

1 **JAST (Journal of Animal Science and Technology) TITLE PAGE**

2 **Upload this completed form to website with submission**

3

ARTICLE INFORMATION	Fill in information in each box below
Article Type	Review article
Article Title (within 20 words without abbreviations)	Enhancing Animal Breeding through Quality Control in Genomic Data - A Review
Running Title (within 10 words)	Quality Control in Genomic Data
Author	Jungjae Lee ¹ , Jong Hyun Jung ² and Sang-Hyon Oh ³
Affiliation	¹ Jenomics Jenetics Company, Pyeongtaek 17869, Korea ² Jung P&C Institute, Yongin 16950, South Korea ³ Division of Animal Science, Gyeongsang National University, Jinju 52725, South Korea
ORCID (for more information, please visit https://orcid.org)	Jungjae Lee https://orcid.org/0000-0002-6145-8862 Jong Hyun Jung https://orcid.org/0000-0003-3667-7710 Sang-Hyon Oh https://orcid.org/0000-0002-9696-9638
Competing interests	No potential conflict of interest relevant to this article was reported.
Funding sources State funding sources (grants, funding sources, equipment, and supplies). Include name and number of grant if available.	This work was carried out with the support of "Cooperative Research Program for Agriculture Science and Technology Development (Project No. RS-2023-00232087)" Rural Development Administration, Republic of Korea.
Acknowledgements	N/A
Availability of data and material	Upon reasonable request, the datasets of this study can be available from the corresponding author.
Authors' contributions Please specify the authors' role using this form.	Conceptualization: LJJ, JJH, OSH. Data curation: LJJ, JJH, OSH. Formal analysis: LJJ, JJH, OSH. Methodology: LJJ, JJH, OSH. Software: N/A Validation: LJJ, JJH, OSH. Investigation: LJJ, JJH, OSH. Writing - original draft: LJJ, JJH, OSH. Writing - review & editing: LJJ, JJH, OSH.
Ethics approval and consent to participate	This article does not require IRB/IACUC approval because there are no human and animal participants.

4

5 **CORRESPONDING AUTHOR CONTACT INFORMATION**

For the corresponding author (responsible for correspondence, proofreading, and reprints)	Fill in information in each box below
First name, middle initial, last name	Prof. Sang-Hyon Oh
Email address – this is where your proofs will be sent	Prof. Sang-Hyon Oh at shoh@gnu.ac.kr
Secondary Email address	
Address	Division of Animal Science, Gyeongsang National University, Jinju 52725, South Korea
Cell phone number	
Office phone number	
Fax number	

6

7

8 **Enhancing Animal Breeding through Quality Control in Genomic Data - A Review**

9 Jungjae Lee^{1†}, Jong Hyun Jung^{2†} and Sang-Hyon Oh^{3*}

10
11
12
13 † Both authors contributed equally to this manuscript.

14
15 *Corresponding Author: Sang-Hyon Oh

16 Tel: +82-55-772-3285, Fax: +82-55-772-3689, E-mail: shoh@gnu.ac.kr

17
18
19 ¹Jenomics Jenetics Company, Pyeongtaek 17869, Korea

20 ²Jung P&C Institute, Yongin 16950, South Korea

21 ³Division of Animal Science, Gyeongsang National University, Jinju 52725, South Korea

22
ACCEPTED

23 **Abstract**

24 High-throughput genotyping and sequencing has revolutionized animal breeding by
25 providing access to vast amounts of genomic data to facilitate precise selection for desirable
26 traits. This shift from traditional methods to genomic selection provides dense marker
27 information for predicting genetic variants. However, the success of genomic selection heavily
28 depends on the accuracy and quality of the genomic data. Inaccurate or low-quality data can
29 lead to flawed predictions, compromising breeding programs and reducing genetic gains.
30 Therefore, stringent quality control (QC) measures are essential at every stage of data
31 processing. Quality control in genomic data involves managing single nucleotide
32 polymorphism (SNP) quality, assessing call rates, and filtering based on minor allele frequency
33 (MAF) and Hardy-Weinberg equilibrium (HWE). High-quality SNP data is crucial because
34 genotyping errors can bias the estimates of breeding values. Cost-effective low-density
35 genotyping platforms often require imputation to deduce missing genotypes. QC is vital for
36 genomic selection, genome-wide association studies (GWAS), and population genetics
37 analyses because it ensures data accuracy and reliability. This paper reviews QC strategies for
38 genomic data and emphasizes their applications in animal breeding programs. By examining
39 various QC tools and methods, this review highlights the importance of data integrity in
40 achieving successful outcomes in genomic selection, GWAS, and population analyses.
41 Furthermore, this review covers the critical role of robust QC measures in enhancing the
42 reliability of genomic predictions and advancing animal breeding practices.

