

KNNClassifier Code:

```
class KNNClassifier:
    def __init__(self, k):
        self.k = k

    def train(self, X_train, y_train):
        self.X_train = X_train
        self.y_train = y_train

    def getNN(self, testEntry):
        distances = {}

        # get the distance between the test Entry and every training entry
        for trainEntryIdx in range(len(self.X_train)):
            dist = distance.euclidean(self.X_train[trainEntryIdx], testEntry)
            distances[trainEntryIdx] = dist

        #list of training entry index to that entry's distance to the test entry
        sorted_distance_map = list(sorted(distances.items(), key=lambda item: item[1], reverse=False))
        sorted_distance_map = sorted_distance_map[0:self.k]
        return [self.y_train[idx] for idx, distance in sorted_distance_map]

    def predict(self, X_test): # returns a vector of the predictions y_predict
        numTestEntries = len(X_test)
        #initialize y_predict
        y_predict = [0 for i in range(numTestEntries)]

        for testEntryIdx in range(numTestEntries):
            #get the values of the K nearest neighbors
            nn = self.getNN(X_test[testEntryIdx])

            nnCounts = Counter(nn)
            val, count = nnCounts.most_common()[0]
            y_predict[testEntryIdx] = val
        return y_predict

    def score(self, X_test, y_test):
        y_predict = self.predict(X_test)
        meanAccuracy = 0
        for i in range(len(y_test)):
            if y_test[i] == y_predict[i]:
                meanAccuracy += 1
        meanAccuracy = meanAccuracy/len(y_test)
        return meanAccuracy
```

For the KNNRegression class, everything is the same, except instead of returning the majority class in the **predict** function, it returns the average of the *y_train* values of its nearest neighbors. Additionally, in the **score()** function, it uses the R² error metric instead of the mean accuracy

