

## Problem 1

- a. **Firstly**, let's investigate the loss function. We need to solve the constrained optimization problem below.

$$\arg \min_{f(x)} E_{Y|x} \exp \left( -\frac{1}{K} (Y_1 f_1 + \dots + Y_K f_K) \right)$$

$$\text{Subject to } f_1 + f_2 + \dots + f_K = 0$$

The Lagrange of this problem can be written as:

$$\exp\left(-\frac{f_1(X)}{K-1}\right) \text{Prob}(c=1|x) + \exp\left(-\frac{f_2(X)}{K-1}\right) \text{Prob}(c=2|x) \dots + \exp\left(-\frac{f_K(X)}{K-1}\right) \text{Prob}(c=K|x) - \lambda(f_1(X) + f_2(X) + \dots + f_K(X)).$$

The  $\lambda$  is the Lagrange multiplier.

**Then**, while taking derivative with respect to  $f_K$  and  $\lambda$ , we can get

$$-\frac{1}{K-1} \exp\left(-\frac{f_1(X)}{K-1}\right) \text{Prob}(c=1|x) - \lambda = 0$$

$$-\frac{1}{K-1} \exp\left(-\frac{f_2(X)}{K-1}\right) \text{Prob}(c=2|x) - \lambda = 0$$

...

$$-\frac{1}{K-1} \exp\left(-\frac{f_K(X)}{K-1}\right) \text{Prob}(c=K|x) - \lambda = 0$$

$$f_1(X) + f_2(X) + \dots + f_K(X) = 0$$

**Finally** solve the set of equations to obtain the population minimizer.

$$f_k^*(x) = (K-1)(\log \text{Prob}(c=k|x) - \frac{1}{K} \sum_{i=1}^K \log \text{Prob}(c=i|x)), \quad k=1,2,\dots,K$$

$$\text{Therefore, } \arg \max_k f_k^*(x) = \arg \max_k \text{Prob}(c=k|x)$$

This is the term of minimized misclassification error by Bayes optimal classification rule.

We can also get the class probability as a function of population minimizer:

$$\text{Prob}(c=k|x) = \frac{e^{\frac{1}{K-1} f_k^*(x)}}{e^{\frac{1}{K-1} f_1^*(x)} + e^{\frac{1}{K-1} f_2^*(x)} + \dots + e^{\frac{1}{K-1} f_K^*(x)}}, \quad k=1,2,\dots,K$$

- b. The multi-class boosting using this loss function can be obtained by an algorithm similar to AdaBoost. The algorithm goes as follow.

1. Initialize the weights  $w_i = 1/n$ ,  $i=1,2,\dots,n$
2. For  $m=1$  to  $M$ :
  - 1) Fit a classifier  $G^{(m)}(x)$  to the training data by weights  $w_i$ .
  - 2) Compute  $\text{err}^{(m)} = \sum_{i=1}^n w_i \text{sign}(C_i \neq G^{(m)}(x_i)) / \sum_{i=1}^n w_i$
  - 3) Compute  $\alpha^{(m)} = \log \frac{1 - \text{err}^{(m)}}{\text{err}^{(m)}} + \log(K-1)$
  - 4) Set  $w_i = w_i * \exp(\alpha^{(m)} * \text{sign}(C_i \neq G^{(m)}(x_i)))$ ,  $i=1,2,\dots,n$
  - 5) Re-normalized  $w_i$ .

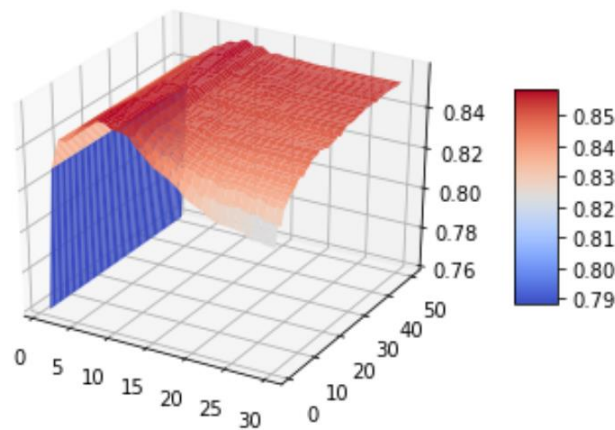
$$3. \text{ Output } C(x) = \arg \max_k \sum_{m=1}^M \alpha^{(m)} \text{sign}(G^{(m)}(x)=k)$$

The most important changes is  $\alpha^{(m)} = \log \frac{1 - \text{err}^{(m)}}{\text{err}^{(m)}} + \log(K-1)$ . In order for  $\alpha^{(m)}$  to be positive, we only need  $(1 - \text{err}^{(m)}) > 1/K$

This new algorithm puts more weight on the misclassified points in step 2 4) than AdaBoost.

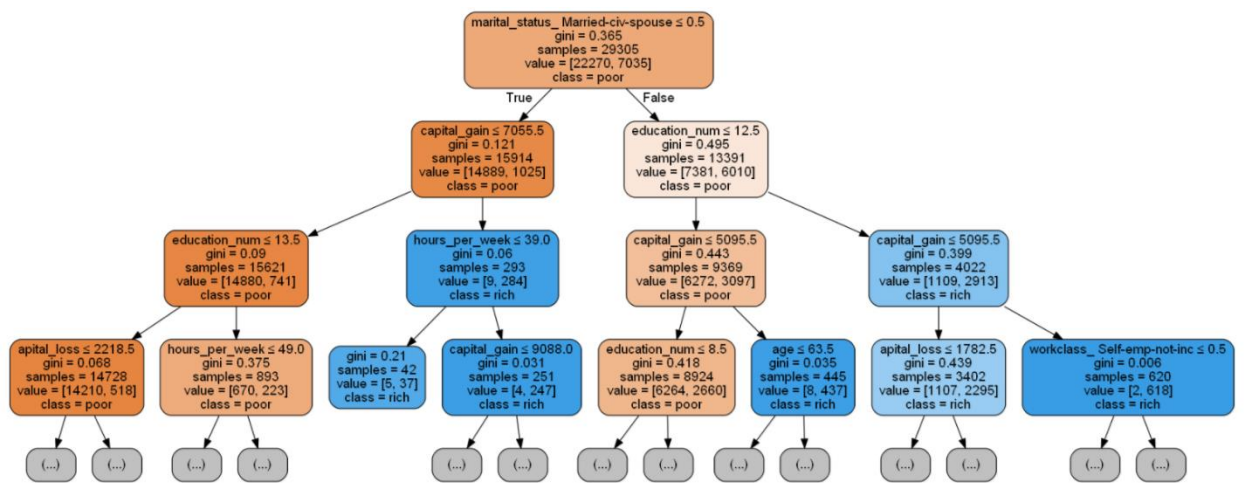
## Problem 2

- I jointed the train set and test set together and use `sklearn.model_selection.train_test_split` split the data into 60-40 train/test data set. I used the train data to calculate the mean and mode of the features and filled in all the missing values. Then I did one-hot preprocessing on the whole data set.
- I used `cv=3` to train the model.  
The best choice of max depth is: 9  
The best choice of minimum number of samples is: 42



The accuracy score when using the optimal max depth and minimum number of samples is 0.855658494139

- The top 3 levels of the decision tree is



### Problem 3

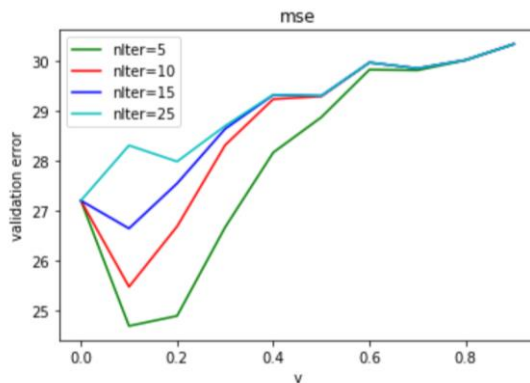
- a. I build the function `gd_tree()` to implement the gradient boosting algorithm using the tree repressors.

`gd_tree()` has parameters as follows:

Input: Xarray, yarray, boosting\_times, shrink, loss\_func, Xnew

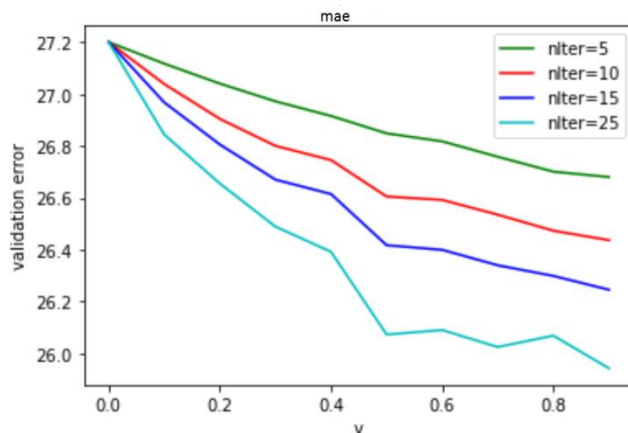
Output: ynew

- b. For the loss type is mean squared error



The optimal model is when  $v$  is 0.1 and  $niter$  is 10. At this time, the validation error comes to the lowest. As we can see, when  $v$  goes up, the RMSE goes down first and then goes up; when  $niter$  goes up, the RMSE goes up. Given the same RMSE (bigger than a certain amount), if the  $v$  is small, the  $niter$  should be big.

For the loss type is mean absolute error



The optimal model is when  $v$  is 0.9 and  $niter$  is 25. At this time, the validation error comes to the lowest. We can deduce from the plots -- Either  $niter$  or  $v$  goes up, the RMSE goes down. So given the same RMSE, if the  $v$  is small, the  $niter$  should be big.

- c. When the model uses the optimal parameters on the test data,  
for the mse loss function, RMSE with is 24.9898888739;  
for the mae loss function, RMSE with is 25.1713995395  
Even though the mae has better performance in validation(with lower RMSE), mse does better  
on test data(with higher RMSE).