

Binary Linear Classifier

Dataset $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, $x_i \in \mathbb{R}^d$, $y \in \{-1, +1\}$
 Weight Vector $w = (w_1, w_2, \dots, w_d) \in \mathbb{R}^d$ Constant $b \in \mathbb{R}$ Given Input $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$
 Compute $\sum_{i=1}^d w_i x_i > b \rightarrow \hat{y} = +1$ $\sum_{i=1}^d w_i x_i < b \rightarrow \hat{y} = -1$ Indicator Function
 Model Parameters $\Omega = \{w, b\}$ \uparrow y_i (True Label) \hat{y}_i (Predicted Label)


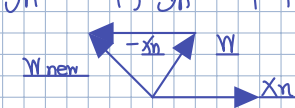
Training Loss Function $J_n(\Omega) = \text{Average \# Of Miss Classified Points} = \frac{1}{N} \sum_{i=1}^N 1\{y_i \neq \hat{y}_i\}$
 Optimal Model Parameters $\Omega^* = \arg\min_{\Omega} J_n(\Omega)$

$x = (x_0=1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$ $w = (w_0=-b, w_1, w_2, \dots, w_d) \in \mathbb{R}^{d+1}$ $w^T x = \sum_{i=0}^d w_i x_i = w_0 x_0 + \sum_{i=1}^d w_i x_i > 0, \hat{y} = +1$ Decision Rule
 $w^T x < 0, \hat{y} = -1$ $\hat{y} = \text{hw}(x) = \text{sign}(w^T x)$

Perceptron Learning Algorithm (Learning Rate Of 1), Input: Linearly Separable Dataset D , Output: $w \in \mathbb{R}^{d+1}$, $J_n(w) = 0$

Initialize w in an arbitrary fashion, $w = (w_0=0, w_1=0, \dots, w_d=0) \in \mathbb{R}^{d+1}$

Step 1: Check If $J_n(w) = 0$, If YES, Output w
 Step 2: Let (x_n, y_n) be a misclassified point in D , If (x_n, y_n) is on the boundary, treat it as misclassified
 If $y_n = +1$ then $w \leftarrow w + x_n$ If $y_n = -1$ then $w \leftarrow w - x_n$ (Go To Step 1)

Geometric Intuition (x_n, y_n) is misclassified, $y_n = +1, \hat{y}_n = -1, w^T x_n \leq 0$ $y_n = -1, \hat{y}_n = +1, w^T x_n \geq 0$



Pocket Algorithm Idea: At each iteration, keep aside the "best weight vector", Run PLA sufficiently many iterations

Linear Regression Data Matrix $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix} (N \times (d+1))$ Observation Vector $y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} (N \times 1)$ Weight Vector $w = \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} ((d+1) \times 1)$ Prediction Vector $\hat{y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{pmatrix} (N \times 1)$
 $\hat{y}_i = x_i^T w$
 $\hat{y} = Xw$
 $y_i = x_{i0}w_0 + x_{i1}w_1 + \dots + x_{id}w_d$
 Loss Function $J_n(w) = \frac{1}{2} \|y - \hat{y}\|^2 = \frac{1}{2} \|y - Xw\|^2$ Regularized Least Square
 $w^* = \arg\min_{w \in \mathbb{R}^{d+1}} J_n(w) = (X^T X)^{-1} X^T y$
 $\arg\min_{w \in \mathbb{R}^{d+1}} \{ \frac{1}{2} \|y - Xw\|^2 + \lambda \|w\|^2 \} = w^*$
 $w^* = (X^T X + \lambda I)^{-1} X^T y$ Regularization Coefficient

Least Square Solution
 Polynomial Fitting, $x \rightarrow z = (1, x, x^2, \dots, x^d)$
 $\nabla f(w) = \begin{pmatrix} \frac{\partial f(w)}{\partial w_1} & \dots & \frac{\partial f(w)}{\partial w_d} \end{pmatrix}^T \in \mathbb{R}^d$ $\nabla_w (w^T A y) = A y$ $\nabla_w (w^T A w) = (A + A^T) w$

