

Regression

- In regression problems, label y is continuous (usually y is scalar)
- y also called response
- Supervised learning problem: given training data $\mathcal{T} = \{\underline{x}_i, y_i\}_{i=1}^N$, find a model f to predict $\hat{y} = f(\underline{x})$

Applications: prediction of house prices based on location attributes, prediction of income based on demographic features, etc.

Loss Function Is "0-1" loss useful? No.

Commonly used loss model for regression,

$$\text{Square loss: } L(f(\underline{x}), y) = (y - f(\underline{x}))^2$$

Linear Regression $f(\underline{x}) = \beta^T \underline{x}$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix} \quad \text{Learned from training data}$$

Example Prediction of home price
 $x_1 = \text{Crime rate}$, $x_2 = \text{pollution level}$

$$f(\underline{x}) = \beta_1 x_1 + \beta_2 x_2 .$$

Linear regression can be used to do more general (non linear) modeling by augmenting \underline{x}

Example:

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$d=2$$

$$\underline{x}' = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \end{bmatrix}$$

$$d'=4$$

Predict response using \underline{x}' as: $\beta_1 + \beta_2 x_1 + \beta_3 x_2 + \beta_4 x_1^2$

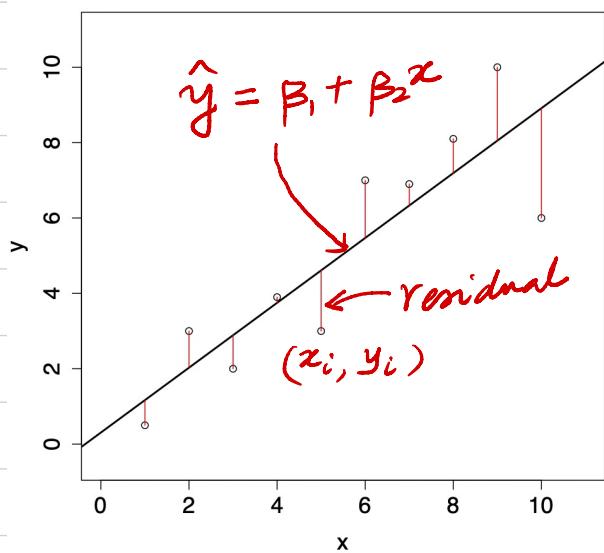
Learning β from $\mathcal{T} = \{\underline{x}_i, y_i\}_{i=1}^N$

For squared loss, define Residual Sum of Squares:

$$RSS(\beta) = \sum_{i=1}^N (\underline{y}_i - \underbrace{\underline{x}_i^T \beta}_{\hat{y}_i})^2$$

Choose β to minimize $RSS(\beta)$

$$\hat{\beta}_{LS} = \arg \min_{\beta} RSS(\beta)$$



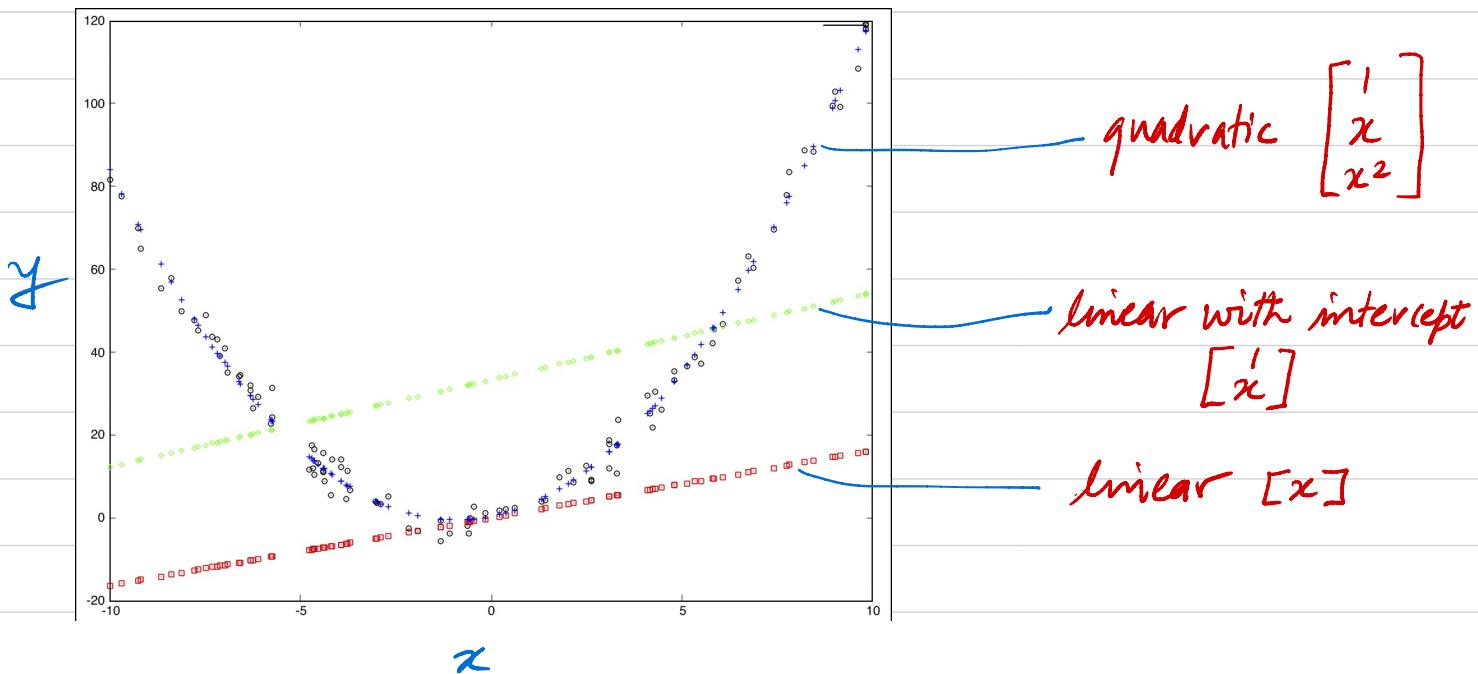
$$RSS(\beta) = \|\underline{y} - X\beta\|^2$$

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}, X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}_{N \times d}$$

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

Model Selection

Train and validate (or cross-validate) several models and choose best one



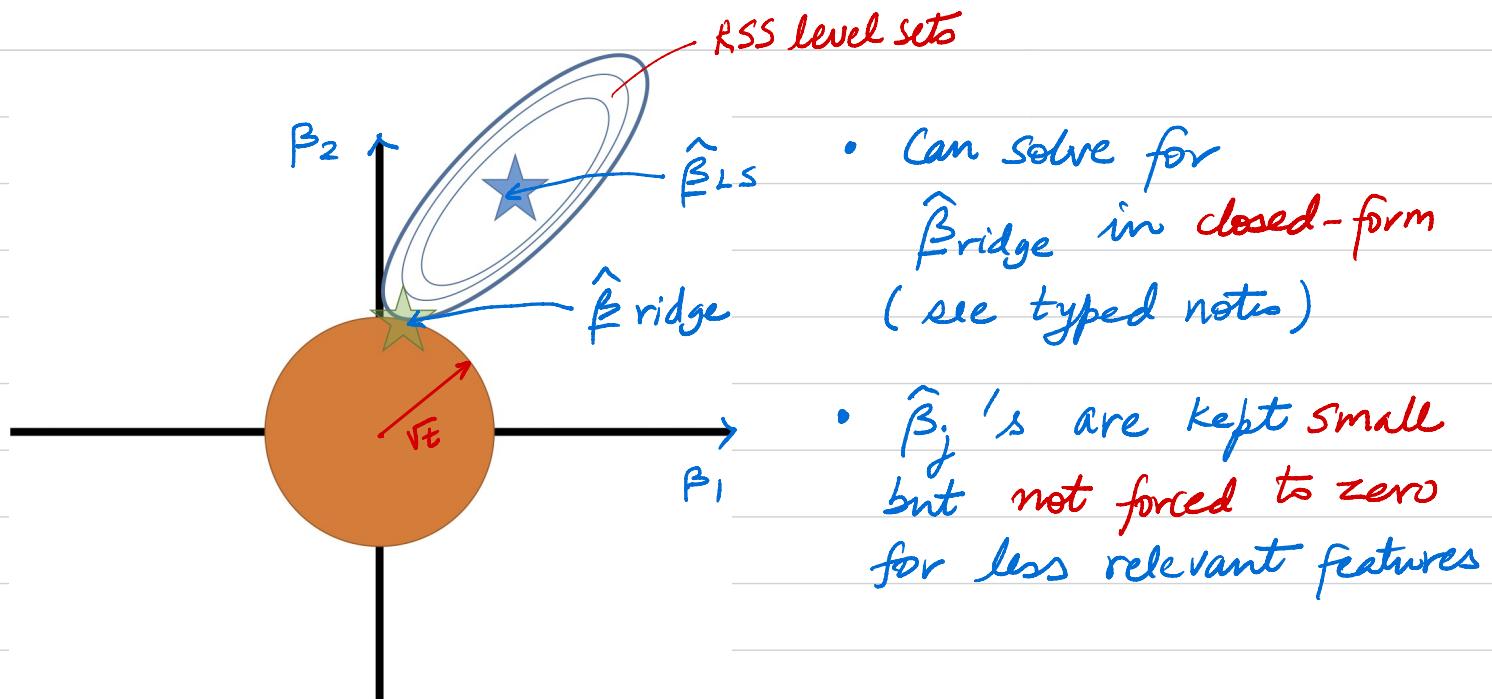
Subset (Feature) Selection

- Often some features are irrelevant or of minor relevance to prediction \hat{y} , e.g. average crime rate in Illinois to home prices in Urbana.
- Goal of feature selection is to automatically give small (or zero) weight β_i to features that are irrelevant. This results in simpler models with higher accuracy.
- Shrinkage Formulation:
Minimize $RSS(\beta)$ subject to $r(\beta) \leq t$
 \uparrow
regularizer

Ridge Regression

$$r(\beta) = \|\beta\|^2 = \sum_i \beta_i^2$$

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



LASSO (Least Absolute Shrinkage and Selection Operator)

$$r(\beta) = \sum_i |\beta_i| \leftarrow L_1\text{-norm}$$

Minimize $\text{RSS}(\beta)$ subject to $\sum_i |\beta_i| \leq t$

