

**ANCESTRY INFORMATIVE MARKERS  
SINGLE NUCLEOTIDE POLYMORPHISMS  
PANEL FOR ANCESTRY ESTIMATION IN THE  
MALAY POPULATION**

**PADILLAH BINTI YAHYA**

**UNIVERSITI SAINS MALAYSIA**

**2021**

**ANCESTRY INFORMATIVE MARKERS  
SINGLE NUCLEOTIDE POLYMORPHISMS  
PANEL FOR ANCESTRY ESTIMATION IN THE  
MALAY POPULATION**

by

**PADILLAH BINTI YAHYA**

**Thesis submitted in fulfilment of the requirements**

**for the Degree of**

**Doctor of Philosophy**

**August 2021**

## ACKNOWLEDGEMENTS

In the name of Allah S.W.T, the Most Gracious and the Most Merciful. All praises to Allah, for giving me strength, patience, perseverance and optimism throughout my postgraduate journey. My biggest gratitude to Him, with His permission this thesis has finally been completed.

My deepest appreciation to my main supervisor, Prof. Dr. Zilfalil Alwi, and my co-supervisors, Associate Prof. Dr. Sarina Sulong and Associate Prof. Dr. Azian Harun for their guidance, support and kind supervision throughout my study. My special thanks to Dr. Sissades Tongsima and his team members, Pongsakorn Wangkumhang, Alisa Wilantho, and Chumpol Ngamphiw at the National Center for Genetic Engineering and Biotechnology (BIOTEC), Thailand Science Park, Pathum Thani Thailand, for their time, guidance, lessons and technical supports in completion this study. Special thanks also to W. Nur Hatin W. Isa, Nur Shafawati Ab Rajab and Nurfazreen Mohd Nasir who have provided me all the SNPs genotype data used in this study and help me a lot in data analyses. Thanks to my great and helpful friends, Nurul Fatimah Azman and Diana Rashid and all students and staff of MyHVP, USM for helping me throughout my study. I am deeply thankful to my families and my best friend, Juriah Mohamed, for standing by my side when times get hard. I couldn't have achieved this without your support. Finally, I would like to thanks Public Service Department Malaysia (JPA) for the HLP scholarship, USM for providing all the facilities to carry out my study and Universiti Sains Malaysia Apex Grant: 1002/PPSP/910343 and NTU Grant (Muhammed Ariff Research Grant (MAS): 304.PPSP.6150148.N119 for supporting this study.

## TABLE OF CONTENTS

ACKNOWLEDGEMENTS .....	ii
TABLE OF CONTENTS .....	iii
LIST OF TABLES .....	viii
LIST OF FIGURES .....	xii
LIST OF SYMBOLS AND ABBREVIATIONS .....	xxxiv
LIST OF APPENDICES .....	xxxix
ABSTRAK .....	xli
ABSTRACT .....	xlii

### CHAPTER 1-INTRODUCTION

1.1	Research background .....	1
1.2	Problem statement .....	7
1.3	Research justifications .....	8
1.4	Objectives.....	10
1.4.1	General Objective .....	10
1.4.2	Specific Objectives .....	11

### CHAPTER 2-LITERATURE REVIEW

2.1	Ancestry .....	12
2.2	Genetic markers used in ancestry studies.....	14
2.2.1	Single Nucleotide Polymorphisms (SNPs).....	16
2.2.1(a)	SNPs as ancestry informative marker.....	20
2.2.1(b)	SNPs genotyping .....	23
2.3	Ancestry informative markers (AIMs).....	24
2.4	Methods for the selection of AIMs .....	29
2.4.1	Informativeness for assignment ( $I_n$ ) .....	31
2.4.2	Principal component analysis (PCA).....	33
2.4.3	Iterative pruning principal component analysis (ipPCA).....	36
2.4.4	Pairwise $F_{ST}$ .....	39
2.4.5	SNPs associated with pigmentation genes .....	40

2.5	Software and algorithm used in AIMs analysis .....	42
2.5.1	STRUCTURE.....	43
2.5.2	ADMIXTURE .....	45
2.5.3	FRAPPE .....	47
2.5.4	PLINK .....	48
2.5.5	Local Ancestry in admixed Populations (LAMP).....	50
2.5.6	SABER .....	51
2.5.7	HAPMIX .....	52
2.5.8	ANCESTRYMAP .....	52
2.5.9	Waikato Environment for Knowledge Analysis (WEKA).....	53
2.5.10	K-nearest-neighbor (KNN).....	55
2.5.11	Naïve Bayes (NB).....	56
2.6	Malay ancestry .....	58
2.7	Theory of the origin of the Malays .....	60
2.8	Genetic Ancestry of the Malays .....	65

## CHAPTER 3-MATERIALS AND METHODS

3.1	Single Nucleotide Polymorphisms (SNPs) genotype datasets .....	69
3.2	Sample size.....	73
3.3	Population datasets.....	73
3.3.1	MyHVP SNP dataset .....	74
3.3.2	Pan-Asian SNP Consortium dataset .....	76
3.3.3	SGVP SNP dataset .....	77
3.3.4	The International HapMap Phase 3 Project SNP dataset .....	78
3.4	Data analysis .....	79
3.4.1	Analysis of Affymetrix GeneChip Mapping Xba 50 K Array data. 83	
3.4.1(a)	SNPs data merging .....	84
3.4.1(b)	SNPs data filtration.....	88
3.4.2	Ancestry informative marker (AIM) SNPs selection from the Affymetrix GeneChip Mapping Xba 50 K Array data.....	97
3.4.2(a)	Selection of AIM SNPs based on $F_{ST}$ .....	97
3.4.2(a)(i)	ipPCA analysis .....	98
3.4.2(a)(ii)	Top $F_{ST}$ calculation.....	105

	3.4.2(a)(iii)	WEKA suite for ancestry-predictive modeling .....	111
	3.4.2(a)(iv)	ADMIXTURE Analysis .....	120
	3.4.2(b)	Selection of AIM SNPs based on $I_n$ .....	122
	3.4.2(b)(i)	Principal Component Analysis (PCA).....	124
	3.4.2(b)(ii)	Informativeness for assignment ( $I_n$ ) .....	129
	3.4.2(b)(iii)	K-nearest-neighbor (KNN).....	132
	3.4.2(c)	Selection of AIM SNPs based on PCA-correlated SNPs (PCAIMs) .....	133
3.5		Analysis of Affymetrix SNP 6 Array and the OMNI 2.5 Illumina data .....	135
	3.5.1	SNPs data merging and filtration .....	139
	3.5.2	Malay AIM SNPs panel selection .....	146

## CHAPTER 4-RESULTS

4.1		Analysis of the Affymetrix GeneChip Mapping Xba 50 K Single Nucleotide Polymorphism (SNP) Array data. ....	150
	4.1.1	Analysis of population stratification using PCA .....	151
	4.1.2	Analysis of population stratification using ADMIXTURE.....	154
	4.1.3	Genetic structure of Malay population as revealed by the ipPCA analysis .....	158
	4.1.4	ADMIXTURE analysis of sub-population SP1 to SP12.....	168
	4.1.5	Pairwise $F_{ST}$ analysis of SP1 to SP12 .....	171
	4.1.6	Selecting Ancestry Informative Marker (AIM) SNPs for the Malays from the Affymetrix GeneChip Mapping Xba 50-K SNPs data ...	173
	4.1.6(a)	AIM SNPs based on ipPCA- $F_{ST}$ .....	173
	4.1.6(b)	ADMIXTURE analysis of all five ancestry Models	181
	4.1.6(c)	AIM SNPs based on $I_n$ .....	192
	4.1.6(d)	AIM SNPs based on PCAIMs .....	204
	4.1.6(e)	Comparison of the performance of PCAIMs and $I_n$ methods.....	214
4.2		Analysis of the Affymetrix SNP-6 Array data.....	216
	4.2.1	Analysis of population stratification using PCA .....	218
	4.2.2	Analysis of population stratification using ADMIXTURE.....	220

4.2.3	Analysis of population stratification using ipPCA .....	224
4.2.4	ADMIXTURE analysis of sub-population SP1 to SP11 .....	232
4.2.5	Pairwise $F_{ST}$ analysis of SP1 to SP11 .....	236
4.2.6	Selecting Ancestry Informative Marker (AIM) SNPs for the Malays from the Affymetrix SNP-6 data .....	237
4.2.6(a)	AIM SNPs based on ipPCA- $F_{ST}$ .....	237
4.2.6(b)	ADMIXTURE analysis of Model 1 to Model 5 .....	244
4.2.6(c)	AIM SNPs based on $I_n$ .....	259
4.2.6(d)	AIM SNPs based on PCA-correlated SNPs .....	272
4.2.6(e)	The combination of the minimal Malay AIM SNPs panel .....	285
4.3	WEKA suite as the predictive-ancestry model .....	289

## CHAPTER 5-DISCUSSION

5.1	Malay AIM SNPs from Affymetrix GeneChip Mapping Xba 50 K and Affymetrix SNP 6 Array SNPs databases .....	296
5.1.1	Malay AIM SNPs from Genome-wide Affymetrix 50K SNPs Array .....	301
5.1.2	Malay AIM SNPs from Genome-wide Affymetrix SNP 6 Array .	304
5.2	Methods of selecting the Malay AIM SNPs .....	306
5.2.1	ipPCA and $F_{ST}$ .....	306
5.2.2	$I_n$ .....	308
5.2.3	PCA-correlated SNPs (PCAIMs) .....	317
5.2.4	Comparison of the three methods ipPCA- $F_{ST}$ , $I_n$ and PCAIMs.....	326
5.2.5	Overlapping AIM SNPs of ipPCA- $F_{ST}$ , $I_n$ and PCAIMs methods	334
5.3	The performance of the panel of Malay AIM SNPs on the eight sub-ethnic groups of Malay population .....	337
5.4	The panel of Malay AIM SNPs related to other published AIM SNPs .....	347
5.5	Limitation of the study .....	350
5.6	Recommendations for future research .....	352

## CHAPTER 6-CONCLUSION

6.1	Outcome of the study .....	356
6.2	Contribution of the study to the medical genetic and forensic communities	358

6.3 Novelty of the study ..... 359

REFERENCES..... 360

APPENDICES

LIST OF ABSTRACTS AND PUBLICATIONS



## LIST OF TABLES

		<b>Page</b>
Table 3.1:	List of populations and number of individual used in this study (Affymetrix GeneChip Mapping Xba 50 K Array). .....	71
Table 3.2:	List of populations and number of individual used in this study (Affymetrix SNP 6 Array ). .....	72
Table 4.1:	Number of individuals left after quality control analysis. ....	151
Table 4.2:	$F_{ST}$ distances between the 12 sub-populations. ....	172
Table 4.3:	The accuracy of each AIM SNPs model selected based on $F_{ST}$ . .....	175
Table 4.4:	The accuracy of classification of individuals to their sub-population using Model 1 (144 AIM SNPs).....	176
Table 4.5:	The accuracy of classification of individuals to their sub-population using Model 2 (229 AIM SNPs).....	177
Table 4.6:	The accuracy of classification of individuals to their sub-population using Model 3 (433 AIM SNPs).....	178
Table 4.7:	The accuracy of classification of individuals to their sub-population using Model 4 (1772 AIM SNPs).....	179
Table 4.8:	The accuracy of classification of individuals to their sub-population using Model 5 (3145 AIM SNPs).....	180
Table 4.9:	The accuracy of each AIM SNPs model selected based on $I_n$ score. ....	200

