

homework 05: a mixture of five

One of the many graduate students in the Holmes group, [Wiggins](#), left the lab suddenly, in protest of his stipend of one shilling a day. He left you a mess of his work in progress. He isn't responding to your emails, so you're doing some detective work of your own.

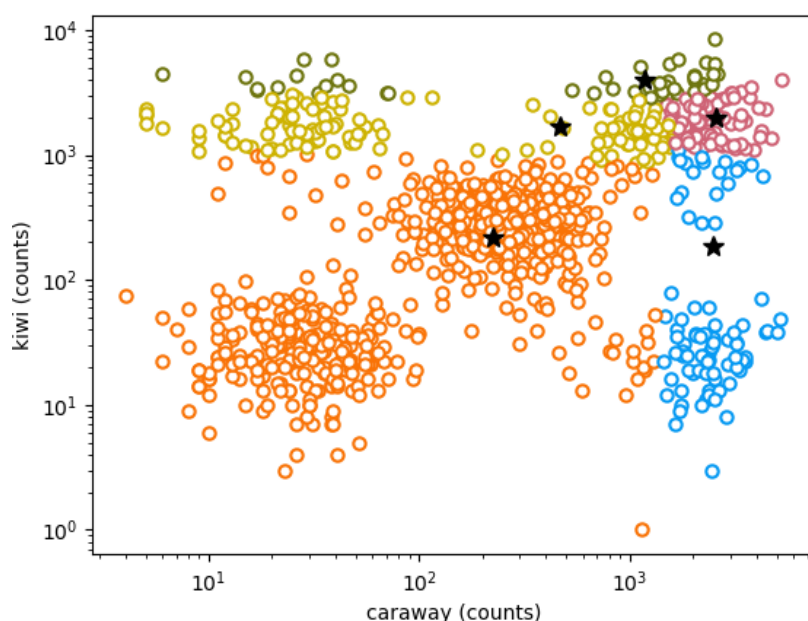
Some of Wiggins' last experiments were single-cell RNA-seqs on differentiated sand mouse embryonic stem cells. You've found records from an experiment in which he was looking at two key early transcription factor genes called *Caraway* and *Kiwi*. Previous work had shown that *Caraway* and *Kiwi* are expressed at intermediate levels in ES cells, but upon differentiation, their mRNA expression patterns break into four different cell types with all four possible combinations of low vs. high expression of these two TFs.

Wiggins' single cell RNA-seq dataset

You've found [a data file](#) where Wiggins had collected mapped read counts for *Caraway* and *Kiwi* in 1000 single differentiated ES cells. Because he expected five "cell types" in the data, he used K-means clustering, with $K=5$, to try to assign each cell to one of the five cell types, and thus estimate the mean expression level (in mapped counts) of the two genes in each cell type, and the relative proportions of the cell types in the population.

However, it's obvious from a figure you found taped into Wiggins' notebook, not to mention the various profanities written on it, that the K-means clustering did not go as hoped. His visualization of

his data does show five clear clusters, but his K-means clustering failed to identify them:



You can see in his notebook that he understood that K-means is prone to local minima, so it's not as if this is a one-off bad solution. His notes indicate that he selected the best of 20 solutions, starting from different random initial conditions. You find the following data table, and a note that this solution had a best "final tot_sqdist = 385706264.1", of 20 solutions with "tot_sqdist" from 385706264.1 to 437847350.4.

cluster	fraction	mean counts: <i>Caraway</i>	<i>Kiwi</i>
0	0.6570	224.2	221.1
1	0.0470	1177.9	3992.5
2	0.0840	2508.2	187.7
3	0.0860	2571.9	2007.5
4	0.1260	469.1	1692.3

Interestingly, it looks like all Wiggins was trying to do was to get K-means clustering to work. He must have also had the ES cells marked with reporter constructs that unambiguously labeled each of their cell types, because [his data file](#) includes a column for the true cell type (0-4), so

Figure 1: Wiggins' figure from his notebook, visualizing his K-means clustering of the 1000 cells in his single cell RNA-seq experiment, for K=5. Black stars indicate the five fitted centroids from his best K-means clustering, and colors indicate the assignments of the 1000 cells to the 5 clusters.

the true clustering is known in these data:

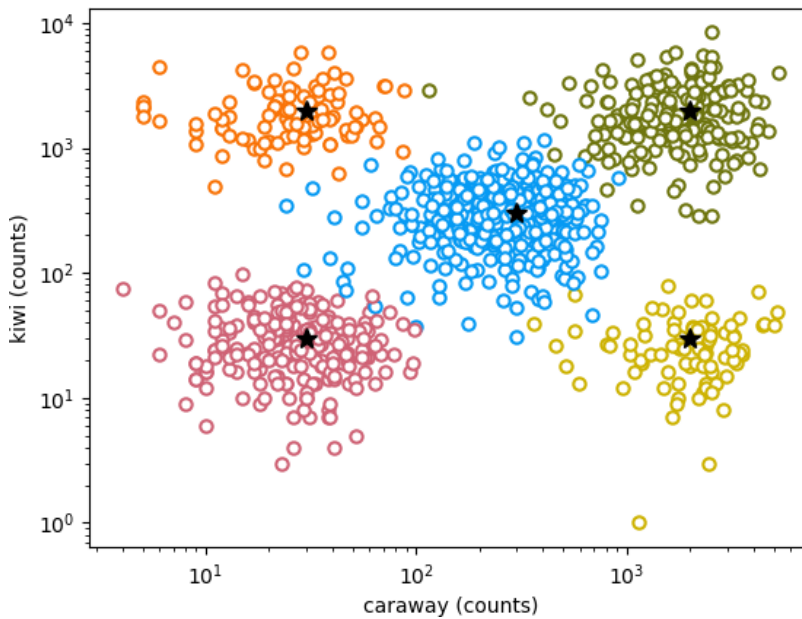


Figure 2: Another figure from Wiggins' notebook, showing the true clustering of the 1000 cells.

1. reproduce Wiggins' K-means result

Write a standard K-means clustering procedure. Use it to cluster [Wiggins' data](#) into $K=5$ clusters. Plot the results, similar to his figure. You should be able to reproduce a similarly bad result.

You'll want to run the K-means algorithm multiple times and choose the best. What is a good statistic for choosing the "best" solution for K-means? You should be able to reproduce Wiggins' "tot_sqdist" measure.

Why is K-means clustering producing this result, when there are clearly five distinct clusters in the data?

2. mixture negative binomial fitting

Now you're going to use what you've learned about mixture models, and about the negative binomial distribution for RNA-seq data.

Write an expectation maximization algorithm to fit

a mixture negative binomial distribution to Wiggins' data, for $Q=5$ components in the mixture.

Assume there is a common dispersion $\phi = 0.3$. This means that all you need to re-estimate in the EM algorithm are the means μ and mixture coefficients π for each mixture component.

Like K-means, EM is a local optimizer, so you will want to run your EM algorithm from multiple initial conditions, and take the best one. What is an appropriate statistic for choosing the "best" fit?

What are the estimated mean expression levels of *Caraway* and *Kiwi* in the five cell types, and the relative proportions of each cell type in the 1000 cells?

Visualize your result in a plot similar to Wiggins'.

3. find a simple fix for K-means

Suggest a simple fix for the problem in applying a standard K-means clustering algorithm to Wiggins' single cell RNA-seq data. Implement the fix, re-run the K-means clustering, pick a "best" solution; report and visualize it.

turning in your work

Submit your Jupyter notebook page (a .ipynb file) to the course Canvas page under the [Assignments tab](#). Please name your file <LastName> <FirstName>_<psetnumber>.ipynb; for example, mine would be EddySean_05.ipynb.

Remember that we grade your psets as if they're lab reports. The quality of your text explanations of what you're doing (and why) are just as important as getting your code to work. We want to see you discuss both your analysis code, and your biological interpretations.

hints

- The true cell type is noted in [the data file](#), so you can check whether your clustering algorithms are working well. (It's sort of obvious where the five clusters are anyway, visually, so it's not like I'm giving much away.)
- K-means (and fitting mixture models by EM) is prone to spurious local optima -- as you'll surely see. A lot of the art is in choosing initial conditions wisely. You may want to try some different strategies.
- The mixture modeling EM algorithm will be eerily parallel to your K-means algorithm (cough that's the point of the pset cough cough). The expectation step, in which you calculate the posterior probability that each data point i belongs to component q , is analogous to the K-means step of assigning a data point i to the current closest centroid q . The maximization step, in which you re-estimate the μ parameters given the posterior-weighted points, corresponds to the K-means step of estimating new centroids from the mean of their assigned points.
- Because the dispersion ϕ is given to you, you only need to estimate μ , the means of each NB component.
- A big hint on part (3): consider how K-means assigns its points to centroids, versus how we're plotting the axes to visualize these data.
- K-means is also prone to artifacts when the variances of the clusters are different, or when the variances aren't uniform in the different directions (dimensions), because it implicitly assumes spherical Gaussian distributions of equal variance. I decided against fitting the NB dispersion parameters

for this pset -- which would be easy enough to do by "method of moment" fitting, though a moderate pain by maximum likelihood, as it turns out.

- You locate [the script Wiggins used to read his data file and produce Figure 2](#). This might help you avoid some of the routine hassles of parsing input and producing output, and focus on the good bit (K-means and EM fitting of a mixture model).

take-home lessons

- On the one hand, I'm using the very intuitive K-means algorithm to illustrate how an EM algorithm works. On the other hand, I'm using EM fitting of a statistical mixture model to illustrate how K-means is a simplified case of mixture modeling that makes strong implicit assumptions.
-