

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

Shijie He

Feb. 23, 2020

## Problem 1

- (a) To get the maximum value and argmax for  $f(\mathbf{w}) = \sum_{i=1}^N \ln p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) + \ln(p(\mathbf{w}))$ , we take the derivative w.r.t  $\mathbf{w}$ . First, we do some linear transformation on the equation so that it is easier to do derivatives.

Since  $w_i \sim N(0, \tau^2)$ , we know  $p(w_i | \tau^2) = \frac{1}{\tau\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{w_i}{\tau})^2}$

$$\begin{aligned} f(\mathbf{w}) &= \sum_{i=1}^N \ln p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) + \ln p(\mathbf{w}) \\ &= N \ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \sum_{i=1}^N \frac{-(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2}{2\sigma^2} + N \ln\left(\frac{1}{\tau\sqrt{2\pi}}\right) + \sum_{i=1}^N -\frac{1}{2}\left(\frac{w_i}{\tau}\right)^2 \\ &= N \ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(\mathbf{y} - \mathbf{x}\mathbf{w})^T (\mathbf{y} - \mathbf{x}\mathbf{w})}{2\sigma^2} + N \ln\left(\frac{1}{\tau\sqrt{2\pi}}\right) - \frac{1}{2\tau^2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

Taking the derivative:

$$\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{\sigma^2} \mathbf{x}^T (\mathbf{y} - \mathbf{x}\mathbf{w}) - \frac{1}{\tau^2} \mathbf{w} = 0$$

We get:

$$\begin{aligned} \mathbf{x}^T \mathbf{y} - \mathbf{x}^T \mathbf{x} \mathbf{w} &= \frac{\sigma^2}{\tau^2} \mathbf{w} \\ \mathbf{x}^T \mathbf{y} &= (\mathbf{x}^T \mathbf{x} + \frac{\sigma^2}{\tau^2} \mathbf{I}) \mathbf{w} \end{aligned}$$

Thus,

$$\mathbf{w}_{\text{MAP}} = (\mathbf{x}^T \mathbf{x} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{x}^T \mathbf{y}$$

We have seen this form before in ridge regression. In ridge,  $\hat{w} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T b$ .

(b) Since  $w_i \sim \text{Laplace}(0, b)$ , we know  $p(w_i | b) = \frac{1}{2b} e^{-\frac{|w_i|}{b}}$ .

$$\begin{aligned}
\mathbf{w}_{\text{MAP}} &= \arg \max_{\mathbf{w}} \sum_{i=1}^N \ln p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) + \ln p(\mathbf{w}) \\
&= \arg \max_{\mathbf{w}} N \ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \sum_{i=1}^N \frac{-(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2}{2\sigma^2} + N \ln\left(\frac{1}{2b}\right) + \sum_{i=1}^N -\frac{|w_i|}{b} \\
&= \arg \max_{\mathbf{w}} - \sum_{i=1}^N \frac{(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2}{2\sigma^2} - \frac{1}{b} \|\mathbf{w}\|_1 \\
&= \arg \min_{\mathbf{w}} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 + \frac{2\sigma^2}{b} \|\mathbf{w}\|_1 \\
&= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 + \frac{2\sigma^2}{Nb} \|\mathbf{w}\|_1
\end{aligned}$$

We have seen this form before in LASSO regression. In LASSO,  $\hat{w} = \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 + \lambda \|\mathbf{w}\|_1$ .

## Problem 2

(a) For Gaussian Discriminant Analysis, we need to find  $\hat{f}(x)$  such that

$$\hat{f}(x) = \arg \max_y p(y | \mathbf{x}) = \arg \max_y p(y)p(\mathbf{x} | y; \mu_y, \Sigma_y)$$

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

Thus,

$$\begin{aligned} \hat{f}(x) &= \mathbb{1}[p(y = 1 | \mathbf{x}) > p(y = 0 | \mathbf{x})] \\ &= \mathbb{1}\left[\ln \frac{p(y = 1 | \mathbf{x})}{p(y = 0 | \mathbf{x})} > 0\right] \\ &= \mathbb{1}\left[\ln \frac{p(y = 1)}{p(y = 0)} - \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} (d(\mathbf{x}, \mu_1, \Sigma_1) - d(\mathbf{x}, \mu_0, \Sigma_0)) > 0\right] \end{aligned}$$

where  $d(\mathbf{x}, \mu, \Sigma) = (\mathbf{x} - \mu)^T (\Sigma)^{-1} (\mathbf{x} - \mu)$ .

The decision bound is a quadratic form of  $\mathbf{x}$ . Thus, Gaussian discriminant analysis induces a quadratic decision boundary.

For  $\Sigma_1 = \Sigma_0$ ,

$$d(\mathbf{x}, \mu_1, \Sigma_1) - d(\mathbf{x}, \mu_0, \Sigma_0) = \mathbf{x}^T (\Sigma_1)^{-1} (\mu_1 - \mu_0) - \frac{1}{2} \mu_1^T (\Sigma_1)^{-1} \mu_1 + \frac{1}{2} \mu_0^T (\Sigma_1)^{-1} \mu_0$$

is linear in  $\mathbf{x}$ . Thus,  $\Sigma_1 = \Sigma_0$  gives a linear decision boundary.

(b) We need  $\Sigma_1$  and  $\Sigma_0$  to be diagonal matrix to make the above classifier a Gaussian Naive Bayes Classifier.