43

44 **Key words:** Animal Breeding, Genomic Selection, Quality Control, Single Nucleotide
45 Polymorphism, Genome-Wide Association Studies

46

47 **Introduction**

48 The rapid evolution of genomic technologies has transformed the landscape of animal
49 breeding. High-throughput genotyping and sequencing provides breeders with access to vast
50 amounts of genomic data and enables the precise selection of desirable traits [1]. These
51 advancements have shifted traditional breeding methods to genomic selection, which leverages
52 dense marker information to predict the genetic variants of individuals [2]. However, the
53 success of genomic selection depends heavily on the accuracy and quality of the genomic data.
54 Inaccurate or low-quality data can lead to inaccurate predictions that can compromise breeding
55 programs and reduce their genetic gains [3]. Therefore, to ensure reliable predictions and
56 maximize the potential of genomic selection, it is essential to implement stringent quality
57 control (QC) measures at every stage of data processing.

58 Genomic data quality control has several key components including the management of
59 single nucleotide polymorphism (SNP) quality, the assessment of call rates, and filtering based
60 on minor allele frequency (MAF) and Hardy-Weinberg equilibrium (HWE) [4]. High-quality
61 SNP data is indispensable because errors in genotyping can lead to biased estimates of breeding
62 values, which decreases the effectiveness of selection strategies [5]. Moreover, cost-effective
63 low-density genotyping platforms often suffer from incomplete marker data so it is necessary
64 to use imputation to deduce the missing genotypes [6].

65 Quality control processes are crucial for genomic selection, genome-wide association
66 studies (GWAS), and population genetics analyses. These processes help ensure that the
67 genomic data is accurate, reliable, and free from biases introduced by genotyping errors,
68 population stratification, or other confounding factors [7, 8]. This paper reviews quality control
69 (QC) strategies for genomic data and their applications in animal breeding programs. By
70 examining various QC tools and methods, this paper aims to show the critical role that data

71 integrity plays in achieving successful outcomes in genomic selection, GWAS, and population
72 analyses [4, 5].

73

74 **Genotyping methods**

75 *Whole-genome Sequencing (WGS)*

76 Whole-genome sequencing is a comprehensive method for analyzing the entire genome.
77 Due to the decreased cost of sequencing and the ability to produce large amounts of genomic
78 data, whole-genome sequencing has become a powerful tool for genomic research. SNP calling
79 from WGS genomic data involves a series of critical steps to ensure accurate identification of
80 genetic variants. The process starts with raw data preprocessing, where tools like FastQC
81 evaluate the read quality [9]. This step is followed by trimming to remove adapters and low-
82 quality bases by using either Trimmomatic or Cutadapt [10, 11].

83 The cleaned reads are then aligned to a reference genome with BWA-MEM or Bowtie2 to
84 generate SAM/BAM files [12, 13]. These files are subsequently sorted, indexed, and processed
85 to mark PCR duplicates with Samtools, while the base quality scores are recalibrated using
86 GATK [14, 15]. Variant calling is performed using tools such as GATK's HaplotypeCaller,
87 FreeBayes, or Bcftools, which identify SNPs based on differences between the sequenced reads
88 and the reference genome [15, 16, 17].

89 In post-calling, variants undergo filtering to remove false positives via GATK's hard
90 filtering or Variant Quality Score Recalibration (VQSR). The filtered SNPs are then annotated
91 with functional information using tools like ANNOVAR or SnpEff [18, 19]. Quality checks
92 include the use of VCFtools for statistical analysis and IGV for visualization, and ensure the
93 reliability of the called SNPs [16, 20]. Joint genotyping across multiple samples and using
94 population-specific reference panels are recommended to enhance the accuracy of SNP calling

95 in WGS.

96

97 *SNP arrays*

98 SNP arrays have significantly advanced genomic research in animal science by enabling
99 the large-scale genotyping of SNPs. The development of SNP arrays began in the early 2000s
100 to meet the demand for efficient and cost-effective methods to genotype large numbers of SNPs
101 across the genome [21, 22]. Early arrays marked a significant advancement by allowing
102 simultaneous genotyping of thousands of SNPs, facilitating genome-wide association studies
103 (GWAS) and the study of genetic variation in populations [21].

104 Over time, these arrays have evolved to include higher-density SNPs to improve coverage
105 and accuracy, as seen in the Illumina BovineSNP50 array which has become a standard tool in
106 cattle genomics [23, 24]. Today, SNP arrays are essential for selecting desirable traits,
107 estimating genetic merit, and managing inbreeding in animal breeding [1, 2]. Quality control
108 of SNP array data is crucial for ensuring accurate and reliable results, and involves assessing
109 call rates, filtering based on minor allele frequency (MAF), and checking for Hardy-Weinberg
110 equilibrium [4]. Tools such as PLINK and GenomeStudio are commonly used in these QC
111 processes [5, 25].

112

113 **QC in animal genomics**

114 *Minor Allele Frequency (MAF)*

115 Minor allele frequency is a key metric in genetic studies. It represents the frequency at
116 which the less common allele occurs in a given population. MAF is important for identifying
117 rare variants which may not significantly contribute to overall genetic variation but can be
118 crucial in specific contexts. MAF is calculated by determining the frequency of both alleles at

119 a locus and taking the minimum of these two values. For example, if allele A has a frequency
120 of 0.8 and allele a has a frequency of 0.2, the MAF would be 0.2. SNPs with very low MAFs,
121 typically below 0.01 or 0.05, are often excluded from analyses because they may represent
122 sequencing errors or lack statistical power in association studies [5].

123 Tools like PLINK and VCFtools [5, 16] are widely used to calculate MAF, with PLINK's
124 --freq command being particularly popular [4]. In animal breeding, many researchers set
125 threshold values for MAF to balance the need for sufficient variation while minimizing noise
126 from rare variants. Typically, MAF thresholds in animal breeding studies range from 0.01 to
127 0.05 depending on the study's objectives and the population structure being analyzed. For
128 instance, a study on dairy cattle by Pryce et al. [26] and Kim et al. [27] used a MAF threshold
129 of 0.01 to ensure that the SNPs included were sufficiently informative for genomic predictions
130 while also minimizing the influence of rare variants that might lead to spurious associations.

131

132 *Call rate*

133 Call rate is another critical quality control metric that measures the proportion of
134 successfully genotyped samples for a specific SNP. A high call rate indicates that a SNP has
135 been consistently detected across the sample population, while a low call rate may suggest
136 issues with the genotyping process, such as poor quality or technical errors [7].

137 The call rate is calculated by dividing the number of successful genotype calls for a SNP
138 by the total number of samples, then multiplying by 100 to express it as a percentage.

139

$$140 \text{ Call Rate} = \frac{\text{Number of successfully genotyped markers (or samples)}}{\text{Total number of markers (or samples)}} \times 100$$

141

142 For instance, if 95 out of 100 samples have a successful genotype call for a SNP, the call
143 rate would be 95% [4]. Normally, markers with a call rate less than 95% are removed, though
144 other studies have set more stringent or lenient thresholds depending on the study design and
145 objectives. For example, some studies have removed markers with a call rate below 99% to
146 ensure extremely high data quality[28], while others have used a more relaxed threshold of 90%
147 when working with larger datasets[29].

148 Tools like PLINK, SNP & Variation Suite (SVS), and GenomeStudio are widely used for
149 calculating and filtering SNPs based on call rates because they offer robust functionalities for
150 quality control in genomic studies. PLINK is particularly popular due to its comprehensive
151 command-line interface, where the --missing command calculates call rates at both the marker
152 and sample levels, allowing researchers to easily filter out SNPs and samples that fall below
153 the desired threshold [5]. SNP & Variation Suite (SVS) offers a user-friendly graphical interface
154 and integrates various statistical tools, making it ideal for complex datasets and large-scale
155 studies [30]. GenomeStudio by Illumina is another powerful tool specifically designed for
156 managing and analyzing genotyping data with features for calculating call rates, identifying
157 low-quality markers, and visualizing data for further inspection [25]. These tools are essential
158 for ensuring that only high-quality data is used in subsequent analyses to improve the reliability
159 of genomic outcomes.

160

161 *Hardy-Weinberg Equilibrium*

162 Hardy-Weinberg equilibrium (HWE) is a fundamental principle in population genetics. It
163 states that allele and genotype frequencies in a population will remain constant from generation
164 to generation in the absence of evolutionary influences [31]. Testing for HWE is an important
165 quality control step because deviations from this equilibrium can indicate issues such as

166 genotyping errors, population stratification, or selection pressures [32]. To test for HWE, the
167 observed genotype frequencies are compared to the expected frequencies under equilibrium
168 conditions. For a biallelic SNP with alleles A and a, the expected genotype frequencies are p^2
169 for AA, $2pq$ for Aa, and q^2 for aa, where p and q represent the allele frequencies [33]. A chi-
170 square test is commonly used to assess whether the differences between the observed and
171 expected frequencies are statistically significant. Tools like PLINK and VCFtools are used to
172 perform HWE tests [34]. SNPs that show significant deviation from HWE, typically with a p-
173 value less than 0.001, are often excluded from analyses to prevent biases that could arise from
174 genotyping errors or other confounding factors [4].

