

Solving Kernel Ridge Regression with Gradient Descent for a Non-Constant Kernel (after an introduction to kernels)

Oskar Allerbo

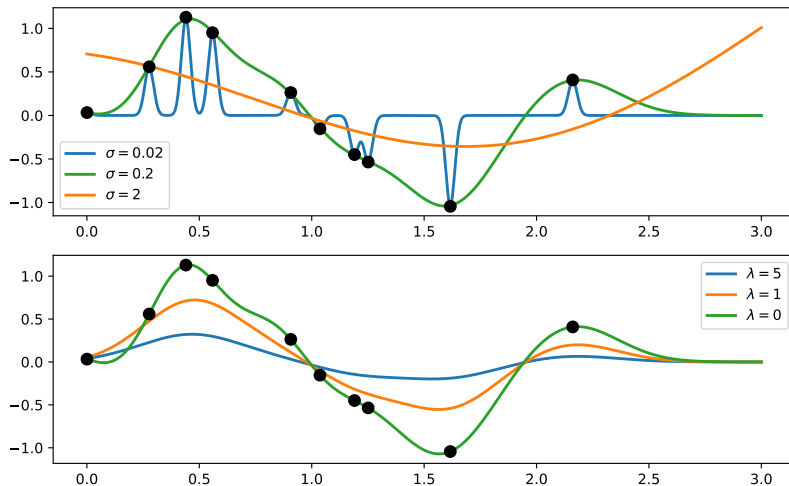
KTH Royal Institute of Technology

2024-03-11

<https://github.com/allerbo/hemavan24/pres.pdf>
<https://github.com/allerbo/hemavan24/slides.pdf>

Kernel Ridge Regression

Kernel Ridge Regression



Kernel Ridge Regression

Kernel ridge regression (KRR) is a generalization of linear ridge regression that is

Kernel Ridge Regression

Kernel ridge regression (KRR) is a generalization of linear ridge regression that is

- non-linear in the data,

Kernel Ridge Regression

Kernel ridge regression (KRR) is a generalization of linear ridge regression that is

- non-linear in the data,
- but linear in the parameters.

Kernel Ridge Regression

Kernel ridge regression (KRR) is a generalization of linear ridge regression that is

- non-linear in the data,
- but linear in the parameters.
- a convex problem, with a closed-form solution.

Outline

- 1 Introduction to Kernels
- 2 Kernel Principal Component Analysis (KPCA)
- 3 Kernel Ridge Regression (KRR)
- 4 Kernel Gradient Descent for Non-constant Kernels (KGD)
- 5 Generalization to Neural Networks (NNs)
- 6 Conclusions

What is a Kernel?

A kernel function:

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:
 $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}. \mathbf{x}, \mathbf{x}' \in \mathbb{R}^p.$

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:

$$k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}. \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^p.$$

- Is symmetric:

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x}).$$

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:

$$k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}. \quad \mathbf{x}, \mathbf{x}' \in \mathbb{R}^p.$$

- Is symmetric:

$$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x}).$$

- Is positive semi-definite:

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0. \quad \text{For all } \mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p, c_i, c_j \in \mathbb{R}.$$

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:
 $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$. $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$.
- Is symmetric:
 $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$.
- Is positive semi-definite:
 $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$. For all $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$, $c_i, c_j \in \mathbb{R}$.
- Is the dot product of the feature expansions of \mathbf{x} and \mathbf{x}' :
 $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^\top \varphi(\mathbf{x}')$. $\varphi(\mathbf{x}), \varphi(\mathbf{x}') \in \mathbb{R}^q$.

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:
 $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$. $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$.
- Is symmetric:
 $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$.
- Is positive semi-definite:
 $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$. For all $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$, $c_i, c_j \in \mathbb{R}$.
- Is the dot product of the feature expansions of \mathbf{x} and \mathbf{x}' :
 $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^\top \varphi(\mathbf{x}')$. $\varphi(\mathbf{x}), \varphi(\mathbf{x}') \in \mathbb{R}^q$.
 $q = \infty$ is a possibility!

What is a Kernel?

A kernel function:

- Takes two arguments and outputs a scalar:
 $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$. $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$.
- Is symmetric:
 $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$.
- Is positive semi-definite:
 $\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$. For all $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$, $c_i, c_j \in \mathbb{R}$.
- Is the dot product of the feature expansions of \mathbf{x} and \mathbf{x}' :
 $k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^\top \varphi(\mathbf{x}')$. $\varphi(\mathbf{x}), \varphi(\mathbf{x}') \in \mathbb{R}^q$.
 $q = \infty$ is a possibility!
- (Is associated to a Reproducing Kernel Hilbert Space, RKHS.)

Examples of Kernels

Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$

Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^q, c \in \mathbb{R}, q \in \mathbb{N}$

Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^q, c \in \mathbb{R}, q \in \mathbb{N}$
- Gaussian: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}}, \sigma > 0$

Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^q, c \in \mathbb{R}, q \in \mathbb{N}$
- Gaussian: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}}, \sigma > 0$
- Laplace: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}}, \sigma > 0$

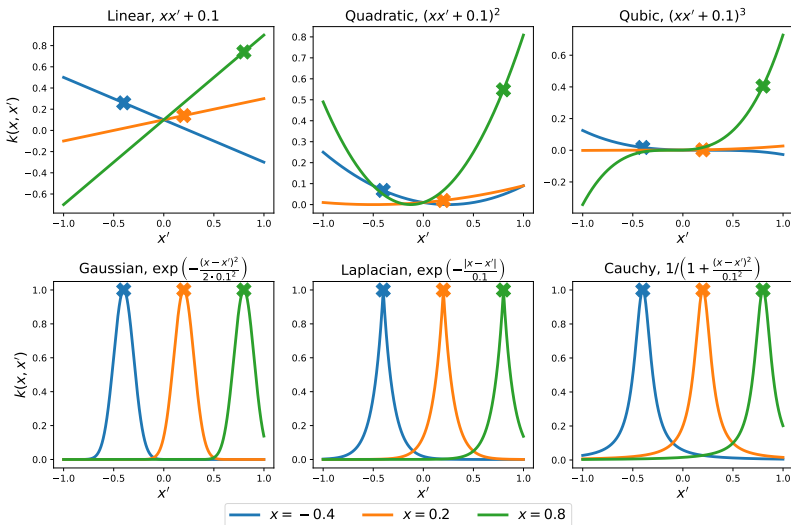
Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^q, c \in \mathbb{R}, q \in \mathbb{N}$
- Gaussian: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}}, \sigma > 0$
- Laplace: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}}, \sigma > 0$
- Cauchy: $k(\mathbf{x}, \mathbf{x}') = \frac{1}{1 + \frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{\sigma^2}}, \sigma > 0$

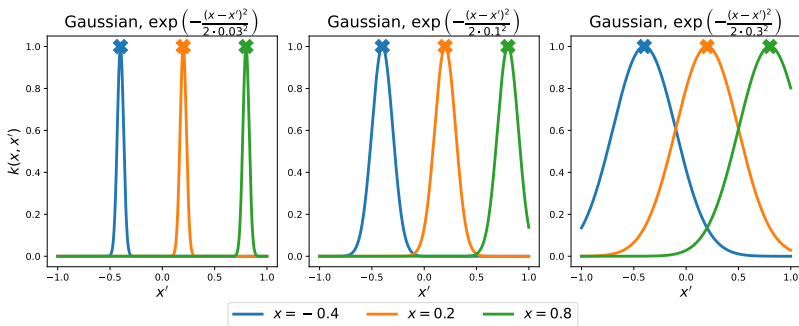
Examples of Kernels

- Linear: $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c, c \in \mathbb{R}$
- Polynomial: $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^q, c \in \mathbb{R}, q \in \mathbb{N}$
- Gaussian: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2}}, \sigma > 0$
- Laplace: $k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}}, \sigma > 0$
- Cauchy: $k(\mathbf{x}, \mathbf{x}') = \frac{1}{1 + \frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{\sigma^2}}, \sigma > 0$
- Matérn: $k(\mathbf{x}, \mathbf{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)} \cdot \left(\sqrt{2\nu} \cdot \frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma} \right)^\nu \cdot K_\nu \left(\sqrt{2\nu} \cdot \frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma} \right)$

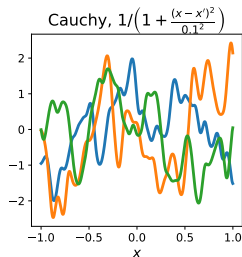
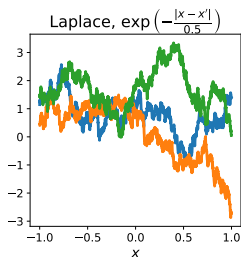
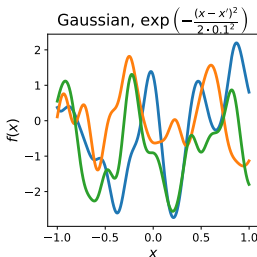
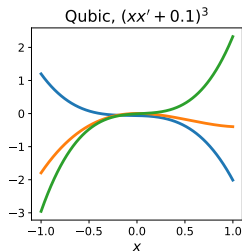
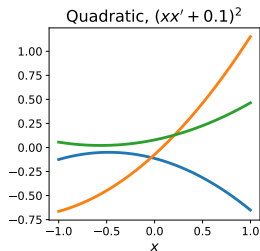
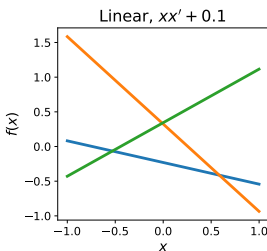
Examples of Kernels (1D Case)



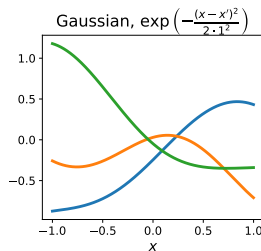
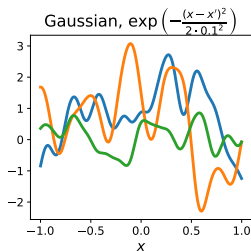
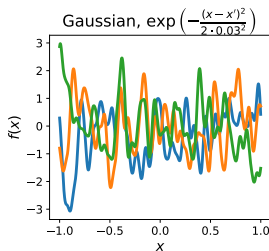
Examples of Kernels (1D Case)



Examples of Functions (1D Case)



Examples of Functions (1D Case)



Examples of Feature Expansions (1D Case)

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$
 $\varphi(x) = [x \quad \sqrt{c}]^T$

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$

$$\varphi(x) = [x \quad \sqrt{c}]^\top$$

$$\varphi(x)^\top \varphi(x') = x \cdot x' + \sqrt{c} \cdot \sqrt{c}$$

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$
 $\varphi(x) = [x \quad \sqrt{c}]^\top$
 $\varphi(x)^\top \varphi(x') = x \cdot x' + \sqrt{c} \cdot \sqrt{c}$
- Quadratic, $k(x, x') = (x \cdot x' + c)^2$
 $\varphi(x) = [x^2 \quad \sqrt{2c}x \quad c]^\top$

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$

$$\varphi(x) = [x \quad \sqrt{c}]^\top$$

$$\varphi(x)^\top \varphi(x') = x \cdot x' + \sqrt{c} \cdot \sqrt{c}$$

- Quadratic, $k(x, x') = (x \cdot x' + c)^2$

$$\varphi(x) = [x^2 \quad \sqrt{2c}x \quad c]^\top$$

$$\varphi(x)^\top \varphi(x') = x^2 \cdot x'^2 + 2c \cdot x \cdot x' + c \cdot c$$

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$
 $\varphi(x) = [x \quad \sqrt{c}]^\top$
 $\varphi(x)^\top \varphi(x') = x \cdot x' + \sqrt{c} \cdot \sqrt{c}$
- Quadratic, $k(x, x') = (x \cdot x' + c)^2$
 $\varphi(x) = [x^2 \quad \sqrt{2c}x \quad c]^\top$
 $\varphi(x)^\top \varphi(x') = x^2 \cdot x'^2 + 2c \cdot x \cdot x' + c \cdot c$
- Gaussian, $k(x, x') = e^{-\frac{(x-x')^2}{2\sigma^2}}$
 $\varphi(x) = e^{-\frac{x^2}{2\sigma^2}} \left[1 \quad \frac{x^1}{\sigma^1 \sqrt{1!}} \quad \cdots \quad \frac{x^k}{\sigma^k \sqrt{k!}} \cdots \right]^\top$

Examples of Feature Expansions (1D Case)

- Linear, $k(x, x') = x \cdot x' + c$

$$\varphi(x) = [x \quad \sqrt{c}]^\top$$

$$\varphi(x)^\top \varphi(x') = x \cdot x' + \sqrt{c} \cdot \sqrt{c}$$

- Quadratic, $k(x, x') = (x \cdot x' + c)^2$

$$\varphi(x) = [x^2 \quad \sqrt{2c}x \quad c]^\top$$

$$\varphi(x)^\top \varphi(x') = x^2 \cdot x'^2 + 2c \cdot x \cdot x' + c \cdot c$$

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

$$\varphi(x) = e^{-\frac{x^2}{2\sigma^2}} \left[1 \quad \frac{x^1}{\sigma^1 \sqrt{1!}} \quad \dots \quad \frac{x^k}{\sigma^k \sqrt{k!}} \dots \right]^\top$$

$$\varphi(x)^\top \varphi(x') = e^{-\frac{1}{2\sigma^2}(x^2 + x'^2)} \cdot \underbrace{\sum_{k=0}^{\infty} \frac{(x \cdot x' / \sigma^2)^k}{k!}}_{=e^{\frac{2 \cdot x \cdot x'}{2\sigma^2}}}$$

Notation

Notation

Training Data: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$.

New Observations: $\mathbf{X}^* \in \mathbb{R}^{n^* \times p}$.

Predictions: $\hat{\mathbf{f}} \in \mathbb{R}^n$, $\hat{\mathbf{f}}^* \in \mathbb{R}^{n^*}$.

Notation

Training Data: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$.

New Observations: $\mathbf{X}^* \in \mathbb{R}^{n^* \times p}$.

Predictions: $\hat{\mathbf{f}} \in \mathbb{R}^n$, $\hat{\mathbf{f}}^* \in \mathbb{R}^{n^*}$.

Feature Expansion Matrices:

$\Phi = \Phi(\mathbf{X}) \in \mathbb{R}^{n \times q}$, $\Phi^* = \Phi^*(\mathbf{X}^*) \in \mathbb{R}^{n^* \times q}$.

Notation

Training Data: $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$.

New Observations: $\mathbf{X}^* \in \mathbb{R}^{n^* \times p}$.

Predictions: $\hat{\mathbf{f}} \in \mathbb{R}^n$, $\hat{\mathbf{f}}^* \in \mathbb{R}^{n^*}$.

Feature Expansion Matrices:

$$\Phi = \Phi(\mathbf{X}) \in \mathbb{R}^{n \times q}, \Phi^* = \Phi^*(\mathbf{X}^*) \in \mathbb{R}^{n^* \times q}.$$

Kernel Matrices:

$$\mathbf{K} = \mathbf{K}(\mathbf{X}) \in \mathbb{R}^{n \times n}, \mathbf{K}^* = \mathbf{K}(\mathbf{X}^*, \mathbf{X}) \in \mathbb{R}^{n^* \times n}.$$

$$\mathbf{K} = \Phi \Phi^\top, \mathbf{K}^* = \Phi^* \Phi^\top.$$

$$k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^\top \varphi(\mathbf{x}').$$

Notation

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1^\top & \text{---} \\ \text{---} & \mathbf{x}_2^\top & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \mathbf{x}_n^\top & \text{---} \end{bmatrix}$$

$$\mathbf{X}^* = \begin{bmatrix} \text{---} & \mathbf{x}_1^{*\top} & \text{---} \\ \text{---} & \mathbf{x}_2^{*\top} & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \mathbf{x}_{n^*}^{*\top} & \text{---} \end{bmatrix}$$

$$\Phi = \begin{bmatrix} \text{---} & \varphi(\mathbf{x}_1)^\top & \text{---} \\ \text{---} & \varphi(\mathbf{x}_2)^\top & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \varphi(\mathbf{x}_n)^\top & \text{---} \end{bmatrix}$$

$$\Phi^* = \begin{bmatrix} \text{---} & \varphi(\mathbf{x}_1^*)^\top & \text{---} \\ \text{---} & \varphi(\mathbf{x}_2^*)^\top & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \varphi(\mathbf{x}_{n^*}^*)^\top & \text{---} \end{bmatrix}$$

$$\Phi^* \Phi^\top = \mathbf{K}^* = \begin{bmatrix} k(\mathbf{x}_1^*, \mathbf{x}_1) & k(\mathbf{x}_1^*, \mathbf{x}_2) & \dots & k(\mathbf{x}_1^*, \mathbf{x}_n) \\ k(\mathbf{x}_2^*, \mathbf{x}_1) & k(\mathbf{x}_2^*, \mathbf{x}_2) & \dots & k(\mathbf{x}_2^*, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_{n^*}^*, \mathbf{x}_1) & k(\mathbf{x}_{n^*}^*, \mathbf{x}_2) & \dots & k(\mathbf{x}_{n^*}^*, \mathbf{x}_n) \end{bmatrix}$$

Kernel Principal Component Analysis

Kernel Principal Component Analysis

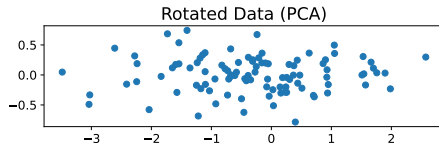
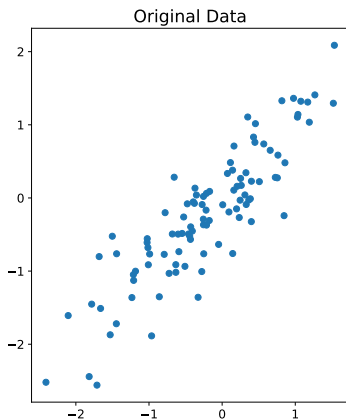
Principal Component Analysis (PCA):

Rotate data to find directions with maximum variance.

Kernel Principal Component Analysis

Principal Component Analysis (PCA):

Rotate data to find directions with maximum variance.



Kernel Principal Component Analysis

Principal Component Analysis (PCA):
Rotate data to find directions with maximum variance.

