
Relation Augmentation: A Gradient Boosting Approach for Detecting Genomic Anomalies

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Genomic anomalies, or variations, are often shared between members of the same
2 species. Although rare, these changes may result in disease or an increase in host
3 fitness. Most approaches for detecting structural variation rely on high quality
4 data and are typically limited to one type of structural variant such as deletions or
5 inversions. We propose a new data augmentation method to mitigate errors in low
6 quality DNA sequencing data by leveraging offspring DNA information to predict
7 genomic variants in their associated parent. To our knowledge, this is the first time
8 such an approach has been proposed to detect multiple structural variation classes,
9 including complex variants, using related individual data. The author(s) of this
10 work identify as Latinx.

11 1 Introduction - Biological Motivation

12 Advances in DNA sequencing have increased the number of large sequencing studies, with the goal of
13 quantifying genomic variation and its influence on both genotypes and phenotypes in species [1, 2, 3].
14 These genomic anomalies may appear as a single basepair change (SNP) or as a rearrangement
15 of a larger region. Detecting these larger regions, known as structural variants (SVs), remains a
16 challenging problem. This is particularly true when incorporating low quality DNA sequencing data
17 [4, 5, 6]. The ability to detect such genomic variation remains an important area of study, as these
18 changes have applications in detecting negative and positive outcomes (e.g. cancer susceptibility and
19 increased fitness) [7, 8, 9, 10].

20 Of particular interest is being able to detect complex structural variants, including deletions, inversion,
21 translocations, or duplications [11]. This typically involves aligning sample DNA to a known refer-
22 ence genome and looking for differences. In certain cases, these variations may occur simultaneously,
23 resulting in complex variants (see Fig. 1) [12]. Our work builds upon previous methods by using
24 related individuals – in this case, offspring data – as a data augmentation method [13, 14, 15, 16, 17].

25 2 Method

26 Our method was designed with the intention of identifying the various structural variants from
27 features provided from aligned sequencing data using `samttools` [18]. Rather than identifying all of
28 the associated structural variants using a single model, each model was trained to identify each variant
29 individually. This required converting a multi-class problem into a series of binary classification
30 problems. By focusing each model on a single class, the capability of the model is more robust to the
31 task and will not suffer because of the severe class imbalances present in other classes. After testing
32 on various ensemble methods, the optimal method for classification across all classes was gradient
33 boosting. Gradient boosting is an additive method that seeks to minimize a given cost function by
34 incorporating new decision trees to compensate for the shortcomings of previous decision trees. New

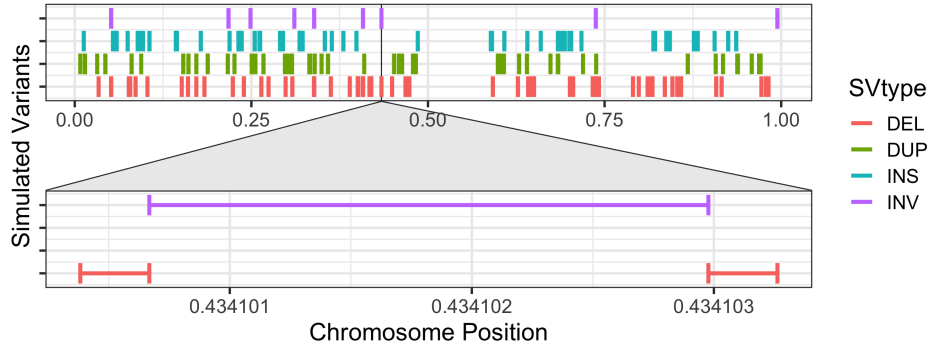


Figure 1: *Top.* We simulate the parent signal by introducing a set of 169 structural variants (deletions, duplications, insertions, and inversions) in Chromosome 1 of the human genome reference (relative position shown). *Bottom.* We illustrate inversion-deletions – known as complex structural variants – where variation occurs simultaneously and an inversion event is flanked by two small deletion events.

35 trees are added to the ensemble using a form of functional gradient descent where the new tree's
 36 parameters are chosen to minimize the loss of the cost function [19].

37 2.1 Simulated Data

38 We first simulate a number of structural variants in an individual and then try to predict the structural
 39 variants introduced in the simulated data. Moreover, we wish to correctly identify when there
 40 are complex rearrangements introduced with low-quality data. The parent signal was simulated
 41 using SURVIVOR with 50 (of each) duplications, insertions, and deletions, 4 inversions, and 5
 42 inversion-deletion events [20]. This tool was also used to simulate 8 offspring based on this parent.

43 3 Results

44 Four models were trained on aligned, low quality* DNA sequencing offspring data and then tested on
 45 the parent sequencing data. Fig. 2 demonstrates the effectiveness of our method and we note that
 46 as the number of offspring increases, the area under the curve (AUC) increases with the number of
 47 offspring. Although our method only uses discordant data, sequencing data that did not match to the
 48 reference, our model was able to correctly classify 13 insertion and 6 inversion (3 regular, 3 complex)
 49 events, while Lumpy detected 0 insertions and the same 6 inversions. For future work, we plan on
 50 incorporating more features with concordant and split-read sequencing data for all individuals. We
 51 also plan on taking advantage of stacking (an ensemble learning technique) to improve classification.

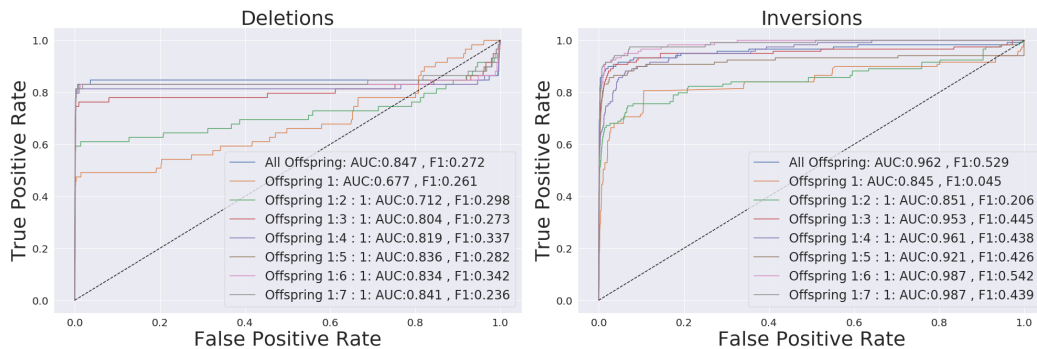


Figure 2: *Left.* Receiver operating characteristics (ROC) curve for deletions for the parent signal using 1 through 8 offspring sequencing data. *Right.* ROC curve for inversions for the parent signal training our model with 1 through 8 offspring sequencing data. Both figures include the AUC and F1 score as performance metrics.

