

COMS W4721 Spring 2020 Homework 2: Linear Classifiers, Decision Trees

Instruction

Please prepare your write-up as a typeset PDF document (which can be generated using LaTex or Word). If you choose to hand-write certain portions of the assignment, make sure your handwriting is legible and the consistency in the format with other pages which are typeset (e.g. page size, page numbering). If we cannot read your handwriting, you may not receive credit for the question. This write-up contains all of your supporting materials which include plots, source code, and proofs for each of the parts in the assignment. Submit your assignment on Gradescope by clearly marking the pages for each part. On the first page of your write-up, please typeset: (1) your name and your UNI; (2) all of your collaborators whom you discussed the assignment with; (3) the parts of the assignment you had collaborated on. The solutions to the problems need to start from the second page. Please write up the solutions by yourself. The academic rules of conduct is found in the course syllabus.

Suggestions

If necessary, please define notations and explain reasoning behind the solutions as concisely as possible. Solutions without explanations when needed may receive no credit. Points can be deducted for solutions with unnecessarily long explanations for lack of clarity. Source code comment can be useful for explaining the logic behind your solutions. Please start early!

Problem 1: Linear Regression (20 points)

In our lecture, we discussed linear regression in which the data $\{(\mathbf{x}^{(i)} \in \mathbb{R}^d, y^{(i)})\}_{i=1}^N$ is generated such that $y^{(i)}$'s independent of others when conditioned on $\mathbf{x}^{(i)}$.

Let us assume that $p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right)$. In this problem we will explore the Bayesian formulation of linear regression called *maximum a posteriori estimate*:

$$\begin{aligned}
\mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N) \\
&= \arg \max_{\mathbf{w}} \frac{p(\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N | \mathbf{w}) p(\mathbf{w})}{\int_{\mathbf{w}'} p(\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N | \mathbf{w}') p(\mathbf{w}') d\mathbf{w}'} \\
&= \arg \max_{\mathbf{w}} p(\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N | \mathbf{w}) p(\mathbf{w}) \\
&= \arg \max_{\mathbf{w}} p(\{\mathbf{x}^{(i)}\}_{i=1}^N | \mathbf{w}) p(\{y^{(i)}\}_{i=1}^N | \{\mathbf{x}^{(i)}\}_{i=1}^N; \mathbf{w}) p(\mathbf{w}) \\
&= \arg \max_{\mathbf{w}} \ln p(\{\mathbf{x}^{(i)}\}_{i=1}^N | \mathbf{w}) + \ln p(\{y^{(i)}\}_{i=1}^N | \{\mathbf{x}^{(i)}\}_{i=1}^N; \mathbf{w}) + \ln p(\mathbf{w}) \\
&= \arg \max_{\mathbf{w}} \ln p(\{y^{(i)}\}_{i=1}^N | \{\mathbf{x}^{(i)}\}_{i=1}^N; \mathbf{w}) + \ln p(\mathbf{w}) \\
&= \arg \max_{\mathbf{w}} \left(\sum_{i=1}^N \ln p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right) + \ln p(\mathbf{w})
\end{aligned}$$

where the integral in the denominator can be dropped since it has no dependence on \mathbf{w} and the natural log can be taken because it is a monotonic transform. For the last equation $\ln p(\{\mathbf{x}^{(i)}\}_{i=1}^N; \mathbf{w})$ is dropped because the data has no dependence on \mathbf{w} . Depending on the prior function on the weight $p(\mathbf{w})$, we can get several flavors of Bayesian linear regression.

- (a) (10 points) Suppose each component of \mathbf{w} is selected such that $w_i \sim N(0, \tau^2)$ i.i.d. What is \mathbf{w}_{MAP} in this case? Have we seen this form before?
- (b) (10 points) Let us assume each component of \mathbf{w} is selected such that $w^{(i)} \sim \text{Laplace}(0, b)$ i.i.d. What is \mathbf{w}_{MAP} in this case? Have we seen this form before?

problem 1

$$a) \quad w_i \sim N(0, \tau^2) \\ p(w | 0, \tau^2) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\tau^2}} e^{-\frac{1}{2\tau^2} (w_i)^2}$$

$$\begin{aligned} w_{map} &= \arg \max_w \left[\sum_{i=1}^N \ln p(y^{(i)} | x^{(i)}, w) + \ln p(w) \right] \\ f &= \sum_{i=1}^N \ln \frac{1}{\sqrt{2\pi\tau^2}} \exp \left(\frac{-(y^{(i)} - w^T x^{(i)})^2}{2\tau^2} \right) + \sum_{i=1}^d \ln \left(\frac{1}{\sqrt{2\pi\tau^2}} \exp \left(-\frac{1}{2\tau^2} (w_i)^2 \right) \right) \\ &= \ln \frac{1}{\sqrt{2\pi\tau^2}} - \frac{1}{2\tau^2} \sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2 + \ln \frac{d}{\sqrt{2\pi\tau^2}} - \frac{1}{2\tau^2} \sum_{i=1}^d (w_i^2) \\ &= \ln \frac{1}{\sqrt{2\pi\tau^2}} - \frac{1}{2\tau^2} (y - xw)^T (y - xw) + \ln \frac{d}{\sqrt{2\pi\tau^2}} - \frac{1}{2\tau^2} w^T w \end{aligned}$$

Take derivative

$$\begin{aligned} \frac{df}{dw} &= \frac{1}{2\tau^2} (x^T y - x^T x w) - \frac{1}{\tau^2} w \\ &= \frac{1}{2\tau^2} x^T (y - xw) - \frac{1}{\tau^2} w \stackrel{\text{set } 0}{=} 0 \\ \frac{1}{2\tau^2} x^T (y - xw) &= \frac{1}{\tau^2} w \\ x^T y - x^T x w &= \frac{2\tau^2}{\tau^2} w \\ w_{map} &= x^T y \cdot (x^T x + \frac{2\tau^2}{\tau^2} I)^{-1} \end{aligned}$$

From ridge regression, we have similar form

$$w_{RR} = (\lambda I + x^T x)^{-1} \cdot x^T y$$

b) $w^{(i)} \sim \text{Laplace}(0, b)$

$$p(w_i | b) = \frac{1}{\pi} \frac{1}{2b} e^{-\frac{|w_i|}{b}}$$

$$\begin{aligned} w_{map} &= \arg \max \left(\sum_{i=1}^N \ln \frac{1}{\sqrt{2\pi\tau^2}} \exp \left(\frac{-(y^{(i)} - w^T x^{(i)})^2}{2\tau^2} \right) + \sum_{i=1}^d \ln \frac{1}{2b} \exp \left(-\frac{|w_i|}{b} \right) \right) \\ &= \arg \max \left(-\frac{1}{2\tau^2} \sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2 - \frac{1}{b} \sum_{i=1}^d |w_i| \right) \\ &= \arg \min \left(\sum_{i=1}^N (y^{(i)} - w^T x^{(i)})^2 + \frac{2b^2}{b} \|w\| \right) \end{aligned}$$

From Lasso regression, we see similar formula

$$\hat{w}_{Lasso} = \arg \min \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \hat{w}^T x^{(i)})^2 + \lambda \|w\|$$

Problem 2: Gaussian Discriminant Analysis (15 points)

- (a) (5 points) Prove that a Gaussian discriminant analysis in \mathbb{R}^d induces a quadratic decision boundary. You can assume a binary classification task for the purpose of this task. What properties of the parameters will ensure a linear decision boundary?
- (b) (10 points) Let's take the Gaussian discriminant analysis in \mathbb{R}^d from part (a). What properties of the parameters can make the classifier a Gaussian Naive Bayes Classifier? Give a short mathematical proof to validate your answer.

problem 2.