Kernel Principal Component Analysis

Principal Component Analysis (PCA):

Rotate data to find directions with maximum variance.

$$\mathbf{X}^\top \mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{P}^\top, \mathbf{Z} = \mathbf{X} \mathbf{P}.$$

Kernel Principal Component Analysis

Principal Component Analysis (PCA):

Rotate data to find directions with maximum variance.

$$\mathbf{X}^\top \mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{P}^\top, \mathbf{Z} = \mathbf{X} \mathbf{P}.$$

Equivalent, dual formulation:

$$\mathbf{X} \mathbf{X}^\top = \mathbf{U} \mathbf{D} \mathbf{U}^\top, \mathbf{Z} = \mathbf{U} \sqrt{\mathbf{D}}.$$

Kernel Principal Component Analysis

Principal Component Analysis (PCA):

Rotate data to find directions with maximum variance.

$$\mathbf{X}^\top \mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{P}^\top, \mathbf{Z} = \mathbf{X} \mathbf{P}.$$

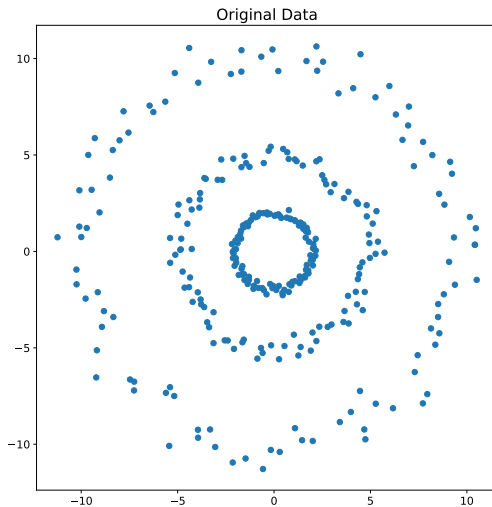
Equivalent, dual formulation:

$$\mathbf{X} \mathbf{X}^\top = \mathbf{U} \mathbf{D} \mathbf{U}^\top, \mathbf{Z} = \mathbf{U} \sqrt{\mathbf{D}}.$$

Kernel PCA:

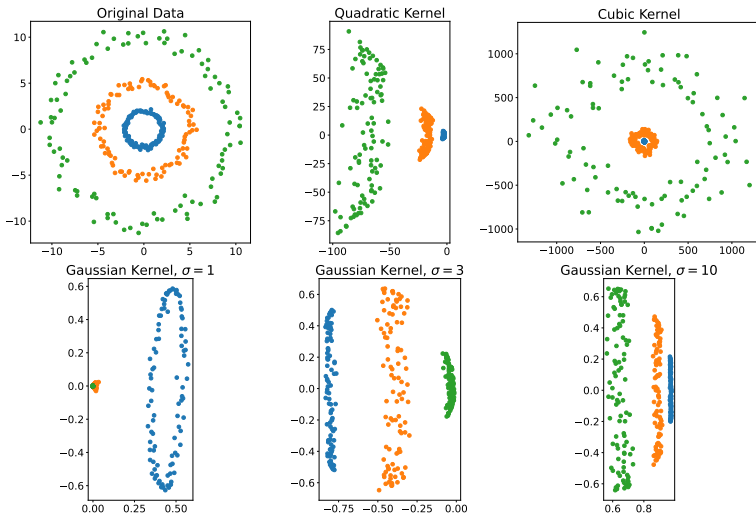
$$\Phi \Phi^\top = \mathbf{K} = \mathbf{U}_K \mathbf{D}_K \mathbf{U}_K^\top, \mathbf{Z}_K = \mathbf{U}_K \sqrt{\mathbf{D}_K}.$$

Kernel Principal Component Analysis

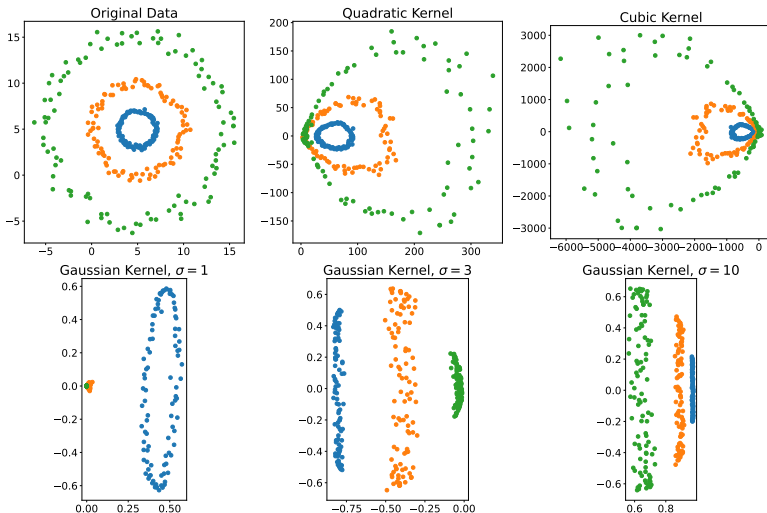


What about this?

Kernel Principal Component Analysis



Kernel Principal Component Analysis



Kernel Ridge Regression

Kernel Ridge Regression

- Linear Ridge Regression

Kernel Ridge Regression

- Linear Ridge Regression
- Dual Formulation of Linear Ridge Regression

Kernel Ridge Regression

- Linear Ridge Regression
- Dual Formulation of Linear Ridge Regression
- Ridge Regression in Feature Space

Kernel Ridge Regression

- Linear Ridge Regression
- Dual Formulation of Linear Ridge Regression
- Ridge Regression in Feature Space
- Dual Formulation of Ridge Regression in Feature Space

Kernel Ridge Regression

- Linear Ridge Regression
- Dual Formulation of Linear Ridge Regression
- Ridge Regression in Feature Space
- Dual Formulation of Ridge Regression in Feature Space
(=Kernel Ridge Regression)

Linear Ridge Regression

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{\lambda}{2} \underbrace{\|\beta\|_2^2}_{=\beta^\top \beta}$$

Linear Ridge Regression

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{\lambda}{2} \underbrace{\|\beta\|_2^2}_{=\beta^\top \beta} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

Linear Ridge Regression

$$\begin{aligned}\hat{\beta} &= \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \frac{\lambda}{2} \underbrace{\|\beta\|_2^2}_{=\beta^\top \beta} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

Predictions are given by

$$\begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{X} \hat{\beta} \\ \mathbf{X}^* \hat{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{X}^* (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \end{bmatrix}.$$

Dual Formulation of Linear Ridge Regression

Dual formulation for $\beta = \mathbf{X}^\top \alpha$:

Dual Formulation of Linear Ridge Regression

Dual formulation for $\beta = \mathbf{X}^\top \alpha$:

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{X}^\top \alpha\|_2^2 + \frac{\lambda}{2} \underbrace{\alpha^\top \mathbf{X}\mathbf{X}^\top \alpha}_{=\|\alpha\|_{\mathbf{X}\mathbf{X}^\top}^2}$$

Dual Formulation of Linear Ridge Regression

Dual formulation for $\beta = \mathbf{X}^\top \alpha$:

$$\begin{aligned}\hat{\alpha} &= \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{X}^\top \alpha\|_2^2 + \frac{\lambda}{2} \underbrace{\alpha^\top \mathbf{X}\mathbf{X}^\top \alpha}_{=\|\alpha\|_{\mathbf{X}\mathbf{X}^\top}^2} \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y}\end{aligned}$$

Dual Formulation of Linear Ridge Regression

Dual formulation for $\beta = \mathbf{X}^\top \alpha$:

$$\begin{aligned}\hat{\alpha} &= \underset{\alpha \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{X}^\top \alpha\|_2^2 + \frac{\lambda}{2} \underbrace{\alpha^\top \mathbf{X}\mathbf{X}^\top \alpha}_{=\|\alpha\|_{\mathbf{X}\mathbf{X}^\top}^2} \\ &= (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y}\end{aligned}$$

Predictions are given by

$$\begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{X} \cdot \mathbf{X}^\top \hat{\alpha} \\ \mathbf{X}^* \cdot \mathbf{X}^\top \hat{\alpha} \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \\ \mathbf{X}^* \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \end{bmatrix}.$$

Linear Ridge Regression

Predictions given by

$$\begin{bmatrix} \mathbf{X}\hat{\boldsymbol{\beta}} \\ \mathbf{X}^*\hat{\boldsymbol{\beta}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{X}^*(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \end{bmatrix}$$

and

$$\begin{bmatrix} \mathbf{X}\mathbf{X}^\top \hat{\boldsymbol{\alpha}} \\ \mathbf{X}^*\mathbf{X}^\top \hat{\boldsymbol{\alpha}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \\ \mathbf{X}^*\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \end{bmatrix}.$$

Linear Ridge Regression

Predictions given by

$$\begin{bmatrix} \mathbf{X}\hat{\boldsymbol{\beta}} \\ \mathbf{X}^*\hat{\boldsymbol{\beta}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{X}^*(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top \mathbf{y} \end{bmatrix}$$

and

$$\begin{bmatrix} \mathbf{X}\mathbf{X}^\top \hat{\boldsymbol{\alpha}} \\ \mathbf{X}^*\mathbf{X}^\top \hat{\boldsymbol{\alpha}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \\ \mathbf{X}^*\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \end{bmatrix}.$$

However,

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_{p \times p})^{-1} \mathbf{X}^\top = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_{n \times n})^{-1}.$$

Ridge Regression in Feature Space

$$\mathbf{x} \in \mathbb{R}^p \mapsto \varphi(\mathbf{x}) \in \mathbb{R}^q.$$

Ridge Regression in Feature Space

$\mathbf{x} \in \mathbb{R}^p \mapsto \varphi(\mathbf{x}) \in \mathbb{R}^q.$

E.g. polynomial regression, $x \mapsto [1, x, x^2, \dots, x^{q-1}]$.

Ridge Regression in Feature Space

$$\mathbf{x} \in \mathbb{R}^p \mapsto \varphi(\mathbf{x}) \in \mathbb{R}^q.$$

E.g. polynomial regression, $x \mapsto [1, x, x^2, \dots, x^{q-1}]$.

$$\mathbf{X} \in \mathbb{R}^{n \times p} \mapsto \mathbf{\Phi} \in \mathbb{R}^{n \times q}, \mathbf{X}^* \in \mathbb{R}^{n^* \times p} \mapsto \mathbf{\Phi}^* \in \mathbb{R}^{n^* \times q}$$

Ridge Regression in Feature Space

$\mathbf{x} \in \mathbb{R}^p \mapsto \varphi(\mathbf{x}) \in \mathbb{R}^q$.

E.g. polynomial regression, $x \mapsto [1, x, x^2, \dots, x^{q-1}]$.

$\mathbf{X} \in \mathbb{R}^{n \times p} \mapsto \Phi \in \mathbb{R}^{n \times q}$, $\mathbf{X}^* \in \mathbb{R}^{n^* \times p} \mapsto \Phi^* \in \mathbb{R}^{n^* \times q}$

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^q} \frac{1}{2} \|\mathbf{y} - \Phi\beta\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \Phi\Phi^\top \alpha\|_2^2 + \frac{\lambda}{2} \|\alpha\|_{\Phi\Phi^\top}^2$$

Ridge Regression in Feature Space

$$\mathbf{x} \in \mathbb{R}^p \mapsto \varphi(\mathbf{x}) \in \mathbb{R}^q.$$

E.g. polynomial regression, $x \mapsto [1, x, x^2, \dots, x^{q-1}]$.

$$\mathbf{X} \in \mathbb{R}^{n \times p} \mapsto \Phi \in \mathbb{R}^{n \times q}, \mathbf{X}^* \in \mathbb{R}^{n^* \times p} \mapsto \Phi^* \in \mathbb{R}^{n^* \times q}$$

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^q} \frac{1}{2} \|\mathbf{y} - \Phi\beta\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \Phi\Phi^\top \alpha\|_2^2 + \frac{\lambda}{2} \|\alpha\|_{\Phi\Phi^\top}^2$$

Predictions are given by

$$\begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \Phi (\Phi^\top \Phi + \lambda I_{p \times p})^{-1} \Phi^\top \mathbf{y} \\ \Phi^* (\Phi^\top \Phi + \lambda I_{p \times p})^{-1} \Phi^\top \mathbf{y} \end{bmatrix} = \begin{bmatrix} \Phi\Phi^\top (\Phi\Phi^\top + \lambda I_{n \times n})^{-1} \mathbf{y} \\ \Phi^*\Phi^\top (\Phi\Phi^\top + \lambda I_{n \times n})^{-1} \mathbf{y} \end{bmatrix}$$

Kernel Ridge Regression

For $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$, $\mathbf{K}^* = \Phi^*\Phi^\top \in \mathbb{R}^{n^* \times n}$,

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \frac{\lambda}{2} \|\alpha\|_{\mathbf{K}}^2$$

$$\begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} \hat{\alpha} = \begin{bmatrix} \mathbf{K}(\mathbf{K} + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \\ \mathbf{K}^*(\mathbf{K} + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \end{bmatrix}$$

Kernel Ridge Regression

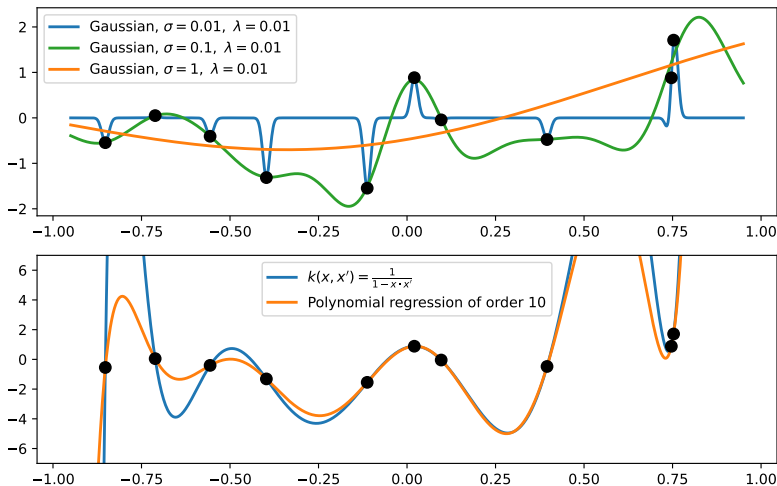
For $\mathbf{K} = \Phi\Phi^\top \in \mathbb{R}^{n \times n}$, $\mathbf{K}^* = \Phi^*\Phi^\top \in \mathbb{R}^{n^* \times n}$,

$$\hat{\alpha} = \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{K}\alpha\|_2^2 + \frac{\lambda}{2} \|\alpha\|_{\mathbf{K}}^2$$

$$\begin{bmatrix} \hat{f} \\ \hat{f}^* \end{bmatrix} = \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} \hat{\alpha} = \begin{bmatrix} \mathbf{K}(\mathbf{K} + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \\ \mathbf{K}^*(\mathbf{K} + \lambda \mathbf{I}_{n \times n})^{-1} \mathbf{y} \end{bmatrix}$$

$q = \infty$ is OK, since Φ and Φ^* are never explicitly calculated.

Kernel Ridge Regression



$$\frac{1}{1-x} = 1 + x + x^2 + x^3 + x^4 + \dots, \quad |x| < 1$$

Kernel Gradient Descent for Non-Constant Kernels

Kernel Gradient Descent for Non-Constant Kernels

Kernel gradient descent in function/prediction space:
(with regularization through early stopping)

$$\begin{bmatrix} \hat{\mathbf{f}}_{t+1} \\ \hat{\mathbf{f}}_{t+1}^* \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}_t \\ \hat{\mathbf{f}}_t^* \end{bmatrix} - \eta \cdot \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} (\hat{\mathbf{f}}_t - \mathbf{y})$$

$$\left(\text{where } \begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} \hat{\alpha} \right)$$

Kernel Gradient Descent for Non-Constant Kernels

Kernel gradient descent in function/prediction space:
(with regularization through early stopping)

$$\begin{bmatrix} \hat{\mathbf{f}}_{t+1} \\ \hat{\mathbf{f}}_{t+1}^* \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}_t \\ \hat{\mathbf{f}}_t^* \end{bmatrix} - \eta \cdot \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} (\hat{\mathbf{f}}_t - \mathbf{y})$$

$$\left(\text{where } \begin{bmatrix} \hat{\mathbf{f}} \\ \hat{\mathbf{f}}^* \end{bmatrix} = \begin{bmatrix} \mathbf{K} \\ \mathbf{K}^* \end{bmatrix} \hat{\alpha} \right)$$

Time dependent kernels:

$$\begin{bmatrix} \hat{\mathbf{f}}_{t+1} \\ \hat{\mathbf{f}}_{t+1}^* \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}_t \\ \hat{\mathbf{f}}_t^* \end{bmatrix} - \eta \cdot \begin{bmatrix} \mathbf{K}_t \\ \mathbf{K}_t^* \end{bmatrix} (\hat{\mathbf{f}}_t - \mathbf{y})$$

Kernel Regression with Gradient Descent and Non-Constant Kernels

Proposition

For a translational invariant kernel with bandwidth σ ,
 $k(\mathbf{x}, \mathbf{x}', \sigma) = k\left(\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}\right)$, and for constant training time, t ,

$$\left\| \nabla_{\mathbf{x}^*} \hat{f}(\mathbf{x}^*, t) \right\|_2 \leq \frac{1}{\sigma} \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

i.e. a larger bandwidth results in a simpler model, and a smaller bandwidth in a more complex model.

Kernel Regression with Gradient Descent and Non-Constant Kernels

Proposition

For a translational invariant kernel with bandwidth σ , $k(\mathbf{x}, \mathbf{x}', \sigma) = k\left(\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}\right)$, and for constant training time, t ,

$$\left\| \nabla_{\mathbf{x}^*} \hat{f}(\mathbf{x}^*, t) \right\|_2 \leq \frac{1}{\sigma} \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

i.e. a larger bandwidth results in a simpler model, and a smaller bandwidth in a more complex model.

We use $1/\sigma$ as a proxy for complexity.

Kernel Regression with Gradient Descent and Non-Constant Kernels

Proposition

For a translational invariant kernel with bandwidth σ ,
 $k(\mathbf{x}, \mathbf{x}', \sigma) = k\left(\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}\right)$, and for constant training time, t ,

$$\left\| \nabla_{\mathbf{x}^*} \hat{\mathbf{f}}(\mathbf{x}^*, t) \right\|_2 \leq \frac{1}{\sigma} \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

i.e. a larger bandwidth results in a simpler model, and a smaller bandwidth in a more complex model.

