# Adding Custom Cluster and Edit Functions

## INTRODUCTION

Occasionally custom Cluster + Edit functions are required to achieve the desired clustering results (typically when there are very specific requirements). This document outlines the high-level steps to integrate a new function.

Screenshots below.
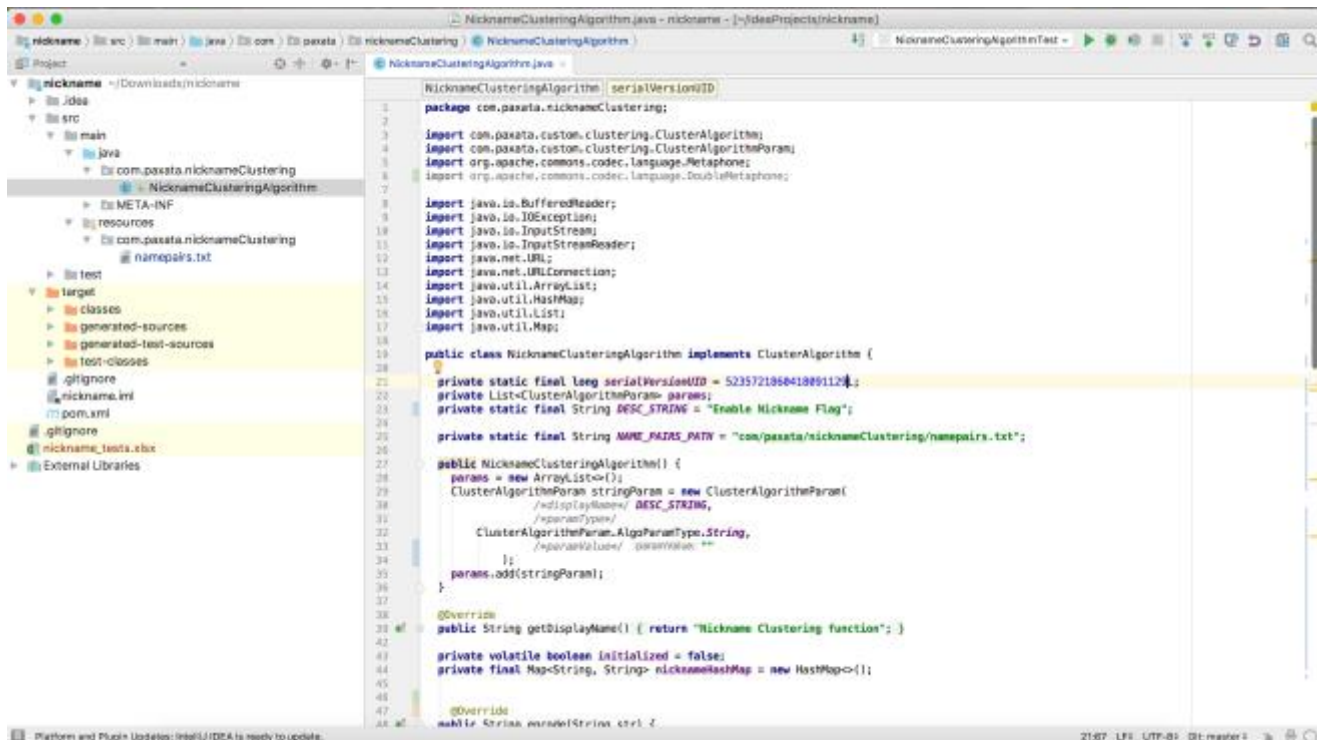
Questions – callum@paxata.com

When a *Project* is created, one of the functions I can perform is a **Cluster + Edit**, as the screenshot shows. Out of the box, I have 3 strategies I can use to perform this (**fingerprint, metaphone** and **ngram**):

You can use a different strategy, by including a custom **Cluster + Edit** function. This is done by writing java code to create a custom JAR (Java ARchive) file. (Paxata provides the sample code on how to do this). Below is an example of some custom code:



Once this is compiled and a JAR file has been created we need to perform 3 things to make this active:

1) Copy the JAR file (in this case "nickname.jar" to the correct location on the pipeline server). It is recommended that we create a new directory called "custom" for our custom cluster + edit JAR files:

2) Navigate the *clustering-algorithms.properties* file to tell Paxata to add a new custom **cluster + edit** function:



Then add a single line in the *clustering-algorithms.properties* file (this will depend on how the JAR file has been packaged). The text I am using is

*"cluster.nickname=com.paxata.nicknameClustering.NicknameClusteringAlgorithm"* (as displayed below):



3) Finally, restart the pipeline and the core server by typing *"service paxata-pipeline restart",* followed by *"service paxata-server restart":*

This will add a new **Cluster + Edit** algorithm to use in your Paxata installation (in this case called *"nickname"*):

Companies around the globe rely on Paxata to get smart about information. Paxata is the pioneer that intelligently empowers all business consumers to transform raw data into ready information, instantly and automatically, with an enterprise-grade, self-service data preparation application and machine learning platform. Our Adaptive Information Platform weaves data into an information fabric from any source and any cloud to create trusted insights. Business consumers use clicks, not code to achieve results in minutes, not months. With Paxata, Be an Information Inspired Business.

Paxata is headquartered in Redwood City, California with offices in New York, Ohio, Washington D.C., and Singapore.

**Paxata**®