(i) From bayes rule :

$$f_{\text{bayes}}(x) = \arg \max p(x|y) \cdot p(y) \quad \text{where } p(x|y) \text{ is class conditional density}$$

and $p(y)$ is class prior

Assume we have data

$$\{(x^{(i)}, y^{(i)})\}_{i=1}^n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$$

$x^{(i)} \in X \subseteq \mathbb{R}^d, y^{(i)} \in Y \subseteq \{0, 1\}$

Then from Gaussian Assumption

$$p(x, y) = p(y)p(x|y) = \begin{cases} p_0 \frac{1}{\sqrt{2\pi}\delta_0} e^{-\frac{(x-\mu_0)^2}{2\delta_0^2}} & \text{if } y=0 \\ p_1 \frac{1}{\sqrt{2\pi}\delta_1} e^{-\frac{(x-\mu_1)^2}{2\delta_1^2}} & \text{if } y=1 \end{cases}$$

The Bayes optimal one under the assumed joint distribution depends on

$$\mathbb{I}(\Pr(y=1|x) \geq \Pr(y=0|x)) \stackrel{\text{bayes}}{\Rightarrow} \mathbb{I}(p(x|y=1)p(y=1) \geq p(x|y=0)p(y=0))$$

$$\Rightarrow \mathbb{I}\left(-\frac{(x-\mu_1)^2}{2\delta_1^2} - \log \sqrt{2\pi}\delta_1 + \log p_1 \geq -\frac{(x-\mu_0)^2}{2\delta_0^2} - \log \sqrt{2\pi}\delta_0 + \log p_0\right)$$

$$\Rightarrow \mathbb{I}(ax^2 + bx + c \geq 0) \quad \text{So the decision boundary is not linear}$$

In matrix form

Gaussian Assumption

$$p(x|y, \mu_y, \Sigma_y) = \frac{1}{(2\pi)^{d/2} |\Sigma_y|^{1/2}} \exp\left(-\frac{(x-\mu_y)^T (\Sigma_y)^{-1} (x-\mu_y)}{2}\right)$$

$$\text{let } d_{\Sigma}(x, \mu) = (x-\mu)^T (\Sigma)^{-1} (x-\mu)$$

$$f(x) = \mathbb{I}(\Pr(y=1|x) > \Pr(y=0|x)) \Rightarrow \mathbb{I}\left(\ln \frac{\Pr(y=1|x)}{\Pr(y=0|x)} > 0\right)$$

$$\Rightarrow \mathbb{I}\left(\ln \frac{p(x|y, \mu_1, \Sigma) p(y=1)}{p(x|y, \mu_0, \Sigma) p(y=0)}\right) + (\Pr(y=1) = \pi_1)$$

$$\Rightarrow \mathbb{I}\left(\ln \frac{\pi_1}{1-\pi_1} - \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} (d_{\Sigma_1}(x, \mu_1) - d_{\Sigma_0}(x, \mu_0)) > 0\right)$$

Since $d_{\Sigma}(x, \mu)$ is in quadratic form, GDA has a quadratic decision boundary.

Further, if we assume $\Sigma_1 = \Sigma_0 = \Sigma$

$$\begin{aligned} d_{\Sigma}(x, \mu_1) - d_{\Sigma}(x, \mu_0) &= (x-\mu_1)^T (\Sigma)^{-1} (x-\mu_1) - (x-\mu_0)^T (\Sigma)^{-1} (x-\mu_0) \\ &= x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \mu_1 \Sigma^{-1} \\ &\quad - x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} x - \mu_0^T \Sigma^{-1} \mu_0 \end{aligned}$$

$$= x^T \Sigma^{-1} (\mu_0 - \mu_1) - \mu_1^T \Sigma^{-1} \mu_1 + \mu_0^T \Sigma^{-1} \mu_0$$

So it is linear in x .

(2) From Gaussian Discriminant Analysis

$$\begin{aligned}
 \hat{f}(x) &= \arg \max_{y \in \{0,1\}} \Pr(y|x) = \arg \max_{y \in \{0,1\}} \underbrace{\Pr(y)}_{\pi_y} \Pr(x|y; \mu_y, \Sigma_y) \\
 &= \arg \max_{y \in \{0,1\}} \frac{1}{(2\pi)^{d/2} |\Sigma_y|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) \right\} \cdot \pi_y \\
 &= \mathbb{I} [\Pr(y=1|x) > \Pr(y=0|x)] \\
 &= \mathbb{I} \left[\ln \frac{\pi_1}{1-\pi_1} - \frac{1}{2} \ln \frac{\Sigma_1}{\Sigma_0} - \frac{1}{2} (\ln \Sigma_1(\mu_1, \mu_0) - \ln \Sigma_0(\mu_1, \mu_0)) > 0 \right] \\
 \text{Since } \Sigma_1, \Sigma_0 \text{ are diagonal.} \quad |\Sigma_1| &= \prod_{k=1}^d \lambda_k, \quad |\Sigma_0| = \prod_{k=1}^d \lambda_k' \\
 d_{\Sigma_1}(x, \mu_1) &= (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = \sum_{k=1}^d \frac{(x_k - \mu_{1,k})^2}{\lambda_k} \\
 &= \mathbb{I} \left[\ln \frac{\pi_1}{1-\pi_1} - \frac{1}{2} \sum_{k=1}^d \left[\ln \frac{\lambda_{k,y=1}}{\lambda_{k,y=0}} \right] - \frac{1}{2} \left(\sum_{k=1}^d \left(\frac{x_k - \mu_{1,k}, y=1}{\lambda_{k,y=1}} \right)^2 - \sum_{k=1}^d \left(\frac{x_k - \mu_{1,k}, y=0}{\lambda_{k,y=0}} \right)^2 \right) > 0 \right]
 \end{aligned}$$

