

RNA-seq Report

Nathan Hughes (JIC)

May 16, 2019

Contents

1	Data setup	3
1.1	Helper funcs for pprinting	3
1.2	Load up counts and DE	3
2	Inspect Samples	4
2.1	Creating a distance map of samples using normalised counts	4
2.1.1	Samples separated	4
2.1.2	Samples together	5
3	Simple Analysis	6
3.1	Largest/Lowest expression sum	6
3.2	PCA on count data	7

1 Data setup

1.1 Helper funcs for pprinting

```

1 import tabulate
2 import IPython
3
4 class OrgFormatter(IPython.core.formatters.BaseFormatter):
5     format_type = IPython.core.formatters.Unicode('text/org')
6     print_method = IPython.core.formatters.ObjectName('_repr_org_')
7
8 def pd_dataframe_to_org(df):
9     return tabulate.tabulate(df, headers='keys', tablefmt='orgtbl', showindex='always')
10
11 ip = get_ipython()
12 ip.display_formatter.formatters['text/org'] = OrgFormatter()
13
14 f = ip.display_formatter.formatters['text/org']
15 f.for_type_by_name('pandas.core.frame', 'DataFrame', pd_dataframe_to_org)
16
17 print('Lets go!')
```

Lets go!

1.2 Load up counts and DE

```

1 import pandas as pd
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 counts = pd.read_csv(
6     "/Users/hughesn/Transcripts/RNA-Seq/Analysis/Data/norml_count_data.csv",
7     index_col=0)
8 xl = pd.ExcelFile(
9     "/Users/hughesn/Transcripts/RNA-Seq/Analysis/Data/diff_from_col0:False_onlyDiff:False.xlsx")
10 sheet_names = xl.sheet_names
11 dfs = []
12 for s in sheet_names:
13     d = xl.parse(s)
14     d['sample'] = s.split("|")[0].replace(" ", "")
15     dfs.append(d)
16
17 DE = pd.concat(dfs)
18 DE = DE.rename_axis('gene').sort_values(by=['gene', 'log2FoldChange'],
19                                         ascending=[False, False])
20 print("Loaded data")
```

Loaded data

2 Inspect Samples

2.1 Creating a distance map of samples using normalised counts

2.1.1 Samples separated

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 from scipy.spatial.distance import pdist, squareform
4
5 distances = pdist(counts.T.values, metric='euclidean')
6 dist_matrix = squareform(distances)
7 dist_df = pd.DataFrame(dist_matrix, columns = counts.columns, index=counts.columns)
8
9 sns.clustermap(dist_df)

```

<Figure size 720x720 with 4 Axes>

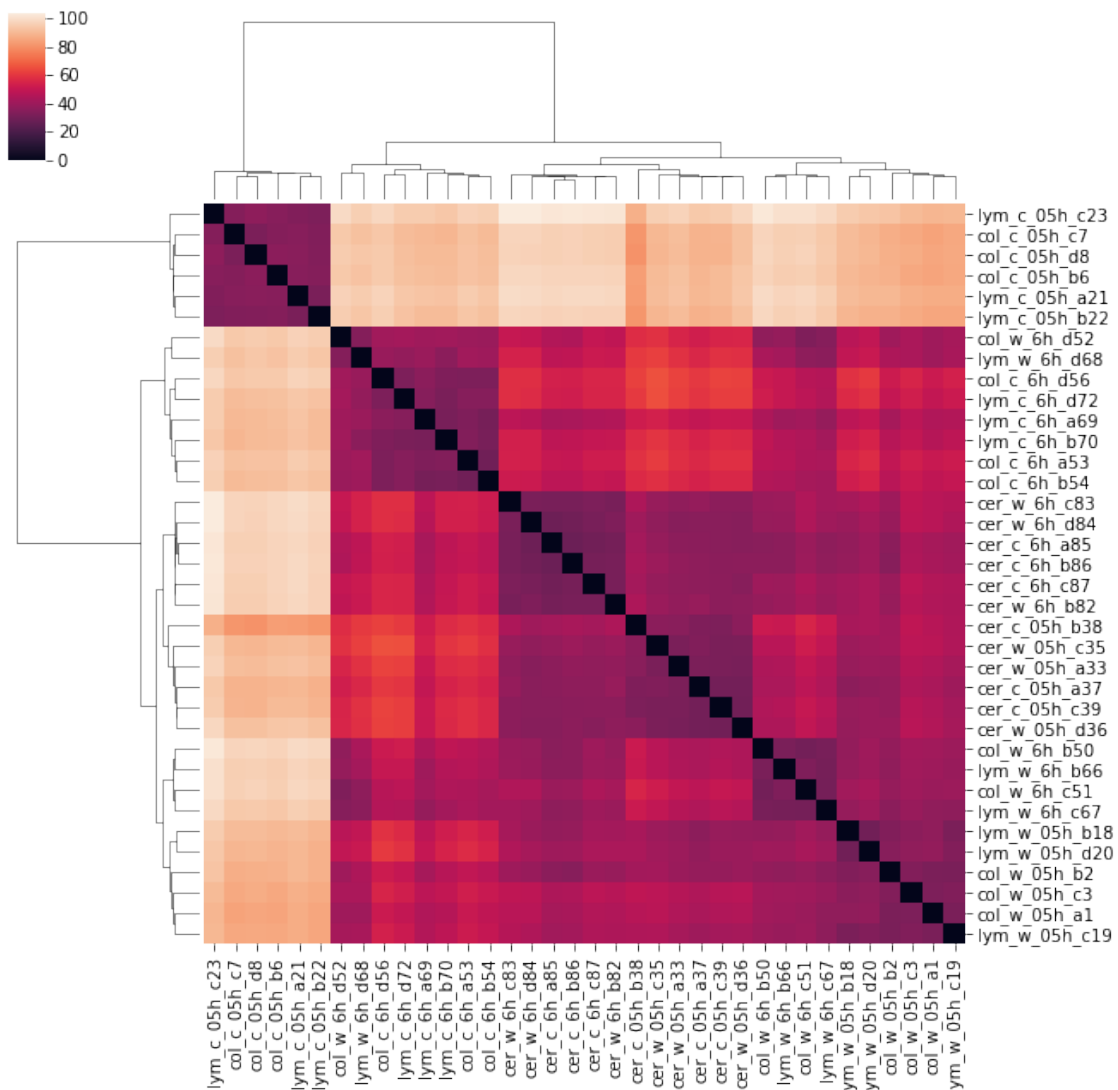


Figure 1: Distance map between samples

2.1.2 Samples together

```

1 def collapse_counts(counts):
2     u_cols = list(set([l.rsplit("_", 1)[0] for l in list(counts.columns)]))
3     cols = list(counts.columns)
4     ss = []
5     for uc in u_cols:
6         cs = [c for c in cols if c.startswith(uc)]
7         ss.append(counts[cs].sum(axis=1).rename(uc))
8     dc = pd.concat(ss, axis=1)
9     return dc
10 collapsed_counts = collapse_counts(counts)
11 distances = pdist(collapsed_counts.T.values, metric='euclidean')
12 dist_matrix = squareform(distances)
13 dist_df = pd.DataFrame(dist_matrix, columns = collapsed_counts.columns, index=collapsed_counts.columns)
14 sns.clustermap(dist_df)

```

<Figure size 720x720 with 4 Axes>

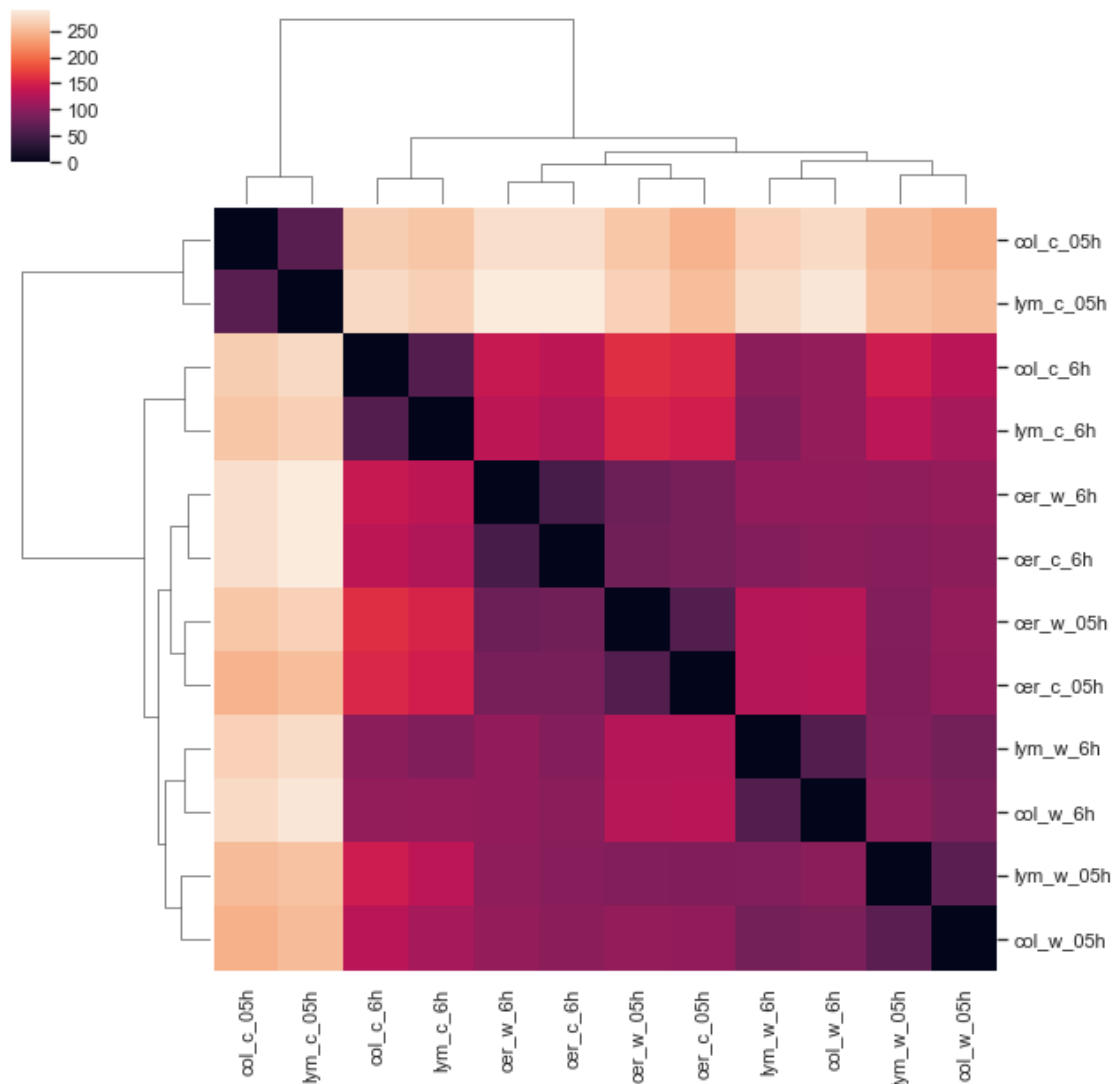


Figure 2: Distance map between samples, pooled together

3 Simple Analysis

3.1 Largest/Lowest expression sum

```

1 #DE.sum(axis=1).sort_values(by=['log2FoldChange'], ascending=[False]).head(3)
2
3 locs = DE[['log2FoldChange']].groupby(['gene']).sum().sort_values(by='log2FoldChange', ascending=False).head(20).index.
4 top = DE.loc[locs]
5 top = top.pivot(columns='sample', values='log2FoldChange')
6
7 locs = DE[['log2FoldChange']].groupby(['gene']).sum().sort_values(by='log2FoldChange', ascending=True).head(20).index.
8 bot = DE.loc[locs]
9 bot = bot.pivot(columns='sample', values='log2FoldChange')
10
11 both = pd.concat([top,bot])
12 both['col_w_05h'] = 0
13
14 sns.clustermap(both, cmap='bwr')

```

<Figure size 720x720 with 4 Axes>

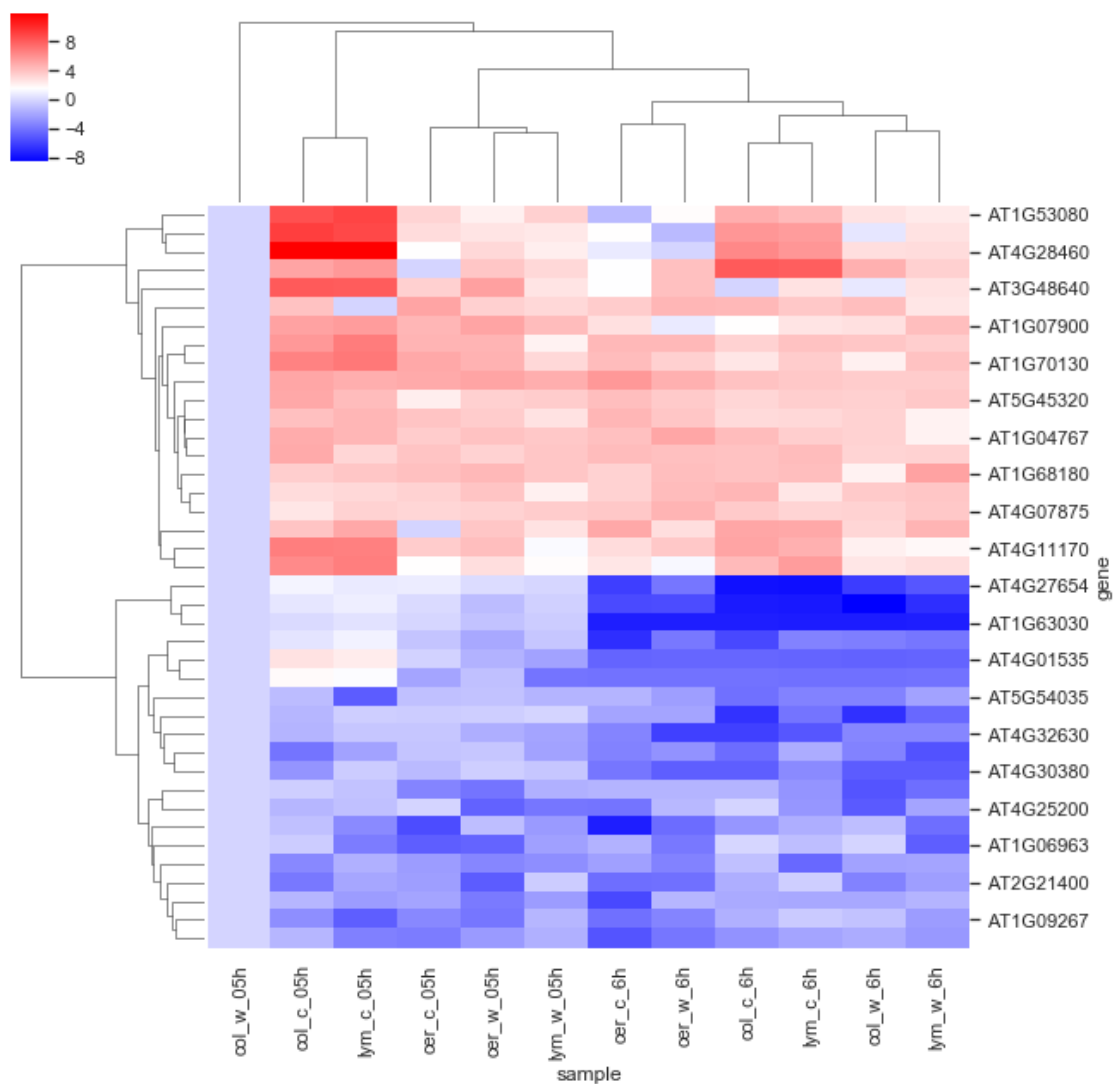


Figure 3: Largest and least DE genes

3.2 PCA on count data

```

1  from sklearn.decomposition import PCA
2  from sklearn.preprocessing import StandardScaler
3  sns.set()
4
5  cols = list(counts.columns)
6
7  counts_genotype = [c.split("_")[0] for c in cols]
8  counts_treatment = [c.split("_")[1] for c in cols]
9  counts_time = [c.split("_")[2] for c in cols]
10
11 x = StandardScaler().fit_transform(counts.T.values)
12
13 pca = PCA(n_components=2)
14 principalComponents = pca.fit_transform(x)
15 principalDf = pd.DataFrame(data=principalComponents, columns=[
16     'principal component 1', 'principal component 2'])
17
18 principalDf['genotype'] = counts_genotype
19 principalDf['treatment'] = counts_treatment
20 principalDf['time'] = counts_time
21
22 g = sns.FacetGrid(principalDf, col='time', row='genotype', hue='treatment')
23
24 g = g.map(plt.scatter, 'principal component 1',
25     'principal component 2').add_legend()
26
27 print("Explained variance from PC1 & 2 respectively:")
28 print(pca.explained_variance_ratio_)

```

Explained variance from PC1 & 2 respectively: [0.21077632 0.14325373]

<Figure size 483.925x648 with 6 Axes>

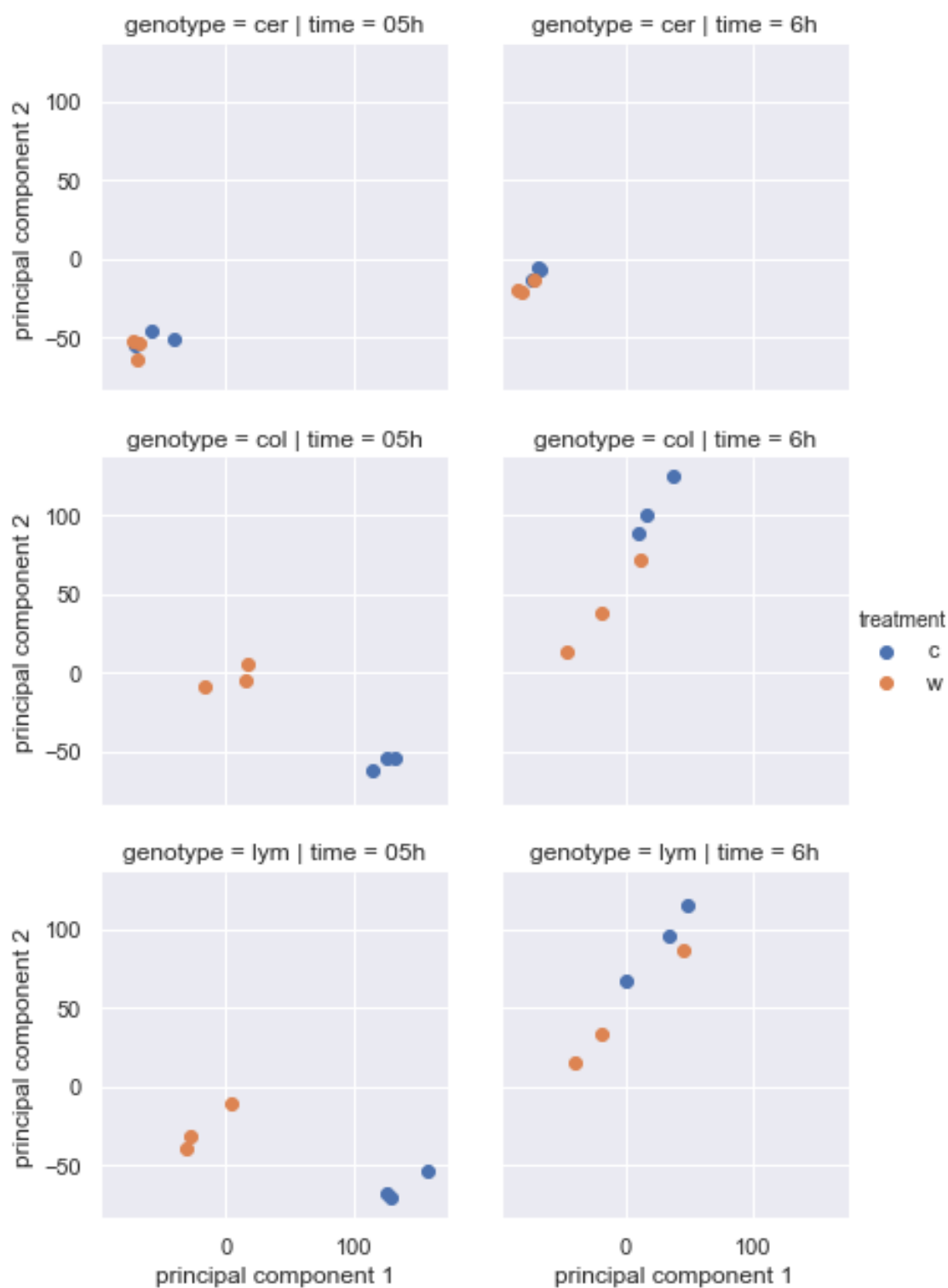


Figure 4: PCA of sample counts