

1 (Conditioning a Gaussian) Note that from Murphy page 113. "Equation 4.69 is of such importance in this book that we have put a box around it, so you can easily find it." That equation is important. Read through the proof of the result. Suppose we have a distribution over random variables $\mathbf{x} = (x_1, x_2)$ that is jointly Gaussian with parameters

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix},$$

where

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \mu_2 = 5, \quad \Sigma_{11} = \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}, \quad \Sigma_{21}^T = \Sigma_{12} = \begin{bmatrix} 5 \\ 11 \end{bmatrix}, \quad \Sigma_{22} = [14].$$

Compute

- The marginal distribution $p(x_1)$.
- The marginal distribution $p(x_2)$.
- The conditional distribution $p(x_1|x_2)$
- The conditional distribution $p(x_2|x_1)$

(a) From Murphy, we know $p(\vec{x}_1) = \mathcal{N}(\vec{x}_1 | \vec{\mu}_1, \Sigma_{11})$ so:

$$p(\vec{x}_1) = \mathcal{N}(\vec{x}_1 | \vec{\mu}_1, \Sigma_{11}) = \left| \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}\right) \right| \rightarrow \text{checked sol to make sure this was it...}$$

(b) Also we know $p(\vec{x}_2) = \mathcal{N}(\vec{x}_2 | \vec{\mu}_2, \Sigma_{22})$ so:

$$p(\vec{x}_2) = \mathcal{N}(\vec{x}_2 | \vec{\mu}_2, \Sigma_{22}) = \left| \mathcal{N}(5, 14) \right|$$

(c) From Murphy eq 4.69, $p(\vec{x}_1 | \vec{x}_2) = \mathcal{N}(\vec{x}_1 | \vec{\mu}_{1|2}, \Sigma_{11|2})$

Note: $\vec{\mu}_{1|2} = \vec{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\vec{x}_2 - \vec{\mu}_2)$

$$= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 5 \\ 11 \end{bmatrix} \left(\frac{1}{14} \right) (\vec{x}_2 - 5)$$

$$\vec{\mu}_{1|2} = \frac{1}{14} \begin{bmatrix} 5 \\ 11 \end{bmatrix} (\vec{x}_2 - 5)$$

$$\Sigma_{11|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \quad \text{since } \begin{bmatrix} 5 \\ 11 \end{bmatrix} = \Sigma_{21}^T$$

$$= \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix} - \begin{bmatrix} 5 \\ 11 \end{bmatrix} \frac{1}{14} \begin{bmatrix} 5 & 11 \end{bmatrix}$$

$$= \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix} - \frac{1}{14} \begin{bmatrix} 25 & 55 \\ 55 & 121 \end{bmatrix} = \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix} - \begin{bmatrix} 25/14 & 55/14 \\ 55/14 & 121/14 \end{bmatrix}$$

$$\Sigma_{11|2} = \begin{bmatrix} 59/14 & 57/14 \\ 57/14 & 61/14 \end{bmatrix}$$

(★ From
Murphy
4.69)

(d) Again, from Murphy (extrapolated), $P(\vec{x}_2 | \vec{x}_1) = \mathcal{N}(\vec{\mu}_{2|1}, \Sigma_{2|1})$

$$\star \vec{\mu}_{2|1} = \vec{\mu}_2 + \Sigma_{21} \Sigma_{11}^{-1} (\vec{x}_1 - \vec{\mu}_1) = 5 + [5 \ 11] \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}^{-1} (\vec{x}_1 - \begin{bmatrix} 0 \\ 0 \end{bmatrix})$$

$$= 5 + [5 \ 11] \frac{1}{14} \begin{bmatrix} 13 & -8 \\ -8 & 6 \end{bmatrix} (\vec{x}_1 - 0)$$

$$\hookrightarrow \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}^{-1} = \frac{1}{14} \begin{bmatrix} 13 & -8 \\ -8 & 6 \end{bmatrix}$$

$$= 5 + \frac{1}{14} [-23 \ 26] (\vec{x}_1)$$

$$= 5 + \begin{bmatrix} -23 & 26 \\ 14 & 14 \end{bmatrix} \vec{x}_1$$

$$\boxed{\vec{\mu}_{2|1} = 5 + \begin{bmatrix} -23 & 13 \\ 14 & 7 \end{bmatrix} \vec{x}_1}$$

$$\star \Sigma_{2|1} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} = 14 - [5 \ 11] \frac{1}{14} \begin{bmatrix} 13 & -8 \\ -8 & 6 \end{bmatrix} \begin{bmatrix} 5 \\ 11 \end{bmatrix}$$

$$= 14 - \frac{1}{14} [-23 \ 26] \begin{bmatrix} 5 \\ 11 \end{bmatrix}$$

$$= 14 - \frac{1}{14} (171)$$

$$\boxed{\Sigma_{2|1} = \frac{25}{14}}$$

2 (MNIST) In this problem, we will use the MNIST dataset, a classic in the deep learning literature as a toy dataset to test algorithms on, to set up a model for logistic regression and softmax regression. In the starter code, we have already parsed the data for you. However, you might need internet connection to access the data and therefore successfully run the starter code.

The problem is this: we have images of handwritten digits with 28×28 pixels in each image, as well as the label of which digit $0 \leq \text{label} \leq 9$ the written digit corresponds to. Given a new image of a handwritten digit, we want to be able to predict which digit it is. The format of the data is label, pix-11, pix-12, pix-13, ... where pix-ij is the pixel in the i th row and j th column.

- (a) (**logistic**) Restrict the dataset to only the digits with a label of 0 or 1. Implement L2 regularized logistic regression as a model to compute $\mathbb{P}(y = 1|x)$ for a different value of the regularization parameter λ . Plot the learning curve (objective vs. iteration) when using Newton's Method *and* gradient descent. Plot the accuracy, precision ($p = \mathbb{P}(y = 1|\hat{y} = 1)$), recall ($r = \mathbb{P}(\hat{y} = 1|y = 1)$), and F1-score ($F1 = 2pr/(p+r)$) for different values of λ (try at least 10 different values including $\lambda = 0$) on the test set and report the value of λ which maximizes the accuracy on the test set. What is your accuracy on the test set for this model? Your accuracy should definitely be over 90%.
- (b) (**softmax**) Now we will use the whole dataset and predict the label of each digit using L2 regularized softmax regression (multinomial logistic regression). Implement this using gradient descent, and plot the accuracy on the test set for different values of λ , the regularization parameter. Report the test accuracy for the optimal value of λ as well as its learning curve. Your accuracy should be over 90%.

(a) Implement L2 reg log. regression to find $\mathbb{P}(y=1|\vec{x})$:

Recall log. model : $\mathbb{P}(y=1|\vec{x};\vec{\theta}) = \sigma(\vec{\theta}^T \vec{x})$ \hat{y}

neg. loglikelihood: $n\ell(\theta) = - \sum_i y_i \log \sigma(\vec{\theta}^T \vec{x}_i) + (1-y_i) \log(1 - \sigma(\vec{\theta}^T \vec{x}_i)) + \frac{\lambda}{2} \|\theta\|_2^2$ L2 norm

$$\text{So, } \nabla_{\theta} \ell = \sum_i y_i (\vec{1} - \sigma(\vec{\theta}^T \vec{x}_i)) \vec{x}_i - (1-y_i) \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i + \lambda \vec{\theta}$$

$$= \sum_i [y_i - y_i \sigma(\vec{\theta}^T \vec{x}_i) - \sigma(\vec{\theta}^T \vec{x}_i) + y_i \sigma(\vec{\theta}^T \vec{x}_i)] \vec{x}_i + \lambda \vec{\theta}$$

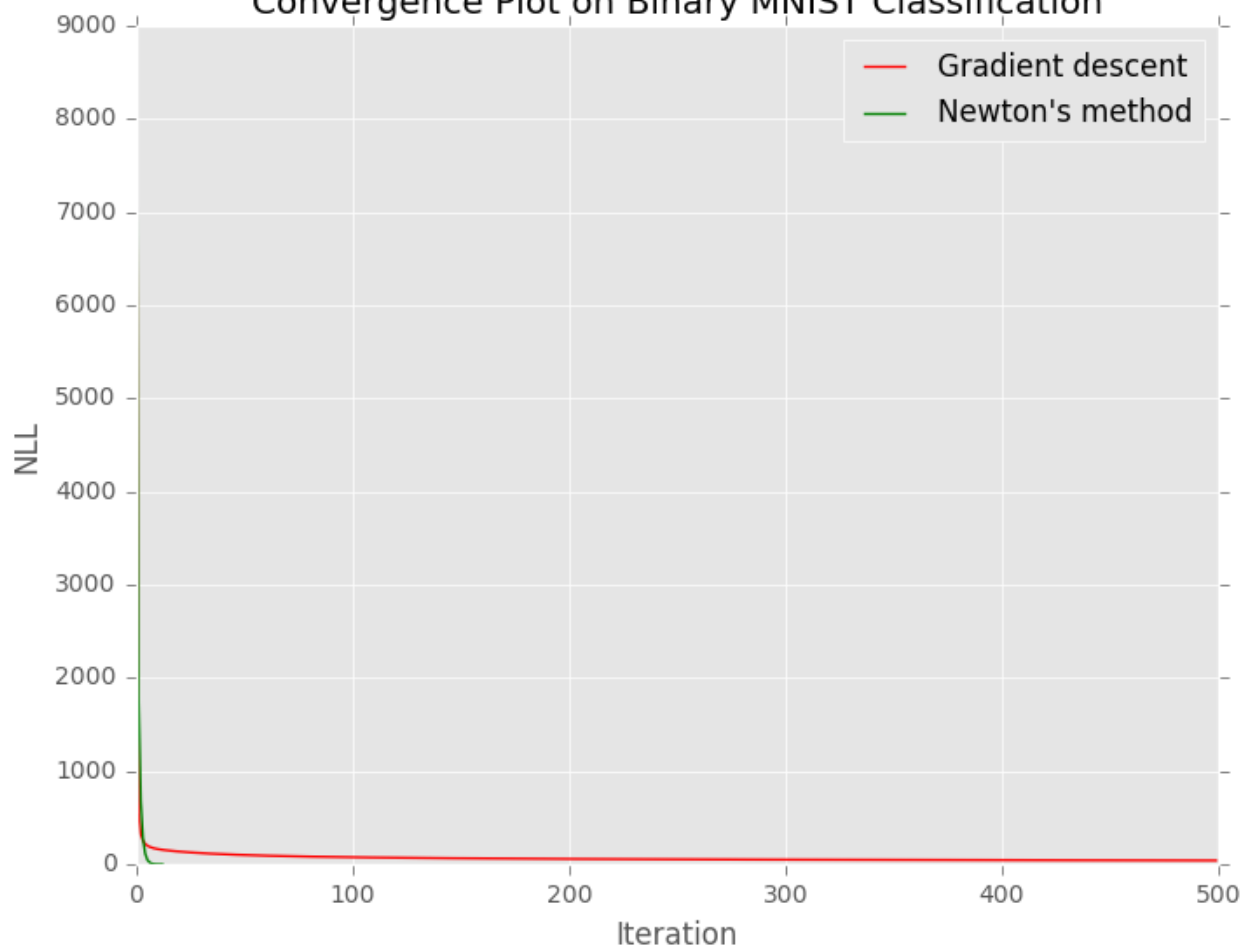
$$\left(\begin{array}{l} \sum_i \vec{x}_i = X \\ \frac{1}{n} A^T B = B^T A \end{array} \right) \Rightarrow = \left[X^T (\sigma(X\vec{\theta}) - \vec{y}) + \lambda \vec{\theta} \right]$$

