

# A Research Paper on Smashing the Curse of Dimensionality Using Variant RF on Genomics Data

Dr. Raman Chadha (Professor, head), Kaljot Shama (Research Scholar, M.Tech, CSE)  
CGC Technical Campus, Jhanjeri, Mohali Punjab

**Abstract-** Gene prediction one in all the foremost difficult tasks within the order analysis, that several tools are developed and area unit still evolving. Machine learning is used to notice the complicated genomic interactions that may cause malady like polygenic disease or cancer. However, current machine learning approaches area unit unable to subsume these knowledge volumes. We have a tendency to introduce variant spark Random Forests (RF) algorithmic rule to affect huge datasets which has associate ensemble of call trees and to be applied to high dimensional datasets.

**Keywords-** Variantspark RF, machine learning (ML), whole genome sequencing (WGS)

## I. INTRODUCTION

We ar within the inside of the digital revolution wherever shoppers and business demand selections be supported proof collected from information. The ensuing agreement regarding of virtually everything produces datasets that don't seem to be solely growing vertically however conjointly horizontally by capturing a lot of info about these events. The Challenge of massive and "wide" information is very pronounces within the health and bio science house wherever, to Illustrate, whole ordering sequencing technology permits researchers to interrogate all three billion base pairs of the human ordering. genetic science is associate knowledge base field of science that specialize in the structure, function, evolution, mapping and writing of genomes. A ordering is associate organism's complete set of deoxyribonucleic acid , together with of its genes. genetic science aims at the collective characterization and quantification of genes. genetic science conjointly involve the sequencing and analysis of genomes through uses of high output deoxyribonucleic acid sequencing. Whole ordering sequencing (WGS) technology permits researchers to interrogate all. three billion base pairs of the human ordering. As such, the analysis of medical genetic science information is at the forefront of this growing got to apply sophisticates machine learning to giant high dimensional datasets .A common task during this field is to spot illness genes, that's wherever tiny errors within the cistron sequence have had a prejudicious health impact reminiscent of cancer. However, Biology is far a lot of difficult than that. thus we tend to highlight the variant spark RF overcome the matter of curse of dimensionality". Machine language could be a branch of AI that offers computers the power to be told new patterns with very little to no human intervention. The machine learning models learn from

previous computations to provide a lot of correct results as a lot of information is fragmentize Straightforward example is Facebook is face detection formula , that uses machine learning techniques to spot the individuals within the photos and gets refined over time. Machine Learning ways, especially Random Forest (RF), on the opposite hand ar similar temperament to spot sets of options (e.g., mutations or a lot of usually variants) that ar prophetic or related to a label (e.g., disease). the most action of Variant Spark RF is to alter — for the primary time — the identification of illness inflicting variants by taking the higher-order genome-wide interactions between genomic loci under consideration. for example the advantage of this new powerful analysis, we tend to created an artificial dataset that simulates the mechanics of a fancy illness or makeup. we tend to decision this artificial affliction "Hipsterism." to form this, we tend to initial known peer-reviewed and revealed traits, reminiscent of propensity for facial hair or higher occasional consumption, that ar ordinarily related to being a reformist. we tend to then score every individual within the one thousand Genomes Project dataset with the formula below, that joins info from these genome-wide locations in a very equally non-purely-additive means as a true complicated makeup would:

$$\text{HipsterIndex} = ((2 + \text{GT}[\text{B6}]) * (1.5 + \text{GT}[\text{R1}])) + ((0.5 + \text{GT}[\text{C2}]) * (1 + \text{GT}[\text{B2}]))$$

Where GT stands for the genotype at this position with homozygous reference encoded as 0, heterozygote as 1, and homozygote alternative as 2.

We then label individuals with a score above 10 as being a Hipster. The genomic information from all individuals with the synthetic Hipster label was then used to train Variant Spark RF to find the features that are most predictive or associated with this synthetic Hipster-phenotype. VariantSpark RF was able to correctly identify the 4 correct locations purely from the given Hipster label. Not only that, it identified the location in order of their exact weighting in the score's formula, which similar tools were not able to achieve (Fig 1).

