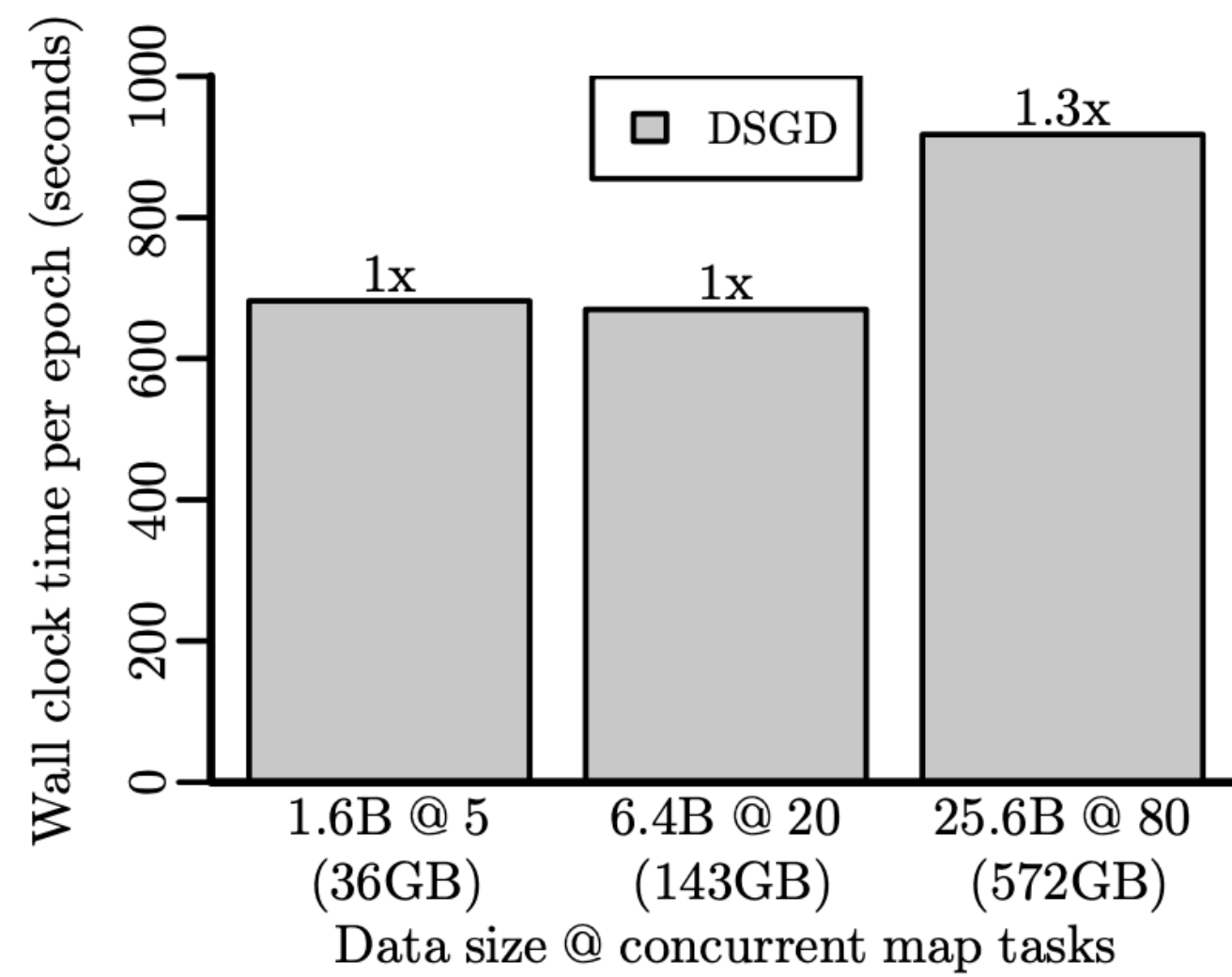
(a) Netflix data (NZSL, R cluster @ 64)



(b) Synth. data (L2, $\lambda = 0.1$, R cluster @ 64)

Experiment

3. Scalability



(c) Scalability (Hadoop cluster)

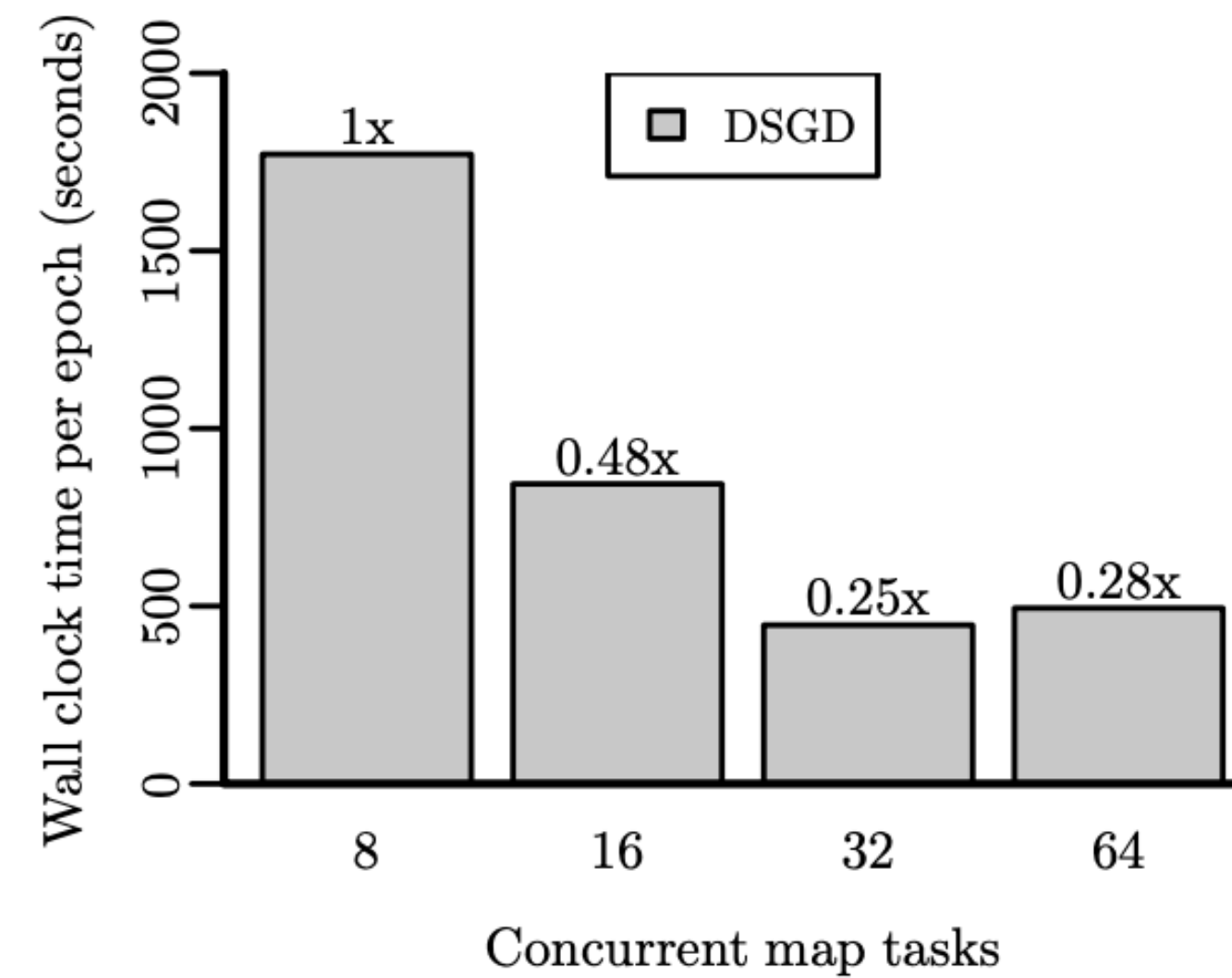


Figure 3: Speed-up experiment (Hadoop cluster, 143GB data)

Conclusion

- Develop stratified version of SGD (**SSGD**)
- Refine SSGD to obtain **DSGD**, a distributed matrix-factorization algorithm that can efficiently handle web-scale matrices
- Future Work
 - Investigate alternative loss functions, regularizations, stratification schemes, emerging distributed-processing platforms
 - Extend to other applications

Thank you