

```
In [1]: import scanpy as sc
import tiledb
import numpy as np
from sklearn.metrics import adjusted_mutual_info_score, adjusted_rand_score
```

```
In [2]: # load anndata (tabula sapiens - only blood, UBERON:0000178)
# https://cellxgene.cziscience.com/collections/e5f58829-1a66-40b5-a624-9046778e74f.
adata = sc.read_h5ad('blood_sapiens.h5ad')

# don't worry about batch effects by selecting one donor & assay
adata = adata[np.logical_and(adata.obs['donor_id']=='TSP7',
                             adata.obs['assay_ontology_term_id']=='EFO:0009922')]

# reset data to raw counts to preprocess from scratch
# using standard workflow
adata.X = adata.layers['decontXcounts']
sc.pp.filter_genes(adata, min_cells=3)
sc.pp.normalize_total(adata, target_sum=1e4)
sc.pp.log1p(adata)

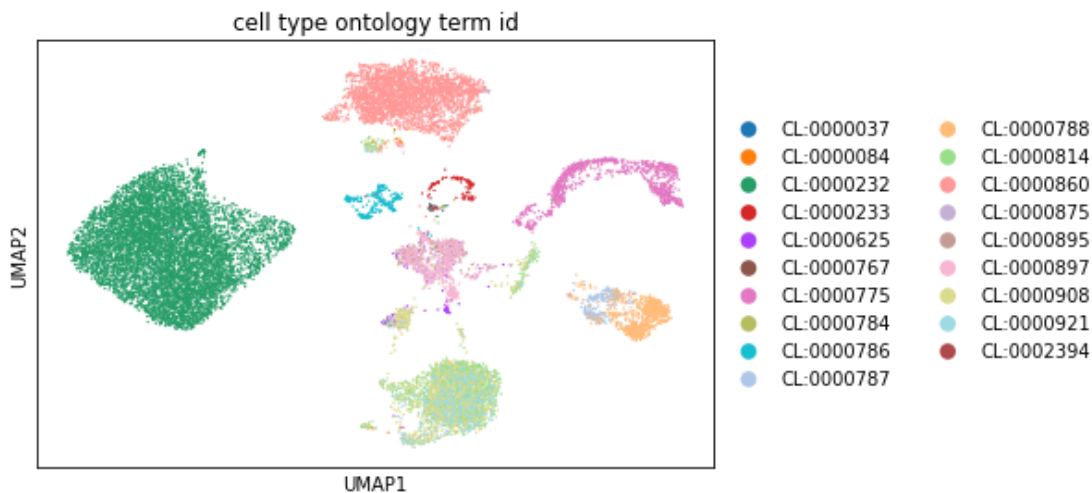
/home/ubuntu/env/lib/python3.9/site-packages/scanpy/preprocessing/_simple.py:25
1: ImplicitModificationWarning: Trying to modify attribute `.var` of view, initializing view as actual.
    adata.var['n_cells'] = number
```

```
In [3]: # get marker genes from gene expression snapshot
X = tiledb.open('prod-cube/marker_genes/')
marker_genes_df = X.df(['UBERON:0000178', 'NCBITaxon:9606'], [])
marker_genes_df = marker_genes_df[marker_genes_df['effect_size_ttest'].notnull()]
```

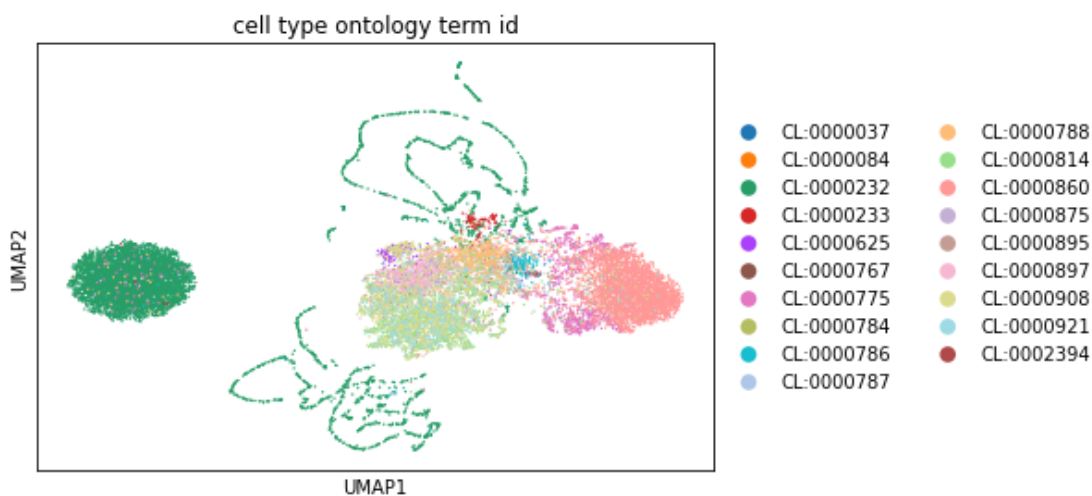
```
In [4]: var_names = np.array(list(adata.var_names))
def agg_func(df):
    g = np.array(list(df['gene_ontology_term_id']))
    df = df[np.in1d(g, var_names)]
    x = df['effect_size_ttest']
    ix = np.argsort(x)[-5:]
    l = list(np.array(list(df['gene_ontology_term_id'])))[ix]
    assert len(set(l)) == len(l)
    return l
marker_genes = list(set(np.concatenate(marker_genes_df.groupby('cell_type_ontology_
print('Found', len(marker_genes), 'unique marker genes.'))
```

Found 293 unique marker genes.

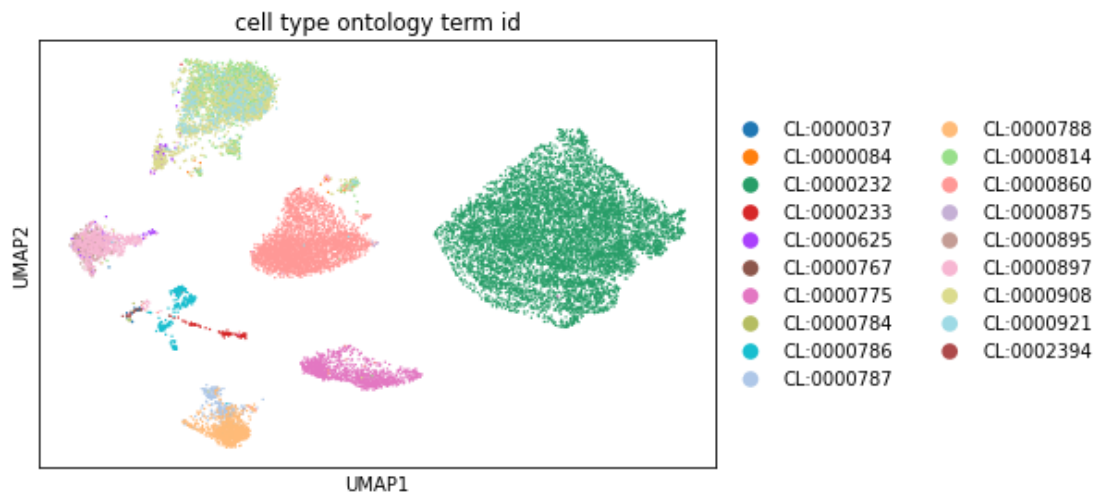
```
In [5]: # analyze using standard workflow
sc.pp.highly_variable_genes(adata,n_top_genes=3000)
adata_orig = adata[:,adata.var['highly_variable']]
sc.tl.pca(adata_orig)
sc.pp.neighbors(adata_orig)
sc.tl.umap(adata_orig)
sc.pl.scatter(adata_orig,basis='umap',color='cell_type_ontology_term_id')
sc.tl.leiden(adata_orig)
```



```
In [6]: # analyze using random genes (same number of genes as the selected markers)
random_genes = np.random.choice(adata.var_names,replace=False,size=len(marker_genes))
adata_rand = adata[:,random_genes]
sc.tl.pca(adata_rand)
sc.pp.neighbors(adata_rand)
sc.tl.umap(adata_rand)
sc.pl.scatter(adata_rand,basis='umap',color='cell_type_ontology_term_id')
sc.tl.leiden(adata_rand)
```



```
In [7]: # analyze using marker genes
adata_sub = adata[:,marker_genes]
sc.tl.pca(adata_sub)
sc.pp.neighbors(adata_sub)
sc.tl.umap(adata_sub)
sc.pl.scatter(adata_sub,basis='umap',color='cell_type_ontology_term_id')
sc.tl.leiden(adata_sub)
```



```
In [8]: ari_orig = adjusted_rand_score(adata.obs['cell_type_ontology_term_id'], adata_orig.obs['cell_type_ontology_term_id'])
ari_sub = adjusted_rand_score(adata.obs['cell_type_ontology_term_id'], adata_sub.obs['cell_type_ontology_term_id'])
ari_rand = adjusted_rand_score(adata.obs['cell_type_ontology_term_id'], adata_rand.obs['cell_type_ontology_term_id'])
```

```
In [9]: nmi_orig = adjusted_mutual_info_score(adata.obs['cell_type_ontology_term_id'], adata_orig.obs['cell_type_ontology_term_id'])
nmi_sub = adjusted_mutual_info_score(adata.obs['cell_type_ontology_term_id'], adata_sub.obs['cell_type_ontology_term_id'])
nmi_rand = adjusted_mutual_info_score(adata.obs['cell_type_ontology_term_id'], adata_rand.obs['cell_type_ontology_term_id'])
```

```
In [10]: print("Adjusted rand score")
print("Default",ari_orig)
print("Markers",ari_sub)
print("Random",ari_rand)
print("\nNormalized mutual information")
print("Default",nmi_orig)
print("Markers",nmi_sub)
print("Random",nmi_rand)
```

```
Adjusted rand score
Default 0.3537439087476417
Markers 0.3066631737299441
Random 0.16314821730360585
```

```
Normalized mutual information
Default 0.6929651763528714
Markers 0.6800783658599662
Random 0.4520186295366571
```