

**Imperial College
London**

Using CNN to estimate
recombination rates in *Anopheles
gambiae* genomes

Liyang Huang
August 2020

A thesis submitted for the partial fulfillment of the requirements for the degree
of Master of Research at Imperial College London

Formatted in the journal style of *Molecular Ecology*.
Submitted for the MRes in Computational Methods in Ecology and Evolution

1 Introduction

Malaria is an acute febrile disease. In non-immune individuals, symptoms usually appear 10 to 15 days after being bitten by a mosquito. *Plasmodium falciparum* malaria can develop into a severe disease, often resulting in death. *Anopheles gambiae* is the main malaria vector in Africa. Researchers from the *Anopheles gambiae* 1000 Genome Project sequenced 765 cases of *Anopheles gambiae* from 15 regions in 8 countries in Africa. Studies have shown that mosquitoes are the most diverse eukaryote, with more than 50 million SNPs identified in accessible genomes [1]. These data reveal complex population structures and gene flow patterns, and prove ancient expansion, recent bottlenecks, and local changes in effective population sizes.

In order to commit to malaria control, the "WHO Global Malaria Technical Strategy 2016-2030" adopted by the World Health Assembly in May 2015 provides a technical framework for all malaria-endemic countries [2]. It aims to guide and support regional and national plans in efforts to control and eliminate malaria. In addition to chemical methods to control disease vectors, there is also the use of genetic engineering. Genome is the center of malaria research. At present, the genomes of *Plasmodium falciparum*, the vector mosquito *Anopheles gambiae*, and humans have all been sequenced. Scientists modify the genes of *Anopheles* mosquitoes to shorten their lifespan or to fight off malaria [3]. The study of recombination rate is inseparable from the genetic modification of the human genome.

Accurately estimating the complete picture of the recombination rate of the entire genome in a natural population is the main goal of genomics because the way of contact affects everything from genetic maps to understanding evolutionary history [4]. However, the genetic basis of many complex phenotypes is still largely unknown, mainly due to the environmental factors involved, the polygenicity of traits and the small impact of each associated mutation. Therefore, it is difficult to estimate the characteristics of natural selection in the genome. When only a limited amount of summary statistics is used, these functions are secret, complex, and difficult to detect. To overcome these problems, I explored the application of deep learning in evolutionary biology.

Deep learning is a part of machine learning, which is a reasoning framework based on artificial neural network (ANN). ANN consists of input (also called features) and output (response), which are connected by a series of nodes in hidden layers [5]. The input unit receives various forms and structures of information based on an internal weighting system, while the neural network tries to understand the presented information to generate an output report. Just as humans need rules and guidelines to arrive at results

44 or outputs, artificial neural networks also use a set of learning rules (short
45 for backpropagation errors) called backpropagation to refine their output.
46 The artificial neural network will initially go through a training phase, in
47 which it will learn to recognize patterns in the data from visual, auditory or
48 textual aspects. In this monitoring phase, the network compares the output
49 of actual production with the output of expected production. The difference
50 between the two results can be adjusted by backpropagation. This means
51 that the network will run in reverse, from the output unit to the input
52 unit, to adjust its connection weights between the units until the difference
53 between the actual result and the expected result produces the smallest pos-
54 sible error. After training, the ANN can predict the response based on any
55 new data input. Andrew Ng from Coursera said that deep learning is the
56 first class of algorithms. If you feed them more data, performance just keeps
57 getting better [6].

58
59 However, when using deep learning algorithms, the information needs
60 to be simplified into summary statistics to reduce the data dimension and
61 capture most of the relevant information. For genomic data, an alterna-
62 tive method is to make full use of genomic information and process vertical
63 genomic data through image representation. Each color, called a discrete
64 value, is used to define the appearance of a specific allele. In this way, the
65 detection of markers can be directly transformed into the problem of pattern
66 recognition in image analysis. Under this data representation, Convolutional
67 Neural Network (CNN) is the most suitable algorithm for feature extraction
68 and prediction.

69
70 CNN is a branch of ANN designed specifically for processing images.
71 Since each pixel will be considered as a unique feature, standard artificial
72 neural networks will become unnecessarily complex. Convolutional neural
73 networks can have tens or hundreds of layers, and each layer learns to detect
74 different features of the image. Apply the filter to each training image at
75 a different resolution and use the output of each convolution image as the
76 input to the next layer. Filters can start with very simple functions (such
77 as brightness and edges) and increase the complexity of the function that
78 uniquely defines the object. CNN uses multiple layers of filtering (called
79 convolution), and each layer processes adjacent pixels grouped in a window
80 and then moves them to cover the entire image [7]. Then the weights asso-
81 ciated with each filter are repeatedly adjusted during the training process
82 to detect information-rich local patterns. Therefore, the convolutional layer
83 also has the additional purpose of automatically extracting information fea-
84 tures, and then passing this information as an input unit to several fully
85 connected layers for prediction.

86
87 In this report, I explored the application of deep learning in evolutionary

88 biology, and implemented a project called ImaGene, which applies convo-
89 lutional neural networks to population genomic data to detect and quan-
90 tify natural selection, estimate recombination rates in *Anopheles gambiae*
91 genomes.

92

93 **2 methods and data**

94 **2.1 Computing tools**

95 In this project, I use three computing tools, python latex and Script. Python
96 is used to filter data from sample data and implement CNN. Script is used
97 to generate simulate data via msms [8]. Latex is used to write report, make
98 the format more normalization.

99

100 **2.2 Data**

101 In this project, I use ag10000G as a sample data. Ag10000G samples are
102 sequenced by the Wellcome Trust Sanger Institute Parasites and Microbes
103 Programme(link is external) using Illumina high-throughput technology [9].
104 The sequence data are then used to discover genetic variation between sam-
105 ples and to make genotype calls. Ag10000G is an international collaboration
106 using whole genome deep sequencing to provide a high-resolution view of
107 genetic variation in natural populations of *Anopheles gambiae*, the princi-
108 pal vector of *Plasmodium falciparum* malaria in Africa.

109

110 **2.3 Methods**

111 The propose of this project is use CNN to estimate recombination rates in
112 *Anopheles gambiae* genomes, I implement it in a user-friendly open source
113 program, ImaGene, available at <https://github.com/mfumagalli/ImaGene> .

114

115 Step1: using vcftools get sample data, which is meaningful from ag10000G,
116 read this data from vcf file and store it to Imagene objects. vcftools is a
117 program designed for working with VCF files, the purpose of vcftools is to
118 provide an easy-to-access method for using complex genetic variation data
119 in the form of VCF files. Because some reason, I cannot download complete
120 data. I just download data of *anopheles gambiae* from Ghana. After filter-
121 ing with vcftools, I found that the number of genomes is 67. However, I
122 found that there are many genes in this genetic data that are meaningless,
123 and I need to delete the genotypes whose flag is. And pass. I use `-max-`
124 `missing-count 1` to exclude sites with missing genotypes over all individuals,

125 and use `--recode-INFO-all--stdout` to get the new vcf file as a sample data
126 for this project. After getting the sample data, read this data and store it
127 to Imagen objects.

128

129 Step2: use txt file to store parameters, use script files to call txt file and
130 generate simulation data via msms. In this step, write a txt file called params
131 to store parameters. This parameters including the path to msms.jar, direc-
132 tory of where I will store the simulation data, reference effective population
133 size, length of the locus in bp, mutation rate in $4 * N_e * L$ scale, number
134 of chromosomal copies to simulate as a locus and sample size. For perform-
135 ing a binary or multiclass classification, it need to define these parameters,
136 RHORANGE is range and step for the recombination rate, NREPL is the
137 number of replicates (simulations) per value of recombination rate to be
138 estimated, this data should more than 1000, it will get the better result.
139 Meanwhile, it needs to define NBATCH which is the number of batches
140 for each simulation. For processing msms, it needs to define NTHREADS
141 which is a thread of msms. Generatedataset.sh is a script which can accepts
142 an input file which specifies the parameters of the simulations. Finally use
143 subprocess.call to perform population genomic data simulation, the script
144 divides the simulation into different batches for later training.

145

146 Step3: run and process simulation data to be used for training the CNN
147 Before training, please parse the encoded simulation file from msms and con-
148 vert it into a binary matrix. The method of processing population genomic
149 data is through image representation. The image is a three-dimensional ma-
150 trix, and the third dimension is color. Since most genetic variations are in
151 dialogue, it is very convenient to convert such full-color images to black and
152 white images. Now reduce the color dimension to length 1, and each pixel
153 encodes the frequency of one of the two alleles in a locus. If the ancestral
154 state of the genome is unknown, recode the alignment so that the allele with
155 the highest frequency in each column is converted to zero, and the allele with
156 the lowest frequency is converted to one. Check the sample allele frequency
157 of the selected allele and choose to impose the middle part of the region.
158 Then, I used different criteria to sort the rows and columns. The data is
159 classified and divided into three categories for subsequent training. In this
160 project, I am divided into 20, 40, 60.