Table 4.10:	The accuracy of classification of individuals to their sub-population using 250 AIM SNPs selected based on $I_n$ score.....	200
Table 4.11:	The accuracy of classification of individuals to their sub-population using 2000 AIM SNPs selected based on $I_n$ score.....	201
Table 4.12:	The accuracy of classification of individuals (after ipPCA pruning) to their sub-population using 250 AIM SNPs selected based on $I_n$ score.....	202
Table 4.13:	The accuracy of classification of individuals (after ipPCA pruning) to their sub-population using 2000 AIM SNPs selected based on $I_n$ score.....	203
Table 4.14:	The accuracy of each AIM SNPs model selected based on PCA-correlated SNPs.....	210
Table 4.15:	The accuracy of classification of individuals to their sub-population using 250 AIM SNPs selected based on PCA-correlated SNPs.....	210
Table 4.16:	The accuracy of classification of individuals to their sub-population using 2000 AIM SNPs selected based on PCA-correlated SNPs.....	211
Table 4.17:	The accuracy of classification of individuals to their sub-population (after ipPCA pruning) using 250 AIM SNPs selected based on PCA-correlated SNPs.....	212

Table 4.18:	The accuracy of classification of individuals to their sub-population (after ipPCA pruning) using 2000 AIM SNPs selected based on PCA-correlated SNPs. ....	213
Table 4.19:	List of the populations and number of individuals after quality control analysis using PLINK. ....	217
Table 4.20:	$F_{ST}$ value between the sub-populations.....	236
Table 4.21:	The accuracy of each AIM SNPs model selected based on $F_{ST}$ .....	238
Table 4.22:	The accuracy of classification of individuals to their sub-population using Model 1 (101 AIM SNPs).....	239
Table 4.23:	The accuracy of classification of individuals to their sub-population using Model 2 (157 AIM SNPs).....	240
Table 4.24:	The accuracy of classification of individuals to their sub-population using Model 3 (294 AIM SNPs).....	241
Table 4.25:	The accuracy of classification of individuals to their sub-population using Model 4 (1250 AIM SNPs).....	242
Table 4.26:	The accuracy of classification of individuals to their sub-population using Model 5 (2240 AIM SNPs).....	243
Table 4.27:	The accuracy of each AIM SNPs model selected based on $I_n$ score. ....	268
Table 4.28:	The accuracy of classification of individuals to their sub-population based on 100 SNPs selected using $I_n$ method.....	269
Table 4.29:	The accuracy of classification of individuals to their sub-population based on 200 SNPs selected using $I_n$ method.....	270

Table 4.30:	The accuracy of classification of individuals to their sub-population based on 2000 SNPs selected using $I_n$ method.....	271
Table 4.31:	The accuracy of each AIM SNPs model selected based on PCAIMs.....	281
Table 4.32:	The accuracy of classification of individuals to their sub-population based on 100 SNPs selected using PCAIMs. ....	282
Table 4.33:	The accuracy of classification of individuals to their sub-population based on 200 SNPs selected using PCAIMs. ....	283
Table 4.34:	The accuracy of classification of individuals to their sub-population based on 2000 SNPs selected using PCAIMs. ....	284
Table 4.35:	The accuracy of classification of individuals to their sub-population using the combination of 555 SNPs selected based on $F_{ST}$ , $I_n$ and PCAIMs. ....	286
Table 5.1:	The classification accuracy of individuals to each sub-ethnic group of Malays using the 250 AIM SNPs selected based on PCAIMs method.....	343
Table 5.2:	The classification accuracy of individuals to each sub-ethnic group of Malays using the 2000 AIM SNPs selected based on PCAIMs method.....	344
Table 5.3:	The classification accuracy of individuals to each sub-ethnic group of Malays using the 250 AIM SNPs selected based on $I_n$ method.....	345
Table 5.4:	The classification accuracy of individuals to each sub-ethnic group of Malays using the 2000 AIM SNPs selected based on $I_n$ method.....	346

## LIST OF FIGURES

	<b>Page</b>
Figure 2.1: Map of the Malay Archipelago. Available from: <a href="https://www.google.com/url">https://www.google.com/url</a> [Accessed 21 July 2019].....	64
Figure 2.2: Map showing the out of Taiwan theory (red line) (adapted from Oppenheimer and Richards, 2001). .....	65
Figure 3.1: Overview of the data analysis carried out on the Affymetrix GeneChip Mapping Xba 50 K Array SNPs database to identify the AIM SNPs panel for ancestry inference of the Malay population.....	81
Figure 3.2: Overview of the data analysis carried out on the Affymetrix SNP 6 Array and Illumina SNPs database to identify the AIM SNPs panel for ancestry inference of the Malay population.....	82
Figure 3.3: An example of SNPs genotype data (MY-KD population). A and B indicates homozygous allele, H for heterozygous allele and U for missing genotype. ....	84
Figure 3.4: Running the PLINK software in command prompt 'C:\' to create a PED and MAP file to be used for further analysis.....	86
Figure 3.5: The PED output file which combined all genotypes of the 493 individuals from 19 populations.....	87
Figure 3.6: The MAP output file of the combined genotypes which listed all 52,501 SNPs on autosomal chromosome 1 to 22 with the SNP identifier and SNP position. ....	87

Figure 3.7:	The quality control process was carried out using PLINK software where the malay_data was trimmed for minor allele frequency (maf) less than 1%, missingness per marker (geno) more than 5% , HWE $p < 10^{-6}$ and individual with missing genotype (mind) more than 10%. The cleaned SNPs data will be saved as malay_data_qc. ....	90
Figure 3.8:	The quality control log showing the filtering processes performed by the PLINK software. In this filtering process, one individual was removed for low genotyping (mind more than 10%) and 2342 markers were excluded based on HWE test $p < 10^{-6}$ , 1427 SNPs failed missingness test of 5% and 1200 SNPs failed minor allele frequency test of 1%. The remaining SNPs left for further quality control processes were 47,649.....	91
Figure 3.9:	One individual was removed from the PED file and saved in malay_data_qc.irem file. The individual was from MY-BG sub-ethnic group.....	91
Figure 3.10:	The MAP file created by PLINK software (malay_data_qc.map) showing the final list of 47,649 SNPs after the filtration process. ....	92
Figure 3.11:	The PED file created by PLINK software (malay_data_qc.ped) showing list of genotypes left (492) after the filtration process. ....	92

Figure 3.12:	PLINK software listed all individuals failed IBD test at 0.185 threshold value and removed them from the PED file.....	93
Figure 3.13:	Running the LD filtration processes using PLINK software. SNPs which failed the set LD criteria were saved in a file, malay_data_qc_ibd.prune.out and further removed from cleaned SNPs data.....	95
Figure 3.14:	SNPs which were not in LD were extracted out using PLINK software and saved in PED and MAP file format for further analysis.....	95
Figure 3.15:	Final cleaned SNPs (35,457) were saved in a MAP file for further analysis. ....	96
Figure 3.16:	Genotypes of the final 478 individuals (in ACTG code) were saved in a PED file format for further analysis. ....	96
Figure 3.17:	The cleaned SNPs data ‘AllsampruneLD_malay.ped’ which was in ACTG code format is converted to ipPCA format using in-house python script and saved as AllsampruneLD.txt.....	99
Figure 3.18:	Five packages are needed to run the ipPCA in the RStudio which including the R stats package (stats), 3D scatter plot, misc functions (e1071), matrix and matrix exponential (expm).....	100
Figure 3.19:	Example of ipPCA_result.html output file. The final node (pink boxes) contained homogenous clusters where the samples in each cluster were identified by their sample id. ....	102

Figure 3.20:	Example of ipPCA_scatter_by_ipPCA.html output file. The final node (pink boxes) contained homogenous clusters.....	103
Figure 3.21:	Example of ipPCA_scatter_by_label.html output file. The final node (pink boxes) contained homogenous clusters where the colors of each group are related to predefined group.....	103
Figure 3.22:	Example of ipPCA_eigenvalue_plots.html output file showing the eigen value plots and the number of final samples in each cluster. ....	104
Figure 3.23:	The PED file of the 12 SPs in ACTG code saved as output.ped. ....	104
Figure 3.24:	Pairwise $F_{ST}$ calculated between the 12 SPs using the python script <b>7_fst.py</b> .....	107
Figure 3.25:	Example of pairwise $F_{ST}$ of all 35,457 SNPs calculated between SP1 and SP2. The $F_{ST}$ values are listed in column sixth. ....	107
Figure 3.26:	The list of SNPs which have high $F_{ST}$ value (Model 1 with 144 SNPs).....	108
Figure 3.27:	Example of the PED file ‘output1_Top3_weka.ped’ of Model 1 (SP2-YRI population) generated in ACTG code. ....	108
Figure 3.28:	The MAP file of Model 1 with 144 top SNPs.....	109
Figure 3.29:	The output1_Top3_weka_recode.raw in ‘012’ code. The missing genotype ‘NA’ need to be changed to ‘-1’ in order to be converted to .csv file and further used in WEKA	



	suite. This file contained FID (family ID), IID (individual ID), paternity and maternity information which is set as '0', sex information set as '3', phenotype set as '-9' and the genotype in '012' code. ....	109
Figure 3.30:	Example of the merged .csv files of all 12 sub-populations (Model 1). The ID of the sub-population is assigned under column 'Class' (example SP1-Indian).....	110
Figure 3.31:	Example of the input file for WEKA which contained the attributes and class information.....	116
Figure 3.32:	The Preprocess screen displaying the class label and number of samples in every class (up right) and the bar chart can be visualized (lower right) if the class is selected. ....	117
Figure 3.33:	The output results of the ancestry-predictive model performance evaluation which showing correctly and incorrectly classified instances and the respective statistical errors based on 144 AIM SNPs panel (Model 1). ....	118
Figure 3.34:	The area under ROC (AUC) plot of the SP3 (Malay-I) indicating good classification accuracy with AUC value of 0.9556.....	119
Figure 3.35:	ADMIXTURE analysis results representing the genetic structure pattern of the eight world populations.....	121
Figure 3.36:	Input SNPs data in 012 coded for PCA analysis and missing data 'N' were converted to numerical value by imputation using mean values. ....	123

Figure 3.37: The calculated eigenvalue and cumulative variance of the SNPs data.....	125
Figure 3.38: The scree plot showing the three first eigenvalues which larger than the others. ....	126
Figure 3.39: Loading values for calculation of score values for every SNPs. ....	127
Figure 3.40: The score values for every SNPs.....	128
Figure 3.41: Rosenberg’s information-theoretic informativeness for assignment ( $I_n$ ) scores was calculated for all 35,457 SNPs in pairwise manner (Malay versus other populations). ....	130
Figure 3.42: The final 250 SNPs which had the largest $I_n$ score were selected to developed the ancestry-predictive model for the Malay population.....	131
Figure 3.43: The test sample represented by the circle shape where its ancestry is determined by the ancestry of its nearest neighbours (the square shapes based on the majority voting) at K=5. ....	133
Figure 3.44: The SNPs genotype of the eight Malaysian Chinese.....	136
Figure 3.45: The SNPs genotype of the 17 Malaysian Indian. ....	136
Figure 3.46: The SNPs genotype of the 10 Malay samples.....	137
Figure 3.47: The SNPs genotype of the 18 Malay samples.....	137
Figure 3.48: The SNPs genotype of 24 Malay Kedah samples generated from OMNI 2.5 Illumina platform. ....	138
Figure 3.49: The SNPs genotype of 24 Malay Kelantan samples generated from OMNI 2.5 Illumina platform.....	138