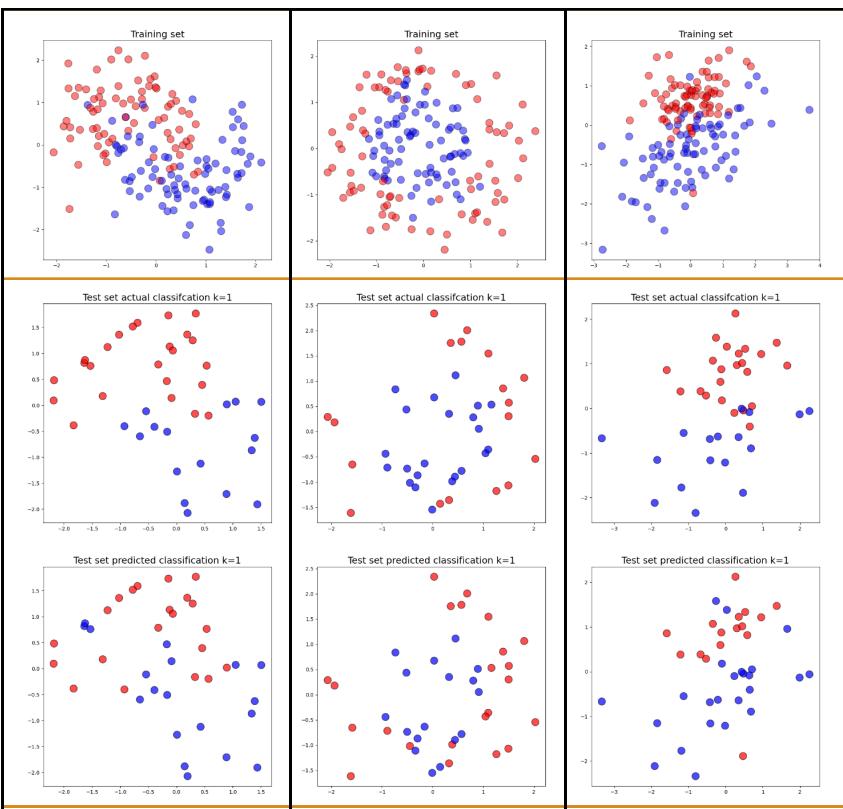
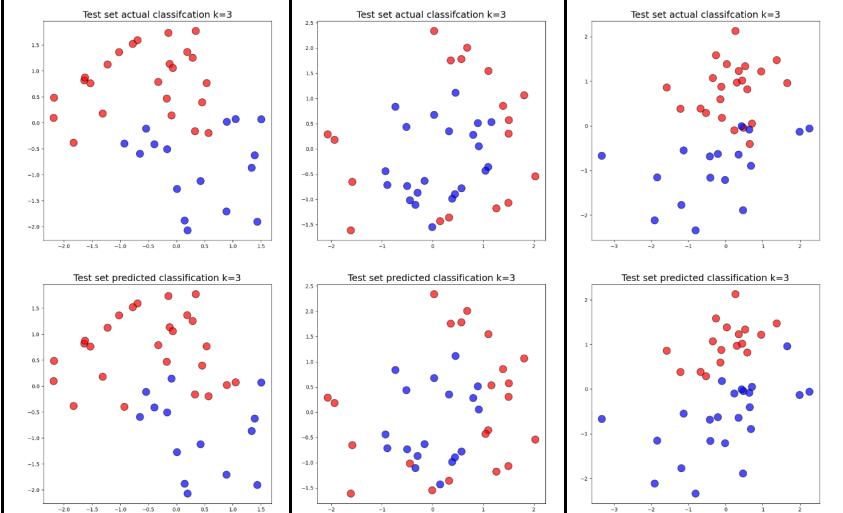
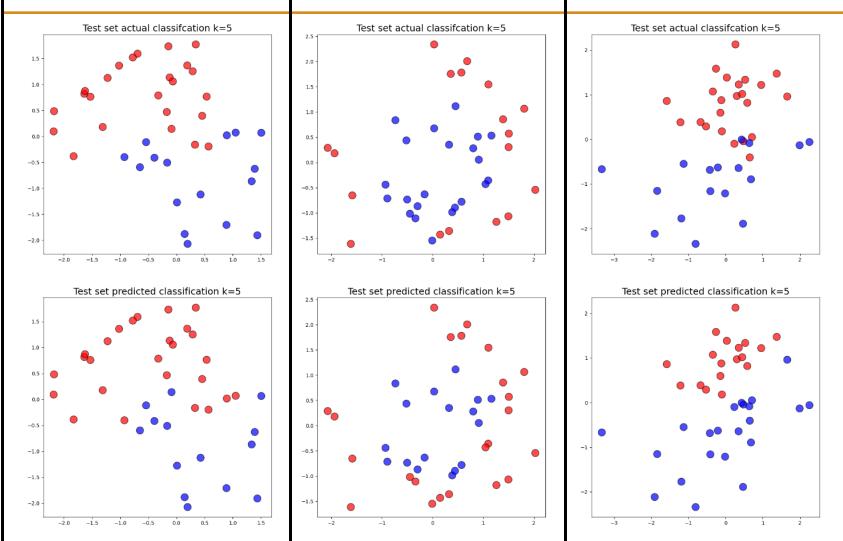
KNNRegression Class

```
class KNNRegression:  
    def __init__(self,k):  
        self.k = k  
  
    def train(self, X_train, y_train):  
        self.X_train = X_train  
        self.y_train = y_train  
  
    def getNN(self, testEntry):  
        distances = {}  
  
        # get the distance between the test Entry and every training entry  
        for trainEntryIdx in range(len(self.X_train)):  
            dist = distance.euclidean(self.X_train[trainEntryIdx], testEntry)  
            distances[trainEntryIdx] = dist  
  
        #list of training entry index to that entry's distance to the test entry  
        sorted_distance_map = list(sorted(distances.items(), key=lambda item: item[1], reverse=False))  
        sorted_distance_map = sorted_distance_map[0:self.k]  
        return [self.y_train[idx] for idx, distance in sorted_distance_map]  
  
    def predict(self, X_test):  
        numTestEntries = len(X_test)  
        #initialize y_predict  
        y_predict = [0 for i in range(numTestEntries)]  
  
        for testEntryIdx in range(numTestEntries):  
            #get the values of the K nearest neighbors  
            nn = self.getNN(X_test[testEntryIdx])  
  
            nnCounts = Counter(nn)  
            val, count = nnCounts.most_common()[0]  
            y_predict[testEntryIdx] = sum(nnCounts.values()) / len(nnCounts)  
        return y_predict  
  
    def score(self, X_test, y_test):  
        y_predict = self.predict(X_test)  
        error = r2_score(y_test, y_predict)  
        return error
```

Data set 1

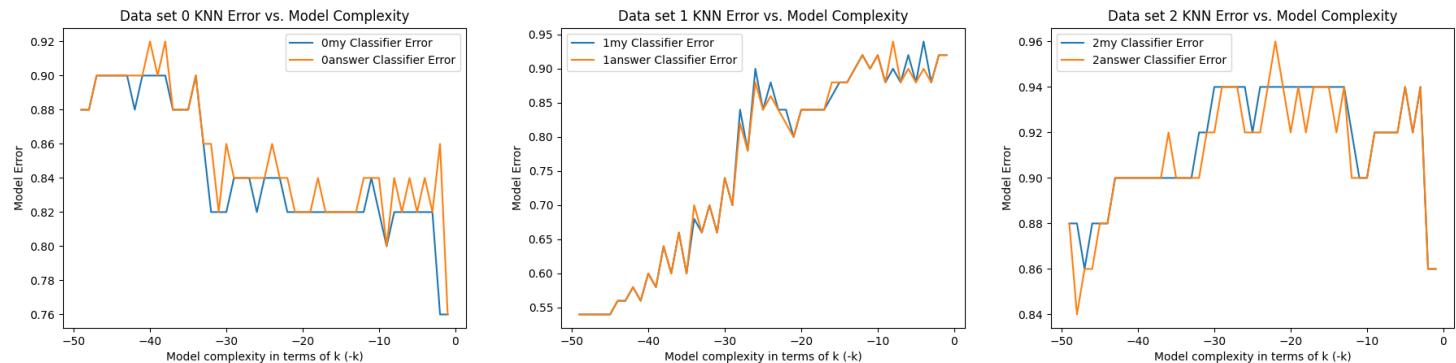
Data set 2

Data set 3

k=1**k=3****k=5**

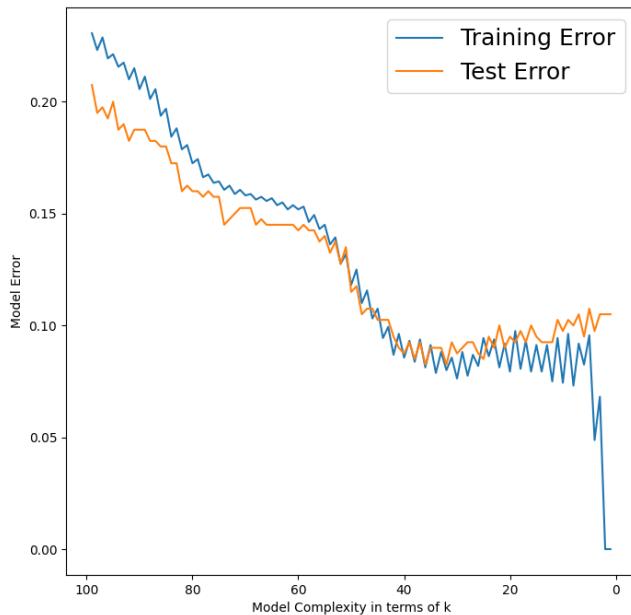
(Model Complexity vs. Error)

My KNN Classification against Correct Solution



The point of these figures is to compare the error values I got using my KNN Classifier to an existing KNN classifier. My results matched up relatively closely with that of the imported KNN classifier. Moreover, they also show model complexity versus the model's error (shown in better detail below).

