

Linear Regression

Derivatives, Gradients, Hessians

$$f(x): \mathbb{R} \rightarrow \mathbb{R}, \quad f'(x) = \frac{df}{dx}: \mathbb{R} \rightarrow \mathbb{R} \quad f(x) = x^2, f'(x) = 2x, f''(x) = 2$$

$$f(w): \mathbb{R}^d \rightarrow \mathbb{R}, \quad \nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad \nabla f = \left(\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right)$$

$$\nabla^2 f: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$$

$$\begin{bmatrix} \frac{\partial^2 f}{\partial w_1 \partial w_1} & & \\ & \ddots & \\ & & \frac{\partial^2 f}{\partial w_i \partial w_j} & \dots \\ & & & \ddots \\ & & & & \dots \\ & & & & & \frac{\partial^2 f}{\partial w_d \partial w_d} \end{bmatrix}$$

Statistical Learning

$$(x_1, y_1), \dots, (x_n, y_n) \quad x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

$l_w(x, y)$: loss function

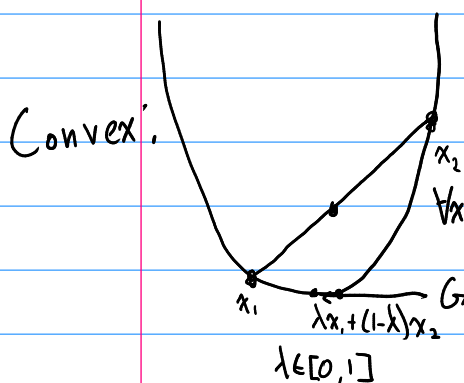
Goal: $\operatorname{argmin}_w E[l_w(x, y)]$, given $(x_i, y_i) \sim D$

$$l_w(x, y) = \begin{cases} 0 & \text{if } \operatorname{sign}(\langle w, x \rangle) = y \\ 1 & \text{else} \end{cases} \Rightarrow \operatorname{argmin}_w \Pr_{(x, y) \sim D} [\operatorname{sign}(\langle w, x \rangle) \neq y]$$

Empirical Risk Minimization (ERM)

$$\operatorname{argmin}_w \frac{1}{n} \sum l_w(x_i, y_i) \rightarrow \operatorname{argmin}_w E[l_w(x, y)]$$

Convexity and Optimization



$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex iff

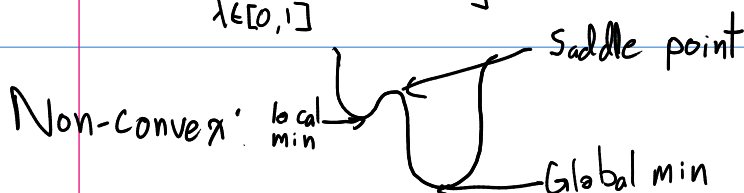
$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

$$\forall x: f''(x) \geq 0$$

(1D fn's)

$$\nabla^2 f(x) \geq 0$$

$M \in \mathbb{R}^{d \times d}$ is Positive Semidefinite (PSD) iff $\forall v Mv \geq 0 \quad \forall v \in \mathbb{R}^d$



Fermat's condition

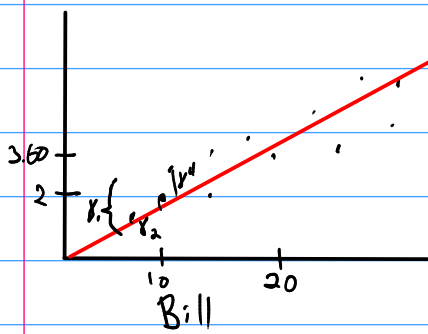
If x is a local extremum of f , then $\nabla f(x) = 0$.

Further, if f convex, converse is true, $\nabla f(x) = 0 \Rightarrow$ global extremum.

Linear Regression

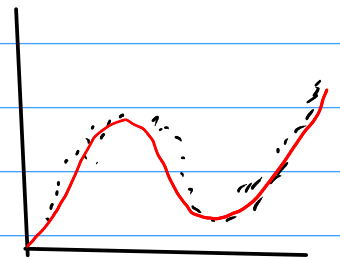
$[X, 1]$ $[w, b]$

Tip



$$l_w(x, y) = (y - \langle w, x \rangle)^2$$

$$\text{Loss} = \sum y_i^2$$



$$\sum (y_i - \langle w, x_i \rangle)^2 \quad A = \begin{bmatrix} -x_1 \\ \vdots \\ -x_n \end{bmatrix} \in \mathbb{R}^{n \times d} \quad z = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$\begin{aligned} \|Aw - z\|_2^2 &= (Aw - z)^T (Aw - z) = (w^T A^T - z^T) (Aw - z) \\ &= w^T A^T A w - z^T A w - w^T A^T z + z^T z \\ &= w^T A^T A w - 2z^T A w + z^T z \end{aligned}$$

Claim: If $f(x) = x^T A x + x^T b + c$, $\nabla f(x) = (A + A^T)x + b$.

$$\begin{aligned} \nabla_w \|Aw - z\|_2^2 &= (A^T A + (A^T A)^T)w - 2A^T z \\ &= 2A^T A w - 2A^T z. \end{aligned}$$

$$\begin{aligned} \nabla_w^2 \|Aw - z\|_2^2 &= 2A^T A. & 2(v^T A^T) A v &\geq 0 \quad \forall v \\ &\geq 0 & &= 2\|Av\|_2^2 \geq 0 \end{aligned}$$

$\therefore \uparrow$ convex

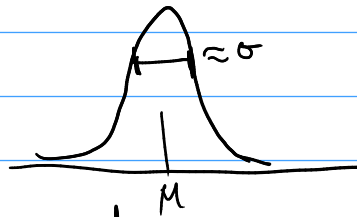
$$\nabla_w \|Aw - z\|_2^2 = 0 \Leftrightarrow 2A^T A w - 2A^T z = 0 \quad \hat{y} = \langle \hat{w}, x \rangle$$

$$A^T A w = A^T z$$

Optimize LS. by finding \hat{w} s.t. $A^T A \hat{w} = A^T z$,
 $\hat{w} = (A^T A)^{-1} A^T z$

Squared loss via Maximum Likelihood Estimation (MLE)

Gaussian $N(\mu, \sigma^2)$
 $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$



$$\operatorname{argmax}_{\text{params}} \Pr[\text{observed data} \mid \text{model params}]$$

$$y_i = \langle w, x_i \rangle + z_i, \quad z \sim N(0, \sigma^2)$$

$$\hat{w} = \operatorname{argmax}_w \Pr[(x_1, y_1) \dots (x_n, y_n) \mid w]$$

$$= \operatorname{argmax}_w \prod_i \Pr[(x_i, y_i) \mid w]$$

$$= \operatorname{argmax}_w \prod_i \Pr[y_i \mid x_i, w] \Pr[x_i \mid w]$$

$$= \operatorname{argmax}_w \prod_i \Pr[y_i \mid x_i, w]$$

$$= \operatorname{argmax}_w \log \left(\prod_i \Pr[y_i \mid x_i, w] \right)$$

$$= \operatorname{argmax}_w \sum \log \Pr[y_i \mid x_i, w]$$

$$y_i \mid x_i, w \sim N(\langle w, x_i \rangle, \sigma^2)$$

$$= \operatorname{argmax}_w \sum \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \langle w, x_i \rangle)^2}{2\sigma^2}\right) \right)$$

$$= \operatorname{argmax}_w -\frac{1}{2\sigma^2} \sum (y_i - \langle w, x_i \rangle)^2$$

$$= \operatorname{argmin}_w \sum (y_i - \langle w, x_i \rangle)^2$$

Regularization



$$X_1 = [0, 1] \quad y_1 = 1$$

$$X_2 = [\varepsilon, 1] \quad y_2 = -1$$

$$\hat{w} = \begin{bmatrix} -2/\varepsilon \\ 1 \end{bmatrix}$$

Tikhonov Reg. Ridge Regression

$$\min \|Aw - z\|_2^2 + \lambda \|w\|_2^2$$

↖ hyperparameter

Lasso

$$\min \|Aw - z\|_2^2 + \lambda \|w\|_1$$

sparse
not closed form solve

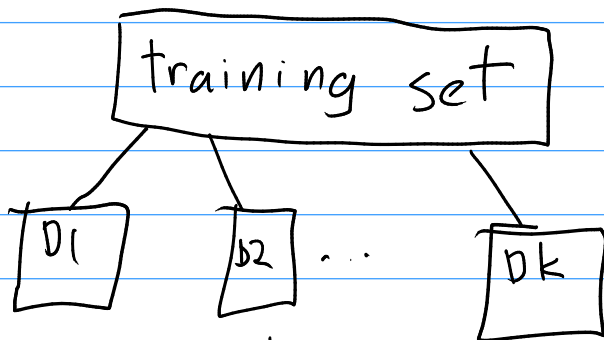
Hyperparameters, Cross Validation

Training dataset

$\lambda \in \{0.01, 0.1, 0.5, 1\}$

~~Validation dataset~~

Test dataset



$k \approx 5? k=10?$

For each λ :

For $i = 1$ to k :

$$w_{\lambda,i} = \text{train} \left(\bigcup_{\substack{j=1 \\ j \neq i}}^k D_j \right)$$

$$\text{perf}_{\lambda,i} = \text{acc}_{w_{\lambda,i}}(D_i)$$

$$\text{perf}_{\lambda} = \sum_{i=1}^k \text{perf}_{\lambda,i}$$

Return $\underset{\lambda}{\text{argmax}} \text{perf}_{\lambda}$

$$\|Aw - z\|_2^2 + \lambda \|w\|_2^2 \leftarrow \text{train}$$

$$\|Aw - z\|_2^2 \leftarrow \text{test, validation}$$

Claim: If $f(x) = x^T A x + x^T b + c$, $\nabla f(x) = (A + A^T)x + b$.

$$\nabla f(x) = \nabla x^T A x + \nabla x^T b + \nabla c$$

$$\nabla c = \vec{0}$$

$$\nabla b^T x, \quad b^T x = \sum_{i=1}^d b_i x_i, \quad \frac{\partial}{\partial x_i} b^T x = b_i \Rightarrow \nabla b^T x = b$$

$$\nabla x^T A x, \quad x^T A x = \sum_{u,v} x_u x_v A_{uv}$$

$$\frac{\partial}{\partial x_i} x^T A x = 2x_i A_{ii} + \sum_{v \neq i} x_v A_{iv} + \sum_{u \neq i} x_u A_{ui} \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$= \underbrace{\sum_v x_v A_{iv}}_{\langle A_{i \cdot}, x \rangle} + \underbrace{\sum_u x_u A_{ui}}_{\sum_u x_u A_{iu}^T}$$

$$(Ax)_i \quad (A^T x)_i$$

$$A = \begin{pmatrix} - & A_{1 \cdot} & - \\ - & A_{i \cdot} & - \\ - & & - \end{pmatrix}$$

$$(Ax)_i = \langle A_{i \cdot}, x \rangle$$

$$\nabla x^T A x = (Ax) + (A^T x) = (A + A^T)x$$