

This homework is due **Monday, October 30 at 10pm.**

1 Getting Started

You may typeset your homework in latex or submit neatly handwritten and scanned solutions. Please make sure to start each question on a new page, as grading (with Gradescope) is much easier that way! Deliverables:

1. Submit a PDF of your writeup to assignment on Gradescope, “HW[n] Write-Up”
2. Submit all code needed to reproduce your results, “HW[n] Code”.
3. Submit your test set evaluation results, “HW[n] Test Set”.

After you’ve submitted your homework, be sure to watch out for the self-grade form.

- (a) Before you start your homework, write down your team. Who else did you work with on this homework? List names and email addresses. In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

- (b) Please copy the following statement and sign next to it:

I certify that all solutions are entirely in my words and that I have not looked at another student’s solutions. I have credited all external sources in this write up.

2 Classification Policy

Suppose we have a classification problem with classes labeled $1, \dots, c$ and an additional “doubt” category labeled $c + 1$. Let $f : \mathbb{R}^d \rightarrow \{1, \dots, c + 1\}$ be a decision rule. Define the loss function

$$L(f(x), y) = \begin{cases} 0 & \text{if } f(x) = y \quad f(x) \in \{1, \dots, c\}, \\ \lambda_s & \text{if } f(x) \neq y \quad f(x) \in \{1, \dots, c\}, \\ \lambda_r & \text{if } f(x) = c + 1 \end{cases} \quad (1)$$

where $\lambda_s \geq 0$ is the loss incurred for making a misclassification and $\lambda_r \geq 0$ is the loss incurred for choosing doubt. Hence, the risk of classifying a new data point x as class $f(x) \in \{1, 2, \dots, c + 1\}$ is

$$R(f(x)|x) = \sum_{i=1}^c L(f(x), i) P(Y = i|x).$$

(a) **Show that the following policy obtains the minimum risk:**

- (a) Choose class i if $P(Y = i|x) \geq P(Y = j|x)$ for all j and $P(Y = i|x) \geq 1 - \frac{\lambda_r}{\lambda_s}$.
- (b) Choose doubt otherwise.

Solution: Let $r : \mathbb{R}^d \rightarrow \{1, \dots, c + 1\}$ be the decision rule which implements the described policy. We will prove that in expectation the rule r is at least as good as the arbitrary rule f . Let $x \in \mathbb{R}^d$ be a data point, which we want to classify.

- Assume that property (1) holds and so $r(x) = i$, $P(Y = i|x) \geq P(Y = j|x)$ for all j , and $P(Y = i|x) \geq 1 - \lambda_r/\lambda_s$. Then:

$$R(r(x) = i|x) = \sum_{j=1}^c \ell(r(x) = i, y = j) P(Y = j|x) = \lambda_s \sum_{j=1, j \neq i}^c P(Y = j|x) = \lambda_s (1 - P(Y = i|x))$$

We are going to consider two cases for $f(x)$:

- Let $f(x) = k$, with $k \neq i, c + 1$. Then:

$$R(f(x) = k|x) = \lambda_s (1 - P(Y = k|x))$$

Hence $R(r(x) = i|x) \leq R(f(x) = k|x)$, since $P(Y = i|x) \leq P(Y = k|x)$.

- Let $f(x) = c + 1$. Then:

$$R(f(x) = k|x) = \lambda_r$$

Hence $R(r(x) = i|x) \leq R(f(x) = k|x)$, since $P(Y = i|x) \geq 1 - \lambda_r/\lambda_s$.

- Assume that property (2) holds and so $r(x) = c + 1$. Then:

$$R(r(x) = i|x) = \lambda_r$$

Let $f(x) = k$, with $k \neq c + 1$. Then:

$$R(f(x) = k|x) = \lambda_s (1 - P(Y = k|x))$$

Property (1) does **not** hold which means that:

$$\max_{j \in \{1, \dots, c\}} P(Y = j|x) < 1 - \lambda_r / \lambda_s$$

hence $P(Y = k|x) < 1 - \lambda_r / \lambda_s$, and so $R(r(x) = i|x) < R(f(x) = k|x)$.

Therefore in every case we proved that the rule r is at list as good as the arbitrary rule f , which proves that r is an optimal rule.

- (b) **What happens if $\lambda_r = 0$? What happens if $\lambda_r > \lambda_s$? Explain why this is consistent with what one would expect intuitively.**

Solution: If $\lambda_r = 0$, then property (1) will hold iff there exists an $i \in \{1, \dots, c\}$ such that $P(r(x) = i|x) = 1$. So we will either classify x in class i if we are 100% sure about this, or else we will choose doubt. Of course this is completely consistent with our intuition, because choosing doubt does not have any penalty at all, since $\lambda_r = 0$.

If $\lambda_r > \lambda_s$, then we will always classify x in the class $i \in \{1, \dots, c\}$ which gives the highest probability of correct classification. Once again this makes sense, since the cost of choosing doubt is higher than classifying x in any of the classes, hence our best option is to classify x in the class which gives the highest probability for a correct classification.

3 LDA and CCA

Consider the following random variable $X \in \mathbb{R}^d$, generated using a *mixture of two Gaussians*. Here, the vectors $\mu_1, \mu_2 \in \mathbb{R}^d$ are arbitrary (mean) vectors, and $\Sigma \in \mathbb{R}^{d \times d}$ represents a positive definite (covariance) matrix. For now, we will assume that we know all of these parameters.

Draw a label $L \in \{1, 2\}$ such that the label 1 is chosen with probability π_1 (and consequently, label 2 with probability $\pi_2 = 1 - \pi_1$), and generate $X = \mathcal{N}(\mu_L, \Sigma)$.

- (a) Now given a particular $X \in \mathbb{R}^d$ generated from the above model, we wish to find its label. **Write out the decision rule corresponding to the following estimates of L :**

- MLE
- MAP

Your decision rule should take the form of a threshold: if some function $f(X) > T$, then choose the label 1, otherwise choose the label 2. **When are these two decision rules the same?**

Solution: Let us use the shorthand $p(X = x|L = 1) = p_1(x)$, and define $p_2(x)$ analogously. In computing the MLE, we need to find which of $p_1(x)$ and $p_2(x)$ is larger. Taking a ratio of the two, we have

$$\frac{p_1(x)}{p_2(x)} = \exp\left\{-\frac{1}{2}(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1) + \frac{1}{2}(x - \mu_2)^\top \Sigma^{-1}(x - \mu_2)\right\}. \quad (2)$$

Let us now use the following simplifying fact. We have

$$\begin{aligned} v^\top Mv - u^\top Mu &= v^\top Mv - v^\top Mu - (u - v)^\top Mu \\ &= v^\top M(v - u) + (v - u)^\top Mu \\ &= (v + u)^\top M(v - u), \end{aligned}$$

where we have used the fact that M is symmetric.

