**Researchers:**
Carlos Mougan
David Masip  - realxdata
Jordi Nin  - ESADE
Oriol Pujol - University of Barcelona

# Quantile Encoder:
Tackling high cardinality categorical features in supervised learning

# Goals of the presentation

- Overview of a Machine Learning **pipeline**.

- Review of the most **common techniques to handle categorical data**.

- **Sktools** library

- State of the art on **features with high cardinality**.

- Introduction to **Quantile Encoder**.

# Quantile Encoder

## Tackling High-cardinality Categorical Features in Regression Problems

**Carlos Mougan**[†] · **David Masip** [†]
**Jordi Nin** · **Oriol Pujol**

**Abstract** In this paper, we provide an analyisis, an implementation, and the results of tackling high cardinality categorical features in tabular data regression datasets with the quantile. The Quantile Encoder outperforms in a consistent way the traditional statistical mean target encoder. To deal with the overfitting for categories with few examples, the Quantile Encoder can benefit from shrinkage in order to avoid it. We give empirical evidence on public datasets of the achievements of this method against state of the art statistical encoding techniques. We also provide support for which metrics yield better results and provide a quantitative analysis of the results. Finally, we create a set of features with different quantiles that provide more information about the categorical feature in question making a performance boost of the models.
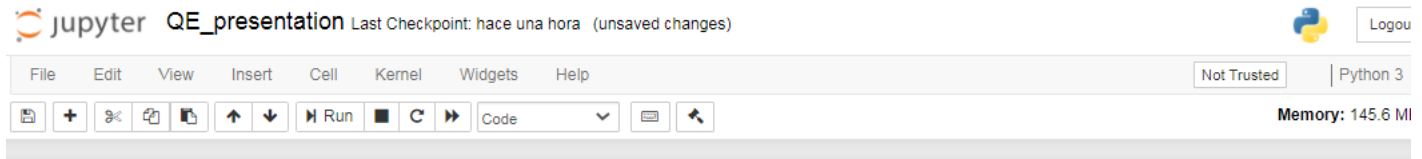
[†] - These two authors contributed equally

Oriol Pujol
Department of Mathematics and Computer Science - Universitat de Barcelona, Spain
E-mail: oriol_pujol@ub.edu

Jordi Nin
Department of Operations, Innovation and Data Sciences - Universitat Ramon Llull, ESADE
E-mail: jordi.nin@esade.edu

- In review for **The European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases –** ECML PKDD 2021

- **Paper** at https://darwin.escb.eu/livelink/livelink/app/nodes/306435319

- **Sktools** library**:** https://sktools.readthedocs.io/

# Presentation on Github

# Presentation takeaways

- How to **deal with categorical features** in supervised learning

- Machine Learning **pipeline** example

- **Sktools & category encoders** python libraries

# Paper contributions

- **Encoding**: Quantile Encoder

- **Optimization**: Not all encodings are optimal for all metrics and loss functions

- **Feature engineering**: Set of features quantile features (Summary Encoder)

# Finish

Feel free to drop any questions