

Use of GAPIT for Genome Wide Association Studies

Presented by Duke Pauli Montana State University



Hosted by Shawn Yarnes Plant Breeding and Genomics





Use of GAPIT for Genome Wide Association Studies

Duke Pauli

Montana State University





United States Department of Agriculture National Institute of Food and Agriculture



Overview

- Part 1
 - Learning objectives
 - Introduction and background
 - Installing RStudio
- Part 2
 - Sample Files
 - Loading GAPIT
- Part 3
 - Loading files, checking
 - Doing the analysis
- Part 4
 - Output
 - Questions



Learning Objectives

- Install RStudio
- Load needed packages and source code for GAPIT
- Load sample files
- Running basic analyses
- Using the output



Disclaimer

- What I am assuming for today's tutorial
 - Have run R before, familiar with basic commands
 - Have an understanding of the principals of GWAS
- What I <u>will not</u> cover today
 - Installing and running R
 - The eXtension website has a great tutorial
 - What GWAS is and the theory behind it
 - Genomic prediction functions of GAPIT



What is GAPIT?

- <u>Genome Association and Prediction Integrated Tool</u>
- Statistical package that is run in the R software environment
- Developed by Alex Lipka and Zhiwu Zhang
- Alexander E. Lipka et al. (2012) GAPIT: Genome Association and Prediction Integrated Tool. Bioinformatics. doi: 10.1093/ bioinformatics/bts444
- Uses statistical tools implemented in other programs like TASSEL

Why use GAPIT?

- Although there are programs that implement the same statistical methods, size of data sets can be a problem
- GAPIT can handle large data sets, file splitting
- Some of these programs have complex user interfaces, most of you are already familiar with R
- DOES NOT REQUIRE THAT MUCH CODE TO RUN AN ANALYSIS
- Do all data manipulations and analysis in one common environment
- Supporting information is good, excellent user manual
- Surprisingly easy to use
- Developers are extremely helpful
- Faster

Installing R studio

- Open source, integrated development environment (IDE)
- Nicer user interface for R
- To install, go to <u>http://www.rstudio.com/</u>, select download now (green button)
- Select "Download RStudio Desktop"
- Select the version that is recommended for your system
- Follow setup wizard
- Has a script window, R console, displays objects that are currently in your workspace, and another display that shows graphs or packages

RStudio



- 0

Running GAPIT

- USE THE MANUAL!!!!
- It is like an R vignette but a thousand times better
- Has all the code, can copy and paste in to R
- Describes the formatting of data for input
- Lists all the options for doing analyses
- Describes the output that is generated
- Explains more advanced uses of the program, more example code



Sample Files

- Barley data set generated for the Barley CAP
- All data available from http://
- Data from replicated yield trial, 768 lines
- Trait is grain protein content, important in malting/brewing
- GAPIT handles multiple traits



Sample Files

- Do an analysis using HapMap formatted genotype data
- Again, check manual for formatting
- Two files
 - HapMap_Genotypes data in HapMap format
 - Protein protein content of grain in %
- 2359 markers
- Missing data was set at 20%
 - (Markers missing more than 20% were removed)

HapMap Format

- Genetic map and marker data in one file
- Make sure there are no # signs in column headers, typically show up in assembly column
- Columns that are required for GAPIT are rs, chrom, pos everything else can be filled in NA with regard to assembly, center, protLSID etc.
- In the import process make sure header=F

