

In this module, we expand our discussion on distance measures and the concept of normalization. As we saw in k-NN, we use distance as a proxy for similarity. As distance calculations may be affected by the units of the measurement, it is also important that we normalize that data so that the measurements can be compared. This is no different from comparing cost of living indices across countries with different currencies where we standardize or normalize the indices to a single currently.

k-Nearest Neighbor (k-NN) Classification

- A distance metric is used to quantify the distance (i.e. dissimilarity) between different observations.
- Calculating the distance between the observation to be predicted and the training samples is the first step of the k-NN algorithm.
- We now introduce some of the common distance metrics that are used by k-NN algorithm. We start with distance metrics for numerical variables and then will consider categorical variables.



Let us now take a look at some distance metrics both for categorical variables, and for quantitative variables.

Common Distance Metrics

- Suppose *p* and *q* are two n-dimensional vectors. The following distance metrics are commonly used in *k*-NN algorithm:
 - ullet City Block or Manhattan distance: $d(p,q) = \sum_{i=1}^n |p_i q_i|$
 - Euclidian distance: $d(p,q) = \sqrt{\sum_{i=1}^{n} (p_i q_i)^2}$
 - lacktriangle Minkowski distance: $d(p,q) = (\sum_{i=1}^n |p_i q_i|^c)^{1/c}$
- As you can see, for c = 1, c = 2, the Minkowski metric becomes equal to the Manhattan and Euclidean metrics respectively.



We noted in the previous module that the Euclidian distance was commonly used with k-NN. Here are two other measures which we use. Note that all three measures are slight variations of each other. Again, any measure we choose should have low computational load as prediction is done as needed, and distances are calculated for each data value against the whole training set.

Sensitivity to Scale I

- Calculation of distance is very sensitive to the scale of individual attributes.
- Suppose we want to calculate the distance between the records of patients using only two attributes (variables) weight and height:
- Patient 1: weight = 66kg height=1.90m Patient 2: weight = 63kg height=1.45m Patient 3: weight = 65kg height=1.50m
- Which of the patients 1 and 2 is the closest match (i.e. nearest neighbor) to patient 3? Assume City Block is used as the distance metric.



If you remember your beginning statistics, especially the calculation of variances, you will realize that distance measures, like variances, are affected by the scale or units of measurement.

Sensitivity to Scale II

- Distance (patient 3, patient 1)= |65-66| + |1.50-1.90| = 1.40 Distance (patient 3, patient 2) = |65-63| + |1.50-1.45| = 2.05 So patient 1 is the closer match to patient 3!
- As we can see, the effect of height differences was completely masked by weight, due to different scales of the variables.
- What if height was expressed in cm instead of m?
 Distance (patient 3, patient 1)=|65-66|+|150-190|=41
 Distance (patient 3, patient 2)=|65-63|+|150-145|=7
 Completely a different result!



As you can see the result is affected by the units of measurement. One approach to solve this problem is the concept of standardization or nomalization where we convert all values of a variable into a unitless range of numbers. It is then possible to directly compare the effect or magnitude differences between variable values.

Normalization

- To avoid dependency of distance calculations to scale of data, in practice we usually normalize the data before calculating distances.
- Normalization helps to prevent attributes with large ranges out-weight attributes with small ranges
- Common Methods:
 - min-max normalization

$$x = \frac{x - x_{min}}{x_{max} - x_{min}}$$

• z-score normalization

$$\chi = \frac{x - \bar{x}}{s}$$



While both the min-max approach and z-score normalization approaches are commonly used, the latter is a more robust measures as the standard deviation In the denominator uses all values in its calculation. Further, once the data has been standardized according to the z-score formula, the data will have a mean of 0 and a variance of 1. This allows direct comparison of values across variables. For example, if you receive a Z-score of 1.4 in your exam, and another student in a totally different class receives a Z-score of 1.2, then we can state that you did better in your class compared to what the other student did in his or her class.

Normalization Example

Age	min-max (0-1)	z-score
44	0.421	0.450
35	0.184	-0.450
34	0.158	-0.550
34	0.158	-0.550
39	0.289	-0.050
41	0.342	0.150
42	0.368	0.250
31	0.079	-0.849
28	0.000	-1.149
30	0.053	-0.949
38	0.263	-0.150
36	0.211	-0.350
42	0.368	0.250
35	0.184	-0.450
33	0.132	-0.649
45	0.447	0.550
34	0.158	-0.550
65	0.974	2.548
66	1.000	2.648
38	0.263	-0.150



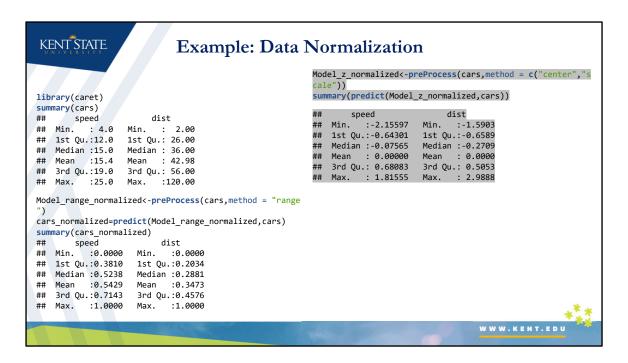
Note that all negative values of z-score indicate scores less than the average. So, age of 35 is less than the average age for that group. This information cannot be determined if we use the min-max approach.

Data Transformation: Normalization

- The preProcess () function that is in the 'caret' package is a powerful method which has implemented a number of data processing and transformation methods.
- The function implements min-max normalization using "range" as the method or z-score scaling when using "center" and "scale" as input method parameters.
- The function, however, creates a model which needs to be applied to the data using the predict () function.



Let us now implement normalization in R using the caret library.		



Code found here can be found in the github account for the class. Note that after normalization, the mean and standard deviation will become 0 and 1 respectively on the data that was used to normalize, as long as we use the z-normalization approach. This is not true if we use the range approach as seen above.

Distance Metrics for Categorical Variables

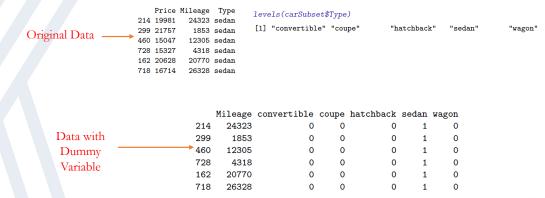
- Numerical operations such as those involved in calculation of distance metrics are not possible for categorical variables.
- Categorical variables, however, can be converted to dummy variables to perform numerical operations.
- Dummy variables are artificially defined variables designed to convert a categorical variables with multiple levels into individual binary variables.
- The process also some times is referred to as "One hot encoding".



How do you calculate distances when dealing with categorical variables? One approach is to code them as numeric. This is a common technique that is frequently used when we discuss Deep Learning. For now, the concept is to code these categorical variables as Dummy Variables that take on values 0 and 1. An example follows.



Dummy Variables: Example





Here note that we have created the same number of dummy variables as there are categories for Car Types. The dummy variables take on a value 1 if a particular car belongs to that subtype. 0 otherwise.

Dummy Variables Implementation



Dummy variables can be easily accommodated as part of the caret package. Remember to load the caret and ISLR packages before using the dummyVars function.

This concludes our module on distance calculation. Make sure to examine the code on the github account.