

ML 1

2. • The reason we subtract the mean before dividing is that we are trying to get a "standard score" or "z-score" of the data. This transformation makes the $\text{mean} = 0$ and the $\text{std} = 1$
- This wouldn't be achieved if we divided by the standard deviation first
 - Consider:

$$\alpha = [105, 101, 95]$$

$$\frac{\alpha - \text{mean}(\alpha)}{\text{std}(\alpha)} = [0.7259 \quad 0.4148 \quad -1.106]$$

$$\frac{\alpha}{\text{std}(\alpha)} - \text{mean}(\alpha) = [-67.2470, -67.5581, -69.1135]$$

Not centered around zero

2. We use the pre-computed mean of the training set, to normalize the test set for 2 reasons:
1. The training set is much larger and thus the sample mean and sample standard deviation are more representative of the population
 2. It is generally best to not derive anything of your model from the test set, so that benchmarking to other algorithms is much easier

3 ai) The derivative of a function and the derivative of the log of a function are zero at the same place → Thus minimizing one is like minimizing the other

$$\text{aii) } f(x) = \log \left[\prod_{i=1}^m \sigma(a_i^T x)^{b_i} (1 - \sigma(a_i^T x))^{1-b_i} \right]$$

$$\text{Log rules: 1) } \log(a \cdot b) = \log(a) + \log(b)$$

$$2) \quad \log(a^b) = b \log(a)$$

$$f(x) = \sum_{i=1}^m \log \left(\sigma(a_i^T x)^{b_i} (1 - \sigma(a_i^T x))^{1-b_i} \right) \\ \sum_{i=1}^m \left[\log(\sigma(a_i^T x)^{b_i}) + \log(1 - \sigma(a_i^T x))^{1-b_i} \right]$$

$$f(x) = \sum_{i=1}^m \left[b_i \log(\sigma(a_i^T x)) + (1-b_i) \log(1 - \sigma(a_i^T x)) \right]$$

$$\frac{\partial}{\partial x_k} (\alpha_i^T x) = \alpha_i^T \cdot \delta$$

$$f(x) = \sum_{i=1}^m [b_i \log(\sigma(\alpha_i^T x)) + (1-b_i) \log(1-\sigma(\alpha_i^T x))]$$

Note: same as homework 1, except missing $\frac{1}{m}$ & -

$$g(h) = b \log(h) + (1-b) \log(1-h)$$

$$h(z) = \frac{1}{1+e^{-z}} \quad z(x) = \alpha^T x$$

$$\frac{dg}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dz} \cdot \frac{dz}{dx}$$

$$\frac{dg}{dh} = \frac{b}{h} - \frac{1-b}{1-h} = \frac{b(1-h)}{h(1-h)} - \frac{h(1-b)}{h(1-h)} = \frac{b-bh-h+hb}{h(1-h)} = \frac{b-h}{h(1-h)}$$

$$\frac{dh}{dz} = h(z)(1-h(z))$$

$$\begin{aligned} \frac{dz}{dx} &= \alpha^T & \frac{dg}{dx} &= \frac{b-h}{h(1-h)} \cdot \cancel{h(1-h)} \cdot \alpha^T \\ &= (b_i - \sigma(\alpha_i^T x)) \alpha_i^T \end{aligned}$$

$$\sum_{i=1}^m (b_i - \sigma(\alpha_i^T x)) \alpha_i^T$$

$$\nabla f(x) = A^T (b - \sigma(Ax))$$

$$\nabla f(x) = A^T (b - \sigma(Ax))$$

$$= A^T b - A^T \sigma(Ax)$$

$$\nabla^2 f(x) = A A^T \sigma(Ax)(I - \sigma(Ax))$$

- Since Hessian for HW1 is convex
then this should be concave

- 4b iv) Overall, the BTLS version is preferred even though it didn't perform much better.
- I do not have to worry about picking a stepsize as much
 - It converges in less iterations
 - Although each iteration of BTLS took a longer time

4 b v). Surprisingly, the linear regression model performed much better than the logistic regression by a large margin.

- This is surprising since logistic regression was much harder to implement.
- I would definitely choose linear regression model for this scenario, since I got much better results and it was easier. However, I would potentially use logistic regression if I had to deal with more classes (All numbers 1-9). Assuming I fixed the bugs...