

# COURS DE STATISTIQUE

## PROJET 1 : RAPPORT

MORMONT ROMAIN  
3<sup>ÈME</sup> BACHELIER INGÉNIEUR CIVIL, ORIENTATION INGÉNIEUR CIVIL  
OPTIONS *informatique* ET *électricité et électronique*  
s110940

ANNÉE ACADÉMIQUE 2013-2014

---

## Table des matières

<b>1</b>	<b>Analyse descriptive</b>	<b>3</b>
1.1	Point (a) . . . . .	3
1.2	Point (b) . . . . .	3
1.3	Point (c) . . . . .	4
1.4	Point (d) . . . . .	4
1.5	Point (e) . . . . .	5
<b>2</b>	<b>Calcul de statistiques sur échantillons</b>	<b>5</b>
2.1	Remarque sur <code>datasample</code> . . . . .	5
2.2	Point (a) . . . . .	6
2.3	Point (b) . . . . .	7
<b>3</b>	<b>Annexe</b>	<b>10</b>
3.1	Code . . . . .	10

# 1 Analyse descriptive

## 1.1 Point (a)

Les histogrammes des fréquences pour les questions 1, 2 et 3 de théorie sont respectivement donnés sur les Figures 1(a), 1(b) et 1(c). On constate qu'une majorité des étudiants ont eu plus de la moitié pour la question 1. La cote apparaissant le plus souvent (le mode) pour cette même question est le 11 (18 étudiants).

Pour la question 2, la majorité des étudiants a une cote dans l'intervalle  $[6, 15]$ . Les notes apparaissant le plus souvent sont 9 et 10 (15 étudiants pour chacune de ces cotes).

Pour la question 3, la répartition des cotes est plus ou moins uniforme et tourne autour de 7 élèves par cote. On constate néanmoins que 17 étudiants ont eu 0 et qu'aucun étudiant n'a eu 9.

De manière générale, la répartition des fréquences pour les différentes questions nous indique que la Question 1 a été, dans l'ensemble, mieux réussie que la Question 2 qui a elle-même été mieux réussie que la Question 3.

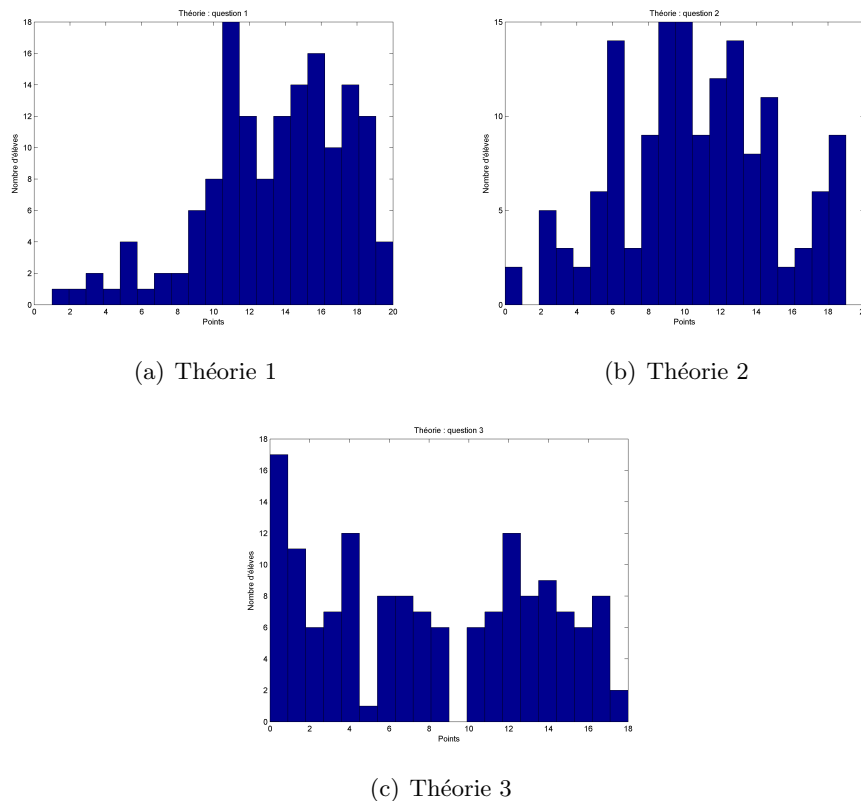


FIGURE 1 – Fréquences pour les questions de théorie

## 1.2 Point (b)

Les moyennes, modes, médianes et écarts types sont donnés dans la Table 1. On constate que l'exercice 2 a été mieux réussi que les deux autres. En effet, plus de la moitié des élèves ont eu 16 sur 20 ou plus et la cote la plus représentée est 20.

L'**exercice 3** a été quant à lui le **moins bien réussi** car plus de la moitié des élèves ont un moins de 5 sur 20 (ou 5 sur 20) avec une majorité d'élèves ayant eu 2.

En ce qui concerne les intervalles *normaux*  $[\mu_i - \sigma_i, \mu_i + \sigma_i]$ , on trouve :

- pour la **question 1** : 64.86% des étudiants dans l'intervalle  $[3.5549, 14.0937]$  (96 étudiants).
- pour la **question 2** : 57.43% des étudiants dans l'intervalle  $[10.3998, 19.6948]$  (85 étudiants)
- pour la **question 3** : 73.65% des étudiants dans l'intervalle  $[1.4892, 9.6730]$  (109 étudiants)

Exercice	Moyenne	Médiane	Mode	Écart-type
<b>Question 1</b>	8.8243	8	6	5.2694
<b>Question 2</b>	15.0473	16	20	4.6475
<b>Question 3</b>	5.5811	5	2	4.0919

TABLE 1 – Exercices d'examen

### 1.3 Point (c)

Les boîtes à moustaches mettent en évidence des résultats aberrants pour les trois projets mais pas pour la question d'examen sur le projet (voir Figure 2). Les valeurs aberrantes sont reprises dans la Table 2.

Les premier et troisième quartiles des résultats des projets et de la question d'examen sur le projet 3 sont donnés dans la Table 3.

Projet	Notes aberrantes
1	0 ( $\times 4$ ), 5.5, 7, 9
2	0 ( $\times 4$ ), 7, 11, 11.5, 12
3	0 ( $\times 4$ ), 8.11, 10.33, 10.5, 12 ( $\times 3$ )

TABLE 2 – Résultats aberrants pour les projets

Quartile	Projet 1	Projet 2	Projet 3	Examen
1 <sup>er</sup>	15.5	16	16.5	0
3 <sup>ème</sup>	19	18.5	19	18

TABLE 3 – Premier et troisième quartiles

### 1.4 Point (d)

Les graphiques des fréquences relatives cumulées pour les moyennes des questions de théorie et des questions d'exercice sont donnés respectivement sur les Figures 3(a) et 3(b). La proportion d'étudiant ayant obtenu une cote dans un certain intervalle  $[a, b]$  est obtenue à l'aide de la fonction des fréquences relatives cumulées  $F$  :

$$F(a \leq x_i \leq b) = F(x_i = b) - F(x_i = a)$$

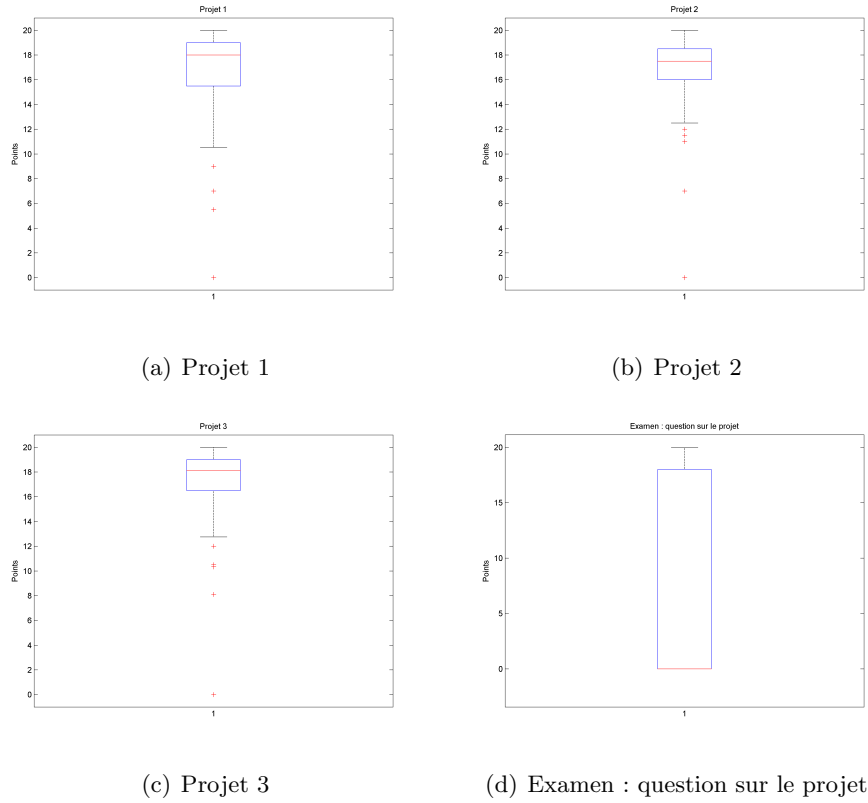


FIGURE 2 – Résultat des projets et de la question sur le projet

Les proportions obtenues pour l'intervalle  $[12, 15]$  pour les moyennes théorie et l'exercice sont respectivement **18.24%** et **22.30%**. On constate, de plus, que la forme des graphes est similaire à la forme du graphe théorique de la fréquence relative cumulée pour une loi normale (surtout pour la moyenne de théorie).

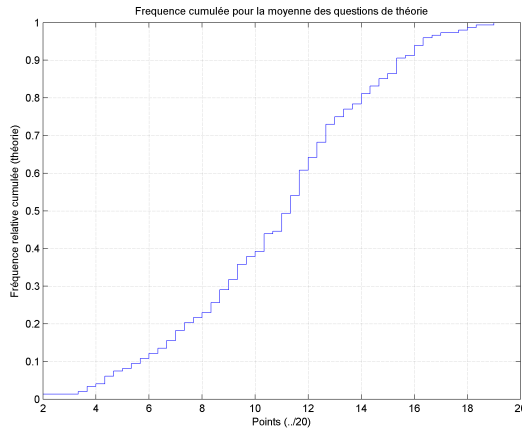
### 1.5 Point (e)

Le scatterplot entre les résultats du projet 3 et de la question sur le projet 3 est donné sur la Figure 4. Le coefficient de corrélation obtenu est **0.2106**. En se basant sur ce coefficient de corrélation, on ne peut pas tirer de conclusion de l'influence de la réussite ou non du projet 3 sur la réussite ou non de la question sur le projet 3. En effet, on observe qu'une même proportion de personne ayant réussi le projet 3 (cote  $> 10$ ) a réussi et a raté la question sur le projet 3.

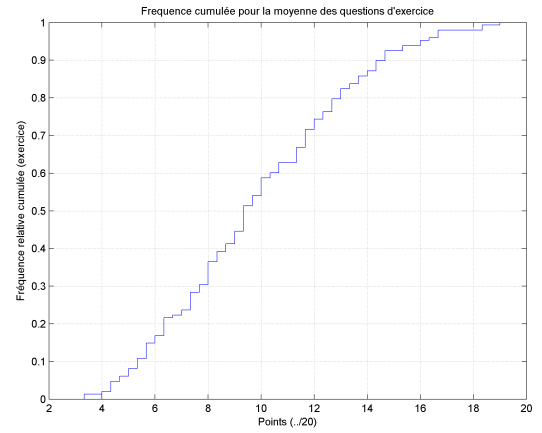
## 2 Calcul de statistiques sur échantillons

### 2.1 Remarque sur datasample

Les échantillons sont calculés avec une implémentation propre de la fonction `datasample`. Le choix d'implémenter cette fonction est du au fait qu'elle n'est pas présente sur la version R2010a de Matlab.



(a) Théorie



(b) Exercice

FIGURE 3 – Fréquences relatives cumulées

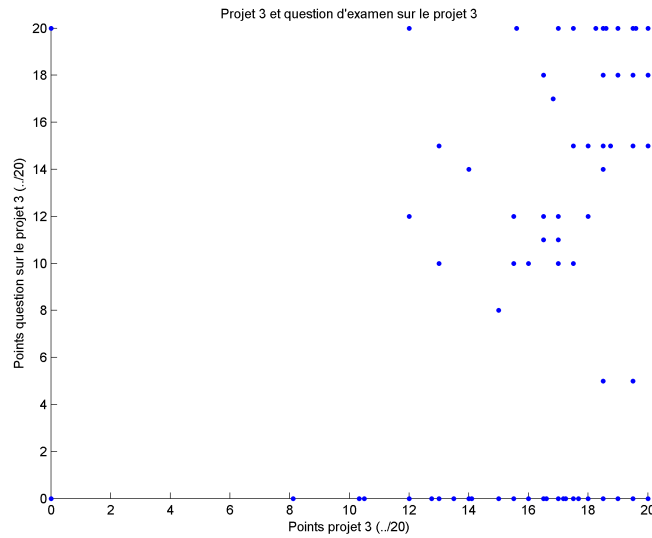


FIGURE 4 – Scatterplot entre les résultats du projet 3 et de la question sur le projet 3

Cette implémentation utilise la fonction `randi` qui effectue bien un tirage *aléatoire*<sup>1</sup> avec remise et permet donc, comme espéré, d’obtenir un échantillon i.i.d.

## 2.2 Point (a)

(i) Un échantillon généré aléatoirement à l’aide de `datasample` est donné dans la Figure 5 et les moyennes, médianes et écart types pour les résultats des exercices sont données dans la Table 4.

On observe une légère imprécision sur les statistiques calculées sur base d’un échantillon comparées à celles calculées sur base de la population. En effet, la sélection d’un échantillon

1. Le tirage étant en fait **pseudo aléatoire** et basé sur un *seed*, il pourrait être nécessaire d’initialiser ce seed avec une autre valeur que la valeur par défaut afin de ne pas générer des séries de valeurs identiques à chaque ouverture de Matlab. Néanmoins, cette opération ne me semble pas nécessaire étant donné l’ampleur de ce projet.

15, 19, 21, 24, 42, 63, 72, 81, 94, 118, 119,  
121, 135, 136, 136, 142, 142, 143, 143, 144

FIGURE 5 – Index des individus de l'échantillon étudié

	Exercice	Moyenne	Médiane	Ecart type
Échantillon	1	9.7500	9.0000	5.6557
	2	15.600	17.000	4.7617
	3	4.6000	3.5000	4.1218
Population	1	8.8243	8.0000	5.2694
	2	15.0473	16.000	4.6475
	3	5.5811	5.0000	4.0919

TABLE 4 – Moyennes, médianes et écart types pour les exercices

provoque une **perte d'information** par rapport à la population. On constate aussi, en comparant aux statistiques de la population, que l'écart-type varie moins que les autres statistiques.

(ii) Dans la continuité de ce qui a été dit au point précédent, on constate que la perte d'information liée à la sélection d'un échantillon entraîne des différences (moins ou plus de données aberrantes, déplacement des quartiles,...) entre les boîtes à moustaches tracées pour la Question 1 et celle donnée sur la Figure 6. On peut néanmoins observer que, malgré ces variations évidentes, les boîtes sont positionnées de la même manière qu'à la question 1. Cette observation n'est pas surprenante étant donné que l'échantillon a été tiré de la population qui a donné les premières boîtes.

(iii) Les courbes des fréquences cumulées pour la moyenne des questions de théorie pour l'échantillon et la population sont données sur la Figure 7. On constate que la fonction relative à l'échantillon contient moins de marches que celle de la population. Encore une fois, cette observation n'est pas surprenante étant donné que le nombre d'individus est beaucoup plus petit dans le premier cas. De ce fait, la fonction est moins proche de la courbe théorique (distribution normale) que la courbe de la population.

La distance de Kolmogorov-Smirnov entre les deux courbes est calculée à l'aide de la fonction `kstest2`. La distance de K-S entre les fonctions pour l'échantillon et la population est **0.0986**. Cette valeur, qui est relativement faible, indique que les deux distributions suivent probablement de la même loi de probabilité.

### 2.3 Point (b)

Les histogrammes pour les sous-questions (i), (ii) et (iii) sont donnés sur la Figure 8.

(i) La moyenne des moyennes (pour l'exercice 1) est **8.8645**. Cette moyenne est **proche** de la moyenne pour la population et bien plus précise que la valeur obtenue avec un seul échantillon. On peut expliquer cette précision accrue par le fait que l'augmentation du nombre d'échantillon a *atténué la perte d'information*. En effet, pour 100 échantillons le nombre d'individus différents sondés a fortement augmenter.

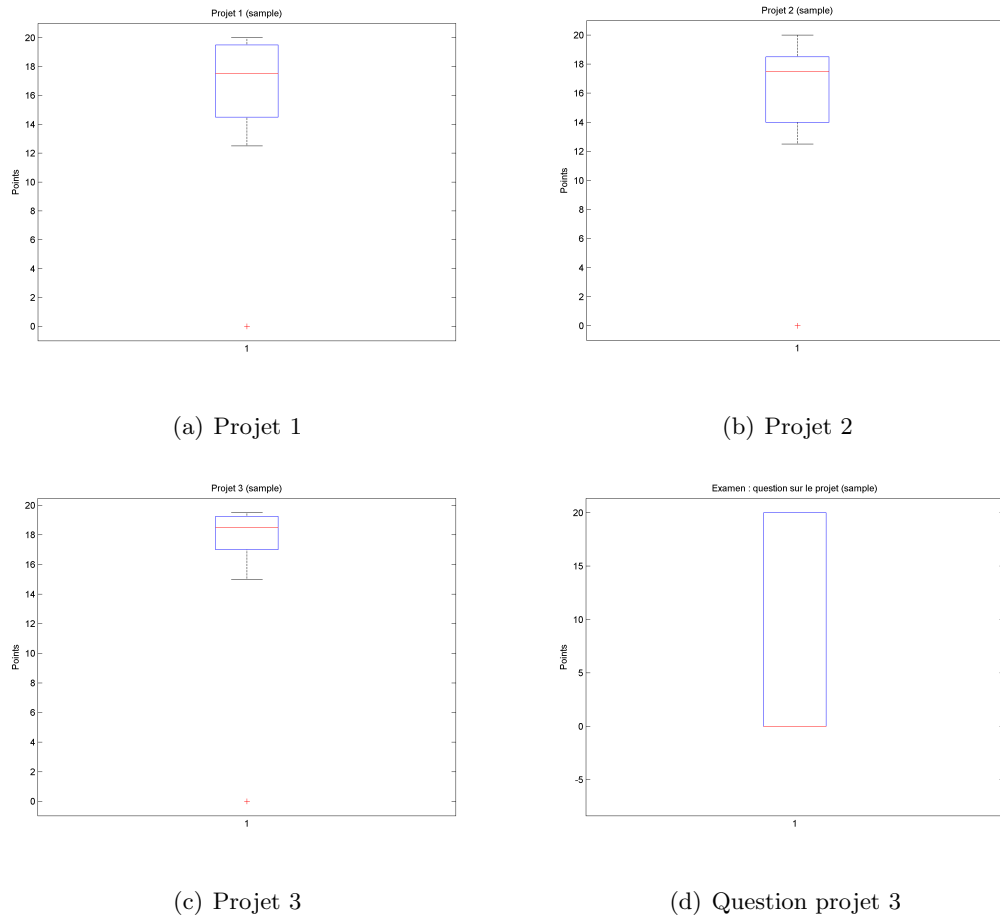


FIGURE 6 – Résultat des projets pour l'échantillon

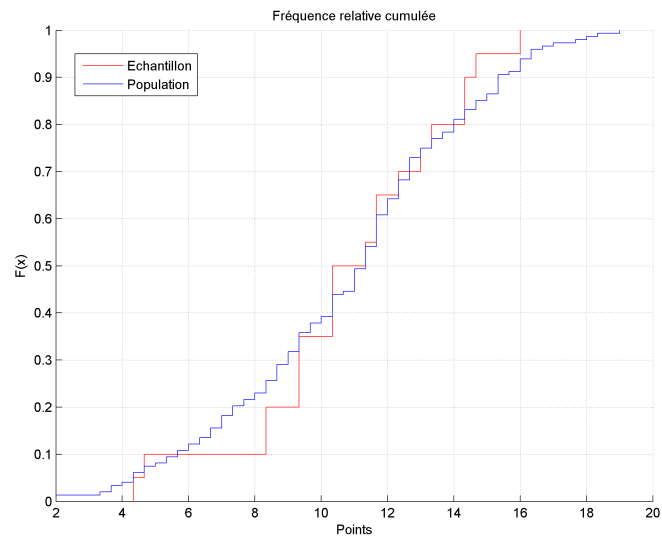


FIGURE 7 – Fréquences cumulées pour la moyenne des questions de théorie

L'allure de l'histogramme rappelle la **loi normale**. En effet, on constate une accumulation d'individus autour de la moyenne de cette variable et une décroissance de la



fonction de part et d'autre de cette moyenne.

(ii) La moyenne des médianes est **8.2250**. Cette valeur est, comme au point 2.b.(i), plus précise que celle trouvée pour un seul échantillon.

On constate, à nouveau, l'accumulation des médianes autour de la moyenne des médianes mais la décroissance de part et d'autre de la moyenne est moins évidente que pour la moyenne des moyennes.

(iii) La moyenne des écart-types est **5.3321**. Encore une fois, cette valeur est plus précise que pour un seul échantillon.

Pour ce qui est du graphe, on retrouve une loi normale sur base des mêmes observations que celles énoncés aux points précédents.

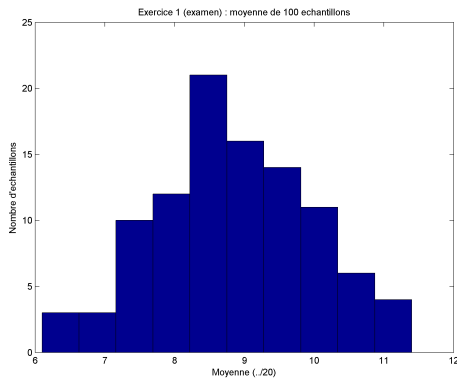
(iv) et (v) Les histogrammes des distances de Kolmogorov-Smirnov entre les fonctions de fréquences cumulées des exercices pour 100 échantillons par rapport à la population sont donnés sur la Figure 9.

L'allure des histogrammes rappelle une loi normale. La distance de K-S pour un échantillon aléatoire et la population dont l'échantillon est extrait évolue selon une loi normale.

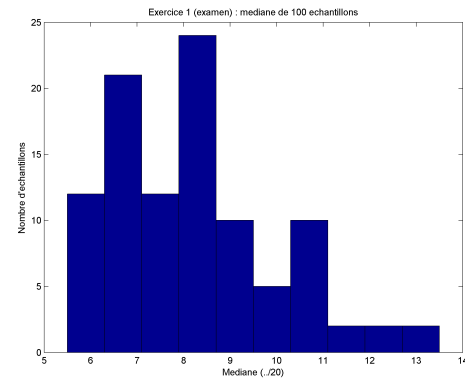
## 3 Annexe

### 3.1 Code

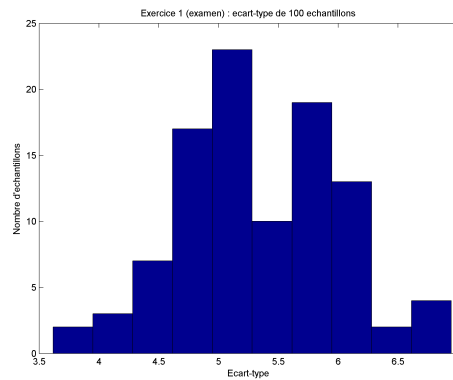
```
1 clear all;
2 close all;
3
4 % Load data
5 s_score = load('stat_data.mat');
6 score = s_score.m;
7 clear s_score;
8 % Splitting different columns into different variables
9 exam_exer = score(:,8:10);
10 exam_theo = score(:,5:7);
11 exam_proj = score(:,4);
12 projects = score(:,1:4);
13 n_students = length(score(:,1));
14
15 %% Question 1
16 % Point a)
17 figure; hist(exam_theo(:,1), 20); % theorie 1
18 title('Theorie : question 1');
19 ylabel('Nombre d\'eleves');
20 xlabel('Points');
21
22 figure; hist(exam_theo(:,2), 20); % theorie 2
23 title('Theorie : question 2');
24 ylabel('Nombre d\'eleves');
25 xlabel('Points');
```



(a) Moyenne



(b) Médiane



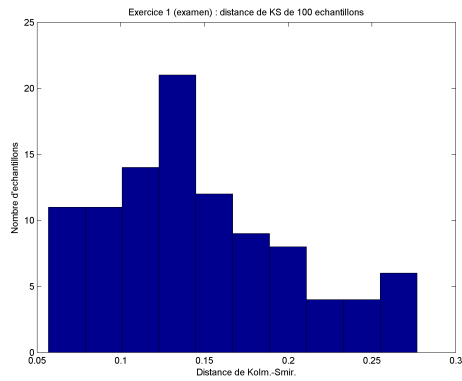
(c) Écart-type

FIGURE 8 – Calculs statistiques pour l'exercice 1 (100 échantillons)

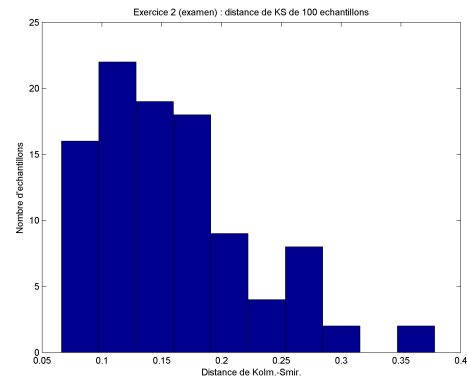
```

26
27 figure; hist(exam_theo(:,3), 20); % theorie 3
28 title('Theorie : question 3');
29 ylabel('Nombre d''eleves');
30 xlabel('Points');
31
32 % Point b)
33 ex_mean = mean(exam_exer);
34 ex_median = median(exam_exer);
35 ex_mode = mode(exam_exer);
36 ex_ectype = std(exam_exer, 1);
37
38 norm_min_threshold = ex_mean - ex_ectype;
39 norm_max_threshold = ex_mean + ex_ectype;
40
41 n_students_normal_exer = [length(find(exam_exer(:,1) >= norm_min_threshold(1) &
42     exam_exer(:,1) <= norm_max_threshold(1) )) ...
43     length(find(exam_exer(:,2) >= norm_min_threshold(2) & exam_exer(:,2) <=
44         norm_max_threshold(2) )) ...
45     length(find(exam_exer(:,3) >= norm_min_threshold(3) & exam_exer(:,3) <=
46         norm_max_threshold(3) ))];
47
48 % Point c)
49 figure; boxplot(projects(:,1));

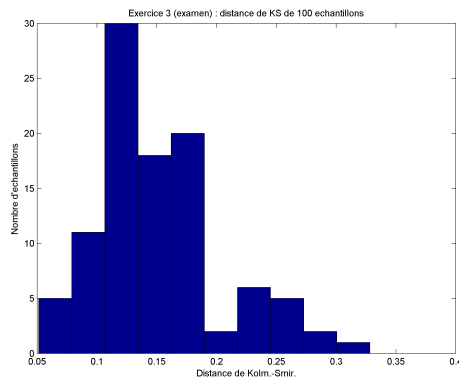
```



(a) Exercice 1



(b) Exercice 2



(c) Exercice 3

FIGURE 9 – Distance de Kolmogorov-Smirnov de 100 échantillons par rapport à la population

```

47 title('Projet 1');
48 ylabel('Points');
49
50 figure; boxplot(projects(:,2));
51 title('Projet 2');
52 ylabel('Points');
53
54 figure; boxplot(projects(:,3));
55 title('Projet 3');
56 ylabel('Points');
57
58 figure; boxplot(exam_proj);
59 title('Examen : question sur le projet');
60 ylabel('Points');
61
62 proj_Q1 = quantile(projects, 0.25);
63 proj_Q3 = quantile(projects, 0.75);
64
65 min_threshold = proj_Q1 - 1.5 * (proj_Q3 - proj_Q1);
66 max_threshold = proj_Q3 + 1.5 * (proj_Q3 - proj_Q1);
67
68 % Finding outliers indexes

```

```

69 outlier_index_p1 = [find(projects(:,1) < min_threshold(1)); find(projects(:,1) >
    max_threshold(1))];
70 outlier_index_p2 = [find(projects(:,2) < min_threshold(2)); find(projects(:,2) >
    max_threshold(2))];
71 outlier_index_p3 = [find(projects(:,3) < min_threshold(3)); find(projects(:,3) >
    max_threshold(3))];
72 outlier_index_exam_p = [find(projects(:,4) < min_threshold(4)); find(projects(:,4) >
    max_threshold(4))];
73
74 outliers_p1 = projects(outlier_index_p1, 1).';
75 outliers_p2 = projects(outlier_index_p2, 2).';
76 outliers_p3 = projects(outlier_index_p3, 3).';
77 outliers_ex_proj = projects(outlier_index_exam_p, 4).';
78
79 clear outlier_index_p1 outlier_index_p2 outlier_index_p3 outlier_index_exam_p;
80
81 % Point d)
82 theo_mean = mean(exam_theo.').';
83 exer_mean = mean(exam_exer.').';
84
85 figure; cdfplot(theo_mean);
86 title('Frequence cumulee pour la moyenne des questions de theorie');
87 xlabel('Points (./20)');
88 ylabel('Frequence relative cumulee (theorie)');
89
90 figure; cdfplot(exer_mean);
91 title('Frequence cumulee pour la moyenne des questions d\'exercice');
92 xlabel('Points (./20)');
93 ylabel('Frequence relative cumulee (exercice)');
94
95 [ecdf_theo x_ecdf_theo] = ecdf(theo_mean);
96 [ecdf_exer x_ecdf_exer] = ecdf(exer_mean);
97
98 %close all;
99 res_min = find(x_ecdf_exer >= 12, 1);
100 res_max = find(x_ecdf_exer <= 15, 1, 'last');
101 score12_15_exer = (ecdf_exer(res_max) - ecdf_exer(res_min));
102
103 res_min = find(x_ecdf_theo >= 12, 1);
104 res_max = find(x_ecdf_theo <= 15, 1, 'last');
105 score12_15_theo = (ecdf_theo(res_max) - ecdf_theo(res_min));
106
107 clear res_min res_max;
108 % Point e)
109 figure; scatter(projects(:,3), exam_proj, 4*pi, 'blue', 'fill');
110 title('Projet 3 et question d\'examen sur le projet 3');
111 xlabel('Points projet 3 (./20)');
112 ylabel('Points question sur le projet 3 (./20)');
113
114 corr_proj3 = corrcoef(projects(:,3), exam_proj);
115 corr_coef_p3 = corr_proj3(1,2);
116
117 clear corr_proj3;

```

Listing 1 – Question1.m

```

1 clear all;

```

```

2 close all;
3
4 % Load data
5 s_score = load('stat_data.mat');
6 score = s_score.m;
7 clear s_score;
8 % Splitting different columns into different variables
9 exam_exer = score(:,8:10);
10 exam_theo = score(:,5:7);
11 exam_proj = score(:,4);
12 projects = score(:,1:4);
13 n_students = length(score(:,1));
14
15 %% Question 2.a
16 [sample_2a sample_index_2a] = datasample(score, 20);
17
18 % Question 2.a.i
19 ex_mean = mean(sample_2a(:,8:10));
20 ex_median = median(sample_2a(:,8:10));
21 ex_ectype = std(sample_2a(:,8:10));
22
23 % Question 2.a.ii
24 figure; boxplot(sample_2a(:,1));
25 title('Projet 1 (sample)');
26 ylabel('Points');
27
28 figure; boxplot(sample_2a(:,2));
29 title('Projet 2 (sample)');
30 ylabel('Points');
31
32 figure; boxplot(sample_2a(:,3));
33 title('Projet 3 (sample)');
34 ylabel('Points');
35
36 figure; boxplot(sample_2a(:,4));
37 title('Examen : question sur le projet (sample)');
38 ylabel('Points');
39
40 % Question 2.a.iii
41 theo_mean_sample_2a = mean(sample_2a(:,5:7).').';
42 theo_mean = mean(exam_theo.').';
43
44 ks = ksdist(theo_mean_sample_2a, theo_mean);
45 [~, ~, ks_f] = kstest2(theo_mean_sample_2a, theo_mean);
46 figure; cdfplot(theo_mean_sample_2a);
47 figure; cdfplot(theo_mean);
48
49 %% Question 2.b
50 n_sample = 100;
51 sample_size = 20;
52 sample_2b = zeros(sample_size, n_sample);
53
54 for i=1:n_sample
55     sample_2b(:,i) = datasample(exam_exer(:,1), sample_size);
56 end
57
58 % Resultat de la question 1 (population)

```

```

59 ex1_mean = mean(exam_exer(:,1));
60 ex1_median = median(exam_exer(:,1));
61 ex1_std = mean(exam_exer(:,1));
62
63 % Question 2.b.i
64 ex1_mean_s = mean(sample_2b).';
65
66 figure; hist(ex1_mean_s);
67 title('Exercice 1 (examen) : moyenne de 100 echantillons');
68 ylabel('Nombre d\'echantillons');
69 xlabel('Moyenne (../20)');
70
71 mean_ex1_mean_s = mean(ex1_mean_s);
72
73 % Question 2.b.ii
74 ex1_median_s = median(sample_2b).';
75
76 figure; hist(ex1_median_s);
77 title('Exercice 1 (examen) : mediane de 100 echantillons');
78 ylabel('Nombre d\'echantillons');
79 xlabel('Mediane (../20)');
80
81 mean_ex1_median_s = mean(ex1_median_s);
82
83 % Question 2.b.iii
84 ex1_std_s = std(sample_2b).';
85
86 figure; hist(ex1_std_s);
87 title('Exercice 1 (examen) : ecart-type de 100 echantillons');
88 ylabel('Nombre d\'echantillons');
89 xlabel('Ecart-type');
90
91 mean_ex1_std_s = mean(ex1_std_s);
92
93 % Question 2.b.iv
94
95 ks_ex1 = zeros(100,1);
96
97 for i = 1:100
98     ks_ex1(i) = ksdist(exam_exer(:,1), sample_2b(:,i));
99 end
100
101 figure; hist(ks_ex1, 10);
102 title('Exercice 1 (examen) : distance de KS de 100 echantillons');
103 ylabel('Nombre d\'echantillons');
104 xlabel('Distance de Kolm.-Smir. ');
105
106 % Question 2.b.v
107
108 ks_ex2 = zeros(100,1);
109 ks_ex3 = zeros(100,1);
110
111 for i = 1:100
112     [~, ~, ks_ex2(i)] = kstest2(exam_exer(:,2), datasample(exam_exer(:,2),
113         sample_size));
114     [~, ~, ks_ex3(i)] = kstest2(exam_exer(:,3), datasample(exam_exer(:,3),
115         sample_size));

```

```

114 end
115
116 figure; hist(ks_ex2, 10);
117 title('Exercice 2 (examen) : distance de KS de 100 echantillons');
118 ylabel('Nombre d''echantillons');
119 xlabel('Distance de Kolm.-Smir.');
```

```

120
121 figure; hist(ks_ex3, 10);
122 title('Exercice 3 (examen) : distance de KS de 100 echantillons');
123 ylabel('Nombre d''echantillons');
124 xlabel('Distance de Kolm.-Smir.');
```

Listing 2 – Question2.m

```

1 function [sample index_perm] = datasample(data, k)
2     % datasample(data, k)
3     %     This function returns a random sample from the matrix
4     %     $data$.
5     %     It randomly selects $k$ lines in the matrix.
6     %     /\ $k$ must be less than the numbers of row in $data$
7     %
8     %     PARAMETERS :
9     %     - data : the matrix from which a sample must be
10    %     extracted
11    %     - k : the size of the sample
12    %
13    %     RETURN :
14    %     - sample : the sample extracted
15    %     - index_perm : the (row) indexes of the values
16    %     extracted in $data$
17
18    [len_x ~] = size(data);
19
20    if(k >= len_x)
21        error('The sample size ''k'' must be less than or equal to the number of
22        line in the matrix');
23    else
24        index_perm = randi(len_x, k, 1);
25        index_perm = sort(index_perm);
26        sample = data(index_perm,:);
27    end
28 end
```

Listing 3 – datasample.m