

$$1. \quad \sigma(x) = \frac{1}{1+e^{-x}}$$

Oscar Scholin
179 HW 2

$$\begin{aligned} a) \quad \sigma'(x) &= -(1+e^{-x})^{-2} \cdot -e^{-x} = \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \left(\frac{1}{1+e^{-x}} \right) \left(1 - \frac{1}{1+e^{-x}} \right) \\ &= \sigma(x) (1 - \sigma(x)) \quad \square \end{aligned}$$

b) Gradient of log likelihood

$$\text{Likelihood: } L(\theta) = P(x|\theta)$$

$$\text{log likelihood: } \ell(\theta) = \ln(L(\theta))$$

assuming binary classification

We know $L(\theta) = \prod_i p_i^{y_i} (1-p_i)^{1-y_i}$, where y_i is the class label, p_i is the probability predicted by the model.

$$\text{So } -\ell(\theta) = -\sum_i [y_i \ln(p_i) + (1-y_i) \ln(1-p_i)]. \text{ Also, } p_i = \sigma(\underbrace{\vec{\theta}^T \vec{x}_i}_{\vec{\theta} \cdot \vec{x}_i})$$

$$\begin{aligned} \therefore \nabla_{\theta}(-\ell(\theta)) &= -\sum_i y_i \frac{\sigma'(\vec{\theta}^T \vec{x}_i)}{\sigma(\vec{\theta}^T \vec{x}_i)} + (1-y_i) \frac{-\sigma'(\vec{\theta}^T \vec{x}_i)}{1-\sigma(\vec{\theta}^T \vec{x}_i)} \\ &= -\sum_i y_i (1 - \sigma(\vec{\theta}^T \vec{x}_i)) \vec{x}_i - (1-y_i) \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i \quad \text{by part a} \\ &= -\sum_i y_i \vec{x}_i - y_i \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i - \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i + y_i \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i \\ &= \sum_i (\sigma(\vec{\theta}^T \vec{x}_i) - y_i) \vec{x}_i \\ &= \sum_i (\xi_i - y_i) \vec{x}_i \quad \text{for } \xi_i = \sigma(\vec{\theta}^T \vec{x}_i) \\ &= X^T (\vec{\xi} - \vec{y}) \quad (\text{since the above is a dot product}) \end{aligned}$$

c) Compute the Hessian:

$$H = \nabla_{\theta} (\nabla_{\theta}(-\ell))^T$$

$$= \nabla_{\theta} (X^T (\vec{\xi} - \vec{y}))^T$$

$$= \nabla_{\theta} (\vec{\xi}^T X - \vec{y}^T X) \quad \text{independent of } \vec{\theta}$$

$$= \nabla_{\theta} (\vec{\xi}^T X)$$

$$\text{using definition of } \vec{\xi} \text{ and generalizing the dot product}$$

$$= \nabla_{\theta} (\sigma(x \cdot \theta))^T X$$

$$= X^T \text{diag}(\vec{\xi} (1 - \vec{\xi})) X = X^T S X$$

$$\therefore H \text{ is PSD} \iff \text{diag}(\vec{\xi}(1-\vec{\xi})) \text{ PSD.}$$

$$\text{Note } \xi_i = \sigma(\vec{\theta}^T \vec{x}_i). \therefore \sigma(\vec{\theta}^T \vec{x}_i)(1-\sigma(\vec{\theta}^T \vec{x}_i)) \geq 0 \text{ since range } \sigma : [0, 1].$$

$$\therefore H \text{ is PSD.}$$

2. Normalize Gaussian.

$$\text{Compute } Z = \int_{-\infty}^{\infty} e^{-x^2/2\sigma^2} dx = \sqrt{\pi \frac{1}{\frac{1}{2\sigma^2}}} = \boxed{\sqrt{2\pi}\sigma}. \quad \square$$

$$\text{using } \int_{-\infty}^{\infty} e^{-ax^2} dx = \sqrt{\frac{\pi}{a}}$$

3. a) Start with

$$\arg \max_{\vec{w}} \sum_{i=1}^N \ln \mathcal{N}(y_i | \underbrace{w_0 + \vec{w}^T \vec{x}_i}_{\text{predicted } y_i}, \sigma^2) + \underbrace{\sum_{j=1}^D \ln \mathcal{N}(w_j | 0, \tau^2)}_{\text{log likelihood of these weights occurring}}.$$

Note the Gaussian prior is $P(\vec{w}) = \prod_j \mathcal{N}(w_j | 0, \tau^2)$
 i.e. assume each weight comes from Gaussian centered @ 0 variance τ^2 by Occam's razor

