

# Kernel Methods: An Infinity Game

COMS21202, Part III

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

1

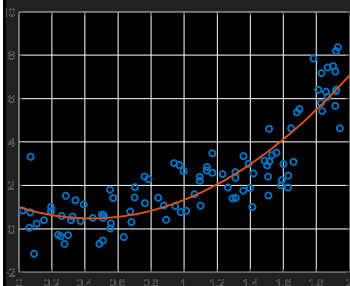
## Objectives

- Applying kernel tricks in LS.
- Knowing common choices of kernel functions.

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

2

**Recall:**  $y = \exp(1.5x - 1) + \epsilon$ ,  $\epsilon \sim N(0,1)$

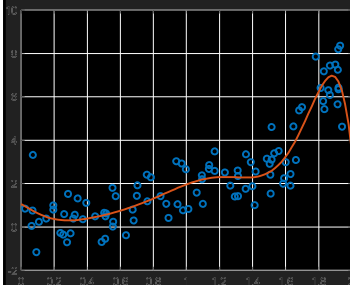


- Polynomial transform with  $b = 2$ .
- Tr. error: 108.97

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

3

**Recall:**  $y = \exp(1.5x - 1) + \epsilon$ ,  $\epsilon \sim N(0,1)$



Polynomial transform with  $b = 8$ .

Tr. error: 78.87

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

4

## Observation

- By increasing output dimension of feature transform  $f(x)$ , we increase the flexibility of  $\hat{y}$ .
- Why don't we keep increasing  $m$  to get a super flexible  $\hat{y}$ ?
  - Do not worry the overfitting now.
- Problem:** large  $m$  causes **numerically issues**.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

5

## Numerical Issues of LS Solution

- Suppose  $f(x) \in R^m$ .
- As we discussed before, if  $m > n$ 
  - $f(X)^T f(X)$  is **singular**.
  - LS solution,  $\hat{\beta} = (f(X)^T f(X))^{-1} f(X)^T y$  cannot be calculated.
  - Shorten  $f(X)$  as  $F$  from now on.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

6

## A Numerical Hack: Regularized LS Solution

- Instead of calculating
  - $\hat{\beta} = (F^T F)^{-1} F^T y$
- We calculate
  - $\hat{\beta} = (F^T F + \lambda I)^{-1} F^T y$
  - Where  $I \in R^{m \times m}$  is identity matrix, improves the invertibility.
  - $\lambda$  is some fixed value, say 0.01.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

7

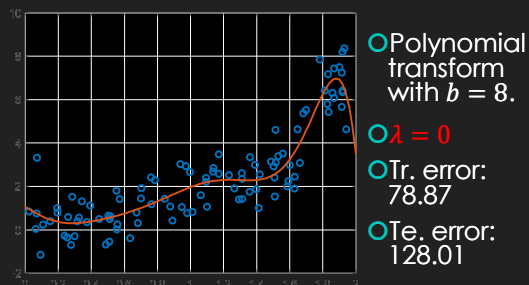
## Regularized LS Solution and Overfitting

- $\lambda I$  helps battle overfitting too (!):
  - Increasing  $\lambda$  decreases the magnitude of  $\hat{\beta}$ , making  $\hat{y}$  approx. a constant 0, which in fact, reduces the flexibility.
  - Show when  $\lambda \rightarrow \infty, \hat{\beta} \approx 0$ .
  - One stone, two birds.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

8

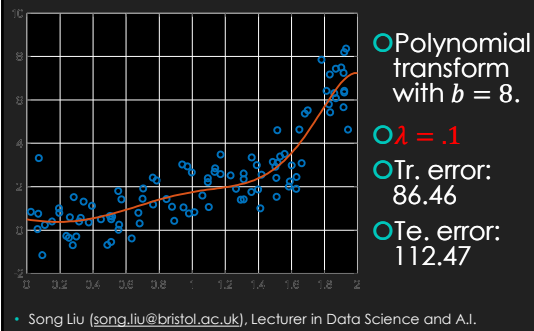
**Example:**  $y = \exp(1.5x - 1) + \epsilon$ ,  
 $\epsilon \sim N(0,1)$



• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

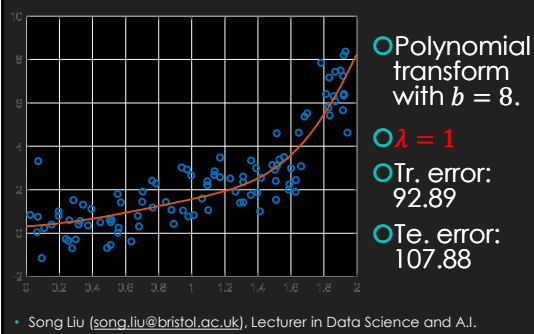
9

**Example:**  $y = \exp(1.5x - 1) + \epsilon$ ,  
 $\epsilon \sim N(0,1)$



10

**Example:**  $y = \exp(1.5x - 1) + \epsilon$ ,  
 $\epsilon \sim N(0,1)$



11

## Regularized LS Solution and Overfitting

- $\lambda$  is called regularization parameter.
  - Should be fixed before fitting.
  - Can be tuned by selecting the value that minimizes testing error.
  - Just like how we select  $b$  for  $f$ .

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

12

### Can we still raise the game?

💡 Can we design  $f(x)$  transforms the original  $x$  into a **infinitely dim. vector**?

- It should create a super flexible  $\hat{y}$ !
- Recall  $\hat{\beta} = (F^T F + \lambda I)^{-1} F^T y$ 
  - Problem: now  $F^T F \in R^{m \times m}$ ,  $m$  is infinity.
  - How do you store  $F$  in computer??

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

13

---

---

---

---

---

---

---

---

### Numerical Hack, #2:

Rewrite Solution using Woodbury identity

- Remarkably,
  - $\hat{\beta} = (F^T F + \lambda I)^{-1} F^T y = F^T (F F^T + \lambda I)^{-1} y$
- Hint, Woodbury identity:
- $(P^{-1} + B^T B)^{-1} B^T = P B^T (B P B^T + I)^{-1}$
- Live demonstration

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

14

---

---

---

---

---

---

---

---

### Numerical Hack, #2:

Rewrite Solution using Woodbury identity

- $\hat{\beta} = F^T (F F^T + \lambda I)^{-1} y$ 
  - Now instead of  $F^T F \in R^{m \times m}$ , we just need to compute  $F F^T \in R^{n \times n}$ .
  - Let  $K := F F^T$ , where
  - $K^{(i,j)} = k(x_i, x_j) = \langle f(x_i), f(x_j) \rangle$ ,
  - i.e.,  $k(x_i, x_j)$  is the inner product of two  $m$  dimensional feature transform.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

15

---

---

---

---

---

---

---

---

### Numerical Hack, #2: Rewrite Solution using Woodbury identity

- $\hat{\mathbf{y}} = \langle \hat{\boldsymbol{\beta}}, \mathbf{f}(\mathbf{x}_0) \rangle$
- $\hat{\mathbf{y}} = \langle \mathbf{f}(\mathbf{x}_0), \mathbf{F}^\top (\mathbf{F}\mathbf{F}^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \rangle$   
 $= \langle \mathbf{f}(\mathbf{x}_0) \mathbf{F}^\top, (\mathbf{F}\mathbf{F}^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \rangle$
- Rewrite  $\mathbf{f}(\mathbf{x}_0) \mathbf{F}^\top$  as  $\mathbf{k} \in \mathbb{R}^n$  we can see
  - $\mathbf{k}^{(i)} = k(\mathbf{x}_0, \mathbf{x}_i) = \langle \mathbf{f}(\mathbf{x}_0), \mathbf{f}(\mathbf{x}_i) \rangle$
  - Verify this your self!

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

16

### Numerical Hack, #3: Evaluating only the Inner Products

- $\hat{\mathbf{y}} := \mathbf{k}(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$
- Note how  $\mathbf{f}(\mathbf{x})$  only appears in the form of inner products!
- 💡 Even if cannot write  $\mathbf{f}(\mathbf{x})$  explicitly, we may still compute its inner product!
  - design "an inner product function  $k$ " mimics behaviour of inner product .
  - Forget about the existence of  $\mathbf{f}$ !

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

17

### Numerical Hack, #3: Evaluating only the Inner Products

- It turns out, you **cannot** pick inner product function  $k$  arbitrarily.
  - Must "behaves like" a inner product.
- However, there are many **known choices** of  $k$  corresponds to inner products of powerful, even infinite dimensional feature transform  $\mathbf{f}$ .
  - **Even** if we cannot write  $\mathbf{f}$  down.

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

18

## Kernel Function

- Our inner product function  $k(x_i, x_j)$  is called **kernel function** in machine learning literatures.
- If an explicit  $f$  can be derived from  $k$ ,
  - We say,  $k$  induces feature transform  $f$ .

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

19

---

---

---

---

---

---

---

---

## The History of Kernel Methods

- Kernel methods were extremely important research topics in machine learning community in the early 2000s.
- It is now referred as "shallow methods", in comparison to deep neural network models.
- It still enjoys great popularity for its simple mathematical expressions and power to represent extremely complex model.

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

20

---

---

---

---

---

---

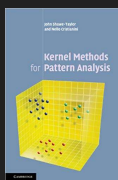
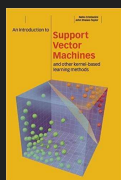
---

---

## Kernel @ Bristol



- Prof. Nello Cristianini at EngMath is one of the world renowned leading scientists in kernel methods.



• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

21

---

---

---

---

---

---

---

---

## Choices of $k$

- Linear kernel function:
  - $k(x_i, x_j) := \langle x_i, x_j \rangle$
  - Implicit feature transform  $f(x) = x$ .
- Polynomial kernel function with degree  $b$ :
  - $k(x_i, x_j) := (\langle x_i, x_j \rangle + 1)^b$
- PC: write down induced  $f(x)$  by polynomial kernels  $b = 2$ .

