# Lecture 16: Kernel Ridge Regression
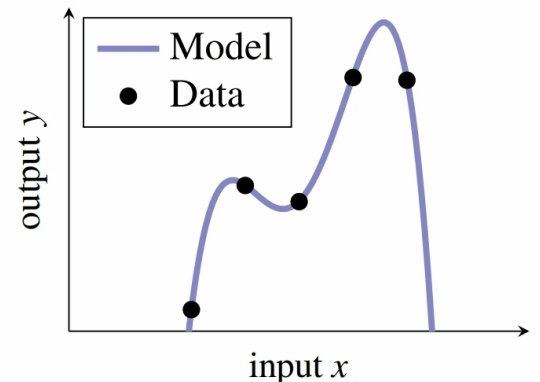
- Previously, in our discussion on polynomial regression, we had shown how could map the input variables to a new feature space of polynomials

$$y = a_0 + a_1 x + \epsilon \quad \overset{\text{Mapping}}{\Longrightarrow} \quad y = b_0 + b_1 x + b_2 x^2 + \cdots + b_p x^p + \epsilon$$

Original input space          New input space

- While we created these non-linear transformations of the original input, we were still using linear regression, since the parameters $b_0, b_1, \ldots, b_p$ appear linearly with $\underline{\phi}(x) = \begin{bmatrix} 1 & x & x^2 & \cdots & x^p \end{bmatrix}^T$ as the new input

$$\boxed{y = \underline{\theta}^T \underline{\phi}(x) + \epsilon}$$

Linear regression with a 4th order polynomial



output $y$ / input $x$ — Model, ● Data

- For vector-valued input $\underline{x}$, the non-linear transformation could be expressed as

$$y = \underbrace{\underline{\phi}^T(\underline{x})}_{1 \times d} \underbrace{\underline{\theta}}_{d \times 1} + \underbrace{\epsilon}_{1 \times 1}$$
$$\underset{1 \times 1}{}$$

$$\underline{x} \in \mathbb{R}^P$$

$$\underline{\phi}(\underline{x}) \in \mathbb{R}^d$$
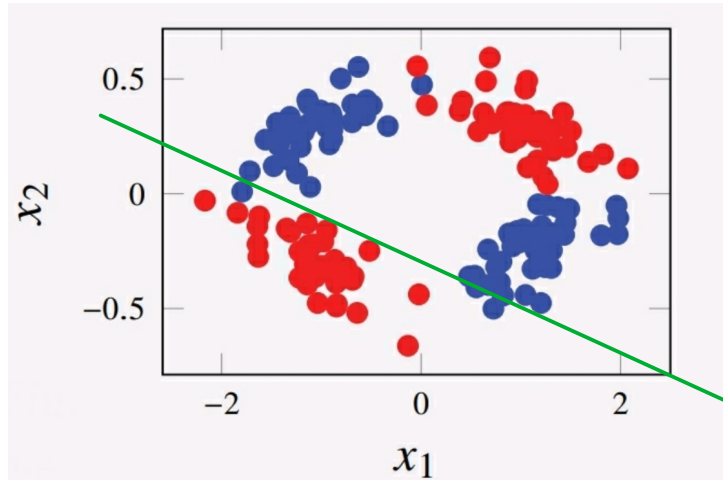
$$\underline{\theta} \in \mathbb{R}^d$$

- Any choice of nonlinear transformation $\underline{\phi}(\underline{x})$ can be used!

- Writing the vectorized linear regression for training data $\{\underline{x}_i, y_i\}_{i=1}^N$

$$\underline{\underline{X}} = \begin{bmatrix} - & \underline{x}_1^T & - \\ - & \underline{x}_2^T & - \\ & \vdots & \\ - & \underline{x}_N^T & - \end{bmatrix}, \qquad \underline{\underline{\Phi}}(\underline{\underline{X}}) = \begin{bmatrix} - & \underline{\phi}(\underline{x}_1)^T & - \\ - & \underline{\phi}(\underline{x}_2)^T & - \\ & \vdots & \\ - & \underline{\phi}(\underline{x}_N)^T & - \end{bmatrix}, \qquad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$N \times P \qquad\qquad\qquad N \times d \qquad\qquad\qquad N \times 1$$
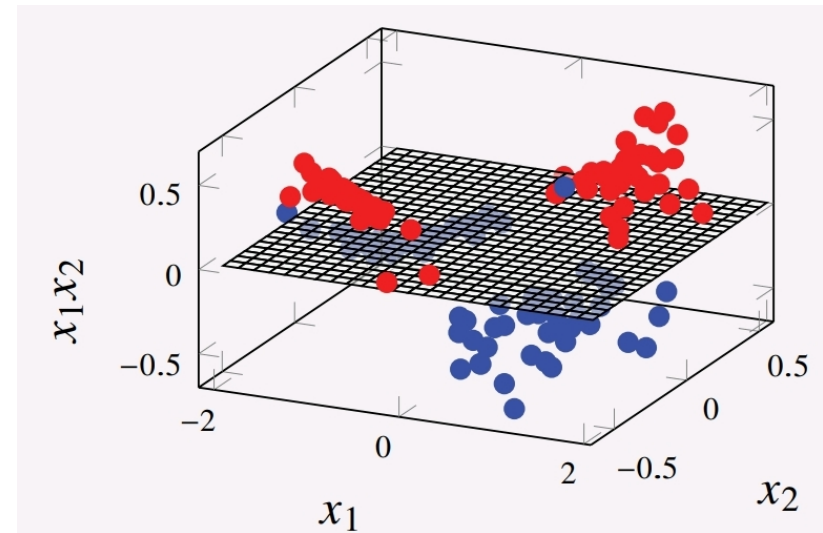
$$\boxed{\underline{y} = \underline{\underline{\Phi}}(\underline{\underline{X}}) \underline{\theta} + \underline{\epsilon}}$$

# Example of non-linear feature transformation for classification



A linear classifier would not work on the original input space

$\left(\begin{array}{l} \text{there is no line that can} \\ \text{separate the two classes} \end{array}\right)$

With an introduction of an extra feature $x_1 x_2$ the problem becomes linearly separable

A carefully engineered transformation $\underline{\phi}(x)$ in linear regression or linear classification may perform very well for a specific ML problem

- We would like a $\underline{\phi}(\underline{x})$ that would work for most problems

- Thus, $\underline{\phi}(\underline{x})$ should contain a lot of transformations that could possibly be of interest to most problems

- Therefore, we should choose $d$, the dimension of $\underline{\phi}(\underline{x})$, really large

  - $d > N$ and eventually let $d \to \infty$

**# of input features** $\nearrow$

**# of training data-points** $\uparrow$

$$
\begin{bmatrix} \vdots \end{bmatrix}_{N \times 1} = \begin{bmatrix} \vdots & \vdots & \vdots & \cdots & \vdots & \vdots \end{bmatrix}_{N \times d} \begin{bmatrix} \vdots \end{bmatrix}_{d \times 1}
$$

- However, increasing the flexibility of a model also means it can overfit the training data

- We will have to use some kind of regularization to prevent overfitting

- Let us use an $L_2$-regularization

- Reformulating the linear regression with transformed features $\underline{\phi}(x)$, we get the estimate of parameters as

$$\hat{\underline{\Theta}} = \underset{\underline{\Theta}}{\arg\min} \; J(\underline{\Theta})$$

$$= \underset{\underline{\Theta}}{\arg\min} \left[ \frac{1}{N} \sum_{i=1}^{N} \left( y_i - \underline{\phi}(x_i)^T \underline{\Theta} \right)^2 + \lambda \|\underline{\Theta}\|_2^2 \right]$$

