

Machine Learning - HW2

Kai Liao

September 2020

1 Logistic Regression

Lemma. Let matrix A be symmetric, there exists $t < \infty$ such that $tI + A$ is positive semi-definite.

Proof. Suppose A is a symmetric matrix. There exists an orthogonal matrix P such that $P'AP = C$ is a diagonal matrix and $P'P = I$. Therefore, $P'(tI + A)P = tI + C$ and $tI + C$ is a diagonal matrix. Choose $t = \max_i \{\text{diag}(-C, i)\} + 1$,¹ then all diagonal elements are positive. Therefore, $tI + A$ must be positive semi-definite. \square

Define

$$\sigma(z) = \frac{1}{1 + \exp(-z)} \text{ where } z = f_{\mathbf{w}}(\mathbf{x}_k) \quad (1)$$

Therefore,

$$\partial L_D(\mathbf{w}) = \sum_{i=1}^m -(y_i \log \sigma(z_i) + (1 - y_i) \log (1 - \sigma(z_i))) \quad (2)$$

Notice that,

$$\frac{d\sigma(z)}{dz} = \exp(-z)(1 + \exp(-z))^{-2} = \frac{1}{1 + \exp(-z)} \frac{\exp(-z)}{1 + \exp(-z)} = \sigma(z)(1 - \sigma(z)) \quad (3)$$

First, I compute the first order derivative of the loss function.

$$\begin{aligned} \frac{\partial L_D(\mathbf{w})}{\partial \mathbf{w}} &= - \sum_{k=1}^n \left[y_k \times \frac{1}{\sigma(z)} \times \sigma(z)(1 - \sigma(z)) \times \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}} + (1 - y_k) \times \frac{1}{1 - \sigma(z)} \times (-\sigma(z))(1 - \sigma(z)) \times \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}} \right] \\ &= \sum_{k=1}^n \left[(\sigma(z) - y_k) \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}} \right] \end{aligned} \quad (4)$$

Then, I compute the second order derivatives (Hessian matrix).

$$\begin{aligned} \frac{\partial^2 L_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} &= \sum_{k=1}^n \left[(\sigma(z) - y_k) \frac{\partial^2 f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i \partial \mathbf{w}_j} + \frac{\partial(\sigma(z) - y_k)}{\partial z} \frac{\partial z}{\partial \mathbf{w}_j} \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i} \right] \\ &= \sum_{k=1}^n \left[(\sigma(z) - y_k) \frac{\partial^2 f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i \partial \mathbf{w}_j} + \sigma(z)(1 - \sigma(z)) \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i} \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_j} \right] \end{aligned} \quad (5)$$

The Hessian matrix H for the L_2 regularized new objective is:

$$\frac{\partial^2 \tilde{L}_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} = \frac{\partial^2 L_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} + \alpha \text{ if } i = j \quad (6)$$

¹Define $\text{diag}(A, i)$ = the diagonal element in the i row

$$\frac{\partial^2 \tilde{L}_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} = \frac{\partial^2 L_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \text{ if } i \neq j \quad (7)$$

The generate a sufficient large α , I have the following two methods:

Method 1:

This Hessian matrix H can be written as

$$H = H^o + \alpha I \quad (8)$$

where H' is,

$$\frac{\partial^2 \tilde{L}_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} = \frac{\partial^2 L_D(\mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \quad (9)$$

By the *lemma*, there exists α such that H is positive semi-definite. To find a value of α , we can follow the process in the proof the lemma. Firstly, I diagonalize the matrix H to $P'H^oP = C$. Then, let $\alpha = \max_i \{\text{diag}(-C, i)\} + 1$. By the lemma, $H = H^o + \alpha I$ is positive semi-definite.

Method 2:

The Hessian matrix H can also be written as

$$H = H^1 + H^2 + \alpha I \quad (10)$$

where

$$H^1[i, j] = \sum_{k=1}^n \left[(\sigma(z) - y_k) \frac{\partial^2 f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \right] \quad (11)$$

$$H^2[i, j] = \sum_{k=1}^n \left[\sigma(z)(1 - \sigma(z)) \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_i} \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial \mathbf{w}_j} \right] \quad (12)$$

By triangle inequality,

$$\|H^1\|_2 \leq u \quad (13)$$

and

$$\|H^2\|_2 \leq Ng \quad (14)$$

Therefore the min eigenvalue for $H^1 + H^2$ is greater than $-u - Ng$. So we can set $\alpha = u + Ng$, the H will be positive definite.

2 SVM with a Squared Loss

Proof. The optimization problem is

$$\min L(w, b, e) = \frac{1}{2} w^\top w + \frac{1}{2} C \sum_{i=1}^n e_i^2, \quad \text{s.t.} \quad y_i - w^\top x_i = e_i \quad (15)$$

We can write the Lagrangian as

$$\mathcal{L} = \frac{1}{2} w^\top w + \frac{1}{2} C \sum_{i=1}^n e_i^2 + \sum_{i=1}^n \lambda_i (y_i - w^\top x_i - e_i) \quad (16)$$

The KKT conditions are

$$\begin{aligned} (1) \quad & \frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^n \lambda_i x_i = 0 \\ (2) \quad & \frac{\partial \mathcal{L}}{\partial e_i} = C e_i - \lambda_i = 0 \\ (3) \quad & \frac{\partial \mathcal{L}}{\partial \lambda_i} = y_i - w^\top x_i - e_i = 0 \end{aligned} \quad (17)$$

We can solve for w as follow

$$\begin{aligned}
w &= C \sum_{i=1}^n e_i x_i \\
&= C \sum_{i=1}^n (y_i - w^T x_i) x_i \\
&= -C \sum_{i=1}^n w^T x_i x_i + C \sum_{i=1}^n y_i x_i \\
&= -C \sum_{i=1}^n x_i x_i^T w + C \sum_{i=1}^n y_i x_i
\end{aligned} \tag{18}$$

Therefore,

$$w + C \sum_{i=1}^n x_i x_i^T w = C \sum_{i=1}^n y_i x_i \tag{19}$$

It is equivalent to

$$w = \left(\frac{1}{C} I + \sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n y_i x_i \tag{20}$$

□

3 Linear SVM

Please find the code and answers after section 4

4 AdaBoost

4.1

Proof. Firstly, I would like to quickly repeat the proof of the theorem we have shown in class.

Theorem *If the weak learning assumption holds, AdaBoost's misclassification error decays exponentially fast:*

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[y_i \neq H(x_i)]} \leq e^{-2\gamma_W^2 L A^T} \tag{21}$$

Proof of the Theorem.

$$\begin{aligned}
R^{\text{train}}(\lambda_{t+1}) &= R^{\text{train}}(\lambda_t + \alpha_t \mathbf{e}_{j_t}) = \frac{1}{n} \sum_{i=1}^n e^{-[\mathbf{M}(\lambda_t + \alpha_t \mathbf{e}_{j_t})]_i} = \frac{1}{n} \sum_{i=1}^n e^{-(\mathbf{M}\lambda_t)_i - \alpha_t M_{ij_t}} \\
&= e^{-\alpha_t} \frac{1}{n} \sum_{i: M_{ij_t}=1} e^{-(\mathbf{M}\lambda_t)_i} + e^{\alpha_t} \frac{1}{n} \sum_{i: M_{ij_t}=-1} e^{-(\mathbf{M}\lambda_t)_i}
\end{aligned} \tag{22}$$

Define $Z_t = \sum_{i=1}^n e^{-(\mathbf{M}\lambda_t)_i}$, we have

$$d_{t,i} = e^{-(\mathbf{M}\lambda_t)_i} / Z_t \tag{23}$$

and

$$\frac{Z_t}{n} d_+ = \frac{Z_t}{n} \sum_{i: M_{ij_t}=1} d_{t,i} = \frac{1}{n} \sum_{i: M_{ij_t}=1} e^{-(\mathbf{M}\lambda_t)_i} \tag{24}$$

and similarly

$$\frac{Z_t}{n} d_- = \frac{1}{n} \sum_{i: M_{ij_t}=-1} e^{-(\mathbf{M}\lambda_t)_i} \tag{25}$$

Therefore,

$$\begin{aligned}
R^{\text{train}}(\boldsymbol{\lambda}_{t+1}) &= e^{-\alpha} \frac{Z_t}{n} d_+ + e^{\alpha} \frac{Z_t}{n} d_- \\
&= R^{\text{train}}(\boldsymbol{\lambda}_t) [e^{-\alpha} d_+ + e^{\alpha} d_-] \\
&= R^{\text{train}}(\boldsymbol{\lambda}_t) [e^{-\alpha} (1 - d_-) + e^{\alpha} d_-] \\
&= R^{\text{train}}(\boldsymbol{\lambda}_t) \left[\left(\frac{d_-}{1 - d_-} \right)^{1/2} (1 - d_-) + \left(\frac{1 - d_-}{d_-} \right)^{1/2} d_- \right] \\
&= R^{\text{train}}(\boldsymbol{\lambda}_t) 2 [d_- (1 - d_-)]^{1/2} \\
&= R^{\text{train}}(\boldsymbol{\lambda}_t) 2 [\epsilon_t (1 - \epsilon_t)]^{1/2}
\end{aligned} \tag{26}$$

Then,

$$\begin{aligned}
R^{\text{train}}(\boldsymbol{\lambda}_T) &= \prod_{t=1}^T 2 \sqrt{\epsilon_t (1 - \epsilon_t)} \\
&= \prod_{t=1}^T 2 \sqrt{\left(\frac{1}{2} - \gamma_t \right) \left(\frac{1}{2} + \gamma_t \right)} \\
&= \prod_t \sqrt{1 - 4\gamma_t^2} \\
&\leq \prod_t \sqrt{e^{-4\gamma_t^2}} \\
&= \prod_t e^{-2\gamma_t^2} \\
&= e^{-2 \sum_{t=1}^T \gamma_t^2}
\end{aligned} \tag{27}$$

Finally,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[y_i \neq H(x_i)]} \leq R^{\text{train}}(\boldsymbol{\lambda}_T) \leq e^{-2 \sum_{t=1}^T \gamma_t^2} \leq e^{-2\gamma_{WLA}^2 T} \tag{28}$$

By the theorem,

$$0 \leq \sum_{i=1}^n \mathbf{1}_{[y_i \neq H(x_i)]} \leq n e^{-2\gamma_{WLA}^2 T} \tag{29}$$

Since

$$\lim_{T \rightarrow \infty} n e^{-2\gamma_{WLA}^2 T} = 0 \tag{30}$$

by sandwich lemma,

$$\lim_{T \rightarrow \infty} \sum_{i=1}^n \mathbf{1}_{[y_i \neq H(x_i)]} = 0 \tag{31}$$

4.2

$$R^{\text{train}}(\boldsymbol{\lambda}) = \sum_{i=1}^n w_i e^{-(M\boldsymbol{\lambda})_i} \tag{32}$$

