# Applied Statistics - Exercise 8

## 1. Central Limit Theorem (T)

Let $X_1, X_2, \ldots$ be a sequence of independent, identically distributed random variables, with probability density function given by

$$f(x) = \begin{cases} x, & \text{if } 0 < x \le 1 \\ -x + 2, & \text{if } 1 < x < 2 \\ 0, & \text{otherwise.} \end{cases} \tag{1}$$

Use central limit theorem to approximate $P(X_1 + X_2 + \ldots + X_{277} > 301)$.

## 2. Simple Statistics (R)

Consider the `firstchi` dataset, available in the `UsingR` package, which you can load using the `library(UsingR)` statement. Using R functions, compute the following numerical statistics for the dataset.

- the sample mean
- the sample variance
- the 30th empirical percentile
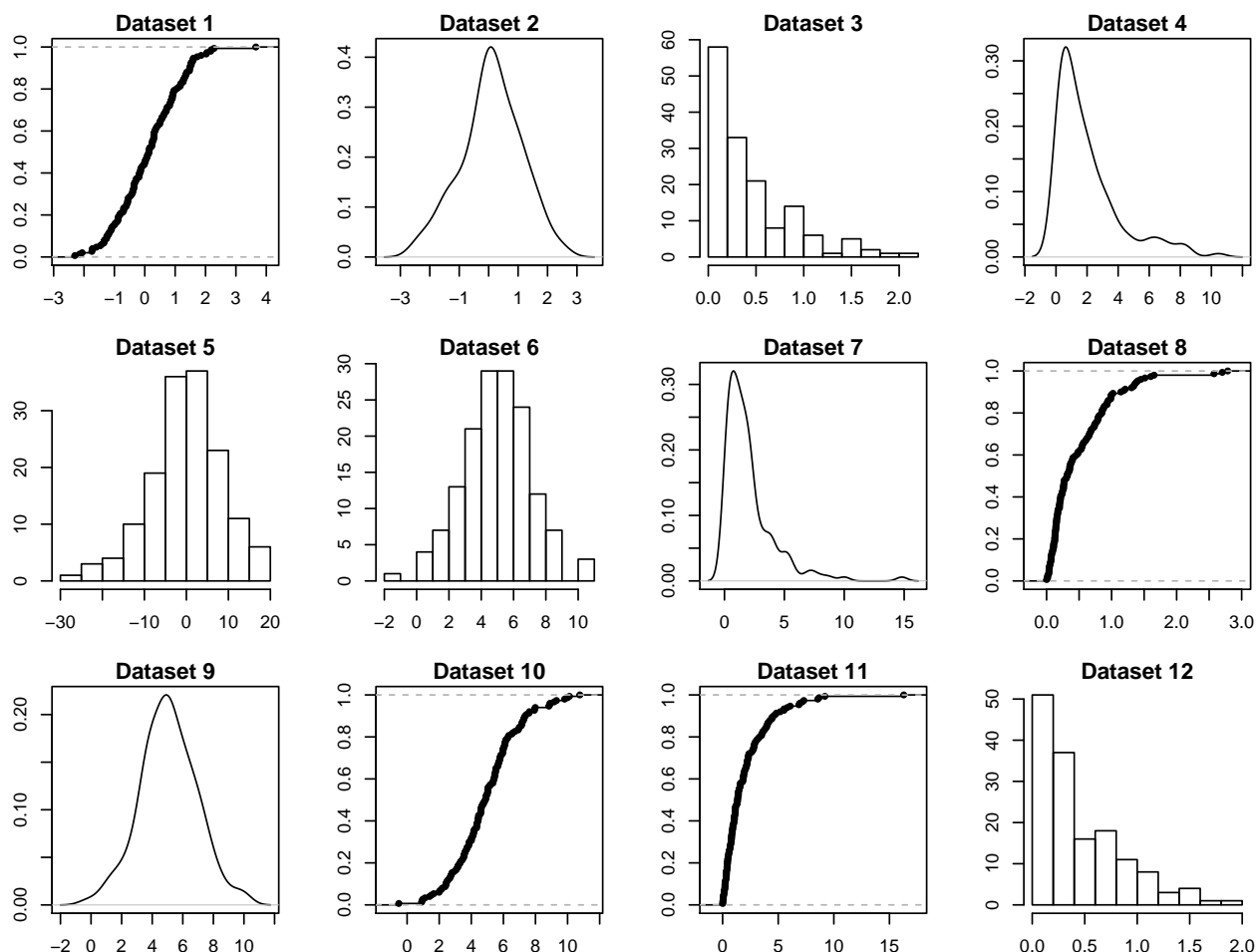- the median
- the MAD

You can refer to Section 2.3 of *Using R for introductory statistics*.

## 5. Recognizing plots (Theory)

Consider the following distributions:

- $N(0, 1)$
- $N(0, 8)$
- $N(5, 2)$
- $Exp(2)$
- $Exp(1/2)$

The following plots report histograms, kernel density estimates, and empirical distribution functions, each for a different dataset of 150 points generated from the above distributions. For each plot, say which type of plot it is (i.e. if it's a histogram, a kernel density estimate or an empirical distribution function), and identify from which of the above distributions it was generated.

**Dataset 1**  **Dataset 2**  **Dataset 3**  **Dataset 4**

**Dataset 5**  **Dataset 6**  **Dataset 7**  **Dataset 8**

**Dataset 9**  **Dataset 10**  **Dataset 11**  **Dataset 12**

# 4. Plotting distributions (R)

The `diamond` dataset of the `UsingR` package contains the price in Singapore dollars of 48 diamond rings, along with their size in carats.

1. Plot the kernel density estimate of prices. Try different bandwidths. How many modes are there? Look also at the empirical cumulative distribution function. Discuss your findings.
2. Plot a scatterplot of prices versus sizes. Does any relation between the two quantities show up?

# 5. Mean and median of two datasets (Theory)

Consider two datasets $x_1, \ldots, x_n$ and $y_1, \ldots, y_m$. Note that they have different lengths. Let $\bar{x}$ be the sample mean of the first, and $\bar{y}$ the sample mean of the second. Consider the combined dataset $x_1, \ldots, x_n, y_1, \ldots, y_m$ with $m + n$ elements, obtained by concatenating the two original datasets.

a. Is it true that the sample mean of the combined dataset is equal to $\frac{\bar{x}+\bar{y}}{2}$? If yes, provide a proof, if no, provide a counterexample.
b. Consider the case where $m = n$, i.e. the two datasets have the same size. In this special case, is the sample mean of the combined dataset equal to $\frac{\bar{x}+\bar{y}}{2}$? If yes, provide a proof, if no, provide a counterexample.

c. Consider now the sample medians $Med_x$ and $Med_x$ of the two datasets, in the general case of $m \neq n$. Is it true that the sample median of the combined dataset is equal to $\frac{Med_x + Med_y}{2}$? If yes, provide a proof, if no, provide a counterexample.

d. In the special case of $m = n$, is the sample median of the combined dataset is equal to $\frac{Med_x + Med_y}{2}$? If yes, provide a proof, if no, provide a counterexample.