

Features: Representing your data

COMS21202, Part III

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

1

Introduction

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

2

Machine Learning Pipeline

```
graph TD; A[Data Acquisition<br/>sampling/quantization] --> B[Feature Engineering<br/>representing your data]; B --> C[Training algorithms<br/>classification/regression/clustering...]; C --> D[Prediction/Inference<br/>Making decisions, cats or dogs?];
```

Part I

Part I, II

We are here, Part III

3

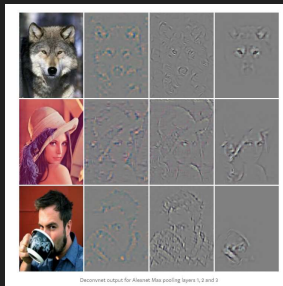
How does machine see the world?

- Machine does not see the world in the same way we do.
 - It does not need to.
 - It only needs the representation of info to perform its task.

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

4

How does machine learning algorithm see the world?



- Visualization of layers in Alexnet.
- Zeiler and Fergus, ECCV 2014

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

5

Turning Data into Features

- Modern machine learning **rarely** uses **raw data** input to perform learning tasks.
- Raw input is usually transformed into a more powerful representation: **features**.
- This procedure of representing data using features is usually referred as **feature engineering** in literatures.

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

6

Feature Engineering

- **Task:** finding a feature **transform function** $f(x)$, which takes a **d -dimensional raw input x** and outputs a **m -dimensional feature vector**.
- Feature function f is the medium through which your learning algorithm interacts with your data.
- Let us put feature engineering in the context of **Least squares**.

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

7

An Appetizer

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

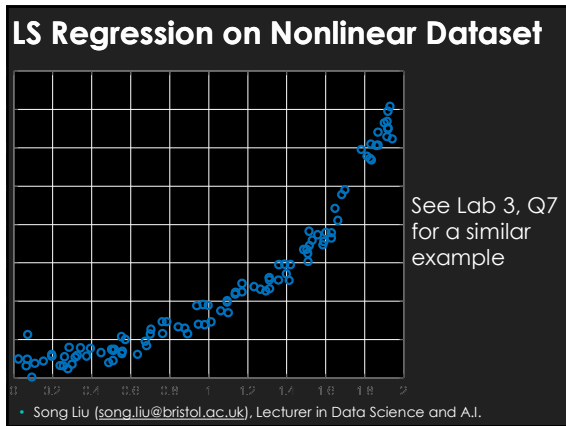
8

Least Squares (LS) + Feature Transform f

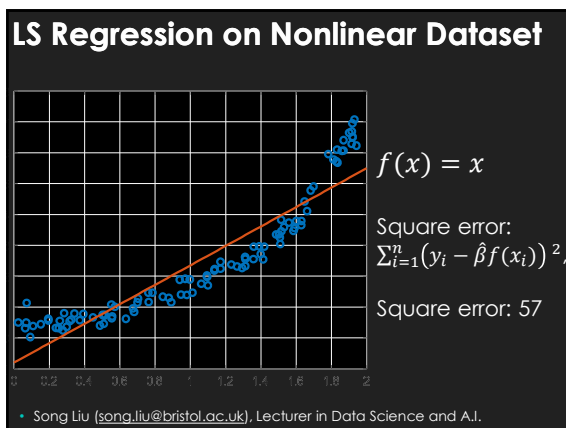
- Recall, given $D = \{(y_i, x_i)\}_i, y_i \in \mathbb{R}$,
- LS solves the following minimization:
 - $\hat{\beta} := \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta x_i)^2$ (1)
- Replace x with $f(x)$, a **feature transform**
 - $\hat{\beta} := \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta f(x_i))^2$ (2)
- (1) and (2) are identical if $f(x) = x$.

