# Randomized SVD, CUR Decomposition, and SPSD Matrix Approximation

**Shusen Wang**

# Outline

- CX Decomposition & Approximate SVD
- CUR Decomposition
- SPSD Matrix Approximation

# CX Decomposition

- Given any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$

- The CX decomposition of $\mathbf{A}$
    1. Sketching: $\mathbf{C} = \mathbf{AP} \in \mathbb{R}^{m \times c}$
    2. Find $\mathbf{X}$ such that $\mathbf{A} \approx \mathbf{CX}$
        - E.g. $\mathbf{X}^{\star} = \mathrm{argmin}_{\mathbf{X}} \left\|\mathbf{A} - \mathbf{CX}\right\|_{\mathrm{F}}^{2} = \mathbf{C}^{\dagger}\mathbf{A}$
        - It costs $O(mnc)$

# CX Decomposition

- Let the sketching matrix $\mathbf{P} \in \mathbb{R}^{n \times c}$ be defined in the table.

- $\min_{\mathrm{rank}(\mathbf{X}) \leq k} \left\| \mathbf{A} - \mathbf{CX} \right\|_F^2 \ \leq \ (1 + \epsilon) \left\| \mathbf{A} - \mathbf{A}_k \right\|_F^2$

| | Uniform sampling | Leverage score sampling | Gaussian projection | SRHT | Count sketch |
|---|---|---|---|---|---|
| $c \geq$ | $O\left( \nu k \left( \log k + \frac{1}{\epsilon} \right) \right)$ | $O\left( k \left( \log k + \frac{1}{\epsilon} \right) \right)$ | $O\left( \frac{k}{\epsilon} \right)$ | $O\left( (k + \log n) \left( \log k + \frac{1}{\epsilon} \right) \right)$ | $O\left( k^2 + \frac{k}{\epsilon} \right)$ |

$\nu$ is the column coherence of $\mathbf{A}_k$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

$$A \approx CX$$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

$$\mathbf{A} \approx \mathbf{CX} = \mathbf{U}_C \boldsymbol{\Sigma}_C \mathbf{V}_C^T \mathbf{X}$$

SVD: $\mathbf{C} = \mathbf{U}_C \boldsymbol{\Sigma}_C \mathbf{V}_C^T \in \mathbb{R}^{m \times c}$

**Time cost**: $O(mc^2)$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

$$\mathbf{A} \approx \mathbf{CX} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^T \mathbf{X} = \mathbf{U}_C \mathbf{Z}$$

Let $\mathbf{\Sigma}_C \mathbf{V}_C^T \mathbf{X} = \mathbf{Z} \in \mathbb{R}^{c \times n}$

SVD: $\mathbf{C} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^T \in \mathbb{R}^{m \times c}$

**Time cost**: $O(mc^2 + nc^2)$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

$$\mathbf{A} \approx \mathbf{CX} = \mathbf{U}_{\mathrm{C}}\boldsymbol{\Sigma}_{\mathrm{C}}\mathbf{V}_{\mathrm{C}}^{\mathrm{T}}\mathbf{X} = \mathbf{U}_{\mathrm{C}}\mathbf{Z} = \mathbf{U}_{\mathrm{C}}\mathbf{U}_{\mathrm{Z}}\boldsymbol{\Sigma}_{\mathrm{Z}}\mathbf{V}_{\mathrm{Z}}^{\mathrm{T}}$$

Let $\boldsymbol{\Sigma}_{\mathrm{C}}\mathbf{V}_{\mathrm{C}}^{\mathrm{T}}\mathbf{X} = \mathbf{Z} \in \mathbb{R}^{c \times n}$

SVD: $\mathbf{C} = \mathbf{U}_{\mathrm{C}}\,\boldsymbol{\Sigma}_{\mathrm{C}}\mathbf{V}_{\mathrm{C}}^{\mathrm{T}} \in \mathbb{R}^{m \times c}$

SVD: $\mathbf{Z} = \mathbf{U}_{\mathrm{Z}}\boldsymbol{\Sigma}_{\mathrm{Z}}\mathbf{V}_{\mathrm{Z}}^{\mathrm{T}} \in \mathbb{R}^{c \times n}$

**Time cost**: $O(mc^2 + nc^2 + nc^2)$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

$m \times s$ matrix with orthonormal columns

$$\mathbf{A} \approx \mathbf{CX} = \mathbf{U}_C \boldsymbol{\Sigma}_C \mathbf{V}_C^T \mathbf{X} = \mathbf{U}_C \mathbf{Z} = \mathbf{U}_C \mathbf{U}_Z \boldsymbol{\Sigma}_Z \mathbf{V}_Z^T$$

$s \times n$ matrix with orthonormal rows

Let $\boldsymbol{\Sigma}_C \mathbf{V}_C^T \mathbf{X} = \mathbf{Z} \in \mathbb{R}^{c \times n}$

diagonal matrix

SVD: $\mathbf{C} = \mathbf{U}_C \boldsymbol{\Sigma}_C \mathbf{V}_C^T \in \mathbb{R}^{m \times c}$
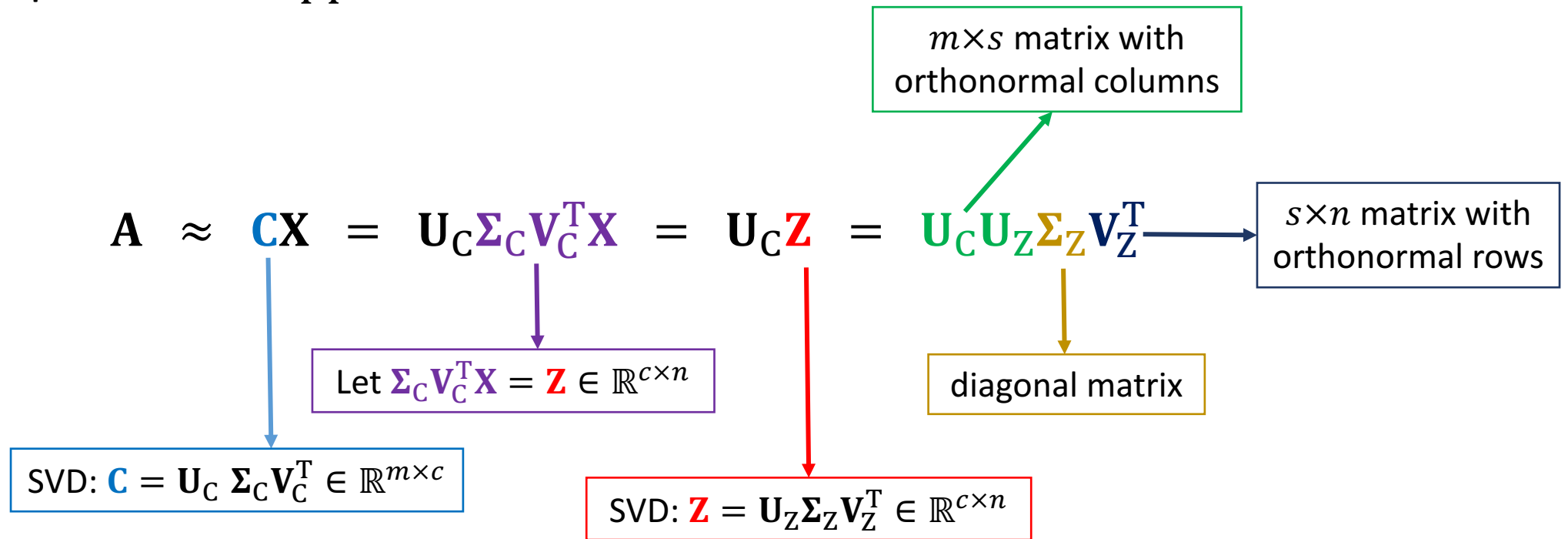
SVD: $\mathbf{Z} = \mathbf{U}_Z \boldsymbol{\Sigma}_Z \mathbf{V}_Z^T \in \mathbb{R}^{c \times n}$

**Time cost**: $O(mc^2 + nc^2 + nc^2 + mc^2)$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD

- Done! Approximate rank $c$ SVD: $\mathbf{A} \approx (\mathbf{U_C U_Z})\mathbf{\Sigma_Z V_Z^T}$

$$\mathbf{A} \approx \mathbf{CX} = \mathbf{U_C \Sigma_C V_C^T X} = \mathbf{U_C Z} = \mathbf{U_C U_Z \Sigma_Z V_Z^T}$$

$m \times s$ matrix with orthonormal columns

$s \times n$ matrix with orthonormal rows

diagonal matrix

**Time cost:** $O(mc^2 + nc^2 + nc^2 + mc^2) = O(mc^2 + nc^2)$

# CX Decomposition ⇔ Approximate SVD

- CX decomposition ⇔ approximate SVD


- Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{C} \in \mathbb{R}^{m \times c}$, the approximate SVD costs
  - $O(mnc)$ time
  - $O(mc + nc)$ memory

# CX Decomposition

- The CX decomposition of $\mathbf{A} \in \mathbb{R}^{m \times n}$
  - Optimal solution: $\mathbf{X}^\star = \text{argmin}_{\mathbf{X}} \left|\left|\mathbf{A} - \mathbf{CX}\right|\right|_{\text{F}}^2 = \mathbf{C}^\dagger \mathbf{A}$
  - How to make it more efficient?

# CX Decomposition

- The CX decomposition of $\mathbf{A} \in \mathbb{R}^{m \times n}$

  - Optimal solution: $\mathbf{X}^\star = \boxed{\operatorname{argmin}_{\mathbf{X}} \left\| \mathbf{A} - \mathbf{CX} \right\|_F^2} = \mathbf{C}^\dagger \mathbf{A}$

  - How to make it more efficient?

A regression problem!

# Fast CX Decomposition

- Fast CX [Drineas, Mahoney, Muthukrishnan, 2008][Clarkson & Woodruff, 2013]
  - Draw another sketching matrix $\boldsymbol{S} \in \mathbb{R}^{m \times s}$
  - Compute $\widetilde{\mathbf{X}} = \mathrm{argmin}_{\mathbf{X}} \left\| \boldsymbol{S}^T (\mathbf{A} - \mathbf{CX}) \right\|_{\mathrm{F}}^2 = \left( \boldsymbol{S}^{\mathrm{T}} \mathbf{C} \right)^{\dagger} \left( \boldsymbol{S}^{\mathrm{T}} \mathbf{A} \right)$
  - Time cost: $O(ncs) + \mathrm{TimeOfSketch}$
  - When $s = \tilde{O}(c/\epsilon)$,
  $$\left\| \mathbf{A} - \mathbf{C}\widetilde{\mathbf{X}} \right\|_{\mathrm{F}}^2 \leq (1 + \epsilon) \cdot \min_{\mathbf{X}} \left\| \mathbf{A} - \mathbf{CX} \right\|_{\mathrm{F}}^2$$

# Outline

- CX Decomposition & Approximate SVD
- CUR Decomposition
- SPSD Matrix Approximation

# CUR Decomposition

- Sketching
    - $\mathbf{C} = \mathbf{A}\mathbf{P_C} \in \mathbb{R}^{m \times c}$
    - $\mathbf{R} = \mathbf{P_R^T}\mathbf{A} \in \mathbb{R}^{r \times n}$
- Find $\mathbf{U}$ such that $\mathbf{CUR} \approx \mathbf{A}$
- CUR $\Leftrightarrow$ Approximate SVD
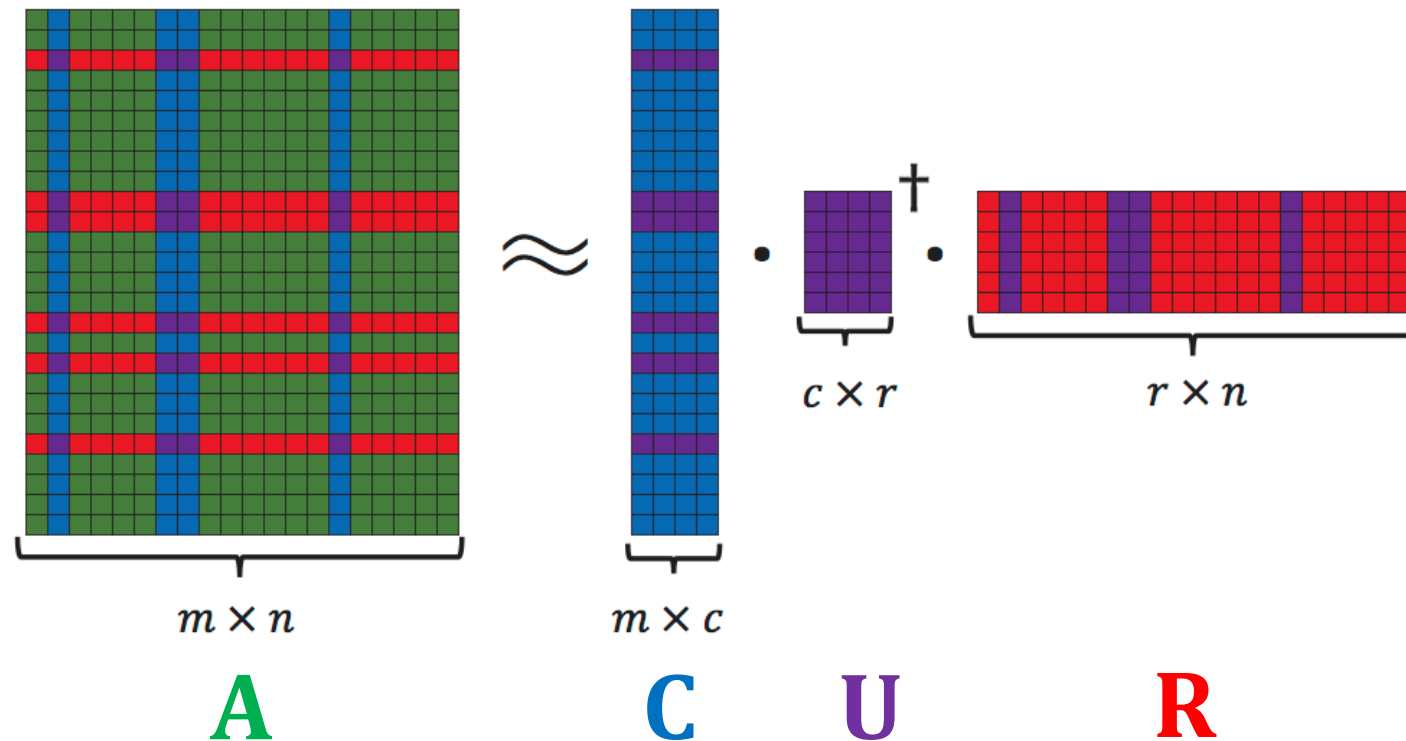    - In the same way as "*CX $\Leftrightarrow$ Approximate SVD*"

# CUR Decomposition

- Sketching
  - $\mathbf{C} = \mathbf{A}\mathbf{P_C} \in \mathbb{R}^{m \times c}$
  - $\mathbf{R} = \mathbf{P_R^T}\mathbf{A} \in \mathbb{R}^{r \times n}$
- Find $\mathbf{U}$ such that $\mathbf{CUR} \approx \mathbf{A}$
- CUR $\Leftrightarrow$ Approximate SVD
  - In the same way as "*CX* $\Leftrightarrow$ *Approximate SVD*"
- 3 types of $\mathbf{U}$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:

$$U = \left(P_R^T A P_C\right)^\dagger$$



$A \qquad C \qquad U \qquad R$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:
$$U = \left( P_R^T A P_C \right)^\dagger$$

- Recall the fast CX decomposition
$$A \approx C\widetilde{X} = C\left( P_R^T C \right)^\dagger \left( P_R^T A \right)$$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:

$$U = \left(P_R^T A P_C\right)^\dagger$$

- Recall the fast CX decomposition

$$A \approx C\widetilde{X} = C\left(P_R^T C\right)^\dagger \left(P_R^T A\right) = CUR$$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:

$$\mathbf{U} = \left(\mathbf{P}_\mathbf{R}^T \mathbf{A} \mathbf{P}_\mathbf{C}\right)^\dagger$$

- Recall the fast CX decomposition

$$\mathbf{A} \approx \mathbf{C}\widetilde{\mathbf{X}} = \mathbf{C}\left(\mathbf{P}_\mathbf{R}^T \mathbf{C}\right)^\dagger \left(\mathbf{P}_\mathbf{R}^T \mathbf{A}\right) = \mathbf{CUR}$$

- They're equivalent: $\mathbf{C}\,\widetilde{\mathbf{X}} = \mathbf{C}\,\mathbf{U}\,\mathbf{R}$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:

$$\mathbf{U} = \left(\mathbf{P}_\mathbf{R}^T \mathbf{A} \mathbf{P}_\mathbf{C}\right)^\dagger$$

- Recall the fast CX decomposition

$$\mathbf{A} \approx \mathbf{C}\widetilde{\mathbf{X}} = \mathbf{C}\left(\mathbf{P}_\mathbf{R}^T \mathbf{C}\right)^\dagger\left(\mathbf{P}_\mathbf{R}^T \mathbf{A}\right) = \mathbf{CUR}$$

- They're equivalent: $\mathbf{C}\,\widetilde{\mathbf{X}} = \mathbf{C}\,\mathbf{U}\,\mathbf{R}$

- Require $c = \tilde{O}\left(\frac{k}{\epsilon}\right)$ and $r = \tilde{O}\left(\frac{c}{\epsilon}\right)$ such that

$$\left|\left|\mathbf{A} - \mathbf{CUR}\right|\right|_\mathbf{F}^2 \leq (1 + \epsilon)\left|\left|\mathbf{A} - \mathbf{A}_k\right|\right|_\mathbf{F}^2$$

# CUR Decomposition

- Type 1 [Drineas, Mahoney, Muthukrishnan, 2008]:

$$\mathbf{U} = \left(\mathbf{P_R^T A P_C}\right)^{\dagger}$$

- Efficient
  - $O(rc^2) + \text{TimeOfSketch}$
- Loose bound
  - Sketch size $\propto \epsilon^{-2}$
- Bad empirical performance

# CUR Decomposition

- Type 2: Optimal CUR

$$\mathbf{U}^{\star} = \min_{\mathbf{U}} \left|\left|\mathbf{A} - \mathbf{CUR}\right|\right|_{F}^{2} = \mathbf{C}^{\dagger}\mathbf{AR}^{\dagger}$$

# CUR Decomposition

- Type 2: Optimal CUR

$$\mathbf{U}^{\star} = \min_{\mathbf{U}} \left\| \mathbf{A} - \mathbf{CUR} \right\|_F^2 = \mathbf{C}^{\dagger} \mathbf{A} \mathbf{R}^{\dagger}$$

- Theory [W & Zhang, 2013], [Boutsidis & Woodruff, 2014]:
  - $\mathbf{C}$ and $\mathbf{R}$ are selected by the adaptive sampling algorithm
  - $c = O\left(\frac{k}{\epsilon}\right)$ and $r = O\left(\frac{k}{\epsilon}\right)$
  - $\left\| \mathbf{A} - \mathbf{CUR} \right\|_F^2 \leq (1 + \epsilon) \left\| \mathbf{A} - \mathbf{A}_k \right\|_F^2$

# CUR Decomposition

- Type 2: Optimal CUR

$$\mathbf{U}^{\star} = \min_{\mathbf{U}} \left\|\mathbf{A} - \mathbf{CUR}\right\|_{F}^{2} = \mathbf{C}^{\dagger}\mathbf{AR}^{\dagger}$$

- Inefficient
  - $\mathrm{O}(mnc) + \mathrm{TimeOfSketch}$

# CUR Decomposition

- Type 3: Fast CUR [W, Zhang, Zhang, 2015]
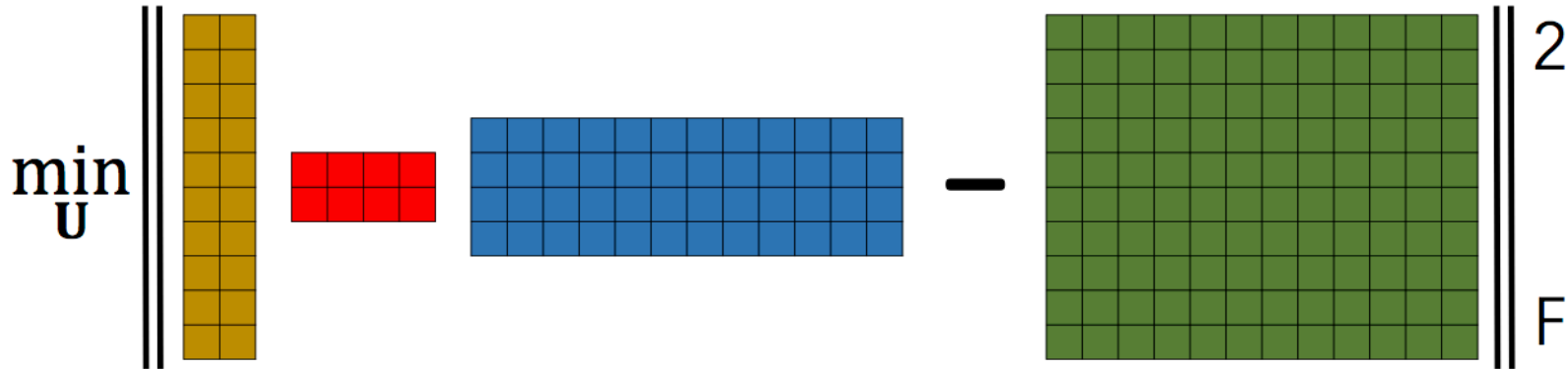  - Draw 2 sketching matrices $\mathbf{S_C}$ and $\mathbf{S_R}$
  - Solve the problem
  $$\widetilde{\mathbf{U}} = \min_{\mathbf{U}} \left\|\mathbf{S}_C^T(\mathbf{A} - \mathbf{CUR})\mathbf{S_R}\right\|_F^2 = \left(\mathbf{S_C^T C}\right)^\dagger \left(\mathbf{S}_C^T \mathbf{A} \mathbf{S_R}\right)\left(\mathbf{R} \mathbf{S_R}\right)^\dagger$$
- Intuition?

# CUR Decomposition
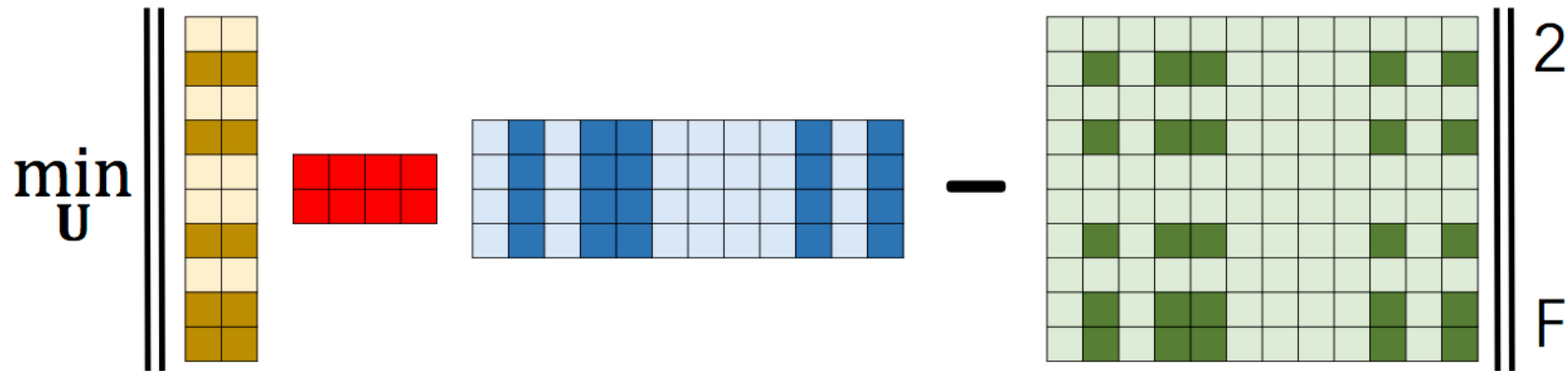
- The optimal $\mathbf{U}$ matrix is obtained by the optimization problem

$$\mathbf{U}^\star = \min_{\mathbf{U}} \left\| \mathbf{C}\mathbf{U}\mathbf{R} - \mathbf{A} \right\|_F^2$$

# CUR Decomposition

- Approximately solve the optimization problem, e.g. by column selection

# CUR Decomposition

- Solve the small scale problem

$$\min_{U} \left\| \quad - \quad \right\|_F^2$$

# CUR Decomposition

- Type 3: Fast CUR [W, Zhang, Zhang, 2015]
  - Draw 2 sketching matrices $\mathbf{S_C} \in \mathbb{R}^{m \times s_c}$ and $\mathbf{S_R} \in \mathbb{R}^{n \times s_r}$
  - Solve the problem

$$\widetilde{\mathbf{U}} = \min_{\mathbf{U}} \left\| \mathbf{S}_C^T (\mathbf{A} - \mathbf{CUR}) \mathbf{S_R} \right\|_F^2 = (\mathbf{S_C^T C})^\dagger (\mathbf{S_C^T A S_R})(\mathbf{R S_R})^\dagger$$

- Theory
  - $s_c = O\left(\frac{c}{\epsilon}\right)$ and $s_r = O\left(\frac{r}{\epsilon}\right)$
  - $\left\| \mathbf{A} - \mathbf{C\widetilde{U}R} \right\|_F^2 \leq (1 + \epsilon) \cdot \min_{\mathbf{U}} \left\| \mathbf{A} - \mathbf{CUR} \right\|_F^2$

# CUR Decomposition

- Type 3: Fast CUR [W, Zhang, Zhang, 2015]
  - Draw 2 sketching matrices $\mathbf{S_C} \in \mathbb{R}^{m \times s_c}$ and $\mathbf{S_R} \in \mathbb{R}^{n \times s_r}$
  - Solve the problem
  $$\widetilde{\mathbf{U}} = \min_{\mathbf{U}} \left\| \mathbf{S}_C^T (\mathbf{A} - \mathbf{CUR}) \mathbf{S_R} \right\|_F^2 = (\mathbf{S_C^T C})^\dagger (\mathbf{S_C^T A S}_R)(\mathbf{R S_R})^\dagger$$
- Efficient
  - $O(s_c s_r (c + r)) + \text{TimeOfSketch}$
- Good empirical performance

**A:**
$m = 1920$
$n = 1168$

**C** and **R:**
- $c = r = 100$
- uniform sampling

**Original**

**Type 2: Optimal CUR**

**Type 1: Fast CX**

**Type 3: Fast CUR**
$s_c = 2c, \qquad s_r = 2r$

**Type 3: Fast CUR**
$s_c = 4c, \qquad s_r = 4r$

# Conclusions

- Approximate truncated SVD
  - CX decomposition
  - CUR decomposition (3 types)
- Fast CUR is the best

# Outline

- CX Decomposition & Approximate SVD

- CUR Decomposition

- SPSD Matrix Approximation

# Motivation 1: Kernel Matrix

- Given $n$ samples $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{R}^d$ and kernel function $\kappa(\cdot,\cdot)$.

- E.g. Gaussian RBF kernel

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\left\|\mathbf{x}_i - \mathbf{x}_j\right\|_2^2}{\sigma^2}\right).$$

# Motivation 1: Kernel Matrix

- Given $n$ samples $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{R}^d$ and kernel function $\kappa(\cdot, \cdot)$.

- E.g. Gaussian RBF kernel

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\left\|\mathbf{x}_i - \mathbf{x}_j\right\|_2^2}{\sigma^2}\right).$$

- Computing the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$
  - where $k_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$
  - costs $O(n^2 d)$ time

# Motivation 2: Matrix Inversion

- Solve the linear system

$$(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$$

  to find $\mathbf{w} \in \mathbb{R}^n$.

- $\mathbf{K} \in \mathbb{R}^{n \times n}$ is the kernel matrix
- $\mathbf{y} = [y_1, \cdots, y_n] \in \mathbb{R}^n$ contains the labels

# Motivation 2: Matrix Inversion

- Solve the linear system
$$(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$$
  to find $\mathbf{w} \in \mathbb{R}^n$.

- Solution: $\mathbf{w}^\star = (\mathbf{K} + \alpha \mathbf{I}_n)^{-1}\mathbf{y}$

# Motivation 2: Matrix Inversion

- Solve the linear system

$$(\mathbf{K} + \alpha\mathbf{I}_n)\mathbf{w} = \mathbf{y}$$

  to find $\mathbf{w} \in \mathbb{R}^n$.

- Solution: $\mathbf{w}^\star = (\mathbf{K} + \alpha\mathbf{I}_n)^{-1}\mathbf{y}$

- It costs
  - $O(n^3)$ time
  - $O(n^2)$ memory.

# Motivation 2: Matrix Inversion

- Solve the linear system

$$(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$$

  to find $\mathbf{w} \in \mathbb{R}^n$.

- Solution: $\mathbf{w}^\star = (\mathbf{K} + \alpha \mathbf{I}_n)^{-1}\mathbf{y}$

- It costs
  - $O(n^3)$ time
  - $O(n^2)$ memory.

- Performed by
  - Kernel ridge regression
  - Least squares kernel SVM

# Motivation 3: Eigenvalue Decomposition

- Find the top $k$ ($\ll n$) eigenvectors of $\mathbf{K}$.

- It costs
  - $\tilde{O}(n^2 k)$ time
  - $O(n^2)$ memory.

# Motivation 3: Eigenvalue Decomposition

- Find the top $k$ ($\ll n$) eigenvectors of $\mathbf{K}$.
- It costs
  - $\tilde{O}(n^2 k)$ time
  - $O(n^2)$ memory.
- Performed by
  - Kernel PCA ($k$ is the target rank)
  - Manifold learning ($k$ is the target rank)

# Computational Challenges

- Time costs
  - Computing kernel matrix: $O(n^2 d)$
  - Matrix inversion: $O(n^3)$
  - Rank $k$ eigenvalue decomposition: $O(n^2 k)$

# Computational Challenges

- Time costs

  - Computing kernel matrix: $O(n^2 d)$
  - Matrix inversion: $O(n^3)$
  - Rank $k$ eigenvalue decomposition: $O(n^2 k)$

**At least quadratic time!**

# Computational Challenges

- Time costs
  - Computing kernel matrix: $O(n^2 d)$
  - Matrix inversion: $O(n^3)$
  - Rank $k$ eigenvalue decomposition: $O(n^2 k)$
- Memory costs
  - Inversion and eigenvalue decomposition: $O(\textcolor{red}{n^2})$

# Computational Challenges

- Time costs
  - Computing kernel matrix: $O(n^2 d)$
  - Matrix inversion: $O(n^3)$
  - Rank $k$ eigenvalue decomposition: $O(n^2 k)$
- Memory costs
  - Inversion and eigenvalue decomposition: $O(n^2)$
  - Because
    - the numerical algorithms are pass-inefficient
    - ➔ form **K** and keep it in memory

# Computational Challenges

- Time costs
  - Computing kernel matrix: $O(n^2 d)$
  - Matrix inversion: $O(n^3)$
  - Rank $k$ eigenvalue decomposition: $O(n^2 k)$

- Memory costs
  - Inversion and eigenvalue decomposition: $O(n^2)$
  - Because
    - the numerical algorithms are pass-inefficient
    - ➔ form **K** and keep it in memory

When $n = 10^5$, the $n \times n$ matrix costs 80GB memory!

# How to Speedup?

- Efficiently form the low-rank approximation

$$\mathbf{K} \approx \mathbf{C}\,\mathbf{U}\,\mathbf{C}^{\mathrm{T}}$$



$$\mathbf{K} \approx \mathbf{C} \cdot \mathbf{U} \cdot \mathbf{C}^{\mathrm{T}}$$

$n \times n$     $n \times c$     $c \times c$     $c \times n$

$\mathbf{K}$     $\mathbf{C}$     $\mathbf{U}$     $\mathbf{C}^{\mathrm{T}}$

# How to Speedup?

- Efficiently form the low-rank approximation
$$\mathbf{K} \approx \mathbf{C} \, \mathbf{U} \, \mathbf{C}^{\mathrm{T}}$$

- Equivalent $\mathbf{K} \approx \mathbf{L} \, \mathbf{L}^{\mathrm{T}}$

# Efficient Matrix Inversion

- Solve the linear system $(\mathbf{K} + \alpha\mathbf{I}_n)\mathbf{w} = \mathbf{y}$:

$$\mathbf{w}^\star = (\mathbf{K} + \alpha\mathbf{I}_n)^{-1}\mathbf{y}$$

# Efficient Matrix Inversion

- Approximately solve the linear system $(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$:

  - Replace $\mathbf{K}$ by $\mathbf{LL}^{\mathrm{T}}$:   $\mathbf{w}^{\star} = (\mathbf{K} + \alpha \mathbf{I}_n)^{-1}\mathbf{y} \approx \left(\mathbf{LL}^{\mathrm{T}} + \alpha \mathbf{I}_n\right)^{-1}\mathbf{y}$

# Efficient Matrix Inversion

- Approximately solve the linear system $(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$

  - Replace $\mathbf{K}$ by $\mathbf{L}\mathbf{L}^{\mathrm{T}}$:   $\mathbf{w}^{\star} = (\mathbf{K} + \alpha \mathbf{I}_n)^{-1}\mathbf{y} \approx \left(\mathbf{L}\mathbf{L}^{\mathrm{T}} + \alpha \mathbf{I}_n\right)^{-1}\mathbf{y}$

  - Expand the inversion by the Woodbury identity

# Efficient Matrix Inversion

- Approximately solve the linear system $(\mathbf{K} + \alpha \mathbf{I}_n)\mathbf{w} = \mathbf{y}$

  - Replace $\mathbf{K}$ by $\mathbf{L}\mathbf{L}^{\mathrm{T}}$:   $\mathbf{w}^{\star} = (\mathbf{K} + \alpha \mathbf{I}_n)^{-1}\mathbf{y} \approx \left(\mathbf{L}\mathbf{L}^{\mathrm{T}} + \alpha \mathbf{I}_n\right)^{-1}\mathbf{y}$

  - Expand the inversion by the Woodbury identity

$$(\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1}$$

# Efficient Matrix Inversion

- Approximately solve the linear system $(\mathbf{K} + \alpha\mathbf{I}_n)\mathbf{w} = \mathbf{y}$

  - Replace $\mathbf{K}$ by $\mathbf{L}\mathbf{L}^{\mathrm{T}}$:   $\mathbf{w}^\star = (\mathbf{K} + \alpha\mathbf{I}_n)^{-1}\mathbf{y} \approx \left(\mathbf{L}\mathbf{L}^{\mathrm{T}} + \alpha\mathbf{I}_n\right)^{-1}\mathbf{y}$

  - Expand the inversion by the Woodbury identity

  $$\mathbf{w}^\star \approx \alpha^{-1}\mathbf{y} + \alpha^{-1}\mathbf{L}(\alpha\mathbf{I} + \mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^{\mathbf{T}}\mathbf{y}$$

# Efficient Matrix Inversion

- Approximately solve the linear system $(\mathbf{K} + \alpha\mathbf{I}_n)\mathbf{w} = \mathbf{y}$

  - Replace $\mathbf{K}$ by $\mathbf{L}\mathbf{L}^{\mathrm{T}}$:   $\mathbf{w}^\star = (\mathbf{K} + \alpha\mathbf{I}_n)^{-1}\mathbf{y} \approx \left(\mathbf{L}\mathbf{L}^{\mathrm{T}} + \alpha\mathbf{I}_n\right)^{-1}\mathbf{y}$

  - Expand the inversion by the Woodbury identity
  $$\mathbf{w}^\star \approx \alpha^{-1}\mathbf{y} + \alpha^{-1}\mathbf{L}(\alpha\mathbf{I} + \mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^{\mathrm{T}}\mathbf{y}$$


  - Time cost: $O(nc^2)$

Linear in $n$, much better than $O(n^3)$

# Efficient Eigenvalue Decomposition

- Approximately compute the $k$-eigenvalue decomposition of $\mathbf{K}$

  - SVD: $\mathbf{L} = \mathbf{U_L \Sigma_L V_L}$

  - $\mathbf{K} \approx \mathbf{LL^T} = \mathbf{U_L \Sigma_L^2 U_L^T}$

# Efficient Eigenvalue Decomposition

- Approximately compute the $k$-eigenvalue decomposition of $\mathbf{K}$

  - SVD: $\mathbf{L} = \mathbf{U_L \Sigma_L V_L}$

  - $\mathbf{K} \approx \mathbf{LL^T} = \mathbf{U_L \Sigma_L^2 U_L^T}$

  - Approximate $k$- eigenvalue decomposition of $\mathbf{K}$

    - eigenvectors: the first $k$ vectors in $\mathbf{U_L}$

  - Time cost: $O(nc^2)$

# Efficient Eigenvalue Decomposition

- Approximately compute the $k$-eigenvalue decomposition of $\mathbf{K}$

  - SVD: $\mathbf{L} = \mathbf{U_L}\mathbf{\Sigma_L}\mathbf{V_L}$

  - $\mathbf{K} \approx \mathbf{L}\mathbf{L^T} = \mathbf{U_L}\mathbf{\Sigma_L^2}\mathbf{U_L^T}$

  - Approximate $k$- eigenvalue decomposition of $\mathbf{K}$

    - eigenvectors: the first $k$ vectors in $\mathbf{U_L}$

  - Time cost: $O(nc^2)$
    - Much lower than $\tilde{O}(n^2 k)$

# Sketching Based Models

- How to find such an approximation?

$$\mathbf{K} \approx \mathbf{C}\,\mathbf{U}\,\mathbf{C}^{\mathrm{T}}$$

# Sketching Based Models

- How to find such an approximation?

$$\mathbf{K} \approx \mathbf{C}\,\mathbf{U}\,\mathbf{C}^{\mathrm{T}}$$

- Sketching based Methods: $\mathbf{C} = \mathbf{K}\mathbf{S} \in \mathbb{R}^{n \times c}$ is a sketch of $\mathbf{K}$.
  - $\mathbf{S} \in \mathbb{R}^{n \times c}$ can be column selection or random projection matrix

# Sketching Based Models

- How to find such an approximation?

$$\mathbf{K} \approx \mathbf{C}\,\mathbf{U}\,\mathbf{C}^{\mathrm{T}}$$

- Sketching based Methods: $\mathbf{C} = \mathbf{K}\boldsymbol{S} \in \mathbb{R}^{n \times c}$ is a sketch of $\mathbf{K}$.
  - $\boldsymbol{S} \in \mathbb{R}^{n \times c}$ can be column selection or random projection matrix

- Three methods:
  - The prototype model [HMT11, WZ13, WLZ16]
  - The fast model [WZZ15]
  - The Nyström method [WS15, GM13]

# The Prototype Model

- Objective: $\mathbf{K} \approx \mathbf{CUC}^{\mathrm{T}}$

- Minimize the approximation error by
$$\mathbf{U}^{\star} = \underset{\mathbf{U}}{\mathrm{argmin}} \left\lVert \mathbf{K} - \mathbf{CUC}^{\mathrm{T}} \right\rVert_{F}^{2} = \mathbf{C}^{\dagger} \mathbf{K} \left( \mathbf{C}^{\dagger} \right)^{\mathrm{T}}.$$

# The Prototype Model

- Objective: $\mathbf{K} \approx \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}}$

- Minimize the approximation error by

$$\mathbf{U}^{\star} = \underset{\mathbf{U}}{\operatorname{argmin}} \left\| \mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}} \right\|_{F}^{2} = \mathbf{C}^{\dagger}\mathbf{K}\left(\mathbf{C}^{\dagger}\right)^{\mathrm{T}}.$$

Extension of the random SVD to SPSD matrix [HMT11]

# The Prototype Model

- Objective: $\mathbf{K} \approx \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}}$

- Minimize the approximation error by

$$\mathbf{U}^{\star} = \operatorname*{argmin}_{\mathbf{U}} \left\|\left| \mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}} \right\|\right|_{F}^{2} = \mathbf{C}^{\dagger}\mathbf{K}(\mathbf{C}^{\dagger})^{\mathrm{T}}.$$

- Time: $O(n^2 c)$

---

- The time complexity is nearly the same to the $k$-eigenvalue decomposition.
- It is much faster than the $k$- eigenvalue decomposition in practice.

# The Prototype Model

- Objective: $\mathbf{K} \approx \mathbf{CUC}^{\mathrm{T}}$

- Minimize the approximation error by
$$\mathbf{U}^\star = \underset{\mathbf{U}}{\operatorname{argmin}} \left\| \mathbf{K} - \mathbf{CUC}^{\mathrm{T}} \right\|_F^2 = \mathbf{C}^\dagger \mathbf{K} (\mathbf{C}^\dagger)^{\mathrm{T}}.$$

- Time: $O(n^2 c)$

- #Passes: one

# The Prototype Model

- Objective: $\mathbf{K} \approx \mathbf{CUC}^{\mathrm{T}}$

- Minimize the approximation error by
$$\mathbf{U}^{\star} = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{K} - \mathbf{CUC}^{\mathrm{T}} \right\|_F^2 = \mathbf{C}^{\dagger}\mathbf{K}(\mathbf{C}^{\dagger})^{\mathrm{T}}.$$

- Time: $O(n^2 c)$

- #Passes: one

- Memory: $O(nc)$
  - Put $k_{ij}$ in memory only when it is visited
  - Keep $\mathbf{C}^{\dagger}$ in memory

# The Prototype Model

- Error Bound
  - $k \ll n$ is arbitrary integer
  - $\mathbf{P}$ samples $c = O\left(\dfrac{k}{\epsilon}\right)$ columns by adaptive sampling
  - $\mathbb{E}\left|\left|\mathbf{K} - \mathbf{C}\mathbf{U}^{\star}\mathbf{C}^{\mathrm{T}}\right|\right|_F^2 \leq (1 + \epsilon)\left|\left|\mathbf{K} - \mathbf{K}_k\right|\right|_F^2$

# The Prototype Model

- Limitations
  - $\mathbf{U}^{\star} = \mathbf{C}^{\dagger}\mathbf{K}(\mathbf{C}^{\dagger})^{\mathrm{T}}$
  - Time cost is $O(n^2 c)$
  - Requires observing the whole of $\mathbf{K}$

# The Prototype Model

- Prototype model: $\mathbf{K} \approx \mathbf{C}\,\mathbf{U}^\star\,\mathbf{C}^\mathrm{T}$, where

$$\mathbf{U}^\star = \operatorname*{argmin}_{\mathbf{U}} \left\|\mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^\mathrm{T}\right\|_\mathrm{F}^2.$$



$$\min_{\mathbf{U}} \left\| \underset{\substack{n \times n \\ \mathbf{K}}}{\Box} - \underset{\substack{n \times c \\ \mathbf{C}}}{\Box} \cdot \underset{\substack{c \times c \\ \mathbf{U}}}{\Box} \cdot \underset{\substack{c \times n \\ \mathbf{C}^\mathrm{T}}}{\Box} \right\|_\mathrm{F}^2$$

# The Fast Model

- Column/row selection
  - Form $\mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P}$ and $\mathbf{P}^{\mathrm{T}}\mathbf{C}$



$$\min_{\mathbf{U}} \left\| \underbrace{\phantom{\mathbf{P}^{\mathrm{T}}\mathbf{K}}}_{\substack{n \times n \\ \mathbf{P}^{\mathrm{T}}\mathbf{K}}} - \underbrace{\phantom{\mathbf{P}^{\mathrm{T}}\mathbf{C}}}_{\substack{n \times c \\ \mathbf{P}^{\mathrm{T}}\mathbf{C}}} \cdot \underbrace{\phantom{\mathbf{U}}}_{\substack{c \times c \\ \mathbf{U}}} \cdot \underbrace{\phantom{\mathbf{C}^{\mathrm{T}}}}_{\substack{c \times n \\ \mathbf{C}^{\mathrm{T}}}} \right\|_{F}^{2}$$

# The Fast Model

- Column/row selection
  - Form $\mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P}$ and $\mathbf{P}^{\mathrm{T}}\mathbf{C}$

$$\min_{\mathbf{U}} \left\| \underbrace{\phantom{XXXXXX}}_{\substack{n \times n \\ \mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P}}} - \underbrace{\phantom{XX}}_{\substack{n \times c \\ \mathbf{P}^{\mathrm{T}}\mathbf{C}}} \cdot \underbrace{\phantom{XX}}_{\substack{c \times c \\ \mathbf{U}}} \cdot \underbrace{\phantom{XXXXX}}_{\substack{c \times n \\ \mathbf{C}^{\mathrm{T}}\mathbf{P}}} \right\|_{\mathrm{F}}^{2}$$

# The Fast Model

- $\mathbf{K} \approx \mathbf{C}\, \widetilde{\mathbf{U}}\, \mathbf{C}^{\mathrm{T}}$, where

$$\widetilde{\mathbf{U}} = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{P}^{\mathrm{T}}\left(\mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}}\right)\mathbf{P} \right\|_{F}^{2}.$$

$$\underset{\mathbf{U}}{\min} \left\| \quad \underbrace{\phantom{XXXX}}_{p \times p} \quad - \quad \underbrace{\phantom{X}}_{p \times c} \cdot \underbrace{\phantom{XX}}_{c \times c} \cdot \underbrace{\phantom{XXXX}}_{c \times p} \quad \right\|_{F}^{2}$$

$$\mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P} \qquad\qquad \mathbf{P}^{\mathrm{T}}\mathbf{C} \quad \mathbf{U} \qquad\qquad \mathbf{C}^{\mathrm{T}}\mathbf{P}$$

# The Fast Model

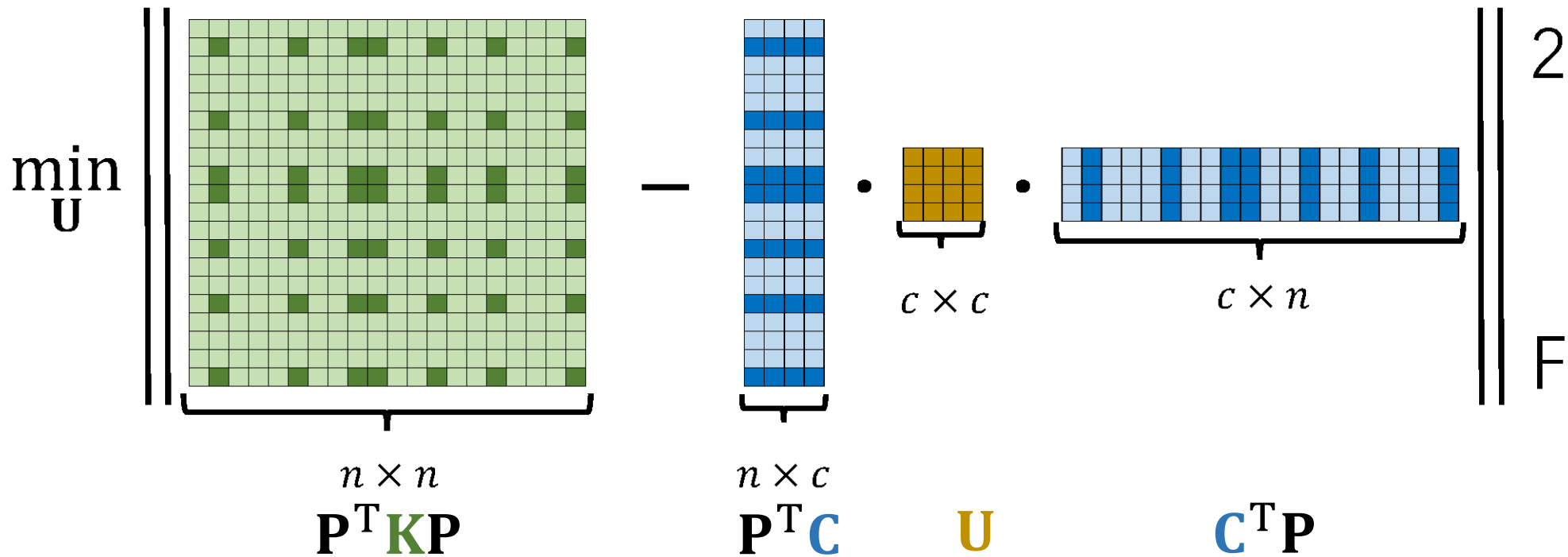- Prototype model: $\mathbf{U}^\star = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^T \right\|_F^2 = \mathbf{C}^\dagger \mathbf{K} (\mathbf{C}^\dagger)^T$

- Fast model: $\widetilde{\mathbf{U}} = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{P}^T (\mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^T) \mathbf{P} \right\|_F^2 = (\mathbf{P}^T \mathbf{C})^\dagger (\mathbf{P}^T \mathbf{K} \mathbf{P}) (\mathbf{C}^T \mathbf{P})^\dagger.$

# The Fast Model

- Prototype model: $\mathbf{U}^{\star} = \underset{\mathbf{U}}{\text{argmin}} \left\| \mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}} \right\|_{\mathrm{F}}^{2} = \mathbf{C}^{\dagger}\mathbf{K}(\mathbf{C}^{\dagger})^{\mathrm{T}}$

- Fast model: $\widetilde{\mathbf{U}} = \underset{\mathbf{U}}{\text{argmin}} \left\| \mathbf{P}^{\mathrm{T}}(\mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}})\mathbf{P} \right\|_{F}^{2} = (\mathbf{P}^{\mathrm{T}}\mathbf{C})^{\dagger}(\mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P})(\mathbf{C}^{\mathrm{T}}\mathbf{P})^{\dagger}.$

- Theory
  - $p = O\left(\sqrt{\dfrac{nc}{\epsilon}}\right)$
  - $\mathbf{P}$ is column selection matrix (according to the row leverage scores of $\mathbf{C}$)
  - Then $\left\| \mathbf{K} - \mathbf{C}\widetilde{\mathbf{U}}\mathbf{C}^{\mathrm{T}} \right\|_{F}^{2} \leq (1 + \epsilon) \left\| \mathbf{K} - \mathbf{C}\mathbf{U}^{\star}\mathbf{C}^{\mathrm{T}} \right\|_{F}^{2}$
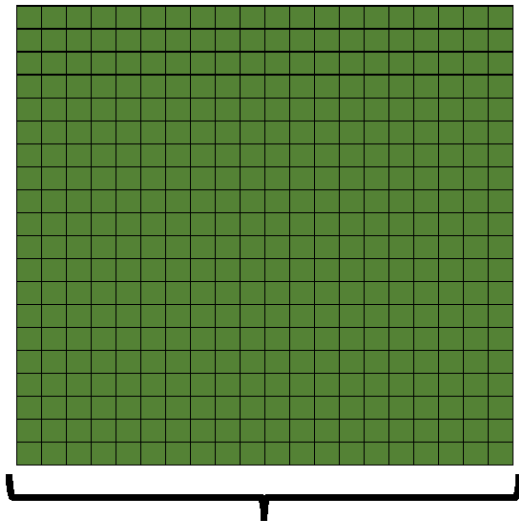
The faster model is nearly as good as the prototype model!

# The Fast Model

- Prototype model: $\mathbf{U}^\star = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^T \right\|_F^2 = \mathbf{C}^\dagger \mathbf{K} (\mathbf{C}^\dagger)^T$

- Fast model: $\widetilde{\mathbf{U}} = \underset{\mathbf{U}}{\mathrm{argmin}} \left\| \mathbf{P}^T(\mathbf{K} - \mathbf{C}\mathbf{U}\mathbf{C}^T)\mathbf{P} \right\|_F^2 = (\mathbf{P}^T\mathbf{C})^\dagger (\mathbf{P}^T\mathbf{K}\mathbf{P})(\mathbf{C}^T\mathbf{P})^\dagger.$

- Theory

  - $p = O\left(\sqrt{\dfrac{nc}{\epsilon}}\right)$

  - $\mathbf{P}$ is column selection matrix (according to the row leverage scores of $\mathbf{C}$)

  - Then $\left\| \mathbf{K} - \mathbf{C}\widetilde{\mathbf{U}}\mathbf{C}^T \right\|_F^2 \leq (1 + \epsilon) \left\| \mathbf{K} - \mathbf{C}\mathbf{U}^\star\mathbf{C}^T \right\|_F^2$

- Overall time cost: $O(p^2 c + nc^2) = O(nc^3/\epsilon)$

$$\boxed{\text{linear in } n}$$
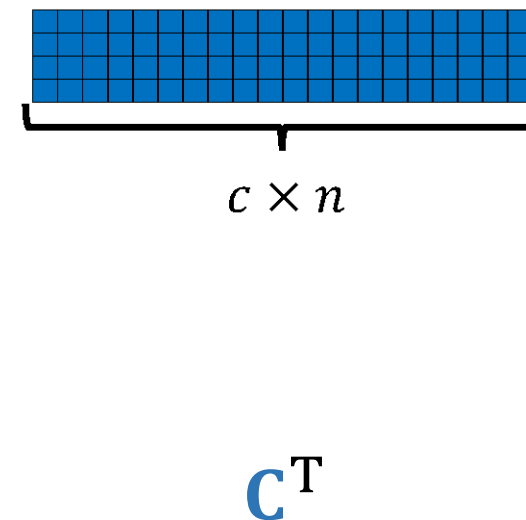
# The Nyström Method



$n \times n$

**K**

# The Nyström Method

- **S** ($n{\times}c$): column selection matrix
- **C** = **K****S** ($n{\times}c$)



$n \times n$

**K**

$n \times c$

**C**

$c \times n$

**C**$^\mathrm{T}$

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix
- $\mathbf{C} = \mathbf{K}\mathbf{S}$ ($n \times c$), $\mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{K}\mathbf{S} = \mathbf{S}^{\mathrm{T}}\mathbf{C}$ ($c \times c$)



$n \times n$

$\mathbf{K}$

$n \times c$

$\mathbf{C}$
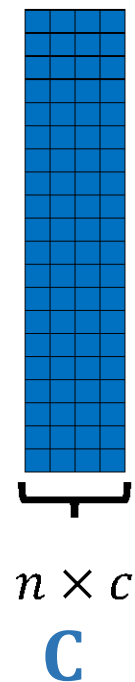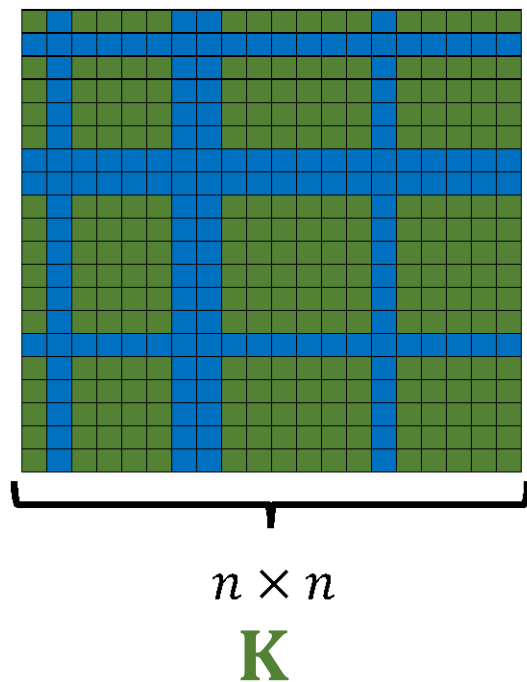
$c \times c$

$\mathbf{W}$

$c \times n$

$\mathbf{C}^{\mathrm{T}}$

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix
- $\mathbf{C} = \mathbf{K}\mathbf{S}$ ($n \times c$), $\mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{K}\mathbf{S} = \mathbf{S}^{\mathrm{T}}\mathbf{C}$ ($c \times c$)
- The Nyström method:  $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^{\dagger}\,\mathbf{C}^{\mathrm{T}}$



$$n \times n \qquad n \times c \qquad c \times c \qquad c \times n$$

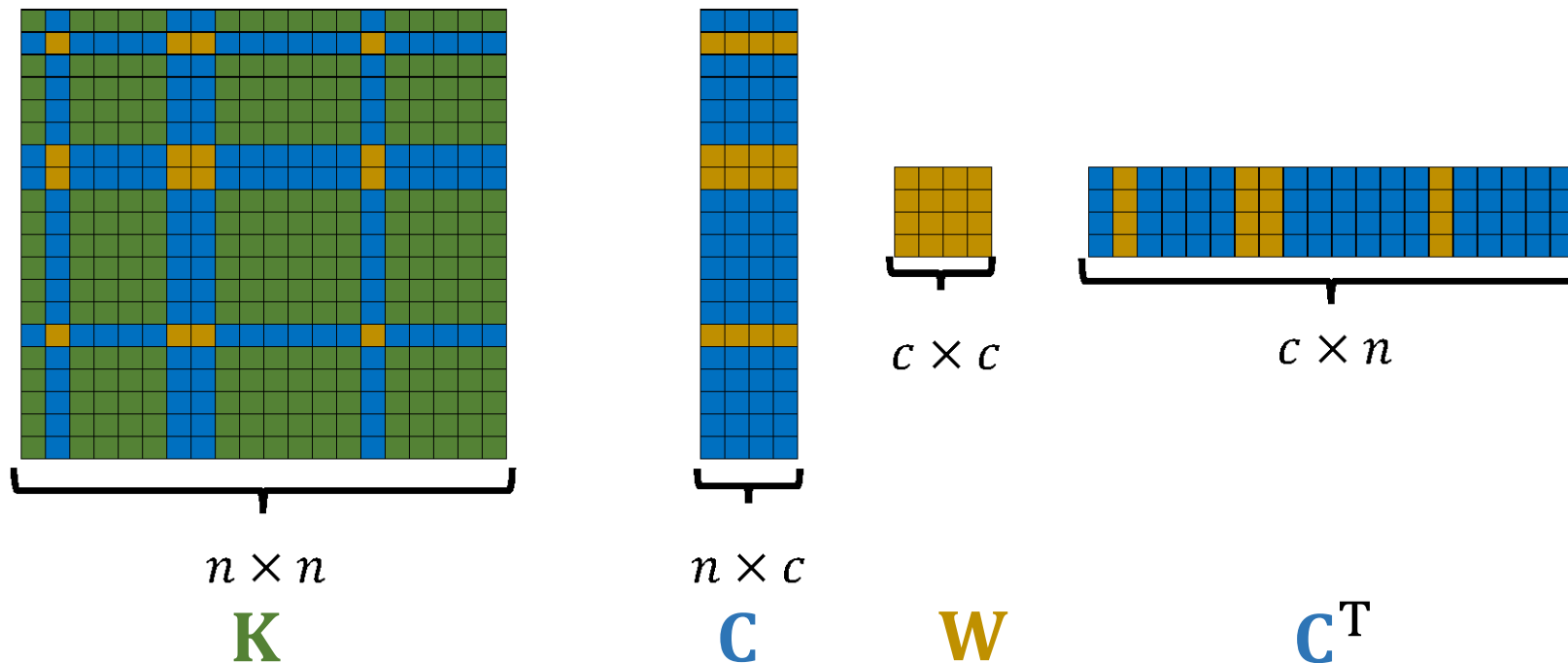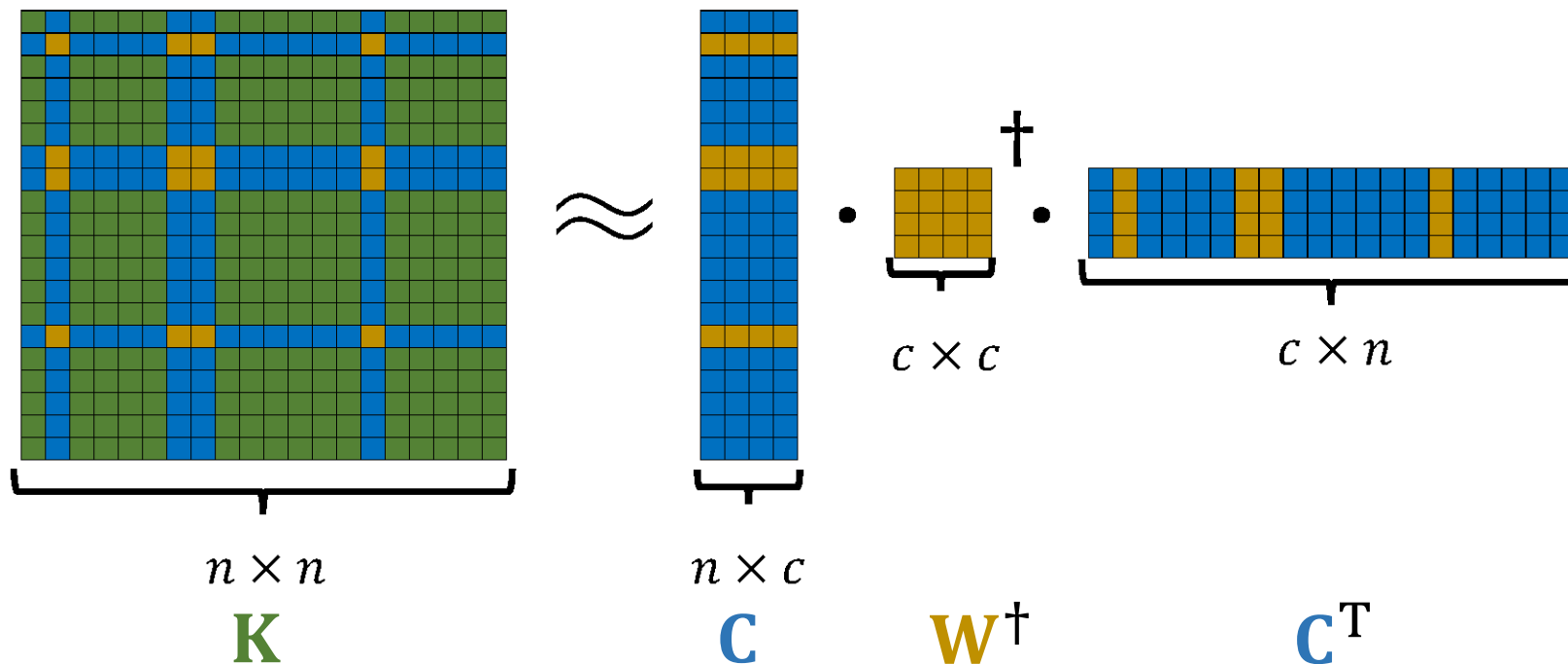$$\mathbf{K} \qquad \mathbf{C} \qquad \mathbf{W}^{\dagger} \qquad \mathbf{C}^{\mathrm{T}}$$

# The Nyström Method

- $\mathbf{S}$ ($n{\times}c$): column selection matrix

- $\mathbf{C} = \mathbf{KS}$ ($n{\times}c$), $\mathbf{W} = \mathbf{S}^\mathrm{T}\mathbf{KS} = \mathbf{S}^\mathrm{T}\mathbf{C}$ ($c{\times}c$)

- The Nyström method: $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^\dagger\,\mathbf{C}^\mathrm{T}$

- New explanation:

  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left|\left|\color{red}{\mathbf{P}^\mathrm{T}}\color{black}{(\mathbf{K} - \mathbf{CXC}^\mathrm{T})}\color{red}{\mathbf{P}}\right|\right|_\mathrm{F}^2$

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix

- $\mathbf{C} = \mathbf{K}\mathbf{S}$ ($n \times c$), $\mathbf{W} = \mathbf{S}^\mathrm{T}\mathbf{K}\mathbf{S} = \mathbf{S}^\mathrm{T}\mathbf{C}$ ($c \times c$)

- The Nyström method: $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^\dagger\,\mathbf{C}^\mathrm{T}$

- New explanation:

  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left|\left| \mathbf{P}^\mathrm{T}(\mathbf{K} - \mathbf{C}\mathbf{X}\mathbf{C}^\mathrm{T})\mathbf{P} \right|\right|_\mathrm{F}^2$

  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left|\left| \mathbf{S}^\mathrm{T}(\mathbf{K} - \mathbf{C}\mathbf{X}\mathbf{C}^\mathrm{T})\mathbf{S} \right|\right|_\mathrm{F}^2$$

# The Nyström Method

- $\mathbf{S}$ ($n{\times}c$): column selection matrix
- $\mathbf{C} = \mathbf{KS}$ ($n{\times}c$), $\mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{KS} = \mathbf{S}^{\mathrm{T}}\mathbf{C}$ ($c{\times}c$)
- The Nyström method: $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^{\dagger}\,\mathbf{C}^{\mathrm{T}}$
- New explanation:
  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}}\left\Vert\mathbf{P}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{P}\right\Vert_{\mathrm{F}}^{2}$
  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}}\left\Vert\mathbf{S}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{S}\right\Vert_{\mathrm{F}}^{2}$$
$$= (\mathbf{S}^{\mathrm{T}}\mathbf{C})^{\dagger}(\mathbf{S}^{\mathrm{T}}\mathbf{KS})(\mathbf{C}^{\mathrm{T}}\mathbf{S})^{\dagger}$$

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix
- $\mathbf{C} = \mathbf{KS}$ ($n \times c$), $\textcolor{red}{\mathbf{W} = \mathbf{S}^\mathrm{T}\mathbf{KS} = \mathbf{S}^\mathrm{T}\mathbf{C}}$ ($c \times c$)
- The Nyström method: $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^\dagger\,\mathbf{C}^\mathrm{T}$
- New explanation:
  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left\| \mathbf{P}^\mathrm{T}(\mathbf{K} - \mathbf{CXC}^\mathrm{T})\mathbf{P} \right\|_\mathrm{F}^2$
  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left\| \mathbf{S}^\mathrm{T}(\mathbf{K} - \mathbf{CXC}^\mathrm{T})\mathbf{S} \right\|_\mathrm{F}^2$$

$$= \left(\textcolor{red}{\mathbf{S}^\mathrm{T}\mathbf{C}}\right)^\dagger \left(\textcolor{red}{\mathbf{S}^\mathrm{T}\mathbf{KS}}\right)\left(\textcolor{red}{\mathbf{C}^\mathrm{T}\mathbf{S}}\right)^\dagger$$

$$= \textcolor{red}{\mathbf{W}^\dagger \mathbf{W} \mathbf{W}^\dagger}$$

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix
- $\mathbf{C} = \mathbf{KS}$ ($n \times c$), $\color{red}{\mathbf{W} = \mathbf{S}^T\mathbf{KS} = \mathbf{S}^T\mathbf{C}}$ ($c \times c$)
- The Nyström method: $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^\dagger\,\mathbf{C}^T$
- New explanation:
  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\operatorname{argmin}} \left\|\mathbf{P}^T(\mathbf{K} - \mathbf{CXC}^T)\mathbf{P}\right\|_F^2$
  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\operatorname{argmin}} \left\|\mathbf{S}^T(\mathbf{K} - \mathbf{CXC}^T)\mathbf{S}\right\|_F^2$$

$$= \color{red}{(\mathbf{S}^T\mathbf{C})^\dagger(\mathbf{S}^T\mathbf{KS})(\mathbf{C}^T\mathbf{S})^\dagger}$$

$$= \color{red}{\mathbf{W}^\dagger\mathbf{W}\mathbf{W}^\dagger = \mathbf{W}^\dagger}$$

# The Nyström Method

- $\mathbf{S}$ ($n{\times}c$): column selection matrix

- $\mathbf{C} = \mathbf{KS}$ ($n{\times}c$), $\mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{KS} = \mathbf{S}^{\mathrm{T}}\mathbf{C}$ ($c{\times}c$)

- The Nyström method:  $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^{\dagger}\,\mathbf{C}^{\mathrm{T}}$

- New explanation:

  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}}\left|\left|\mathbf{P}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{P}\right|\right|_{\mathrm{F}}^{2}$

  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}}\left|\left|\mathbf{S}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{S}\right|\right|_{\mathrm{F}}^{2}$$

$$= (\mathbf{S}^{\mathrm{T}}\mathbf{C})^{\dagger}(\mathbf{S}^{\mathrm{T}}\mathbf{KS})(\mathbf{C}^{\mathrm{T}}\mathbf{S})^{\dagger}$$

$$= \mathbf{W}^{\dagger}\mathbf{WW}^{\dagger} = \mathbf{W}^{\dagger}$$

- The Nystrom method is special instance of the fast model.

# The Nyström Method

- $\mathbf{S}$ ($n \times c$): column selection matrix
- $\mathbf{C} = \mathbf{KS}$ ($n \times c$), $\mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{KS} = \mathbf{S}^{\mathrm{T}}\mathbf{C}$ ($c \times c$)
- The Nyström method:  $\mathbf{K} \approx \mathbf{C}\,\mathbf{W}^{\dagger}\,\mathbf{C}^{\mathrm{T}}$
- New explanation:
  - Recall the fast model: $\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left\| \mathbf{P}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{P} \right\|_{\mathrm{F}}^{2}$
  - Setting $\mathbf{P} = \mathbf{S}$, then

$$\widetilde{\mathbf{X}} = \underset{\mathbf{X}}{\mathrm{argmin}} \left\| \mathbf{S}^{\mathrm{T}}(\mathbf{K} - \mathbf{CXC}^{\mathrm{T}})\mathbf{S} \right\|_{\mathrm{F}}^{2}$$
$$= (\mathbf{S}^{\mathrm{T}}\mathbf{C})^{\dagger}(\mathbf{S}^{\mathrm{T}}\mathbf{KS})(\mathbf{C}^{\mathrm{T}}\mathbf{S})^{\dagger}$$
$$= \mathbf{W}^{\dagger}\mathbf{W}\mathbf{W}^{\dagger} = \mathbf{W}^{\dagger}$$

- The Nystrom method is special instance of the fast model.
- It is approximate solution to the prototype model

# The Nyström Method

- Cost
  - Time: $O(nc^2)$
  - Memory: $O(nc)$

# The Nyström Method

- Cost
  - Time: $O(nc^2)$
  - Memory: $O(nc)$

  Very efficient!

# The Nyström Method

- Cost
  - Time: $O(nc^2)$
  - Memory: $O(nc)$

  Very efficient!

- Error bound: weak

# Comparisons

- $\mathbf{C} = \mathbf{KS} \in \mathbb{R}^{n \times c}, \mathbf{W} = \mathbf{S}^{\mathrm{T}} \mathbf{KS} = \mathbf{S}^{\mathrm{T}} \mathbf{C} \in \mathbb{R}^{c \times c}$
- SPSD matrix approximation: $\mathbf{K} \approx \mathbf{CUC}^{\mathrm{T}}$
  - The prototype model: $\mathbf{U} = \mathbf{C}^{\dagger} \mathbf{K} \left(\mathbf{C}^{\dagger}\right)^{\mathrm{T}}$
  - The fast model: $\mathbf{U} = \left(\mathbf{P}^{\mathrm{T}} \mathbf{C}\right)^{\dagger} \left(\mathbf{P}^{\mathrm{T}} \mathbf{KP}\right) \left(\mathbf{C}^{\mathrm{T}} \mathbf{P}\right)^{\dagger}$
  - The Nyström method: $\mathbf{U} = \mathbf{W}^{\dagger}$

# Comparisons

- $\mathbf{C} = \mathbf{KS} \in \mathbb{R}^{n \times c}, \mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{KS} = \mathbf{S}^{\mathrm{T}}\mathbf{C} \in \mathbb{R}^{c \times c}$

- SPSD matrix approximation: $\mathbf{K} \approx \mathbf{CUC}^{\mathrm{T}}$

  - The prototype model: $\mathbf{U} = \mathbf{C}^{\dagger}\mathbf{K}\left(\mathbf{C}^{\dagger}\right)^{\mathrm{T}}$

  - The fast model: $\mathbf{U} = \left(\mathbf{P}^{\mathrm{T}}\mathbf{C}\right)^{\dagger}\left(\mathbf{P}^{\mathrm{T}}\mathbf{KP}\right)\left(\mathbf{C}^{\mathrm{T}}\mathbf{P}\right)^{\dagger}$

  - The Nyström method: $\mathbf{U} = \mathbf{W}^{\dagger}$

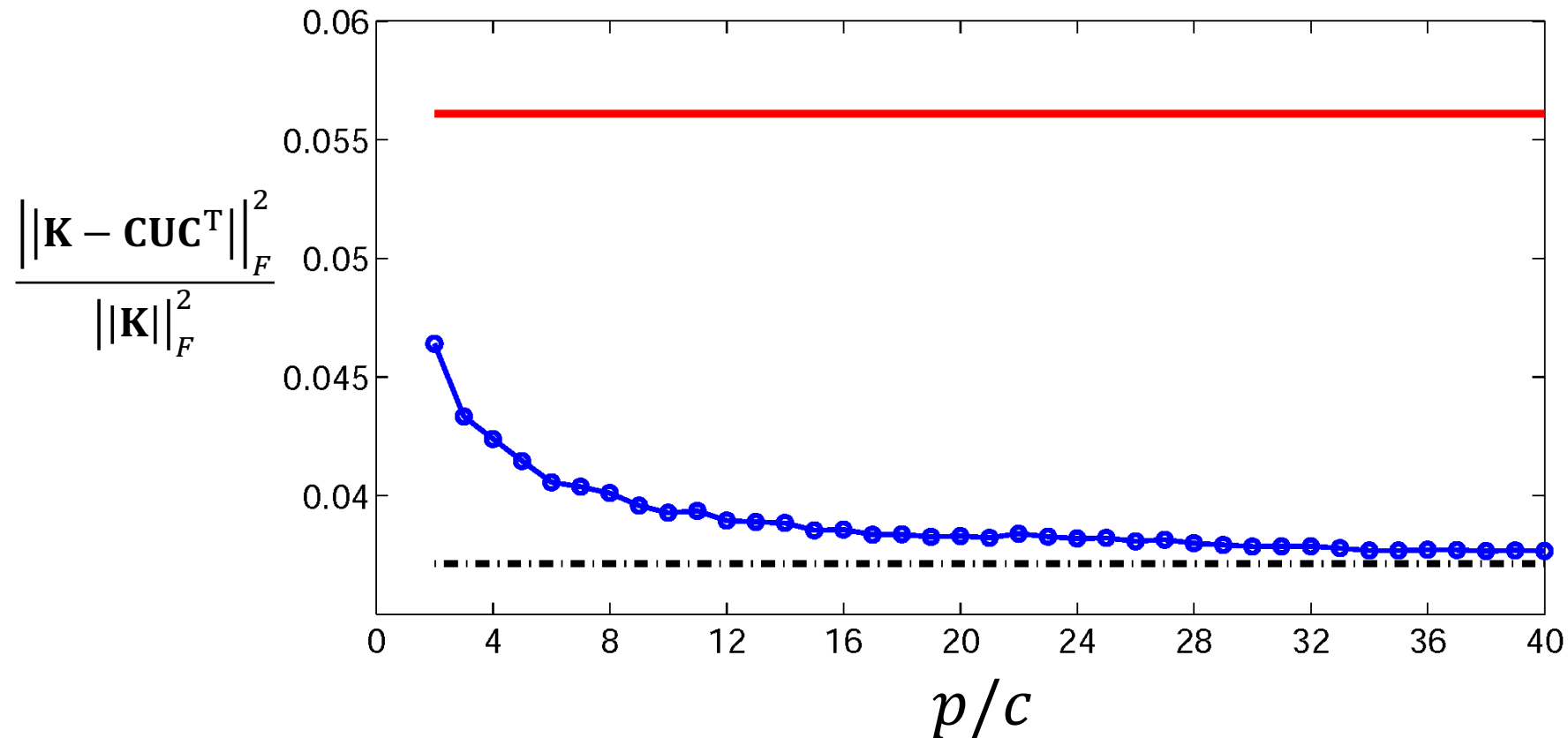> When $\mathbf{P} = \mathbf{I}_n$, the prototype model $\Longleftrightarrow$ the fast model

# Comparisons

- $\mathbf{C} = \mathbf{K}\mathbf{S} \in \mathbb{R}^{n \times c}, \mathbf{W} = \mathbf{S}^{\mathrm{T}}\mathbf{K}\mathbf{S} = \mathbf{S}^{\mathrm{T}}\mathbf{C} \in \mathbb{R}^{c \times c}$

- SPSD matrix approximation: $\mathbf{K} \approx \mathbf{C}\mathbf{U}\mathbf{C}^{\mathrm{T}}$

  - The prototype model: $\mathbf{U} = \mathbf{C}^{\dagger}\mathbf{K}(\mathbf{C}^{\dagger})^{\mathrm{T}}$

  - The fast model: $\mathbf{U} = (\mathbf{P}^{\mathrm{T}}\mathbf{C})^{\dagger}(\mathbf{P}^{\mathrm{T}}\mathbf{K}\mathbf{P})(\mathbf{C}^{\mathrm{T}}\mathbf{P})^{\dagger}$

  - The Nyström method: $\mathbf{U} = \mathbf{W}^{\dagger}$

When $\mathbf{P} = \mathbf{S}$, the Nyström method $\Leftrightarrow$ the fast model

# Comparisons

- $c = 150, n = 100c$, vary $p$ from $2c$ to $40c$



The Nyström Method
$O(nc^2)$ time

The Fast Model
$O(nc^2 + p^2c)$ time

The Prototype Model
$O(n^2c)$ time

# Conclusions

- Motivations
  - Avoid forming the kernel matrix
  - Avoid inversion/decomposition
- Prototype model, fast model, Nystrom
  - They have connections
  - The fast model and Nystrom are practical