Let  $d$  be the length of  $\mathbf{x}$ , and the entries of diagonal matrix  $\Sigma_0$  be  $\sigma_{k,y=1}$  for  $k \in \{1, 2, 3, \dots, d\}$ .

We know  $d(\mathbf{x}, \mu, \Sigma) = (\mathbf{x} - \mu)^T (\Sigma)^{-1} (\mathbf{x} - \mu) = \sum_{k=1}^d \left(\frac{x_k - \mu_k}{\sigma_k}\right)^2$  for diagonal matrix  $\Sigma$ .

Hence,

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \mathbb{1}\left[\ln \frac{p(y = 1)}{p(y = 0)} - \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} (d(\mathbf{x}, \mu_1, \Sigma_1) - d(\mathbf{x}, \mu_0, \Sigma_0)) > 0\right] \\ &= \mathbb{1}\left[\ln \frac{p(y = 1)}{p(y = 0)} - \frac{1}{2} \ln \frac{\prod_{k=1}^d \sigma_{k,y=1}}{\prod_{k=1}^d \sigma_{k,y=0}} - \frac{1}{2} \left(\sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=1}}{\sigma_{k,y=1}}\right)^2 - \sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2\right) > 0\right] \\ &= \mathbb{1}\left[\ln \frac{p(y = 1)}{p(y = 0)} - \frac{1}{2} \sum_{k=1}^d \ln \frac{\sigma_{k,y=1}}{\sigma_{k,y=0}} - \frac{1}{2} \left(\sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=1}}{\sigma_{k,y=1}}\right)^2 - \sum_{k=1}^d \left(\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}}\right)^2\right) > 0\right] \end{aligned}$$

On the other hand for Gaussian Naive Bayes Classifier:

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \arg \max_y \prod_{k=1}^d p(x_k \mid y; \mu_k, \sigma_k^2) p(y) \\
&= \mathbb{1}[p(y=1) \prod_{k=1}^d p(x_k \mid y=1, \mu_{k,y=1}, \sigma_{k,y=1}) > p(y=0) \prod_{k=1}^d p(x_k \mid y=0, \mu_{k,y=0}, \sigma_{k,y=0})] \\
&= \mathbb{1}[\ln \frac{p(y=1) \prod_{k=1}^d p(x_k \mid y=1, \mu_{k,y=1}, \sigma_{k,y=1})}{p(y=0) \prod_{k=1}^d p(x_k \mid y=0, \mu_{k,y=0}, \sigma_{k,y=0})} > 0] \\
&= \mathbb{1}[\ln \frac{p(y=1)}{p(y=0)} - \frac{1}{2} \sum_{k=1}^d \ln \frac{\sigma_{k,y=1}}{\sigma_{k,y=0}} - \frac{1}{2} (\sum_{k=1}^d (\frac{x_k - \mu_{k,y=1}}{\sigma_{k,y=1}})^2 - \sum_{k=1}^d (\frac{x_k - \mu_{k,y=0}}{\sigma_{k,y=0}})^2) > 0]
\end{aligned}$$

Thus, Gaussian Naive Bayes Classifier Gaussian discriminant Classifier when the covariance matrix is diagonal.

### Problem 3: Logistic Regression

- (a) To maximize the optimization function and for large penalty coefficient  $\lambda$ , the optimal solution for the penalized parameter would be closed to 0. That can be shown on the plot that if the penalized parameter is  $w_0$ , the separation line would be a line with no intercept. By choosing the penalized parameter to be  $w_1$ , the separation line would be a horizontal line and if  $w_2$  is chosen, it would be a vertical line.

As we can tell from the plot, horizontal separation line would give the largest training error. Regularization of  $w_1$  and  $w_2$  would give similar training error.

- (b) No regularization:

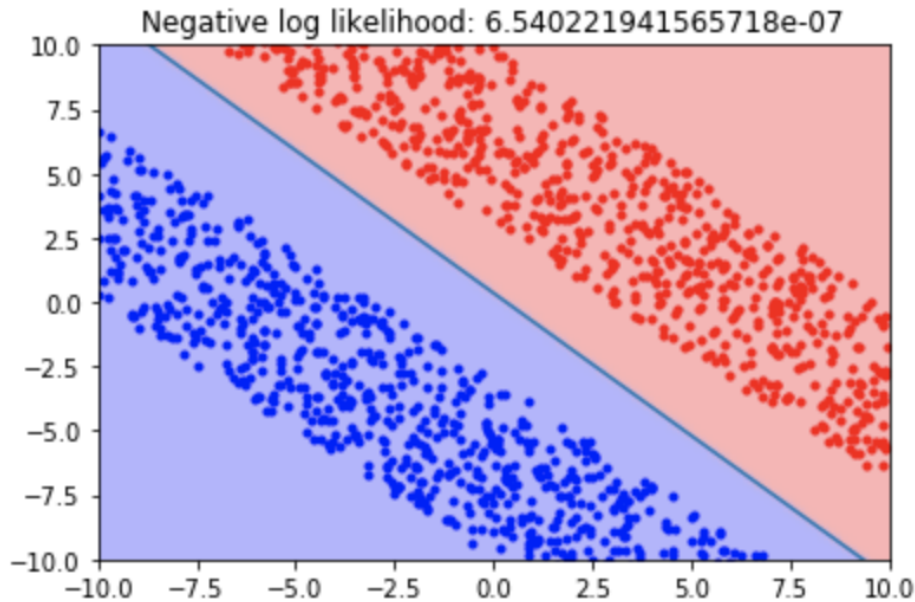


Figure 1: No regularization

We can see that the points are separated in two parts. There exist lines that can perfectly separate the data with no training error. As we can see above, there is no training error with the model without regularization.

