# SOMMAIRE

**MK**
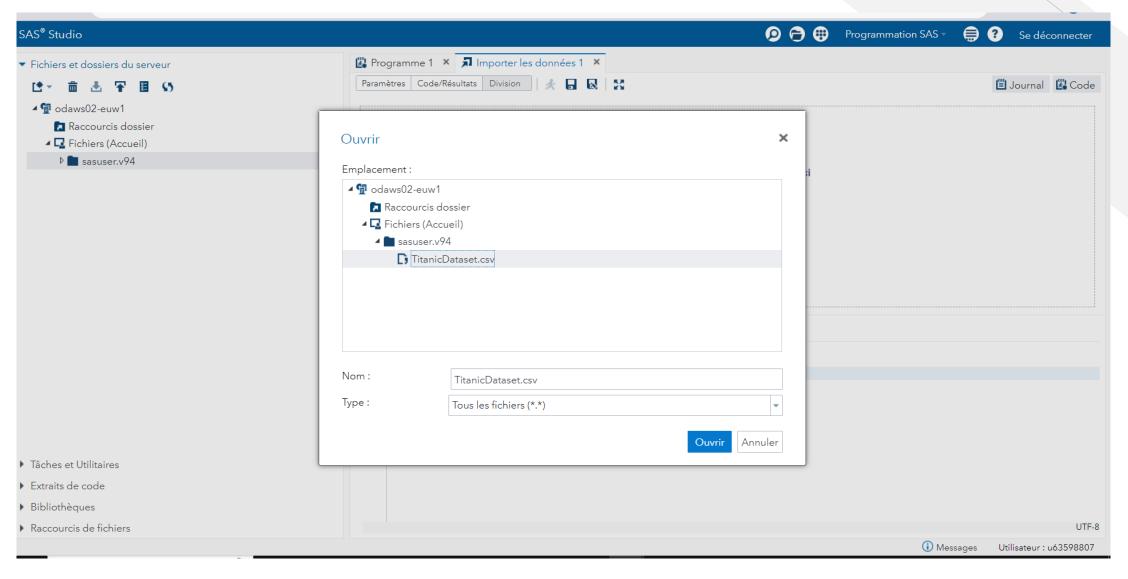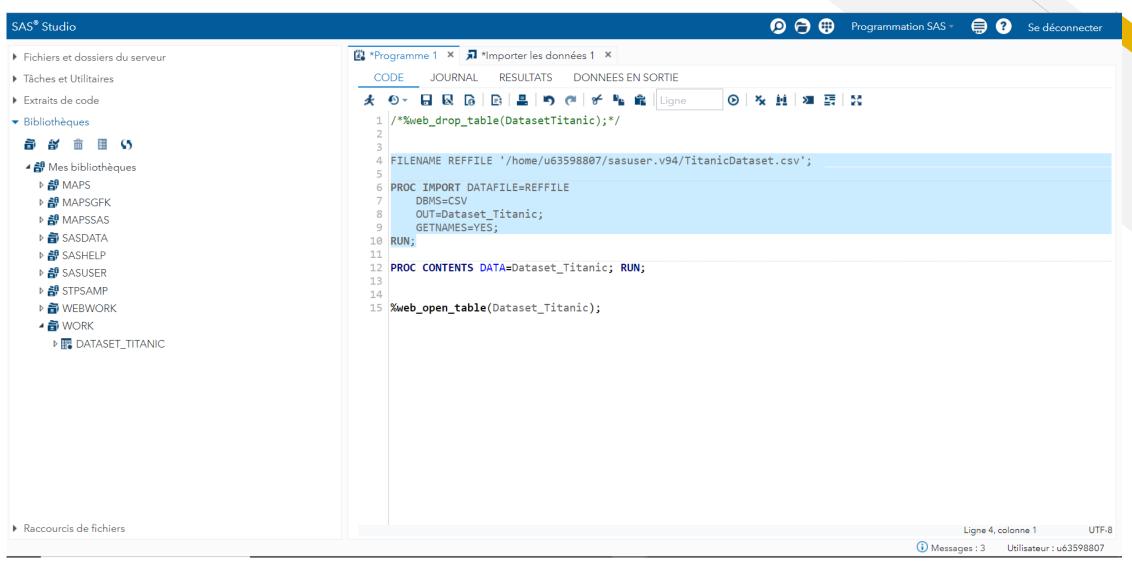
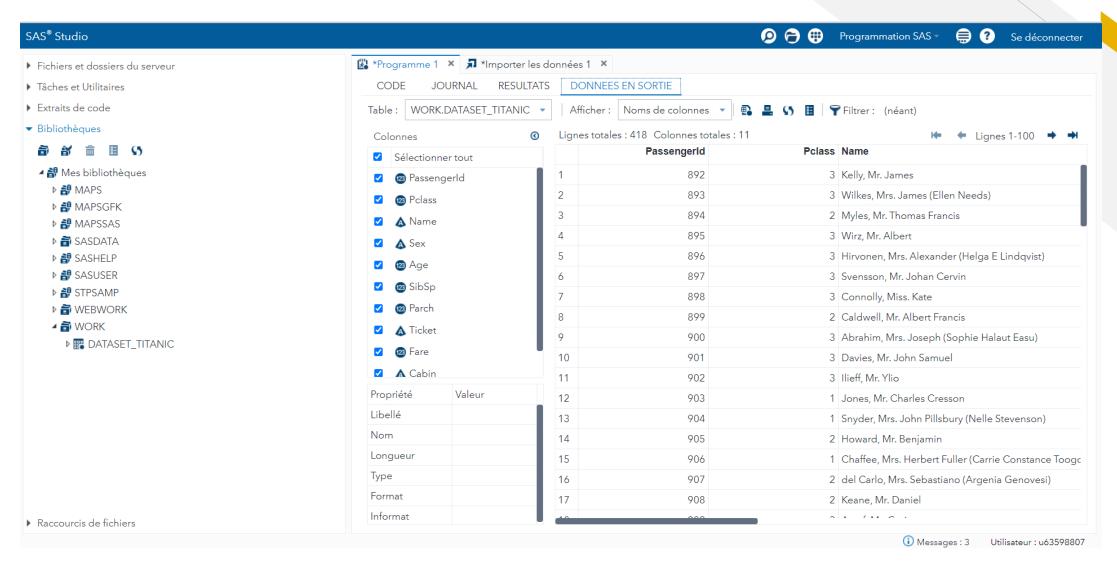# Importation des données

# Importation des données (1/3)

# Importation des données (2/3)

# Importation des données (3/3)

# Étude des données

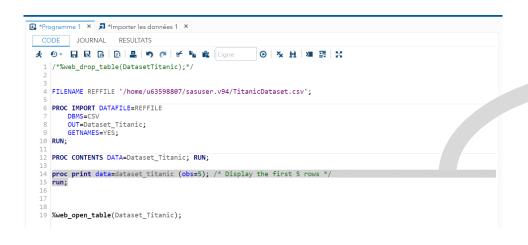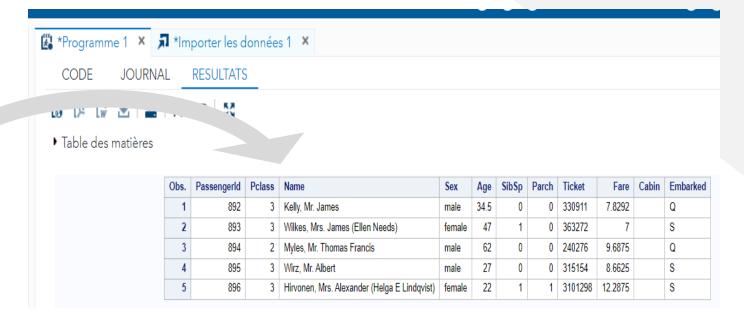# Étude des données (1/3)

**MK**

# Étude des données – Statistiques (2/3)

```
 6  PROC IMPORT DATAFILE=REFFILE
 7      DBMS=CSV
 8      OUT=Dataset_Titanic;
 9      GETNAMES=YES;
10  RUN;
11
12  PROC CONTENTS DATA=Dataset_Titanic; RUN;
13
14  /*First Few Rows of the Dataset*/
15  proc print data=dataset_titanic (obs=5); /* Display the first 5 r    */
16  run;
17
18  /* Summary Statistics*/
19  proc means data=dataset_titanic;
20      var age fare Pclass ;
21  run;
22
```
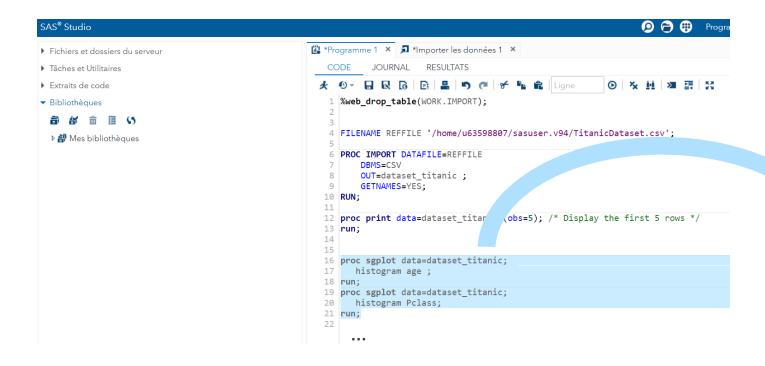
La procédure MEANS

