

Each row in `pRating` corresponds to a food, and each column to a possible rating (column 1 corresponds to rating 1, column 2 to a rating of 2, and so on). To make it easier to read, here is the same information in a table:

Food	Rating						
	1	2	3	4	5	6	7
Ants	70%	20%	0%	0%	0%	0%	10%
Brains	10%	10%	0%	0%	0%	70%	10%
Cucumber	10%	10%	20%	30%	20%	10%	0%
Dandelions	10%	20%	20%	20%	10%	10%	10%
Emu	10%	10%	20%	30%	10%	10%	10%
Food	20%	30%	20%	0%	30%	0%	0%

- (a) One day you overhear the Klinklefigs talking about the meal they were given that day (which consisted of one of the six food items). One of them says “Ugh. It was so horrible. I would only rate it as a 1.” For each food, what is the probability that it was served that meal? In other words, what is your posterior distribution over food items given this data?
- (b) The other Klinklefigs now all weigh in with their ratings of the meal. Here is their data:

klinklefig	a	b	c	d	e	f	g	h	i
rating	1	4	3	5	2	5	3	5	1

Note that “klinklefig a” is the same one who gave you the rating in part (a). The others are all different Klinklefigs. Now what does your posterior distribution over foods look like?

- (c) Klinklefig b suddenly says “Wait! I changed my mind! I think I would give it a 2, not a 4.” How does this change your posterior distribution? Why does it change so sharply?

3. Learning to categorise zombies (30 marks)

As is traditional for an alien invasion, when the Klinklefigs devour a human brain, the corpse rises from the dead as a zombie. Ghastly as this is, it has opened up a new applied problem that needs to be solved: learning to classify zombies. Our agents in the field have been able to capture and dissect a number of these zombies, and we have learned that they come in four basic types, unimaginatively named Zombie Type 1, 2, 3 and 4. They don’t *look* much different to each other, but they seem to have different characteristics. Observationally, our field agents are able to assess zombies in terms of five characteristics: an estimate of the *health* of the zombie, an *armour* rating, a *speed*, an assessment of its *aggression* and its overall *damage* output.

You have recently been appointed as the resident Zombie Scientist, and have been asked to develop a tool that can infer the underlying type of a zombie based on the five observable ratings that a field agent provides. Your tool will be a *k*-nearest neighbours classifier, and your data can be found either in the `zombieSightings.Rdata` file or the `zombieSightings.csv` file released with this problem set.

For this exercise you can either use the `kNN` function in the `classifiers.R` file released in class, or you can write your own *k*-NN classifier if you prefer. But please note: if you’re using the code from the `classifiers.R` file, **make sure you have the most recent version before you start the problem set!** The original release of this file had a tiny bug in the `kNN` function that I’ve now fixed. If you use an older version you’ll get slightly incorrect results.

As for the data set, in the `zombieSightings.Rdata` file you’ll find a single matrix called `zombies`. Each of the 300 rows corresponds to a single zombie sighting. Here’s the first 8 rows:

```
> zombies[1:8,]
  health armour speed aggression damage type
[1,]    76  1194   235     134    267    1
[2,]   121  1260   569     131    238    1
[3,]   137  1278   108     131    301    3
[4,]   100  1155     1     115    208    4
[5,]   142  1221   236     125    226    1
[6,]   168  1242     1     130    299    2
[7,]   140  1114    63     103    158    2
[8,]   186  1297   511     131    287    3
```

As you can see, the first 5 columns are the *features*, and the 6th column is the *label*. The goal overall will be to use the feature information (i.e., `zombies[,1:5]`) to predict the label (i.e., `zombies[,6]`).

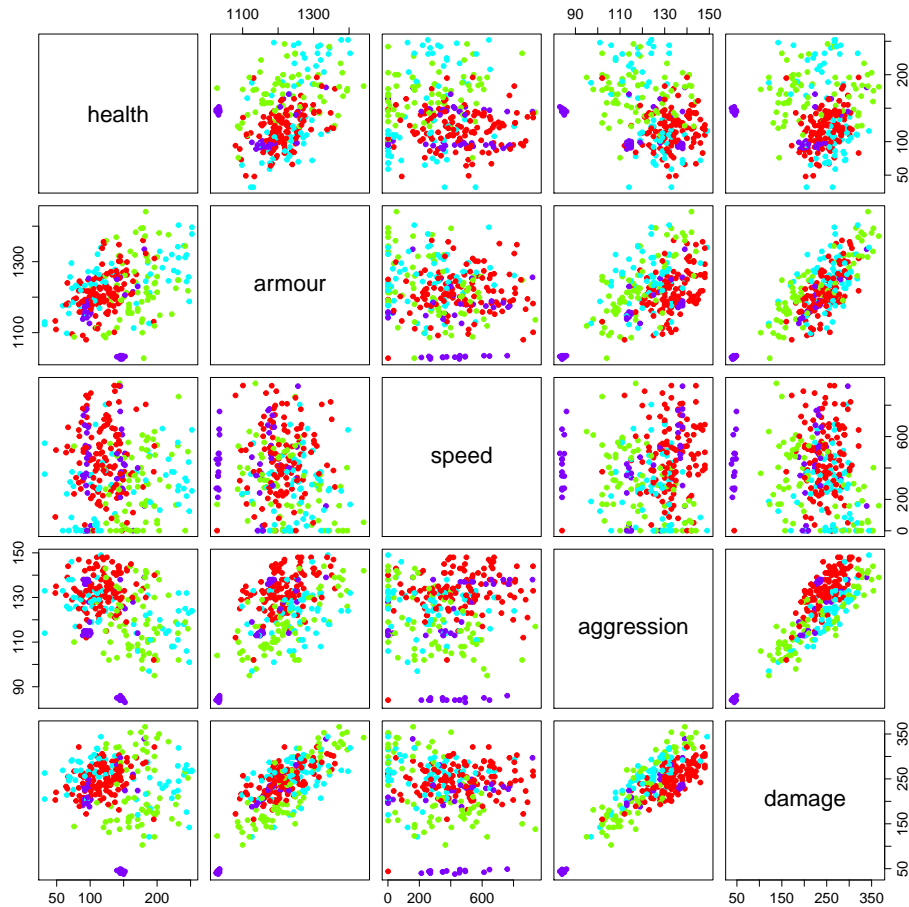


Figure 1: Scatter plots showing every pair of feature values. The red dots are Type 1 zombies, the green dots are Type 2 zombies, the light blue dots are Type 3 zombies and the purple dots are Type 4 zombies.

You can immediately tell that the data are going to be kind of weird-looking. Notice that some of the zombies have speed 1, whereas the others have quite large numbers. That's not an accident. Sometimes a zombie is so badly mangled that it can barely move, but most zombies are actually pretty quick. Clearly, we have messy data to work with! To give you a feel for the data set as a whole, Figure 1 draws a scatter plot showing every possible pair of features, plotting the four types of zombies in different colours.

- The first thing to do is get a classifier running on a small data set. Run a k nearest neighbours classifier using only the first 20 zombies (i.e., `zombies[1:20,]`) as the training data, and using the *last* 100 zombies (i.e., `zombies[201:300,]`) as your test data. Run the classifier for all values of k from 1 to 10 nearest neighbours. For all values of k , what proportion of the zombies in the test set are correctly classified? How good is this classifier relative to what you'd expect by chance?
- Now that we've got a classifier working, let's give it a bit more data to work with. Repeat the exercise for part (a), but using the first 200 zombies as training data and the last 100 zombies as test data. Compare your new results to your previous ones. Was one version better than the other? If so, why?
- One possible problem with the exercises that we've just run is that we're using the raw data as our input, but all the features seem to be on different scales. The table below shows you the maximum value and the minimum value for each of the different features, along with the difference between the two (i.e., the range):

	health	armour	speed	aggression	damage
minimum	32	1026	1	83	38
maximum	252	1441	939	149	366
range	220	415	938	66	328

It is sometimes a good idea to try rescaling the raw data. Specifically, rescale all of the features (not the label!) so that the minimum value is 0 and the maximum value is 1. Include the rescaled feature values for the first 8 zombies in your answer (similar to the way I did with the raw data at the start of this question). For the rest of this question, we'll work with the rescaled data and not with the raw data.

- (d) Repeat the exercise from part (b) using the rescaled data. That is, run a k -nearest neighbours classifier for all values of k between 1 and 10, using the first 200 observations as training data and the last 100 as test data. Does the classifier work better on the rescaled data?
- (e) Why is there a difference between the results in parts (b) and (d)? For this question you don't need to prove anything or do any calculations. Just describe in words *why* you think the classifier behaves differently when you rescale the data.
- (f) The overall accuracy can be a bit misleading. Not every category is guessed correctly with the same effectiveness. Repeat the exercise from part (d), but this time report the how accurately each of the zombie types is classified. For example, if there are 20 Type 1 zombies in the test set, how often are those correctly classified as Type 1 zombies? Do this for all four Types, and for all values of k between 1 and 10.
- (g) With respect to the results in part (f), take a close look at the pattern of accuracy for Type 1 zombies as k increases. Compare that to the pattern for Type 4 zombies. Are they the same or different? Try to explain why you're seeing the pattern that you do. Hint: try calculating how many of each type of zombie exist in the data. What does that mean when you calculate k nearest neighbours for different values of k for those types?
- (h) When training classifiers in real life, it's not a good idea to rely solely on a single train-test split. A better approach is to randomly divide the data into training and test sets several times (with each time containing different training and testing data), and then average the results. That's what we'll cover in this exercise. You'll need to randomly divide the data into 10 train-test splits (the `sample` function in R may be helpful to you in doing this). Each training set should contain 200 observations and each test set should contain 100 observations. How well does the classifier perform, for all values of k between 1 and 20? In your answer you should report the overall accuracy as well as the category-specific performance in (f).

4. Sub-types of Klinklefigs (20 marks)

After extensive experience with Klinklefigs, you realise that there might be several different races. Although it's hard to tell the difference between Klinklefigs, you have been able to identify two dimensions that appear to differentiate: WRINKLINESS and HEIGHT. You send many scouts out to gather data about the distribution of these features in the population: for every Klinklefig they see, they rate its WRINKLINESS and measure its HEIGHT. This is shown in Figure 2 and the information can be found in the file `klinklefigs.RData` available with this problem set.

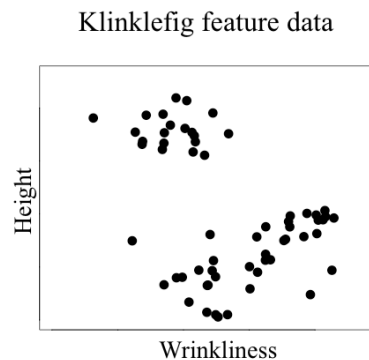


Figure 2: A graph of all of the klinklefigs your scouts have classified.

- (a) Download the code called `mixtureofgaussians.R`, which corresponds to the code presented in lecture. It is complete except for the line(s) of code that perform the E-step and M-step; locations where this code should go are clearly specified in the file. Fill in the missing code so that it now fully works. (Hint: you may use – but don't have to use – the function `getgaussiandistance`, which is in the `kmeanshelperfns.R` file). Submit the code corresponding to the modifications you made.

- (b) Based on this data, you want to figure out how many distinct races of Klinklefigs there are. To do this, we'll have you run your Mixture of Gaussians (MoG) code on this dataset many times for different possible numbers of clusters k , and for each k calculate the amount of agreement between different solutions. The file `mogsupportfns.R` contains a function called `getrandindex` which will allow you to calculate the overlap between any two different solutions based on a measure called Adjusted Rand Index (*adjRand*). It takes as arguments two r_{curr} matrices (from two different runs of `mixtureofgaussians`) and returns a value that is 1 if the matrices/clusters are identical, 0 if they are no more similar than would be expected by chance, and mid-range values to reflect some degree of overlap. Run your MoG code 10 times each for $k = 2$, $k = 3$, and $k = 4$. Then use `getrandindex` to calculate the *adjRand* for every possible pairing of datasets for each k . That is, you'll be comparing all of your results for $k = 2$ to each other, all your results at $k = 3$ to each other, and so on; this will be 45 pairs for each k . You need not compare results from $k = 2$ to $k = 3$ or $k = 4$. Doing this will yield a measure of the amount of agreement or similarity between different runs of `mixtureofgaussians`.
- What is the mean *adjRand* at each k ? The standard deviation? (You can use the `mean` and `sd` function for this).
- (c) What do you think your result in (b) reflects about how many distinct races there really are? In general, do you think estimating the "true" value of k by looking at the amount of agreement between different runs makes sense? Why or why not? (Note: you do not need to have succeeded in coding the `mixtureofgaussians` in order to answer this "in general" question; it is a thought question only).
- (d) Another idea for determining how many races there really are is to compare the likelihoods for each of the best clusterings found for each possible number of races. I did this for 2, 3, and 4 races and found that two races had the lowest likelihood, three races the next highest, and four races the highest likelihood. Does this suggest that there are actually four distinct races? Why or why not? (Note: you do not need to have succeeded in coding the `mixtureofgaussians` in order to answer this question; it is a thought question only).