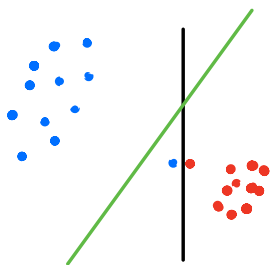


Today Least Squares
Regression



$$\min_{w, b} \frac{1}{2} \|w\|_2^2$$

$$y_i (w^T x_i + b) \geq 1 \quad \forall i$$

Hard-Margin SVM

$$\min \frac{1}{2} \|w\|_2^2 + C \cdot \frac{1}{n} \sum_{i=1}^n s_i$$

Soft-Margin SVM

$$y_i (w^T x_i + b) \geq 1 - s_i \quad \forall i$$

$$s_i \geq 0 \quad \forall i$$

$$\Downarrow$$

$$s_i \geq \max(0, 1 - y_i \cdot (w^T x_i + b))$$

$$= (1 - y_i \cdot (w^T x_i + b))_+$$

$s_i \geq a$
 $s_i \geq 0$

$$\max(0, a) = (a)_+$$

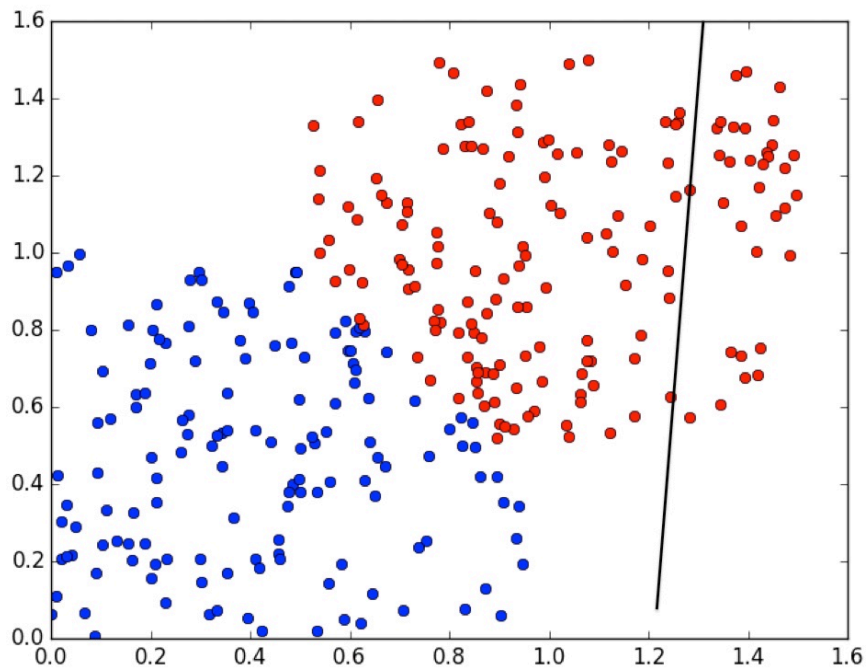
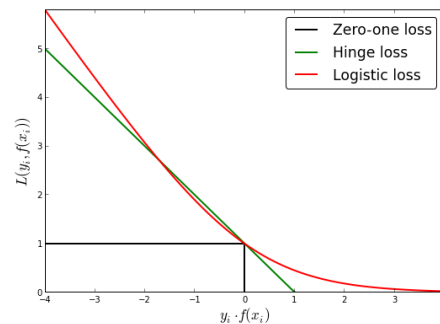
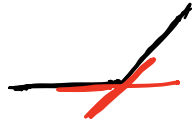
$$\min_{w, b} \quad \frac{1}{2} \|w\|_2^2 + c \frac{1}{n} \sum_{i=1}^n (1 - y_i (w^T x_i + b))_+ \quad \text{Soft Margin SVM.}$$

$$\min_{w, b} \quad \frac{1}{2} \|w\|_2^2 + c \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

- Hinge Loss: $L(y, f(x)) = (1 - yf(x))_+$ where $f(x) = wx + b$
(Soft-Margin SVM)

- Logistic Loss $L(y, f(x)) = \log(1 + e^{-yf(x)})$
(Logistic Regression)

We can use Gradient Descent (sub-gradient Descent)



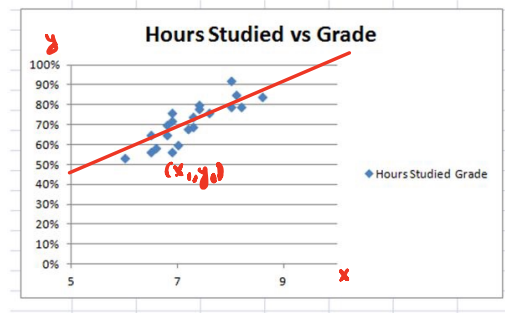
Least Squares Regression

Predict the value of a continuous target variable y

$(x_1, y_1), \dots, (x_n, y_n)$ training data

$x \in \mathbb{R}^d$, $y \in \mathbb{R}$ random.