Choosing the direction j :

$$\begin{aligned}
j_t &\in \operatorname{argmax}_j \left[- \frac{\partial R^{\text{train}}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_j)}{\partial \alpha} \Big|_{\alpha=0} \right] \\
&= \operatorname{argmax}_j \left[- \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n w_i e^{-(\mathbf{M}(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_j))_i} \right] \Big|_{\alpha=0} \right] \\
&= \operatorname{argmax}_j \left[- \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n w_i e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha (\mathbf{M}_j)_i} \right] \Big|_{\alpha=0} \right] \\
&= \operatorname{argmax}_j \left[- \frac{\partial}{\partial \alpha} \left[\frac{1}{n} \sum_{i=1}^n w_i e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha M_{ij}} \right] \Big|_{\alpha=0} \right] \\
&= \operatorname{argmax}_j \left[\frac{1}{n} \sum_{i=1}^n w_i M_{ij} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} \right]
\end{aligned} \tag{33}$$

Define $Z_t = \sum_{i=1}^n e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i}$, we have

$$d_{t,i} = e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} / Z_t \tag{34}$$

$$j_t \in \operatorname{argmax}_j \sum_{i=1}^n w_i M_{ij} d_{t,i} \tag{35}$$

Choosing the step α :

$$\begin{aligned}
0 &= \frac{\partial R(\boldsymbol{\lambda}_t + \alpha \mathbf{e}_{j_t})}{\partial \alpha} \Big|_{\alpha_t} \\
&= -\frac{1}{n} \sum_{i=1}^n w_i M_{ij_t} e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i - \alpha_t M_{ij_t}} \\
&= -\frac{1}{n} \sum_{i: M_{ij_t}=1} w_i e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} e^{-\alpha_t} - \frac{1}{n} \sum_{i: M_{ij_t}=-1} -w_i e^{-(\mathbf{M}\boldsymbol{\lambda}_t)_i} e^{\alpha_t}
\end{aligned} \tag{36}$$

Define $d_+ = \sum_{i: M_{ij_t}=1} w_i d_{t,i}$ and $d_- = \sum_{i: M_{ij_t}=-1} w_i d_{t,i}$

$$\begin{aligned}
0 &= \sum_{i: M_{ij_t}=1} w_i d_{t,i} e^{-\alpha_t} - \sum_{i: M_{ij_t}=-1} w_i d_{t,i} e^{\alpha_t} \\
&= d_+ e^{-\alpha_t} - d_- e^{\alpha_t}
\end{aligned} \tag{37}$$

Therefore,

$$\alpha_t = \frac{1}{2} \ln \frac{d_+}{d_-} = \frac{1}{2} \ln \frac{1 - d_-}{d_-} \tag{38}$$

So the coordinate descent algorithm is:

```

 $d_{1,i} = 1/n$  for  $i = 1 \dots n$ 
 $\lambda_1 = 0$ 
loop  $t = 1 \dots T$ 
 $j_t \in \operatorname{argmax}_j \sum_{i=1}^n w_i M_{ij} d_{t,i}$ 
 $d_- = \sum_{M_{ij_t}=-1} w_i d_{t,i}$ 
 $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - d_-}{d_-} \right)$ 
 $\lambda_{t+1} = \lambda_t + \alpha_t e_{j_t}$ 
 $d_{t+1,i} = e^{-(M\lambda_{t+1})_i} / Z_{t+1}$  for each  $i$ , where  $Z_{t+1} = \sum_{i=1}^n w_i e^{-(M\lambda_{t+1})_i}$ 
end

```

This is also adaboost:

$$\begin{aligned}
j_t &\in \operatorname{argmin}_j \sum_i w_i d_{t,i} \mathbf{1}_{[h_j(x_i) \neq y_i]} \\
&= \operatorname{argmax}_j \left[- \sum_{i: M_{ij}=-1} w_i d_{t,i} \right] \\
&= \operatorname{argmax}_j \left[\left[\sum_{i: M_{ij}=1} w_i d_{t,i} + \sum_{i: M_{ij}=-1} w_i d_{t,i} \right] - 2 \sum_{i: M_{ij}=-1} w_i d_{t,i} \right] \\
&= \operatorname{argmax}_j \sum_{i: M_{ij}=1} w_i d_{t,i} - \sum_{i: M_{ij}=-1} w_i d_{t,i} \\
&= \operatorname{argmax}_j \sum_{i=1}^n w_i M_{ij} d_{t,i}
\end{aligned} \tag{39}$$

$$\epsilon_t = \sum_i w_i d_{t,i} \mathbf{1}_{[h_{j_t}(x_i) \neq y_i]} = \sum_{i: h_{j_t}(x_i) \neq y_i} w_i d_{t,i} = \sum_{i: M_{ij_t}=-1} w_i d_{t,i} = d_- \tag{40}$$

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t} = \frac{1}{2} \ln \frac{1 - d_-}{d_-} \tag{41}$$

So AdaBoost minimizes the exponential loss by coordinate descent.

4.3

Please find the code and answer behind

LINEAR SVM

Due Date : 9/28 Monday 10:15 PM EST

```
In [1]: pip install libsvm
```

Requirement already satisfied: libsvm in c:\users\kaike\anaconda3\lib\site-packages (3.23.0.4)

Note: you may need to restart the kernel to use updated packages.

```
In [2]: import numpy as np
import matplotlib.pyplot as plt
import scipy.io as io
import libsvm
import math
from sklearn.svm import SVC
from libsvm.svmutil import *

%matplotlib inline
```

3.1 Linear Support Vector Machine on toy data

3.1.1

Generate a training set of size 100 with 2D features (X) drawn at random as follows:

- $X_{\text{neg}} \sim \mathcal{N}([-5, -5], 5 \cdot I_2)$ and correspond to negative labels (-1)
- $X_{\text{pos}} \sim \mathcal{N}([5, 5], 5 \cdot I_2)$ and correspond to positive labels (+1)

Accordingly, $X = [X_{\text{neg}}, X_{\text{pos}}]$ is a 100×2 array, Y is a 100×1 array of values $\in \{-1, 1\}$.

Draw a scatter plot of the full training dataset with the points colored according to their labels.

```
In [3]: import random

# Generate binary class dataset
np.random.seed(0)

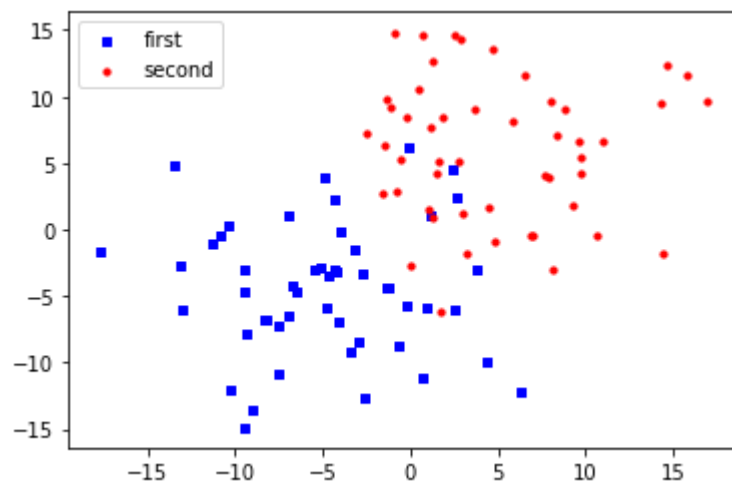
n_samples = 100
center_1 = [-5, -5]
center_2 = [5, 5]

# random number of X is Xneg
random_seperation = 50
#int(n_samples*random.uniform(0, 1))

# Generate Data:
Y = []
Xneg = np.random.normal(-5,5,size = (int(random_seperation),2))
for i in range(int(random_seperation)):
    Y.append(-1)
Xpos = np.random.normal(5,5,size = (int(100 - random_seperation),2))
for i in range(int(100 - random_seperation)):
    Y.append(1)
X = np.concatenate((Xneg, Xpos))

# Scatter plot:
X_neg_plot_1 = []
X_neg_plot_2 = []
X_pos_plot_1 = []
X_pos_plot_2 = []
for i in range(len(X)):
    if Y[i] == -1 :
        X_neg_plot_1.append(X[i][0])
        X_neg_plot_2.append(X[i][1])
    elif Y[i] == 1 :
        X_pos_plot_1.append(X[i][0])
        X_pos_plot_2.append(X[i][1])

plt.scatter(X_neg_plot_1, X_neg_plot_2, s=10, c='b', marker="s", label='first'
)
plt.scatter(X_pos_plot_1, X_pos_plot_2, s=10, c='r', marker="o", label='second')
plt.legend(loc='upper left');
plt.show()
```

3.1.2

Train a linear support vector machine on the data with $C = 1$ and draw the decision boundary line that separates o and x. Mark the support vectors separately (ex.circle around the point).

Note: You can use the `libsvm.svmutil` functions with the `kernel_type` set to 0, indicating a linear kernel and `svm_type` set to 0 indicating C-SVC. Also note that the `support_vector` coefficients returned by the LIBSVM model are the dual coefficients.

```

In [4]: # Define the SVM problem
problem = svm_problem(Y,X)
# Define the hyperparameters
param = svm_parameter('-t 0 -s 0')
# Train the model
model = svm_train(problem, param)

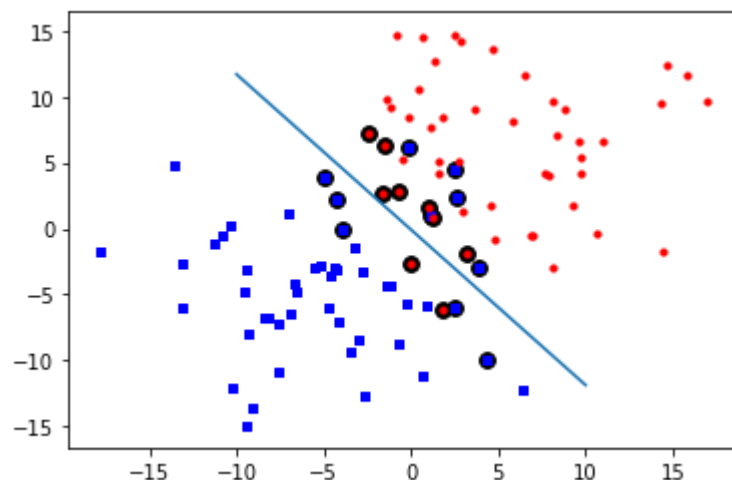
# Compute the slope and intercept of the separating line/hyperplane with the
# use of the support vectors
# and other information from the LIBSVM model.
sv = model.get_sv()
sv_coef = model.get_sv_coef()
w = np.matmul(np.array(X)[np.array(model.get_sv_indices()) - 1].T, sv_coef)
b = -model.rho.contents.value
if model.get_labels()[1] == -1:
    w = -w
    b = -b

# Draw the scatter plot, the decision boundary line, and mark the support vect
ors.

# Plot the line
x_bdd = [-10,10]
y_bdd = []
y_bdd.append(-(x_bdd[0] * w[0] + b) / w[1])
y_bdd.append(-(x_bdd[1] * w[0] + b) / w[1])
plt.plot(x_bdd, y_bdd)
# Plot supporting vector
for i in model.get_sv_indices():
    plt.scatter(X[i - 1][0], X[i - 1][1], color='black', s=60)
# Plot all training point
plt.scatter(X_neg_plot_1, X_neg_plot_2, s=10, c='b', marker="s")
plt.scatter(X_pos_plot_1, X_pos_plot_2, s=10, c='r', marker="o")

plt.show()

```



3.1.3

Draw a line that separates the data for 8 different C ($10^{-5} \sim 10^7$). Plot the number of support vectors vs. C (plot x-axis on a log scale). How does the number of support vectors change as C increases and why does it change like that?

