

Compte Rendu

Sonny Klotz - Jean-Didier Pailleux - Malek Zemni

*Interface de chargement, de contrôle
et d'analyse statistique des données
pour la constitution d'un graphe de flux*

27/05/2017



Module *Projet*

Table des matières

1	Architecture de l'application	1
1.1	Organigramme et données échangées	1
1.2	Format du fichier CSV	3
1.3	Fonctionnalités des modules	3
2	Contraintes et Langages de programmation	5
2.1	Contraintes	5
2.2	Choix langages et frameworks	6
3	Partie technique du projet	7
4	Points délicats	7
5	Organisation interne du groupe	7
6	Coût	8

Introduction

Ce document est le compte-rendu final de notre travail qui s'inscrit dans le cadre du module *Projet* de la licence informatique de l' **UVSQ**. Le sujet de ce projet a été proposé par l'entreprise **DCbrain**.

Notre client, **DCbrain**, développe des outils qui permettent de visualiser le comportement des **réseaux physiques** afin de prédire et trouver les problèmes de ces réseaux dans le but de les optimiser. Ces réseaux physiques sont principalement les réseaux industriels, de fluide ou de distribution, tels les réseaux électriques. Des données sont collectées à partir de ces réseaux puis analysées grâce aux technologies du **Big Data**. Le but est d'en extraire une valeur ajoutée.

Notre tâche correspond au travail indispensable d'analyse préliminaire de ces données, afin de faciliter leur exploitation. Ce projet a donc pour objectif de fournir aux utilisateurs une application web qui permettra de charger les données collectées, les visualiser et enfin les analyser.

Dans une première partie de ce document, on présentera l'architecture de notre application, illustrée par un organigramme qui a été préalablement établi. Cet organigramme représente un découpage de l'application en modules avec les différentes informations qui circulent entre ces modules.

Dans une deuxième partie, on parlera des différents langages de programmation choisis pour réaliser l'application et des contraintes qui ont justifié ces choix.

Ensuite, dans une troisième partie, on traitera la partie technique de notre projet, c'est à dire, l'application, son fonctionnement et ses problèmes.

Finalement, dans les deux dernières parties, on établira un bilan quant à l'organisation interne au sein du groupe et un bilan comparatif des coûts présumés et des coûts finaux de notre produit.

1 Architecture de l'application

1.1 Organigramme et données échangées

Cet organigramme a été préalablement établi dans le cahier des charges du produit. Il représente la décomposition en modules de l'application, ainsi que les informations qui circulent entre ces modules :

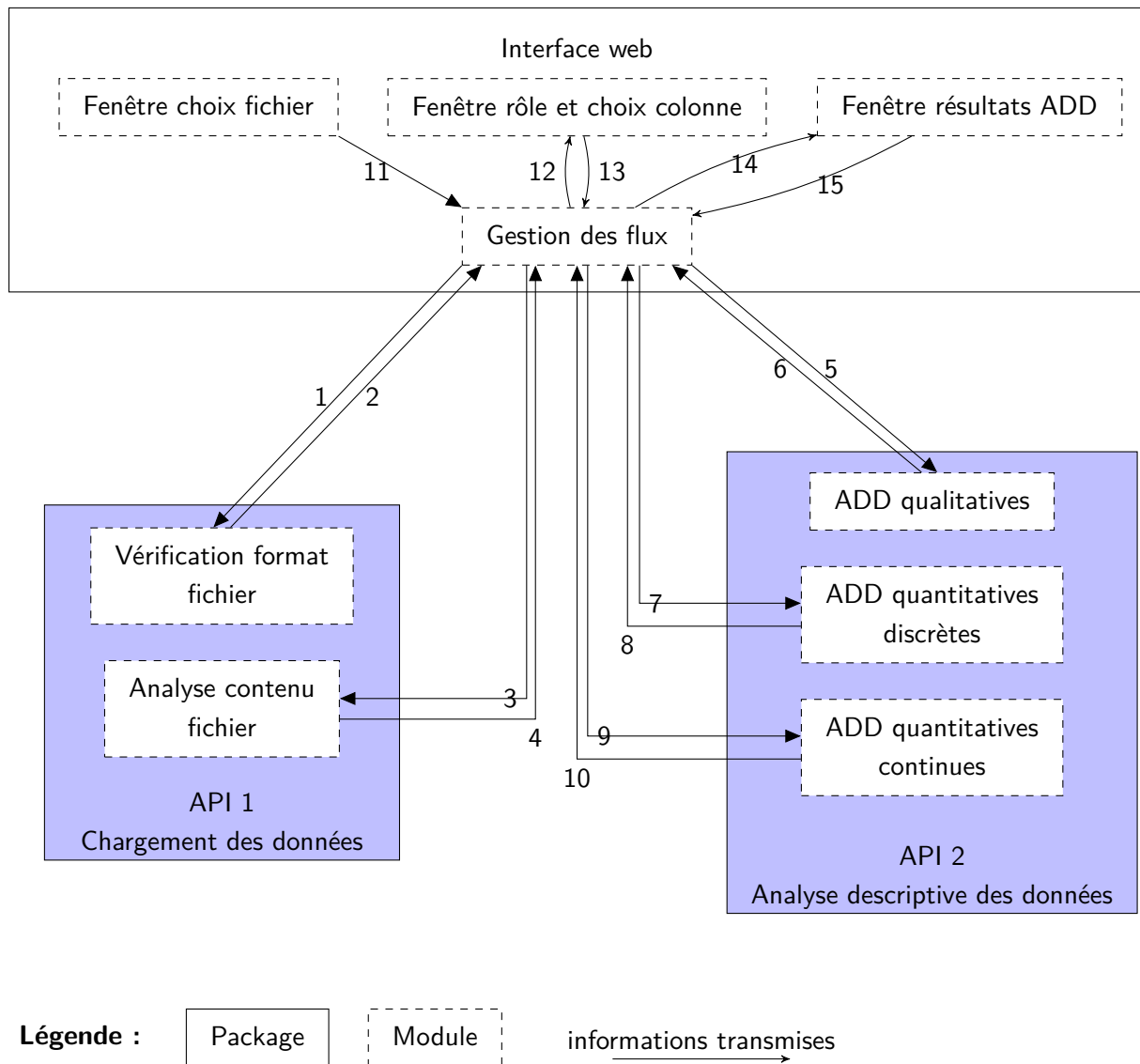


FIGURE 1 – Organigramme des différents modules du logiciel

Notes :

- (1) Fichier CSV : lancement de la vérification de son format
- (2) Code d'erreur : fichier OK ou ERREUR
- (3) Fichier CSV : lancement de l'analyse de son contenu
- (4) **structure 1** : contenant les données du fichier, le nombre de lignes et le nombre de colonnes (connus à partir de la taille de la structure)
- structure 2** : contenant 3 informations sur chaque colonne : le type, le rôle et les positions des données erronées ou manquantes
- (5) Ensemble de données de type qualitatif
- (6) Erreur ou effectifs, effectifs cumulés, fréquences, fréquences cumulées, diagramme en secteur,

histogramme

(7) Ensemble de données de type quantitatif discret

(8) Erreur ou indicateurs de tendance central, de dispersion et de forme, les anomalies, la distribution des données, un diagramme à moustaches

(9) Ensemble de données de type quantitatif continu

(10) Même données que (8)

(11) Chemin du fichier CSV importé

(12) Informations du (4) et ensemble de données contenu dans le fichier CSV

(13) Signal de validation du choix de la colonne, et noms des colonnes

(14) Envoi des résultats d'analyses de (6), (8) et (10)

(15) Signal de contrôle : demande d'exportation des résultats de l'ADD, analyse d'une autre colonne, ou importation d'un autre fichier

1.2 Format du fichier CSV

Le format du fichier a été établi par le client **DCbrain**. Son contenu est assez structuré et est décrit par des colonnes aux types prédéfinies :

timestamp	parent	enfant	mesure 1	mesure 2
January 1st 2017, 15 :00 :00.000	102	95	26644.235	176.253

Remarque : le graphe de flux utilisé par **DCbrain** pour analyser les réseaux de ses clients est orienté, d'où l'utilisation des nœuds *parent* - *enfant* numérotés. Ceux-ci permettent d'identifier une connexion précise.

Les colonnes *mesure* représentent des données mesurées sur une connexion à un temps donné. Le nombre de ces colonne n'est pas limité et en pratique il ne dépassera pas la dizaine.

1.3 Fonctionnalités des modules

1. API 1 - Chargement des données

1. Module Vérification format fichier :

- Vérification de la possibilité d'ouverture du fichier.
- Vérification du contenu : le fichier est un CSV contenant du texte brut non formaté. (pas de mise en forme avec des balises ou autres marqueurs)
- Vérification de l'accessibilité du fichier en lecture.

