# Appendix for the Paper
# "Integrating Independent Layer-Wise Rank Selection with Low-Rank SVD Training for Model Compression: A Theory-Driven Approach"

## A  Proof of Theorems 1, 2, and 3

**Proposition 1** (Theorem 4.2 [Wright and Ma, 2022]). *Let $W \in \mathbb{R}^{m \times n}$ be a matrix, and $r = \text{rank}(W)$. $W$ can be decomposed as $U\Sigma V^T$, where $U \in \mathbb{R}^{m \times r}$ and $V \in \mathbb{R}^{n \times r}$, such that $UU^T = I$ and $VV^T = I$, $\Sigma \in \mathbb{R}^{r \times r}$ is a is diagonal matrix, i.e., $\Sigma = \text{diag}(\sigma)$, $\sigma = [\sigma_1, \sigma_2, \ldots, \sigma_r]$, and $\sigma_k (k \in [r])$ are singular values of $W$, where $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$. Then, we have $W = \sum_{i=1}^{r} \sigma_i U_{:,i} V_{i,:}^T$.*

**Lemma 1** (Eckart–Young–Mirsky Theorem [Golub *et al.*, 1987]). *Let $W \in \mathbb{R}^{m \times n}$ be a matrix, and $r = \text{rank}(W)$. Following the same settings in Proposition 1, we define $W_k$ to be the best rank-k approximation to $W$ in the spectral norm, i.e., $W_k = U_k \Sigma_k V_k^T = \sum_{i=1}^{k} \sigma_i U_{:,i} V_{i,:}^T$, where $U_k, \Sigma_k, V_k$ are top-k vectors truncated from $U, \Sigma, V$. Then, we have $||W - W_k||_2 = \sigma_{k+1}$, where $|| \cdot ||_2$ stands for the spectral norm.*

*Proof.* We have

$$||W - W_k||_2 = ||\sum_{i=1}^{r} \sigma_i V_{i,:} U_{:,i}^T - \sum_{i=1}^{k} \sigma_i V_{i,:} U_{:,i}^T||_2 \qquad \text{(Proposition 1)}$$

$$= ||\sum_{i=k+1}^{r} \sigma_i V_{i,:} U_{:,i}^T||_2$$

$$= \sigma_{k+1} \qquad \text{(The definition of spectral norm)}$$

$\square$

**Theorem 1** (The output difference bound for rank-$k$ approximation over $L$-layer neural networks). *We denote $a^l$ to be the activation function for the l-th layer, and assume $a^l$ is $\rho_l$-Lipschitz and $a^l(0) = 0$ for all $l \in [1, L]$. Let $X^0$ be the initial input vector, $X^l$ and $X_k^l$ be the output vectors as a result of passing the full-rank matrix $W^l$ and low-rank matrix $W_k^l$ through the l-th layer, respectively, and $\sigma_i^l$ be the i-th singular value of $W^l$. We define $k^l$ such that the top $k^l$ largest singular values of the full-rank matrix $W^l$ are kept in the corresponding low-rank SVD approximated matrix $W_k^l$ in layer l. Then, the output difference from rank-k approximation over L-layer feed-forward networks $||X^L - X_k^L||_2$ is upper-bounded by $\left( \prod_{l=1}^{L} \rho_l \sigma_1^l \right) \left( \sum_{l=1}^{L} \frac{\sigma_{k^l+1}^l}{\sigma_1^l} \right) ||X^0||_2$.*

*Proof.* For the output difference at layer $l + 1$, we have

$$||X^{l+1} - X_k^{l+1}||_2 = ||a^{l+1}(W^{l+1}X^l) - a^{l+1}(W_k^{l+1}X_k^l)||_2$$

$$\leq \rho_{l+1}||W^{l+1}X^l - W_k^{l+1}X_k^l||_2$$

$$\leq \rho_{l+1}||W^{l+1}X^l - W_k^{l+1}X^l + W_k^{l+1}X^l - W_k^{l+1}X_k^l||_2$$

$$\leq \rho_{l+1}||W^{l+1}X^l - W_k^{l+1}X^l||_2 + \rho_{l+1}||W_k^{l+1}X^l - W_k^{l+1}X_k^l||_2$$

$$\leq \rho_{l+1}||W^{l+1} - W_k^{l+1}||_2 \cdot ||X^l||_2 + \rho_{l+1}||W_k^{l+1}||_2 \cdot ||X^l - X_k^l||_2$$

$$\leq \rho_{l+1}\sigma_{k^{l+1}+1}^{l+1}||X^l||_2 + \rho_{l+1}\sigma_1^{l+1} \cdot ||X^l - X_k^l||_2. \qquad \text{(Lemma 1)}$$