My KNNClassifier Training and Testing Error vs Model Complexity



How do Results change with different values of k?

k is the **hyperparameter** for this model which controls the model's complexity. When k is large ($k = n$), the model complexity is low as the model will just predict the majority group every time. When k is small ($k=1$), the model complexity is high.

In the figures above, I plotted the model's error (comparing $y_{predict} = KNN(X_test)$ with y_{test}) against the model's complexity (since small K means high model complexity and large k means small model complexity, I plotted it against $-k$ so as the x-axis increases, so does the model complexity). Additionally, I compared my results to a correct implementation of KNNClassifier and KNNRegressor (my results are in **blue**, the correct implementation results are in **orange**).

When comparing the training and testing error, we see that when model complexity is low, both the training and testing errors are similar. When model complexity is high, although training error is low, the testing error is high. This makes sense as if we overfit our model to our training data with high model complexity, it's unlikely to generalize well on future/test data.

Linear Regression:

2. Consider the regression model $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{X} \in \mathbb{R}^{n \times p}$ has rank $p < n$, $\beta \in \mathbb{R}^p$, and $\epsilon \in \mathbb{R}^n$ a random vector. Recall that the ordinary least squares (OLS) estimator is given by $\hat{\beta}_{OLS} = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$. Assume ϵ is such that $E(\epsilon) = \mathbf{0}$ and $\operatorname{cov}(\epsilon) = \sigma^2 \mathbf{I}$.

- (a) Provide an expression for $\hat{\beta}_{OLS}$ by solving the optimization problem above.

$$\begin{aligned}\hat{\beta}_{OLS} &= \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 = \underset{\beta}{\operatorname{argmin}} ((\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)) \\ &= \underset{\beta}{\operatorname{argmin}} ((\mathbf{y}^T - \beta^T \mathbf{X}^T)(\mathbf{y} - \mathbf{X}\beta)) \\ &= \underset{\beta}{\operatorname{argmin}} ((\mathbf{y}^T \mathbf{y} - 2(\mathbf{X}\beta)^T \mathbf{y} + (\mathbf{X}\beta)^T (\mathbf{X}\beta)))\\ \text{To find } \hat{\beta} \text{ (the argmin),} \\ \text{take the derivative wrt } \beta \text{ and} \\ \text{set to 0 to find the min.} \quad 0 = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \hat{\beta}_{OLS} \\ \mathbf{X}^T \mathbf{X} \hat{\beta}_{OLS} &= \mathbf{X}^T \mathbf{y} \\ \hat{\beta}_{OLS} &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})\end{aligned}$$

- (b) Calculate $E(\hat{\beta}_{OLS})$ and $\operatorname{cov}(\hat{\beta}_{OLS})$.

$$\begin{aligned}E(\hat{\beta}_{OLS}) &= E((\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})) \\ &= E(\beta + \mathbf{X}^{-1} \epsilon) \\ &= E(\beta) + \mathbf{X}^{-1} E(\underbrace{\epsilon}_0) \quad \text{linearity of expectation} \\ &= \beta + 0 \\ &= \boxed{\beta}\end{aligned}$$

sidenote

$$\begin{aligned}&(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T (\mathbf{X}\beta + \epsilon)) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}\beta + \mathbf{X}^T \epsilon) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})\beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\ &= \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\ &= \beta + \mathbf{X}^{-1} \epsilon\end{aligned}$$

$$\begin{aligned}\operatorname{cov}(\hat{\beta}_{OLS}) &= E[(\hat{\beta}_{OLS} - E(\hat{\beta}_{OLS})) (\hat{\beta}_{OLS} - E(\hat{\beta}_{OLS}))^T] \\ &= E((\beta + \mathbf{X}^{-1} \epsilon) - \beta)((\beta + \mathbf{X}^{-1} \epsilon) - \beta)^T) \quad \text{from above, } E(\hat{\beta}_{OLS}) = \beta \\ &= E((\mathbf{X}^{-1} \epsilon)(\mathbf{X}^{-1} \epsilon)^T) \quad \text{and } \hat{\beta}_{OLS} = \beta + \mathbf{X}^{-1} \epsilon \\ &= E(\mathbf{X}^{-1} \epsilon \epsilon^T (\mathbf{X}^{-1})^T) \\ &= \mathbf{X}^{-1} E(\epsilon \epsilon^T) (\mathbf{X}^{-1})^T \\ &= \sigma^2 \mathbf{X}^{-1} (\mathbf{X}^{-1})^T \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

\Leftarrow since $E(\epsilon) = 0$, $\operatorname{cov}(\epsilon) = E(\epsilon \epsilon^T) - E(\epsilon)^T$
 $\Rightarrow \sigma^2 \mathbf{I} = E(\epsilon \epsilon^T)$

(c) Show that the OLS predictions, $\hat{y} = \mathbf{X}\hat{\beta}_{OLS}$, can be written as $\hat{y} = \mathbf{H}\mathbf{y}$ for some matrix \mathbf{H} .

$$\hat{y} = \mathbf{X}\hat{\beta}_{OLS}$$

$$\hat{y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}(\mathbf{X}^\top \mathbf{y})$$

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top \text{ satisfies the equation}$$

Matrix Analysis:

2. Answer the following:

- (a) Let $\mathbf{x} \in \mathbb{R}^p$ be a random vector with known mean $E(\mathbf{x}) = \boldsymbol{\mu}$ and known covariance $\text{cov}(\mathbf{x}) = \boldsymbol{\Sigma}$, where $\boldsymbol{\Sigma}$ is positive definite. Find an affine transformation of \mathbf{x} such that the transformed variable, $\tilde{\mathbf{x}}$, is such that $E(\tilde{\mathbf{x}}) = \mathbf{0}$ and $\text{cov}(\tilde{\mathbf{x}}) = \mathbf{I}$

$$\tilde{\mathbf{x}} = A\mathbf{x} + b$$

- $E(\tilde{\mathbf{x}}) = E(A\mathbf{x} + b) = AE[\mathbf{x}] + b = \mathbf{0} \Rightarrow b = -A\boldsymbol{\mu}$
- $\text{Cov}(\tilde{\mathbf{x}}) = E((\tilde{\mathbf{x}} - E(\tilde{\mathbf{x}}))(\tilde{\mathbf{x}} - E(\tilde{\mathbf{x}}))^\top)$

$$= E((\tilde{\mathbf{x}})(\tilde{\mathbf{x}})^\top)$$

$$= E((A\mathbf{x} + b)(A\mathbf{x} + b)^\top)$$

$$= E((A\mathbf{x} - A\boldsymbol{\mu})(A\mathbf{x} - A\boldsymbol{\mu})^\top)$$

$$= E(A(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top A^\top)$$

$$\stackrel{\text{substitute}}{=} E((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) A^\top$$

$$\stackrel{\mathbf{I}}{=} A E((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) A^\top$$

$$\stackrel{\mathbf{I}}{=} A \sum \mathbf{A}^\top$$

$\boldsymbol{\Sigma}$ is positive definite \exists matrix $R \in \mathbb{R}^{p \times p}$ such that $\boldsymbol{\Sigma} = R^\top R$

$$\begin{aligned} \mathbf{I} &= A R^\top R A^\top \\ (A^\top A)^{-1} &= R^\top R \\ A^\top A &= R^{-1} (R^\top)^{-1} = R^{-1} (R^{-1})^\top \\ \Rightarrow A^\top &= R^{-1} \Rightarrow A = (R^\top)^{-1} \end{aligned}$$

$$\therefore \text{Let } A = (R^\top)^{-1} \text{ and } b = - (R^\top)^{-1} \boldsymbol{\mu}$$

$$\tilde{\mathbf{x}} = A\mathbf{x} + b$$

(b) Let $x, y \in \mathbb{R}^n$. Find the eigenvalues and eigenvectors of xy^T .

$$\underset{n \times 1}{xy^T} \underset{n \times 1}{v} = \underset{n \times 1}{\lambda v}$$

notice how $y^T x \in \mathbb{R}$ is a scalar.

We see that $(xy^T)x = x(y^T x) = (y^T x)x$ since $y^T x \in \mathbb{R}$

$\Rightarrow \lambda_1 = y^T x$ is an eigenvalue of xy^T with corresponding eigenvector x .

Moreover, it's the only eigenvalue since if λ_2 with v_2 s.t $xy^T v_2 = \lambda_2 v_2$

$$\text{then, } \underset{\epsilon \mathbb{R}}{x(y^T v_2)} = \underset{\epsilon \mathbb{R}}{(y^T v_2)x} = \underset{\epsilon \mathbb{R}}{\lambda_2 v_2}$$

$$v_2 = c x \text{ for some } c \in \mathbb{R}$$

$\Rightarrow v_2 = cx$ for some $c \in \mathbb{R} \Rightarrow x$ and v_2 correspond to the same eigenvalue

$$\Rightarrow \lambda_1 = \lambda_2.$$

$\therefore y^T x$ is the only eigen value of xy^T with corresponding vector x .

(c) Let A be a real symmetric $p \times p$ matrix with eigenvalues $\lambda_1, \dots, \lambda_p$. Show that $\sum_{i=1}^p \sum_{j=1}^p a_{ij}^2 = \sum_{i=1}^p \lambda_i^2$.

Note how :

$$\begin{aligned} (A^2)_{ii} &= a_{1i} \cdot a_{1i} + a_{2i} \cdot a_{2i} + a_{3i} \cdot a_{3i} + \dots + a_{ni} \cdot a_{ni} \\ &= a_{1i}^2 + a_{2i}^2 + \dots + a_{ni}^2 \quad \text{since } A \text{ is symmetric} \end{aligned}$$

$$\begin{aligned} (a) \operatorname{tr}(A^2) &= \sum_{i=1}^p (A^2)_{ii} \\ &= a_{11}^2 + a_{21}^2 + \dots + a_{n1}^2 \\ &\quad + a_{12}^2 + a_{22}^2 + \dots + a_{nn}^2 \\ &\quad \vdots \\ &\quad + a_{1n}^2 + a_{2n}^2 + \dots + a_{nn}^2 \\ &= \sum_{i=1}^p \sum_{j=1}^p a_{ij}^2 \end{aligned}$$

(b) Since A is symmetric $\Rightarrow \exists U \in \mathbb{R}^n$ s.t $U \Lambda U^T = A$

$$\begin{aligned} \operatorname{tr}(A^2) &= \operatorname{tr}(U \Lambda U^T U \Lambda U^T) \\ &= \operatorname{tr}(U \Lambda U U^T) \\ &= \operatorname{tr}(\Lambda^2) \\ &= \sum_{i=1}^p \lambda_i^2 \quad \text{since } \Lambda \text{ is diagonal} \end{aligned}$$

$$\text{Combining (a) \& (b) } \Rightarrow \sum_{i=1}^p \sum_{j=1}^p a_{ij}^2 = \sum_{i=1}^p \lambda_i^2 \quad \blacksquare$$

Multivariate Statistics:

3. Suppose $\{\mathbf{x}_i\}_{i=1}^n$ are n IID samples from a p -variate normal distribution with known mean $\mathbf{0}$ and unknown covariance Σ (positive definite). Show that

$$\hat{\Sigma} = \frac{\mathbf{X}^T \mathbf{X}}{n}$$

is the maximum likelihood estimate of Σ , where

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$$

is the *data matrix*.

Here are some suggested steps to get you started:

- (i) Write down the (joint) likelihood for this problem.
- (ii) Calculate the log-likelihood (dropping constant terms as appropriate).
- (iii) Take the gradient of the log-likelihood to get the score equations.
- (iv) Solve for the maximum likelihood estimate.

Hint: It may be easier to parametrize the problem in terms of the precision (inverse covariance) matrix $\Theta = \Sigma^{-1}$, find $\hat{\Theta}$, and then use the invariance of the maximum likelihood estimate to find $\hat{\Sigma} = (\hat{\Theta})^{-1}$.