2. Module Analyse contenu fichier :

Ce module comprend deux fonctionnalités principales :

- Lecture du contenu du fichier CSV :
On initialise une première structure (**structure 1**) pour y sauvegarder le contenu du fichier.
On lit ligne par ligne des caractères du fichier. A chaque fois qu'on détecte un caractère de séparation (une virgule, un point-virgule ou une tabulation), on stocke les caractères lus (la donnée) dans la structure.
Cette fonctionnalité fournit la **structure 1** contenant les données du fichier, le nombre de lignes et le nombre de colonnes (connus à partir de la taille de la structure).
 - Descriptions préliminaires des données de chaque colonne du fichier CSV :
On initialise une deuxième structure (**structure 2**) pour y stocker des informations sur chaque colonne.
On lit une par une les données de chaque colonne. On vérifie le type de la donnée en le comparant au type attendu. Si le type ne correspond pas, on signale dans la structure que la colonne contient une donnée erronée ou manquante (repérée par sa position dans la colonne). A la fin du parcours d'une colonne, on pourra lui attribuer un rôle/nom.
Cette fonctionnalité fournit donc la **structure 2** contenant 3 informations sur chaque colonne : le type, le rôle et les positions des données erronées ou manquantes.
- Ce module va donc fournir les deux structures **structure 1** et **structure 2** décrites ci-dessus.

2. API 2 - Analyse descriptive de données

1. Module ADD qualitatives :
 - Calcul des effectifs, effectifs cumulés, fréquences et fréquences cumulées des données
 - Préparation des représentations graphiques : diagramme en secteur et histogramme
2. Module ADD quantitatives discrètes :
 - Statistiques : moyenne, quantiles, variance, écart-type, coefficients de symétrie et d'aplatissement de Fisher
 - Anomalies de Tukey : toute donnée hors de l'intervalle $[Q1 - 1.5*IQ; Q3 + 1.5*IQ]$, où $Q1$ et $Q3$ sont les quartiles, et IQ l'écart inter-quartiles.
 - Préparation des représentations graphiques : fonction de distribution, fonction de distribution cumulative, boîte à moustaches
3. Module ADD quantitatives continues :
 - Discrétisation de l'étendue (découpage en classe d'intervalles)
 - Calcul des quantiles pour le cas de données continues
 - Préparation des représentations graphiques : fonction de distribution cumulative

3. Interface web

1. Module Gestion des flux :
 - Gestion des branchements : exécution normale ou arrêts pour cause d'erreur.
 - Interface entre les différentes fonctionnalités : communique les données nécessaires entre les modules.
2. Module Fenêtre choix fichier :
 - Récupération un fichier CSV en renseignant son chemin en parcourant l'arborescence de fichiers, ou de la manière d'un Drag & Drop.
3. Module Fenêtre rôle et choix colonne :
 - Affichera le nombre de lignes/colonnes contenu dans le CSV.
 - Affichera le titre du fichier.
 - Affichera un échantillon du contenu du CSV (environ les 1000 premières lignes) avec un système de scroll.
 - Affichage des lignes erronées (numéro de la ligne + contenu + type d'erreur).
 - Mise en place d'un système de navigation sous forme d'onglet (Onglet erreurs, onglet échantillon,...). Cela permettra d'éviter que la fenêtre contienne trop d'informations.
 - L'utilisateur devra sélectionner la colonne avec un clic, puis pourra lancer l'analyse sur celle-ci.
4. Module Fenêtre résultats ADD :
 - Affichage des résultats d'analyse descriptive : informations statistiques de l'API 2 et représentations graphiques pour visualiser les données dans leur ensemble.
 - Une fonctionnalité de retour en arrière permet de sélectionner une nouvelle colonne sans relancer l'applet.

2 Contraintes et Langages de programmation

2.1 Contraintes

Le développement de l'application s'accompagne de contraintes imposées. Celles-ci établie par DCbrain au commencement du projet. On peut en énumérer trois, qui sont les suivantes :

Contrainte	Fiche
------------	-------

1. Le produit doit fournir une application web	<p>Description : notre produit sera une application fonctionnant sur un navigateur web, appelée applet.</p> <p>Justification : assure une très grande portabilité et fournit à l'utilisateur une interface interactive.</p> <p>Critère de satisfaction : on peut lancer l'application sur un navigateur web.</p>
2. Le produit doit être développé avec un langage de programmation compatible avec l'analyse de données	<p>Description : le langage de programmation choisi doit inclure des bibliothèques qui permettent d'analyser les données (tâches de data mining et de machine Learning).</p> <p>Justification : permettre une analyse de données efficace.</p> <p>Critère de satisfaction : on peut effectuer une analyse descriptive de données.</p>
3. Le produit doit fournir une API d'analyse de données en sortie	<p>Description : l'application doit intégrer des API d'analyse descriptive de données qui pourront être livrées en sortie au client.</p> <p>Justification : permettre une réutilisabilité des fonctionnalités majeures du produit.</p> <p>Critère de satisfaction : on peut exporter une API d'analyse de données en sortie.</p>

2.2 Choix langages et frameworks

Les contraintes appliquées au projet ainsi que les exigences définies nous ont permis de fixer notre choix sur un langage de programmation : **Python**.

