

# Machine Learning

## Linear Regression 2

Hiroshi Shimodaira and Hao Tang

March 2026

*Ver. 1.0.1*

# Topics

- Linear regression with feature transformation
- Polynomial regression
- Overfitting in linear regression
- Linear regression with regularisation
- Ridge regression
- Kernel ridge regression

## Augmented vector representation

$$y = \mathbf{w}^\top \mathbf{x} + b = [\mathbf{w}^\top \quad b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}^\top \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \mathbf{w}'^\top \mathbf{x}' \quad (1)$$

- Fitting  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$  is equivalent to appending 1 to  $\mathbf{x}$  and fitting  $f(\mathbf{x}') = \mathbf{w}'^\top \mathbf{x}'$ .
- The 1 can be seen as a feature independent of the input.

## Feature functions

- Suppose we have a data point  $\mathbf{x} = [x_1 \ x_2 \ x_3]^T$ .
- The data point after appending 1 becomes

$$[1 \ x_1 \ x_2 \ x_3]^T \quad (2)$$

- The data point after appending 1 and quadratic terms becomes

$$\phi(\mathbf{x}) = [1 \ x_1 \ x_2 \ x_3 \ x_1x_2 \ x_2x_3 \ x_1x_3 \ x_1^2 \ x_2^2 \ x_3^2]^T \quad (3)$$

- We call  $\phi(\cdot)$  a *feature function*.
- Linear regression with polynomial feature transformation is called *polynomial regression*.

# Linear regression with feature transformation

- In general,  $\phi : \mathbb{R}^d \mapsto \mathbb{R}^D$  can be any function.
- Instead of  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ , we now have  $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ .

- Instead of  $\mathbf{X} = \underbrace{\begin{bmatrix} \dot{\mathbf{x}}_1^\top \\ \vdots \\ \dot{\mathbf{x}}_N^\top \end{bmatrix}}_{N \times d}$ , we have  $\Phi = \underbrace{\begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{bmatrix}}_{N \times D}$

- The mean-squared error can be written as

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i)^2 \quad (4)$$

$$= \frac{1}{N} \|\Phi \mathbf{w} - \mathbf{y}\|_2^2. \quad (5)$$

## Linear regression with feature transformation *(cont.)*

- The cost function can be written compactly as

$$L = \|\Phi \mathbf{w} - \mathbf{y}\|_2^2 \quad (6)$$

$$= (\Phi \mathbf{w} - \mathbf{y})^\top (\Phi \mathbf{w} - \mathbf{y}) \quad (7)$$

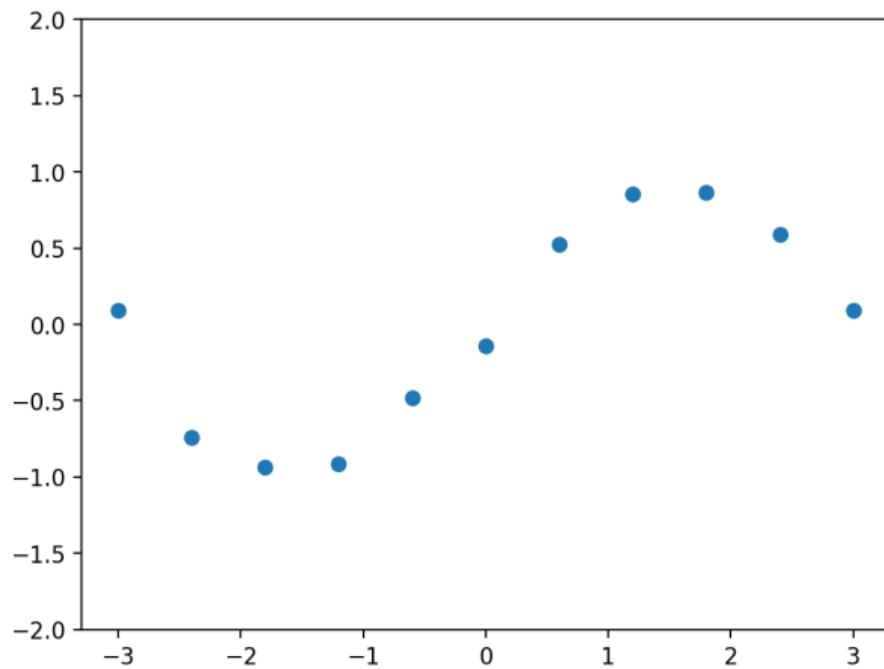
$$= \mathbf{w}^\top \Phi^\top \Phi \mathbf{w} - 2\mathbf{w}^\top \Phi^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}. \quad (8)$$