(i) $\mathbf{x}_i \sim N_p(\mathbf{0}, \Sigma)$ Since each $\mathbf{x}_i^{(1)}$ are IID, the joint likelihood:

$$L(\Sigma | \mathbf{X}) = \prod_{i=1}^n f_{\mathbf{x}_i}(\mathbf{x}_i) = \prod_{i=1}^n \underbrace{\left(\frac{1}{(2\pi)^{\frac{p}{2}} \cdot |\Sigma|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(\mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i)} \right)}_{f_{\mathbf{x}_i}(\mathbf{x}_i)}$$

(ii) log-likelihood:

$$\begin{aligned} l(\Sigma | \mathbf{X}) &= \log(L(\Sigma | \mathbf{X})) \\ &= \sum_{i=1}^n \log \left(\frac{1}{(2\pi)^{\frac{p}{2}} \cdot |\Sigma|^{\frac{1}{2}}} \cdot e^{-\frac{1}{2}(\mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i)} \right) \quad \text{properties of logs} \\ &= \sum_{i=1}^n \log(2\pi)^{-\frac{p}{2}} + \sum_{i=1}^n \log|\Sigma|^{-\frac{1}{2}} + \sum_{i=1}^n -\frac{1}{2}(\mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i) \\ &= C + \frac{n}{2} \log(\frac{1}{|\Sigma|}) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i) \quad \text{where } C \text{ is a constant WRT. } \Sigma \end{aligned}$$

(iii) Following the hint, take the gradient wrt $\Theta = \Sigma^{-1}$

$$\begin{aligned} \nabla_{\Theta} l(\Theta | \mathbf{X}) &= \nabla_{\Theta} \left[\frac{n}{2} \log(|\Theta|) + \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \Theta \mathbf{x}_i) \right] \\ &= \frac{n}{2} (\Theta^{-1})^T - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \quad (\mathbf{x}^T \mathbf{x})^T = \mathbf{x}^T \mathbf{x} \\ 0 &= \frac{n}{2} (\hat{\Theta}^{-1})^T - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \end{aligned}$$

$$\Rightarrow (\hat{\theta}^{-1})^T = \frac{1}{n} X^T X \Rightarrow \hat{\theta}^{-1} = \frac{1}{n} X^T X$$

$$\Rightarrow \hat{\Sigma} = (\hat{\theta})^{-1} \Rightarrow \hat{\Sigma} = \frac{X^T X}{n}$$

Convex Optimization:

4. For each of the following, is the function convex in $\mathbf{x} \in \mathbb{R}^n$? Strictly convex? Concave? None of the above? Briefly justify your answer.

(a) $\max(0, 1 - y * \alpha^T \mathbf{x})$ for $y \in \{\pm 1\}$. - convex: a hyperplane that becomes 0 

(b) $\log(1 + e^{\alpha^T \mathbf{x}})$ - strictly convex:

(c) $-y \log(\alpha^T \mathbf{x}) - (1 - y) \log(1 - \alpha^T \mathbf{x})$ for $y \in \{0, 1\}$ - equation of a line?

(d) $\|\mathbf{y} - \mathbf{Ax}\|_1$ - convex:  visualizing the graph, there are constant slope lines, so not strictly convex

(e) $\|\mathbf{y} - \mathbf{Ax}\|_2^2 + \|\mathbf{x}\|_2^2$ - convex: both $\|\mathbf{y} - \mathbf{Ax}\|_2^2$ and $\|\mathbf{x}\|_2^2$ are convex \Rightarrow their sum is convex
(we optimize this form in Ridge Regression, so makes sense that it is convex)

In each of the above, $\alpha \in \mathbb{R}^n$ is a fixed vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a full-rank matrix with $m < n$.

5. ~~For parts (c) and (e) above, calculate the gradient of the function and write out a gradient descent (or ascent) method in pseudo-code to minimize (or maximize) the function.~~

didn't get to this