

Kernelization of matrix updates, when and how?

Manfred Warmuth¹ Wojciech Kotłowski² Shuisheng Zhou³

¹University of California, Santa Cruz

²Poznań University of Technology, Poland

³Xidian University, Xian, China

ALT 2012, Lyon

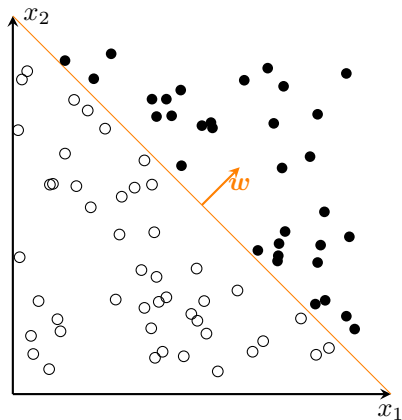
Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates

Outline

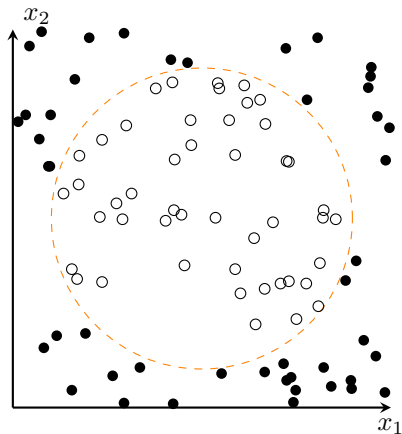
- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates

Simple start – linear binary classification



- Examples $\mathcal{S} = \{(\mathbf{x}_t, \ell_t)\}_{t=1}^T$.
- Linear hypothesis $\mathbf{w} = \mathbf{w}(\mathcal{S})$
- Prediction on instance \mathbf{x} : $\mathbf{w} \cdot \mathbf{x}$

What if data not close to linear?

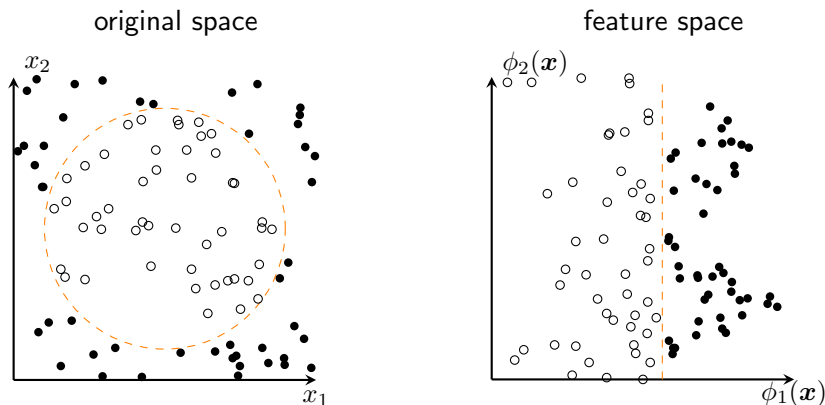


- Simply invent new features

Close to linear in feature space

Embed instances into a feature space:

$$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^m$$



■ $\phi(x_1, x_2) = \left(\sqrt{x_1^2 + x_2^2}, \arctan \frac{x_2}{x_1} \right)$.

The Kernel Trick [BGV92]

- If w linear combination of expanded instances, then

$$w \cdot x = \underbrace{\sum_{t=1}^T c_t \phi(x_t) \cdot \phi(x)}_w = \sum_{t=1}^T c_t \underbrace{\phi(x_t) \cdot \phi(x)}_{K(x_t, x)}$$

- Kernel function $K(x, \tilde{x})$ often efficient to compute.

Which algorithms can be kernelized?

- Necessary and sufficient conditions for kernelizability.
 - Vectors [Warmuth & Vishwanathan, 2005]
 - Asymmetric matrices.
 - Symmetric matrices.
- Prove new versions of the Representer Theorem.
- Build kernelizable algorithms with matrix parameters.
 - Multiplicative updates.

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates

- Labeled sequence of examples $\mathcal{S} = \{(\mathbf{x}_t, \ell_t)\}_{t=1}^T$.
- Instance matrix \mathbf{X} .
Contains instances $\mathbf{x}_1, \dots, \mathbf{x}_T$ as columns.
- Kernel matrix $\mathbf{X}'\mathbf{X}$.
Contains pairwise dot products between the instances \mathbf{x}_t :

$$(\mathbf{X}'\mathbf{X})_{ij} = \mathbf{x}_i \cdot \mathbf{x}_j$$

(we assume instances \mathbf{x}_t are already expanded).

- Augmented kernel matrix $\widehat{\mathbf{X}}'\widehat{\mathbf{X}}$, where $\widehat{\mathbf{X}}$ contains instances $\mathbf{x}_1, \dots, \mathbf{x}_T$ and an unlabeled instance \mathbf{x} as columns.

Learning algorithm:

$$\mathcal{A}: \mathcal{S} \times \mathbf{x} \rightarrow \mathbb{R}.$$

- A mapping from example sequence $\mathcal{S} = \{(\mathbf{x}_t, \ell_t)\}_{t=1}^T$ followed by a next instance \mathbf{x} to some output range.

Definition: kernelizable algorithm

\mathcal{A} is **kernelizable**, if for all \mathcal{S}, \mathbf{x} with the same labels and the same augmented kernel matrix $\widehat{\mathbf{X}}' \widehat{\mathbf{X}}$, algorithm \mathcal{A} maps to the same output.

Kernelizable algorithms: rotational invariance

- For any orthogonal matrix U , let:

$$US = \{(Ux_t, \ell_t)\}_{t=1}^T.$$

- Transformation does not affect the kernel matrix:

$$(UX)'UX = X'U'UX = X'X.$$

- **Lemma:**

Two sequences are orthogonal transformations of each other if and only if the kernel matrices associated with the sequences are the same.

Kernalizable algorithms: definition

Original definition

\mathcal{A} is **kernelizable**, if for all \mathcal{S}, \mathbf{x} with the same labels and the same augmented kernel matrix $\widehat{\mathbf{X}}' \widehat{\mathbf{X}}$, algorithm \mathcal{A} maps to the same output.

Rewritten definition (based on Lemma)

\mathcal{A} is **kernelizable** if and only if for all sequences \mathcal{S} , next instances \mathbf{x} , and orthogonal matrices \mathbf{U} , $\mathcal{A}(\mathcal{S}, \mathbf{x}) = \mathcal{A}(\mathbf{U}\mathcal{S}, \mathbf{U}\mathbf{x})$.

- 1 \mathcal{A} produces a **weight vector**:

$$w: \mathcal{S} \rightarrow \mathbf{R}^n.$$

- 2 \mathcal{A} predicts with:

$$\mathcal{A}(\mathcal{S}, \mathbf{x}) = \mathbf{w} \cdot \mathbf{x}.$$

Theorem

A linear algorithm \mathcal{A} is kernelizable if and only if for every input sequence $\mathcal{S} = \{(\mathbf{x}_t, \ell_t)\}_{t=1}^T$:

- 1 \mathbf{w} is a linear combination of the instances of \mathcal{S} ,

$$\mathbf{w} = \sum_{t=1}^T c_t \mathbf{x}_t,$$

- 2 The coefficients $\{c_t\}_{t=1}^T$ depend on \mathcal{S} only via the kernel matrix $\mathbf{X}'\mathbf{X}$.

Which algorithms are kernelizable?

- Algorithm that makes linear predictions with the parameter vector which is a linear combination of instances?

- Algorithms to which the Representer Theorem [Kimeldorf & Wahba, 71] applies?

$$\mathbf{w} = \arg \min_{\tilde{\mathbf{w}}} \left(\|\tilde{\mathbf{w}}\|^2 + \eta \sum_t \text{loss}(\tilde{\mathbf{w}} \cdot \mathbf{x}_t) \right)$$

Solution \mathbf{w} is a linear combination of the \mathbf{x}_t .

Which algorithms are kernelizable?

- Algorithm that makes linear predictions with the parameter vector which is a linear combination of instances?

⇒ Necessary, but not sufficient.

- Algorithms to which the Representer Theorem [Kimeldorf & Wahba, 71] applies?

$$\mathbf{w} = \arg \min_{\tilde{\mathbf{w}}} \left(\|\tilde{\mathbf{w}}\|^2 + \eta \sum_t \text{loss}(\tilde{\mathbf{w}} \cdot \mathbf{x}_t) \right)$$

Solution \mathbf{w} is a linear combination of the \mathbf{x}_t .

⇒ Sufficient, but not necessary.

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case**
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates

Asymmetric matrix data

- Instances are **outer products** $\mathbf{x}_t \mathbf{y}'_t$, where $\mathbf{x}_t \in \mathbf{R}^n$, $\mathbf{y}_t \in \mathbf{R}^m$

$$\begin{array}{|c|} \hline \mathbf{x} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \mathbf{y}' \\ \hline \end{array} = \begin{array}{|c|} \hline \mathbf{x}\mathbf{y}' \\ \hline \end{array}$$

$$\mathcal{S} = \{(\mathbf{x}_t \mathbf{y}'_t, \ell_t)\}_{t=1}^T.$$

- Instance matrices:**

\mathbf{X} containing instances $\mathbf{x}_1, \dots, \mathbf{x}_T$ as columns,

\mathbf{Y} containing instances $\mathbf{y}_1, \dots, \mathbf{y}_T$ as columns.

- Two kernel matrices $\mathbf{X}'\mathbf{X}$ and $\mathbf{Y}'\mathbf{Y}$.
- Corresponding augmented kernel matrices $\widehat{\mathbf{X}}'\widehat{\mathbf{X}}$ and $\widehat{\mathbf{Y}}'\widehat{\mathbf{Y}}$.
- Application:** collaborative filtering:
 \mathbf{x}_t describes user,
 \mathbf{y}_t describes movie.

- **Definition:** \mathcal{A} is **kernelizable**, if for all $\mathcal{S}, \mathbf{xy}'$ with the same labels and the same augmented kernel matrices $\widehat{\mathbf{X}}'\widehat{\mathbf{X}}, \widehat{\mathbf{Y}}'\widehat{\mathbf{Y}}$, \mathcal{A} maps to the same output.
- For any orthogonal matrices $\mathbf{U} \in \mathbb{R}^{n \times n}, \mathbf{V} \in \mathbb{R}^{m \times m}$, let:

$$\mathbf{USV}' = \{(\mathbf{U}\mathbf{x}_t\mathbf{y}_t'\mathbf{V}', \ell_t)\}_{t=1}^T.$$

- **Rewritten definition:** \mathcal{A} is **kernelizable** if and only if for all sequences \mathcal{S} , next instances \mathbf{x} , and orthogonal matrices \mathbf{U}, \mathbf{V} ,

$$\mathcal{A}(\mathcal{S}, \mathbf{xy}') = \mathcal{A}(\mathbf{USV}', \mathbf{Uxy}'\mathbf{V}').$$

- 1 \mathcal{A} produces a **weight matrix**:

$$\mathbf{W} : \mathcal{S} \rightarrow \mathbf{R}^{m \times n}.$$

- 2 \mathcal{A} predicts with:

$$\mathcal{A}(\mathcal{S}, \mathbf{x}\mathbf{y}') = \text{tr}(\mathbf{W}' \mathbf{x}\mathbf{y}') = \mathbf{y}'\mathbf{W}'\mathbf{x}.$$

Theorem

A linear algorithm \mathcal{A} is kernelizable if and only if for every input sequence $\mathcal{S} = \{(\mathbf{x}_t \mathbf{y}'_t, \ell_t)\}_{t=1}^T$:

- 1 \mathbf{W} is a linear combination of the instances of \mathcal{S} ,

$$\mathbf{W} = \sum_{i,j=1}^T c_{ij} \mathbf{x}_i \mathbf{y}'_j = \mathbf{X} \mathbf{C} \mathbf{Y}',$$

- 2 The coefficients $\{c_{ij}\}_{i,j=1}^T$ depend on \mathcal{S} only via the kernel matrices $\mathbf{X}'\mathbf{X}$ and $\mathbf{Y}'\mathbf{Y}$.

Theorem

A linear algorithm \mathcal{A} is kernelizable if and only if for every input sequence $\mathcal{S} = \{(\mathbf{x}_t \mathbf{y}'_t, \ell_t)\}_{t=1}^T$:

- 1 \mathbf{W} is a linear combination of the instances of \mathcal{S} ,

$$\mathbf{W} = \sum_{i,j=1}^T c_{ij} \mathbf{x}_i \mathbf{y}'_j = \mathbf{X} \mathbf{C} \mathbf{Y}',$$

- 2 The coefficients $\{c_{ij}\}_{i,j=1}^T$ depend on \mathcal{S} only via the kernel matrices $\mathbf{X}'\mathbf{X}$ and $\mathbf{Y}'\mathbf{Y}$.

- Compare with vector case: $\mathbf{w} = \sum_{t=1}^T c_t \mathbf{x}_t$.

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case**
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates

Symmetric matrix data

- Instances are **outer products** $\mathbf{x}_t \mathbf{x}'_t$, where $\mathbf{x}_t \in \mathbf{R}^n$.

The diagram shows a vertical rectangular box on the left containing the variable \mathbf{x} . To its right is a horizontal rectangular box containing the variable \mathbf{x}' . An equals sign is placed between these two boxes and a larger square box on the right. Inside the square box is the expression $\mathbf{x}\mathbf{x}'$, representing the result of the outer product.

$$\mathcal{S} = \{(\mathbf{x}_t \mathbf{x}'_t, \ell_t)\}_{t=1}^T.$$

- **Instance matrix:**
 \mathbf{X} containing instances $\mathbf{x}_1, \dots, \mathbf{x}_T$ as columns,
- Kernel matrix $\mathbf{X}'\mathbf{X}$.
- Corresponding augmented kernel matrix $\widehat{\mathbf{X}}'\widehat{\mathbf{X}}$.
- **Application:** PCA.

Kernelizable algorithm

- **Definition:** \mathcal{A} is **kernelizable**, if for all $\mathcal{S}, \mathbf{x}\mathbf{x}'$ with the same labels and the same augmented kernel matrix $\widehat{\mathbf{X}}'\widehat{\mathbf{X}}$, \mathcal{A} maps to the same output.
- For any orthogonal matrix U , let:

$$USU' = \{(U\mathbf{x}_t\mathbf{x}'_tU', \ell_t)\}_{t=1}^T.$$

- **Rewritten definition:** \mathcal{A} is **kernelizable** if and only if for all sequences \mathcal{S} , next instances $\mathbf{x}\mathbf{x}'$, and orthogonal matrices U ,

$$\mathcal{A}(\mathcal{S}, \mathbf{x}\mathbf{x}') = \mathcal{A}(USU', U\mathbf{x}\mathbf{x}'U').$$

- 1 \mathcal{A} produces a **symmetric weight matrix**:

$$\mathbf{W} : \mathcal{S} \rightarrow \mathbf{R}^{n \times n}.$$

- 2 \mathcal{A} predicts with:

$$\mathcal{A}(\mathcal{S}, \mathbf{x}\mathbf{x}') = \text{tr}(\mathbf{W} \mathbf{x}\mathbf{x}') = \mathbf{x}'\mathbf{W}\mathbf{x}.$$

Theorem

A linear algorithm \mathcal{A} is kernelizable if and only if for every input sequence $\mathcal{S} = \{(\mathbf{x}_t \mathbf{x}'_t, \ell_t)\}_{t=1}^T$:

- 1 \mathbf{W} is a linear combination of the instances of \mathcal{S} and the identity matrix,

$$\mathbf{W} = \sum_{i,j=1}^T c_{ij} \mathbf{x}_i \mathbf{x}'_j + c \mathbf{I} = \mathbf{X} \mathbf{C} \mathbf{X}' + c \mathbf{I}.$$

- 2 The coefficients $\{c_{ij}\}_{i,j=1}^T, c$ depend on \mathcal{S} only via the kernel matrix $\mathbf{X}' \mathbf{X}$.

Theorem

A linear algorithm \mathcal{A} is kernelizable if and only if for every input sequence $\mathcal{S} = \{(\mathbf{x}_t \mathbf{x}'_t, \ell_t)\}_{t=1}^T$:

- 1 \mathbf{W} is a linear combination of the instances of \mathcal{S} and the identity matrix,

$$\mathbf{W} = \sum_{i,j=1}^T c_{ij} \mathbf{x}_i \mathbf{x}'_j + c \mathbf{I} = \mathbf{X} \mathbf{C} \mathbf{X}' + c \mathbf{I}.$$

- 2 The coefficients $\{c_{ij}\}_{i,j=1}^T, c$ depend on \mathcal{S} only via the kernel matrix $\mathbf{X}' \mathbf{X}$.

- Compare with asymmetric case: $\mathbf{W} = \sum_{i,j=1}^T c_{ij} \mathbf{x}_i \mathbf{y}'_j$.

	vectors	asym. matrices	sym. matrices
instances	\mathbf{x}_t	$\mathbf{x}_t \mathbf{y}'_t$	$\mathbf{x}_t \mathbf{x}'_t$
kernel matrix	$\mathbf{X}' \mathbf{X},$	$\mathbf{X}' \mathbf{X}, \mathbf{Y}' \mathbf{Y}$	$\mathbf{X}' \mathbf{X}$
transformation	$\mathcal{S} \mapsto \mathbf{U} \mathbf{S}$	$\mathcal{S} \mapsto \mathbf{U} \mathbf{S} \mathbf{V}'$	$\mathcal{S} \mapsto \mathbf{U} \mathbf{S} \mathbf{U}'$
weights	$\mathbf{w} \in \mathbb{R}^n$	$\mathbf{W} \in \mathbb{R}^{n \times m}$	$\mathbf{W} \in \mathbb{R}^{n \times n}, \mathbf{W}' = \mathbf{W}$
prediction	$\mathbf{w} \cdot \mathbf{x}$	$\text{tr}(\mathbf{W}' \mathbf{x} \mathbf{y}')$	$\text{tr}(\mathbf{W} \mathbf{x} \mathbf{x}')$
kernelizability	$\mathbf{w} = \sum_i c_i \mathbf{x}_i$	$\mathbf{W} = \sum_{ij} c_{ij} \mathbf{x}_i \mathbf{y}'_j$	$\mathbf{W} = \sum_{ij} c_{ij} \mathbf{x}_i \mathbf{x}'_j + c \mathbf{I}$

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem**
- 6 Efficiency
- 7 Example: Multiplicative updates

$$\inf_{\mathbf{W}} \Omega(\mathbf{W}) + \eta \sum_t \text{loss}(\text{tr}(\mathbf{W} \mathbf{x}_t \mathbf{y}'_t))$$

- [Abernethy et al, 09] Kernelizable, if $\Omega(\mathbf{W})$ is:
 - spectrally invariant,
 - non-decreasing in singular values of \mathbf{W} .

Our version of the Representer Theorem

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}).$$

- Kernelizable, if:
 - \mathbf{W}^* is unique,
 - $\mathcal{L}(\mathbf{W}, \mathcal{S})$ is rotationally invariant in a sense that for any orthogonal \mathbf{U}, \mathbf{V} :

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}) \implies \mathbf{U}\mathbf{W}^*\mathbf{V}' = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathbf{U}\mathcal{S}\mathbf{V}')$$

- Conditions incomparable to [Abernethy et al, 09].
- E.g.,

$$\inf_{\mathbf{W} \in \mathcal{W}} \Omega(\mathbf{W}) + \eta \sum_t \text{loss}(\text{tr}(\mathbf{W} \mathbf{x}_t \mathbf{y}_t')),$$

for some spectrally invariant Ω and \mathcal{W} .

The Representer Theorem for symmetric matrices

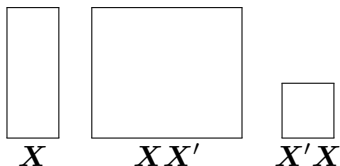
$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}).$$

- Kernelizable, if:
 - \mathbf{W}^* is unique,
 - $\mathcal{L}(\mathbf{W}, \mathcal{S})$ is rotationally invariant in a sense that for any orthogonal \mathbf{U} :

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}) \implies \mathbf{U}\mathbf{W}^*\mathbf{U}' = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathbf{U}\mathcal{S}\mathbf{U}')$$

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency**
- 7 Example: Multiplicative updates



- Eigendecomposition of **large** Gram Matrix XX' can be computed from eigendecomposition of **small** kernel matrix $X'X$.
- Key to kernelization of PCA and Fisher Linear Discriminant Functions. [Schölkopf et al., 98; Mika et al., 99]

- Need decomposition of **large** Gram matrix $\mathbf{XCY}' \in \mathbf{R}^{n \times m}$ in terms of **small** kernel matrices $\mathbf{X}'\mathbf{X}$ and $\mathbf{Y}'\mathbf{Y} \in \mathbf{R}^{T \times T}$

Lemma

SVD of generalized Gram Matrix $\mathbf{X} \mathbf{C} \mathbf{Y}'$ can be computed from the SVD of $\sqrt{\mathbf{X}'\mathbf{X}} \mathbf{C} \sqrt{\mathbf{Y}'\mathbf{Y}}$

$n \times T$ $T \times T$ $T \times m$

- Symmetric case covered by setting $\mathbf{Y} = \mathbf{X}$.

Outline

- 1 Introduction
- 2 The vector case
- 3 The asymmetric matrix case
- 4 The symmetric matrix case
- 5 Representer Theorem
- 6 Efficiency
- 7 Example: Multiplicative updates**

Matrix Exponentiated Gradient

- Offline problem:

$$\min_{\mathbf{W} \in \mathcal{W}} \sum_t \ell_t(\text{tr}(\mathbf{W} \mathbf{x}_t \mathbf{x}_t')),$$

where \mathcal{W} is a set of positive definite matrices with unit trace.

- Online algorithm derived from:

$$\begin{aligned} \mathbf{W}_{t+1} = \operatorname{argmin}_{\mathbf{W} \in \mathcal{W}} & \text{tr}(\mathbf{W} (\mathbf{log} \mathbf{W} - \mathbf{log} \mathbf{W}_t)) \\ & + \eta \text{tr} \left(\frac{\partial \ell(\text{tr}(\mathbf{W}'_t \mathbf{x}_t \mathbf{x}_t'))}{\partial \mathbf{W}_t} \mathbf{W} \right). \end{aligned}$$

Rotationally invariant problem with unique solution

⇒ kernelizable.

Matrix Exponentiated Gradient

- MEG update:

$$\mathbf{W}_{t+1} = \frac{\exp\left(-\eta \sum_{i=1}^t \partial_i(\mathbf{W}_i) \mathbf{x}_i \mathbf{x}_i'\right)}{Z_t},$$

where $\partial_i = \frac{\partial \ell_i(z)}{\partial z} \Big|_{z=\mathbf{W}_i \mathbf{x}_i \mathbf{x}_i'}$ and Z_t is such that $\text{tr}(\mathbf{W}_{t+1}) = 1$.

- has the form $\mathbf{X} \mathbf{C} \mathbf{X}' + c \mathbf{I}$ where \mathbf{C} is diagonal.
⇒ SVD of $\mathbf{X} \mathbf{C} \mathbf{X}'$ can be computed efficiently.
- Applications: online PCA.

- Necessary and sufficient conditions for kernelizability for asymmetric and symmetric matrix data.
- New versions of the Representer Theorem.
- Kernelization of multiplicative updates and other online algorithms.

- Necessary and sufficient conditions for kernelizability for asymmetric and symmetric matrix data.
- New versions of the Representer Theorem.
- Kernelization of multiplicative updates and other online algorithms.

Thank you!