

Organizing Data Science Projects

NIH BYOB, 2019-08-22

Justin M Fear

- IRTA Fellow (NIDDK/NIH)

-  Genomics
-  Gene Regulation
-  Drosophila

- Contact:

-  @jfear
-  justin.fear@nih.gov
-  Justin Fear@nih-byob.slack.com

- Slides are at:

- <http://talks.geneticsunderground.com>



Why Project Organization





Find things quickly

(≤ 10 min)

- Find the exact code used to generate result
- Tweak a plot
- Jump back into a project picking up where you left off

Reuse code

Automate the boring stuff

- Standardized workflows
- Modularize code for testing, maintenance, and usability
- Save brain power for fun tasks





Share code and results

- Quickly show a collaborator what you did
- Summarize side analyses for archive



Recover from data disasters

- Oops we swapped sample names now you need to re-run
- We forgot to tell you about these addition 10 samples
- The file we sent you was truncated



Reproducible Research



Have you seen this?

Please don't do this

Poor uses of file names

```
.  
├── deg_jmf_final_v2.sh  
├── deg_lmm_final.sh  
├── deg_lmm_v1.sh  
├── deg_lmm_v2.sh  
├── deg_step1_jmf_v1.sh  
└── deg_step2a_jmf_v1.sh  
├── deg_step2_jmf_v1.sh  
├── deg_step3_jmf_v1.sh  
├── deg_step3_jmf_v2.sh  
├── deg_step3_lmm_v1.sh  
└── download_data_v1.sh  
└── download_data_v2.sh
```

- File names are not great for version control
 - If you name something final, you will always have another version.
- File names are not great at ordering steps in a workflow
 - Adding or re-ordering steps is confusing at best.

