

1. We handled missing attributes by calculating the mode of every attribute for the data set and storing them in an array. Then, we assigned the corresponding modal value to each missing attribute.

2. (restricted to 16 leaf nodes)

$(\text{oppnuminjured} \geq 3.0 \wedge \text{opprundifferential} < 18.0 \wedge \text{rundifferential} \geq 39.0 \wedge \text{opprundifferential} \geq 16.0 \wedge \text{rundifferential} \geq 42.0 \wedge \text{homeaway} = 0 \wedge \text{rundifferential} < 59.0 \wedge \text{winpercent} < 0.404196846425) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{oppnuminjured} < 0.0 \wedge \text{winpercent} < 0.590582828588 \wedge \text{oppdayssincegame} \geq 2.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{oppnuminjured} < 0.0 \wedge \text{winpercent} \geq 0.590582828588 \wedge \text{oppwinningpercent} \geq 0.494600828567 \wedge \text{winpercent} \geq 0.600132770207) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{opprundifferential} < 18.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{oppnuminjured} \geq 0.0 \wedge \text{opprundifferential} \geq 18.0 \wedge \text{oppdayssincegame} < 1.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{oppnuminjured} \geq 0.0 \wedge \text{opprundifferential} \geq 18.0 \wedge \text{oppdayssincegame} \geq 1.0 \wedge \text{oppwinningpercent} < 0.0806987298889) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge \text{oppnuminjured} \geq 0.0 \wedge \text{opprundifferential} \geq 18.0 \wedge \text{oppdayssincegame} \geq 1.0 \wedge \text{oppwinningpercent} \geq 0.0806987298889 \wedge \text{rundifferential} < 3.0 \wedge \text{winpercent} \geq 0.699154368204) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} < 0.125358171942 \wedge \text{oppwinningpercent} < 0.0068583423205) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} < 0.125358171942 \wedge \text{oppwinningpercent} \geq 0.0068583423205 \wedge \text{opprundifferential} < 23.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} < 0.125358171942 \wedge \text{oppwinningpercent} \geq 0.0068583423205 \wedge \text{opprundifferential} \geq 23.0 \wedge \text{winpercent} < 0.387374577582 \wedge \text{oppwinningpercent} \geq 0.0651108904444) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.125358171942 \wedge \text{oppwinningpercent} < 0.270328014607 \wedge \text{winpercent} < 0.289124011761) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.125358171942 \wedge \text{oppwinningpercent} < 0.270328014607 \wedge \text{winpercent} \geq 0.289124011761 \wedge \text{oppwinningpercent} < 0.189465166099 \wedge \text{winpercent} < 0.371677431261) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.125358171942 \wedge \text{oppwinningpercent} < 0.270328014607 \wedge \text{winpercent} \geq 0.289124011761 \wedge \text{oppwinningpercent} \geq 0.189465166099 \wedge \text{opprundifferential} < 36.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.125358171942 \wedge \text{oppwinningpercent} < 0.270328014607 \wedge \text{winpercent} \geq 0.289124011761 \wedge$

$\text{oppwinningpercent} \geq 0.189465166099 \wedge \text{opprundifferential} \geq 36.0 \wedge$
 $\text{winpercent} < 0.528333144456) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.125358171942 \wedge$
 $\text{oppwinningpercent} \geq 0.270328014607 \wedge \text{oppwinningpercent} < 0.307734761131 \wedge$
 $\text{winpercent} < 0.539602893417 \wedge \text{winpercent} < 0.539007439959 \wedge \text{opppdayssincegame} < 2.0) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge$
 $\text{oppwinningpercent} \geq 0.125358171942 \wedge \text{oppwinningpercent} \geq 0.270328014607 \wedge$
 $\text{oppwinningpercent} < 0.307734761131 \wedge \text{winpercent} < 0.539602893417 \wedge$
 $\text{winpercent} < 0.539007439959 \wedge \text{opppdayssincegame} \geq 2.0)$

3. $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge$
 $\text{oppnuminjured} < 0.0 \wedge \text{winpercent} \geq 0.590582828588 \wedge$
 $\text{oppwinningpercent} \geq 0.494600828567 \wedge \text{winpercent} \geq 0.600132770207)$

“If the opponent has no injured players (also means less than 3 and less than 1 from the other splits) and we have no injured players (also means less than 1 from other split) and our win percentage is equal to or higher than 60.01% (also means higher than or equal to 59.06% from other split) and the opponent winning percentage is equal to or higher than 49.46%, then we are predicted to win the game”

4. We implemented top-down pruning. We first get the mode classification value for the validation set (whether winner is 1 or 0 the most). Then, we look at the accuracy of the current node. If the accuracy of just predicting the mode value for winner (blind guessing) is greater than the accuracy of the current node, we prune that entire branch and replace it with a leaf node with the mode winner value (0 or 1). As nodes higher up the tree are the ones with the best information gain ratio, no node below it will be better, and thus we can prune that entire branch. If the current node's accuracy is higher, we recurse through its children.

5. (restricted to 16 leaf nodes)

$(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} < 0.0 \wedge$
 $\text{oppnuminjured} < 0.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge$
 $\text{numinjured} < 0.0 \wedge \text{oppnuminjured} \geq 0.0 \wedge \text{opprundifferential} < 30.0) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge$
 $\text{oppwinningpercent} \geq 0.110564199499 \wedge \text{oppwinningpercent} < 0.277250882063 \wedge$
 $\text{winpercent} < 0.452110687309) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} < 1.0 \wedge \text{numinjured} \geq 0.0 \wedge \text{oppwinningpercent} \geq 0.110564199499 \wedge$
 $\text{oppwinningpercent} \geq 0.277250882063) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} < 19.0 \wedge \text{rundifferential} < 12.0 \wedge$
 $\text{rundifferential} < 3.0 \wedge \text{weather} = 1) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} < 19.0 \wedge \text{rundifferential} < 12.0 \wedge$
 $\text{rundifferential} \geq 3.0 \wedge \text{opprundifferential} < 10.0 \wedge \text{winpercent} < 0.667310699054) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} < 19.0$
 $\wedge \text{rundifferential} \geq 12.0 \wedge \text{rundifferential} < 30.0 \wedge \text{opprundifferential} < 8.0 \wedge$
 $\text{oppwinningpercent} < 0.379750315096 \wedge \text{winpercent} < 0.573437483861) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} < 19.0$

$\wedge \text{rundifferential} \geq 12.0 \wedge \text{rundifferential} < 30.0 \wedge \text{opprundifferential} < 8.0 \wedge$
 $\text{oppwinningpercent} \geq 0.379750315096) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} < 19.0 \wedge \text{rundifferential} \geq 12.0 \wedge$
 $\text{rundifferential} < 30.0 \wedge \text{opprundifferential} \geq 8.0 \wedge \text{oppwinningpercent} \geq 0.527969069748 \wedge$
 $\text{rundifferential} \geq 22.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} < 19.0 \wedge \text{rundifferential} \geq 12.0 \wedge \text{rundifferential} \geq 30.0) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} < 23.0 \wedge \text{rundifferential} \geq 71.0) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} \geq 23.0 \wedge \text{opprundifferential} < 43.0 \wedge$
 $\text{rundifferential} < 53.0 \wedge \text{winpercent} \geq 0.409365261398 \wedge \text{temperature} < 74.4658041955 \wedge$
 $\text{rundifferential} < 14.0 \wedge \text{rundifferential} < 2.0) \vee (\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge$
 $\text{oppnuminjured} \geq 1.0 \wedge \text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} \geq 23.0 \wedge$
 $\text{opprundifferential} < 43.0 \wedge \text{rundifferential} \geq 53.0 \wedge \text{rundifferential} < 104.0 \wedge$
 $\text{rundifferential} < 87.0 \wedge \text{opprundifferential} < 27.0 \wedge \text{oppwinningpercent} \geq 0.257787028218 \wedge$
 $\text{oppwinningpercent} \geq 0.324607506889 \wedge \text{rundifferential} \geq 54.0 \wedge \text{startingpitcher} = 1) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} \geq 23.0 \wedge \text{opprundifferential} < 43.0 \wedge$
 $\text{rundifferential} \geq 53.0 \wedge \text{rundifferential} < 104.0 \wedge \text{rundifferential} < 87.0 \wedge$
 $\text{opprundifferential} < 27.0 \wedge \text{oppwinningpercent} \geq 0.257787028218 \wedge$
 $\text{oppwinningpercent} \geq 0.324607506889 \wedge \text{rundifferential} \geq 54.0 \wedge \text{startingpitcher} = 2) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} \geq 23.0 \wedge \text{opprundifferential} < 43.0 \wedge$
 $\text{rundifferential} \geq 53.0 \wedge \text{rundifferential} < 104.0 \wedge \text{rundifferential} < 87.0 \wedge$
 $\text{opprundifferential} < 27.0 \wedge \text{oppwinningpercent} \geq 0.257787028218 \wedge$
 $\text{oppwinningpercent} \geq 0.324607506889 \wedge \text{rundifferential} \geq 54.0 \wedge \text{startingpitcher} = 3) \vee$
 $(\text{oppnuminjured} < 3.0 \wedge \text{numinjured} < 1.0 \wedge \text{oppnuminjured} \geq 1.0 \wedge$
 $\text{opprundifferential} \geq 19.0 \wedge \text{opprundifferential} \geq 23.0 \wedge \text{opprundifferential} < 43.0 \wedge$
 $\text{rundifferential} \geq 53.0 \wedge \text{rundifferential} < 104.0 \wedge \text{rundifferential} < 87.0 \wedge$
 $\text{opprundifferential} < 27.0 \wedge \text{oppwinningpercent} \geq 0.257787028218 \wedge$
 $\text{oppwinningpercent} \geq 0.324607506889 \wedge \text{rundifferential} \geq 54.0 \wedge \text{startingpitcher} = 4)$