Ce choix est justifié par le fait que **Python** soit compatible pour l'Analyse Descriptives de Données, puisqu'il est doté de nombreux modules permettant d'effectuer du calcul scientifique sur des données, donc l'analyse de données. De même, il est compatible avec le développement d'applications web, grâce aux nombreux frameworks disponibles.

Pour ce dernier, notre choix s'est porté sur **Flask**, qui s'efforce d'être le plus simple possible pour une prise en main et le plus petit possible. À l'opposé de Django ou autres, il est donc plus abordable pour une première approche sur le développement d'applications web.

L'application devant être une applet intégrée dans une page web pour s'exécuter, le recours à des langages de balisage comme **HTML** pour la présentation et **CSS** pour la mise en forme de ces pages web étaient nécessaires.

Ensuite, pour permettre à l'application de faire communiquer les pages web avec un serveur web sans occasionner le rechargement de la page, nous avons eu recours à **Ajax**. Il s'agit d'un concept de programmation web reposant sur plusieurs technologies comme le **Javascript**, qui lui se charge d'établir la connexion entre la page web et serveur. **Javascript** est également employé pour rendre les pages web interactives. **Jquery** est quant à lui utilisé pour faire le lien entre **HTML** et **Ajax**.

Enfin, pour la génération de la documentation des fonctionnalités du programme, on utilisera le logiciel **Sphinx** qui se charge de générer la docstring à partir des commentaires présents dans le code python.

3 Partie technique du projet

Quoi marche (comment ça marche - à l'aide de la partie avant) QUoi marche pas (pb et et resolution si trouvee)

1 page, plus tard

4 Points délicats

Justification des points qui divergent par rapport aux specs et organigramme de départ

1 page, plus tard à moins que t'ai déjà des idées - (Peut-être rajouté les points délicats de flask/application web sur l'approche et tout ce qui va avec ?

5 Organisation interne du groupe

Ci-dessous se trouve le tableau de la répartition des modules pour chaque membre du groupe établie avant le commencement du développement :

Module	Malek	Sonny	Jean-Didier	Total
Gestion des flux			x	1
Fenêtre choix fenêtre			x	1
Fenêtre rôle et choix colonne	x			1
Fenêtre résultats ADD		x		1

ADD qualitatives			x	1
ADD quantitatives discrètes		x		1
ADD quantitative continues		x		1
Vérification format fichier	x			1
Analyse contenu fichier	x			1

1 page On expliquera plus tard comment ça s'est déroulé concrètement par rapport à ce qui était prévu.

6 Coût

Tableau indiquant l'estimation et le coût final en nombre de lignes de code, pour chaque module :

Module	Nombre de lignes	Justification	Coût final
Gestion des flux	15	Mise en forme du main et appel de l'application	
Fenêtre choix fichier	10 + 20	Fonctions pour : Drag& Drop + Système de fichiers	
Fenêtre rôle et choix colonne	5 + 20 + 10 + 10 + 20	Communication avec le module application + Affichage de l'interface+ lecteurs des valeurs + affichage des valeurs	
Fenêtre résultats ADD	10 + 3* 30	Envoie d'informations au module application + Construction des graphe pour l'ADD pour les 3 types d'analyse	
ADD qualitatives	20 + 20*3	Application des formules pour les calculs de fréquences et d'effectifs + calcul des valeurs pour la construction de 3 graphes	
ADD quantitatives discrètes	60 + 20*2	Application des formules attaché à l'analyse quantitative discret + calcul des valeurs pour la construction de 2 graphes	
ADD quantitatives continu	20 + 10 + 10 + 5 + 20*2	Parcours + choix précision classe d'intervalle + écriture + communication avec les modules+ calcul des valeurs pour la construction de 2 graphes	
Vérification format fichier	30	Ouverture fichier + vérification si ouverture en lecture + présence de texte formaté ou non	

Analyse contenu fichier	20 + 5 + 25 + 10	Recopie et vérification + initialisation de la structure+ Parcours du fichiers avec condition + Fonction pour donner nom et type de colonne	
Coût Total	565	Estimation totale du coût	

Conclusion

1 demi-page, plus tard