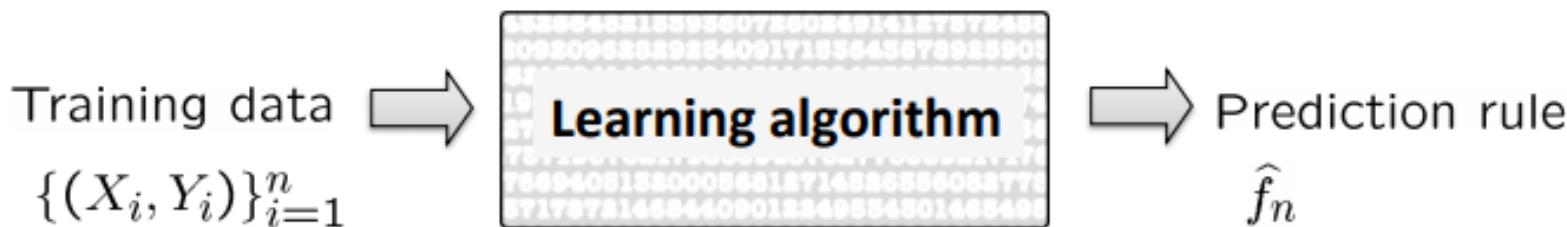


Linear Regression

Xiangli Chen

Regression algorithm



- Linear regression
- Regularized linear regression – ridge, lasso regression
- Polynomial regression
- Kernel regression
- Regression trees, splines, wavelet estimators, ...

Restrict class of predictors

Optimal predictor: $f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$

Empirical Minimizer: $\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$

Class of predictors

Law of Large Numbers:

$$\frac{1}{n} \sum_{i=1}^n [\text{loss}(Y_i, f(X_i))] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

Restrict class to avoid overfitting

- Class of linear functions
- Class of polynomial functions
- Class of nonlinear functions

Linear regression

least square Estimator

$$\hat{f}_{Lm} = \underset{f \in \mathcal{F}_L}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

— class of linear functions

Uni-variate case:

$$f(x) = \beta_0 + \beta_1 x \quad \beta_0: \text{intercept} \quad \beta_1: \text{slope}$$

Multi-variate case:

$$f(x) = \beta^T x \quad \text{where } x = [1, x_1, \dots, x_p]^T, \beta = [\beta_0, \dots, \beta_p]^T$$

Least squares estimator

$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad f(X_i) = X_i \beta$$



$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (X_i \beta - Y_i)^2 \quad \hat{f}_n^L(X) = X \hat{\beta}$$

$$= \arg \min_{\beta} \frac{1}{n} (X \beta - Y)^T (X \beta - Y)$$

$$X = \begin{bmatrix} X_1^T \\ \vdots \\ X_n^T \end{bmatrix} = \begin{bmatrix} 1, X_1, \dots, X_P \\ \vdots \\ 1, X_n, \dots, X_P \end{bmatrix} \quad Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

Least Square Estimator

Let $J(\beta) = (X\beta - y)^T (X\beta - y)$

so $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} J(\beta)$

$$\frac{\partial J(\beta)}{\partial \beta} \Big|_{\hat{\beta}=0} \Rightarrow (X^T X) \hat{\beta} = X^T y$$

If $X^T X$ is invertible $\hat{\beta} = X(X^T X)^{-1} X^T y$

- $X^T X$ is invertible, when $X^T X$ has full rank
- If $X^T X$ is not invertible, try regularization, gradient descent learning method

Geometric Interpretation of LS

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

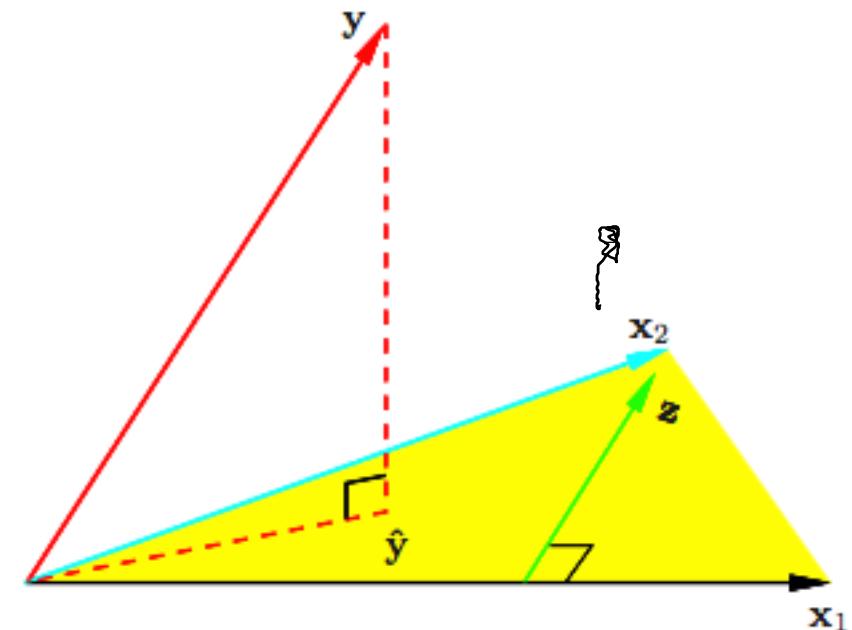
\hat{y} is a linear combination of columns of X
on the space spanned by those columns

Note that

$$X^T(y - \hat{y}) = X^T y - X^T X (X^T X)^{-1} X^T y = 0$$

so $y - \hat{y}$ is orthogonal to the space

\hat{y} is the orthogonal project of y into the space



Gradient Descent

When $X^T X$ is invertible or computationally expensive if X is huge

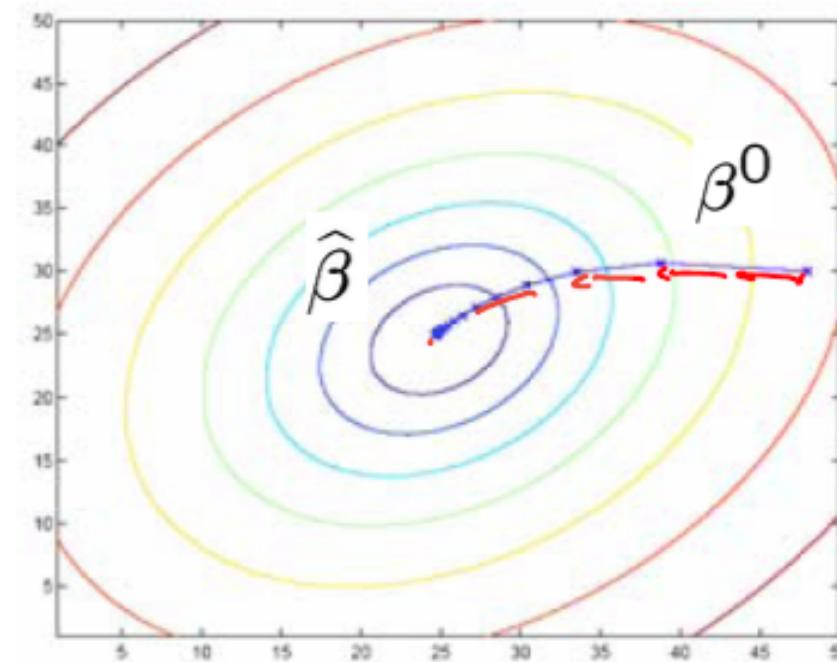
$$\hat{\beta} = \arg \min_{\beta} J(\beta)$$

$J(\beta)$ is convex with respect to β

Gradient descent step

$$\begin{aligned}\beta^{t+1} &= \beta^t - \alpha \frac{\partial J(\beta)}{\partial \beta} \\ &= \beta^t - 2\alpha X^T(X\beta^t - y)\end{aligned}$$

α is the step size



Gaussian Assumption and Likelihood

$$y = \beta^T x + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad Y | X=x \sim N(\beta^T x, \sigma^2)$$

Given random sample (i.i.d) $\{x_i, y_i\}_{i=1}^N$, conditional likelihood

$$f(y_{i=1}^N | x_{i=1}^N, \beta, \sigma^2) = \prod_{i=1}^N f(y_i | x_i, \beta, \sigma^2)$$

Maximize conditional log likelihood (minimize empirical conditional log loss)

$$\hat{\beta}_{MLE} = \arg \max_{\beta} \sum_{i=1}^N \log f(y_i | x_i, \beta, \sigma^2) \quad (\text{discriminative approach})$$

$$= \arg \min_{\beta} \sum_{i=1}^N (\beta^T x_i - y_i)^2 = \hat{\beta}$$

Least square estimation is same as maximize conditional log likelihood estimate under a Gaussian model

Shrinkage Methods

Ridge regression – regularized by L_2 penalty

$$\hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{arg\,min}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad \lambda \geq 0$$

Note that β_0 has been left out of the penalty term

- Adding c to y_i 's, β_0 increasing by c as well and also predicted values \hat{y}_i 's. It is not true if β_0 is penalized
- β_0 affects the shift, not the complexity

Same Problem

Re-express the ridge regression

$$\sum_{i=1}^N \left(y_i - \beta_0^C - \sum_{j=1}^P (x_{ij} - \bar{x}_j) \beta_j^C \right)^2 + \lambda \sum_{j=1}^P \beta_j^{C^2}$$

$$\text{where } \beta_0^C = \beta_0 + \sum_{j=1}^P \bar{x}_j \beta_j \quad \beta_j^C = \beta_j \quad j=1, 2, \dots, P$$

$$\text{Note that } \hat{\beta}_0^C = \bar{y} \quad \hat{\beta}_j^C = \hat{\beta}_j$$

(at $y_i = y_i - \bar{y}$ and $x_{ij} = x_{ij} - \bar{x}_j$ gives ridge regression without intercept)

$$\sum_{i=1}^N (y_i - \sum_{j=1}^P x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^P \beta_j^2 \Leftrightarrow (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y \quad (X^T X + \lambda I) \text{ always has full rank}$$

SVD Decomposition

MAP Interpretation

Maximize a posterior (MAP)

$$f(\beta | y_{i=1}^N, x_{i=1}^N, \sigma^2) \propto f(y_{i=1}^N | \beta, x_{i=1}^N, \sigma^2) f(\beta | x_{i=1}^N, \sigma^2)$$

Gaussian Prior

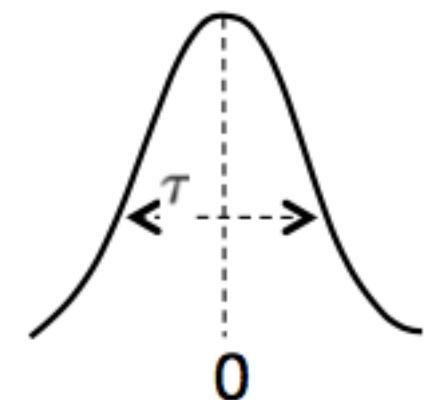
$$\beta \sim N(0, \tau^2 I) \quad f(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$

$$\hat{\beta}_{MAP} = \arg \max_{\beta} \log f(y_{i=1}^N | \beta, x_{i=1}^N, \sigma^2) f(\beta)$$

$$= \arg \min_{\beta} \sum_{i=1}^N (y_i - \beta^T x_i)^2 + (\frac{\sigma}{\tau})^2 \|\beta\|_2^2$$

(at $\lambda = (\frac{\sigma}{\tau})^2$, the expression is the same to ridge regression)

$$\hat{\beta}_{MAP} = \hat{\beta}_{ridge}$$



Ridge Regression

Lasso

Laplace Prior

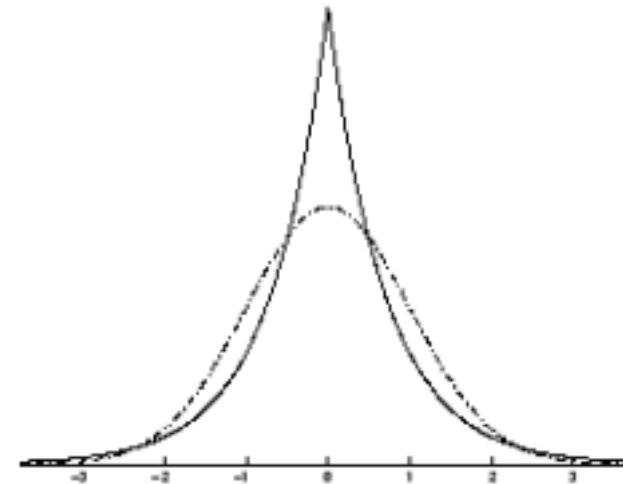
$$\beta \sim \text{Laplace}(0, t) \quad f(\beta) \propto e^{-|\beta|/t}$$

MAP estimator

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^N (\gamma_i - \beta^T x_i)^2 + \frac{\sigma^2}{t} \|\beta\|_1$$

Let $\lambda = \frac{\sigma^2}{t}$, the expression is Lasso (regularized with l1 term)

No closed form solution — a quadratic programming problem



Lasso

Prime Problem

Prime problem

$$\hat{\beta} = \arg \min_{\beta} J(\beta) \Rightarrow \hat{\beta} = \arg \min_{\beta} J(\beta) + \lambda \text{pen}(\beta)$$

subject to $\text{pen}(\beta) \leq \eta$

Lagrangian form

$$\text{Ridge: } \|\beta\|_2^2 \leq \eta \quad \text{Lasso: } \|\beta\|_1 \leq \eta$$

$\|\beta\|_p \leq \eta$ gives convex set of β

- There is a one to one correspondence between λ and η .
- When λ increases, η decreases, an increasing number of coordinates are driven to zero.

Ridge regression vs lasso

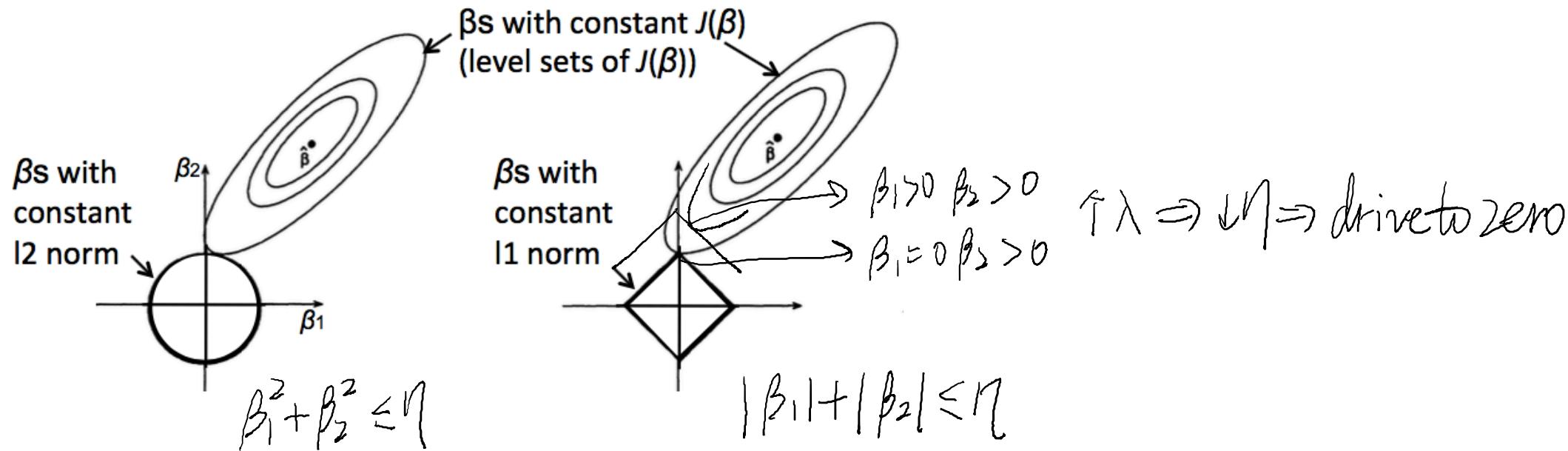
Ridge Regression:

$$\text{pen}(\beta) = \|\beta\|_2^2$$

Lasso:

$$\text{pen}(\beta) = \|\beta\|_1$$

HOT!



Lasso results in sparse solution – more zero coordinates,
good for high dimension

Polynomial Regression

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_M x^M$$

$$X = \begin{bmatrix} 1 & x_1 & \dots & x_1^M \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^M \end{bmatrix} \quad \hat{\beta} = (X^T X)^{-1} X^T y$$