6. Splits for the **unpruned** tree: 2,799
 Split difference: 2,360

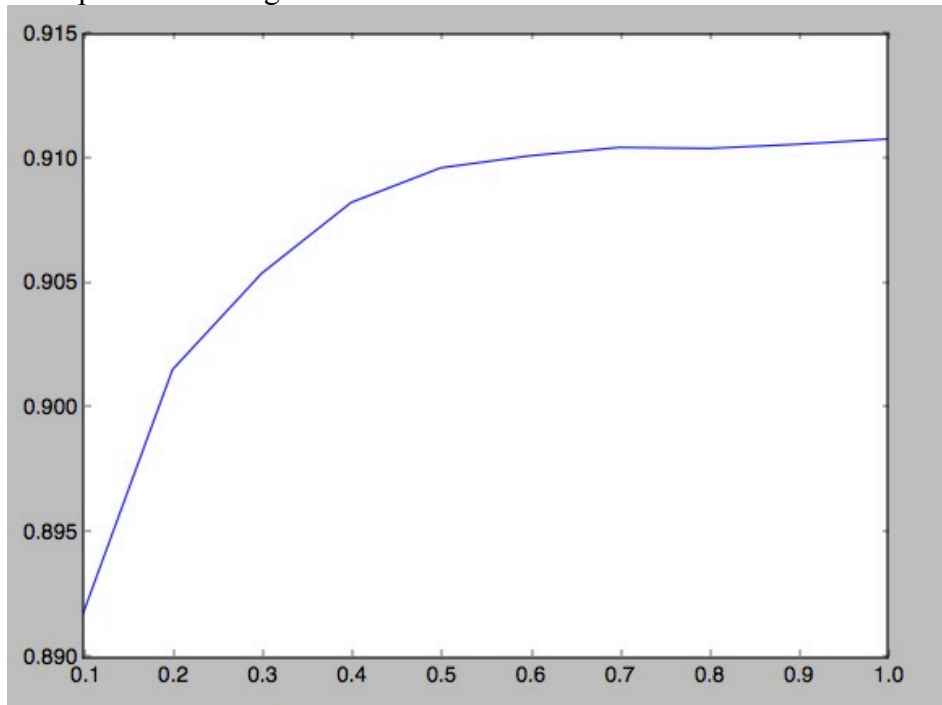
Splits for the **pruned** tree: 439

7. **Unpruned** accuracy: 0.916176470588 **Pruned** accuracy: 0.924474789916

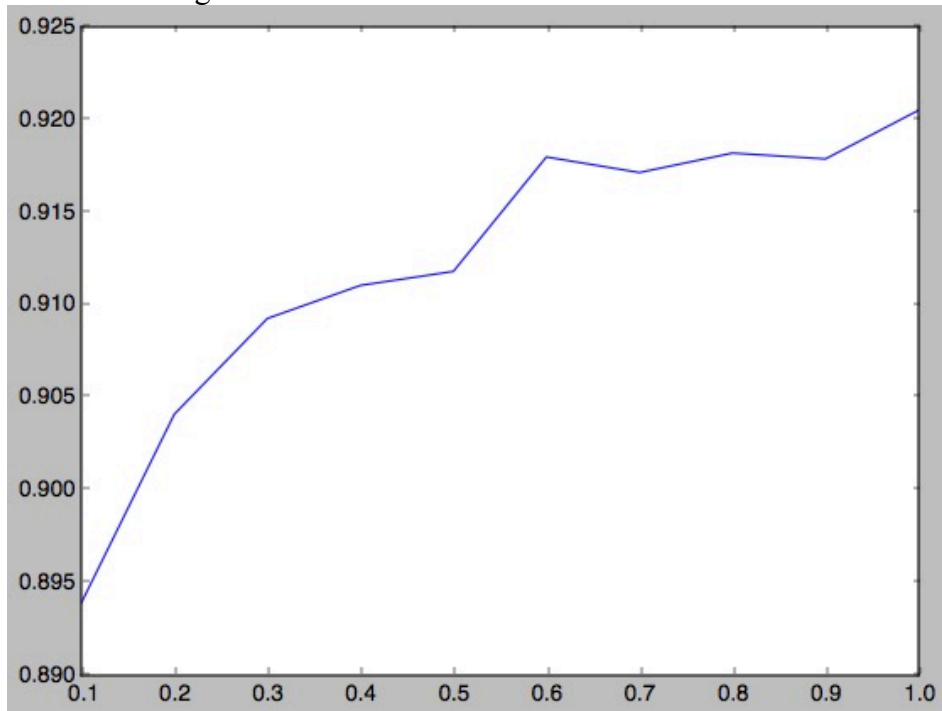
There is a slight difference in validation accuracy between the two, with the pruned tree being more accurate. This is because the depth and numeric split count values allowed for overfitting in the unpruned tree, and thus caused a decrease in validation accuracy.

Pruning the tree helped to eliminate the overfitting. However, the pruned accuracy isn't much higher than the unpruned accuracy. This is because we didn't allow the unpruned tree to way overfit the training data due to time constraints (by increasing depth or split counts).

8. Unpruned learning curve:



Pruned learning curve:



The pruned learning curve seems to have higher accuracy at every point as compared to the unpruned learning curve. It seems that amount of information doesn't change the fact that our pruned tree has higher validation accuracy than the unpruned tree.

9. We think the pruned tree will perform better on the unlabeled test set. Pruning helps deal with overfitting, which makes it more applicable to an unseen example, as shown with the validation accuracy difference. Whatever has the higher validation accuracy should do better on the test data. This is backed up by the graph on slide 43 of the Decision Trees lecture. The graph shows pruned trees consistently performing better on test data than unpruned trees.

10. Most of the work done for the assignment was worked on together with pair programming. However, we can try to split it up based on who did the majority of each part.

I (Jonathan) did most of:

ID3 functions

Alex did most of:

Node functions

Pruning functions

Graphing functions

11. First, using IGR instead of using just IG was a good model decision. With pruning, the validation accuracy using IGR was 0.924474789916 . Changing this to using just IG (still with pruning) actually dropped the validation accuracy to 0.913865546218 . This difference in accuracy is bigger than the accuracy difference between pruning and not pruning. Additionally, IGR isn't much more computationally expensive, meaning it is clearly the better option for the model. However, changing step size to include every possible threshold value in the data set is computationally prohibitive. While it may very well increase accuracy, the large difference in training time means using larger step sizes might be the better choice for the model.