

Distributed Random Forests: Resilient?

Raksha Kumaraswany
and
Katherine Metcalf

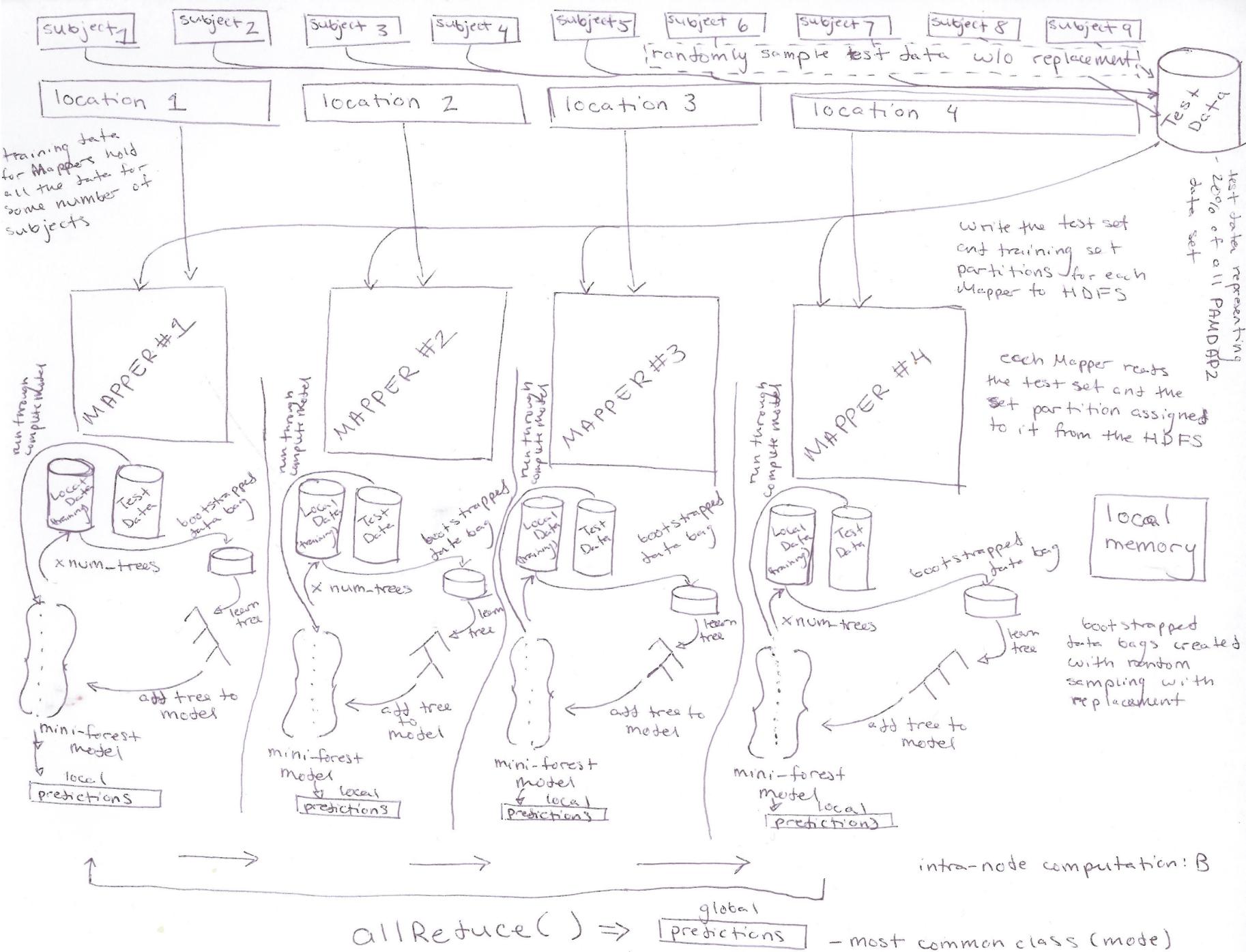
Problem Description

- Explore the effect of learning distributed mini-forests on local data sets that vary according to feature value and label distributions
- Important question b/c passing around data to build partitions with globally representative distributions is computationally expensive and puts the data at higher security risks
- The PAMDAP2 is a time series, activity classification data set:
 - binarized the labels into sedentary and non-sedentary actions
 - 53 continuous feature values from 3 IMU sensors and 1 heart rate monitor
 - > 3M data points
- Split the PAMDAP2 data set among the HARP Mappers by person to create local data sets with distinct feature value and label distributions

