

Logistic Regression and Cross Entropy Loss

Binary Classification

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i \in \mathbb{R}^{d+1}$$

$$y_i \in \{-1, +1\}$$

$$\text{Output}(\hat{P}_n(-1 | x), \hat{P}_n(1 | x))$$

$$\hat{P}_n(y | x) = \frac{e^{y w^T x}}{1 + e^{y w^T x}}$$

Loss Function: Log Loss

$$e_n(w) = -\log \hat{P}_n(y | x_n)$$

$$E_n(w) = \frac{1}{N} \sum_{n=1}^N e_n(w)$$

$$D' = \{(x_1, p_1), (x_2, p_2), \dots, (x_N, p_N)\}$$

$$p_i = (p_i(1), p_i(2))$$

$$\text{if } y_i = 1, p_i = (1, 0)$$

$$\text{if } y_i = -1, p_i = (0, 1)$$

$$\begin{aligned} CE(p_i, \hat{P}_w) &= -(p_i(1) \log \hat{P}_w(1 | x_n) + \\ &\quad p_i(2) \log \hat{P}_w(-1 | x_n)) \\ &= e_n(w) \end{aligned}$$

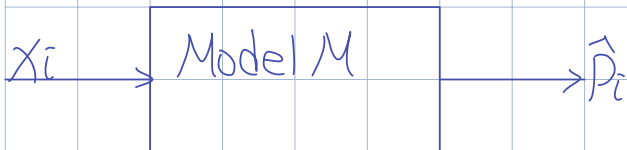
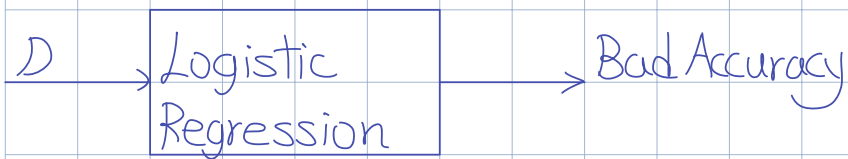
$$E_n(w) = \frac{1}{N} \sum_{n=1}^N CE(p_n, \hat{P}_w(x_n))$$

Knowledge Distillation

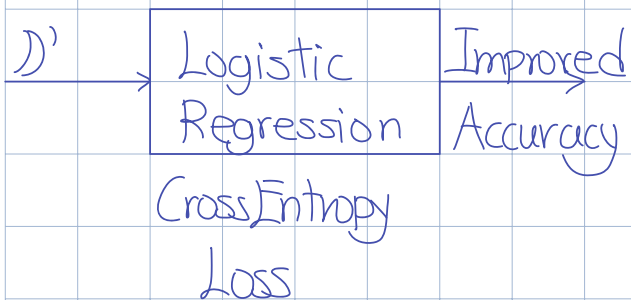


Limited Data

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$



$$D' = \{(x_1, \hat{p}_1), (x_2, \hat{p}_2), \dots, (x_N, \hat{p}_N)\}$$

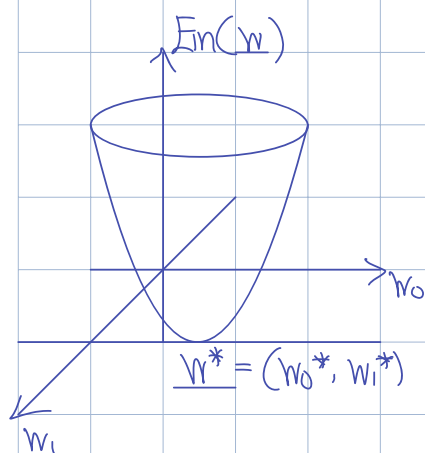


Training for Logistic Regression

Minimize training error

$$E_n(w) = \frac{1}{N} \sum_{n=1}^N e_n(w)$$

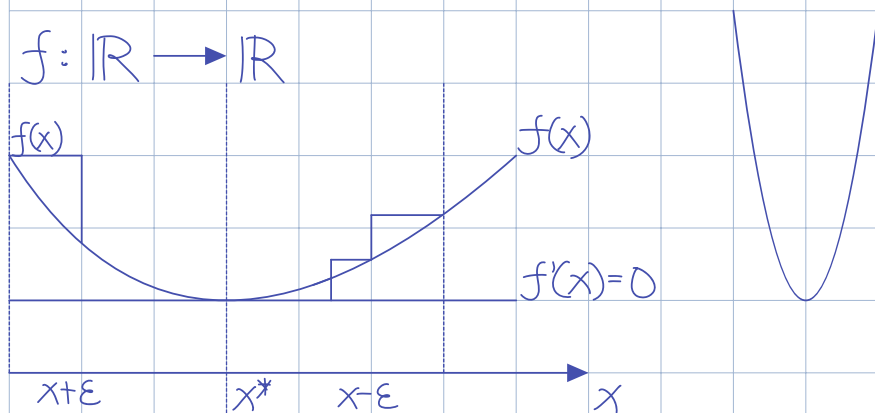
$$e_n(w) = \log(1 + e^{-y_n w^T x_n})$$



$\nabla E_n(w) = 0 \rightarrow$ No analytical solution

To Solve For w^*

Solve numerically using Gradient Descent



If x^* minimizes $f(x)$
 $f'(x^*) = 0$

Analytically computing x^* may not be tractable

Consider any $x \in \mathbb{R}$

if $x < x^*$, $f(x)$ is decreasing, $f'(x) < 0$

if $x > x^*$, $f(x)$ is increasing, $f'(x) > 0$

if $x = x^*$, $f'(x) = 0$

Gradient Descent for $f: \mathbb{R} \rightarrow \mathbb{R}$

1. Initialize $x = x_0$

2. If $f'(x) \approx 0$ then stop & output x

3. If $f'(x) > 0$ then $x = x - \epsilon$

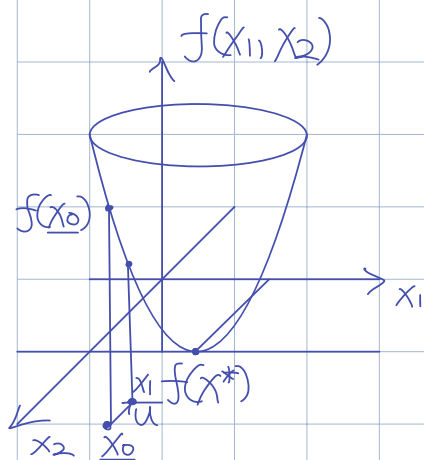
4. If $f'(x) < 0$ then $x = x + \epsilon$

5. Go to step 2

$\epsilon = \text{Step Size}$

Typically decrease ϵ as the iterations progress

$f(x_1, x_2)$



Let x_0 be the initial point

Update Rule

$$x_1 = x_0 + \delta u$$

$u = \text{direction vector}$

$\delta = \text{step size along } u$

δ will be fixed to a small constant
(eg 10^{-3})

We will focus on selecting \underline{u}

$$f(\underline{x}_1) = f(\underline{x}_0 + \delta \underline{u})$$

$$\|\underline{u}\| = 1$$

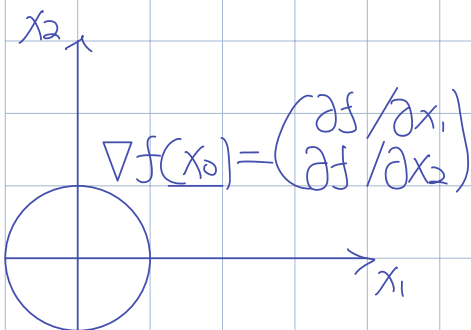
pick \underline{u} among all direction vectors that minimize $f(\underline{x}_1)$

$$\underline{u}^* = \underset{\underline{u} \in \mathbb{R}^2}{\operatorname{argmin}} f(\underline{x}_0 + \delta \underline{u})$$

Since δ is small, we can apply the Taylor series approximation
 $f(\underline{x}_0 + \delta \underline{u}) \approx f(\underline{x}_0) + \delta \underline{u}^T \nabla f(\underline{x}_0)$

$$\underset{\underline{u}, \|\underline{u}\|=1}{\operatorname{argmin}} \{ \underbrace{f(\underline{x}_0)}_{\text{Constant}} + \delta \underline{u}^T \nabla f(\underline{x}_0) \}$$

$$= \underset{\underline{u}, \|\underline{u}\|=1}{\operatorname{argmin}} (\underline{u}^T \nabla f(\underline{x}_0))$$



$$\underline{u}^* = \frac{-\nabla f(\underline{x}_0)}{\|\nabla f(\underline{x}_0)\|}$$

Update Rule

$$\underline{x_1} = \underline{x_0} - \delta \frac{\nabla f(x_0)}{\|\nabla f(x_0)\|}$$

$$\underline{x_t} = \underline{x_{t-1}} - \delta_t \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}$$