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

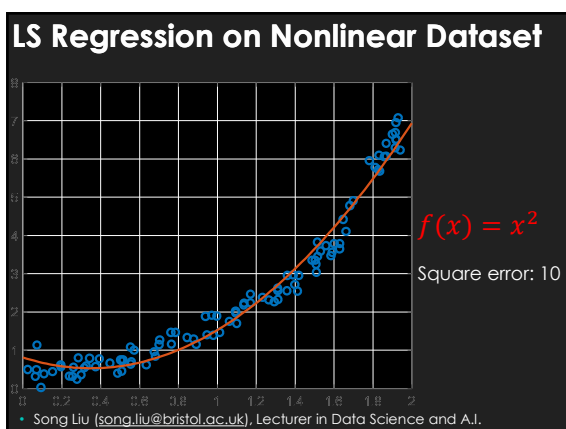
9



10



11



12

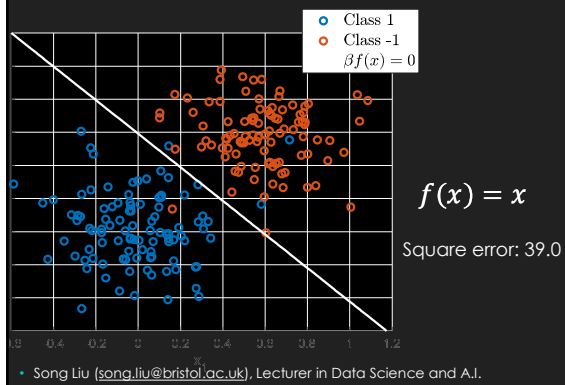
LS Classification + Feature Transform

- Classification dataset: $D = \{(y_i, x_i)\}_{i=1}^n, y \in \{-1, 1\}$.
- Now y only takes two discrete values -1 or 1 as **class labels**.
 - If $y_i = 1/-1$, x_i belongs to pos/neg class.
- Solving LS on D using feature transform f :
 - $\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n (y_i - \beta f(x_i))^2$
- $\hat{\beta} f(x) = 0$ indicates the **classification boundary**.
 - Why?

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

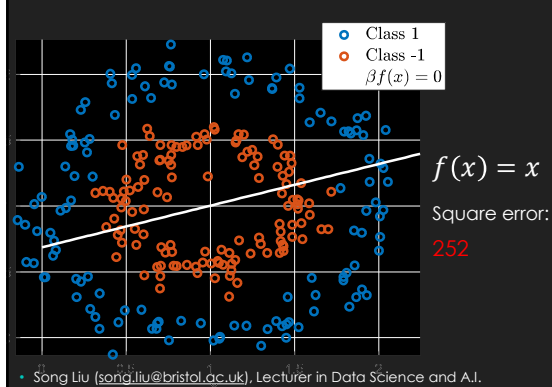
13

LS Classification

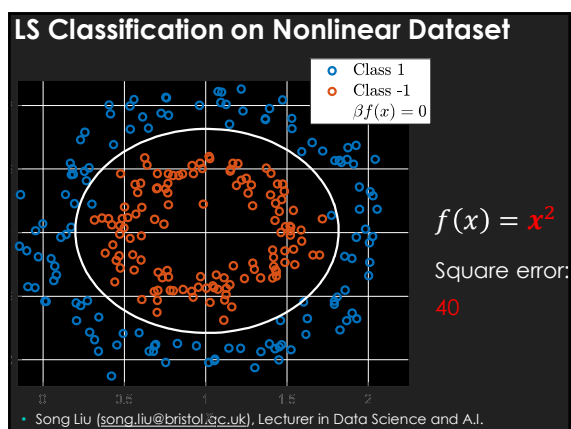


14

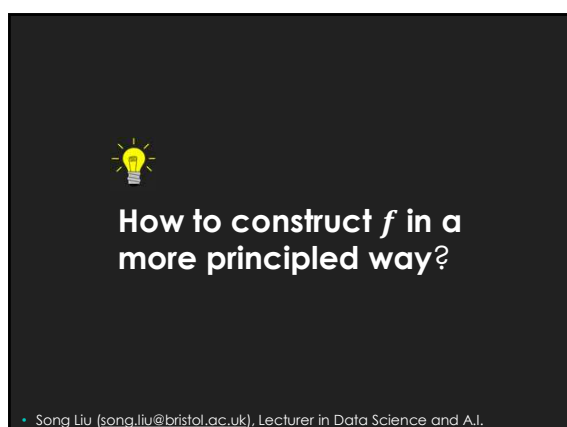
LS Classification on Nonlinear Dataset



15



16



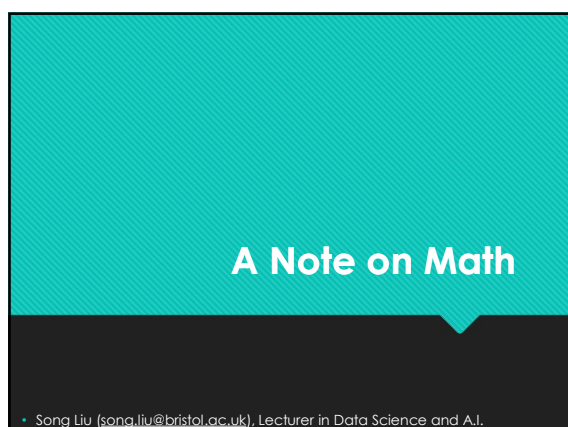
17

Two schools of thoughts:

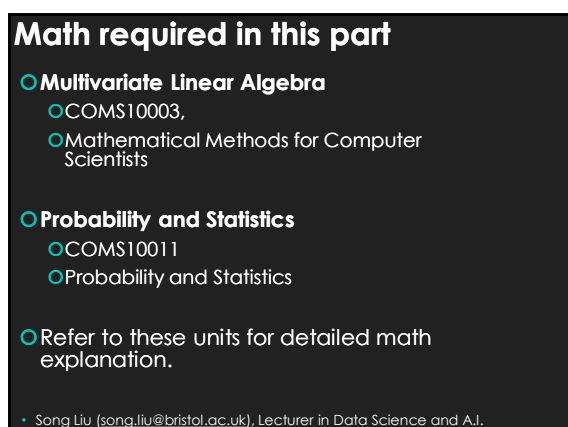
- Choosing f **manually** (Week 20,22)
 - **Pros:**
 - Efficient, require little computational effort.
 - Works well if you have domain knowledge.
 - **Cons:** Less flexible, requires tuning on different datasets.
- Choosing f **automatically** (Week 21)
 - **Pros:** Adaptive, automatically done on different datasets
 - **Cons:**
 - Extra computational burden.
 - Hard to integrate your domain knowledge.
- **Real-world problem solving involves a bit of both!!**

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

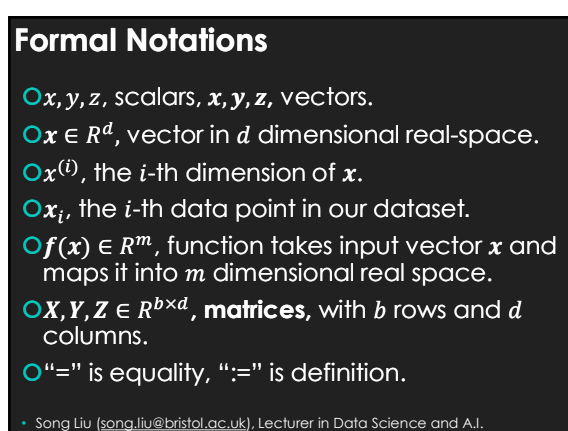
18



19



20



21

Polynomial Transform

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

22

A Generic Model

- We introduce a generic model.
- $\hat{y} := \langle \boldsymbol{\beta}, \mathbf{f}(\mathbf{x}) \rangle = \sum_i \beta^{(i)} f^{(i)}(\mathbf{x})$.
 - Inner product between $\boldsymbol{\beta}$ and \mathbf{f} .
 - \hat{y} is linear w.r.t. parameter $\boldsymbol{\beta}$.
- Special case:
 - when $f(x), \beta \in R, \hat{y} = \beta f(x)$.

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

23

Polynomial Transform

- Let $\mathbf{f}(\mathbf{x})$ be polynomial functions:
- When $x \in R, \mathbf{f}(x) := [x^0, x^1, x^2, \dots, x^b]$.
 - b is called the degree of \mathbf{f} .
 - $\mathbf{f}(x) = [0, x, x^2]$ is called a degree 2 polynomial trans. on x .

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

24

Polynomial Transform

- When $\mathbf{x} \in \mathbb{R}^d$,
 - $\mathbf{f}(\mathbf{x}) := [\mathbf{h}(x^{(1)}), \mathbf{h}(x^{(2)}), \dots, \mathbf{h}(x^{(d)})]$.
 - $\mathbf{h}(t) := [t^0, t^1, t^2, \dots, t^b] \in \mathbb{R}^{b+1}$.
 - $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^{d(b+1)}$, which means $\boldsymbol{\beta} \in \mathbb{R}^{d(b+1)}$.
 - PC: Write down $f^{(i)}(\mathbf{x})$ given i, b and d .

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

25

Polynomial Transform on Data Matrix

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ is data matrix with n observations and d dimensions.

$$\mathbf{f}(\mathbf{X}) := \begin{bmatrix} \mathbf{f}(\mathbf{x}_1) \\ \mathbf{f}(\mathbf{x}_2) \\ \vdots \\ \mathbf{f}(\mathbf{x}_n) \end{bmatrix} \in \mathbb{R}^{n \times d(b+1)}.$$

- We expanded our data matrix.
 - from d to $d(b+1)$

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

26

Pairwise Polynomial Transform

- So far, the polynomial transform is applied on each dimension:
 - i.e., $\mathbf{f}(\mathbf{x}) = [\mathbf{h}(x^{(1)}), \mathbf{h}(x^{(2)}), \dots, \mathbf{h}(x^{(d)})]$.
- It does **not** consider the dependencies between features.
 - Can be solved by appending cross terms i.e., $\mathbf{f}(\mathbf{x}) := [\mathbf{h}(x^{(1)}), \dots, \mathbf{h}(x^{(d)}), \forall_{u < v} x^{(u)} x^{(v)}]$

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

27

LS Solution

○ $\hat{\beta} = \arg \min \sum_{i=1}^n (y_i - \beta^T f(x_i))^2$

○ $\hat{\beta} = (f(X)^T f(X))^{-1} f(X)^T y$

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

28

Questions

○ At least, how many observations are needed to compute $\hat{\beta}$ with $f \in R^{d(b+1)}$ using the formula above?

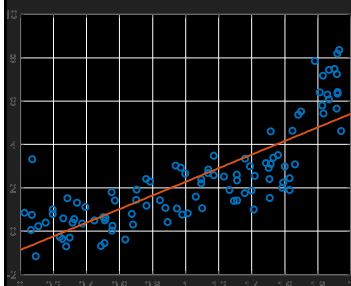
○ <https://pollev.com/songliu644>

○ PC: what is the computational complexity?

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

29

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

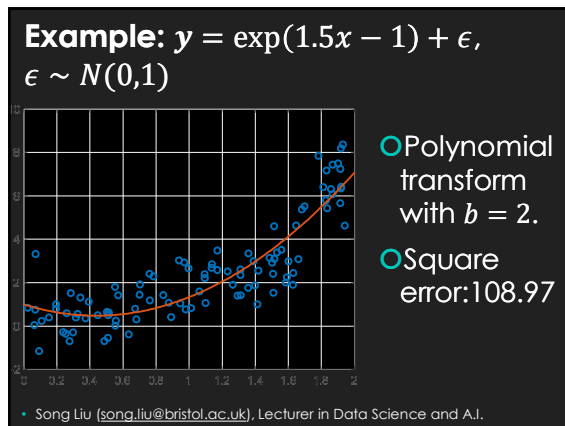


○ Polynomial transform with $b = 1$.

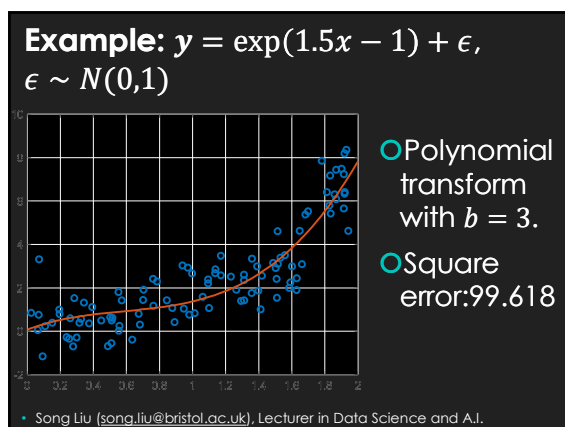
○ Square error: 171.0

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

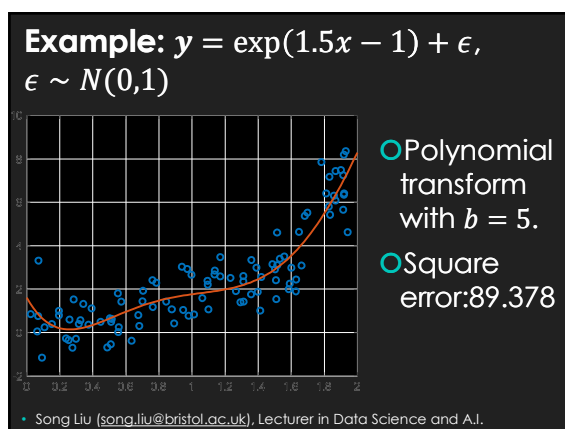
30



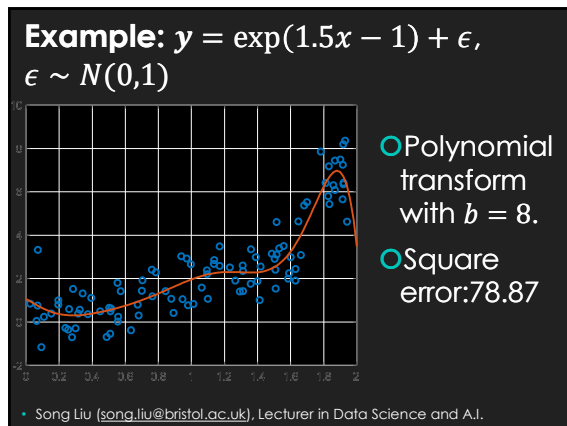
31



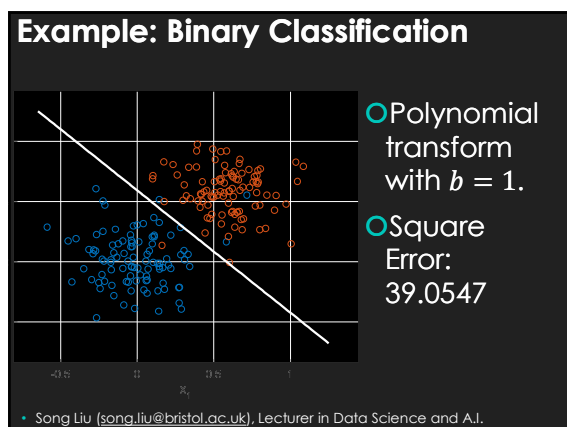
32



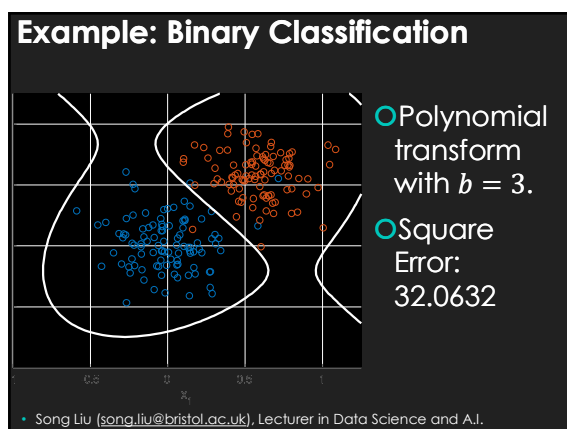
33



34



35



36

Observations:

- Pay attention on
 - how square error keeps **dropping** when **increasing** degree b .
 - how \hat{y} becomes more **flexible** when **increasing** b .
- We will revisit this point in the next lecture.

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

37

Why it works?

- 1-dimensional intuition: Taylor Series.
- Taylor Series of $g(x)$ at 0:
 - $g(x) = g(0)(x-0)^0 + g'(0)(x-0)^1 + \frac{g''(0)}{2!}(x-0)^2 + \frac{g'''(0)}{3!}(x-0)^3 + \dots$
- You can approximate a **smooth** function using polynomial terms (at some cost).

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

38

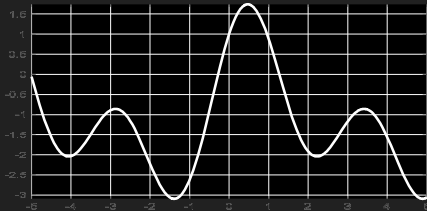
Fourier Series

- What are **other ways** of decomposing a function?
- Suppose we have a periodic signal $g(x)$ over the time domain.
 - e.g. a sound wave or a stock price
 - $g(x) = a_0 + \sum_{i=1}^{\infty} [a_i \sin(ix) + b_i \cos(ix)]$
 - This decomposition is called Fourier Series.

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

39

Fourier Series



○ $g(x) = \sin(x) + \cos(x) + \sin(2x) + \cos(2x)$

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

40

Trigonometric Transform

○ Trigonometric Transform are used to approximate function over **time domain**.

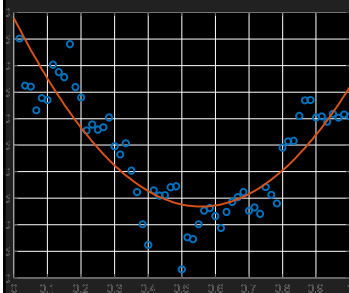
○ $f(x) := [1, \sin(x), \cos(x), \sin(2x), \cos(2x) \dots \sin(bx), \cos(bx)]$

○ $f(x) \in \mathbb{R}^{2b+1}$

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

41

Example: Apple Stock Price, Feb 2019



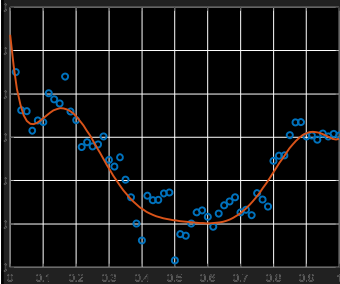
○ Trigonometric transform with $b = 1$.

○ Squared error: $1.5681e+03$

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

42

Example: Apple Stock Price, Feb 2019



Trigonometric transform with $b = 4$.

Squared error: 699.9117

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

43

Linear Expansion of Basis Functions

Polynomial and Trigonometric transforms based on the idea a function can be approximated by:

- $y \approx \hat{y} = \sum_{i=1}^m \beta^{(i)} f^{(i)}(x)$
- called a linear basis expansion of y
- $f^{(i)}$ are called **basis function**
 - Polynomial basis, Trigonometric basis...

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

44

Radial Basis Function (RBF)

RBF is another widely used basis function for function approximation.

$$f^{(i)}(x) := \exp\left(-\frac{\|x - x_i\|^2}{\sigma^2}\right)$$

- $\sigma > 0$ is called width
- σ is determined **before** fitting
- A practice is setting σ as the median of all pairwise distances of x in your data.

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

45

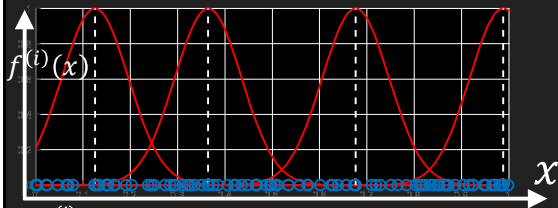
Radial Basis Function (RBF)

- x_i are called **RBF centroids**.
- x_i can be **randomly chosen** from the x in your dataset
- $f(x) := [1, f^{(1)}(x), f^{(2)}(x), \dots, f^{(b)}(x)]$
 - Do not forget 1!

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

46

Radial Basis Function (RBF)

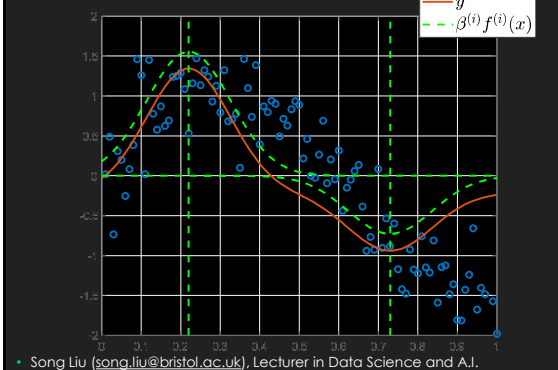


- $f^{(i)}(x)$ are visualized in red at random 4 centroids among 100 uniformly drawn x .
- At each "bump",
 - If $\beta^{(i)} > 0$, basis at $x^{(i)}$ gives \hat{y} a "lift".
 - If $\beta^{(i)} < 0$, basis at $x^{(i)}$ gives \hat{y} a "push".

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

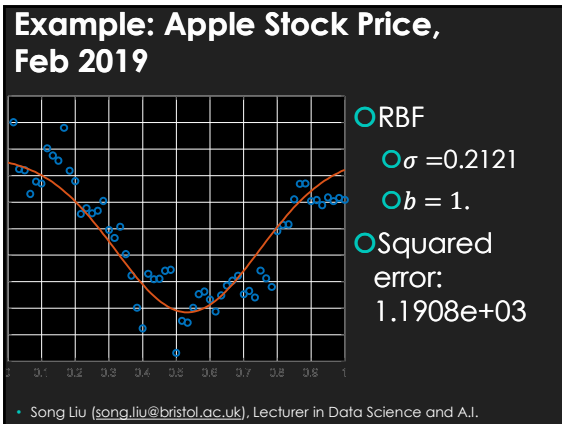
47

RBF feature, $b = 2$

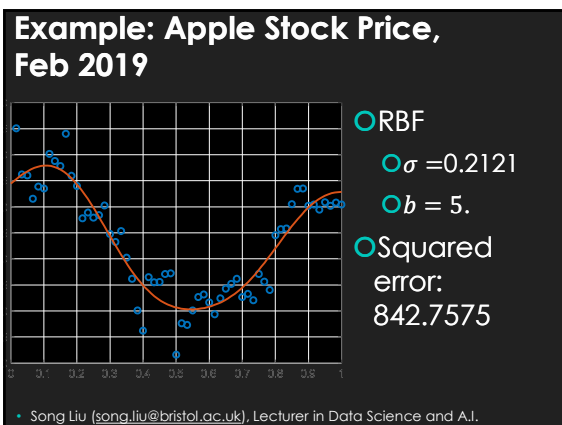


• Song Liu (song.liu@bristol.ac.uk), Lecturer in Data Science and A.I.

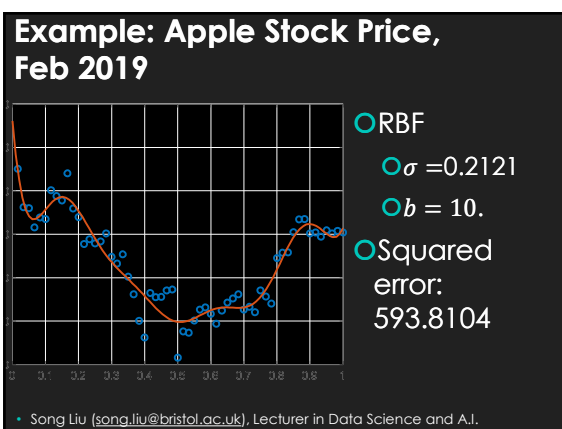
48



49

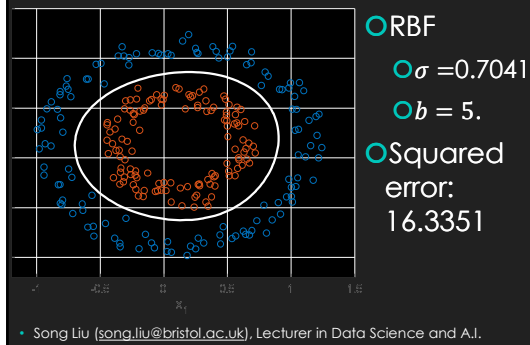


50



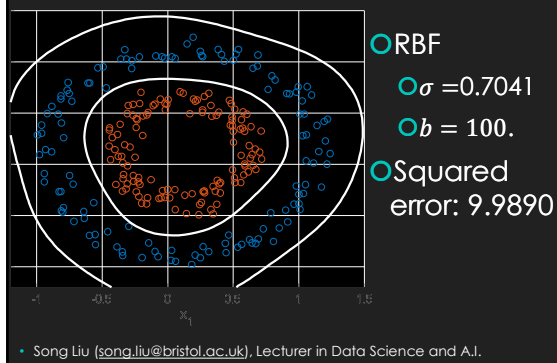
51

Example: Double Ring Classification



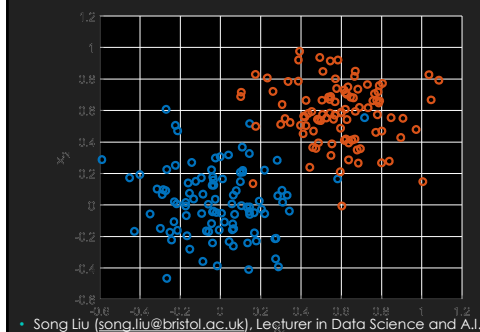
52

Example: Double Ring Classification



53

Selecting Features using Prior Knowledge



54

Question

Seeing your dataset above, what f should you use for classification? Hint: consider computational cost and overfitting

- Polynomial, $b = 1$
- Polynomial, $b = 2$
- Polynomial, $b = 3$
- RBF, $b = 100$

○ <https://bit.ly/2FnjryC>

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

55

Conclusion

Feature transform can be crucial to regression and classification tasks.

Three useful feature transform:

- Polynomial
- Trigonometric (on time series)
- RBF

As b increases, \hat{y} become more flexible, squared error is lowered.

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

56

Unanswered Questions

Increasing b drops squared error.

- How do you select number of basis b ?

Knowing an f with a larger b makes \hat{y} more flexible, can we make $b = \infty$?

Next two lectures, The selection of number of basis b and Kernel methods.

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

57

To know more...

○ The Elements of Statistical Learning:
Data Mining, Inference, and
Prediction, Hastie et al., 2009

○ 2.3.1 Linear Models and Least Squares

○ 2.6.3 Function Approximation

○ 2.8.3 Basis Functions

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