Note: You might prefer to use the command-line style of `svm_parameter` initialization such as: `svm_parameter('-s 0 -t 0')` to indicate a linear kernel and C-SVC as the SVM type.

```

In [5]: C_range = [10**-5, 10**-3, 1, 10, 100, 10**3, 10**5, 10**7]
num_sv = []

# Loop over a similar setup to that in the previous code block.
for C in C_range:
    param = svm_parameter('-s 0 -t 0 -c ' + str(C))
    model = svm_train(problem, param)
    num_sv.append(model.get_nr_sv())
    sv = model.get_SV()
    sv_coef = model.get_sv_coef()

    # Calculate w and b
    w = np.matmul(np.array(X)[np.array(model.get_sv_indices()) - 1].T, sv_coef)

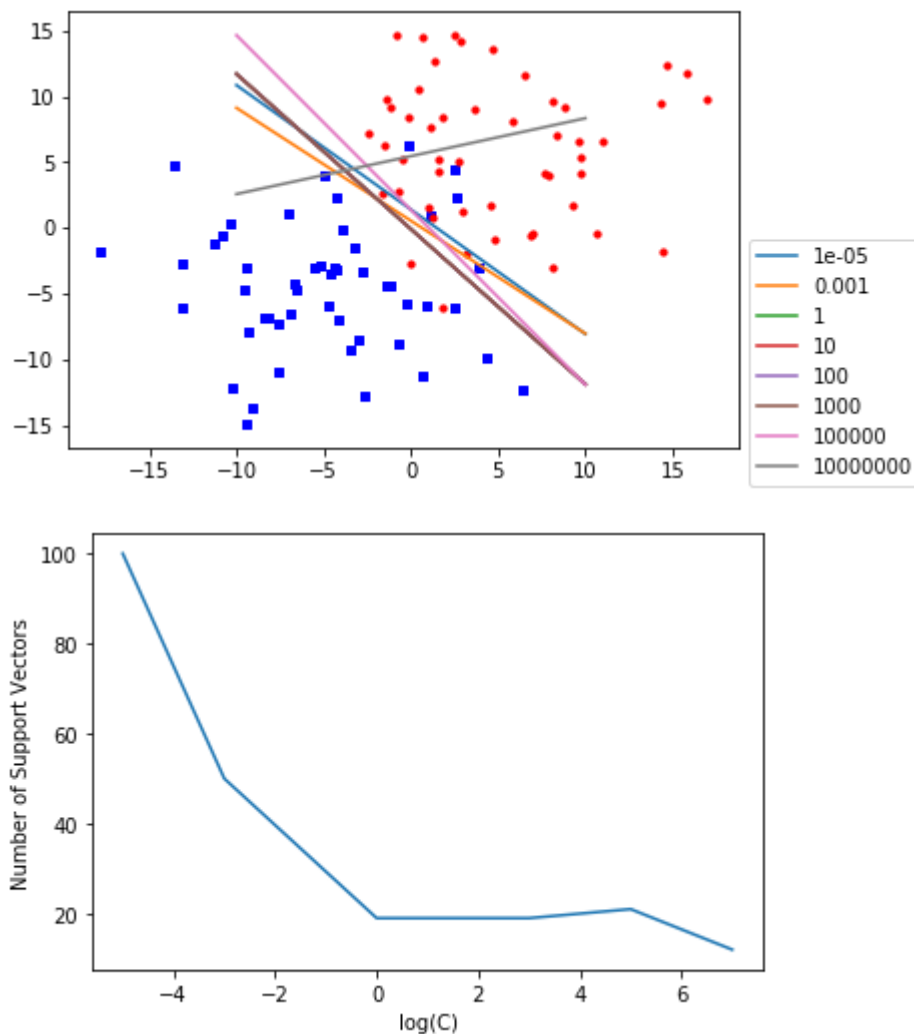
    b = -model.rho.contents.value
    if model.get_labels()[1] == -1:
        w = -w
        b = -b

    # Draw the scatter plot with multiple decision lines on top (one for each
    value of C)
    x_bdd = [-10,10]
    y_bdd = []
    y_bdd.append(-(x_bdd[0] * w[0] + b) / w[1])
    y_bdd.append(-(x_bdd[1] * w[0] + b) / w[1])
    plt.plot(x_bdd, y_bdd, label = str(C))

# Plot all training point
plt.scatter(X_neg_plot_1, X_neg_plot_2, s=10, c='b', marker="s")
plt.scatter(X_pos_plot_1, X_pos_plot_2, s=10, c='r', marker="o")
plt.legend(loc='upper left', bbox_to_anchor=(1, 0.5))
plt.show()

# Draw the num_sv vs. C plot
log_C = []
for C in C_range:
    log_C.append(math.log(C,10))
plt.xlabel('log(C)')
plt.ylabel('Number of Support Vectors')
plt.plot(log_C, num_sv)
plt.show()

```



How does the number of support vectors change as C increases and why does it change like that?

C governs the importance of avoiding misclassifying each training sample. A high C suggests that the SVM penalizes the misclassification very badly. Support vectors are training sample which the output of SVM has a value between $[0, 1]$. If I increase C , the margin will be narrower because the main goal is to reduce misclassification. Therefore, less vectors will be support vectors.

3.1.4

Now try rescaling the data to the $[0, 1]$ range and repeat the steps of the previous question (3.1.3) and over the same range of C values. Are the decision boundaries different from those in the previous question? What does this imply about (a) the geometric margin and (b) the relative effect of each feature on the predictions of the trained model ?

Solution below:

SVM tries to maximize the distance between the separating plane and the support vectors. If one feature (i.e. one dimension in this space) has very large values, it will dominate the other features when calculating the distance. If you rescale all features (e.g. to [0, 1]), they all have the same influence on the distance metric.

In this case, the decision boundary will shift and shirk the same way with the changes of features. But the change is propotional, because the we rescale two features by very similar proportion ($\max(\text{feature1}) - \min(\text{feature1})$ is approx equal to $\max(\text{feature2}) - \min(\text{feature2})$), the geometric margin will decrease proportionally and the relative effect of each feature remains pretty much unchanged. However if we rescale each feature differently, the relative effect of each feature will be changed.

```
In [6]: import sklearn
        from sklearn import preprocessing
        min_max_scaler = preprocessing.MinMaxScaler()

        # Single line below:
        X_train_minmax = min_max_scaler.fit_transform(X)
```

```

In [7]: C_range = [10**-5, 10**-3, 1, 10, 100, 10**3, 10**5, 10**7]
num_sv = []

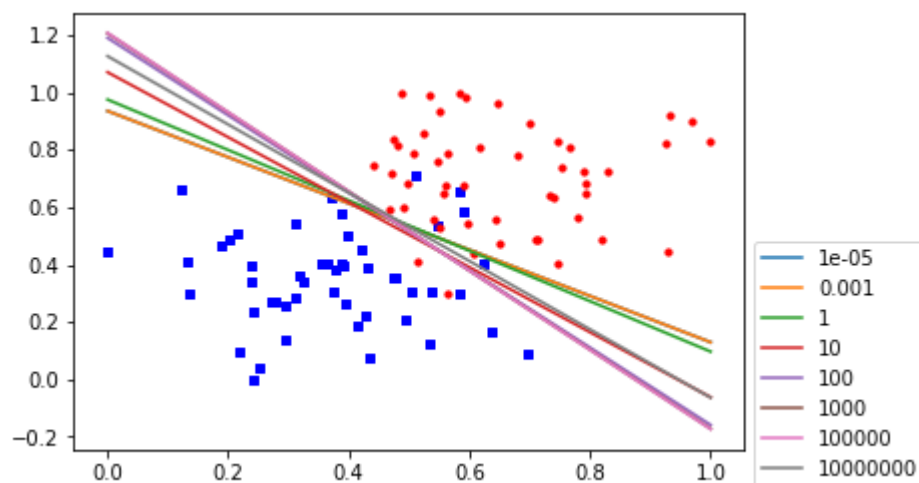
# Repeat the loop from 3.1.3
problem = svm_problem(Y,X_train_minmax)
for C in C_range:
    param = svm_parameter('-s 0 -t 0 -c ' + str(C))
    model = svm_train(problem, param)
    num_sv.append(model.get_nr_sv())
    sv = model.get_SV()
    sv_coef = model.get_sv_coef()

    # Calculate w and b
    w = np.matmul(np.array(X_train_minmax)[np.array(model.get_sv_indices()) -
1].T, sv_coef)
    b = -model.rho.contents.value
    if model.get_labels()[1] == -1:
        w = -w
        b = -b

    # Draw the scatter plot with multiple decision lines on top (one for each
value of C)
    x_bdd = [0,1]
    y_bdd = []
    y_bdd.append(-(x_bdd[0] * w[0] + b) / w[1])
    y_bdd.append(-(x_bdd[1] * w[0] + b) / w[1])
    plt.plot(x_bdd, y_bdd, label = str(C))

# Plot all training point
X_neg_plot_1_minmax = []
X_neg_plot_2_minmax = []
X_pos_plot_1_minmax = []
X_pos_plot_2_minmax = []
for i in range(len(X)):
    if Y[i] == -1 :
        X_neg_plot_1_minmax.append(X_train_minmax[i][0])
        X_neg_plot_2_minmax.append(X_train_minmax[i][1])
    elif Y[i] == 1 :
        X_pos_plot_1_minmax.append(X_train_minmax[i][0])
        X_pos_plot_2_minmax.append(X_train_minmax[i][1])
plt.scatter(X_neg_plot_1_minmax, X_neg_plot_2_minmax, s=10, c='b', marker="s")
plt.scatter(X_pos_plot_1_minmax, X_pos_plot_2_minmax, s=10, c='r', marker="o")
plt.legend(loc='upper left', bbox_to_anchor=(1, 0.5))
plt.show()

```



3.2 MNIST

Multiclass kernel SVM. In this problem, we'll use support vector machines to classify the MNIST data set of handwritten digits.

