

# Kernel methods

## BIFOLD Aqtivate workshop

**Stefan Blücher**

10.02.2024



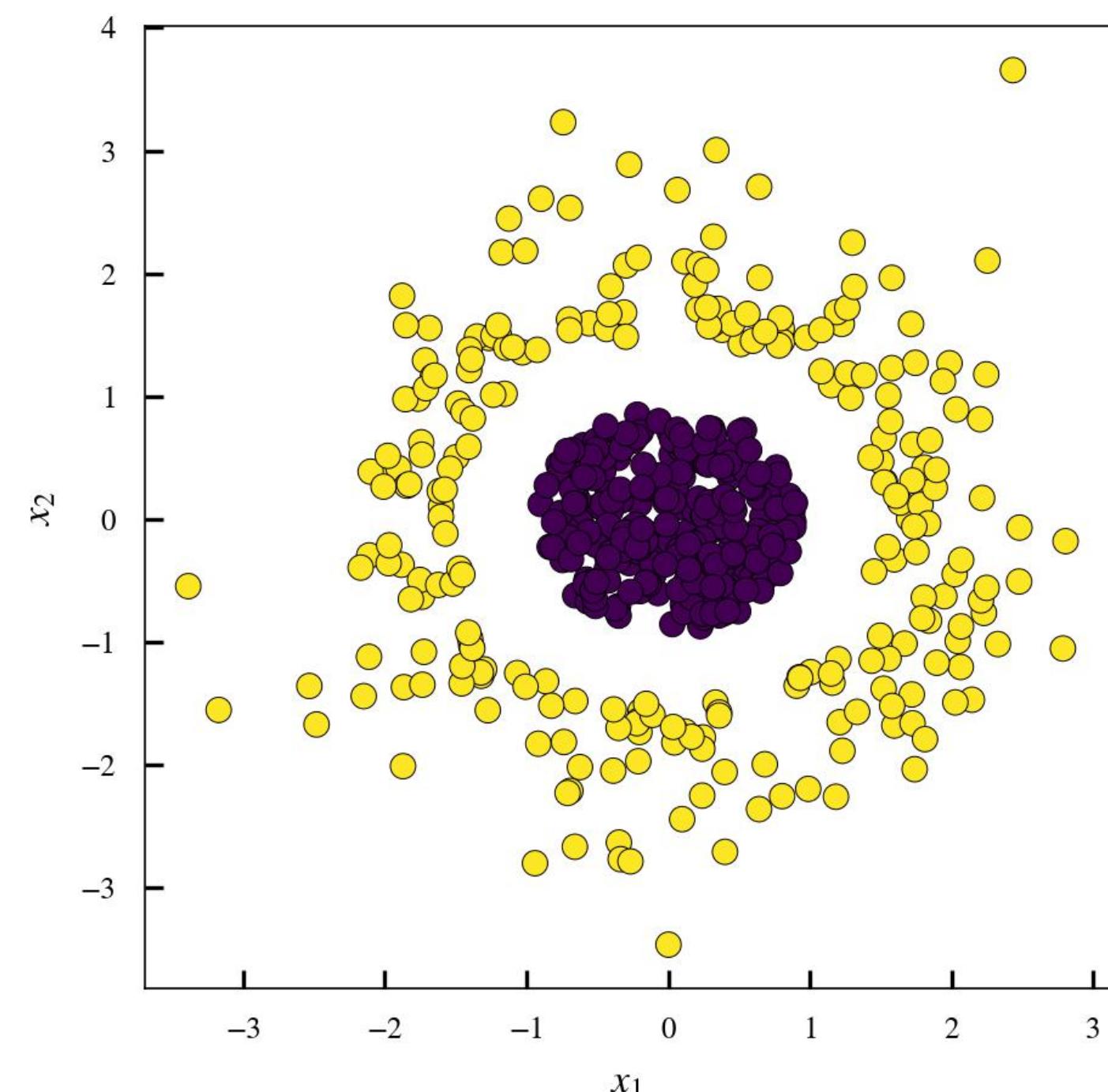
Credits to Lorenz Vaitl  
for slides

# Motivation

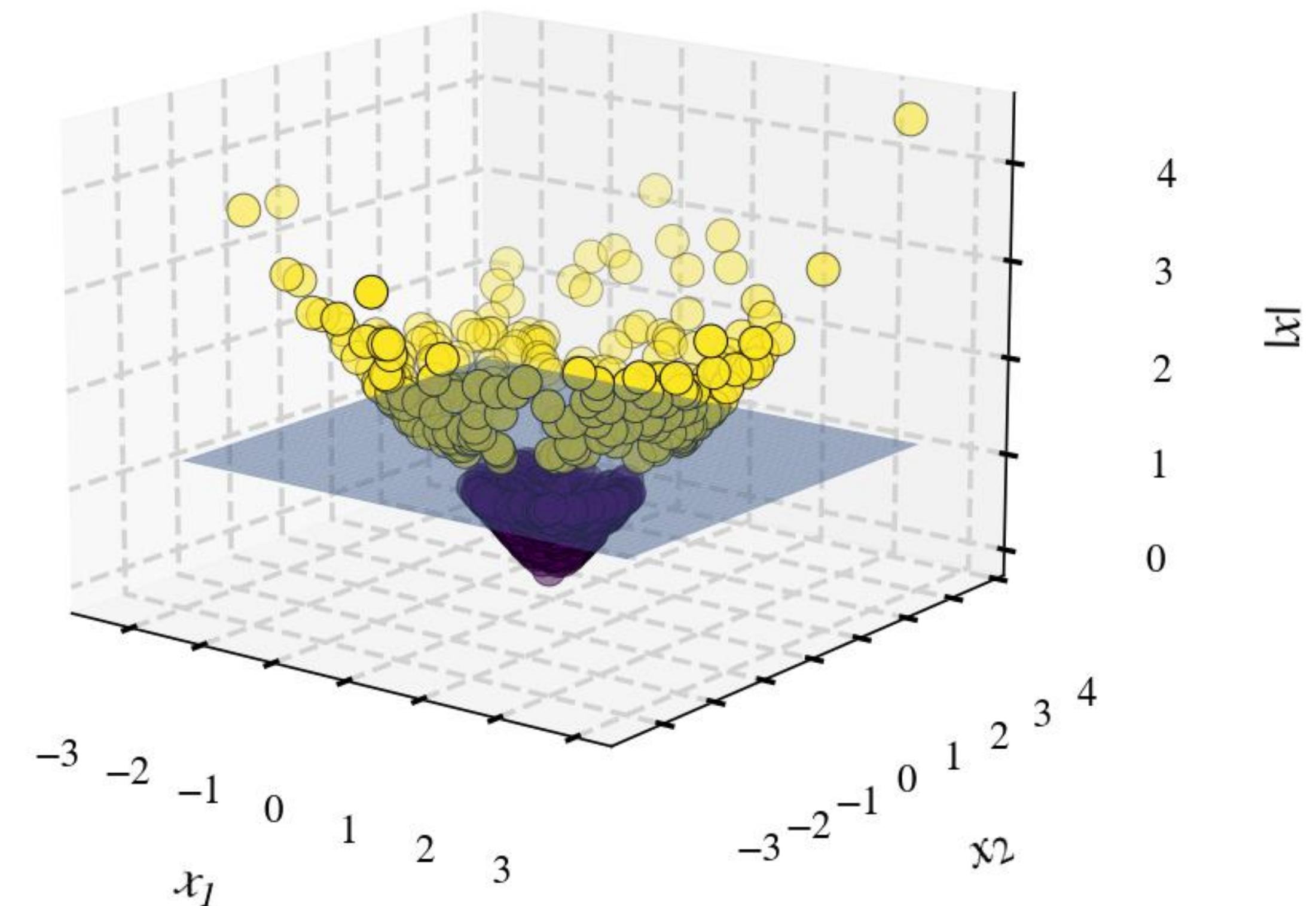
From linear to non-linear models

# Better features for linear models

Binary classification problem



Polynomial kernel



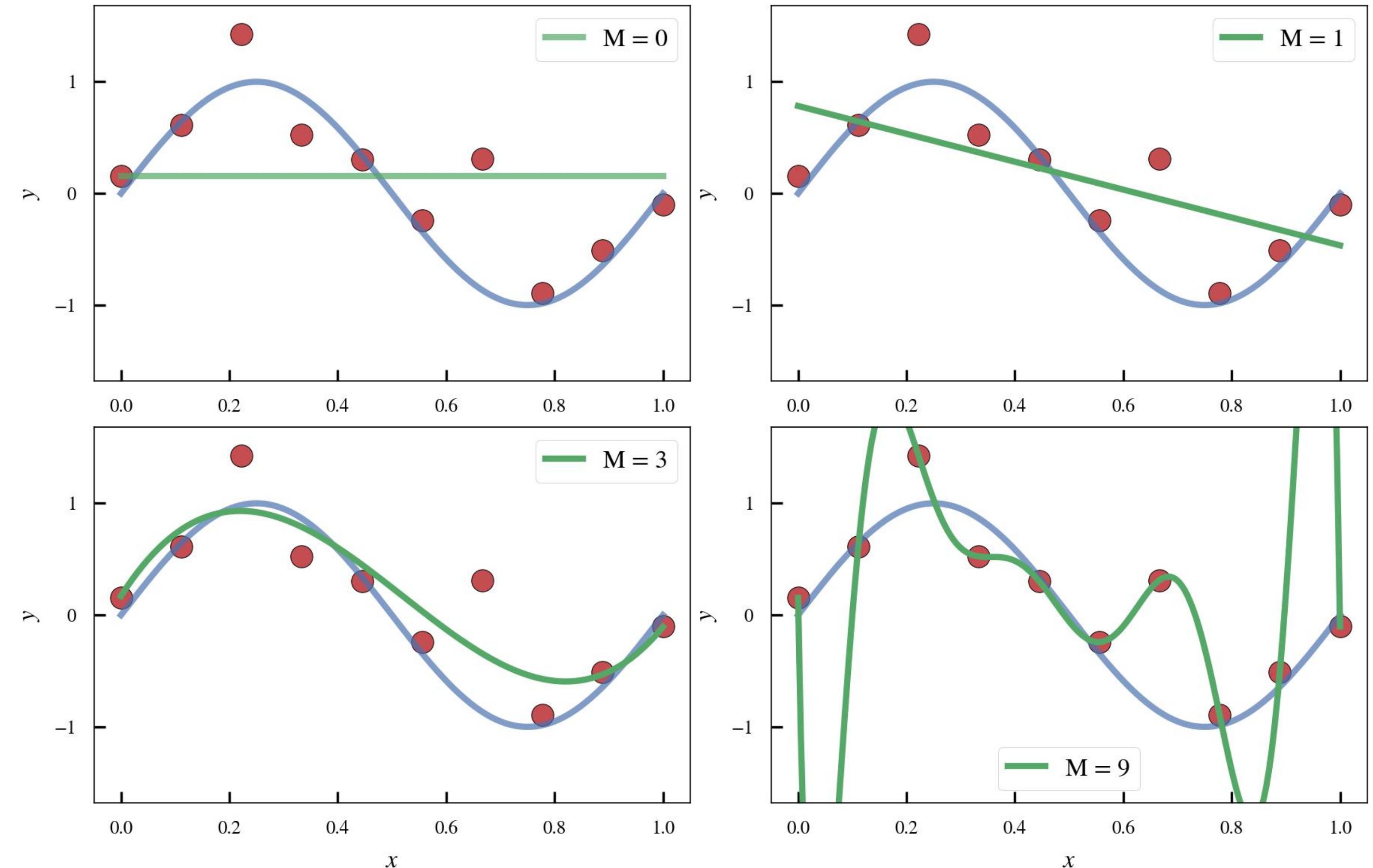
Kernel describing degree- $d$  polynomial  
 $k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^\top \mathbf{x}_2 + c)^d$

# Linear regression with polynomial features

$$\hat{y}(x) = w_0 + w_1 \cdot x^1 + \dots + w_M \cdot x^M$$

→ Here we are using the basis function

$\phi_M(x) = (x^0, x^1, \dots, x^M)$ ,  
which has a  $(M + 1)$ -dimensional  
feature space  $\mathcal{F}$

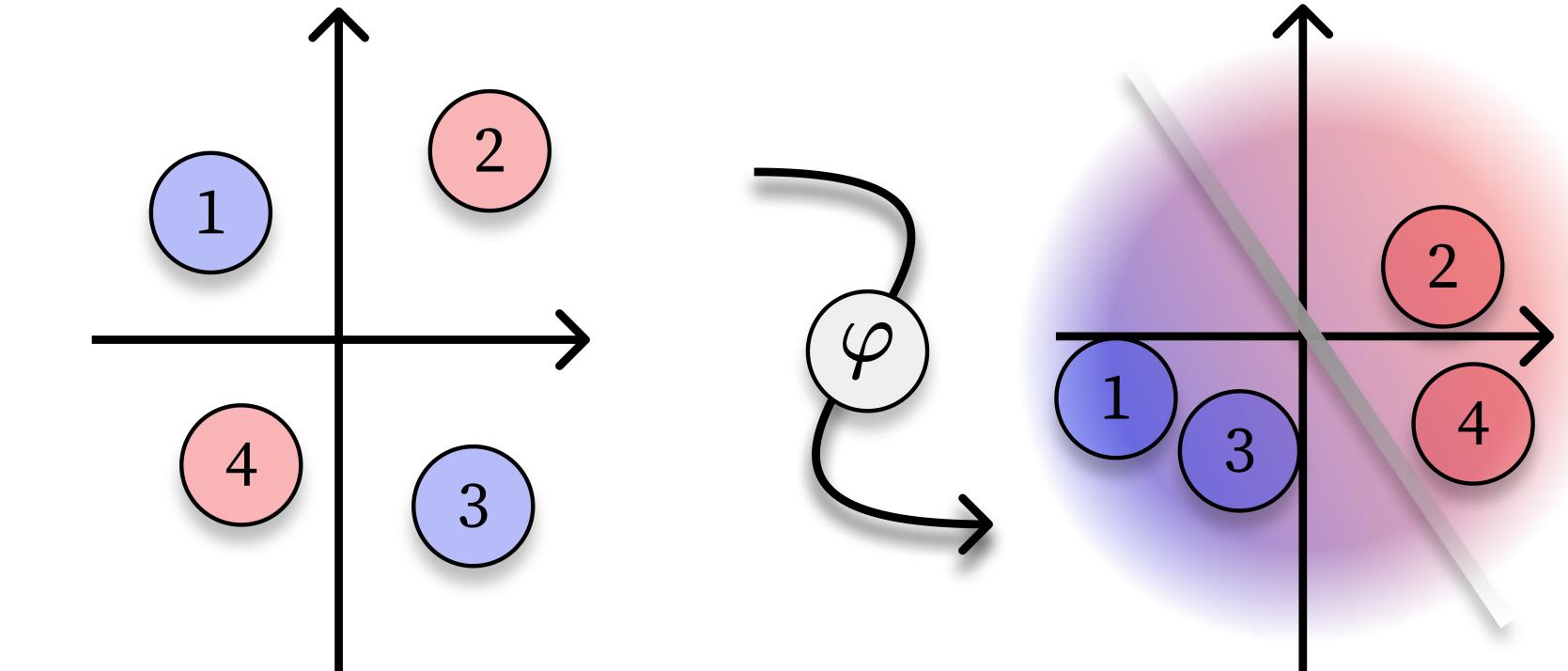
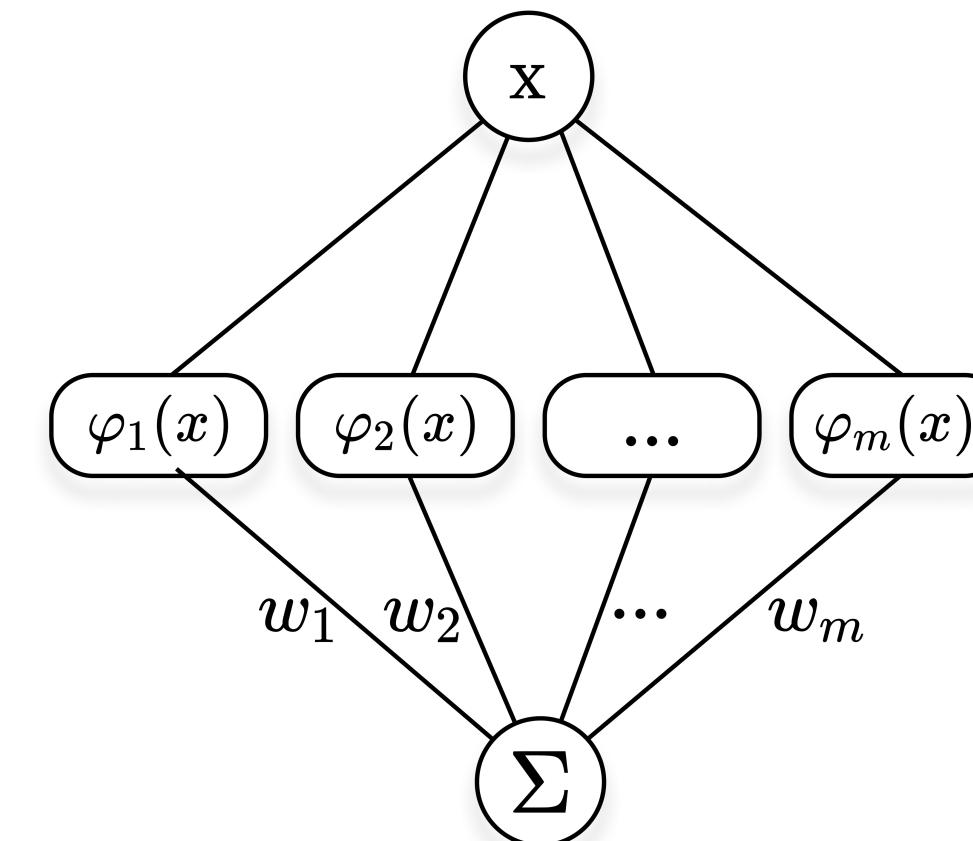


[adapted from Bishop, 2007]

# **Kernel ridge regression**

**Kernel trick in practise**

# Kernel trick



[adapted from Jäkel et al., 2009]

1. Map the data into a (high dimensional) feature space,  $x \mapsto \varphi(x)$
2. Look for linear relations in the feature space

We call this form **dual representation**

Only scalar products  $\rightarrow$  we can put in kernels:

$$x_i^\top x_j \rightarrow \varphi(x_i)^\top \varphi(x_j) = k(x_i, x_j) = K_{ij}$$

# Ridge regression

$$\hat{y}(x) = w_0 + w_1 \cdot x^1 + \dots + w_M \cdot x^M$$

$$\begin{matrix} N \rightarrow \\ \textcolor{red}{\hat{y}} \end{matrix} = \begin{matrix} D \rightarrow \\ \textcolor{gray}{w}^T \end{matrix} \cdot \begin{matrix} N \rightarrow \\ \textcolor{blue}{X} \end{matrix}$$

Linear ridge regression looks for a linear combination of features  $w$  that minimizes the prediction error and has a small norm

We can write this term in several equivalent ways:

$$\begin{aligned}\mathcal{E}_{RR}(w) &= \sum_n (y_n - w^\top x_n)^2 + \lambda \sum_d w_d^2 \\ &= \|\mathbf{y} - \mathbf{w}^\top \mathbf{X}\|^2 + \lambda \|\mathbf{w}\|^2 \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X} \mathbf{y}^\top + \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w} + \lambda \mathbf{w}^\top \mathbf{w}\end{aligned}$$

# Ridge regression

$$\begin{array}{c} N \rightarrow \\ \hat{\mathbf{y}} \end{array} = \begin{array}{c} D \rightarrow \\ \mathbf{w}^T \end{array} \cdot \begin{array}{c} N \rightarrow \\ \mathbf{x} \end{array}$$

$$\mathcal{E}_{RR}(\mathbf{w}) = \mathbf{y}\mathbf{y}^T - 2\mathbf{w}^T \mathbf{X} \mathbf{y}^T + \mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w} + \lambda \mathbf{w}^T \mathbf{w}$$

Computing the derivative w.r.t.  $\mathbf{w}$  yields

$$\frac{\partial \mathcal{E}_{RR}(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X} \mathbf{y}^T + 2\mathbf{X} \mathbf{X}^T \mathbf{w} + 2\lambda \mathbf{w}$$

Setting the gradient to zero and rearranging terms the optimal  $\mathbf{w}$  is

$$2\mathbf{X} \mathbf{X}^T \mathbf{w} + \lambda \mathbf{w} = 2\mathbf{X} \mathbf{y}^T$$

$$(\mathbf{X} \mathbf{X}^T + \lambda I) \mathbf{w} = \mathbf{X} \mathbf{y}^T$$

$$\mathbf{w} = (\mathbf{X} \mathbf{X}^T + \lambda I)^{-1} \mathbf{X} \mathbf{y}^T$$

D<sup>D</sup>  
Compute  
inverse in  
feature space

# Kernelize ridge regression

Setting the gradient to zero and rearranging terms the optimal  $\mathbf{w}$  is

$$2\mathbf{X}\mathbf{X}^\top \mathbf{w} + \lambda 2\mathbf{w} = 2\mathbf{X}\mathbf{y}^\top$$

$$(\mathbf{X}\mathbf{X}^\top + \lambda I)\mathbf{w} = \mathbf{X}\mathbf{y}^\top$$

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1} \mathbf{X}\mathbf{y}^\top$$

 Compute inverse in feature space

Setting the gradient to 0 and rearranging terms the optimal  $\mathbf{w}$  satisfies

$$\mathbf{w} = \underbrace{\mathbf{X} \frac{1}{\lambda} (\mathbf{y}^\top - \mathbf{X}^\top \mathbf{w})}_{:= \boldsymbol{\alpha} (\in \mathbb{R}^{n \times 1})} = \sum_i^n \alpha_i \mathbf{x}_i$$