Also, we notice that,

$$||X^l||_2 \leq \left( \prod_{i=1}^{l} \rho_i ||W^i||_2 \right) ||X^0||_2 = \left( \prod_{i=1}^{l} \rho_i \sigma_1^i \right) ||X^0||_2.$$

If we let $c_{l+1} = \rho_{l+1}\sigma_{k^{l+1}+1}^{l+1} \left( \prod_{i=1}^{l} \rho_i \sigma_1^i \right)$ and $d_{l+1} = \rho_{l+1}\sigma_1^{l+1}$, then we have,

$$||X^{l+1} - X_k^{l+1}||_2 \leq c_{l+1}||X^0||_2 + d_{l+1}||X^l - X_k^l||_2.$$

By induction over $L$ layers, we have,

$$\|X^L - X_k^L\|_2 \leq c_L\|X^0\|_2 + d_L\|X^{L-1} - X_k^{L-1}\|_2$$

$$\leq c_L\|X^0\|_2 + d_L\left(c_{L-1}\|X^0\|_2 + d_{L-1}\|X^{L-2} - X_k^{L-2}\|_2\right)$$

$$= (c_L + d_L c_{L-1})\|X^0\|_2 + d_L d_{L-1}\|X^{L-2} - X_k^{L-2}\|_2$$

$$\leq (c_L + d_L c_{L-1} + d_L d_{L-1} c_{L-2})\|X^0\|_2 + d_L d_{L-1} d_{L-2}\|X^{L-3} - X_k^{L-3}\|_2$$

$$\cdots$$

$$\leq \left(c_L + d_L c_{L-1} + d_L d_{L-1} c_{L-2} + \cdots + \left(\prod_{l=0}^{L-2} d_{L-l}\right)c_1\right)\|X^0\|_2 + \left(\prod_{l=1}^{L-1} d_{L-l}\right)\|X^0 - X^0\|_2$$

$$= \left(c_L + d_L c_{L-1} + d_L d_{L-1} c_{L-2} + \cdots + \left(\prod_{l=0}^{L-2} d_{L-l}\right)c_1\right)\|X^0\|_2$$

$$= \left(\prod_{l=1}^{L}\rho_l\sigma_1^l\right)\left(\sum_{l=1}^{L}\frac{\sigma_{k^l+1}^l}{\sigma_1^l}\right)\|X^0\|_2$$

$\square$

**Theorem 2** (The loss error bound from rank-$k$ approximation in classification problems)**.** *Following the settings in Theorem 1, we consider a $C$-class classification problem. Let $X_i^L \in \mathbb{R}^C$ and $X_{i,k}^L \in \mathbb{R}^C$ be the output logits when feeding a input $X_i^0$ sampled from the training dataset $\mathcal{D}^{tr}$, e.g., $\mathcal{D}^{tr} = \{X_i^0, y_i\}_{i=1}^{R}$, from the full-rank parameter space $\mathcal{W}$ and low-rank parameter space $\mathcal{W}_k$, respectively. Particularly, $\mathcal{W} = \{W^1, W^2, \ldots, W^L\}$, $\mathcal{W}_k = \{W_k^1, W_k^2, \ldots, W_k^L\}$, $X_i^L = f_{\mathcal{W}}(X_i^0)$, $X_{i,k}^L = f_{\mathcal{W}_k}(X_{i,k}^0)$, and $\|X_i^0\|_2 \leq B$, for $\forall\ i \in [1,R]$. Let $z_i = \texttt{softmax}(X_i^L)$ and $z_{i,k} = \texttt{softmax}(X_{i,k}^L)$, where* \texttt{softmax} *is the softmax function. We consider the cross-entropy function as the loss function, i.e., $g(z,y) = -y^T\log(z)$. Let $L(\mathcal{W}; X_i^0) = g(z_i, y_i)$ and $L(\mathcal{W}_k; X_i^0) = g(z_{i,k}, y_i)$. Now, we set $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta, \forall\ l \in [1, L]$. Then, we have, for $\forall\ \epsilon > 0$, $\exists\ \delta = \frac{\epsilon}{\sqrt{2}BL\left(\prod_{l=1}^{L}\rho_l\sigma_1^l\right)}$, s.t. $|L(\mathcal{W};\mathcal{D}^{tr}) - L(\mathcal{W}_k;\mathcal{D}^{tr})| < \epsilon$.*