Substituting $v = x - \mu_2$, $u = x - \mu_1$, we have

$$\log \frac{p_1(x)}{p_2(x)} = \left(x - \frac{\mu_1 + \mu_2}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_2).$$

Notice that we are interested in comparing this quantity with 0; if it is greater than zero, then we prefer 1 to 2, and we prefer 2 otherwise. The MLE decision rule therefore takes the form

$$\hat{L}_{\text{MLE}} = \begin{cases} 1 & \text{if } \left(X - \frac{\mu_1 + \mu_2}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_2) > 0 \\ 2 & \text{otherwise.} \end{cases}$$

For the MAP estimator, we are interested in comparing the quantities $p(L = 1|X = x)$ and $p(L = 2|X = x)$. By an application of Bayes rule, we are comparing $\frac{1}{p(X=x)}p(X = x|L = 1)p(L = 1)$ with $\frac{1}{p(X=x)}p(X = x|L = 2)p(L = 2)$, and so it suffices to compute the ratio

$$\frac{p(X = x|L = 1)p(L = 1)}{p(X = x|L = 2)p(L = 2)} = \frac{\pi_1 p_1(x)}{\pi_2 p_2(x)}.$$

Since we are going to compute the logarithm of this expression and set it to zero, it suffices to check the condition

$$\log \frac{p_1(x)}{p_2(x)} > \log \frac{\pi_2}{\pi_1}.$$

The thresholding is therefore identical to the calculation above up to an offset, and we have

$$\hat{L}_{\text{MAP}} = \begin{cases} 1 & \text{if } \left(X - \frac{\mu_1 + \mu_2}{2} \right)^\top \Sigma^{-1} (\mu_1 - \mu_2) > \log \frac{\pi_2}{\pi_1} \\ 2 & \text{otherwise.} \end{cases}$$

Clearly, when $\pi_1 = \pi_2 = 1/2$, these two estimators are exactly the same.

- (b) You should have noticed that the function f is *linear* in its argument X , and takes the form $w^\top (X - v)$. We will now show that CCA defined on a suitable set of random variables leads to precisely the same decision rule.

Let $Y \in \mathbb{R}^2$ be a one hot vector denoting the realization of the label ℓ , i.e. $Y_\ell = 1$ if $L = \ell$, and zero otherwise. Let $\pi_1 = \pi_2 = 1/2$. **Compute the covariance matrices Σ_{XX} , Σ_{XY} and Σ_{YY} as a function of μ_1, μ_2, Σ .** Recall that the random variables are not zero-mean.

Solution: Let us first compute the covariance matrix of X . Recall that $X \sim N(\mu_1, \Sigma)$ when $L = 1$ (which happens with probability $1/2$), and $X \sim N(\mu_2, \Sigma)$ otherwise.

Also note that the covariance matrix of a non-zero mean RV can be written as $\mathbb{E}[(Z - \mu)(Z - \mu)^\top] = \mathbb{E}[ZZ^\top] - 2\mu\mathbb{E}[Z^\top] + \mu\mu^\top$.

Let us now use $\mu = \frac{\mu_1 + \mu_2}{2}$ to denote the mean of X . Hence, we have

$$\Sigma_{XX} = \mathbb{E}[XX^\top] - \mu\mu^\top.$$

We also have

$$\begin{aligned} \mathbb{E}[XX^\top] &= \frac{1}{2}\mathbb{E}[XX^\top | L = 1] + \frac{1}{2}\mathbb{E}[XX^\top | L = 2] \\ &= \frac{1}{2}(\Sigma + \mu_1\mu_1^\top) + \frac{1}{2}(\Sigma + \mu_2\mu_2^\top) \\ &= \Sigma + \frac{1}{2}(\mu_1\mu_1^\top + \mu_2\mu_2^\top). \end{aligned}$$

Combining this with our calculation above, we have

$$\begin{aligned} \Sigma_{XX} &= \Sigma + \frac{1}{2}(\mu_1\mu_1^\top + \mu_2\mu_2^\top) - \mu\mu^\top \\ &= \Sigma + \frac{1}{4}(\mu_1 - \mu_2)(\mu_1 - \mu_2)^\top. \end{aligned}$$

(Go ahead and verify the last line.)

Let us now compute Σ_{YY} . Using the same idea, and noting that $\mathbb{E}[Y] = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$, we have

$$\Sigma_{YY} = \mathbb{E}[YY^\top] - \begin{bmatrix} 1/4 & 1/4 \\ 1/4 & 1/4 \end{bmatrix}.$$

Also note that $\mathbb{E}[Y_i Y_j] = 0$ if $i \neq j$ (since one of the random variables will be zero by the one-hot property), and $\mathbb{E}[Y_1^2] = \pi_1 = 1/2$ and $\mathbb{E}[Y_2^2] = \pi_2 = 1/2$.

Thus, the covariance matrix is given by

$$\begin{aligned} \Sigma_{YY} &= \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix} - \begin{bmatrix} 1/4 & 1/4 \\ 1/4 & 1/4 \end{bmatrix} \\ &= \begin{bmatrix} 1/4 & -1/4 \\ -1/4 & 1/4 \end{bmatrix} \\ &= \frac{1}{4} \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}^\top \end{aligned}$$

To conclude, we must compute the covariance matrix Σ_{XY} . Note that we can again condition

on the labels to accomplish this, and using $\mathbf{1}/2$ to denote the all half vector, we have

$$\begin{aligned}
 \Sigma_{XY} &= \frac{1}{2}\mathbb{E}[(X - \mu)(Y - \mathbf{1}/2)^\top | L = 1] + \frac{1}{2}\mathbb{E}[(X - \mu)(Y - \mathbf{1}/2)^\top | L = 2] \\
 &= \frac{1}{2}\mathbb{E}[(X - \mu) \begin{bmatrix} 1/2 & -1/2 \end{bmatrix} | L = 1] + \frac{1}{2}\mathbb{E}[(X - \mu) \begin{bmatrix} -1/2 & 1/2 \end{bmatrix} | L = 2] \\
 &= \frac{1}{4}\mathbb{E}[(X - \mu) \quad -(X - \mu)] | L = 1 + \frac{1}{4}\mathbb{E}[\quad -(X - \mu) \quad (X - \mu)] | L = 2 \\
 &= \frac{1}{4} \begin{bmatrix} \mu_1 - \mu_2 & \mu_2 - \mu_1 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} \mu_1 - \mu_2 & \mu_2 - \mu_1 \end{bmatrix} \\
 &= \frac{1}{4} \begin{bmatrix} \mu_1 - \mu_2 & \mu_2 - \mu_1 \end{bmatrix} \\
 &= \frac{1}{4}(\mu_1 - \mu_2) \begin{bmatrix} 1 & -1 \end{bmatrix}^\top.
 \end{aligned}$$

- (c) Let us now perform CCA on the two random variables. Recall that in order to find the first canonical directions, we look for vectors $u \in \mathbb{R}^d$ and $v \in \mathbb{R}^2$ such that $\rho(u^\top X, v^\top Y)$ is maximized.