x≣	II ☐														?	• – 1	8 ×
FI	LE HO	OME INS	ERT PAGE	LAYOUT	FORMULA	S DATA	REVIEW	VIEW							W	illiam Pauli	- 0
Pas	te ∛	Calibri B I U	- 11 -			≫ - ₽ € = ₽	Wrap Text Merge & Ce	Ge enter - \$	neral	 .00 Co .00 →.0 For 	nditional Formatting ▼	ormat as C Table - Sty	ell les •	sert • \sum elete • $$	Sort & F Filter • Se	ind & elect •	
Clipboard 🖬 Font 🖬 Alignment 🖬 Number 🖬 Styles Cells Editing																^	
A1 \checkmark : $\times \checkmark f_x$ rs															*		
	Α	В	С	D	E	F	G	Н	I	J	К	L	М	N	0	Р	
1	rs	alleles	chrom	pos	strand	assembly	center	protLSID	assayLSID	panelLSID	QCcode	HARRING	ROBUST	TRADITIO	AC_METC	BARONES	S ND
2	11_10002	A/G	1	6670	NA	NA	NA	NA	NA	NA	NA	AA	GG	GG	AA	AA	GG
3	11_10003	A/G	6	5275	NA	NA	NA	NA	NA	NA	NA	AA	AA	AA	AA	AA	AA
4	11_10006	5 T/C	1	7394	NA	NA	NA	NA	NA	NA	NA	TT	CC	TT	CC	TT	TT
5	11_10008	B T/G	3	5640	NA	NA	NA	NA	NA	NA	NA	GG	GG	GG	GG	TT	GG
6	11_10010	T/C	4	6600	NA	NA	NA	NA	NA	NA	NA	CC	CC	TT	CC	CC	CC
7	11_10011	C/G	3	5640	NA	NA	NA	NA	NA	NA	NA	GG	GG	GG	GG	CC	GG
8	11_10012	2 T/C	2	5890	NA	NA	NA	NA	NA	NA	NA	TT	CC	CC	TT	TT	CC
9	11_10013	s T/C	6	4544	NA	NA	NA	NA	NA	NA	NA	CC	TT	TT	CC	CC	TT
10	11_10015	6 C/G	6	9739	NA	NA	NA	NA	NA	NA	NA	GG	CC	CC	GG	GG	CC
11	11_10017	7/C	2	0	NA	NA	NA	NA	NA	NA	NA	TT	CC	CC	CC	TT	CC
12	11_10023	C/G	6	2235	NA	NA	NA	NA	NA	NA	NA	CC	GG	GG	CC	GG	GG
13	11_10024	T/C	5	10759	NA	NA	NA	NA	NA	NA	NA	CC	CC	CC	CC	CC	CC
14	11_10025	6 A/G	7	2113	NA	NA	NA	NA	NA	NA	NA	GG	GG	GG	AA	AA	GG
15	11_10026	5 T/C	3	3283	NA	NA	NA	NA	NA	NA	NA	TT	TT	TT	TT	TT	TT
16	11_10028	B A/G	4	0	NA	NA	NA	NA	NA	NA	NA	GG	GG	GG	AA	AA	GG
17	11_10030	A/G	1	1805	NA	NA	NA	NA	NA	NA	NA	AA	GG	GG	GG	GG	AA
10	•	HapMap	_genotypes	+						:	4	66	~~	<u> </u>	<u> </u>	~~~	

+ 100%

Phenotype File

- Straightforward text file
- Taxa used for header of entries
- GAPIT can run multiple phenotypes, i.e. analyze multiple traits in one analysis
- One phenotype file with multiple columns for different phenotypes
- Missing data coded as NA or NaN
- During data import use header=T

X≣	l 🖬 🕤	- @						pro	otein - Excel							? 🛧 -	- 8	×
E	ILE	HOME INS	SERT PAG	E LAYOUT	FORMUL	AS DAT	A REVIEW	VIEW								William Pa	auli 👻	-
Pas	ste	Calibri B I U	• 11 • • □ • ♪ Font	A A • <u>A</u> •	= = =		Wrap Text	enter 👻	General \$ • % • Number		Conditional Formatting ▼	Format as Ce Table - Style Styles		Insert • Delete • Format • Cells	∑ · A Z v Sort & Filter	Find & Select -		~
A1		- : X	√ <i>f</i> x ⊺	аха														*
		Α	В	С	D	E	F	G	Н	Ι	J	K	L	Μ	N	0		
1	Taxa		protein															
2	03WA-1	.05.4HW-1	13.3															
3	03WA-1	.68.14W-3	15.6															
4	03WA-2	03.18HW-1	14.8															
5	03WA-2	03.1HW-1	14.2															
6	03WA-2	03.9HW-1	14.9															
7	03WA-2	04.19W-1	15.8															
8	03WA-2	04.22W-4	13.7															
9	03WA-2	04.4W-2	11.8															
10	03WA-2	04.9W-1	13.5															
11	06WA-4	06.11	11.9															
12	06WA-4	06.18	12.5															
13	06WA-4	06.21	11.7															
14	06WA-4	06.3	12.4															
15	06WA-4	06.6	12.3															
16	06WA-4	06.9	12.6															
17	06WA-4	08.14	12.6															
10		protein	+								: •]	
REA	.DY												Ħ	B	• – —		- + 1	.00%

Running GAPIT

- Need to install some other packages that GAPIT uses
- GAPIT uses code from the Buckler lab website for running the package
- Using source code allows them to update the package continuously without having to constantly install new versions of the package
- Uses EMMA which also must be loaded
- Install.packages puts them into your library
- Library loads them for use in your current session
- ***Important if GAPIT isn't loading check website for updated code
 - http://www.maizegenetics.net/gapit

```
RStudio
File Edit Code View Plots Session Project Build Tools Help
🝳 📲 🥣 🖌 🔝 🔁 😓 🖉 🕗 🖓 😔 💭 😟
 Code for GAPIT webinar.R* * pheno *
                                                                                                                   -6
        📊 🔲 Source on Save 🛛 🔍 🎢 🗸
                                                                                                📑 Run 📑 🕞 Source 👻 📗
      # to run GAPIT we need to access some tools that are located in other packages. The following lines of
      # code will install the required packages in to your R library and then will import them so that GAPIT
      # may access them. The easiest way to do this is to highlight all of the following code and select run.
   3
      # Once you have done this, you will only need to use the library functions to access them at future dates.
   4
   5
   6
      source("http://www.bioconductor.org/biocLite.R")
      biocLite("multtest")
   7
   8
      install.packages("gplots")
   9
      install.packages("LDheatmap")
  10
  11
      install.packages("genetics")
  12
  13
      library(multtest)
      library("gplots")
  14
      library("LDheatmap")
  15
  16
      library("genetics")
  17
      library("compiler") #this library is already installed in R
  18
  19
      #these lines of code install the actual GAPIT functions from the maize genetics website
      source("http://www.maizegenetics.net/images/stories/bioinformatics/GAPIT/gapit_functions.txt")
  20
      source("http://www.maizegenetics.net/images/stories/bioinformatics/GAPIT/emma.txt")
  21
  22
```

Upload Files

- A few ways...
- Use the import wizard in RStudio
 - Under Tools tab, then import data set, will walk you through it
- Manually, need to know the address of the file
- Use file.choose()
 - Opens a file browser, navigate to the file, double click it
 - Displays file name in console
 - Copy that and use in the read.table statement
 - E.g. Example_data<-read.table("D;\\GAPIT\\protein.txt",header=T)

Uploading Files

RStud	dio	
File Ec	dit Code View Plots Session Project Build Tools Help	
•.•	達 🗝 🔚 🔚 🗁 Go to file/function	획 Pro
Cod	e for GAPIT webinar.R* × pheno ×	
	🕞 🔲 Source on Save 🛛 🔍 🎢 🗸	📑 Run 📑 📑 Source 🗸 🎚
21	Source(help.//www.matzegeneeres.nee/ mages/storres/oronnormatics/oniti/emma.exe /	
23	# Now we need to import the data files that we will be using for doing the analysis. Best to have	
24	# data saved in a text format. If you are using R studio you can import the files using the import	
25	# wizard located under the tools tab at the top. This may be the easiest as it walks you through the	
26	# process. You may perfer the command line method as it gives you more controll over the import of the	
28	# the file in the approviate directory. Select the file and the name of the file will appear in the	-
29	# consol window. Copy and paste this name into the read.table command.	
30		
31		
32	# If our data is in the HapMap format, you will only need the two files-the phenotypic and the hapmap	
33	# file. We will do this analysis first for the sake of time.	
25	TITE.CNOOSE() hanman genok-read table("D:\\CARIT\\HanMan genotynes tyt" header-E)	
36	pheno<-read.table("D:\\GAPIT\\protein.txt".header=T)	
37	human can can be a film a film occurrence file and - 1	

Doing the analyses

- Know your data
- Do a quick check to make sure files are properly loaded
 - Use the str() command to get information about the data object
- Plot the phenotype data
- Look at basic descriptive statistics
- Helps insure that data is properly loaded

Checking phenotype data

```
RStudio
File Edit Code View Plots Session Project Build Tools Help
된 📲 🥣 🖌 🔝 🛛 📥 🛛 🥕 Go to file/function
Code for GAPIT webinar.R* * pheno *
       📊 🔲 Source on Save 🛛 🔍 🎢 🗸
                                                                                                                          Run 🐤 Run
  37
  38
      # good idea to check our pehontype data to make sure the file strucutre is correct and how the data
     # is distributed, checking for outliers
  39
  40
     str(pheno) # gives us information on the object, in this case a data frame and other information
  41
      hist(pheno$protein) # creates a histogram plot of our data, things look pretty good, we have a
  42
       #lines that are bit extreme but with 775 lines we would expect about 27 to be at least 3 sd's away from
  43
  44
       # the mean.
  45
  46
      #some basic statistics to look at
  47
      mean(pheno$protein) # 12.23
      range(pheno$protein) # 9.5 to 15.8
  48
      sd(pheno$protein) # .88
  49
     which(is.na(pheno$protein)) # look for lines with missing data, there should be none which is confimred
  сн.
```