Linear combination of samples

# Kernelize ridge regression

$$\mathcal{E}_{RR}(\mathbf{w}) = \mathbf{y}\mathbf{y}^\top - 2\mathbf{w}^\top \mathbf{X}\mathbf{y}^\top + \mathbf{w}^\top \mathbf{X}\mathbf{X}^\top \mathbf{w} + \lambda \mathbf{w}^\top \mathbf{w}$$

$$\mathbf{w} = \mathbf{X}\boldsymbol{\alpha}$$

Let's look at the error function  $\mathcal{E}_{RR}$  w.r.t.  $\boldsymbol{\alpha}$

$$\mathcal{E}_{RR}(\boldsymbol{\alpha}) = \mathbf{y}\mathbf{y}^\top - 2\boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_K \mathbf{y}^\top + \boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_K \underbrace{\mathbf{X}\mathbf{X}^\top}_K \mathbf{X} \boldsymbol{\alpha} + \lambda \boldsymbol{\alpha}^\top \underbrace{\mathbf{X}^\top \mathbf{X}}_K \boldsymbol{\alpha}$$

Now we can optimize for the new parameters

$$\frac{\partial}{\partial \boldsymbol{\alpha}} \mathcal{E}_{RR}(\boldsymbol{\alpha}) = 0$$

$$\boldsymbol{\alpha} = (K + \lambda I_N)^{-1} K \mathbf{y}$$



Compute  
inverse in  
data space

# Kernelize ridge regression

$$\mathcal{E}_{RR}(\alpha) = \mathbf{y}\mathbf{y}^\top - 2\alpha^\top \underbrace{\mathbf{X}^T \mathbf{X}}_K \mathbf{y}^\top + \alpha^\top \underbrace{\mathbf{X}^T \mathbf{X}}_K \underbrace{\mathbf{X} \mathbf{X}^\top}_K \mathbf{X} \alpha + \lambda \alpha^\top \underbrace{\mathbf{X}^T \mathbf{X}}_K \alpha$$

We call this form **dual representation**

Only scalar products  $\rightarrow$  we can put in kernels:

$$\mathbf{x}_i^\top \mathbf{x}_j \rightarrow \varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) = K_{ij}$$

We write  $k(X, X)$  as  $K$

# Compare training strategies

**Kernel Ridge Regression:**

$$\alpha = \frac{1}{\lambda} (\mathbf{y}^\top - \varphi(\mathbf{X}_{train})^\top \mathbf{w})$$

$$\lambda \alpha = \mathbf{y}^\top - \varphi(\mathbf{X}_{train})^\top \varphi(\mathbf{X}_{train}) \alpha$$

$$\mathbf{y}^\top = (\varphi(\mathbf{X}_{train})^\top \varphi(\mathbf{X}_{train}) + \lambda I) \alpha$$

$$\alpha = (\varphi(\mathbf{X}_{train})^\top \varphi(\mathbf{X}_{train}) + \lambda I)^{-1} \mathbf{y}^\top$$

$$\alpha = (K + \lambda I)^{-1} \mathbf{y}^\top$$

This  $\alpha$  minimizes  $\mathcal{E}_{RR}(\alpha)$ .

**Ridge Regression**

Train  $\mathbf{w}$  which minimizes  $\mathcal{E}_{RR}(\mathbf{w})$ :

$$\mathbf{w} = (\varphi(\mathbf{X}) \varphi(\mathbf{X})^\top + \lambda I)^{-1} \varphi(\mathbf{X}) \mathbf{y}^\top$$

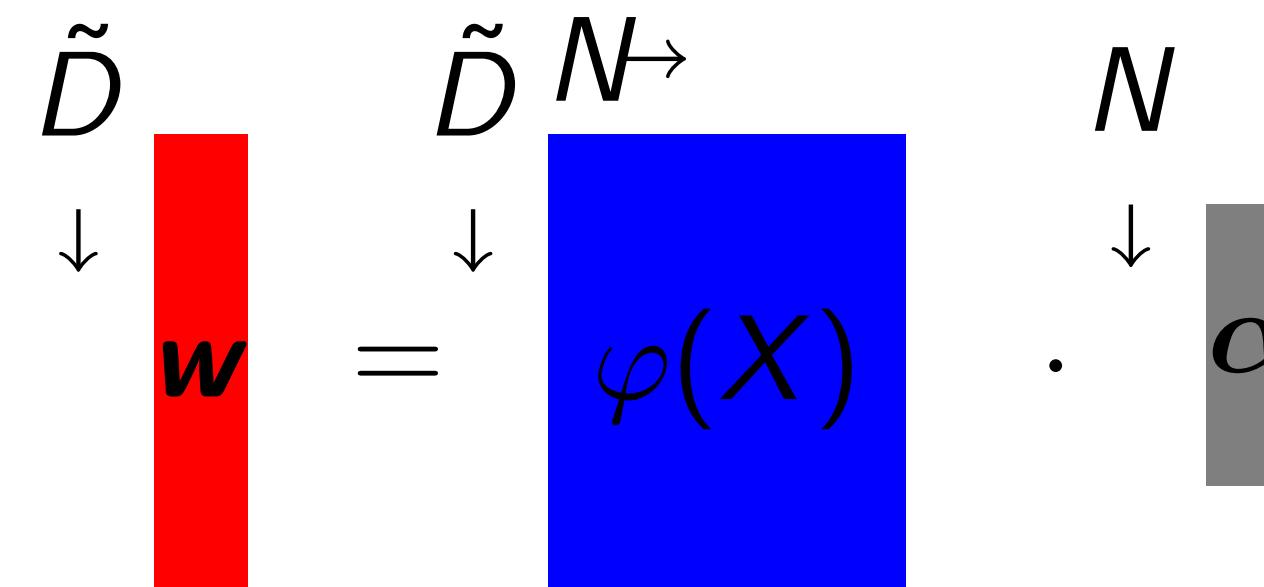
For KRR we write  $\varphi(\mathbf{x}) \in \mathbb{R}^{\tilde{D}}$  instead of  $\mathbf{x}$

We defined  $\alpha_i = \frac{1}{\lambda} (y_i - \varphi(\mathbf{x}_i)^\top \mathbf{w})$

Optimal  $\mathbf{w} = \varphi(\mathbf{X}_{train}) \alpha$

KRR trains  $\alpha \in \mathbb{R}^n$ , RR trains  $\mathbf{w} \in \mathbb{R}^{\tilde{D}}$

For  $\varphi(\mathbf{x})^\top \varphi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$  the two models are equivalent (but not the runtime complexity)



# Predictions for new data

$$\begin{aligned}y_{new} &= \mathbf{w}^T \varphi(\mathbf{x}_{new}) \\&= (\varphi(\mathcal{X}_{train})\boldsymbol{\alpha})^T \varphi(\mathbf{x}_{new}) \\&= \boldsymbol{\alpha}^T \varphi(\mathcal{X}_{train})^T \varphi(\mathbf{x}_{new}) \\&= \boldsymbol{\alpha}^T k(\mathcal{X}_{train}, \mathbf{x}_{new}) \\&= \mathbf{y}_{train}(K + \lambda I)^{-1} k(\mathcal{X}_{train}, \mathbf{x}_{new})\end{aligned}$$

Optimal  $\boldsymbol{\alpha} = (K + \lambda I)^{-1} \mathbf{y}^T$

$\mathbf{w} = \varphi(\mathcal{X}_{train})\boldsymbol{\alpha}$

$K(a, b) = K(b, a)^T$

We call  $K = k(\mathcal{X}_{train}, \mathcal{X}_{train})$   
the Gram matrix

# Summary

## Kernel ridge regression

- Input: kernel function  $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}')$ , regularization hyperparameter  $\lambda$ , training dataset  $(\mathbf{X}_{train}, \mathbf{y}_{train})$  and test datapoints  $\mathbf{X}_{test}$

- Training <sup>2</sup>:

$$\boldsymbol{\alpha} = (k(\mathbf{X}_{train}, \mathbf{X}_{train}) + \lambda I)^{-1} \mathbf{y}_{train}^T$$

- Predicting:

$$\hat{\mathbf{y}}_{test} = \boldsymbol{\alpha}^T k(\mathbf{X}_{train}, \mathbf{X}_{test})$$

---

<sup>2</sup>Since we need  $\mathbf{X}_{train}$  again for predictions, we cannot forget the data like for RR

# **Theory**

**Why all this works**

# Constructing a kernel from features

## Definition

We define a Kernel function  $k$  as a dot product in feature space  $\mathcal{F}$  where

$$\begin{aligned}\varphi : \mathcal{X} &\rightarrow \mathcal{F} \\ x &\mapsto \varphi(x)\end{aligned}$$

and so

$$\begin{aligned}k : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ k(x_i, x_j) &= \varphi(x_i)^T \cdot \varphi(x_j)\end{aligned}$$

## Kernel Trick (Kernel Substitution)

For any algorithm that can be formulated such that the input vector  $\mathbf{x}$  enters only in terms of scalar products  $\mathbf{x}^T \cdot \mathbf{x}'$ , we can replace each scalar product by a kernel  $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \cdot \varphi(\mathbf{x}')$ .

Why do it?

For  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{d \times 1} \Rightarrow \mathbf{x}^T \cdot \mathbf{x}' \in \mathbb{R}^{1 \times 1}$  regardless of  $d$

For  $\varphi(\mathbf{x}) \in \mathbb{R}^{\tilde{d} \times 1}$  instead of  $\mathbf{x}$  with typically  $\tilde{d} \gg d$

$\Rightarrow$  we get more complex (powerful) models

By using kernel  $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^{1 \times 1}$

We do not need to explicitly calculate high-dimensional  $\varphi(\mathbf{x})$

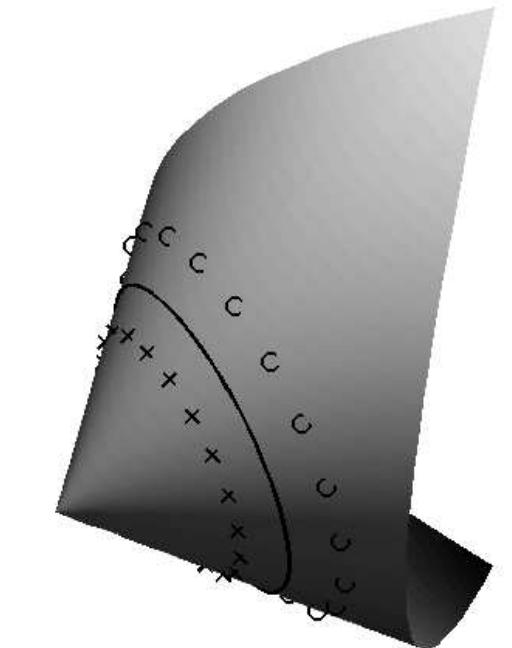
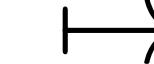
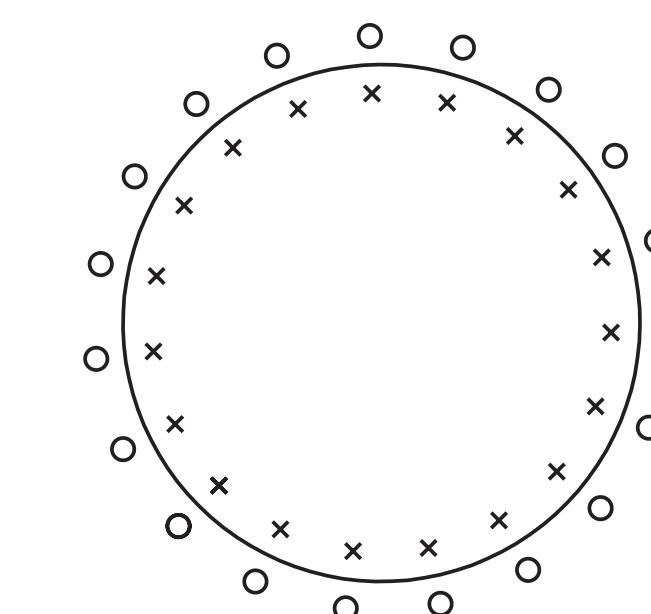
# Example for a kernel

$$\varphi : \mathbf{x} = (x_1, x_2)^\top \mapsto \varphi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

The corresponding kernel:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \varphi(\mathbf{x})^\top \cdot \varphi(\mathbf{y}) \\ &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2)^\top \\ &= x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \\ &= (x_1y_1 + x_2y_2)^2 = (\mathbf{x}^\top \mathbf{y})^2 \end{aligned}$$

Visualizing  $\varphi(\mathbf{x})$



→ With  $k(\cdot, \cdot)$  we implicitly work in  $\mathbb{R}^3$ , but only operate in  $\mathbb{R}^2$

## Definition (Positive semi-definite symmetric kernels)

A kernel

$$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

is said to be positive semi-definite symmetric

if for any  $\{x_1, \dots, x_n\} \subseteq \mathcal{X}$ , the matrix  $K = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{n \times n}$  is symmetric positive semidefinite.<sup>1</sup>

*Symmetric Matrix:*

$$A = A^\top$$

*Positive semi-definite Matrix:*

$$x^\top A x \geq 0 \quad \forall x$$

alternatively, all eigenvalues  $\lambda$  of  $A$  non-negative.

---

<sup>1</sup>For ease of notation we may use  $K(X, X') = [K(x_i, x'_j)]_{ij}$  for describing the matrix of the kernel-function evaluated on all sample-pairs.  $K(X, X)$  is often called the *Gram matrix* of  $X$ .

## Definition

We define a Kernel function  $k$  as a dot product in feature space  $\mathcal{F}$  where

$$\varphi : \mathcal{X} \rightarrow \mathcal{F}$$

$$x \mapsto \varphi(x)$$

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

$$k(x_i, x_j) = \varphi(x_i)^T \cdot \varphi(x_j)$$

## Definition (Positive semi-definite symmetric kernels)

if for any  $\{x_1, \dots, x_n\} \subseteq \mathcal{X}$ , the matrix  $K = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{n \times n}$  is symmetric positive semidefinite.<sup>1</sup>

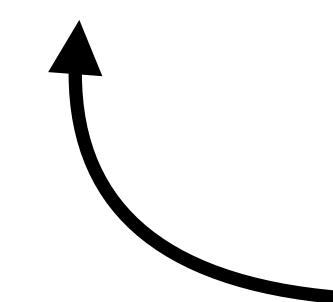
$$K = \kappa(X, X) = \varphi(X)^T \cdot \varphi(X)$$

Feature space

$$= (U\Sigma V^T)^T \cdot U\Sigma V^T = V\Sigma (U^T U) \Sigma V^T$$

Singular value decomposition

$$= V\Sigma^2 V^T$$



Positive spectrum

## Mercer's Theorem [Mercer, 1909]

*Non-technical formulation:*

If kernel is positive semi-definite symmetric then one can find a map  $\varphi$  such that  $k(x, x') = \varphi(x)^T \varphi(x')$  for almost all  $x \in \mathcal{X}$

- To see if your kernel is valid, show it is positive semi-definite symmetric!
- You can construct kernels from other kernels, e.g. by sum, product or exponentiating

Alternatively show  $k(x, x') = \varphi(x)^T \varphi(x')$

# Popular kernel functions

Linear Kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

Polynomial Kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + c)^p$$

Gaussian Kernel/ RBF (more on this later)

$$k(\mathbf{x}_i, \mathbf{x}_j) = e^{\|\mathbf{x}_i - \mathbf{x}_j\|^2 / -2\sigma^2}$$

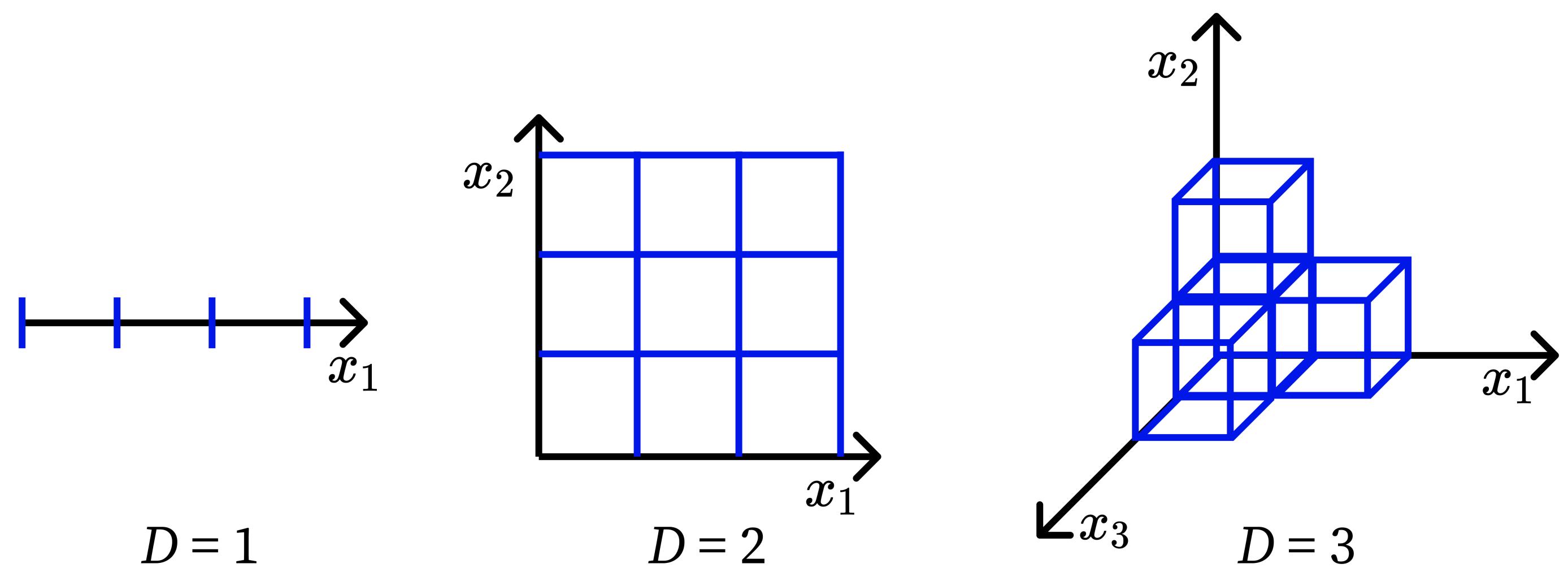
Note that we are never directly operating in feature space

# Curse of dimensionality

A big problem with high dimensional features spaces

When the dimensionality increases, the volume of the space increases so fast that the available data becomes sparse

The amount of data needed for a reliable result often grows exponentially with the dimensionality



[Adapted from Bishop, 2007]

# Why should this work?

How do we find the optimal weight

$$\mathbf{w} \in \mathbb{R}^{\tilde{d}} ?$$

Seems impossible with  $\tilde{d} \gg n$



# Why it works!

Representer Theorem [Kimeldorf and Wahba, 1971]

The minimizing function  $f^*$  of a (regularized) error function on some training data  $\mathbf{x}_i$  can be written in terms of the kernel  $k$  as

$$f^*(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}_i).$$

I.e. only scalar products between  $\varphi(\mathbf{x})$  and the training data  $\varphi(\mathbf{x}_i)$  are relevant.  
For the model  $f(\mathbf{x}) = \mathbf{w}^\top \varphi(\mathbf{x})$  this means that  $\mathbf{w}$  can be written as

$$\mathbf{w} = \sum_{i=1}^N \alpha_i \varphi(\mathbf{x}_i).$$

# Kernels as similarity measures

$$f^*(\mathbf{x}_{\text{new}}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}_{\text{new}}, \mathbf{x}_i)$$

Kernel methods are *memory-based* methods:

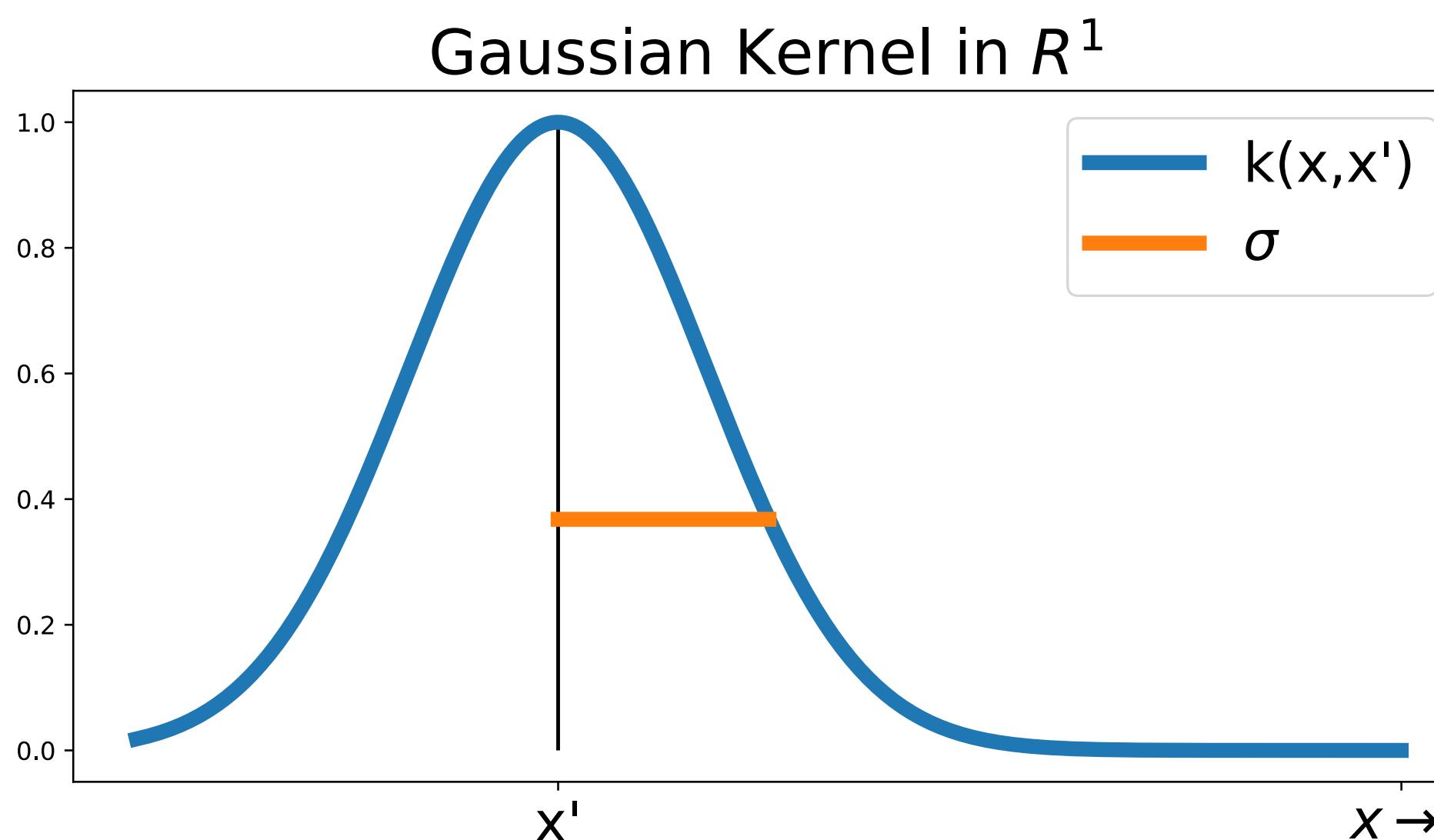
- store the entire training set
- define similarity of data points by kernel function

$$k(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

- new predictions require comparison with previously learned examples

# Law of universal generalization

- Roger Shepard argued that the **Perceptual Similarity** of new data  $x$  decays exponentially with distance from prototype  $x'$  [Shepard, 1987]
  - motivation for using the Gaussian Kernel



Gaussian Kernel (RBF Kernel)

$$k(x', x) = e^{\frac{-(x'-x)^2}{2\sigma^2}}$$

Because  $\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$ , the basis function  $\varphi(x)$  is infinite dimensional

# Kernels as similarity measures

**Kernel Ridge Regression:**

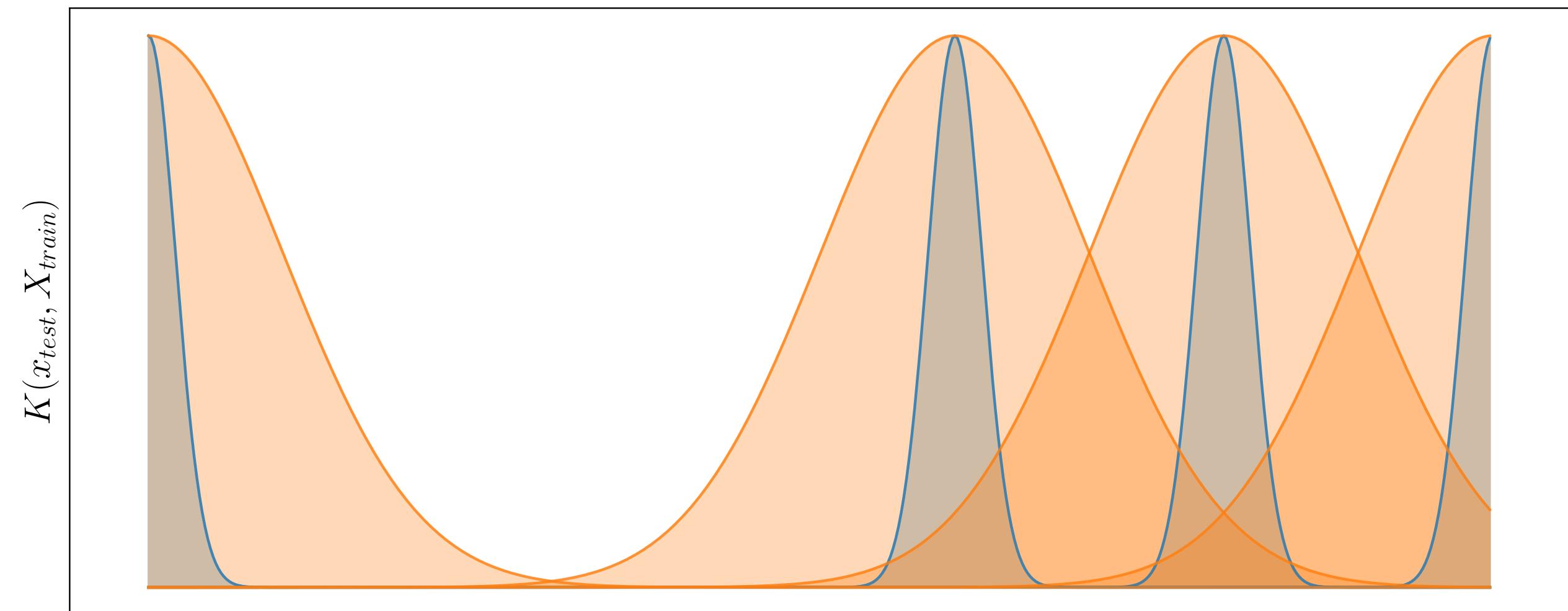
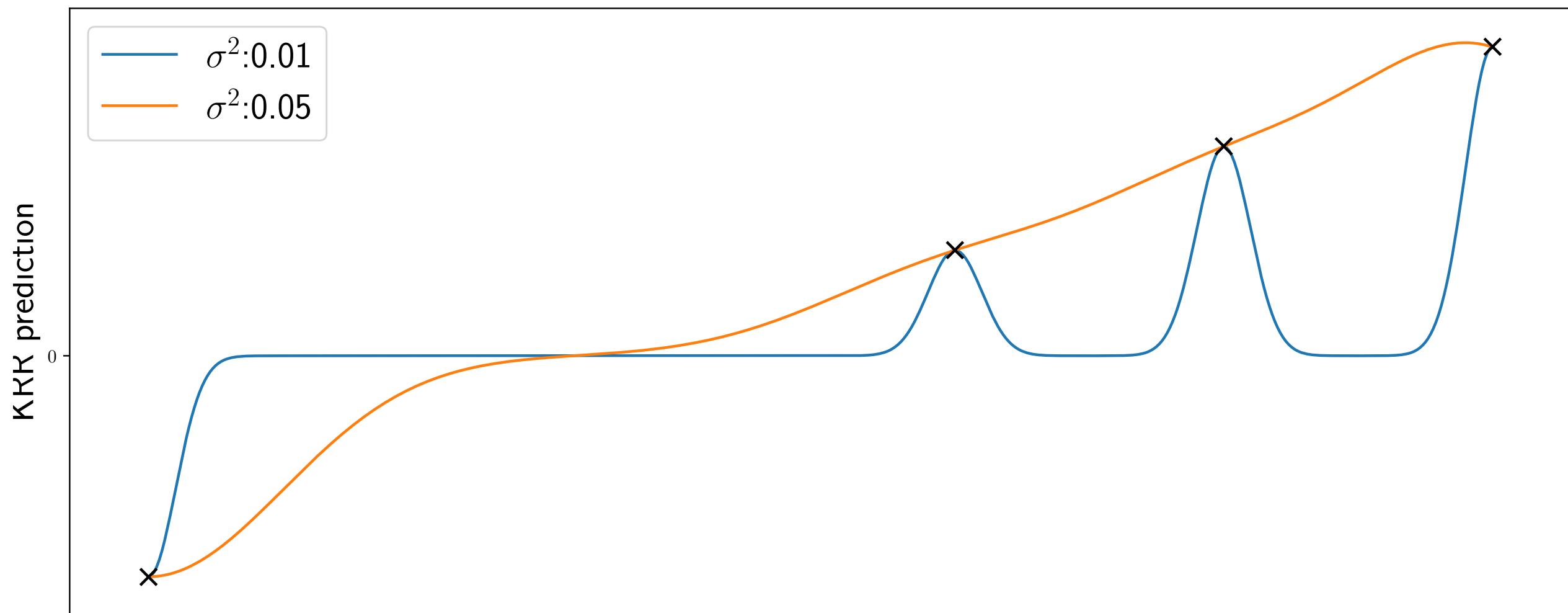
**Example**

KRR with Gaussian Kernels:

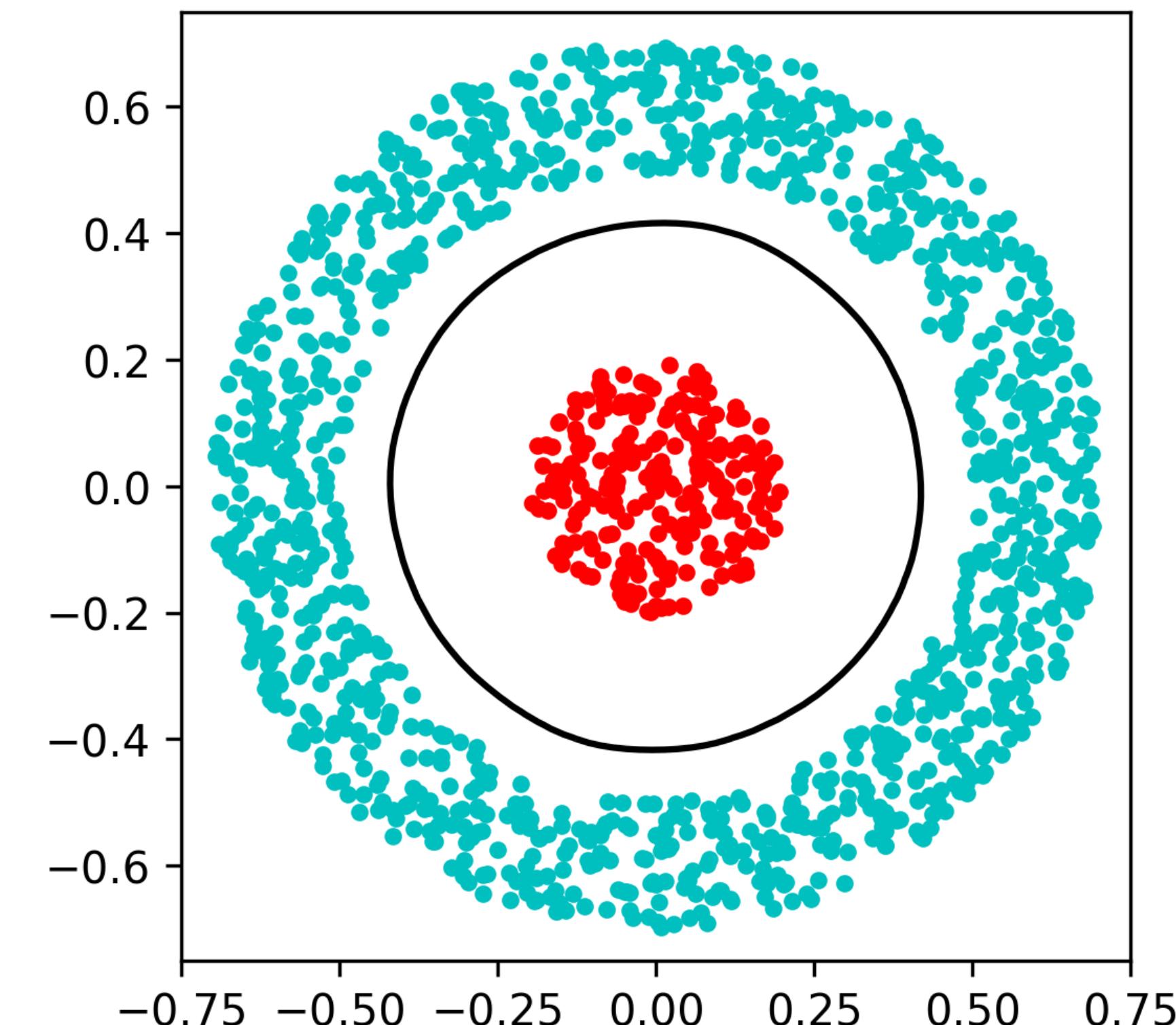
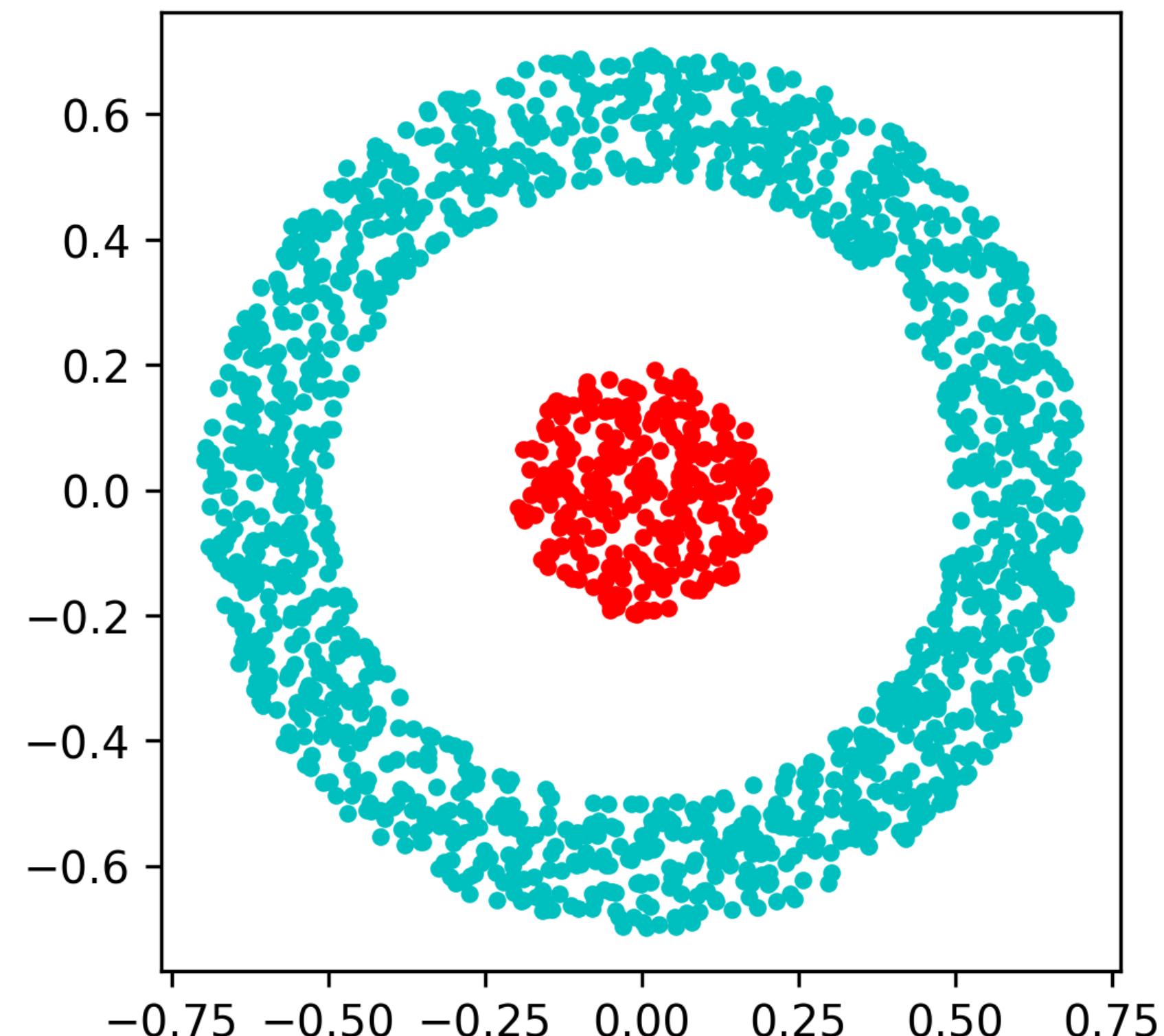
$$k(x, x') = \exp \left\{ -\frac{\|x - x'\|^2}{2\sigma^2} \right\}$$

Predictions:

$$y_{new} = y_{train}(K + \lambda I)^{-1} k(X_{train}, x_{new})$$



We can also kernelize LDA with e.g. a Gaussian Kernel ( $\sigma = 0.1$ )



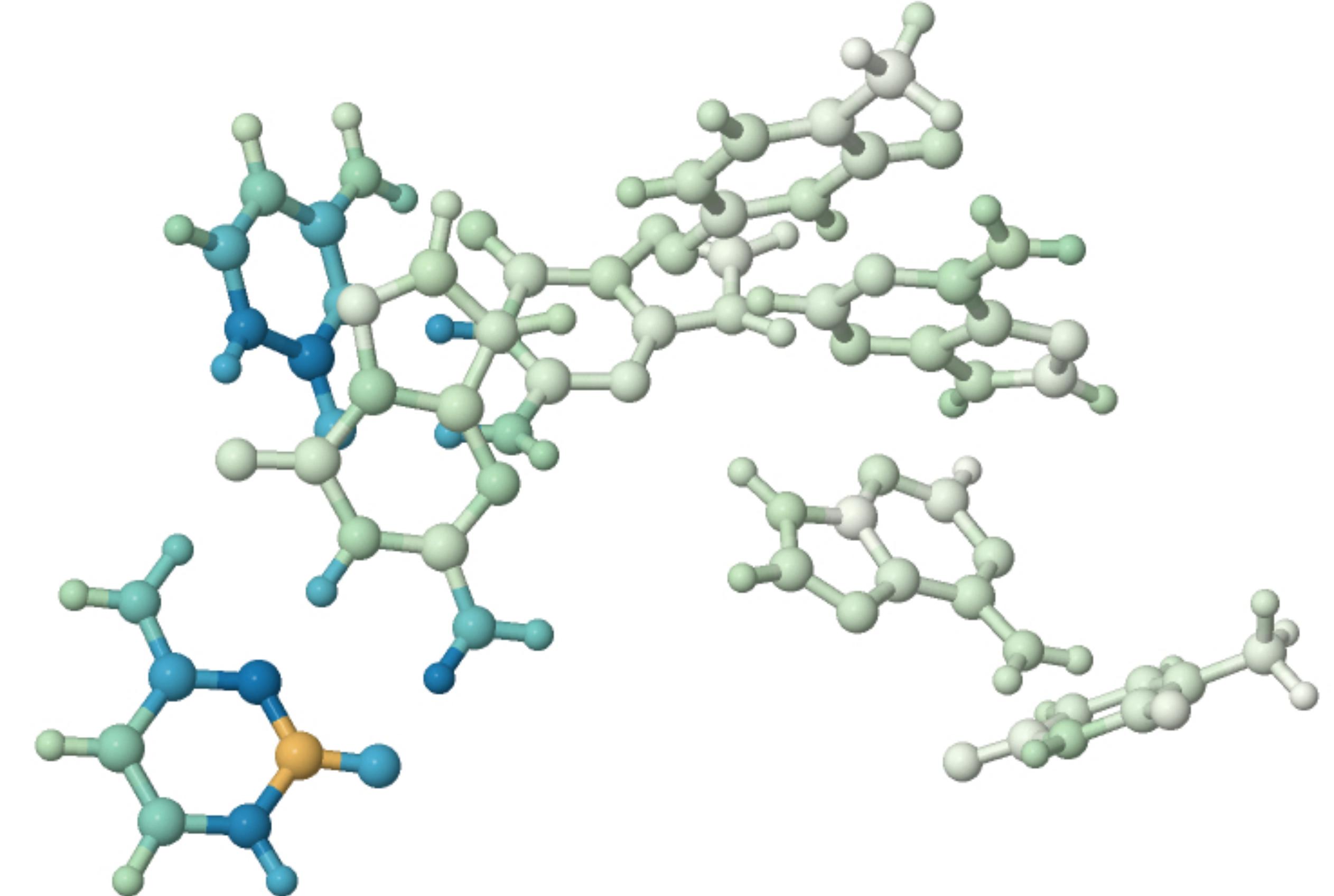
→ We are able to do non-linear classification!

# **Practical examples**

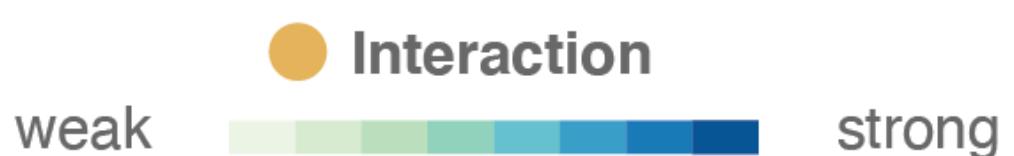
**[Gaussian processes], NTK and sGDML**



# Machine learning force fields with symmetrized kernels

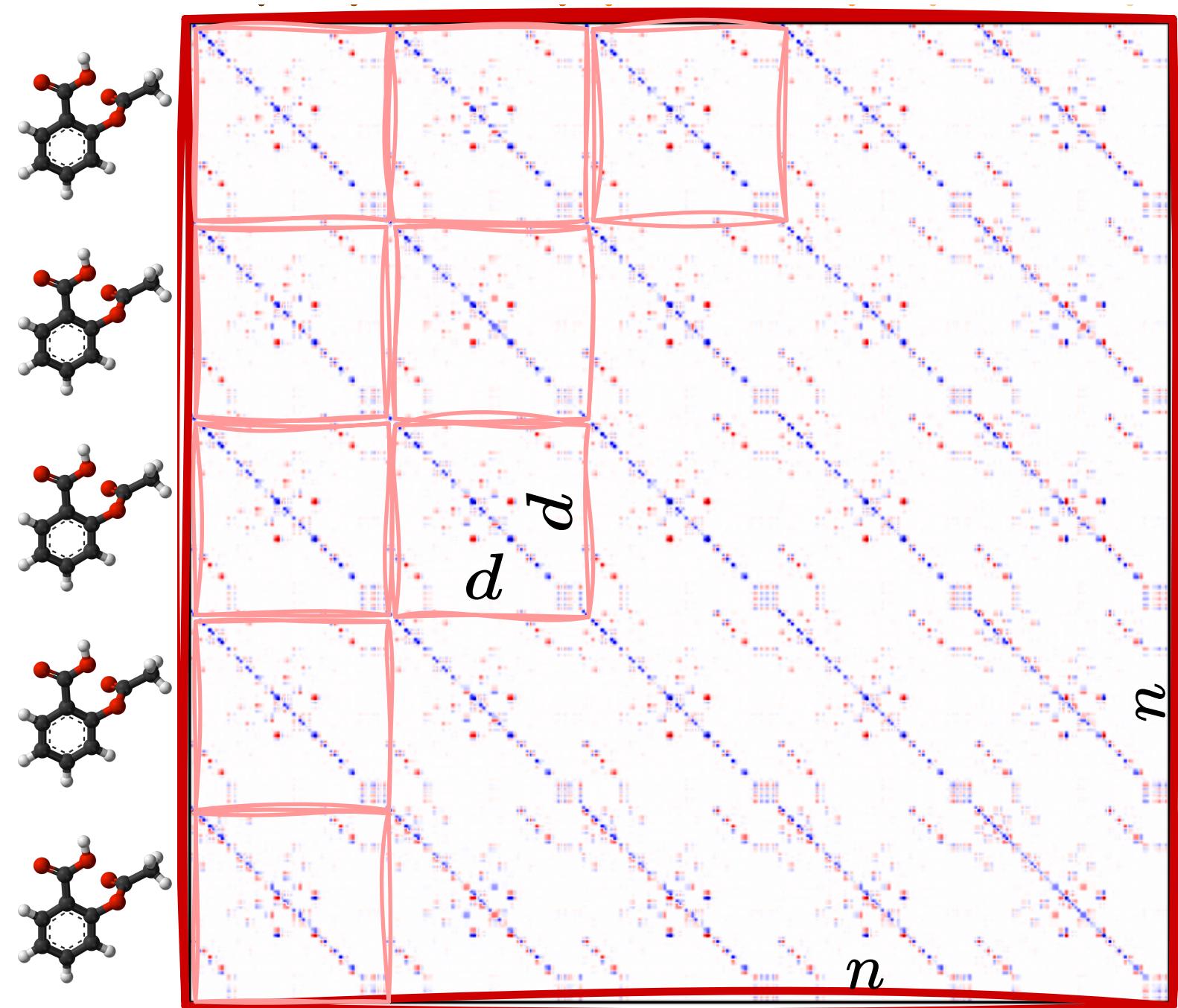


All atoms interact in quantum  
many-body systems.

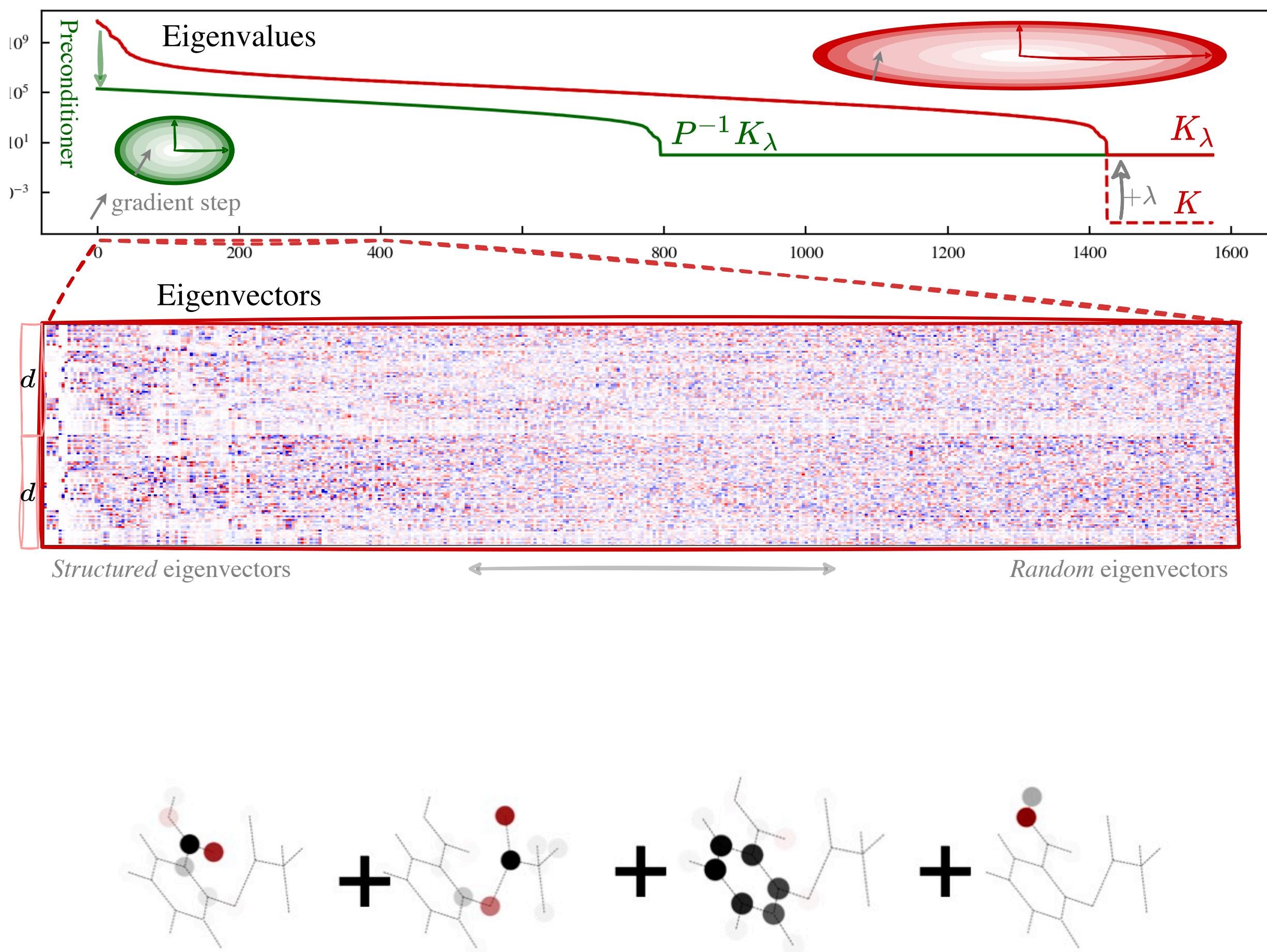


# Singular value decomposition

Kernel between molecules

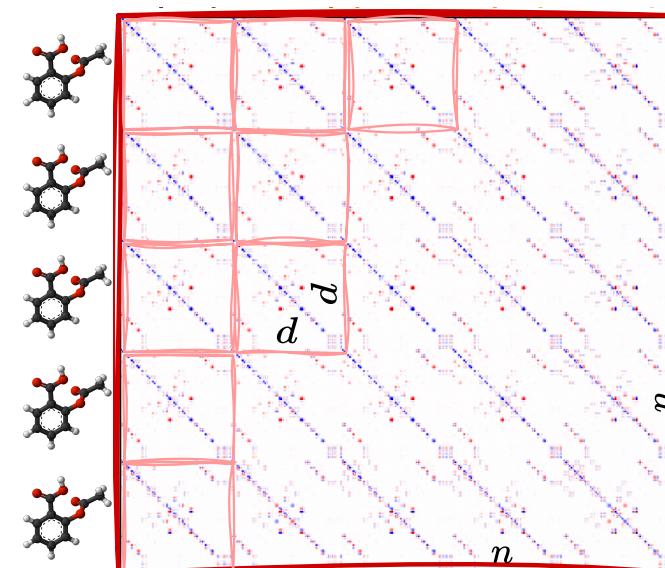


Dynamics  
(positions)



**$k$  prototype molecules**

Kernel matrix



Training solution

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

Invert kernel matrix

Cholesky decomposition

VS.

Iterative solver

Closed form

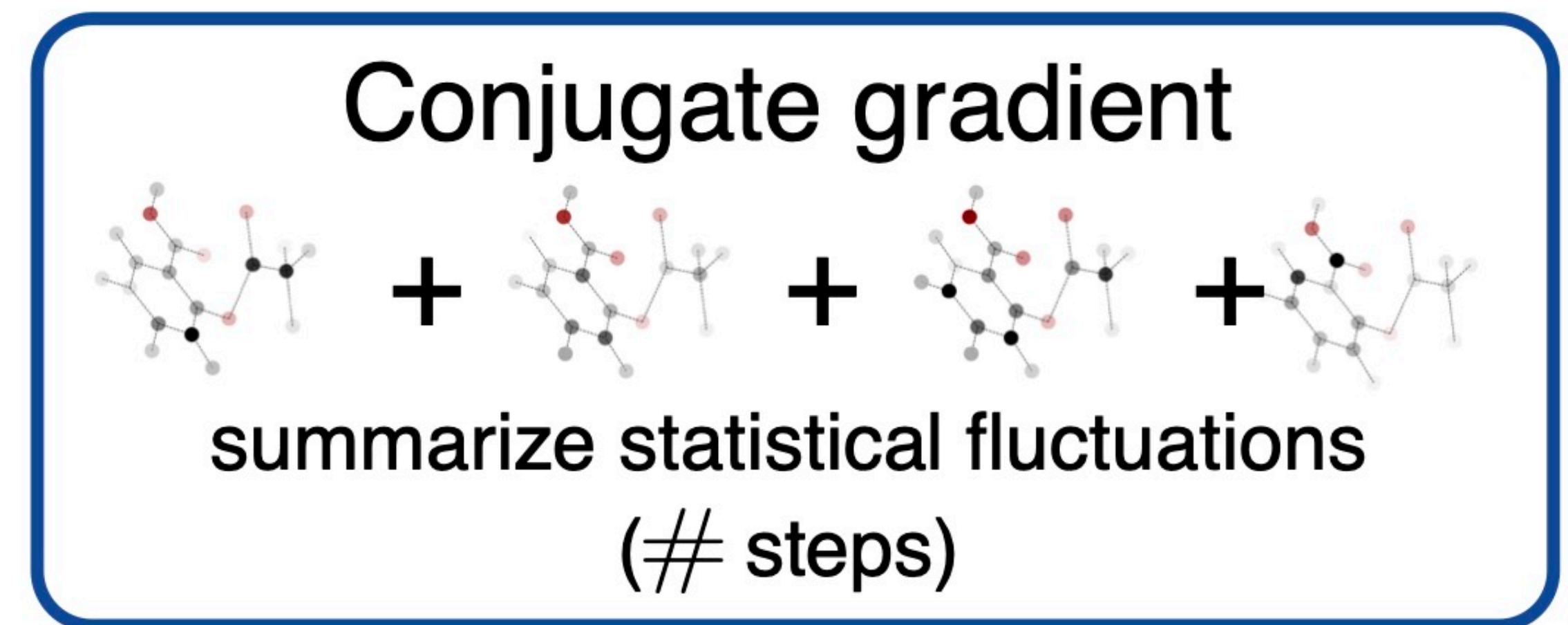
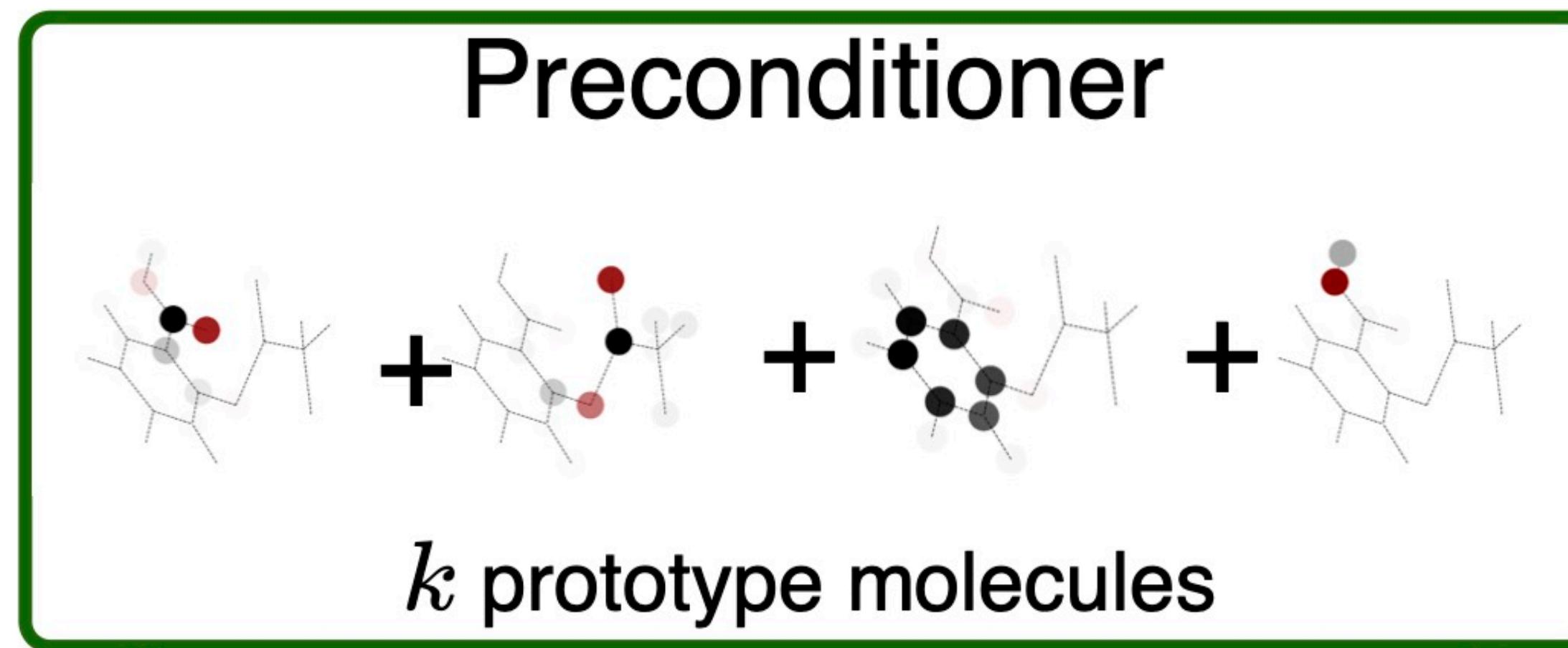
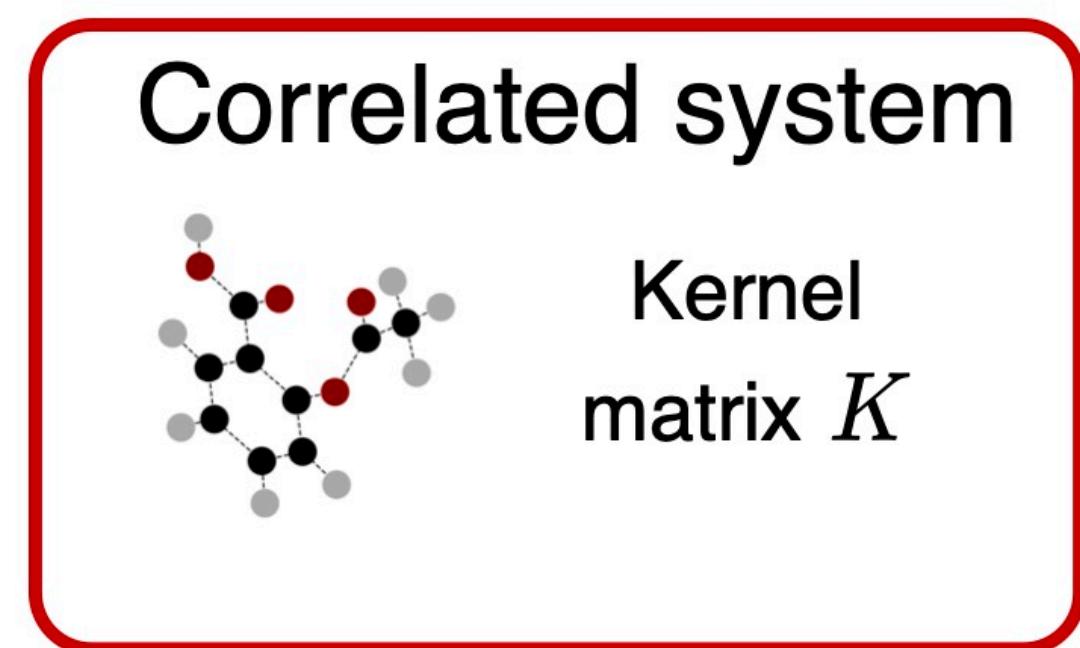
$$\mathbf{K} = \mathbf{L}^T \times \mathbf{L}$$

Requires storing full matrix

Convergence on problem

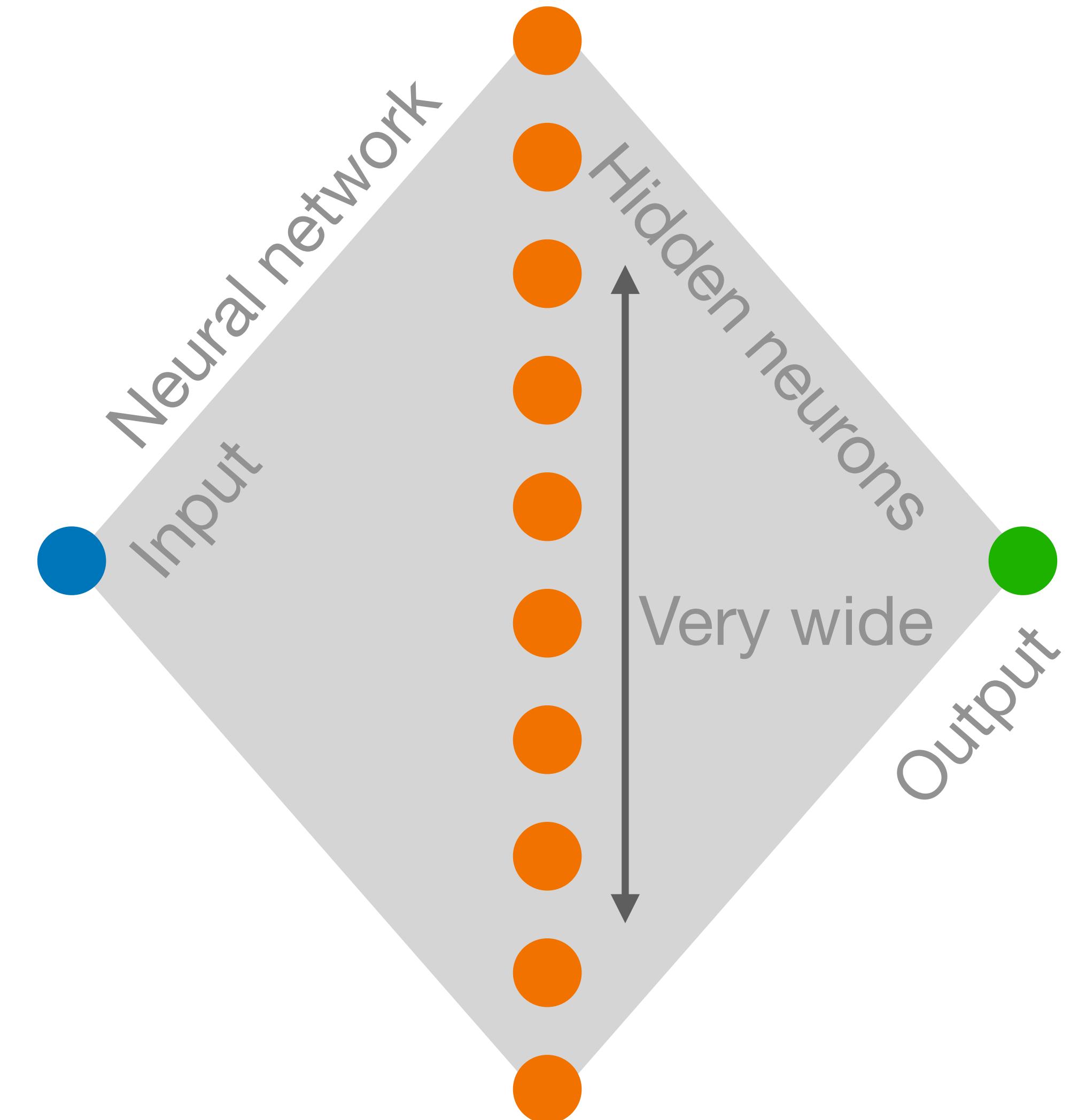


Correlations induce bad conditioning



# Neural Tangent Kernels (NTK)

- Theory for deep (wide) neural networks
- Converge to kernel model
- <https://www.youtube.com/watch?v=raT2ECrvbag>



# **Summary**

# Kernel methods – Pros and Cons

- + Powerful Modeling Tool  
(non-linear problems become linear in kernel space)
- + Omnipurpose Kernels  
(Gaussian works well in many cases)
- + Kernel methods can handle symbolic objects
- + When you have less data points than your data has dimensions kernel methods can offer a dramatic speedup
  
- Difficult to understand what's happening in kernel space
- Model complexity increases with number of data points
- If you have too much data, kernel methods can be slow

# References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2007.
- C. F. Gauß. Theoria motus corporum coelestium in sectionibus conicis solem ambientium. *Göttingen*, 1809.
- A. E. Hoerl and R. W. Kennar. Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 12(1):69–82, 1970.
- F. Jäkel, B. Schölkopf, and F. A. Wichmann. Does cognitive science need kernels? *Trends Cogn Sci*, 13(9):381–388, 2009. doi: 016/j.tics.2009.06.002.
- G. Kimeldorf and G. Wahba. Some results on tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1):82–95, 1971.
- A.-M. Legendre. *Nouvelles méthodes pour la détermination des orbites des comètes*, chapter Sur la methode des moindres quarres. Firmin Didot, <http://imgbase-scd-ulp.u-strasbg.fr/displayimage.php?pos=-141297>, 1805.
- J. Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909.
- R. N. Shepard. Toward a universal law of generalization for psychological science. *Science*, 237(4820):1317–23, 1987.
- A. N. Tychonoff. On the stability of inverse problems. *Doklady Akademii Nauk SSSR*, 39(5):195–198, 1943.