Make file names descriptive and concise.

```
myscript_100.sh myscript_158.sh myscript_214.sh myscript_271.sh myscript_328.sh myscript_385.sh myscript_441.sh myscript_499.sh  
myscript_101.sh myscript_159.sh myscript_215.sh myscript_272.sh myscript_329.sh myscript_386.sh myscript_442.sh myscript_49.sh  
myscript_102.sh myscript_15.sh myscript_216.sh myscript_273.sh myscript_32.sh myscript_387.sh myscript_443.sh myscript_4.sh  
myscript_103.sh myscript_160.sh myscript_217.sh myscript_274.sh myscript_330.sh myscript_388.sh myscript_444.sh myscript_500.sh  
myscript_104.sh myscript_161.sh myscript_218.sh myscript_275.sh myscript_331.sh myscript_389.sh myscript_445.sh myscript_50.sh  
myscript_105.sh myscript_162.sh myscript_219.sh myscript_276.sh myscript_332.sh myscript_38.sh myscript_446.sh myscript_51.sh  
myscript_106.sh myscript_163.sh myscript_221.sh myscript_277.sh myscript_333.sh myscript_390.sh myscript_447.sh myscript_52.sh  
myscript_107.sh myscript_164.sh myscript_220.sh myscript_278.sh myscript_334.sh myscript_391.sh myscript_448.sh myscript_53.sh  
myscript_108.sh myscript_165.sh myscript_221.sh myscript_279.sh myscript_335.sh myscript_392.sh myscript_449.sh myscript_54.sh  
myscript_109.sh myscript_166.sh myscript_222.sh myscript_27.sh myscript_336.sh myscript_393.sh myscript_44.sh myscript_55.sh  
myscript_10.sh myscript_167.sh myscript_223.sh myscript_280.sh myscript_337.sh myscript_394.sh myscript_450.sh myscript_56.sh  
myscript_110.sh myscript_168.sh myscript_224.sh myscript_281.sh myscript_338.sh myscript_395.sh myscript_451.sh myscript_57.sh  
myscript_111.sh myscript_169.sh myscript_225.sh myscript_282.sh myscript_339.sh myscript_396.sh myscript_452.sh myscript_58.sh  
myscript_112.sh myscript_170.sh myscript_226.sh myscript_283.sh myscript_33.sh myscript_397.sh myscript_453.sh myscript_59.sh  
myscript_113.sh myscript_171.sh myscript_227.sh myscript_284.sh myscript_340.sh myscript_398.sh myscript_454.sh myscript_5.sh  
myscript_114.sh myscript_172.sh myscript_228.sh myscript_285.sh myscript_341.sh myscript_399.sh myscript_455.sh myscript_60.sh  
myscript_115.sh myscript_173.sh myscript_229.sh myscript_286.sh myscript_342.sh myscript_39.sh myscript_456.sh myscript_61.sh  
myscript_116.sh myscript_174.sh myscript_230.sh myscript_287.sh myscript_343.sh myscript_3.sh myscript_457.sh myscript_62.sh  
myscript_117.sh myscript_175.sh myscript_231.sh myscript_288.sh myscript_344.sh myscript_400.sh myscript_458.sh myscript_63.sh  
myscript_118.sh myscript_176.sh myscript_232.sh myscript_289.sh myscript_345.sh myscript_401.sh myscript_459.sh myscript_64.sh  
myscript_119.sh myscript_177.sh myscript_233.sh myscript_290.sh myscript_347.sh myscript_402.sh myscript_45.sh myscript_65.sh  
myscript_110.sh myscript_178.sh myscript_234.sh myscript_291.sh myscript_348.sh myscript_403.sh myscript_460.sh myscript_66.sh  
myscript_120.sh myscript_179.sh myscript_235.sh myscript_292.sh myscript_349.sh myscript_404.sh myscript_461.sh myscript_67.sh  
myscript_121.sh myscript_179.sh myscript_236.sh myscript_293.sh myscript_34.sh myscript_405.sh myscript_462.sh myscript_68.sh  
myscript_122.sh myscript_179.sh myscript_236.sh myscript_293.sh myscript_34.sh myscript_406.sh myscript_463.sh myscript_69.sh  
myscript_123.sh myscript_180.sh myscript_237.sh myscript_294.sh myscript_350.sh myscript_407.sh myscript_464.sh myscript_6.sh  
myscript_124.sh myscript_181.sh myscript_238.sh myscript_295.sh myscript_351.sh myscript_408.sh myscript_465.sh myscript_70.sh  
myscript_125.sh myscript_182.sh myscript_239.sh myscript_296.sh myscript_352.sh myscript_409.sh myscript_466.sh myscript_71.sh  
myscript_126.sh myscript_183.sh myscript_239.sh myscript_297.sh myscript_353.sh myscript_40.sh myscript_467.sh myscript_72.sh  
myscript_127.sh myscript_184.sh myscript_240.sh myscript_298.sh myscript_354.sh myscript_410.sh myscript_468.sh myscript_73.sh  
myscript_128.sh myscript_185.sh myscript_241.sh myscript_299.sh myscript_355.sh myscript_411.sh myscript_469.sh myscript_74.sh  
