

# MULTIPLE-KERNEL REGRESSION WITH SPARSITY CONSTRAINTS

---

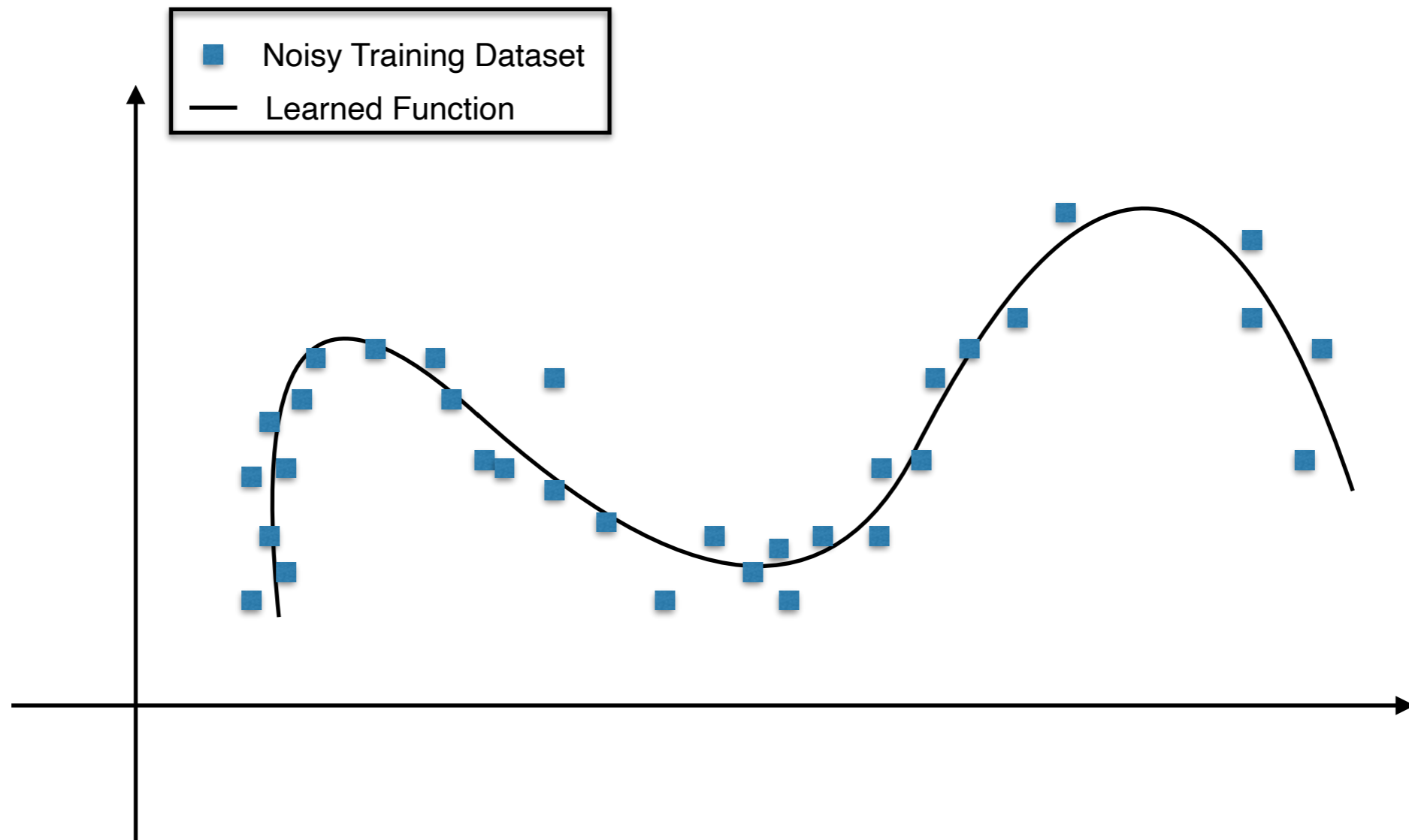
[Shayan Aziznejad](#) and Michael Unser

Biomedical Imaging Group, École Polytechnique Fédérale de Lausanne

- **Supervised Learning:** Determining an unknown function given its nonuniform samples
- The goal is to recover  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  based on its noisy samples  $y_m \approx f(\mathbf{x}_m)$ , where  $\mathbf{x}_m \in \mathbb{R}^d$  and  $m = 1, 2, \dots, M$ .
- General formulation as a minimization problem:

$$\min_{f \in \mathcal{F}} \underbrace{\sum_{m=1}^M \mathbb{E}(f(\mathbf{x}_m), y_m)}_{\text{Data Fidelity}} + \underbrace{\lambda \mathcal{R}(f)}_{\text{Regularization}}$$

# EXAMPLE



# REPRODUCING KERNEL HILBERT SPACES (RKHS)

- $\mathcal{H}(\mathbb{R}^d)$ : **Hilbert** space of functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- $\mathcal{H}(\mathbb{R}^d)$  is an RKHS, if the **sampling functional** is **continuous**.

Equivalently,  $\delta(\cdot - \mathbf{x}_0) \in \mathcal{H}'(\mathbb{R}^d)$

- **Unique** reproducing kernel of  $\mathcal{H}(\mathbb{R}^d)$ :  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that

$$k(\mathbf{x}, \cdot) \in \mathcal{H}(\mathbb{R}^d), \quad \forall f \in \mathcal{H}(\mathbb{R}^d) : f(\mathbf{x}) = \langle k(\mathbf{x}, \cdot), f(\cdot) \rangle_{\mathcal{H}}$$

# REPRESENTER THEOREM

- Supervised learning in an **RKHS**:

$$\min_{f \in \mathcal{H}(\mathbb{R}^d)} \sum_{m=1}^M \mathbb{E}(f(\mathbf{x}_m), y_m) + \lambda \|f\|_{\mathcal{H}}^2$$

- **Representer theorem** [Scholkopf et al. 2001]: If the solution exists, it is in the form of

$$f(\mathbf{x}) = \sum_{m=1}^M a_m \mathbf{k}(\mathbf{x}, \mathbf{x}_m) \quad a_m \in \mathbb{R}, \mathbf{x}_m \in \mathbb{R}^d.$$

- Linear combination of M kernels on **data points**

# APPLICATION OF REP. THEOREM

- Optimal solution:  $f(\mathbf{x}) = \sum_{m=1}^M a_m \mathbf{k}(\mathbf{x}, \mathbf{x}_m)$   $a_m \in \mathbb{R}, \mathbf{x}_m \in \mathbb{R}^d$ .
- Gram matrix:  $G \in \mathbb{R}^{M \times M}$  such that  $[G]_{m,n} = \mathbf{k}(\mathbf{x}_m, \mathbf{x}_n) \Rightarrow \|f\|_{\mathcal{H}}^2 = \mathbf{a}^T G \mathbf{a}$
- **Reduced** minimization problem:  $\min_{\mathbf{a} \in \mathbb{R}^M} \sum_{m=1}^M \mathbb{E}([G\mathbf{a}]_m, y_m) + \lambda \mathbf{a}^T G \mathbf{a}$
- **Closed-form** solution if  $\mathbb{E}(y, z) = (y - z)^2$  (Tikhonov 1963)

$$\mathbf{a} = (G^T G + \lambda G)^{-1} G^T \mathbf{y}$$

# SPARSE KERNEL EXPANSION

- **Reducing complexity**: crucial for **large datasets**
- Sparsity-enforcing loss: support-vector machines [Vapnik 1998]
- **Sparsity-promoting** regularizer: **generalized LASSO** [Roth 2004]

$$\min_{\mathbf{a} \in \mathbb{R}^M} \sum_{m=1}^M \mathbb{E}([\mathbf{G}\mathbf{a}]_m, y_m) + \lambda \|\mathbf{a}\|_{\ell_1}$$

- Banach-space formulations: gTV regularization  
[Unser et al. 2017, Bach 2017]