ACCEPTED

175 Table 1. Tool list for quality control processes

Tools	Function	Reference
GEMMA	Application of linear mixed models and related models to GWAS	[4]
PLINK	Run association analyses and perform QC and regression steps	[5]
FastQC	Quality control checks on raw sequence data	[9]
Trimmomatic	Trim and crop FASTQ data	[10]
Cutadapt	finds and removes adapter sequences, primers, poly-A tails	[11]
BWA-MEM	produce multiple primary alignments for different part of a query sequence	[12]
Bowtie2	aligning sequencing reads to long reference sequences	[13]
Samtools	Manipulate alignments in the SAM, BAM, and CRAM formats	[14]
GATK	Variant calling using sequencing data	[15]
VCFTools	Summarize, filter out, convert data into other file formats	[16]
FreeBayes	Bayesian genetic variant detector designed to fine SNPs	[17]
SnEff	Annotation on genetic variants and predicts their effects on genes	[18]
ANNOVAR	Generate gene-based annotation	[19]
IGV	Visualization tool to simultaneously integrate and analyze multiple types of genomic data	[20]
GenomeStudio	Normalize, cluster, and call genotypes	[25]
SVS	Perform analyses and visualizations on genomic and phenotypic data	[33]
BEAGLE	Genotype calling, phasing, and genotype imputation	[39]
Fimpute	Haplotype estimation or phasing and genotype imputation	[40]
Impute2	Genotype imputation and haplotype phasing	[47]
Minimac	performs imputation with pre-phased haplotypes	[48]

176

177 **Application**

178 *Population analysis*

179 Population analysis is invaluable for genomic studies in animal science because it enables
180 researchers to assess the genetic structure, diversity, and evolutionary dynamics within and
181 between populations. Accurately characterizing population structures is crucial for identifying
182 subpopulations, measuring inbreeding levels, and understanding the genetic background of
183 breeding populations, all of which are essential for maintaining genetic diversity and improving
184 selection outcomes [35]. Tools such as PLINK, ADMIXTURE, and STRUCTURE are
185 commonly employed to detect key characteristics for understanding the genetic landscape of
186 animal populations, such as population stratification, admixture, and genetic differentiation [5,
187 36]. For example, ADMIXTURE provides estimates of individual ancestry proportions. These
188 estimates allow researchers to detect mixed genetic backgrounds that could influence trait
189 analysis [36]. Quality control measures, such as filtering based on MAF, HWE, and genotyping
190 call rates ensure the data used for population analysis is reliable [4,37]. MAF filtering helps
191 exclude rare alleles that may introduce noise or result from genotyping errors [5]. Similarly,
192 HWE filtering removes SNPs that deviate from expected frequencies due to selection or
193 population substructures in order to prevent potential biases in the analysis [37]. Proper quality
194 control improves the accuracy of population structure analyses and mitigates the risk of
195 confounding in subsequent analyses such as GWAS and genomic selection [4]. By accurately
196 characterizing population structures, researchers can identify unique genetic markers and
197 enhance their understanding of trait inheritance, and then design breeding strategies that
198 optimize genetic gain and preserve diversity to support sustainable livestock production [35,
199 36].

200

201 *GWAS*

202 Genome-wide association studies (GWAS) are powerful tools for identifying genetic
203 variants associated with complex traits in animal breeding such as growth traits, disease
204 resistance, reproductive traits, and carcass traits [2, 4]. The reliability of GWAS findings hinges
205 on rigorous quality control (QC) procedures that ensure high-quality data throughout the
206 process. This begins with careful study design and population selection, where potential
207 confounders like population stratification are addressed through methods such as Principal
208 Component Analysis (PCA) and linear mixed models to correct for genetic structure within the
209 population [38]. Phenotype data must be accurately collected and screened for outliers to
210 minimize noise. Genotype data undergoes thorough QC, including filtering SNPs based on call
211 rates, MAF, and deviations from HWE [4, 5]. For instance, SNPs with low call rates are
212 excluded to avoid unreliable data that could lead to false-positive associations, while MAF
213 filtering focuses the analysis on common variants that are more likely to have sufficient
214 statistical power to detect true associations. HWE filtering is employed to remove SNPs that
215 significantly deviate from expected allele frequencies because such deviations may indicate
216 genotyping errors or underlying selection pressures [5]. To reduce redundancy and
217 computational burden, linkage disequilibrium (LD) pruning is performed and missing
218 genotypes are often imputed via reference panels using Fimpute or BEAGLE [39, 40]. Tools
219 like PLINK and GEMMA are widely used to implement QC measures and conduct association
220 tests because they offer a robust framework for analyzing large genomic datasets [4]. Statistical
221 analysis in GWAS is carried out using models appropriate for the trait under study, and
222 corrections for multiple testing to mitigate the risk of false positives and meta-analysis may be
223 employed when integrating results from multiple studies [41]. To ensure the robustness and
224 high accuracy of the GWAS models, a 5-fold cross-validation is often used. In this method, the
225 datasets are divided into five subsets. The model is iteratively trained on four subsets and tested

