

Exercise 4

Sunday, 18 May 2025 10:32 PM

* Exercise 4 : Logistic Regression

- Task 1 :

For a single observation (x_i, y_i) , the Bernoulli model gives

$$P(y_i | x_i; \omega) = p_i^{y_i} (1-p_i)^{1-y_i}$$

where we will set $p_i = \sigma(\omega^T x_i)$

Joint likelihood for all N i.i.d obs.,

$$L(\omega) = \prod_{i=1}^N P(y_i | x_i; \omega) = \prod_{i=1}^N (p_i)^{y_i} (1-p_i)^{1-y_i}$$

Taking log,

$$\ell(\omega) = \log L(\omega) = \sum_{i=1}^N [y_i \log p_i + (1-y_i) \log (1-p_i)]$$

As $p_i = \sigma(\omega^T x_i)$

$$\ell(\omega) = \sum_{i=1}^N [y_i \log[\sigma(\omega^T x_i)] + (1-y_i) \log(1-\sigma(\omega^T x_i))]$$

This is log-likelihood for log. regres. on Bernoulli.

- Task 2 :

a] Convexity of negative log-likelihood

We know that log-likelihood is, its negative would be,

$$-\ell(\omega) = \sum_{i=1}^N [y_i \log[\sigma(\omega^T x_i)] - (1-y_i) \log(1-\sigma(\omega^T x_i))]$$

The grad. of each data pt. is,

$$\nabla_{\omega} [-\ell_i(\omega)] = [\sigma(\omega^T x_i) - y_i] x_i$$

Diff. again & we get hessian,

$$\nabla^2 [-\ell(\omega)] = \sum_{i=1}^N \sigma(\omega^T x_i) [1 - \sigma(\omega^T x_i)] x_i x_i^T$$

As $\sigma(\omega^T x_i) [1 - \sigma(\omega^T x_i)] \geq 0$ & $x_i x_i^T$ is positive semi-definite, their sum is also semi-definite.

$$\therefore \nabla^2 [-\ell(\omega)] \geq 0 \Rightarrow \ell(\omega) \text{ is convex.}$$

I believe convexity guarantees there are local minima.
 Any local min. is global. Thus, grad. descent will converge to global optimum (unique).

b) Binary cross-entropy loss :

We can see that,

$$\max_w l(w) \Leftrightarrow \min_w (-l(w)) \Leftrightarrow \min_w -\frac{1}{N} l(w)$$

Dividing by N turns negative log-likelihood into exactly the average binary cross-entropy:

$$\begin{aligned} -\frac{1}{N} l(w) &= -\frac{1}{N} \sum_{i=1}^N [y_i \log \sigma(w^T x_i) + (1-y_i) \log(1-\sigma(w^T x_i))] \\ &= L_{BCE}(w) \end{aligned}$$

Thus,

$$\arg \max_w l(w) = \arg \min_w L_{BCE}(w)$$

showing that max.-likelihood estimation under Bernoulli GLM is equal. to minimizing the binary cross-entropy loss.

* Task 3 :

a) Divergence of the MLE under linear separability.

The log-likelihood for $\{(x_i, y_i)\}_{i=1}^N$,

$$l(w) = \sum_{i=1}^N [y_i \log \sigma(w^T x_i) + (1-y_i) \log(1-\sigma(w^T x_i))]$$

$$\text{where } \sigma(z) = 1/(1+e^{-z})$$

By assumption there exists some weight w' s.t.,

$$w'^T x_i \begin{cases} > 0 & \text{if } y_i = 1, \\ < 0 & \text{if } y_i = 0 \end{cases}$$

Let's take seq. $w_k = k w' \quad - [k=1, 2, \dots]$

Then for each i ,

- if $y_i = 1$, $\omega_k^T \mathbf{x}_i = k (\omega^T \mathbf{x}_i) \rightarrow +\infty$ so
 $\sigma(\omega_k^T \mathbf{x}_i) \rightarrow 1$ & $\log \sigma(\omega_k^T \mathbf{x}_i) \rightarrow 0$
- If $y=0$, $\omega_k^T \mathbf{x}_i = k (\omega^T \mathbf{x}_i) \rightarrow -\infty$ so
 $\sigma(\omega_k^T \mathbf{x}_i) \rightarrow 0$ & $\log(1 - \sigma(\omega_k^T \mathbf{x}_i)) \rightarrow 0$
- ∴ Every term in $\ell(\omega_k)$ approaches zero &

$$\ell(\omega_k) \rightarrow 0 \text{ as } k \rightarrow \infty$$

For finite ω , $\ell(\omega) < 0$. Thus supremum of log-likelihood is 0, and it attained only in the limit $\|\omega\| \rightarrow \infty$

No finite max. exists, ∴ MLE fails to converge under linear separability.

b) Mitigation:

In order to obtain optimization with a finite sol'n one typically penalizes large weights.

Introducing L_2 or L_1 penalty to neg. log likelihood.

$$\min_{\omega} [-\ell(\omega) + \lambda \|\omega\|_2^2] \text{ or } \min_{\omega} [-\ell(\omega) + \lambda \|\omega\|_1]$$

Also if we place a gaussian (for L_2) or Laplace (L_1) prior on ω yields a MAP estimate that remains finite even when the data are linearly separable. Regularization prevents weight mag. from diverging & often improves generalization.