From Gaussian Naive Bayes classifier

$$\begin{aligned}
 \hat{f}(x) &= \arg \max_{y \in \{0,1\}} \prod_{k=1}^d \Pr(x_k | y; \mu_k, \lambda_k) \cdot \underbrace{\Pr(y)}_{\pi_y} \quad x_k \in \{0,1\} \\
 &= \mathbb{I} \left[\prod_{k=1}^d \frac{1}{\sqrt{2\pi\lambda_k}} \exp \left\{ -\frac{1}{2} \left(\frac{x_k - \mu_{k,y=1}}{\lambda_{k,y=1}} \right)^2 \right\} \cdot \pi_1 > \prod_{k=1}^d \frac{1}{\sqrt{2\pi\lambda_k}} \exp \left\{ -\frac{1}{2} \left(\frac{x_k - \mu_{k,y=0}}{\lambda_{k,y=0}} \right)^2 \right\} \cdot \pi_0 \right] \\
 &= \mathbb{I} \left[\ln \frac{\prod_{k=1}^d \frac{1}{\sqrt{2\pi\lambda_{k,y=1}}} \exp \left\{ -\frac{1}{2} \left(\frac{x_k - \mu_{k,y=1}}{\lambda_{k,y=1}} \right)^2 \right\} \cdot \pi_1}{\prod_{k=1}^d \frac{1}{\sqrt{2\pi\lambda_{k,y=0}}} \exp \left\{ -\frac{1}{2} \left(\frac{x_k - \mu_{k,y=0}}{\lambda_{k,y=0}} \right)^2 \right\} \cdot \pi_0} > 0 \right] \\
 &= \mathbb{I} \left[\ln \frac{\pi_1}{1-\pi_1} - \frac{1}{2} \sum_{k=1}^d \ln \frac{\lambda_{k,y=1}}{\lambda_{k,y=0}} - \frac{1}{2} \left(\sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=1}}{\lambda_{k,y=1}} \right)^2 - \sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=0}}{\lambda_{k,y=0}} \right)^2 \right) > 0 \right]
 \end{aligned}$$

Therefore, when the covariance matrix Σ_y is diagonal
Gaussian NBC is a special case of Gaussian Discriminant Analysis.

Problem 3: Logistic Regression (30 points)

- (a) (15 points) Logistic Regression can be used to solve the binary classification task as depicted in the figure. A simple logistic regression model is given below:

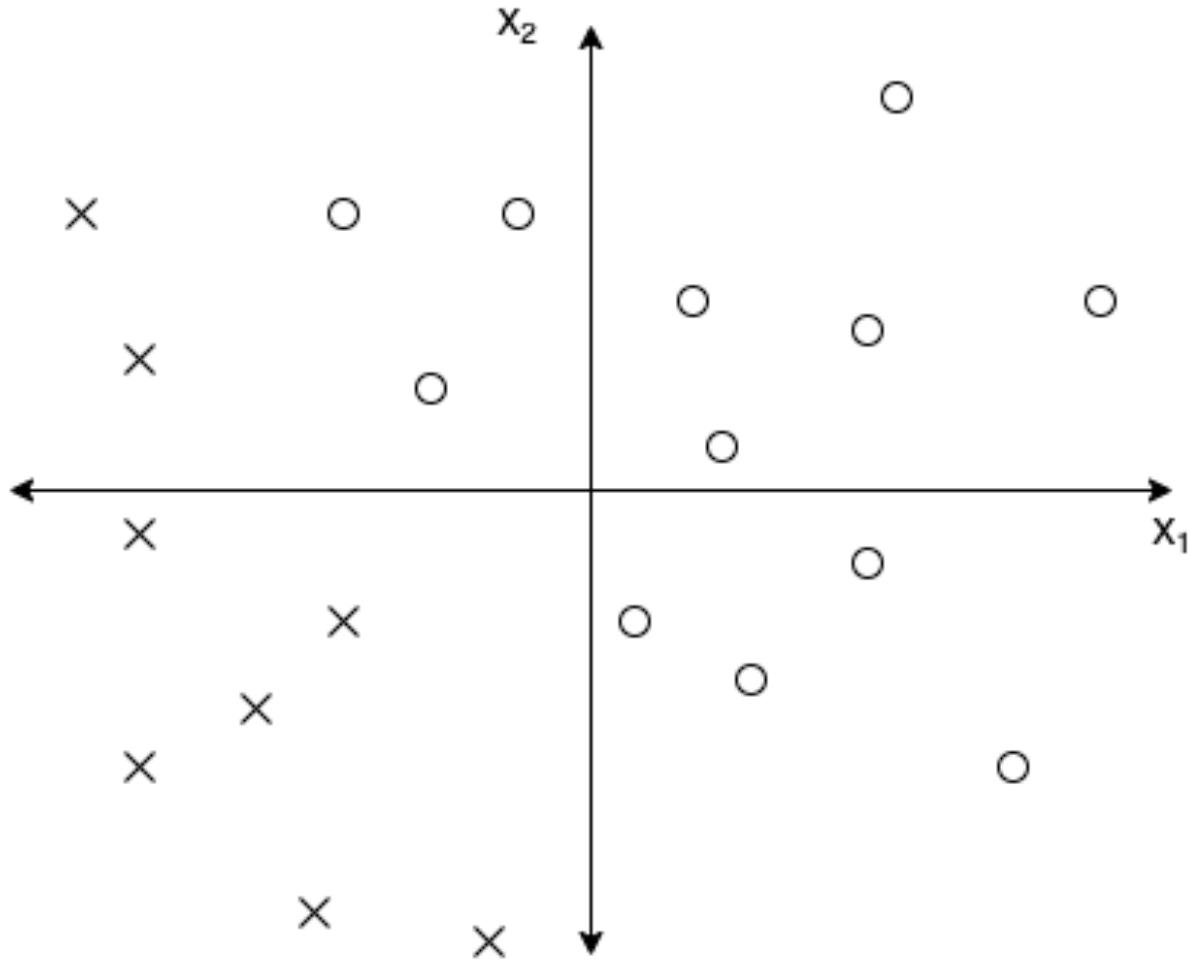


Figure 1: Binary classification dataset

$$P(y = \text{class1} | \vec{x}, \vec{w}) = \frac{1}{1 + \exp(-w_0 - w_1 x_1 - w_2 x_2)}$$

The given data is linearly separable at this point and thus a logistic regression model can be fit to separate the data with zero training error.

Consider a regularized logistic regression model where our optimization (maximization in this case) function becomes:

$$\sum_{i=1}^N \log(P(y^{(i)} | x^{(i)}; w_0, w_1, w_2)) - \lambda \cdot w_j^2$$

for a very large λ . The regularization penalties used penalize one parameter at a time, ie. either one of $\{w_0, w_1, w_2\}$ are penalized at a given time. How does the training error change with regularization of each parameter? Provide a brief mathematical justification for each of your answers.