Figure 3.50: The Affymetrix SNPs data of the 28 Malays, 17 Indians and 8 Chinese were merged into one PED file before filtration process. ....	141
Figure 3.51: Quality control performed on the Affymetrix SNPs data using PLINK software. ....	141
Figure 3.52: Quality control performed on the Illumina SNPs data using PLINK software. ....	142
Figure 3.53: The PED file of the filtered Affymetrix SNPs data comprising 53 genotypes. ....	142
Figure 3.54: The MAP file of the filtered Affymetrix SNPs data comprising 248,615 SNPs. ....	143
Figure 3.55: The PED file of the filtered Illumina SNPs data comprising 48 genotypes. ....	143
Figure 3.56: The MAP file of the filtered Illumina SNPs data comprising 249,243 SNPs. ....	144
Figure 3.57: The PED file of the merged SNPs data comprising 1766 genotypes. ....	144
Figure 3.58: The MAP file of the merged SNPs data comprising 37,487 SNPs. ....	145
Figure 3.59: The 1766 individuals were clustered into 11 final nodes (pink boxes) which correspond to 11 sub-populations by ipPCA analysis. All (160) Malay samples except five were clustered in Node 21. ....	148

Figure 3.60:	Results of the WEKA analysis of Model 4 (1250 AIM SNPs selected based on $F_{ST}$ ) showing classification accuracy of more than 90%.	149
Figure 4.1:	PCA plot showing four main clusters; Yoruba, Indian, Semang and mixture of Malays, Proto-Malay, Indonesian and Chinese Yunnan individuals.	153
Figure 4.2:	PCA plot showing the heterogeneous clusters are not well separated albeit slightly separation of the Proto-Malay individuals from the others (orange circles).	153
Figure 4.3:	The genetic structure pattern of the admixed Malay population revealed by the ADMIXTURE analysis (K=2 to K=5) based on the 35,457 cleaned SNPs.	156
Figure 4.4:	The genetic structure pattern of the admixed Malay population revealed by the ADMIXTURE analysis (K=5 to K=9) based on the 35,457 cleaned SNPs. It demonstrated genetic ancestry shared with the Indonesian (more than 50%), Chinese (~20%), Indian (~20%) and Proto-Malay (MY-TM) populations.	157
Figure 4.5:	Eigen value plots of sub populations SP1 to SP13 (correspond to their respective final nodes) and the eigendev values of each assigned sub populations.	163
Figure 4.6:	Scattered plots of sub populations SP1 to SP13 (correspond to their respective final nodes). Individuals in SP4, SP5, SP6, SP7, SP9 and SP13 are seen more	

	scattered than other sup-populations indicating the	
	heterogeneous genetic components. ....	166
Figure 4.7:	ipPCA based on 35,457 SNP data divided the 19	
	populations into 13 sub-populations where Malay	
	individuals were seen clustered into five different sub-	
	populations (SP3, SP4, SP8, SP12 and SP13) with most of	
	them grouped in SP3, SP4 and SP8.....	167
Figure 4.8:	ADMIXTURE analysis at K=2 to K=5 after ipPCA	
	pruning (387 individuals and 35,457 SNPs).....	169
Figure 4.9:	ADMIXTURE analysis at K=6 to K=9 after ipPCA	
	pruning (387 individuals and 35,457 SNPs).....	170
Figure 4.10:	Neighbour-joining (NJ) tree of the 12 sub-populations	
	based on pairwise $F_{ST}$ distances.....	172
Figure 4.11:	ADMIXTURE analysis using 144 AIM SNPs (Model 1) at	
	K=2 to K=5 failed to reveal genetic pattern differences	
	between the Malay population and the Indonesian, Proto-	
	Malay and Chinese populations.....	182
Figure 4.12:	ADMIXTURE analysis of 144 AIM SNPs (Model 1) at	
	K=6 to K=9 showed spurious results which failed to give	
	any conclusive results.....	183
Figure 4.13:	ADMIXTURE analysis using the 229 AIM SNPs ( Model	
	2) at K=2 to K=5 failed to reveal genetic pattern	
	differences between the Malay population and the	
	Indonesian, Proto-Malay and Chinese populations.....	184

Figure 4.14:	ADMIXTURE analysis of the 229 AIM SNPs (Model 2) revealed spurious results at K=6 to K=9. ....	185
Figure 4.15:	ADMIXTURE analysis results of Model 3 (433 AIM SNPs) at K=2 to K=5 revealed genetic structure pattern differences of the Malay population from the Yoruba, Indian, Semang, Proto-Malay and Chinese populations. However Malay and Indonesian populations cannot be differentiated at K=5.....	186
Figure 4.16:	ADMIXTURE analysis results of Model 3 (433 AIM SNPs) at K=6 to K=9. Using this AIM SNPs Yoruba, Indian, Jahai, Kensui, Proto-Malay, Chinese Wa and Chinese Jinuo genetic patterns could be differentiated to each other at K=8 albeit mask with spurious alleles. ....	187
Figure 4.17:	ADMIXTURE analysis results of Model 4 (1772 AIM SNPs) at K=2 to K=5. Five distinct populations were revealed at K=5; Yoruba, Indian, Semang, Proto-Malay, and Chinese. ....	188
Figure 4.18:	ADMIXTURE analysis results of Model 4 (1772 AIM SNPs) at K=6 to K=9. The results were comparable to the full 35,457 SNPs. At K=8 all 12 sub-populations demonstrated distinctive genetic pattern. ....	189
Figure 4.19:	ADMIXTURE analysis results of Model 5 (3145 AIM SNPs) at K=2 to K=5. Five distinct populations were revealed at K=5; Yoruba, Indian, Semang, Proto-Malay, and Chinese. Indonesian and Malay populations can also	

	be differentiated based on the differences amount of contribution of yellow colour block (Indian genetic component).....	190
Figure 4.20:	ADMIXTURE analysis results of Model 5 (3145 AIM SNPs) at K=6 to K=9. The results were comparable to the full 35,457 SNPs. At K=8 all 12 sub-populations demonstrated distinctive genetic pattern. ....	191
Figure 4.21:	The performance of AIM SNPs in assigning the Malay individuals into correct Malay cluster using $I_n$ ranking.....	194
Figure 4.22:	Comparison of PCA score using full set of SNPs (a) and 250 AIM SNPs (b). Both sets of SNPs enable the clustering of individuals according to their genetic ancestry. The 250 AIM SNPs was also enable the separation of individuals from all four closely related populations namely MY-TM, Indonesian, Malays and Chinese (Yunnan) (c). ....	195
Figure 4.23:	The genetic structure pattern of the admixed Malays population revealed by the ADMIXTURE analysis based on the 250 SNPs data at K=2 to K=5. ....	196
Figure 4.24:	The genetic structure pattern of the admixed Malays population revealed by the ADMIXTURE analysis based on the 250 SNPs data at K=6 to K=9. ....	197
Figure 4.25:	The genetic structure pattern of the admixed Malays population revealed by the ADMIXTURE analysis based on the 2000 SNP data at K=2 to K=5. ....	198

Figure 4.26:	The genetic structure pattern of the admixed Malays population revealed by the ADMIXTURE analysis based on the 2000 SNP data at K=6 to K=9. ....	199
Figure 4.27:	ADMIXTURE analysis of 250 AIM SNPs selected based on PCA-correlated SNPs at K=2 to K=5. ....	206
Figure 4.28:	ADMIXTURE analysis of 250 AIM SNPs selected based on PCA-correlated SNPs at K=6 to K=9. ....	207
Figure 4.29:	ADMIXTURE analysis of 2000 AIM SNPs selected based on PCA-correlated SNPs at K=2 to K=5. ....	208
Figure 4.30:	ADMIXTURE analysis of 2000 AIM SNPs selected based on PCA-correlated SNPs at K=5 to K=9. ....	209
Figure 4.31:	Comparison of number of overlapping AIM SNPs selected using both $I_n$ and PCAIMs methods. ....	215
Figure 4.32:	Comparison of the performance of AIM SNPs selected between $I_n$ and PCAIMs approaches. ....	215
Figure 4.33:	PCA analysis demonstrated five distinct clusters where populations demonstrated close genetic ancestry were grouped together in one cluster. ....	218
Figure 4.34:	Independent PCA analysis of Group 4 and 5 showing cluster of Malay population separated from Japanese and Chinese group (a). Independent PCA analysis of Group 2 and 3 also enable the separation Indian population from the European and Mexican populations (b). Individuals of Malaysian and Singaporean Malays were seen clustered	



	together albeit demonstrated several individuals' outliers (c) and (d). .....	219
Figure 4.35:	ADMIXTURE analysis of all 17 populations (Indian representing by MAL-IN and sgvp-ins; Malay representing by MAL-MY and sgvp-mas; Chinese representing by MY-CH, sgvp-chs, CHB, and CHD) using 37 487 SNPs at K=2 to K=10. Unique genetic patterns of all populations demonstrated at K=8, albeit at this number of K, TSI and CEU showed homogenous similar genetic patterns and GIH and Indian populations demonstrated slightly different genetic patterns from each other.....	223
Figure 4.36:	Eigen value plots of sub populations SP1 to SP11 (correspond to their respective final nodes) and the eigendev values of each assigned sub populations.....	227
Figure 4.37:	Scattered plots of sub populations SP1 to SP11 (correspond to their respective final nodes). Individuals in SP2, SP4 and SP10 are seen more scattered than other sup-populations indicating the heterogeneous genetic components.....	230
Figure 4.38:	ipPCA analysis re-clustered individuals into 11 sub- populations where Malay individuals were all clustered in SP6 except five Singaporean Malay individuals which were clustered with Chinese ancestor in SP5.....	231
Figure 4.39:	ADMIXTURE analysis of the all 11 SPs at K=2 to K=9 based on 37 487 SNPs. Genetic pattern of Malay	

	population can be differentiated from all other world population at K=7. Contribution of Indian (purple colour block) and Chinese (green colour block) ancestor genetic components can be seen in the Malays.....	235
Figure 4.40:	ADMIXTURE analysis of Model 1 (101 AIM SNPs) at K=2 to K=6. Malay population hardly differentiated from the Chinese and Japanese populations.....	246
Figure 4.41:	ADMIXTURE analysis of Model 2 (157 AIM SNPs) at K=2 to K=8. Malay population demonstrated admixed genetic pattern which can be differentiated from other world populations at K=7.....	249
Figure 4.42:	ADMIXTURE analysis of Model 3 (294 AIM SNPs) at K=2 to K=8.....	252
Figure 4.43:	ADMIXTURE analysis of Model 4 (1250 AIM SNPs) at K=2 to K=9. Malay population demonstrated unique genetic pattern at K=7. All other world populations were also demonstrated distinctive genetic pattern at K=8.....	255
Figure 4.44:	ADMIXTURE analysis of Model 5 (2240 AIM SNPs) at K=2 to K=9. The results are comparable with the full number of SNPs.....	258
Figure 4.45:	The performance of small number of SNPs selected based on $I_n$ .....	261
Figure 4.46:	ADMIXTURE analysis of 100 AIM SNPs selected based on $I_n$ algorithm at K=2 to K=5. Higher K gave spurious results.....	262