- Linear regression with $L_2$-regularization has closed-form solution

$$\hat{\underline{\Theta}} = \left( \underline{\underline{X}}^T \underline{\underline{X}} + N\lambda \underline{\underline{I}} \right)^{-1} \underline{\underline{X}}^T \underline{y} \quad (\text{recall}!)$$

$$\underset{d \times 1}{\hat{\underline{\Theta}}} = \left( \underbrace{\underline{\underline{\Phi}}(\underline{\underline{X}})^T}_{d \times N} \underbrace{\underline{\underline{\Phi}}(\underline{\underline{X}})}_{N \times d} + \underbrace{N\lambda \underline{\underline{I}}}_{d \times d} \right)^{-1} \underbrace{\underline{\underline{\Phi}}(\underline{\underline{X}})^T}_{d \times N} \underbrace{\underline{y}}_{N \times 1}$$

- Linear regression with $L_2$-regularization has closed-form solution

$$\underbrace{\hat{\underline{\Theta}}}_{d \times 1} = \left( \underbrace{\underline{\Phi}(\underline{X})^T}_{d \times N} \underbrace{\underline{\Phi}(\underline{X})}_{N \times d} + \underbrace{N \lambda \underline{I}}_{d \times d} \right)^{-1} \underbrace{\underline{\Phi}(\underline{X})^T}_{d \times N} \underbrace{\underline{Y}}_{N \times 1}$$

- The downside of choosing a very large number of features, $d$, is that $\text{size}(\hat{\underline{\Theta}}) = d \to \infty$ and storing $\hat{\underline{\Theta}}$ digitally becomes a problem!

- During prediction, we need to use the $d$-dimensional parameter vector $\hat{\underline{\Theta}}$

$$\hat{y}(\underline{x}_*) = \underbrace{\underline{\phi}(\underline{x}_*)^T}_{1 \times d} \underbrace{\hat{\underline{\Theta}}}_{d \times 1}$$

- But if $d \to \infty$, how to scale computations or meet storage demands ??

- Let's try to reformulate the prediction

$$\hat{y}(\underline{x}_*) = \underline{\phi}(\underline{x}_*)^T \underset{d\times 1}{\underset{1\times d}{\hat{\Theta}}} = \hat{\Theta}^T \underline{\phi}(\underline{x}_*)$$

$$= \left[ \underbrace{\left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda\underline{\underline{I}} \right)^{-1}}_{d\times d} \underbrace{\underline{\underline{\Phi}}(\underline{\underline{x}})^T}_{d\times N} \underbrace{\underline{y}}_{N\times 1} \right]^T \underbrace{\underline{\phi}(\underline{x}^*)}_{d\times 1}$$

$$= \underbrace{\underline{y}^T}_{1\times N} \underbrace{\underline{\underline{\Phi}}(\underline{\underline{x}})}_{N\times d} \underbrace{\left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda\underline{\underline{I}} \right)^{-1}}_{d\times d} \underbrace{\underline{\phi}(\underline{x}^*)}_{d\times 1}$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{N\times 1}$$

This entire expression is independent of 'd' and if we could compute this N-dimensional vector directly, then it would be great

- However, $\left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda\underline{\underline{I}} \right)^{-1}$ still requires inverting

  $\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxx}}_{d\times d}$ a $d\times d$ matrix !!

$$\hat{y}(\underline{x}_*) = \underline{Y}^T \underline{\underline{\Phi}}(\underline{\underline{x}}) \left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda\underline{\underline{I}} \right)^{-1} \underline{\phi}(\underline{x}_*)$$

$$\underbrace{\hphantom{\left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda\underline{\underline{I}} \right)}}_{d \times d}$$

- To prevent inverting a $d \times d$ matrix, where $d$ is very large, lets use a matrix identity: $\underline{\underline{A}} \left( \underline{\underline{A}}^T \underline{\underline{A}} + \underline{\underline{I}} \right)^{-1} = \left( \underline{\underline{A}}\,\underline{\underline{A}}^T + \underline{\underline{I}} \right)^{-1} \underline{\underline{A}}$

$$\hat{y}(\underline{x}_*) = \underbrace{\underline{Y}^T}_{1 \times N} \underbrace{\left( \underline{\underline{\Phi}}(\underline{\underline{x}}) \underline{\underline{\Phi}}(\underline{\underline{x}})^T + N\lambda\underline{\underline{I}} \right)^{-1}}_{N \times N} \underbrace{\underbrace{\underline{\underline{\Phi}}(\underline{\underline{x}})}_{N \times d} \underbrace{\underline{\phi}(\underline{x}_*)}_{d \times 1}}_{N \times 1}$$

- We can now compute $\hat{y}(\underline{x}_*)$ without having to deal with any $d$-dimensional quantities, if we can compute $\underline{\underline{\Phi}}(\underline{\underline{x}})\,\underline{\underline{\Phi}}(\underline{\underline{x}})^T$ & $\underline{\underline{\Phi}}(\underline{\underline{x}})\,\underline{\phi}(\underline{x}^*)$ directly

$$\underline{\underline{\Phi}}(\underline{\underline{x}}) = \begin{bmatrix} \text{---}\ \underline{\phi}(\underline{x}_1)^T\ \text{---} \\ \text{---}\ \underline{\phi}(\underline{x}_2)^T\ \text{---} \\ \vdots \\ \text{---}\ \underline{\phi}(\underline{x}_N)^T\ \text{---} \end{bmatrix}$$

$N \times d$

$$\underline{\underline{\Phi}}(\underline{\underline{x}})^T = \begin{bmatrix} | & & | \\ | & & | \\ \underline{\phi}(\underline{x}_1) & \cdots & \underline{\phi}(\underline{x}_N) \\ | & & | \\ | & & | \end{bmatrix}$$

$d \times N$

$$\hat{y}(\underline{x}_*) = \underline{y}^T \left( \underline{\underline{\Phi}}(\underline{x})\, \underline{\underline{\Phi}}(\underline{x})^T + N\lambda \underline{\underline{I}} \right)^{-1} \underline{\underline{\Phi}}(\underline{x})\, \underline{\phi}(\underline{x}_*)$$

$\underbrace{\phantom{\underline{y}^T}}_{1 \times N}$ $\underbrace{\phantom{\left( \underline{\underline{\Phi}}(\underline{x})\, \underline{\underline{\Phi}}(\underline{x})^T + N\lambda \underline{\underline{I}} \right)}}_{N \times N}$ $\underbrace{\phantom{\underline{\underline{\Phi}}(\underline{x})\, \underline{\phi}(\underline{x}_*)}}_{N \times 1}$

- Lets look at the two matrix multiplications

$$\underline{\underline{\Phi}}(\underline{x})\, \underline{\underline{\Phi}}(\underline{x})^T = \begin{bmatrix} \underline{\phi}(\underline{x}_1)^T \underline{\phi}(\underline{x}_1) & \underline{\phi}(\underline{x}_1)^T \underline{\phi}(\underline{x}_2) & \cdots & \underline{\phi}(\underline{x}_1)^T \underline{\phi}(\underline{x}_N) \\ \underline{\phi}(\underline{x}_2)^T \underline{\phi}(\underline{x}_1) & \underline{\phi}(\underline{x}_2)^T \underline{\phi}(\underline{x}_2) & \cdots & \underline{\phi}(\underline{x}_2)^T \phi(x_N) \\ \vdots & \vdots & \vdots & \vdots \\ \underline{\phi}(\underline{x}_N)^T \underline{\phi}(\underline{x}_1) & \underline{\phi}(\underline{x}_N)^T \underline{\phi}(\underline{x}_2) & \cdots & \phi(\underline{x}_N)^T \phi(x_N) \end{bmatrix}$$

$\underbrace{\underline{\underline{\Phi}}(\underline{x})}_{N \times d}\ \underbrace{\underline{\underline{\Phi}}(\underline{x})^T}_{d \times N}$    $\underbrace{\phantom{xx}}_{N \times N}$

(under $\underline{\phi}(\underline{x}_2)^T$: $\underbrace{\phantom{xx}}_{1 \times d}$, under $\underline{\phi}(\underline{x}_1)$: $\underbrace{\phantom{xx}}_{d \times 1}$)

$N \times N$

$$\underline{\underline{\Phi}}(\underline{x})\, \underline{\phi}(\underline{x}_*) = \begin{bmatrix} \underline{\phi}(\underline{x}_1)^T \underline{\phi}(\underline{x}_*) \\ \underline{\phi}(\underline{x}_2)^T \underline{\phi}(\underline{x}_*) \\ \vdots \\ \underline{\phi}(\underline{x}_N)^T \underline{\phi}(\underline{x}_*) \end{bmatrix}$$

$\underbrace{\underline{\underline{\Phi}}(\underline{x})}_{N \times d}\ \underbrace{\underline{\phi}(\underline{x}_*)}_{d \times 1}$    $\underbrace{\phantom{xx}}_{N \times 1}$

- $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$ is an inner product between two d-dimensional vectors $\underline{\phi}(\underline{x})$ & $\underline{\phi}(\underline{x}')$

- $\underline{\phi}(\underline{x})$ enters the prediction $y(\underline{x}_*)$ only as these inner products

Lets take an example of polynomial transformation

$$x \leftarrow \text{scalar } (p=1)$$
$$\underset{1 \times 1}{}$$

$\underline{\phi}(x)$ is say a third-order scaled polynomial

of the form :
$$\begin{bmatrix} 1 \\ \sqrt{3}\, x \\ \sqrt{3}\, x^2 \\ x^3 \end{bmatrix}$$
$$4 \times 1$$

$$\underline{\phi}(x)^T \underline{\phi}(x') = \begin{bmatrix} 1 & \sqrt{3}x & \sqrt{3}x^2 & x^3 \end{bmatrix} \begin{bmatrix} 1 \\ \sqrt{3}\, x' \\ \sqrt{3}\, x'^2 \\ x'^3 \end{bmatrix}$$

$$= 1 + 3xx' + 3x^2 x'^2 + x^3 x'^3 = (1 + xx')^3$$

In general, if $\underline{\phi}(x)$ is a suitably scaled polynomial of order 'd', then

$$\underline{\phi}(x)^T \underline{\phi}(x') = (1 + xx')^d$$

inner product

- Usually to compute $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$

  - One has to first compute $d$-dimensional vectors $\underline{\phi}(\underline{x})$ and $\underline{\phi}(\underline{x}')$, and

  - then compute their inner product $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$

- However, for the previous example, we found that we could have just evaluated the expression $(1 + x x')^d$ directly

- Important point: If we make the choice of $\underline{\phi}(\underline{x})$ s.t. the inner product $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$ can be computed without first calculating $\underline{\phi}(\underline{x})$, we can let $d \to$ very large

- Infact, if you don't really care about $\underline{\phi}(\underline{x})$ explicitly sometimes, then the need of deriving $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$ can be bypassed by using the concept of kernels

# Introducing the idea of kernels

- In simple terms, a kernel $K(\underline{x}, \underline{x}')$ $\leftarrow$ kappa is any function that takes two arguments $\underline{x}$ and $\underline{x}'$ from the same space $\mathbb{R}^p$ and returns a scalar

- We will mostly limit ourselves to kernels that are real-valued and symmetric

  i.e. $K(\underline{x}, \underline{x}') = K(\underline{x}', \underline{x}) \in \mathbb{R}$ for all $\underline{x}$ and $\underline{x}'$

  For example, $K(x, x') = (1 + x x')^d$ is such a kernel

- In fact, the inner product of two non-linear input transformation is also an example of a kernel:

$$K(\underline{x}, \underline{x}') = \underline{\phi}(\underline{x})^\top \underline{\phi}(\underline{x}')$$

- So instead of choosing $\underline{\phi}(\underline{x})$ and deriving its inner product $\underline{\phi}(\underline{x})^\top \underline{\phi}(\underline{x}')$ sometimes one can choose a kernel $K(\underline{x}, \underline{x}')$ directly $\leftarrow$ KERNEL TRICK