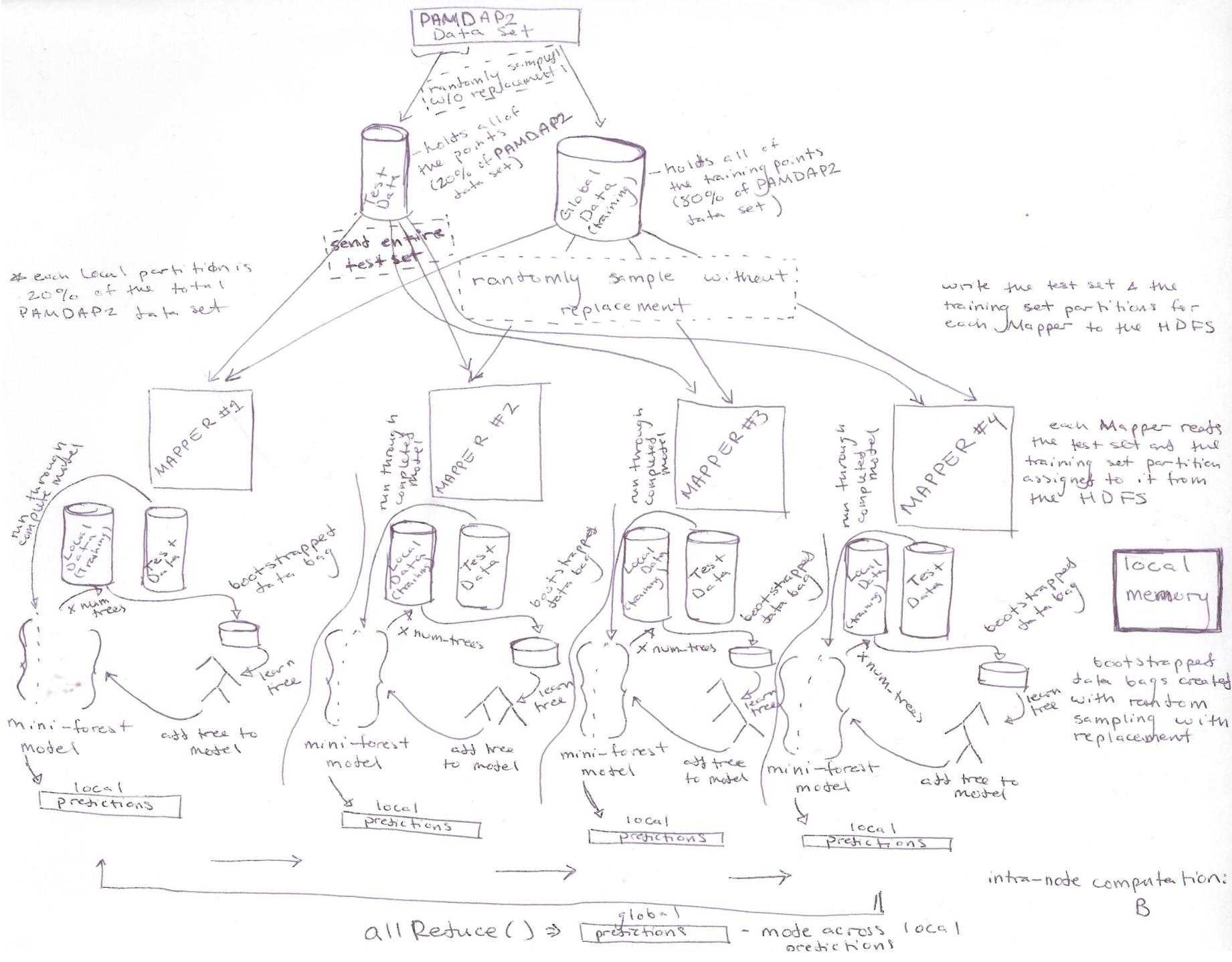
Random Forest

- ensemble classifier
- idea - divide & conquer
- divide - the dataset into bootstrap samples
- conquer - by learning tree for each sample
- gather all trees - a forest, the new model
- prediction - mode, or stochastic prediction based on distribution
- Advantage:
 - these independently learned models, when aggregated create a resilient model
 - not prone to bias
- Disadvantage:
 - computational complexity scaled with data size: many avenues for parallelism
 - bootstrap sample creation, forest creation, tree creation, prediction.
 - prone to variance, but can be controlled with sampling strategies: strategic bag-of-little bootstraps
- What we do - parallel forest creation, with strategic bag-of-little bootstraps