Figure 4.47:	ADMIXTURE analysis of 200 AIM SNPs selected based on $I_n$ algorithm at K=2 to K=7. Higher K gave spurious results.....	264
Figure 4.48:	ADMIXTURE analysis of 2000 AIM SNPs selected based on $I_n$ algorithm at K=2 to K=9. Malay genetic structure can be differentiated from other world population including the Chinese and Japanese. ....	267
Figure 4.49:	Comparison of the performance of the PCAIMs methods at $k=2$ and $k=3$ with the $I_n$ method. In this study PCAIMs method demonstrated slightly better performance compared to $I_n$ method.....	273
Figure 4.50:	Comparison of shared or overlap SNPs selected using PCAIMs methods at $k=2$ and $k=3$ with the $I_n$ method. Percentage of overlap SNPs selected between the PCAIMs and $I_n$ method are quite low.....	273
Figure 4.51:	ADMIXTURE analysis of 100 AIM SNPs selected based on PCAIMs algorithm at K=2 to K=6. Higher K gave spurious results. Genetic pattern of Malay population slightly different from the other world population, however Chinese and Japanese population cannot be differentiated using this 100 AIM SNPs. ....	275
Figure 4.52:	ADMIXTURE analysis of 200 AIM SNPs selected based on PCAIMs algorithm at K=2 to K=7. Spurious results were seen at K=6 and K=7. Genetic pattern of Malay population slightly different from the other world	

	population, however Chinese and Japanese population cannot be differentiated using this 200 AIM SNPs. ....	277
Figure 4.53:	ADMIXTURE analysis of 2000 AIM SNPs selected based on PCAIMs algorithm at K=2 to K=9. The results were comparable with the full SNPs (37,487). ....	279
Figure 4.54:	ADMIXTURE analysis of the 555 SNPs showing the genetic pattern of the Malay population compared to the other 11 world populations at K=2 to K=8. Malay genetic structure distinctly differentiated from the Chinese and Japanese using this 555 set of AIM SNPs at K=7. ....	288
Figure 4.55:	Results of the kappa statistic consistently increases with the increases of correctly classified percentage. Kappa statistic of Model 4 is 0.88 which indicates nearly complete agreement with the prediction of true class. ....	292
Figure 4.56:	The performance of AIM SNPs selected by $I_n$ method better than PCAIMs (478 instances) with high kappa statistic. ....	292
Figure 4.57:	The performance of classification of all ancestry-predictive model for the Malay-1 sub-population. ....	293
Figure 4.58:	The performance of classification of all ancestry-predictive models is quite low for the Malay-2 sub-population which could be contributed by the small number of samples studied. ....	293

Figure 4.59:	Comparison of the classification performance of ancestry-predictive models developed based on $F_{ST}$ , $I_n$ and PCAIMs methods.....	295
Figure 4.60:	Comparison of classification performance of all ancestry-predictive models which revealed the poor performance of the small number of AIM SNPs where the TP and $F$ -measure were seen dropped below 70%. .....	295
Figure 5.1:	Comparison of 250 AIM SNPs selected using $I_n$ approach with the Full number of SNPs and 2000 AIM SNPs. It is demonstrated that 250 AIM SNPs can be used to differentiate Malay population with the other six populations at $K=5$ . Increasing the number of SNPs to 2000 enable a clearer plot of admixed genetic structure of the Malays and the distinct genetic pattern of the Jahai (brown colour) and Kensiu (blue colour) of sub-ethnic Semang, the Proto-Malay (orange colour) as well as the Chinese Jinuo (pink colour) and Wa (purple colour) of sub-ethnic Yunnan.....	312
Figure 5.2:	Comparison of ADMIXTURE results analysis of the 250 AIM SNPs selected using $I_n$ approach with 2000 AIM SNPs and full SNPs set for the populations re-assigned by ipPCA algorithm (SP1 to SP12).....	313
Figure 5.3:	Classification performance of 250 AIM SNPs is comparable to the 2000 AIM SNPs with the TP value of all above 0.8, except for sub-populations Malay-1 and	

	Malay-2. This is due to the small number of individual representing the sub-populations which contributed to the lower TP value.....	314
Figure 5.4:	Overall classification performance of 2000 AIM SNPs is better than 250 AIM SNPs. Assigning of individuals to their sub-populations (SP1 to SP12 respectively) reducing the number of individuals representing each population hence classification performance of all sub-population dramatically reduced. ....	314
Figure 5.5:	Comparison of AIM SNPs selected using $I_n$ approach with the full SNPs set. It was noted that Malay population can only be differentiated from the other Asian populations using 2000 AIM SNPs.....	316
Figure 5.6:	Classification performance of 2000 AIM SNPs selected using $I_n$ method is remarkable for the Malay population (SP6) and for almost all the others population except for SP5, SP10 and SP11.....	317
Figure 5.7:	Comparison of PCAIMs AIM SNPs panel for the 50K SNPs data. The 2000 AIM SNPs panel demonstrated comparable ADMIXTURE results as the full SNPs at K=8. The 250 AIM SNPs unable to differentiate Malays from the Proto-Malay, Indonesian and Chinese (Yunnan) populations at K=5 (higher K gave spurious results). The 2000 AIM SNPs panel clearly revealed the genetic characteristic of the Malays (shared major genetic	

characteristic with Indonesian (green colour block) with the combination of traces of Indian (yellow colour block), Chinese (purple/pink colour block) and the Proto-Malay (orange colour block) genetic component)..... 320

Figure 5.8: Comparison of PCAIMs AIM SNPs panel for the 50K SNPs data. Genetic characteristic of the 12 sub-populations which assigned by the ipPCA analysis comparable to the full SNPs set when analyse using the 2000 AIM SNPs panel (at K=8). However both Malay-1 and Malay-2 sub-populations did not form distinct genetic characteristic when ADMIXTURE analysis was run using 250 AIM SNPs panel at K=5 (higher K gave spurious results). The 250 AIM SNPs panel enable the differentiation of the Malays from the Yoruba, Indian, Kensiou and Jahai populations. .... 321

Figure 5.9: The 2000 AIM SNPs demonstrated a better classification performance than 250 AIM SNPs for the Malay-2, Indonesian and Proto-Malay population but comparable or slightly lower for the pooled Malays, Malay-1, Chinese, Indian, Yoruba and Semang populations..... 322

Figure 5.10: Overall performance of 2000 AIM SNPs panel is better than 250 AIM SNPs for all sub-populations assigned by the ipPCA. However 250 AIM SNPs panel successfully classified the Malay-1 individuals at almost 0.8 TP value..... 322

Figure 5.11: Comparison of ADMIXTURE results of AIM SNPs panel with full SNPs set. The small number of SNPs (100 and 200 AIM SNPs) selected using <i>PCAIMs</i> enables the differentiation of Malays genetic structure from the other 11 populations. A distinct Malay genetic structure can be seen at K=8 when 2000 AIM SNPs were used to run the ADMIXTURE analysis. ....	325
Figure 5.12: Classification performance of three AIM SNPs panel selected using <i>PCAIMs</i> . Overall performance of the 2000 AIM SNPs panel is better than the 100 and 200 AIM SNPs panel. All AIM SNPs panel enable the classification of individuals to their correct sub-populations with TP value more than 0.4 except for sub-populations SP5, SP10 and SP11.....	326
Figure 5.13: Comparison of ADMIXTURE plots at K=6 developed from panel of 200 AIM SNPs selected by the $I_n$ and <i>PCAIMs</i> method. The Malay population hardly differentiated from the others Asian population. Heavy background noise was also seen in both $I_n$ and <i>PCAIMs</i> plots. ....	329
Figure 5.14: ADMIXTURE plot developed from the 157 AIM SNPs selected by <i>ipPCA-F<sub>ST</sub></i> method. Genetic structure of the Malay population could be differentiated from the others population including the Asian population albeit with	



	heavy background noise interrupted their genetic characteristic plot as compared to the full SNPs set.....	329
Figure 5.15:	Comparison of ADMIXTURE plots at K=8 developed from panel of 2000 AIM SNPs selected by the $I_n$ and $PCAIMs$ method. The Malay population could be differentiated from the others population albeit slight background noise was observed in both $I_n$ and $PCAIMs$ plots. ....	330
Figure 5.16:	ADMIXTURE plot developed from the 1250 AIM SNPs selected by $ipPCA-F_{ST}$ method. Malay population formed distinct genetic structure which could be differentiated from the others population. Genetic structure of others population consistently re-produced using the 1250 AIM SNPs except for SP10.....	330
Figure 5.17:	Comparison of genetic structure revealed by a minimal number of AIM SNPs panel (at K=5) selected by $F_{ST}$ , $I_n$ and $PCAIMs$ method. Overall performance of the AIM SNPs panel is comparable albeit more noise can be seen in $I_n$ and $PCAIMs$ plots.....	333
Figure 5.18:	Comparison of genetic structure revealed by 1772 and 2000 AIM SNPs panel (at K=8) selected by $F_{ST}$ , $I_n$ and $PCAIMs$ method. The performance of all three methods is comparable to each other. The admixed genetic structure of the Malay population is successfully revealed using these numbers of SNPs.....	334

Figure 5.19: Schematic diagram of Malay ancestry customized chip  
(<http://www.affymetrix.com>). Probes complementary to  
the 1772 AIM SNPs loci and the 1250 AIM SNPs loci  
could be incorporated onto a chip and further be validated  
on the Malay population. .... 355

## LIST OF SYMBOLS AND ABBREVIATIONS

%	: Percentage
+	: Plus
±	: Plus minus
®	: Registered Sign
Σ	: Sum of
≤	: Less than and equal to
≥	: More than and equal to
>	: More than
<	: Less than
~	: Approximately
β	: Beta
α	: Alfa
δ	: Absolute allele frequency differences
λ <sup>2</sup>	: Chi-squared test
r <sup>2</sup>	: Regression value
K	: Number of cluster
P	: Allele frequencies
Q	: ancestral component or membership coefficient
Pr(C <sub>k</sub> /d)	: Posterior probabilities of data
A	: Adenine
aAIMs	: Ancient ancestry informative markers
ABI	: Applied Biosystems
AIM	: Ancestry informative marker

ANOVA	: Analysis of variance
ASW	: Africa America
AUC	: Area under curve
bp	: Basepair
C	: Cytosine
CEU	: European
CHB	: Han Chinese Beijing
CHD	: Chinese Colorado
CN-JN	: Chinese Jinuo
CN-WA	: Chinese Wa
dbSNP	: Database Single Nucleotide Polymorphism
ddNTPs	: Dideoxynucleotides triphosphates
DNA	: Deoxyribonucleic Acid
EM	: Expectation maximization
$F_{ST}$	: Allele frequency differences between populations
FN	: False negative
FP	: False positive
FRET	: Fluorescence resonance energy transfer
G	: Guanine
GIH	: Gujarati Indian
GWAS	: Genome-wide association studies
HGDP-CEPH:	Human Genome Diversity Cell Line Panel
HLA	: Human Leukocyte Antigen
HMM	: Hidden markov model
HNB	: Hidden Naïve Bayes

<i>H.Pylori</i>	: Helicobacter pylori
HVR	: Hyper variable region
HWE	: Hardy Weinberg Equilibrium
IBS	: Identical-by state
IBD	: Identical-by descent
$I_n$	: Informativeness for assignment
ID-JV	: Indonesian Java
ID-ML	: Indonesian Malay
ID-TR	: Indonesian Toraja
IN-DR	: Indian Telugu
IN-WL	: Indian Marathi
INDEL	: Insertion deletion
ipPCA	: Iterative pruning principal component analysis
JPT	: Japanese
KNN	: K-nearest neighbor
LAMP	: Local ancestry in admixed population
LD	: Linkage disequilibrium
LSBL	: Locus specific branch length
LWK	: Luhya
MAF	: Minor allele frequency
MAL-CH	: Chinese Malaysia
MAL-IN	: Indian Malaysia
MAL-MY	: Melayu Malaysia
MC1R	: Melanocortin 1 reseptor
MCMC	: Markov Chain Monte Carlo

MDS	: Multidimensional scale
MEX	: Mexican
MHC	: Major histocompatibility complex
MHMM	: Markov-hidden markov model
MKK	: Maasai
mtDNA	: Mitochondrial DNA
MY-BG	: Melayu Bugis
MY-BJ	: Melayu Banjar
MY-CH	: Melayu Champa
MY-JH	: Jahai
MY-JV	: Melayu Java
MY-KD	: Melayu Kedah
MY-KN	: Melayu Kelantan
MY-KS	: Kensui
MY-MN	: Melayu Minang
MY-PT	: Melayu Patani
MY-TM	: Temuan
MyHVP	: Malaysian node of the human variome project
NB	: Naïve Bayes
NCBI	: National Center for Biotechnology Information
PCA	: Principal Component Analysis
PCAIMs	: PCA-correlated ancestry informative markers
RFLP	: Restriction fragment length polymorphisms
RMSE	: Root mean squared error
ROC	: Receiver operating characteristic

SBE	: Single-base extension
SGVP	: Singapore Genome Variation Project
sgvp-chs	: Singapore Chinese
sgvp-ins	: Singapore Indian
sgvp-mas	: Singapore Melayu
SINEs	: Short interspersed nuclear elements
SP	: Sub-population
STRs	: Short Tandem Repeats
SNPs	: Single Nucleotide Polymorphisms
SVD	: Singular value decomposition
T	: Thymine
TB	: Tuberculosis
TP	: True positive rate
TSI	: Toscani
WEKA	: Waikato Environment for Knowledge Analysis
YRI	: Yoruba

## LIST OF APPENDICES

APPENDIX A	Ethical approval
APPENDIX B	Steps to run python script in command prompt
APPENDIX C	Python script to covert SNPs data to MAP and PED file format
APPENDIX D	Running quality control (QC) filtering using PLINK
APPENDIX E	Steps to convert SNPs data to ipPCA format
APPENDIX F	Steps to calculate the $F_{ST}$ using python script
APPENDIX G	Steps to run WEKA Suite the ancestry-predictive modeling
APPENDIX H	Step to convert PED and MAP data to ADMIXTURE format
APPENDIX I	List of Malay AIM SNPs panel
APPENDIX J	Reviewer comments



**PENANDA MAKLUMAT Keturunan POLIMORFISME NUKLEOTIDA  
TUNGGAL UNTUK ANGGARAN Keturunan DALAM POPULASI  
MELAYU**

**ABSTRAK**

Penanda maklumat keturunan (AIM) dapat digunakan untuk menyimpulkan (inferens) keturunan seseorang individu bagi meminimumkan ketidaktepatan maklumat keturunan yang dilaporkan sendiri yang digunakan pada masa ini dalam penyelidikan bioperubatan. Dalam kajian ini, tiga kaedah digunakan dalam membina panel SNP AIM Melayu, iaitu analisis *iterative pruning principal component* yang digabungkan dengan *pairwise  $F_{ST}$  (ipPCA- $F_{ST}$ )*, *informativeness for assignment ( $I_n$ )* dan *PCA-correlated SNPs (PCAIMs)*. Dua set data genotip SNP Melayu yang diperolehi daripada Projek Variome Manusia Malaysia (MyHVP) telah digunakan untuk mengekstrak panel SNP AIM Melayu iaitu 135 genotip SNP Melayu yang dihasilkan melalui platform Affymetrix GeneChip Mapping Xba 50K array dan 76 genotip SNP Melayu yang dihasilkan melalui platform Affymetrix SNP-6 array dan SNP OMNI2.5 Illumina array. Tambahan 89 lagi genotip SNP Melayu yang dihasilkan melalui platform Affymetrix SNP-6 array diperolehi dari Projek Variasi Genom Singapura (SGVP). Pangkalan data SNP Pan-Asian digunakan sebagai populasi rujukan untuk pangkalan data genotip SNP Melayu pertama manakala pangkalan data HapMap Fasa 3 dan pangkalan data SGVP digunakan sebagai populasi rujukan untuk pangkalan data genotip SNP Melayu kedua. Ketepatan setiap panel SNP AIM Melayu yang dihasilkan dinilai dengan menggunakan *machine learning "ancestry-predictive model"* yang dibina dengan menggunakan WEKA, platform *machine learning* komprehensif yang ditulis dalam Java. Seterusnya, corak genetik bangsa Melayu

dianalisis menggunakan program ADMIXTURE berdasarkan kepada panel SNP AIM tersebut. Hasil analisis menunjukkan bahawa model SNP AIM 144, 299 dan 433 yang dipilih dari data Affymetrix 50K SNP menggunakan kaedah *ipPCA-F<sub>ST</sub>*, mengklasifikasikan individu Melayu masing-masing, dengan ketepatan 76.5%, 70.6% dan 82.4%. Ketepatan meningkat kepada 88.2% menggunakan model SNP AIM 1772 dan 3145. Model SNP AIM 250 dan 2000 yang dipilih menggunakan kaedah  $I_n$ , berjaya mengklasifikasikan individu Melayu dengan ketepatan masing-masing, 89.8% dan 96.1%. Walaubagaimanapun, ketepatan sedikit lebih rendah, masing-masing 78% dan 92.1% untuk bilangan SNP yang sama yang dipilih dengan kaedah *PCAIM*. Analisis ADMIXTURE menunjukkan corak genetik populasi Melayu dapat dibezakan dengan jelas dari populasi rujukan menggunakan panel SNP AIM 1772 dan 3145 yang dipilih oleh *ipPCA-F<sub>ST</sub>* dan 2000 SNP yang dipilih oleh  $I_n$  dan *PCAIMs*. Panel SNP AIM 101, 157 dan 294 yang dipilih dari data Affymetrix SNP-6 menggunakan kaedah *ipPCA-F<sub>ST</sub>*, menunjukkan ketepatan klasifikasi masing-masing 88.8%, 94.4% dan 96.9%. Ketepatan meningkat kepada 100%. menggunakan panel SNP AIM 1250 dan 2240. Model 100, 200 dan 2000 SNPs dipilih menggunakan  $I_n$ , menunjukkan ketepatan klasifikasi masing-masing 67.5%, 80% dan 100%. *PCAIM* menunjukkan ketepatan 68.8%, 81.9% dan 99.4%, masing-masing untuk bilangan panel SNP AIM yang sama. Corak genetik populasi Melayu dapat dibezakan dengan populasi dunia lain yang digunakan dalam kajian ini menggunakan panel SNP AIM 1250 dan 2240 yang dipilih menggunakan kaedah *ipPCA-F<sub>ST</sub>* dan 2000 AIM SNP yang dipilih menggunakan kaedah  $I_n$  dan *PCAIM*.

**ANCESTRY INFORMATIVE MARKERS SINGLE NUCLEOTIDE  
POLYMORPHISMS PANEL FOR ANCESTRY ESTIMATION IN THE  
MALAY POPULATION**

**ABSTRACT**

Ancestry-informative markers (AIMs) can be used to infer an individual's ancestry to minimize the inaccuracy of self-reported ethnicity in biomedical research. The AIM-SNP panels for the Malay population were developed using three methods in this study: iterative pruning principal component analysis (ipPCA) combined with pairwise  $F_{ST}$  (ipPCA- $F_{ST}$ ), informativeness for assignment ( $I_n$ ), and PCA-correlated SNPs (or PCA-informative markers; PCAIMs). The Malay AIM-SNP panels were designed using two sets of Malay SNP genotype datasets stored in SNP arrays hosted by the Malaysian Node of the Human Variome Project (MyHVP). The first dataset contained 135 Malay SNPs genotypes generated from Affymetrix GeneChip Mapping Xba 50K array platform. The second dataset contained 76 Malay SNPs genotypes generated from Affymetrix SNP-6 array and OMNI2.5 Illumina SNP array platforms. In addition, 89 Malay SNP genotypes from the Singapore Genome Variation Project (SGVP) using the Affymetrix SNP-6 array platform were also included in the second set of the Malay SNP genotype datasets. The Pan-Asian SNP dataset was used as a reference population for the first Malay SNP genotype dataset, whereas the International HapMap Phase 3 project and SGVP datasets served as the reference populations for the second Malay SNP genotype datasets. The accuracy of each resulting Malay AIM-SNP panel was evaluated using machine learning “ancestry-predictive model” constructed using WEKA, a comprehensive machine learning

platform written in Java. The ADMIXTURE program was used to explore the genetic pattern of Malays based on the selected AIM-SNP panels. The results showed that models with 144, 299, and 433 AIM-SNP panels selected from the Affymetrix 50K SNPs dataset using the ipPCA-F<sub>ST</sub> method, correctly classified Malay individuals with an accuracy of 76.5, 70.6, and 82.4%, respectively. The accuracy further increased to 88.2% when using models with 1772 and 3145 AIM-SNP panels. Models with 250 and 2000 SNPs ranked by  $I_n$ , correctly classified Malay individuals with an accuracy of 89.8 and 96.1%, respectively. However, the accuracy was slightly lower, 78 and 92.1%, respectively, for the same number of SNPs selected by the PCAIM method. ADMIXTURE analysis showed that the genetic structure of the Malay population can be distinctly differentiated from the reference populations using 1772 and 3145 AIM-SNP panels selected by ipPCA-F<sub>ST</sub>, and 2000 SNPs selected by  $I_n$  and PCAIM. Models with 101, 157, and 294 AIM-SNP panels selected from the Affymetrix SNP-6 dataset using the ipPCA-F<sub>ST</sub> method demonstrated classification accuracy of 88.8, 94.4, and 96.9%, respectively. Remarkable results were obtained using 1250 and 2240 AIM-SNP panels, where the accuracy increased to 100%. Models with 100, 200, and 2000 ranked by  $I_n$ , correctly classified Malay individuals with an accuracy of 67.5, 80, and 100%, respectively. For the same number of AIM-SNP panels, the PCAIM showed an accuracy of 68.8, 81.9, and 99.4%, respectively. The genetic structure of the Malay population can be differentiated from the other world populations used in this study using the 1250 and 2240 AIM-SNP panels selected by ipPCA-F<sub>ST</sub> and 2000 AIM-SNP panels selected by  $I_n$  and PCAIM.

# CHAPTER 1

## INTRODUCTION

### 1.1 Research background

More than 100,000 years ago, modern humans were thought to have migrated from Africa to other parts of the world (Cavalli-Sforza, 2007). They formed local communities influenced by their surrounding environment and tended to mate in proximity, contributing to genetic drift and natural selection (Cavalli-Sforza, 2007). Mutations have also arisen, resulting in biological differences, though they retain some of their ancestor's genomic information. The migration, genetic drift, mutation and natural selection operated in parallel with demographic and historical events resulted in variants that are rare in some population but not in others which are likely arisen recently and contributed to the population differences (Cavalli-Sforza, 2007, Paschou et al., 2010). The differences of the genetic patterns were portrayed in their genetic ancestry and population structure carried in the genome of each individual (Cavalli-Sforza, 2007, Paschou et al., 2010).

Microsatellite markers were used to examine patterns of human genetic variation and population genetic structure across the entire genome (Paschou et al., 2010). Studies of population genetic structure based on Single Nucleotide Polymorphisms (SNPs) genome-wide data have successfully revealed the clines of genetic diversity around the world, especially with the advent of modern technologies and the realization of the HapMap project (Paschou et al., 2010). More recent studies of genetic ancestry to infer individual membership down to a population within a continent has attracted

considerable attention because of their value in biomedical, population genetics, anthropological, and forensic applications (Bryc et al., 2015, Byun et al., 2017, Das and Upadhyai, 2018, Vongpaisarnsin et al., 2017, Zeng et al., 2016).

To reveal the genetic variation within and among populations, the genetic distance between them can be measured by calculating the frequencies of the variant allele. Wright's F-statistics ( $F_{ST}$ ) is commonly used to measure genetic differentiation. Small  $F_{ST}$  revealed similar allele frequencies, whereas large  $F_{ST}$  indicated that allele frequencies within each population differed (Holsinger and Weir, 2009). This variant/allele can be chosen as a candidate marker to infer ancestry for an individual of that population, and it is called ancestry-informative markers (AIMs).

AIMs are DNA markers that show high allele frequency differences between populations from different geographic regions; thus, this marker could infer an individual's biogeographic ancestry (Kayser and de Knijff, 2011). AIMs can be found on any DNA polymorphisms such as short tandem repeats (STRs), Alu elements, insertion-deletion polymorphisms (INDELs) or single nucleotide polymorphisms (SNPs) (Algee-Hewitt et al., 2016, Esposito et al., 2018, Gómez-Pérez et al., 2010, Inácio et al., 2016). Both haploid and diploid genetic markers can be used to study genetic ancestry or biogeographic ancestry. Mitochondrial DNA (mtDNA) and Y-chromosome polymorphisms are the two haploid markers that have been frequently used to study biogeography ancestry (Chaitanya et al., 2014, Dulik et al., 2012, Salas et al., 2006). However, mtDNA (maternal lineage) and Y-chromosome (paternal lineage) markers do not provide comprehensive information on individual ancestry due to their haplotype nature. Comparatively, autosomal markers can provide more

information about an individual's genetic ancestry because they represent a much greater proportion of genome history (both maternal and paternal ancestry) (both maternal and paternal ancestry) (Royal et al., 2010).

Autosomal markers commonly studied for genetic ancestry were STRs (Kutanan et al., 2014, Nunez et al., 2010, Phillips et al., 2013a), Alu elements (Gómez-Pérez et al., 2010, Hormozdiari et al., 2011, Krishnaveni and Prabhakaran, 2015), INDELs (Moriot et al., 2018, Pereira et al., 2012, Tao et al., 2019, Zaumsegel et al., 2013) and SNPs (Fondevila et al., 2013, Hwa et al., 2017, Kidd et al., 2014, Poetsch et al., 2013). Nonetheless, SNPs was the marker of choice for the ancestry studies due to its stability, abundance in human genome, demonstrated pronounced frequency variation among populations and thousands of SNPs can be assayed simultaneously using high throughput platform such as microarray chips (Royal et al., 2010). Furthermore, autosomal SNPs are almost entirely used to estimate genetic ancestry in epidemiological applications.

Single base pair substitution or SNPs play an important role to the variation among individuals including susceptibility to disease and reactions to drug (Kruglyak and Nickerson, 2001). SNPs play an important role in an individual susceptibility to most diseases and drugs metabolism and it is also said to be directly involved in determining the phenotype of an individual such as eye color, hair and skin; facial morphology and height (Butler, 2012). Millions of SNPs have been identified and made available at dbSNP homepage at NCBI and HapMap, however a small subset of SNPs (10-100s) should be enough to accurately infer an individual ancestry (Fondevila et al., 2013, Gettings et al., 2014, Sampson et al., 2011). AIM-SNPs are a small set of informative

SNPs. The advantages of AIM-SNPs over a random set of autosomal SNP markers include its ability to offer increased power for ancestry inference and to define admixed population (to determine the relative percentages of descendants). Simultaneously, a smaller set of markers can reduce genotyping costs while increasing throughput (Royal et al., 2010).

Several approaches have been used by researchers to select SNPs for ancestry studies. Rosenberg et al., (2003) suggested that SNPs can be ranked individually based on their ability to distinguish ancestry by calculating the  $F_{ST}$  value, allele frequency and the informativeness for assignment ( $I_n$ ).  $I_n$  is the measure of information provided by multiallelic markers about individual ancestry (Rosenberg et al., 2003). Paschou et al., (2010) chose SNPs that are strong contributors to the principal component analysis (PCA) or PCA-correlated SNPs (PCAIMs). PCA is a multivariate analysis that provides a new coordinate system (Kayser and de Knijff, 2011). Kersbergen et al., (2009) used a model-based clustering approach STRUCTURE program to estimate genetic diversity among multiple groups of individuals and then used the pairwise  $F_{ST}$  ranking procedure to identify AIM-SNPs.

Iterative pruning PCA (ipPCA) is another algorithm suggested for searching AIM-SNPs. This algorithm assigned individuals to sub-populations and calculated the total number of sub-populations present, and the STRUCTURE program was then used to select the appropriate AIM-SNPs (Intarapanich et al., 2009). Galanter et al., (2012) used Locus Specific Branch Length (LSBL) approach to discover AIM-SNPs. Based on  $F_{ST}$  values, LSBL is a measure of population structure in one population sample relative to two other population samples. The AIM-SNPs were selected based on the



highest LSBL for that population. The selected AIM-SNPs were tested for linkage disequilibrium (LD), physical distance, and heterogeneity. The AIMS were excluded if the markers were in LD or the alleles showed significant allele frequency heterogeneity between the samples representing each ancestral group (Galanter et al., 2012).

ADMIXTURE program is another approach to observe the genetic structure of studied groups/clusters while simultaneously selecting AIM-SNPs (Vongpaisarnsin et al., 2015). Other approaches include selecting SNPs from databases that exhibit marked allele frequency differences between populations (high  $F_{ST}$  value) and strongest contributors to the PCA (Harrison et al., 2008, Phillips et al., 2013b, Sampson et al., 2011); exploring existing SNP panels with hundreds of genetic markers using commercially available SNP genotyping arrays (Kidd et al., 2014) and selecting SNPs with genes involved in melanin syntheses such as MC1R, OCA2, ASIP, or SLC45A2 (Gettings et al., 2014, Poetsch et al., 2013, Soejima and Koda, 2007).

Individual AIM-SNPs can then be genotyped using common SNP genotyping platforms such as allelic-specific hybridization, primer extension, oligonucleotide ligation, or invasive cleavage (Sobrinho et al., 2005). Further, allelic products of these methods can be detected with several detection systems such as fluorescence-electrophoresis (Bouakaze et al., 2009, Fondevila et al., 2013, Mosquera-Miguel et al., 2009, Poetsch et al., 2013), fluorescence resonance energy transfer (FRET) (Lareu et al., 2001, Nicklas and Buel, 2008), fluorescence arrays (Divne and Allen, 2005, Zeng et al., 2012), and mass spectrometry (Li et al., 1999, Shi et al., 2011). Primer extension combined with the fluorescence-electrophoresis allelic detection method, also known

as mini-sequencing technology, is one of the most suitable methods for analyzing a small number of AIM-SNPs. This is because most laboratories have an automatic capillary electrophoresis instrument, which is also used for STR genotyping. To facilitate the detection of the SNP allelic products, the commercialized multiplex single-base extension reaction SNaPshot® kit (Applied Biosystems, USA), which utilizes fluorescent ddNTPs, is also available on the market (Daniel et al., 2009, Phillips et al., 2013a, Rogalla et al., 2015, Wei et al., 2016).

High throughput technology such as microarray is more powerful for analyzing hundreds or thousands of AIM-SNPs simultaneously. Keating et al., (2013) developed the Identitas v1 Forensic Chip, a diagnostic tool comprising of 201,173 genome-wide autosomal, X-chromosomal, Y-chromosomal, and mitochondrial SNPs for simultaneously inferring biogeographic ancestry, appearance, relatedness, and gender. The chip, which was manufactured by Illumina, uses the well-established Infinium technology (Keating et al., 2013). Galanter et al., (2012) genotyped 446 AIM-SNPs using both Affymetrix and Illumina platforms, and Krjutskov et al., (2009) used a microarray platform to analyze a 124-plex SNP comprising of 49 mtDNA SNPs, 29 Y-chromosomal SNPs, and 46 autosomal SNPs (Krjutskov et al., 2009).

According to previous studies, certain diseases are more prevalent in one ethnic group than others, such as hypertension, end-stage renal disease, tuberculosis, lung function, and prostate cancer (Daya et al., 2014, Menezes et al., 2015, Royal et al., 2010). Cappetta et al., (2015) studied the effect of genetic ancestry on leukocyte global DNA methylation in cancer patients and suggested that genetic ancestry should be considered as a modifying factor in epigenetic association studies, especially in

admixed populations. Thus, understanding ethnicity/ancestry and the substructure is vital for properly designing case-control association studies and identifying disease predisposing alleles that may differ across ethnic groups (Cappetta et al., 2015, Cavalli-Sforza, 2007, Royal et al., 2010, Tishkoff and Kidd, 2004). The current practice of using the self-identified ethnicity/ancestry approach in association diseases or medical genetics studies may result in false-positive or false-negative results, especially in studies of admixed populations, because this approach cannot account for the percentage of admixture in admixed cases (Liu et al., 2013b). Furthermore, understanding one's ancestry background can help in the proper diagnosis and subsequent treatment of diseases.

## **1.2 Problem statement**

Self-reported ethnicity has limitations when conducting genetic studies (Al-Alem et al., 2014, Al-Naamani et al., 2017, Mersha and Abebe, 2015). A study participant's ethnicity is frequently misreported. This contributes to the spurious association with false-positive or false-negative results. Self-reported ethnicity errors may occur when subjects are unaware of their true ethnicity, only know their recent ancestors, or identify with one ethnic group despite their admixed background. Therefore, the accuracy of biomedical studies will be affected, and the study will fail to provide novel insights into variation in disease susceptibility and adverse drug reactions in the studied population or individual. As a result, analyzing the human genome can provide information about an individual's ancestry, especially AIM-SNPs, which can be used to describe actual genetic variation in studied individuals or populations.

Many studies on AIM-SNPs have been carried out in various populations worldwide, including the Brazilian population (Lins et al., 2010), the United States population (African American, East Asian, European American, and Hispanic American/Native American) (Getting et al., 2014), European-descendent populations (Huckins et al., 2014), the Han Chinese population (Qin et al., 2014), the Thai population (Vongpaisarnsin et al., 2015), Australia and Pacific region populations (Santos et al., 2016) and Singapore population (Ramani et al., 2017). To date, however, no study has been carried out to identify Malay AIM-SNPs for the Malay population. The designing of Malay AIM-SNP panels can be used to ascertain Malay ancestry and reveal the extent of admixture in Malay individual subjects. This information may improve the accuracy and precision of biomedical and pharmacological studies carried out in the Malay population.

### **1.3 Research justifications**

Genetic epidemiology shows that many Mendelian diseases are concentrated in a few, usually small social or ethnic groups, especially for the rarer diseases (Cavalli-Sforza, 2007). Ethnicity/ancestry of individuals has also been proven to affect their response to some specific therapeutic agents or drugs (Al-Naamani et al., 2017). The study of biogeographical ancestry using AIM-SNPs is an important tool in the human medical research for gaining a better understanding of the ancestry-associated variations in human diseases without denying the contribution of other external factors such as the socio-economic, dietary and environmental background. Moreover, as the world becomes multi-ethnic, mixed marriage is becoming more common in some populations, making it difficult to assign a single ethnicity to an individual.