*Proof.* First, we define $\frac{\partial L(\mathcal{W}, X_i^0)}{\partial x} \in \mathbb{R}^C$ to be the partial derivative of $L(\mathcal{W}, X_i^0)$ with respect to the output layer $X_k^L$, for the given input $X_i^0$. Then, we have

$$\left\|\frac{\partial L(\mathcal{W}, X_i^0)}{\partial x}\right\|_2 = \|z_i - y_i\|_2.$$

Without loss of generality, for $z_i = [z_{i,1}, \ldots, z_{i,C}] \in \mathbb{R}^C, y_i = [y_{i,1}, \ldots, y_{i,C}] \in \mathbb{R}^C$, we assume $y_{i,c} = 0$, where $c \in [1, C-1]$ and $y_{1,C} = 1$ and $z_{i,1} + \cdots + z_{i,C} = 1$, where $z_{i,j} \in [0,1], \forall\ j \in [1, C]$, Then

$$\left\|\frac{\partial L(\mathcal{W}, X_i^0)}{\partial x}\right\|_2 = \|z_i - y_i\|_2$$

$$= \sqrt{(z_{i,1} - y_{i,1})^2 + (z_{i,2} - y_{i,2})^2 + \cdots + (z_{i,C} - y_{i,C})^2}$$

$$= \sqrt{z_{i,1}^2 + z_{i,2}^2 + \cdots + z_{i,C-1}^2 + (z_{i,C} - 1)^2}$$

$$= \sqrt{z_{i,1}^2 + \cdots + z_{i,C-1}^2 + z_{i,C}^2 - 2z_{i,C} + 1}$$

$$= \sqrt{z_{i,1}^2 + \cdots + z_{i,C-1}^2 + z_{i,C}^2 - 2(z_{i,1} + \cdots + z_{i,C})z_{i,C} + 1}$$

$$= \sqrt{(z_{i,1} - z_{i,C})^2 + \cdots + (z_{i,C-1} - z_{i,C})^2 - Cz_{i,C}^2 + 1}$$

$$\leq \sqrt{(z_{i,1} + \cdots + z_{i,C-1} - (C-1)z_{i,C})^2 - Cz_{i,C}^2 + 1}$$

(The equality case holds when $(z_{i,j} - z_{i,C})(z_{i,j'} - z_{i,C}) = 0$ for $j, j' \in [1, C-1]$)

$$= \sqrt{(1 - Cz_{i,C})^2 - Cz_{i,C}^2 + 1}$$

$$\leq \sqrt{2} \qquad\qquad \text{(The equality case holds when } z_{i,C} = 0)$$

Overall, the equality case occurs when $z_{i,j} = 1$ if $j \neq C$ and all rest $z_{i,j'}$ are 0 if $j' \neq j$ and $j' \in [1, C]$.

Let $\delta = \frac{\epsilon}{\sqrt{2}BL\left(\prod_{l=1}^{L} \rho_l \sigma_1^l\right)}$, where $||X_i^0||_2 \leq B$ and $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta, \forall l \in [L]$. With Lagrange's mean value theorem, we have

$$|L(\mathcal{W}; X_i^0) - L(\mathcal{W}_k; X_i^0)| \leq \max_{\substack{x \in tX_k^L + (1-t)X^L \\ t \in [0,1]}} \left\{ \left|\left| \frac{\partial L}{\partial x} \right|\right|_2 \right\} ||X^L - X_k^L||_2$$

$$\leq \sqrt{2}||X^L - X_k^L||_2$$

$$\leq \sqrt{2}||X_i^0||_2 \sum_{l=1}^{L} \frac{\sigma_{k^l+1}^l}{\sigma_1^l} \left( \prod_{l=1}^{L} \rho_l \sigma_1^l \right)$$

$$< \sqrt{2}B\delta L \left( \prod_{l=1}^{L} \rho_l \sigma_1^l \right) = \epsilon.$$

Finally, $|L(\mathcal{W}; \mathcal{D}^{tr}) - L(\mathcal{W}_k; \mathcal{D}^{tr})| \leq \frac{1}{R} \sum_{i=1}^{R} |L(\mathcal{W}; X_i^0) - L(\mathcal{W}_k; X_i^0)| < \epsilon.$ $\qquad\square$

**Theorem 3** (The loss error bound from rank-$k$ approximation in regression problems)**.** *Following the settings in Theorem 1, we consider a regression problem. Let $X_i^L$ and $X_{i,k}^L$ be the output at layer $L$ when feeding a input $X_i^0$ sampled from the training dataset $\mathcal{D}^{tr}$, e.g., $\mathcal{D}^{tr} = \{X_i^0, y_i\}_{i=1}^R$, from the full-rank parameter space $\mathcal{W}$ and low-rank parameter space $\mathcal{W}_k$, respectively. Particularly, $\mathcal{W} = \{W^1, W^2, \ldots, W^L\}$, $\mathcal{W}_k = \{W_k^1, W_k^2, \ldots, W_k^L\}$, $X_i^L = f_{\mathcal{W}}(X_i^0)$, $X_{i,k}^L = f_{\mathcal{W}_k}(X_{i,k}^0)$, and $||X_i^0||_2 \leq B$, for $\forall i \in [1, R]$. We consider the loss function as $g(z, y) = ||z - y||_2$. Let $L(\mathcal{W}; X_i^0) = g(X_i^L, y_i)$ and $L(\mathcal{W}_k; X_i^0) = g(X_{i,k}^L, y_i)$. Now, we set $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta, \forall l \in [1, L]$. Then, we have, for $\forall \epsilon > 0$, $\exists \delta = \frac{\epsilon}{BL\left(\prod_{l=1}^{L} \rho_l \sigma_1^l\right)}$, s.t. $|L(\mathcal{W}; \mathcal{D}^{tr}) - L(\mathcal{W}_k; \mathcal{D}^{tr})| < \epsilon.$*

*Proof.* Note that

$$|L(\mathcal{W}; X_i^0) - L(\mathcal{W}_k; X_i^0)| = \left| \, ||X_i^L - y_i||_2 - ||X_{i,k}^L - y_i||_2 \, \right|$$

$$\leq ||(X_i^L - y_i) - (X_{i,k}^L - y_i)||_2$$

$$= ||X_i^L - X_{i,k}^L||_2.$$

Let $\delta = \frac{\epsilon}{BL\sqrt{2}\left(\prod_{l=1}^{L} \rho_l \sigma_1^l\right)}$, where $||X_i^0||_2 \leq B$ and $\frac{\sigma_{k^l+1}^l}{\sigma_1^l} < \delta, \forall l \in [L]$. We have

$$|L(\mathcal{W}; X_i^0) - L(\mathcal{W}_k; X_i^0)| \leq ||X^L - X_k^L||_2$$

$$\leq ||X_i^0||_2 \sum_{l=1}^{L} \frac{\sigma_{k^l+1}^l}{\sigma_1^l} \left( \prod_{l=1}^{L} \rho_l \sigma_1^l \right)$$

$$< B\delta L \left( \prod_{l=1}^{L} \rho_l \sigma_1^l \right) = \epsilon.$$

Finally, $|L(\mathcal{W}; \mathcal{D}^{tr}) - L(\mathcal{W}_k; \mathcal{D}^{tr})| \leq \frac{1}{R} \sum_{i=1}^{R} |L(\mathcal{W}; X_i^0) - L(\mathcal{W}_k; X_i^0)| < \epsilon.$ $\qquad\square$

# B Statistical Analysis of Our Pilot Study in Section 4.1

|  | Full-Rank | Low-Rank | | | | |
|---|---|---|---|---|---|---|
| $\epsilon$ | 0 | 0.17 | 0.23 | 0.28 | 0.33 | 0.56 |
| $\delta$ | 0 | 0.015 | 0.021 | 0.025 | 0.03 | 0.047 |
| $k^1$ (Layer 1) | 2 | 2 | 2 | 2 | 2 | 2 |
| $k^2$ (Layer 2) | 100 | 9 | 8 | 7 | 6 | 3 |
| $k^3$ (Layer 3) | 3 | 3 | 3 | 3 | 3 | 3 |

Table 4: Statistics of our pilot study.