Methods: Local-Label and Local-Label- Feature data

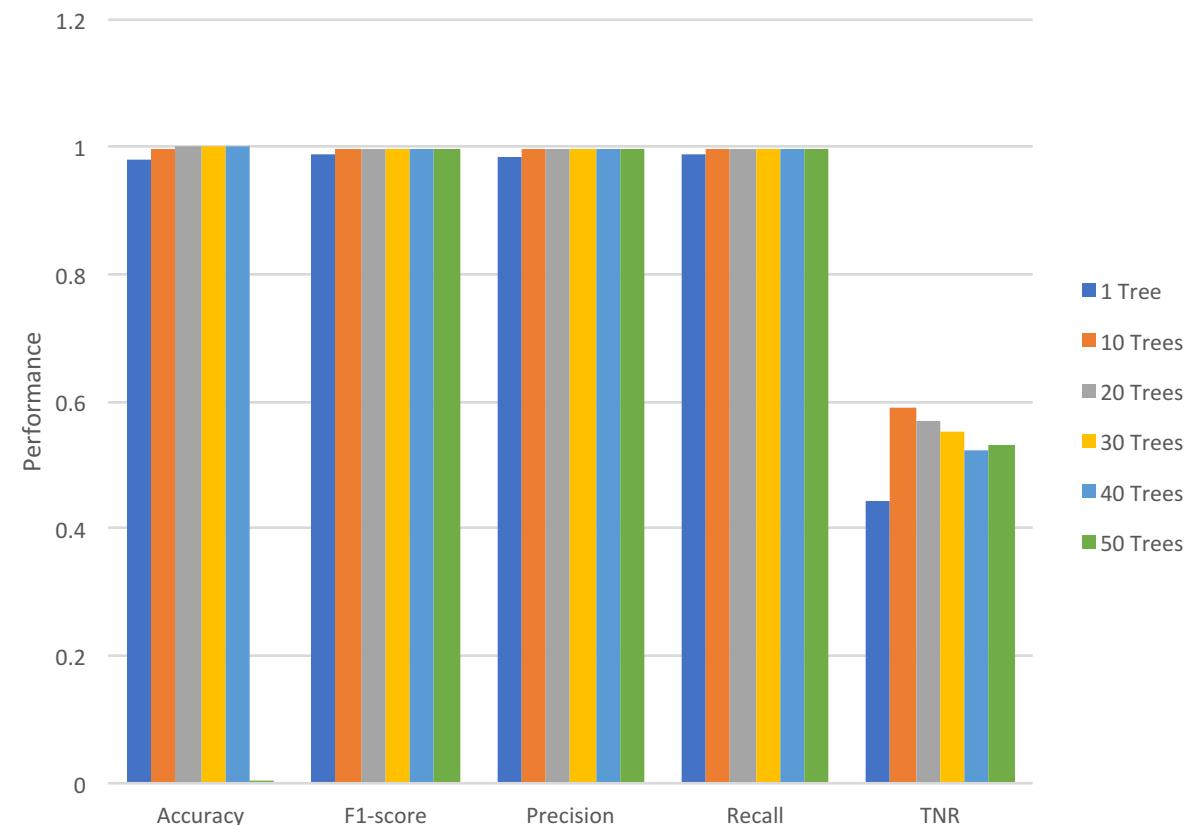


Methods: Global data

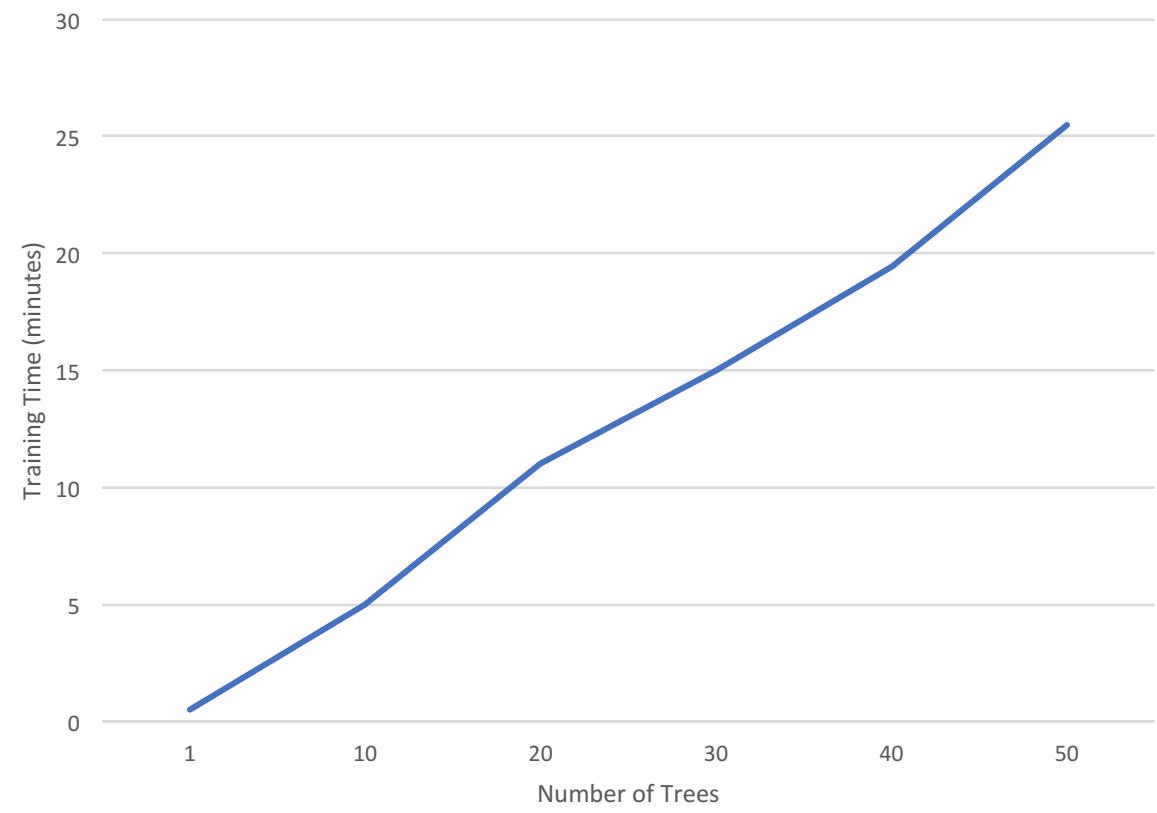