| Variable | N | Moyenne | Ec-type | Minimum | Maximum |
|---|---|---|---|---|---|
| Age | 332 | 30.2725904 | 14.1812092 | 0.1700000 | 76.0000000 |
| Fare | 417 | 35.6271885 | 55.9075762 | 0 | 512.3292000 |
| Pclass | 418 | 2.2655502 | 0.8418376 | 1.0000000 | 3.0000000 |

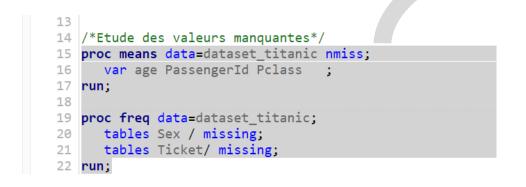# Étude des données – Distribution (2/3)

# Étude des données – Valeurs manquantes (2/3)

MK

# Étude des données – Corrélation des variables (3/3)

# Data Processing

# Data Processing – Gestion des Valeurs manquantes (méthode de la moyenne)

**MK**

La procédure MEANS

| Variable | Nbre manquant |
|----------|---------------|
| Age | 86 |
| PassengerId | 0 |
| Pclass | 0 |
| Gender | 0 |

La procédure MEANS

| Variable | Nbre manquant |
|----------|---------------|
| Age | 0 |
| PassengerId | 0 |
| Pclass | 0 |
| Gender | 0 |

*Importer les données 1 ✕   WORK.AGE_MEAN ✕   WORK.MEAN_AGE ✕

...nes ▾   | Filtrer : (néant)

Lignes totales : 1 Colonnes totales : 3

| | _TYPE_ | _FREQ_ | mean_age |
|---|--------|--------|----------|
| 1 | 0 | 418 | 30.272590361 |

```
74
75  /*Handling the Age missing values*/
76  /* Calculate the mean of Age */
77  proc means data=dataset_titanic mean noprint;
78     var Age;
79     output out=age_mean mean=mean_age;
80  run;
81  /*Storing the mean in a new column and in each row"*/
82  data dataset_titanic;
```

# Data Processing – Gestion des Valeurs (Encodage)

# Data Science - Exemple

**MK**

```
149
150
151   /* Train a logistic regression model and save it */
152   proc logistic data=train_titanic_dataset outmodel=your_trained_model;
153      /* Target variable: Survived (1 for survived, 0 for not survived) */
154      /* Definition of predictor variables */
155      model Survived(event='1') = Age Fare Pclass Gender;
156
157      /* Specify options (e.g., selection methods, interactions, etc.) */
158      /* selection=stepwise; */
159
160      /* Output the results, including parameter estimates and model fit statistics */
161      ods output ParameterEstimates=LogRegParams FitStatistics=ModelFitStats;
162   run;
163
164
```

# RÉGRESSION LOGISTIQUE – Evaluation du modèle

```
179  /* Load the scored dataset (containing predicted probabilities) */
180  data scored;
181      set scored; /* Replace with the actual name of your scored dataset */
182  run;
183
184  /* Calculate the Mean Squared Error (MSE) */
185  data squared_error;
186      set scored;
187      /* Calculate the squared error for each observation */
188      squared_error = (Survived - P_1) ** 2;
189  run;
190
191  /* Calculate the overall Mean Squared Error (MSE) */
192  proc means data=squared_error mean;
193      var squared_error;
194      output out=mse_results mean=MSE;
195  run;
196
```

## La procédure MEANS

| Variable d'analyse : squared_error |
| --- |
| **Moyenne** |
| 0.1429785 |