# GENERALIZED TOTAL VARIATION (gTV)

- **gTV norm** :  $\|\mathbf{L}\{f\}\|_{\mathcal{M}}$ , where  $\mathbf{L}$  is an invertible LSI operator
- $(\mathcal{M}(\mathbb{R}^d), \|\cdot\|_{\mathcal{M}})$  is a **generalization** of  $(L_1(\mathbb{R}^d), \|\cdot\|_{L_1})$ :

$$f \in L_1(\mathbb{R}^d) \Rightarrow \|f\|_{\mathcal{M}} = \|f\|_{L_1}$$

$$\delta(\cdot - \mathbf{x}_0) \in \mathcal{M}(\mathbb{R}^d), \quad \|\delta(\cdot - \mathbf{x}_0)\|_{\mathcal{M}} = 1$$

- The corresponding **native space**:

$$\mathcal{M}_{\mathbf{L}}(\mathbb{R}^d) = \{f : \mathbb{R}^d \rightarrow \mathbb{R}; \|\mathbf{L}\{f\}\|_{\mathcal{M}} < +\infty\}$$

- The shift-invariant kernel associated to  $\mathbf{L}$ :

$$\mathbf{k} = \mathbf{L}^{-1}\{\delta\}$$



# REPRESENTER THEOREM

- **Supervised learning** with **gTV** regularization:

$$\min_{f \in \mathcal{M}_L(\mathbb{R}^d)} \sum_{m=1}^M \mathbb{E}(f(\mathbf{x}_m), y_m) + \lambda \|\mathbf{L}\{f\}\|_{\mathcal{M}}$$

- **Representer theorem** [Unser et al. 2017]: Under certain assumptions, there exists an optimal solution that admits the kernel expansion

$$f(\mathbf{x}) = \sum_{n=1}^N a_n k(\mathbf{x} - \mathbf{z}_n) \quad a_n \in \mathbb{R}, \mathbf{z}_n \in \mathbb{R}^d, N \leq M.$$

- gTV regularization  $\Rightarrow \ell_1$  penalty on the kernel coefficients
- **Adaptive** positions  $\Rightarrow$  **Sparse** representation

# APPLICATION OF REP. THEOREM

- Optimal solution:  $f(\cdot) = \sum_{n=1}^N a_n k(\cdot - \mathbf{z}_n)$ ,
- Gram matrix:  $G_Z \in \mathbb{R}^{M \times N}$  such that  $[G_Z]_{m,n} = k(\mathbf{x}_m - \mathbf{z}_n)$
- **Reduced** minimization problem:  $\min_{\mathbf{a} \in \mathbb{R}^M, Z \in \mathbb{R}^{d \times N}} \sum_{m=1}^M \mathbb{E}([G_Z \mathbf{a}]_m, y_m) + \lambda \|\mathbf{a}\|_{\ell_1}$
- Grid-based algorithms: [Gupta et al. 2017, Debarre et al. 2019]

Suitable for **low dimensional problems**

# MULTIPLE KERNEL LEARNING

- Motivation: Increasing the model flexibility to have a more accurate regression [Lanckriet et al. 2004][Bach et al. 2004].
- $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_N$  are  $N$  RKHS with the kernel functions  $k_1, k_2, \dots, k_N$
- Learn a new positive-definite kernel  $k_\mu = \sum_{n=1}^N \mu_n k_n$
- Joint minimization problem:

$$\min_{\mathbf{a} \in \mathbb{R}^M, \boldsymbol{\mu} \in \mathbb{R}^N} \sum_{m=1}^M \mathbb{E}([G_\mu \mathbf{a}]_m, y_m) + \lambda \mathbf{a}^T G_\mu \mathbf{a}$$

# MULTIPLE KERNEL REGRESSION WITH gTV

- **Multi-component** model for the target function
- **Sparsity**-promoting regularization

$$\min_{f_n \in \mathcal{M}_{L_n}(\mathbb{R}^d)} \sum_{m=1}^M \mathbb{E}(f(\mathbf{x}_m), y_m) + \lambda \sum_{n=1}^N \|\mathbf{L}_n\{f_n\}\|_{\mathcal{M}} \quad s.t. \quad f = \sum_{n=1}^N f_n$$

- **Representer theorem:** Under certain assumptions, there exists an optimal solution that admits the kernel expansion

$$f(\mathbf{x}) = \sum_{n=1}^N \sum_{j=1}^{M_n} a_{n,j} k_n(\mathbf{x} - \mathbf{z}_{n,j}), \quad a_{n,j} \in \mathbb{R}, \mathbf{z}_{n,j} \in \mathbb{R}^d, \quad \sum_{n=1}^N M_n \leq M$$

- The number of **active kernels** is upper bounded by  $M$   
**Independent** of the number of components!!