3.2.1

Load in the MNIST data using from the provided mnist-original.mat file on sakai. First split the data into training and testing by simply taking the first 60k points as training and the rest as testing. Then sample 500 data points for each of the 10 categories (for a total of 5000 training points) from the 60k training photos. These 5k points are now our training set. Finally, sample 500 data points for each of the 10 categories from the 10k testing photos. These 5k points are now our testing set.

Note: For data loading, you might want to use `scipy.io.loadmat`.


```

In [8]: import scipy
import scipy.io
from scipy.io import loadmat

np.random.seed(0)

minist_data = loadmat('mnist-original.mat')
X = minist_data['data'].transpose()
y = minist_data['label'].reshape(-1)

X_train_pool = X[:60000,:]
y_train_pool = y[:60000]
X_test_pool = X[60000:,:]
y_test_pool = y[60000:]

# get the index for sample selection
index = []

for label in set(y_train_pool):
    index_temp = []
    index_temp = np.argwhere(y_train_pool == label).reshape(-1)
    index += list(np.random.choice(index_temp, size = 500, replace = False))

X_train = np.take(X_train_pool, index, axis = 0)
y_train = np.take(y_train_pool, index, axis = 0)

index = []

for label in set(y_test_pool):
    index_temp = []
    index_temp = np.argwhere(y_test_pool == label).reshape(-1)
    index += list(np.random.choice(index_temp, size = 500, replace = False))

X_test = np.take(X_test_pool, index, axis = 0)
y_test = np.take(y_test_pool, index, axis = 0)

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

```

```

In [9]: np.unique(y_train, return_counts=True) #ensure each Label has 500 examples.

```

```

Out[9]: (array([0., 1., 2., 3., 4., 5., 6., 7., 8., 9.]),
array([500, 500, 500, 500, 500, 500, 500, 500, 500, 500], dtype=int64))

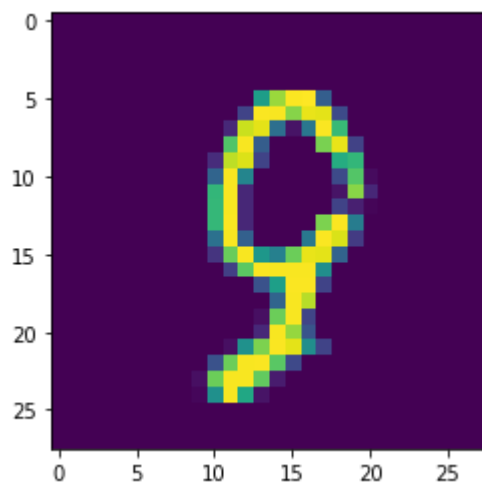
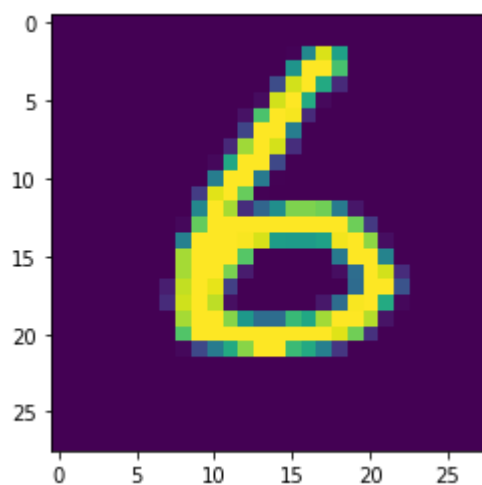
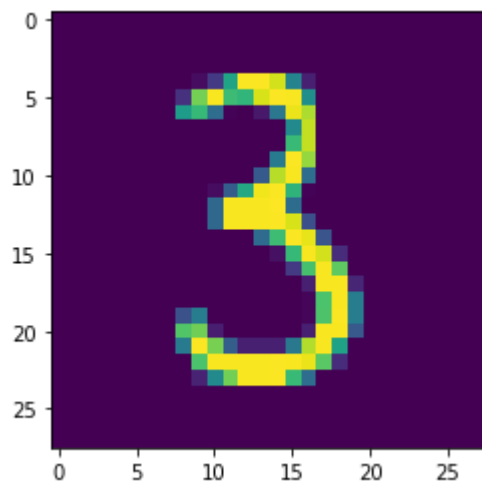
```

3.2.2

Draw 3 different digits using `pyplot.imshow()`.

```
In [10]: from random import randrange

digits = [3,6,9]
for digit in digits:
    index_temp = []
    index_temp = np.argwhere(y_train == digit).reshape(-1)
    plt.imshow(X_train[index_temp[randrange(500)],:].reshape(28,28))
    plt.show()
```



3.2.3

For each value $C = 10^{-12} \sim 10^{12}$ train a support vector machine with a linear kernel and compute its accuracy on the test set subsampled previously. Plot test accuracy and the number of support vectors (two separate plots) vs. C for $C = 10^{-12} \sim 10^{12}$ (plot 7 points or more with the x-axis on a log scale).

```
In [11]: C_range = []
for i in range(-12, 13):
    C_range.append(10**i)

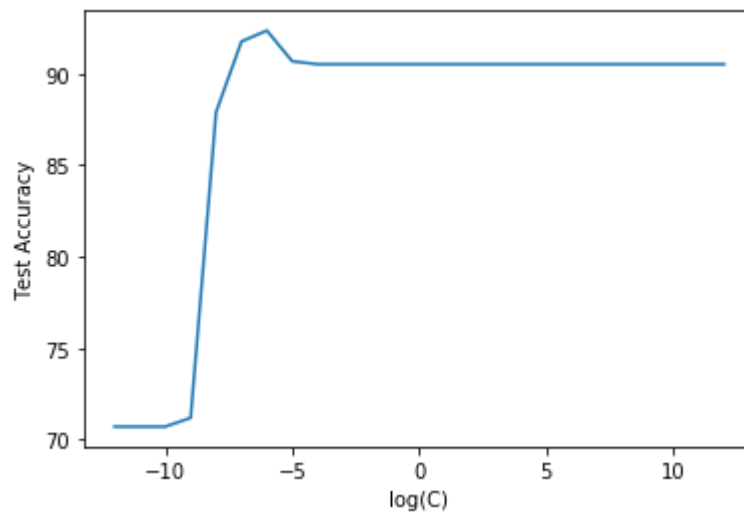
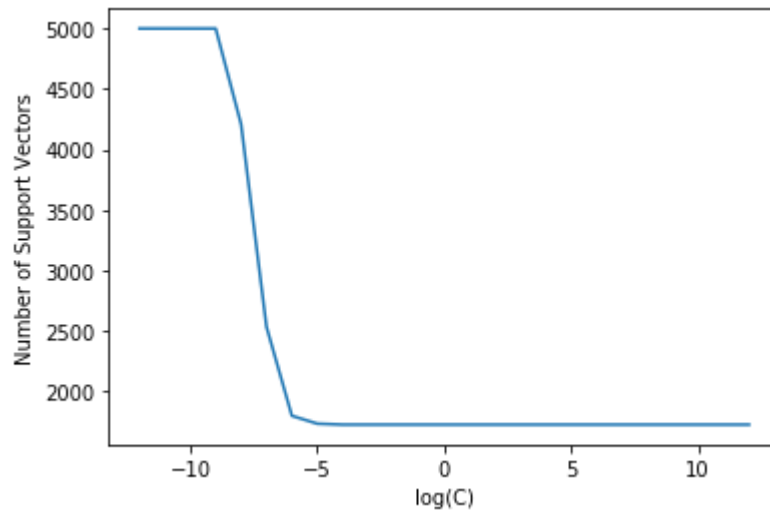
accuracy = []
num_sv = []
for C in C_range:
    # Define the SVM problem
    problem = svm_problem(y_train, X_train)
    # Define the hyperparameters
    param = svm_parameter('-s 0 -t 0 -c ' + str(C))
    # Train the model
    model = svm_train(problem, param)
    # make prediction
    p_label, p_acc, p_val = svm_predict(y_test, X_test, model)
    accuracy.append(p_acc[0])
    # number of sv
    num_sv.append(model.get_nr_sv())
```

```
Accuracy = 70.72% (3536/5000) (classification)
Accuracy = 70.72% (3536/5000) (classification)
Accuracy = 70.72% (3536/5000) (classification)
Accuracy = 71.2% (3560/5000) (classification)
Accuracy = 87.9% (4395/5000) (classification)
Accuracy = 91.76% (4588/5000) (classification)
Accuracy = 92.36% (4618/5000) (classification)
Accuracy = 90.68% (4534/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
Accuracy = 90.52% (4526/5000) (classification)
```

```
In [12]: log_C = []
for C in C_range:
    log_C.append(math.log(C,10))

plt.plot(log_C, num_sv)
plt.xlabel('log(C)')
plt.ylabel('Number of Support Vectors')
plt.show()

plt.plot(log_C, accuracy)
plt.xlabel('log(C)')
plt.ylabel('Test Accuracy')
plt.show()
```



In []:

Boosting a decision stump

The goal of this notebook is to implement your own boosting module.

- Go through an implementation of decision trees.
- Implement Adaboost ensembling.
- Use your implementation of Adaboost to train a boosted decision stump ensemble.
- Evaluate the effect of boosting (adding more decision stumps) on performance of the model.
- Explore the robustness of Adaboost to overfitting.

This file is adapted from course material by Carlos Guestrin and Emily Fox.

Let's get started!

Import some libraries

```
In [1]: ## please make sure that the packages are updated to the newest version.

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

Getting the data ready

We will be using a subset of the [LendingClub \(https://www.kaggle.com/wendykan/lending-club-loan-data\)](https://www.kaggle.com/wendykan/lending-club-loan-data) dataset.

```
In [2]: loans = pd.read_csv('loan_small.csv')
```

Recoding the target column

We re-assign the target to have +1 as a safe (good) loan, and -1 as a risky (bad) loan. In the next cell, the features are also briefly explained.

```
In [3]: features = ['grade',          # grade of the loan
                   'term',           # the term of the loan
                   'home_ownership', # home ownership status: own, mortgage or re
nt                                     nt
                   'emp_length',     # number of years of employment
                   ]

loans['safe_loans'] = loans['loan_status'].apply(lambda x : +1 if x=='Fully Pa
id' else -1)

## please update pandas to the newest version in order to execute the followin
g line
loans.drop(columns=['loan_status'], inplace=True)

target = 'safe_loans' # this variable will be used later
```

Transform categorical data into binary features

In this assignment, we will work with **binary decision trees**. Since all of our features are currently categorical features, we want to turn them into binary features using 1-hot encoding.

We can do so with the following code block:

```
In [4]: loans = pd.get_dummies(loans)
```

Let's see what the feature columns look like now:

```
In [5]: features = list(loans.columns)
features.remove('safe_loans') # Remove the response variable
features
```

```
Out[5]: ['term_ 36 months',
'term_ 60 months',
'grade_A',
'grade_B',
'grade_C',
'grade_D',
'grade_E',
'grade_F',
'grade_G',
'home_ownership_MORTGAGE',
'home_ownership_NONE',
'home_ownership_OTHER',
'home_ownership_OWN',
'home_ownership_RENT',
'emp_length_1 year',
'emp_length_10+ years',
'emp_length_2 years',
'emp_length_3 years',
'emp_length_4 years',
'emp_length_5 years',
'emp_length_6 years',
'emp_length_7 years',
'emp_length_8 years',
'emp_length_9 years',
'emp_length_< 1 year']
```

Train-test split

We split the data into training and test sets with 80% of the data in the training set and 20% of the data in the test set. We use `seed=1` so that everyone gets the same result.

