# Herramientas Computacionales para Ciencias
## Homework 12b

Mauricio Sevilla*

06/05/2019

## Kullback–Leibler divergence

> The Kullback–Leibler divergence, also called *Relative Entropy*, allows us to quantify how different two distributions of probability are. Even though it cannot be called a *metric*, it can measure the amount of information lost when using a distribution $Q$ as an approximation of $P$. The divergence $D_{KL}(P\|Q) = 0$ when the two distributions are exactly the same.
> The Kullback–Leibler divergence, on the discrete case, is calculated as follows,
>
> $$D_{\mathrm{KL}}(P\|Q) = \sum_i P(i) \log\left(\frac{P(i)}{Q(i)}\right). \tag{1}$$
>
> We are only going to need the discrete case because we are using a discretization in order the use the computer.

We are going to see *How different* is our description of a randomly generated set of points to a certain distribution by calculating the *The Kullback–Leibler divergence* between then using (1).

- First, we are going to use some ideas previously developed, out specific case of the central limit theorem, so we have to generate some sets (100000) of uniform distributed points (Only 10 points) and save only the averages on an array.

- Plot the histogram and save the bins and frequencies, use as bins `np.linspace(0,1,100)` and the option `density=True`. Use labels and titles!.

- Define a Gaussian function to make a fit, use three variables: amplitude, deviation and centroid.

- Separate the values saved for the histogram, so that you consider the center of the bin for each $x$.

- Graphically, estimate the initial parameters to do a fit of the histogram (Using the separation you just did).

- Calculate the Kullback–Leibler divergence.

- Calculate the Kullback–Leibler divergence using the Gaussian model, but for the number of sets taking the values from $10$ to $10000$ and plot the results. You may use logarithmic scale on $y$ and $x$.

   **Hint** You can calculate the histograms using the function `np.histogram` instead of `plt.hist`

---

*email=j.sevillam@uniandes.edu.co