

Measuring Culinary Diversity of Cities

Girish Ganesan

[Introduction](#)

[Entropy](#)

[Examples of Entropy](#)

[Culinary Diversity Index](#)

[The Battle for Culinary Diversity](#)

[Data](#)

Introduction

When discussing food in a town, one often hears comments about the diversity of cuisines available. Certainly Manhattan / NYC tops the list of the places where one can find authentic cuisines from most parts of the world. When coming to LA or San Francisco, one wonders whether we have the same level of diversity available there. How about Toronto? Or London?

As a Data Scientist I have always wondered whether we could develop a quantitative metric to measure the diversity of cuisines; an absolute metric that could be used to compare culinary diversity across cities.

With that in mind, I set about to develop a Culinary Diversity Index that would measure the diversity of cuisines available in a city. The key idea is simple: we calculate the percentage of restaurants in a city for each type of cuisine. The most diverse city will have an equal percentage for each cuisine. So if we consider a list of 100 cuisines, the ideal city will have 1% of restaurants with each cuisine.

That brings us the next question: How about a city which has “almost 1%” for this number? How would we measure this error? To answer this, we leverage the concept of *Entropy*. In the next section we will motivate the discussion of Entropy and see how it can be applied to our problem of constructing a Culinary diversity index.

Entropy

In the previous section we looked at equitable distribution of cuisines as a possible measure of diversity. Instead of percentages let us look at the fraction of the total restaurants that serve a particular cuisine. If we pick a random restaurant at a location from Foursquare, this fraction will be indicative of the probability that we pick that particular cuisine. Thus if 10% of the restaurants serve Indian cuisine, the probability that a randomly picked restaurant will have Indian Cuisine is 0.1.

The thoughts we had in the previous section can now be cast in a probability framework. Diversity is maximum when all cuisines are equally probable. The more a city deviates from this equiprobable distribution of cuisines, the less diverse it is.

Even when cuisines are equiprobable, the larger the number of cuisines the larger should be the diversity. That is a city with 100 cuisines (each with probability $1/100$) should be more diverse than a city with 10 cuisines (each with probability $1/10$).

Thankfully, the science of *Information Theory* has given us a tool to measure this type of randomness: [Entropy](#). Consider a city with cuisines $c_1, c_2, c_3, \dots c_n$. Let the fraction of each cuisine be $p_1, p_2, p_3, \dots p_n$. If we pick a restaurant at random the *event* that we will end up picking a cuisine c_k occurs with a probability p_k . If we consider the set of all possible events, the Entropy of this collection of events is:

$$E = -\sum p_k \log_2(p_k)$$

Where \log_2 represents the base 2 logarithm.

Examples of Entropy

Consider an unbiased coin. The possible events that result from the coin toss are Head and Tail, each with a probability $\frac{1}{2}$. The entropy of this event collection is:

$$E = -\frac{1}{2} \log_2(\frac{1}{2}) + -\frac{1}{2} \log_2(\frac{1}{2}) = -\frac{1}{2} * (-1) + -\frac{1}{2} * (-1) = 1$$

However if the coin is biased, say the probability of Heads is 0.7, then the entropy of the events is only 0.88. And if the coin has two heads (or two tails) then only one event is possible and the Entropy is 0.

Thus we see that the more the distribution is closer to an equal probability distribution, the higher the entropy.

Now consider a fair six sided die instead of a fair coin. The cast of a die has six outcomes: 1 - 6. The Entropy of this collection of events is:

$$E = \sum \frac{1}{6} * \log_2(\frac{1}{6}) = 6 * \frac{1}{6} * \log_2(\frac{1}{6}) = 2.585$$

This is higher than the Entropy of a fair coin. Thus we see that even among equiprobable event sets, the set with higher number of events has higher Entropy. This ties in with our diversity requirement. A city with 100 equiprobable cuisines should be ranked higher than one with 10 equiprobable cuisines. With the Entropy measure in hand, we are now ready to define the Culinary Diversity Index.

Culinary Diversity Index

Consider a city with cuisines $c_1, c_2, c_3, \dots c_n$. Let the fraction of each cuisine be $p_1, p_2, p_3, \dots p_n$. The *Culinary Diversity Index* (CDI) of the city is defined as:

$$CDI = -\sum p_k \log_2(p_k)$$

The higher the CDI the more diverse the cuisines. This metric gives us a measure to compare different cities in terms of culinary diversity.

The Battle for Culinary Diversity

For the battle of Neighborhoods project I compare the culinary diversity of three major cities: New York, San Francisco and Toronto. I plan to calculate the CDI for each of the cities and see which one is more diverse in terms of CDI ranking.

Data

To calculate the CDI for a city, we need a list of cuisines and the number of restaurants serving that cuisine. This data is obtained from Foursquare. The cuisines are obtained from the Restaurant subsection of the [venue category hierarchy](#) of Foursquare. Not all cuisines are used. For example, I consider Indian cuisine but do not dive into other sub-categories. The list of cuisines considered and their Foursquare codes will be provided in the Python code.

The data for a given location / cuisine is found by querying the Foursquare places API with the location and cuisine code.