Sequential Results

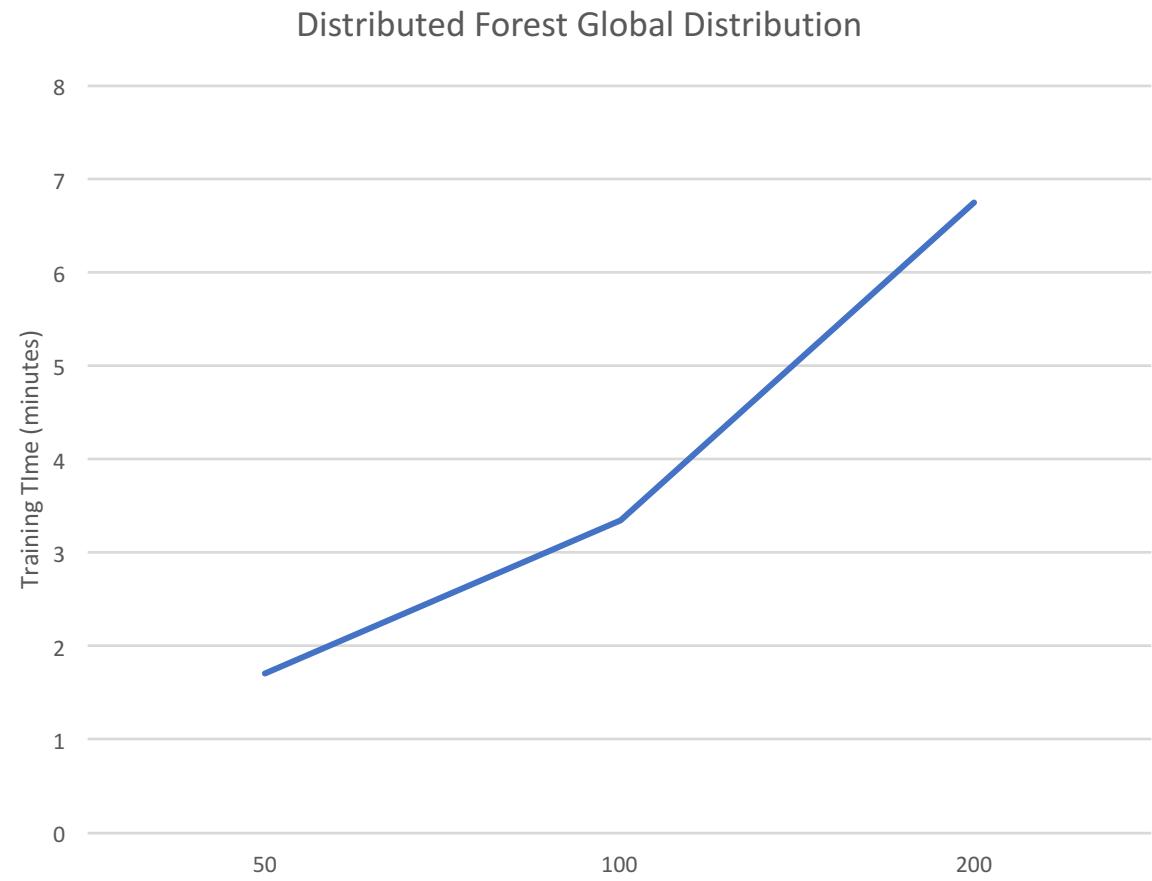
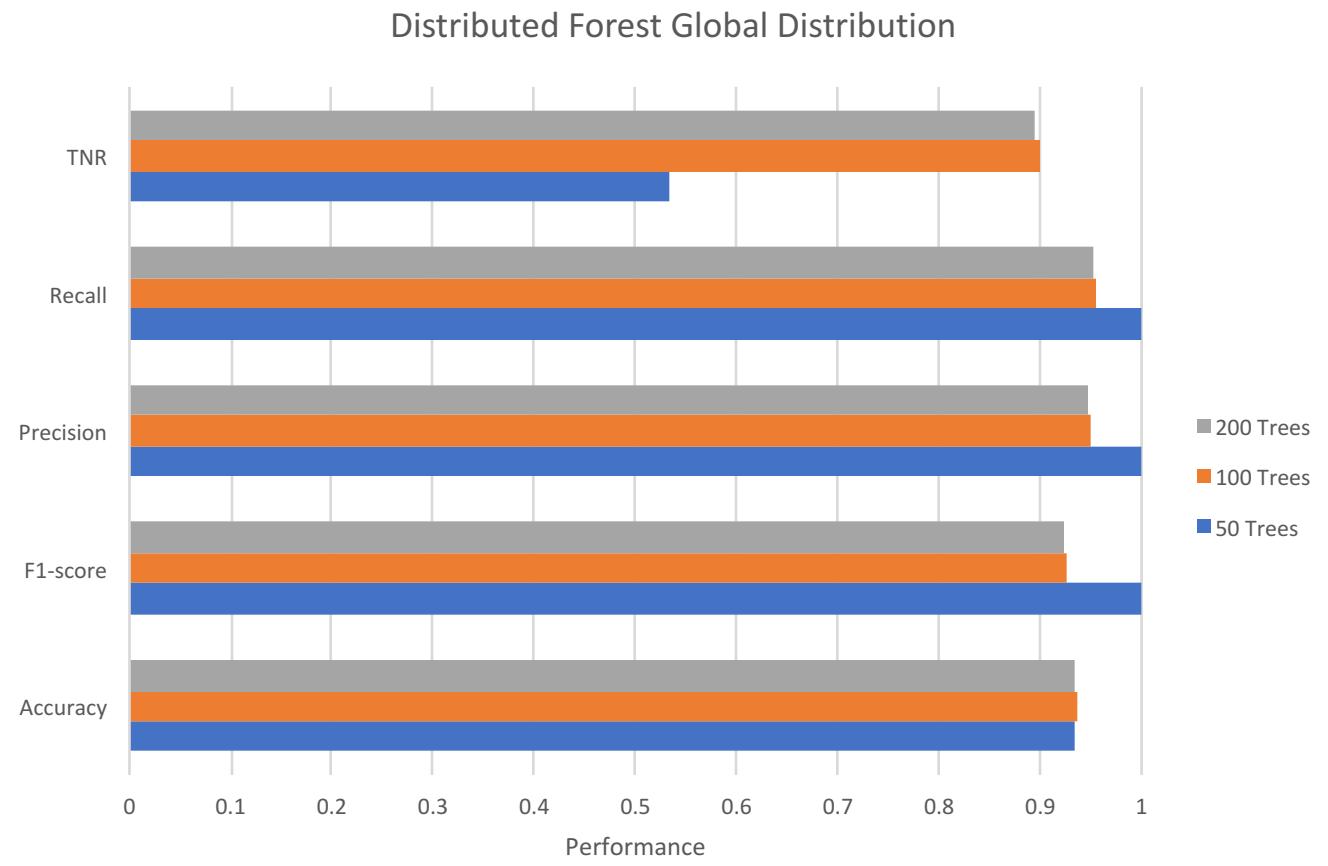
Sequential Random Forest Performance