Regularization of  $x_1$ :

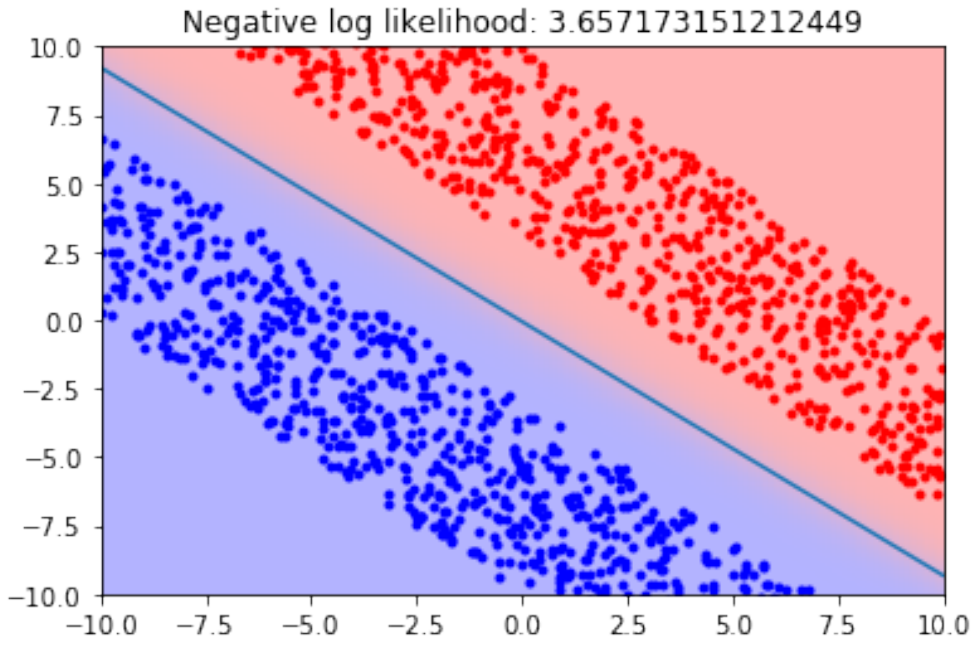


Figure 2: Regularization of  $x_1$ ,  $\lambda = 1$

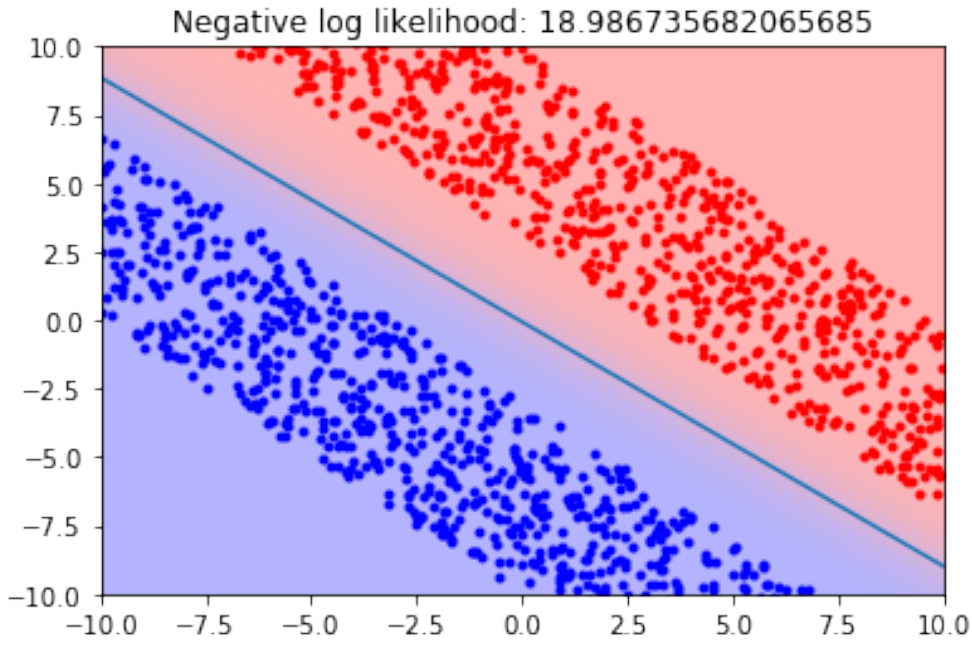


Figure 3: Regularization of  $x_1$ ,  $\lambda = 10$

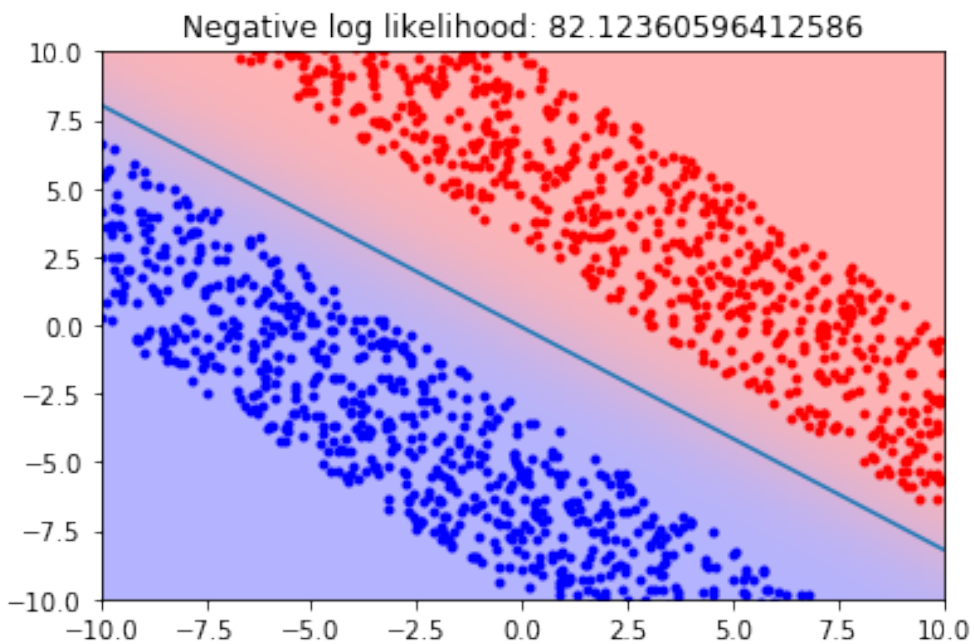


Figure 4: Regularization of  $x_1$ ,  $\lambda = 100$