161

162 Step4: implement, train and evaluate the CNN In this step, ImaGene
163 interacts directly with Keras [10] model. ImaGene provides utilities to mon-
164 itor and evaluate training, and it uses standard methods to perform binary
165 and multi-class classification tasks. Use iteration methods to build the CNN
166 model. In this model, all data needs to be classified and trained in batches.
167 The total number of training samples is 300, but the data of $300 * 67 * 205$
168 is very large. In order to reduce the time and memory consumption of the

169 computer when running, select batch training. A batchsize of 32 means
170 that 32 samples are a group. For this project, I implemented a CNN with
171 two 64-unit and one 128-unit 2D convolutional layers with a kernel size of
172 33 and a stride of 1 1. The convolution layer extracts local features in the
173 picture through the filtering of the convolution kernel. The pooling layer is
174 implemented every time the 2D convolutional layer is completed, and the
175 kernel size is 22. A total of three pooling layers have been implemented.
176 The pooling layer can reduce the data dimension more effectively than the
177 convolutional layer. This can not only greatly reduce the amount of calcu-
178 lation, but also effectively avoid overfitting. Finally, a fully connected layer
179 with 64 units is applied to output the result. After each training, I will get
180 the updated score, including the accuracy and loss. When the number of
181 simulations is all completed training, I will get the final generated model.
182 Next, I need to optimize the model, reduce the loss value and increase the
183 value of accuracy by changing the parameters, and plot a confusion matrix.
184 The confusion matrix can more intuitively show the relationship between
185 the predicted value and the true value. If the color is darker, it means the
186 match is higher.

187

188 Step5: substitute sample data into the CNN model and estimate the
189 recombine rate In the end, I used the trained model to estimate the re-
190 combination rate of the natural selection category. First, adjust the size of
191 the actual data to match the data used for training. After that, the data
192 is converted into the required format. The model will output a probability
193 vector for the given data, and the sum of this vector is 1

194

195 **3 Results**

196 **3.1 Image representation of *Anopheles gambiae* genome data**

197 After read sample data from vcf file and store it to Image objects, genomic
198 data has been transformed into an image representation. I have one image
199 with 134 rows (equivalent to the number of sampled chromosomal copies)
200 and 2200 columns representing all genomic positions reported.

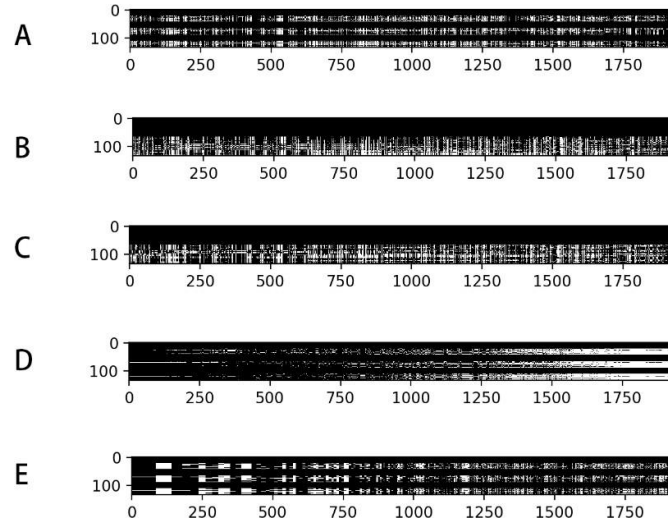


Figure 1

201

202 Figure 1 shows the genome data of anopheles gambiae. it has 134 rows
 203 and 1926 columns. Every row means one chromosome, because anopheles
 204 gambiae is diploid, the number of genomes is 67, the number of chromosomes
 205 is 134. Every column means genetic locus, it means in this chromosome has
 206 1926 locus. A is the data without sort. The data shown in this picture
 207 has no rules to follow, and it is messy. Since the order is arbitrary, no
 208 meaning can be seen. So, I can sort rows (and columns) according to several
 209 conditions to help feature extraction. B is the data with sort only rows by
 210 their distance, arrange from top to bottom according to the distance between
 211 the longest appearing lines. C is the data with sort only rows by their
 212 frequency, arrange from top to bottom according most frequent haplotypes.
 213 With these two sorts, the picture can express the genome more intuitively.
 214 On the other hand, each column means a locus and contains information
 215 about the relative position of the polymorphism along the locus. The sorting
 216 of the columns contains information about linkage disequilibrium, which can
 217 provide information for detecting selective sweeps. However, this order is
 218 also affected by mutation and recombination events. So, we sort the column.
 219 D is the data with sort only columns by their distance, arrange from left to
 220 right according to the distance between the longest appearing lines. E is the
 221 data with sort only columns by their frequency, arrange from left to right
 222 according most frequent haplotypes.

3.2 Evaluate pipelines under various data and learning configurations

We evaluate whether to change a single parameter to test the choice. In this project, the purpose is to evaluate the accuracy of using CNN to detect and quantify recombination rate events under different learning and data manipulation settings. In this program, 300 pictures can be obtained through msms simulation, each row represents a haplotype randomly sampled from the population. Processing images by sorting rows and columns can improve detection.

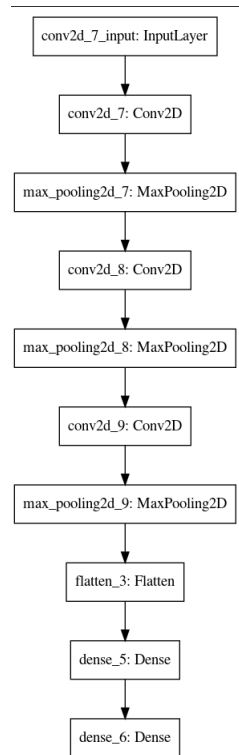


Figure 2

Like figure 2, in this project, I implemented a CNN with two 64-unit and a 128-unit 2D convolutional layer with a kernel size of 33 and a stride of 11. The pooling layer is implemented every time the 2D convolutional layer is completed, and the kernel size is 22. A total of three pooling layers have been implemented. Finally, a fully connected layer with 64 units is applied to output the result. At the same time, I deleted the column corresponding to the allele frequency less than 0.015. After sorting, we adjust the size of all images to 67205 pixels. In order to prevent over-fitting, the "real-time simulation" method is used. The algorithm will train the newly generated

242 data in each cycle. At the same time, it retains the complete training data
243 set. Each cycle separates 10

244

245 In order to make my results more convincing, I use the controlled variable
246 method. The controlled variable method is a variance reduction technique
247 used in the Monte Carlo method. It uses the information about the error in
248 the estimation of the known quantity to reduce the error in the estimation of
249 the unknown quantity. This method can eliminate interference and directly
250 reveal the influence of a single factor on the change of the research object.
251 Therefore, in this project, I will make four variable changes to detect and
252 quantify the accuracy of the recombination rate event.

253

254 The initialization parameters are the reference effective population size
255 equals 100000, the length of the locus in bp is 100000, mutation rate in
256 $4 \times N_e \times \text{LEN}$ equals 40, the number of chromosomal copies to simulate equals
257 100, the range and step for the recombination rate is 20 (in $4N \times \text{length}$ units,
258 weak), 40 (medium), 60 (strong). The number of replicates per value of re-
259 combination rate to be estimated equals 1000, the number of batches for
260 each simulation equals 10, thread of msms is 4 and every simulation trained
261 on 1-epoch model.

262

263 In the first time, the parameter which is number of batches for each sim-
264 ulation be changed. 10,20,30. In these learning, the total number of images
265 is 600.

266

267 The number of batches for each simulation equals 10(model 1):

268 The testing results (loss:0.96794 accuracy:0.63833)

269 The value of predict (0.73869,0.24316,0.01815)

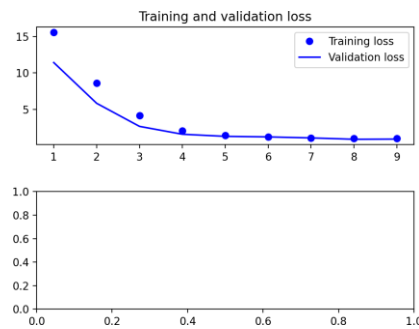


Figure 3

270

271

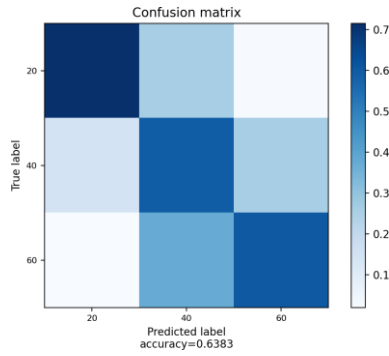


Figure 4

272

273 Figure 3 shows that the loss value during training and verification in 10
 274 simulations, it shows the value of loss from 15 reduce to 0.9. and Figure
 275 4 shows the confusion matrix. The darkest color means this piece has the
 276 highest matching degree. From confusion matrix, I found if choose 20, the
 277 accuracy of selection is the highest. When I use this model in the real data,
 278 I found the result is as predicted by the model. This means that the model
 279 can be used to identify fragments of the genome with a weak recombination
 280 rate.

281

282 The number of batches for each simulation equals 20(model 2):

283 The testing results (loss:0.86217 accuracy:0.66833)

284 The value of predict (0.82825,0.15678,0.01497)

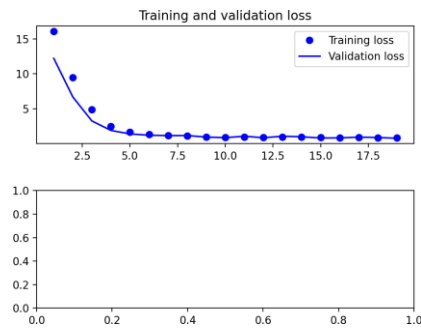


Figure 5

285

286

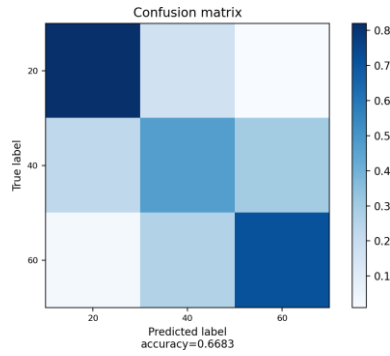


Figure 6

287

288 Figure 5 shows that the loss value during training and verification in 20
 289 simulations, it shows the value of loss from 15 reduce to 0.86. and Figure 6
 290 shows the confusion matrix. From confusion matrix, I found if choose 20, the
 291 accuracy of selection is the highest. When I use this model in the real data, I
 292 found the result is as predicted by the model. This means that the model can
 293 be used to identify fragments of the genome with a weak recombination rate

294

295 The number of batches for each simulation equals 30(model 3):

296 The testing results (loss:0.92432 accuracy:0.63833)

297 The value of predict (0.69016,0.27488,0.03496)

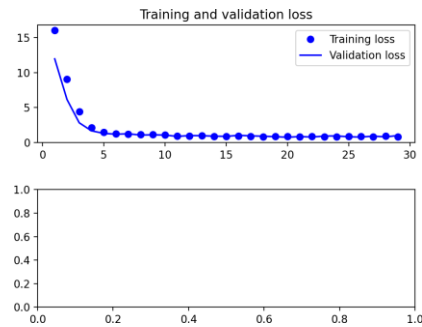


Figure 7

298

299

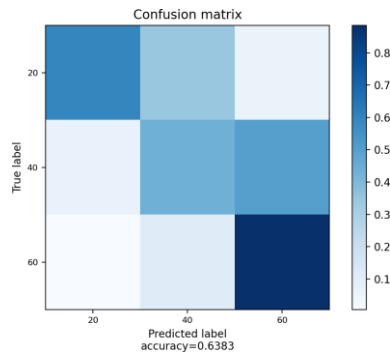


Figure 8

300

301 Figure 7 shows that the loss value during training and verification in 30
 302 simulations, it shows the value of loss from 15 reduce to 0.92. and Figure 8
 303 shows the confusion matrix. From confusion matrix, I found if choose 60, the
 304 accuracy of selection is the highest. When I use this model in the real data,
 305 I found that the results are not the same as predicted by the model. This
 306 means that the model cannot be used to identify fragments of the genome
 307 with a recombination rate

308

309 Next, the parameter which is the length of the locus in bp is be changed,
 310 100000, 200000, 300000. Cause that, mutation rate in $4*Ne*LEN$ is be
 311 changed, 40, 80, 120.

312

313 Length of the locus in bp equals 100000, mutation rate in $4*Ne*LEN$
 314 equals 40(model 4):

315 The testing results (loss:0.96794 accuracy:0.63833)

316 The value of predict (0.73869,0.24316,0.01815)

317 Same as model 1

318

319 Length of the locus in bp equals 200000, mutation rate in $4*Ne*LEN$
 320 equals 80(model 5):

321 The testing results (loss:1.15424 accuracy:0.49500)

322 The value of predict (0.83262,0.13596,0.03141)

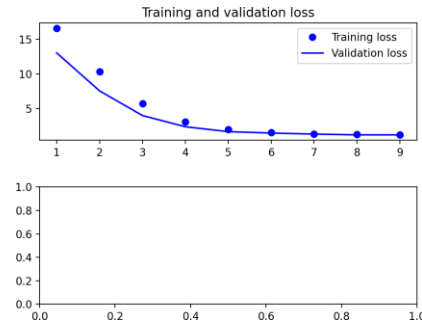


Figure 9

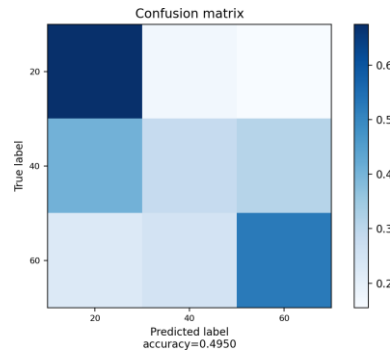


Figure 10

Figure 9 shows that the loss value during training and verification when length of the locus in bp equals 200000, mutation rate in $4 \times Ne \times LEN$ equals 80, it shows the value of loss from 15 reduce to 0.83. and Figure 10 shows the confusion matrix. From confusion matrix, I found if choose 20, the accuracy of selection is the highest. When I use this model in the real data, I found result is as predicted by the model. This means that the model can be used to identify fragments of the genome with a weak recombination rate.