myscript_129.sh myscript_186.sh myscript_242.sh myscript_299.sh myscript_356.sh myscript_412.sh myscript_470.sh myscript_75.sh  
myscript_12.sh myscript_187.sh myscript_243.sh myscript_299.sh myscript_357.sh myscript_413.sh myscript_470.sh myscript_76.sh  
myscript_130.sh myscript_188.sh myscript_244.sh myscript_300.sh myscript_358.sh myscript_414.sh myscript_471.sh myscript_77.sh  
myscript_131.sh myscript_189.sh myscript_245.sh myscript_301.sh myscript_359.sh myscript_415.sh myscript_472.sh myscript_78.sh  
myscript_132.sh myscript_189.sh myscript_246.sh myscript_302.sh myscript_35.sh myscript_416.sh myscript_473.sh myscript_79.sh  
myscript_133.sh myscript_190.sh myscript_247.sh myscript_303.sh myscript_360.sh myscript_417.sh myscript_474.sh myscript_7.sh  
myscript_134.sh myscript_191.sh myscript_248.sh myscript_304.sh myscript_361.sh myscript_418.sh myscript_475.sh myscript_80.sh  
myscript_135.sh myscript_192.sh myscript_249.sh myscript_305.sh myscript_362.sh myscript_419.sh myscript_476.sh myscript_81.sh  
myscript_136.sh myscript_193.sh myscript_24.sh myscript_306.sh myscript_363.sh myscript_41.sh myscript_477.sh myscript_82.sh  
myscript_137.sh myscript_194.sh myscript_250.sh myscript_307.sh myscript_364.sh myscript_420.sh myscript_478.sh myscript_83.sh  
myscript_138.sh myscript_195.sh myscript_251.sh myscript_308.sh myscript_365.sh myscript_421.sh myscript_479.sh myscript_84.sh  
myscript_139.sh myscript_196.sh myscript_252.sh myscript_309.sh myscript_366.sh myscript_422.sh myscript_47.sh myscript_85.sh  
myscript_13.sh myscript_197.sh myscript_253.sh myscript_30.sh myscript_367.sh myscript_423.sh myscript_480.sh myscript_86.sh  
myscript_140.sh myscript_198.sh myscript_254.sh myscript_310.sh myscript_368.sh myscript_424.sh myscript_481.sh myscript_87.sh  
myscript_141.sh myscript_199.sh myscript_255.sh myscript_311.sh myscript_369.sh myscript_425.sh myscript_482.sh myscript_88.sh  
myscript_142.sh myscript_19.sh myscript_256.sh myscript_312.sh myscript_36.sh myscript_426.sh myscript_483.sh myscript_89.sh  
myscript_143.sh myscript_1.sh myscript_257.sh myscript_313.sh myscript_370.sh myscript_427.sh myscript_8.sh  
myscript_144.sh myscript_200.sh myscript_258.sh myscript_314.sh myscript_371.sh myscript_428.sh myscript_485.sh myscript_90.sh  
myscript_145.sh myscript_201.sh myscript_259.sh myscript_315.sh myscript_372.sh myscript_429.sh myscript_486.sh myscript_91.sh  
myscript_146.sh myscript_202.sh myscript_25.sh myscript_316.sh myscript_373.sh myscript_42.sh myscript_487.sh myscript_92.sh  
myscript_147.sh myscript_203.sh myscript_260.sh myscript_317.sh myscript_374.sh myscript_430.sh myscript_488.sh myscript_93.sh  
myscript_148.sh myscript_204.sh myscript_261.sh myscript_318.sh myscript_375.sh myscript_431.sh myscript_489.sh myscript_94.sh  
myscript_14.sh myscript_206.sh myscript_263.sh myscript_31.sh myscript_377.sh myscript_433.sh myscript_490.sh myscript_95.sh  
myscript_149.sh myscript_207.sh myscript_264.sh myscript_320.sh myscript_378.sh myscript_434.sh myscript_491.sh myscript_97.sh  
myscript_150.sh myscript_208.sh myscript_265.sh myscript_321.sh myscript_379.sh myscript_435.sh myscript_492.sh myscript_98.sh  
myscript_151.sh myscript_209.sh myscript_266.sh myscript_322.sh myscript_37.sh myscript_436.sh myscript_493.sh myscript_99.sh  
myscript_152.sh myscript_209.sh myscript_267.sh myscript_323.sh myscript_380.sh myscript_437.sh myscript_494.sh myscript_9.sh  
myscript_153.sh myscript_20.sh myscript_268.sh myscript_324.sh myscript_381.sh myscript_438.sh myscript_495.sh myscript_95.sh  
myscript_154.sh myscript_210.sh myscript_269.sh myscript_325.sh myscript_382.sh myscript_439.sh myscript_496.sh myscript_96.sh  
myscript_155.sh myscript_211.sh myscript_269.sh myscript_326.sh myscript_383.sh myscript_443.sh myscript_497.sh myscript_97.sh  
myscript_156.sh myscript_212.sh myscript_26.sh myscript_327.sh myscript_384.sh myscript_440.sh myscript_498.sh myscript_98.sh  
(base) dataripper ~ /tmp/lots_of_files $
```

Poor uses of folders

- One folder to rule them all

For small projects you can get by with putting everything in a single folder. But, once you have more than 20 files it is hard to find things. Using search is only useful if you know what you are searching for.

- Too many folders

As above, once you get more than 20 or so folders is becomes hard to find things.

Make your own folder hierarchy and stick to it.

Poor uses of scripts

```
# Run this part first  
# bunch of lines of code  
  
# Run this part second  
# bunch of lines of code  
  
# Run this part third  
bunch of lines of code  
  
# Run this part fourth  
# bunch of lines of code
```

- Don't comment and uncomment parts of script

Your-future-self will have no clue what was actually run

- Don't copy and paste from a script

Beginners often has scripts with lots of comments describing each step. They copy and paste from the script onto the command line. Like commenting and uncommenting you may not know exactly what you did.

A photograph of an open closet. The top section shows hanging clothes on a wooden rod. The bottom section features several white drawers. The leftmost drawer contains folded items in various colors like red, white, and grey. The middle drawer is filled with folded towels or linens in colors such as pink, white, and blue. The rightmost drawer is filled with folded dark-colored items like black and dark blue. The overall image conveys a sense of organization and cleanliness.

How to get organized

Branch: master ▾

Commits on Aug 19, 2019
munge biomarkers for miriam and sharvani. jfear committed 2 days ago View commit 48f748e Copy
Merge pull request #190 from jfear/munge_data_for_demas ... jfear committed 2 days ago View commit dc4eb5d Copy
Commits on Aug 16, 2019
Removes old bulk demasculinization plot. jfear committed 5 days ago View commit bad4077 Copy
Adds number of genes to bars. jfear committed 5 days ago View commit 7ef4533 Copy
Adds basic plotting scripts. jfear committed 5 days ago View commit 417de53 Copy
Adds munging for demasculinization. jfear committed 5 days ago View commit a8aff18 Copy
Adjusts gene_set titles for plots. jfear committed 5 days ago View commit 292a601 Copy
Merge pull request #189 from jfear/pull_out_ns ... jfear committed 5 days ago View commit Verified 0e30146 Copy
Runs x-to-a on larval bulk gene sets. jfear committed 5 days ago View commit 8106b86 Copy
Outputs larval bulk deg gene lists. jfear committed 5 days ago View commit 0492659 Copy
Adds ns gene subsets to x-to-a-wf. jfear committed 5 days ago View commit 0c43c44 Copy
Renames male/testis biased fbgns and outputs ns gene lists. jfear committed 5 days ago View commit 26fcfa78 Copy
Merge pull request #188 from jfear/adult_carcass ... jfear committed 5 days ago View commit Verified 6d76a5e Copy
Renames figure output to soma. jfear committed 5 days ago View commit d041e4d Copy
Adds a sex biased comparison of just soma. jfear committed 5 days ago View commit da034f8 Copy
Renamed male biased to sex biased for clarity. jfear committed 5 days ago View commit 6463d6a Copy
Renamed testis biased to gonad biased for clarity. jfear committed 5 days ago View commit a75a2f7 Copy
Merge pull request #122 from jfear/demas_stack ... jfear committed 5 days ago View commit Verified 5197931 Copy
Quick notebook to prototype a few plots. jfear committed 5 days ago View commit ca50d31 Copy
Commits on Aug 15, 2019

Tools of the trade (VCS)

Version Control Software

- Use VCS such as git

VCS is like MS Words track changes

- Use the cloud to store to your code

Github, GitLab, and Bitbucket are common places to store your version controlled code.

You can make code private or public, and it is available anywhere with internet.

Tools of the trade (Workflows)



- Orchestrate analysis using workflow software
 - Galaxy
 - Make
 - Snakemake

```
my_project/
├── config
├── data -> ../remote_store/data
├── docs
├── environment.yaml
└── envs
    ├── deseq2.yaml
    └── scrublet.yaml
├── example1-wf
    ├── config
    ├── scripts
    └── Snakefile
├── example2-wf
    ├── config
    ├── scripts
    └── Snakefile
├── geo-wf
    ├── config
    ├── scripts
    └── Snakefile
├── lcdb-references -> ../remote_store/lcdb-references
├── notebook
    ├── 2019-07-01_play_with_deseq_setting.Rmd
    └── 2019-08-01_explore_howard_et_al.ipynb
├── output -> ../remote_store/output
└── paper_submission-wf
    ├── config
    ├── scripts
    └── Snakefile
├── scripts
└── Snakefile
└── src
    └── my_project
```

Project Folder Organization

- Folder structure is personal preference
- Folder names are personal preference

There are general *best practices*

```
my_project/
├── config
├── data -> ../remote_store/data
├── docs
├── environment.yaml
└── envs
    ├── deseq2.yaml
    └── scrublet.yaml
├── example1-wf
│   ├── config
│   ├── scripts
│   └── Snakefile
├── example2-wf
│   ├── config
│   ├── scripts
│   └── Snakefile
├── geo-wf
│   ├── config
│   ├── scripts
│   └── Snakefile
├── lcdb-references -> ../remote_store/lcdb-references
├── notebook
│   ├── 2019-07-01_play_with_deseq_setting.Rmd
│   └── 2019-08-01_explore_howard_et_al.ipynb
└── output -> ../remote_store/output
├── paper_submission-wf
│   ├── config
│   ├── scripts
│   └── Snakefile
├── scripts
└── Snakefile
└── src
    └── my_project
```

1. Use the same folder structure and names across projects

- Everything has its place
- You can quickly move between projects

2. Separate original data, generated data, and scripts

Not stored in version control

```
└── data # original and external  
    ├── lcdb-references # multi-project  
    └── output # generated output  
        ├── example1-wf  
        ├── example2-wf  
        └── paper_submission-wf
```

- Improves mobility
- Delineates what you generated
- Allows reuse of common data across projects

```
data # original and external  
└── external  
    ├── DroID_DPiM_2018-03-29.txt # website  
    ├── Ferrari_et_al_2006.tsv # paper  
    ├── Ferrari_et_al_2006.readme # paper details  
    └── FlyBase/ # community  
        └── maria/ # collaborator  
    └── rnaseq_samples # our data  
        ├── ...  
        └── w1118_LG_m_r4_B_C12.fastq.gz  
    └── singleCellSeqData # out data  
        ├── ...  
        └── SV_9_10X_Te/
```

3. Uses workflows to orchestrate

```
./example1-wf
```

```
├── config
│   └── config.yaml
│       └── samptable.tsv
└── scripts/
    └── Snakefile
```

```
from larval_gonad.config import read_config

common_config = read_config('../config/common.yaml')
sample_table = pd.read_csv('../expression-atlas-wf/config/samptable.tsv', sep='\t')
sample_attrs = sample_table.samplename.str.extract('(?P<species>\w+)_(?P<tissue>\w+)_(?P<sex>\w+)_(?P<re
SPECIES = sample_attrs.species.unique().tolist()
TISSUE = sample_attrs.tissue.unique().tolist()

rule all:
    input:
        '../output/neox-wf/sturgill_2007.xls',
        '../output/neox-wf/sturgill_2007_mullerD.feather',
        '../output/neox-wf/sturgill_2007_mullerE.feather',
        expand(
            '../output/neox-wf/{species}_muller_{suffix}.feather',
            species=['dpse', 'dwil'],
            suffix=['A', 'D', 'E']
        )

rule download_sturgill:
    output: '../output/neox-wf/sturgill_2007.xls'
    shell: """
        curl -o {output[0]} https://media.nature.com/original/nature-assets/nature/journal/v450/n7167/ex
"""

rule sturgill_data_prep:
    input: rules.download_sturgill.output[0]
    output:
        mullerD='../output/neox-wf/sturgill_2007_mullerD.feather',
        mullerE='../output/neox-wf/sturgill_2007_mullerE.feather'
    script: 'scripts/parse_sturgill.py'

rule muller_arm_movement_classes:
    input: '../output/expression-atlas-wf/muller_arm_assignment.feather'
    output:
        muller_a = '../output/neox-wf/{species}_muller_A.feather',
        muller_d = '../output/neox-wf/{species}_muller_D.feather',
        muller_e = '../output/neox-wf/{species}_muller_E.feather'
    script: 'scripts/muller_arms_classes.py'
```

4. Modularize reusable code

```
lcdb-wf@56c948d #submodules  
src/   # project level package  
└── my_project  
    ├── io.py  
    ├── plotting.py  
    └── stats.py  
└── tests/  
    ├── test_io.py  
    └── test_stats.py  
└── setup.py
```

```
def decompress_seq(x: int, length=16):  
    """ Un-pack a DNA sequence from a 2-bit format  
  
    Based on code from: https://github.com/10XGenomics/cellranger  
    cellranger/lib/python/cellranger/utils.py  
  
    Parameters  
    -----  
    x : int  
        Number sequence to be decoded.  
    length : int  
        Length of the barcode. This can be found in the molecular info hdf5  
        file from 10x genome.  
        molInfo.get_node_attr('/metrics', 'chemistry_barcode_read_length')  
  
    """  
    bits = 64  
    x = np.uint64(x)  
    assert length <= (bits / 2 - 1)  
    if x & (1 << (bits - 1)):  
        return "N" * length  
    result = bytearray(length)  
    for i in range(length):  
        result[(length - 1) - i] = bytearray(NUCS[x & np.uint64(0b11)].encode())[0]  
        x = x >> np.uint64(2)  
    return result.decode()  
  
    NUCS: list  
  
def two_bit_mapper(iterable):  
    """Return a dictionary mapping 2bit encoded Seqs.  
  
    Parameters  
    -----  
    iterable : list-like  
        Unique list of 2bit encoded sequences.  
  
    Returns  
    -----  
    dict : Mapper from encoded to decoded  
    """
```

5. Use a style guide and linters

- Consistent style improves readability
- Just google my language and style guide
- Linters catch syntax errors and point out style problems.
 - flake8 # python
 - lintr # R
- Fix ugly code with software
 - black # python
 - styler # R

```
for (i in seq(10)) {  
  for (j in seq(100)) {  
    if (i == j) {  
      print(TRUE)  
    } else if (i %% j == 0) {  
      print("modulo")  
    } else {  
      print(FALSE)  
    }  
  }  
}
```

```
for (i in seq(10)) {  
  for (j in seq(100)) {  
    if (i == j) {  
      print(TRUE)  
    } else if (i %% j == 0) {  
      print("modulo")  
    } else {  
      print(FALSE)  
    }  
  }  
}
```

6. Split out configuration for consistency

```
./config # Project config
├── common.yaml
├── gene_sets.yaml
└── colors.yaml

./example1-wf # Workflow config
├── config
│   ├── config.yaml
│   └── samptable.tsv
```

Project config

Contains info that is needed across the project.

- Project name and github url
- Assembly and Annotation
- alpha level

Workflow config

Anything you may tweak in the future.

- Various thresholds
- Workflow specific references
- Various Mappings (i.e. file name to title)

7. Use containers and environments (portability and

One of the hardest problems in data science is managing software.

```
./environment.yaml # project env  
  
./envs # specific tools conda envs  
└── deseq2.yaml  
└── scrublet.yaml  
└── seurat2.yaml  
└── seurat3.yaml
```

Containers (Docker, Singularity)

- Completely reproducible system
 - Kernel and Software

Environments (Conda, pipenv)

- Install and manage software versions
- Different versions of software can be installed in different environments



8., 9., 10. Documentation

What is not documented, stays not documented

What to document (Everything!)

- How was the data generated
- Record all “experiments”
 - failed attempts
 - comparing different methods
- Record the reasoning for any decision points
- Clearly describe how to get final results

Where to document (Everywhere!)

- Sample/Resource Table
- README
- Top of scripts
- Function/Class Docstrings
- Code comments (but not too many)
- Literate Programming (i.e. notebooks)
- Project Blog

Sample Table

```
./example1-wf
  └── config
      ├── config.yaml
      └── samptable.tsv
```

A1_OCPA1_OCP_1.fastq.gz	OCP	A1	A	1	25		f	0	1	0	0	A
A6_TCPA6_TCP_2.fastq.gz	TCP	A6	A	6	15		m	1	0	0	0	B

Add as much information about your samples.

Top of Scripts

- Describe what the script does
- Any major decisions that you made
- Anything to help you remember

```
"""Decide what cells to use.          Justin Fear, 2 months ago • Refactored cellselection.

Before doing downstream analysis, I need to decide what cells to use. There are four different steps of cell selection.

1. Remove empty cells
2. Remove heterotypic doublets
3. Remove homotypic doublets
4. Remove low complexity cells (or high mitochondrial expression, or high rRNA expression)

STEP 1: There are several tools available for identifying empty cells. All of these tools model total UMI and find a threshold to separate GEMs with cells and GEMs without cells. Here I am three different tools (cell ranger v2, cell ranger v3, dropLetUtils).

DECISION: I will go forward using the consensus of cell ranger v3 and dropLetUtils.

STEP 2: There are several tools available to identify heterotypic doublets. All of these methods generate synthetic doublets by mixing cells from different clusters. They then re-cluster and call cells that cluster with the synthetic cells doublets. I have decided to only use scrublet because it uses the raw count data directly while other method require a pre-clustered datasets.

DECISION: I remove cells scrublet calls as doublets

STEP 3 + 4: Homotypic doublets are impossible to identify bioinformatically. The only signal would be cells with high UMI counts and a large number of expressed genes. Similarly, cells with low complexity I need to remove cells with low gene search: str ntially cells with high mito or rRNA expression. To look at this problem I use a grid search approach where I try the combination of different filtering thresholds and compare clustering of these cells.

DECISION: Filtering of mitochondrial and rRNA had no affect on clustering so will be ignored. The low end cutoff of 200 and 500 expressed genes had very similar results, while a more extreme 1,000 genes behaved very differently. Because we have potentially quisent cell types I will use the 200 gene cutoff. The high cutoffs of 5,000 and 6,000 were similar while no cutoff performed very differently suggesting these high end genes were driving clustering. I will use 5,000 because these high expressing cells are potentially doublets.

"""

import os
import yaml
from itertools import chain
from pathlib import Path

import pandas as pd
from snakemake.shell import shell

from larval_gonad.config import read_config
from larval_gonad.io import pickle_dump

configfile: 'config/config.yaml'
common_config = read_config('../config/common.yaml')
ASSEMBLY = common_config['assembly']
TAG = common_config['tag']

SAMPLES = pd.read_csv('../config/scrnaseq-sampletable.tsv', sep='\t', index_col=0).index.tolist()
```

Functions and Classes

- Any function you will call from another script.
- Add type hints if it is confusing what goes in.
- Add examples to clearly show what the function does.

```
def decompress_seq(x: int, length=16):
    """ Un-pack a DNA sequence from a 2-bit format

    Based on code from: https://github.com/10XGenomics/cellranger
    cellranger/lib/python/cellranger/utils.py

    Parameters
    -----
    x : int
        Number sequence to be decoded.
    length : int
        Length of the barcode. This can be found in the molecular info hdf5
        file from 10x genome.
        molInfo.get_node_attr('/metrics', 'chemistry_barcode_read_length')

    """
    bits = 64
    x = np.uint64(x)
    assert length <= (bits / 2 - 1)
    if x & (1 << (bits - 1)):
        return "N" * length
    result = bytearray(length)
    for i in range(length):
        result[(length - 1) - i] = bytearray(NUCS[x & np.uint64(0b11)].encode())[0]
        x = x >> np.uint64(2)
    return result.decode()

def two_bit_mapper(iterable):
    """Return a dictionary mapping 2bit encoded Seqs.

    Parameters
    -----
    iterable : list-like
        Unique list of 2bit encoded sequences.

    Returns
    -----
    dict : Mapper from encoded to decoded
    """

NUCS: list
```

Literate Programming

```
./notebook
└── 2019-08-01_bulk_deg.Rmd
└── 2019-08-10_bulk_ma.ipynb

./docs
└── cell_number_counts.ipynb
└── permutation_summary.ipynb
```

- Jupyter Notebooks
- R Notebooks and Rmarkdown

<https://github.com/markusschanta/awesome-jupyter>

```
In [2]: # Imports
import numpy as np
import pandas as pd

from IPython.display import HTML, Image
import qgrid

import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns

from lcdblib.images.SVG import nb_svg

# Module settings
pd.options.display.max_columns = 999
sns.set_context('poster')
```

I have have been exploring the TaDa data and playing with four different analysis methods.

- Clough, Emily, Erin Jimenez, Yoo-Ah Kim, Cale Whitworth, Megan C. Neville, Leonie U. Hempel, Hania J. Pavlou, et al. 2014. "Sex- and Tissue-Specific Functions of Drosophila Doublesex Transcription Factor Target Genes." *Developmental Cell* 31 (6): 761–73. doi:10.1016/j.devcel.2014.11.021.
- Marshall, Owen J., and Andrea H. Brand. 2015. "Damidseq_pipeline: An Automated Pipeline for Processing DamID Sequencing Datasets." *Bioinformatics (Oxford, England)* 31 (20): 3371–73. doi:10.1093/bioinformatics/btv386.
- Carl, Sarah H., and Steven Russell. 2015. "Common Binding by Redundant Group B Sox Proteins Is Evolutionarily Conserved in Drosophila." *BMC Genomics* 16: 292. doi:10.1186/s12864-015-1495-3.
- Maksimov, Daniil A., Petr P. Laktionov, and Stepan N. Belyakin. 2016. "Data Analysis Algorithm for DamID-Seq Profiling of Chromatin Proteins in Drosophila Melanogaster." *Chromosome Research: An International Journal on the Molecular, Supramolecular and Evolutionary Aspects of Chromosome Biology*, October. doi:10.1007/s10577-016-9538-4.

Each pipeline has its advantages and disadvantages. The Marshall method is not statistically rigorous and can only handle duplicates, but it has an easy to run pipeline. The Clough and Carl are more statistically satisfying, but their implementation requires deciphering the papers. The main difference between these two methods are how reads are counted (Clough uses 500bp windows while Carl extends reads and counts overlapping GATC sites).

Both Carl or Clough methods reasonable, but in the end I felt both counting methods are wrong (discussed later). I have devised my own counting method, which is similar to the new Maksimov method, but use the same analysis approach as Carl and Clough.

For me to effectively argue for my method over other published methods, it is important to understand these data. Next I visualize select steps of the DamID protocol ([DamID in Living Color](#)). Then I show some real data ([Example Data](#)) at the genome level. Finally I show some comparisons of the different methods and hopefully convince you that my proposed method seems the most reasonable.

DamID in Living Color

```
In [3]: # Get dam ID SVG
s = nb_svg('.../figs/drawings/dam_cuts.slides.svg')
#s.getLayers()
```

Getting a gut feeling for what to expect from the data is important for data science. It was extremely valuable to participate with Miriam in DamID and library construction. I also needed to visualize the DamID process, which lead to some insight. Here I show some key steps in a DamID experiment.

Dedicated Project Blog

- Aggregate notebooks
 - bookdown # R
 - jupyter webbook # python
- Various static site generators
 - Pelican
 - Nikola
 - jekyll
 -

Larval Testis X Inactivation

Data Prep and Quality Control

[Cell Selection](#)

Selecting Clustering Algorithm

Seurat2 Testis 1 Clustering

Clustering

Integrated Clustering (n = 3)

Testis 1 Individual Clustering

Testis 2 Individual Clustering

Testis 3 Individual Clustering

Testis 4 Individual Clustering

Intronless Gene Expression

Neo X Gene Movement

D. Pseudoobscura (SP vs Cytes)

D. Pseudoobscura (SP vs EPS)

D. Pseudoobscura (SP vs PS)

D. Pseudoobscura (EPS vs PS)

D. willistoni (SP vs Cytes)

D. willistoni (SP vs EPS)

D. willistoni (SP vs PS)

D. willistoni (EPS vs PS)

cellranger3-wf uses cell ranger v3 with default settings. Recently, 10X Genomics released version 3 of cell ranger.

This version adds an additional step to the cell calling algorithm which improves calling of low RNA content cells.

Interestingly, this method called more cells than **cellranger-force-wf** in all replicates except **Testis 1**.

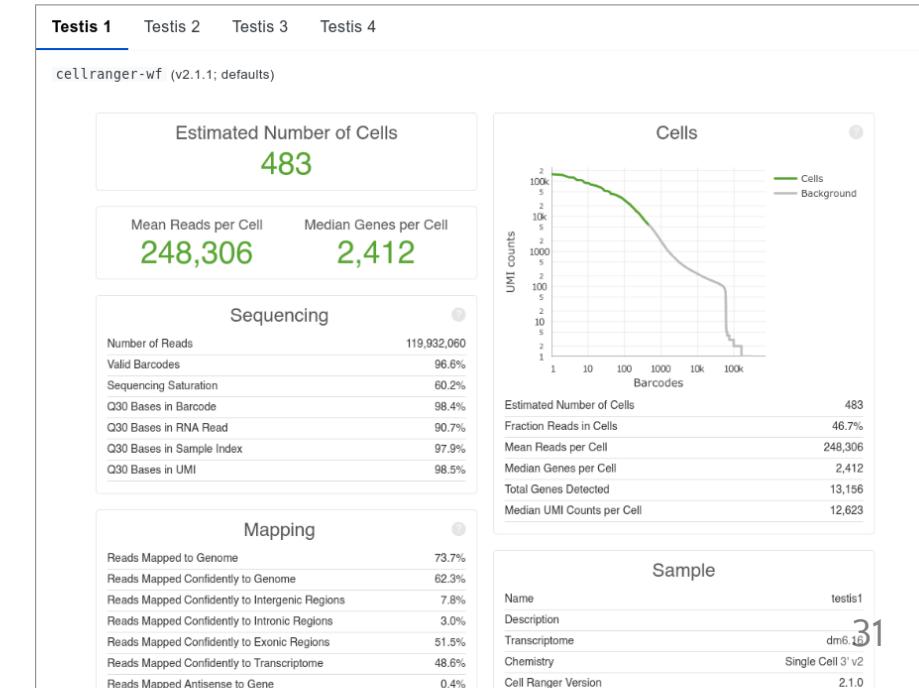
droputils uses DropLetUtils v3.0.1 with default settings. DropLetUtils is an R packaged with several algorithms for scRNA-Seq. I used the `emptydrops` function which models UMI to classify empty GEMs. This method performs very different for each replicate. This method seems to over call cells in **Testis 1**, **Testis 3**, and **Testis 4** while performing similarly to **cell ranger** on **Testis 2**.

Without a ground truth data set it is impossible to know which method is best approximating real cell calls. Consensus is often used in the absence of truth. However, since all of these methods use total UMI in their models, the 4-way consensus would be the same as the most conservative cell calls (**cellranger-wf**). Given the known limitation of **cellranger-wf** and the *ad hoc* nature of **cellranger-force-wf** I decided to use the consensus of **cellranger3-wf** and **droputils**.

Table 2. Cell count after removing empty GEMs.

Replicate	cellranger-wf	cellranger-force-wf	cellranger3-wf	droputils	2-way Consensus [‡]
Testis 1	483	3,000	2,826	13,884	2,790
Testis 2	550	3,000	6,385	5,765	4,801
Testis 3	423	8,000	12,485	39,872	12,515
Testis 4	349	8,000	15,033	34,761	15,033

[‡]Intersection of **cellranger3-wf** and **droputils**.



10 Best Practices

1. Use the same structure and names across projects
2. Separate original data, generated data, and scripts
3. Use workflows to orchestrate
4. Split out configuration for consistency
5. Modularize reusable code
6. Use a style guide and linters
7. Use containers and environments
8. Document as you go
9. Document as you go

Resources

- Cookiecutter
- Github
- Bitbucket
- Snakemake
- Data Science git-flow