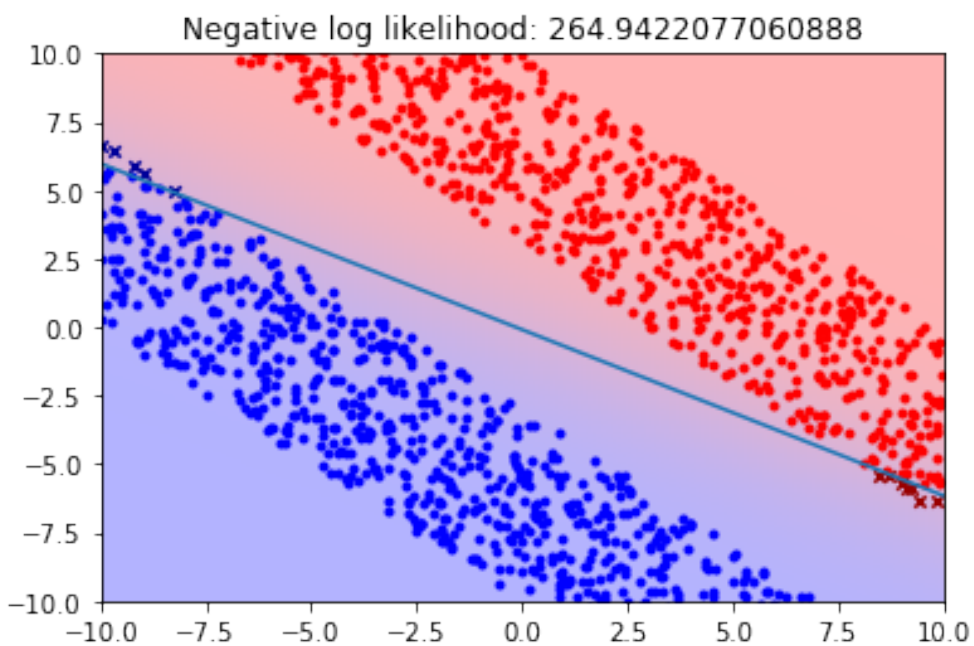


Figure 5: Regularization of  $x_1$ ,  $\lambda = 1000$

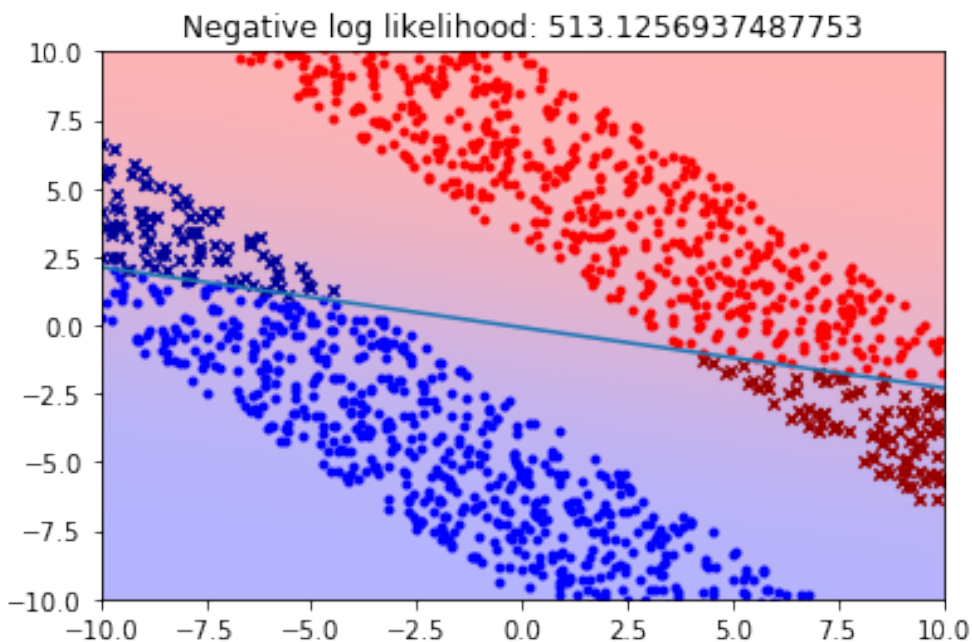


Figure 6: Regularization of  $x_1$ ,  $\lambda = 10000$

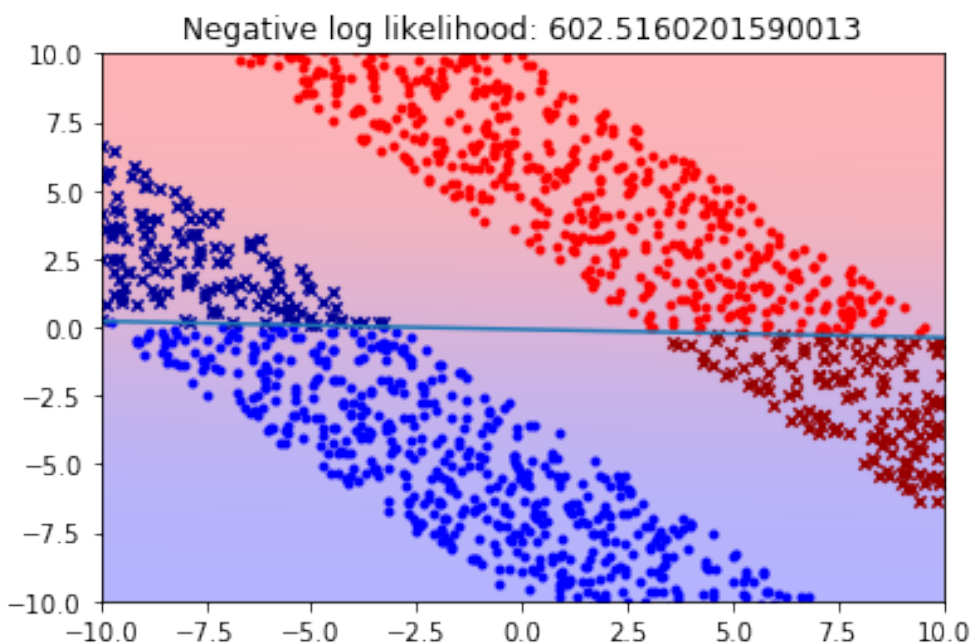


Figure 7: Regularization of  $x_1$ ,  $\lambda = 100000$

As discussed in the first part of this problem, a large regularization parameter will result in a horizontal or vertical separation line. As we can see from the above graphs, as the regularization parameter  $\lambda$  becomes larger, the separation line becomes "more and more" horizontal.



(c) Regularization of  $x_2$ :

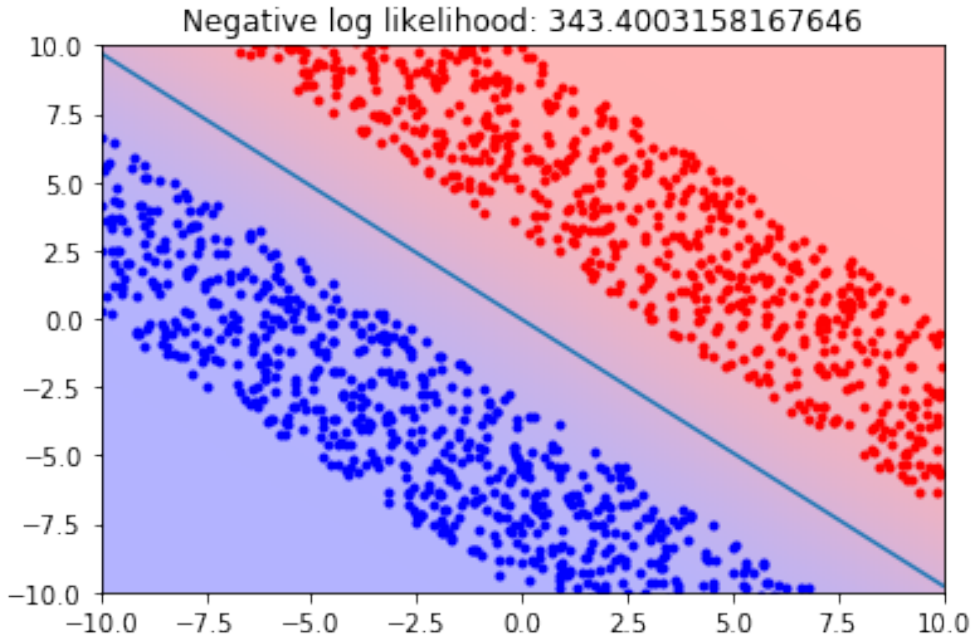


Figure 8: Regularization of  $x_2$ ,  $\lambda = 10^3$

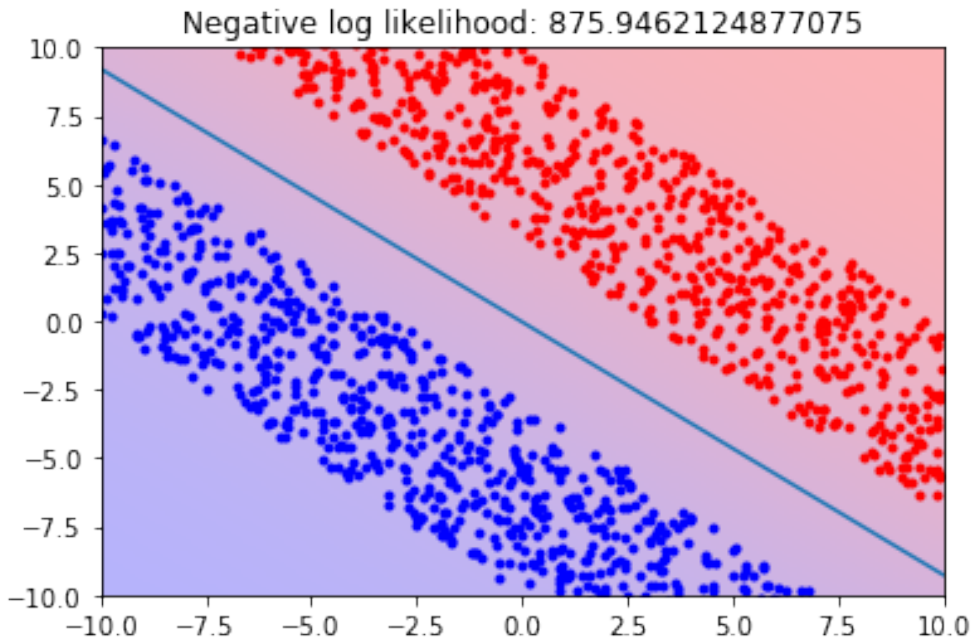


Figure 9: Regularization of  $x_2$ ,  $\lambda = 10^4$

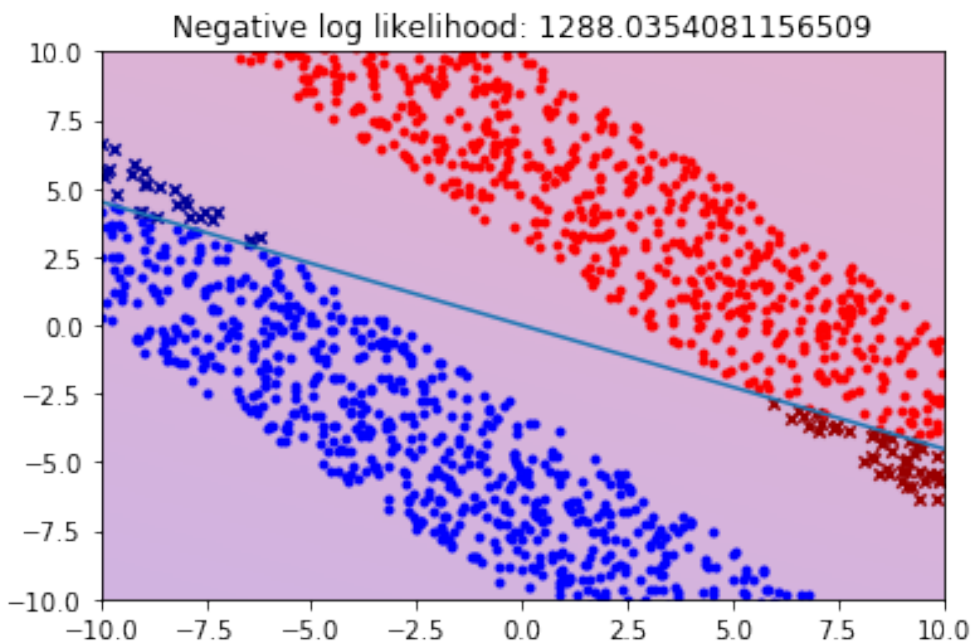


Figure 10: Regularization of  $x_2$ ,  $\lambda = 10^5$

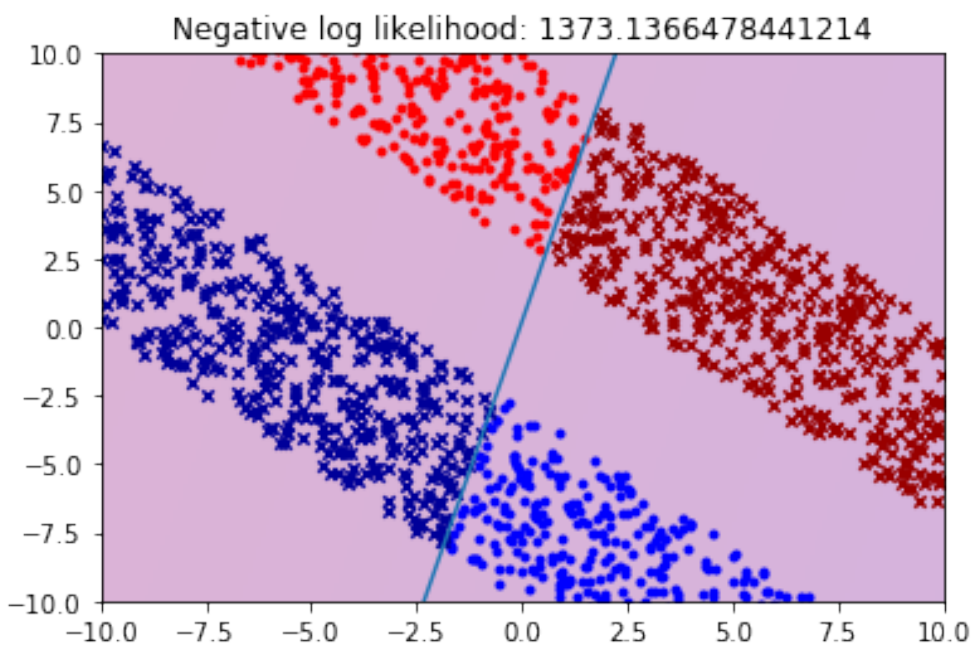