Linear Regression $f(x) = \vec{w}^T x + b$



Performance measure mean squared error (MSE)

$$R(f) = \mathbb{E}_{x,y} [(y - f(x))^2]$$

$f(x)$ linear $f(x) = \vec{w}^T x + b$

$$R(w, b) = \mathbb{E}_{x,y} [(y - \vec{w}^T x - b)^2]$$

$P_{x,y}$ unknown. We can estimate R directly

$$\hat{R}(w, b) = \frac{1}{n} \sum_{i=1}^n (y_i - \vec{w}^T x_i - b)^2$$

Add a regularization term $\|w\|_2^2$, \dots

$$\min_{w, b} \overbrace{\frac{1}{n} \sum (y_i - \bar{w}^T x_i - b)^2}^{L(w, b)} + \lambda \|w\|_2^2$$

$\lambda = 0 \Rightarrow$ LS Regression

$\lambda > 0 \Rightarrow$ Ridge Regression

$$\frac{\partial}{\partial b} L(w, b) = 2 \frac{1}{n} \sum (y_i - \bar{w}^T x_i - b) (-1) = 0$$

$$b^* = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{w}^T x_i) = \bar{y} - \bar{w}^T \bar{x}$$

where $\bar{y} = \frac{1}{n} \sum y_i$ $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Plugging b^* in L

$$\min_w \frac{1}{n} \sum (y_i - \bar{y} - \bar{w}^T (x_i - \bar{x}))^2 + \lambda \|w\|_2^2$$

Define $\tilde{y}_i = y_i - \bar{y}$ $\tilde{x}_i = x_i - \bar{x}$

$$\min \frac{1}{n} \|\tilde{y} - \tilde{X} w\|_2^2 + \lambda \|w\|_2^2 \quad \tilde{y} = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix} \quad \tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix}$$

$$w^* = (\tilde{X}^T \tilde{X} + n \lambda I)^{-1} \tilde{X}^T \tilde{y} \quad b^* = \bar{y} - (w^*)^T \bar{x}$$

Directly: $O(nd^2)$

\rightarrow Gradient Descent $w_{t+1} = w_t - \alpha_t \cdot 2(\tilde{X}^T \tilde{X} + n \lambda I) w_t - \tilde{X}^T \tilde{y})$

One iteration: $O(nd)$
 $\tilde{x}^T (\tilde{X} w_t)$, $\tilde{x}^T \tilde{y}$

$$= w_t - \alpha_t \sum_{i=1}^n 2 \tilde{x}_i (\tilde{w}^T \tilde{x}_i - \tilde{y}_i) + 2 \lambda w_t$$

- n is large
- dataset does not fit into memory.

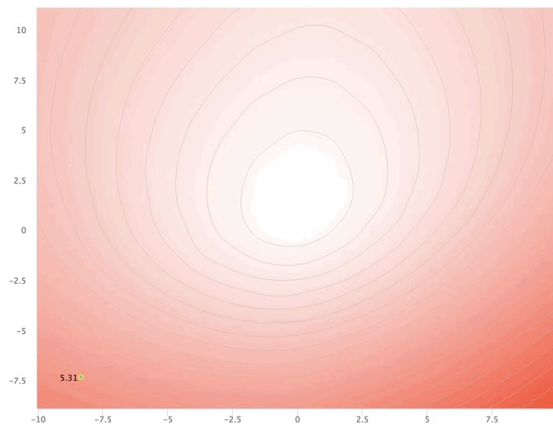
→ Stochastic Gradient Descent

For $i=1, \dots, n$ in random order

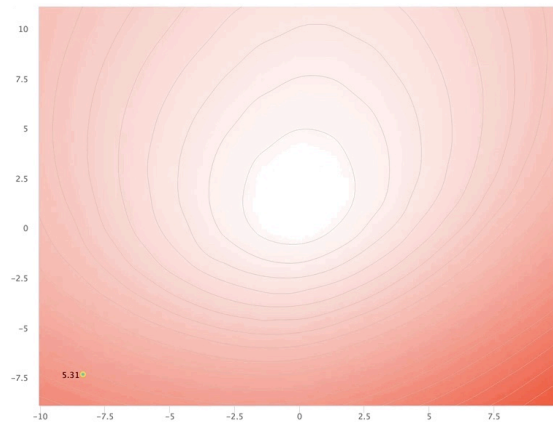
$$w_{t+1} = w_t - 2 \alpha_t (2 \tilde{x}_i (\tilde{w}^T \tilde{x}_i - \tilde{y}_i) + 2 \lambda w_t)$$

One iteration: $O(d)$

- n is extremely large
- online data



GD



SGD

Probability Model: Suppose

$$p(y | x, \theta) = N(\tilde{w}^T x + b, \sigma^2) \Rightarrow E[y_i | x_i] = \tilde{w}^T x_i + b$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$= \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(y - \tilde{w}^T x - b)^2}{2\sigma^2}}$$

$$\arg \max_{\theta} \log \prod_{i=1}^n P(y_i | x_i, \theta)$$

$$\theta_{ML} = \arg \max_{\theta} -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{w}^T x_i - b)^2 - \frac{n}{2} \log 2\pi\sigma^2$$

Maximum Likelihood Estimate = Linear Least Squares

MAP estimation: $p(y | x, \theta) = N(\bar{w}^T x + b, \sigma^2)$

Suppose we have a prior dist on w

$$p(w) = N(0, \tau^2 I) \propto e^{-\frac{\|w\|_2^2}{2\tau^2}}$$

$$\hat{w}_{MAP} = \arg \max_w \log P(y | x, \theta) \cdot P(w)$$

$$= \arg \max_w -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{w}^T x_i - b)^2 - \frac{1}{2\tau^2} \|w\|_2^2$$

$$= \arg \min \sum_{i=1}^n (y_i - \bar{w}^T x_i - b)^2 + \underbrace{\left(\frac{\sigma^2}{\tau^2} \right)}_{=\lambda} \|w\|_2^2$$

$$\arg \min \sum_{i=1}^n (y_i - \bar{w}^T x_i - b)^2 + \lambda \cdot \|w\|_2^2$$