Show that the maximizing u^* is proportional to $\Sigma^{-1}(\mu_1 - \mu_2)$. Recall that u^* is that “direction” of X that contributes most to predicting Y . **What is the relationship between u^* and the function $f(X)$ computed in part (a)?**

Hint: The Sherman-Morrison formula for matrix inversion may be useful:

Suppose $A \in \mathbb{R}^{d \times d}$ is an invertible square matrix and $a, b \in \mathbb{R}^d$ are column vectors. Then,

$$(A + ab^\top)^{-1} = A^{-1} - \frac{A^{-1}ab^\top A^{-1}}{1 + b^\top A^{-1}a}.$$

Solution: Recall the definition of the correlation coefficient. We are interested in vectors u, v such that we maximize

$$\begin{aligned}
 \rho(u^\top X, v^\top Y) &= \frac{\text{cov}(u^\top X, v^\top Y)}{\sqrt{\text{var}(u^\top X) \text{var}(v^\top Y)}} \\
 &= \frac{u^\top \Sigma_{XY} v}{\sqrt{u^\top \Sigma_{XX} u} \sqrt{v^\top \Sigma_{YY} v}}.
 \end{aligned}$$

Now using the expressions we computed above, we have

$$\begin{aligned}
 u^\top \Sigma_{XY} v &= \frac{1}{4} u^\top (\mu_1 - \mu_2) \begin{bmatrix} 1 & -1 \end{bmatrix}^\top v \\
 v^\top \Sigma_{YY} v &= \frac{1}{4} \left(\begin{bmatrix} 1 & -1 \end{bmatrix}^\top v \right)^2.
 \end{aligned}$$

Consequently, we see that the factor $\begin{bmatrix} 1 & -1 \end{bmatrix}^\top v$ cancels in both the numerator and denominator, yielding

$$\rho(u^\top X, v^\top Y) = \frac{1}{2} \frac{u^\top (\mu_1 - \mu_2)}{\sqrt{u^\top \Sigma_{XX} u}},$$

and we are interested in maximizing this expression over all u . Performing the change of variables $w = \Sigma_{XX}^{-1/2}u$, we are interested in the expression

$$\max_w \frac{w^\top \Sigma_{XX}^{-1/2}(\mu_1 - \mu_2)}{\|w\|_2}.$$

Since the above expression does not change when w is multiplied by a constant, we may assume that it has unit ℓ_2 norm, and so the optimization problem becomes

$$\max_{w: \|w\|_2=1} w^\top \Sigma_{XX}^{-1/2}(\mu_1 - \mu_2).$$

To find the maximizer, we use Cauchy Schwarz inequality to obtain

$$w^\top \Sigma_{XX}^{-1/2}(\mu_1 - \mu_2) \leq \|w\|_2 \|\Sigma_{XX}^{-1/2}(\mu_1 - \mu_2)\|_2,$$

and equality is attained by the vector $w^* = \alpha \Sigma_{XX}^{-1/2}(\mu_1 - \mu_2)$ for a scalar α that normalizes w^* to be unit norm.

We can then obtain $u^* = \Sigma_{XX}^{-1/2}w^* = \alpha \Sigma_{XX}^{-1}(\mu_1 - \mu_2)$.

We are now in a position to apply the hint. We know that

$$\Sigma_{XX}^{-1} = \Sigma^{-1} - \frac{\Sigma^{-1}(\mu_1 - \mu_2)(\mu_1 - \mu_2)^\top \Sigma^{-1}}{4 + (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2)};$$

notice now that $(\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2)$ is a scalar β . We also have $u^* \propto \Sigma_{XX}^{-1}(\mu_1 - \mu_2) = \Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{4+\beta} \Sigma^{-1}(\mu_1 - \mu_2)\beta$.

In other words, we have $u^* = \gamma \Sigma^{-1}(\mu_1 - \mu_2)$. Note that this is precisely the same direction w predicted by LDA in part (a)! What does this mean?

(d) **Repeat the above part for arbitrary probabilities π_1, π_2 .**

Solution: We now show that essentially the same relations hold for the covariance matrices in this case by repeating part (b); part (c) then follows unchanged for the more general case.

Now denote $\mu = \pi_1\mu_1 + \pi_2\mu_2$. We have

$$\Sigma_{XX} = \mathbb{E}[XX^\top] - \mu\mu^\top.$$

We also have

$$\begin{aligned} \mathbb{E}[XX^\top] &= \pi_1 \mathbb{E}[XX^\top | L=1] + \pi_2 \mathbb{E}[XX^\top | L=2] \\ &= \pi_1 (\Sigma + \mu_1\mu_1^\top) + \pi_2 (\Sigma + \mu_2\mu_2^\top) \\ &= \Sigma + \pi_1\mu_1\mu_1^\top + \pi_2\mu_2\mu_2^\top. \end{aligned}$$

Combining this with our calculation above, we have

$$\begin{aligned} \Sigma_{XX} &= \Sigma + \pi_1\mu_1\mu_1^\top + (1 - \pi_1)\mu_2\mu_2^\top - \mu\mu^\top \\ &= \Sigma + \pi_1\pi_2(\mu_1 - \mu_2)(\mu_1 - \mu_2)^\top. \end{aligned}$$

Let us now compute Σ_{YY} . Using the same idea, and noting that $\mathbb{E}[Y] = \begin{bmatrix} \pi_1 \\ \pi_2 \end{bmatrix}$, we have

$$\Sigma_{YY} = \mathbb{E}[YY^\top] - \begin{bmatrix} \pi_1^2 & \pi_1\pi_2 \\ \pi_1\pi_2 & \pi_2^2 \end{bmatrix}.$$

Also note that $\mathbb{E}[Y_i Y_j] = 0$ if $i \neq j$ (since one of the random variables will be zero by the one-hot property), and $\mathbb{E}[Y_1^2] = \pi_1$ and $\mathbb{E}[Y_2^2] = \pi_2$.

