



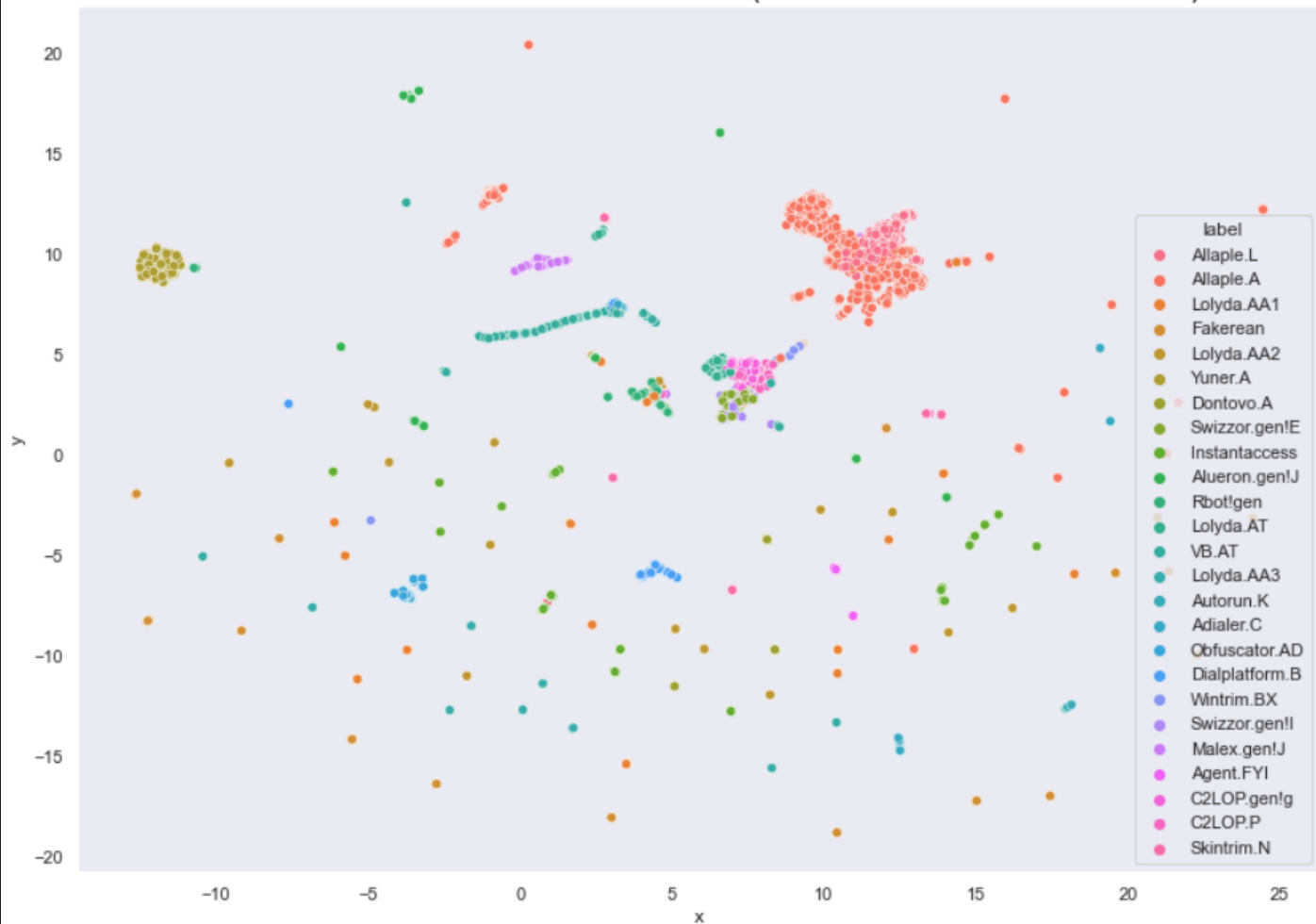
# TP4

Gabriel Semorile

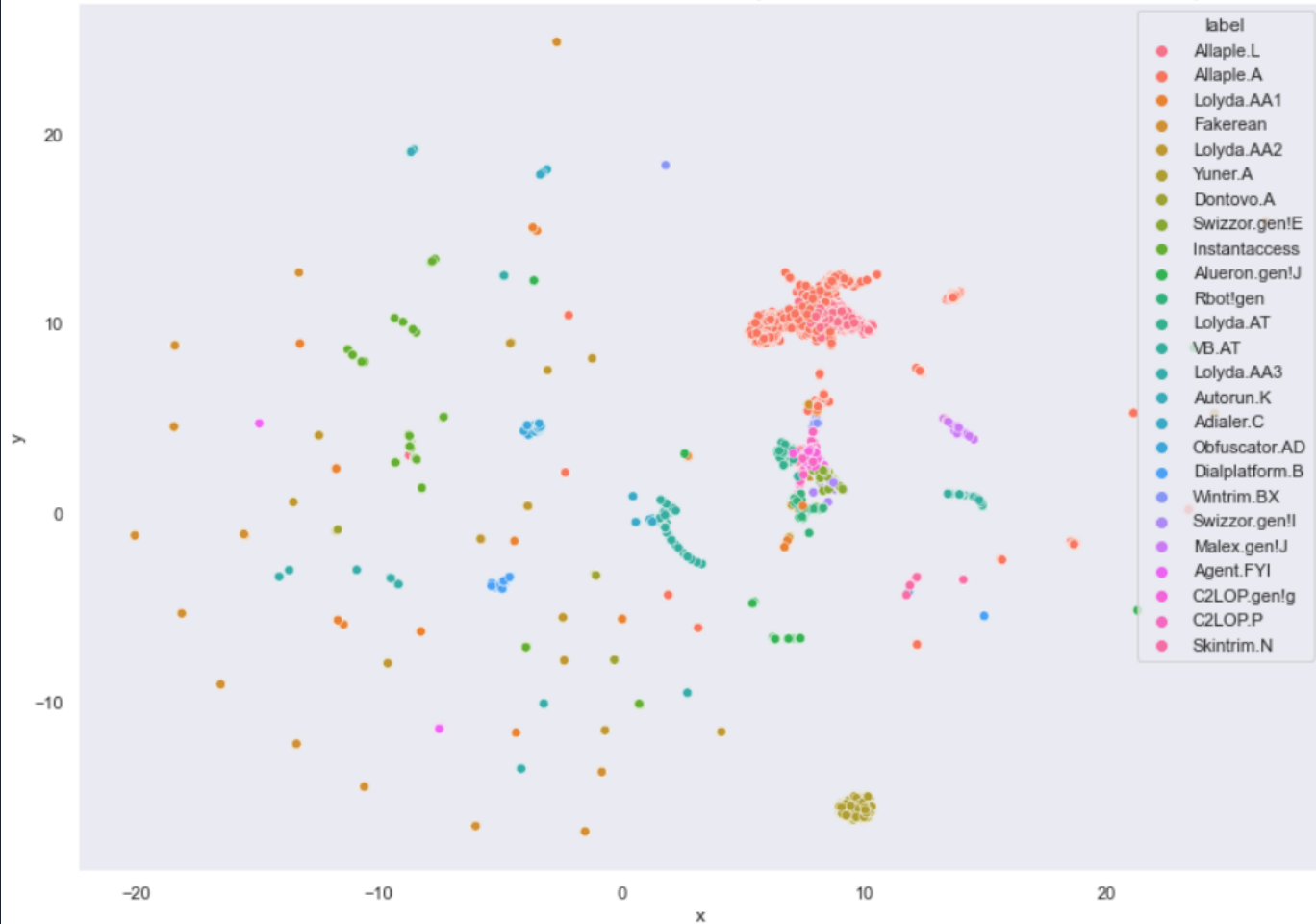
# Preparación del Dataset

- 1) Separo de forma aleatoria las imagenes entre train, validation y test.
- 2) Redimensiono las imágenes.
- 3) Utilizo UMAP para la reducción de dimensiones previa al entrenamiento.
- 4) Primero reduzco a 2 dimensiones para graficar.

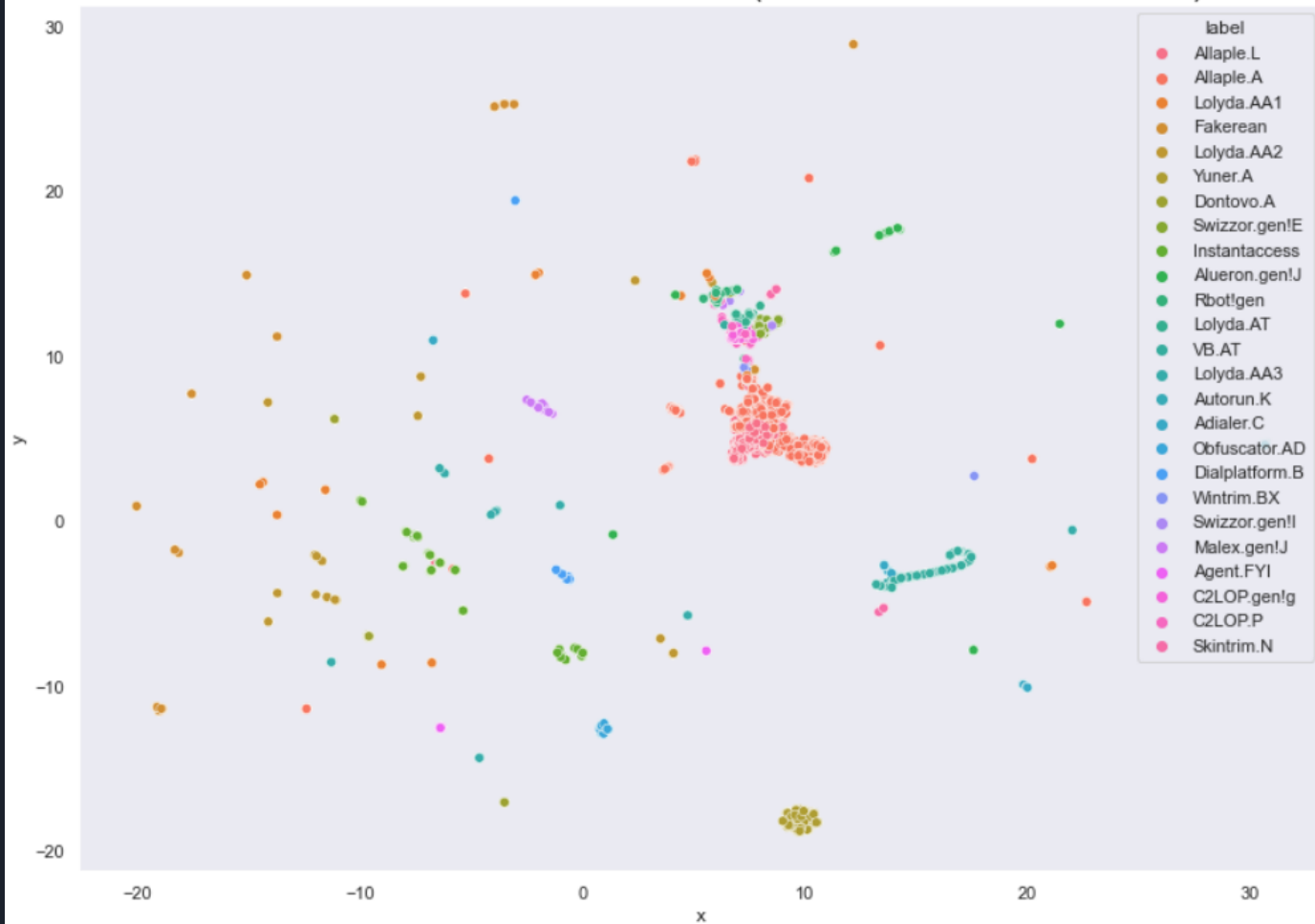
Reducción de dimensiones con UMAP (cant de vecinos cercanos = 10)



Reducción de dimensiones con UMAP (cant de vecinos cercanos = 15)



Reducción de dimensiones con UMAP (cant de vecinos cercanos = 20)