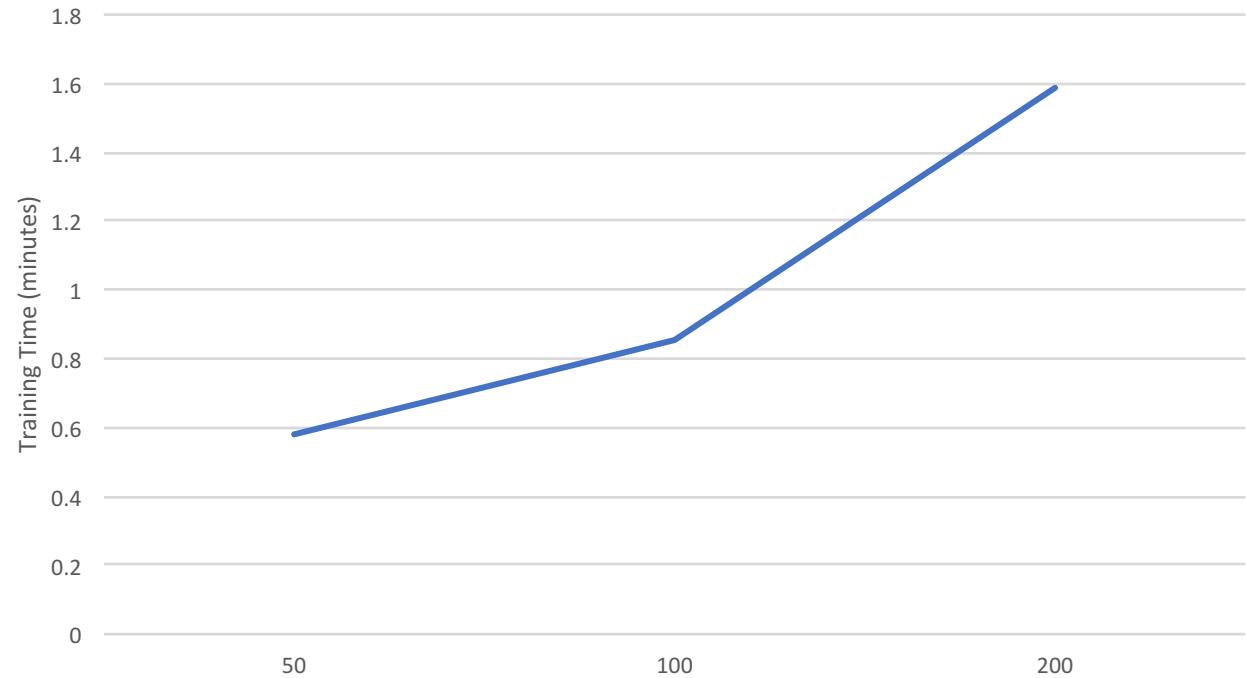
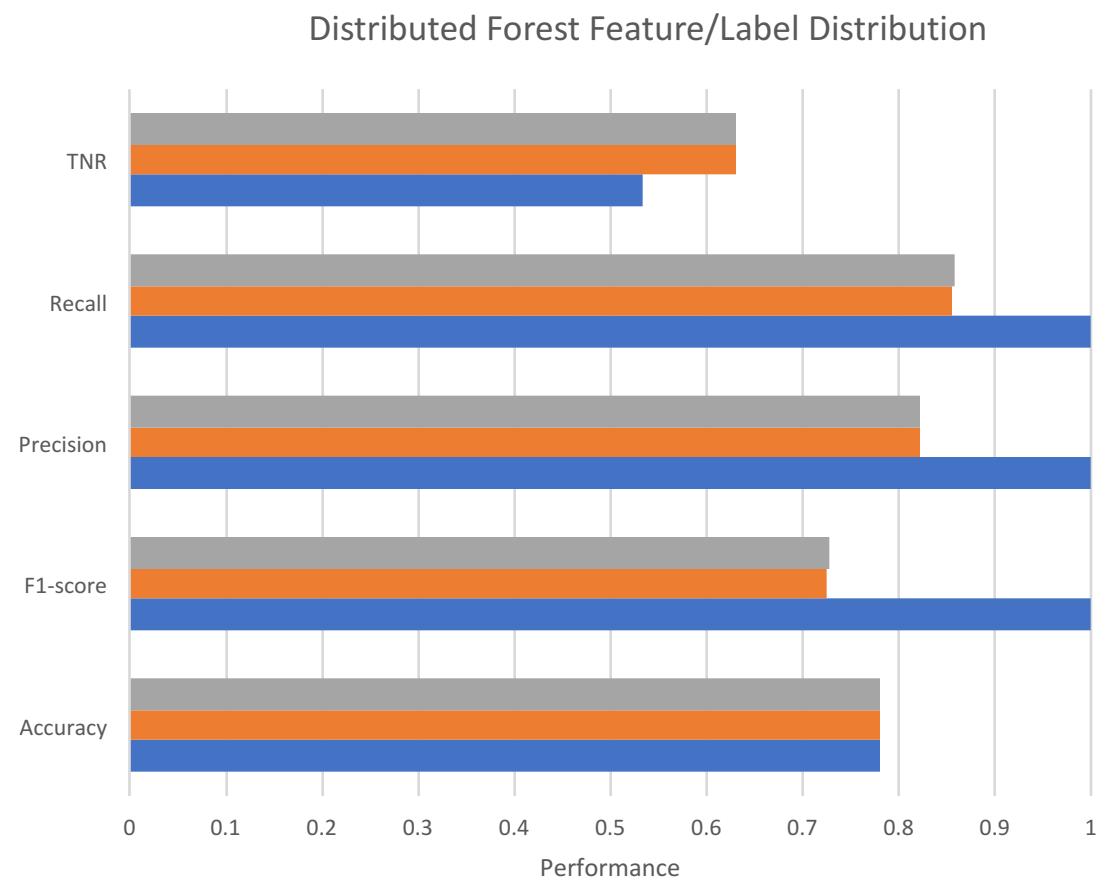
Sequential Random Forest Train Time



