## 2D-data

$$L = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda m^2$$

$$= \sum_{i=1}^{n} (y_i - mx_i - b)^2 + \lambda m^2$$

$$\frac{\partial L}{\partial b}$$

$$L = \sum_{i=1}^{n} (y_i - mx_i - b)^2 + \lambda m^2$$

$$\frac{\partial L}{\partial b} = -2\sum_{i=1}^{n}(y_i - mx_i - b) = 0$$

$$\sum_{i=1}^{n}(y_i - mx_i - b) = 0$$

$$\sum_{i=1}^{n} y_i - m\sum_{i=1}^{n}x_i - \sum_{i=1}^{n}b = 0$$

$$\bar{y} - m\bar{x} - b = 0$$

$$-b = \bar{y} - m\bar{x}$$

$$\boxed{b = \bar{y} - m\bar{x}}$$

$$\frac{\partial L}{\partial m}$$

$$L = \sum_{i=1}^{n}(y_i - mx_i - b)^2 + \lambda(m^2)$$

$$= \sum_{i=1}^{n}(y_i - mx_i - \bar{y} + m\bar{x})^2 + \lambda(m^2)$$

$$= \sum_{i=1}^{n}[(y_i - \bar{y}) + m(\bar{x} - x_i)]^2 + \lambda m^2$$

$$\frac{\partial L}{\partial m} = 2\sum_{i=1}^{n}(\bar{x} - x_i)[(y_i - \bar{y}) + m(\bar{x} - x_i)] + 2\lambda m = 0$$

$$2\left[\sum_{i=1}^{n}[(\bar{x}-x_i)(y_i-\bar{y}) + m(\bar{x}-x_i)^2] + \lambda m\right] = 0$$

$$\sum_{i=1}^{n}(\bar{x}-x_i)(y_i-\bar{y}) + m\sum_{i=1}^{n}(\bar{x}-x_i)^2 + \lambda m = 0$$

$$m\left[\sum_{i=1}^{n}(\bar{x}-x_i)^2 + \lambda\right] = -\sum_{i=1}^{n}(\bar{x}-x_i)(y_i-\bar{y})$$

$$\boxed{m = \frac{\sum_{i=1}^{n}(x_i-\bar{x})(y_i-\bar{y})}{\sum_{i=1}^{n}(x_i-\bar{x})^2 + \lambda}}$$

### SLR vs RIDGE

$$\boxed{m = \frac{\sum(x_i-\bar{x})(y_i-\bar{y})}{\sum(x_i-\bar{x})^2}} \; > \; \boxed{m = \frac{\sum(x_i-\bar{x})(y_i-\bar{y})}{\sum(x_i-\bar{x})^2 + \lambda}} \to \text{hyperparameter } [0, \infty)$$

$$\uparrow \lambda \downarrow m$$

$$\boxed{b = \bar{y} - m\bar{x}} \; < \; \boxed{b = \bar{y} - m\bar{x}}$$

## nD-data

| $x_{11}$ | $x_{12}$ | $x_{13}$ | ... | $x_{1m}$ | $y_1$ |
|---|---|---|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ | ... | $x_{2m}$ | $y_2$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ | ... | $x_{3m}$ | $y_3$ |
| $x_{n1}$ | $x_{n2}$ | $x_{n3}$ | ... | $x_{nm}$ | $y_n$ |

$$\boxed{\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_m x_m}$$

$$\hat{y} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & x_{23} & \cdots & x_{2m} \\ \vdots & & & & & \\ 1 & x_{n1} & x_{n2} & x_{n3} & \cdots & x_{nm} \end{bmatrix}_{n \times (m+1)} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}_{(m+1) \times 1}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$X \qquad \beta$$

① $\boxed{\hat{y} = X\beta}$

$\lambda(\beta_0^2 + \beta_1^2 + \cdots + \beta_n^2)$

② $\boxed{E = (y - \hat{y})^2(y + \hat{y}) + \lambda\|\beta\|^2}$

$$= (y - \hat{y})^T(y - \hat{y}) + \lambda \beta^T \beta$$

$$= y^T y - y^T \hat{y} - \hat{y}^T y + \hat{y}^T \hat{y} + \lambda \beta^T \beta$$

Let prove that $\hat{y}^T y$ is a symmetric matrix

$$\hat{y}^T_{1 \times n} y_{n \times 1} = [\quad]_{1 \times 1} \to \text{scalar} \; \therefore \; \hat{y}^T y \text{ is a symmetric matrix}$$

③ $\boxed{E(\beta) = y^T y - 2\hat{y}^T y + \hat{y}^T \hat{y} + \lambda \beta^T \beta}$

↳ find such value for $\beta$ matrix for which $E(\beta)$ is minimum

$$\boxed{\frac{dE}{d\beta} = 0}$$

$$E = y^T y - 2\hat{y}^T y + \hat{y}^T \hat{y} + \lambda \beta^T \beta$$

$$= y^T y - 2\beta^T X^T y + \beta^T X^T X \beta + \lambda \beta^T \beta$$

$$\frac{dE}{d\beta} = -2 X^T y + 2 X^T X \beta + 2\lambda \beta = 0 \qquad \to \text{let } X^T X = A, \text{ given } A \text{ is symmetric}$$
$$\frac{d(x^T A x)}{dx} = 2 A x$$

$$-X^T y + A\beta + \lambda\beta = 0$$

$$A\beta + \lambda\beta = X^T y$$

$$\beta^T A + \beta^T \lambda = y^T X$$

$$\beta^T (A + \lambda) = y^T X$$

$$\beta^T (A + \lambda)(A + \lambda)^{-1} = y^T X (A + \lambda)^{-1}$$

$$\beta^T = y^T X (A + \lambda)^{-1}$$

$$\beta = [y^T X (A + \lambda)^{-1}]^T$$

$$\beta = [(A+\lambda)^{-1}]^T X^T y \quad \to \text{need to prove } (A+\lambda)^{-1} \text{ is symmetric matrix}$$
$$\text{i.e. } [(A+\lambda)^{-1}]^T = (A+\lambda)^{-1}$$

$$= (A+\lambda)^{-1} X^T y \qquad \text{we know } A \text{ is a symmetric matrix}$$
$$\text{so } A+\lambda \text{ is also a symmetric matrix}$$

$$\boxed{\beta = (X^T X + \lambda I)^{-1} X^T y} \qquad \text{let } A + \lambda = C$$

(from sklearn.linear_model import Ridge
reg = Ridge(alpha, linear_solvers))

$\text{let } A + \lambda = C$

$$C C^{-1} = I$$
$$[C^{-1}]^T C = I$$
$$[C^{-1}]^T C C^{-1} = I C^{-1}$$
$$[C^{-1}]^T I = I C^{-1}$$
$$[C^{-1}]^T = C^{-1} \quad \therefore [(A+\lambda)^{-1}]^T = (A+\lambda)^{-1}$$

## With Gradient Descent

$$\beta_m = \beta_{old} - \eta \frac{dL}{d\beta_m}$$

Assume only one feature, with two records.

| $X_1$ | $Y$ |
|---|---|
| $x_1$ | $y_1$ |
| $x_2$ | $y_2$ |

$$L = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2 + \lambda \beta_1^2$$

$$= (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \lambda \beta_1^2$$

$$= (y_1 - \beta_0 - \beta_1 x_1)^2 + (y_2 - \beta_0 - \beta_1 x_2)^2 + \lambda \beta_1^2$$

How?

$$\frac{dL}{dw} = X^T X w - X^T y + \lambda w$$

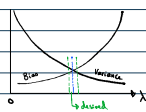$$\frac{dL}{d\beta_0} = (-2)(y_1 - \hat{y}_1) + (-2)(y_2 - \hat{y}_2)$$

$$= -2\sum_{i=1}^{n}(y_i - \hat{y}_i)$$

$$\frac{dL}{d\beta_1} = -2x_1(y_1 - \hat{y}_1) - 2x_2(y_2 - \hat{y}_2) + 2\lambda \beta_1$$

$$= -2\sum_{i=1}^{n}x_i(y_i - \hat{y}_i) + 2\lambda \beta_1$$

## Points to Remember

1. As $\lambda \uparrow$, coefficients
   - approaches 0
   - shrink
   - very close to 0 **But Never 0**
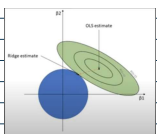   - converge toward 0  ↳ very important

2. Higher coefficient values are impacted more in comparison to lower coefficients values.

3. Bias Variance Tradeoff : when we apply regularization we aim to ↑ bias ↓ variance i.e. generalize our model. This is subjected to the hyperparameter $\lambda$.



4. Effect on loss function :
   → As $\lambda \uparrow$
   - loss function tends to shift towards the origin
   - loss function is translated upwards
   - it also shrinks

5. Why is it called Ridge — [study Hard constraint Ridge Regression]

6. Apply Ridge when coef >= 2 & when we don't want to remove any features

7. Can deal with multicollinearity

class `sklearn.linear_model.`**`Ridge`**(*alpha=1.0, *, fit_intercept=True, copy_X=True, max_iter=None, tol=0.0001, solver='auto', positive=False, random_state=None*)　[source]

1. Can deal with multicollinearity