

# CS189–Spring 2016 — Solutions to Homework 3

Shaun Singh, TA

March 13, 2016

## Problem 1: Independence vs. Correlation

- (a) Essentially, there are 4 possible points  $(X, Y)$  can be, all with equal probability ( $\frac{1}{4}$ ):  $\{(0, 1), (0, -1), (1, 0), (-1, 0)\}$ . If graphed onto the Cartesian Plane, these point form "crosshairs".

To show that  $X$  and  $Y$  are uncorrelated, we need to prove:

$$E[(X - \mu_X)(Y - \mu_Y)] = E[X - \mu_X]E[Y - \mu_Y]$$

$$E[XY] = E[X]E[Y] = 0$$

Since for  $\mu_X$  and  $\mu_Y$ , we see that

$$E[X] = E[Y] = \frac{1}{2} * 0 + \frac{1}{2} * \left(\frac{1}{2} + \frac{-1}{2}\right) = 0$$

Notice for that whenever  $X$  is nonzero,  $Y$  is zero (vice versa). Thus,  $E[XY] = 0$  since one of the terms is always zero, and we have shown that  $X$  and  $Y$  are uncorrelated. However, to show that  $X$  and  $Y$  are independent, we must show that:

$$P(X|Y) = P(X)$$

Unfortunately, this is not the case.  $P(X = 0) = \frac{1}{2}$ , but  $P(X = 0|Y = 1) = 0$ . Thus,  $X$  and  $Y$  are not independent.

- (b) To prove independence for  $X$  and  $Y$ , we need to show that  $P(X, Y) = P(X)P(Y)$ . We know that  $P(X|Y) = P(X)$  since  $B_1$  is random 0-1 from the perspective of  $Y$ . Similar arguments can be made for all pairs of random variables. Thus,  $X$ ,  $Y$ , and  $Z$  are pairwise independent.

However,  $P(X|Y, Z) \neq P(X)$  since given the information in Y and Z, we can predict the value in X by the following relation. Thus, it is not mutually independent

$$Y \oplus Z = (B_2 \oplus B_3) \oplus (B_1 \oplus B_3) = B_1 \oplus B_2 \oplus 0 = B_1 \oplus B_2 = X$$

## **Problem 2: Isocontours of Normal Distributions**

- (a) See appendix for graphs

### **Problem 3: Visualizing Eigenvectors of Gaussian Covariance Matrix**

- (a) See appendix for graphs

## Problem 4: Covariance Matrices and Decompositions

- (a) We know that without loss of generality that the covariance matrix  $\Sigma_X \in \mathbb{R}^{N,N}$  corresponding to random variable  $X \in \mathbb{R}^N$  is positive semidefinite, which means that  $\forall x \in \mathbb{R}^N, x^\top \Sigma x \geq 0$ . Unfortunately, we require that  $\Sigma$  be positive definite ( $\forall x \in \mathbb{R}^N, x^\top \Sigma x > 0$ ) in order for it to be invertible, since invertible matrices cannot have any eigenvalues be 0.

When might our covariance matrix  $\Sigma$  have an eigenvalue of 0? The most general case would be when one or more of the RV's are dependent on the others. We require our basis of eigenvectors to be independent in order for the covariance matrix to have no eigenvalues of 0. Another case would be when one or more of the RV's are deterministic.

For example, consider a random variable  $X \in \mathbb{R}^N$  where the first  $n - 1$  elements are standard normal RV's, but the  $n$ -th element is deterministic (for example  $X_n = 15$  always). Thus, all elements for the final row and column of  $\Sigma$  will be 0 since  $\text{cov}(X_i, X_n) = 0 \forall i$ , implying that  $\Sigma$  has an eigenvalue of 0

We can easily convert random variable  $X$  into a new random variable  $X' \in \mathbb{R}^{N-1}$  by deleting the  $N$ th dependent term, resulting in a positive definite covariance matrix  $\Sigma_{X'}$  and invertible. This transformation does not result in any loss of information.

- (b) Use the Spectral Decomposition Theorem to convert  $\Sigma$  into the following, where  $U$  is a unitary matrix of orthonormal eigenvectors  $\vec{e}_i \forall i \in [0 \dots N]$  and  $D$  is a diagonal matrix with eigenvalues  $\lambda_i \forall i \in [0 \dots N]$  located at indices corresponding to eigenvectors in  $U$ . Note: all eigenvalues  $> 0$  since  $\Sigma$  is positive definite

$$\Sigma = UDU^\top \Rightarrow \Sigma^{-1} = (UDU^\top)^{-1} = (U^\top)^{-1}D^{-1}U^{-1} = UD^{-1}U^\top \quad (1)$$

This is since a unitary matrix  $U$  is such that  $U^{-1} = U^\top$ . Note that if diagonal matrix  $D$  has values  $d_{i,i} \forall i$ , then  $D^{-1}$  has value  $\frac{1}{d_{i,i}} \forall i$ . Once again, since  $\Sigma$  was positive definite, the value  $\frac{1}{d_{i,i}}$  exists.

Now, we decompose  $D^{-1}$  into its square-root by defining  $Q$  as a diagonal matrix with diagonal values  $\frac{1}{\sqrt{d_{i,i}}}$ . Verify that  $QQ = D^{-1}$  and that  $Q^\top = Q$ . Thus, we have:

$$\Sigma^{-1} = UD^{-1}U^\top = UQQU^\top = UQQ^\top U^\top \quad (2)$$

$$\Sigma^{-1} = A^\top A \quad (3)$$

Where we defined  $(UQ)^\top = A$ . Therefore

$$x^\top \Sigma^{-1} x = x^\top A^\top A x = (Ax)^\top (Ax) = \|Ax\|_2^2 \quad (4)$$

Note: This process is closely related to Cholesky Decomposition, which will require one to use QR Decomposition to show that  $\Sigma^{-1} = LL^\top$  for all invertible covariance matrices  $\Sigma$  where  $L$  is a lower triangular matrix.

- (c)  $x^\top \Sigma^{-1} x$  is a scalar written in vector quadratic form. It looks like an incomprehensible value, but when we convert it to  $\|Ax\|_2^2$ , we see that in reality its just the squared L2 norm of  $Ax$ , which measures the squared distance from the data vector  $x$  from the mean (in this case 0). Note that we can change the mean to be any arbitrary value without loss of generality.
- (d) Recall from Part B our decomposition for  $\Sigma^{-1}$ , which was as follows where  $U$  is a unitary matrix,  $D$  is a diagonal matrix.

$$\Sigma^{-1} = UD^{-1}U^\top = A^\top A \quad (5)$$

Note that  $\|x\|_2 = 1$  and  $\|Ux\|_2 = 1$  since unitary matrices are orthonormal and preserve magnitude. Define  $q = Ux$ , we have

$$\|Ax\|_2^2 = x^\top A^\top A x = x^\top UD^{-1}U^\top x = q^\top D^{-1}q \quad (6)$$

We can choose our  $x$  such that  $q$  will be any Euclidean Basis Vector  $\vec{e}_i$  such that the  $i$ th element is 1 and all other elements are 0. Therefore, the maximum value that  $\|Ax\|_2^2$  is  $\frac{1}{\lambda_i}$ , where  $\lambda_i$  is the minimum eigenvalue of  $\Sigma$ . The minimum value that  $\|Ax\|_2^2$  is  $\frac{1}{\lambda_j}$ , where  $\lambda_j$  is the maximum eigenvalue of  $\Sigma$ .

If we have  $X_i \perp X_j \forall i, j$ , then  $\text{cov}(X_i, X_j) = 0 \forall i, j$  meaning that off diagonal terms for  $\Sigma$  are 0. Thus, we can find  $\Sigma^{-1}$  directly, where

$$\Sigma_{i,j}^{-1} = \begin{cases} \frac{1}{\sigma_i^2} & \text{if } i == j \\ 0 & \text{else} \end{cases}$$

Therefore, if we have  $X_i \perp X_j \forall i, j$ , the maximum value that  $\|Ax\|_2^2$  is  $\frac{1}{\sigma_i^2}$ , where  $\sigma_i^2$  is the minimum variance. The minimum value of  $\|Ax\|_2^2$  is  $\frac{1}{\sigma_j^2}$ , where  $\sigma_j^2$  is the maximum variance.