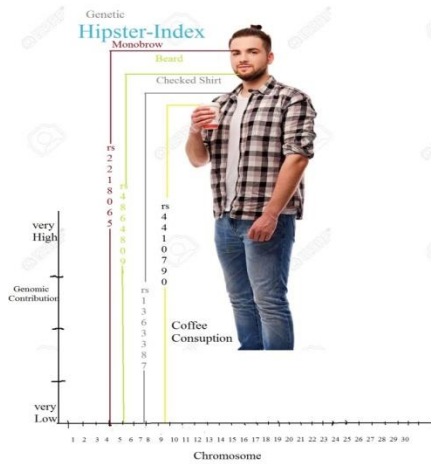
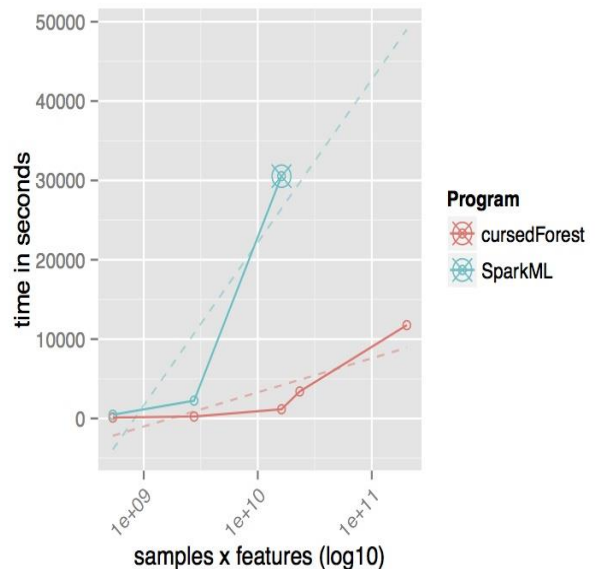
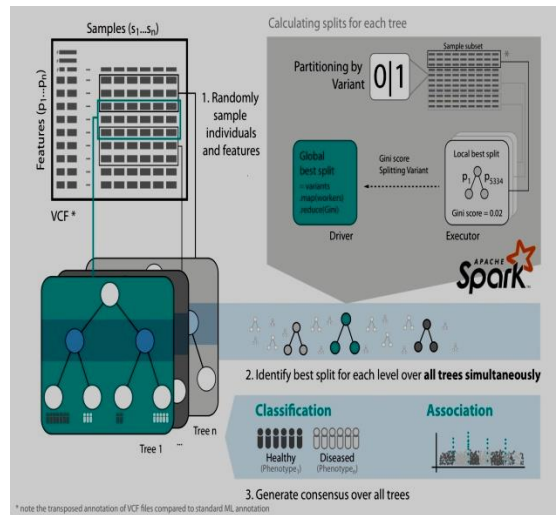


Fig.1:Synthetic phenotype demonstrating the ability of VariantSpark RF to identify sets of variants contributing to the phenotype even in a non-additive way.

We needed to re-evaluate this dataset victimisation the random forest implementation in Spark MLlib to enhance the accuracy victimisation supervised learning. However, MLlib’s RF wasn’t ready to method the whole dataset Associate in Nursingd ran out of memory on an on-premise Hadoop-cluster with twelve Executors (16 Intel Xeon E5-2660@2.20GHz processor cores and 128GB of RAM) for even atiny low set of the initial knowledge (2,504 samples half dozen,450,364 features). RF has conjointly a reduced risk of over fitting compared to different machine learning ways, that is crucial for things wherever the dataset has more options than samples. These things suffer from the “curse of spatiality,” and RF overcomes this by building freelance. VariantSparkRF starts by indiscriminately distribution subsets of the info to Spark Executors for call tree building (Fig 2). It then calculates the simplest split over all nodes and trees at the same time. This implementation avoids communication bottlenecks between Spark Driver and Executors as data exchange is least, permitting it to create massive numbers of trees with efficiency. This surveys the answer area fitly to cater for immeasurable options and thousands of samples. Furthermore, Variant Spark RF has memory efficient representation of genomics data, optimized communication patterns and computation batching. It also provides efficient implementation of Out-Of-Bag (OOB) error, which substantially simplifies parameter tuning over the computationally more costly alternative of cross-validation. We implemented Variant Spark RF in scala as it is the most performance interface languages to Apache Spark. Also, new updates to Spark and the interacting APIs will be deployed in scala first, which has been important when working on top of a fast evolving framework.