Logistic Regression $P_r(y=1|x) = \theta(s) = \theta(w^T x)$
 Sigmoid Activation Function $\hat{P}_w(1|x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$ $\hat{P}_w(-1|x) = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$ $\hat{P}_w(y|x) = \frac{e^{y w^T x}}{1 + e^{y w^T x}}$
 $\theta(s) = \frac{e^s}{1 + e^s}$ (NonLinear)
 Input D , Model Parameter $w \in \mathbb{R}^{d+1}$, Output $[\hat{P}_w(1|x), \hat{P}_w(-1|x)]$
 Log Loss Function, $-\log \hat{P}_w(y|x)$, $E_n(w) = -\log \hat{P}_w(y_n|x_n) = \log(1 + e^{-y_n w^T x_n})$
 $J_n(w) = \frac{1}{N} \sum_{n=1}^N E_n(w) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n w^T x_n})$ $w^* = \arg\min_{w \in \mathbb{R}^{d+1}} J_n(w)$ $w^* = \arg\min_{w \in \mathbb{R}^{d+1}} (J_n(w) + \lambda \|w\|^2)$

Maximum Likelihood Viewpoint $\Pr(y_1, y_2, \dots, y_N | x_1, x_2, \dots, x_N) = \prod_{n=1}^N \Pr(y_n | x_n)$ Independent Event
 $w^* = \arg\max_{w \in \mathbb{R}^{d+1}} \prod_{n=1}^N \Pr(y_n | x_n) = \arg\max_{w \in \mathbb{R}^{d+1}} \log \left(\prod_{n=1}^N \Pr(y_n | x_n) \right) = \arg\max_{w \in \mathbb{R}^{d+1}} \sum_{n=1}^N \log(\Pr(y_n | x_n)) = \arg\min_{w \in \mathbb{R}^{d+1}} J_n(w)$

Cross Entropy Viewpoint $S = \{s_1, s_2, \dots, s_M\}$, $P = (p(s_1), p(s_2), \dots, p(s_M))$, $Q = (q(s_1), q(s_2), \dots, q(s_M))$
 $E(P, Q) = -\sum_{i=1}^M P(s_i) \log q(s_i)$ $E_n(w) = -\{ \frac{1}{2} \{y_n = +1\} \log \hat{P}_w(1|x_n) + \frac{1}{2} \{y_n = -1\} \log \hat{P}_w(-1|x_n) \}$
 $P_n = (\frac{1}{2} \{y_n = +1\}, \frac{1}{2} \{y_n = -1\})$ $Q_n = (\log \hat{P}_w(1|x_n), \log \hat{P}_w(-1|x_n))$ $J_n(w) = \frac{1}{N} \sum_{n=1}^N E(P_n, \hat{P}_w(x_n))$

Knowledge Distillation Transfer knowledge from a large model to a smaller one
 Training For Logistic Regression (No Analytical Solution) (Gradient Descent)
 1. Initialize $x = x_0$
 2. If $f'(x) \approx 0$, stop and output x
 3. If $f'(x) > 0$, $x = x - \epsilon$
 4. If $f'(x) < 0$, $x = x + \epsilon$
 5. Go to step 2
 6. (Step Size) $\epsilon \in 10^{-3}$

$x_t = x_0 + \delta u$ | $|u| = 1$ | $u^* = \arg \min_u f(x_0 + \delta u)$ Taylor Approximation $f(x_0 + \delta u) \approx f(x_0) + \delta u^T \nabla f(x_0)$
 $u^* = -\frac{\nabla f(x_0)}{\|\nabla f(x_0)\|}$, $x_t = x_{t-1} + \delta_t u^*$ $u \in \mathbb{R}^2$ $x_t = x_{t-1} - \epsilon_t \nabla f(x_{t-1})$ (Step Size) $\epsilon_t \propto \epsilon_t \|\nabla f(x_{t-1})\|$