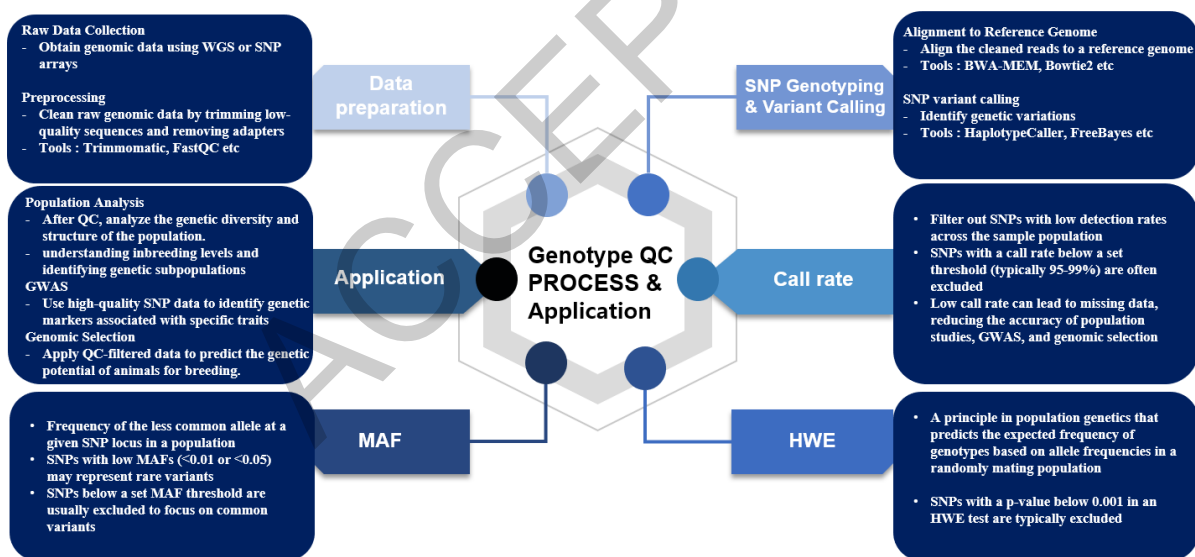
226 on the remaining one to help validate the model's accuracy and mitigate overfitting [42]. The
227 results from GWAS offer valuable genetic variants for traits which can be targeted in marker-
228 assisted selection and genomic selection programs. Genomic selection aims to ultimately
229 improve the genetic merit of livestock populations [2].

230

231 *Genomic Selection*

232 Genomic selection (GS) allows for the selection of animals based on SNP markers [43].
233 With the introduction of genomic selection, animal breeding has dramatically advanced by
234 overcoming the limitations of traditional selection methods like best linear unbiased prediction
235 (BLUP) and marker-assisted selection [43, 44]. GS relies on dense SNP data to estimate
236 genomic breeding values, which are used to predict an individual's genetic potential for
237 economically important traits [2]. The accuracy of genomic selection models is dependent upon
238 the quality of the genomic data and the reliability of GS models can be enhance significantly
239 by the inclusion of imputation methods to handle missing or low-density SNP data [45].
240 Imputation is beneficial in low-density platforms because it allows for the cost-effective use of
241 genotyping while still leveraging the power of high-density SNP information. Imputation
242 increases the accuracy of genomic predictions by inferring missing genotypes in order to
243 improve the reliability of estimated breeding values even with fewer markers [46]. Several
244 imputation tools, including FImpute [40], Beagle [39], Impute2 [47], and Minimac [48] are
245 widely used in animal breeding to enhance the accuracy of genomic selection models.
246 Therefore, strict quality control is essential [49]. Quality control methods, such as filtering
247 SNPs based on call rates, MAF, and HWE, is critical to ensuring that the data is vigorous and
248 reliable. High call rates are important because missing data can introduce bias and reduce the
249 reliability of genomic estimated breeding values. Similarly, excluding SNPs with low MAF

250 helps to avoid the noise associated with rare variants that may have little impact on prediction
 251 accuracy. Ensuring that SNPs conform to HWE expectations also prevents the inclusion of
 252 markers affected by selection, mutation, or other factors that could bias the genomic selection
 253 models [4, 5]. Advanced computational tools, such as GBLUP (Genomic Best Linear Unbiased
 254 Prediction) and ssBLUP (single-step BLUP), and Bayesian methods (BayesA, BayesB, BayesC)
 255 integrate SNP effects across the genome to enhance the precision of breeding value predictions
 256 [50, 51]. By using high-quality genomic data, genomic selection enables breeders to make more
 257 accurate decisions that lead to faster genetic gains and the improvement of traits such as milk
 258 yield, growth rate, and carcass weight in livestock. This approach not only enhances the
 259 efficiency of breeding programs but also contributes to the long-term sustainability and
 260 productivity of animal populations [35].



261

262 Figure 1. Overall flowchart from data preparation to Application in animal breeding

263

264 **Conclusion**

265 High-throughput genotyping and sequencing has significantly advanced the field of
266 animal breeding by enabling precise selection for desirable traits. However, the success of
267 genomic selection hinges on the accuracy and quality of the genomic data used. Rigorous
268 quality control (QC) measures are essential to ensure data integrity. These measures include
269 SNP quality management, call rate assessment, and filtering based on minor allele frequency
270 (MAF) and Hardy-Weinberg equilibrium (HWE). These QC processes are crucial for genomic
271 selection, genome-wide association studies (GWAS), and population genetics analyses.
272 Implementing stringent QC strategies enhances the reliability of genomic predictions, which
273 improves breeding programs and genetic gains. By maintaining high standards of data quality,
274 researchers and breeders can make informed decisions that lead to sustainable and productive
275 advancements in animal breeding.

276

277

ACCEPTED

278 **Reference**

- 279 1. Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME. Invited review: Genomic selection
280 in dairy cattle: Progress and challenges. *Journal of dairy science*. 2009;92(2):433-43.
281 <https://doi.org/10.3168/jds.2008-1646>
- 282 2. Meuwissen TH, Hayes BJ, Goddard M. Prediction of total genetic value using genome-
283 wide dense marker maps. *genetics*. 2001;157(4):1819-29.
284 <https://doi.org/10.1093/genetics/157.4.1819>
- 285 3. Wiggans G, VanRaden P, Cooper T. The genomic evaluation system in the United States:
286 Past, present, future. *Journal of dairy science*. 2011;94(6):3202-11.
287 <https://doi.org/10.3168/jds.2010-3866>
- 288 4. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan KT. Data
289 quality control in genetic case-control association studies. *Nature protocols*.
290 2010;5(9):1564-73. <https://doi.org/10.1038/nprot.2010.116>
- 291 5. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a
292 tool set for whole-genome association and population-based linkage analyses. *The*
293 *American journal of human genetics*. 2007;81(3):559-75. <https://doi.org/10.1086/519795>
- 294 6. García-Ruiz A, Cole JB, VanRaden PM, Wiggans GR, Ruiz-López FJ, Van Tassell CP.
295 Changes in genetic selection differentials and generation intervals in US Holstein dairy
296 cattle as a result of genomic selection. *Proceedings of the National Academy of Sciences*.
297 2016;113(28):E3995-E4004. <https://doi.org/10.1073/pnas.1519061113>
- 298 7. Laurie CC, Doheny KF, Mirel DB, Pugh EW, Bierut LJ, Bhangale T, et al. Quality control
299 and quality assurance in genotypic data for genome-wide association studies. *Genetic*
300 *epidemiology*. 2010;34(6):591-602. <https://doi.org/10.1002/gepi.20516>
- 301 8. Marees AT, De Kluiver H, Stringer S, Vorspan F, Curis E, Marie-Claire C, et al. A tutorial
302 on conducting genome-wide association studies: Quality control and statistical analysis.
303 *International journal of methods in psychiatric research*. 2018;27(2):e1608.
304 <https://doi.org/10.1002/mp.1608>
- 305 9. Andrews S. *FastQC: a quality control tool for high throughput sequence data*. Cambridge,
306 United Kingdom; 2010.

- 307 10. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
308 data. *Bioinformatics*. 2014;30(15):2114-20.
309 <https://doi.org/10.1093/bioinformatics/btu170>
- 310 11. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
311 *EMBnet journal*. 2011;17(1):10-2. <https://doi.org/10.14806/ej.17.1.200>
- 312 12. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
313 *bioinformatics*. 2009;25(14):1754-60. <https://doi.org/10.1093/bioinformatics/btp324>
- 314 13. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature methods*.
315 2012;9(4):357-9. <https://doi.org/10.1038/nmeth.1923>.
- 316 14. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence
317 alignment/map format and SAMtools. *bioinformatics*. 2009;25(16):2078-9.
318 <https://doi.org/10.1093/bioinformatics/btp352>
- 319 15. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
320 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
321 sequencing data. *Genome research*. 2010;20(9):1297-303.
322 <https://doi.org/10.1101/gr.107524.110>
- 323 16. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant
324 call format and VCFtools. *Bioinformatics*. 2011;27(15):2156-8.
325 <https://doi.org/10.1093/bioinformatics/btr330>
- 326 17. Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing.
327 *arXiv preprint arXiv:12073907*. 2012. <https://doi.org/10.48550/arXiv.1207.3907>
- 328 18. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for
329 annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs
330 in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *fly*. 2012;6(2):80-
331 92. <https://doi.org/10.4161/fly.19695>
- 332 19. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants
333 from high-throughput sequencing data. *Nucleic acids research*. 2010;38(16):e164-e.
334 <https://doi.org/10.1093/nar/gkq603>
- 335 20. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.

- 336 Integrative genomics viewer. Nature biotechnology. 2011;29(1):24-6.
337 <https://doi.org/10.1038/nbt.1754>
- 338 21. Kennedy GC, Matsuzaki H, Dong S, Liu W-m, Huang J, Liu G, et al. Large-scale
339 genotyping of complex DNA. Nature biotechnology. 2003;21(10):1233-7.
340 <https://doi.org/10.1038/nbt869>
- 341 22. Matsuzaki H, Dong S, Loi H, Di X, Liu G, Hubbell E, et al. Genotyping over 100,000
342 SNPs on a pair of oligonucleotide arrays. Nature Methods. 2004;1(2):109-11.
343 <https://doi.org/10.1038/nmeth718>
- 344 23. Matukumalli LK, Lawley CT, Schnabel RD, Taylor JF, Allan MF, Heaton MP, et al.
345 Development and characterization of a high density SNP genotyping assay for cattle. PloS
346 one. 2009;4(4):e5350. <https://doi.org/10.1371/journal.pone.0005350>
- 347 24. Illumina, Inc. BovineSNP50 Genotyping BeadChip. Technical Note; 2009.
- 348 25. Illumina, Inc. GenomeStudio Software. Technical Note; 2010.
- 349 26. Pryce JE, Haile-Mariam M, Goddard ME, Hayes BJ. Identification of genomic regions
350 associated with inbreeding depression in Holstein and Jersey dairy cattle. Genetics
351 Selection Evolution. 2014;46:1-14. <https://doi.org/10.1186/s12711-014-0071-7>
- 352 27. Kim S, Lim B, Cho J, Lee S, Dang C-G, Jeon J-H, et al. Genome-wide identification of
353 candidate genes for milk production traits in Korean Holstein cattle. Animals.
354 2021;11(5):1392. . <https://doi.org/10.3390/ani11051392>
- 355 28. Verma SS, de Andrade M, Tromp G, Kuivaniemi H, Pugh E, Namjou-Khales B, et al.
356 Imputation and quality control steps for combining multiple genome-wide datasets.
357 Frontiers in genetics. 2014;5:370. . <https://doi.org/10.3389/fgene.2014.00370>
- 358 29. Lee J, Kim Y, Cho E, Cho K, Sa S, Kim Y, et al. Genomic analysis using Bayesian
359 methods under different genotyping platforms in Korean Duroc pigs. Animals.
360 2020;10(5):752. <https://doi.org/10.3390/ani10050752>
- 361 30. *SNP & Variation Suite*™ (Version 8.x) [Software]. Bozeman, MT: Golden Helix, Inc.
362 Available from <http://www.goldenhelix.com>.
- 363 31. Edwards A. GH Hardy (1908) and hardy–weinberg equilibrium. Genetics.