- (b) (15 points) You are given a dataset *logistic_regression.csv* which has two features x_1, x_2 and the corresponding *class* label. Please write a small program in a language of your choice to optimize the logistic regression function to

- Fit the data in the csv file without regularization
- Regularize the squared weight of w_1 associated with the feature x_1 .
For each value of $\lambda \in [10^0, 10^1, 10^2, 10^3, 10^4, 10^5]$: plot the decision boundary along with the points in the dataset. You should have 6 plots.
- Regularize the squared weight of w_2 associated with the feature x_2 .
For each value of $\lambda \in [10^3, 10^4, 10^5, 10^6, 10^7, 10^8]$: plot the decision boundary along with the points in the dataset. You should have 6 plots.

Please attach all plots for all of the three subparts with your analysis. You may use `scipy.optimize.minimize` for implementing your code. You do not need to include the code. You may find the Python notebook for problem 3 useful for the starter code (please see section with "Modify Me"'s).

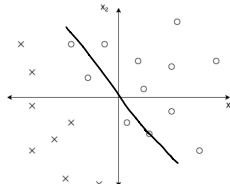
problem 3

a) optimization function :

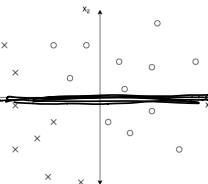
$$\sum_{i=1}^N \log(P(y^{(i)} | x^{(i)}, w_0, w_1, w_2)) - \lambda w_j^2$$

when λ is very large, the penalty term has large impact. To minimize the cost function, w_j has to be 0.

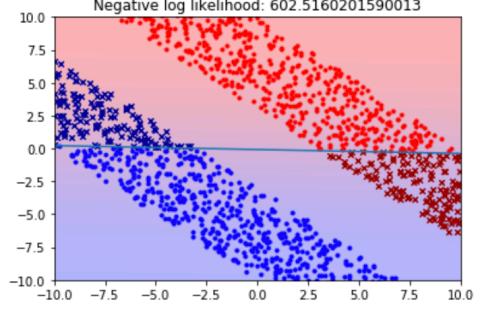
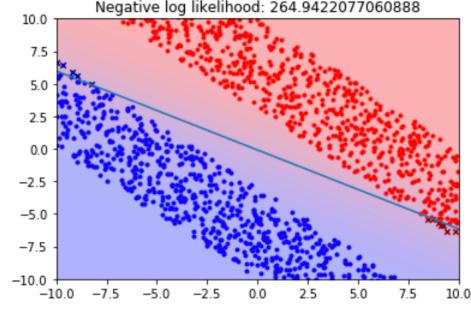
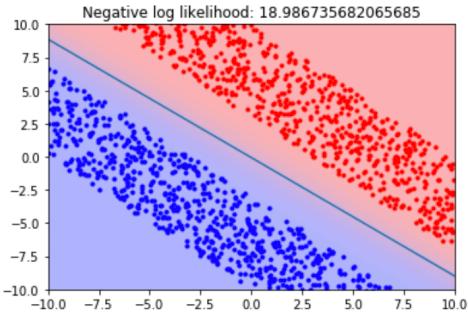
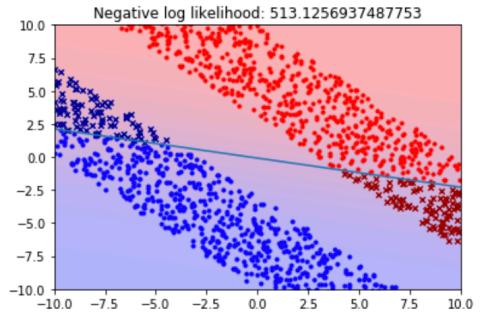
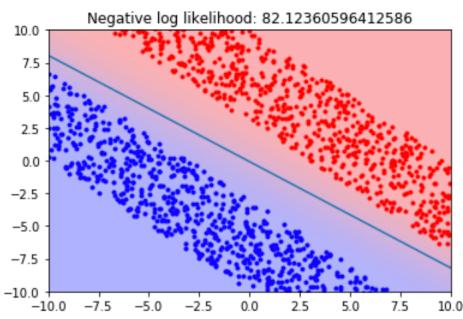
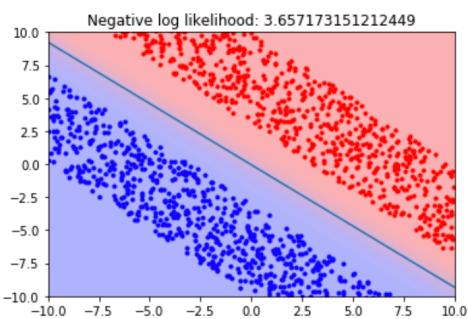
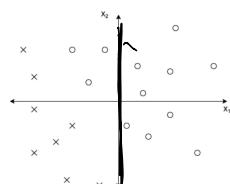
when $w_0 = 0$,
train error increase,
around 3-4 points
misclassified.

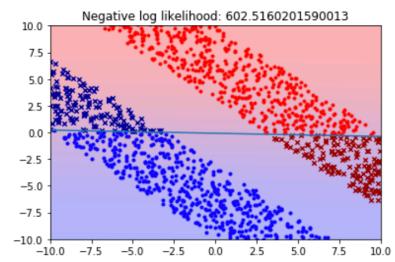
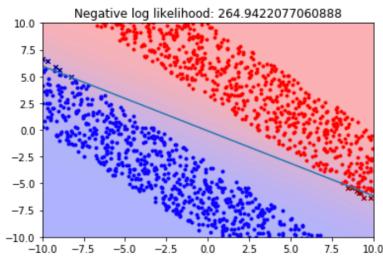
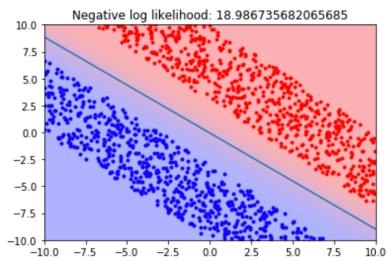
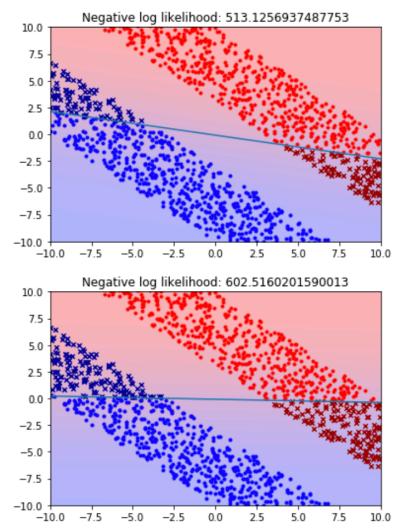
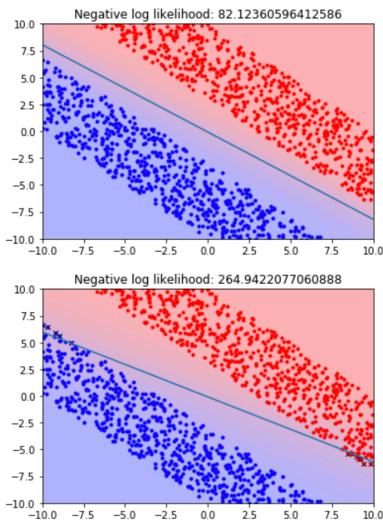
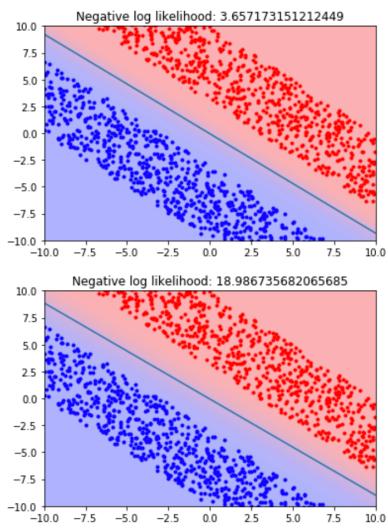


