

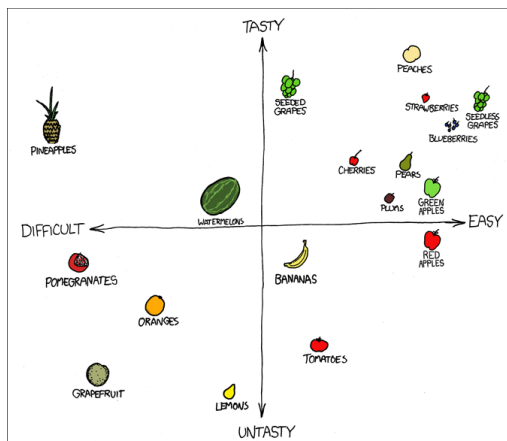
School of Computing and Information Systems
The University of Melbourne
COMP30027

MACHINE LEARNING (Semester 1, 2019)

Tutorial sample solutions: Week 2

1. Revise the definitions of **instances** and **attributes** (or **features**). For the following problems, identify what the instances and attributes might consist of:
 - (a) Building a system that guesses what the weather (temperature, precipitation, etc.) will be like tomorrow
 - (b) Predicting products that a customer would be interested in buying, based on other purchases that customer has previously made
 - (c) Automatically identifying the author of a given piece of literature
 - (d) Finding the best burrito in the United States of America
2. What are the main differences between **supervised** and **unsupervised** machine learning? What kinds of “concepts” do we typically attempt to “learn” in Machine Learning — for each, identify whether they are primarily supervised or unsupervised.
 - (a) Based on the problems in the previous question, identify the “concept” for each one, and conjecture whether a typical strategy is likely to use supervised or unsupervised Machine Learning.
 - (b) (Extension) Brainstorm some mechanisms for solving supervised problems, and some mechanisms for solving unsupervised problems. In what ways are they the same? In what ways are they different?
 - For (a), it seems fairly clear that each instance will be a day; depending on how we construe the problem, various properties could be attributes — the most logical is probably the corresponding data (temperature, precipitation, humidity, wind speed, etc.) from the previous day(s). The question suggests that there are multiple concepts here — corresponding to the various weather features of the particular day that we are trying to predict; assuming that we can access historical data for the particular location, (supervised) regression seems like the most plausible ML strategy.
 - For (b), there are a couple of different ways of construing the problem:
 - If we attempt to exhaustively label every product for every customer as either “interested” or “not interested”, then we have a classification problem, where we might try to predict these labels based on some properties of the product and customer;
 - If we instead construe a customer as an instance, we might then try to find a single product (or set of products) that the customer would be interested in. Whether this would be a supervised problem (probably classification) or unsupervised problem (probably association rule mining, or perhaps clustering) would depend on how likely we are to be able to access labelled data for a sufficient number of customers, so that we could build a sensible model.
 - For (c), again there is some question about the problem domain, for example:
 - If we have a single unknown piece of literature and a fixed set of authors who may have written it — and a collection of their previous writing — then this is probably a classification problem, where we might associate each piece of writing with the words (or grammatical structure, and perhaps metadata) contained within it;
 - If we have an open-domain problem — that potentially anybody could have written it — then collecting labelled data would be possible (i.e. classification), albeit obnoxious. We might instead prefer to use a clustering approach based on the document’s linguistic properties (although this is unlikely to identify a single author);

- We might instead have a situation like plagiarism detection, where we don't have access to very much data for any individual author. In that case, simple classification is unlikely to be very effective (because our model might be insufficient to represent each author), but we could try something like outlier detection or *semi-supervised learning* (which we'll talk about later in semester).
 - The example in (d) might seem whimsical, but this was actually attempting somewhat seriously (<https://fivethirtyeight.com/tag/burrito-bracket/> — you might like to examine their study design and features). The key question here is what the instances are: it probably isn't the case that we are looking for a single unique burrito that we can hold in our hands and say that it is truly the “best” one (whatever that means), but rather a particular restaurant (or product from a restaurant) that is consistent “better” than comparable products from other restaurants.
 - Perhaps these questions seem obtuse, but the first stage in many ML problems is working out what we are trying to achieve, and then asking ourselves whether we can collect or interpret the data in such a way that applying a Machine Learning strategy could plausibly help us to find the information that we are looking for.
3. Based on the following dataset¹ representation, would you consider “bananas” to be more similar to “apples” or “oranges”?



- (a) Identify the attributes in this dataset. What **types** of attribute are they (implicitly)? What other attributes might be relevant, and what are their types?
- As with the previous questions, we would first need to establish what it is that we are trying to do:
 - the obvious interpretation of this data is that the “easy” and “tasty” axes are recording the two attributes for the various instances (of each kind of fruit); however, at that point, it isn't clear what we could possibly wish to learn.
 - An alternative explanation might be that we wish to be able to plot other fruits (avocados, kiwifruit, dragonfruit, ...) in this space. In this case, “easy” and “tasty” are the concepts, but then the attributes that the author used to plot the given fruits are completely unknown to us — so we might need to brainstorm a set of measurable properties of fruits for ourselves, and see if they allow us to discern the given information correctly (probably using some kind of regression method).
- (b) (Extension) Critically assess the data: do you agree with how the problem is framed; do you consider the instances to be fairly located in the feature space?

¹By Randall Munroe, <https://xkcd.com/388/>, used under a Creative Commons Attribution-NonCommercial 2.5 License.