

main

August 21, 2021

## 1 GO analysis using GOATOOLS

```
[1]: import functools
import pandas as pd
import collections as cx
from pybiomart import Dataset
# GO analysis
from goatools.base import download_go_basic_obo
from goatools.base import download_ncbi_associations
from goatools.obo_parser import GODag
from goatools.anno.genetogo_reader import Gene2GoReader
from goatools.goea.go_enrichment_ns import GOEnrichmentStudyNS
```

```
[2]: @functools.lru_cache()
def get_database():
    dataset = Dataset(name="hsapiens_gene_ensembl",
                      host="http://www.ensembl.org",
                      use_cache=True)
    db = dataset.query(attributes=["ensembl_gene_id",
                                  "external_gene_name",
                                  "entrezgene_id"],
                       use_attr_names=True).dropna(subset=['entrezgene_id'])
    return db

@functools.lru_cache()
def get_deg():
    fn = '../_m/genes/diffExpr_EAvsAA_FDR05.txt'
    return pd.read_csv(fn, sep='\t')

@functools.lru_cache()
def convert2entrez():
    df = get_deg()
    if 'EntrezID' in df.columns:
        return df.rename(columns={'EntrezID': 'entrezgene_id'})
    else:
```

```

        return df.merge(get_database(), left_on='ensemblID',
                        right_on='ensembl_gene_id')

@functools.lru_cache()
def get_upregulated():
    df = convert2entrez()
    return df.loc[(df['t'] > 0)]

@functools.lru_cache()
def get_downregulated():
    df = convert2entrez()
    return df.loc[(df['t'] < 0)]

```

```

[3]: def obo_annotation(alpha=0.05):
    # database annotation
    fn_obo = download_go_basic_obo()
    fn_gene2go = download_ncbi_associations() # must be gunzip to work
    obodag = GODag(fn_obo) # downloads most up-to-date
    anno_hs = Gene2GoReader(fn_gene2go, taxids=[9606])
    # get associations
    ns2assoc = anno_hs.get_ns2assoc()
    for nspc, id2gos in ns2assoc.items():
        print("{NS} {N:,} annotated human genes".format(NS=nspc, N=len(id2gos)))
    goeaobj = GOEnrichmentStudyNS(
        get_database()['entrezgene_id'], # List of human genes with entrez IDs
        ns2assoc, # geneid/GO associations
        obodag, # Ontologies
        propagate_counts = False,
        alpha = alpha, # default significance cut-off
        methods = ['fdr_bh'])
    return goeaobj

def run_goea(direction):
    if direction == "Up":
        df = get_upregulated()
    elif direction == "Down":
        df = get_downregulated()
    else:
        df = convert2entrez()
    geneids_study = {z[0]:z[1] for z in zip(df['entrezgene_id'], df['Symbol'])}
    goeaobj = obo_annotation()
    goea_results_all = goeaobj.run_study(geneids_study)
    goea_results_sig = [r for r in goea_results_all if r.p_fdr_bh < 0.05]

```

```

ctr = cx.Counter([r.NS for r in goea_results_sig])
print('Significant results[{TOTAL}] = {BP} BP + {MF} MF + {CC} CC'.format(
    TOTAL=len(goea_results_sig),
    BP=ctr['BP'], # biological_process
    MF=ctr['MF'], # molecular_function
    CC=ctr['CC'])) # cellular_component

if direction == "Up":
    label = "upregulated"
elif direction == "Down":
    label = "downregulated"
else:
    label = "allDEG"
goeobj.wr_xlsx("GO_analysis_%s.xlsx" % label, goea_results_sig)
goeobj.wr_txt("GO_analysis_%s.txt" % label, goea_results_sig)

```

## 1.1 Gene ontology

```

[4]: for direction in ["All", "Up", "Down"]:
    run_goea(direction)

```

```

requests.get(http://purl.obolibrary.org/obo/go/go-basic.obo, stream=True)
WROTE: go-basic.obo

```

```

FTP RETR ftp.ncbi.nlm.nih.gov gene/DATA gene2go.gz -> gene2go.gz
gunzip gene2go.gz
go-basic.obo: fmt(1.2) rel(2021-08-18) 47,217 GO Terms
HMS:0:00:04.253571 330,313 annotations, 20,685 genes, 18,684 GOs, 1 taxids READ:
gene2go
CC 19,422 annotated human genes
BP 18,505 annotated human genes
MF 18,190 annotated human genes

```

```

Load BP Gene Ontology Analysis ...
70% 20,236 of 29,107 population items found in association

```

```

Load CC Gene Ontology Analysis ...
74% 21,428 of 29,107 population items found in association

```

```

Load MF Gene Ontology Analysis ...
70% 20,354 of 29,107 population items found in association

```

```

Run BP Gene Ontology Analysis: current study set of 2759 IDs ... 81% 1,719 of
2,111 study items found in association
77% 2,111 of 2,759 study items found in population(29107)
Calculating 12,429 uncorrected p-values using fisher
12,429 GO terms are associated with 17,848 of 29,107 population items
4,694 GO terms are associated with 1,719 of 2,759 study items

```

```

METHOD fdr_bh:
    11 GO terms found significant (< 0.05=alpha) ( 8 enriched + 3
purified): statsmodels fdr_bh
    159 study items associated with significant GO IDs (enriched)
    37 study items associated with significant GO IDs (purified)

Run CC Gene Ontology Analysis: current study set of 2759 IDs ... 86% 1,820 of
2,111 study items found in association
    77% 2,111 of 2,759 study items found in population(29107)
Calculating 1,753 uncorrected p-values using fisher
    1,753 GO terms are associated with 18,711 of 29,107 population items
    838 GO terms are associated with 1,820 of 2,759 study items
METHOD fdr_bh:
    39 GO terms found significant (< 0.05=alpha) ( 39 enriched + 0
purified): statsmodels fdr_bh
    1,696 study items associated with significant GO IDs (enriched)
    0 study items associated with significant GO IDs (purified)

Run MF Gene Ontology Analysis: current study set of 2759 IDs ... 84% 1,768 of
2,111 study items found in association
    77% 2,111 of 2,759 study items found in population(29107)
Calculating 4,420 uncorrected p-values using fisher
    4,420 GO terms are associated with 17,838 of 29,107 population items
    1,596 GO terms are associated with 1,768 of 2,759 study items
METHOD fdr_bh:
    15 GO terms found significant (< 0.05=alpha) ( 13 enriched + 2
purified): statsmodels fdr_bh
    1,452 study items associated with significant GO IDs (enriched)
    3 study items associated with significant GO IDs (purified)
Significant results[65] = 11 BP + 15 MF + 39 CC
    65 items Wrote: GO_analysis_allDEG.xlsx
    65 GOEA results for 1826 study items. Wrote: GO_analysis_allDEG.txt
EXISTS: go-basic.obo
EXISTS: gene2go
go-basic.obo: fmt(1.2) rel(2021-08-18) 47,217 GO Terms
HMS:0:00:04.541304 330,313 annotations, 20,685 genes, 18,684 GOs, 1 taxids READ:
gene2go
CC 19,422 annotated human genes
BP 18,505 annotated human genes
MF 18,190 annotated human genes

Load BP Gene Ontology Analysis ...
    70% 20,236 of 29,107 population items found in association

Load CC Gene Ontology Analysis ...
    74% 21,428 of 29,107 population items found in association

Load MF Gene Ontology Analysis ...

```

70% 20,354 of 29,107 population items found in association

Run BP Gene Ontology Analysis: current study set of 1337 IDs ... 79% 760 of 967 study items found in association

72% 967 of 1,337 study items found in population(29107)

Calculating 12,429 uncorrected p-values using fisher

12,429 GO terms are associated with 17,848 of 29,107 population items

2,805 GO terms are associated with 760 of 1,337 study items

METHOD fdr\_bh:

5 GO terms found significant (< 0.05=alpha) ( 4 enriched + 1 purified): statsmodels fdr\_bh

48 study items associated with significant GO IDs (enriched)

2 study items associated with significant GO IDs (purified)

Run CC Gene Ontology Analysis: current study set of 1337 IDs ... 84% 813 of 967 study items found in association

72% 967 of 1,337 study items found in population(29107)

Calculating 1,753 uncorrected p-values using fisher

1,753 GO terms are associated with 18,711 of 29,107 population items

530 GO terms are associated with 813 of 1,337 study items

METHOD fdr\_bh:

16 GO terms found significant (< 0.05=alpha) ( 16 enriched + 0 purified): statsmodels fdr\_bh

602 study items associated with significant GO IDs (enriched)

0 study items associated with significant GO IDs (purified)

Run MF Gene Ontology Analysis: current study set of 1337 IDs ... 82% 793 of 967 study items found in association

72% 967 of 1,337 study items found in population(29107)

Calculating 4,420 uncorrected p-values using fisher

4,420 GO terms are associated with 17,838 of 29,107 population items

942 GO terms are associated with 793 of 1,337 study items

METHOD fdr\_bh:

8 GO terms found significant (< 0.05=alpha) ( 7 enriched + 1 purified): statsmodels fdr\_bh

630 study items associated with significant GO IDs (enriched)

1 study items associated with significant GO IDs (purified)

Significant results[29] = 5 BP + 8 MF + 16 CC

29 items Wrote: GO\_analysis\_upregulated.xlsx

29 GOEA results for 770 study items. Wrote: GO\_analysis\_upregulated.txt

EXISTS: go-basic.obo

EXISTS: gene2go

go-basic.obo: fmt(1.2) rel(2021-08-18) 47,217 GO Terms

HMS:0:00:04.466198 330,313 annotations, 20,685 genes, 18,684 GOs, 1 taxids READ: gene2go

CC 19,422 annotated human genes

BP 18,505 annotated human genes

MF 18,190 annotated human genes

```

Load BP Gene Ontology Analysis ...
  70% 20,236 of 29,107 population items found in association

Load CC Gene Ontology Analysis ...
  74% 21,428 of 29,107 population items found in association

Load MF Gene Ontology Analysis ...
  70% 20,354 of 29,107 population items found in association

Run BP Gene Ontology Analysis: current study set of 1423 IDs ... 84%    960 of
1,145 study items found in association
  80% 1,145 of 1,423 study items found in population(29107)
Calculating 12,429 uncorrected p-values using fisher
  12,429 GO terms are associated with 17,848 of 29,107 population items
  3,212 GO terms are associated with    960 of 1,423 study items
METHOD fdr_bh:
  10 GO terms found significant (< 0.05=alpha) ( 7 enriched + 3
purified): statsmodels fdr_bh
  81 study items associated with significant GO IDs (enriched)
  15 study items associated with significant GO IDs (purified)

Run CC Gene Ontology Analysis: current study set of 1423 IDs ... 88% 1,008 of
1,145 study items found in association
  80% 1,145 of 1,423 study items found in population(29107)
Calculating 1,753 uncorrected p-values using fisher
  1,753 GO terms are associated with 18,711 of 29,107 population items
  657 GO terms are associated with 1,008 of 1,423 study items
METHOD fdr_bh:
  28 GO terms found significant (< 0.05=alpha) ( 28 enriched + 0
purified): statsmodels fdr_bh
  910 study items associated with significant GO IDs (enriched)
  0 study items associated with significant GO IDs (purified)

Run MF Gene Ontology Analysis: current study set of 1423 IDs ... 85%    976 of
1,145 study items found in association
  80% 1,145 of 1,423 study items found in population(29107)
Calculating 4,420 uncorrected p-values using fisher
  4,420 GO terms are associated with 17,838 of 29,107 population items
  1,101 GO terms are associated with    976 of 1,423 study items
METHOD fdr_bh:
  14 GO terms found significant (< 0.05=alpha) ( 13 enriched + 1
purified): statsmodels fdr_bh
  768 study items associated with significant GO IDs (enriched)
  0 study items associated with significant GO IDs (purified)
Significant results[52] = 10 BP + 14 MF + 28 CC
  52 items WROTE: GO_analysis_downregulated.xlsx
  52 GOEA results for 1007 study items. WROTE: GO_analysis_downregulated.txt

```

[ ]: