

Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

The starter files can be found under the Resource tab on course website. The graphs for problem 2 generated by the sample solution could be found in the corresponding zipfile. These graphs only serve as references to your implementation. You should generate your own graphs for submission. Please print out all the graphs generated by your own code and submit them together with the written part, and make sure you upload the code to your Github repository.

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 $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ that is jointly Gaussian with parameters

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

where

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

Compute

- (a) The marginal distribution $p(\mathbf{x}_1)$.
- (b) The marginal distribution $p(\mathbf{x}_2)$.
- (c) The conditional distribution $p(\mathbf{x}_1|\mathbf{x}_2)$
- (d) The conditional distribution $p(\mathbf{x}_2|\mathbf{x}_1)$

(a) The marginal distribution is:

$$p(x_1) = \mathcal{N}(\mu_1, \Sigma_1 1) = \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}\right)$$

(b) The marginal distribution is:

$$p(x_2) = \mathcal{N}(\mu_2, \Sigma_2 2) = \mathcal{N}(5, 14)$$

(c) We are given

$$p(x_1|x_2) = \mathcal{N}(\mu_{1|2}, \Sigma_{1|2}), \text{ where}$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2) = \frac{1}{14} \begin{bmatrix} 5 \\ 11 \end{bmatrix} (x_2 - 5),$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} = \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix} - \frac{1}{14} \begin{bmatrix} 5 \\ 11 \end{bmatrix} \begin{bmatrix} 5 & 11 \end{bmatrix} = \begin{bmatrix} \frac{59}{14} & \frac{57}{14} \\ \frac{57}{14} & \frac{61}{14} \end{bmatrix}$$

(d) Similarly to part (c), we have

$$p(x_2|x_1) = \mathcal{N}(\mu_{2|1}, \Sigma_{2|1}), \text{ where}$$

$$\mu_{2|1} = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(x_1 - \mu_1) = 5 + \begin{bmatrix} 5 & 11 \end{bmatrix} \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}^{-1} (x_1 - \mu_1) = 5 + \begin{bmatrix} -\frac{23}{14} & \frac{13}{7} \end{bmatrix} x_1,$$

$$\Sigma_{2|1} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} = 14 - \begin{bmatrix} 5 & 11 \end{bmatrix} \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}^{-1} \begin{bmatrix} 5 \\ 11 \end{bmatrix} = \frac{25}{14}$$

■

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 $\mathbb{P}(y = 1|\mathbf{x})$ 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) Recalling the logistic model, we have $P(y = 1|x; \theta) = \sigma(\theta^T x)$, combined with a Gaussian prior on the weights, we obtain the negative log likelihood

$$nll(\theta) = - \sum_i y_i \log \sigma(\theta^T x) + (1 - y_i) \log(1 - \sigma(\theta^T x)) + \frac{\lambda}{2} \|\theta\|_2^2$$

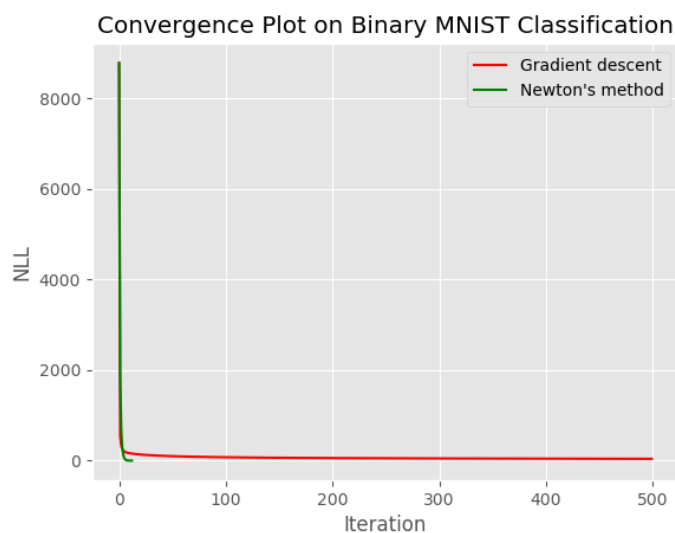
Taking the gradients:

$$\begin{aligned} \nabla_{\theta} l &= \sum_i y_i (1 - \sigma(\theta^T x)) x - (1 - y_i) \sigma(\theta^T x) x + \lambda \theta \\ &= \sum_i [y_i - \sigma(\theta^T x_i)] x_i + \lambda \theta \\ &= X^T (\sigma(X\theta) - y) + \lambda \theta \end{aligned}$$

Thus, we get the Hessian:

$$\begin{aligned}
 \nabla^2 l &= \frac{d}{d\theta} \nabla l^\top \\
 &= \sum_i \nabla_{\theta} \sigma(\theta^\top x) x_i^\top + \lambda I \\
 &= X^\top \text{diag}[\sigma(X\theta)(1 - \sigma(X\theta))]X + \lambda I
 \end{aligned}$$

We plot the learning curves for both gradient descent and Newton's method.



From this figure, we can tell that Newton's method is much faster than gradient descent. Now we plot the accuracy, precision, recall, and F1 score for different values of λ .



For a regularization parameter of 5.0 or above, accuracy is maximized for all parameters.

- (b) Softmax regression provides us with $\mathbb{P}(y = c|x, W) = \frac{1}{Z} \exp(w_c^\top x) = \frac{\exp(w_c^\top x)}{\sum_i \exp(s_i^\top x)}$. Assuming a Gaussian prior on each column of W produces the negative log likelihood:

$$nll(W) = -\log \prod_i \prod_c \mu_{ic}^{y_{ic}} - \lambda \text{tr}(W^\top W)$$

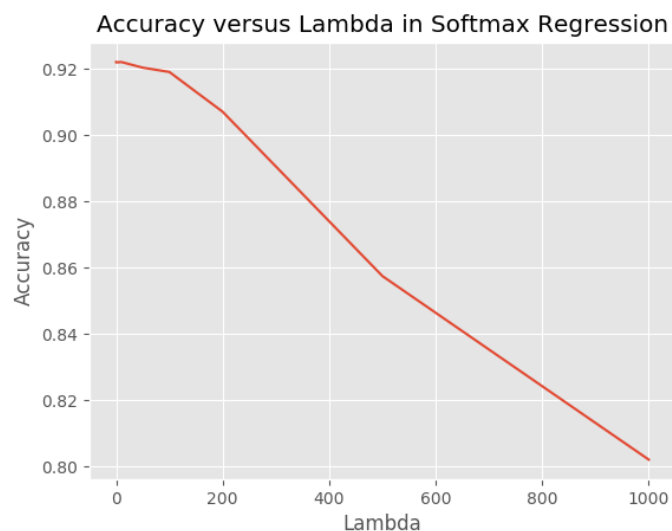
$$= \prod_i \prod_c y_{ic} \log \mu_{ic} + \lambda \text{tr}(W^\top W)$$

$$\nabla_W nll = X^\top (\mu - y) + \lambda W,$$

$y \in \{0,1\}^{n \times c}$ is the one-hot encoding of output y . Also, $\mu \in \{0,1\}^{n \times c}$ is defined as

$$\mu_i = S(x_i) = \frac{\exp(W^\top x)}{1^\top \exp(W^\top x)}$$

Now, we plot the accuracy on the test set for different values of λ using gradient descent.



The maximum test accuracy was 0.9221 when $\lambda = 0.01$. Using this optimal regularization parameter, we plot the convergence plot:

■

