

Intersecciones de Bases de Datos

Archivo: Taller-2_LPrentt-20-03-2021.py **Sólo para análisis de intersección. No hay limpieza de base de datos**

No realicé merge de los dataframe = ColombiaCB-5March21.csv, Top100Startups- Colombia.xlsx y Empresas Unicorn - Contactos.xlsx.

Convertí los nombres de las empresas en los 3 dataframe a minúsculas (lowercase) y realice 3 intersecciones:

1. CrunchBase vs Top100, salieron 60 coincidencias
2. CrunchBase vs Unicorn, salieron 19 coincidencias
3. CrunchBase vs Top100 vs Unicorn, salieron 12 coincidencias

Intersecciones de Bases de Datos

Código para opción 1

1. CrunchBase vs Top100, salieron 60 coincidencias. Lo implemente con 2 For

```
df["y"]=0  
df["Flag"]=0
```

```
intersect=set(df['Organization Name']).intersection(set(dftop['Organization']))  
len(intersect)
```

Forma L. Prentt

```
k = df.shape[1]  
for i in range(df.shape[0]):  
    for j in range(dftop.shape[0]):  
        if df["Organization Name"][i]==dftop["Organization"][j]:  
            df.iloc[i:i+1, k-2:k-1] = 1  
            break
```

Forma de Miller Quiroga

```
for j in intersect:  
    df.loc[df['Organization Name'] == j, ['Flag']] = 1
```

Las 2 formas producen los mismos resultados, sin embargo la Miller Quiroga es mas eficiente

Intersecciones de Bases de Datos

Código para opción 2

2. CrunchBase vs Unicorn, salieron 19 coincidencias. Lo implemente con 2 For

```
df1["y"]=0  
df1["Flag"]=0
```

```
intersect2=set(df1['Organization Name']).intersection(set(dfunicorn['Name']))  
len(intersect2)
```

Forma L. Prentt

```
k = df1.shape[1]  
for i in range(df1.shape[0]):  
    for j in range(dfunicorn.shape[0]):  
        if df1["Organization Name"][i]==dfunicorn["Name"][j]:  
            df1.iloc[i:i+1, k-2:k-1] = 1  
            break
```

Forma Miller Quiroga

```
for t in intersect2:  
    df1.loc[df1['Organization Name'] == t, ['Flag']] = 1
```

Las 2 formas producen los mismos resultados, sin embargo la Miller Quiroga es mas eficiente

Intersecciones de Bases de Datos

Código para opción 3

3. CrunchBase vs Top100 vs Unicorn, salieron 12 coincidencias. Lo implemente con 2 For

```
df2["y"]=0
df2["Flag"]=0

intersect3=set(df['Organization Name']).intersection(set(dftop['Organization'])).intersection(set(dfunicorn['Name']))
len(intersect3)

# conversión del conjunto intersect3 a numpy.array
int3=np.array(list(intersect3))

##### Forma L. Prentt
k = df2.shape[1]
for i in range(df2.shape[0]):
    for j in range(len(intersect3)):
        if df2["Organization Name"][i]==int3[j]:
            df2.iloc[i:i+1, k-2:k-1] = 1
            break

##### Forma Miller Quiroga
for l in intersect3:
    df2.loc[df2['Organization Name'] == l, ['Flag']] = 1
```

