
Chapter 1

NIPTeR: an R package for fast and accurate trisomy prediction in non-invasive prenatal testing

BMC Bioinformatics 2018;19:531.

DOI: 10.1186/s12859-018-2557-8

PubMed ID: 30558531

L.F. Johansson^{1,2}, H.A. de Weerd^{1,2,3}, E.N. de Boer¹, F. van Dijk^{1,2},
G.J. te Meerman¹, R.H. Sijmons¹, B. Sikkema-Raddatz¹, M.A. Swertz^{1,2}

1. University of Groningen, University Medical Center Groningen, Department of Genetics, Groningen, The Netherlands

2. University of Groningen, University Medical Center Groningen, Genomics Coordination Center, Groningen, The Netherlands

3. School of Bioscience, Systems biology research center, University of Skövde, Skövde, Sweden

Received 2018 Oct 2; Accepted 2018 Dec 4; Published online 2018 Dec 17.

Abstract

Background Various algorithms have been developed to predict fetal trisomies using cell-free DNA in non-invasive prenatal testing (NIPT). As basis for prediction, a control group of non-trisomy samples is needed. Prediction accuracy is dependent on the characteristics of this group and can be improved by reducing variability between samples and by ensuring the control group is representative for the sample analyzed.

Results NIPTer is an open-source R Package that enables fast NIPT analysis and simple but flexible workflow creation, including variation reduction, trisomy prediction algorithms and quality control. This broad range of functions allows users to account for variability in NIPT data, calculate control group statistics and predict the presence of trisomies.

Conclusion NIPTer supports laboratories processing next-generation sequencing data for NIPT in assessing data quality and determining whether a fetal trisomy is present. NIPTer is available under the GNU LGPL v3 license and can be freely downloaded from <https://github.com/molgenis/NIPTer> or CRAN.

1.1 Background

Non-invasive prenatal testing (NIPT) is rapidly becoming the new standard in prenatal screening for fetal aneuploidy [1]. In NIPT, cell-free DNA from the pregnant woman's blood plasma, which consists of both maternal and fetal DNA fragments, is analysed. Next to SNP-based methods [6], low-coverage whole genome next-generation sequencing (NGS) is often used [4, 14], and various algorithms, software programs and packages have been developed to analyse this type of data [3, 15, 16, 13, 12]. In literature, many methods have been described that depend on a statistical comparison between a sample of interest and a reference set of non-trisomy control samples [4, 14, 5, 7]. The RAPIDR and DASAF R packages, for instance, have been described [11, 10] and they made several of these algorithms available, including GC-correction, the standard Z-score and the Normalized Chromosome Value (NCV), to create an analysis workflow in R. However, those packages lack features like chi-squared-based variation reduction (χ^2 VR), regression-based Z-score (RBZ) and Match QC. These are all algorithms that we have extensively discussed before [7]. In short, χ^2 VR detects chromosomal regions that have a higher variability than expected by chance and reduces their weight so that, after correction, they have less impact on the fraction of reads mapped to the different chromosomes. The RBZ is an alternative Z-score calculation based on stepwise regression with forward selection. In the RBZ positive or negative correlation between chromosomal fractions is used to predict the number of reads to map onto the chromosome of interest if no trisomy is present. The Match QC score is a sum-of-squares-based approach to compare chromosomal fractions between the test sample and controls, and it provides a measure by which to determine whether a control group is representative for a specific sample. Here we report NIPTeR, an R package that provides fast NIPT analysis for research and diagnostics and provides users with multiple methods for variation reduction, prediction and quality control based upon comparison of a sample with a set of negative control samples.

1.2 Implementation

NIPTeR users can create different workflows for variation reduction and aneuploidy prediction using thirteen functions as building blocks (Fig. ??). A stepwise practical example for using these building blocks is presented as a case report in Additional file 1.

NIPTeR analysis uses two core objects. The first object is NIPTSample, which contains the counts of aligned sequence reads in 50,000 bp bins for a specific sample. The second object is NIPTControlGroup, which contains a series of NIPTSamples for comparison. Users generate NIPTSample using the function `bin_bam_sample`, which needs a BAM file [9] as input. The user can optionally select to count reads mapped to the forward and reverse strands separately, so that they can each be used as a separate predictor. The `as_control_group` function converts a series of NIPTSample objects into a NIPTControlGroup. Within NIPTeR, users can manage an existing NIPTControlGroup using the `add_samples_controlgroup`, `remove_sample_controlgroup` and `remove_duplicates_controlgroup` functions.

Both NIPTSample and NIPTControlGroup can undergo one or more variation reduction steps to adjust the bin read counts, either using the `gc_correct` function for weighted bin GC correction [5] or LOESS GC correction [2] or the `chi_correct` function for χ^2 VR. Each NIPTSample object shows the correction status for the autosomes and the sex chromosomes separately and indicates which variation reduction methods have been performed (or that they are 'uncorrected'). χ^2 VR can be applied to uncorrected or GC-corrected samples, and makes use of a NIPTSample and a NIPTControlGroup having an identical correction status.

Using the fractions of reads mapped to the different chromosomes, trisomy prediction can be generated for a given NIPTSample based on the NIPTControlGroup using three different prediction algorithms: (1) `calculate_z_score`, which uses a standard Z-score [4]; (2) `calculate_ncv_score`, which uses an NCV [14]; and (3) `perform_regression`, which uses RBZ. All three trisomy prediction functions use NIPTControlGroup to calculate the expected fraction of reads on the chromosome of interest. For NCV, this calculation is done in a separate function, `prepare_ncv`, because the calculation is time-intensive and only has to

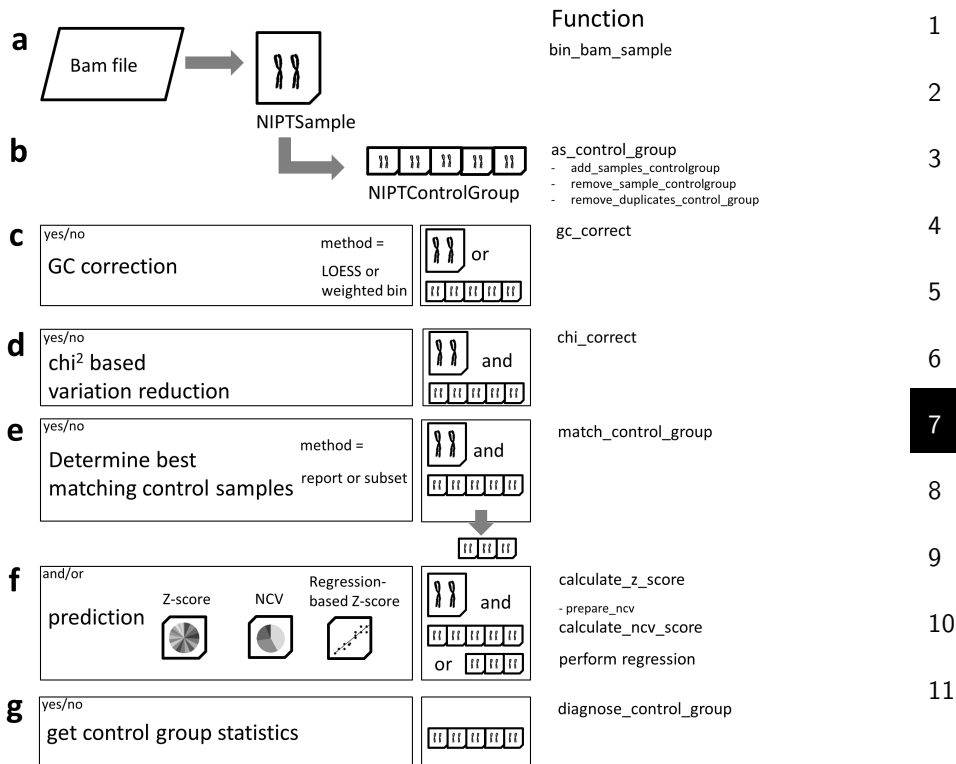


Figure 1.1: Workflow and functions of NIPTeR. **a** A BAM file is transformed into an NIPTSample object; **b** a series of NIPTSample objects can then be transformed into an NIPTControlGroup object; **c** optional LOESS or weighted bin GC correction; **d** optional chi-squared-based variation reduction; **e** optional comparison of NIPTSample and NIPTControlGroup and possible selection of a subset that best-matches the control group samples; **f** three different prediction methods: Z-score, normalized chromosome value or regression-based Z-score; **g** optional check of control group statistics

be performed once for each `NIPTControlGroup`. The prediction functions then compare the observed fraction of reads of the chromosome of interest in the `NIPTSample` with the expected fraction. In NCV and RBZ calculations, users have the option of excluding selected chromosomes as predictors. Since chromosomes 13, 18 and 21 are the most likely candidates for a trisomy, these are excluded by default, but users do have the option of including them. The functions `prepare_ncv` and `perform_regression` provide users the option of using a train and test set to prevent over-fitting the models they create.

In addition to providing Z-scores, the functions also produce control group statistics. The function `match_control_group` provides a Match QC score, a calculation that shows how well the sample fits within the control group based on the fraction of reads mapped to the different chromosomes, a measure that can be shown in a report. Alternately, users can select a subset of best-matching control samples as a sample-specific control group using the arguments `mode = "report"` or `"subset"`. When a sample has an anomalously high Match QC score, the control samples being used are not suitable as a control group for the sample being analyzed. A second quality control function, `diagnose_control_group`, calculates Z-scores for all samples and chromosomes in a `NIPTControlGroup` as well as the mean, standard deviation and Shapiro-Wilk test of those Z-scores. This information can be used to curate the control group as explained in detail in Additional file 1.

1.3 Results

1.3.1 Workflow

All these NIPTeR building blocks can be combined into an analysis workflow. For example, the NIPTeR workflow for the Fan & Quake analysis [5], using a weighted bin GC correction and a standard Z-score prediction for trisomy 21, and given a GC-corrected control group is:

```
> NIPTsample <- bin_bam_sample(bam_filepath =
  "/Path/to/bam/sample.bam")
> NIPTsample_gc <- gc_correct(nipt_object = NIPTsample,
  method = "bin")
> Zscore21NIPTsample <- calculate_z_score(nipt_sample =
  NIPTsample_gc, nipt_control_group = NIPTControlGroup_gc,
  chromo_focus = 21)
```

In addition, control group statistics and the match control of the sample to the control group can be performed:

```
> NIPTcontrol_diagnose - diagnose_control_group(nipt_control_group
  = NIPT_control_group_gc)
> MatchQC <- match_control_group(nipt_sample = NIPTsample_gc,
  nipt_control_group = NIPT_control_group_gc, mode = "report")
```

1.3.2 Prediction and control group statistics

The output formats of the `calculate_z_score` and `calculate_ncv_score` functions are similar. An example result of the main output reads:

```
Zscore21NIPTsample$sample_Zscore
[1] 0.4575612

Zscore21NIPTsample$control_group_statistics
mean          SD          Shapiro_P_value
1.380646e-02   7.184378e-05   9.498096e-01
```

Here, the Z-score is 0.45, which falls within the -3 to 3 range and leads to the conclusion that this sample does not have a trisomy 21. The `control_group_statistics` show the mean fraction of sequence

reads mapping to chromosome 21 and the standard deviation (SD) of the fractions between the control samples. The Shapiro_P_value tests for control group normality, and control groups with a value above 0.05 can be considered to be normally distributed.

The output of `perform_regression` is slightly different and gives four predictions based on different models when set to the default setting:

	Prediction_set_1	Prediction_set_2	Prediction_set_3	Prediction_set_4
Z_score_sample	0.695389767405796	0.436463271170429	0.437555582217223	-0.268842730284741
CV	0.00536568258297721	0.00502335300817695	0.00483989627449594	0.00486660271957713
cv_types	Practical_CV	Practical_CV	Practical_CV	Practical_CV
P_value_shapiro	0.430190936876808	0.844844184734285	0.478810106756347	0.606229054979589
Pred_chrom ¹	3F 1F 2R 7F	3R 22F 1R 5R	6R 10F 8R 17F	20F 12F 19R 14F
Mean_test_set	0.998406705791639	0.997692920712523	0.998044728541847	0.997802000172399
CV_train_set	0.00441576466562767	0.004609720864648	0.00479265227193279	0.00492160650642337

Here, in addition to the RBZ, the coefficient of variation (CV) of the test set is given as a measure of control group variability. The type of CV is given as well, in which “Practical CV” is the true CV. If there is a risk of over-fitting the model on the control set, a theoretical CV is used. In addition to the Shapiro P value, `perform_regression` reports the mean of the test set (which should be close to one) and the CV of the training set (based on which the chromosomes used to create the prediction model are selected), where reads mapped to the forward and reverse strands are used as separate entities.

1.3.3 Quality control

Using the `diagnose_control_group` function, control samples that have outliers that could hamper prediction can be detected.

¹In practice `Pred_chrom` is written in full as: `Predictor_chromosomes`. For layout purposes a shorthand is used here.

```
> NIPTcontrol_diagnose$abberant_scores
```

	Chromosome	Sample_name	Z_score
1	17F	sample21	3.13281485801102
2	1R	sample21	3.1290608434065
3	17R	sample21	3.33995848430216
4	22R	sample24	3.08496372975161
...			
19	8F	sample21	-3.85723794269498
20	5R	sample21	-3.16594249087773
21	16R	sample21	-3.5467264109158

This example shows that, for many chromosomes in sample 21 one or both of the strands have a Z-score higher than 3. This means that there is more variability in this sample than expected, pointing to a low quality sample. As explained in more detail in Additional file 1, we recommend that users remove samples that have more than one aberrant score (Z-score outside the -3 to 3 range) from the control group.

When looking at the individual Match QC scores of the GC corrected NIPTSample compared to the GC corrected NIPTControlGroup, the list of sum of squares of differences in chromosomal fractions of the test sample compared to each control sample is shown:

	Sum_of_squares
sample86	1.919715e-07
sample74	2.155461e-07
...	
sample40	1.089867e-06
sample21	2.028651e-06

In general, the lower the sum of squares, the more representative a control sample is for the test sample. The average of all sum of squares for an NIPTSample is the Match QC score. A Match QC score for a specific sample that falls outside 3 SD of the control group Match QC, indicates that the control group is not suitable for

analysis of the sample.

Further examples and results can be found in the NIPTeR package vignette [8] and the case report provided in Additional file 1. A demonstration of the NIPTeR GC-correction methods is given in Additional file 2 and a comparison of NIPTeR results with manual calculations is available for the χ^2 VR in Additional file 3 and for the prediction methods and Match QC score in Additional file 4.

The NIPTeR package requires R 3.1.0 or higher, the stats and sets packages as available on CRAN, and the RSamtools and S4Vectors Bioconductor packages.

1.3.4 Performance

NIPTeR performance was tested on three different machines and operating systems (Additional file 5). Given a pre-processed control group of 100 samples, one sample was processed in 3 to 4 min (on average), including both GC correction and χ^2 VR and using the Z-score and RBZ as prediction algorithms for chromosomes 13, 18 and 21. NCV analysis was performed in an additional 1 to 6 min using a maximum number of 6 to 9 chromosomes as denominator.

1.4 Conclusion

NIPTeR allows for fast NIPT analysis and flexible workflow creation and includes variation correction and prediction algorithms as well as QC control. Algorithms used in NIPTeR are validated as described in Johansson and de Boer et al. (2017) [7]. NIPTeR is available under the GNU GPL open source license and can be freely downloaded from <https://github.com/molgenis/NIPTeR> or CRAN.

1.5 Availability and requirements

Project name: NIPTeR. Project home page: <https://CRAN.R-project.org/package=NIPTeR> Source page: <https://github.com/molgenis/NIPTeR> Operating system(s): Linux, MacOS, Windows. Programming language: R. Other requirements:

R (3.1.0 or higher), RSamtools, sets, stats, S4Vectors. Licence: GNU Lesser General Public License v3.0. Any restrictions to use by non-academics: none

Acknowledgments

We thank Kate Mc Intyre for editorial advice.

Authors' contributions

LJ is the main author. LJ and HdW conceived and designed the NIPTeR package. Together with FvD they developed and implemented the application. LJ, HdW, EdB and GtM designed and validated algorithms and implementation. RS, BS and MS were responsible for project administration and supervision. All authors read and approved the final version of this manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

1.6 Additional files

Additional files can be accessed online:

