

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
- $f(\mathbf{x}) = w_0 + \mathbf{x}^T \mathbf{w}$ ($f_{\mathbf{w}}(\mathbf{x}) = \mathbf{X} \mathbf{w}$)
- $f(\mathbf{x}) = w_0 + \sum_{j=1}^p w_j x_j$

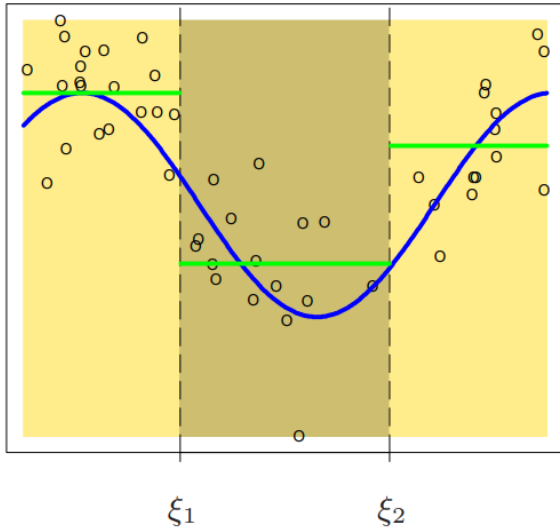
- Have advantage when the data size is small: avoid overfitting
- But it imposes significant limitations on the model

Feature Transformation with Basis Functions

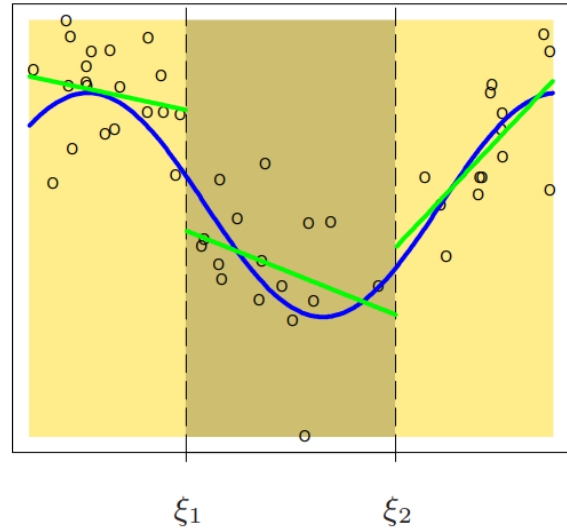
- $f_{\{\mathbf{w}, b\}}(\mathbf{x}) = \sum_{j=1}^p w_j \phi(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$
- $(\mathbf{w} = (w_1, \dots, w_n)^T \quad \boldsymbol{\phi} = (\phi_1, \dots, \phi_p))$
- Polynomial Regressions
- *Gaussian Basis Function*: $\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$
- *Sigmoid Basis Function*: $\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$
(*logistic sigmoid function*: $\sigma(a) = \frac{1}{1 + \exp(-a)}$)
- *Splines (piecewise polynomials)*
- $f(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4(x - \xi_1)_+^3 + w_5(x - \xi_2)_+^3$

- Splines:

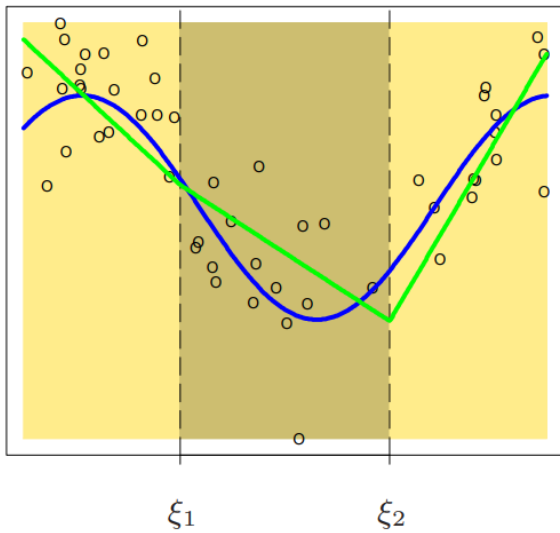
Piecewise Constant



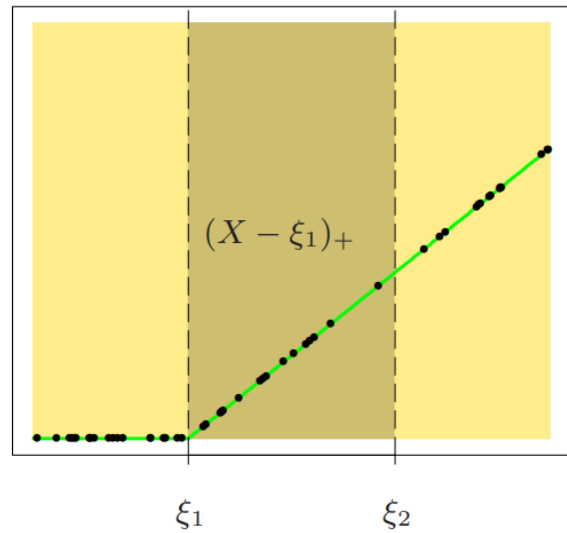
Piecewise Linear



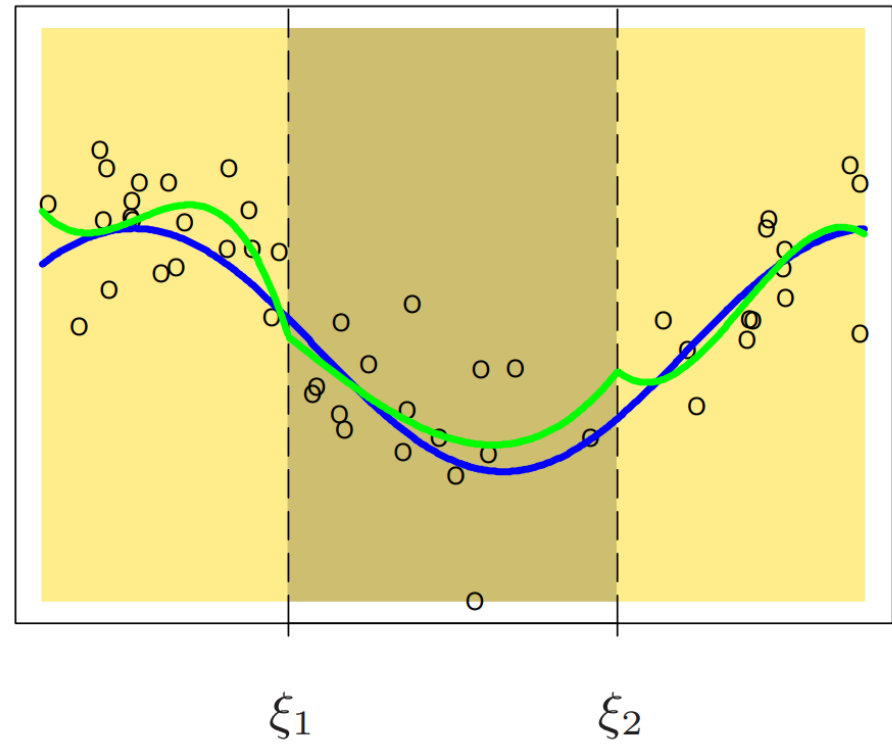
Continuous Piecewise Linear



Piecewise-linear Basis Function



Cubic Splines



Least Squares Regression

- Minimizing the squared error:
- $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}}$
- $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} RSS = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (f_{\mathbf{w}}(\mathbf{x}) - y)^2$
 $= \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$
- Take derivative w.r.t \mathbf{w} and set the derivative to be 0 \rightarrow
- $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- $\hat{\mathbf{y}} = \mathbf{X}_{new} \hat{\mathbf{w}} = \mathbf{X}_{new} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Properties of the LS estimator: $\hat{\mathbf{w}}$

For $\mathbf{y} = \mathbf{f}_{\mathbf{w}}(\mathbf{x}) + \boldsymbol{\epsilon} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$

- Assumptions:

- $E(\boldsymbol{\epsilon}) = 0$ (Mean of the errors are zeros)
- $Cov(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$ (errors are uncorrelated with equal variance)
(usually hard to satisfy)

- Conclusions:

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

Ridge Regression

- $\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j + \lambda \sum_{j=1}^p \beta_j^2 \right\}$
 $= \underset{\beta}{\operatorname{argmin}} \left\{ (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \lambda ||\boldsymbol{\beta}||^2 \right\}$

- Equivalently:

- $\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right\}$
(= $\underset{\beta}{\operatorname{argmin}} \left\{ (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) \right\}$
subject to $\sum_{j=1}^p \beta_j^2 \leq t$
(subject to $||\boldsymbol{\beta}||^2 \leq t$)

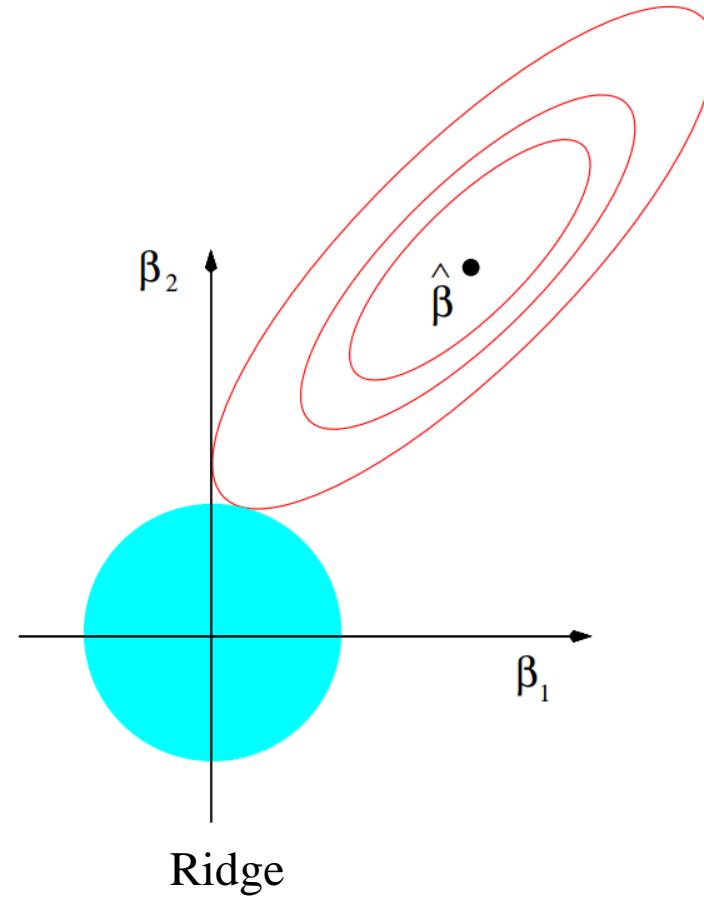
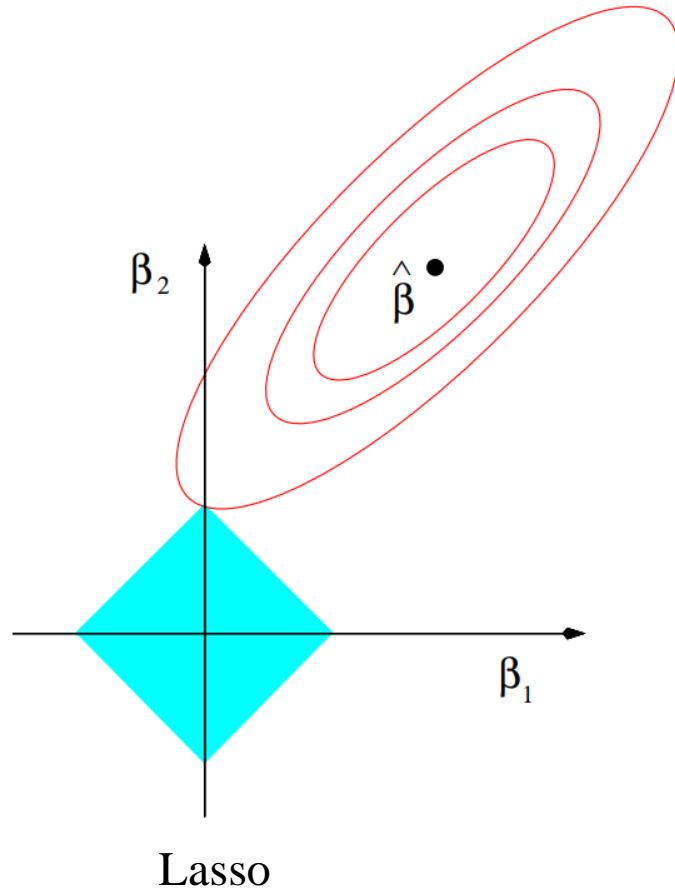
Lasso

- $\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j + \lambda \sum_{j=1}^p |\beta_j| \right\}$
(= $\underset{\beta}{\operatorname{argmin}} \{ (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \lambda |\boldsymbol{\beta}|_1 \}$)

- Equivalently:

- $\hat{\beta}^{ridge} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right\}$
(= $\underset{\beta}{\operatorname{argmin}} \{ (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) \}$
subject to $|\boldsymbol{\beta}|_1 \leq t$
(subject to $|\boldsymbol{\beta}|_1 \leq t$)

Geometry of Ridge and Lasso regression



Code Demo

Data Processing

Ridge Regression and Lasso

