

Applications of Optimization to Machine Learning

Linear Regression :

Dataset :

Training data.

$$\begin{Bmatrix} x_1, x_2, \dots, x_n \\ y_1, y_2, \dots, y_n \end{Bmatrix}$$

$$x_i \in \mathbb{R}^d \quad \forall i$$
$$\underline{y_i} \in \underline{\mathbb{R}}$$

Goal :

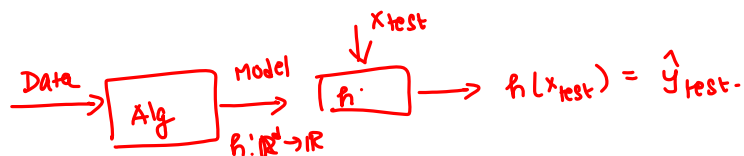
$$\underline{f_2} : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\hat{y}_{\text{test}} = f_2(x_{\text{test}})$$

\uparrow
 $\in \mathbb{R}$

\nwarrow $\in \mathbb{R}^d$

given input during test time.



h is linear $h_w(x) = \underline{w^T x}$ for some $w \in \mathbb{R}^d$.

how to get "Best" h from training data?

Performance measure : Sum of squares error.

$$\min_{w \in \mathbb{R}^d} \underbrace{\sum_{i=1}^n \underbrace{(w^T x_i - y_i)^2}_{\substack{\text{error of } x_i \text{ made by } w.}}}_{f(w)}$$

Specific
Goal
of linear
regression

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$f(w) = \sum_{i=1}^n \frac{(w^T x_i - y_i)^2}{\phantom{f(w) = \sum_{i=1}^n}} \quad ?$$

$$\downarrow$$

$$\sum_{i=1}^n \underline{h_i(w)}$$

If $h_i(w)$ is convex for all i , then $f(w)$ is convex [\because Property 1
Sum of convex
fns is convex]

$$\underline{h_i(w)} = (w^T x_i - y_i)^2$$

$$= f(g(w))$$

\Rightarrow Convex [Composition of
Linear + const and
inward
Convex]

$$g(w) = \underline{w^T x_i - y_i} \quad \begin{array}{l} \text{Linear + Constant} \\ g: \mathbb{R}^d \rightarrow \mathbb{R} \end{array}$$

$$f(z) = \underline{z^2} \quad \begin{array}{l} \text{Convex} \\ f: \mathbb{R} \rightarrow \mathbb{R} \end{array}$$

Conclusion: f is convex.

$$f(w) = \frac{1}{2} \sum_{i=1}^n \underbrace{(w^T x_i - y_i)^2}_{\text{residual}} \\ = \frac{1}{2} \left\| \begin{bmatrix} (w^T x_1 - y_1) \\ \vdots \\ (w^T x_n - y_n) \end{bmatrix} \right\|_2^2$$

$$f(w) = \frac{1}{2} \|Xw - y\|_2^2$$

$$f(w) = \frac{1}{2} (Xw - y)^T (Xw - y) = \frac{1}{2} \left(\underbrace{w^T X^T X w}_{d \times d} - 2w^T \underline{(X^T y)} + \underline{y^T y} \right)$$

$$\nabla f(w) = \frac{1}{2} \left(2 \underline{(X^T X)} w - 2 \underline{(X^T y)} \right) = \underline{(X^T X)} w - \underline{(X^T y)} \quad \begin{matrix} \nearrow \\ w^T A w \\ 2Aw \end{matrix}$$

$$X = \begin{bmatrix} - & x_1 & - \\ - & x_2 & - \\ & \vdots & \\ - & x_n & - \end{bmatrix}_{n \times d} \quad w = \begin{bmatrix} 1 \\ w \end{bmatrix}_{d \times 1} \quad \begin{bmatrix} 1 \\ y \end{bmatrix}_{n \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$z^T z = \|z\|^2 = \sum_{i=1}^n z_i^2$$

$$\nabla f(w) = (X^T X) w - X^T y$$

$$(X^T X) w^* = X^T y \Rightarrow$$

$$w^* = (X^T X)^+ (X^T y)$$

$$(X^T X)^{-1} (X^T y)$$

But $X^T X$ may not be invertible.

Advantages:-

Analytical solution.

$$(X^T X)^+ (X^T y)$$

issues/cons:

Needs an "inverse" computation $\sim O(d^3)$

↳ Can use iterative procedures
[Gradient descent !!]

$$w^{t+1} = w^t - \eta_t \nabla f(w^t)$$

$(X^T X) w - X^T y$ [does not involve an inverse]

$$\underline{\nabla f(w^*)} = \boxed{(\underline{x^T x})w - \underline{x^T y}}$$

Approximation of Gradient

- Stochastic Gradient descent

- Samples a small set of data points uniformly at random.

- Pretend the points sampled from the new dataset and compute gradient w.r.t w

- Can show that $\frac{1}{T} \sum_{t=1}^T w_t \rightarrow w^*$.