Consequently, understanding one's ancestry is a vital step in ensuring the proper designing of biomedical studies.

The development of AIM-SNP panels will help to address this problem, and the genetic admixture of an admixed individual will be fully discovered. This genomic knowledge will also be extremely useful in rapidly extending our knowledge of Mendelian diseases and individual responses to drugs that may be available for a specific disease. The accuracy of disease association studies can be enhanced, and false-positive and false-negative associations can be avoided. Individual genetic ancestry estimation can also bring us closer to more personalized/individualized genetic-based medicine.

In Malaysia, the Malays are the main ethnic group residing mostly in Peninsular Malaysia and some in Sabah and Sarawak. Malays are said to be an admixed population and demonstrated heterogeneous genetic structure (Hatin et al., 2011, Hatin et al., 2014, Deng et al., 2014, Deng et al., 2015, Hoh et al., 2015). The availability of Malay SNPs database hosted by the Malaysian Node of the Human Variome Project (MyHVP) and the Singapore Genome Variation Project (SGVP) enable the selection of AIM SNPs for the Malay population.

The Malay AIM-SNP panels can be used by the medical community to examine their subjects (who self-identified themselves as Malay) before recruiting them in their disease association studies, particularly for research involving ancestry/race as a vital contributor to the severity of the diseases. Pharmacology researchers may use these Malay AIM-SNPs to confirm their Malay ancestry subjects to gain an accurate

response to their drug research testing. These Malay AIM-SNPs are also useful in forensic investigations. The unknown samples found at the crime scene could be subjected to ancestry screening tests using these AIM-SNP panels, facilitating police investigations.

The MyHVP Malay SNP databases will be used in this study to select suitable SNPs that will be used to develop the Malay AIM-SNP panels using three approaches: pairwise  $F_{ST}$  (combined with ipPCA),  $I_n$  and PCA-correlated SNPs (or PCAIMs). It is hypothesized that hundreds of SNPs from these databases will be selected as AIM candidates for Malay ancestry. The accuracy of each resulting Malay AIM panel will be evaluated using a machine learning “ancestry-predictive model” constructed using WEKA, a comprehensive machine learning platform written in Java, and the genetic pattern based on the selected AIM-SNPs will be evaluated using the ADMIXTURE program. Consequently, the selected Malay AIM-SNPs can differentiate them from their closely related genetic ancestry counterpart and other world populations, and they can be used to describe admixture in the Malay population.

## **1.4 Objectives**

### **1.4.1 General Objective**

To select a subset of ancestry informative marker (AIM) SNPs from the Genome-wide Affymetrix GeneChip Mapping Xba 50 K array and Genome-wide Affymetrix SNP-6 array of the Malay SNPs datasets to design a panel of AIM SNPs that can estimate the ancestry of Malay population.

### 1.4.2 Specific Objectives

1. To identify a subset of autosomal AIM-SNPs from Malay SNP datasets, a genome-wide Affymetrix GeneChip Mapping Xba 50 K array with Pan-Asian SNP datasets, as reference samples.
2. To identify a subset of autosomal AIM-SNPs from Malay SNPs datasets, Genome-wide Affymetrix SNP-6 array platform with SGVP, and International HapMap Phase 3 project (HapMap) datasets as reference samples.
3. To compare the performance of the three AIM-SNP identification techniques: ipPCA combined with pairwise  $F_{ST}$  (ipPCA- $F_{ST}$ ),  $I_n$  and PCA-correlated SNPs (or PCAIMs).
4. To evaluate the genetic structure of the Malay population using the ipPCA and the ADMIXTURE program, as well as to validate the accuracy of the selected Malay AIM-SNP panels using ADMIXTURE and WEKA suite as the ancestry-predictive model.
5. To identify common SNPs between AIM-SNPs obtained from the Affymetrix GeneChip Mapping Xba 50 K Array and the Affymetrix SNP-6 array.

## CHAPTER 2

### LITERATURE REVIEW

#### 2.1 Ancestry

Ethnicity, race and ancestry are three terms that used interchangeably to describe a group of persons that show different appearance, comes from different background and geographic origins and practising different cultural and social. However, there are differences between the terms race, ethnic and ancestry. Race is a term that was used during the slaving era as a way to differentiate people based on their appearance such as skin colour, feature, social, geographic origin and cultural practice (National Human Genome Research Institute, 2005, Tishkoff and Kidd, 2004). Race is often equated with *continental ancestry* which assumes the existence of the five races which correspond well to geographic regions that are Africa, Europe, East Asia, Oceania and the Americas (Tishkoff and Kidd, 2004, Royal et al., 2010). Race was replaced by ethnicity in 20<sup>th</sup> centuries because of this term is more appropriate to differentiate peoples come from different biogeographic origin (National Human Genome Research Institute, 2005).

Ethnicity emphasized more on characteristic to group people such as cultural, languages, religion, custom, dress, diet, as well as their historical or territorial identity hence not just their genetic ancestry. However by using the ethnicity to group people it may suffers several shortcomings especially ethnic group contains subgroups for example, “Hispanic” ethnic group has subgroups Cuban Americans, Mexican Americans and Puerto Ricans and the Malay ethnic group has various subgroups such



as Melayu Bugis, Melayu Banjar, Melayu Minang and Melayu Kelantan. Combining these groups into a single category does not result in a better understanding of these groups (National Human Genome Research Institute, 2005, Tishkoff and Kidd, 2004). Ethnicity and racial which both used to cluster population by broad geographic origin is not adequate to represent genetic variation in human (Tishkoff and Kidd, 2004). Categorizing people in term of ancestry is more appropriate since it may recognize a single predominant source or multiple sources (National Human Genome Research Institute, 2005).

Self-identified race or ethnicity to identify individual in biomedical studies is not sufficient (Mersha and Abebe, 2015). Knowledge of ancestry can be important clinically and in biomedical studies due to the differential distribution of normal genetic variation and of genetic variation affecting diseases which were resulted from the genetic drift, natural selection, mutation and migration processes (Tishkoff and Kidd, 2004). Knowing the genetic ancestry of the subject in biomedical research studies will be important to avoid spurious results in biomedical studies related to population stratification or may be used to map susceptibility variants that might be differentially distributed with respect to ancestry (Royal et al., 2010).

According to US Census Bureau, ancestry is “*a person’s ethnic origin, the person’s parents or descent or the ‘roots’ where the person comes from*” (Royal et al., 2010). Ancestry may be referred to a group of persons who are geographically different (biogeographic ancestry) for example Asian, sub-Saharan African, African American, European; geopolitically different for example Vietnamese, Zambian, or Norwegian

or culturally different for example Brahmin, Lemba or Apache (National Human Genome Research Institute, 2005).

We are somewhat related to each other by varying degree because of our common origin as a species (Royal et al., 2010). Mutation and recombination process contributed to differences in an individual's genome and various segments of an individual's genome may have ancestral histories that trace to different populations (Royal et al., 2010). Ancestry of an individual can be inferred from the genetic data based on the differences of the allele frequencies of the loci geographically. So ancestry can be used as alternative of the race and ethnicity in genetic research to describe a group of people and to study the population histories and relationships. The inferring of ancestry from the genetic data can be used in genealogical, anthropological and epidemiological research (Royal et al., 2010). Self-identified race or ethnicity which is commonly used in pharmacological and biomedical studies can be accompanied by the information of the genetic ancestral and considering together the education, environmental, diet, socioeconomic and culture factors which may contributed to the health disparities (Mersha and Abebe, 2015, Royal et al., 2010).

## **2.2 Genetic markers used in ancestry studies**

The study of individual ancestry is useful in biomedical studies especially in the era of personalized medicine. Studies carried out by Mersha and Abebe (2015), showed limitation of self-reported race or ethnicity in biomedical research. Phenotype-based race information often disagrees with the genetic ancestry inferred using ancestry informative markers based on genetic or genomic data (Mersha and Abebe, 2015).

Estimation of the ancestry using genetic markers is not impossible by using the allele frequency differences between populations (Tishkoff and Kidd, 2004). Markers that demonstrate large differences in allele frequencies are good candidates to be used as a tool to estimate genetic ancestry (Barnes, 2010, Tishkoff and Kidd, 2004). Genetic markers such as Single Nucleotide Polymorphisms (SNPs), Y chromosomal haplotypes, mitochondrial DNA (mtDNA) haplotypes, X chromosomal, Short Tandem Repeats (STRs), Insertion deletion polymorphisms (INDEL), as well as Alu elements and Human Leukocyte Antigen polymorphisms (HLA) have been used to study the ancestry of the world populations (Mersha and Abebe, 2015).

Most of the ancestry studies have been carried out using haploid markers (mtDNA and Y chromosomal) (Cai et al., 2011, Corach et al., 2010, Dulik et al., 2012, Nunez et al., 2010) or multiple unlinked autosomal markers which are diploid where priority selected to be ancestry informative (Bryc et al., 2015, Gettings et al., 2014, Nakaoka et al., 2013). The selection of genetic markers in ancestry studies mostly rely on the purposes of the study where in genealogical and anthropological study haploid markers may be preferred whereas for the epidemiological research, it rely mostly on allele frequencies of the autosomal SNPs. Genetic epidemiologists employ the methods of ancestry inference using autosomal SNPs to identify genetic associations with diseases either to control for statistical biases related to population stratification among cases and controls or as a strategy to map susceptibility variants that might be differentially distributed with respect to ancestry (Royal et al., 2010).

Recently most of the studies of ancestry concentrated on the SNP marker on autosomal DNA. This is due to the haploid marker lack recombinant and autosomal give more information regarding the paternal and maternal ancestry which produce more accurate data. Moreover the studies of SNP data on autosomal DNA are more benefited to the epidemiological research especially the studies of relationships of the diseases susceptibility and protected SNPs.

On the other hand, combination of both haploid and diploid markers analysis might reveal more information about the genetic history of studied populations. Studies carried out by (Simonson et al., 2011) combined mitochondrial, Y-chromosomes and SNPs on autosomal chromosome to show that the Iban population from the Malaysian state of Sarawak exhibited genetic similarity with Indonesian and mainland Southeast Asian populations. Nunez et al., (2010) analysed the mitochondrial control region, Y-chromosome STRs and autosomal STRs to ascertain the origin of the Nicaraguan ancestors whereas Corach et al., (2010) used selected SNP marker on autosomal chromosome, mitochondrial DNA and Y-chromosome to study the genetic admixture of Argentineans.

### **2.2.1 Single Nucleotide Polymorphisms (SNPs)**

Single Nucleotide Polymorphisms (SNPs) is a bi-allelic genetic marker. It refers to a single base sequence variation at a particular point in the genome and can be found abundantly throughout the human genome with a frequency of about one in 1,000 bp (Brookes, 1999). SNP markers are mostly bi-allelic markers which usually have two alleles per marker; A/G, C/T, A/T, T/G, C/G or A/C (Butler, 2012). SNPs can be found

in coding regions of gene, non-coding regions or in the intergenic regions (Syvänen, 2001).

SNPs found in coding region can be further divided into two types; the synonymous and non-synonymous SNPs. The synonymous SNPs do not affect the protein sequence whereas the non-synonymous SNPs alter the function or structure of the encoded proteins resulted in the recessively or dominantly inherited monogenic disorder (Syvänen, 2001). SNPs is the simplest form of DNA variation among individuals but yet known to be a very important genetic marker which responsible to the phenotypic differences between individuals. This bi-allelic marker plays an important role to the differences of individual's drug responses as well as the progression and development of many genetic diseases and said to be directly involved in determining the phenotype of an individual such as eye colour, hair and skin; facial morphology and height (Butler, 2012). SNP has been extensively studied not only for the biomedical research, but also in the field of anthropological and forensic research.