We use $1/\sigma$ as a proxy for complexity.

$$\begin{bmatrix} \hat{\mathbf{f}}_{t+1} \\ \hat{\mathbf{f}}_{t+1}^* \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}_t \\ \hat{\mathbf{f}}_t^* \end{bmatrix} - \eta \cdot \begin{bmatrix} \mathbf{K}(\sigma_t) \\ \mathbf{K}^*(\sigma_t) \end{bmatrix} (\hat{\mathbf{f}}_t - \mathbf{y})$$

Kernel Regression with Gradient Descent and Non-Constant Kernels

Proposition

For a translational invariant kernel with bandwidth σ ,
 $k(\mathbf{x}, \mathbf{x}', \sigma) = k\left(\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\sigma}\right)$, and for constant training time, t ,

$$\left\| \nabla_{\mathbf{x}^*} \hat{\mathbf{f}}(\mathbf{x}^*, t) \right\|_2 \leq \frac{1}{\sigma} \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

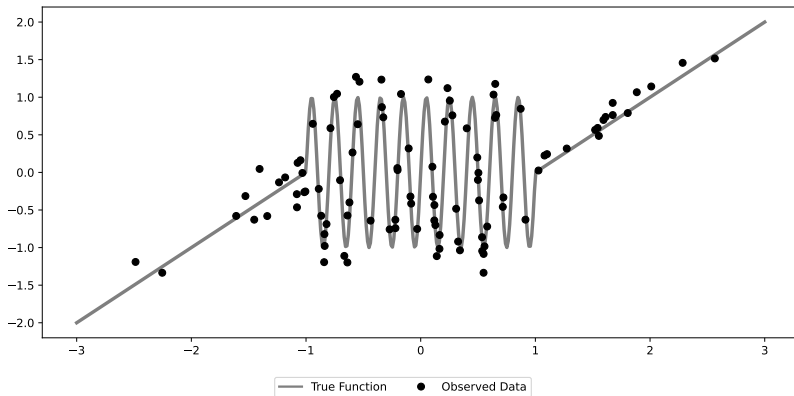
i.e. a larger bandwidth results in a simpler model, and a smaller bandwidth in a more complex model.

We use $1/\sigma$ as a proxy for complexity.

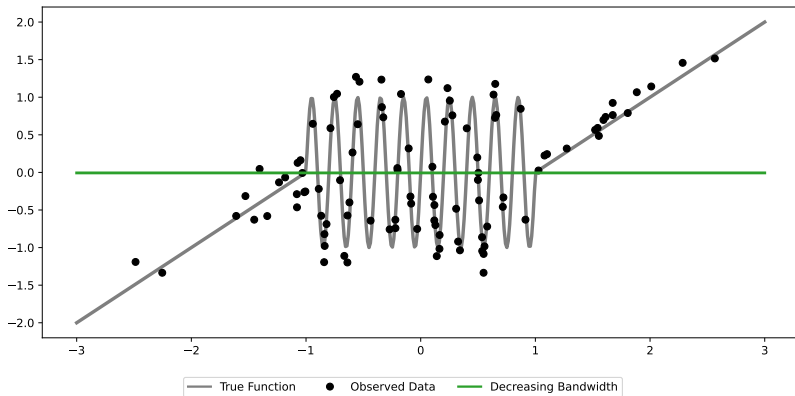
$$\begin{bmatrix} \hat{\mathbf{f}}_{t+1} \\ \hat{\mathbf{f}}_{t+1}^* \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}_t \\ \hat{\mathbf{f}}_t^* \end{bmatrix} - \eta \cdot \begin{bmatrix} \mathbf{K}(\sigma_t) \\ \mathbf{K}^*(\sigma_t) \end{bmatrix} (\hat{\mathbf{f}}_t - \mathbf{y})$$

Idea: Start with a kernel with large bandwidth (a simple model).
Gradually decrease the bandwidth towards zero during training.

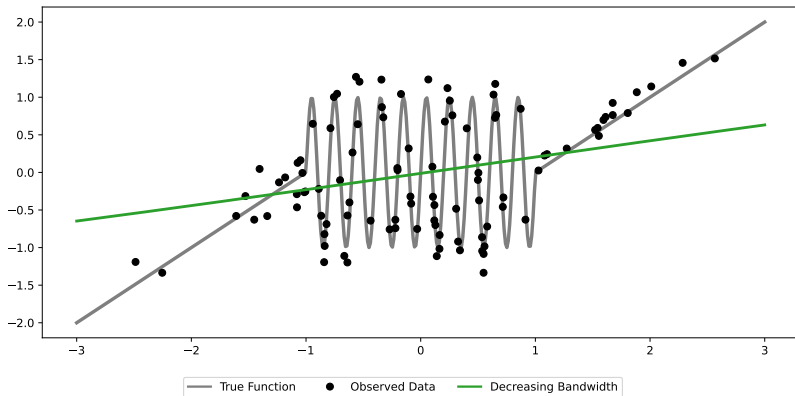
Kernel Regression with Gradient Descent and Non-Constant Kernels



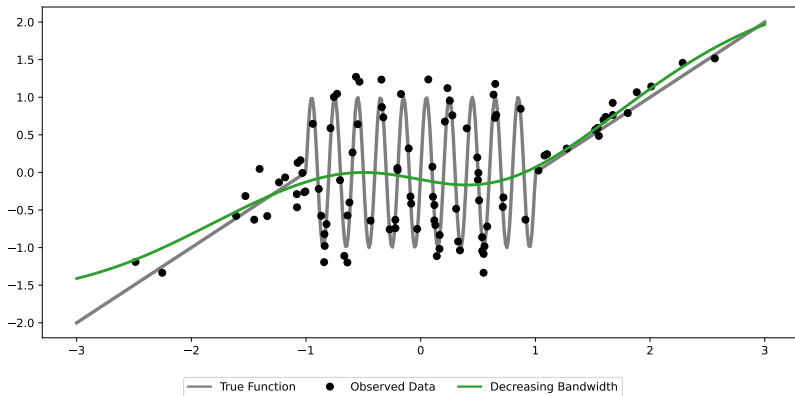
Kernel Regression with Gradient Descent and Non-Constant Kernels



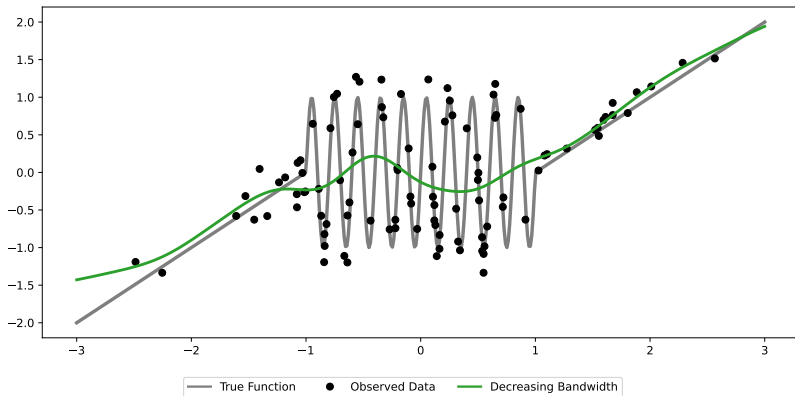
Kernel Regression with Gradient Descent and Non-Constant Kernels



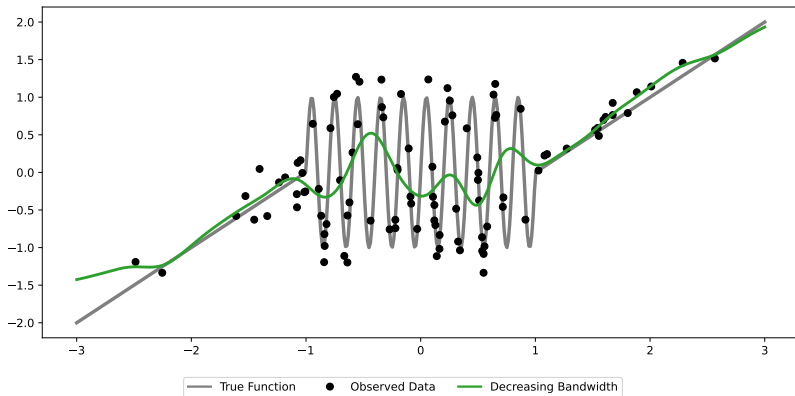
Kernel Regression with Gradient Descent and Non-Constant Kernels



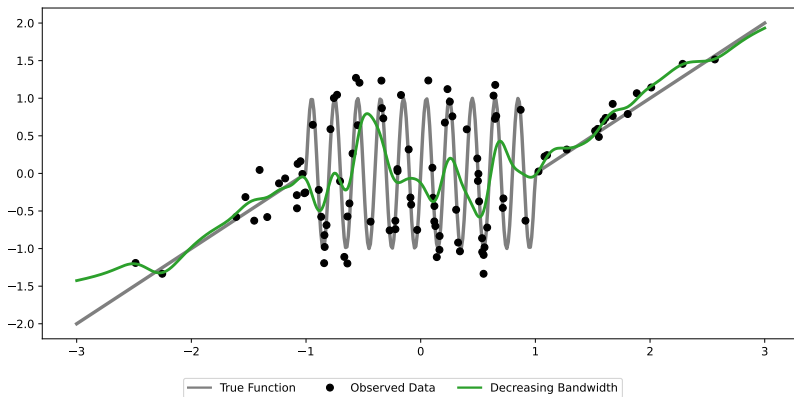
Kernel Regression with Gradient Descent and Non-Constant Kernels



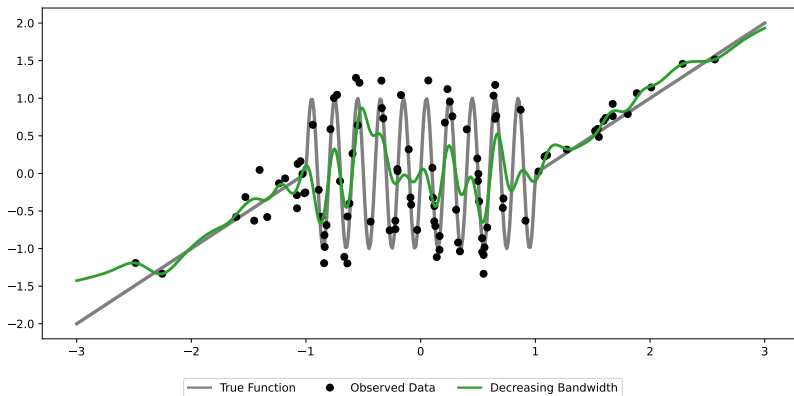
Kernel Regression with Gradient Descent and Non-Constant Kernels



