

# Cardiac ACMG-ClinVar Penetrance Estimation

*James Diao, under the supervision of Arjun Manrai*

*June 26, 2017*

## Contents

<b>1</b>	<b>Download, Transform, and Load Data</b>	<b>2</b>
1.1	Collect ACMG Gene Panel . . . . .	2
1.2	Download ClinVar VCF . . . . .	3
1.3	Download 1000 Genomes VCFs . . . . .	3
1.4	Import and Process 1000 Genomes VCFs . . . . .	4
1.5	Import and Process gnomAD/ExAC VCFs . . . . .	4
1.6	Collect 1000 Genomes Phase 3 Populations Map . . . . .	5
<b>2</b>	<b>Common Pathogenic Variants by Ancestry</b>	<b>6</b>

**Working Directory:** /Users/jamesdiao/Documents/Kohane\_Lab/2017-ACMG-penetrance/ACMG\_Penetrance

# 1 Download, Transform, and Load Data

## 1.1 Collect ACMG Gene Panel

<http://www.ncbi.nlm.nih.gov/clinvar/docs/acmg/>

## Table from ACMG SF v2.0 Paper 60 x 8 (selected rows):

	Phenotype	MIM_disorder	PMID_Gene_Reviews_entry
<b>N1</b>	Hereditary breast and ovarian cancer	604370 612555	20301425
<b>N2</b>	Hereditary breast and ovarian cancer	604370 612555	20301425
<b>N3</b>	Li-Fraumeni syndrome	151623	20301488
<b>N4</b>	Peutz-Jeghers syndrome	175200	20301443
<b>N5</b>	Lynch syndrome	120435	20301390

Table continues below

	Typical_age_of_onset	Gene	MIM_gene	Inheritance	Variants_to_report
<b>N1</b>	Adult	BRCA1	113705	AD	KP&EP
<b>N2</b>	Adult	BRCA2	600185	AD	KP&EP
<b>N3</b>	Child/Adult	TP53	191170	AD	KP&EP
<b>N4</b>	Child/Adult	STK11	602216	AD	KP&EP
<b>N5</b>	Adult	MLH1	120436	AD	KP&EP

## ACMG-59 Genes:

```
## [1] BRCA1 BRCA2 TP53 STK11 MLH1 MSH2 MSH6 PMS2
## [9] APC MUTYH BMPR1A SMAD4 VHL MEN1 RET PTEN
## [17] RB1 SDHD SDHAF2 SDHC SDHB TSC1 TSC2 WT1
## [25] NF2 COL3A1 FBN1 TGFBR1 TGFBR2 SMAD3 ACTA2 MYH11
## [33] MYBPC3 MYH7 TNNT2 TNNI3 TPM1 MYL3 ACTC1 PRKAG2
## [41] GLA MYL2 LMNA RYR2 PKP2 DSP DSC2 TMEM43
## [49] DSG2 KCNQ1 KCNH2 SCN5A LDLR APOB PCSK9 ATP7B
## [57] OTC RYR1 CACNA1S
```

## 1.2 Download ClinVar VCF

[ftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf\\_GRCh37/clinvar.vcf.gz](ftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf_GRCh37/clinvar.vcf.gz)

ClinVar is the central repository for variant interpretations. Relevant information from the VCF includes:

(a) CLNSIG = “Variant Clinical Significance, 0 - Uncertain, 1 - Not provided, 2 - Benign, 3 - Likely benign, 4 - Likely pathogenic, 5 - Pathogenic, 6 - Drug response, 7 - Histocompatibility, 255 - Other”

(b) CLNDBN = “Variant disease name”

(c) CLNDSDBID = “Variant disease database ID”

(d) CLNREVSTAT = “Review Status, no\_assertion, no\_criteria, single - criterion provided single submitter, mult - criteria provided multiple submitters no conflicts, conf - criteria provided conflicting interpretations, exp - Reviewed by expert panel, guideline - Practice guideline”

(e) INTERP = Pathogenicity (likely pathogenic or pathogenic; CLNSIG = 4 or 5)

## Processed ClinVar data frame 224657 x 20 (selected rows/columns):

VAR_ID	CHROM	POS	ID	REF	ALT	CLNSIG	INTERP
1_955619_G_C	1	955619	.	G	C	Likely_benign	FALSE
1_957568_A_G	1	957568	.	A	G	Uncertain_significance	FALSE
1_957605_G_A	1	957605	.	G	A	Likely_benign	TRUE
1_957640_C_T	1	957640	.	C	T	Uncertain_significance	FALSE

Table continues below

GOLD_STARS	pathogenic	benign	conflicted	MSID	CLNREVSTAT	CLNDSDBID
1	FALSE	TRUE	FALSE	210112	1	1
1	FALSE	TRUE	FALSE	263166	1	1
0	TRUE	FALSE	FALSE	243036	1	1
1	FALSE	TRUE	FALSE	128296	1	1

## 1.3 Download 1000 Genomes VCFs

[ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/ALL.\[chrom\].phase3\\_\[version\].20130502.genotypes.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/ALL.[chrom].phase3_[version].20130502.genotypes.vcf.gz)

Downloaded 1000 Genomes VCFs are saved in: /Users/jamesdiao/Documents/Kohane\_Lab/2017-ACMG-penetrance/1000G/

gene	name	chrom	start	end	downloaded
BRCA1	NM_007294	17	41196311	41277500	TRUE
BRCA2	NM_000059	13	32889616	32973809	TRUE
TP53	NM_000546	17	7571719	7590868	TRUE
STK11	NM_000455	19	1205797	1228434	TRUE
MLH1	NM_000249	3	37034840	37092337	TRUE

## 1.4 Import and Process 1000 Genomes VCFs

- Unnest the data frames to 1 row per variant\_ID key (CHROM\_POSITION\_REF\_ALT).
- Remove all insertions, deletions, CNV, etc, and keep only missense variants (1 REF, 1 ALT)
- For 1000 Genomes: convert genomes to allele counts. For example: (0|1) becomes 1, (1|1) becomes 2. Multiple alleles are unnested into multiple counts. For example: (0|2) becomes 0 for the first allele (no 1s) and 1 for the second allele (one 2).

## Processed 1000 Genomes VCFs: 141467 x 2516 (selected rows/columns):

GENE	AF_1000G	VAR_ID	CHROM	POS	ID	REF	ALT
BRCA1	0.004193290	17_41196363_C_T	17	41196363	rs8176320	C	T
BRCA1	0.008386580	17_41196368_C_T	17	41196368	rs184237074	C	T
BRCA1	0.000998403	17_41196372_T_C	17	41196372	rs189382442	T	C
BRCA1	0.342252000	17_41196408_G_A	17	41196408	rs12516	G	A
BRCA1	0.000399361	17_41196409_G_C	17	41196409	rs548275991	G	C

Table continues below

HG00096	HG00097	HG00099	HG00100	HG00101	HG00102
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
1	0	1	1	0	2
0	0	0	0	0	0

## 1.5 Import and Process gnomAD/ExAC VCFs

- Unnest the data frames to 1 row per variant\_ID key (CHROM\_POSITION\_REF\_ALT).
- Remove all insertions, deletions, CNV, etc, and keep only missense variants (1 REF, 1 ALT)
- Collect superpopulation-level allele frequencies: African = AFR, Latino = AMR, European (Finnish + Non-Finnish) = EUR, East.Asian = EAS, South.Asian = SAS.

## Processed gnomAD VCFs: 96742 x 48 (selected rows/columns):

	GENE	AF_GNOMAD	VAR_ID
<b>8337</b>	MLH1	0.00007430	3_37092125_A_G
<b>37134</b>	FBN1	0.00001190	15_48725034_C_T
<b>66525</b>	DSP	0.00003190	6_7583743_G_A
<b>32042</b>	TSC2	0.00000524	16_2136877_G_A
<b>17246</b>	MUTYH	0.00003310	1_45804261_G_A

## Processed ExAC VCFs: 59883 x 45 (selected rows/columns):

	GENE	AF_EXAC	VAR_ID
<b>9547</b>	APC	0.003545000	5_112174456_A_T
<b>21545</b>	WT1	0.000008336	11_32450117_G_A
<b>37094</b>	RYS2	0.000009841	1_237755051_A_G
<b>41240</b>	DSP	0.000008281	6_7581713_C_G
<b>44501</b>	KCNQ1	0.000008297	11_2608792_G_T

## 1.6 Collect 1000 Genomes Phase 3 Populations Map

This will allow us to assign genotypes from the 1000 Genomes VCF to ancestral groups.

From: [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/integrated\\_call\\_samples\\_v3.20130502.ALL.panel](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/integrated_call_samples_v3.20130502.ALL.panel)

## Phase 3 Populations Map Table: 2504 x 4 (selected rows)

sample	pop	super_pop	gender
HG02757	GWD	AFR	female
HG03175	ESN	AFR	male
NA19780	MXL	AMR	male
HG02265	PEL	AMR	male
NA18619	CHB	EAS	female
HG00534	CHS	EAS	female

## 2 Common Pathogenic Variants by Ancestry