when $w_1 = 0$
train error increase most
around 6 points
misclassified.



when $w_2 = 0$
train error increase
around 3 points
misclassified.





Problem 4: Decision Tree (35 points)

Decision trees are useful for classification. They follow a tree structure by splitting the data among the feature at each level which is most descriptive (gives largest information gain). For more information on decision trees see [this chapter of Mitchell's machine learning book](#).

We will be classifying whether a banknote is fraudulent or not using this [Banknotes dataset](#). The dataset provides signal-based features of the banknotes and we will predict whether the banknote is fraudulent (1) or not (0).

We will walk you through writing some key functions of a decision tree from scratch.

- (a) (3 points) In the decision tree iPython notebook, start by importing the data (numpy array or pandas dataframe both work) and splitting the dataset by train and test (usually 80% train, 20% test). If you'd like to shuffle the data before splitting, you can use `numpy.random.shuffle`
- (b) (8 points) Next, you will implement some functions in the `class DecisionTree`, which uses the `Node` class to build a decision tree. For each level, we calculate the current uncertainty, then split upon the feature at the threshold which gives us the highest information gain. First, implement the function `get_uncertainty` which implements the entropy and gini index uncertainty depending on the keyword `metric` that is given.
- (c) (8 points) Using the `get_uncertainty` method, implement the `getInfoGain` method which calculates the information gain when splitting the data on `node.data` on `split_index`. `node.data` stores the relevant data at each node of the decision tree (root note will have all the data, then its children will split the data in 2 parts on the `split_index`, etc.)
- (d) (10 points) Lastly, in `get_feature_threshold` find the feature that gives the largest information gain and `assign` the feature number (column number) to `node.feature` while updating the threshold values (check the docstrings for more detailed information).
- (e) (6 points) After you implement the above steps, the `buildTree` function will recursively build the decision tree. Then when you run `homework_evaluate` which will call `self.predict` to evaluate the accuracy on the test set. Try different values of K (depth of tree a.k.a. number of features the tree will split on) and compare the performance. Which feature gives the largest information gain? Which feature is the least useful for the decision tree? How does varying the uncertainty metric (gini versus entropy) vary the performance (and why)? (For these questions we're looking for some analysis of the resulting model and some thought. Does not have to be a long paragraph.)

Congrats! You built a decision tree! Now (optional) try it on other data for fun. We've included a breast cancer dataset with both cardinal and categorical variables for you to try. With categorical variables, we one-hot encode them (e.g. for a 3-class categorical variable, 2 is represented as [0, 1, 0]) to allow the decision tree to classify them correctly (think about why this is the case, hint: does the magnitude of the number give us useful information here?).

Hello and welcome to decision trees!

Decision trees are often pretty effect learning algorithms, and certainly serve as an interesting technical exercise in data preprocessing and recursion. This notebook will walk you through some of the basic notions of how the implementation should be executed, and also give you a chance to write some of your own code.

Suggested reading before you start:

<https://www.cs.princeton.edu/courses/archive/spring07/cos424/papers/mitchell-dectrees.pdf>

(<https://www.cs.princeton.edu/courses/archive/spring07/cos424/papers/mitchell-dectrees.pdf>)

Throughout the notebook you'll see TODO tags in the comments. This is where you should insert your own code to make the functions work! If you get stuck, we encourage you to come to office hours. You can also try to look at APIs and documentation online to try to get a sense how certain methods work. If you take inspiration from any source online other than official documentation, please be sure to cite the resource! Good luck!

In [3]:

```
import pandas as pd
import scipy.io
import numpy as np
import sklearn.model_selection as model_selection
```

We will be using decision trees to classify if a banknote is fradulent (class 1) or not fradulent (class 0). Download data from <https://archive.ics.uci.edu/ml/datasets/banknote+authentication#>
(<https://archive.ics.uci.edu/ml/datasets/banknote+authentication#>)

Import data

In [4]:

```
# TODO YOUR CODE HERE
def import_data(split = 0.8, shuffle=False):
    """Read in the data, split it by split percentage into train and test data,
    and return X_train, y_train, X_test, y_test as numpy arrays"""
    df = pd.read_csv("data_banknote_authentication.txt", header = None)
    dfv = df.values
    if shuffle == True:
        np.random.shuffle(dfv)
        df = pd.DataFrame(dfv)

    y = np.array(df[4])
    X = np.array(df.drop([4], axis=1))
    X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, train_size=split, random_state=0)
    print("data imported")

    return X_train, y_train, X_test, y_test
```

In [5]:

```
class Node:  
    """Each node of our decision tree will hold values such as left and right children,  
    the data and labels being split on, the threshold value & index in the datafram  
e for a particular feature,  
    and the uncertainty measure for this node"""  
    def __init__(self, data, labels, depth):  
        """  
        data: X data  
        labels: y data  
        depth: depth of tree  
        """  
        self.left = None  
        self.right = None  
  
        self.data = data  
        self.labels = labels  
        self.depth = depth  
  
        self.threshold = None # threshold value  
        self.threshold_index = None # threshold index  
        self.feature = None # feature as a NUMBER (column number)  
        self.label = None # y label  
        self.uncertainty = None # uncertainty value
```

In [10]:

```

class DecisionTree:
    def __init__(self, K=5, verbose=False):
        """
        K: number of features to split on
        """

        self.root = None
        self.K = K
        self.verbose = verbose

    def buildTree(self, data, labels):
        """Builds tree for training on data. Recursively called _buildTree"""
        self.root = Node(data, labels, 0)
        if self.verbose:
            print("Root node shape: ", data.shape, labels.shape)
        self._buildTree(self.root)

    def _buildTree(self, node):
        # get uncertainty measure and feature threshold
        node.uncertainty = self.get_uncertainty(node.labels)
        self.get_feature_threshold(node)

        index = node.data[:, node.feature].argsort() # sort feature for return
        node.data = node.data[index]
        node.labels = node.labels[index]

        # check label distribution.
        label_distribution = np.bincount(node.labels)
        majority_label = node.labels[0] if len(label_distribution) == 1 else np.argmax(label_distribution)

        if self.verbose:
            print("Node uncertainty: %f" % node.uncertainty)

        # Split left and right if threshold is not the min or max of the feature or
        # every point has the
        # same label.
        if node.threshold_index == 0 or node.threshold_index == node.data.shape[0]
        or \
            len(label_distribution) == 1:
            node.label = majority_label
        else:
            node.left = Node(node.data[:node.threshold_index], node.labels[:node.threshold_index], node.depth + 1)
            node.right = Node(node.data[node.threshold_index:], node.labels[node.threshold_index:], node.depth + 1)
            node.data = None
            node.labels = None

        # If in last layer of tree, assign predictions
        if node.depth == self.K:
            if len(node.left.labels) == 0:
                node.right.label = np.argmax(np.bincount(node.right.labels))
                node.left.label = 1 - node.right.label
            elif len(node.right.labels) == 0:

```

```

        node.left.label = np.argmax(np.bincount(node.left.labels))
        node.right.label = 1 - node.left.label
    else:
        node.left.label = np.argmax(np.bincount(node.left.labels))
        node.right.label = np.argmax(np.bincount(node.right.labels))
    return

else: # Otherwise continue training the tree by calling _buildTree
    self._buildTree(node.left)
    self._buildTree(node.right)

def predict(self, data_pt):
    return self._predict(data_pt, self.root)

def _predict(self, data_pt, node):
    feature = node.feature
    threshold = node.threshold
    if node.label is not None:
        return node.label
    elif data_pt[node.feature] < node.threshold:
        return self._predict(data_pt, node.left)
    elif data_pt[node.feature] >= node.threshold:
        return self._predict(data_pt, node.right)

def get_feature_threshold(self, node, metric="entropy"):
    """ TODO Find the feature that gives the largest information gain. Update node.threshold,
    node.threshold_index, and node.feature (a number representing the feature.
    e.g. 2nd column feature would be 1)
    Make sure to sort the columns of data before you try to find the threshold
    index (look at numpy argsort) and set the values
    for node.threshold, node.threshold_index, and node.feature
    return: None
    """
    node.threshold = 0
    node.threshold_index = 0
    node.feature = 0
    # TODO YOUR CODE HERE
    gain = 0
    feature = 0
    for col in range(node.data.shape[1]):
        node.feature = col
        for i in range(node.labels.shape[0]):
            new_gain = self.getInfoGain(node, i, metric = metric)
            if new_gain > gain:
                gain = new_gain
                node.threshold_index = i
                feature = col
    node.feature = feature
    node.threshold = sorted(node.data[:,node.feature])[node.threshold_index]

def getInfoGain(self, node, split_index, metric = "entropy"):
    """
    TODO Get information gain using the variables in the parameters,
    split_index: index in the feature column that you are splitting the classes
    on
    """

```

```

    return: information gain (float)
    """
    # TODO YOUR CODE HERE
    index = node.data[:,node.feature].argsort()

        leftInfo = self.get_uncertainty(node.labels[index[:split_index]] , metric=metric)
        rightInfo = self.get_uncertainty( node.labels[index[split_index:]], metric=metric)
        selfInfo = self.get_uncertainty(node.labels, metric=metric)
        n = len(node.labels)
        result = selfInfo - (leftInfo*(split_index/n) + rightInfo*((n-split_index)/n))

    return result

def get_uncertainty(self, labels, metric="entropy"):
    """
    TODO Get uncertainty. Implement entropy AND gini index metrics.
    np.bincount(labels) and labels.shape might be useful here
    return: uncertainty (float)
    """

    if labels.shape[0] == 0:
        return 1

    p = np.bincount(labels)/labels.shape[0]
    entropy = -sum(p*np.log(p))
    gini = sum(p*(1-p))

    result = 1
    if metric == "entropy":
        result = entropy
    elif metric == "gini":
        result = gini

    return result

def printTree(self):
    """Prints the tree including threshold value and feature name"""
    self._printTree(self.root)

def _printTree(self, node):
    if node is not None:
        if node.label is None:
            print("\t" * node.depth, "(%d, %d)" % (node.threshold, node.feature))
        else:
            print("\t" * node.depth, node.label)
            self._printTree(node.left)
            self._printTree(node.right)

    def homework_evaluate(self, X_train, labels, X_test, y_test):
        n = X_train.shape[0]

```

```
n = X_train.shape[0]

count = 0
for i in range(n):
    if self.predict(X_train[i]) == labels[i]:
        count += 1

print("The decision tree is %d percent accurate on %d training data" % ((count / n) * 100, n))

n = X_test.shape[0]

count = 0
for i in range(n):
    if self.predict(X_test[i]) == y_test[i]:
        count += 1

print("The decision tree is %d percent accurate on %d test data" % ((count / n) * 100, n))

return count / n
```

Run the tree

Try different values of K (depth of tree a.k.a. number of features the tree will split on) and compare the performance. Which feature gives the largest information gain? Which feature is the least useful for the decision tree?

In [103]:

```
def featureList(node, features):
    if node is not None and node.feature is not None:
        features.append(node.feature)
        featureList(node.left, features)
        featureList(node.right, features)
```

In [104]:

```
X_train, y_train, X_test, y_test = import_data(split=0.8)

tree = DecisionTree(K=10, verbose=False)
tree.buildTree(X_train, y_train)
tree.homework_evaluate(X_train, y_train, X_test, y_test)
features = []
featureList(tree.root,features)
print(features)

data imported

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: divide by zero encountered in log
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: invalid value encountered in multiply

The decision tree is 99 percent accurate on 1097 training data
The decision tree is 98 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 0, 0, 0, 0, 1, 2,
0, 0, 0]
```

In [105]:

```
k_list = range(1, 7)
for i in k_list:
    tree = DecisionTree(K=i, verbose=False)
    tree.buildTree(X_train, y_train)
    print('K={}'.format(i))
    tree.homework_evaluate(X_train, y_train, X_test, y_test)
    features = []
    featureList(tree.root, features)
#tree.printTree()
print(features)
print('The feature gives largest information gain is {}'.format(max(set(features), key = features.count)))
print('The feature gives least information gain is {}'.format(min(set(features), key = features.count)))
```

```
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: divide by zero encountered in log
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: invalid value encountered in multiply

K=1:
The decision tree is 89 percent accurate on 1097 training data
The decision tree is 89 percent accurate on 275 test data
[0, 1, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 1

K=2:
The decision tree is 95 percent accurate on 1097 training data
The decision tree is 93 percent accurate on 275 test data
[0, 1, 2, 0, 0, 2, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 1

K=3:
The decision tree is 98 percent accurate on 1097 training data
The decision tree is 96 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 2, 1, 3, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=4:
The decision tree is 98 percent accurate on 1097 training data
The decision tree is 97 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 1, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=5:
The decision tree is 99 percent accurate on 1097 training data
The decision tree is 98 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 0, 0, 1, 2, 0, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=6:
The decision tree is 99 percent accurate on 1097 training data
The decision tree is 98 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 0, 0, 0, 0, 1, 2, 0, 0, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3
```