Length of the locus in bp equals 300000, mutation rate in $4 \times Ne \times LEN$ equals 120(model 6):

The testing results (loss:1.12191 accuracy:0.51833)

The value of predict (0.65387,0.27168,0.07445)

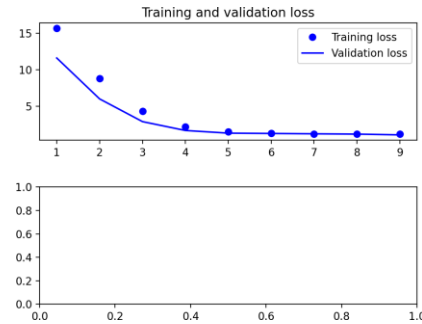


Figure 11

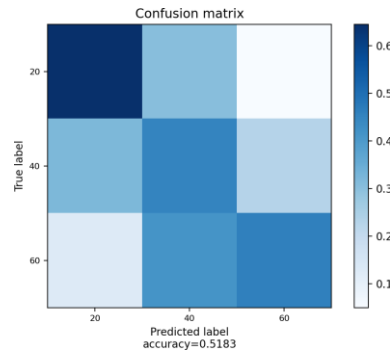


Figure 12

Figure 11 shows that the loss value during training and verification when length of the locus in bp equals 300000, mutation rate in $4 \times N_e \times L \times \mu$ equals 120, it shows the value of loss from 15 reduce to 1.12. and Figure 12 shows the confusion matrix. From confusion matrix, I found if choose 20, the accuracy of selection is the highest. When I use this model in the real data, I found result is as predicted by the model. This means that the model can be used to identify fragments of the genome with a weak recombination rate.

Then, the parameter which is the number of replicates per value of recombination rate to be estimated is be changed. 1000, 3000, 5000.

The number of replicates per value of recombination rate equals 1000(model 7)

Same as model 1

The number of replicates per value of recombination rate equals 3000(model

357 8)
 358 The testing results (loss:1.12190 accuracy:0.51833)
 359 The value of predict (0.65387,0.27168,0.07444)

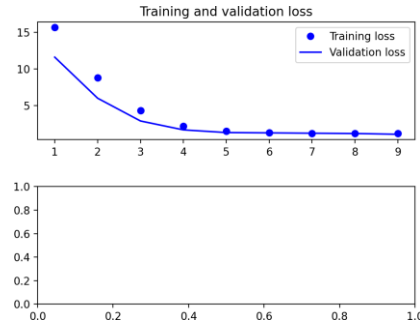


Figure 13

360
 361

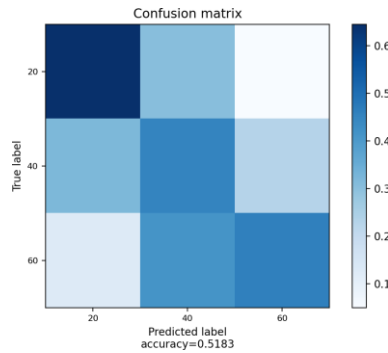


Figure 14

362
 363 Figure 13 shows that the loss value during training and verification when the
 364 number of replicates per value of recombination rate equals 3000, it shows
 365 the value of loss from 15 reduce to 1.12. and Figure 14 shows the confusion
 366 matrix. From confusion matrix, I found if choose 20, the accuracy of selec-
 367 tion is the highest. When I use this model in the real data, I found result
 368 is as predicted by the model. This means that the model can be used to
 369 identify fragments of the genome with a weak recombination rate.

370
 371 The number of replicates per value of recombination rate equals 5000 (model
 372 9)
 373 The testing results (loss:0.99253 accuracy:0.63667)
 374 The value of predict (0.80850,0.17148,0.02001)

