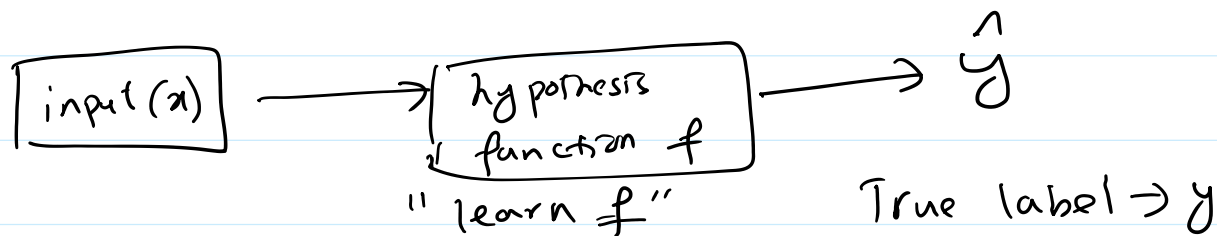


(A) Derivation of logistic Regression



$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

or generally \sim Exponential family (η)

PMF $P(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$

$$h(x) = E[y|x]$$

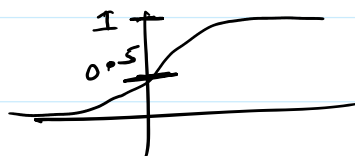
$$\text{or } h_{\theta}(x) = E[y|x; \theta] = \phi = \frac{1}{1 + e^{-x}}$$

$$= \frac{1}{1 + e^{-\theta^T x}}$$

$g(\eta) = E[T(y); \eta]$ is canonical response function

In our case Sigmoid.

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad \text{when } g \text{ is sigmoid}$$



PMF

$$P(y|x; \theta) = \phi^y (1 - \phi)^{1-y}$$

$$= (h_0(x))^y (1 - h_0(x))^{1-y}$$

Likelihood (function of θ)

$$\begin{aligned} L(\theta) &= P(y|x; \theta) \\ &= \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \end{aligned} \left[\begin{array}{l} y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \\ \text{assumption,} \\ \epsilon^{(i)} \text{ are IID} \\ \text{distributed} \end{array} \right]$$

$$= \prod_{i=1}^n (h_0(x^{(i)})^{y^{(i)}} (1 - h_0(x^{(i)}))^{(1-y^{(i)})})$$

Since \log is strictly increasing function,
we could find θ that maximize $\log(L(\theta))$

So,

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^n y^{(i)} \log h_0(x^{(i)}) + (-y^{(i)}) \log (1 - h_0(x^{(i)}))$$

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = (y - h_0(x)) x_j$$

note $[g'(x) = g(x)(1 - g(x))]$

$$\text{So, } \theta := \theta + \alpha (y^{(i)} - h_0(x^{(i)})) x$$

③ Problem set 1(A) and 1(B)

(a) [10 points] In lecture we saw the average empirical loss for logistic regression:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})),$$

where $y^{(i)} \in \{0, 1\}$, $h_{\theta}(x) = g(\theta^T x)$ and $g(z) = 1/(1 + e^{-z})$.

Find the Hessian H of this function, and show that for any vector z , it holds true that

$$z^T H z \geq 0.$$

Hint: You may want to start by showing that $\sum_i \sum_j z_i x_i x_j z_j = (x^T z)^2 \geq 0$. Recall also that $g'(z) = g(z)(1 - g(z))$.

Remark: This is one of the standard ways of showing that the matrix H is positive semi-definite, written " $H \succeq 0$." This implies that J is convex, and has no local minima other than the global one. If you have some other way of showing $H \succeq 0$, you're also welcome to use your method instead of the one above.

Answer:

For simplicity assume we have only one training example on the set. So we have

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_0} J(\theta) \\ \frac{\partial}{\partial \theta_1} J(\theta) \end{bmatrix}$$

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta} (y \log(h_{\theta}(x)) + (1-y) \log(1-h_{\theta}(x)))$$

$$= \frac{y}{h_{\theta}(x)} \frac{\partial}{\partial \theta} h_{\theta}(x) + \frac{(1-y)}{1-h_{\theta}(x)} \frac{\partial}{\partial \theta} (1-h_{\theta}(x))$$

$$= g'(\theta^T x) \left(\frac{y}{h_{\theta}(x)} + \frac{y-1}{1-h_{\theta}(x)} \right)$$

$$\begin{pmatrix} \because h_{\theta}(x) = g(\theta^T x) \\ y(\theta^T x) = g(\theta^T x)(1-g(\theta^T x)) \end{pmatrix} = g'(\theta^T x) x \left(\frac{y - h_{\theta}(x)}{h_{\theta}(x)(1-h_{\theta}(x))} \right)$$

$$= \cancel{g(\theta^T x)} x \frac{(y - h_{\theta}(x))}{\cancel{g(\theta^T x)}} = (y - h_{\theta}(x)) x$$

\Rightarrow st $\frac{\partial}{\partial \theta_j} J(\theta) = (y - h_{\theta}(x)) x_j$

Now

$$H_{ij} = x_i h_{\theta}(x) (h_{\theta}(x) + 1) x_j$$

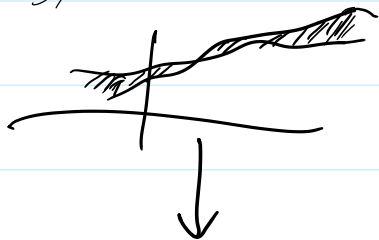
$$\mathbb{E} H = \sum_i \sum_j H_{ij} z_i z_j = \sum_i \sum_j h_{\theta}(x) (h_{\theta}(x) + 1) z_i x_i z_j x_j$$

$$= \log(h(n+1)) \sum_{i,j} z_i z_j$$

$$= \log(h(n+1)) (\sum z)^2$$

≥ 0

So $\sum H z \geq 0$ is true, $H \geq 0$ is this PSD of H is convex



Use input sample to estimate error.



$$E[\text{div}(f(x; \theta), g(x))] \approx \frac{1}{n} \sum_{i=1}^n \text{div}(f(x_i; \theta), \hat{y}_i)$$

Implementation Notes (C)

Instead of iteratively going down hill by ' α ' learn rate (gradient descent),

we used quadratic approximation using

Newton's method. For vector input,

$$\theta := \theta - H^{-1} \nabla_{\theta} l(\theta) \quad (\text{Newton-Raphson method})$$

'numerical analysis'

H is Hessian, $H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\theta)$

Algorithm

Training

① Initialize the parameter θ

② Define $h_0(x) = g(\theta^T x)$ where g is sigmoid
 $y \in \{0, 1\}$ (in ps '1')

② Compute $\nabla_{\theta} J(\theta)$
empirical loss, $\frac{1}{n} \sum (y - h_0(x))x$

③ Compute Hessian H ,
average, $\frac{1}{n} \sum h_0(x)(1-h_0(x))x^T x$

④ Use Newton-Raphson, Repeat until
Convergence

⑤ To predict, feed the input x to
hypothesis, $h_0(x)$ to produce
prediction.

(b) [5 points] **Coding problem.** Follow the instructions in `src/p01b_logreg.py` to train a logistic regression classifier using Newton's Method. Starting with $\theta = \vec{0}$, run Newton's Method until the updates to θ are small: Specifically, train until the first iteration k such that $\|\theta_k - \theta_{k-1}\|_1 < \epsilon$, where $\epsilon = 1 \times 10^{-5}$. Make sure to write your model's predictions to the file specified in the code.

Answer: