	# Import the libraimport scanpy as import harmonypy import leidenalg	
2]:	anndata 0.9.1 scanpy 1.9.3	versions() n["HOME"]+"/public/lab6" _10x_mtx(DATADIR, prefix="GSM5114461_S6_A11_", cache=True)
	PIL anyio attr babel backcall beta_ufunc binom_ufunc bottleneck brotli certifi	8.3.1 NA 21.2.0 2.9.1 0.2.0 NA NA 1.3.2 NA 2022.12.07
	cffi chardet charset_normalize cloudpickle colorama cycler cython_runtime cytoolz dask	1.14.6 4.0.0
	dateutil debugpy decorator defusedxml entrypoints fastjsonschema fsspec google h5py	2.8.2 1.6.7 5.0.9 0.7.1 0.3 NA 2021.07.0 NA 3.3.0
	harmonypy idna igraph ipykernel ipython_genutils jedi jinja2 joblib json5	NA 3.1 0.10.4 6.17.1 0.2.0 0.18.0 3.0.1 1.0.1 NA
	jsonschema jupyter_server jupyterlab_server kiwisolver leidenalg llvmlite markupsafe matplotlib matplotlib_inline	1.3.1 0.9.1 0.36.0 2.0.1 3.4.2
	mpl_toolkits natsort nbclassic nbformat nbinom_ufunc numba numexpr numpy packaging	NA 8.3.1 NA 5.8.0 NA 0.53.1 2.7.3 1.21.1 21.0
	pandas parso pexpect pickleshare pkg_resources plotly prometheus_client prompt_toolkit psutil ptyprocess	1.5.3 0.8.2 4.8.0 0.7.5 NA 5.14.1 NA 3.0.19 5.8.0 0.7.0
	pvectorc pydev_ipython pydevconsole pydevd pydevd_file_utils pydevd_plugins pydevd_tracing pyexpat pygments	NA NA NA NA 2.9.5 NA
	pyparsing pyrsistent pytz requests rfc3339_validator rfc3986_validator ruamel scipy send2trash session_info	
	six sklearn sniffio socks storemagic tblib terminado texttable threadpoolctl	1.16.0 0.24.2 1.2.0 1.7.1 NA 1.7.0 0.10.1 1.6.7 2.2.0
	tlz toolz tornado traitlets typing_extensions urllib3 wcwidth websocket yaml	1.26.6 0.2.5 0.57.0 5.4.1
		25.0.2 7.25.0 7.4.9 4.12.0 3.0.16 6.4.0 ckaged by conda-forge (default, Jun 19 2021, 00:32:32) [GCC 9.3.0] .88.1.el7.x86_64-x86_64-with-glibc2.31
3]:	DATADIR=os.enviro dsets = ["GSM5114 adatas = {} for ds in dsets:	on updated at 2023-05-15 23:47 Is into one AnnData object by "concatenating" multiple anndata objects. In ["HOME"]+"/public/lab6" 461_S6_A11", "GSM5114464_S7_D20", "GSM5114474_M3_E7"]
	combined = ad.con combined.obs_name GSM5114461_S6_A11 GSM5114464_S7_D20 GSM5114474_M3_E7 /opt/conda/lib/py	
5]:	utils.warn_name combined # will p AnnData object wi obs: 'dataset adatas["GSM511446	s_duplicates("obs") rint out the dimensions of the combined dataset loaded th n_obs × n_vars = 12357 × 20621 ' 1_S6_A11"] # will print out dimensions of one of the individual datasets
6]: 6]:	var: 'gene_id adatas["GSM511446 AnnData object wi var: 'gene_id	th n_obs × n_vars = 4793 × 20621 s', 'feature_types' 4_S7_D20"] # print dimensions of second dataset th n_obs × n_vars = 4910 × 20621 s', 'feature_types' 4_M3_E7"] # print dimensions of third dataset
	var: 'gene_id combined.obs # wi # Yo	th n_obs × n_vars = 2654 × 20621 s', 'feature_types' ll print out info about each cell. u should see a "dataset" column indicating which dataset each cell came from dataset A GSM5114461_S6_A11
	AAACCTGTCCTAAGTC AAACCTGTCGCCTGT	C GSM5114461_S6_A11 A GSM5114461_S6_A11 G GSM5114461_S6_A11 T GSM5114461_S6_A11 GSM5114474_M3_E7
	TTTGTCATCCCGGAT	G GSM5114474_M3_E7 G GSM5114474_M3_E7 G GSM5114474_M3_E7
9]:	Question 2: Initial Filterscellsc.pp.filter_cell	tering and normalizing your dataset ering Code s(combined, min_counts = 1000) # min 1000 counts (> 1000 reads) s(combined, min_genes = 200) # min 200 genes (> 200 genes expressed) s(combined, min_counts = 15) # min 15 counts (> 15 total count)
9]:	<pre>combined # check # n_obs are cells AnnData object wi</pre>	s(combined, min_cells = 5) # min 5 cells (> 5 cells) number of cells and genes after filtering and n_vars are genes th n_obs × n_vars = 10133 × 15779 ', 'n_counts', 'n_genes'
	/opt/conda/lib/py d in a future ver plot_data = [np	I Filtering r_genes(combined, n_top=20) # plot 20 highest expressed genes thon3.9/site-packages/seaborn/categorical.py:82: FutureWarning: iteritems is deprecated and will be remained sion. Use .items insteadasarray(s, float) for k, s in iter_data]
	MT-CO1 - MT-CO2 - MT-CO3 - MT-CYB - MT-ND4 - TTR - CHGA - MT-ND2 - FTL - MT-ND3 - MT	
	MT-ND3	20 30 40 50 60 % of total counts
	<pre>combined.var['hig sc.pp.calculate_q # Visualize violi sc.pl.violin(comb sc.pl.scatter(com sc.pl.scatter(com</pre>	ercent of counts in each cell that are from mitochondrial genes. hPercent'] = combined.var_names.str.startswith('MT') c_metrics(combined, qc_vars = ['highPercent'], percent_top = None, log1p = False, inplace = True) n and scatter plots of QC metrics including the percent mitochondria per cell, the count number per cell ined, ['pct_counts_highPercent', 'total_counts', 'n_genes_by_counts'], jitter = 0.4, multi_panel = True bined, x = 'total_counts', y = 'pct_counts_highPercent') bined, x = 'total_counts', y = 'n_genes_by_counts')
	pct 100 - 80 -	bined, x = 'n_genes_by_counts', y = 'pct_counts_highPercent') counts_highPercent 120000 -
	60 - 9	9 6000 - 4000 - 4000 - 2000 - 2000 -
	100 - 100 - 80 - 40 -	
	40 - 0 - 0 20000	40000 60000 80000 100000 120000 total_counts
	- 0000 0000 0000 0000 0000 0000	
	2000	00 40000 60000 80000 100000 120000 total_counts
	20 - 20 - 0 - 0 - 0 - 0 - 0 - 0 - 0 - 0	
4]:		<pre>n_genes_by_counts th a high percentage of counts from mitochondrial genes. The paper we got the data from suggested using ined[(combined.obs.pct_counts_highPercent < 25) & (combined.obs.total_counts < 70000), :] # Also filter</pre>
	sc.pp.normalize_p	er_cell(combined, counts_per_cell_after=1e4) # normalize to 10,000 reads/cell ned) # log transform
	<pre>sc.pp.normalize_p sc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP'</pre>	<pre>er_cell(combined, counts_per_cell_after=1e4) # normalize to 10,000 reads/cell ned) # log transform highest dispersion - use batch_key = dataset and n_top_genes = 500 able_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500) var.sort_values('dispersions_norm', ascending = False).index[0:5]) #sort and print Y', 'LYZ', 'KRT17', 'NTS'], dtype='object')</pre>
	sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP' adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0	er_cell(combined, counts_per_cell_after=1e4) # normalize to 10,000 reads/cell ned) # log transform highest dispersion - use batch_key = dataset and n_top_genes = 500 able_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500) var.sort_values('dispersions_norm', ascending = False).index[0:5]) #sort and print Y', 'LYZ', 'KRT17', 'NTS'], dtype='object') Look at Dataset n_cells highPercent n_cells_by_counts mean_counts pct_dropout_by_counts total_counts highly_variable means dispersion 1470 False 1470 0.366130 85.492944 3710.0 False 0.183694 0.8048 3514 False 3514 0.478239 65.321228 4846.0 False 0.380966 0.5673 417 False 417 0.043324 95.884733 439.0 False 0.030544 0.3873 24 False 24 0.002467 99.763150 25.0 False 0.002365 -0.0893
7]:	sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP' adata_filt.var #	er_cell(combined, counts_per_cell_after=1e4) # normalize to 10,000 reads/cell ned) # log transform highest dispersion - use batch_key = dataset and n_top_genes = 500 abble_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500) var.sort_values('dispersions_norm', ascending = False).index[0:5]) #sort and print Y', 'LYZ', 'KRT17', 'NTS'], dtype='object') Look at Dataset n_cells highPercent n_cells_by_counts mean_counts pct_dropout_by_counts total_counts highly_variable means dispersion 1470 False
7]: [7]: _	sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP' adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0	### dataset notes highPercent notes note
7]: [7]: _	sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP'] adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0 AC007325.4 2477.0 AC007325.2 52.0 AC004556.3 29.0 AC240274.1 560.0 15779 rows × 13 column # We'll manually # dataset for the genes = ["GCG", "CPA1", "CLPS] adata_var = adata Question 5: Batch Co	### dispersion - use batch_key = dataset and n_top_genes = 500 abble_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500 abble_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500 abble_genes(adata_filt, batch_key = 'dataset', n_top_genes = 500) var.sort.values('dispersions_norm', ascending = False).index(8:5)) #sort and print Y', 'LYZ', 'KRT17', 'NTS'], dtype='object') Look at Dataset n_cells highPercent n_cells_by_counts
7]: [7]: _	sc.pp.normalize_p sc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP' adata_filt.var #	### Properties of the properti
7]: [7]: _	sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0 AC007325.2 52.0 AC004556.3 29.0 AC240274.1 560.0 15779 rows × 13 column # We'll manually # dataset for the genes = ["GCG", " "CPA1", "CLPS adata_var = adata adata_var = ada	### Properties of the properti
7]: [7]: _ 8]: [sc.pp.normalize_psc.pp.log1p(combi Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.00 NOC2L 4846.00 KLHL17 439.00 PLEKHN1 25.00 HES4 7191.00	# Log transform Alignest dispersion - use batch key = dataset and n top genes = 589
7]: [7]: _ 8]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_vari.print(adata_filt. Index(['PPY', 'NP'] adata_filt.var #	### Appart dispersion - use birch_lepy = distance and n_top_penes = 500
7]: [7]: [8]: [Sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP'] adata_filt.var #	### Card (Combined) - Counts per card after 16 / 2 mornal for to 10,000 reado/cell ### Card (Combined) - Counts (Combined) - Combined Com
7]: [7]: [8]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt Index(['PPY', 'NP adata_filt.var #	### Caching Counts_per_cell_after=tee) = normalize to 10,000 reads/cell ### 2 for crossform - one backs way - declared con o con penes = 500 ### 2 for crossform - one backs way - declared con o con penes = 500 ### 2 for controlled (inspercious ward, secondary = balacy).make(0:1) #back and print ### 2 for controlled (inspercious ward, secondary = balacy).make(0:1) #back and print ### 2 for controlled (inspercious ward, secondary = balacy).make(0:1) #back and print ### 2 for controlled (inspercious ward) ward counts per dispose by counts total counts balaby variable means depends ### 2 for controlled (inspercious ward) ward counts per dispose by counts total counts balaby variable means depends ### 2 for controlled (inspercious ward) ward counts per dispose by counts total counts balaby variable means depends ### 2 for counts ward of counts and counts per dispose by counts total counts by the counts of counts by the cou
7]: [7]: [8]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt. Index(['PPY', 'NP) adata_filt.var #	### CALL CONTROLS, COURTS per CALL afformation of the control of t
7]: [7]: [9]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0 AC007325.2 52.0 AC240274.1 560.0 15779 rows × 13 colurn # We'll manually # dataset for the genes = ["GCG", ""CPA1", "CLPS adata_var = a	### CALL CONTROLS, COURTS per CALL afformation of the control of t
7]: [7]: [9]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0 AC007325.2 52.0 AC240274.1 560.0 15779 rows × 13 colurn # We'll manually # dataset for the genes = ["GCG", ""CPA1", "CLPS adata_var = a	### # Journal Community Secretary # Additional and Languagement 2000
7]: [7]: [9]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.4 2477.0 AC007325.2 52.0 AC240274.1 560.0 15779 rows × 13 colurn # We'll manually # dataset for the genes = ["GCG", ""CPA1", "CLPS adata_var = a	## CALL PROPERTY OF THE CONTROL OF T
7]: [7]: [9]: [sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_variprint(adata_filt.* Index(['PPY', 'NP adata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC01043.1 10144.0 AC007325.4 2477.0 AC007325.2 52.0 AC04556.3 29.0 AC240274.1 560.0 15779 rows × 13 colur # we'll manually # dataset for the genes = ["GCG", ""CPA1", "CLPS adata_var = adata_war = adata_war = adata_war.obsm['.exc.pp.pca(adata_#sc.pl.pca(adata_#sc.pl.pca(adata_#sc.pl.pca(adata_#sc.pl.pca(adata_war).pl.	### Commonwest Commonwest (1997 - Inchesion and Commonwest Commonw
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP) adata_filt. Var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.2 52.0 AC024526.3 29.0 AC240274.1 560.0 15779 rows × 13 colur # We'11 manually # dataset for the genes = ["GCC", " "CPA1", "CLPS adata_var = adata_var. obsm[' 2023-05-15 23:49: Characters are actual #sc.pp.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_vsc.pl.p	### Committee 1.00
7]: [7]: [9]: [9]: [SC.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_vari ript(adata_filt.var #	
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP) adata_filt. Var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC011043.1 10144.0 AC007325.2 52.0 AC024526.3 29.0 AC240274.1 560.0 15779 rows × 13 colur # We'11 manually # dataset for the genes = ["GCC", " "CPA1", "CLPS adata_var = adata_var. obsm[' 2023-05-15 23:49: Characters are actual #sc.pp.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_vsc.pl.p	### Commence of the Commence o
7]: [7]: [9]: [9]: [SC.pp.normalize_psc.pp.log1p(combi) Question 4: # Find genes with sc.pp.highly_vari ript(adata_filt.var #	### CANSON CONTROL OF THE PROPERTY OF THE PROP
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.logip(combi) Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP adata_filt.var # n.counts SAMD1 3710.0 NC2L 4846.0 KLHL17 499.0 PLEKHN1 25.0 HESA 7191.0 AC01043.1 10144.0 AC07325.2 52.0 AC240274.1 560.0 15779 rows × 13 colur # we'll manually # denes = ["GCG", "" "CPA", "CLPs adata_var =	### CASISTON CONTROL OF THE TIME AND CONTROL OF THE TI
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.logip(combi) Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP adata_filt.var # n.counts SAMD1 3710.0 NC2L 4846.0 KLHL17 499.0 PLEKHN1 25.0 HESA 7191.0 AC01043.1 10144.0 AC07325.2 52.0 AC240274.1 560.0 15779 rows × 13 colur # we'll manually # denes = ["GCG", "" "CPA", "CLPs adata_var =	### ACCIDENT CONTROL OF SIGNATURE OF SIGNATU
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.logip(combi) Question 4: # Find genes with sc.pp.highly_vari print(adata_filt. Index(['PPY', 'NP adata_filt.var # n.counts SAMD1 3710.0 NC2L 4846.0 KLHL17 499.0 PLEKHN1 25.0 HESA 7191.0 AC01043.1 10144.0 AC07325.2 52.0 AC240274.1 560.0 15779 rows × 13 colur # we'll manually # denes = ["GCG", "" "CPA", "CLPs adata_var =	### ACCIDENT CONTROL OF SIGNATURE OF SIGNATU
7]: [7]: [9]: [9]: [Sc.pp.normalize_psc.pp.log1(combi) Question 4: # Find genes with sc.pp.highly_vari print(a(ata_filt.var # n_counts SAMD11 3710.0 NOC2L 4846.0 KLHL17 439.0 PLEKHN1 25.0 HES4 7191.0 AC007325.4 2477.0 AC007325.2 52.0 AC004556.3 29.0 AC240274.1 560.0 15779 rows × 13 colur # We'll manually # detaset ["GCC", "CPA", "CLPS adata_var adata Question 5: Batch Co Every code below is a characters are actual #Sc.pp.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_wsc.pl.pca(adata_vsc.pca(adata_	### ACCIDENT CONTROL OF SIGNATURE OF SIGNATU
7]: [# Find genes with sc.pp. nightly—variprinte (dispays) and the sc.pp. neghbors (msc.pp. neghbors (msc.p	Compared
7]: [# Find genes with sc.pp. nightly—variprinte (dispays) and the sc.pp. neghbors (msc.pp. neghbors (msc.p	Control of the c
7]: [# Find genes with sc.pp. nightly—variprinte (dispays) and the sc.pp. neghbors (msc.pp. neghbors (msc.p	Control of the c
7]: [Sc.pp. normalize, post post post post post post post post	Compared
7]: [CPUP. normalize, per	### CONTROL CO
	# SAMO WITH PROPERTY OF THE MERCAND TO THE MERCAND	### Control of the Co
	# SNE WITH Period # SNE WITH Period # SNE WITH Period # SAMDII 3710.0 **NOC1 4846.0 **NOC1 4846.0 **NOC1 4846.0 **NOC1 4846.0 **NOC1 4946.0 **ACO0325.4 **ACO0325.4 **SCO0325.4 **SCO0325.	Committee Comm
	# CSNE with datase # CSNE with d	Compared Continues Compare
	CONTRACTOR AND	Comparison Com