<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-018-2557-8>

1

2

3

4

5

6

7

8

9

10

11

Bibliography

- [1] Megan Allyse, Mollie Minear, Margaret Rote, Anthony Hung, Subhashini Chandrasekharan, Elisa Berson, and Shilpa Sridhar. Non-invasive prenatal testing: a review of international implementation and challenges. *International Journal of Women's Health*, page 113, Jan 2015.
- [2] Eric Z. Chen, Rossa W. K. Chiu, Hao Sun, Ranjit Akolekar, K. C. Allen Chan, Tak Y. Leung, Peiyong Jiang, Yama W. L. Zheng, Fiona M. F. Lun, Lisa Y. S. Chan, and et al. Noninvasive prenatal diagnosis of fetal trisomy 18 and trisomy 13 by maternal plasma dna sequencing. *PLoS ONE*, 6(7):e21791, Jul 2011.
- [3] Z. Chen. Development of Bioinformatics Algorithms for Trisomy 13 and 18 Detection by Next Generation Sequencing of Maternal Plasma DNA. The Chinese University of Hong Kong, 2011.
- [4] R. W. K. Chiu, K. C. A. Chan, Y. Gao, V. Y. M. Lau, W. Zheng, T. Y. Leung, C. H. F. Foo, B. Xie, N. B. Y. Tsui, F. M. F. Lun, and et al. Non-invasive prenatal diagnosis of fetal chromosomal aneuploidy by massively parallel genomic sequencing of dna in maternal plasma. *Proceedings of the National Academy of Sciences*, 105(51):20458–20463, Dec 2008.
- [5] H. Christina Fan and Stephen R. Quake. Sensitivity of noninvasive prenatal detection of fetal aneuploidy from maternal plasma using shotgun sequencing is limited only by counting statistics. *PLoS ONE*, 5(5):e10439, May 2010.
- [6] Megan P. Hall, Matthew Hill, Bernhard Zimmermann, Styrmir Sigurjónsson, Margaret Westemeyer, Jennifer Saucier, Zachary Demko, and Matthew Rabinowitz. Non-invasive prenatal detection of trisomy 13 using a single nucleotide polymorphism- and informatics-based approach. *PLoS ONE*, 9(5):e96677, May 2014.

BIBLIOGRAPHY

- [7] L. F. Johansson, E. N. de Boer, H. A. de Weerd, F. van Dijk, M. G. Elferink, G. H. Schuring-Blom, R. F. Suijkerbuijk, R. J. Sinke, G. J. te Meerman, R. H. Sijmons, and et al. Novel algorithms for improved sensitivity in non-invasive prenatal testing. *Scientific Reports*, 7(1), May 2017.
- [8] Johansson L.F. and de Weerd H.A. Nipter vignette., 2016.
- [9] H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, and R. Durbin. The sequence alignment/map format and samtools. *Bioinformatics*, 25(16):2078–2079, Jun 2009.
- [10] Baohong Liu, Xiaoyan Tang, Feng Qiu, Chunmei Tao, Junhui Gao, Mengmeng Ma, Tingyan Zhong, JianPing Cai, Yixue Li, and Guohui Ding. Dasaf: An r package for deep sequencing-based detection of fetal autosomal abnormalities from maternal cell-free dna. *BioMed Research International*, 2016:1–7, 2016.
- [11] Kitty K. Lo, Christopher Boustred, Lyn S. Chitty, and Vincent Plagnol. Rapidr: an analysis package for non-invasive prenatal testing of aneuploidy. *Bioinformatics*, 30(20):2965–2967, Jul 2014.
- [12] Minh-Duy Phan, Thong V. Nguyen, Huong N. T. Trinh, Binh T. Vo, Truc M. Nguyen, Nguyen H. Nguyen, Tho T. Q. Nguyen, Thuy T. T. Do, Tuyet T. D. Hoang, Kiet D. Truong, and et al. Establishing and validating noninvasive prenatal testing procedure for fetal aneuploidies in vietnam. *The Journal of Maternal-Fetal & Neonatal Medicine*, page 1–7, Jul 2018.
- [13] Martin Sauk, Olga Žilina, Ants Kurg, Eva-Liina Ustav, Maire Peters, Priit Paluoja, Anne Mari Roost, Hindrek Teder, Priit Palta, Nathalie Brison, and et al. Niptmer: rapid k-mer-based software package for detection of fetal aneuploidies. *Scientific Reports*, 8(1), Apr 2018.
- [14] A. J. Sehnert, B. Rhee, D. Comstock, E. de Feo, G. Heilek, J. Burke, and R. P. Rava. Optimal detection of fetal chromosomal abnormalities by massively parallel dna sequencing of cell-free fetal dna from maternal blood. *Clinical Chemistry*, 57(7):1042–1049, Apr 2011.
- [15] Roy Straver, Erik A. Sistermans, Henne Holstege, Allerdien Visser, Cees B. M. Oudejans, and Marcel J. T. Reinders. Wisecondor: detection of fetal aberrations from shallow sequencing maternal plasma based on a within-sample comparison scheme. *Nucleic Acids Research*, 42(5):e31–e31, Oct 2013.
- [16] Jianfeng Yang, Xiaofan Ding, and Weidong Zhu. Improving the calling of non-invasive prenatal testing on 13-/18-/21-trisomy by support vector machine discrimination. Nov 2017.

List of Tables

LIST OF TABLES

List of Figures

1.1	Workflow and functions of NIPTeR	7
-----	--	---