To maximize  $f(X)$ , we want the superscript above the exponent to be minimal since there is a negative sign. Thus, for  $\|Ax\|_2^2$  to be minimal, we want to choose  $x$  to be the vector corresponding to the eigenvector corresponding to the maximal eigenvector  $\lambda_j$  or maximum variance  $\sigma_j^2$  if independent.

## Problem 5: Gaussian Classifiers For Digits

- (a) Say we have i.i.d observations  $X_1 \dots X_n$ , then the MLEs for a multivariate Gaussian distribution are as follows

$$\text{MLE of } \hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i, \text{ MLE of } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

Note that the MLEs are just the sample mean and covariances. The MLE of  $\hat{\mu}$  is unbiased, but the MLE of  $\hat{\Sigma}$  is biased, with  $E[\hat{\Sigma}] = \frac{n-1}{n} \Sigma$ . In order to get an unbiased sample, use  $\frac{1}{n-1}$  instead of  $\frac{1}{n}$ .

- (b) For each training set, count all the labels belonging to set  $\theta_i$ . Note that there are 10 sets in total (numbers 0...9). Divide the count for a particular set  $\theta_i$  by the total across all sets. One can thus view the distribution as the sample prior (empirically, how many of a class have we seen in the total data set).

$$P(\theta_i) = \frac{\theta_i}{\sum_{j=1}^{10} \theta_j}$$

- (c) I visualized the covariance matrix for the "zero" class, for the 10000 image set from train-small.mat (See the appendix for details). Upon visualization through imagesc, the 784 x 784 matrix appears to show an X of positive values in the center, with negative values along the edges of the X. The X is enclosed by a 150 zeros on every side. Note, the values themselves are extremely close to zero, about  $+/- 10^{-4}$ . Also, the non-zero values in the visualization appear to be pixelated, as in the X seems to be formed by dots of non-zero values on a zero canvas.

The interpretation for the band of zeros surrounding the X (and for zeros dispersed the non-zero area) is that pixel values at those locations from 1 to 784 never vary, in which case it is likely that one of the pixels is always zero. The reason the image looks pixelated is because we had to convert a 28 x 28 image to 1 x 784: the resulting row vector of pixels will have slices of non-zero values representing the original digits non-zero pixels (the rest of the vectors will be zeros). Thus, taking a covariance of these row vectors will result in pixelation patches within the overall covariance matrix. Finally, the non-zero indices represent the density/most likely areas for the image to have non-zero pixels.

- (d) (i) Since the multivariate Gaussians share the same covariance matrix  $\Sigma_{avg}$ , the decision boundary is linear since the covariance matrices will cancel out. With the 100 data point training set, the algorithm achieved a 51.1% error rate. With the 10000 data point training set, the algorithm achieved a 11.72% error rate. See appendix for full graph.

- (ii) In this case, the multivariate Gaussians do not share the same covariance matrix  $\Sigma_i$ , thus the decision boundary is quadratic since the covariance matrices will not cancel out. Essentially, the cross terms will remain (with maximum degree of two), and this will curve the decision boundary.

With the 100 data point training set, the algorithm achieved a 21.2% error rate. With the 10000 data point training set, the algorithm achieved a 7.67% error rate. See appendix for full graph.

- (iii) Part b generates a significantly superior error rate. This is because it does not make the assumption that all classes share the same covariance matrix, and allows the decision boundary to properly mimic the optimal class choice (instead of approximating using a linear boundary). This comes at the cost of calculation efficiency, as now the computer has to calculate cross terms.



## Problem 6: Linear Regression

1. Implement a linear regression model with least squares. Include your code in the submission.

You should add a constant term to the training data (e.g. add another dimension to each data point, with the value of 1). This is same as adding the bias term to linear regression (see discussion 4 question 1). Note that each data point  $\mathbf{x}^{(i)T}$  is a row of the training data matrix  $X$ .

**Solution:** In the code there should be some equation close to the form:

$$X^T X \mathbf{w} = X^T \mathbf{y} \text{ or } \mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$$

Although it would be weird if you did this, other valid solutions include linear programming and gradient descent.

2. Test your trained model on the validation set. What is the residual sum of squares (RSS) on the validation set? What is the range of predicted values (min, max)? Do they make sense?

**Solution:**

$$RSS = 5.7950e12.$$

The range of predicted values is  $-5.6563e04$  to  $7.1080e05$ . They don't really make sense because negative values are predicted.

3. Plot the regression coefficients  $\mathbf{w}$  (plot the value of each coefficient against the index of the coefficient). Be sure to exclude the coefficient corresponding to the constant offset you added earlier.

**Solution:** Plot should show that the 1st, 7th, and 8th coefficients are the most significant. There should not be a 9th coefficient that dominates all (this would be the bias coefficient).

4. Plot a histogram of the residuals of the training data (the residual corresponding to point  $i$  is  $f(\mathbf{x}^{(i)}) - y^{(i)}$ ). What distribution does this resemble?

**Solution:** Normal distribution.

## Appendix

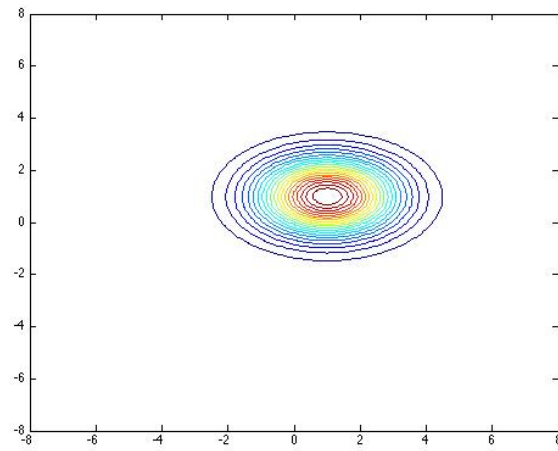


Figure 1: Problem 2a

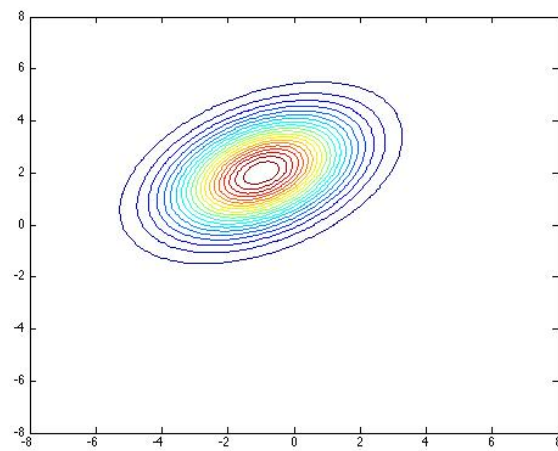


Figure 2: Problem 2b

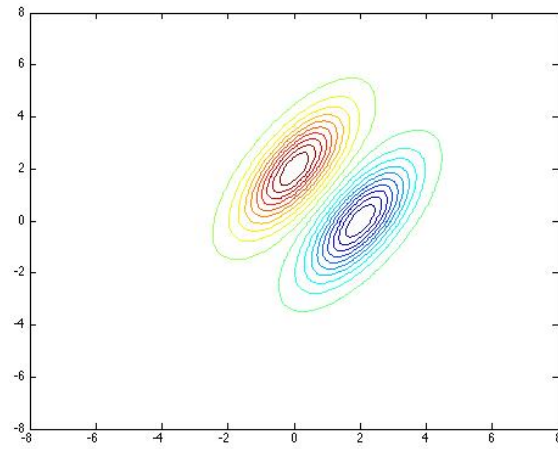


Figure 3: Problem 2c

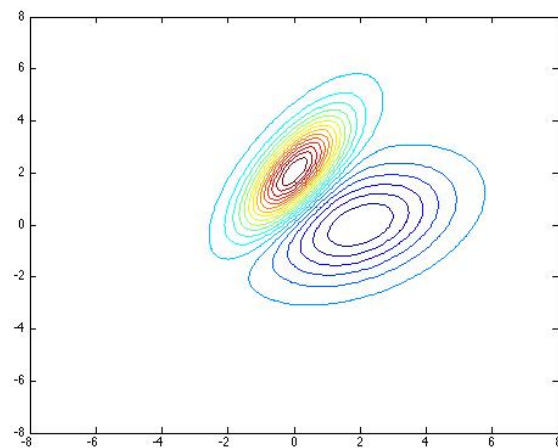


Figure 4: Problem 2d

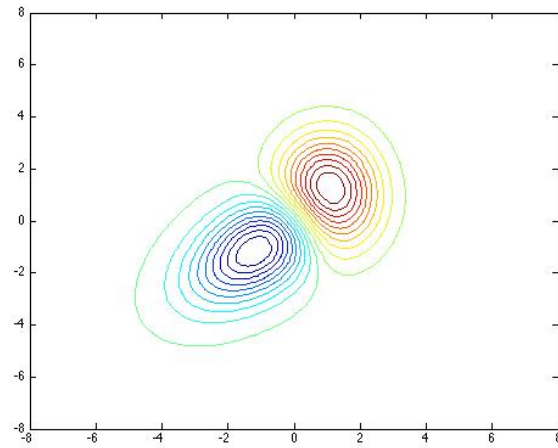


Figure 5: Problem 2e

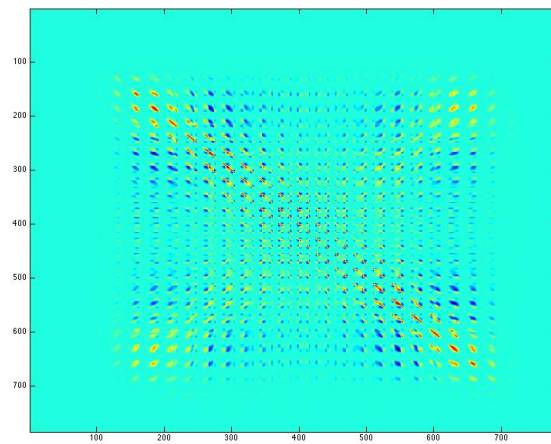


Figure 6: Problem 3c

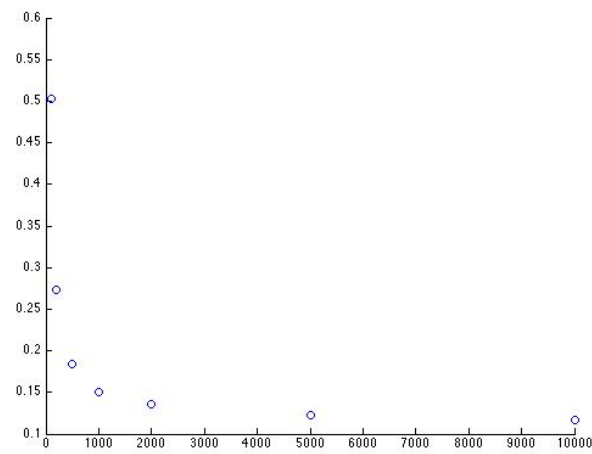


Figure 7: Problem 3d-i: Using an average covariance matrix across classes per data set

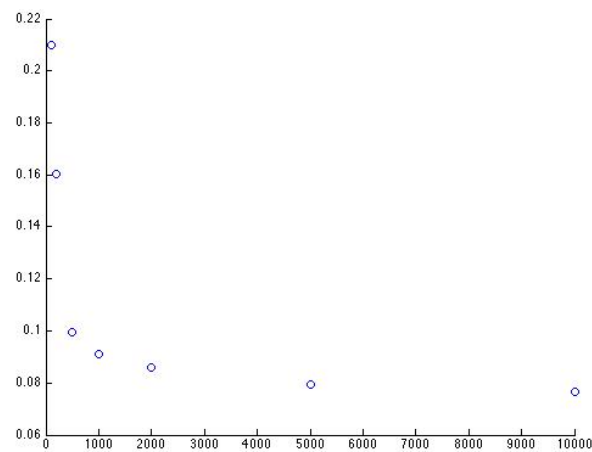


Figure 8: Problem 3d-ii: Using specific covariance matrices