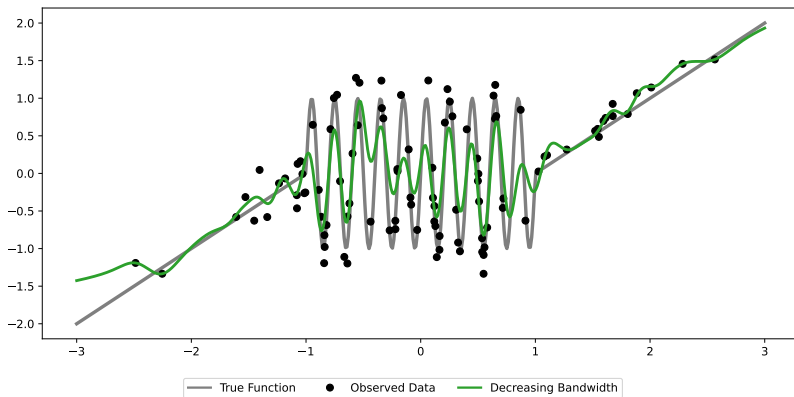
Kernel Regression with Gradient Descent and Non-Constant Kernels



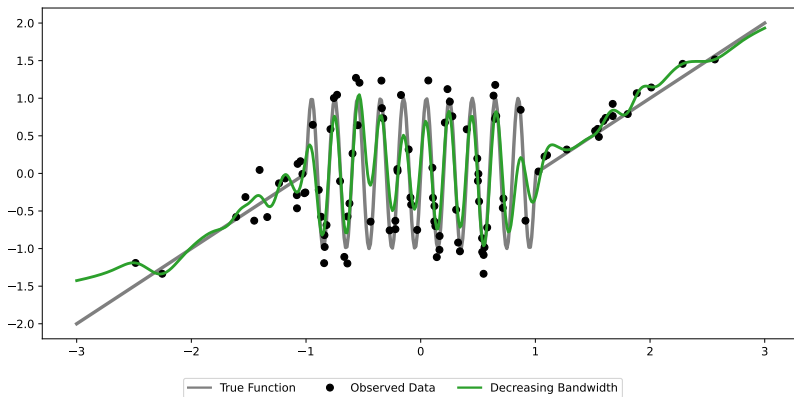
Kernel Regression with Gradient Descent and Non-Constant Kernels



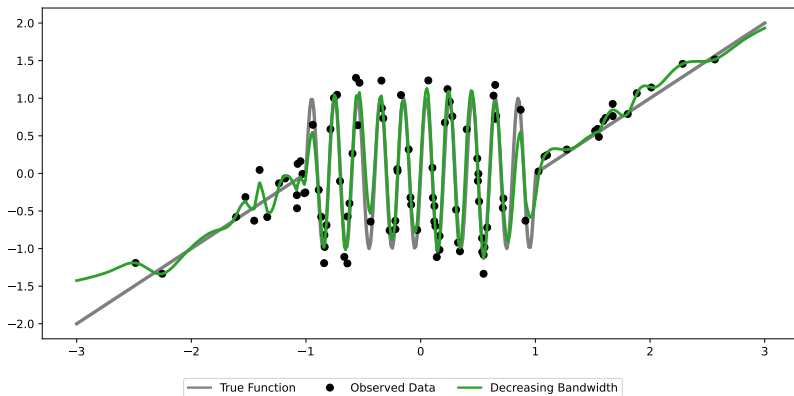
Kernel Regression with Gradient Descent and Non-Constant Kernels



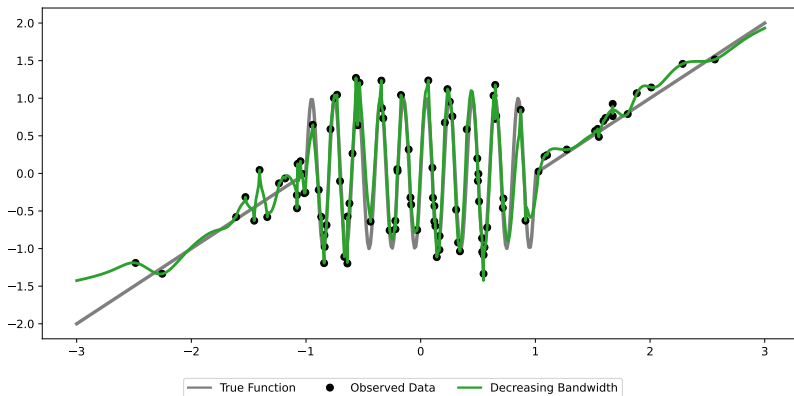
Kernel Regression with Gradient Descent and Non-Constant Kernels



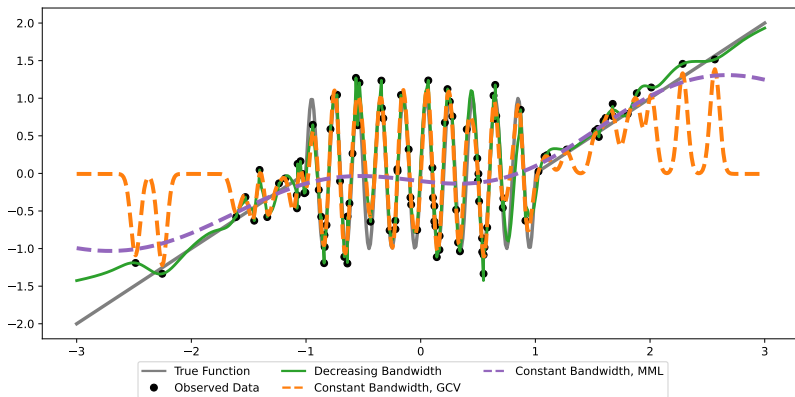
Kernel Regression with Gradient Descent and Non-Constant Kernels



Kernel Regression with Gradient Descent and Non-Constant Kernels



Kernel Regression with Gradient Descent and Non-Constant Kernels



Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

A too simple model generalizes poorly.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

- A too simple model generalizes poorly.

- A model of appropriate complexity generalizes well.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

- A too simple model generalizes poorly.

- A model of appropriate complexity generalizes well.

- A too complex model generalizes poorly (overfitting).

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

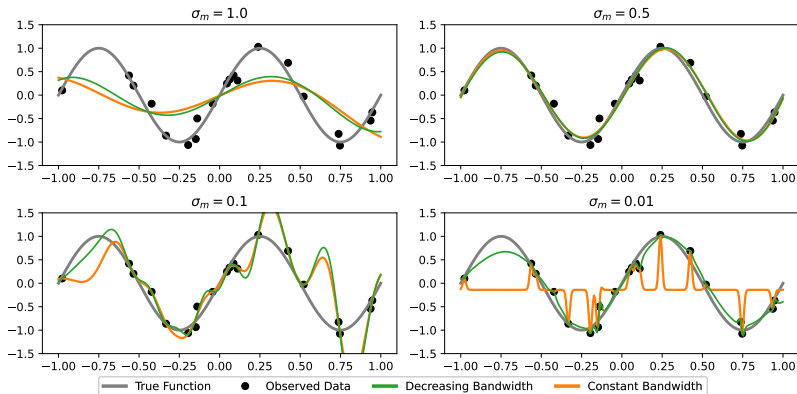
A too simple model generalizes poorly.	}	Classical statistical knowledge
A model of appropriate complexity generalizes well.		
A too complex model generalizes poorly (overfitting).		

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Generalization as function of model complexity:

A too simple model generalizes poorly.	}	Classical statistical knowledge
A model of appropriate complexity generalizes well.		
A too complex model generalizes poorly (overfitting).		
An extremely complex model generalizes well (double descent).		

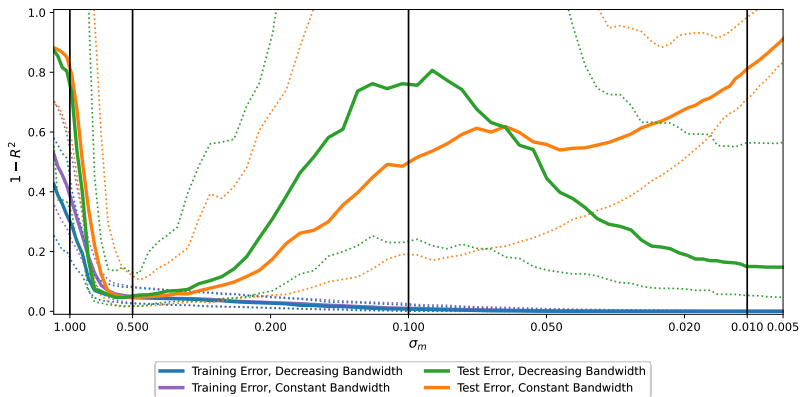
Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent



σ_m : Minimum allowed bandwidth when decreasing the bandwidth.

Employed bandwidth when using a constant bandwidth.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent



Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Proposition (Simplified)

$$|\hat{f}(\mathbf{x}^*, t, \sigma_m)| \leq \min \left(\overline{\sigma^{-1}}(\sigma_m), \overline{k_2^*}(\sigma_m) \right) \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

where $\overline{\sigma^{-1}}(\sigma_m)$ increases with model complexity and $\overline{k_2^*}(\sigma_m)$ decreases with model complexity.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Proposition (Simplified)

$$|\hat{f}(\mathbf{x}^*, t, \sigma_m)| \leq \min \left(\overline{\sigma^{-1}}(\sigma_m), \overline{k_2^*}(\sigma_m) \right) \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

where $\overline{\sigma^{-1}}(\sigma_m)$ increases with model complexity and $\overline{k_2^*}(\sigma_m)$ decreases with model complexity.

Low complexity: $\overline{\sigma^{-1}}(\sigma_m)$ is small $|\hat{f}|$ is too small.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Proposition (Simplified)

$$|\hat{f}(\mathbf{x}^*, t, \sigma_m)| \leq \min \left(\overline{\sigma}^{-1}(\sigma_m), \overline{k}_2^*(\sigma_m) \right) \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

where $\overline{\sigma}^{-1}(\sigma_m)$ increases with model complexity and $\overline{k}_2^*(\sigma_m)$ decreases with model complexity.

Low complexity: $\overline{\sigma}^{-1}(\sigma_m)$ is small $|\hat{f}|$ is too small.

Moderate complexity: $\overline{\sigma}^{-1}(\sigma_m)$ is moderate $|\hat{f}|$ is appropriate.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Proposition (Simplified)

$$|\hat{f}(\mathbf{x}^*, t, \sigma_m)| \leq \min \left(\overline{\sigma}^{-1}(\sigma_m), \overline{k}_2^*(\sigma_m) \right) \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

where $\overline{\sigma}^{-1}(\sigma_m)$ increases with model complexity and $\overline{k}_2^*(\sigma_m)$ decreases with model complexity.

Low complexity: $\overline{\sigma}^{-1}(\sigma_m)$ is small $|\hat{f}|$ is too small.

Moderate complexity: $\overline{\sigma}^{-1}(\sigma_m)$ is moderate $|\hat{f}|$ is appropriate.

High complexity: $\overline{\sigma}^{-1}(\sigma_m)$ is large $|\hat{f}|$ is too large.

Kernel Regression with Gradient Descent and Non-Constant Kernels, Double Descent

Proposition (Simplified)

$$|\hat{f}(\mathbf{x}^*, t, \sigma_m)| \leq \min \left(\overline{\sigma^{-1}}(\sigma_m), \overline{k_2^*}(\sigma_m) \right) \cdot t \cdot C(k(\cdot), \mathbf{X}, \mathbf{y})$$

where $\overline{\sigma^{-1}}(\sigma_m)$ increases with model complexity and $\overline{k_2^*}(\sigma_m)$ decreases with model complexity.

Low complexity:	$\overline{\sigma^{-1}}(\sigma_m)$ is small	$ \hat{f} $ is too small.
Moderate complexity:	$\overline{\sigma^{-1}}(\sigma_m)$ is moderate	$ \hat{f} $ is appropriate.
High complexity:	$\overline{\sigma^{-1}}(\sigma_m)$ is large	$ \hat{f} $ is too large.
Very high complexity:	$\overline{k_2^*}(\sigma_m)$ is moderate	$ \hat{f} $ is appropriate.

Neural Tangent Kernel Gradient Flow

Neural Tangent Kernel Gradient Flow

$$\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y})$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{f}(t + \Delta t) \\ \hat{f}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{f}(t) \\ \hat{f}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{f}(t) - y) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{f}(t + \Delta t) \\ \hat{f}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{f}(t) \\ \hat{f}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{f}(t) - y) \end{aligned}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial t} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t} \end{bmatrix}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} K(t) \\ K^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial t} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t} \end{bmatrix} \stackrel{\text{C.R.}}{=} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial \hat{\boldsymbol{\theta}}(t)}{\partial t}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial t} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t} \end{bmatrix} \stackrel{\text{C.R.}}{=} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial \hat{\boldsymbol{\theta}}(t)}{\partial t} \stackrel{\text{G.D.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\boldsymbol{\theta}}(t)}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \left[\frac{\frac{\partial \hat{\mathbf{f}}(t)}{\partial t}}{\frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t}} \right] &\stackrel{\text{C.R.}}{=} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial \hat{\boldsymbol{\theta}}(t)}{\partial t} \stackrel{\text{G.D.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\boldsymbol{\theta}}(t)} \\ &\stackrel{\text{C.R.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \left(\frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \right)^\top \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)} \end{aligned}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \left[\frac{\frac{\partial \hat{\mathbf{f}}(t)}{\partial t}}{\frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t}} \right] &\stackrel{\text{C.R.}}{=} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial \hat{\boldsymbol{\theta}}(t)}{\partial t} \stackrel{\text{G.D.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\boldsymbol{\theta}}(t)} \\ &\stackrel{\text{C.R.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \left(\frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \right)^\top \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)} \\ &\stackrel{\text{def}}{=} - \begin{bmatrix} \boldsymbol{\Phi}(t) \boldsymbol{\Phi}(t)^\top \\ \boldsymbol{\Phi}^*(t) \boldsymbol{\Phi}(t)^\top \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)} \end{aligned}$$

Neural Tangent Kernel Gradient Flow

$$\begin{aligned} \begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} &= \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} - \Delta t \cdot \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \\ \Leftrightarrow \frac{\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix}}{\Delta t} &= - \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\hat{\mathbf{f}}(t) - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial t} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial t} \end{bmatrix} &\stackrel{\text{C.R.}}{=} \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial \hat{\boldsymbol{\theta}}(t)}{\partial t} \stackrel{\text{G.D.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\boldsymbol{\theta}}(t)} \\ &\stackrel{\text{C.R.}}{=} - \begin{bmatrix} \frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \\ \frac{\partial \hat{\mathbf{f}}^*(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \end{bmatrix} \cdot \left(\frac{\partial \hat{\mathbf{f}}(t)}{\partial \hat{\boldsymbol{\theta}}(t)} \right)^\top \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)} \\ &\stackrel{\text{def}}{=} - \begin{bmatrix} \boldsymbol{\Phi}(t) \boldsymbol{\Phi}(t)^\top \\ \boldsymbol{\Phi}^*(t) \boldsymbol{\Phi}(t)^\top \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)} = - \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} \cdot \frac{\partial L(\hat{\mathbf{f}}(t))}{\partial \hat{\mathbf{f}}(t)}. \end{aligned}$$

Improved Generalization by Neural Tangent Control

- Kernel gradient descent performs best when the bandwidth decreases toward zero,

Improved Generalization by Neural Tangent Control

- Kernel gradient descent performs best when the bandwidth decreases toward zero,
- that is, when $\begin{bmatrix} K \\ K^* \end{bmatrix}$ goes to $\begin{bmatrix} I_{n \times n} \\ \mathbf{0}_{n^* \times n} \end{bmatrix}$.

Improved Generalization by Neural Tangent Control

- Kernel gradient descent performs best when the bandwidth decreases toward zero,
- that is, when $\begin{bmatrix} K \\ K^* \end{bmatrix}$ goes to $\begin{bmatrix} I_{n \times n} \\ \mathbf{0}_{n^* \times n} \end{bmatrix}$.
- Can this behaviour be enforced for the neural tangent kernel?

Improved Generalization by Neural Tangent Control

- Kernel gradient descent performs best when the bandwidth decreases toward zero,
- that is, when $\begin{bmatrix} K \\ K^* \end{bmatrix}$ goes to $\begin{bmatrix} I_{n \times n} \\ \mathbf{0}_{n^* \times n} \end{bmatrix}$.
- Can this behaviour be enforced for the neural tangent kernel?
- Work in progress, but it seems so...

Conclusions

Kernels models

Conclusions

Kernels models

- are non-linear in data.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.
- Kernel Ridge Regression.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.
- Kernel Ridge Regression.

Kernel regression with gradually increased model complexity

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.
- Kernel Ridge Regression.

Kernel regression with gradually increased model complexity

- reduces the need for hyper parameter selection.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.
- Kernel Ridge Regression.

Kernel regression with gradually increased model complexity

- reduces the need for hyper parameter selection.
- exhibits a double descent behavior.

Conclusions

Kernels models

- are non-linear in data.
- are linear in parameters (and thus convex with closed-form solutions).
- correspond to a (possibly infinite dimensional) feature expansion.

Examples:

- Kernel Principal Component Analysis.
- Kernel Ridge Regression.

Kernel regression with gradually increased model complexity

- reduces the need for hyper parameter selection.
- exhibits a double descent behavior.
- is generalizable to any parametric model trained with gradient descent.

The End

Thank you!

Alternative Interpretation

- In forward stagewise additive modelling, model m is used to fit the residuals of model $m - 1$.

$$\hat{\mathbf{f}}'_m(\mathbf{X}) = \operatorname{argmin}_{\mathbf{f}'_m \in \mathcal{F}_m} L\left(\mathbf{f}'_m(\mathbf{X}), \mathbf{y} - \hat{\mathbf{f}}_{m-1}(\mathbf{X})\right)$$
$$\hat{\mathbf{f}}_m(\mathbf{X}) = \hat{\mathbf{f}}_{m-1}(\mathbf{X}) + \hat{\mathbf{f}}'_m(\mathbf{X}),$$

- If we interpret every bandwidth update during training as a model change, then kernel gradient descent with decreasing bandwidth is forward stagewise additive modelling with increasing model complexity.

