

Linear Regression

ML Pipeline: fit Class F, Training loss L, Optim. Method

Multiple Lin Reg: loss fct: $\ell_{\text{sq}}(f(x), y) = (y - f(x))^2$

$$\ell_{\text{abs}}(f(x), y) = |f(x) - y|, \ell_{\text{huber}}(f(x), y) = \begin{cases} \frac{1}{2}(f(x) - y)^2 & \text{if } |f(x) - y| \leq \delta \\ \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$$

$$\ell_{\text{exp}}(f(x), y) = \mathbb{I}\{\text{max}\{0, f(x)\} + (1-\gamma)\text{max}\{0, f(x)\} - \gamma, \gamma\} \cdot \ell_{\text{sq}}(f(x), y)$$

$\ell_{\text{ce}}(f(x), y) = \text{f underestimates}, \ell_{\text{ce}}(f, y) \rightarrow \text{approximates}$
 $\text{dim } X \times \text{inversible}$

Lin Reg Estimator: $\hat{w} = \underset{w \in \mathbb{R}^n}{\text{argmin}} \|y - Xw\|_2^2 = (X^T X)^{-1} X^T y$

$$L \cdot \hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y = T_x y \quad (\text{projection onto span}(X))$$

$L \cdot \hat{w} \rightarrow X \cdot \hat{w}$ not invertible $\rightarrow \hat{w}$ not unique b/c $\text{dim}(\text{ker}(X^T X)) = 2$

$\hat{w}_{\text{minNorm}} \text{ Solution: } \hat{w} = \underset{w \in \mathbb{R}^n}{\text{argmin}} \|w\|_2 \text{ s.t. } y = Xw$

$$\Rightarrow \hat{w} = X^T X^{-1} y = (X^T X)^{-1} X^T y \quad \text{where } (X^T X)^{-1} X^T X = I_n$$

Nonlinear LS: $\hat{w} = \underset{w \in \mathbb{R}^n}{\text{argmin}} \|y - f(w)\|_2^2 = (\hat{F}^T \hat{F})^{-1} \hat{F}^T y$

$$\hat{F}^T \hat{F} \text{ invertible} \Leftrightarrow \text{rank}(\hat{F}) = p, \hat{F}^T \hat{F} \text{ full rank} \Rightarrow \text{proj.} \rightarrow \text{uniqueness}$$

Remarks: $\hat{w} = (\hat{w}_1, \dots, \hat{w}_n) \Rightarrow w_i = \bar{y} - \bar{w}_i \bar{x}_i, w_i^2 = \frac{\text{cov}(x_i, y)}{\text{var}(x_i)}$

Optimization: Iterative Approach: ① Comp. step: ℓ_{sq} , ② might not converge

$$\text{GD: } w^{t+1} = w^t - \eta \cdot \nabla \ell(w)$$

Conv. of GD for LinReg: GD converges for $\min_w \|y - Xw\|_2^2$

$$\text{if } y \in \text{span}(X), \eta_{\text{opt}} = \frac{2}{\lambda_{\text{min}} + \lambda_{\text{max}}} \text{ well conditioned if } \lambda_{\text{min}} \approx \lambda_{\text{max}}$$

Other Grad. Based Methods: 2nd Order: use Hessian of ℓ

$$\text{Momentum: } w^{t+1} = w^t + \alpha(w^t - w^{t-1}) - \beta \nabla \ell(w^t)$$

$$\text{Adaptive: } w^{t+1} = w^t - \frac{1}{\sqrt{\eta_t}} \frac{\partial \ell}{\partial w}(w^t), \eta_t = \frac{1}{\eta_{\text{min}} + \eta_{\text{max}}} \cdot \eta_t$$

Stochastic GO: random select 1 datapoint x_i for GO ✓

$$(\text{Mini})\text{Batch GO: } S \subset \{1, \dots, n\}: \ell_S(w) = \frac{1}{|S|} \sum_{i \in S} \ell(f_i(w), y_i)$$

L_1 in practice: D-S, u.S, 1k train Model on S's

L_1 is larger \Rightarrow Var. & Comp. Comp. ↑

Remarks: $w^k = (1 - \eta/1)V^T \rightarrow \text{Proj. Mat. to Span}(X), B = V^T(1 - \eta/1)V^T Y$

Convexity: $f((1 + \eta/1)x) \geq f(x) + (1 + \eta/1)f'(x) \Leftrightarrow \ell^2(f(x), y) \geq \ell(f(x), y) + \ell(f'(x), y)$