- 364 2008;179(3):1143-50. <https://doi.org/10.1534/genetics.104.92940>
- 365 32. Nielsen R, Slatkin M. An introduction to population genetics: theory and applications:
366 Sinauer Associates Sunderland, MA; 2013.
- 367 33. Mayo O. A century of Hardy–Weinberg equilibrium. *Twin Research and Human*
368 *Genetics*. 2008;11(3):249-56. <https://doi.org/10.1375/twin.11.3.249>
- 369 34. Graffelman J, Weir B. Testing for Hardy–Weinberg equilibrium at biallelic genetic
370 markers on the X chromosome. *Heredity*. 2016;116(6):558-68.
371 <https://doi.org/10.1038/hdy.2016.20>
- 372 35. Hill WG. Understanding and using quantitative genetic variation. *Philosophical*
373 *Transactions of the Royal Society B: Biological Sciences*. 2010;365(1537):73-85.
374 <https://doi.org/10.1098/rstb.2009.0203>
- 375 36. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in
376 unrelated individuals. *Genome research*. 2009;19(9):1655-64.
377 <https://doi.org/10.1101/gr.094052.109>
- 378 37. Minelli C, Thompson JR, Abrams KR, Thakkinstian A, Attia J. How should we use
379 information about HWE in the meta-analyses of genetic association studies?. *International*
380 *journal of epidemiology*. 2008;37(1):136-46. <https://doi.org/10.1093/ije/dym234>
- 381 38. McVean G. A genealogical interpretation of principal components analysis. *PLoS*
382 *genetics*. 2009;5(10):e1000686. <https://doi.org/10.1371/journal.pgen.1000686>
- 383 39. Browning BL, Browning SR. A unified approach to genotype imputation and haplotype-
384 phase inference for large data sets of trios and unrelated individuals. *The American*
385 *Journal of Human Genetics*. 2009;84(2):210-23.
386 <https://doi.org/10.1016/j.ajhg.2009.01.005>
- 387 40. Sargolzaei M, Chesnais JP, Schenkel FS. A new approach for efficient genotype
388 imputation using information from relatives. *BMC genomics*. 2014;15:1-12.
389 <https://doi.org/10.1186/1471-2164-15-478>
- 390 41. Pearson TA, Manolio TA. How to interpret a genome-wide association study. *Jama*.
391 2008;299(11):1335-44. <https://doi.org/10.1001/jama.299.11.1335>

- 392 42. Kohavi R, editor A study of cross-validation and bootstrap for accuracy estimation and
393 model selection. Ijcai; 1995: Montreal, Canada.
- 394 43. Meuwissen T, Hayes B, Goddard M. Genomic selection: a paradigm shift in animal
395 breeding. *Animal frontiers*. 2016;6(1):6-14. <https://doi.org/10.2527/af.2016-0002>
- 396 44. Zhang Z, Zhang Q, Ding X. Advances in genomic selection in domestic animals. *Chinese*
397 *science bulletin*. 2011;56:2655-63. <https://doi.org/10.1007/s11434-011-4632->
- 398 45. Berry D, Kearney J. Imputation of genotypes from low-to high-density genotyping
399 platforms and implications for genomic selection. *Animal*. 2011;5(8):1162-9.
400 <https://doi.org/10.1017/S1751731111000309>
- 401 46. García-Ruiz A, Cole JB, VanRaden PM, Wiggans GR, Ruiz-López FJ, Van Tassell CP.
402 Changes in genetic selection differentials and generation intervals in US Holstein dairy
403 cattle as a result of genomic selection. *Proceedings of the National Academy of Sciences*.
404 2016;113(28):E3995-E4004. <https://doi.org/10.1073/pnas.1519061113>
- 405 47. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method
406 for the next generation of genome-wide association studies. *PLoS genetics*.
407 2009;5(6):e1000529. <https://doi.org/10.1371/journal.pgen.1000529>
- 408 48. Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR. Fast and accurate
409 genotype imputation in genome-wide association studies through pre-phasing. *Nature*
410 *genetics*. 2012;44(8):955-9. <https://doi.org/10.1038/ng.2354>
- 411 49. VanRaden P. Symposium review: How to implement genomic selection. *Journal of Dairy*
412 *Science*. 2020;103(6):5291-301. <https://doi.org/10.3168/jds.2019-17684>
- 413 50. Habier D, Fernando RL, Kizilkaya K, Garrick DJ. Extension of the Bayesian alphabet for
414 genomic selection. *BMC bioinformatics*. 2011;12:1-12. [https://doi.org/10.1186/1471-](https://doi.org/10.1186/1471-2105-12-186)
415 [2105-12-186](https://doi.org/10.1186/1471-2105-12-186)
- 416 51. Misztal I, Legarra A, Aguilar I. Computing procedures for genetic evaluation including
417 phenotypic, full pedigree, and genomic information. *Journal of dairy science*.
418 2009;92(9):4648-55. <https://doi.org/10.3168/jds.2009-2064>