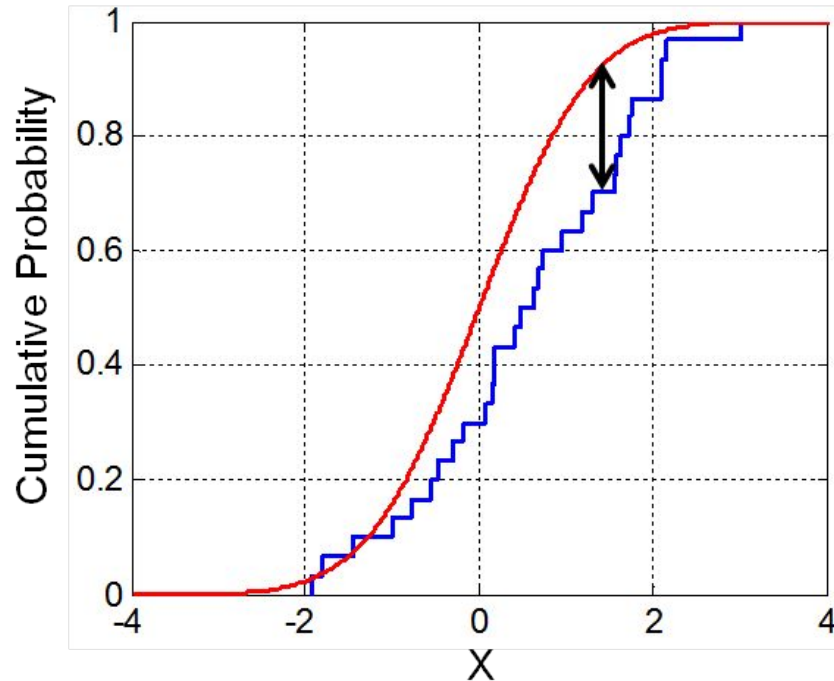# Kolmogorov-Smirnov test behaviour using Montecarlo Simulation

Dragos Tanasa - Statistical Learning

# Kolmogorov-Smirnov test



This is a statistical test that can be used to determine both if two sample came from the same distribution or if a sample is drawn from a particular probability distribution. In both cases this is done by confronting the cumulative distribution functions

$$F_n(x) = \frac{\text{number of (elements in the sample} \leq x)}{n} = \frac{1}{n} \sum_{i=1}^{n} 1_{(-\infty, x]}(X_i)$$

The statistics of interest is D_n defined as:

$$D_n = \sup_x |F_n(x) - F(x)|$$

# Code implementation

```python
33  def main():
34
35      # SIMULATION I: data is sampled from the distributions we are testing.
36      mu = 10
37      sigma = 3
38      n = 100
39      experiments = 500
40
41      d = []
42
43      for i in range(experiments):
44          data = data_generating_process(mu, sigma, n)
45          d.append(d_calculation(data, mu, sigma))
46
47      threshold = 1.36 / (pow(n, 0.5))
48
49      d = np.sort(d)
50      p = np.searchsorted(d, threshold, side='right') / d.size
51      print(p)
52
```

For every iteration in the MM simulation a normal distributed sample is generated.
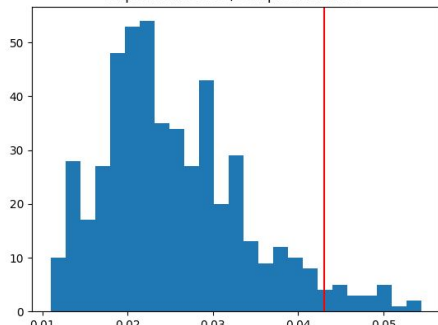
The d_calculation function outputs the value of D_n calculated between our sample and a Normal distribution with parameters given as an input to the function.

The threshold is set using tabulated values
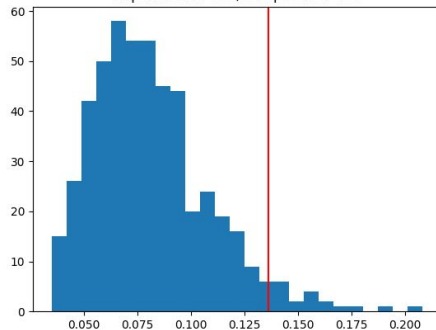
# Observations

Experiments: 500, sample size: 1000
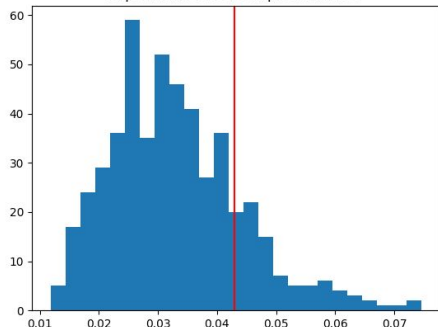


Case I:

Same distribution

Experiments: 500, sample size: 100



Case II:

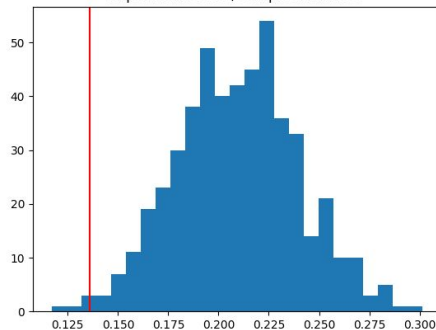Slightly different distribution, small sample size

Experiments: 500, sample size: 1000



Case III:

Slightly different distribution, bigger sample size

Experiments: 500, sample size: 1000



Case IV:

Different distributions

# References:

- https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Smirnov_test

- https://oak.ucc.nau.edu//rh83/Statistics/ks1/