$\rightarrow f(y) \geq f(x) + \ell(f(x), y - x) \Leftrightarrow$

Lemma: (i) f.g conv. \Leftrightarrow f.conv. & g.conv., (ii) max{f, g} conv. if f convex & g affine

(iii) f.g convex if f convex & g decreasing, g convex

(iv) f convex, diff., x stat point \rightarrow x global min (conv. \Leftrightarrow unique min)

Strongly Convex with w.r.t. \Leftrightarrow f conv. & $\nabla^2 f(x) \succcurlyeq L$

Strongly Convex \Leftrightarrow strictly convex, $\nabla^2 f(x) \succcurlyeq L$

Model Evaluation & Selection:

Estim. Err.: $\ell(\hat{f}(x), f(x))$, Pred. Err.: $\ell(f(x), y)$

Gener. Err.: Average Pred. Error: $\mathbb{E} \sum_i \ell(\hat{f}_i(x_i), y_i)$

Estim. of GE: Empirical Risk: $\ell(\hat{f}(x), f(x)) = \frac{1}{n} \sum_i \ell(f(x_i), y_i)$

$\Rightarrow \hat{f}_0 = \underset{f \in \mathcal{F}}{\text{argmin}} \ell(f, \hat{f}(x))$, GE too optimistic if train & test on same

$D = \mathbb{D}_{\text{train}} \cup \mathbb{D}_{\text{test}}, |\mathbb{D}_{\text{train}}| = |\mathbb{D}_{\text{test}}| \rightarrow$ better for model training,

but test error is not very accurate, $|\mathbb{D}_{\text{test}}| = |\mathbb{D}_{\text{train}}| \rightarrow$ worse for testing

$\ell_{\text{abs}}(f(x), y) = |f(x) - y|, \ell_{\text{huber}}(f(x), y) = \begin{cases} \frac{1}{2}(f(x) - y)^2 & \text{if } |f(x) - y| \leq \delta \\ \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$

$\ell_{\text{exp}}(f(x), y) = \mathbb{I}\{\text{max}\{0, f(x)\} + (1-\gamma)\text{max}\{0, f(x)\} - \gamma, \gamma\} \cdot \ell_{\text{sq}}(f(x), y)$

$\ell_{\text{ce}}(f(x), y) = \text{f underestimates}, \ell_{\text{ce}}(f, y) \rightarrow \text{approximates}$
 $\text{dim } X \times \text{inversible}$

Lin Reg Estimator: $\hat{w} = \underset{w \in \mathbb{R}^n}{\text{argmin}} \|y - Xw\|_2^2 = (X^T X)^{-1} X^T y$

$$\hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y = T_x y \quad (\text{projection onto span}(X))$$

$\hat{w} \rightarrow X \cdot \hat{w}$ not invertible $\rightarrow \hat{w}$ not unique b/c $\text{dim}(\text{ker}(X^T X)) = 2$

$\hat{w}_{\text{minNorm}} \text{ Solution: } \hat{w} = \underset{w \in \mathbb{R}^n}{\text{argmin}} \|w\|_2 \text{ s.t. } y = Xw$

$$\Rightarrow \hat{w} = X^T X^{-1} y = (X^T X)^{-1} X^T y \quad \text{where } (X^T X)^{-1} X^T X = I_n$$

Cross Validation: If $D = D_{\text{train}} \cup D_{\text{test}}$, $D = \mathbb{D}_{\text{train}} \cup \mathbb{D}_{\text{test}} \cup \mathbb{D}_{\text{val}}$ \rightarrow f still depends on D_{val}

\rightarrow Cross Validation: $\text{CV}_k(M) = \frac{1}{k} \sum_{i=1}^k \ell_{\text{sq}}(f_i(x_i), y_i)$

Model Selection w.r.t. set: $D = D_{\text{train}} \cup D_{\text{test}} \cup D_{\text{val}}$

$$\hat{f}_0 = \underset{f \in \mathcal{F}}{\text{argmin}} \ell(f, \hat{f}(x)) = \hat{f}_0(D_{\text{train}})$$

$\hat{f}_0(D_{\text{test}}) \rightarrow$ f calc. on D_val

<

Convolutional Neural Networks

Images can be represented as $\{x_i\}$ matrix of grayscale

$\{x_i\}$ matrix of color

Idea: Units share same weights, θ

& each unit connected to few

Units in the next layer $\begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

1D Conv: $w \in \mathbb{R}^k, x \in \mathbb{R}^d \rightarrow w \cdot x \in \mathbb{R}^{d-k+1}$

CNN: ANN with $\tilde{v}^{(l)} = \varphi(\tilde{w}^{(l)}) = v^{(l-1)}$ layers,

but may also contain fully connected layers

Remark: Weights are called filters.

Multidimensional Convolution:

Consider $n \times m \times 3$

RGB image with

m filters of size $f \times g$.

Padding: Add a frame (couple of white rectangles)

to the image \Rightarrow new dim: $(n+2p) \times (m+2p)$

Stride: Stepsize of filters. Must match dim

of filters: $(m+2p-f) \% S = 0$

Output Dim: $\left(\frac{n+2p-f}{S} + 1\right) \times \left(\frac{m+2p-f}{S} + 1\right) \times m$

Clustering (Unsupervised Learning)

Unsupervised learning: learning without y 's.

(Goal: extract meaningful patterns) \rightarrow Regression

(Goal: clustering \rightarrow Classification)

Standard Clustering: group D into clusters based on similarity

Hierarchical Clust: 1) Form individual clusters $\forall i \in D$ points

2) Iteratively merge pairs of clusters that are closest

K-Means Clustering & Lloyd's Heuristic

Each K-means iter. requires amount of memory polynomial in $m \cdot K$

Each Cluster represented by its center $\mu_j \in \mathbb{R}^d$

each K-means iter. performs one of addition & multiplication operations with $m \cdot K$

Def: Cluster assignment: $z_i = \text{argmax}_j \|x_i - \mu_j\|_2^2$

Goal: learn μ_j in training \rightarrow assign z_i in testing based on μ_j

$\mu_j = \frac{1}{m_j} \sum_i x_i, \forall j=1..K$, non-convex &

$\mu_j = \text{argmin}_{\mu_j} \sum_i \|x_i - \mu_j\|_2^2$ np-hard

Lloyd's Heuristic: Iteratively update $z, \mu_j = \frac{1}{m_j} \sum_i x_i$

Def: $\tilde{z}^{(t)} = (\tilde{x}_1^{(t)}, \tilde{x}_2^{(t)}, \dots, \tilde{x}_m^{(t)})$, $\mu_j^{(t)} = \text{argmin}_{\mu_j} \sum_i \|x_i - \mu_j^{(t)}\|_2^2$

1) Init: $\tilde{x}_i^{(0)} = x_i$, 2) Repeat: $\tilde{x}_i^{(t+1)} = \text{argmin}_{\mu_j} \sum_i \|x_i - \tilde{x}_i^{(t)}\|_2^2$

$\mu_j^{(t+1)} = \frac{1}{m_j} \sum_i \tilde{x}_i^{(t)}, \forall j=1..K$, \rightarrow f+1 until conv.

Convergence Analysis & Limitations:

Then: $\hat{R}(\mu, \tilde{z}) = \sum_i \|x_i - \mu\|_2^2 \rightarrow \hat{R}(\mu^{(t)}, \tilde{z}^{(t)}) \geq \hat{R}(\mu^{(t+1)}, \tilde{z}^{(t)})$

\hookrightarrow Lloyd's Algorithm guaranteed to converge.

But only to local min. Θ : $\text{Convexity} \rightarrow \text{Local min.} \rightarrow \text{Global min.}$

Init Schemes & K-Means++

ML method: efficiently find $f \in F$ \rightarrow Model f

Optimization method \rightarrow Training loss L \rightarrow Function class F

i.i.d. samples \rightarrow Distribution P^* \rightarrow ML inference method: efficiently find $\hat{P} \in P$ \rightarrow Model \hat{P}

Training data D \rightarrow Optimization method \rightarrow Loss on distributions \rightarrow Distribution class P

ideally close \rightarrow e.g. $\hat{P} \neq P^*$ not equally likely

strongly deg. on init. \rightarrow \hat{P} $\neq P^*$ \rightarrow suboptimal sol.

if \hat{P} is not close to P^* \rightarrow \hat{P} $\neq P^*$ \rightarrow suboptimal sol.

$\hat{P} = \text{argmin}_{\hat{P}} \sum_i \hat{P}(x_i) \log \frac{\hat{P}(x_i)}{P^*(x_i)}$

$\hat{P}(x_i) = \text{exp} \left(\frac{\log \hat{P}(x_i)}{\log \hat{P}(x_i)} \right)$

$\hat{P}(x_i) = \text{argmax}_{\hat{P}} \sum_i \hat{P}(x_i) \log \frac{\hat{P}(x_i)}{P^*(x_i)}$

<math