Distributed Results

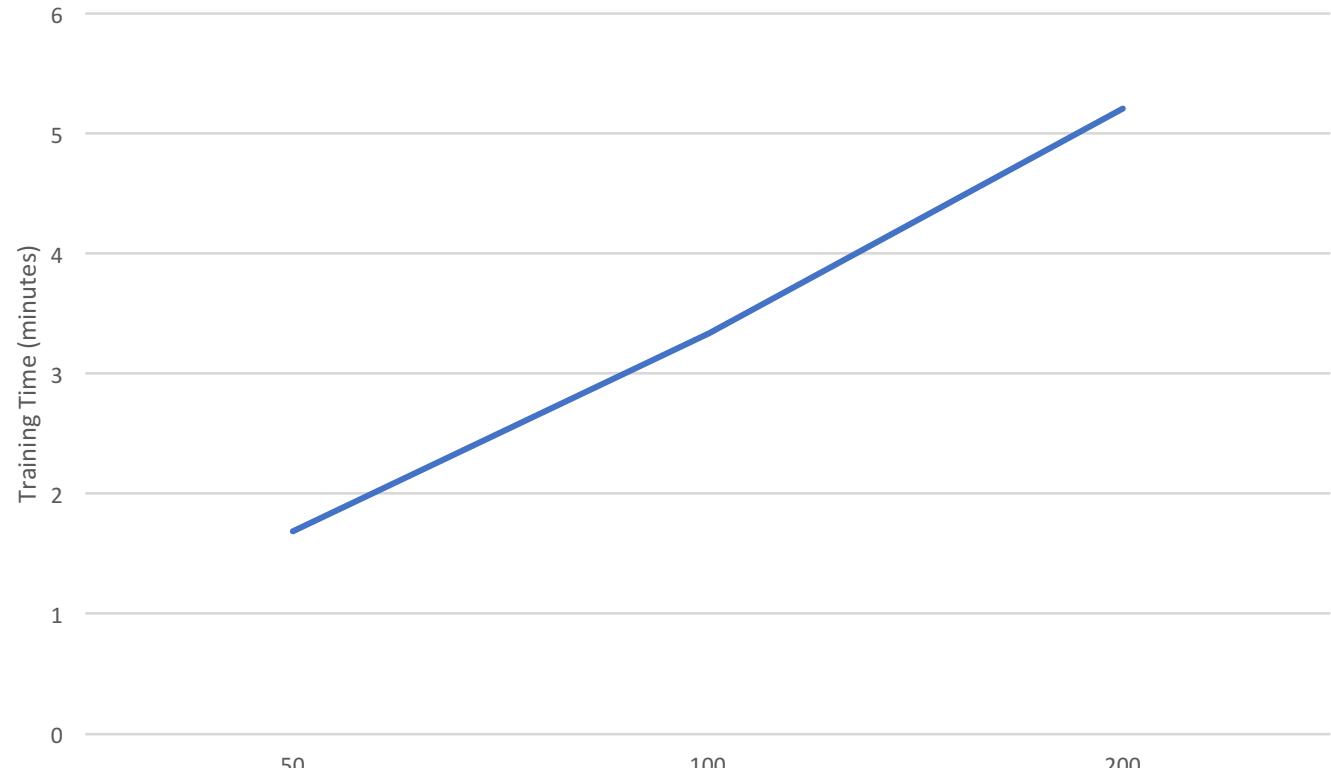


Distributed Forest Feature/Label Distribution

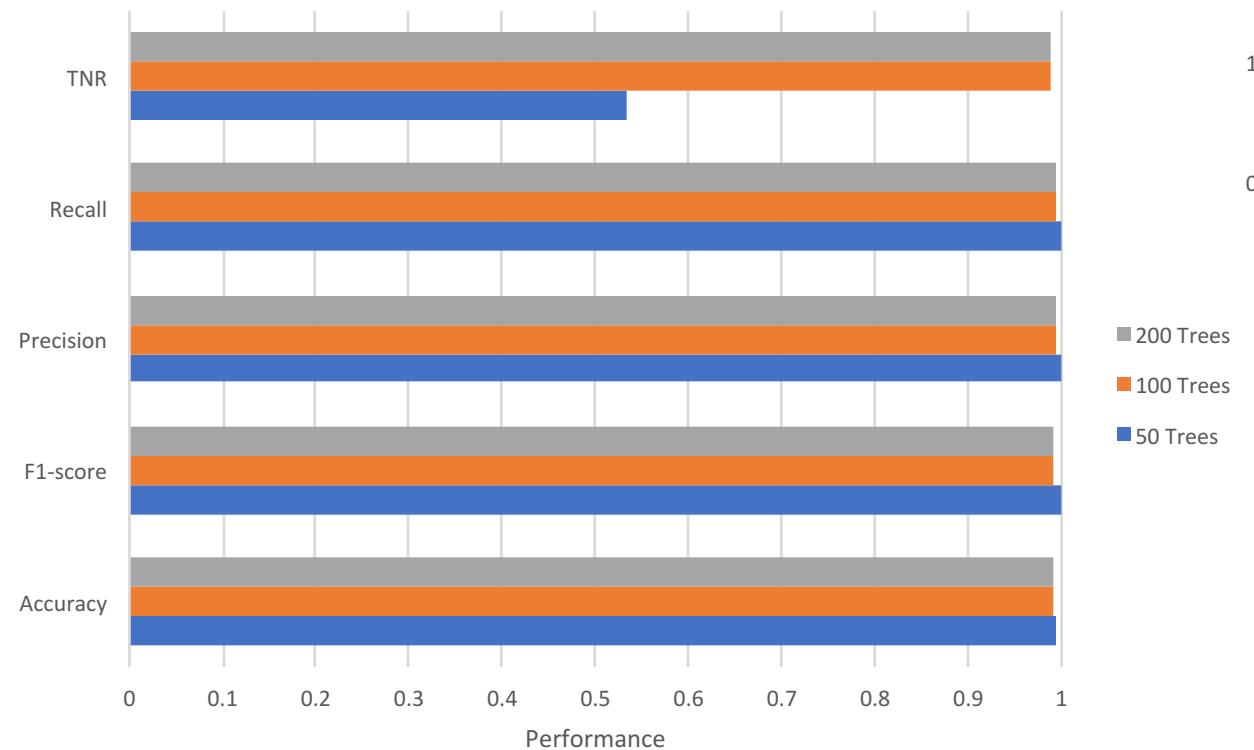


■ 200 Trees
■ 100 Trees
■ 50 Trees

Distributed Forest Label Distribution



Distributed Forest Label Distribution



Conclusions

- Random Forest is robust enough to be able to handle large subsets of trees that are built from different label distributions, however, performance drops dramatically when both feature value and label distributions differ
- There is a significant difference in the time it takes to learn the distributed random forests versus the sequential random forests. The differences that can be observed in the results is even more pronounced when you notice that, at most, the sequential code was only run for 50 trees total while the distributed code was run for, 50, 100, and 200 trees per mapper.
- Therefore, this project shows that in cases where the distribution over the feature values does not drastically differ between the data partitions assigned to each mapper, it is not necessary to maintain a global data model to bootstrap from. However, if both the feature value and label distributions differ significantly, then the degradation in performance indicates that it is necessary to take the extra step of maintaining a global data model to bootstrap each mapper's data partition from.

Distributed Random Forests: Resilient?

Problem Description

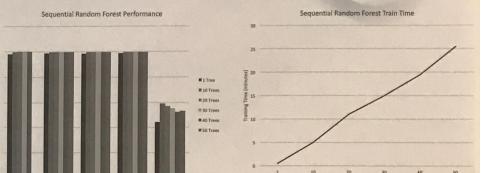
- Explore the effect of learning distributed mini-forests on local data sets that vary according to feature value and label distributions
- Important question b/c passing around data to build partitions with globally representative distributions is computationally expensive and puts the data at higher security risks
- The PAMDAP2 is a time series, activity classification data set:
 - binarized the labels into sedentary and non-sedentary actions
 - 53 continuous feature values from 3 IMU sensors and 1 heart rate monitor
 - > 3M data points
- Split the PAMDAP2 data set among the HARP Mappers by person to create local data sets with distinct feature value and label distributions

Raksha Kumaraswamy
and
Katherine Metcalf

Random Forest

- ensemble classifier
- idea - divide & conquer
- divide - the dataset into bootstrap samples
- conquer - by learning tree for each sample
- gather all trees - a forest, the new model
- prediction - mode, or stochastic prediction based on distribution
- Advantage:**
 - these independently learned models, when aggregated create a resilient forest
 - not prone to bias
- Disadvantage:**
 - computational complexity scaled with data size: many avenues for parallelization
 - bootstrap sample creation, forest creation, tree creation, prediction, pruning
 - prone to variance, but can be controlled with sampling strategies: stratified bootstraps
- What we do** - parallel forest creation, with strategic bag-of-forest bootstraps

Sequential Results



Distributed Results



Conclusions

- Random Forest is robust enough to be able to handle large subsets of trees that are built from different label distributions, however, performance drops dramatically when both feature value and label distributions differ
- There is a significant difference in the time it takes to learn the distributed random forests versus the sequential random forests. The differences that can be observed in the results is even more pronounced when you notice that, at most, the sequential code was only run for 50 trees total while the distributed code was run for, 50, 100, and 200 trees per mapper.
- Therefore, this project shows that in cases where the distribution over the feature values does not drastically differ between the data partitions assigned to each mapper, it is not necessary to maintain a global data model to bootstrap from. However, if both the feature value and label distributions differ significantly, then the degradation in performance indicates that it is necessary to take the extra step of maintaining a global data model to bootstrap each mapper's data partition from.

Distributed Random Forests: Resilient?

Raksha Kumaraswamy
and
Katherine Metcalf

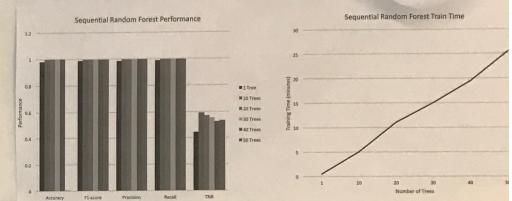
Problem Description

- Explore the effect of learning distributed mini-forests on local data sets that vary according to feature value and label distributions
- Important question b/c passing around data to build partitions with globally representative distributions is computationally expensive and puts the data at higher security risks
- The PAMDAP2 is a time series, activity classification data set:
 - binarized the labels into sedentary and non-sedentary actions
 - 53 continuous feature values from 3 IMU sensors and 1 heart rate monitor
 - > 3M data points
- Split the PAMDAP2 data set among the HARP Mappers by person to create local data sets with distinct feature value and label distributions

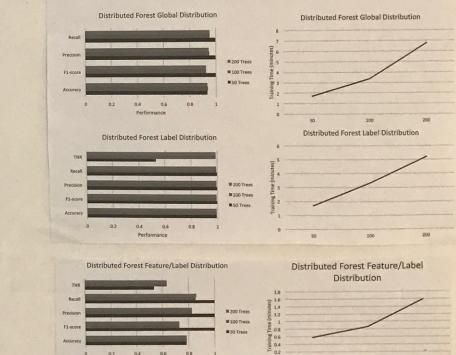
- ensemble classifier
- idea - divide & conquer
- divide - the dataset into bootstrap samples
- conquer - by learning tree for each sample
- gather all trees - a forest, the new model
- prediction - mode, or stochastic prediction based on distribution