RStudio		X					
File Edit Code View Plots Session Project Build Tools Help							
된 🕶 🚽 🔜 📾 🛛 🚔 🖉 The Go to file/function	🔕 Project:	(None) 🝷					
Code for GAPIT webinar.R × pheno ×	Workspace History						
🗇 🗇 📄 🖸 Source on Save 🔍 🎽 - 👘 🕀 🕀 💼 🔂 🕀 Source - 🗐	🞯 🕞 🖙 Import Dataset 🕶 🧹	C					
27 # data. Using the following command, file.choose(), opens a browser allowing	Data						
28 # the file in the approviate directory. Select the file and the name of the	hapmap_geno 2360 obs. of 779 variables						
30	pheno 768 obs. of 2 variables						
31	Functions						
32 # If our data is in the HapMap format, you will only need the two files-the (=	GAPIT(Y = NULL, G = NULL, GD = NULL, GM = NULL, KI = NULL,	Z _					
34 file.choose()	= NULL CV = NULL CV Inheritance = NULL GP = NULL GK =						
<pre>35 hapmap_geno<-read.table("D:\\GAPIT\\HapMap_genotypes.txt",header=F)</pre>	Files Plots Packages Help						
<pre>36 pheno<-read.table("D:\\GAPIT\\protein.txt",header=T) 37</pre>	🖕 🧼 🔎 Zoom 🛛 🗷 Export 👻 👰 🗹 Clear All	C					
38 # good idea to check our pehontype data to make sure the file strucutre is c							
<pre>39 # is distributed, checking for outliers</pre>	Histogram of pheno\$protein						
40 41 str(pheno) # gives us information on the object, in this case a data frame a	riiotograii oi priototprotoin						
42 hist(pheno\$protein) # creates a histogram plot of our data, things look pret	8						
43:2 S (Top Level) ≎							
Console ~/ 🔗							
<pre>> source("http://www.maizegenetics.net/images/stories/bioinformatics/GAPIT/emma.tx > hapmap_geno<-read_table("D:\\GAPIT\\HapMap_genotypes_txt" header=E)</pre>							
<pre>> pheno<-read.table("D:\\GAPIT\\protein.txt",header=T)</pre>							
> str(pheno) # gives us information on the object, in this case a data frame and c							
'data.frame': 768 obs. of 2 variables:							
\$ Taxa : Factor w/ 768 levels "03WA-105.4HW-1",: 1 2 3 4 5 6 7 8 9 10							
<pre>\$ protein: num 13.3 15.6 14.8 14.2 14.9 15.8 13.7 11.8 13.5 11.9 > bist(phone\$protein) # creates a bistogram plot of our data things look protty s</pre>	10 11 12 13 14 15 16						
>	phono@protoin						
	phenosprotein						

GAPIT Function

- To run GAPIT using the HapMap formatted data need 2 files
 - G = HapMap formatted genotype file
 - Y = phenotype file
- The number of input parameters is large, refer to the manual
- We will do a basic analysis accounting for population structure, with and without using compression, using P3D
- When data is in HapMap format GAPIT can impute missing data using
 - Major allele (SNP.impute="Major")
 - Minor allele (SNP.impute="Minor")
 - Heterozygote (SNP.impute="Middle")

Code for analysis

Analysis1<- GAPIT(

Y = name of phenotype file,

G = name of genotype HapMap File,

SNP.impute = method of imputation,

PCA.total = number of principal components to use to control for population structure,

Major.allele.zero=T reports allele effect with respect to the minor allele,

Group.from = 768,

Group.to = 768,

Group.by = 1)

- Last three lines turn compression off
- Set working directory where you want files to be saved (Lots of output)
 - In RStudio Session tab, Set working directory, Choose directory
 - In R, File, Change directory
 - Manually with setwd("folder location")

21 # the following is a very basic analysis. Make sure to change directory to where the results 52 # will be saved. In R studio go to session at the top and select "set working directory" 53 # and select "choose directory"- this will allow you to browse to the correct folder in which you 54 # would like to save the results. Once you locate the correct folder, highlight it and hit select. 55 56 57 # first analysis where compression is not used in the model 58 analysis1<-GAPIT(59 Y=pheno, 60 G=hapmap_geno, 61 SNP.impute="Major", PCA.total=3, 62 63 Major.allele.zero=T, 64 group.from=768, 65 group.to=768, 66 group.by=1) 67 68 # now we use compression to see how that effects the outcome of the analysis 69 analysis2<-GAPIT(70 Y=pheno, 71 G=hapmap_geno, 72 SNP.impute="Major",

- 73 PCA.total=3,
- 74 Major.allele.zero=T)

file GDNULLUserThe common name of file for genotype map for numeric formatfile, GMNULLUserThe common name of file for genotype data in numeric formatfile, pathNULLUserPath for genotype filesfile, form0 >0 The first genotype files named sequentiallyfile, to0 >0 The first genotype files named sequentiallygroup, form1 >0 The starting number of groups of compressiongroup, form1 >1 The ending number of groups of compressiongroup, form1 >1 The ending number of groups of compressiongroup, form1 >1 The ending number of groups of compressionknabip algorithmvanRadenConselle and EMMAAlgorithm to derive kinship from genotypeknabip chotervanRadencomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipknabip choterNULLUserChornosome for LD analysisLD chornosomeNULLUserChornosome Scale and or center and scale the SNPs for LD analysisPCA scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs for Conducting PCASNP FDRAddDomGenetic modelSNP ADFAddDomGenetic modelSNP FDRAddDomGenetic modelSNP FDRFRUEFALSELogic variable to test SNPs sampled to estimate kinship and PCSSNP faction1 $>$ and <1Fraction of SNPs sampl	file.G	NULL	User	The common name of file for genotype in hapmap format
file.GMNULLUserThe common name of file for genotype data in numeric formatfile.pathNULLUserPath for genotype filesfile.form0>0The first genotype files named sequentiallyfile.to0>0The last genotype files named sequentiallygroup.fom10>0The last genotype files named sequentiallygroup.fom11The attring number of groups of compressiongroup.fom1000000>1The attring number of groups of compressiongroup.fom10000001The attring number of groups of compressiongroup.fom1000000>1The attring number of groups of compressiongroup.fom10000001Clustering algorithm to group individuals based on their kinshipkinship algorithmvaragecomplet, ward, single, mequitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship groupMarafiWarageClustering algorithm to group individuals based on their kinshipLD chromosomeNULLUserClustoring algorithm to group individuals based on their kinshipLD chromosomeNULLUserClustoring algorithm to group individuals based on their kinshipSNP FDRNoneScale_Centerd and scaledScale and/or center and scale the SNPs before conducting PCASNP FDRNoneScale_Centerd and scaledScale and/or center and scale the SNPs before conducting PCASNP FDRAddScale Add Centerd and scaledGenetic modelSNP FAITAddDom	file.GD	NULL	User	The common name of file for genotype map for numeric format
file pathNULLUserPath for genotype filesfile form0>0The first genotype files named sequentiallyfile form0>0The last genotype files named sequentiallygroup by10>0The group in interval of compressiongroup ho1>1The starting number of groups of compressiongroup to1000000>1The ending number of groups of compressiongroup to1000000>1Algorithm to derive kinship from genotypekinship.glorithmVanRadenLoiselle and EMMAAlgorithm to group individuals based on their kinshipkinship.glorithmvaragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship.groupMeanMax, Min, and MedianClustering algorithm to group individuals based on their kinshipLD chromosomeNULLUserLocationChromosome for LD analysisLD chromosomeNULLUserLocationScale and/or center and scale the SNPs before conducting PCASNP FDR1>0 and <1	file.GM	NULL	User	The common name of file for genotype data in numeric format
fileform0>0Inferiorm <td>file.path</td> <td>NULL</td> <td>User</td> <td>Path for genotype files</td>	file.path	NULL	User	Path for genotype files
file.to0>0The last genotype files named sequentiallygroup by10>0The grouping interval of compressiongroup.from1>1The starting number of groups of compressiongroup.to1000000>1The ending number of groups of compressionkinship.algorithmVaRadenLoiselle and EMMAAlgorithm to derive kinship from genotypekinship.chusteraveragecomplet, ward, single, mcquitty, median, and centrolClustering algorithm to group individuals based on their kinshipkinship.groupMultUserChronosome for LD analysisLD.choronosomeNULLUserCacaion (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.fDR1OmScale AlformationSNP.fARFAddDomGenetic modelSNP.fARFAddDomGenetic modelSNP.faction1Sola <1	file.from	0	>0	The first genotype files named sequentially