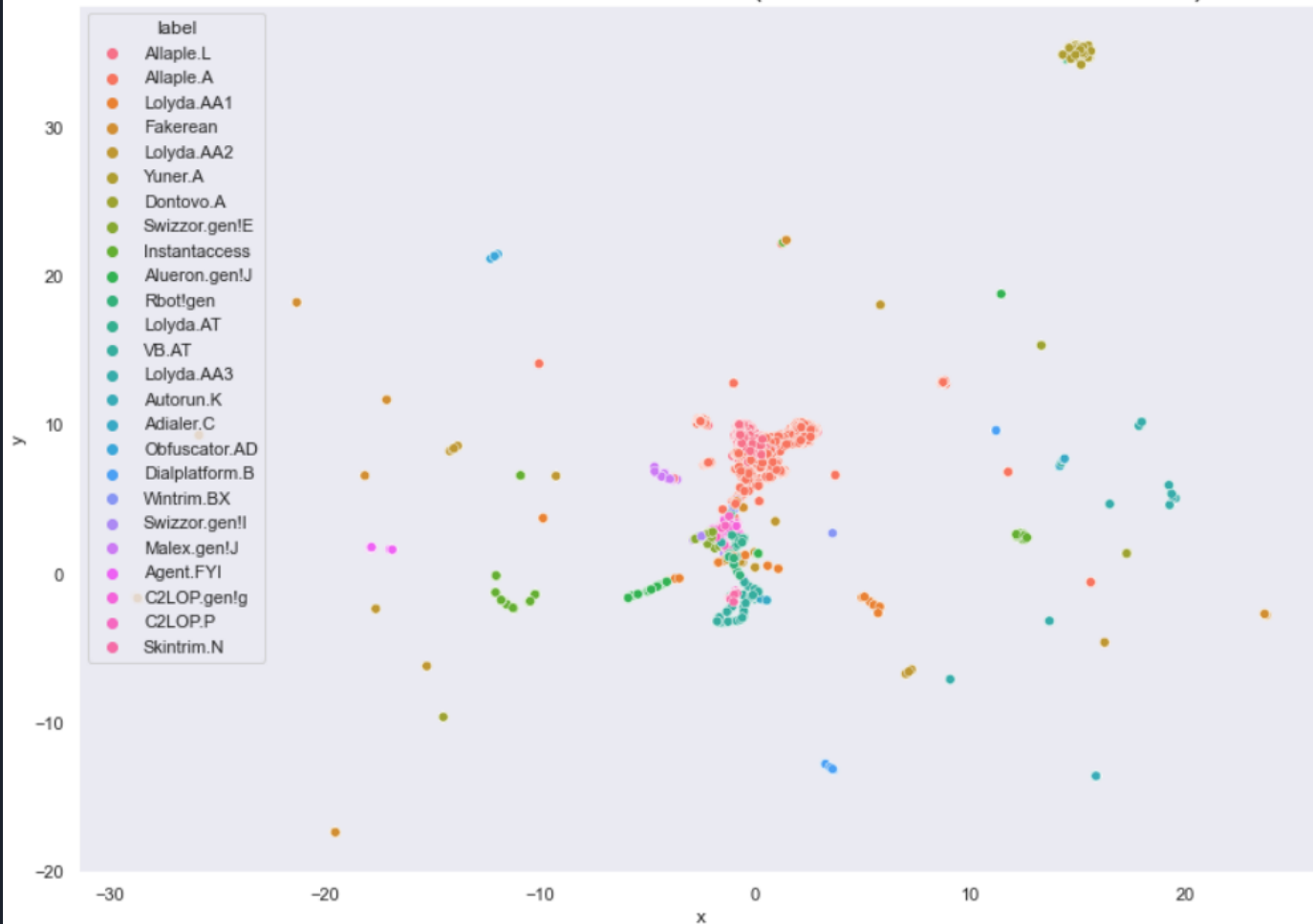
## Reducción de dimensiones con UMAP (cant de vecinos cercanos = 30)



Reducción de dimensiones con UMAP (cant de vecinos cercanos = 40)