```
In [6]: from sklearn.model_selection import train_test_split

train_data, test_data = train_test_split(loans, test_size = 0.2, random_state=
1)
```

Weighted decision trees

Since the data weights change as we build an AdaBoost model, we need to first code a decision tree that supports weighting of individual data points.

Weighted error definition

Consider a model with N data points with:

- Predictions $\hat{y}_1 \dots \hat{y}_n$
- Target $y_1 \dots y_n$
- Data point weights $\alpha_1 \dots \alpha_n$.

Then the **weighted error** is defined by:

$$E(\alpha, \hat{y}) = \frac{\sum_{i=1}^n \alpha_i \times 1[y_i \neq \hat{y}_i]}{\sum_{i=1}^n \alpha_i}$$

where $1[y_i \neq \hat{y}_i]$ is an indicator function that is set to 1 if $y_i \neq \hat{y}_i$.

Write a function to compute weight of mistakes

Write a function that calculates the weight of mistakes for making the "weighted-majority" predictions for a dataset. The function accepts two inputs:

- `labels_in_node` : Targets $y_1 \dots y_n$
- `data_weights` : Data point weights $\alpha_1 \dots \alpha_n$

We are interested in computing the (total) weight of mistakes, i.e.

$$WM(\alpha, \hat{y}) = \sum_{i=1}^n \alpha_i \times 1[y_i \neq \hat{y}_i].$$

This quantity is analogous to the number of mistakes, except that each mistake now carries different weight. It is related to the weighted error in the following way:

$$E(\alpha, \hat{y}) = \frac{WM(\alpha, \hat{y})}{\sum_{i=1}^n \alpha_i}$$

The function **intermediate_node_weighted_mistakes** should first compute two weights:

- WM_{-1} : weight of mistakes when all predictions are $\hat{y}_i = -1$ i.e $WM(\alpha, -\mathbf{1})$
- WM_{+1} : weight of mistakes when all predictions are $\hat{y}_i = +1$ i.e $WM(\alpha, +\mathbf{1})$

where $-\mathbf{1}$ and $+\mathbf{1}$ are vectors where all values are -1 and +1 respectively.

After computing WM_{-1} and WM_{+1} , the function **intermediate_node_weighted_mistakes** should return the lower of the two weights of mistakes, along with the class associated with that weight. We have provided a skeleton for you with `YOUR CODE HERE` to be filled in several places.


```
In [7]: def intermediate_node_weighted_mistakes(labels_in_node, data_weights):
# Sum the weights of all entries with label +1
total_weight_positive = sum(data_weights[labels_in_node == +1])

# Weight of mistakes for predicting all -1's is equal to the sum above
### YOUR CODE HERE
...
WM_neg = total_weight_positive

# Sum the weights of all entries with label -1
### YOUR CODE HERE
...
total_weight_negative = sum(data_weights[labels_in_node == -1])

# Weight of mistakes for predicting all +1's is equal to the sum above
### YOUR CODE HERE
...
WM_pos = total_weight_negative

# Return the tuple (weight, class_label) representing the lower of the two
weights
# class_label should be an integer of value +1 or -1.
# If the two weights are identical, return (weighted_mistakes_all_positiv
e,+1)
### YOUR CODE HERE
...
if WM_pos > WM_neg :
    weight = WM_neg
    flag = -1
elif WM_pos <= WM_neg:
    weight = WM_pos
    flag = 1
return (weight, flag)
```

Checkpoint: Test your **intermediate_node_weighted_mistakes** function, run the following cell:

```
In [8]: example_labels = pd.Series([-1, -1, 1, 1, 1])
example_data_weights = pd.Series([1., 2., .5, 1., 1.])
if intermediate_node_weighted_mistakes(example_labels, example_data_weights) =
= (2.5, -1):
    print('Test passed!')
else:
    print('Test failed... try again!')
    print(intermediate_node_weighted_mistakes(example_labels, example_data_w
eights))
```

Test passed!

Recall that the **classification error** is defined as follows:

$$\text{classification error} = \frac{\# \text{ mistakes}}{\# \text{ all data points}}$$

Function to pick best feature to split on

The next step is to pick the best feature to split on.

The **best_splitting_feature** function takes the data, the features, the targetm and the data weights as input and returns the best feature to split on.

Complete the following function.

```

In [9]: # If the data is identical in each feature, this function should return None

def best_splitting_feature(data, features, target, data_weights):

    # These variables will keep track of the best feature and the corresponding error
    best_feature = None
    best_error = float('+inf')
    num_points = float(len(data))

    # Loop through each feature to consider splitting on that feature
    for feature in features:

        # The left split will have all data points where the feature value is 0
        # The right split will have all data points where the feature value is 1

        left_split = data[data[feature] == 0]
        right_split = data[data[feature] == 1]

        # Apply the same filtering to data_weights to create left_data_weights, right_data_weights
        ## YOUR CODE HERE
        ...
        left_weight_split = data_weights[data[feature] == 0]
        right_weight_split = data_weights[data[feature] == 1]

        # Calculate the weight of mistakes for left and right sides
        ## YOUR CODE HERE
        ...
        WM_left, sign_left = intermediate_node_weighted_mistakes(left_split[target], left_weight_split)
        WM_right, sign_right = intermediate_node_weighted_mistakes(right_split[target], right_weight_split)

        # Compute weighted error by computing
        # ( [weight of mistakes (left)] + [weight of mistakes (right)] ) / [total weight of all data points]
        ## YOUR CODE HERE
        ...
        error = (WM_left + WM_right)/sum(data_weights)
        # If this is the best error we have found so far, store the feature and the error
        if error < best_error:
            best_feature = feature
            best_error = error

    # Return the best feature we found
    return best_feature

```

Checkpoint: Now, we have another checkpoint to make sure you are on the right track.

```
In [10]: example_data_weights = np.array(len(train_data)* [1.5])
if best_splitting_feature(train_data, features, target, example_data_weights)
== 'term_36 months':
    print('Test passed!')
else:
    print('Test failed... try again!')
```

Test passed!

Aside. Relationship between weighted error and weight of mistakes:

By definition, the weighted error is the weight of mistakes divided by the weight of all data points, so

$$E(\alpha, \hat{y}) = \frac{\sum_{i=1}^n \alpha_i \times 1[y_i \neq \hat{y}_i]}{\sum_{i=1}^n \alpha_i} = \frac{WM(\alpha, \hat{y})}{\sum_{i=1}^n \alpha_i}.$$

In the code above, we obtain $E(\alpha, \hat{y})$ from the two weights of mistakes from both sides, $WM(\alpha_{\text{left}}, \hat{y}_{\text{left}})$ and $WM(\alpha_{\text{right}}, \hat{y}_{\text{right}})$. First, notice that the overall weight of mistakes $WM(\alpha, \hat{y})$ can be broken into two weights of mistakes over either side of the split:

$$\begin{aligned} WM(\alpha, \hat{y}) &= \sum_{i=1}^n \alpha_i \times 1[y_i \neq \hat{y}_i] = \sum_{\text{left}} \alpha_i \times 1[y_i \neq \hat{y}_i] + \sum_{\text{right}} \alpha_i \times 1[y_i \neq \hat{y}_i] \\ &= WM(\alpha_{\text{left}}, \hat{y}_{\text{left}}) + WM(\alpha_{\text{right}}, \hat{y}_{\text{right}}) \end{aligned}$$

We then divide through by the total weight of all data points to obtain $E(\alpha, \hat{y})$:

$$E(\alpha, \hat{y}) = \frac{WM(\alpha_{\text{left}}, \hat{y}_{\text{left}}) + WM(\alpha_{\text{right}}, \hat{y}_{\text{right}})}{\sum_{i=1}^n \alpha_i}$$

Building the tree

With the above functions implemented correctly, we are now ready to build our decision tree. A decision tree will be represented as a dictionary which contains the following keys:

```
{
    'is_leaf'           : True/False.
    'prediction'        : Prediction at the leaf node.
    'left'              : (dictionary corresponding to the left tree).
    'right'             : (dictionary corresponding to the right tree).
    'features_remaining': List of features that are possible splits.
}
```

Let us start with a function that creates a leaf node given a set of target values:

```
In [11]: def create_leaf(target_values, data_weights):  
  
    # Create a leaf node  
    leaf = {'splitting_feature' : None,  
            'is_leaf': True}  
  
    # Computed weight of mistakes.  
    weighted_error, best_class = intermediate_node_weighted_mistakes(target_val  
ues, data_weights)  
    # Store the predicted class (1 or -1) in leaf['prediction']  
    ## YOUR CODE HERE  
    ...  
    leaf['prediction'] = best_class  
  
    return leaf
```

We provide a function that learns a weighted decision tree recursively and implements 3 stopping conditions:

1. All data points in a node are from the same class.
2. No more features to split on.
3. Stop growing the tree when the tree depth reaches **max_depth**.

```

In [12]: def weighted_decision_tree_create(data, features, target, data_weights, current_depth = 1, max_depth = 10):
    remaining_features = features[:] # Make a copy of the features.
    target_values = data[target]
    print("-----")
    print("Subtree, depth = %s (%s data points)." % (current_depth, len(target_values)))

    # Stopping condition 1. Error is 0.
    if intermediate_node_weighted_mistakes(target_values, data_weights)[0] <= 1e-15:
        print("Stopping condition 1 reached.")
        return create_leaf(target_values, data_weights)

    # Stopping condition 2. No more features.
    if remaining_features == []:
        print("Stopping condition 2 reached.")
        return create_leaf(target_values, data_weights)

    # Additional stopping condition (limit tree depth)
    if current_depth > max_depth:
        print("Reached maximum depth. Stopping for now.")
        return create_leaf(target_values, data_weights)

    # If all the datapoints are the same, splitting_feature will be None. Create a leaf
    splitting_feature = best_splitting_feature(data, features, target, data_weights)
    remaining_features.remove(splitting_feature)

    left_split = data[data[splitting_feature] == 0]
    right_split = data[data[splitting_feature] == 1]

    left_data_weights = data_weights[data[splitting_feature] == 0]
    right_data_weights = data_weights[data[splitting_feature] == 1]

    print("Split on feature %s. (%s, %s)" % (splitting_feature, len(left_split), len(right_split)))

    # Create a leaf node if the split is "perfect"
    if len(left_split) == len(data):
        print("Creating leaf node.")
        return create_leaf(left_split[target], data_weights)
    if len(right_split) == len(data):
        print("Creating leaf node.")
        return create_leaf(right_split[target], data_weights)

    # Repeat (recurse) on left and right subtrees
    ## YOUR CODE HERE
    left_tree = weighted_decision_tree_create(left_split, remaining_features, target, left_data_weights, current_depth + 1, max_depth)
    right_tree = weighted_decision_tree_create(right_split, remaining_features, target, right_data_weights, current_depth + 1, max_depth)

    return {'is_leaf': False,

```

```

'prediction'      : None,
'splitting_feature': splitting_feature,
'left'           : left_tree,
'right'          : right_tree}

```

Here is a recursive function to count the nodes in your tree:

```

In [13]: def count_nodes(tree):
          if tree['is_leaf']:
              return 1
          return 1 + count_nodes(tree['left']) + count_nodes(tree['right'])

```

Run the following test code to check your implementation. Make sure you get **'Test passed'** before proceeding.