To validate the feasibility of identifying $\delta$ based on our derived $\epsilon$-$\delta$ correlation and determining the optimal $k^l$, we conduct a pilot study using a simple 3-layer feed-forward neural network for a ternary classification problem. Particularly, the input is in 2 dimensions; the output is in 3 dimensions; $W^1 \in \mathbb{R}^{2 \times 100}$, $W^2 \in \mathbb{R}^{100 \times 100}$, and $W^1 \in \mathbb{R}^{100 \times 3}$. Here, $W^1$, $W^2$,
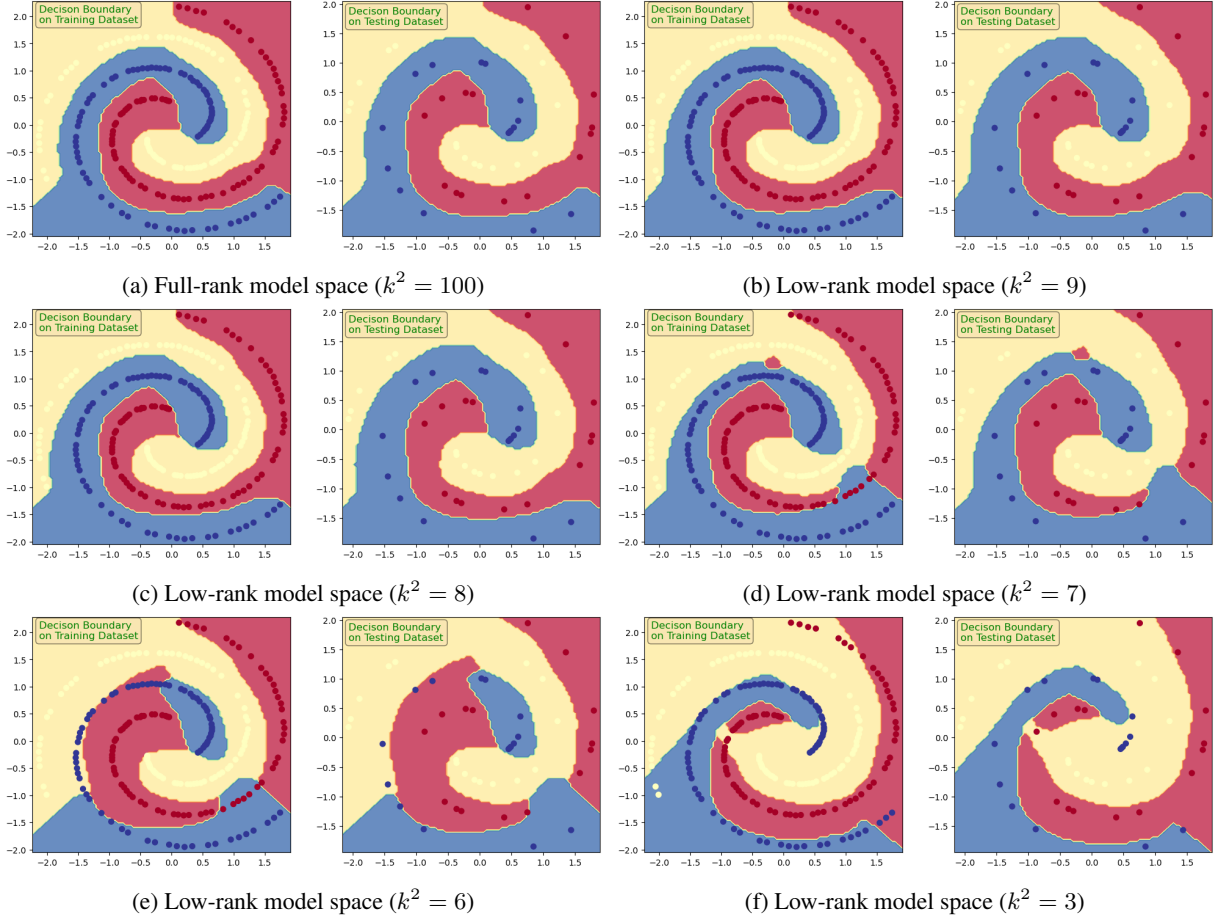
Figure 7: A visualization of the decision boundaries on the training dataset (left column) and testing dataset (right column) through different rank selections.

and $\boldsymbol{W}^3$ are either column full-rank or row full-rank. It means that $\mathrm{rank}(\boldsymbol{W}^1) = 2$, $\mathrm{rank}(\boldsymbol{W}^2) = 100$, and $\mathrm{rank}(\boldsymbol{W}^3) = 3$. Since there is not much room to tune $k$ in $\boldsymbol{W}^1$ and $\boldsymbol{W}^3$, our major focus is on tuning $k$ in $\boldsymbol{W}^2$. The activation functions in the first and second layers are both the ReLU functions. The output from the third layer will be fed into a softmax function to get the prediction probabilities. The loss function is the cross-entropy function. Fig. 7 visualizes the decision boundary on the training (left column in each sub-figure) and testing dataset (right column in each sub-figure) through different rank selections with more details. Table 4 lists the $\epsilon$, $\delta$, and selected $k$ for each layer. As we can find, as the loss error bound $\epsilon$ increases, the corresponding $\delta$ also increases, indicating that more information is truncated as smaller values of $k^l$ are assigned to each layer, respectively. This truncation leads to underfitting in the low-rank model, preventing it from effectively capturing the patterns of the original full-rank model. Particularly, according to the results in Fig. 7, the key turning point in selecting $k^2$ happens from 8 to 7. We can find the decision boundary's shape has appeared different patterns compared with that in larger $k^2$. To some extent, when $k^2 = 8$, it represents the smallest $k^l$, which is sufficient to retain the full-rank model's representational capacity. Further searching would significantly degrade the performance.