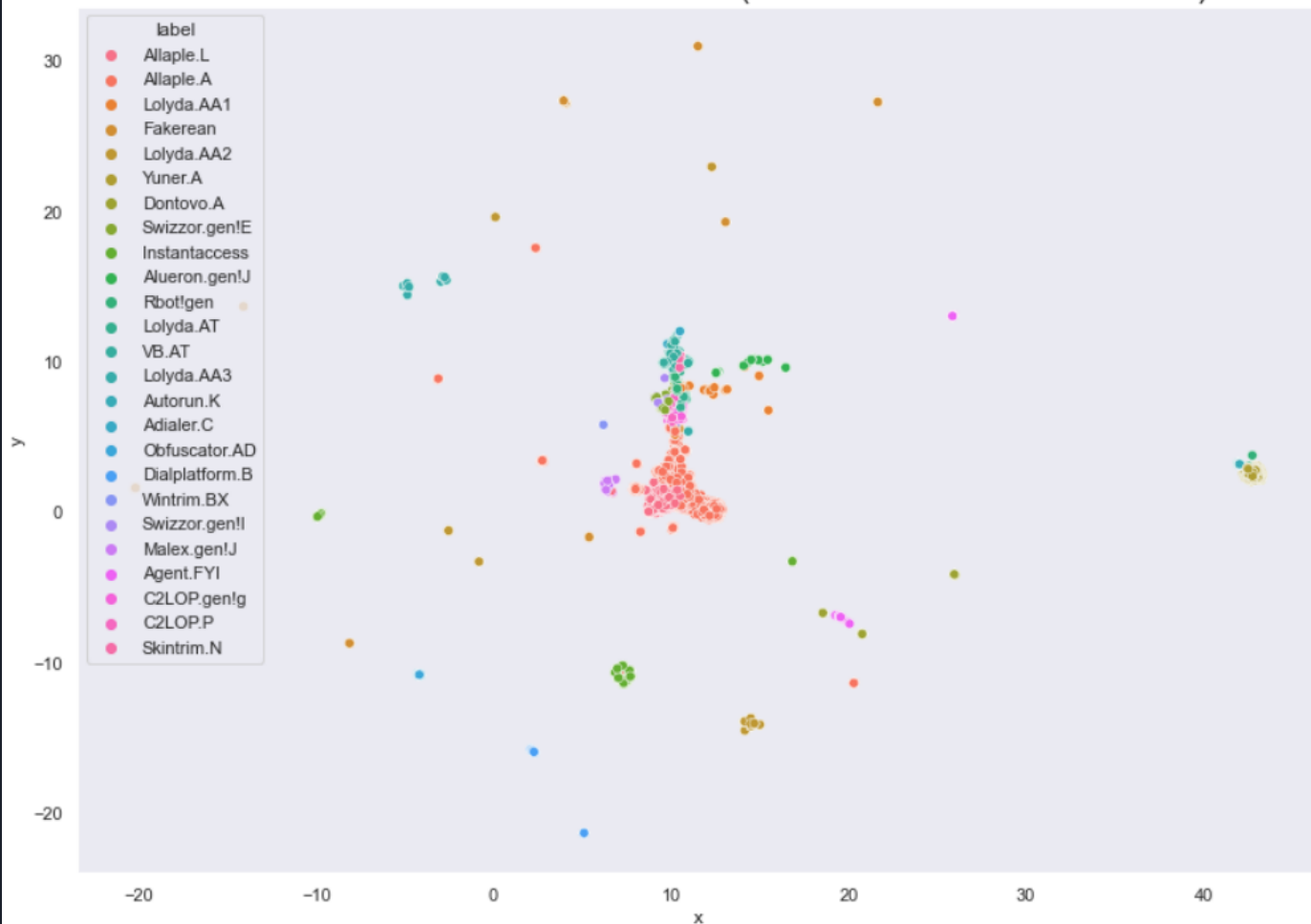


## Reducción de dimensiones con UMAP (cant de vecinos cercanos = 45)





# Reducción de dimensiones con UMAP (cant de vecinos cercanos = 100)



# Baseline con Kmeans

- Estandarizo la data para entrenar con StandardScaler
- Utilizo UMAP para reducir a 300 dimensiones.
- Entreno Kmeans con  $K = N$  (siendo  $N$  la cantidad de clases de virus).

Obtengo los siguientes Scores...

Datasets \ Scores	Purity Score	NMI
Validation	0.6323	0.7058
Test	0.6420	0.7005

# Métricas a utilizar para Clustering

## Purity:

Se basa en calcular para cada punto que proporción de los compañeros del cluster son del mismo tipo.

## Normalized Mutual Info (NMI):

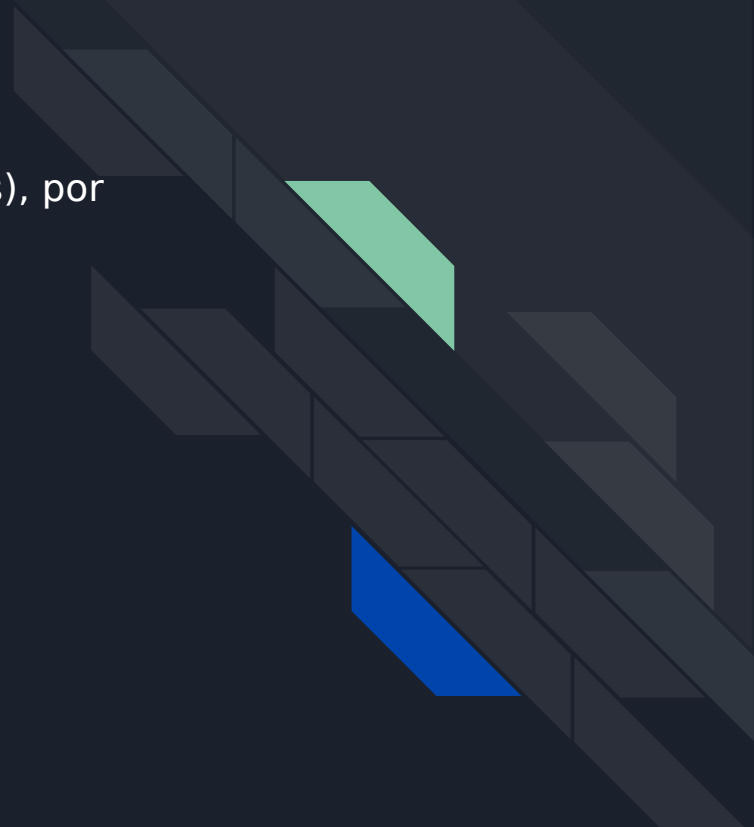
Se basa en normalizar Mutual Info, el cual es 0 si el clustering es aleatorio respecto a la clase de cada punto.

Características \ Scores	Purity Score	NMI
Ventajas	Muy Simple	Balance entre cantidad de clusters y calidad del clustering.
Desventajas	El score aumenta a medida que aumenta la cantidad de clusters.	Más Complejo

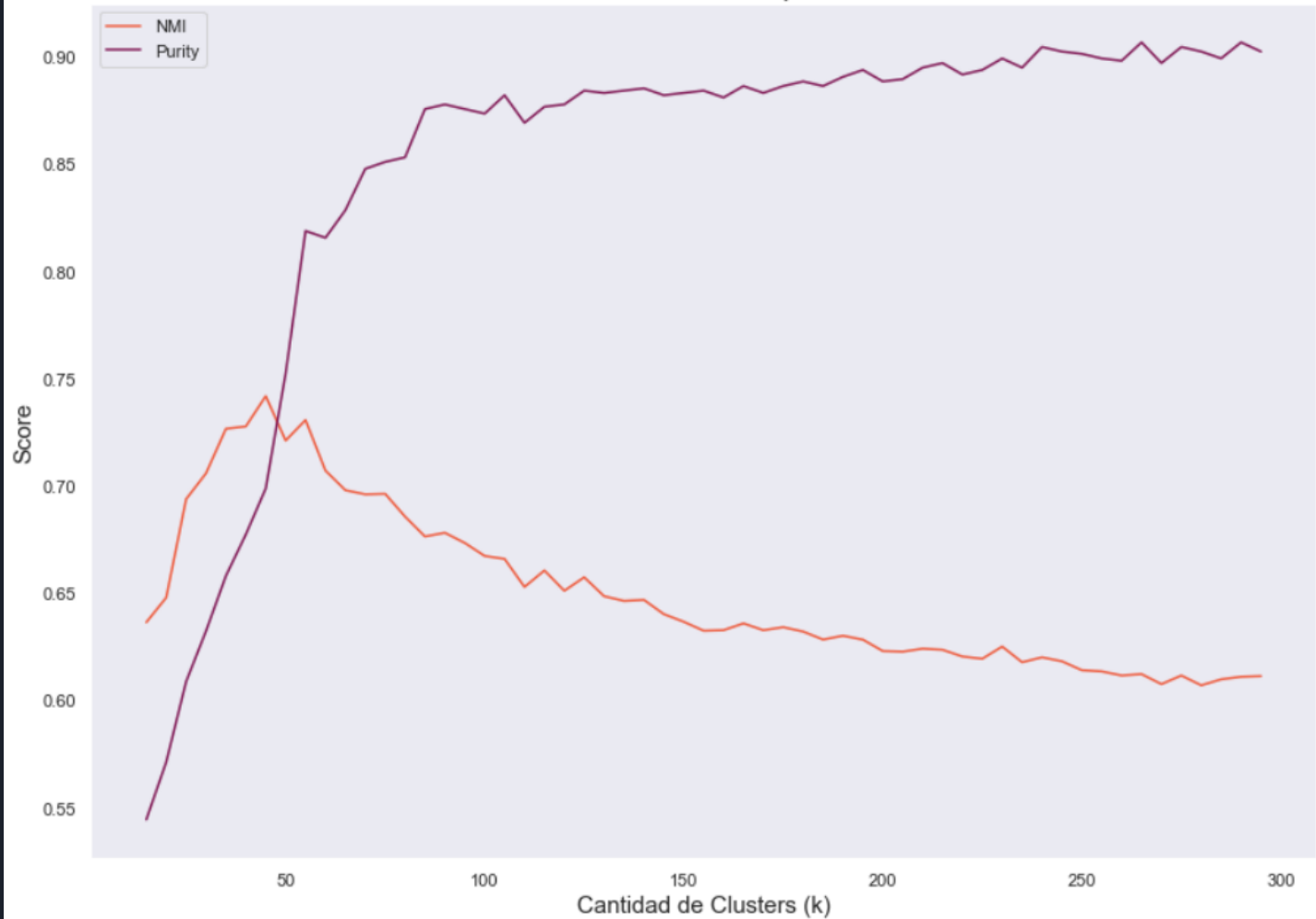
# Kmeans con Búsqueda de Hiper-parámetros