- Advantage:**
 - these independently learned models, when aggregated create a resilient forest
 - not prone to bias
- Disadvantage:**
 - computational complexity scaled with data size: many avenues for parallelization
 - bootstrap sample creation, forest creation, tree creation, prediction, pruning
 - prone to variance, but can be controlled with sampling strategies: stratified bootstraps
- What we do** - parallel forest creation, with strategic bag-of-forest bootstraps

Sequential Results



Distributed Results



Conclusions

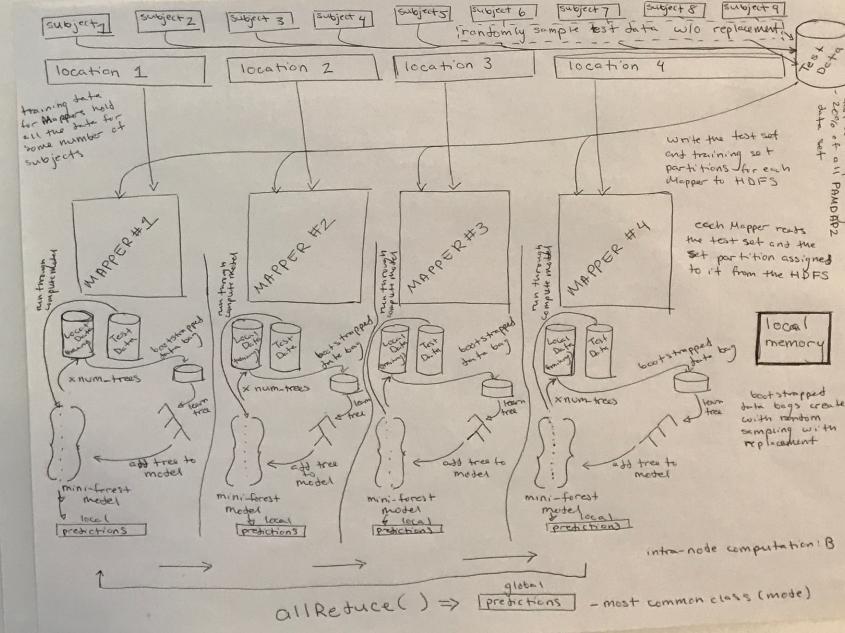
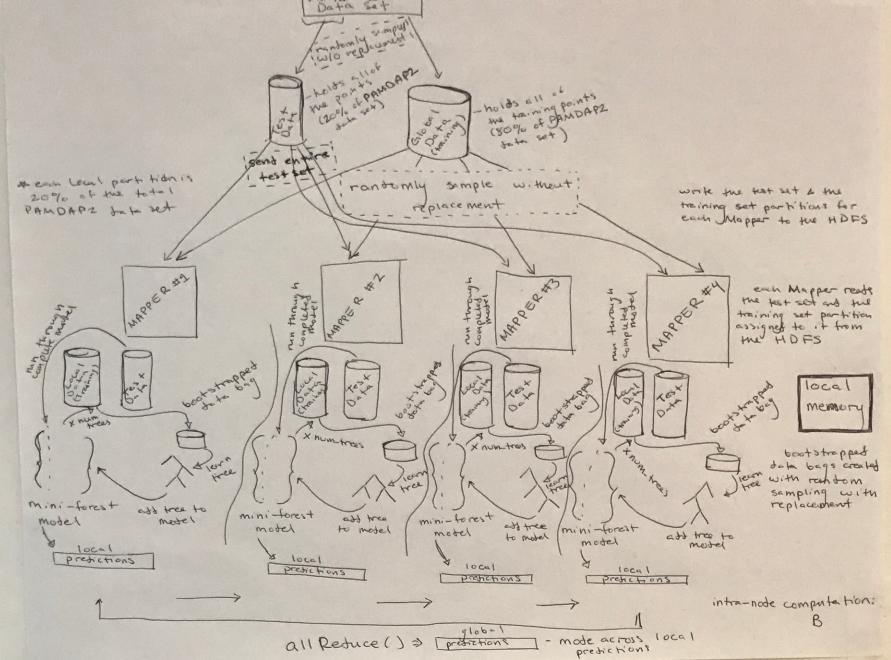
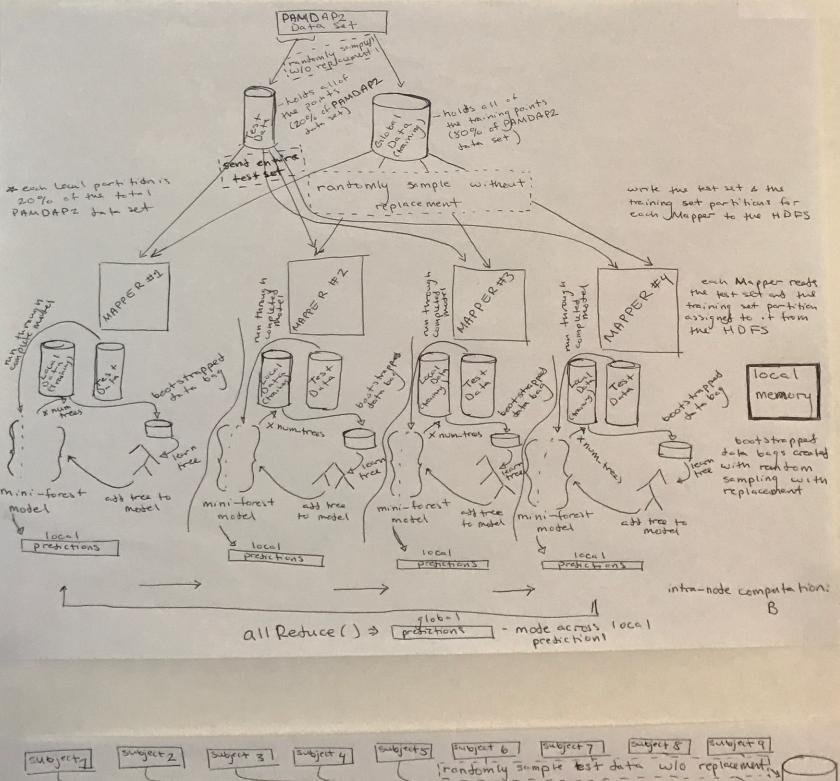
- Random Forest is robust enough to be able to handle large subsets of trees that are built from different label distributions, however, performance drops dramatically when both feature value and label distributions differ
- There is a significant difference in the time it takes to learn the distributed random forests versus the sequential random forests. The differences that can be observed in the results is even more pronounced when you notice that, at most, the sequential code was only run for 50 trees total while the distributed code was run for, 50, 100, and 200 trees per mapper.
- Therefore, this project shows that in cases where the distribution over the feature values does not drastically differ between the data partitions assigned to each mapper, it is not necessary to maintain a global data model to bootstrap from. However, if both the feature value and label distributions differ significantly, then the degradation in performance indicates that it is necessary to take the extra step of maintaining a global data model to bootstrap each mapper's data partition from.

