

Exercise 1

(a) Η maximum likelihood συνάρτηση είναι:

$$P(x_1, \dots, x_n; \theta) = \prod_{i=1}^n (P(\omega_1 | x_i)^{y_i} \cdot P(\omega_2 | x_i)^{1-y_i})$$

Τις όποιες:

$$P(\omega_1 | x_i) = \sigma(\theta^T \cdot x_i)$$

και $P(\omega_1 | x_i) + P(\omega_2 | x_i) = 1$

$$\text{Οπου } \sigma(\theta^T \cdot x) = \frac{1}{1 + e^{-\theta^T \cdot x}}$$

Και επειδή θέλουμε να κάνουμε minimize την Αρνητική log-likelihood:

$$\ln(L(\theta)) = - \sum_{i=1}^n y_i \ln(\sigma(\theta^T x)) + (1-y_i) \ln(1-\sigma(\theta^T x))$$

$$\text{Άρα } \frac{\partial \ln(L(\theta))}{\partial \theta} =$$

$$- \sum_{i=1}^n y_i \frac{1}{\sigma(\theta^T x)} (\sigma(\theta^T x))' + (1-y_i) \cdot \frac{1}{1-\sigma(\theta^T x)} \cdot ((1-\sigma(\theta^T x)))'$$

$$= - \sum_{i=1}^n \left(y_i \frac{1}{\sigma(\theta^T x)} + (1-y_i) \frac{1}{1-\sigma(\theta^T x)} \right) (\sigma(\theta^T x))'$$

$$\begin{bmatrix} \frac{\partial \sigma_1}{\partial \theta_1} \\ \vdots \\ \frac{\partial \sigma_n}{\partial \theta_n} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

$$= - \sum_{i=1}^n \left(y_i \frac{1}{\sigma(\theta^T x)} + (1-y_i) \frac{1}{1-\sigma(\theta^T x)} \right) (\sigma(\theta^T x) (1-\sigma(\theta^T x)))' \frac{\partial \sigma^T x}{\partial \theta}$$

$$= - \sum_{i=1}^n \left(y (1-\sigma(\theta^T x)) - (1-y_i) \cdot \sigma(\theta^T x) \right) x$$

$$= \sum_{i=1}^n (s_i - y_i) \cdot x_i = X^T (S - y)$$

Αρα Gradient Descent:

$$\theta^i = \theta^{i-1} - \mu_i \nabla L(\theta) |_{\theta=\theta^{i-1}} = \theta^{i-1} - \mu_i X^T (S - y)$$

Exercise 2

(a) .ipyub

(b) Θα χρησιμοποιήσουμε gradient descent για να κάνουμε minimise την cost function:

$$J(\theta) = \sum_{i=1}^N (y_i - f(\theta^T x_i))^2$$

$$\boxed{f'(z) = af(z)(1-f(z))}$$

Επίσης μπορούμε να πολλαπλασιάσουμε την cost function με το $\frac{1}{2}$ καθώς δεν θα επηρεάσει το σημείο ελαχίστου που ψάχνουμε. (Για να φύγουν τα δυνάμια)

$$\text{Αρα } J(\theta) = \frac{1}{2} \sum_{i=1}^N (y_i - f(\theta^T x_i))^2$$

$$\frac{\partial J(\theta)}{\partial \theta} = \sum_{i=1}^N (y_i - f(\theta^T x_i)) \cdot (f(\theta^T x_i))'$$

$$= - \sum_{i=1}^N (y_i - f(\theta^T x_i)) \cdot f(\theta^T x_i) \cdot (1 - f(\theta^T x_i)) \cdot \frac{\partial \theta^T x_i}{\partial \theta}$$

$$= - \sum_{i=1}^N ((y_i - f(\theta^T x_i)) \cdot af(\theta^T x_i) \cdot (1 - f(\theta^T x_i)) \cdot x_i$$

$$= -a(y - f) \cdot f \cdot (1 - f) \cdot X^T$$

όπου a = ελεύθερη παράμετρος

$$y = [y_1 \dots y_n]^T$$

$$X^T = [x_1 \dots x_n]^T$$

$$f = [f_1 \dots f_n]^T$$

Αρα το gradient descent σχήμα που θα κάνει minimize την cost function θα είναι:

$$\theta^i = \theta^{i-1} - \mu_i \nabla J(\theta) |_{\theta = \theta^{i-1}} = \theta^{i-1} - \mu_i \cdot (-a)(y-f) \cdot f \cdot (1-f) \cdot X^T$$

(c) Το μοντέλο είναι αδύνατο να επιστρέψει ακριβώς τις τιμές 0 και 1 γιατί παίρνει αυτές τις τιμές στο $\rightarrow +\infty$.

Συνεπώς για να επιστρέφει το 0 θα πρέπει $\theta^T x = -\infty$ και το οποίο είναι πρακτικά αδύνατο.

$$\lim_{\theta^T x \rightarrow +\infty} \sigma(\theta^T x) = 1 \quad \text{και} \quad \lim_{\theta^T x \rightarrow -\infty} \sigma(\theta^T x) = 0$$

(d) Το αποτέλεσμα εκφράζει την πιθανότητα ενός συγκεκριμένου αποτελέσματος. (π.χ. να ανήκει σε κλάση 1)

(e) Όπως παρατηρούμε και από τα plots που δημιουργήσαμε στο ερώτημα (a), αλλάζοντας το a αλλάζει και η κλίση της συνάρτησης. Μικραίνοντας το a βλέπουμε ότι η σιγμοειδής συνάρτηση γίνεται πιο "smooth" και το αντίθετο όταν το αυξάνουμε. Για ένα πάρα πολύ μεγάλο a , η συνάρτηση θα έδινε τιμές 0 και 1.