In [106]:

```
k_list = range(1, 7)
for i in k_list:
    tree = DecisionTree(K=i, verbose=False)
    tree.buildTree(X_train, y_train, metric = "gini")
    print('K={}'.format(i))
    tree.homework_evaluate(X_train, y_train, X_test, y_test)
    features = []
    featureList(tree.root, features)
    print(features)
    print('The feature gives largest information gain is {}'.format(max(set(features), key = features.count)))
    print('The feature gives least information gain is {}'.format(min(set(features), key = features.count)))
```

```
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: divide by zero encountered in log
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: invalid value encountered in multiply

K=1:
The decision tree is 89 percent accurate on 1097 training data
The decision tree is 89 percent accurate on 275 test data
[0, 1, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 1

K=2:
The decision tree is 95 percent accurate on 1097 training data
The decision tree is 93 percent accurate on 275 test data
[0, 1, 2, 0, 0, 2, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 1

K=3:
The decision tree is 98 percent accurate on 1097 training data
The decision tree is 96 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 2, 1, 3, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=4:
The decision tree is 98 percent accurate on 1097 training data
The decision tree is 97 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 1, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=5:
The decision tree is 99 percent accurate on 1097 training data
The decision tree is 98 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 0, 0, 1, 2, 0, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3

K=6:
The decision tree is 99 percent accurate on 1097 training data
The decision tree is 98 percent accurate on 275 test data
[0, 1, 2, 0, 1, 0, 0, 0, 0, 0, 2, 1, 0, 0, 3, 2, 0, 0, 0, 0, 1, 2, 0, 0, 0]
The feature gives largest information gain is 0
The feature gives least information gain is 3
```

As we can see from the above result, the feature 0 has the largest information gain, and the feature 3 has the least. Using Gini and entropy does not have a significant difference. Increase the k would not cause overfit since the test error didn't increase much.

Optional Decision Tree Exercise

This section is designed to give you some exposure to typical preprocessing and allow you to run your decision tree code on another example.

All of the preprocessing has been done for you — there's nothing you need to fill in, but it may be worthwhile to tinker with some of the pieces to make sure you understand how everything fits together.

Otherwise, if your decision tree code works on the bank notes example, you should be able to run through this straight away.

Step 1: Data Preprocessing

To start, you'll need to download the data files from <https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer/> (<https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer/>). The file `breast-cancer.data` contains the actual data you'll need and the file `breast-cancer.names` gives some information about the researchers and the data types (it's probably worth looking at to give you a sense of what's going on).

Note: If you try to open `breast-cancer.data` or `breast-cancer.names` directly, your computer might not know how to handle the file format. To view them, you need to change the file types from `.data` and `.names` to `.txt` — this can be accomplished by simply changing the file name from `breast-cancer.data` to `breast-cancer.txt`

Now let's load the data file into the notebook. All you need to do is put it in the same directory as the notebook and the cell below should locate the appropriate file.

In [12]:

```
import os
current_directory = os.listdir()

try:
    cancer_files = [file for file in current_directory if 'cancer' in file]
    data_file = [file for file in cancer_files if 'names' not in file][0]
except IndexError:
    print('The breast cancer files were not found. Please upload again.')
    data_file = None

if data_file:
    print(f'The data file has been located as {data_file}.')
```

The data file has been located as `breast-cancer.data`.

!! If the cell above returned a file that you don't recognize as the correct data file, you must go back to the upload step and ensure you have successfully uploaded your files before proceeding. !!

Now we'll try to read the sample into our environment:

In [13]:

```
samples = []

with open(data_file, 'r') as file:
    samples = file.readlines()

samples = [sample.split(',') for sample in samples]

print(f'There are {len(samples)} samples in the dataset')
```

There are 286 samples in the dataset

We can now start defining how we want to map our sample values to numeric features that can be used in the decision tree algorithm below. Let's first store our samples in a pandas DataFrame so that it'll be easier to view and work with.

In [14]:

```
import pandas as pd
import numpy as np

columns = ['class', 'age', 'menopause', 'tumor-size', 'inv-nodes',
           'node-caps', 'deg-malig', 'breast', 'breast-quad', 'irradiat']

samples_df = pd.DataFrame(samples, columns=columns)
samples_df['irradiat'] = samples_df['irradiat'].str.replace('\n', '')

samples_df
```

Out[14]:

	class	age	menopause	tumor-size	inv-nodes	node-caps	deg-malig	breast	breast-quad	irradiat
0	no-recurrence-events	30-39	premeno	30-34	0-2	no	3	left	left_low	no
1	no-recurrence-events	40-49	premeno	20-24	0-2	no	2	right	right_up	no
2	no-recurrence-events	40-49	premeno	20-24	0-2	no	2	left	left_low	no
3	no-recurrence-events	60-69	ge40	15-19	0-2	no	2	right	left_up	no
4	no-recurrence-events	40-49	premeno	0-4	0-2	no	2	right	right_low	no
...
281	recurrence-events	30-39	premeno	30-34	0-2	no	2	left	left_up	no
282	recurrence-events	30-39	premeno	20-24	0-2	no	3	left	left_up	yes
283	recurrence-events	60-69	ge40	20-24	0-2	no	1	right	left_up	no
284	recurrence-events	40-49	ge40	30-34	3-5	no	3	left	left_low	no
285	recurrence-events	50-59	ge40	30-34	3-5	no	3	left	left_low	no

286 rows × 10 columns

Our DataFrame looks great! Now that we have the raw data stored in an interpretable way, it's good practice to make a copy before trying to encode everything — if something goes wrong, it'll be easy to just come back up here and reset the encoded DataFrame.

In [15]:

```
encoded_samples_df = samples_df.copy()
```

In the following few cells, we're going to define the different types of variables that exist in our dataset and make sure we assign the appropriate type of encoding.

To start, we'll define our binary variables (0,1) and create Python dictionaries to map the string labels to binary integer values.

It's also helpful to store all of the columns and dictionaries in two lists so that we can easily reference them later on.

In [16]:

```
binary_cols = ['class', 'breast', 'irradiat']

class_map = {'no-recurrence-events': 0, 'recurrence-events': 1}
breast_map = {'left': 0, 'right': 1}
irrad_map = {'yes': 1, 'no': 0}

binary_maps = [class_map, breast_map, irrad_map]
```

Next, we've defined our ordinal variables (those that have obvious ordered structure). Instead of hard-coding the maps here, it's nice to have a function that can take any number of ordinal columns and instantly create all of the corresponding maps. The code below does that exactly, relying on the integer value of the first number in the range (i.e. '50-65') as the sorting key.

In [17]:

```
ordinal_cols = ['age', 'tumor-size', 'inv-nodes', 'deg-malig']

def first_val(x):
    return eval(x.split('-')[0].replace('"', ''))

ordinal_maps = {}

for col in ordinal_cols:
    uniques = sorted(list(encoded_samples_df[col].unique()), key=first_val)
    val_map = {}
    i = 1

    for val in uniques:
        val_map[val] = i
        i += 1

    ordinal_maps[col] = val_map
```

Finally, we have several columns that take $n > 2$ discrete values. Binary encoding won't work here, so we'll use one-hot encoding. Just like in the cell immediately above, we'll use a function here to take an arbitrary number of columns and encode their unique values as one-hot vectors.

In [18]:

```
one_hot_cols = ['menopause', 'breast-quad', 'node-caps']
one_hot_maps = {}

def one_hot_map(col, df):
    one_hot = {}
    unique_vals = list(df[col].unique())
    one_hot_vecs = np.identity(len(unique_vals))

    for i in range(len(unique_vals)):
        one_hot[unique_vals[i]] = one_hot_vecs[i]

    return one_hot

for col in one_hot_cols:
    one_hot_maps[col] = one_hot_map(col, encoded_samples_df)
```

Now that we have all of our encoding maps set up, we'll zip them all up into a master dictionary so that we can easily iterate through them on our `encoded_samples_df` DataFrame.

In [19]:

```
all_maps = dict(one_hot_maps, **ordinal_maps)
for col, val_map in zip(binary_cols, binary_maps):
    all_maps[col] = val_map
```

Before proceeding, let's double check to make sure you've captured all of the columns:

In [20]:

```
if len(all_maps) != len(columns):
    print("You're missing a map! Go back and double check that you didn't miss a column.")
else:
    print("Looks like you successfully created all the maps! Nice work!")
```

Looks like you successfully created all the maps! Nice work!

And now the moment of truth! Let's apply all of our encoding maps on the DataFrame to turn everything into useable features for our decision trees algorithm.

In [21]:

```
for col in columns:
    encoded_samples_df[col] = encoded_samples_df[col].map(all_maps[col])

encoded_samples_df
```

Out[21]:

	class	age	menopause	tumor-size	inv-nodes	node-caps	deg-malig	breast	breast-quad	irradiat
0	0	2	[1.0, 0.0, 0.0]	7	1	[1.0, 0.0, 0.0]	3	0	[1.0, 0.0, 0.0, 0.0, 0.0, 0.0]	0
1	0	3	[1.0, 0.0, 0.0]	5	1	[1.0, 0.0, 0.0]	2	1	[0.0, 1.0, 0.0, 0.0, 0.0, 0.0]	0
2	0	3	[1.0, 0.0, 0.0]	5	1	[1.0, 0.0, 0.0]	2	0	[1.0, 0.0, 0.0, 0.0, 0.0, 0.0]	0
3	0	5	[0.0, 1.0, 0.0]	4	1	[1.0, 0.0, 0.0]	2	1	[0.0, 0.0, 1.0, 0.0, 0.0, 0.0]	0
4	0	3	[1.0, 0.0, 0.0]	1	1	[1.0, 0.0, 0.0]	2	1	[0.0, 0.0, 0.0, 1.0, 0.0, 0.0]	0
...
281	1	2	[1.0, 0.0, 0.0]	7	1	[1.0, 0.0, 0.0]	2	0	[0.0, 0.0, 1.0, 0.0, 0.0, 0.0]	0
282	1	2	[1.0, 0.0, 0.0]	5	1	[1.0, 0.0, 0.0]	3	0	[0.0, 0.0, 1.0, 0.0, 0.0, 0.0]	1
283	1	5	[0.0, 1.0, 0.0]	5	1	[1.0, 0.0, 0.0]	1	1	[0.0, 0.0, 1.0, 0.0, 0.0, 0.0]	0
284	1	3	[0.0, 1.0, 0.0]	7	2	[1.0, 0.0, 0.0]	3	0	[1.0, 0.0, 0.0, 0.0, 0.0, 0.0]	0
285	1	4	[0.0, 1.0, 0.0]	7	2	[1.0, 0.0, 0.0]	3	0	[1.0, 0.0, 0.0, 0.0, 0.0, 0.0]	0

286 rows × 10 columns

Double check your `encoded_samples_df` here to make sure everything passes the sniff test!

Assuming it all looks good, the final step is to break up our DataFrame into individual samples, unpack the nested arrays into their constituent values, and concatenate everything together into a final sample vector...

In [22]:

```
sample_list = list(encoded_samples_df.to_numpy())
```

In [23]:

```
def unpack_nested_arrays(vec):
    final_vec = np.array([])

    for e in vec:
        if isinstance(e, int):
            e = np.array([e])
        final_vec = np.concatenate((final_vec, e))

    return final_vec
```

In [24]:

```
final_vecs = []
for e in sample_list:
    final_vecs.append(unpack_nested_arrays(e))

final_vecs = np.array(final_vecs, dtype=int)
```

Great! The `final_vecs` variable should now point to a 286x19 numpy.ndarray containing all of our samples. Let's finally move on to the machine learning!

Step 2: Decision Tree

We'll run the decision tree process here again.

In [25]:

```
def import_cancer_data(split=0.8, shuffle=True, CUTOFF=0, bins = 256):

    if shuffle:
        np.random.shuffle(final_vecs)

    num_samples = final_vecs.shape[0]
    num_train_samples = int(num_samples*split)

    train_data, test_data = final_vecs[:num_train_samples, :], final_vecs[num_train_samples:, :]

    X_train = train_data[:, 1:]
    y_train = train_data[:, 0]
    X_test = test_data[:, 1:]
    y_test = test_data[:, 0]

    print(X_train.shape)
    print(y_train.shape)

    return X_train, y_train, X_test, y_test
```

Here's the final fit/evaluation code again, but this time using the breast cancer data. You should get somewhere around 0.75 accuracy for the decision tree on this dataset — not 100%, but not too bad for the humble decision tree!

In [26]:

```
X_train, y_train, X_test, y_test = import_cancer_data(split=0.8)

tree = DecisionTree(K=3, verbose=False)
tree.buildTree(X_train, y_train)

tree.homework_evaluate(X_train, y_train, X_test, y_test)

(228, 18)
(228,)

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: divide by zero encountered in log
//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:130: RuntimeWarning: invalid value encountered in multiply
```

The decision tree is 77 percent accurate on 228 training data
The decision tree is 75 percent accurate on 58 test data

Out[26]:

0.7586206896551724

In []: