

CNV GDRC Scripts

Local Environment

We recommend having four `conda` environments on your local machine to run our analysis. These environments are stored in the `envs/` directory. The following commands will create environments called `gdrc_data`, `gdrc_probprog`, `gdrc_snakemake`, and `gdrc_dx` using `conda`.

```
conda env create --file envs/data.yaml --name gdrc_data
conda env create --file envs/probprog.yaml --name gdrc_probprog
conda env create --file envs/snakemake.yaml --name gdrc_snakemake
conda env create --file envs/dx.yaml --name gdrc_dx
```

Analysis of Individual-Level Data

Analysis of the individual-level data requires access to the UK Biobank (UKB). Successful application to the data will provide access to the UKB [research analysis platform](#) (RAP).

Some scripts will need to be run on the RAP using the JupyterLab tool. The rest can be executed using the command line on your local computer. These scripts will spawn jobs on the RAP. Note that memory storage and computation on the RAP incurs expenses. We estimate that the total expenses to perform our analysis is less than £500.

The `dx` [command line utility](#) is useful for interacting with the UKB RAP. It can be used to upload/download files and execute jobs. It will be installed in the `gdrc_dx` environment but needs to be [configured by each user](#).

For scripts that are executed locally using the `dx` command line utility provided by DNAnexus, there is a `project` variable defined at the top of the script. This needs to be modified to use your UKB RAP project string. This will look something like this and can be found on the UKB RAP project information pane.

```
project="project-[24 character hash]"
```

Some of our scripts require software that is not available on the UKB RAP. We use Docker containers to package such software and execute it.

1. [FlashPCA2](#)

1. Use `docker build FlashPCA2/ -t 'flashpca2'` to build the container (on Apple silicon, the command is `docker buildx build --platform=linux/amd64 FlashPCA2 -t 'flashpca2'`).
2. Use `docker save flashpca2 | gzip > flashpca2.tar.gz` to save the image.
3. Upload the image to the UKB RAP at `/Software/`.

2. [VEP](#) with [loftee](#)

1. Upload data required by loftee to the UKB RAP at `/Annotations/LoFs/`.
2. Pull the VEP docker container using `docker pull ensemblorg/ensembl-vep`
3. Use `docker save ensembl-vep | gzip > loftee.tar.gz` to save the image.
4. Upload the image to the UKB RAP at `/Software/`.

Preparing for Association Analysis

Scripts should be executed in the following order.

Order	Script	Environment
1	<code>field_names.ipynb</code>	Local (<code>gdrc_data</code>)
2	<code>extract_covariates_and_phenotypes.ipynb</code>	UKB RAP
3	<code>extract_cnv_burden_genotypes.sh</code>	Local (<code>gdrc_dx</code>)*
4	<code>combine_cnv_burden_genotypes.sh</code>	Local (<code>gdrc_dx</code>)
5	<code>individuals_for_analysis.ipynb</code>	UKB RAP
6	<code>filter_exome_rare_variants_for PCs.sh</code>	Local (<code>gdrc_dx</code>)

Order	Script	Environment
7	<code>estimate_fine_structure_PCs.sh</code>	Local (<code>gdrc_dx</code>)
8	<code>filter_genotypes_for_wgr.sh</code>	Local (<code>gdrc_dx</code>)
9	<code>process_covariates_and_phenotypes.ipynb</code>	UKB RAP
10	<code>combine_nmr_phases.ipynb</code>	Local (<code>gdrc_dx</code>)*
11	<code>wes_allele_frequencies.sh</code>	Local (<code>gdrc_dx</code>)
12	<code>lof_loftee.sh</code>	Local (<code>gdrc_dx</code>)
13	<code>process_loftee_annotation.sh</code>	Local (<code>gdrc_dx</code>)*
14	<code>pick_most_severe_loftee_consequence.sh</code>	Local (<code>gdrc_dx</code>)*
15	<code>extract_cnv_burden_cov_matrices.sh</code>	Local (<code>gdrc_dx</code>)*
16	<code>create_annotation_files_lof.ipynb</code>	UKB RAP
17	<code>create_annotation_files_cnv.ipynb</code>	UKB RAP
18	<code>ecdf_cnv_burden_length.ipynb</code>	Local (<code>gdrc_data</code>)

*Specific instructions below the table.

`combine_nmr_phases.ipynb` : This requires early access to the Phase 3 data from Nightingale. Please contact Nightingale PLC directly for requests to the data. Phase 3 data will eventually join the Phase 1 and Phase 2 data in the UKB.

`extract_cnv_burden_genotypes.sh` :

1. Upload `extract_cnv_burden_genotypes.py` to `/Genotypes/CNVs/` in the UKB RAP.
2. Run `extract_cnv_burden_genotypes.sh` on the local machine in the `gdrc_dx` environment.

`process_loftee_annotation.sh` :

1. Upload `process_loftee_annotation.py` to `/Annotations/LoFs/` in the UKB RAP.
2. Run `process_loftee_annotation.sh` on the local machine in the `gdrc_dx` environment.

`pick_most_severe_loftee_consequence.sh` :

1. Upload `pick_most_severe_loftee_consequence.py` to `/Annotations/LoFs/` in the UKB RAP.
2. Run `pick_most_severe_loftee_consequence.sh` on the local machine in the `gdrc_dx` environment.

`extract_cnv_burden_cov_matrices.sh` :

1. Upload `extract_cnv_burden_cov_matrices.py` to `/Genotypes/CNVs/` in the UKB RAP.
2. Run `extract_cnv_burden_cov_matrices.sh` on the local machine in the `gdrc_dx` environment.

Association Analysis

Scripts should be executed in the following order.

Order	Script	Environment
1	<code>regenie_step_1_lofs.sh</code>	Local (<code>gdrc_data</code>)
2	<code>regenie_step_1_lofs_nmr.sh</code>	Local (<code>gdrc_data</code>)
3	<code>regenie_step_1_lofs_stratification.sh</code>	Local (<code>gdrc_data</code>)
4	<code>regenie_step_1_cnvs.sh</code>	Local (<code>gdrc_data</code>)
5	<code>regenie_step_1_cnvs_nmr.sh</code>	Local (<code>gdrc_data</code>)
6	<code>regenie_step_2_lofs.sh</code>	Local (<code>gdrc_data</code>)
7	<code>regenie_step_2_lofs_nmr.sh</code>	Local (<code>gdrc_data</code>)
8	<code>regenie_step_2_lofs_stratification.sh</code>	Local (<code>gdrc_data</code>)
9	<code>regenie_step_2_cnvs.sh</code>	Local (<code>gdrc_data</code>)
10	<code>regenie_step_2_cnvs_nmr.sh</code>	Local (<code>gdrc_data</code>)

Simulations

Scripts should be executed in the following order.

Order	Script	Environment
1	<code>extract_cnv_burden_cov_matrices_for_simulations.sh</code>	Local (<code>gdrc_data</code>)
2	<code>simulations_fixed.sh</code>	Local (<code>gdrc_data</code>)*
3	<code>simulations_random.sh</code>	Local (<code>gdrc_data</code>)*

*Specific instructions below the table.

`simulations_fixed.sh`:

1. Upload `simulations_fixed_1.ipynb` and `simulations_fixed_2.ipynb` to `/Simulations/` in the UKB RAP.
2. Run `simulations_fixed.sh` on the local machine in the `gdrc_dx` environment.

`simulations_random.sh`:

1. Upload `simulations_random_1.ipynb` and `simulations_random_2.ipynb` to `/Simulations/` in the UKB RAP.
2. Run `simulations_random.sh` on the local machine in the `gdrc_dx` environment.

Analysis of Summary Statistics

These analyses can be performed on summary statistics. You can either use data that you generate yourself using the scripts above, [our summary statistics](#), or any other burden summary statistics. If you are using other duplication burden summary statistics, you will need to obtain in-sample linkage disequilibrium (LD) estimates for the burden genotypes.

In the rest of the sections, we will assume that you are using our summary statistics. The data deposited in Zenodo is expected to be in `sum_stats/data/`. Here is an example of downloading some of the data.

```
mkdir -p sum_stats/data/
curl -L 'https://zenodo.org/records/13852455/files/LoF_Burden_Tests.tar.gz?download=1' \
  -o sum_stats/data/LoF_Burden_Tests.tar.gz
curl -L 'https://zenodo.org/records/13852455/files/CNV_Burden_Tests.tar.gz?download=1' \
  -o sum_stats/data/CNV_Burden_Tests.tar.gz
tar -xvf sum_stats/data/LoF_Burden_Tests.tar.gz -C sum_stats/data/
tar -xvf sum_stats/data/CNV_Burden_Tests.tar.gz -C sum_stats/data/
rm sum_stats/data/LoF_Burden_Tests.tar.gz
rm sum_stats/data/CNV_Burden_Tests.tar.gz
```

Use the following command to download s_{het} estimates from [Zeng et al. 2024](#).

```
curl -L 'https://zenodo.org/records/7939768/files/s_het_estimates.genebayes.tsv?download=1' \
  -o sum_stats/data/s_het_estimates.genebayes.tsv
```

Monotonicity

We use `snakemake` to distribute the workload across a cluster. We have provided the `Snakefile` to execute the workflow. You will have to configure the execution to your specific cluster or local machine.

A basic command to execute the workflow can be found in `run.sh`, which should be executed in the `gdrc_snakemake` environment. This command can be [configured extensively](#) to match your environment.

We also provide a notebook (`monotonicity.ipynb`) that can be executed using the `gdrc_probprog` environment to perform the full analysis for a single trait. This can be used to explore individual traits in depth.

RSS MASH

We use `snakemake` to distribute the workload across a cluster. We have provided the `Snakefile` to execute the workflow. You will have to configure the execution to your specific cluster or local machine.

A basic command to execute the workflow can be found in `run.sh`, which should be executed in the `gdrc_snakemake` environment. This command can be [configured extensively](#) to match your environment.

Figures

Jupyter notebooks to generate all figures and tables included in the manuscript can be found in the `sum_stats/` directory. The resulting figures and tables are included in the `results/` directory. The data that is required for each notebook is mentioned in the first cell.

Some notebooks use a Python kernel, while others use an R kernel. The scripts and their environments are described below. The following order of execution is recommended. For example, results from `figure_1_analysis.ipynb` are used in `figure_1.ipynb` to generate plots.

Order	Notebook	Environment
1	<code>figure_1_analysis.ipynb</code>	Local (<code>gdrc_data</code>), Python
2	<code>figure_1.ipynb</code>	Local (<code>gdrc_data</code>), R
3	<code>figure_2_analysis.ipynb</code>	Local (<code>gdrc_data</code>), Python
4	<code>figure_2.ipynb</code>	Local (<code>gdrc_data</code>), R
5	<code>figure_3_analysis.ipynb</code>	Local (<code>gdrc_probprog</code>), Python
6	<code>figure_3.ipynb</code>	Local (<code>gdrc_data</code>), R
7	<code>figure_4_analysis.ipynb</code>	Local (<code>gdrc_data</code>), Python
8	<code>figure_4.ipynb</code>	Local (<code>gdrc_data</code>), R
9	<code>figure_5_analysis.ipynb</code>	Local (<code>gdrc_data</code>), Python
10	<code>figure_5.ipynb</code>	Local (<code>gdrc_data</code>), R
11	<code>figure_apx_genome_wide_burden.ipynb</code>	Local (<code>gdrc_data</code>), R
12	<code>figure_apx_simulations.ipynb</code>	Local (<code>gdrc_data</code>), R
13	<code>figure_apx_confounding_analysis.ipynb</code>	Local (<code>gdrc_data</code>), Python
14	<code>figure_apx_confounding.ipynb</code>	Local (<code>gdrc_data</code>), R
15	<code>figure_apx_duplication_ld.ipynb</code>	Local (<code>gdrc_data</code>), R
16	<code>table_apx_burden_effects.ipynb</code>	Local (<code>gdrc_data</code>), R
17	<code>gwas_trumpet_plots.sh</code>	Shell
18	<code>figure_apx_gwas_trumpet_plots.ipynb</code>	Local (<code>gdrc_data</code>), R
19	<code>figure_apx_monotonicity_analysis.ipynb</code>	Local (<code>gdrc_probprog</code>), Python
20	<code>figure_apx_monotonicity.ipynb</code>	Local (<code>gdrc_data</code>), R
21	<code>figure_apx_genetic_correlation_analysis.ipynb</code>	Local (<code>gdrc_probprog</code>), Python
22	<code>figure_apx_genetic_correlation.ipynb</code>	Local (<code>gdrc_data</code>), R

Documentation

To generate the documentation, run the following locally using the `gdrc_data` environment.

```
cd doc/  
python compile_doc.py
```