Distributed Random Forest Resilient?

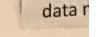
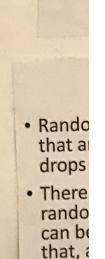
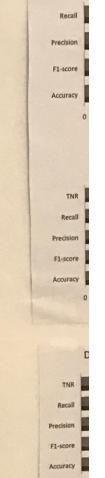
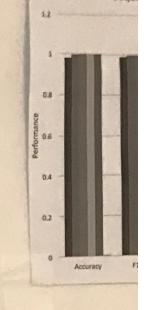
Problem Description

- Explore the effect of learning distributed mini-forests on local data sets that varying according to feature value and label distributions
- Important question b/c passing around data to build partitions with globally representative distributions is computationally expensive and puts the data at higher security risks
- The PAMDAP2 is a time series, activity classification data set:
 - binarized the labels into sedentary and non-sedentary actions
 - 53 continuous feature values from 3 IMU sensors and 1 heart rate monitor
 - > 3M data points
- Split the PAMDAP2 data set among the HARP Mappers by person to create local data sets with distinct feature value and label distributions

Raksha Kumaraswamy
and
Katherine Metcalf



- Rando that all drops
- There rando can be that, a distrik
- There featur assign mode distrik indica data r



nt?

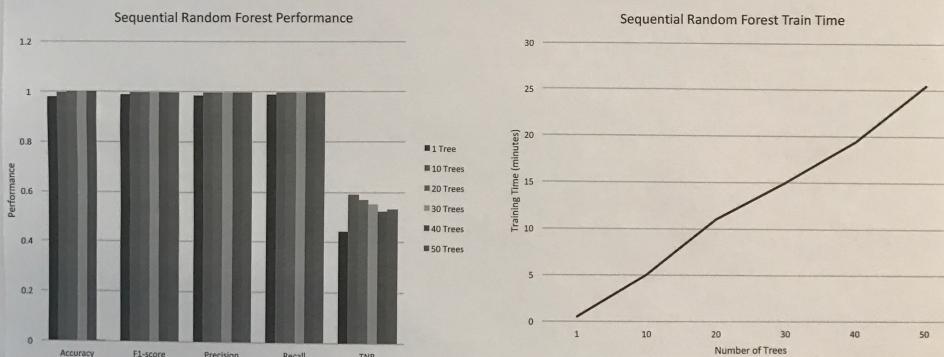
wany

calf

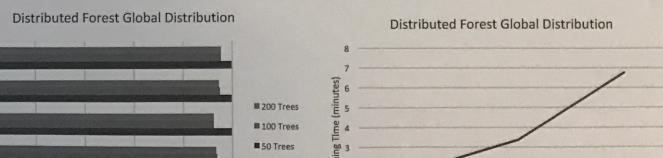
Random Forest

- ensemble classifier
- idea - divide & conquer
- divide - the dataset into bootstrap samples
- conquer - by learning tree for each sample
- gather all trees - a forest, the new model
- prediction - mode, or stochastic prediction based on distribution
- Advantage:
 - these independently learned models, when aggregated create a resilient model
 - not prone to bias
- Disadvantage:
 - computational complexity scaled with data size: many avenues for parallelism
 - bootstrap sample creation, forest creation, tree creation, prediction.
 - prone to variance, but can be controlled with sampling strategies: strategic bag-of-little bootstraps
- What we do - parallel forest creation, with strategic bag-of-little bootstraps

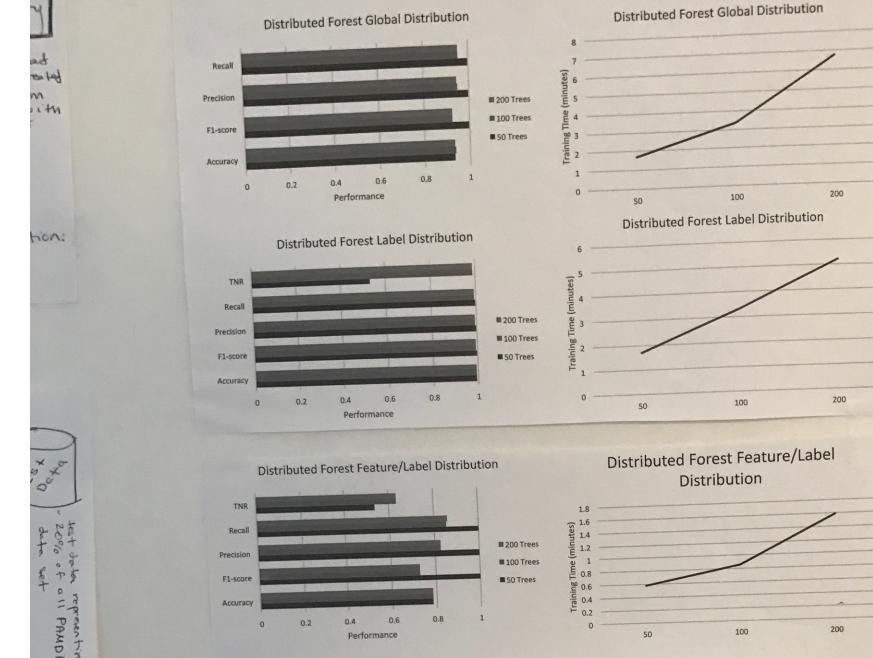
Sequential Results



Distributed Results



Distributed Results



Conclusions

- Random Forest is robust enough to be able to handle large subsets of trees that are built from different label distributions, however, performance drops dramatically when both feature value and label distributions differ
- There is a significant difference in the time it takes to learn the distributed random forests versus the sequential random forests. The differences that can be observed in the results is even more pronounced when you notice that, at most, the sequential code was only run for 50 trees total while the distributed code was run for, 50, 100, and 200 trees per mapper.
- Therefore, this project shows that in cases where the distribution over the feature values does not drastically differ between the data partitions assigned to each mapper, it is not necessary to maintain a global data model to bootstrap from. However, if both the feature value and label distributions differ significantly, then the degradation in performance indicates that it is necessary to take the extra step of maintaining a global data model to bootstrap each mapper's data partition from.