Algorithm

Algorithm 1 Kernel Gradient Descent with Decreasing Bandwidth

Input: Training data, (\mathbf{X}, \mathbf{y}) . Prediction covariates, \mathbf{X}^* . Initial bandwidth, σ_0 . Prior $\mu(\mathbf{x})$.
 Minimum allowed bandwidth, σ_m . Step-size, Δt . Minimum R^2 speed, v_{R^2} . Maximum R^2 , R_{\max}^2 .

Output: Vector of predictions $[\hat{\mathbf{f}}(t)^\top, \hat{\mathbf{f}}^*(t)^\top]^\top$.

- 1: Initialize $[\hat{\mathbf{f}}(0)^\top, \hat{\mathbf{f}}^*(0)^\top]^\top = [\mu(\mathbf{X})^\top, \mu(\mathbf{X}^*)^\top]^\top$.
- 2: Initialize $\sigma(0) = \sigma_0$, $\mathbf{K} = \mathbf{K}(\sigma_0)$ and $\mathbf{K}^* = \mathbf{K}^*(\sigma_0)$.
- 3: **repeat**
- 4: Calculate $R^2(t) = 1 - \frac{\|\mathbf{y} - \hat{\mathbf{f}}(t)\|_2^2}{\|\mathbf{y} - \bar{\mathbf{y}}\|_2^2}$ and $\frac{\partial R^2(t)}{\partial t} = 1 - \frac{\|\mathbf{y} - \hat{\mathbf{f}}(t)\|_{\mathbf{K}(t)}^2}{\|\mathbf{y} - \bar{\mathbf{y}}\|_2^2}$.
- 5: **if** $\frac{\partial R^2(t)}{\partial t} < v_{R^2}$ **then**
- 6: **repeat**
- 7: Decrease $\sigma(t)$ and recalculate $\mathbf{K}(\sigma(t))$ and $\frac{\partial R^2(t)}{\partial t}$
- 8: **until** $\frac{\partial R(t)^2}{\partial t} \geq v_{R^2}$ **or** $\sigma(t) \leq \sigma_m$.
- 9: Recalculate $\mathbf{K}^*(\sigma(t))$.
- 10: **end if**
- 11: Update $[\hat{\mathbf{f}}(t)^\top, \hat{\mathbf{f}}^*(t)^\top]^\top$ according to

$$\begin{bmatrix} \hat{\mathbf{f}}(t + \Delta t) \\ \hat{\mathbf{f}}^*(t + \Delta t) \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{f}}(t) \\ \hat{\mathbf{f}}^*(t) \end{bmatrix} + \Delta t \begin{bmatrix} \mathbf{K}(t) \\ \mathbf{K}^*(t) \end{bmatrix} (\mathbf{y} - \hat{\mathbf{f}}(t))$$

- 12: **until** $R^2(t) \geq R_{\max}^2$.
-

Two Propositions

Proposition 1.

Let $\overline{k_2^*}$ denote the weighted average of $\|\mathbf{k}^*(\tau)\|_2$ during training, with weight function $\|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2$, and let $\overline{s_{\min}}$ denote the time average of the smallest eigenvalue of $\mathbf{K}(\tau)$ during training, i.e.

$$\overline{k_2^*} := \frac{\int_0^t \|\mathbf{k}^*(\tau)\|_2 \cdot \|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2 d\tau}{\int_0^t \|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2 d\tau} \quad \text{and} \quad \overline{s_{\min}} := \frac{1}{t} \int_0^t s_{\min}(\mathbf{K}(\tau)) d\tau.$$

Then, for the KGF estimate, $\hat{f}_\mu(\mathbf{x}^*, t)$,

$$|\hat{f}_\mu(\mathbf{x}^*, t)| \leq \overline{k_2^*} \cdot \min\left(t, \frac{1}{\overline{s_{\min}}}\right) \cdot \|\mathbf{y}_\mu\|_2. \quad (7)$$

Proposition 3.

Let $\overline{\sigma^{-1}}$ denote the weighted average of the inverse bandwidth during training, with weight function $\|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2$, i.e.,

$$\overline{\sigma^{-1}} := \frac{\int_0^t \frac{1}{\sigma(\tau)} \cdot \|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2 d\tau}{\int_0^t \|\mathbf{y} - \hat{\mathbf{f}}(\tau)\|_2 d\tau}.$$

Then, for a kernel $k(\mathbf{x}, \mathbf{x}', \sigma(\tau)) = k\left(\frac{1}{\sigma(\tau)} \cdot \|\mathbf{x} - \mathbf{x}'\|_\Theta\right)$, with bounded derivative, $|k'(u)| \leq k'_{\max} \forall u$, for the KGF estimate, $\hat{f}_\mu(\mathbf{x}^*, t)$,

$$\left\| \frac{\partial \hat{f}_\mu(\mathbf{x}^*, t)}{\partial \mathbf{x}^*} \right\|_2 \leq \overline{\sigma^{-1}} \cdot \min\left(t, \frac{1}{\overline{s_{\min}}}\right) \cdot \|\mathbf{y}_\mu\|_2 \cdot k'_{\max} \cdot \sqrt{n \cdot \|\Theta\|_2}.$$

Assuming the data is centered, so that $\overline{\mathbf{y}_\mu} = 0$, we further obtain

$$\left| \hat{f}_\mu(\mathbf{x}^*, t) \right| \leq \overline{\sigma^{-1}} \cdot \min\left(t, \frac{1}{\overline{s_{\min}}}\right) \cdot \|\mathbf{y}_\mu\|_2 \cdot k'_{\max} \cdot \sqrt{n \cdot \|\Theta\|_2} \cdot \|\mathbf{x}^* - \mathbf{x}_m\|_2 \quad (10)$$

where $\mathbf{x}_m \in \mathbf{X}$ is the observation furthest away from \mathbf{x}^* .

Double Descent

Table 1: Conceptual sketch of how the bound on $|\hat{f}_\mu|$ from Equation 15, and thus the generalization properties, changes with model complexity, σ_m . The active elements in the minimum functions are marked in bold. In the third and fourth rows, it is not obvious which element is smaller, but this uncertainty does not affect the bound. The constants C_1 and C_2 are omitted to improve readability. For lower model complexities, the bound on the deviation from the prior grows with model complexity, but for very complex models it starts to shrink again.

σ_m	Model Complexity	$\overline{\sigma^{-1}}$	$\overline{k_2^*}$	t	$1/\overline{s_{\min}}$	Bound on $ \hat{f}_\mu $	Generalization
Large	Low	Small	Large	Moderate	Large	$\overline{\sigma^{-1}} \cdot t$ Small	Poor, due to basically predicting the prior
Moderate	Moderate	Moderate	Large	Moderate	Large	$\overline{\sigma^{-1}} \cdot t$ Moderate	Good, due to moderate deviations from the prior
Small	High	Large	Large	Moderate	Large	$\overline{\sigma^{-1}} \cdot t$ Large	Poor, due to extreme predictions
Very small	Very high	Very large	Moderate	Moderate	Moderate	$\overline{k_2^*}/\overline{s_{\min}}$ Moderate	Good, due to moderate deviations from the prior

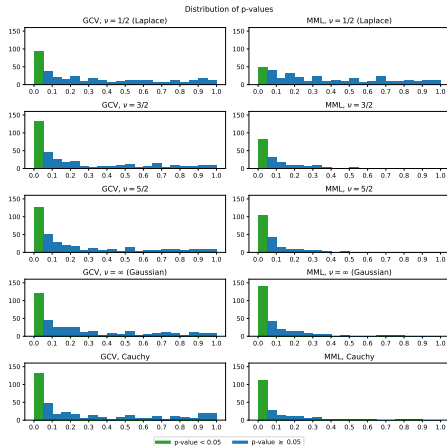
performs excellently on training data but tends to generalize poorly, something that is often referred to as overfitting. However, the wisdom from double descent is that if the model is made even more complex, it may generalize well despite excellent performance on training data.

If we constrain the complexity of the final model by introducing a minimum allowed bandwidth, σ_m , we may, for a constant, fairly long, training time, obtain a double descent behavior in the complexity (i.e. in σ_m). This can qualitatively be seen from the bound on $|\hat{f}_\mu(\mathbf{x}^*)|$ obtained by combining Equations 7 and 10:

$$|\hat{f}_\mu(\mathbf{x}^*, t, \sigma_m)| \leq \min\left(\overline{\sigma^{-1}}(\sigma_m) \cdot C_1, \overline{k_2^*}(\sigma_m)\right) \cdot \min\left(t, \frac{1}{\overline{s_{\min}}(\sigma_m)}\right) \cdot C_2, \quad (15)$$

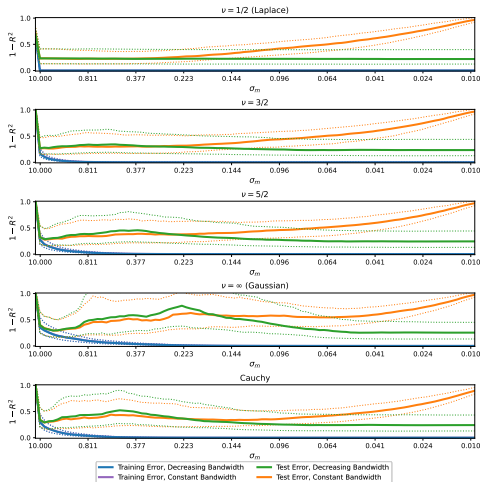
Additional Results

Daily mean temperature in the U.K. during one year, as function of location. 45000 observations, 366 days.
 Wilcoxon Signed-Rank Test, testing whether $R_{\text{test}}^2(\text{decr}) > R_{\text{test}}^2(\text{const})$.



Additional Results

U.K. temperature data



Additional Results

Synthetic data

