# Bias-variance tradeoff

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \tilde{y}_i)^2 = E\left[(y-\tilde{y})^2\right]$$

$$= E\left[(y - E[\tilde{y}])^2\right] + E\left[(\tilde{y} - E[\tilde{y}])^2\right]$$

Bias $\qquad\qquad$ variance

$$+ \sigma^2$$





- ordinary least squares (OLS)

  - Design matrix $X \in \mathbb{R}^{n \times p}$
  - output $y \in \mathbb{R}^n$
  - Estimator $\beta \in \mathbb{R}^p$

$$\beta = \left(X^T X\right)^{-1} X^T y$$

$$\min_{\beta \in \mathbb{R}^D} \frac{1}{n} \left\{ (y - X\beta)^T (y - X\beta) \right\}$$

OLS

$$= \frac{1}{n} \| (y - X\beta) \|_2^2 = C(\beta)$$

$$\| x \|_2 = \sqrt{\sum_i x_i^2}$$

Ridge

$$\min_{\beta \in \mathbb{R}^P} = \boxed{\frac{1}{n} \| y - X\beta \|_2^2 + \lambda \beta^T \beta}$$

$$\underline{\lambda \| \beta \|_2^2}$$

For $\lambda > 0$   $MSE(\lambda) < MSE(OLS)$
$MSE(Ridge)$

OLS

$$\frac{\partial C(\beta)}{\partial \beta} \Rightarrow \beta = \underline{\left(X^T X\right)^{-1} X^T y}$$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 = \left( X^T (y - X\beta) \right) \times (-2)$$

Ridge

$$\frac{\partial C(\beta, \lambda)}{\partial \beta} = 0 = X^T (y - X\beta) - \lambda \beta$$

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

$$\beta = \left( X^T X + \boxed{\lambda I} \right)^{-1} X^T y$$

with finite $\lambda > 0$, no problem with divergencies

the matrix $X^T X + \lambda I$

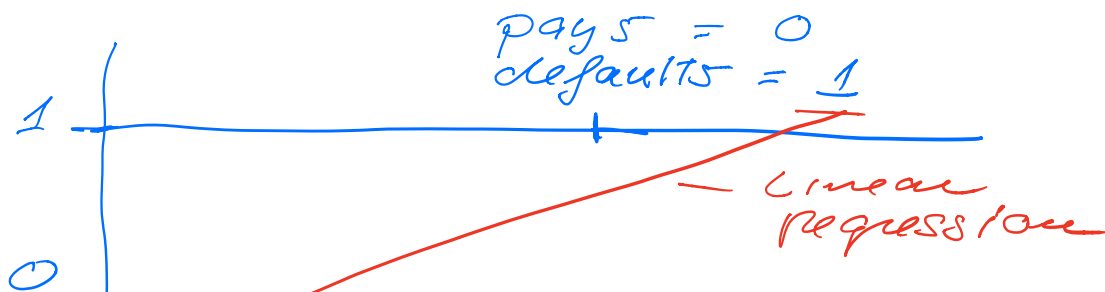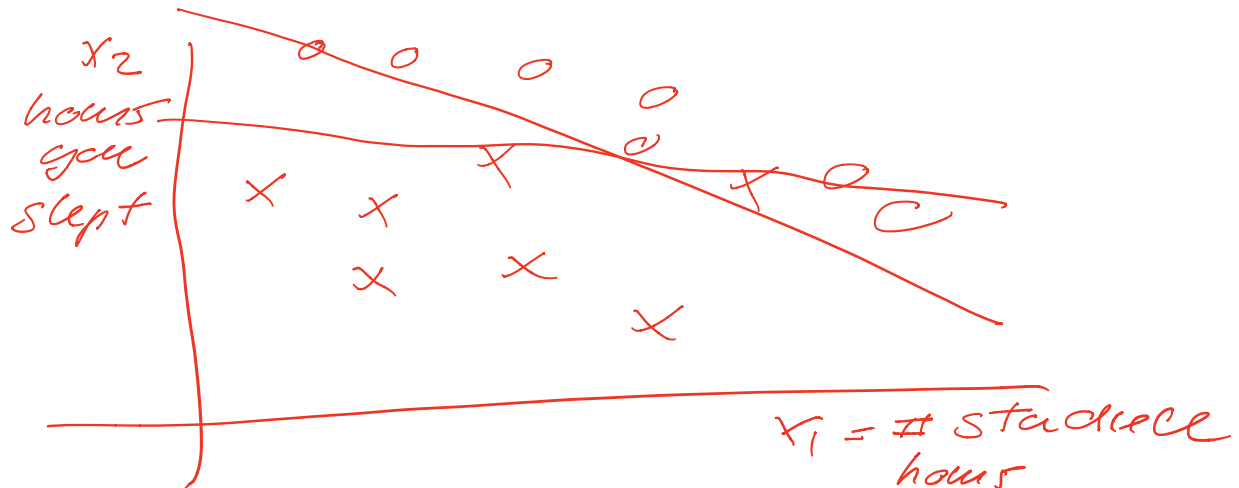$$X^T X \in \mathbb{R}^{P \times P} \quad I \in \mathbb{R}^{P \times P}$$

Lasso optimization :
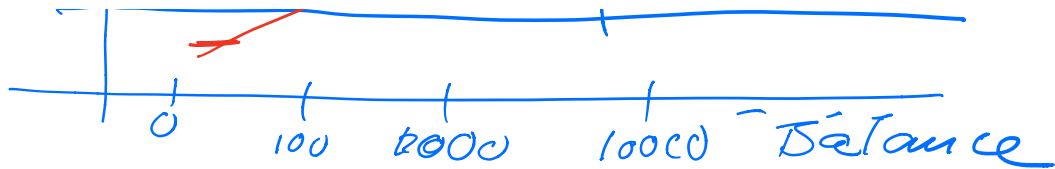
$$C(\beta, \lambda) = \frac{1}{n} \| y - X\beta \|_2^2$$
$$+ \lambda \| \beta \|_1$$

$$\| X \|_1 = \sum_i |x_i|$$

---

classification problem
( Logistic regression )

Example credit card data

pays = 0
defaults = 1

Balance axis:

$0 \quad 100 \quad 8000 \quad 10000$ — Balance

$x_2$ hours you slept

$x_1 = \#$ studied hours

# Logistic Regression

$$p(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}$$

$$1 - p(t) = p(-t)$$

$$t_i = \beta_0 + \beta_1 x_i \qquad \text{Model}$$

Probability $\quad \beta_0 \, \beta_1$

$$P(y_i = 1 \mid x_i \beta) = \frac{e^{t_i}}{1+e^{t_i}}$$

$$P(y_i = 0 \mid x_i \beta) = 1 - P(y_i = 1 \mid x_i \beta)$$

Maximum likelihood approx (iid)

$$D = \{ (y_0, x_0), (y_1, x_1) \cdots (y_{n-1}, x_{n-1}) \}$$

Probability

$$P(D|\beta) =$$

$$\prod_{i=0}^{n-1} \left( P(y_i = 1 | x_i \beta) \right)^{y_i} \left( 1 - P(y_i = 1 | x_i \beta) \right)^{1 - y_i}$$

How do we find $\beta$?

Define cost function as the log (–negative) of $P(D|\beta)$

$$C(X, \beta) = - \sum_i \left\{ y_i \log\left[ P(y_i = 1) \right] \right.$$

$$\left. + (1 - y_i) \log\left[ 1 - P(y_i = 1) \right] \right\}$$

$$P(t_i) = P(\beta_0 + \beta_1 x_i) = \frac{e^{t_i}}{1 + e^{t_i}}$$

$$\frac{\partial C(X\beta)}{\partial \beta_j} = 0$$

$$\frac{\partial C}{\partial \beta} = - X^T (y - P) \in \mathbb{R}^P$$

$$P = [P_0, P_1 \cdots P_{n-1}]$$

$$y, p \in R^n \quad X \in R^{n \times p}$$

Define $W$ with diagonal matrix elements only

$$p_i(1 - p_i)$$

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = \underbrace{X^T W X}_{} \in R^{p \times p}$$

always positive definite,

Need to solve

$$\boxed{X^T(y - p)} = 0 \qquad \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

non-linear in the unknown parameters $\beta$

OLS $\qquad \beta = (X^T X)^{-1} X^T y$

$\longrightarrow$ Family of gradient descent methods,

$$f(\beta) = 0 \qquad \text{Newton's method,}$$

$$\frac{\partial C(\beta)}{}$$

$$\partial \beta$$

#Iterations $k$

$$\beta_{k+1} = \beta_k - \frac{f(\beta_k)}{f'(\beta_k)}$$

$$\|\beta_{k+1} - \beta_k\|_2 < \varepsilon \approx 10^{-10}$$

$$\beta_{k+1} = \beta_k - \left(\frac{\partial^2 c}{\partial\beta\,\partial\beta^T}\right)^{-1}_{\beta=\beta_k} \times \left(\frac{\partial c}{\partial\beta}\right)_{\beta=\beta_k}$$

$$\beta_{k+1} = \beta_k - \gamma_k \left(\frac{\partial c}{\partial\beta}\right)_{\beta=\beta_k}$$

Learning rate