- Hago la búsqueda en torno a K (cantidad de clusters), por lo tanto uso NMI para conseguir el mejor modelo.
- Busco K entre 15 y 300 con un intervalo de 5.

El mejor NMI obtenido es: 0.7418



Distribución de métricas para Kmeans



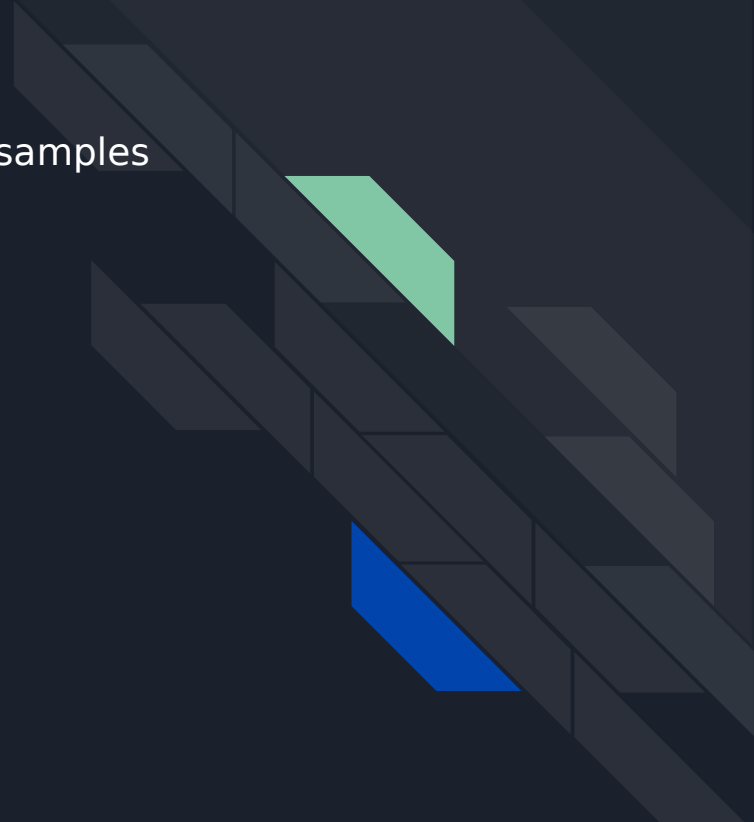
# HDBSCAN

- Hago un Grid Search entre `min_cluster_size`, `min_samples` y `cluster_selection_epsilon`

Los puntajes NMI obtenidos son:

Validation: 0.6705

Test: 0.6752



Muchas  
Gracias!