```

In [14]: example_data_weights = np.array([1.0 for i in range(len(train_data))])
          small_data_decision_tree = weighted_decision_tree_create(train_data, features,
                                                                    target,
                                                                    example_data_weights, max_depth=2)
          if count_nodes(small_data_decision_tree) == 7:
              print('Test passed!')
          else:
              print('Test failed... try again!')
              print('Number of nodes found:', count_nodes(small_data_decision_tree))
              print('Number of nodes that should be there: 7')

```

```

-----
Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)
-----

```

```

Subtree, depth = 2 (8850 data points).
Split on feature grade_A. (8775, 75)
-----

```

```

Subtree, depth = 3 (8775 data points).
Reached maximum depth. Stopping for now.
-----

```

```

Subtree, depth = 3 (75 data points).
Reached maximum depth. Stopping for now.
-----

```

```

Subtree, depth = 2 (23150 data points).
Split on feature grade_D. (19331, 3819)
-----

```

```

Subtree, depth = 3 (19331 data points).
Reached maximum depth. Stopping for now.
-----

```

```

Subtree, depth = 3 (3819 data points).
Reached maximum depth. Stopping for now.
Test passed!

```

Let us take a quick look at what the trained tree is like. You should get something that looks like the following

```
{'is_leaf': False,
  'left': {'is_leaf': False,
    'left': {'is_leaf': True, 'prediction': -1, 'splitting_feature': None},
    'prediction': None,
    'right': {'is_leaf': True, 'prediction': 1, 'splitting_feature': None},
    'splitting_feature': 'grade_A'
  },
  'prediction': None,
  'right': {'is_leaf': False,
    'left': {'is_leaf': True, 'prediction': 1, 'splitting_feature': None},
    'prediction': None,
    'right': {'is_leaf': True, 'prediction': -1, 'splitting_feature': None},
    'splitting_feature': 'grade_D'
  },
  'splitting_feature': 'term. 36 months'
}
```

In [15]: `small_data_decision_tree`

```
Out[15]: {'is_leaf': False,
  'prediction': None,
  'splitting_feature': 'term_ 36 months',
  'left': {'is_leaf': False,
    'prediction': None,
    'splitting_feature': 'grade_A',
    'left': {'splitting_feature': None, 'is_leaf': True, 'prediction': -1},
    'right': {'splitting_feature': None, 'is_leaf': True, 'prediction': 1}},
  'right': {'is_leaf': False,
    'prediction': None,
    'splitting_feature': 'grade_D',
    'left': {'splitting_feature': None, 'is_leaf': True, 'prediction': 1},
    'right': {'splitting_feature': None, 'is_leaf': True, 'prediction': -1}}
```

Making predictions with a weighted decision tree

We give you a function that classifies one data point. It can also return the probability if you want to play around with that as well.


```
In [16]: def classify(tree, x, annotate = False):
# If the node is a leaf node.
if tree['is_leaf']:
    if annotate:
        print("At leaf, predicting %s" % tree['prediction'])
    return tree['prediction']
else:
    # Split on feature.
    split_feature_value = x[tree['splitting_feature']]
    if annotate:
        print("Split on %s = %s" % (tree['splitting_feature'], split_feature_value))
    if split_feature_value == 0:
        return classify(tree['left'], x, annotate)
    else:
        return classify(tree['right'], x, annotate)
```

Evaluating the tree

Now, we will write a function to evaluate a decision tree by computing the classification error of the tree on the given dataset.

Again, recall that the **classification error** is defined as follows:

$$\text{classification error} = \frac{\# \text{ mistakes}}{\# \text{ all data points}}$$

The function called **evaluate_classification_error** takes in as input:

1. tree (as described above)
2. data (a dataframe)

The function does not change because of adding data point weights.

```
In [17]: def evaluate_classification_error(tree, data):
# Apply the classify(tree, x) to each row in your data
# YOUR CODE HERE
...
prediction = []
for i in range(len(data.index)):
    prediction.append(classify(tree, data.iloc[i]))
# Once you've made the predictions, calculate the classification error
return (prediction != data[target]).sum() / float(len(data))
```

```
In [18]: evaluate_classification_error(small_data_decision_tree, test_data)
```

```
Out[18]: 0.390875
```

Example: Training a weighted decision tree

To build intuition on how weighted data points affect the tree being built, consider the following:

Suppose we only care about making good predictions for the **first 10 and last 10 items** in `train_data`, we assign weights:

- 1 to the last 10 items
- 1 to the first 10 items
- and 0 to the rest.

Let us fit a weighted decision tree with `max_depth = 2`.

```
In [19]: # Assign weights
example_data_weights = np.array([1.] * 10 + [0.]*(len(train_data) - 20) + [1.]
                                * 10)

# Train a weighted decision tree model.
small_data_decision_tree_subset_20 = weighted_decision_tree_create(train_data,
                                                                    features, target,
                                                                    example_data_weights, max_depth=2)
```

```
-----
Subtree, depth = 1 (32000 data points).
Split on feature emp_length_10+ years. (22413, 9587)
```

```
-----
Subtree, depth = 2 (22413 data points).
Split on feature grade_A. (19673, 2740)
```

```
-----
Subtree, depth = 3 (19673 data points).
Reached maximum depth. Stopping for now.
```

```
-----
Subtree, depth = 3 (2740 data points).
Stopping condition 1 reached.
```

```
-----
Subtree, depth = 2 (9587 data points).
Stopping condition 1 reached.
```

Now, we will compute the classification error on the `subset_20`, i.e. the subset of data points whose weight is 1 (namely the first and last 10 data points).

```
In [20]: subset_20 = train_data.head(10).append(train_data.tail(10))
evaluate_classification_error(small_data_decision_tree_subset_20, subset_20)
```

```
Out[20]: 0.15
```

Now, let us compare the classification error of the model `small_data_decision_tree_subset_20` on the entire test set `train_data`:

```
In [21]: evaluate_classification_error(small_data_decision_tree_subset_20, train_data)
```

```
Out[21]: 0.445625
```

The model `small_data_decision_tree_subset_20` performs **a lot** better on `subset_20` than on `train_data`.

So, what does this mean?

- The points with higher weights are the ones that are more important during the training process of the weighted decision tree.
- The points with zero weights are basically ignored during training.

Implementing your own Adaboost (on decision stumps)

Now that we have a weighted decision tree working, it takes only a bit of work to implement Adaboost. For the sake of simplicity, let us stick with **decision tree stumps** by training trees with `max_depth=1`.

Recall from the lecture notes the procedure for Adaboost:

1. Start with unweighted data with $\alpha_j = 1$

2. For $t = 1, \dots, T$:

- Learn $f_t(x)$ with data weights α_j
- Compute coefficient \hat{w}_t :

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - E(\alpha, \hat{y})}{E(\alpha, \hat{y})} \right)$$

- Re-compute weights α_j :

$$\alpha_j \leftarrow \begin{cases} \alpha_j \exp(-\hat{w}_t) & \text{if } f_t(x_j) = y_j \\ \alpha_j \exp(\hat{w}_t) & \text{if } f_t(x_j) \neq y_j \end{cases}$$

- Normalize weights α_j :

$$\alpha_j \leftarrow \frac{\alpha_j}{\sum_{i=1}^N \alpha_i}$$

Complete the skeleton for the following code to implement **adaboost_with_tree_stumps**. Fill in the places with YOUR CODE HERE.

```

In [22]: from math import log
         from math import exp

def adaboost_with_tree_stumps(data, features, target, num_tree_stumps):
    # start with unweighted data (uniformly weighted)
    alpha = np.array([1.]*len(data))
    weights = []
    tree_stumps = []
    target_values = data[target]

    for t in range(num_tree_stumps):
        print('=====')
        print('Adaboost Iteration %d' % t)
        print('=====')
        # Learn a weighted decision tree stump. Use max_depth=1
        # YOUR CODE HERE
        ...
        tree_stumps.append(weighted_decision_tree_create(data, features, target, alpha, max_depth=1))

        # Make predictions
        ## YOUR CODE HERE
        ...
        predictions = []
        for i in range(len(data.index)):
            predictions.append(classify(tree_stumps[t], data.iloc[i]))

        print(len(predictions))

        # Produce a Boolean array indicating whether
        # each data point was correctly classified
        is_correct = predictions == target_values
        is_wrong = predictions != target_values

        # Compute weighted error
        ## YOUR CODE HERE
        ...
        weighted_error = sum(alpha * is_wrong) / sum(alpha)

        # Compute model coefficient using weighted error
        ## YOUR CODE HERE
        ...
        weight = 0.5 * log((1 - weighted_error)/weighted_error)

        weights.append(weight)

        # Adjust weights on data point
        ## YOUR CODE HERE
        adjustment = is_correct.apply(lambda is_correct : exp(-weight) if is_correct else exp(weight))

        # Scale alpha by multiplying by adjustment
        # Then normalize data points weights
        ## YOUR CODE HERE
        ...

```

```

        alpha = alpha * adjustment
        alpha = alpha / sum(alpha)

    return weights, tree_stumps

```

Checking your Adaboost code

Train an ensemble of **two** tree stumps and see which features those stumps split on. We will run the algorithm with the following parameters:

- train_data
- features
- target
- num_tree_stumps = 2

```
In [23]: stump_weights, tree_stumps = adaboost_with_tree_stumps(train_data, features, target, num_tree_stumps=2)
```

```

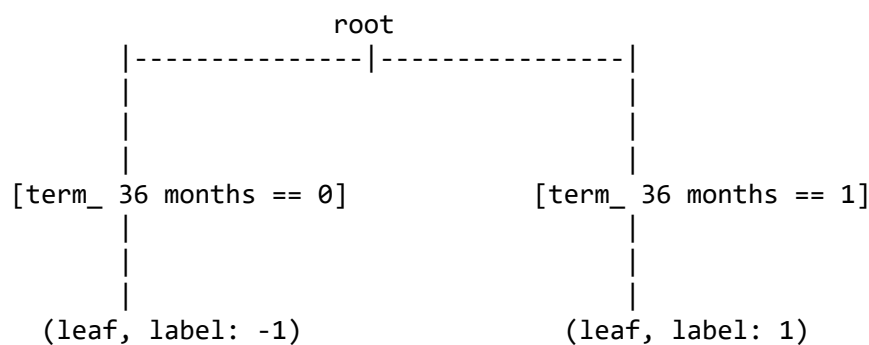
=====
Adaboost Iteration 0
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)
-----
Subtree, depth = 2 (8850 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (23150 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 1
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)
-----
Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000

```

