

```
In [1]: import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import scipy.stats as st
from scipy.stats import ttest_ind
import os
%matplotlib inline

os.chdir("/Users/yolandatiao/GSuite Scripps/MLC_Screen")
import fit_nbinom
```

0. Define Screen Data Class and Transformations

```

In [3]: class screenData():
    def __init__(self, raw_data, col_data, row_data):
        self.raw_df = pd.read_csv(raw_data, index_col="name")
        self.col_df = pd.read_csv(col_data)
        self.row_df = pd.read_csv(row_data, index_col='name')
        self.qcflag = False
        conditions = list(set(list(self.col_df['condition'])))
        self.cond_dict = {}
        for cond in conditions:
            self.cond_dict[cond] = [str(x) for x in list(self.col_df[self.c
raw_cond_df = pd.DataFrame({'name': list(self.raw_df.index.values)})
        for i in self.cond_dict.keys():
            raw_cond_df[i] = list(self.raw_df[self.cond_dict[i]].sum(axis=1
raw_cond_df = raw_cond_df.set_index('name')
        self.raw_df = raw_cond_df

    def qc(self, input_name, count_cutoff, input_vs_output_cutoff):
        self.raw_df['target'] = list(self.row_df['target'])
        self.raw_df['sum'] = list(self.raw_df.sum(axis=1))
        self.raw_df['out_sum'] = list(self.raw_df['sum'] - self.raw_df[inpu
        self.raw_df['input_vs_output'] = list(self.raw_df[input_name] / sel
        self.raw_df = self.raw_df[self.raw_df['sum'] > count_cutoff]
        self.raw_df = self.raw_df[self.raw_df['input_vs_output'] < input_vs
        # update row metadata
        self.row_df = pd.DataFrame({'name': list(self.raw_df.index.values)},
        self.row_df.set_index('name')

        del self.raw_df['sum']
        del self.raw_df['out_sum']
        del self.raw_df['input_vs_output']
        del self.raw_df['target']
        self.qcflag = True

    def norm_count(self):
        if (not self.qcflag):
            print("Proceed without quality control...")
            print("Perform normalization for each sample...")
        self.norm_df = self.raw_df / self.raw_df.sum() * 1000000
        self.norm_df = self.norm_df.astype('int32')
        self.all_normalized_counts = []
        for col in self.norm_df.columns:
            self.all_normalized_counts += list(self.norm_df[col])
        self.all_normalized_counts.sort()

    def nb_pctl(self):
        if not hasattr(self, "norm_df"):
            self.norm_count()
        self.nbinom_parameters = fit_nbinom.fit_nbinom(np.array(self.all_no
        self.nbpctl_df = pd.DataFrame(st.nbinom.cdf(self.norm_df, self.nbin
        self.nbpctl_df.columns = list(self.raw_df.columns)
        self.nbpctl_df.index = list(self.raw_df.index.values)

    def nb_pctl_shift(self):
        if not hasattr(self, "nbpctl_df"):
            self.nb_pctl()
        self.nbpctl_shift_df = pd.DataFrame({"name": list(self.nbpctl_df.in

```

```

conditions = list(self.nbpctl_df.columns)
for i in range(0, len(conditions)):
    for j in range((i+1), len(conditions)):
        i_name = conditions[i]
        j_name = conditions[j]
        comp = "__vs__".join([i_name, j_name])
        self.nbpctl_shift_df[comp] = list(self.nbpctl_df[i_name] -
self.nbpctl_shift_df = self.nbpctl_shift_df.set_index('name')

def target_eff(self):
    if not hasattr(self, "nbpctl_shift"):
        self.nb_pctl_shift()

self.nbpctl_shift_df['target'] = list(self.row_df['target'])
self.target_shift_df = self.nbpctl_shift_df.groupby('target').mean()

self.target_shift_z_df = pd.DataFrame({'target': list(self.target_s
for col in self.target_shift_df.columns:
    self.target_shift_z_df[col] = st.zscore(np.array(self.target_sh
self.target_shift_z_df = self.target_shift_z_df.set_index('target')

del self.nbpctl_shift_df['target']

def target_ttest(self):
    if not hasattr(self, "nbpctl_shift_df"):
        self.nb_pctl_shift()
    comparisons = list(self.nbpctl_shift_df.columns)
    targets = list(set(list(self.row_df['target'])))
    self.target_shift_p_df = pd.DataFrame({'target': targets})
    self.target_shift_st_df = pd.DataFrame({'target': targets})
    self.nbpctl_shift_df['target'] = list(self.row_df['target'])
    for cp in comparisons:
        cp_p = []
        cp_st = []
        cp_nums = list(self.nbpctl_shift_df[cp])
        for tg in targets:
            tg_nums = list(self.nbpctl_shift_df[self.nbpctl_shift_df['t
            cp_tg_tt = st.ttest_ind(tg_nums, cp_nums)
            cp_p.append(cp_tg_tt.pvalue)
            cp_st.append(cp_tg_tt.statistic)
        self.target_shift_p_df[cp] = cp_p
        self.target_shift_st_df[cp] = cp_st
    self.target_shift_p_df = self.target_shift_p_df.set_index('target')
    self.target_shift_st_df = self.target_shift_st_df.set_index('target')
    self.target_shift_p_df = self.target_shift_p_df.fillna(1)
    self.target_shift_st_df = self.target_shift_st_df.fillna(0)
    del self.nbpctl_shift_df['target']

def adj_z(self):
    if not hasattr(self, 'target_shift_p_df'):
        self.target_ttest()
    if not hasattr(self, 'target_shift_z_df'):
        self.target_eff()
    self.target_shift_z_df = self.target_shift_z_df.sort_index()
    self.target_shift_p_df = self.target_shift_p_df.sort_index()
    self.target_shift_z_adj_df = pd.DataFrame({'target': list(self.targ
    comparisons = list(self.target_shift_z_df.columns)

```

```

for cp in comparisons:
    cp_z = list(self.target_shift_z_df[cp])
    cp_p = [(x + 0.01) for x in list(self.target_shift_p_df[cp])]
    cp_z_adj = cp_z / (np.sqrt(cp_p))
    self.target_shift_z_adj_df[cp] = cp_z_adj
self.target_shift_z_adj_df = self.target_shift_z_adj_df.set_index('

def to_csv(self, out_dir):
    not_write = ['col_df', 'row_df']
    items = list(vars(mlc_data).keys())
    items = [x for x in items if ('df' in x and x not in not_write)]
    if (not os.path.exists(out_dir)):
        os.mkdir(out_dir)
    for i in items:
        i_name = "%s/%s.csv"%(out_dir, i)
        i_name = i_name.replace("_df", "")
        vars(mlc_data)[i].to_csv(i_name)

```

1. Read data and apply transformations

```

In [4]: mlc_raw = "/Users/yolandatiao/GSuite Scripps/MLC_Screen/input/sundrud_repor
mlc_col = "/Users/yolandatiao/GSuite Scripps/MLC_Screen/input/sundrud_repor
mlc_row = "/Users/yolandatiao/GSuite Scripps/MLC_Screen/input/sundrud_repor

# Read input data
mlc_data = screenData(mlc_raw, mlc_col, mlc_row)

# Apply quality control, filter: low reads and drop outs
mlc_data.qc('pre_inj', 100, 3)

# Calculate adjusted z score
mlc_data.adj_z()

# Write output
mlc_data.to_csv('mlc_screen')

```

Perform normalization for each sample...

```

/Users/yolandatiao/anaconda3/envs/stats/lib/python3.8/site-packages/scipy/optimize/optimize.py:697: RuntimeWarning: invalid value encountered in double_scalars
    df = (f*((xk + d,) + args)) - f0) / d[k]
/Users/yolandatiao/anaconda3/envs/stats/lib/python3.8/site-packages/numpy/core/fromnumeric.py:3583: RuntimeWarning: Degrees of freedom <= 0 for s
lice
    return _methods._var(a, axis=axis, dtype=dtype, out=out, ddof=ddof,
/Users/yolandatiao/anaconda3/envs/stats/lib/python3.8/site-packages/numpy/core/_methods.py:209: RuntimeWarning: invalid value encountered in double_scalars
    ret = ret.dtype.type(ret / rcount)

```

2. Check output

- Focus on small intestine Mdr1 positive (**si_Mp**) versus small intestine Mdr1 negative (**si_Mn**)

```
In [10]: # Select and sort data
shift_sorted_df = pd.DataFrame(mlc_data.target_shift_df.sort_values(by='si_
shift_p_df = pd.DataFrame(mlc_data.target_shift_p_df.sort_values(by='si_Mp_
shift_st_df = pd.DataFrame(mlc_data.target_shift_st_df.sort_values(by='si_M
shift_z_sorted_df = pd.DataFrame(mlc_data.target_shift_z_df.sort_values(by=
shift_z_adj_sorted_df = pd.DataFrame(mlc_data.target_shift_z_adj_df.sort_va
```

2.1 Adjusted Z Score (of percentile shift)

```
In [32]: shift_z_adj_sorted_df.head(n=10)
```

Out[32]:

si_Mp_vs_si_Mn	
target	
Abcb1a	-13.891919
Nr2f6	-11.480406
Scand1	-8.995634
Thra	-5.575579
Nr1i3	-4.515784
Ppargc1b	-4.167345
Hnf4g	-3.553863
Esr1	-2.414341
Nr2f1	-2.267800
Abcb1b	-2.213765

```
In [31]: shift_z_adj_sorted_df.tail(n=10)
```

Out[31]:

si_Mp_vs_si_Mn	
target	
Thrb	1.286672
Nr4a1	1.745649
Ar	2.496645
Nr5a2	2.945634
Hnf4a	3.324292
Esr2	3.714146
Nr2e3	4.009031
Nr1d1	6.168899
Rxra	16.148703
Ncoa1	19.526014

2.1 Z Score (of percentile shift)

```
In [30]: shift_z_sorted_df.head(n=10)
```

Out[30]:

si_Mp_vs_si_Mn	
target	
Abcb1a	-2.701878
Scand1	-2.062608
Nr2f6	-1.882859
Thra	-1.726952
Ppargc1b	-1.540401
Nr1i3	-1.444528
Hnf4g	-1.317952
Nr2f1	-1.258917
Abcb1b	-1.239469
Esr1	-1.188065

```
In [29]: shift_z_sorted_df.tail(n=10)
```

Out[29]:

si_Mp_vs_si_Mn	
target	
Thrb	0.843061
Nr4a1	0.954861
Hnf4a	1.227635
Ar	1.231297
Nr2e3	1.333887
Nr5a2	1.335143
Esr2	1.483393
Nr1d1	1.817280
Ncoa1	2.327917
Rxra	2.890425

2.2 Percentile shift

```
In [35]: shift_st_df.head(n=10)
```

```
Out[35]:
```

si_Mp_vs_si_Mn	
target	
Nr2f6	-2.405790
Abcb1a	-2.213397
Scand1	-2.038860
Thra	-1.724421
Nr1i3	-1.690060
Ppargc1b	-1.532914
Hnf4g	-1.529252
Esr1	-1.197908
Nr5a1	-1.172956
Rorb	-1.106841

```
In [36]: shift_st_df.tail(n=10)
```

```
Out[36]:
```

si_Mp_vs_si_Mn	
target	
Ppard	0.943491
Nr4a1	1.062240
Ar	1.195139
Nr5a2	1.298300
Esr2	1.445950
Hnf4a	1.533901
Nr2e3	1.647823
Nr1d1	1.777419
Rxra	2.304973
Ncoa1	2.889457

2.4 By target t-test p-value (of percentile shift)

```
In [38]: shift_p_df.head(n=10)
```

Out[38]:

si_Mp_vs_si_Mn	
target	
Ncoa1	0.004214
Nr2f6	0.016898
Rxra	0.022037
Abcb1a	0.027827
Scand1	0.042574
Nr1d1	0.076782
Thra	0.085936
Nr1i3	0.092326
Nr2e3	0.100703
Hnf4a	0.126377

```
In [39]: shift_p_df.tail(n=10)
```

Out[39]:

si_Mp_vs_si_Mn	
target	
Nr1h3	0.898327
Ncoa5	0.901452
Rxrg	0.912657
Nr1i2	0.949010
Esrrg	0.949011
Vdr	0.949358
Rara	0.952286
Ahr	0.976446
Nr2c2	0.991820
Blimp	1.000000

2.4 By target t-test statistics (of percentile shift)


```
In [40]: shift_st_df.head(n=10)
```

Out[40]:

si_Mp_vs_si_Mn	
target	
Nr2f6	-2.405790
Abcb1a	-2.213397
Scand1	-2.038860
Thra	-1.724421
Nr1i3	-1.690060
Ppargc1b	-1.532914
Hnf4g	-1.529252
Esr1	-1.197908
Nr5a1	-1.172956
Rorb	-1.106841

```
In [41]: shift_st_df.tail(n=10)
```

Out[41]:

si_Mp_vs_si_Mn	
target	
Ppard	0.943491
Nr4a1	1.062240
Ar	1.195139
Nr5a2	1.298300
Esr2	1.445950
Hnf4a	1.533901
Nr2e3	1.647823
Nr1d1	1.777419
Rxra	2.304973
Ncoa1	2.889457

```
In [ ]:
```