- Solving the optimal solution gives

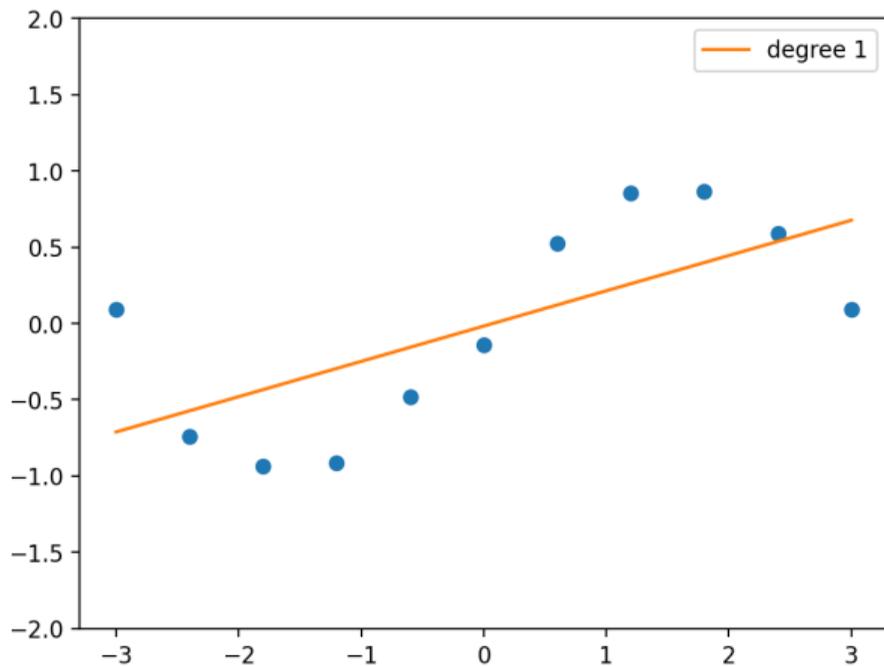
$$\frac{\partial L}{\partial \mathbf{w}} = (\Phi^\top \Phi + (\Phi^\top \Phi)^\top) \mathbf{w} - 2\Phi^\top \mathbf{y} = \mathbf{0}. \quad (9)$$

$$\implies \mathbf{w} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}. \quad (10)$$

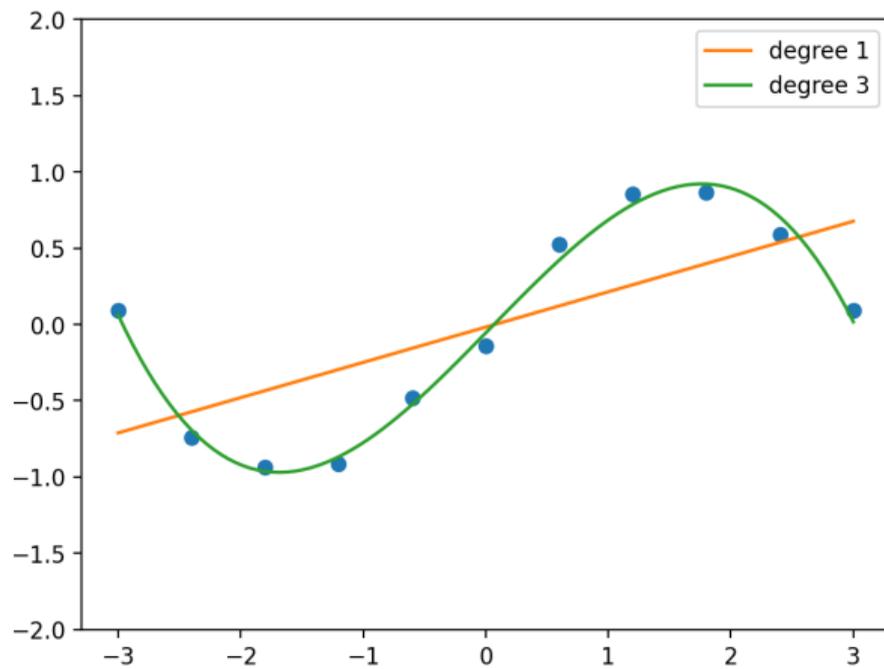
# Examples



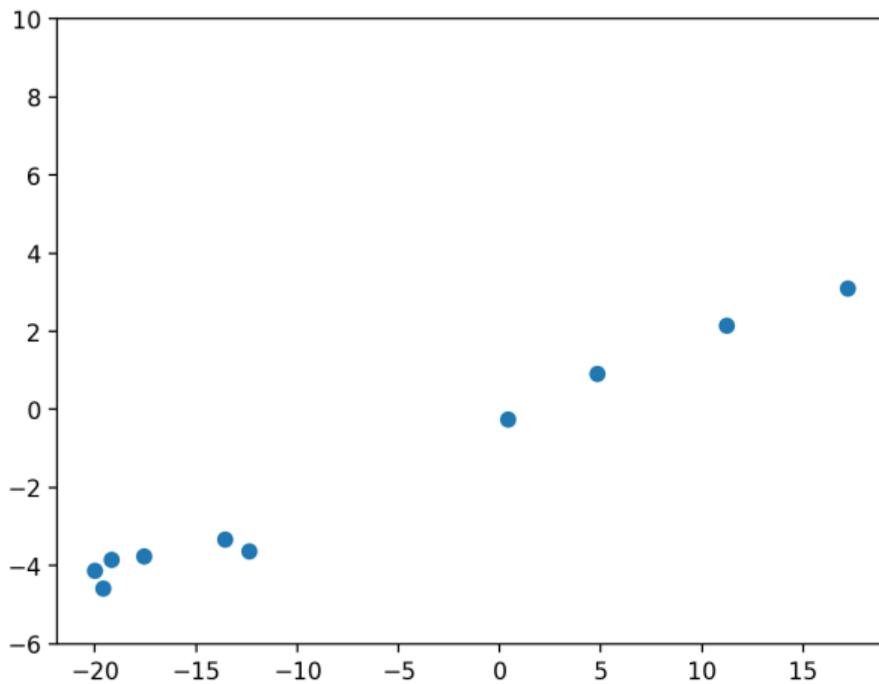
# Examples



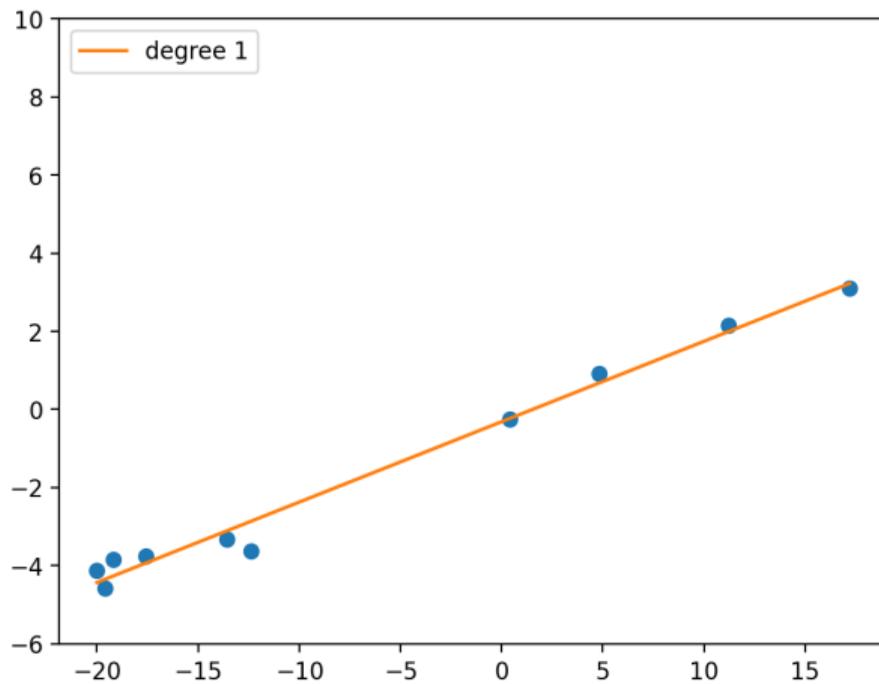
# Examples



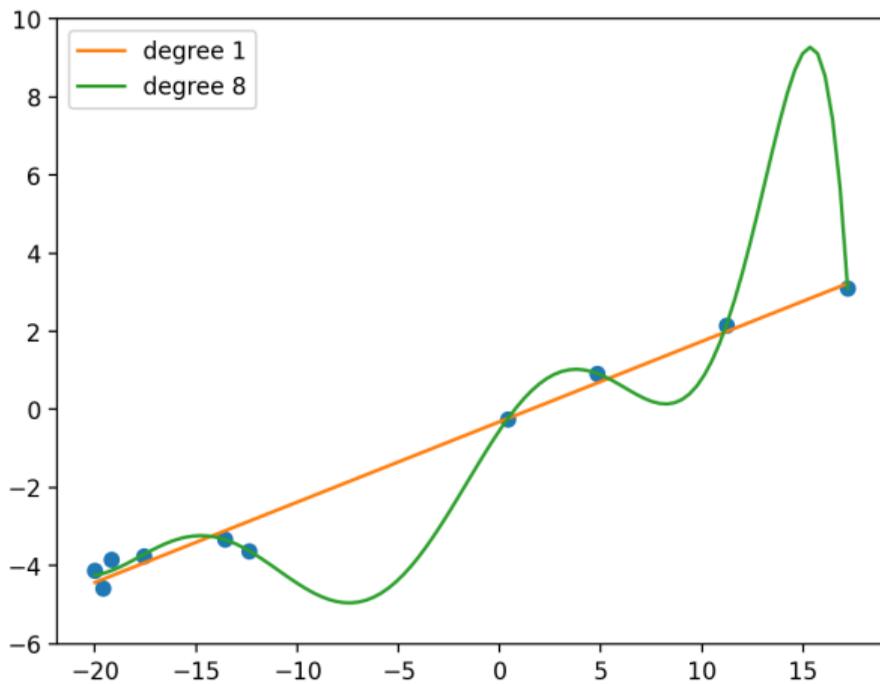
# Examples



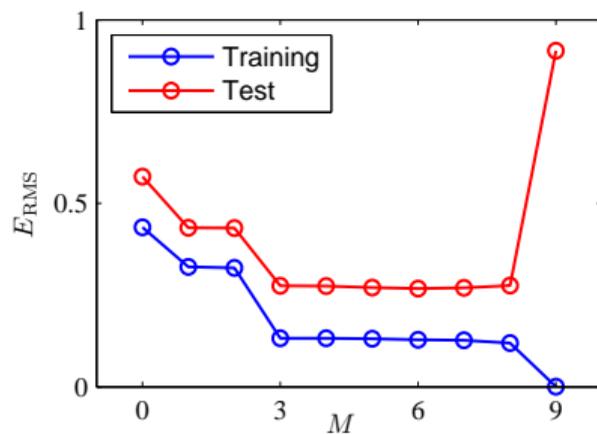
# Examples



# Examples



# Overfitting in linear regression



Credit: C. Bishop, PRML, Figure 1.5

- See slides on generalisation

# Linear regression with regularisation

- Optimisation problem without regularisation

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i)^2 \quad (11)$$

- Optimisation problem with regularisation

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|\mathbf{w}\|^2, \quad \lambda > 0 \quad (12)$$

This is called '*ridge regression*'.

The second term is an  $L^2$ -regulariser.

## Ridge regression: training

$$L = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}^\top \phi(\mathbf{x}_i) - y_i)^2 + \lambda \|\mathbf{w}\|^2 \quad (13)$$

$$= \frac{1}{N} (\Phi \mathbf{w} - \mathbf{y})^\top (\Phi \mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^\top \mathbf{w} \quad (14)$$

$$= \frac{1}{N} \mathbf{w}^\top (\Phi^\top \Phi + N\lambda \mathbf{I}) \mathbf{w} - 2\mathbf{w}^\top \Phi^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \quad (15)$$

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{2}{N} (\Phi^\top \Phi + N\lambda \mathbf{I}) \mathbf{w} - \Phi^\top \mathbf{y} = \mathbf{0} \quad (16)$$

$$\implies (\Phi^\top \Phi + N\lambda \mathbf{I}) \mathbf{w} = \Phi^\top \mathbf{y} \quad (17)$$

$$\mathbf{w} = (\Phi^\top \Phi + N\lambda \mathbf{I})^{-1} \Phi^\top \mathbf{y} \quad (18)$$

## Kernel ridge regression

Using the previous result, the prediction of  $\hat{y}$  given  $\mathbf{x}$  is expressed as

$$\hat{y} = \mathbf{w}^\top \phi(\mathbf{x}) = \mathbf{y}^\top \Phi (\Phi^\top \Phi + N\lambda \mathbf{I})^{-1} \phi(\mathbf{x}) \quad (19)$$

Using the "push-through identity"<sup>(†)</sup> yields:

$$= \mathbf{y}^\top (\Phi \Phi^\top + N\lambda \mathbf{I})^{-1} \Phi \phi(\mathbf{x}) \quad (20)$$

Note that

$$\Phi \Phi^\top = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_N) \\ & \ddots & \\ \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}_N) \end{bmatrix}, \quad \Phi \phi(\mathbf{x}) = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}) \\ \vdots \\ \phi(\mathbf{x}_N)^\top \phi(\mathbf{x}) \end{bmatrix} \quad (21)$$