Thus, the covariance matrix is given by

$$\begin{aligned} \Sigma_{YY} &= \begin{bmatrix} \pi_1 & 0 \\ 0 & \pi_2 \end{bmatrix} - \begin{bmatrix} \pi_1^2 & \pi_1\pi_2 \\ \pi_1\pi_2 & \pi_2^2 \end{bmatrix} \\ &= \begin{bmatrix} \pi_1(1-\pi_1) & -\pi_1\pi_2 \\ -\pi_1\pi_2 & \pi_2(1-\pi_2) \end{bmatrix} \\ &= \pi_1\pi_2 \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}^\top \end{aligned}$$

where the last step uses the fact that $\pi_2 = 1 - \pi_1$.

To conclude, we must compute the covariance matrix Σ_{XY} . Note that we can again condition on the labels to accomplish this, and using π to denote the vector (π_1, π_2) , we have (Note that we again use $\pi_1 = 1 - \pi_2$ below)

$$\begin{aligned} \Sigma_{XY} &= \pi_1 \mathbb{E}[(X - \mu)(Y - \pi)^\top | L = 1] + \pi_2 \mathbb{E}[(X - \mu)(Y - \pi)^\top | L = 2] \\ &= \pi_1 \mathbb{E}[(X - \mu) \begin{bmatrix} \pi_2 & -\pi_2 \end{bmatrix} | L = 1] + \pi_2 \mathbb{E}[(X - \mu) \begin{bmatrix} -\pi_1 & \pi_1 \end{bmatrix} | L = 2] \\ &= \pi_1 \pi_2 \mathbb{E}[(X - \mu) \quad -(X - \mu) | L = 1] + \pi_1 \pi_2 \mathbb{E}[\begin{bmatrix} -(X - \mu) & (X - \mu) \end{bmatrix} | L = 2] \\ &= \pi_1 \pi_2 \begin{bmatrix} \frac{\mu_1 - \mu_2}{2} & \frac{\mu_2 - \mu_1}{2} \end{bmatrix} + \pi_1 \pi_2 \begin{bmatrix} \frac{\mu_1 - \mu_2}{2} & \frac{\mu_2 - \mu_1}{2} \end{bmatrix} \\ &= \pi_1 \pi_2 \begin{bmatrix} \mu_1 - \mu_2 & \mu_2 - \mu_1 \end{bmatrix} \\ &= \pi_1 \pi_2 (\mu_1 - \mu_2) \begin{bmatrix} 1 & -1 \end{bmatrix}^\top. \end{aligned}$$

We therefore see that all the calculation that we did in part (c) carries over without modification, since the covariance matrices have precisely the same factorizations as above.

- (e) Now assume you don't know the parameters μ_1, μ_2, Σ , but are given samples of training data $(x_1, \ell_1), (x_2, \ell_2), \dots, (x_n, \ell_n)$ drawn from the above distribution. You are then given a new X_{test} , which you wish to classify. Using your intuition from CCA, **write down a procedure that predicts the label of X_{test} .**

Solution: We have two options, as a consequence of the work we did above. One option is to perform CCA regression with the random variables X and Y . In particular, each training data point gives us a realization of this pair of random variables, and we have samples x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_n , so this is a straightforward application of the techniques we already know.

The second option is to perform LDA from samples. Let us sketch out this solution here.

We have to compute the covariance matrices and means from data. As we have seen, it suffices to compute just the covariance matrix of X and the means, which we can write from samples $x_1, x_2, \dots, x_n; y_1, y_2, \dots, y_n$. Note that knowing y_i implies that we know the label ℓ_i .

$$\hat{\mu}_1 = \frac{1}{\#\{i : \ell_i = 1\}} \sum_{i:\ell_i=1} x_i$$

$$\hat{\mu}_2 = \frac{1}{\#\{i : \ell_i = 2\}} \sum_{i:\ell_i=2} x_i.$$

We also know that the covariance matrices of both classes are the same, so we can compute the covariance matrix estimates

$$\hat{\Sigma}_1 = \frac{1}{\#\{i : \ell_i = 1\}} \sum_{i:\ell_i=1} (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^\top$$

$$\hat{\Sigma}_2 = \frac{1}{\#\{i : \ell_i = 2\}} \sum_{i:\ell_i=2} (x_i - \hat{\mu}_2)(x_i - \hat{\mu}_2)^\top.$$

In addition, we also have a notion of the label probabilities π_1 and π_2 . In particular, we have

$$\hat{\pi}_1 = \frac{1}{n} \#\{i : \ell_i = 1\},$$

$$\hat{\pi}_2 = \frac{1}{n} \#\{i : \ell_i = 2\}.$$

The final covariance can finally be found by computing

$$\hat{\Sigma} = \hat{\pi}_1 \hat{\Sigma}_1 + \hat{\pi}_2 \hat{\Sigma}_2.$$

We then use the statistic $\hat{v} = \hat{\Sigma}^{-1}(\hat{\mu}_1 - \hat{\mu}_2)$, and when a new point X_{test} comes in, we compute $\hat{v}^\top X_{\text{test}}$ and compare it with a threshold $\log \frac{\hat{\pi}_1}{\hat{\pi}_2}$ (or 0) depending on whether we want to compute an MLE or MAP estimate.

4 Sensors, Objects, and Localization (Part 2)

Let us say there are n objects and m sensors located in a $2d$ plane. The n objects are located at the points $(x_1, y_1), \dots, (x_n, y_n)$. The m sensors are located at the points $(a_1, b_1), \dots, (a_m, b_m)$. We have measurements for the distances between the objects and the sensors: D_{ij} is the measured distance from object i to sensor j . The distance measurement has noise in it. Specifically, we model

$$D_{ij} = \|(x_i, y_i) - (a_j, b_j)\| + Z_{ij},$$

where $Z_{ij} \sim N(0, 1)$. The noise is independent across different measurements.

Starter code has been provided for data generation to aid your explorations. All Python libraries are permitted, though you must use your own implementation of a Neural Network if you want to

get bonus points for implementing stochastic gradient descent. For others parts, you are free to use commercial libraries. You will need to plot ten figures in total in this problem.

Assume we observe $D_{ij} = d_{ij}$ with $(X_i, Y_i) = (x_i, y_i)$ for $i = 1, \dots, n$ and $j = 1, \dots, m$. Here, $m = 7$. **Our goal is to predict $(X_{i'}, Y_{i'})$ from newly observed $D_{i'1}, \dots, D_{i'7}$.** For a data set with n test data, the error is measured by the average distance between the predicted object locations and the true object locations,

$$\frac{1}{n} \sum_{i=1}^n \sqrt{(\hat{x}_i - x_i)^2 + (\hat{y}_i - y_i)^2},$$

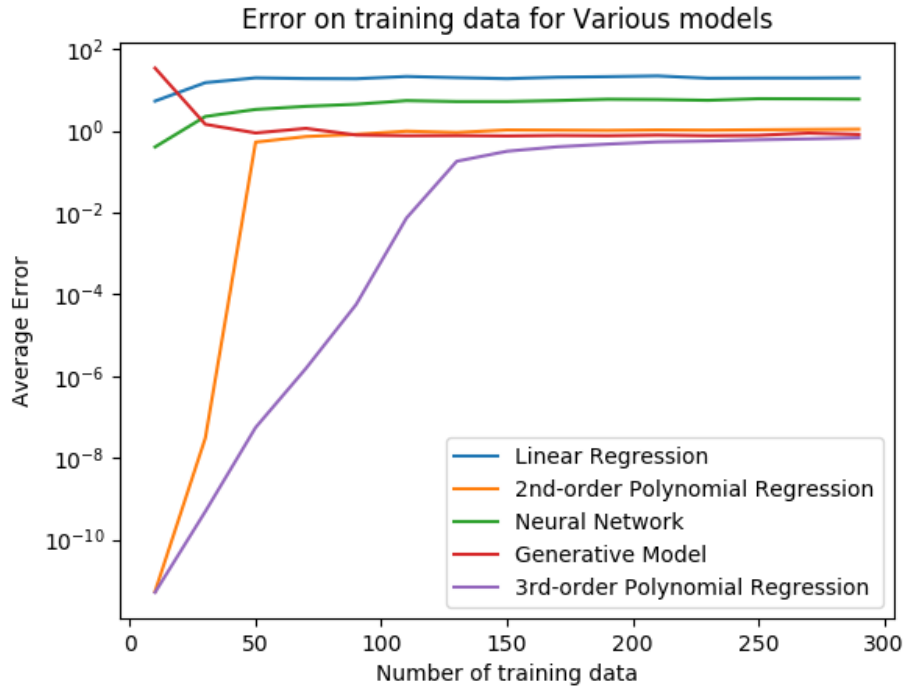
where (\hat{x}_i, \hat{y}_i) is the location of objects predicted by a model.

We are going to consider five models in this problem and compare their performance:

- In an earlier assignment, we first estimated sensor locations and then used the estimated sensor locations to estimate the new object locations. This is called a *Generative Model*.
 - A *Linear Model*. Using the training set, the linear model attempts to learn (X_i, Y_i) from (D_{i1}, \dots, D_{i7}) . Then it predicts $(X_{i'}, Y_{i'})$ from $(D_{i'1}, \dots, D_{i'7})$.
 - A *Second-Order Polynomial Regression Model*. The set-up is similar to the linear model, but including second-order polynomial features.
 - A *Third-Order Polynomial Regression Model*. The set-up is similar to the linear model, but including third-order polynomial features.
 - A *Neural Network Model*. The Neural Network should have two hidden layers, each with 100 neurons, and use `ReLU` as the non-linearity. (You are encouraged to explore on your own beyond this however. These parameters were chosen to teach you a hype-deflating lesson.)
- (a) **Implement the five models listed above. Submit your code on gradescope and include it in your write-up.**
- (b) **Fix a set of 7 sensors and generate the following data sets:**
- 15 training sets where n_{train} varies from 10 to 290 in increments of 20.
 - A “regular” testing data set where $n_{\text{test}} = 1000$.
 - A “shifted” testing data set where $n_{\text{test}} = 1000$. You can do this by setting `original_dist` to `False` in the function `generate_data` in `starter.py`.

The difference between the “regular” testing data and the “shifted” testing data is that the “regular” testing data is drawn from the same distribution as the training data, whereas the “shifted” testing data is farther away. **Train each of the five models on each of the fifteen training sets. Use your results to generate three figures. Each figure should include all of the models on the same plot so that you can compare them:**

- A plot of *training error* versus n_{train} (the amount of data used to train the model) for all of the models.



- A plot of *testing error* on the “regular” test set versus n_{train} (the amount of data used to train the model) for all of the models.
- A plot of *testing error* on the “shifted” test set versus n_{train} (the amount of data used to train the model) for all of the models.

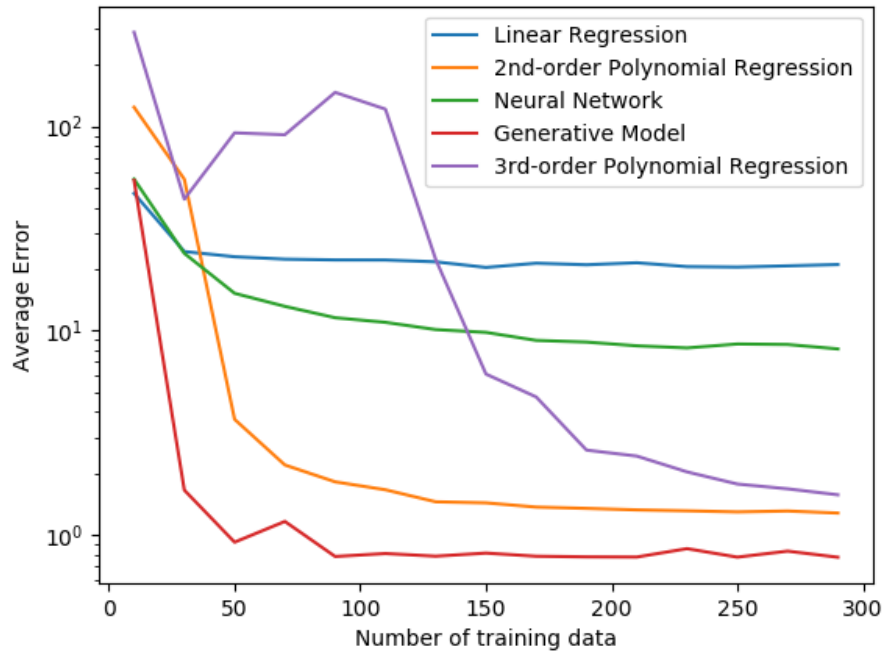
Briefly describe your observations. What are the strengths and weaknesses of each model?