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

22

---

---

---

---

---

---

---

---

## Choices of $k$

- RBF (or Gaussian) kernel:
  - $k(x_i, x_j) := \exp\left(-\frac{\|x_i - x_j\|^2}{\sigma^2}\right)$
  - $f(x)$  induced by  $k$  is **infinitely dimensional**!
  - $\sigma$  is chosen before fitting.
  - Best  $\sigma$  is chosen by minimizing testing error.
- Déjà vu?

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

23

---

---

---

---

---

---

---

---

## Choices of $k$

- How do I pick  $k$ ?
  - Depending on your learning task.
    - e.g., linear/poly kernels are frequently used in natural language processing.
  - Depending on your dataset.
    - e.g., some kernels are even defined for structural inputs, such as strings or graphs.
  - Domain knowledge matters!!
- RBF kernel is a good all-rounded choice for  $x \in \mathbb{R}^d$ .

• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

24

---

---

---

---

---

---

---

---



## Implementation of Kernel LS

- Recall:  $\hat{\mathbf{y}} := \mathbf{k}(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$
- Computational cost
  - $\mathbf{K}$ :  $O(n^2)$
  - $(\mathbf{K} + \lambda \mathbf{I})^{-1}$ : Usually  $O(n^3)$
  - Kernel methods though flexible, is computationally demanding for large  $n$ .

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

25

---

---

---

---

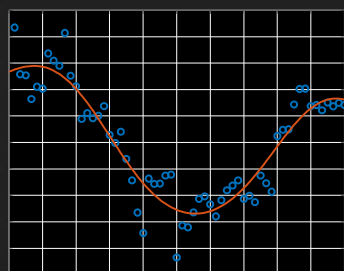
---

---

---

---

## Example: Apple Stock Price, Feb 2019



- RBF kernel
- $\sigma = 0.2121$
- $\lambda = 0.1$ .
- Tr error: 833.58

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

26

---

---

---

---

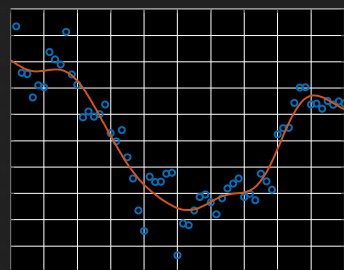
---

---

---

---

## Example: Apple Stock Price, Feb 2019



- RBF kernel
- $\sigma = 0.106$
- $\lambda = 0.1$ .
- Tr error: 666.20

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

27

---

---

---

---

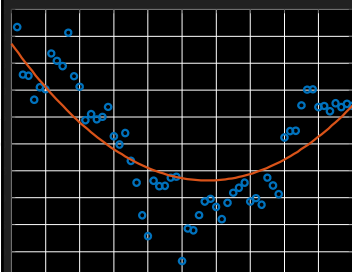
---

---

---

---

### Example: Apple Stock Price, Feb 2019



- Poly. kernel
- $b = 2$
- $\lambda = 0.1$ .
- Tr error: 2068.1

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

28

---

---

---

---

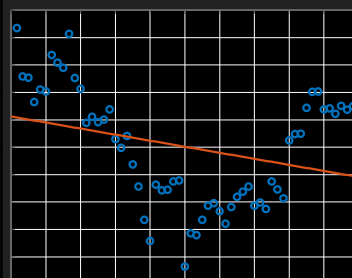
---

---

---

---

### Example: Apple Stock Price, Feb 2019



- Linear kernel
- $\lambda = 0.1$ .
- Tr error: 5964

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

29

---

---

---

---

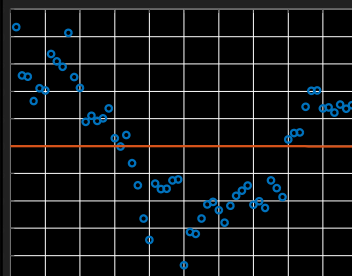
---

---

---

---

### Example: Apple Stock Price, Feb 2019



- Linear kernel
- $\lambda = 1000$ .
- Tr error: 6597

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

30

---

---

---

---

---

---

---

---

## Conclusion

- Kernel methods transform original data point into higher dimensional (potentially **infinitely dim.**) feature vectors.
  - We get super flexible  $\hat{y}$ .
  - Regularization can ease the overfitting caused by flexibility.
- Computation of inf. dimensional features is made possible by kernel trick.
- Important kernel functions:
  - Linear, polynomial, RBF.

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

31

## Proper Names

- Numerical Hack #1 is called **Regularization** in statistics, usually used when handling high dimensional data.
- Numerical Hack #2,3 are called **"kernel tricks"**, usually used for hiding  $f(x)$  inside inner products.
  - Other types of kernel tricks exist.

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

32

- We will swap next week **Wednesday** class with **Monday** class (?)

• Song Liu ([song.liu@bristol.ac.uk](mailto:song.liu@bristol.ac.uk)), Lecturer in Data Science and A.I.

33