Gradient Descent Algorithm Initialize x_0 in some arbitrary fashion, $t = 0, 1, 2, \dots$, Compute $g_t = \nabla f(x_t)$
 Select Direction $u_t = -g_t$, Update $x_{t+1} = x_t + \epsilon_t u_t$, Continue until $\|\nabla f(x_t)\| \approx 0$

Linear Regression $\nabla E_n(w) = \frac{1}{N} \sum_{n=1}^N \nabla E_n(w)$ $\nabla_w E_n(w) = 2(w^T x_n - y_n) x_n$ $w_{t+1} = w_t - \epsilon_t \frac{2}{N} \sum_{n=1}^N (w^T x_n - y_n) x_n$ w_t will converge to w^*

Stochastic Gradient Descent $w_{t+1} = w_t - \epsilon_t \nabla E_n(w_t)$ $n \in \text{Uniform}\{1, 2, \dots, N\}$

Mini Batch Gradient Descent Draw M examples at random without replacement $w_{t+1} = w_t - \frac{\epsilon_t}{M} \sum_{n=1}^M \nabla E_{\text{min}}(w_t)$

SGD Update Rule $w_{k+1} = w_k + \epsilon_k \frac{y_n x_n}{1 + e^{w^T x_n}}$ (Logistic Regression)

Multiple Classes Logistic Regression $\Omega = \{w(1), w(2), w(3)\}$ $\hat{P}_\Omega(1|x), \hat{P}_\Omega(2|x), \hat{P}_\Omega(3|x)$

$\hat{P}_\Omega(i|x) = \frac{e^{w(i)x}}{e^{w(1)x} + e^{w(2)x} + e^{w(3)x}}$ $E_n(\Omega) = -\log \hat{P}_\Omega(y_n|x_n)$ $w(1) - w(2) = w$

$\Omega_{k+1} = \Omega_k - \epsilon_k \nabla_{\Omega_k} E_n(\Omega_k)$ Momentum Coefficient

SGD + Momentum Velocity Vector $v_0 = 0$, $g_t = \nabla E_n(w_t)$, $v_t = -\epsilon_t g_t + \mu v_{t-1}$, $w_{t+1} = w_t + v_t$

Nestor Momentum $v_t = \mu v_{t-1} - \epsilon_t \nabla E_n(w_t + \mu v_{t-1})$, $w_{t+1} = w_t + v_t$ GD: $O(1/k)$, Nestor: $O(1/k^2)$

Ada Grad $g_t = \nabla E_n(w_t) \in \mathbb{R}^d$, $v_t = v_{t-1} + g_t \odot g_t$, $v_t(j) = \sum_{q=1}^t g_q^2(j)$ $w_{t+1}(j) = w_t(j) - \frac{\epsilon_t}{\sqrt{v_t(j)}} g_t(j)$

RMS - PROP $v_t = \beta v_{t-1} + (1 - \beta) g_t \odot g_t$

ADAM $v_t = \mu v_{t-1} - \frac{\epsilon_t}{\sqrt{v_t}} \odot g_t$ $w_{t+1} = w_t + v_t$ $w_{t+1} = w_t - \frac{\epsilon_t}{\sqrt{v_t}} \odot g_t$

For Classification, $\sigma \rightarrow$ Sigmoid, For Regression, $\sigma \rightarrow$ Identity Function $\hat{y}_1 = \text{OR}(h_1, h_2)$ $\hat{y}_2 = \text{AND}(h_1, h_2)$

$\text{OR}(h_1, h_2) = \text{sign}(h_1 + h_2 - 1.5)$ $\text{AND}(h_1, h_2) = \text{sign}(h_1 h_2 - 1.5)$ $\text{XOR}(h_1, h_2) = \text{AND}(\hat{y}_1, \hat{y}_2)$ (M = 0.5)

Layer 1, Implement w_1, w_2, w_3 to output h_1, h_2, h_3 Layer 2, Implement y_1, y_2, y_3 Layer 3, $y = \text{OR}(y_1, y_2, y_3)$

Neural Network Input Layer $l=0$, Hidden Layers $1 \leq l \leq L-1$, Output Layer L , $w_{ij}^{(l)}$ weight between node i in $l-1$ and j in l

$d^{(l)}$ Number of nodes excluding bias node in layer l , $x_j^{(l)}$ Output from node j in layer l , $s_j^{(l)}$ Input to node j in layer l

$s^{(l)} = \begin{bmatrix} s_1^{(l)} \\ \vdots \\ s_{d^{(l)}}^{(l)} \end{bmatrix}$ $\theta(s^{(l)}) = \begin{bmatrix} \theta(s_1^{(l)}) \\ \vdots \\ \theta(s_{d^{(l)}}^{(l)}) \end{bmatrix}$ $x^{(l)} = \begin{bmatrix} 1 \\ s^{(l)} \end{bmatrix}$ $w^{(l)} = \begin{cases} w_{ij}^{(l)} \end{cases}$ Row $0 \leq i \leq d^{(l-1)}$ Column $1 \leq j \leq d^{(l)}$

Forward Propagation Algorithm

Input $x^{(0)} = [1 \ x_1 \ x_2 \ \dots \ x_d]^T$

for $l=1, 2, \dots, L$ do
 $s^{(l)} = w^{(l)T} x^{(l-1)}$ $x^{(l)} = \theta(s^{(l)})$ $\sum_{l=1}^L (d^{(l-1)} + 1) d^{(l)} + \sum_{l=1}^L d^{(l)}$

Output $x^{(L)}$

Linear Regression Minimizer Weights

$E_n(w) = \frac{1}{N} \|y - Xw\|^2 = \frac{1}{N} (y - Xw)^T (y - Xw)$

$= \frac{1}{N} [y^T y - y^T Xw - w^T X^T y + w^T X^T X w]$

$\frac{\partial E_n(w)}{\partial w} = \frac{1}{N} [-X^T y - X^T y + (X^T X + X^T X) w]$

$0 = \frac{1}{N} [2X^T X w - 2X^T y]$ $X^T X w = X^T y$ $w = (X^T X)^{-1} X^T y$

Backward Propagation Algorithm

Backward Message

Input (x, y) , $\Omega = \{w^{(1)}, \dots, w^{(L)}\}$ Output $\delta^{(l)}$ for $l=1, 2, \dots, L$

$\frac{\partial e}{\partial w_{ij}^{(l)}} \forall i, j, l$ 1. Run forward propagation to compute $\{s^{(l)}, x^{(l)}\}$ and $e(\Omega) = g(x^{(L)}, y)$

2. $\delta^{(L)} = \nabla_{x^{(L)}} g(x^{(L)}, y) \odot \theta'(s^{(L)})$

3. For $l=L-1, \dots, 1$

$\delta^{(l)} = [w^{(l+1)T} \delta^{(l+1)}] \odot \theta'(s^{(l)})$

$\frac{\partial e}{\partial w_{ij}^{(l)}} = x_i^{(l-1)} \delta_j^{(l)}$ $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \epsilon_k \frac{\partial e}{\partial w_{ij}^{(l)}}$

Matrix Inverse $\text{Adj}(A) = \text{Cofactor}(A)^T$

$\frac{1}{\det(A)} \text{Adj}(A)$ $\text{Sign}(x) = \begin{cases} -1, x \leq 0 \\ 1, x > 0 \end{cases}$

Sigmoid $\theta(s) = \frac{e^s}{1 + e^s}$, $\theta(s) = \theta(s)(1 - \theta(s))$

$\tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$ $\theta'(s) = (1 - \tanh^2(s))$

ReLU $\theta(s) = \max(0, s)$

Leaky ReLU $\theta(s) = \begin{cases} s, s \geq 0 \\ \alpha s, s < 0 \end{cases}$

Data Augmentation

Weights Initialization