In 2015, CSIRO developed the first version of VariantSpark, which was limited to unsupervised clustering and was built on top of Spark MLlib. we clustered

individuals from the 1000 Genomes Project to identify their ethnicity. This dataset contained ~2,500 individuals with ~80 Million genomic variants and we achieved a correct prediction rate of 82% (accuracy). We wanted to re-evaluate this dataset using the random forest implementation in Spark MLlib to improve the accuracy using supervised learning. However, MLlib’s RF was not able to process the entire dataset and ran out of memory on an on-premise Hadoop-cluster with 12 Executors (16 Intel Xeon E5-2660@2.20GHz CPU cores and 128GB of RAM) for even a small subset of the original data (2,504 samples 6,450,364 features).In contrast, the new version of VariantSpark, which implements the RF with a novel parallelization algorithm built on Spark Core directly, was able to process the entire dataset using the same cluster setup. It is processing over 15 Million variants per second from the 202 Billion variants in the dataset and was finishing in 3 hours. Being the only method to use the whole dataset, VariantSpark RF achieved a higher accuracy 0.96 (OOB=0.02) (Fig 3).



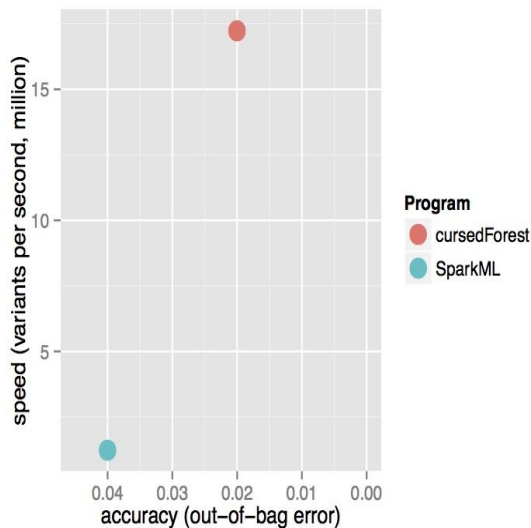


Fig.3: Performance comparison of VariantSpark RF and Spark MLlib on 1000 Genomes Project dataset. **Top** Runtime over different dataset sizes, cross marks the last dataset analyzed successfully. **Bottom** Accuracy achieved at termination point and number of features processed per second.

## II. RESULTS INTERPRETATION

**Variant Spark** is a machine learning library for real-time genomic data analysis (for thousands of samples and millions of variants) and is...

- Built on top of Apache Spark and written in Scala
  - Uses a custom machine learning **random forest** implementation to find the most *important* variants attributing to a phenotype of interest
  - Includes a dataset with a subset of the samples and variants (in VCF format) from the 1000 Genomes Project
1. **chr2\_223034082** (rs2218065) encoding for monobrow is the most important feature.
  2. A group of SNPs encoding for the MEGF10 gene (**chr5\_126626044**), which is involved in Retina horizontal cell formation as the second most important marker, explaining why hipsters prefer checked shirts.
  3. **chr7\_17284577** (rs4410790) the marker for increased coffee consumption is ranked third. **chr4\_54511913** (rs4864809) the marker for beards is fourth.

## III. CONCLUSION

Variant Spark RF, a brand new library that enables random forest to be applied to high-dimensional datasets. The novel Spark-based parallelization permits an oversized range of trees to be designed at the same time, thence enabling the answer area to be searched a lot of thoroughly than different strategies. whereas genetic science is presently the discipline manufacturing the most important volumes of complicated information, the continuing information fiction can bring similar analysis challenges to different disciplines. VariantSpark RF might thence be capable of changing these challenges to opportunities on those disciplines also.

## IV. REFERENCES

- [1]. A, Katal, Wazid M, and Goudar R.H. "Big data: Issues, challenges, tools and Good practices." Noida: 2013, pp. 404 – 409, 8-10 Aug. 2013.
- [2]. Golfarelli, M., & Rizzi, S. (2009). Data warehouse design: modern principles and methodologies. Columbus: McGraw-Hill
- [3]. Almeida, F., and Calistru, C, "The Main Challenges and Issues of Big Data Management", International Journal of Research Studies in Computing, 2(1), 2013, pp. 11-20.
- [4]. Chen, S. Mao, and Y. Liu, "Big data: a survey," *Mobile Networks and Applications*, vol. 19, no. 2, pp.
- [5]. Frost & Sullivan: Global Precision Market Growth Opportunities, Forecast to 2015 2017
- [6]. O'Brien et al. BMC Genomics 2015
- [7]. <https://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-015-2269-7>
- [8]. <https://github.com/bigdatagenomics/adam>
- [9]. [https://aehrc.github.io/VariantSpark/notebook-examples/VariantSpark\\_HipsterIndex.html](https://aehrc.github.io/VariantSpark/notebook-examples/VariantSpark_HipsterIndex.html)