



Shiraz University  
Computer Science and Engineering Department  
Machine Learning Lab

خواجہ نصیرالدین طوسی:  
”هر کس چیزی را بجوید و در راهش کوشش نماید؛ آن را می یابد؛  
هر کس دری را بکوبد و پایداری نماید؛ به درون خانه راه می یابد.“

# Positive Semi-Definite (PSD) Kernel Learning for Supervised and Unsupervised Problems

Dr. Sattar Hashemi

Fatemeh Alavi

alavi.s.fatemeh@gmail.com



# Outline



## Introduction

## PSD Kernels & Indefinite Kernels

## PSD Kernel Learning

- **Parametric Kernel Learning**
- **Data-Dependent Kernel Learning**

## Multiple Kernel Learning

- **Classification and Clustering Tasks**

## Neighborhood Kernel Learning

- **Classification and Clustering Tasks**

## Conclusion



# Introduction

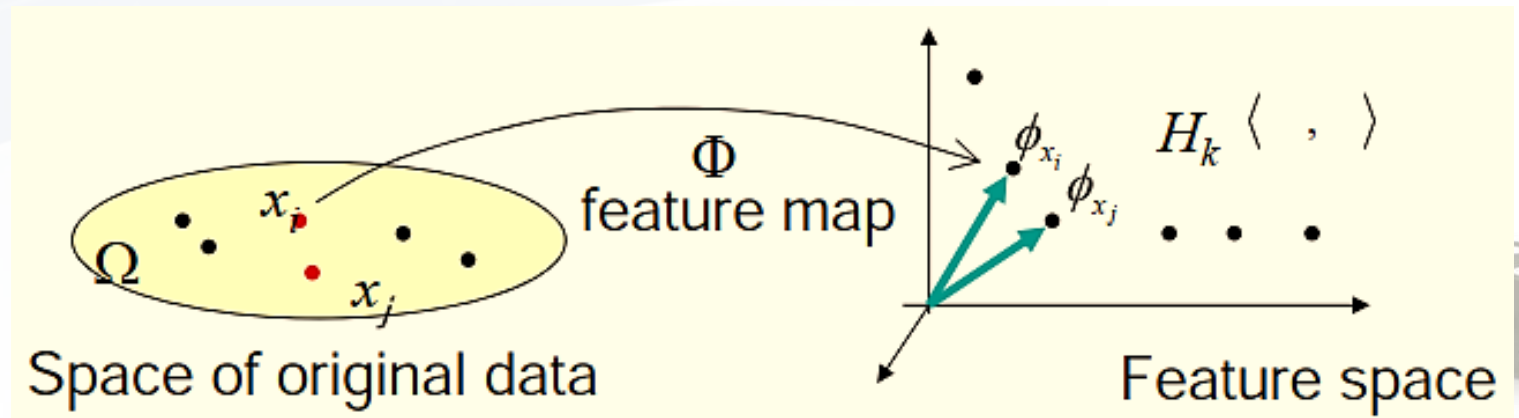


- Traditional kernel methods **implicitly** embed data in some Hilbert spaces, and search for linear relations in the Hilbert spaces.
- The kernel matrix induces a notion of data **similarity** and data relationships.
- Kernels are considered as similarity measures that arise from a particular **representation** of patterns.

**Kernel trick:**

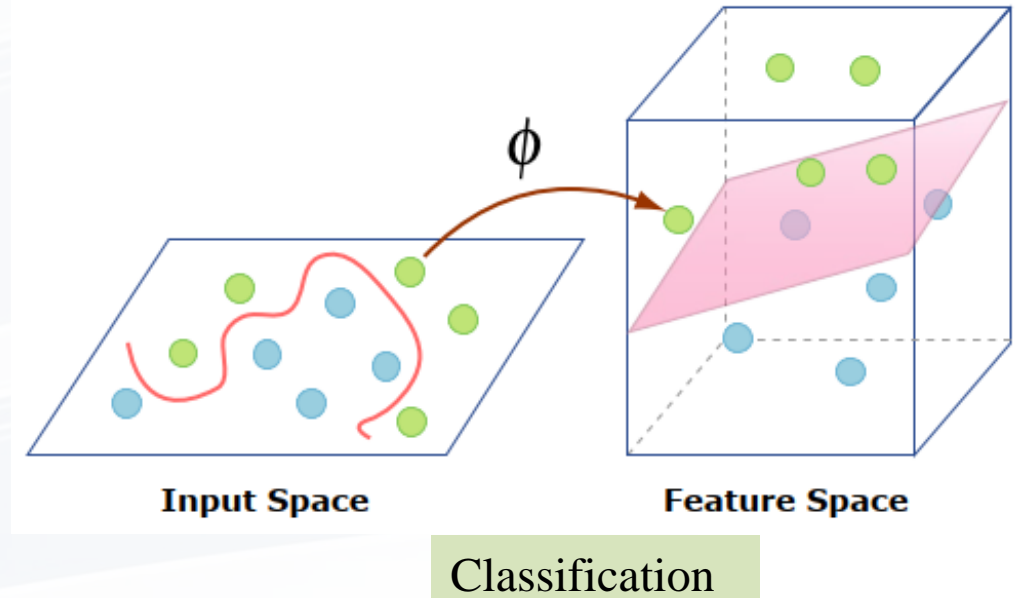
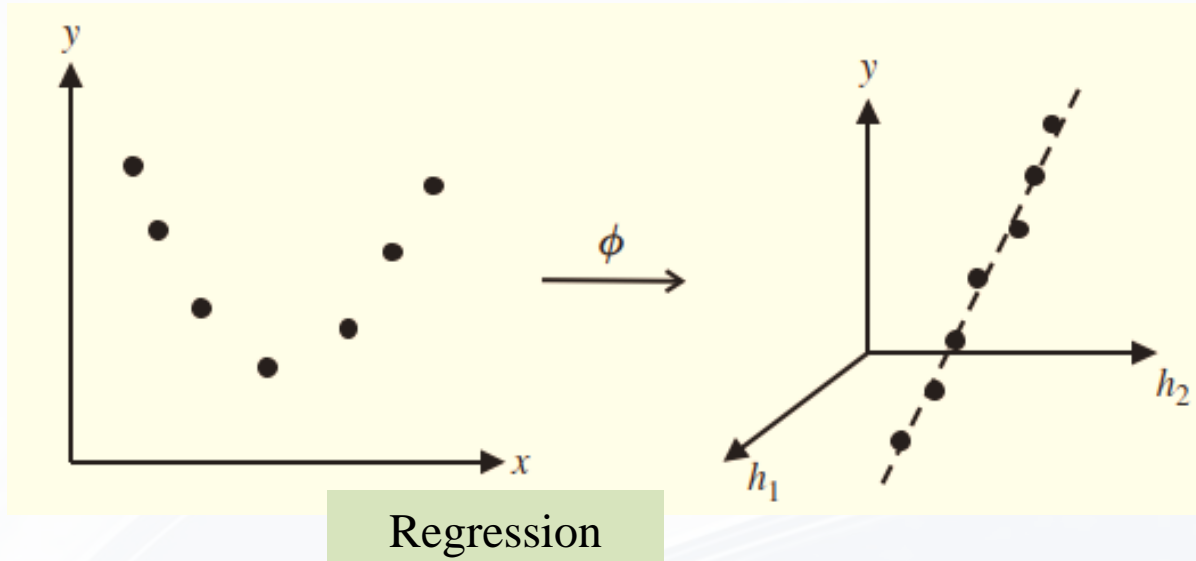
$$K(x, x') = \langle \phi(x), \phi(x') \rangle_H$$

$$\text{Cos}(\cdot) = \frac{\langle \phi(x), \phi(x') \rangle_H}{\|x\| \|x'\|}$$

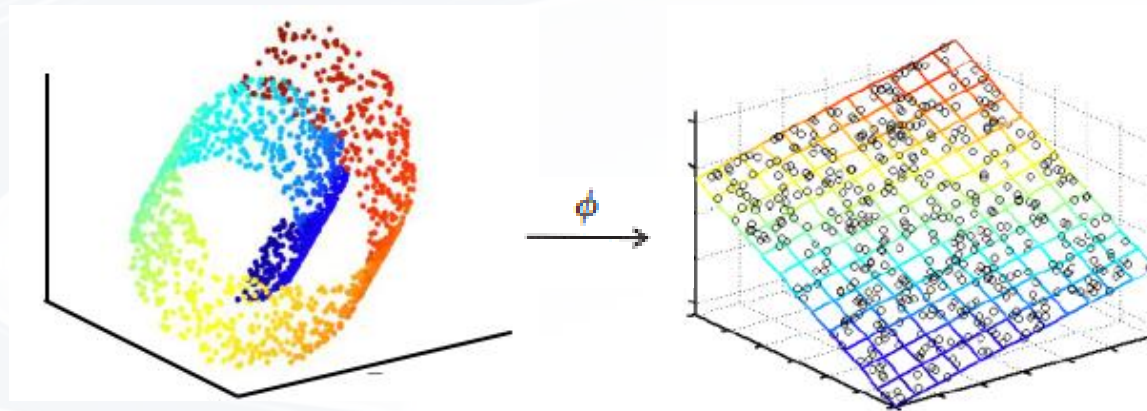


*"Learning with kernels: support vector machines, regularization, optimization, and beyond." MIT press 2018.*

# Introduction (cont.)



Unsupervised  
tasks

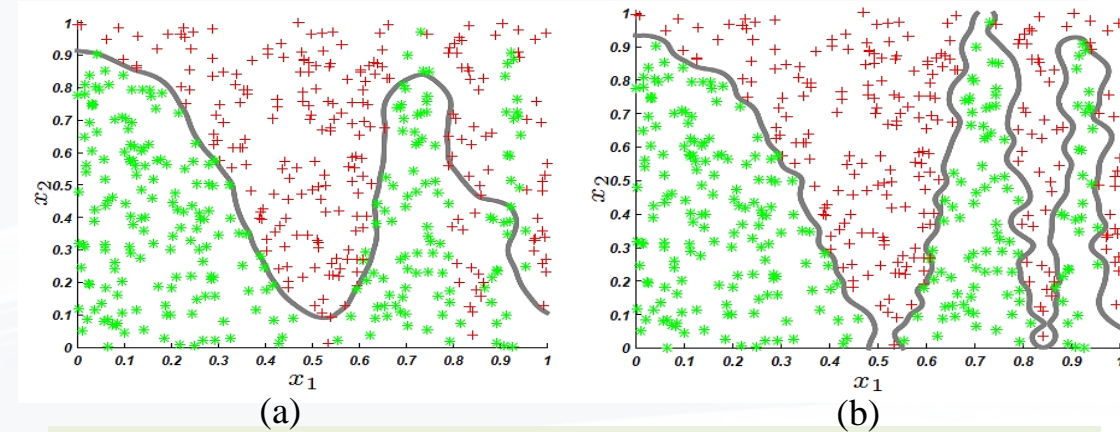




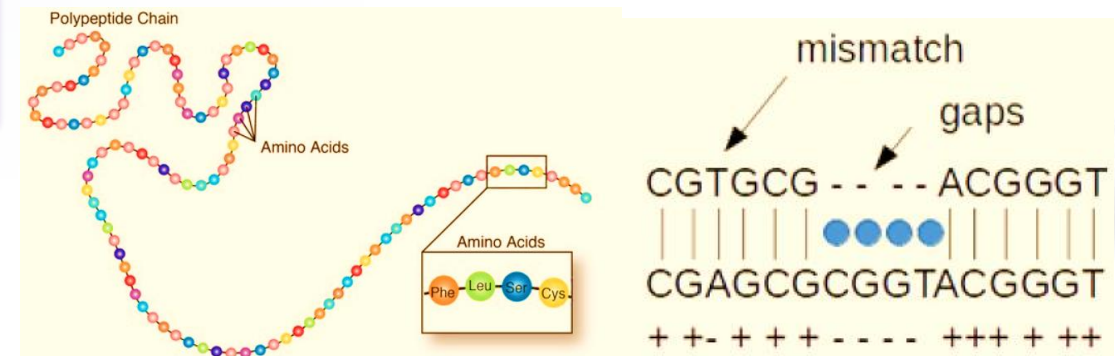
# Introduction (cont.)



- The **performance** of kernel-based methods broadly relies on selecting an appropriate kernel.
- Traditional kernels **globally** perform on given data, and cannot **locally** capture the different relationships of data in different regions.
- Classical kernel matrixes are unable to capture **complex similarities**.