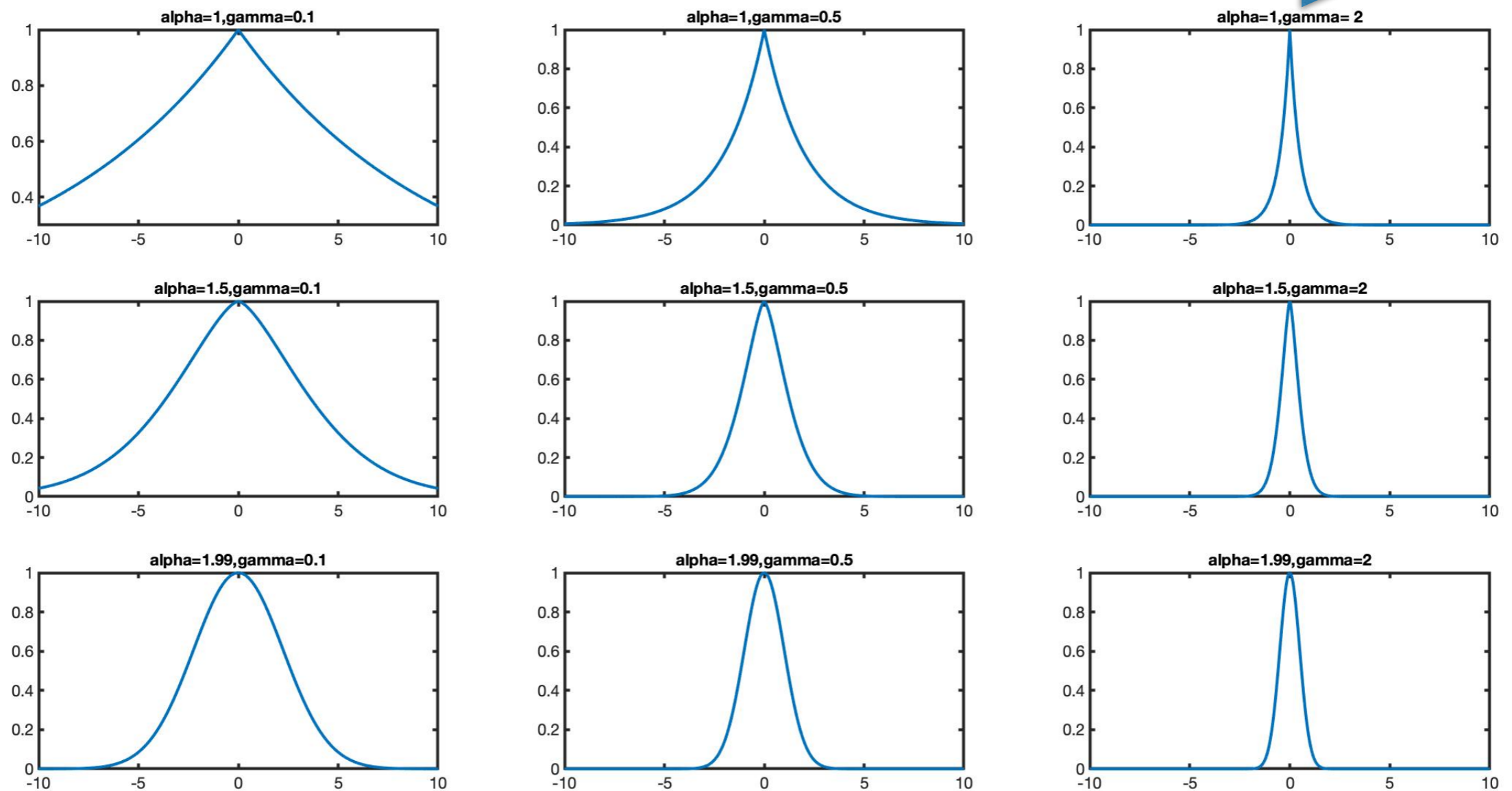
- Proposition: Any function  $k : \mathbb{R}^d \rightarrow \mathbb{R}$  with the following properties is an **admissible** kernel with respect to the gTV regularization:
  1. Non vanishing, smooth and slowly growing frequency response  $\hat{k}(\boldsymbol{\omega})$
  2. Heavy-tailed frequency response;  $\hat{k}(\boldsymbol{\omega}) \geq C(\|\boldsymbol{\omega}\| + 1)^{-\alpha}$
- Example: Characteristic function of heavy-tailed distributions.
- Example: Bessel potentials [Aronszajn 1961]

