# PyPandas: A scalable data cleaning library

Pei-Lun Liao
New York University
pll273@nyu.edu

Chia-Hsien Lin
New York University
chl566@nyu.edu

Shang-Hung Tsai
New York University
st3127@nyu.edu

## ABSTRACT

This is abstract. This is abstract. This is abstract. This is abstract. This is abstract. This is abstract.

## 1 INTRODUCTION

### 1.1 Problem description and goal

In recent days, as storage devices become cheaper, storing big data in thousands of computers become easier. At that time people start to think how to mine big value from the huge dataset. The term, Big Data, was created to describe the phenomenon. As the result, data mining [5] [13] and machine learning get popular in recent years. In these fields the quality of data is the key point to mine good value in the huge dataset. Hence, data cleaning become a challenge in the first step of data management and analysis [2] [6] [14].

To process huge dataset efficiently, distributed systems and frameworks [3] [4] [16] are introduced. Spark [17] is one of the popular open source project used in industry and academia. It is built on Hadoop [16] and provide a way to manage distributed memory. PySpark [10] is an extended library for Python users to use Spark.

In the Python machine learning and data mining environment, Pandas [9] is the most popular data management library. Pandas dataframe provides a way to collect data and the data can be transform into Python primitive data structures or Numpy [11] array easily. Moreover, Pandas provides several data management methods such as missing value filling, data indexing and data profiling.

However, Pandas can only be used in a single machine and could not be scaled to manage huge dataset. In this paper, we propose PyPandas, a library built on PySpark, which support basic data management features, missing value handling or replication removing. Moreover, we add more data processing features into our library to provide user friendly functions such as outlier detection and numerical data scaling. These features are important for the machine learning tasks and are supported in Scikit-Learn [12], a mainstream machine learning library.

Furthermore, text data is another common data type in our storage. Text data are hard to clean because of multiple languages, new word, typo, abbreviation, url and emoji etc. Although we have NLTK [1], a popular text cleaning library, there is no scalable text cleaning library. In this paper, we also provide tools for users to clean text easily.

The features of our library is summarized below.

- Data management library runs on Spark
- Support missing value handling, replicated data removing, outlier detection and numerical data scaling
- Text processing such as url detection, punctuation removing, pattern searching and replacement.

## 2 RELATED WORK

### 2.1 Previous work

There exists a number of open-source library for data cleaning and parallel data frames. Optimus [7] is framework that can perform distributed data cleaning and preprocessing. This framework works well with Spark and its DataFrame can scale in big cluster. It comes with helpful tools such as removing special characters and replacing null values. Dask [15] is another open-source library that supports parallel analytic computing. It provides scalable parallel data structures that extend interfaces like pandas. At the same time, it offers low latency and high responsiveness. SparklingPandas[8] library attempts to combine the power of spark and pandas to scale data analysis. It provides a Pandas-like API that is built using Spark?s DataFrame class. Unfortunately, it only support spark v1.4 and Python 2.7, and its development has ended.

## 3 PROPOSED METHODS

## 4 EXPERIMENT

### 4.1 Dataset

This is the dataset we will use. This is the dataset we will use. This is the dataset we will use. This is the dataset we will use. This is the dataset we will use. This is the dataset we will use.

### 4.2 Evaluation

We will evaluate our project by comparing it against other existing open-source libraries, including Optimus[7], Dask[15], and SparkingPandas[8]. We will compare our data cleaning features with Optimus, particularly those frequently used functionality in data science such as null value handling, duplicate detection, etc. Furthermore, we will evaluate the performance of our implementation of parallel data management by measuring the time and space consumed because it is important to have an efficient implementation to make this library useful in production. We will also examine and compare the overall software architecture of the project. We might perform a survey to assess the usability of the library.

## 5 CONCLUSIONS

This is our conclusion. This is our conclusion. This is our conclusion. This is our conclusion. This is our conclusion. This is our conclusion.

## 6 REFERENCE

## REFERENCES

[1] Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python* (1st ed.). O'Reilly Media, Inc.

[2] Xu Chu, Ihab F. Ilyas, Sanjay Krishnan, and Jiannan Wang. 2016. Data Cleaning: Overview and Emerging Challenges. In *Proceedings of the 2016 International Conference on Management of Data (SIGMOD '16)*. ACM, New York, NY, USA, 2201–2206. https://doi.org/10.1145/2882903.2912574

[3] Jeffrey Dean and Sanjay Ghemawat. 2008. MapReduce: Simplified Data Processing on Large Clusters. *Commun. ACM* 51, 1 (Jan. 2008), 107–113. https://doi.org/10.1145/1327452.1327492

[4] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. 2003. The Google File System. In *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles (SOSP '03)*. ACM, New York, NY, USA, 29–43. https://doi.org/10.1145/945445.945450

[5] Jiawei Han, Micheline Kamber, and Jian Pei. 2011. *Data Mining: Concepts and Techniques* (3rd ed.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.

[6] Ihab F. Ilyas and Xu Chu. 2015. Trends in Cleaning Relational Data: Consistency and Deduplication. *Foundations and TrendsÂő in Databases* 5, 4 (2015), 281–393. https://doi.org/10.1561/1900000045

[7] Iron. 2018. Optimus. (2018). https://github.com/ironmussa/Optimus

[8] Holden Karau and Juliet Hougland. 2015. SparklingPandas. (2015). https://github.com/sparklingpandas/sparklingpandas

[9] Wes McKinney. [n. d.]. pandas: a Foundational Python Library for Data Analysis and Statistics. ([n. d.]).

[10] Amit Nandi. 2015. *Spark for Python Developers.* Packt Publishing.

[11] Travis E. Oliphant. 2015. *Guide to NumPy* (2nd ed.). CreateSpace Independent Publishing Platform, USA.

[12] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.

[13] Anand Rajaraman and Jeffrey David Ullman. 2011. *Mining of Massive Datasets.* Cambridge University Press, New York, NY, USA.

[14] Vijayshankar Raman and Joseph M. Hellerstein. 2001. Potter's Wheel: An Interactive Data Cleaning System. In *Proceedings of the 27th International Conference on Very Large Data Bases (VLDB '01)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 381–390. http://dl.acm.org/citation.cfm?id=645927.672045

[15] Matthew Rocklin. 2015. Dask: Parallel Computation with Blocked algorithms and Task Scheduling. In *Proceedings of the 14th Python in Science Conference*, Kathryn Huff and James Bergstra (Eds.). 130 – 136.

[16] Tom White. 2009. *Hadoop: The Definitive Guide* (1st ed.). O'Reilly Media, Inc.

[17] Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2010. Spark: Cluster Computing with Working Sets. In *Proceedings of the 2Nd USENIX Conference on Hot Topics in Cloud Computing (HotCloud'10)*. USENIX Association, Berkeley, CA, USA, 10–10. http://dl.acm.org/citation.cfm?id=1863103.1863113