In particular, note $\mathcal{N}(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$. Substituting in,

$$\begin{aligned} \arg \max_{\vec{w}} & \sum_{i=1}^N \ln \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(w_0 + \vec{w}^T \vec{x}_i)^2}{2\sigma^2}\right) \right] + \sum_{j=1}^D \ln \left[\frac{1}{\sqrt{2\pi}\tau} \exp\left(-\frac{w_j^2}{2\tau^2}\right) \right] \\ &= \arg \max_{\vec{w}} \sum_{i=1}^N -\frac{(w_0 + \vec{w}^T \vec{x}_i)^2}{2\sigma^2} + \ln \frac{1}{\sqrt{2\pi}\sigma} + \sum_{j=1}^D \frac{-w_j^2}{2\tau^2} + \ln \frac{1}{\sqrt{2\pi}\tau} \\ &= \arg \max_{\vec{w}} -\sum_{i=1}^N \frac{(w_0 + \vec{w}^T \vec{x}_i)^2}{2\sigma^2} + \ln \sqrt{2\pi}\sigma - \sum_{j=1}^D \frac{w_j^2}{2\tau^2} + \ln \sqrt{2\pi}\tau \\ &= \arg \max_{\vec{w}} - \left[\underbrace{N \ln \sqrt{2\pi}\sigma + D \ln \sqrt{2\pi}\tau}_{\text{this term is constant and so doesn't contribute to the optimization}} + \sum_{i=1}^N \frac{(w_0 + \vec{w}^T \vec{x}_i)^2}{2\sigma^2} + \sum_{j=1}^D \frac{w_j^2}{2\tau^2} \right] \end{aligned}$$

Note $\arg \max_{\vec{w}} -(\) = \arg \min_{\vec{w}} (\)$:

$$= \arg \min_{\vec{w}} \sum_{i=1}^N \frac{(w_0 + \vec{w}^T \vec{x}_i)^2}{2\sigma^2} + \sum_{j=1}^D \frac{w_j^2}{2\tau^2}$$

Rescale the problem by $2\sigma^2$ since this will also not affect optimization:

$$= \arg \min_{\vec{w}} \sum_{i=1}^N (w_0 + \vec{w}^T \vec{x}_i)^2 + \sum_{j=1}^D \underbrace{\frac{\sigma^2}{\tau^2}}_{\equiv \lambda} w_j^2,$$

which is Ridge Regression. \square

b) Find solution \vec{x}^* to Ridge Regression.

$$f \equiv \|A\vec{x} - \vec{b}\|^2 + \|\Gamma\vec{x}\|^2$$

Take gradient wrt \bar{x} and set to 0:

$$\begin{aligned}\nabla_{\bar{x}} f &= \nabla_{\bar{x}} \left[(A\bar{x} - \bar{b})^T (A\bar{x} - \bar{b}) + (\Gamma\bar{x})^T (\Gamma\bar{x}) \right] \quad \text{writing out the norms} \\ &= \nabla_{\bar{x}} \left[(\bar{x}^T A^T - \bar{b}^T) (A\bar{x} - \bar{b}) + \bar{x}^T \Gamma^T \Gamma \bar{x} \right] \\ &= \nabla_{\bar{x}} \left[\bar{x}^T A^T A \bar{x} - \underbrace{\bar{x}^T A^T \bar{b} - \bar{b}^T A \bar{x}}_{-2\bar{x}^T A^T \bar{b}} + \bar{b}^T \bar{b} + \bar{x}^T \Gamma^T \Gamma \bar{x} \right] \\ &= 2A^T A \bar{x} - 2A^T \bar{b} + 2\Gamma^T \Gamma \bar{x}\end{aligned}$$

Setting to 0:

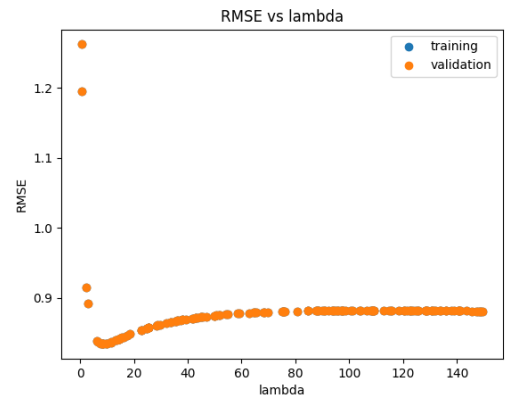
$$2A^T A \bar{x}^* - 2A^T \bar{b} + 2\Gamma^T \Gamma \bar{x}^* = 0$$

$$\begin{aligned}(A^T A + \Gamma^T \Gamma) \bar{x}^* &= A^T \bar{b} \\ \therefore \bar{x}^* &= (A^T A + \Gamma^T \Gamma)^{-1} A^T \bar{b}\end{aligned}$$

$$\text{Let } \Gamma = \sqrt{\lambda} I \Rightarrow \text{minimizes } f = \|A\bar{x} - \bar{b}\|^2 + \lambda \bar{x}^T \bar{x}$$

c) We have the following numerical results:

$$\begin{aligned}\lambda^* &= 8.6497 \\ \text{RMSE on val} &: 0.8340 \\ \text{test} &: 0.8628\end{aligned}$$



d) Solve: minimize $f = \|A\bar{x} + b\mathbf{1} - \bar{y}\|^2 + \|\Gamma\bar{x}\|^2$

$$\begin{aligned}\therefore f &= (A\bar{x} + b\mathbf{1} - \bar{y})^T (A\bar{x} + b\mathbf{1} - \bar{y}) + (\Gamma\bar{x})^T (\Gamma\bar{x}) \\ &= (\bar{x}^T A^T + b^T \mathbf{1}^T - \bar{y}^T) (A\bar{x} + b\mathbf{1} - \bar{y}) + \bar{x}^T \Gamma^T \Gamma \bar{x} \\ &= \bar{x}^T A^T A \bar{x} + 2b^T \mathbf{1}^T A \bar{x} - 2\bar{y}^T A \bar{x} - 2b^T \mathbf{1}^T \bar{y} + \underbrace{b^T \mathbf{1}}_{\mathbf{1}^T \mathbf{1}} + \bar{y}^T \bar{y} + \bar{x}^T \Gamma^T \Gamma \bar{x}\end{aligned}$$

We must find gradient wrt \bar{x} & b since these are fit params:

$$\nabla_{\bar{x}} f = 2A^T A \bar{x} + 2b^T A^T \mathbf{1} - 2\bar{y}^T A + 2\Gamma^T \Gamma \bar{x} \stackrel{!}{=} 0$$

$$\nabla_b f = 2\mathbf{1}^T A \bar{x} - 2\mathbf{1}^T \bar{y} + 2b \mathbf{1}^T \mathbf{1} \stackrel{!}{=} 0$$

$$\therefore b^* = \frac{\mathbf{1}^T (\bar{y} - A\bar{x})}{n}$$

Now solve for \bar{x}^* :

$$A^T A \bar{x} + \Gamma^T \Gamma \bar{x} + \left(\frac{1^T (\bar{y} - A \bar{x})}{n} \right) A^T \mathbf{1} - \bar{y} A^T \mathbf{1} = 0$$

$$\left(A^T A + \Gamma^T \Gamma - \frac{1^T A A^T \mathbf{1}}{n} \right) \bar{x} + \frac{1^T \bar{y} A^T \mathbf{1}}{n} - \bar{y} A^T \mathbf{1} = 0$$

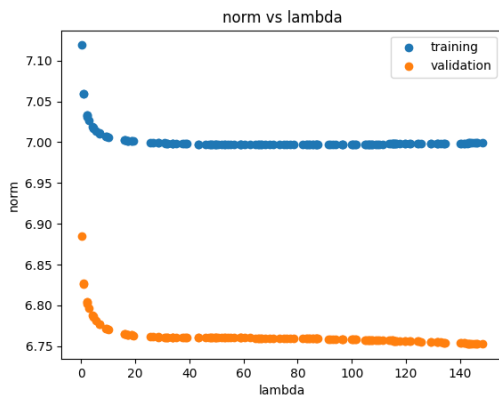
scalar, so can take transpose

$$\left(A^T A + \Gamma^T \Gamma - \frac{1}{n} A^T \mathbf{1} \mathbf{1}^T A \right) \bar{x} + \left(A^T \mathbf{1} \mathbf{1}^T - A^T \right) \bar{y} = 0$$

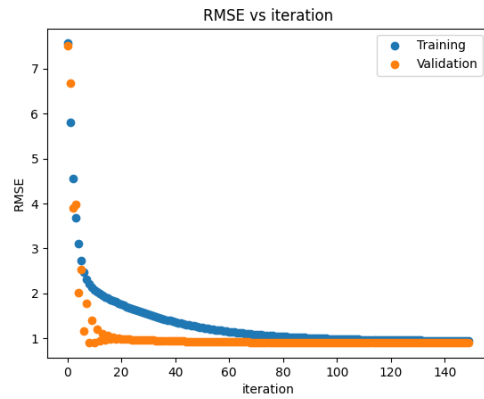
$$A^T \left[\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) A + \Gamma^T \Gamma \right] \bar{x} = A^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \bar{y}$$

$$\therefore \bar{x}^* = \left\{ A^T \left[\left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) A + \Gamma^T \Gamma \right] \right\}^{-1} A^T \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \bar{y}$$

e)



Diff in b : 2.909×10^{-11}
 w : 7.996×10^{-1}



Gradient descent ↑
 Diff in b : 1.539×10^{-1}
 w : 7.996×10^{-1}