# DDA3020 Tutorial 4
# Linear Regression

Rongxiao Qu

School of Data Science

Email: rongxiaoqu@link.cuhk.edu.cn

Office hour: Tue 10:30 - 11:30, by appointment

Date: 2022.10.11

# Contents

- Definition of linearity
- Feature Transformation with Basis Functions
- Solving linear least squares
- Properties of LS estimator
- Generalized Linear Regressions
  - Ridge Regression
  - Lasso
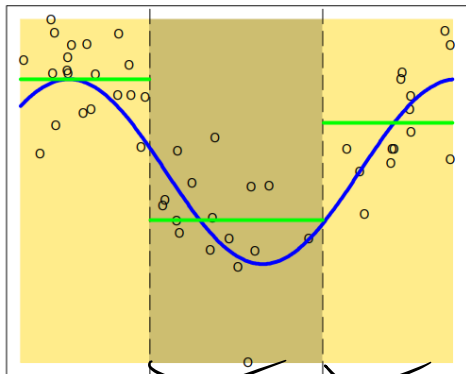- Code Demo

# "Linear" Regression

- A linear combination of the input features
- $f(x) = \underline{w_0} + \underline{x^T}\,\underline{w}$   ($f_w(x) = X\,w$)
- $f(x) = w_0 + \sum_{j=1}^{p} w_j x_j$

$$\begin{bmatrix} \vdots & \overset{x_1}{\vdots} & \overset{x_2}{\vdots} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

- Have advantage when the data size is small: avoid overfitting
- But it imposes significant limitations on the model

# Feature Transformation with Basis Functions

- $f_{\{\mathbf{w},b\}}(\mathbf{x}) = \sum_{j=1}^{p} w_j \phi(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$

- $(\mathbf{w} = (w_1, \ldots, w_n)^T \; \boldsymbol{\phi} = (\phi_1, \ldots, \phi_p))$

- Polynomial Regressions $\quad \phi(x) = x_1^2 \; x_1^3 \quad x_1 x_2$

- *Gaussian Basis Function*: $\phi_j(x) = \exp\left\{-\dfrac{(x-\mu_j)^2}{2s^2}\right\}$

- *Sigmoid Basis Function*: $\phi_j(x) = \sigma\left(\dfrac{x-\mu_j}{s}\right)$
  ( *logistic sigmoid function*: $\sigma(a) = \dfrac{1}{1+\exp(-a)}$ )

- *Splines (piecewise polynomials)*
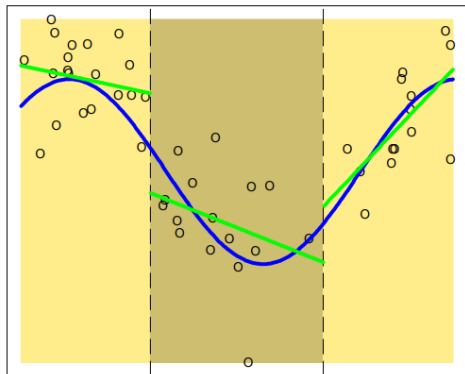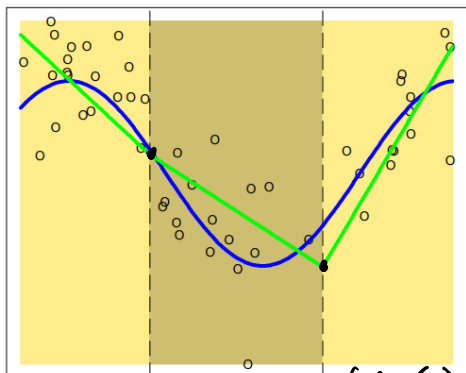- $f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 (x - \xi_1)_+^3 + w_5 (x - \xi_2)_+^3$

- Splines:

$$(x - \xi_1) = \begin{cases} x - \xi_1 & x > \xi_1 \\ 0 & x \leq \xi_1 \end{cases}$$

$$f_w(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 (x - \xi_1)_+^3 + w_5 (x - \xi_2)_+^3)$$
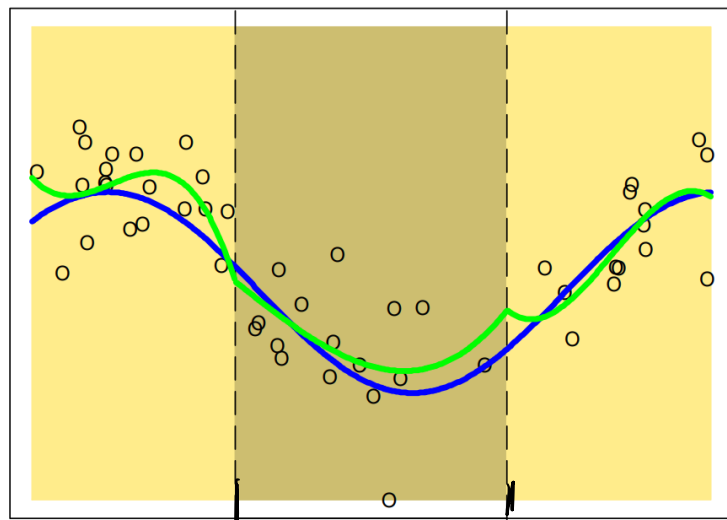