Solution: Code is provided in separate files. `models.py` contains solutions to Part a,b,c,d and e. `plot1.py` contains solutions to Part f. `plot2.py` contains solutions to Part g. `plot3.py` contains solutions to Part h. `plot4.py` contains solutions to Part i. `plot5.py` contains solutions to Part j.

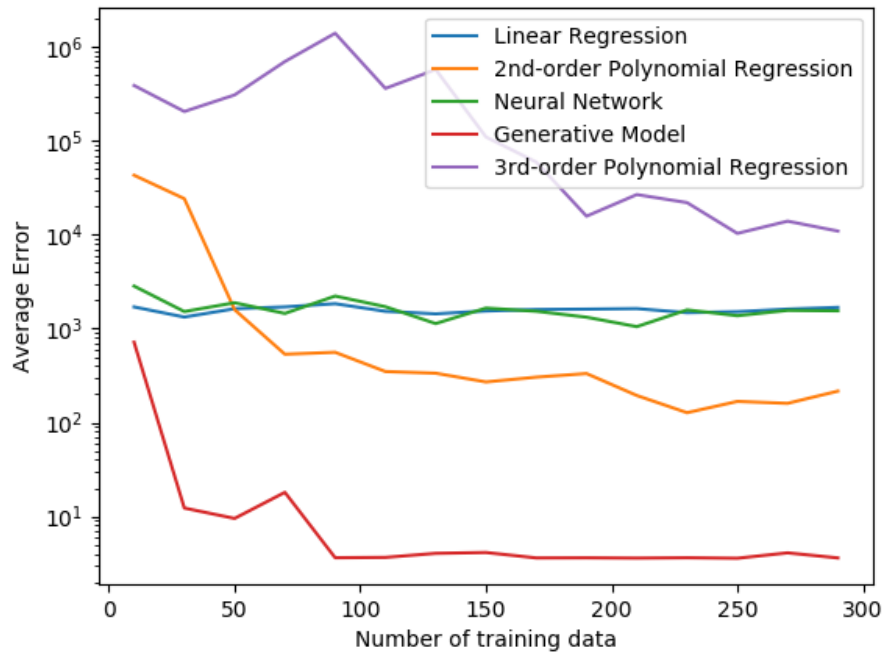
The plot by varying the numbers of data points n_{train} from 10 to 290 by 20 is shown. We observe that the generative model achieves the best performance as expected. Why? When data is drawn from a different distribution, we see that although all the models are performing worse, the generative model is still generalizing at some reasonable level. The rest of the models are behaving quite poorly indeed. This is because extrapolation is fundamentally challenging for any universal-function-approximation type model. Second-order polynomial has the second better performance. Third order polynomial performs well when the data is generated from the same distribution, followed by neural networks. The linear model fails to perform well even with ample training data.

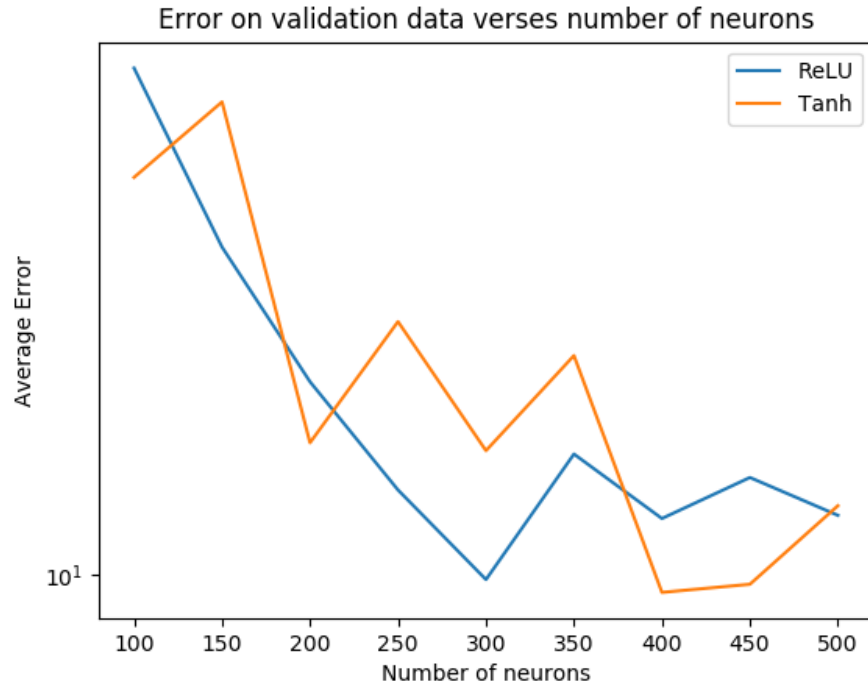
- (c) We are now going to do some hyper-parameter tuning on our neural network. Fix the number of hidden layers to be two and let ℓ be the number of neurons in each of these two layers. Try values for ℓ between 100 and 500 in increments of 50. Use data sets with $n_{\text{train}}=200$ and

Error on test data from the same distribution for Various models



Error on test data from a different distribution for Various models





$n_{\text{test}} = 1,000$. **What is the best choice for ℓ (the number of neurons in the hidden layers)? Justify your answer with plots.**

Solution: The best performance is achieved at the case when number of neurons equals to 300 – 500. (Any reasonable answer that is supported with a plot is acceptable.)

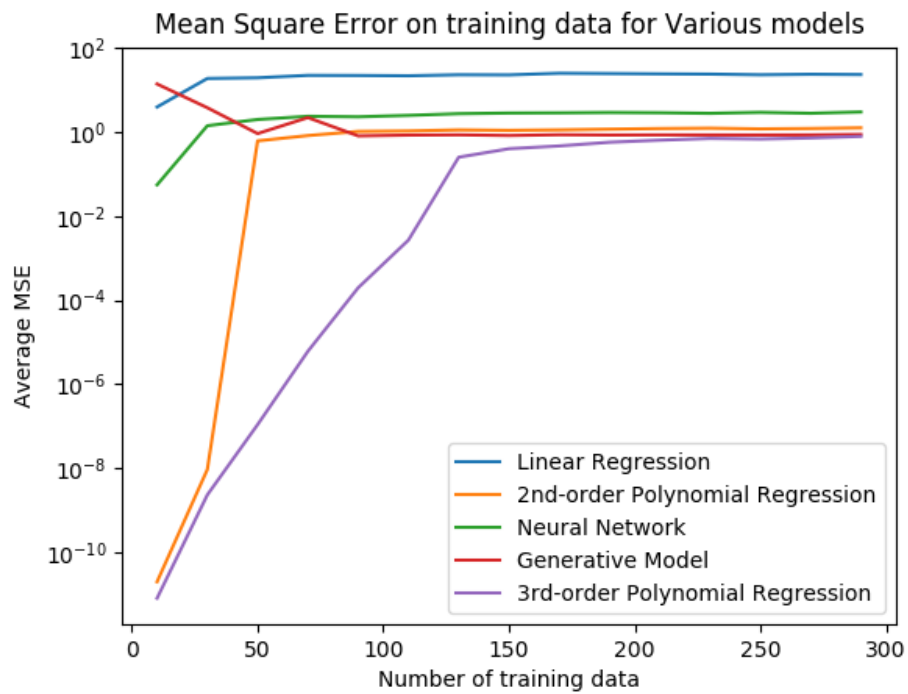
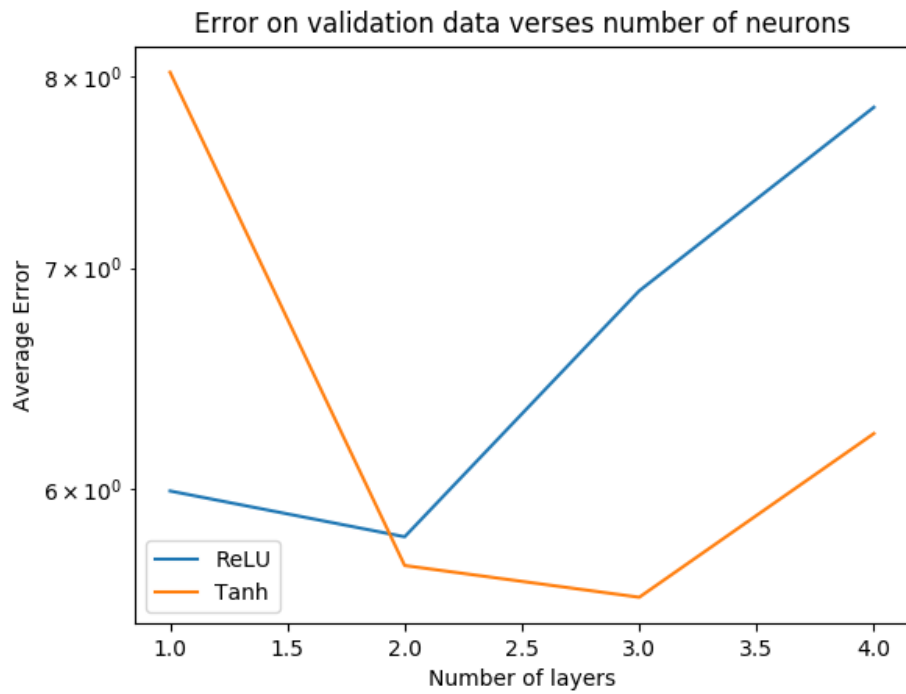
- (d) We are going to do some more hyper-parameter tuning on our neural network. Let k be the number of hidden layers and let ℓ be the number of neurons in each hidden layer. **Write a formula for the total number of weights in our network in terms of ℓ and k . If we want to keep the total number of weights in the network approximately equal to 10000, find a formula for ℓ in terms of k .** Try values of k between 1 and 4 with the appropriate implied choice for ℓ . Use data sets with $n_{\text{train}} = n_{\text{test}} = 200$. **What is the best choice for k (the number of layers)? Justify your answer with plots.**

Solution: The best performance is achieved at the case when number of layers equals to 3 with Tanh and 2 with ReLU. (Any reasonable answer that is supported with a plot is acceptable.)

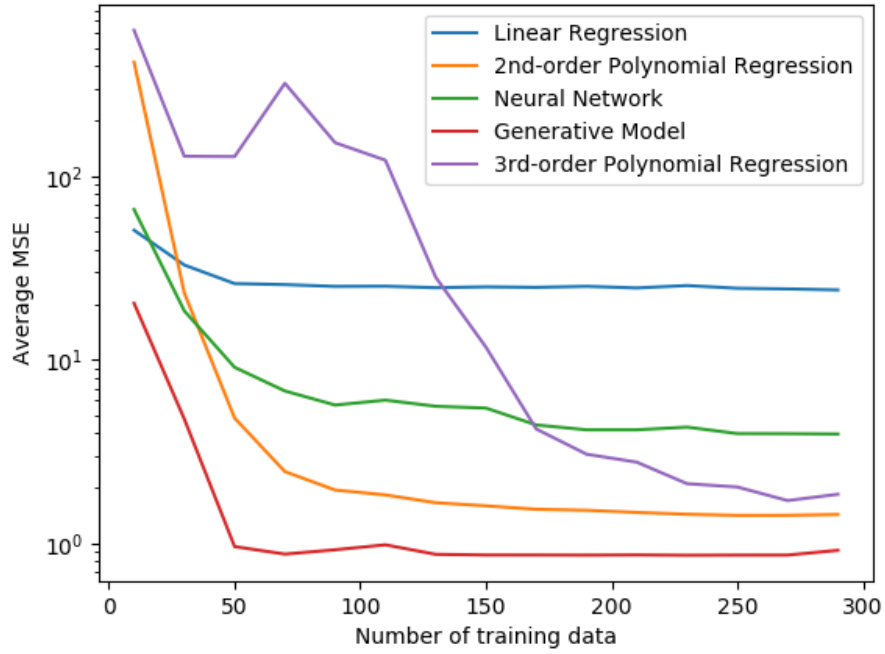
- (e) **Repeat part ((b)), but for the neural network you should use the hyper-parameters you found through hyper-parameter tuning.**

Solution: Performance is shown in Figures on Page 13.