If $\underline{x}$ enters the model as $\underline{\phi}(\underline{x})^T \underline{\phi}(\underline{x}')$ only, we can choose a kernel $K(\underline{x}, \underline{x}')$ directly, instead of choosing $\underline{\phi}(\underline{x})$ <span style="color:green">**KERNEL TRICK**</span>

- Mathematically, we can rewrite

$$\hat{y}(\underline{x}_*) = \underbrace{\underline{Y}^T}_{1 \times N} \underbrace{\left( \underline{\underline{\Phi}}(\underline{x}) \underline{\underline{\Phi}}(\underline{x})^T + N\lambda \underline{\underline{I}} \right)^{-1}}_{N \times N} \underbrace{\underline{\underline{\Phi}}(\underline{x}) \underline{\phi}(\underline{x}_*)}_{N \times 1}$$

as

$$\boxed{\hat{y}(\underline{x}_*) = \underline{Y}^T \left( \underbrace{\underline{\underline{K}}(\underline{X}, \underline{X})}_{N \times N} + N\lambda \underline{\underline{I}} \right)^{-1} \underbrace{\underline{\underline{K}}(\underline{X}, \underline{x}_*)}_{N \times 1}}$$

where

$$\underline{\underline{K}}(\underline{X}, \underline{X}) = \begin{bmatrix} K(\underline{x}_1, \underline{x}_1) & K(\underline{x}_1, \underline{x}_2) & \cdots & K(\underline{x}_1, \underline{x}_N) \\ K(\underline{x}_2, \underline{x}_1) & K(\underline{x}_2, \underline{x}_2) & \cdots & K(\underline{x}_2, \underline{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ K(\underline{x}_N, \underline{x}_1) & K(\underline{x}_N, \underline{x}_2) & \cdots & K(\underline{x}_N, \underline{x}_N) \end{bmatrix}, \quad \underline{\underline{K}}(\underline{X}, \underline{x}_*) = \begin{bmatrix} K(\underline{x}_1, \underline{x}_*) \\ K(\underline{x}_2, \underline{x}_*) \\ \vdots \\ K(\underline{x}_N, \underline{x}_*) \end{bmatrix}$$

- Recall, linear regression with $L_2$-regularization was called as
   ridge regression

- $\hat{y}(\underline{x}_*) = \underline{y}^T \left( \underline{\underline{K}}(\underline{\underline{X}}, \underline{\underline{X}}) + N\lambda\underline{\underline{I}} \right)^{-1} \underline{\underline{K}}(\underline{\underline{X}}, \underline{x}_*)$

   ↳ This equation describes linear regression with $L_2$-regularization using
      a kernel, hence is called **kernel ridge regression**

$$\underline{\underline{K}}(\underline{\underline{X}}, \underline{\underline{X}}) \longleftarrow \text{Gram matrix (or Gramian matrix)}$$

- The design choice is now to select a kernel $K(\underline{x}, \underline{x}')$ instead of $\underline{\phi}(\underline{x})$

- In practice, choosing $K(\underline{x}, \underline{x}')$ is much easier than choosing an
   appropriate $\underline{\phi}(\underline{x})$ especially when the number of transformed
   features (i.e. $d$) is very large

- From computation point of view, we can choose $K(\underline{x}, x')$ arbitrarily, as long as we can compute

$$\hat{y}(\underline{x}_*) = \underline{y}^T \left( \underline{\underline{K}}(\underline{X}, \underline{X}) + N\lambda \underline{\underline{I}} \right)^{-1} \underline{\underline{K}}(\underline{X}, \underline{x}_*)$$

Gram matrix

this must be invertible

- For the inverse $\left( \underline{\underline{K}}(\underline{X}, \underline{X}) + N\lambda \underline{\underline{I}} \right)^{-1}$ to exist, we will restrict ourselves to kernels for which the Gram matrix $\underline{\underline{K}}(\underline{X}, \underline{X})$ is always **PSD**

positive semi-definite

$$\left[ \begin{array}{l} \text{A matrix } \underline{\underline{M}} \text{ is said to be PSD if} \\ \quad \circ \quad \underline{v}^T \underline{\underline{M}} \, \underline{v} \geqslant 0 \quad \text{for all } \underline{v} \\ \quad \bullet \quad \text{equivalently, all eigenvalues of } \underline{\underline{M}} \geqslant 0 \end{array} \right]$$

- Kernels $K(\underline{x}, \underline{x}')$ that leads to a PSD $\underline{\underline{K}}(\underline{X}, \underline{X})$ are called PSD kernels

# Examples of positive semi-definite kernels

- **Squared exponential** kernel (also known as radial basis function, RBF exponentiated quadratic, Gaussian kernel)

$$K(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|_2^2}{2\ell^2}\right)$$

where $\ell > 0$ is a hyperparameter to be chosen by the user (by cross-validation)

- Polynomial kernel

$$K(\underline{x}, \underline{x}') = \left(c + x^\top x'\right)^{d-1}$$

← order of polynomial

- You will see more examples of symmetric PSD kernels later

$$\hat{y}(\underline{x}_*) = \underline{y}^T \underbrace{\left( \underline{\underline{K}}(\underline{X}, \underline{X}) + N\lambda\underline{\underline{I}} \right)^{-1}}_{N \times N \text{ matrix}} \underline{\underline{K}}(\underline{X}, \underline{x}_*)$$

- Inversion of a high-dimensional matrix is a very heavy operation

- Do we need to invert the matrix $\left( \underline{\underline{K}}(\underline{X}, \underline{X}) + N\lambda\underline{\underline{I}} \right)$ every time we predict for a new test input $\underline{x}_*$?

  — Not necessary

- We can introduce an N-dimensional vector $\hat{\underline{\alpha}}$ ← (Dual parameter)

$$\underbrace{\hat{\underline{\alpha}}^T}_{1 \times N} = \underbrace{\begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \\ \vdots \\ \hat{\alpha}_N \end{bmatrix}^{\textcircled{T}}}_{1 \times N} = \underbrace{\underline{y}^T}_{1 \times N} \underbrace{\left( \underline{\underline{K}}(\underline{X}, \underline{X}) + N\lambda\underline{\underline{I}} \right)^{-1}}_{N \times N} \Rightarrow$$

Test prediction

$$\boxed{\hat{y}(\underline{x}_*) = \hat{\underline{\alpha}}^T \underline{\underline{K}}(\underline{X}, \underline{x}_*)}$$

So now, we only need to compute and store $\hat{\underline{\alpha}}$ and $\underline{X}$

# Summary of Kernel Ridge Regression (KRR)

## Training

Input: Training data $T = \{\underline{x}_i, y_i\}_{i=1}^{N}$, a kernel $K$,

regularization parameter $\lambda$

Output: Dual parameter $\hat{\underline{\alpha}}$

- Assemble $\underline{\underline{X}}$ and compute $\underline{\underline{K}}(\underline{\underline{X}}, \underline{\underline{X}})$

- Compute $\hat{\underline{\alpha}}$ as

$$\hat{\underline{\alpha}}^T = \underline{y}^T \left( \underline{\underline{K}}(\underline{\underline{X}}, \underline{\underline{X}}) + N\lambda \underline{\underline{I}} \right)^{-1}$$

## Prediction with kernel ridge regression

Input: Learned dual parameter $\hat{\underline{\alpha}}$ and test input $\underline{x}_*$

Output: Prediction $\hat{y}(\underline{x}_*) = \hat{\underline{\alpha}}^T \underline{\underline{K}}(\underline{\underline{X}}, \underline{x}_*)$