$$G_s(\mathbf{x}, \mathbf{y}) = \mathcal{F}^{-1} \left\{ \frac{1}{(1 + \|\boldsymbol{\omega}\|_2^2)^{\frac{\alpha}{2}}} \right\} (\mathbf{x} - \mathbf{y})$$

# EXAMPLES: SYMMETRIC-ALPHA-STABLES

$$k_{\alpha,\gamma}(\mathbf{x}, \mathbf{y}) = \exp(-\gamma\|\mathbf{x} - \mathbf{y}\|_{\alpha}^{\alpha}), \quad \alpha \in (0, 2), \gamma \in \mathbb{R}^{+}$$

Higher  $\gamma$



Higher  $\alpha$



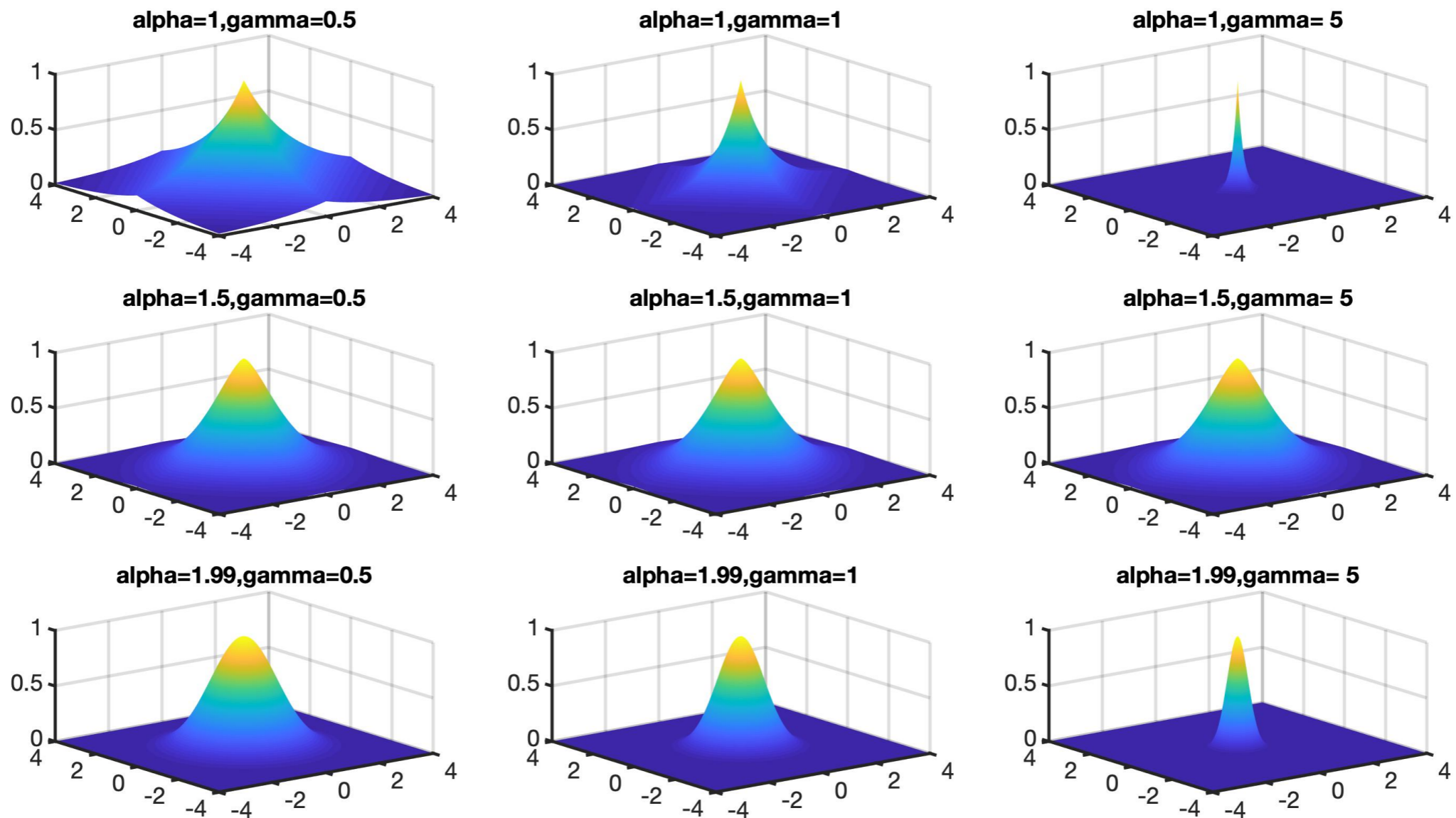
# 2D KERNELS

$$k(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^d k(x_i - y_i).$$

Higher  $\gamma$



Higher  $\alpha$



# CONCLUSION

- We considered a Banach space framework for supervised learning.
- Our framework suggests an adaptive kernel expansion for the learned function.
- We proposed a multi-component model for the target function with gTV regularization.
- Our representer theorem shows that the number of active kernels in the multi-component model is upper-bounded by the number of data points.
- We identified some classes of kernel functions that are admissible to our theory.

**Challenge: Lack of algorithms for high dimensional data**



# REFERENCES

- Schölkopf, Bernhard, Ralf Herbrich, and Alex J. Smola. "A generalized representer theorem." *International conference on computational learning theory*. Springer, Berlin, Heidelberg, 2001.
- Tikhonov, Andrei Nikolaevich. "On the solution of ill-posed problems and the method of regularization." *Doklady Akademii Nauk*. Vol. 151. No. 3. Russian Academy of Sciences, 1963.
- Vladimir Vapnik. "Statistical learning theory Wiley." *New York* (1998): 156-160.
- Roth, Volker. "The generalized LASSO." *IEEE transactions on neural networks* 15.1 (2004): 16-28.
- Unser, Michael, Julien Fageot, and John Paul Ward. "Splines are universal solutions of linear inverse problems with generalized TV regularization." *SIAM Review* 59.4 (2017): 769-793.
- Bach, Francis. "Breaking the curse of dimensionality with convex neural networks." *The Journal of Machine Learning Research* 18.1 (2017): 629-681.
- Gupta, Harshit, Julien Fageot, and Michael Unser. "Continuous-domain solutions of linear inverse problems with Tikhonov versus generalized TV regularization." *IEEE Transactions on Signal Processing* 66.17 (2018): 4670-4684.
- Debarre, Thomas, et al. "B-Spline-Based Exact Discretization of Continuous-Domain Inverse Problems with Generalized TV Regularization." *IEEE Transactions on Information Theory* (2019).
- Lanckriet, Gert RG, et al. "Learning the kernel matrix with semidefinite programming." *Journal of Machine learning research* 5.Jan (2004): 27-72.
- Bach, Francis R., Gert RG Lanckriet, and Michael I. Jordan. "Multiple kernel learning, conic duality, and the SMO algorithm." *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004.
- Aronszajn, Nachman, and Kennan T. Smith. "Theory of Bessel potentials. I." *Annales de l'institut Fourier*. Vol. 11. 1961.

THANKS FOR YOUR ATTENTION!