SNPs have been the marker of choice in biomedical research due to its vital contributions to the function of the regulation and expression of a protein. The studies of SNPs not only can help us to predict the response of individual to various type of drugs or environmental toxin, and risk of developing particular diseases but also to track the inheritance of diseases genes within families (Kim and Misra, 2007). Case-control association studies are the most common application of SNPs in biomedical studies. Large SNPs genotyping data of both the patient and healthy control groups which contributes to the changes in cellular biological processes inducing diseased states usually studied in case-control association studies utilizing SNPs (Kim and

Misra, 2007). The establishment of the relationship between a genotype and a phenotype is based on the comparison of the differences of genotypes for all phenotypic characteristic demonstrated by the groups being studied (Kim and Misra, 2007). The information can be used to characterise the susceptibility genes associated with a disease, hence the encoded protein can be determined for prevention or treatment of the disease (Kim and Misra, 2007).

Pharmacogenomics studies are another discipline that utilizes SNPs as a tool to study the effects of genetic polymorphisms on drug response (Kim and Misra, 2007). This study is becoming popular due to the strong demand of the personalized medication. In pharmacogenomics study where SNPs are utilized as the markers, the aim is to elucidate effects of genetic polymorphisms on drug responses. Patients who have been administered a specific drug are the targeted groups and a large scale SNPs genotyping data is needed to ensure the accuracy and the effectiveness of the study (Kim and Misra, 2007).

In Malaysia, the study of SNPs related to *Helicobacter pylori* (*H.pylori*) and Thalassemia are amongst two common diseases that intensively studied in biomedical research. This is due to the significant differences of *H.pylori* infection prevalence rates among major ethnic groups in Malaysia (ie Malays, Chinese and Indians). The highest infected was observed among Indians adults whereas Malays exhibited low prevalence of *H.pylori* (Goh, 2018, Kumar et al., 2015, Lee et al., 2013, Sasidharan et al., 2011). On the other hand, Malays demonstrated high prevalence of Thalassemia (HbE  $\beta$ -thalassemia) compared to the Chinese and Indians populations (George, 2013).

*H.pylori* is a major gastric bacterial pathogen which has been said to be distributed through the routes of human migration resulted in the division of six ancestral populations; three from Africa, two from Asia and one from Europe (Tay et al., 2009). Study carried out by He et al. (2015) identified several SNPs related to gene that involved in the process of gastric carcinogenesis. Three genes; PGC, PTPN11, and IL1B said to be associated with the susceptibility to gastric carcinogenesis (He et al., 2015). They found the interactions of the SNPs of PGC (rs6912200 and rs4711690), PTPN11 (rs12229892) and IL1B (rs1143623) modified the risks of gastric cancer.

Thalassaemias are autosomal recessive disorders which are caused by the defective synthesis of the globin chain or faulty synthesis of haemoglobin (Yatim et al., 2014). Two common types of thalassemia are  $\alpha$  and  $\beta$ -thalassemia, which resulted from the defective synthesis of alpha and beta chains respectively (Yatim et al., 2014). Common type of thalassemia observed in Malays is HbE  $\beta$ -thalassemia (George, 2013). Nuinoon et al., (2010) reported three SNPs that have high association with  $\beta$ -thalassemia; SNPs of gene HBBP1 (rs2071348), gene HBS1L-MYB (rs9376092) and gene BCL11A (rs766432). Recent studies carried out by Cyrus et al., (2017) revealed another six SNPs located on chromosome 6 related to gene HBS1L-MYB (rs9376090, rs9399137, rs4895441, rs9389269, rs9402686, rs9494142, rs9376090) which has high association with the severity of the  $\beta$ -thalassemia.

### **2.2.1(a) SNPs as ancestry informative marker**

SNPs is a useful biological marker in the anthropological genetics research to study the variations among different groups of humans, reconstruct evolutionary history, human physical traits and it can reveal the history of modern human migration and their adaptation to different environments. SNPs have played an important role in genetic anthropological studies because this polymorphism are believed to be stable and not deleterious to organisms and can be population specific. Most SNPs are located in non-coding regions of the genome hence not known to influence phenotype of an individual and can be used in evolutionary studies. Numerous studies of SNPs marker on mitochondrial DNA (mtDNA) and Y-chromosomal as well as autosomal have been reported (Bryc et al., 2015, Dulik et al., 2012, Elhaik et al., 2013, Kivisild, 2015). Mitochondrial and Y-chromosomal uni-parentally markers have been predominantly used to study human migration in past decades; however the trend has recently shifted to autosomal SNPs. This is due to the breakthrough of the whole genome autosomal SNPs research and the development of miniaturized and automated procedure for analysing thousands of SNPs simultaneously (Kundu and Ghosh, 2015).

Autosomal SNPs marker has been widely used in the study of human migration and tracing the ancestors of human populations. This marker utilizes DNA from the 22 pairs of autosomal chromosomes contributed by both parents; hence more information regarding the history of the ancestors can be obtained compared to the haploid marker such as mitochondrial DNA and Y-chromosomal DNA. Numerous studies have been carried out on autosomal SNPs to infer ancestry across diploid genome (Hou et al., 2014, Huckins et al., 2014, Hwa et al., 2017, Galanter et al., 2012, Kersbergen et al., 2009, Rogalla et al., 2015, Phillips et al., 2013b, Sampson et al., 2011, Santos et al.,



2016, Santos et al., 2011, Vongpaisarnsin et al., 2015, Vongpaisarnsin et al., 2017, Wei et al., 2016).

SNPs distributed throughout the human genome that occur at very different frequencies in different world populations are good candidates as ancestry informative markers (Budowle and Daal, 2008). The use of SNPs as ancestry-informative marker has been numerous published recently (Hwa et al., 2017, Das and Upadhyai, 2018, Esposito et al., 2018, Setser et al., 2020, Vongpaisarnsin et al., 2017). Yang et al., (2005) have identified 199 ancestry informative markers which were distributed throughout the human genome. The SNPs were selected based on allele frequency differences in the ABI database comprising USA Caucasian, African American, Chinese and Japanese (Yang et al., 2005). Using the ancestry informative markers, they successfully demonstrated that both continental and sub-continental populations can be readily distinguished. Furthermore, the contribution of the putative parental population can be examined in admixed population using the SNPs ancestry informative markers (Yang et al., 2005).

Kidd et al., (2014) have developed a panel of 55 highly informative SNPs. The panel has been used to analyse 73 world populations and said to be very robust and efficient to provide excellent information on ancestry especially for forensic application (Kidd et al., 2014). In developing the SNPs ancestry panel, they used several sources of SNPs databases including the ABI databases, the HGDP-CEPH SNPs databases and their own laboratory databases. They successfully identified the largest pairwise allele frequencies differences between the studied populations to develop the panel of SNPs for ancestry inference. The set of 55 ancestry SNPs is opened for improvement because

not all populations were represented and tested and the SNPs selected were less efficient in estimating the admixed populations (Kidd et al., 2014). Subsequently, in 2017 they have added another 14 reference populations allele frequency to enhance the use of the 55 ancestry informative SNPs (Pakstis et al., 2017). The allele frequencies of all 55 SNPs for a total of 139 population samples are available publicly and have been incorporated in commercial kits by ThermoFisher Scientific and Illumina (Pakstis et al., 2017).

Recent studies carried out by Esposito et al. (2018) revealed the usefulness of a panel of ancient ancestry informative markers (aAIMs) in identifying fine-scale ancient population structure in Eurasians. They utilized more than 150 thousand autosomal SNPs from 302 ancient genome classified to 21 populations recovered from Europe, the Middle East and North Eurasia (Esposito et al., 2018). They demonstrated that principal component analysis (PCA)-based approach outperforms other methods such as the Infocalc and the Wright's  $F_{ST}$  in capturing ancient population structure and identifying admixed individuals. Their finding can also be used to improve the accuracy of genetic studies utilizing ancient DNA. On the other hand, Das and Upadhyai, (2018) have approved the robustness and efficiency of autosomal ancestry informative SNPs in analysing the fine genetic structure of the highly admixed population of South Asian genetic origins. Comparison of the three methods; Infocalc,  $F_{ST}$  and the Smart PCA based on their whole genome data of Indian subcontinent, shows that the Infocalc method gave the best results compared to the Smart PCA and  $F_{ST}$ .

### **2.2.1(b) SNPs genotyping**

SNPs genotyping protocol can be divided into two main parts; the biochemical reaction and the detection procedures (Chen and Sullivan, 2003). The biochemical reaction is basically the determination of the allele-specific products of the SNPs (Kim and Misra, 2007). Reviewed carried out by Sobrino et al., (2005) have listed a number of SNP genotyping chemistries such as allelic-specific hybridization, primer extension, oligonucleotide ligation and invasive cleavage. Primer extension involves the incorporation of nucleotides to the DNA template using specific enzyme to discriminate the SNPs alleles. A primer will be designed to anneal to the 3' end of the DNA template (SNPs) and nucleotides will be added to the template by the polymerase enzyme (Kim and Misra, 2007, Sobrino et al., 2005).

The allelic-specific hybridization involves the use of differences in thermal stability of the double-stranded DNA to discriminate the SNPs allele (Kim and Misra, 2007). Target-probe pairs must perfectly complementary to each other thus the effectiveness of the hybridization likely depending on the length and sequence of the probe, location of the SNPs and hybridization condition. This approach is suitable for the high throughput microarray platforms such as incorporated in the GeneChip<sup>®</sup> array technology (Affymetrix, CA) (Kim and Misra, 2007).

Oligonucleotide ligation is a technique where ligase enzymes are used to discriminate the SNPs allele. In these approach three oligonucleotides probes are involved, where the first two oligonucleotides are hybridized to the single stranded DNA template, adjacent to each other. Subsequently, the third probe binds to the template adjacent to the SNP immediately next to the allele-specific probe. The ligation product is detected

by various methods (Kim and Misra, 2007). Invasive cleavage involves the cleaving of the targeted DNA sequence by restriction enzyme and the product can be detected using gel electrophoresis. Invader<sup>®</sup> assay has adopted this technique by using two allele-specific probes attached with two types of dye at either end, the reporter (R) and the quencher (Q) and a common invader probe. The products can be detected using fluorescence analysis.

Allelic product of those methods can be detected with several detection systems such as fluorescence-electrophoresis, fluorescence resonance energy transfer (FRET), fluorescence polarization, fluorescence arrays, mass spectrometry and luminescence (Fondevila et al., 2013, Nicklas and Buel, 2008, Poetsch et al., 2013, Shi et al., 2011, Zeng et al., 2012). The fluorescence-electrophoresis detection method is perhaps the most accessible method due to the availability of the technique in most of the forensic laboratory around the world. The high throughput technology such as microarray (for example Affymetrix and Illumina platforms) is more powerful for analysing hundreds or thousands of SNPs simultaneously which had been used in many ancestry studies (Galanter et al., 2012, Keating et al., 2013, Krjutskov et al., 2009).

### **2.3 Ancestry informative markers (AIMs)**

Historically, human dispersed from East Africa to other part of the world more than 100,000 years ago resulted in the genetic diversity of modern human due to adaptation to new environment, geographical and climate change (Cavalli-Sforza, 2007). The migration contributed to the genetic variation because of the founder populations usually carried a portion of their most immediate ancestral population and at the same