# Primal vs Dual formulation

$$\hat{\underline{\Theta}} = \left( \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{\underline{\Phi}}(\underline{\underline{x}}) + N\lambda \underline{\underline{I}} \right)^{-1} \underline{\underline{\Phi}}(\underline{\underline{x}})^T \underline{y}$$

$$\hat{y}(\underline{x}_*) = \underline{\phi}(\underline{x}_*)^T \hat{\underline{\Theta}}$$

Primal formulation
of linear regression

$$\hat{\underline{\alpha}} = \underline{y}^T \left( \underline{\underline{K}}(\underline{\underline{x}}, \underline{\underline{x}}) + N\lambda \underline{\underline{I}} \right)^{-1}$$

$$\hat{y}(\underline{x}_*) = \hat{\underline{\alpha}}^T \underline{\underline{K}}(\underline{\underline{x}}, \underline{x}_*)$$

Dual formulation
of linear regression

$-\ \hat{\underline{\Theta}} \in \mathbb{R}^d, \ d \to \infty$

$-\ \hat{\underline{\alpha}} \in \mathbb{R}^N, \ N \to \#$ of data pts
$\qquad\qquad\qquad\qquad$ (finite)

- By comparing the two formulation, we can find a relation between $\hat{\underline{\Theta}}$ and $\hat{\underline{\alpha}}$

$$\hat{y}(\underline{x}_*) = \hat{\underline{\Theta}}^T \underline{\phi}(\underline{x}_*) = \hat{\underline{\alpha}}^T \underbrace{\underline{\underline{\Phi}}(\underline{\underline{x}}) \underline{\phi}(\underline{x}_*)}_{\underline{\underline{K}}(\underline{\underline{x}}, \underline{x}_*)}$$

$$\Rightarrow \quad \boxed{\underset{d \times 1}{\hat{\underline{\Theta}}} = \underset{d \times N}{\underline{\underline{\Phi}}(\underline{\underline{x}})^T} \ \underset{N \times 1}{\hat{\underline{\alpha}}}}$$

← this is a general result of Representer theorem

# Simplied version of Representer's Theorem

**Theorem:** Let $\hat{y}(x) = \underline{\Theta}^T \underline{\phi}(x)$ with a fixed nonlinear transform $\underline{\phi}(x)$, with $\underline{\Theta}$ learned from training data $\{\underline{x}_i, y_i\}_{i=1}^N$

(The dimensionality of $\underline{\Theta}$ and $\underline{\phi}(x)$ need not be finite)

Furthermore, let $L(y, \hat{y})$ be any arbitrary loss function & $h: [0, \infty] \longrightarrow \mathbb{R}$ be a strictly monotonically increasing function

Then, the estimate $\hat{\underline{\Theta}}$ which is the argmin of the cost function $J(\Theta)$, i.e.

$$\hat{\underline{\Theta}} = \underset{\underline{\Theta}}{\arg\min} \; \frac{1}{N} \sum_{i=1}^{N} L\left(y_i, \underbrace{\underline{\Theta}^T \underline{\phi}(x_i)}_{\hat{y}(x_i)}\right) + h\left(\|\underline{\Theta}\|_2^2\right)$$

can be written as

$$\hat{\underline{\Theta}} = \underline{\underline{\Phi}}(\underline{x})^T \underline{\alpha}, \quad \text{with some } N\text{-dimensional vector } \underline{\alpha}$$

- What does the representer theorem mean?

  - It suggests that if $\hat{y}(x) = \Theta^T \phi(x)$, and $\Theta$ is to be learned using any loss function and $L_2$-regularization, then $\hat{\Theta}$ can be learned also from its dual parameter $\hat{a}$, using: $\hat{\Theta} = \Phi(x) \hat{\alpha}$

  - An important implication of the representer theorem is that $L_2$-regularization is crucial in order to obtain the dual formalism, and we could not have obtained KRR with say $L_1$-regularization

  - Representer theorem is very important for most kernel methods. It tell us that we can express some models in terms of dual parameters $\alpha$ which are of finite length $N$, and a kernel $K(x, x')$, instead of the primal parameters $\Theta$ (maybe of infinite length $d$) and a $\phi(x)$