```
In [24]: def print_stump(tree):
    split_name = tree['splitting_feature'] # split_name is something like 'term_36 months'
    if split_name is None:
        print("(leaf, label: %s)" % tree['prediction'])
        return None
    split_feature, split_value = split_name.split('_')
    print('
                                root')
    print('
    |-----|-----|')
    print('
    |
    |
    |
    |
    |')
    print(' [{0} == {0}]{1}[{0} == {1}]'.format(split_name, ' '*(27-len(split_name))))
    print('
    |
    |
    |
    |')
    print(' (%s)
    |
    |
    |
    |')
    print(' (%s)' \
    % (('leaf, label: ' + str(tree['left']['prediction']) if tree['left']['is_leaf'] else 'subtree'),
    ('leaf, label: ' + str(tree['right']['prediction']) if tree['right']['is_leaf'] else 'subtree'))
```

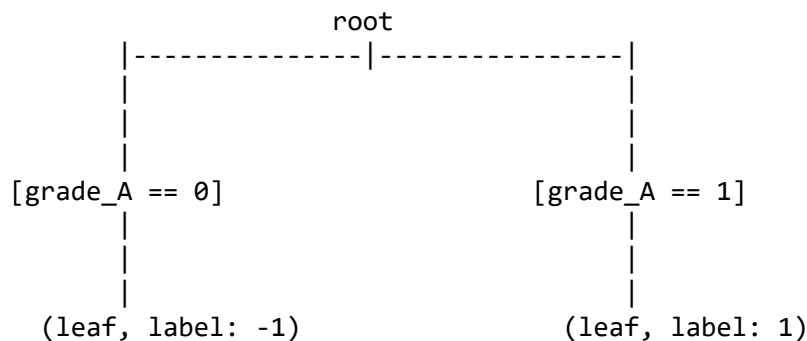
Here is what the first stump looks like:

```
In [25]: print_stump(tree_stumps[0])
```



Here is what the next stump looks like:

```
In [26]: print_stump(tree_stumps[1])
```



```
In [27]: print(stump_weights)
```

```
[0.17198848113764034, 0.1772878063726963]
```

If your Adaboost is correctly implemented, the following things should be true:

- `tree_stumps[0]` should split on **term. 36 months** with the prediction -1 on the left and +1 on the right.
- `tree_stumps[1]` should split on **grade.A** with the prediction -1 on the left and +1 on the right.
- Weights should be approximately `[0.17, 0.18]`

Reminders

- Stump weights (\hat{w}) and data point weights (α) are two different concepts.
- Stump weights (\hat{w}) tell you how important each stump is while making predictions with the entire boosted ensemble.
- Data point weights (α) tell you how important each data point is while training a decision stump.

Training a boosted ensemble of 10 stumps

Let us train an ensemble of 10 decision tree stumps with Adaboost. We run the **adaboost_with_tree_stumps** function with the following parameters:

- `train_data`
- `features`
- `target`
- `num_tree_stumps = 10`

```
In [28]: stump_weights, tree_stumps = adaboost_with_tree_stumps(train_data, features,  
                                                                target, num_tree_stumps=10)
```



```
=====
Adaboost Iteration 0
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)
-----
Subtree, depth = 2 (8850 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (23150 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 1
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)
-----
Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 2
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_D. (26027, 5973)
-----
Subtree, depth = 2 (26027 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (5973 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 3
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_B. (23457, 8543)
-----
Subtree, depth = 2 (23457 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (8543 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 4
=====
-----
Subtree, depth = 1 (32000 data points).
```

Split on feature grade_E. (28766, 3234)

Subtree, depth = 2 (28766 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3234 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 5

=====

Subtree, depth = 1 (32000 data points).
Split on feature home_ownership_MORTGAGE. (16870, 15130)

Subtree, depth = 2 (16870 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (15130 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 6

=====

Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)

Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 7

=====

Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)

Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 8

=====

Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)

Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.

```
Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000
```

```
=====
Adaboost Iteration 9
=====
```

```
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_E. (28766, 3234)
-----
```

```
Subtree, depth = 2 (28766 data points).
Reached maximum depth. Stopping for now.
-----
```

```
Subtree, depth = 2 (3234 data points).
Reached maximum depth. Stopping for now.
32000
```

Making predictions

Recall from the lecture that in order to make predictions, we use the following formula:

$$\hat{y} = \text{sign} \left(\sum_{t=1}^T \hat{w}_t f_t(x) \right)$$

We need to do the following things:

- Compute the predictions $f_t(x)$ using the t -th decision tree
- Compute $\hat{w}_t f_t(x)$ by multiplying the `stump_weights` with the predictions $f_t(x)$ from the decision trees
- Sum the weighted predictions over each stump in the ensemble.

Complete the following skeleton for making predictions:

```
In [29]: def predict_adaboost(stump_weights, tree_stumps, data):
    scores = np.array([0.]*len(data))

    for i, tree_stump in enumerate(tree_stumps):
        predictions = data.apply(lambda x: classify(tree_stump, x), axis = 1)

        # Accumulate predictions on scores array
        # YOUR CODE HERE
        ...
        for j in range(len(data.index)):
            scores[j] = scores[j] + predictions.iloc[j] * stump_weights[i]

    for j in range(len(data.index)):
        if scores[j] > 0:
            scores[j] = 1
        else:
            scores[j] = -1

    return scores
```

```
In [30]: predictions = predict_adaboost(stump_weights, tree_stumps, test_data)

from sklearn.metrics import accuracy_score
accuracy = accuracy_score(test_data[target], predictions)
print('Accuracy of 10-component ensemble = %s' % accuracy)
```

Accuracy of 10-component ensemble = 0.62825

Now, let us take a quick look what the `stump_weights` look like at the end of each iteration of the 10-stump ensemble:

```
In [31]: stump_weights
```

```
Out[31]: [0.17198848113764034,
          0.1772878063726963,
          0.10308067697010909,
          0.08686702058336851,
          0.07220085937793974,
          0.07438562925258671,
          0.05834552873244469,
          0.04545487026475097,
          0.0319454846001187,
          0.023305292432239024]
```

Question i: Are the weights monotonically decreasing, monotonically increasing, or neither?

Reminder: Stump weights (\hat{w}) tell you how important each stump is while making predictions with the entire boosted ensemble.

Stump weights are overall monotonically decreasing.

The reason is that the weighted classification error is converging to 50%. Therefore, the alphas is decreasing, because the weighted misclassification error is in the formula for the alphas. So the stump weights are decreasing.

Performance plots

In this section, we will try to reproduce some performance plots.

How does accuracy change with adding stumps to the ensemble?

We will now train an ensemble with:

- `train_data`
- `features`
- `target`
- `num_tree_stumps = 30`

Once we are done with this, we will then do the following:

- Compute the classification error at the end of each iteration.
- Plot a curve of classification error vs iteration.

First, lets train the model.

```
In [32]: # this may take a while...  
stump_weights, tree_stumps = adaboost_with_tree_stumps(train_data,  
                                                         features, target, num_tree_stumps=30)
```

```
=====
Adaboost Iteration 0
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)
-----
Subtree, depth = 2 (8850 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (23150 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 1
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)
-----
Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 2
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_D. (26027, 5973)
-----
Subtree, depth = 2 (26027 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (5973 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 3
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_B. (23457, 8543)
-----
Subtree, depth = 2 (23457 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (8543 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 4
=====
-----
Subtree, depth = 1 (32000 data points).
```

Split on feature grade_E. (28766, 3234)

Subtree, depth = 2 (28766 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3234 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 5

Subtree, depth = 1 (32000 data points).
Split on feature home_ownership_MORTGAGE. (16870, 15130)

Subtree, depth = 2 (16870 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (15130 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 6

Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)

Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 7

Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)

Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.
32000

=====

Adaboost Iteration 8

Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)

Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.


```
Subtree, depth = 2 (3919 data points).  
Reached maximum depth. Stopping for now.  
32000
```

```
=====
```

Adaboost Iteration 9

```
=====
```

```
-----
```

Subtree, depth = 1 (32000 data points).
Split on feature grade_E. (28766, 3234)

```
-----
```

Subtree, depth = 2 (28766 data points).
Reached maximum depth. Stopping for now.

```
-----
```

Subtree, depth = 2 (3234 data points).
Reached maximum depth. Stopping for now.

```
-----
```

32000

```
=====
```

Adaboost Iteration 10

```
=====
```

```
-----
```

Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)

```
-----
```

Subtree, depth = 2 (8850 data points).
Reached maximum depth. Stopping for now.

```
-----
```

Subtree, depth = 2 (23150 data points).
Reached maximum depth. Stopping for now.

```
-----
```

32000

```
=====
```

Adaboost Iteration 11

```
=====
```

```
-----
```

Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)

```
-----
```

Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.

```
-----
```

Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.

```
-----
```

32000

```
=====
```

Adaboost Iteration 12

```
=====
```

```
-----
```

Subtree, depth = 1 (32000 data points).
Split on feature emp_length_10+ years. (22413, 9587)

```
-----
```

Subtree, depth = 2 (22413 data points).
Reached maximum depth. Stopping for now.

```
-----
```

Subtree, depth = 2 (9587 data points).
Reached maximum depth. Stopping for now.

```
-----
```

32000

```
=====
```

Adaboost Iteration 13

```
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_B. (23457, 8543)
-----
Subtree, depth = 2 (23457 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (8543 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 14
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)
-----
Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 15
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_D. (26027, 5973)
-----
Subtree, depth = 2 (26027 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (5973 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 16
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)
-----
Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 17
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)
-----
```

Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.

32000

=====

Adaboost Iteration 18

=====

Subtree, depth = 1 (32000 data points).
Split on feature grade_E. (28766, 3234)

Subtree, depth = 2 (28766 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (3234 data points).
Reached maximum depth. Stopping for now.

32000

=====

Adaboost Iteration 19

=====

Subtree, depth = 1 (32000 data points).
Split on feature grade_C. (23388, 8612)

Subtree, depth = 2 (23388 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (8612 data points).
Reached maximum depth. Stopping for now.

32000

=====

Adaboost Iteration 20

=====

Subtree, depth = 1 (32000 data points).
Split on feature home_ownership_MORTGAGE. (16870, 15130)

Subtree, depth = 2 (16870 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (15130 data points).
Reached maximum depth. Stopping for now.

32000

=====

Adaboost Iteration 21

=====

Subtree, depth = 1 (32000 data points).
Split on feature term_36 months. (8850, 23150)

Subtree, depth = 2 (8850 data points).
Reached maximum depth. Stopping for now.

Subtree, depth = 2 (23150 data points).
Reached maximum depth. Stopping for now.

```
32000
=====
Adaboost Iteration 22
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_F. (30624, 1376)
-----
Subtree, depth = 2 (30624 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (1376 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 23
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_B. (23457, 8543)
-----
Subtree, depth = 2 (23457 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (8543 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 24
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature emp_length_2 years. (29104, 2896)
-----
Subtree, depth = 2 (29104 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (2896 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 25
=====
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_G. (31657, 343)
-----
Subtree, depth = 2 (31657 data points).
Reached maximum depth. Stopping for now.
-----
Subtree, depth = 2 (343 data points).
Reached maximum depth. Stopping for now.
32000
=====
Adaboost Iteration 26
=====
-----
```

```
Subtree, depth = 1 (32000 data points).
Split on feature grade_A. (28081, 3919)
```

```
-----
Subtree, depth = 2 (28081 data points).
Reached maximum depth. Stopping for now.
```

```
-----
Subtree, depth = 2 (3919 data points).
Reached maximum depth. Stopping for now.
32000
```

```
=====
Adaboost Iteration 27
```

```
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_G. (31657, 343)
```

```
-----
Subtree, depth = 2 (31657 data points).
Reached maximum depth. Stopping for now.
```

```
-----
Subtree, depth = 2 (343 data points).
Reached maximum depth. Stopping for now.
32000
```

```
=====
Adaboost Iteration 28
```

```
-----
Subtree, depth = 1 (32000 data points).
Split on feature home_ownership_OWN. (29204, 2796)
```

```
-----
Subtree, depth = 2 (29204 data points).
Reached maximum depth. Stopping for now.
```

```
-----
Subtree, depth = 2 (2796 data points).
Reached maximum depth. Stopping for now.
32000
```

```
=====
Adaboost Iteration 29
```

```
-----
Subtree, depth = 1 (32000 data points).
Split on feature grade_G. (31657, 343)
```

```
-----
Subtree, depth = 2 (31657 data points).
Reached maximum depth. Stopping for now.
```

```
-----
Subtree, depth = 2 (343 data points).
Reached maximum depth. Stopping for now.
32000
```

Computing training error at the end of each iteration

Now, we will compute the classification error on the **train_data** and see how it is reduced as trees are added.

```
In [33]: error_all = []
         for n in range(1, 31):
             predictions = predict_adaboost(stump_weights[:n], tree_stumps[:n], train_data)
             error = 1.0 - accuracy_score(train_data[target], predictions)
             error_all.append(error)
         print("Iteration %s, training error = %s" % (n, error_all[n-1]))
```

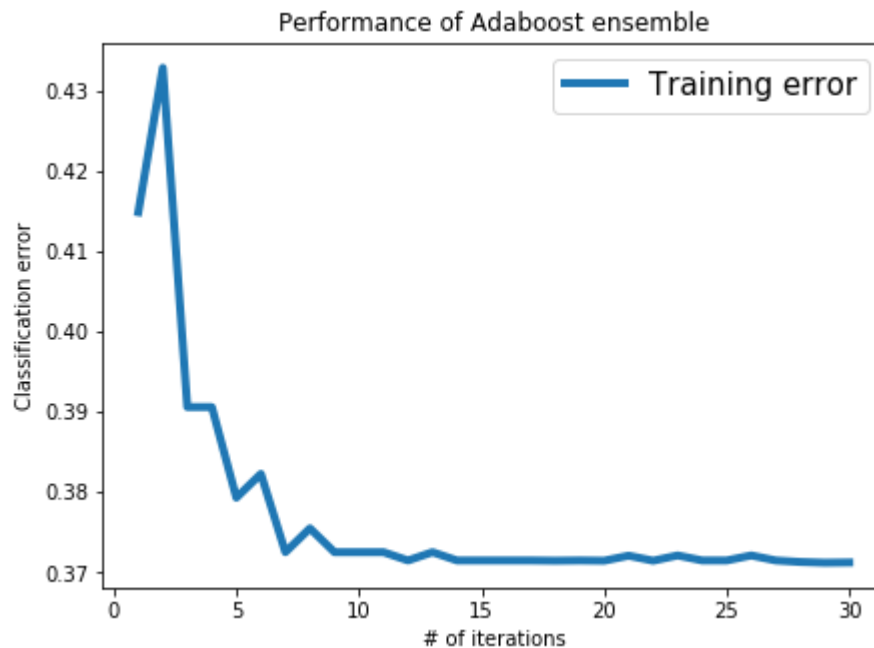
```
Iteration 1, training error = 0.41484374999999996
Iteration 2, training error = 0.43281250000000004
Iteration 3, training error = 0.39059374999999996
Iteration 4, training error = 0.39059374999999996
Iteration 5, training error = 0.37931250000000005
Iteration 6, training error = 0.38228125
Iteration 7, training error = 0.37253125
Iteration 8, training error = 0.37549999999999994
Iteration 9, training error = 0.37253125
Iteration 10, training error = 0.37253125
Iteration 11, training error = 0.37253125
Iteration 12, training error = 0.37150000000000005
Iteration 13, training error = 0.37253125
Iteration 14, training error = 0.37150000000000005
Iteration 15, training error = 0.37150000000000005
Iteration 16, training error = 0.37150000000000005
Iteration 17, training error = 0.37150000000000005
Iteration 18, training error = 0.37146875
Iteration 19, training error = 0.37150000000000005
Iteration 20, training error = 0.37146875
Iteration 21, training error = 0.37209375
Iteration 22, training error = 0.37146875
Iteration 23, training error = 0.37212500000000004
Iteration 24, training error = 0.37150000000000005
Iteration 25, training error = 0.37150000000000005
Iteration 26, training error = 0.37212500000000004
Iteration 27, training error = 0.37150000000000005
Iteration 28, training error = 0.37131250000000005
Iteration 29, training error = 0.37121875000000004
Iteration 30, training error = 0.37124999999999997
```

Visualizing training error vs number of iterations

We have provided you with a simple code snippet that plots classification error with the number of iterations.

```
In [34]: plt.rcParams['figure.figsize'] = 7, 5
plt.plot(list(range(1,31)), error_all, '-', linewidth=4.0, label='Training error')
plt.title('Performance of Adaboost ensemble')
plt.xlabel('# of iterations')
plt.ylabel('Classification error')
plt.legend(loc='best', prop={'size':15})

plt.rcParams.update({'font.size': 16})
```



Evaluation on the test data

Performing well on the training data is cheating, so let's make sure it works on the `test_data` as well. Here, we will compute the classification error on the `test_data` at the end of each iteration.

```
In [35]: test_error_all = []
for n in range(1, 31):
    predictions = predict_adaboost(stump_weights[:n], tree_stumps[:n], test_data)
    error = 1.0 - accuracy_score(test_data[target], predictions)
    test_error_all.append(error)
    print("Iteration %s, test error = %s" % (n, test_error_all[n-1]))
```

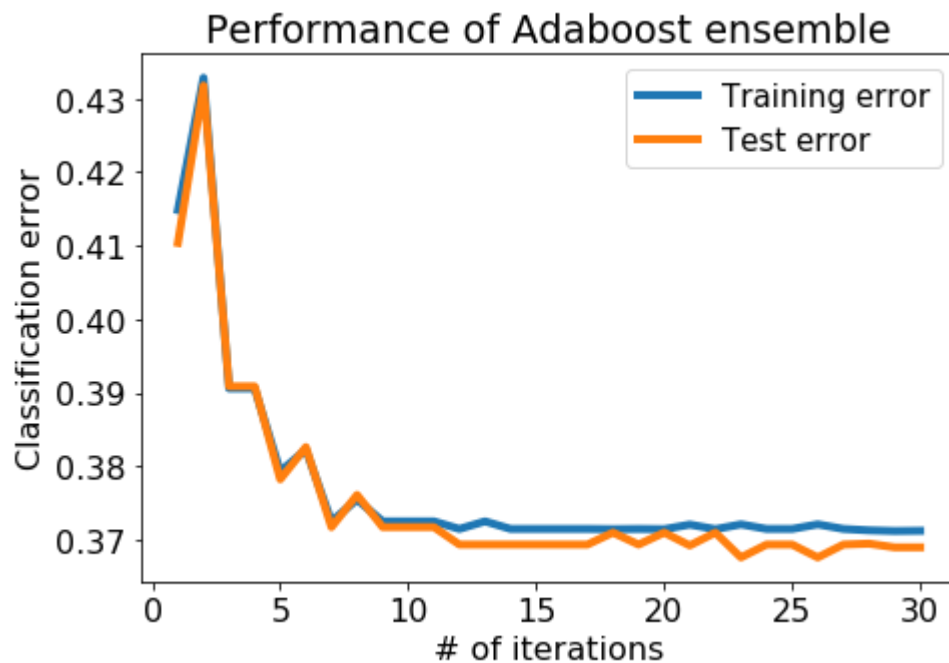
```
Iteration 1, test error = 0.41037500000000005
Iteration 2, test error = 0.43174999999999997
Iteration 3, test error = 0.390875
Iteration 4, test error = 0.390875
Iteration 5, test error = 0.37825
Iteration 6, test error = 0.382625
Iteration 7, test error = 0.37175
Iteration 8, test error = 0.37612500000000004
Iteration 9, test error = 0.37175
Iteration 10, test error = 0.37175
Iteration 11, test error = 0.37175
Iteration 12, test error = 0.369375
Iteration 13, test error = 0.369375
Iteration 14, test error = 0.369375
Iteration 15, test error = 0.369375
Iteration 16, test error = 0.369375
Iteration 17, test error = 0.369375
Iteration 18, test error = 0.371
Iteration 19, test error = 0.369375
Iteration 20, test error = 0.371
Iteration 21, test error = 0.36924999999999997
Iteration 22, test error = 0.371
Iteration 23, test error = 0.367625
Iteration 24, test error = 0.369375
Iteration 25, test error = 0.369375
Iteration 26, test error = 0.367625
Iteration 27, test error = 0.369375
Iteration 28, test error = 0.36950000000000005
Iteration 29, test error = 0.369
Iteration 30, test error = 0.369
```

Visualize both the training and test errors

Now, let us plot the training & test error with the number of iterations.


```
In [36]: plt.rcParams['figure.figsize'] = 7, 5
plt.plot(list(range(1,31)), error_all, '-', linewidth=4.0, label='Training error')
plt.plot(list(range(1,31)), test_error_all, '-', linewidth=4.0, label='Test error')

plt.title('Performance of Adaboost ensemble')
plt.xlabel('# of iterations')
plt.ylabel('Classification error')
plt.rcParams.update({'font.size': 16})
plt.legend(loc='best', prop={'size':15})
plt.tight_layout()
```



Question ii: From this plot (with 30 trees), is there massive overfitting as the # of iterations increases?

No. The Classification error is relatively stable as the number of iteration increase. Therefore, there is no massive overfitting.

In []: