



MAX PLANCK INSTITUTE
FOR DYNAMICS OF COMPLEX
TECHNICAL SYSTEMS
MAGDEBURG



COMPUTATIONAL METHODS IN
SYSTEMS AND CONTROL THEORY

Optimization of Tensor Train Canonical Decomposition in the Support Tensor Machine

Peter Benner (Max Planck Institute DCTS and TU Chemnitz)

Sergey Dolgov (University of Bath)

Kirandeep Kour (Max Planck Institute DCTS)

Martin Stoll (TU Chemnitz)

May 17-21, 2021

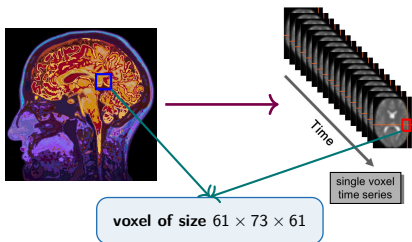
SIAM Applied Linear Algebra, 2021

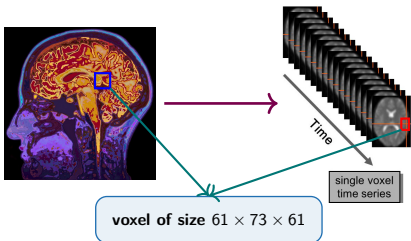
Supported by:



Partners:







Dataset

- + Small Sample Size (SSS)
- + Multi Dimensional (MD)

- Binary Classification
- Learning Model
- Nonlinear Boundary
- + SSS + MD

Efficient Model

- Robustness
 - State-of-the-art
- for SSS + MD

Model for Small Sample Size

- Binary Classification $\rightarrow X$ - input (all the data points), y - output $\{-1, 1\}$
- Learning Model \rightarrow Support Vector Machine (SVM)
- Nonlinear Boundary \rightarrow projection onto higher-dimensional space

Support Vector Machine

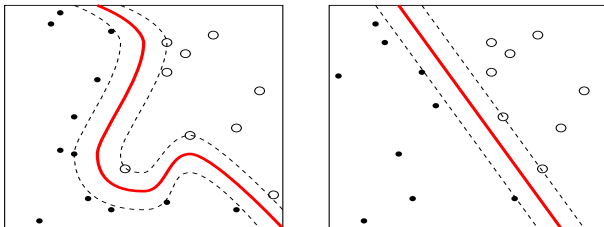


Figure: Non-linear(kernel) and linear SVM.



SVM: It is a discriminative classifier formally defined by a separating hyperplane. Generally, this hyperplane has a margin from both side of the labeled data.



SVM: It is a discriminative classifier formally defined by a separating hyperplane. Generally, this hyperplane has a margin from both side of the labeled data.

- X - input (all the data points),
- y - output $\{-1, 1\}$,
- L - hyperplane,

SVM: It is a discriminative classifier formally defined by a separating hyperplane. Generally, this hyperplane has a margin from both side of the labeled data.

- X - input (all the data points),
- y - output $\{-1, 1\}$,
- L - hyperplane,

we initialize $w \perp L$ and $\text{sign}(\langle x, w \rangle)$ helps in classification. The hypothesis is,

$$h_w(x) = \begin{cases} \langle x, w \rangle + b = 0, & y = 1, \\ \langle x, w \rangle + b \neq 0, & y = -1. \end{cases}$$

SVM: It is a discriminative classifier formally defined by a separating hyperplane. Generally, this hyperplane has a margin from both side of the labeled data.

- X - input (all the data points),
- y - output $\{-1, 1\}$,
- L - hyperplane,

we initialize $w \perp L$ and $\text{sign}(\langle x, w \rangle)$ helps in classification. The hypothesis is,

$$h_w(x) = \begin{cases} \langle x, w \rangle + b = 0, & y = 1, \\ \langle x, w \rangle + b \neq 0, & y = -1. \end{cases}$$

Taking maximum margin to the nearest point for selecting w , b and L .

$$\begin{aligned} \max_w \min_{x_i} & \quad \left| \left\langle x_i, \frac{w}{\|w\|} \right\rangle + b \right| \\ \text{subject to} & \quad \text{sign}(\langle x_i, w \rangle + b) = \text{sign}(y_i). \end{aligned}$$



Linear Boundary Optimization for SVM

$$\min_{w, b, \xi} J(w, b, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{subject to } \underbrace{y_i (w^T x_i + b)}_L \geq 1 - \xi_i, \quad \xi_i \geq 0, i = 1, \dots, N.$$

Linear Boundary Optimization for SVM

$$\min_{w, b, \xi} J(w, b, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{subject to } y_i \underbrace{(w^T x_i + b)}_L \geq 1 - \xi_i, \quad \xi_i \geq 0, i = 1, \dots, N.$$

(Dual) Linear Boundary Optimization for SVM

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0, \quad i = 1, \dots, N.$$



Projection onto higher-dimensional space

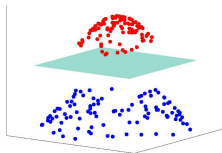
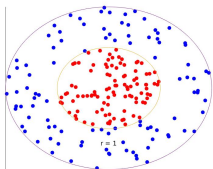


Figure: Nonlinear classification of data in (\mathbb{R}^2) .

Figure: Linear classification in higher-dimension (\mathbb{R}^3) .

Projection onto higher-dimensional space

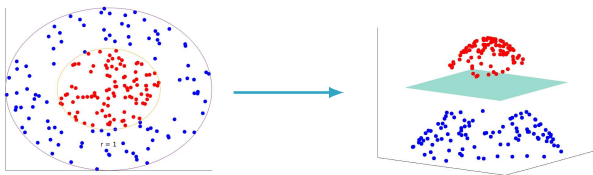


Figure: Nonlinear classification of data in (\mathbb{R}^2) .

Figure: Linear classification in higher-dimension (\mathbb{R}^3) .

Non-linear Boundary Optimization for SVM

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0, \quad i = 1, \dots, N.$$



Feature map: $\phi : \mathcal{X}(\text{input}) \longrightarrow \mathcal{V}$ (feature space)

$$\langle \phi(x), \phi(x') \rangle = k(x, x')$$



Feature map: $\phi : \mathcal{X}(\text{input}) \longrightarrow \mathcal{V}$ (feature space) $\langle \phi(x), \phi(x') \rangle = k(x, x')$

Non-linear Boundary Optimization for SVM

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0, \quad i = 1, \dots, N.$$

Feature map: $\phi : \mathcal{X}(\text{input}) \longrightarrow \mathcal{V}$ (feature space)

$$\langle \phi(x), \phi(x') \rangle = k(x, x')$$

Non-linear Boundary Optimization for SVM

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0, \quad i = 1, \dots, N.$$

Including Multi-dimensional Constraint

Classification for different datatypes

(SVM) $k(x, x') \rightarrow$ RBF

$k(\mathcal{X}, \mathcal{X}') \rightarrow$ How??



- SVM has input as a vector,
 - Storing memory for multidimensional data can be explosive,
 - Structural data can lose information during tensor-vector conversion,



- SVM has input as a vector,
 - Storing memory for multidimensional data can be explosive,
 - Structural data can lose information during tensor-vector conversion,
- Extension of SVM is available and it is called **Support Tensor Machine (STM)**,



- SVM has input as a vector,
 - Storing memory for multidimensional data can be explosive,
 - Structural data can lose information during tensor-vector conversion,
- Extension of SVM is available and it is called **Support Tensor Machine (STM)**,
- STM works directly with tensor as an input.

- SVM has input as a vector,
 - Storing memory for multidimensional data can be explosive,
 - Structural data can lose information during tensor-vector conversion,
- Extension of SVM is available and it is called **Support Tensor Machine (STM)**,
- STM works directly with tensor as an input.

Dual Non-linear Boundary Optimization Problem for STM

$$\max_{\alpha_1, \dots, \alpha_N} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0, \quad i = 1, \dots, N.$$



Tensor and Low-rank Approximation

Tensor-Train Decomposition

- Manipulating a tensor is often prone to the curse of dimensionality $\mathcal{O}(\mathcal{I}^M)$.
- A tensor decomposition works efficiently with curse of dimensionality.

- Manipulating a tensor is often prone to the curse of dimensionality $\mathcal{O}(\mathcal{I}^M)$.
- A tensor decomposition works efficiently with curse of dimensionality.
- Specially, Tensor-Train decomposition provides a compact form with storage $\mathcal{O}(MIR^2)$

Tensor-Train Decomposition [I.V. OSELEDETS; E. TYRTYSHNIKOV, 2009]

$$x_{i_1 i_2 \dots i_M} \cong \sum_{r_0, \dots, r_M} \mathfrak{g}_{r_0, i_1, r_1}^{(1)} \mathfrak{g}_{r_1, i_2, r_2}^{(2)} \dots \mathfrak{g}_{r_{M-1}, i_M, r_M}^{(M)},$$

$$\mathbf{x} \cong \langle\langle \mathfrak{G}^{(1)}, \mathfrak{G}^{(2)}, \dots, \mathfrak{G}^{(M)} \rangle\rangle,$$

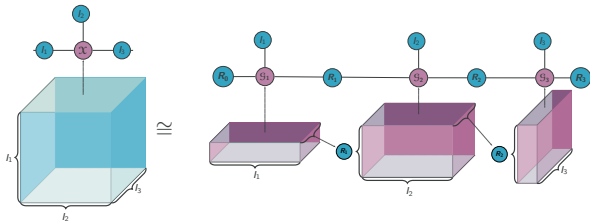


Figure: TT decomposition of a 3-way tensor.

Approximation of Kernel Computation

- Uses Canonical Polyadic (CP) decomposition [F.L. HITCHCOCK, 1927]

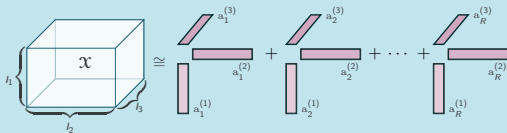


Figure: CP decomposition of a 3-way tensor.

- Dual Structure-preserving Kernel [L. HE; X. KONG; P.S. YU; A.B. RAGIN; Z. HAO; X. YANG, 2014]

(Feature map) $\Psi: \mathbf{X} \times \mathbf{Y} \times \mathbf{Z} \mapsto \mathbb{R}^{\mathcal{H}_1 \times \mathcal{H}_2 \times \mathcal{H}_3}$

$$\Psi: \sum_{r=1}^R a_r^{(1)} \otimes a_r^{(2)} \otimes a_r^{(3)} \mapsto \sum_{r=1}^R \phi(a_r^{(1)}) \otimes \phi(a_r^{(2)}) \otimes \phi(a_r^{(3)})$$

$$\langle \Psi(\mathbf{X}), \Psi(\mathbf{Y}) \rangle = k(\mathbf{X}, \mathbf{Y}) = \sum_{i,j=1}^R k(a_i^{(1)}, b_j^{(1)}) k(a_i^{(2)}, b_j^{(2)}) k(a_i^{(3)}, b_j^{(3)})$$



- Using TT decomposition
 - Stable algorithm
 - Depends on matrix SVD
 - Leads to over-fitting
 - Blocks might not be unique
- Using CP decomposition
 - Simplicity
 - Finding best rank: NP-hard [J. HÅSTAD, 1989]

TT-CP decomposition: [K. KOUR; S. DOLGOV; M. STOLL; P. BENNER, 2020]

Simplicity: Exact TT-CP Expansion

Given TT decomposition of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}$

$$x_{i_1 i_2 \dots i_M} \cong \sum_{r_1, \dots, r_{M-1}}^{R_1, \dots, R_{M-1}} \mathcal{G}_{i_1, r_1}^{(1)} \mathcal{G}_{r_1, i_2, r_2}^{(2)} \dots \mathcal{G}_{r_{M-1}, i_M}^{(M)},$$

CP expansion can be written,

$$x_{i_1 i_2 \dots i_M} \approx \sum_{r=1}^R H_{i_1, r}^{(1)} H_{i_2, r}^{(2)} \dots H_{i_M, r}^{(M)},$$

The rank r can be merged $r = r_1 + (r_2 - 1)R_1 + \dots + (r_m - 1) \prod_{\ell=1}^{m-1} R_\ell$,

where the CP factors are defined as

$$H_{i_m, r}^{(m)} = \mathcal{G}_{i_m, r}^{(m)}, \quad m = 1, \dots, M$$

Algorithm 1: Uniqueness Enforcing TT-SVD

Input: M -dimensional tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}$, prescribed accuracy ϵ .

Output: Return tensor \mathcal{X}' in **unique** TT decomposition with cores $\mathcal{G}^{(1)}, \mathcal{G}^{(2)}, \dots, \mathcal{G}^{(M)}$.

The computed approximation satisfies $\|\mathcal{X} - \mathcal{X}'\|_F \leq \epsilon \|\mathcal{X}\|_F$.

Temporary tensor: $\mathcal{Z} = \mathcal{X}, r_0 = 1$.

for $m = 1$ to $M - 1$ **do**

$\mathbf{Z} := \text{reshape}(\mathcal{Z}, [R_{m-1} I_m, I_{m+1} \dots I_M])$

Compute δ -truncated SVD: $\mathbf{Z} = \mathbf{USV}^T + E, \|E\|_F \leq \delta$.

$\mathbf{U} = [u_1, u_2, \dots, u_{R_m}], \mathbf{S} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{R_m}), \mathbf{V} = [v_1, v_2, \dots, v_{R_m}]$

$\mathbf{Z} = \sigma_1 u_1 v_1^T + \dots + \sigma_{I_m} u_{I_m} v_{I_m}^T$

for $r_m = 1$ to R_m **do**

$i_{r_m}^* = \arg \max_{i=1, \dots, I_m} |u_{r_m}(i)|$,

$\bar{u}_{r_m} := u_{r_m} / \text{sign}(u_{r_m}(i_{r_m}^*)), \bar{v}_{r_m} := v_{r_m} \cdot \text{sign}(u_{r_m}(i_{r_m}^*))$.

$\mathcal{G}_{r_{m-1}, i_m, r_m}^{(m)} = \bar{u}_{r_m}(i_m)$

end for

New core: $\mathcal{G}^{(m)} := \text{reshape}(\mathbf{U}, [R_{m-1}, I_m, R_m])$

$\mathbf{Z} := \mathbf{SV}^T$

end for

$\mathcal{G}^{(m)} = \mathbf{Z}$



Stability: Uniqueness Enforcing TT-SVD

- This ensures that “close” tensor produce “close” TT-cores
- Makes model stable w.r.t. different rank truncation

State-of-the-art: Norm Equilibration

Given a rank- r TT-CP expansion $\llbracket H^{(1)}, H^{(2)}, \dots, H^{(M)} \rrbracket$, we compute the total norm of each of the rank-1 tensors

$$\|\mathcal{N}_r\| = \|H_r^{(1)}\| * \|H_r^{(2)}\| * \dots * \|H_r^{(M)}\|,$$

and distribute this norm equally among the factors,

$$H_r^{(k)} := \frac{H_r^{(k)}}{\|H_r^{(k)}\|} * \|\mathcal{N}_r\|^{1/M}, \quad k = 1, 2, \dots, M.$$

Algorithm 2: TT-MMK: STM Kernel Approximation via TT-CP

Input: data $\{\mathcal{X}_n\}_{n=1}^N \mathbb{R}^{l_1 \times l_2 \times \dots \times l_M}$, TT-rank R .

Output: Kernel Approximation $k \in \mathbb{R}^{N \times N}$

for $n = 1$ **to** N **do**

 Compute TT decomposition $\mathcal{X}_n \cong \langle\langle \mathcal{G}^{(1,n)}, \mathcal{G}^{(2,n)}, \dots, \mathcal{G}^{(M,n)} \rangle\rangle$,

 Compute TT-CP decomposition

$\llbracket H^{(1,n)}, H^{(2,n)}, \dots, H^{(M,n)} \rrbracket = \langle\langle \mathcal{G}^{(1,n)}, \mathcal{G}^{(2,n)}, \dots, \mathcal{G}^{(M,n)} \rangle\rangle$,

end for

for $u, v = 1$ **to** N **do**

$k(\mathcal{X}_u, \mathcal{X}_v) \approx \sum_{i,j=1}^R k(h_i^{(1,u)}, h_j^{(1,v)}) k(h_i^{(2,u)}, h_j^{(2,v)}) \dots k(h_i^{(M,u)}, h_j^{(M,v)})$,

end for

Datasets:

- 1 Resting-state fMRI Brain Images
 - Alzheimer Disease (ADNI¹) - 33 data points
 - Attention Deficit Hyperactivity Disorder (ADHD²) - 200 data points

¹<http://adni.loni.usc.edu/>

²<http://neurobureau.projects.nitrc.org/ADHD200/Data.html>

³<https://aviris.jpl.nasa.gov/>

Datasets:

① Resting-state fMRI Brain Images

- Alzheimer Disease (ADNI¹) - 33 data points
- Attention Deficit Hyperactivity Disorder (ADHD²) - 200 data points

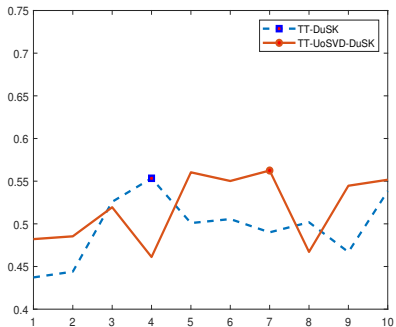
② Hyperspectral Images

- Indian Pines: The HSI images was collected via Aviris Sensor³ over Indian Pines test site - 100 data points
- Salinas: This HSI images data was collected by 224 band Aviris Sensor over Salinas vally, California - 100 data points

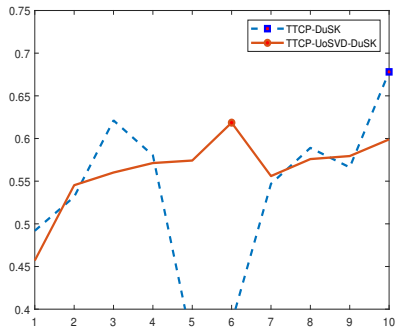
¹<http://adni.loni.usc.edu/>

²<http://neurobureau.projects.nitrc.org/ADHD200/Data.html>

³<https://aviris.jpl.nasa.gov/>

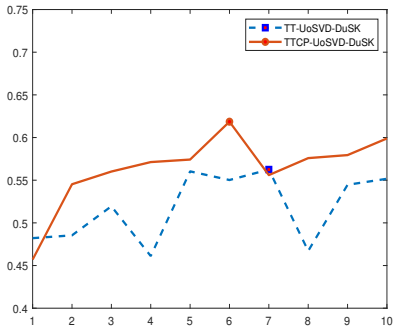


(a) TT with and without enforced uniqueness

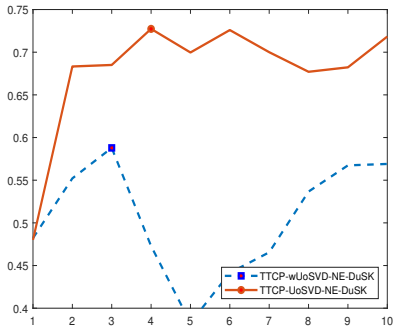


(b) TTCP with and without enforced uniqueness

Figure: Effect of each step of Algorithm for ADNI data set(Rank vs Accuracy).



(a) TT vs TTCP with enforced uniqueness



(b) TTCP norm equilibrium with and without enforced uniqueness

Figure: Effect of each step of Algorithm for ADNI data set (Rank vs Accuracy).

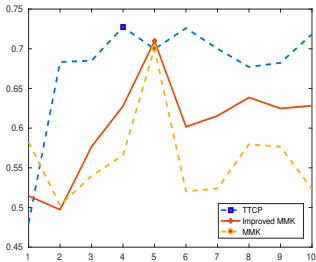


Figure: Accuracy for ADNI dataset w.r.t. truncated TT-CP rank

Table: Best average classification accuracy for different methods and datasets

Methods	ADNI	ADHD	I. Pines	Salinas
SVM	49	52	46	47
STuM	51	54	57	74
DuSK	55	58	60	92
MMK	69	60	93	98
TT-MMK	73	64	99	99

Summary

- Low rank-method “TT-CP decomposition”
- Binary classification kernel model for less data with tensor input
- Approximating kernel not only reduced cost, also increases accuracy
- Stability and reduced computational cost with higher accuracy

Summary

- Low rank-method “TT-CP decomposition”
- Binary classification kernel model for less data with tensor input
- Approximating kernel not only reduced cost, also increases accuracy
- Stability and reduced computational cost with higher accuracy

<https://arxiv.org/abs/2002.05079>

Thankful to SIAM Student Travel Award