We can use the kernel trick  $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^\top \phi(\mathbf{v})$  !

$$(\dagger) : (AB + I)^{-1}A = A(BA + I)^{-1}$$

[https://en.wikipedia.org/wiki/Woodbury\\_matrix\\_identity](https://en.wikipedia.org/wiki/Woodbury_matrix_identity)

## Kernel ridge regression (cont.)

$$\hat{y} = \underbrace{\mathbf{y}^\top}_{1 \times N} \underbrace{(\Phi \Phi^\top + N\lambda \mathbf{I})^{-1}}_{N \times N} \underbrace{\Phi}_{N \times D} \underbrace{\phi(\mathbf{x})}_{D \times 1} \quad (22)$$

$$= \underbrace{\mathbf{y}^\top}_{1 \times N} \underbrace{(\mathbf{K}(\mathbf{X}, \mathbf{X}) + N\lambda \mathbf{I})^{-1}}_{N \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (23)$$

where

$$\mathbf{K}(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ & \ddots & \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}, \quad \mathbf{K}(\mathbf{X}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}) \\ \vdots \\ k(\mathbf{x}_N, \mathbf{x}) \end{bmatrix} \quad (24)$$

NB: the computational complexity of matrix multiplication and matrix inversion for  $n \times n$  matrices is approximately  $O(n^3)$ .

## Kernel ridge regression (*cont.*)

- With the kernel trick,

$$\hat{y} = \underbrace{\mathbf{y}^\top}_{1 \times N} \underbrace{(\mathbf{K}(\mathbf{X}, \mathbf{X}) + N\lambda \mathbf{I})^{-1}}_{N \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (25)$$

$$= \underbrace{\boldsymbol{\alpha}^\top}_{1 \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (26)$$

## Kernel ridge regression (*cont.*)

- With the kernel trick,

$$\hat{y} = \underbrace{\mathbf{y}^\top}_{1 \times N} \underbrace{(\mathbf{K}(\mathbf{X}, \mathbf{X}) + N\lambda \mathbf{I})^{-1}}_{N \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (25)$$

$$= \underbrace{\boldsymbol{\alpha}^\top}_{1 \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (26)$$

- Without the kernel trick,

$$\mathbf{w} = \underbrace{(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + N\lambda \mathbf{I})^{-1}}_{D \times D} \underbrace{\boldsymbol{\Phi}^\top}_{D \times N} \underbrace{\mathbf{y}}_{N \times 1} \quad (27)$$

$$\hat{y} = \underbrace{\mathbf{w}^\top}_{1 \times D} \underbrace{\phi(\mathbf{x})}_{D \times 1} \quad (28)$$

## Kernel ridge regression (*cont.*)

- With the kernel trick,

$$\hat{y} = \underbrace{\mathbf{y}^\top}_{1 \times N} \underbrace{(\mathbf{K}(\mathbf{X}, \mathbf{X}) + N\lambda \mathbf{I})^{-1}}_{N \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (25)$$

$$= \underbrace{\boldsymbol{\alpha}^\top}_{1 \times N} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x})}_{N \times 1} \quad (26)$$

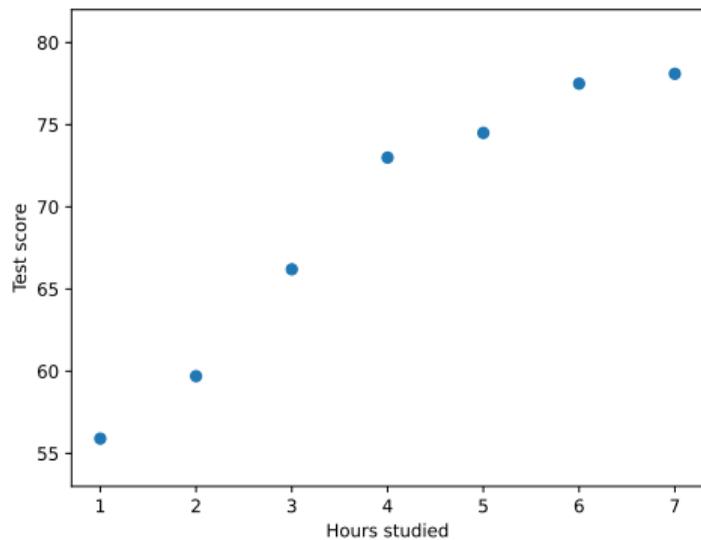
- Without the kernel trick,

$$\mathbf{w} = \underbrace{(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + N\lambda \mathbf{I})^{-1}}_{D \times D} \underbrace{\boldsymbol{\Phi}^\top}_{D \times N} \underbrace{\mathbf{y}}_{N \times 1} \quad (27)$$

$$\hat{y} = \underbrace{\mathbf{w}^\top}_{1 \times D} \underbrace{\phi(\mathbf{x})}_{D \times 1} \quad (28)$$

- If  $D \gg N$ , it's faster and more efficient to use (26) than (28).

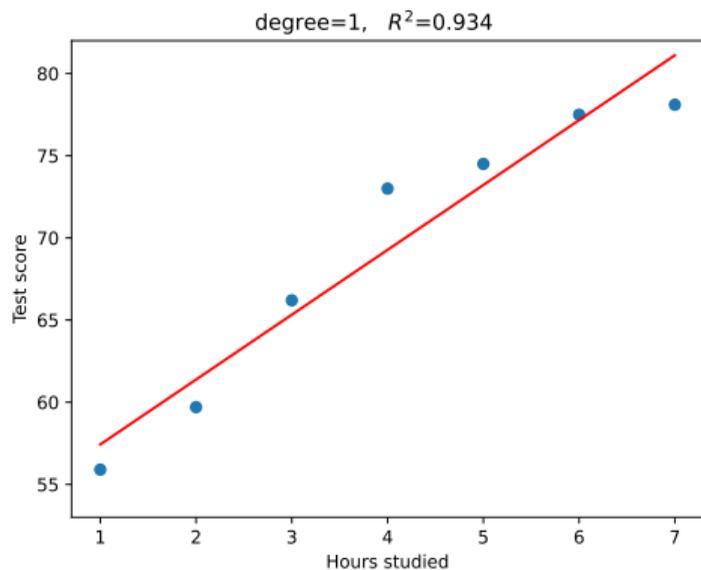
## Example – impact of outliers



Clean data

Data with an outlier

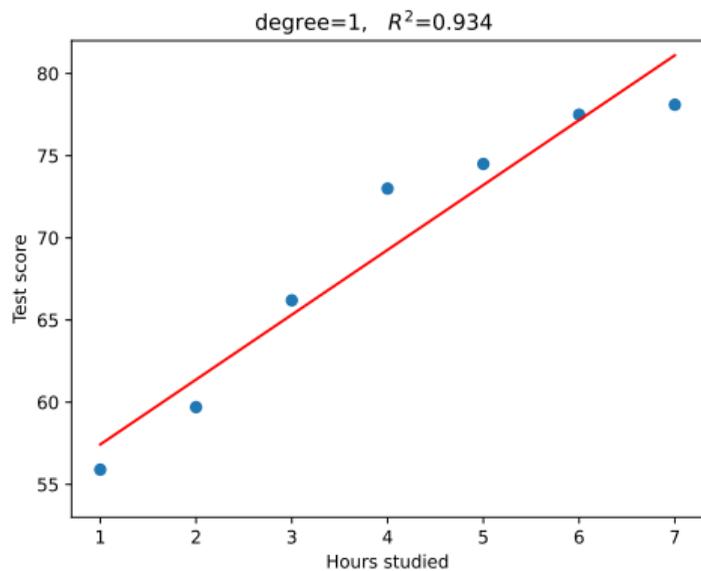
## Example – impact of outliers



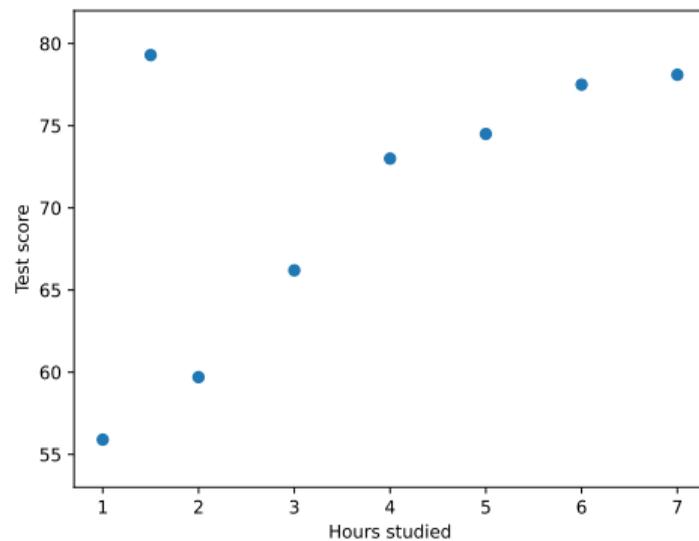
Clean data

Data with an outlier

## Example – impact of outliers

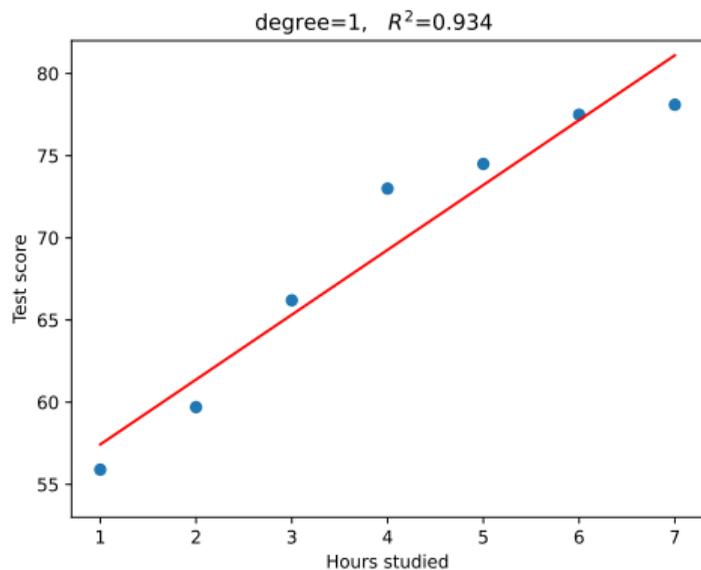


Clean data

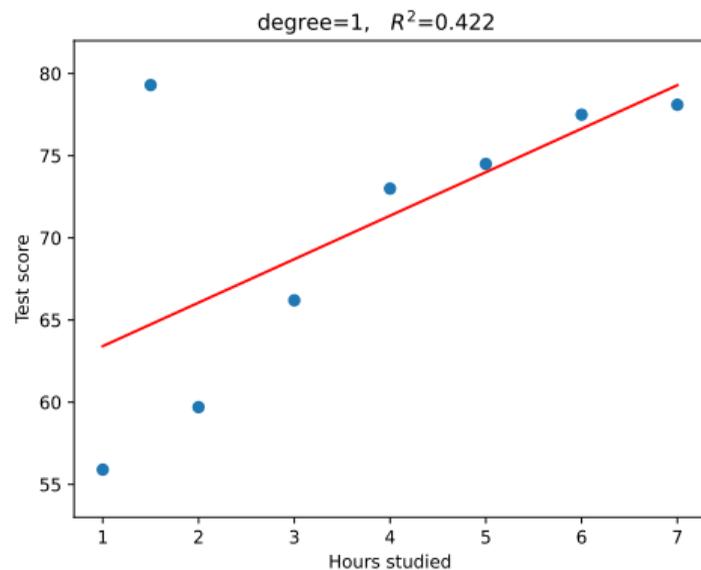


Data with an outlier

## Example – impact of outliers

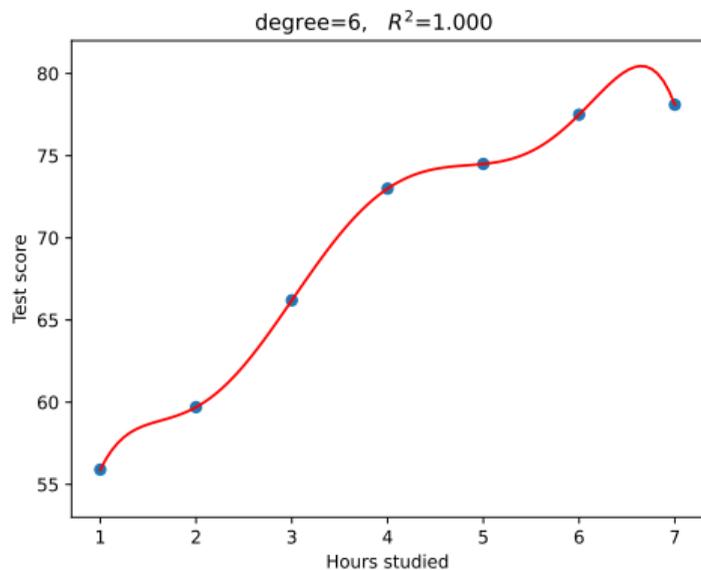


Clean data

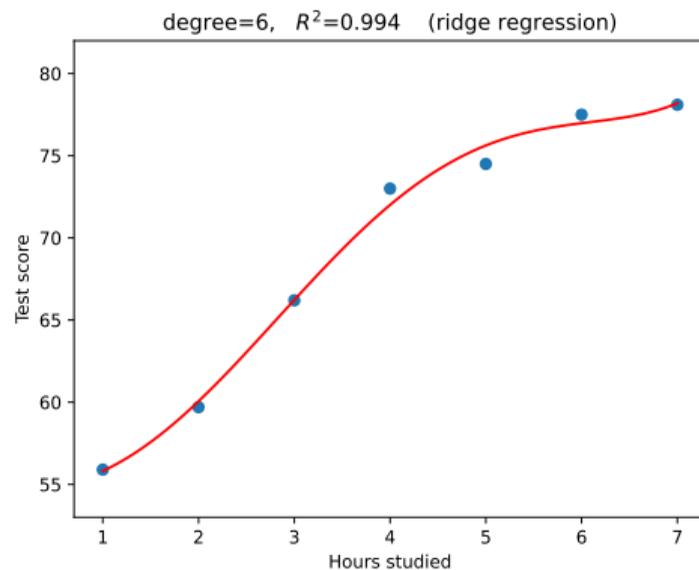


Data with an outlier

## Example – impact of outliers (cont.)

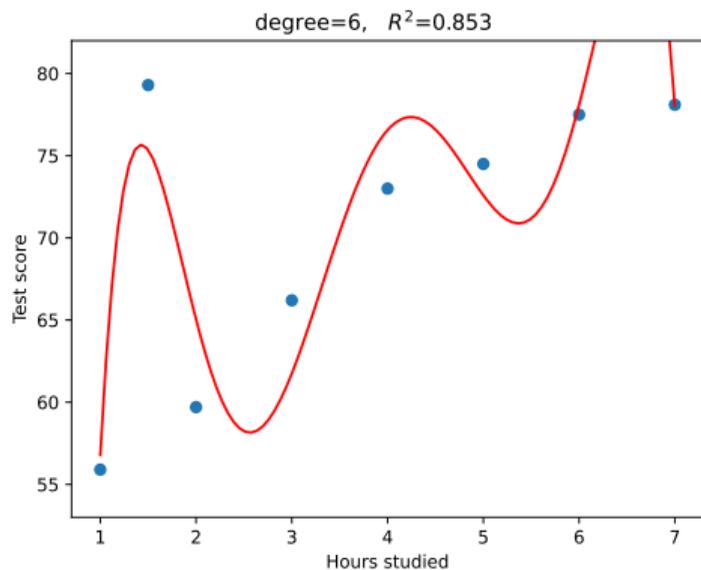


Polynomial regression

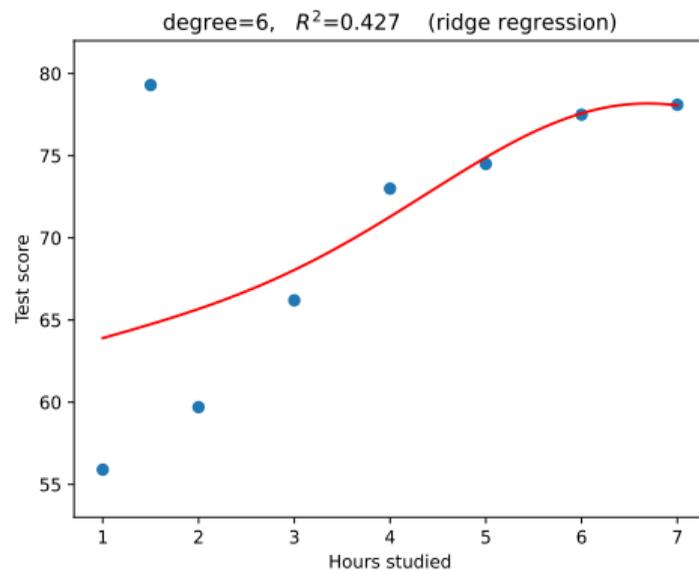


Ridge regression

## Example – impact of outliers (cont.)



Polynomial regression



Ridge regression

## Practical issues

- The runtime is not particularly suitable for large data sets.
- Is it right to evaluate the performance of linear regression on training data alone?
- What if there are outliers?
- What if the features are highly correlated?

## Topics not covered

- Choices of features  $\mathbf{x}$  (feature selection)
- Interpretations of the model parameters  $\mathbf{w}$
- Collinearity
- Heteroscedasticity
- Other linear regression models (e.g., LASSO, Bayesian linear regression)
- Multiple linear regression
- Relationships with neural networks
- Bias-variance trade-off
- Relationships with principal component analysis (PCA)

## Quizzes

1. Show the derivation of (4) to (5).
2. Show the derivation of (8).
3. Show the derivation of (9) to (10).
4. Show the derivation of (19) to (20).