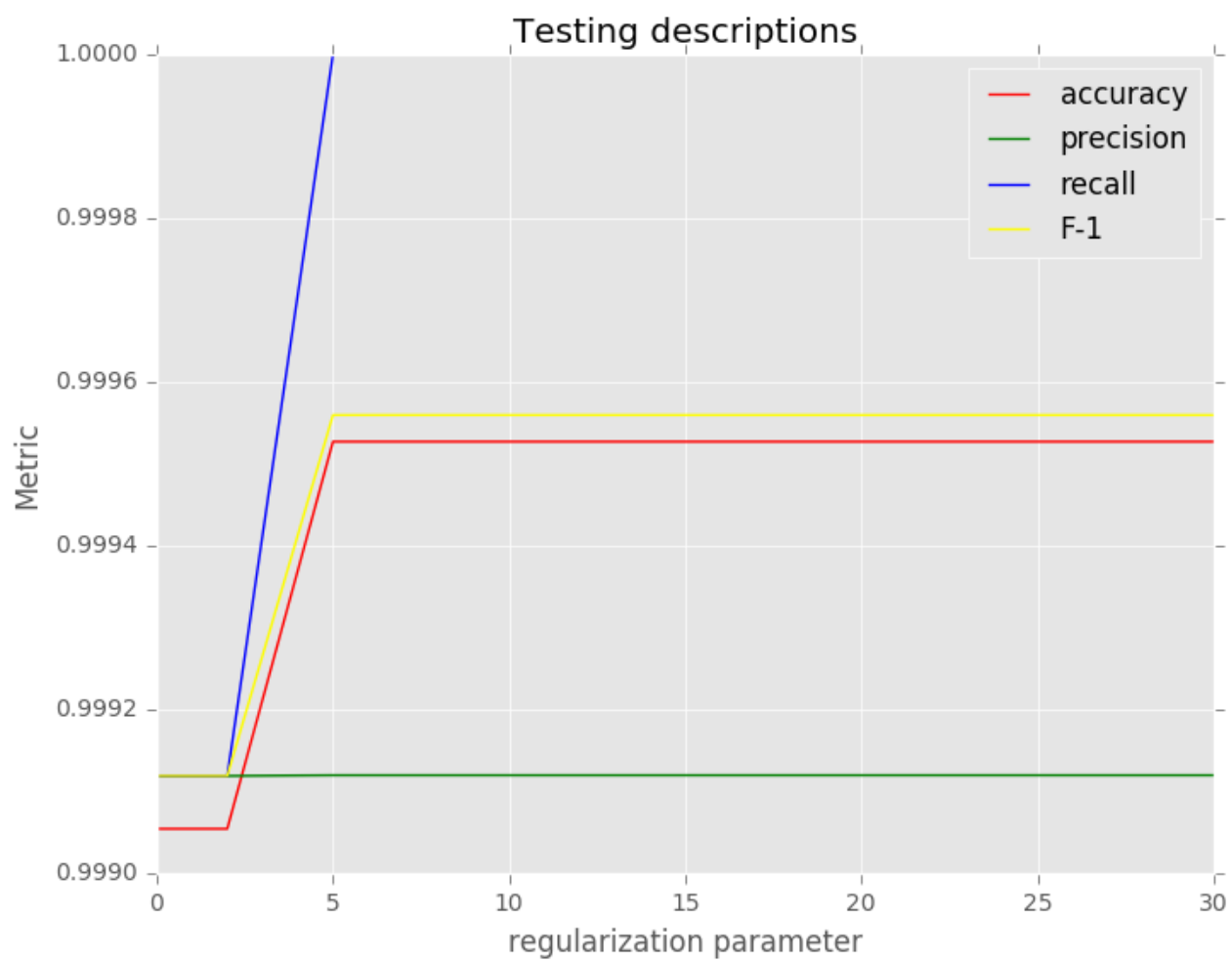
And Hessian: $\nabla^2 \ell = \frac{d}{d\theta} \nabla \ell^T$

$$\begin{aligned} &= \sum_i \nabla_{\theta} \sigma(\vec{\theta}^T \vec{x}_i) \vec{x}_i^T + \lambda \mathbf{I} \\ \text{checked sol'n here} \Rightarrow &= \left[X^T \text{diag}[\sigma(X\vec{\theta})(1-\sigma(X\vec{\theta}))] X + \lambda \mathbf{I} \right] \end{aligned}$$

(see plots attached)

Convergence Plot on Binary MNIST Classification





(b) Softmax $p(y=c | \vec{x}, W) = \frac{1}{Z} \exp(\vec{w}_c^T \vec{x}) = \frac{\exp(\vec{w}_c^T \vec{x})}{\sum_i \exp(\vec{w}_i^T \vec{x})}$ from last wk.

neg log likelihood = $nl(W) = -\log \prod_i \prod_c \mu_{ic}^{y_{ic}} - \lambda \text{tr}(W^T W)$

$\left(\begin{matrix} b/c \log(ab) \\ = \log a + \log b \end{matrix} \right) \hookrightarrow = \sum_i \sum_c y_{ic} \log \mu_{ic} + \lambda \text{tr}(W^T W)$

Then $\nabla_W nl = X^T (\vec{\mu} - \vec{y}) + 2W$ a little similar to (a)

for $y \in \{0, 1\}^{n \times c}$ as one-hot encoding of output y , so

checked
sol'n +
Murphy

$y_{ij} = \begin{cases} 1 & \text{if datapoint } i \text{ is } j \\ 0 & \text{otherwise} \end{cases} \rightarrow \text{predicting label}$

