

Table des matières

1	Introduction	2
---	--------------	---

1 Introduction

Lors de la création d'un dataset à partir de différentes sources, l'un des principaux problèmes consiste à rendre les données suffisamment homogènes - le niveau de suffisance dépendant des objectifs pour lesquels le dataset est construit. Dans le cas des données textuelles, cela s'applique non seulement à un certain nombre d'éléments quantitatifs ou facilement reconnaissables – longueur du texte, langue, etc. – mais aussi à leur propre signification.

Les pages suivantes présentent une approche possible du *pre-processing* des données textuelles collectées à partir de diverses sources, qui permet également de vérifier l'homogénéité du sens des textes contenus dans l'ensemble de données.