

Project Documentation

Team Duck Duck Go

Contents

Introduction	3
Prerequisites:	3
Functions used throughout:	3
Structure	3
GWAS	3
Functional information and Gene Ontology	4
Linkage Disequilibrium	7

Introduction

This web application prototype is designed to retrieve information on Single Nucleotide Polymorphisms (SNPs) seen in Type 1 Diabetes patients identified by Genome wide association studies (GWAS). The database will use information from the GWAS catalogue, along with population data from Ensembl and the 1000 Genomes Project and functional information and Gene Ontology information obtained through Ensembl's VEP tool which is all is retrievable through a user friendly interface through the input of an rsID, Chromosome position or a Gene name. The site also allows the user to calculate Linkage Disequilibrium (LD) of SNPs selected for each population producing a text file containing the LD values and plot these values as a LD heatmap. The user is also able to enter multiple SNPs and return a Manhattan plot of the p-values.

Prerequisites:

Functions used throughout:

A number of functions were used throughout the code such as...

some python code will go here...

Structure

GWAS

This information was downloaded from the GWAS catalogue where a TSV file was downloaded and then trimmed:

some python code will go here...

This code uses pandas to open the TSV file, creates a dataframe called data, and removes any special characters from the column names, as SQL does not interact with special characters very well.

some python code will go here...

The data frame is then filtered further so that it only has data that references T1D in the "disease_trait" column so that only T1D data remains in the dataframe. Next all SNPs that don't have rsIDs are removed, as some cells had incompatible data in this column.

some python code will go here...

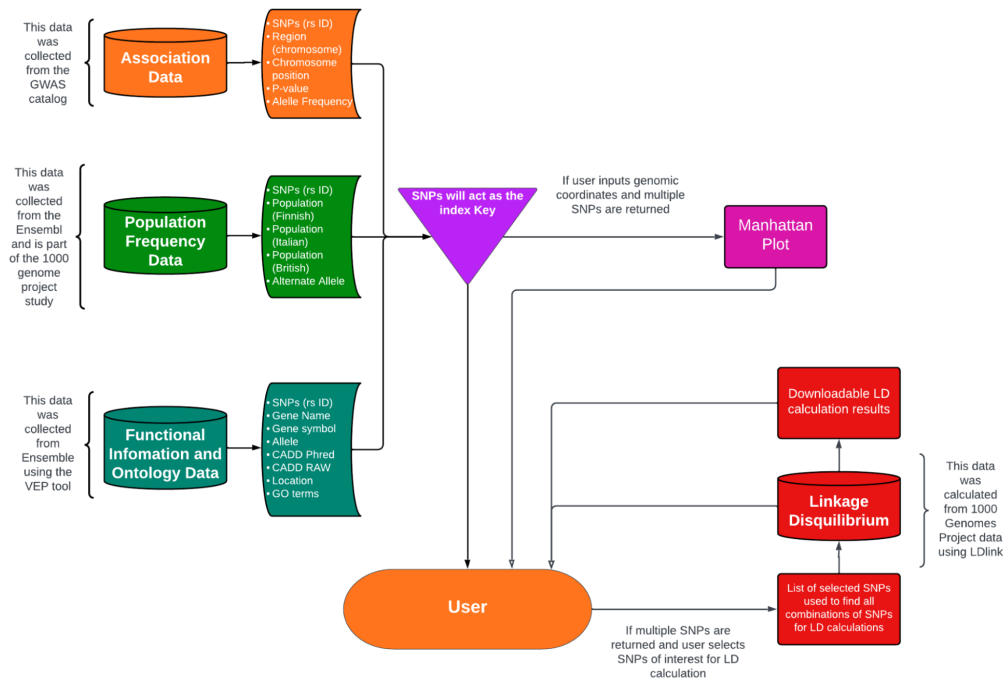


Figure 1: Structure of the data flow

The dataframe is then trimmed again so that it contains only the columns of interest "snps", "region", "chr_pos", "chr_id", "p_value", "mapped_gene".

some python code will go here...

Duplicate mapped_gene information is then removed.

some python code will go here...

Next the column 'snps' was renamed to 'rsid' and made the change directly to the dataframe by setting inplace=True.

Functional information and Gene Ontology

We have used Ensembl's Variant Effect Predictor web tool to gather the Functional and Ontology data by submitting a job with the rsIDs as input from the above GWAS TSV file. After running the job, we get an output of text file with columns like:

- rsid: the reference SNP identifier for the variant.
- Allele: the alternative allele observed at the variant site.

- impact: the impact of the variant on the affected gene
- consequence: the consequence of the variant on the affected gene
- location: the location of the variant within the affected gene
- Gene: the Ensembl gene ID of the affected gene.
- Symbol: the gene symbol or name.
- Feature_Type: the type of genomic feature the variant is located in
- Feature: the Ensembl ID of the specific feature the variant is located in
- Exon: the exon number(s) affected by the variant.
- Intron: the intron number(s) affected by the variant.
- HGVS_c: the HGVS nomenclature for the variant at the cDNA .
- HGVS_p: the HGVS nomenclature for the variant at the protein level.
- cDNA_position: the position of the variant within the cDNA sequence of the affected gene.
- CDS_position: the position of the variant within the coding sequence of the affected gene.
- Protein_position: the position of the variant within the protein sequence of the affected gene.
- Amino_acids: the amino acid change resulting from the variant.
- Codons: the DNA codon change resulting from the variant.
- Existing_variation: additional identifiers for the variant in other databases.
- Distance: the distance to the nearest feature in the same or opposite strand.
- Strand: the genomic strand the variant is located on.
- FLAGS: additional information about the variant .
- SYMBOL_SOURCE: the source of the gene symbol or name.
- HGNC_ID: the HGNC gene ID of the affected gene.
- MANE_SELECT: indication of whether this transcript is the MANE (Matched Annotation from NCBI and EMBL-EBI) Select transcript.

- MANE_PLUS_CLINICAL: indication of whether this transcript is the MANE Select Plus Clinical (MPC) transcript.
- TSL: transcript support level (a measure of transcript annotation confidence).
- APPRIS: annotation of principal isoforms for each gene.
- ENSP: the Ensembl protein ID of the affected protein.
- SIFT: prediction of the effect of the variant on protein function.
- PolyPhen: prediction of the effect of the variant on protein function.
- CLIN_SIG: clinical significance of the variant.
- SOMATIC: indication of whether the variant is somatic or germline.
- PHENO: phenotype association of the variant.
- PUBMED: PubMed ID of publications reporting functional evidence of the variant.
- MOTIF_NAME: name of the DNA motif affected by the variant.
- MOTIF_POS: position of the variant within the affected DNA motif.
- HIGH_INF_POS: indication of whether the variant falls in a highly conserved position within the DNA motif.
- MOTIF_SCORE_CHANGE: the effect of the variant on the score of the affected DNA motif.
- TRANSCRIPTION_FACTORS: transcription factors that bind the affected DNA motif.
- CADD_PHRED: Phred-scaled CADD score (Combined Annotation-Dependent Depletion), which predicts the deleteriousness of variants.
- CADD_RAW: the raw CADD score, which is a measure of the deleteriousness of variants.
- GO Terms: Gene Ontology (GO) terms associated with the affected gene.

The VEP file provides detailed information about the functional and ontological consequences of genetic variants, including their impact on genes, proteins, and pathways.

We have converted the text file to a tsv file and trimmed down the file to include only the rsID, Alleles, CADD_PHRED and CADD_RAW scores columns for the functional data and the rsID, location, gene, symbol, GO terms columns for the Ontology data. We have further used these TSV files in our database.

CADD (Combined Annotation Dependent Depletion) is a tool used for predicting the potential harm caused by genetic variants. The tool generates a score that indicates the likelihood of a variant being deleterious.

One advantage of CADD over SIFT and PolyPhen is that CADD integrates a larger and more diverse set of functional annotations. It also considers the effects of variants on non-coding regions of the genome, which can be important for understanding the functional consequences of variants that are not in protein-coding regions.

Linkage Disequilibrium

Linkage disequilibrium (LD) is the degree of non-random association of the allele of one SNP with the allele of another SNP within a population. LD is typically measured by two metrics: D' and r^2 .

D' is the normalised values of D , the coefficient of linkage disequilibrium, where A and B are alleles of two SNPs in different loci:

$$D = P_A - P_A P_B$$

r^2 is the correlation coefficient between two loci:

$$r^2 = \frac{D^2}{p_A(1 - p_A)p_B(1 - p_B)}$$

Collecting data

Linkage disequilibrium data was obtained from LDlink using the LDmatrix tool. D' and r^2 values for SNPs were calculated using 1000 Genomes Project data for all three populations. LD data was obtained by inputting a list of SNPs from the same chromosome and selecting the population which would be used for allele frequency data for LD calculations. LDmatrix would produce two text files containing a matrix of results for D' and r^2 values calculated between all SNPs pair combinations in the input list. This was performed

separately for each population. Some SNPs did not have any LD data due to a lack of allele frequency data for those SNPs in the 1000 Genomes Project.

LD datasets containing D' and r^2 values for Finnish, Toscani and British populations are loaded in with pandas as separate dataframes. Each dataframe has their index set to the first column which contains SNP rsIDs.

some python code will go here...

This function uses the itertools combination function to create a list of tuples containing all unique pairs of SNPs possible from a list of SNPs. The list is then separated into two lists containing the first and second element of each tuple.

some python code will go here...

An empty dataframe is created to be filled with rows containing data from all six dataframes. This loop uses the two lists created from the SNP list to index each dataframe and extract the respective LD value. These are used to create a list which is converted into a single row pandas dataframe which is added to the empty dataframe using pandas concat until data for all relevant pairwise LD calculations have been added. The completed dataframe is then outputted as a TSV file.

some python code will go here...

Outputting LD results

When a user searches by gene name or chromosomal coordinates, if multiple SNPs are returned, a list of SNPs is used to filter the LD dataset for all rows with entries for all pairwise LD calculations of SNPs in the list and output a results dataframe.

Before filtering, the list is checked for any SNPs which are not in the LD dataset due to lack of LD data and any offending SNPs are removed from the list.

some python code will go here...

The SNP list is then used to create two lists containing the first and second element of each tuple using the `SNP_pair_lists()` function defined earlier.

some python code will go here...

The LD dataset containing all available data for pairwise LD calculations is loaded in with pandas and an empty dataframe is created for the filtered

data. The pair of SNP lists are then used to index the LD dataset dataframe for all rows with pairs of SNPs relevant to the user's search query which are added to the LD results dataframe using pandas concat.

some python code will go here...

LD heatmap plots

When a user searches by gene name or chromosomal coordinates, if multiple SNPs are returned, a list of SNPs is also used to extract LD values for all relevant pairwise SNP calculations to create a dataframe used to create a heatmap plot of LD values.

The LD dataset is loaded in with pandas and SNPs not present in the dataset are removed from the list of SNPs passed from the user query. The SNP list is then used to create two lists containing the first and second element of each tuple using the `SNP_pair_lists()` function defined earlier.

some python code will go here...

An empty dataframe is created to be filled with LD values used to create the LD plot. The pair of SNP lists are then used to index the LD dataset dataframe and extract the LD value for all possible pairwise LD calculations from the SNP list. A row of LD values is created for each SNP where each column corresponds with the pairwise LD calculation with one of the SNPs from the list. Each row is added to the empty dataframe using pandas concat.

some python code will go here...

The LD matrix dataframe is passed to the `ld_plot` function. The number of rows (n) is used to create a mask which will hide half of the heatmap to create a triangular plot. A coordinate matrix is also created to rotate the heatmap plot. The SNP list is used to create the axis labels located at the bottom of the plot. The function's title parameter passes a string which is used to determine the plot title.

some python code will go here...

Manhattan Plot

Flask

Navigation

References