Piecewise Constant

Piecewise Linear

Continuous Piecewise Linear

Piecewise-linear Basis Function

$(X - \xi_1)_+$

$$f_w(x) = w_0 + w_1 x + w_2 (x - \xi_1)_+ + w_3 (x - \xi_2)_+$$

Cubic Splines

# Least Squares Regression

- Minimizing the squared error:

- $\hat{y} = X\hat{w}$

- $\hat{w} = \underset{w}{\arg\min}\, RSS = \underset{w}{\arg\min}\, \sum_{i=1}^{n}(f_w(x) - y)^2$

$$= \underset{w}{\arg\min}(Xw - y)^T(Xw - y)$$

- Take derivative w.r.t $\mathbf{w}$ and set the derivative to be $0 \rightarrow$

- $\hat{w} = (X^T X)^{-1} X^T y$

- $\hat{y} = X_{new}\hat{w} = X_{new}(X^T X)^{-1} X^T y$

# Properties of the LS estimator: $\widehat{w}$

For $y = f_w(x) + \epsilon = Xw + \epsilon$

- Assumptions:

  ① $f_w(x) \to$ true    ② $\epsilon \sim$ not correlated with $y / x$

  - $E(\epsilon) = 0$      (Mean of the errors are zeros)
  - $Cov(\epsilon) = \sigma^2 I$     (errors are uncorrelated with equal variance)

    (usually hard to satisfy)

- Conclusions:
  - $E(\widehat{w}) = w$
  - $Cov(\widehat{w}) = \sigma^2 (X^T X)^{-1}$   $\to$ conduct tests of significance for $w_i's$
  - ($\widehat{w}$ is the best linear unbiased estimator (BLUE) of $w$)

$$E(\hat{w}) = E\left( (X^TX)^{-1}X^Ty_\downarrow \right)$$

$$= E\left( (X^TX)^{-1}X^T(Xw+e) \right)$$

$$= E\left( (X^TX)^{-1}X^TXw + (X^TX)^{-1}X^Te \right)$$

$$= E(w_\downarrow) + (X^TX)^{-1}X^T \underline{\underline{E(e)}}_{=0}$$

$$= w$$

$$\text{Cov}(\hat{w}) = E\left( (\hat{w} - E(\hat{w}))(\hat{w} - E(\hat{w}))^T \right)$$

$$= E\left( ((X^TX)^{-1}X^T(Xw+e)-w)((X^TX)^{-1}X^T(Xw+e)-w)^T \right)$$

$$= E\left( ((X^TX)^{-1}X^Te)((X^TX)^{-1}X^Te)^T \right)$$

$$= E\left( (X^TX)^{-1}X^Tee^TX(X^TX)^{-1} \right)$$

$$= (X^TX)^{-1}X^T \underline{\underline{E(ee^T)}} X(X^TX)^{-1}$$
$$\phantom{= (X^TX)^{-1}X^T} = \text{Cov}(e)$$

$$= \sigma^2 (X^TX)^{-1}\underline{X^TX}(X^TX)^{-1}$$

$$= \sigma^2 (X^TX)^{-1}$$

# Ridge Regression

- $\hat{\beta}^{ridge} = \underset{\beta}{\text{argmin}} \left\{ \sum_{i=1}^{n} y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j + \lambda \sum_{j=1}^{p} \beta_j^2 \right\}$

  $= \underset{\beta}{\text{argmin}} \left\{ (X\beta - y)^T(X\beta - y) + \lambda ||\beta||^2 \right\})$

- Equivalently:

- $\hat{\beta}^{ridge} = \underset{\beta}{\text{argmin}} \left\{ \sum_{i=1}^{n} y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right\}$

  $(= \underset{\beta}{\text{argmin}} \left\{ (X\beta - y)^T(X\beta - y) \right\}$

  $\text{subject to } \sum_{j=1}^{p} \beta_j^2 \leq t$ )

  $(\text{subject to } ||\beta||^2 \leq t)$

$$RSS = (Xw - y)^T (Xw - y) + \lambda w^T w$$

$$= w^T X^T X w - w^T X^T y - y^T X w + y^T y + \lambda w^T w$$

$$\frac{\partial RSS}{\partial w} = 2 X^T X w - X^T y - X^T y + 2\lambda w = 0$$

$$X^T X w + \lambda w = X^T y$$

$$(X^T X + \lambda I) w = X^T y$$

$$\hat{w}_{ridge} = (X^T X + \lambda I)^{-1} X^T y$$

# Lasso

- $\hat{\beta}^{ridge} = \underset{\beta}{\arg\min} \left\{ \sum_{i=1}^{n} y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j + {\color{red}\lambda \sum_{j=1}^{p} |\beta_j|} \right\}$

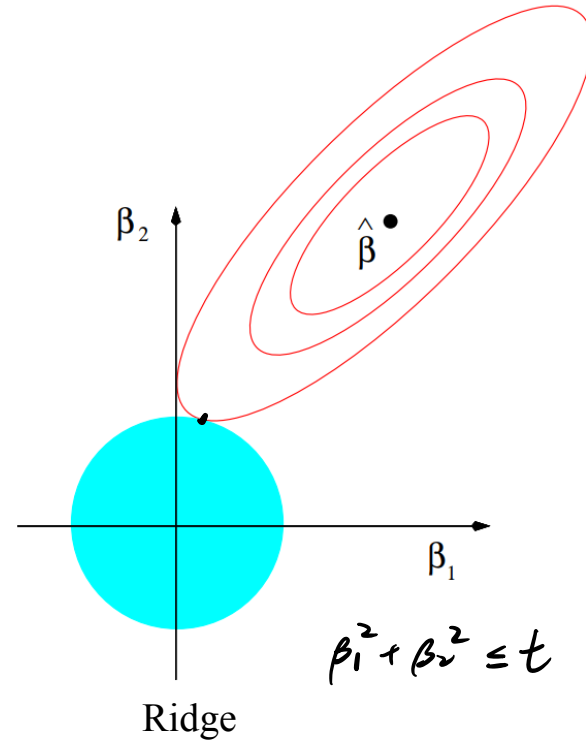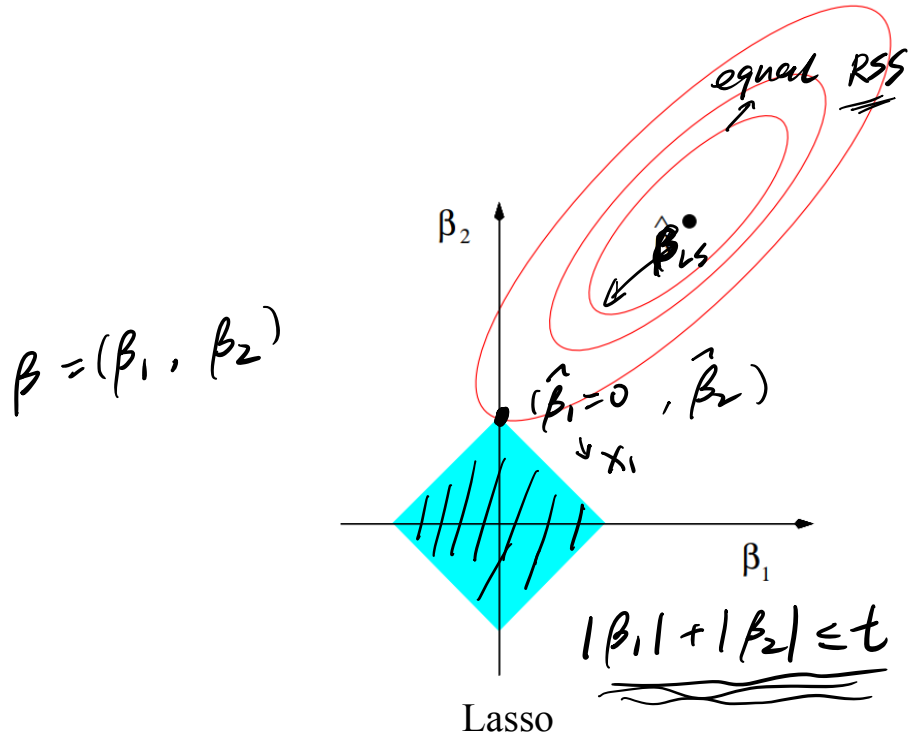$$\left( = \underset{\beta}{\arg\min} \left\{ (X\beta - y)^T(X\beta - y) + \lambda|\beta|_1 \right\} \right)$$

- Equivalently:
- $\hat{\beta}^{ridge} = \underset{\beta}{\arg\min} \left\{ \sum_{i=1}^{n} y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j \right\}$

$$\left( = \underset{\beta}{\arg\min} \left\{ (X\beta - y)^T(X\beta - y) \right\} \right.$$

$$\text{subject to } {\color{red}|\beta|_1 \leq t}$$

$$\left. (\text{subject to } |\beta|_1 \leq t) \right)$$

# Geometry of Ridge and Lasso regression



$\beta = (\beta_1, \beta_2)$

equal RSS

$\hat{\beta}_{LS}$

$(\hat{\beta}_1 = 0, \hat{\beta}_2)$

$\downarrow x_1$

$|\beta_1| + |\beta_2| \le t$

Lasso

$\hat{\beta}$

$\beta_1^2 + \beta_2^2 \le t$

Ridge

# Code Demo

Data Processing

Ridge Regression and Lasso