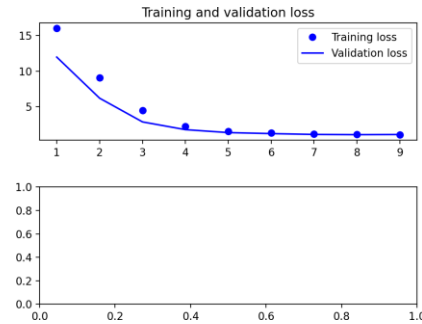


Figure 15

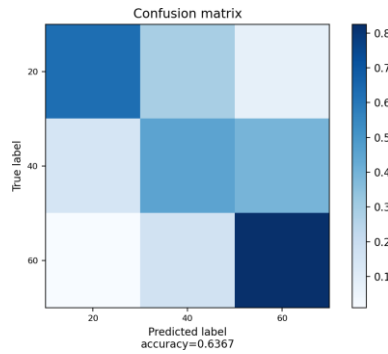


Figure 16

Figure 15 shows that the loss value during training and verification when the number of replicates per value of recombination rate equals 5000, it shows the value of loss from 15 reduce to 0.99. and Figure 16 shows the confusion matrix. From confusion matrix, I found if choose 60, the accuracy of selection is the highest. When I use this model in the real data, I found that the results are not the same as predicted by the model. This means that the model cannot be used to identify fragments of the genome with a recombination rate.

Finally, I changed the number of training sessions per simulation. 1-epoch model. 2-epoch model. 3-epoch model.

Every simulation trained on 1-epoch model (model 10): same as model 1

Every simulation trained on 2-epoch model (model 11):

394 The testing results (loss:0.92612 accuracy:0.64667)
 395 The value of predict (0.76125,0.21737,0.02139)

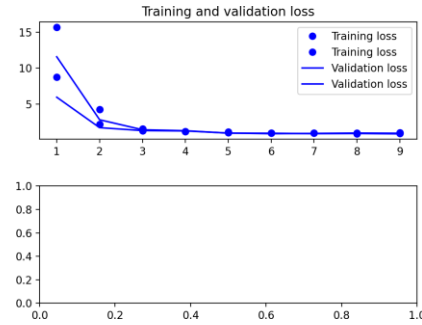


Figure 17

396

397

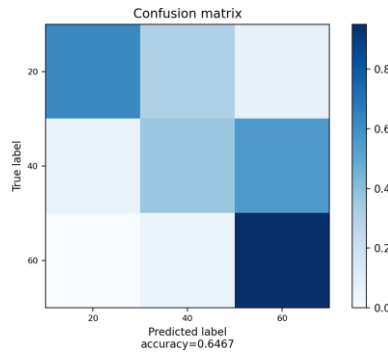


Figure 18

398

399 Figure 17 shows that the loss value during training and verification when
 400 every simulation trained on 2-epoch model, it shows in the first training,
 401 the value of loss from 15 reduce to 0.92. In the second training, the value of
 402 loss from 8 to 0.92. Although the initial loss values of the two trainings are
 403 different, they will eventually fall to the same loss value. Figure 18 shows
 404 the confusion matrix. From confusion matrix, I found if choose 60, the ac-
 405 curacy of selection is the highest. When I use this model in the real data,
 406 I found that the results are not the same as predicted by the model. This
 407 means that the model cannot be used to identify fragments of the genome
 408 with a recombination rate.

409

410 Every simulation trained on 3-epoch model (model 12):

411 The testing results (loss:0.85148 accuracy:0.67167)

412 The value of predict (0.73038,0.24809,0.02154)

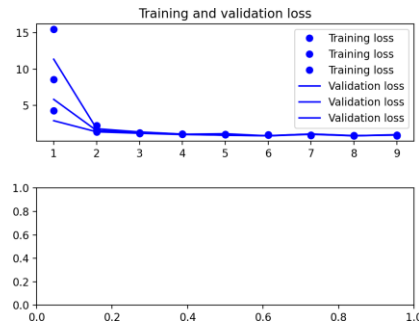


Figure 19

413

414

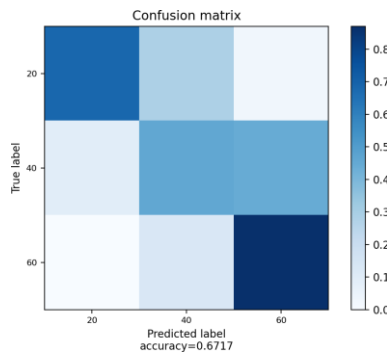


Figure 20

415

416 Figure 19 shows that the loss value during training and verification when
 417 every simulation trained on 3-epoch model, it shows in the first training,
 418 the value of loss from 15 reduce to 0.85. In the second training, the value
 419 of loss from 8 to 0.85. In the last training, the value of loss from 4 to 0.85.
 420 Although the initial loss values of the two trainings are different, they will
 421 eventually fall to the same loss value. Figure 18 shows the confusion matrix.
 422 From confusion matrix, I found if choose 60, the accuracy of selection is the
 423 highest. When I use this model in the real data, I found that the results are
 424 not the same as predicted by the model. This means that the model cannot
 425 be used to identify fragments of the genome with a recombination rate.

426

427 According to the above model, I think the most model is model 2. Ex-
 428 cept model 3, model 9, model 11 and model 12, other model can used to
 429 identify fragments of the genome with a recombination rate, but if I choose

430 model 2, the value of accuracy is highest.

431

432 **4 Discussion**

433 In this research, I used ImaGene program, which uses keras to train the
434 CNN model. In this project, I proved that convolutional neural networks
435 have a strong advantage in detecting and quantifying recombination rates of
436 *Anopheles gambiae* genomes. The use of convolutional neural networks can
437 effectively learn the corresponding features from many samples, avoiding the
438 complex feature extraction process, and more efficiently estimating the gene
439 recombination rate.

440

441 However, the project still has many parts that need to be improved and
442 expanded to make its predictions more accurate and reliable than the pre-
443 dictions introduced in this article. The random initialization method we
444 choose for setting the initial network parameters before training may not be
445 optimal. Although various tests have been carried out, it is proved that the
446 initial value will optimize the model with modifying the number of simula-
447 tions. However, the number of tests performed is too few, and more tests
448 are applied to achieve maximum verification accuracy.

449

450 At the same time, if the units of the convolutional layer are increased, it
451 may have obvious computational advantages. However, further research is
452 needed to evaluate the effect of adjusting the unit size of the convolutional
453 layer, and the trade-off between computational speed and accuracy when
454 increasing the unit size.

455 Finally, CNN can also be improved, such as using OctConv. OctConv is
456 like the "compressor" of Convolutional Neural Network (CNN) [11]. Us-
457 ing it to replace traditional convolution can improve the effect while saving
458 the consumption of computing resources. For example, for a classic image
459 recognition algorithm, replacing the traditional convolution, the recognition
460 accuracy on ImageNet can be improved by 1.2

461

462 **5 Conclusions**

463 In this research, I used ImaGene, used msms to implement simulation data,
464 and used kares to implement CNN training, detection and quantification
465 of recombination rates. I also learned how data processing and learning
466 settings affect prediction accuracy. The discoveries of these efforts will help
467 to change the genes of *Anopheles gambiae* genomes, better solve the African
468 Malaria, and reveal new associations with complex diseases.

469 References

- 470 [1] Anopheles gambiae 1000 Genomes Consortium et al. Genetic diversity
471 of the african malaria vector anopheles gambiae. *Nature*, 552(7683):96–
472 100, 2017.
- 473 [2] Roll Back Malaria et al. World malaria report 2005. *World Health
474 Organization and UNICEF*, 2005.
- 475 [3] Junitsu Ito, Anil Ghosh, Luciano A Moreira, Ernst A Wimmer, and
476 Marcelo Jacobs-Lorena. Transgenic anopheline mosquitoes impaired in
477 transmission of a malaria parasite. *Nature*, 417(6887):452–455, 2002.
- 478 [4] Luis Torada, Lucrezia Lorenzon, Alice Beddis, Ulas Isildak, Linda Pat-
479 tini, Sara Mathieson, and Matteo Fumagalli. Imagenet: a convolutional
480 neural network to quantify natural selection from genomic data. *BMC
481 bioinformatics*, 20(9):337, 2019.
- 482 [5] S Agatonovic-Kustrin and R Beresford. Basic concepts of artificial
483 neural network (ann) modeling and its application in pharmaceutical
484 research. *Journal of pharmaceutical and biomedical analysis*, 22(5):717–
485 727, 2000.
- 486 [6] Jason Brownlee. What is deep learning. *Machine learning mastery*, 16,
487 2016.
- 488 [7] Andrej Karpathy. Convolutional neural networks (cnns/convnets).
489 *CS231n Convolutional Neural Networks for Visual Recognition*, 2016.
- 490 [8] Gregory Ewing and Joachim Hermisson. Msms: a coalescent simulation
491 program including recombination, demographic structure and selection
492 at a single locus. *Bioinformatics*, 26(16):2064–2065, 2010.
- 493 [9] Anopheles gambiae 1000 Genomes Consortium et al. Ag1000g phase 1
494 ar2 data release. malariagen, 2014.
- 495 [10] Antonio Gulli and Sujit Pal. *Deep learning with Keras*. Packt Publish-
496 ing Ltd, 2017.
- 497 [11] Yunpeng Chen, Haoqi Fan, Bing Xu, Zhicheng Yan, Yannis Kalantidis,
498 Marcus Rohrbach, Shuicheng Yan, and Jiashi Feng. Drop an octave:
499 Reducing spatial redundancy in convolutional neural networks with oc-
500 tave convolution. In *Proceedings of the IEEE International Conference
501 on Computer Vision*, pages 3435–3444, 2019.