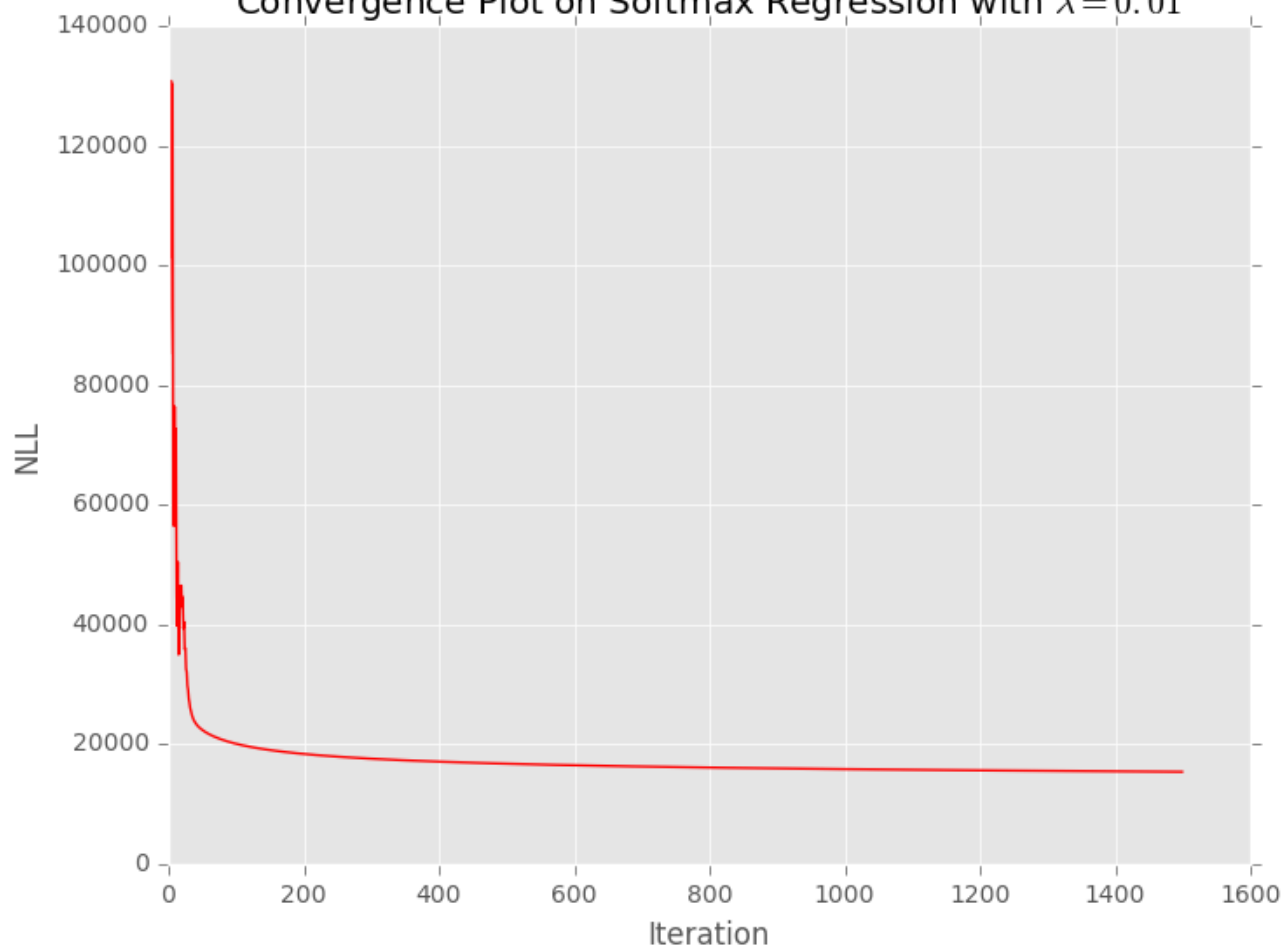
and $\vec{y} \vec{1}_c^T = \vec{1}_n$

Also for $\vec{\mu} \in [0, 1]^{n \times c}$, $\vec{\mu}_i = S(\vec{x}_i) = \frac{\exp(W^T \vec{x})}{\vec{1}^T \exp(W^T \vec{x})}$

checked
sol'n & asked prof. Gu
in office hrs

\rightarrow Accuracy 0.9221 for $\lambda = 0.01$

Convergence Plot on Softmax Regression with $\lambda = 0.01$



Accuracy versus Lambda in Softmax Regression