## C   Our Proposed Rank Selection Enabled Low-Rank SVD Training Algorithm

We present our proposed rank selection enabled low-rank SVD training algorithm in detail, as illustrated in Algorithms 3, 4, and 5.

---

**Algorithm 3:** Proposed Rank Selection Enabled Low-Rank SVD Training Algorithm

---
**Input:** full-rank parameters $\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V}$, loss error tolerance $\epsilon$, stop searching precision $\Delta\delta$, training epoch $E$
**Output:** low-rank parameters $\mathcal{U}_k$, $\boldsymbol{\Sigma}_k$, $\mathcal{V}_k$

**1** Initialize parameters $\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V}$ ;
**2** $e \leftarrow 1$ ;
**3** **while** $e \leq E$ **do**
**4**    Update $\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V}$ based on loss function $L(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V})$ with an appropriate optimizer and extract the learning loss
        $L_T(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V})$; `// `$L_O(\mathcal{U}, \mathcal{V})$` and `$L_R(\boldsymbol{\Sigma})$` are not used in the next-step rank`
        `selection`
**5**    $\mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k \leftarrow$ `RankSelection`$(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V}, \epsilon, \Delta\delta, L_T(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V}))$;
**6**    $\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V} \leftarrow \mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k$; `// Use truncated models for the next round of training`
**7**    $e \leftarrow e + 1$;
**8** **return** $\mathcal{U}_k$, $\boldsymbol{\Sigma}_k$, $\mathcal{V}_k$ ;

---

---

**Algorithm 4:** Rank Selection Function

---
**1** **Function** `RankSelection` ($\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V}$, $\epsilon$, $\Delta\delta$, $L_T(\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V})$)**:**
**2**    $floss \leftarrow L_T(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V})$;
**3**    $l \leftarrow 0; u \leftarrow 1; \delta \leftarrow (l+u)/2$;
**4**    **while** $|l - u| \geq \Delta\delta$ *or* $|floss - loss| \geq \epsilon$ **do**
**5**       $\mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k \leftarrow$ `Truncation`$(\mathcal{U}, \boldsymbol{\Sigma}, \mathcal{V}, \delta)$;
**6**       $loss \leftarrow L_T(\mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k)$;
**7**       **if** $|floss - loss| < \epsilon$ **then**
**8**          $l \leftarrow \delta; \delta \leftarrow (l+u)/2$;
**9**       **else**
**10**         $u \leftarrow \delta; \delta \leftarrow (l+u)/2$;
**11**    **return** $\mathcal{U}_k$, $\boldsymbol{\Sigma}_k$, $\mathcal{V}_k$ ;

---

---

**Algorithm 5:** Truncation Function

---
**1** **Function** `Truncation` ($\mathcal{U}$, $\boldsymbol{\Sigma}$, $\mathcal{V}$, $\delta$)**:**
**2**    $\mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k \leftarrow \emptyset$;
**3**    **for** *each* $\boldsymbol{U}^l \in \mathcal{U}$, $\boldsymbol{\Sigma}^l \in \boldsymbol{\Sigma}$, $\boldsymbol{V}^l \in \mathcal{V}$ **do**
**4**       $k^l \leftarrow \arg\max_k \{k | \sigma_k^l / \sigma_1^l \geq \delta\}$;
**5**       Truncate top-$k^l$ vectors from $\boldsymbol{U}^l, \boldsymbol{\Sigma}^l, \boldsymbol{V}^l$ to get $\boldsymbol{U}_k^l, \boldsymbol{\Sigma}_k^l, \boldsymbol{V}_k^l$, and add into $\mathcal{U}_k, \boldsymbol{\Sigma}_k, \mathcal{V}_k$ ;
**6**    **return** $\mathcal{U}_k$, $\boldsymbol{\Sigma}_k$, $\mathcal{V}_k$

---