group.by10>0The grouping interval of compressiongroup.from1>1The starting number of groups of compressiongroup.fo1000000>1The ending number of groups of compressionkinship algorithmVanRadenLoiselle and EMMAAlgorithm to derive kinship from genotypekinship chyteraveragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship group.byMeanMax, Min, and MedianChromosome for LD analysisLD chromosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	file.to	0	>0	The last genotype files named sequentially
group.fom1>1The starting number of groups of compressiongroup.to1000000>1The ending number of groups of compressionkinship.algorithmVanRadenLoiselle and EMMAAlgorithm to derive kinship from genotypekinship.clusteraveragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship.groupMaanMax, Min, and MedianMethod to derive kinship among groupsLD.chromosomeNULLUserChromosome for LD analysisLD.chromosomeNULLUserCacian (center) of SNPs for LD analysisPCA.scalingNoneScale, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	group.by	10	>0	The grouping interval of compression
group.to1000000>1The ending number of groups of compressionkinship.algorithmVanRadenLoiselle and EMMAAlgorithm to derive kinship from genotypekinship.clusteraveragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship.groupMeanMax, Min, and MedianMethod to derive kinship among groupsLD.chromosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	group.from	1	>1	The starting number of groups of compression
kinship.algorithmVanRadenLoiselle and EMMAAlgorithm to derive kinship from genotypekinship.chusteraveragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship.groupMeanMax, Min, and MedianMethod to derive kinship among groupsLD.chormosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	group.to	1000000	>1	The ending number of groups of compression
kinship.clusteraveragecomplete, ward, single, mcquitty, median, and centroidClustering algorithm to group individuals based on their kinshipkinship.groupMeanMax, Min, and MedianMethod to derive kinship among groupsLD.chromosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	kinship.algorithm	VanRaden	Loiselle and EMMA	Algorithm to derive kinship from genotype
kinship.groupMeanMax, Min, and MedianMethod to derive kinship among groupsLD.chromosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered.and.scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	kinship.cluster	average	complete, ward, single, mcquitty, median, and centroid	Clustering algorithm to group individuals based on their kinship
LD.chromosomeNULLUserChromosome for LD analysisLD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	kinship.group	Mean	Max, Min, and Median	Method to derive kinship among groups
LD.locationNULLUserLocation (center) of SNPs for LD analysisPCA.scalingNoneScaled, Centered and scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	LD.chromosome	NULL	User	Chromosome for LD analysis
PCA scalingNoneScaled, Centered, and, scaledScale and/or center and scale the SNPs before conducting PCASNP.FDR1>0 and <1	LD.location	NULL	User	Location (center) of SNPs for LD analysis
SNP.FDR1>0 and <1Threshold to filter SNP on FDRSNP.MAF0>0 and <1	PCA.scaling	None	Scaled, Centered.and.scaled	Scale and/or center and scale the SNPs before conducting PCA
SNP.MAF0>0 and <1Minor Allele Frequency to filter SNPs in GWAS reportsSNP.effectAddDomGenetic modelSNP.P3DTRUEFALSELogic variable to use P3D or not for testing SNPsSNP.fraction1>0 and <1	SNP.FDR	1	>0 and <1	Threshold to filter SNP on FDR
SNP.effectAddDomGenetic modelSNP.P3DTRUEFALSELogic variable to use P3D or not for testing SNPsSNP.fraction1>0 and <1	SNP.MAF	0	>0 and <1	Minor Allele Frequency to filter SNPs in GWAS reports
SNP.P3DTRUEFALSELogic variable to use P3D or not for testing SNPsSNP.fraction1>0 and <1	SNP.effect	Add	Dom	Genetic model
SNP.fraction 1 >0 and <1 Fraction of SNPs sampled to estimate kinship and PCs SNP.test TRUE FALSE Logic variable to test SNPs or not	SNP P3D			
SNP.test TRUE FALSE Logic variable to test SNPs or not	5111.1.515	TRUE	FALSE	Logic variable to use P3D or not for testing SNPs
	SNP.fraction	TRUE 1	FALSE >0 and <1	Logic variable to use P3D or not for testing SNPs Fraction of SNPs sampled to estimate kinship and PCs

Output

- QQ plot
- PCA graph
- Allelic Effects Estimate
- BLUPs
- GWAS results
- Manhattan Plots
 - Genome-wide
 - Chromosome specific
- There is more output than this
- Model Selection (BIC)

QQ Plot

.protein

Population Str.

GWAS Results

	Α	В	Insert	Function	E	F	G	Н	
1	SNP	Chromoso	Position	P.value	maf	nobs	Rsquare.of.Model.without.SNP	Rsquare.of.Model.with.SNP	FDR_Adjusted_P-values
2	12_10811	6	4544	7.96E-10	0.160156	768	0.453308974	0.481035038	1.60E-06
3	12_10199	6	4544	2.03E-09	0.172526	768	0.453308974	0.479661509	1.60E-06
4	12_10575	6	4544	2.03E-09	0.172526	768	0.453308974	0.479661509	1.60E-06
5	12_20685	7	9179	4.00E-09	0.049479	768	0.453308974	0.478672367	2.36E-06
6	12_11437	7	9179	2.97E-08	0.048177	768	0.453308974	0.475758998	1.40E-05
7	12_30301	7	9179	1.52E-06	0.041667	768	0.453308974	0.470114965	0.000599332
8	11_20340	2	8592	3.99E-06	0.435547	768	0.453308974	0.46875265	0.001345092
9	12_11353	6	5275	1.12E-05	0.061198	768	0.453308974	0.467302678	0.002932349
10	12_30032	6	5275	1.12E-05	0.061198	768	0.453308974	0.467302678	0.002932349
11	11_10697	4	11466	1.56E-05	0.013021	768	0.453308974	0.466836667	0.003680907
12	11_10003	6	5275	0.000146	0.221354	768	0.453308974	0.463730878	0.031397203
13	11_10256	7	7785	0.000171	0.03776	768	0.453308974	0.463519898	0.033560523
14	11_21239	5	6934	0.000282	0.227865	768	0.453308974	0.462832961	0.051170938
15	11_20714	6	6704	0.000383	0.132161	768	0.453308974	0.462415252	0.060326767
16	11_10095	5	13716	0.000403	0.231771	768	0.453308974	0.462345234	0.060326767
17	12 10938	1	5060	0.000409	0.029297	768	0.453308974	0.462325767	0.060326767

Allelic Effects Estimate

	A	В	С	D	Ł	F
1	SNP	Chromoso	Position	Allelic Effe	ct Estimate	;
2	11_10002	1	6670	0.002848		
3	11_10003	6	5275	-0.23801		
4	11_10006	1	7394	-0.03255		
5	11_10008	3	5640	0.129514		
6	11_10010	4	6600	-0.04677		
7	11_10011	3	5640	0.100203		
8	11_10012	2	5890	-0.27811		
9	11_10013	6	4544	0.155822		
10	11_10015	6	9739	0.016424		
11	11_10017	2	0	-0.09578		
12	11_10023	6	2235	0.012274		
13	11 10024	5	10759	0.045846		

Genome-wide Manhattan Plot

.protein

Chromosome Specific Manhattan Plot

Optimal Compression Parameters

- Optimal method to calculate group kinship is "Mean"
- Optimal clustering method is "Average"
- Number of groups is 758
- -2*log likelihood is 1487.65
- Heritability is 0.824

Your data

- Most of the time the data is not in HapMap form, what do you do?
- Build your own HapMap file, code SNPs as AA, AC, CC
- GAPIT does not impute marker data for numerically formatted genotype files
 - Impute in another program like TASSEL and import that genotype file
 - Impute with R, use marker average for missing data etc.
- If you use a numeric format you need to provide a genetic map file that has marker information
- In the GAPIT function the G = genotype file is changed to GD = genotype file and you must also include GM = genetic map file

Code for numeric data

83	
84	analysis2<-GAPIT(
85	Y=pheno,
86	GD=numeric_geno,
87	GM=genetic_map,
88	SNP.impute="Major",
89	PCA.total=3,
90	Major.allele.zero=T
91	
02	

Genetic Map

۵		alibri	· 11 ·	A A =		»-	Wrap Text	
Pas	te 💉 I	в <u>I</u> <u>U</u> -	- 5	• <u>A</u> • I		∉ ≇ 🛱	Merge & Ce	nter 👻
Clip	board 🗔		Font	F2		Alignment		5
A1	•	: 🗙 🗸	fx N	ame				
	Α	В	С	D	E	F	G	F
1	Name	Chromoso	Position					
2	11_10002	1	66.7					
3	11_10003	6	52.75					
4	11_10006	1	73.94					
5	11_10008	3	56.4					
6	11_10010	4	66					
7	11_10011	3	56.4					
8	11_10012	2	58.9					
9	11_10013	6	45.44					
10	11_10015	6	97.39					
11	11_10017	2	0					
12	11_10023	6	22.35					
13	11_10024	5	107.59					
14	11_10025	7	21.13					
15	11_10026	3	32.83					
16	11_10028	4	0					

Numeric Genotype Format

Pas	∎	Calibri B I U -	- 11			⊗∙ ≩≀ ∉ ∉ ₫।	Wrap Text Merge & Ce	nter - \$	eneral	▼ .00 .00 .00 →.0	Conditional Formatting ▼	Format as Ce Table - Style	Inser Dele Formes	rt ▼ Σ ete ▼ ↓ ▼ nat ▼	Sort & F	ind & elect •	
Clip	board 🗔		Font	E.		Alignment		E.	Number	F <u>a</u>		Styles	Cell	s	Editing		^
A1	-	: ×	f _x	Гаха													٣
	Α	В	С	D	E	F	G	Н	Ι	J	K	L	Μ	Ν	0	Р	
1	Taxa	11_10002	11_10003	11_10006	11_10008	11_10010	11_10011	11_1001	2 11_10013	11_100	015 11_1001	7 11_10023	11_10024 1	1_10025	11_10026	11_10028	3 11_
2	03WA-10	5 2	2	2 2	0	0	0		2 0		2	2 2	0	2	2	2	2
3	03WA-168	8 2	2	2 0	2	2	2		2 0		2	0 0	0	0	2	(C
4	03WA-203	3 2	2	2 0	0	0	2		2 0		0	2 0	0	2	2	2	2
5	03WA-203	3 2	2	2 2	0	0	0		2 0		0	2 0	0	0	2	2	2
6	03WA-20	3 2	2	2 0	0	0	0		2 0		2	2 0	0	0	2	2	2
7	03WA-204	4 2	2	2 2	0	0	0		2 0		0	2 0	1	2	2	2	2
8	03WA-204	4 2	2	2 2	2	0	2		2 0		0	2 0	2	0	2	2	2
9	03WA-204	4 2	2	2 2	0	0	0		2 0		2	2 0	0	0	2	1	2
10	03WA-204	4 2	2	2 2	2	0	2		2 0		0	2 0	2	0	2	1	2
11	06WA-40	6 2	2	2 0	0	0	0		2 0		0	0 2	0	0	2	()
12	06WA-40	6 2	2	2 2	0	0	0		2 0		0	0 1	0	0	2	()
13	06WA-40	6 2	2	2 0	0	0	0		2 0		0	0 2	0	0	2	()
14	06WA-40	6 2	2	2 2	0	0	0		2 0		0	0 2	0	0	2	1	1
15	06WA-40	6 2	2	2 2	0	0	0		2 0		0	0 1	0	2	2	2	2
16	06WA-40	6 2	2	2 2	0	0	0		2 0		0	0 0	0	2	2	()
17	06WA-408	8 2	2	2 2	0	0	0		2 0		0	0 0	0	2	2	1	2
10		Numeric_	genotype	+			0				: •		2	2			

Problems

- Consult the manual, example code for nearly all situations that might arise
- Check the website, check the date on the user manual because that is updated as well as the code. The updated manual reflects changes to the code
- Tutorial and example data sets
- Contact the makers very friendly and EXTREMELY helpful

Getting the most out of GAPIT

- Adjust analysis parameters
- Run several analyses
- Again, consult the model for more advanced features
- There is more example code in the provided script file that uses more of the options for running an analysis

Thank you

- Alex Lipka and Zhiwu Zhang
- eXtension Team
- Participants thanks for showing up

Please fill out the survey evaluation. You will be contacted via email.

Today's Presentation, Sample Data, and Links Available http://www.extension.org/pages/68355

> Sign up for PBG News http://pbgworks.org

Sign up for Future Webinars and View Archive http://www.extension.org/pages/60426