Figure 11: Regularization of  $x_2$ ,  $\lambda = 10^6$

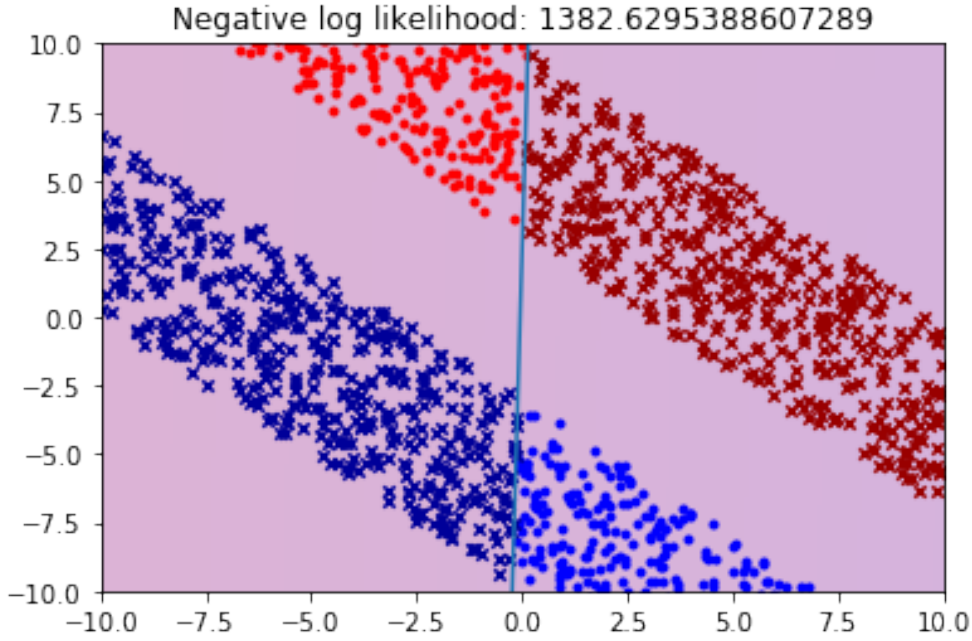


Figure 12: Regularization of  $x_2$ ,  $\lambda = 10^7$

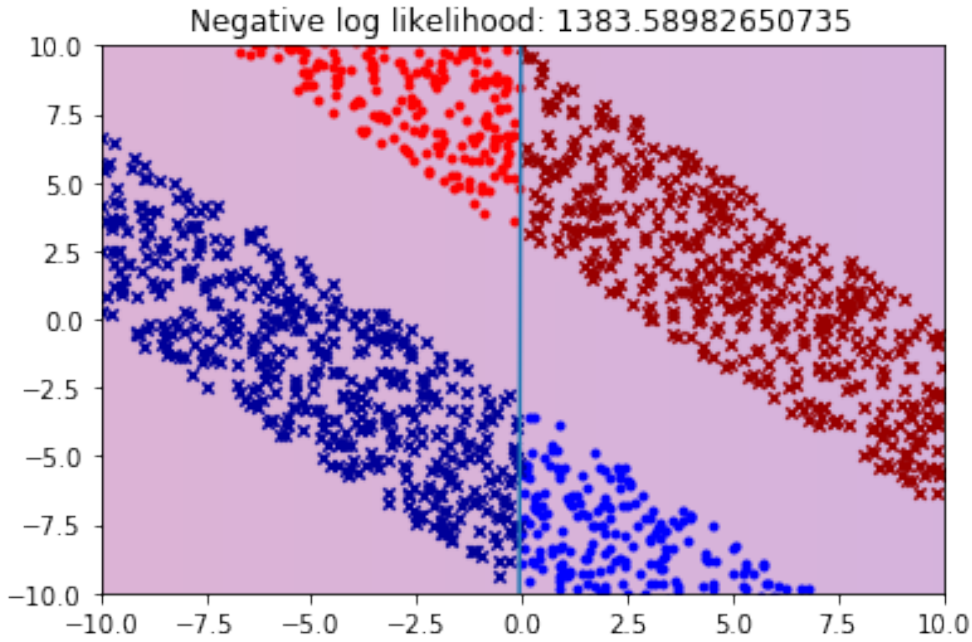


Figure 13: Regularization of  $x_2$ ,  $\lambda = 10^8$

Similar to the regularization of  $x_1$ , the separation line becomes more and more vertical as  $\lambda$  increases.

# HW2 Problem 4

February 29, 2020

## 1 Hello and welcome to decision trees!

Decision trees are often pretty effective 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/mitchdectrees.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!

```
[1]: import pandas as pd
import scipy.io
import numpy as np
from sklearn.model_selection import train_test_split
```

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

### 1.1 Import data

```
[2]: # 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"""
    # TODO

    data = pd.read_csv('data_banknote_authentication.txt', header = None)
    data = np.array(data)
    if shuffle == True:
        data = np.random.shuffle(data)

    X = data[:, :-1]
    y = data[:, -1]
```

```

y = y.astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
↳random_state=1)
print("data imported")

return X_train, y_train, X_test, y_test

```

```

[3]: 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
↳dataframe 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

```

```

[4]: 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, metric = "entropy"):
        """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, metric = metric)

def _buildTree(self, node, metric = "entropy"):

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

    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
    n = node.labels.shape[0]
    gain = 0
    fea = 0
    for k in range(node.data.shape[1]):
        node.feature = k
        for i in range(n):
            new_gain = self.getInfoGain(node, i, metric = metric)
            if new_gain > gain:
                gain = new_gain
                node.threshold_index = i
                fea = k
    node.feature = fea
    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
        n = len(node.labels)
        nl = split_index + 1
        nr = n - nl
        index = node.data[:, node.feature].argsort()
        qm = self.get_uncertainty(node.labels, metric = metric)
        ql = self.get_uncertainty(node.labels[index[:split_index]], metric = ↪
        ↪ metric)
        qr = self.get_uncertainty(node.labels[index[split_index:]], metric = ↪
        ↪ metric)
        return qm - (nl / n * ql) - (nr / n * qr)

    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
        # TODO YOUR CODE HERE
        else:
            p = np.bincount(labels) / labels.shape[0]
            if 0 not in p:
                if metric == "entropy":
                    return - sum(p * np.log(p))
                elif metric == "gini":
                    return sum(p * (1-p))
                else:
                    return 1
            else:
                return 1

    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]

    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

```

## 1.2 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?

```
[5]: X_train, y_train, X_test, y_test = import_data(split=0.8)
```

```

def getfeature(node, features):
    if node != None:
        if node.feature != None :
            features.append(node.feature)
            getfeature(node.left, features)
            getfeature(node.right, features)

# Use Entropy

```

```

for i in range(1, 7):
    print("For k=", i, ":")
    tree = DecisionTree(K=i, verbose=False)
    tree.buildTree(X_train, y_train)
    tree.homework_evaluate(X_train, y_train, X_test, y_test)
    features = []
    getfeature(tree.root, features)
    print("Largest information gain feature is:", np.argmax(np.
↪bincount(features)))
    print("Least information gain feature is:", np.argmin(np.
↪bincount(features)), "\n")

```

data imported

For k= 1 :

The decision tree is 91 percent accurate on 1097 training data

The decision tree is 88 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 0

For k= 2 :

The decision tree is 94 percent accurate on 1097 training data

The decision tree is 92 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 1

For k= 3 :

The decision tree is 96 percent accurate on 1097 training data

The decision tree is 95 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 1

For k= 4 :

The decision tree is 99 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

For k= 5 :

The decision tree is 99 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

For k= 6 :

The decision tree is 99 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

```
[6]: # Use Gini
for i in range(1, 7):
    print("For k=", i, ":")
    tree = DecisionTree(K=i, verbose=False)
    tree.buildTree(X_train, y_train, metric = "gini")
    tree.homework_evaluate(X_train, y_train, X_test, y_test)
    features = []
    getfeature(tree.root, features)
    print("Largest information gain feature is:", np.argmax(np.
    ↳bincount(features)))
    print("Least information gain feature is:", np.argmin(np.
    ↳bincount(features)), "\n")
```

For k= 1 :

The decision tree is 89 percent accurate on 1097 training data

The decision tree is 89 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 1

For k= 2 :

The decision tree is 96 percent accurate on 1097 training data

The decision tree is 95 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 1

For k= 3 :

The decision tree is 98 percent accurate on 1097 training data

The decision tree is 97 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

For k= 4 :

The decision tree is 98 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

For k= 5 :

The decision tree is 99 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data

Largest information gain feature is: 0

Least information gain feature is: 3

For k= 6 :

The decision tree is 99 percent accurate on 1097 training data

The decision tree is 98 percent accurate on 275 test data  
Largest information gain feature is: 0  
Least information gain feature is: 3

To get to know which feature has the largest information gain and which has the least, I collect the selected feature column for each node and try to find out which feature are selected most. As we can see from the above result, the first feature has the largest information gain, and the fourth feature has the least. Using Gini and entropy does not have a significant difference regarding the training and testing accuracy since the way of calculating uncertainty is quite similar.

## 2

---

### 3 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.

#### 3.0.1 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/>. 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.

```
[ ]: 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}.')
```

**!! 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:

```
[ ]: 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')
```

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.

```
[ ]: 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
```

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.

```
[ ]: 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.

```
[ ]: 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.

```
[ ]: 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.

```
[ ]: 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.

```
[ ]: 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:

```
[ ]: 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!")
```

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.

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

encoded_samples_df
```

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...

```
[ ]: sample_list = list(encoded_samples_df.to_numpy())
```

```
[ ]: 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
```

```
[ ]: 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!

### 3.0.2 Step 2: Decision Tree

We'll run the decision tree process here again.

```
[ ]: 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!

```
[ ]: 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)
```