Gaussian kernel:  $e^{-\gamma \|x-x'\|_2^2}$  (a)  $\gamma = 1$  (b)  $\gamma = 25$

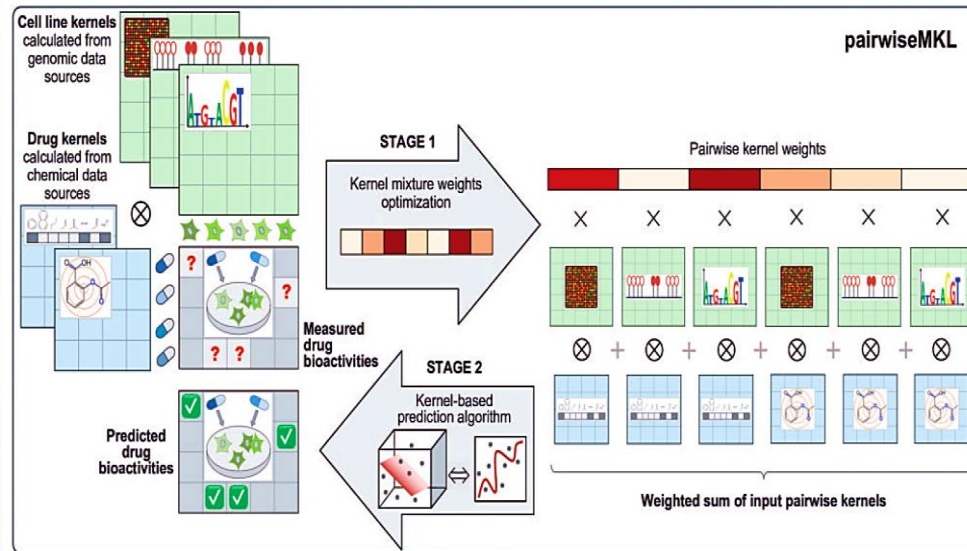


*"Classification with Truncated  $\ell_1$  Distance Kernel.", TPAMI 2015.*

# Applications of Kernel Learning

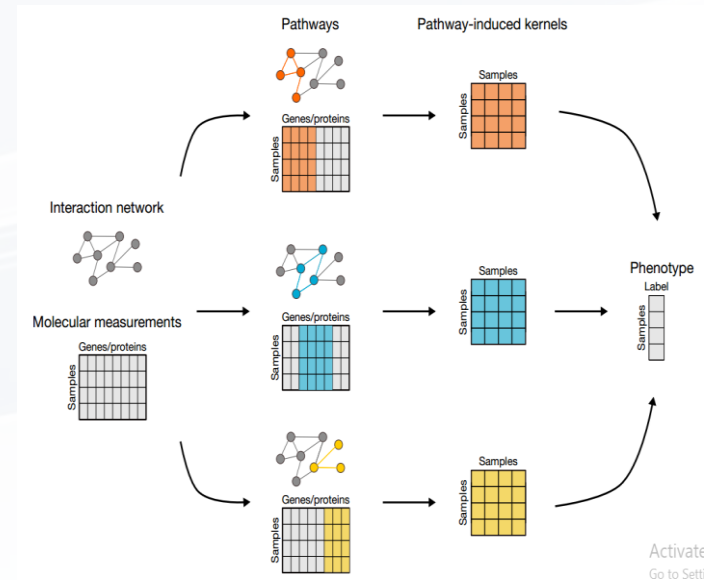


## Drug Bioactivity Prediction



*"Learning with multiple pairwise kernels for drug bioactivity prediction.", Bioinformatics 2018.*

## Bioinformatics: Patient Stratification Pathway Induced Multiple Kernel Learning



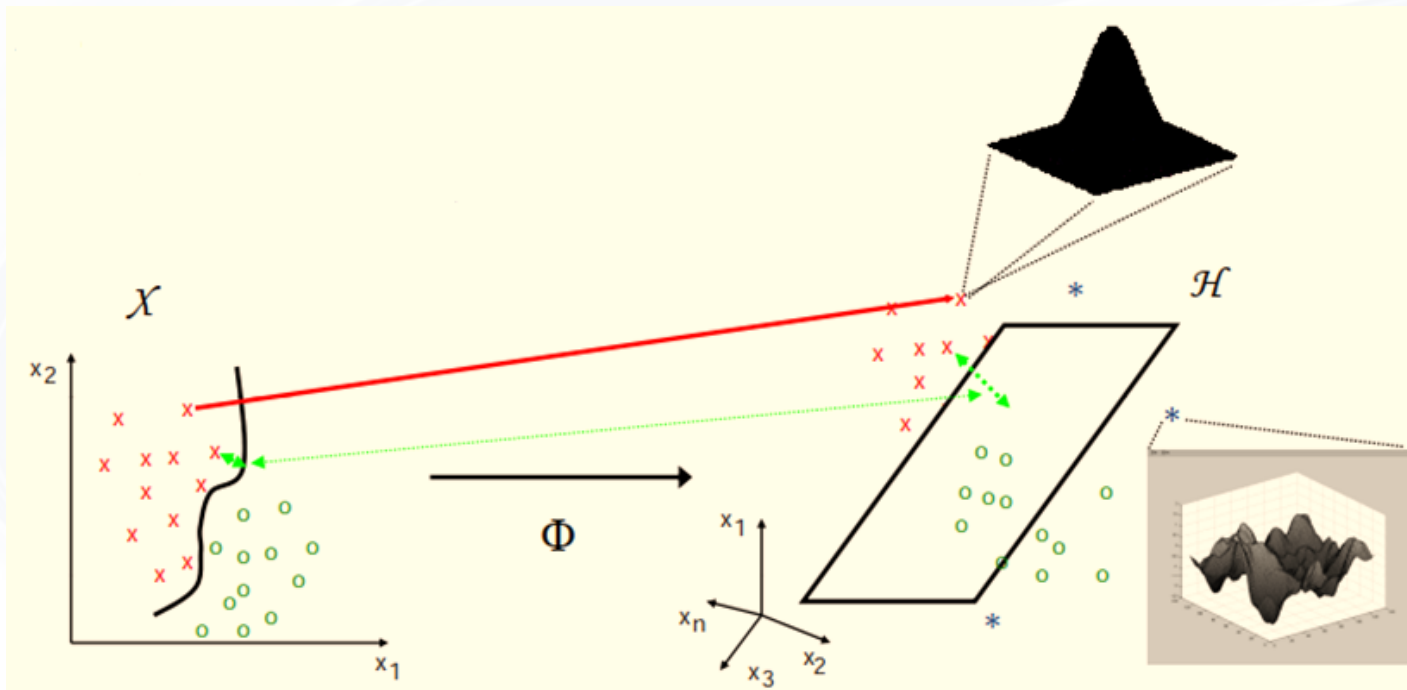
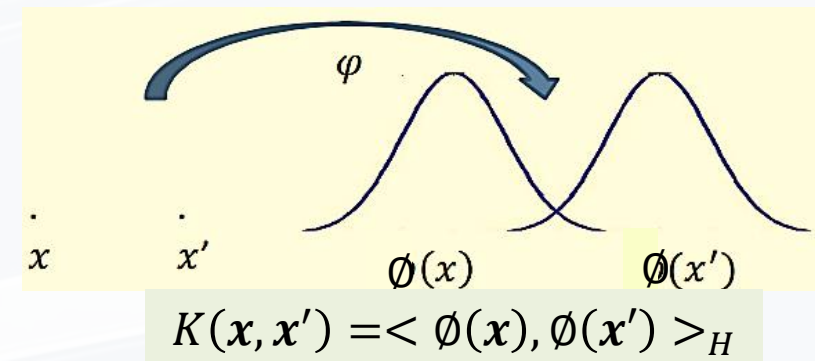
*"PIMKL: Pathway-Induced Multiple Kernel Learning.", NPJ Systems Biology and Applications, 2019.*



# Deep Viewpoint



- **Bochner's theorem:** a kernel can be represented (in dual form) as a probability distribution, and so the search for a kernel becomes a search over distributions.



$$\Phi(x) = K(., x) = e^{-\gamma \|.-x\|^2}$$

$$\Phi(x') = K(., x') = e^{-\gamma \|.-x'\|^2}$$

$$\langle \Phi(x), \Phi(x') \rangle = K(x, x') = e^{-\gamma \|x-x'\|^2}$$

*"Learning with kernels: support vector machines, regularization, optimization, and beyond." MIT press 2018.*

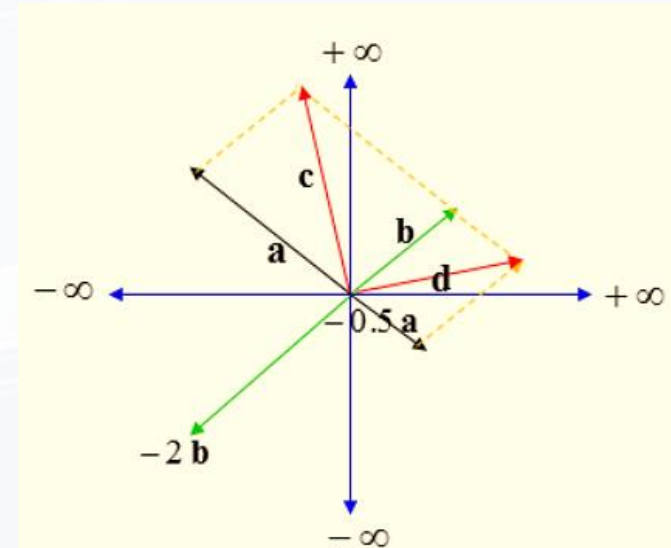
# Hilbert Space (Functional Space)



- **Vector Space:**

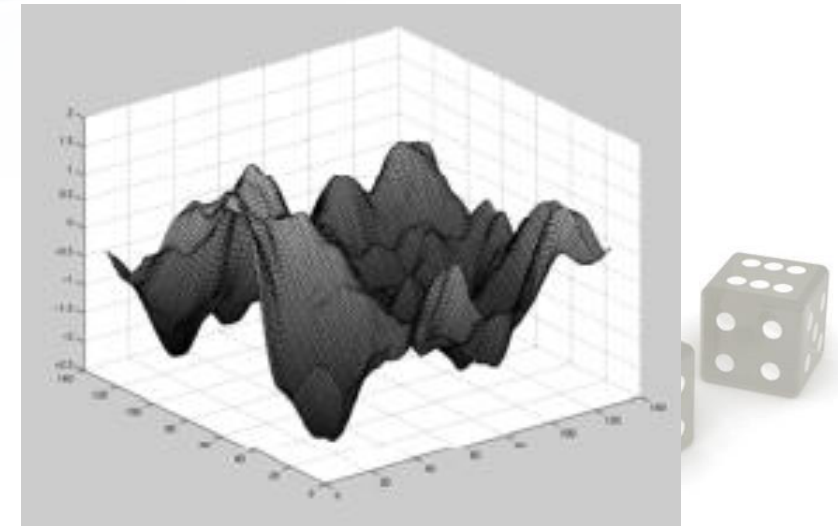
- V set of objects: vectors & functions
- + vector addition:  $\mathbf{a} \in V, \mathbf{b} \in V \rightarrow \mathbf{a} + \mathbf{b} \in V$
- . scalar multiplication:  $\mathbf{b} \in V, \alpha \in R \rightarrow \alpha \mathbf{b} \in V$

$$V = \{v | v = \sum_{i=1}^n \alpha_i x_i \quad \forall x_i \in X\}$$



- **Hilbert space:** a complete vector space with dot product and a norm
- Functional space:

$$H = \{f(.) | f(.) = \sum_{i=1}^n \alpha_i K(., x_i) = \sum_{i=1}^n \alpha_i \Phi(x_i) \quad \forall x_i \in X\}$$





# Representer Theorem



## Representer Theorem:

$$f^* = \underset{f(x) \in H}{\operatorname{argmin}} \operatorname{loss}((x_1, y_1, f(x_1)), \dots, (x_n, y_n, f(x_n))) + \Omega(\|f\|_H^2)$$

although the optimization problem seems to be in an infinite-dimensional space, the solution only lies in the span of  $n$  particular kernels centered on  $n$  training points.

$$w = \sum_{i=1}^n \alpha_i \Phi(x_i) \rightarrow$$
$$f(x) = \langle w, \Phi(x) \rangle + w_0 = \sum_{i=1}^n \alpha_i \underbrace{\langle \Phi(x_i), \Phi(x) \rangle}_{K(x_i, x)} + w_0$$

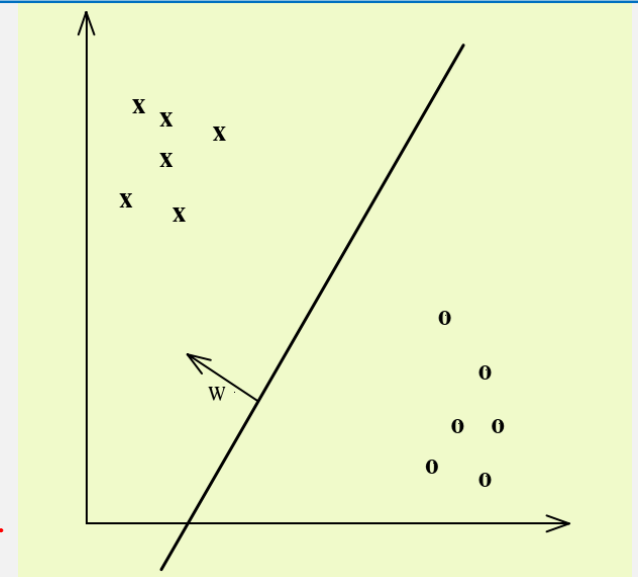
*“Kernel methods for pattern analysis.”, 2004, Cambridge university press.*

*“Learning with kernels: support vector machines, regularization, optimization, and beyond.”, 2018, the MIT Press.*

*“Learning kernel classifiers: theory and algorithms.”, 2001, MIT press.*

$$w = \sum_{i=1}^n \alpha_i x_i \rightarrow$$
$$f(x) = \langle w, x \rangle + w_0 = \sum_{i=1}^n \alpha_i \langle x_i, x \rangle + w_0$$

$$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$
$$(x, x') \mapsto k(x, x'),$$



# Valid Kernel



Without constructing  $\phi(\mathbf{x})$ , Necessary and sufficient condition for  $k(\mathbf{x}, \mathbf{x}')$  to be a kernel is:

Matrix kernel  $\mathbf{K}$  is positive semi-definite (PSD)

**Mercer's theorem:** any continuous, symmetric, positive semi-definite kernel function  $k(\mathbf{x}, \mathbf{x}')$  can be expressed as a dot product in a high-dimensional space

$$\mathbf{K} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

- $k(x_i, x_i) \geq 0 \quad \forall x_i \in X$
- All eigenvalues of  $\mathbf{K}$  satisfy  $\lambda_i \geq 0$



# PSD Kernels vs. Indefinite Kernels



## Indefinite Kernels RKKS Space

Representer theorem

Non-convex problems

No Mercer's condition

No kernel trick

## PSD Kernels RKHS Space

Representer theorem

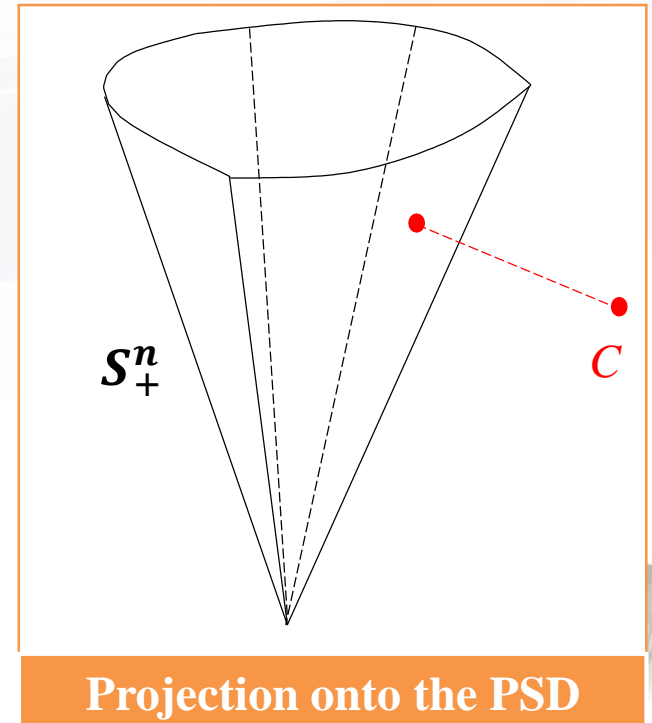
Convex problems

Mercer's condition

Kernel trick

$\mathcal{S}^n$  : the space of symmetric matrices

$\mathcal{S}_+^n$  : cone of PSD matrices



*"Classification with Truncated  $\ell_1$  Distance Kernel.", TPAMI 2015.*

*"Indefinite Kernel Logistic Regression with Concave-Inexact-Convex Procedure.", TNNLS 2018.*

# PSD Kernel Learning Methods



## Parametric Kernel Learning (Inductive Learning)

- Multiple Kernel Learning
  - Linear & Nonlinear Kernel Learning
  - Finite & Infinite Kernel Learning
- Deep Kernel Learning
- Bayesian Multiple Kernel Learning

## Data-dependent Kernel Learning (Transductive Learning)

- Optimal Neighborhood Kernel Learning
- Low-rank Kernel Learning





# Parametric Kernel Learning: Multiple KL

## Finite kernel learning

- Base kernels should be specified **in advance**.
- Optimal kernel is a weighted combination of predefined base kernels over a parametrized **discrete** set.
- **Gaussian kernel**:  $e^{-\gamma\|x-x'\|_2^2}$ ,  $\gamma = 2^{-15}, \dots, 2^{15}$

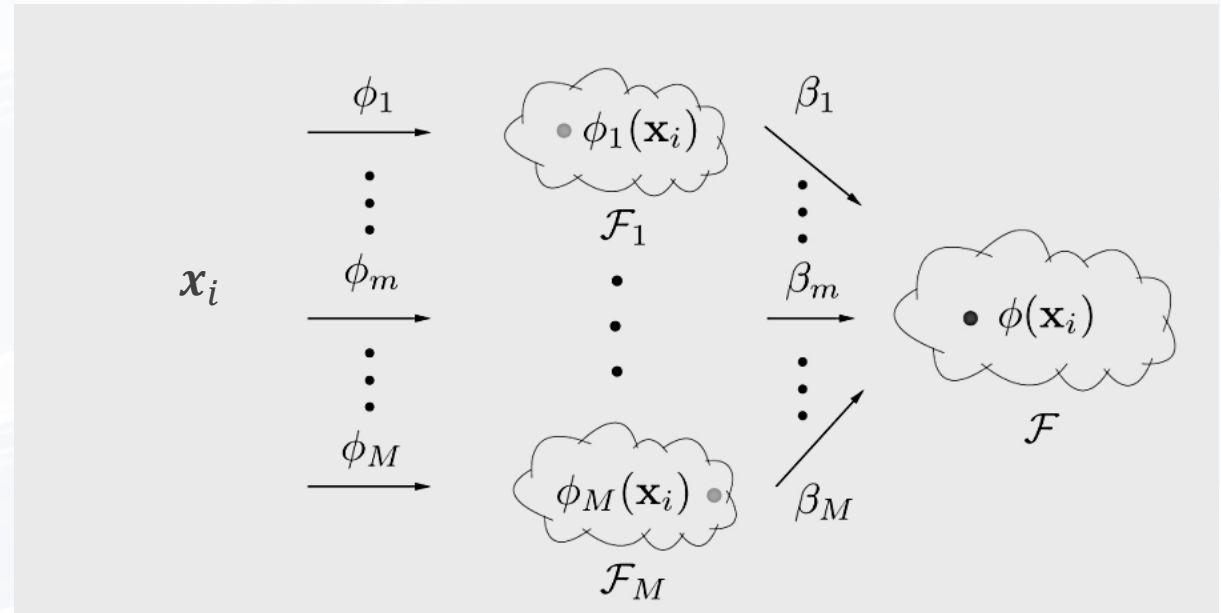
## Infinite kernel learning

$$K^{infinite} = \left\{ \int_{\Omega} K_{\theta} dp(\theta) : p \in \mathcal{M}(\Omega) \right\}; \theta \in \Omega$$

- Base kernels can be learned **automatically**.
- Infinite kernel learning improves the accuracy.
- Optimal kernel is a weighted combination of learned base kernels over a **continuous** set.
- **Gaussian kernel**:  $e^{-\mu\|x-x'\|_2^2}$ ,  $\mu \in [2^{-15}, 2^{15}]$

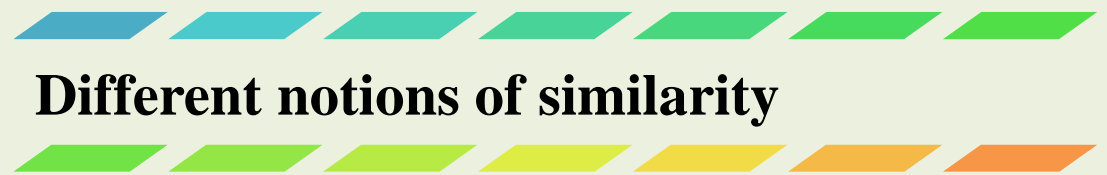
*"Multiple kernel learning algorithms.", JMLR 2011.*

*"Infinite kernel learning: generalization bounds and algorithms.", AAAI 2017.*



**Different representations**

**Different notions of similarity**



# Deep Kernel Learning



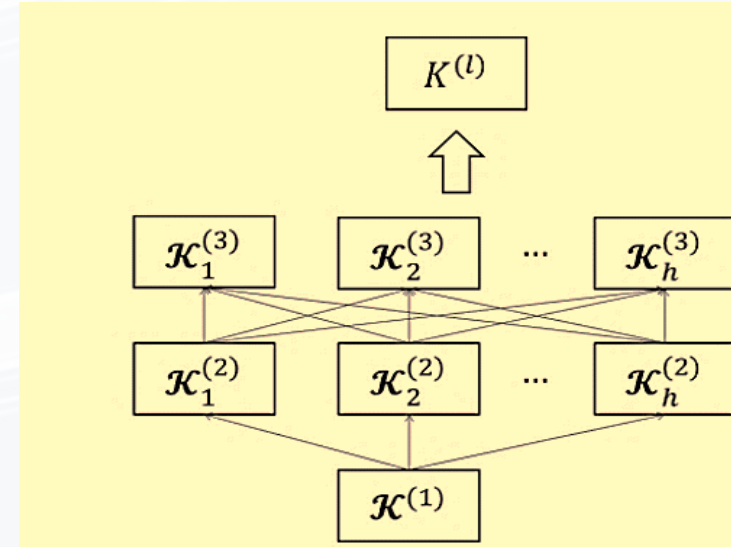
## Multi-layer Composition of kernels

$$\mathcal{K}^{(l)}(\mathbf{x}_i, \mathbf{x}_j) = \phi^{(l)}(\dots \phi^{(1)}(\mathbf{x}_i)) \cdot \phi^{(l)}(\dots \phi^{(1)}(\mathbf{x}_j))$$

$$K^{(2)} = \phi^{(2)}(\phi^{(1)}(x_i)) \cdot \phi^{(2)}(\phi^{(1)}(x_j)) = e^{-2\gamma} e^{2\gamma K(x_i x_j)}$$

## Nested Kernel:

$$K^{(l)} = \{w_{1,1}^{(l)} K_{1,1}^{(l)} (w_{1,1}^{(l-1)} K_{1,1}^{(l-1)} + \dots) + \dots w_{h,m}^{(l)} K_{h,m}^{(l)} (\dots)\}$$



## Deep Learning (Advantages):

- End-to-end learning
- “richness” of representations

## Deep Learning (Limitations):

- Too many hyper-parameters
- a single type of kernel
- Tuning of several hyper-parameters in the discrete space

*“Two-Layer Multiple Kernel Learning.”, AISTATS 2011.*

*“A Representer Theorem for Deep Kernel Learning.” JMLR 2019.*



# MKL (Classification Tasks)



**Kernel Target Alignment (KTA)** shows that in the feature space the data distribution is somehow correlated to the label distribution.

$$KTA(K_1, K_2) = \frac{\langle K_1, K_2 \rangle_F}{\|K_1\|_F \|K_2\|_F} \longrightarrow \text{Ideal Kernel: } K_y = yy^T \begin{cases} K_{ij}^y = 1 \text{ if } y_i = y_j \\ K_{ij}^y = -1 \text{ if } y_i \neq y_j \end{cases}$$

*"On kernel-target alignment. ", NIPS 2001.*

*Alignf Algorithm*

$$a = [\langle K_1, K_y \rangle_F \quad \dots \quad \langle K_m, K_y \rangle_F], \quad M_{ij} = \langle K_i, K_j \rangle_F = \text{Tr}(K_i K_j)$$

$$\text{QP: } \max_w \frac{\langle K_w, K_y \rangle_F}{\sqrt{\langle K_w, K_w \rangle_F}}, \text{ s. t. } \|w\|_2 = 1, w \geq 0, K_w = \sum_{p=1}^m w_p K_p \rightarrow \min_{\|w\|_2=1, w \geq 0} w^T M w - 2w^T a$$

*"Algorithms for Learning Kernels Based on Centered Alignment. ", JMLR 2012.*

# MKL (Classification Tasks)



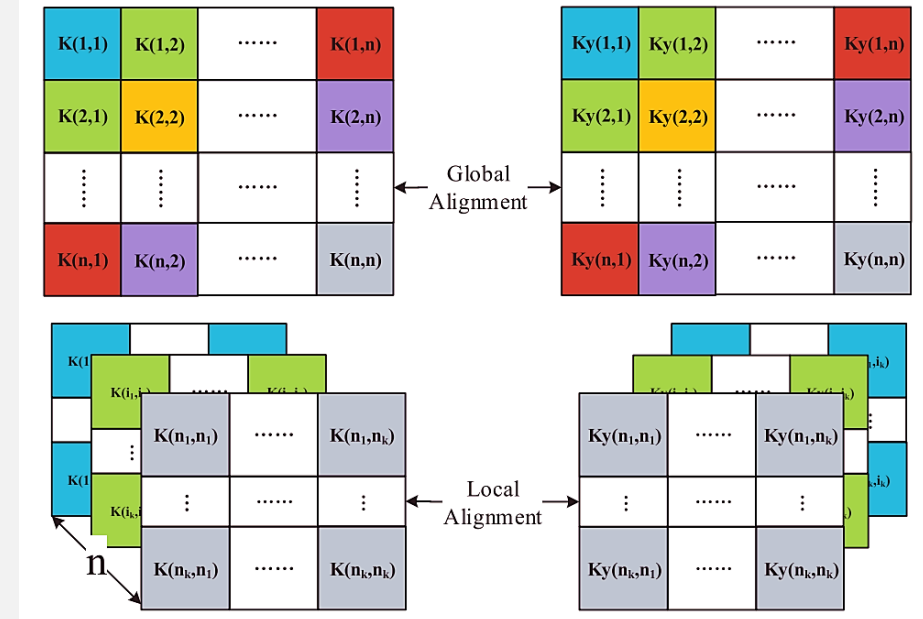
- MKL with Hybrid Kernel Alignment (HKTA)

$$KTA_{hybrid}(K_1, K_2) = (1 - \lambda)KTA_{local}(K_1, K_2) + \lambda KTA_{global}(K_1, K_2)$$

$$\max_{w \geq 0, \|w\|_2=1} (1 - \lambda) \left( \frac{1}{n} \sum_{i=1}^n \frac{\langle K_w^{(i)}, K_y^{(i)} \rangle_F}{\|K_w^{(i)}\|_F} \right) + \lambda \frac{\langle K_w, K_y \rangle_F}{\|K_w\|_F}$$

Local Kernel:

$$K^{(i)} \in R^{k \times k}$$



"Multiple kernel learning with hybrid kernel alignment maximization .", Pattern Recognition 2017.



# Multiple Kernel K-Means (Clustering Tasks)



## Kernel k-means (KKM)

$$\min_{\mathbf{Z} \in \{0,1\}^{n \times k}} \sum_{i=1, c=1}^{n,k} Z_{ic} \|\phi(\mathbf{x}_i) - \boldsymbol{\mu}_c\|_2^2 \quad s.t. \quad \sum_{c=1}^k Z_{ic} = 1, \boldsymbol{\mu}_c = \frac{1}{n_c} \sum_{i=1}^n Z_{ic} \phi(\mathbf{x}_i)$$



Reformulated equation

$$\min_{\mathbf{H} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{K}(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)) \quad s.t. \quad \mathbf{H}^T \mathbf{H} = \mathbf{I}_k$$

## Multiple kernel k-means (MKKM)

$$K_w(\mathbf{x}_i, \mathbf{x}_j) = \phi_w(\mathbf{x}_i)^T \phi_w(\mathbf{x}_j) = \sum_{p=1}^m w_p^2 K_p(\mathbf{x}_i, \mathbf{x}_j)$$

$$\phi_w(\mathbf{x}) = [w_1 \phi_1(\mathbf{x})^T, w_2 \phi_2(\mathbf{x})^T, \dots, w_m \phi_m(\mathbf{x})^T]^T$$



Multiple kernels

$$\min_{\mathbf{H} \in \mathbb{R}^{n \times k}, \mathbf{w} \in \mathbb{R}^m} \text{Tr}(\mathbf{K}_w(\mathbf{I}_n - \mathbf{H}\mathbf{H}^T)) \\ s.t. \quad \mathbf{H}^T \mathbf{H} = \mathbf{I}_k, \quad \mathbf{w}^T \mathbf{1}_m = 1, \mathbf{w} \geq \mathbf{0}$$

"Multiple Kernel k-Means Clustering with Matrix-Induced Regularization.", AAAI 2016.

# Kernel K-Means (Clustering Tasks)



Multiple kernel k-means with matrix-induced regularization (MKKM-MR)

$$\min_{H \in \mathbb{R}^{n \times k}, \mathbf{w} \in \mathbb{R}^m} \text{Tr}(K_w(I_n - HH^T)) + \frac{\lambda}{2} \mathbf{w}^T M \mathbf{w}$$
$$\text{s. t. } H^T H = I_k, \mathbf{w}^T \mathbf{1}_m = 1, \mathbf{w} \geq 0$$

*"Multiple Kernel k-Means Clustering with Matrix-Induced Regularization.", AAAI 2016.*

Optimal neighborhood kernel clustering with multiple kernels

$$\min_{H \in \mathbb{R}^{n \times k}, G \in \mathcal{S}_+^n, \mathbf{w}} \text{Tr}(G(I_n - HH^T)) + \frac{\rho}{2} \|G - K_w\|_F^2 + \frac{\lambda}{2} \mathbf{w}^T M \mathbf{w}$$
$$\text{s. t. } H^T H = I_k, \mathbf{w}^T \mathbf{1}_m = 1, \mathbf{w} \geq 0, G \succcurlyeq 0$$

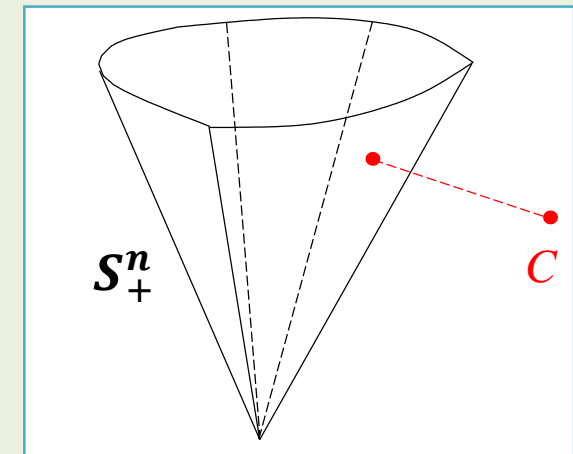
where

$$K_w(\mathbf{x}_i, \mathbf{x}_j) = \phi_w(\mathbf{x}_i)^T \phi_w(\mathbf{x}_j) = \sum_{p=1}^m w_p^2 K_p(\mathbf{x}_i, \mathbf{x}_j)$$

*"Optimal neighborhood kernel clustering with multiple kernels.", AAAI 2017.*

$\mathcal{S}^n$  : the space of symmetric matrices

$\mathcal{S}_+^n$  : cone of PSD matrices



Projection onto the PSD Cone

# References



- ✓ Schölkopf, B., Smola, A. J., & Bach, F. (2018). *”Learning with kernels: support vector machines, regularization, optimization, and beyond.”* MIT press.
- ✓ Shawe-Taylor, J., & Cristianini, N. (2004). *”Kernel methods for pattern analysis.”*, Cambridge university press.
- ✓ Moeller, J., Srikumar, V., Swaminathan, S., Venkatasubramanian, S., & Webb, D. (2016, September). **”Continuous kernel learning.”** In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases* (pp. 657-673). Springer, Cham.
- ✓ Huang, X., Suykens, J. A., Wang, S., Hornegger, J., & Maier, A. (2017). **”Classification With Truncated  $l_1$  Distance Kernel.”**, *IEEE transactions on neural networks and learning systems (TNNLS)*, 29(5), 2025-2030.
- ✓ Liu, F., Huang, X., Gong, C., Yang, J., & Suykens, J. A. (2018). **”Indefinite kernel logistic regression with concave-inexact-convex procedure.”**, *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 30(3), 765-776.
- ✓ Gönen, M., & Alpaydın, E. (2011). **”Multiple kernel learning algorithms.”** *Journal of machine learning research (JMLR)*, 12(Jul), 2211-2268.
- ✓ M. Gonen (2012), **”Bayesian efficient multiple kernel learning”**, ICML, 1206.6465.
- ✓ E.V. Strobl, S. Visweswaran, (2013) **”Deep multiple kernel learning”**, in: 12th International Conference on Machine Learning and Applications, vol. 1, IEEE, 2013, pp. 414–417.
- ✓ Aioli, M. Donini, (2015) **”EasyMKL: a scalable multiple kernel learning algorithm”**, Neurocomputing 169 215–224.
- ✓ Y. Wang, X. Liu, Y. Dou, Q. Lv, Y. Lu, (2017) **”Multiple kernel learning with hybrid kernel alignment maximization”**, Pattern Recognit. 70, 104–111 .
- ✓ N. Cristianini, J. Kandola, A. Elisseeff, J. Shawe-Taylor, (2006) **”On kernel target alignment”**, in: Innovations in Machine Learning, Springer, pp. 205–256.
- ✓ M.P. Kumar, B. Packer, D. Koller, (2010), **”Self-paced learning for latent variable models”**, in: NIPS, vol. 1, p. 2.



# References



- ✓ A.R. Meenakshi, C. Rajian, (1999) “**On a product of positive semidefinite matrices**”, Linear Algebra Appl. 295 (1–3) 3–6.
- ✓ A.R. Conn, N.I.M. Gould, P.L. Toint, (2000) “**Trust Region Methods**”, SIAM.
- ✓ Du, L., Zhou, P., Shi, L., Wang, H., Fan, M., Wang, W., Shen, Y.D., (2015) “**Robust multiple kernel k-means using l21-norm**”, in: Twenty-fourth international joint conference on artificial intelligence, pp. 3476–3482.
- ✓ Gonen, M., Margolin, A.A., (2014) “**Localized data fusion for kernel k-means clustering with application to cancer biology**”. Advances in Neural Information Processing Systems 27, 1305–1313.
- ✓ Huang, J., Nie, F., Huang, H., (2015) “**A new simplex sparse learning model to measure data similarity for clustering**”, in: Twenty-fourth international joint conference on artificial intelligence.
- ✓ Yao, Y., Li, Y., Jiang, B., Chen, H., (2020) “**Multiple kernel k-means clustering by selecting representative kernels**”. IEEE Transactions on Neural Networks and Learning Systems.
- ✓ Zhou, S., Liu, X., Li, M., Zhu, E., Liu, L., Zhang, C., Yin, J., (2019) “**Multiple kernel clustering with neighbor-kernel subspace segmentation**”. IEEE transactions on neural networks and learning systems 31, 1351–1362.
- ✓ Shen, H.T., Zhu, Y., Zheng, W., Zhu, X., (2020) “**Half-quadratic minimization for unsupervised feature selection on incomplete data**”. IEEE transactions on neural networks and learning systems 32, 3122–3135.
- ✓ Boyd, S., Parikh, N., Chu, E., (2011) “**Distributed optimization and statistical learning via the alternating direction method of multipliers**”. Now Publishers Inc.
- ✓ Charbonnier, P., Blanc-F'eraud, L., Aubert, G., Barlaud, M., (1997). “**Deterministic edge-preserving regularization in computed imaging**”. IEEE Transactions on image processing 6, 298–311.