52 References

- 53 [1] . G. P. Consortium *et al.*, “An integrated map of genetic variation from 1,092 human genomes,” *Nature*,
54 vol. 491, no. 7422, pp. 56–65, 2012.
- 55 [2] J.-Y. Li, J. Wang, and R. S. Zeigler, “The 3,000 rice genomes project: new opportunities and challenges for
56 future rice research,” *GigaScience*, vol. 3, no. 1, pp. 1–3, 2014.
- 57 [3] Z. R. Chalmers, C. F. Connelly, D. Fabrizio, L. Gay, S. M. Ali, R. Ennis, A. Schrock, B. Campbell,
58 A. Shlien, J. Chmielecki, *et al.*, “Analysis of 100,000 human cancer genomes reveals the landscape of
59 tumor mutational burden,” *Genome medicine*, vol. 9, no. 1, p. 34, 2017.
- 60 [4] E. P. Consortium *et al.*, “Identification and analysis of functional elements in 1% of the human genome by
61 the encode pilot project,” *Nature*, vol. 447, no. 7146, p. 799, 2007.
- 62 [5] C. Alkan, B. P. Coe, and E. E. Eichler, “Genome structural variation discovery and genotyping,” *Nature
63 Reviews Genetics*, vol. 12, no. 5, p. 363, 2011.
- 64 [6] J. R. MacDonald, R. Ziman, R. K. C. Yuen, L. Feuk, and S. W. Scherer, “The database of genomic variants:
65 a curated collection of structural variation in the human genome,” *Nucleic acids research*, vol. 42, no. D1,
66 pp. D986–D992, 2013.
- 67 [7] J. Weischenfeldt, O. Symmons, F. Spitz, and J. O. Korb, “Phenotypic impact of genomic structural
68 variation: insights from and for human disease,” *Nature Reviews Genetics*, vol. 14, no. 2, p. 125, 2013.
- 69 [8] I. Martincorena and P. J. Campbell, “Somatic mutation in cancer and normal cells,” *Science*, vol. 349,
70 no. 6255, pp. 1483–1489, 2015.
- 71 [9] C. Jeong, D. B. Witonsky, B. Basnyat, M. Neupane, C. M. Beall, G. Childs, S. R. Craig, J. Novembre,
72 and A. Di Rienzo, “Detecting past and ongoing natural selection among ethnically tibetan women at high
73 altitude in nepal,” *PLOS Genetics*, vol. 14, pp. 1–30, 09 2018.
- 74 [10] G. A. Gnecci-Ruscone, P. Abondio, S. De Fanti, S. Sarno, M. G. Sherpa, P. T. Sherpa, G. Marinelli,
75 L. Natali, M. Di Marcello, D. Peluzzi, *et al.*, “Evidence of polygenic adaptation to high altitude from
76 tibetan and sherpa genomes,” *Genome biology and evolution*, vol. 10, no. 11, pp. 2919–2930, 2018.
- 77 [11] P. Stankiewicz and J. R. Lupski, “Structural variation in the human genome and its role in disease,” *Annual
78 review of medicine*, vol. 61, pp. 437–455, 2010.
- 79 [12] Z. Stephens, C. Wang, R. K. Iyer, and J.-P. Kocher, “Detection and visualization of complex structural
80 variants from long reads,” *BMC bioinformatics*, vol. 19, no. 20, p. 508, 2018.
- 81 [13] T. Rausch, T. Zichner, A. Schlattl, A. M. Stütz, V. Benes, and J. O. Korb, “Delly: structural variant
82 discovery by integrated paired-end and split-read analysis,” *Bioinformatics*, vol. 28, no. 18, pp. i333–i339,
83 2012.
- 84 [14] R. M. Layer, C. Chiang, A. R. Quinlan, and I. M. Hall, “Lumpy: a probabilistic framework for structural
85 variant discovery,” *Genome biology*, vol. 15, no. 6, p. R84, 2014.
- 86 [15] S. S. Sindi and B. J. Raphael, “Identification of structural variation,” *Genome Analysis: Current Procedures
87 and Applications*, p. 1, 2014.
- 88 [16] M. Banuelos, R. Almanza, L. Adhikari, S. Sindi, and R. F. Marcia, “Sparse signal recovery methods
89 for variant detection in next-generation sequencing data,” 2016. Proceedings of the *IEEE International
90 Conference on Acoustics, Speech and Signal Processing*.
- 91 [17] M. Spence, M. Banuelos, R. F. Marcia, and S. Sindi, “Detecting inherited and novel structural variants in
92 low-coverage parent-child sequencing data,” *Methods*, 2019.
- 93 [18] H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, and R. Durbin,
94 “The sequence alignment/map format and samtools,” *Bioinformatics*, vol. 25, no. 16, pp. 2078–2079, 2009.
- 95 [19] J. H. Friedman, “Greedy function approximation: a gradient boosting machine,” *Annals of statistics*,
96 pp. 1189–1232, 2001.
- 97 [20] D. C. Jeffares, C. Jolly, M. Hoti, D. Speed, L. Shaw, C. Rallis, F. Balloux, C. Dessimoz, J. Bähler, and
98 F. J. Sedlazeck, “Transient structural variations have strong effects on quantitative traits and reproductive
99 isolation in fission yeast,” *Nature communications*, vol. 8, p. 14061, 2017.