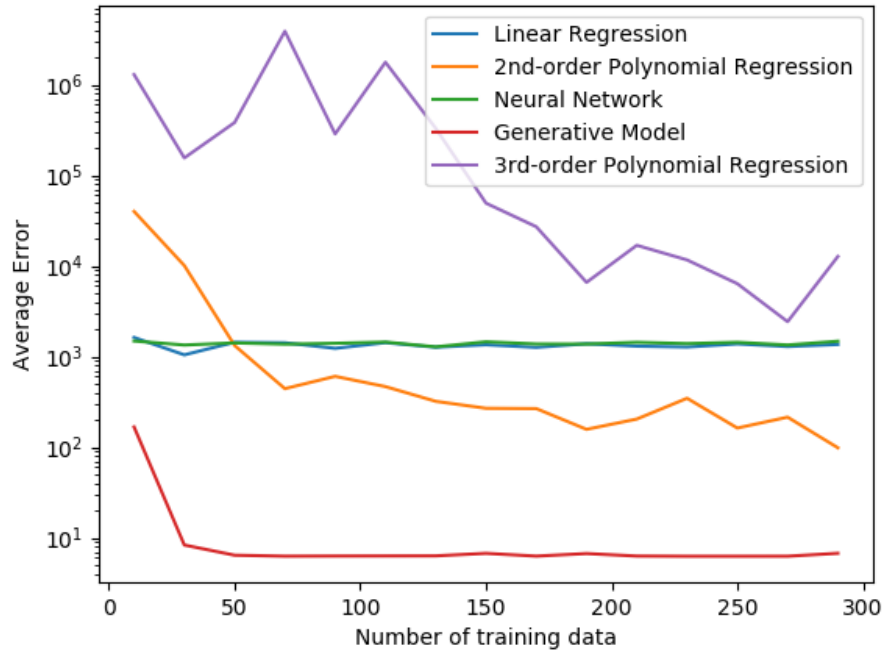
- (f) (Bonus) Instead of training your neural network model using gradient descent, use batch stochastic gradient descent. **How does the neural network trained with SGD compare to the neural network trained with gradient descent? Justify your answer with plots.** (For this part, you must use your own implementation.)



MSE on test data from the same distribution for Various models



Error on test data from a different distribution for Various models





Solution: The learning rate and batchsize is chosen based on the following formula and the results are shown in figures.

$$\text{Number of iterations} = \frac{\text{number of data}}{\text{Batchsize}} \cdot \text{Number of iterations of Gradient Descent} \quad (3)$$

and the learning rate to be

$$\text{Number of iterations} = \frac{\text{Batchsize}}{\text{number of data}} \cdot \text{Learning rate of Gradient Descent} \quad (4)$$

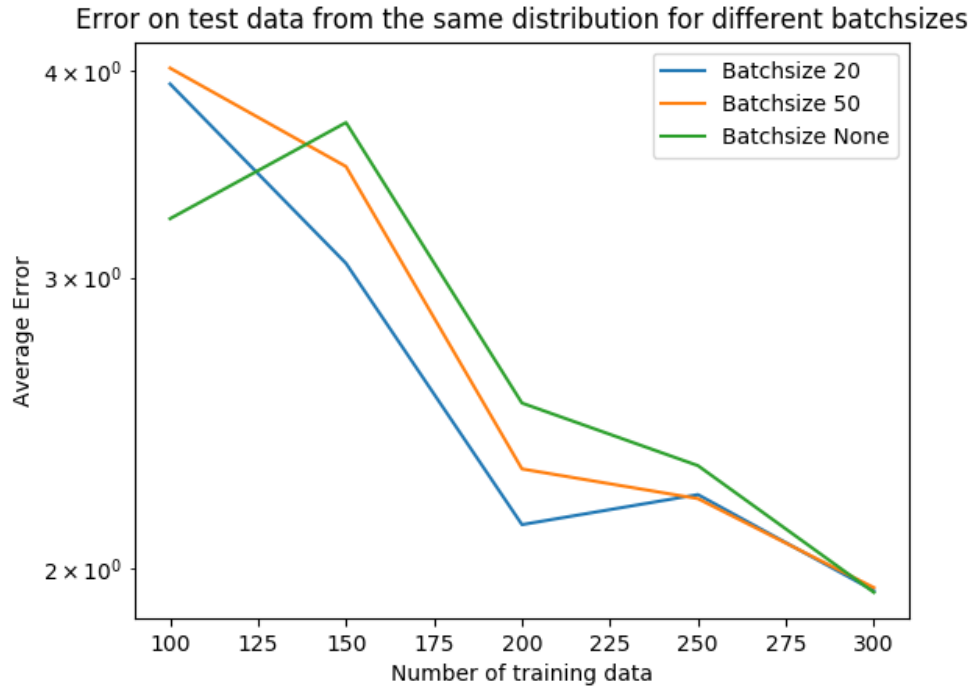
Code lies in plot5.py. SGD achieves similar performance to GD.

- (g) (Bonus) You might have seen that the neural network performance was dissapointing. Explore increasing the training data as well as playing with the hyper-parameters on your own to improve neural network performance for samples drawn from the “regular” data set. Can you get anything to generalize to the “shifted” test data? **Report your parameters and present your evidence.**

Solution: Neither increasing the training data nor playing with the hyper-parameters on your own can improve neural network performance for samples drawn from the “regular” data set. The reason is that the data are generated from a different distribution and neural network cannot transfer its performance to that, as mentioned above. (Performance of the tuned neural network on data generated from a different distribution can be seen in the figures above.)

5 Tensorflow Installation

Install Tensorflow and run the tutorial.



https://www.tensorflow.org/get_started/mnist/beginners

By now, hopefully you have a reasonably good understanding of how Backpropagation and SGD work to train Neural Networks. Implementing these on your own in a good learning experience, but when it comes time for practical applications you are probably going to use existing packages. Starting with the next homework, you will no longer be required to use your own implementation of Neural Networks. Instead, you may use Tensorflow, where SGD and Backpropagation are implemented for you. This week, you should get a head start by installing Tensorflow.

Submit your code (it is fine to copy/paste the tutorial, as long as you get it running) and the exact value you get for the training error on MNIST (it should be approximately 92%).

6 Your Own Question

Write your own question, and provide a thorough solution.

Writing your own problems is a very important way to really learn material. The famous “Bloom’s Taxonomy” that lists the levels of learning is: Remember, Understand, Apply, Analyze, Evaluate, and Create. Using what you know to create is the top-level. We rarely ask you any HW questions about the lowest level of straight-up remembering, expecting you to be able to do that yourself. (e.g. make yourself flashcards) But we don’t want the same to be true about the highest level.

As a practical matter, having some practice at trying to create problems helps you study for exams much better than simply counting on solving existing practice problems. This is because thinking about how to create an interesting problem forces you to really look at the material from the perspective of those who are going to create the exams.

Besides, this is fun. If you want to make a boring problem, go ahead. That is your prerogative. But it is more fun to really engage with the material, discover something interesting, and then come up with a problem that walks others down a journey that lets them share your discovery. You don't have to achieve this every week. But unless you try every week, it probably won't happen ever.