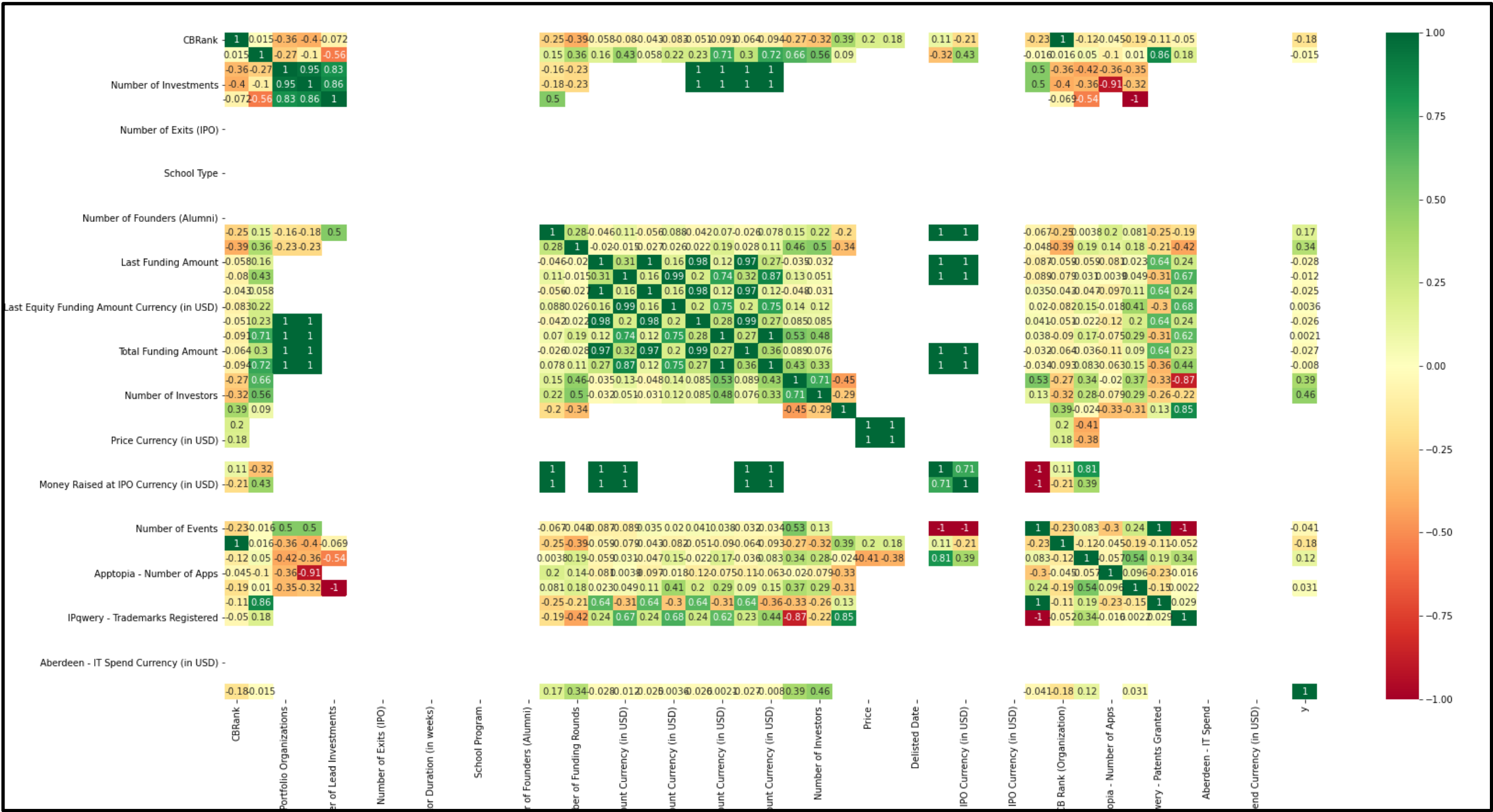
Para Análisis de Regresión se usa el archivo **Taller-2_LPrentt.20A-03-2021.py** :

Se eliminaron todas las empresas que no pertenecen a Colombia

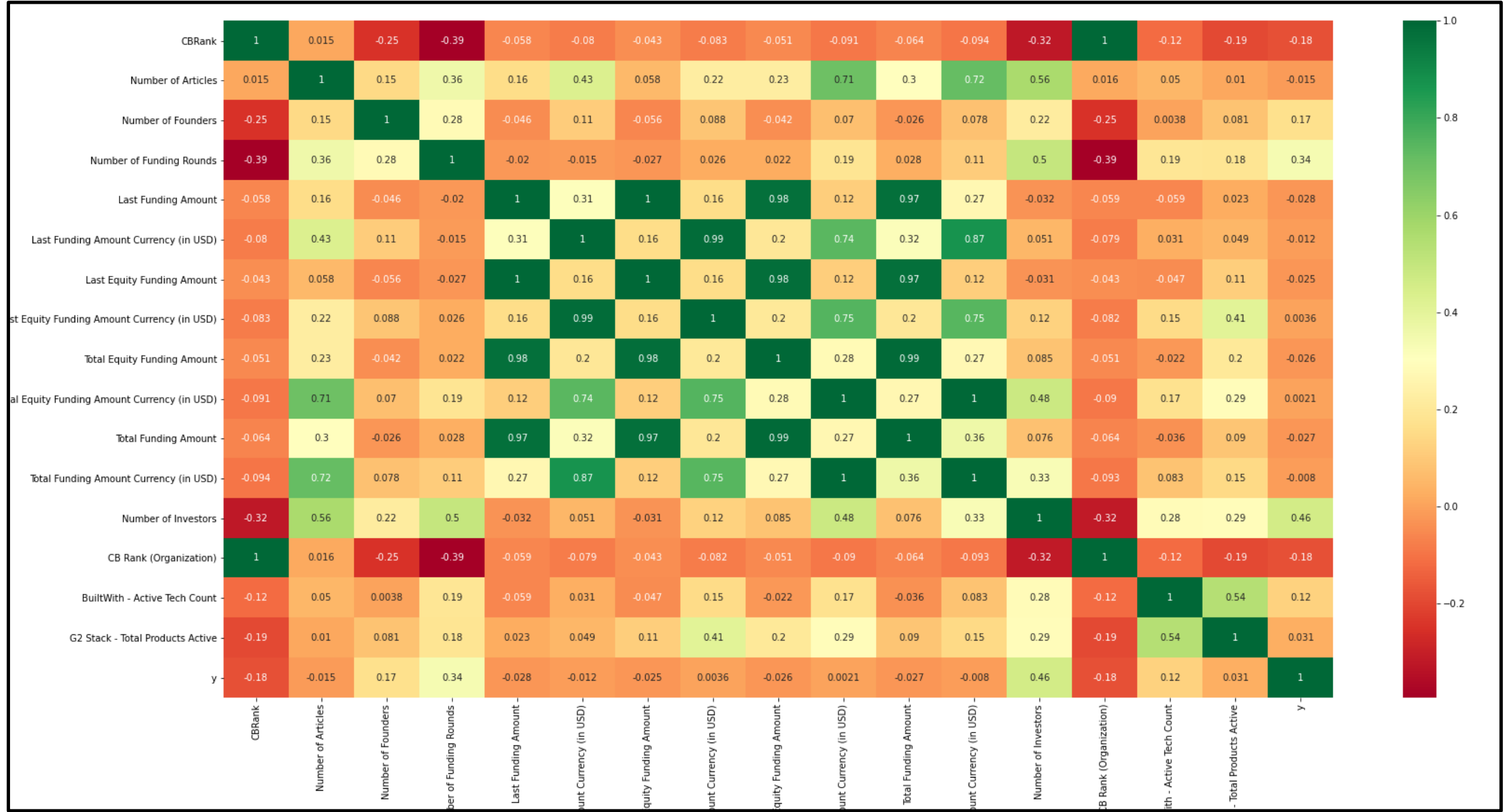
El cual realiza las intersecciones con el código de Miller Quiroga y realiza limpieza de base de datos crunchbase. La variable objetivo a modelar que contiene las intersecciones es “y”

La base de datos final para trabajar es el dataframe **df2**

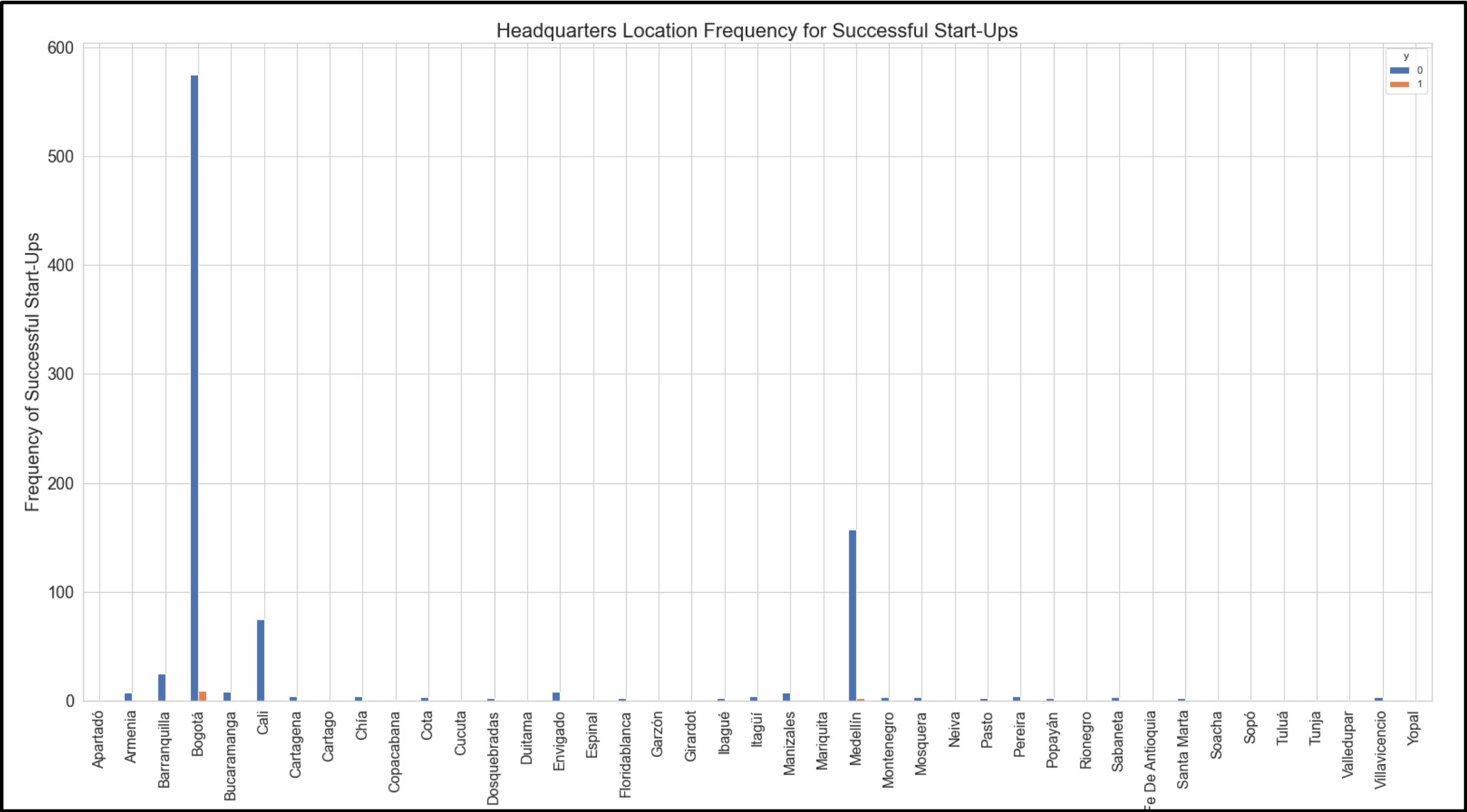
Matriz de Correlación, sin filtro de columnas con mas 80% de nan



Matriz de Correlación, con filtro de columnas con mas 80% de nan

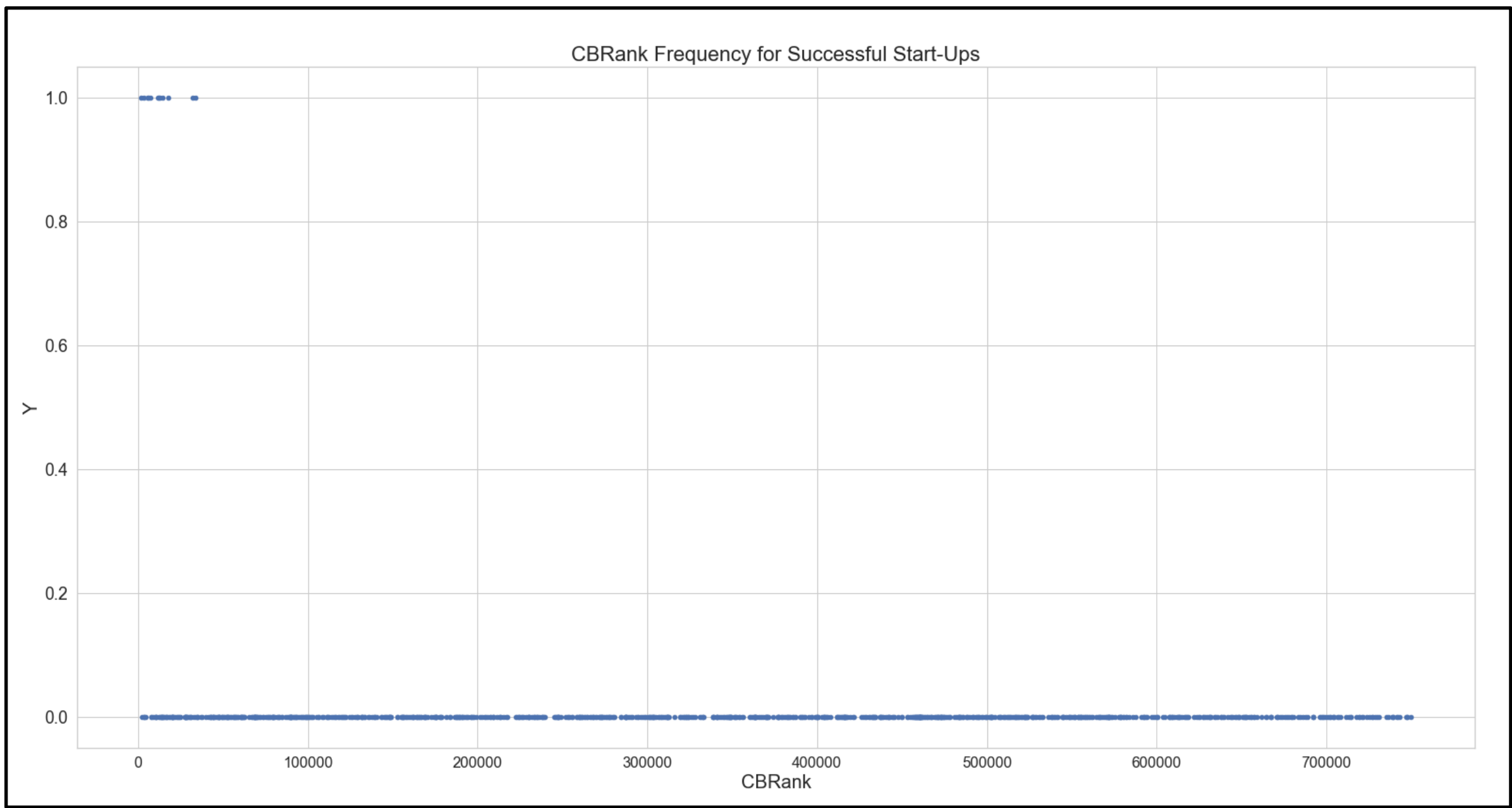


Headquarters Location Frequency for Successful Start-Ups



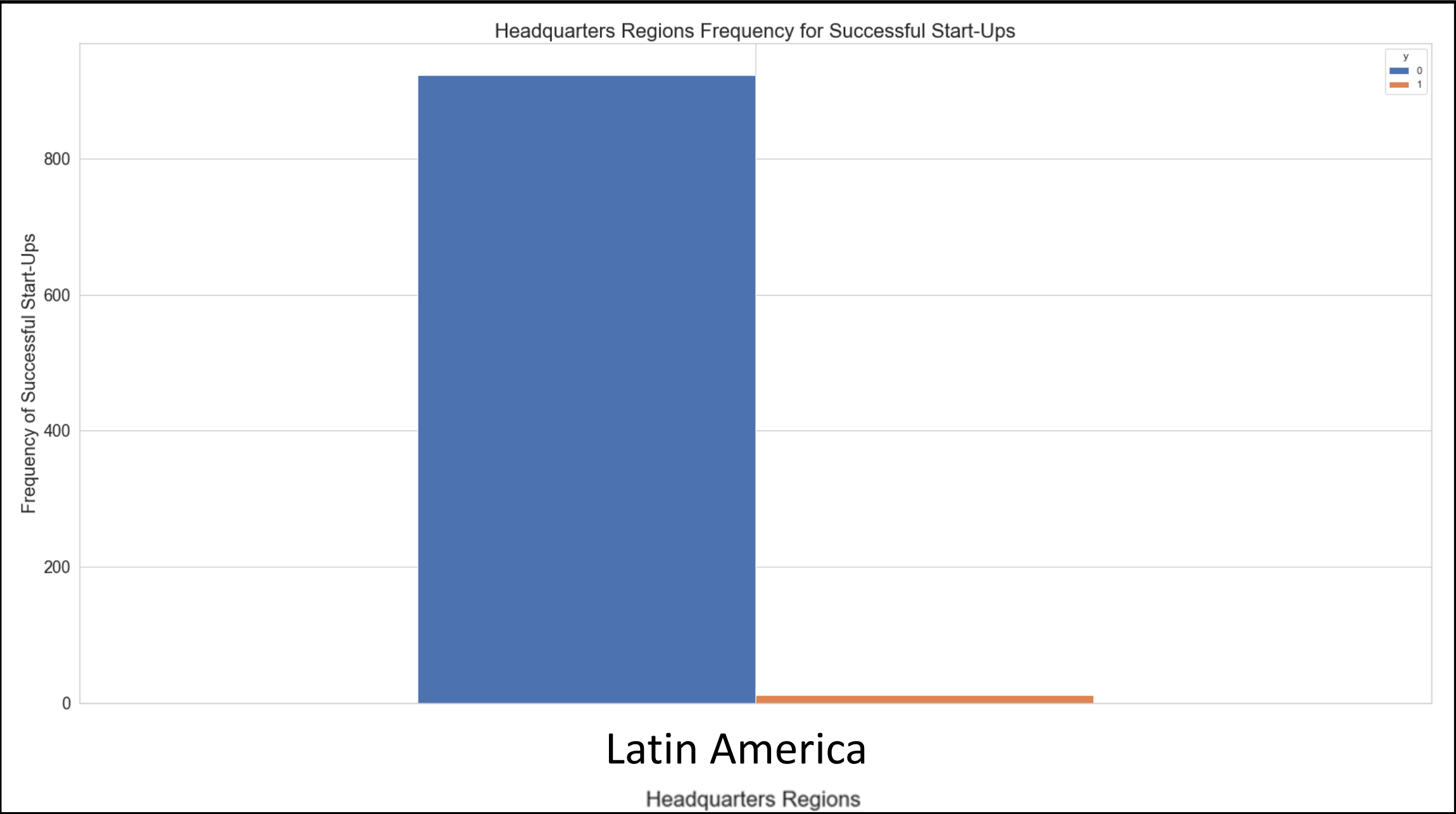
Si Importa la localización de la startup para su éxito. Bogotá y Medellín tienen las Start-ups con éxito, sin embargo también se puede eliminar para evitar sesgo de que sólo hay éxito en Bogotá y Medellín

CBRank vs Successful Start-Ups



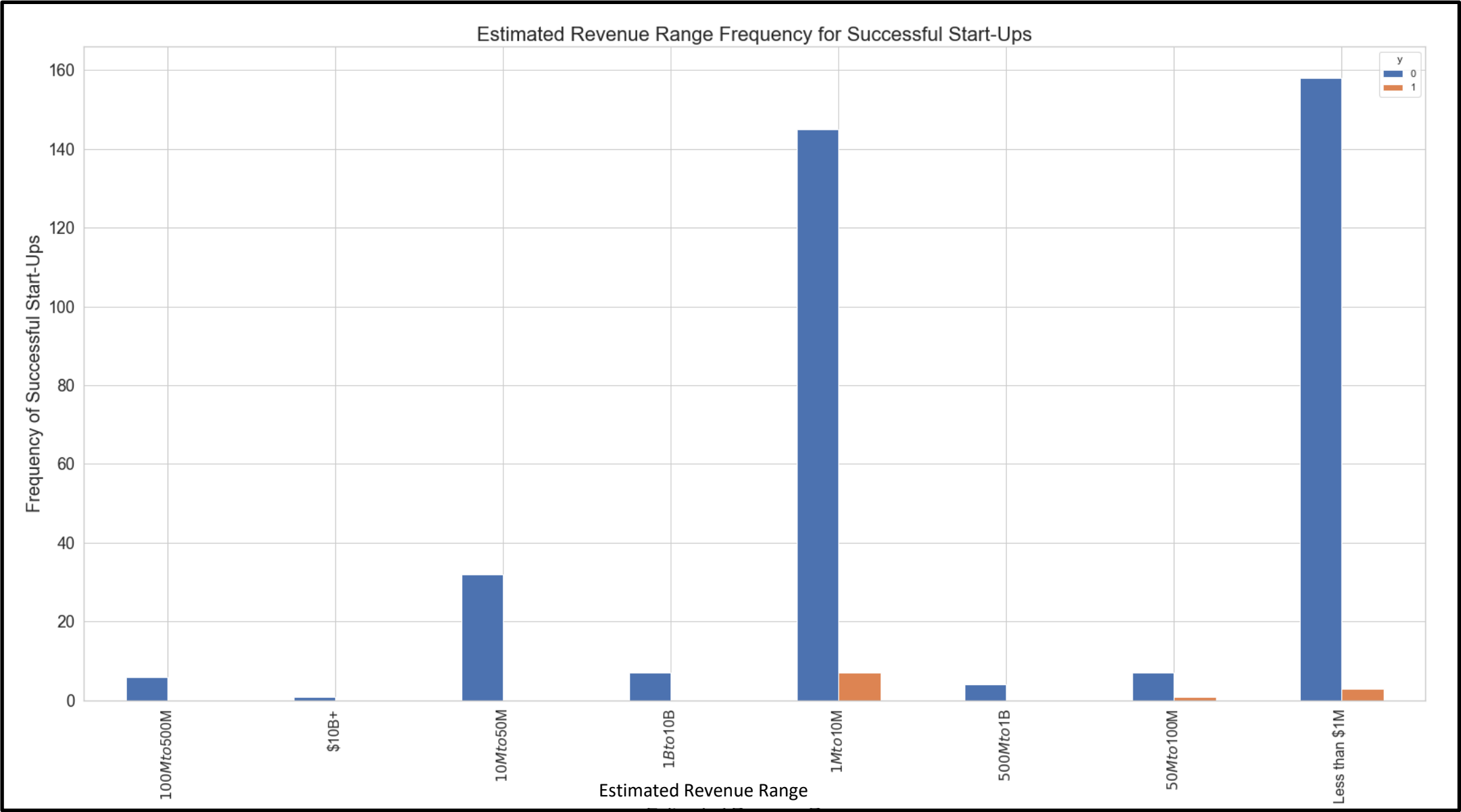
Si/No Importa Cbrank para éxito de Startups. (Tengo duda)

Headquarters Regions vs Successful Start-Ups



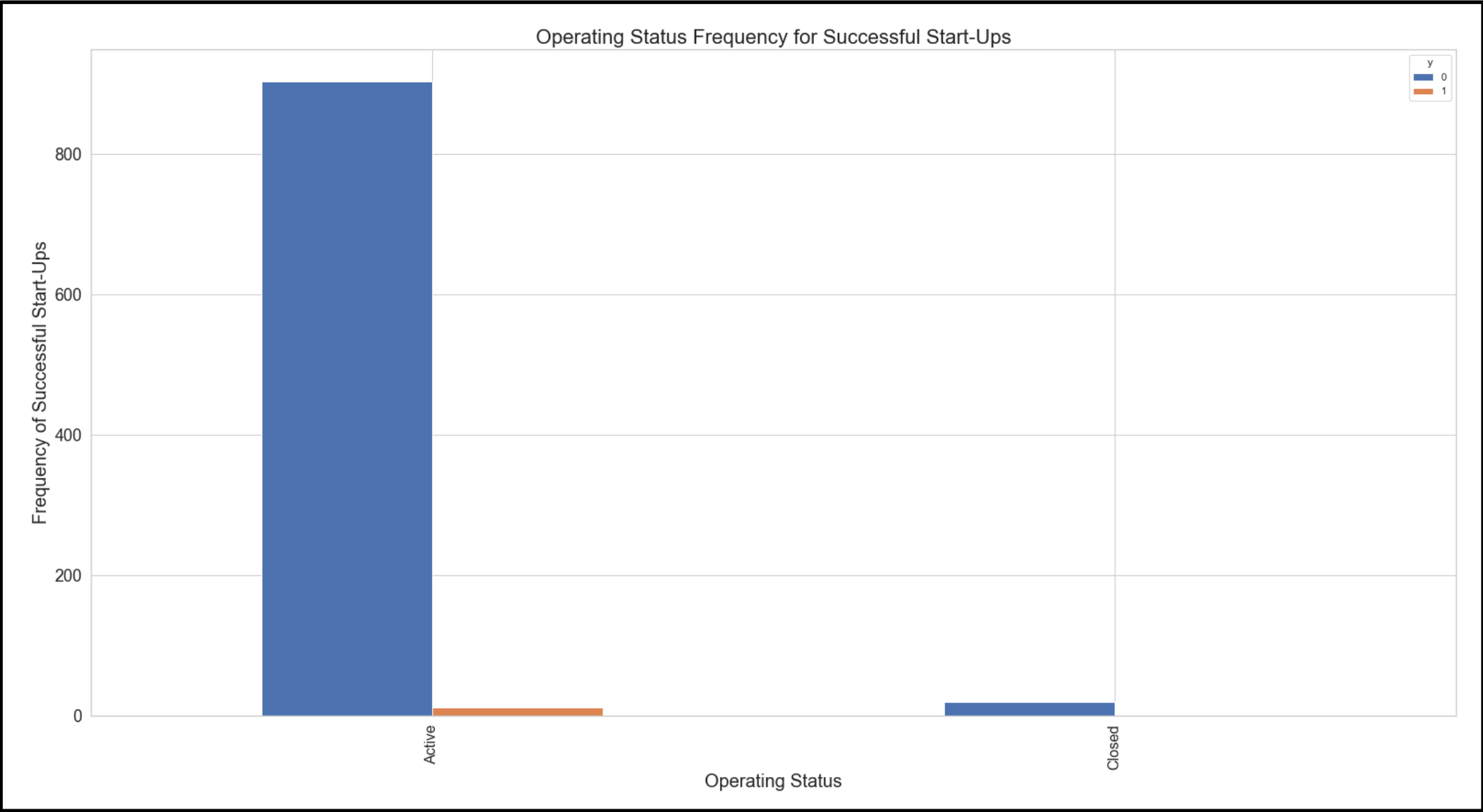
No importa Headquarters Regions para éxito de Startups

Estimated Revenue Range vs Successful Start-Ups



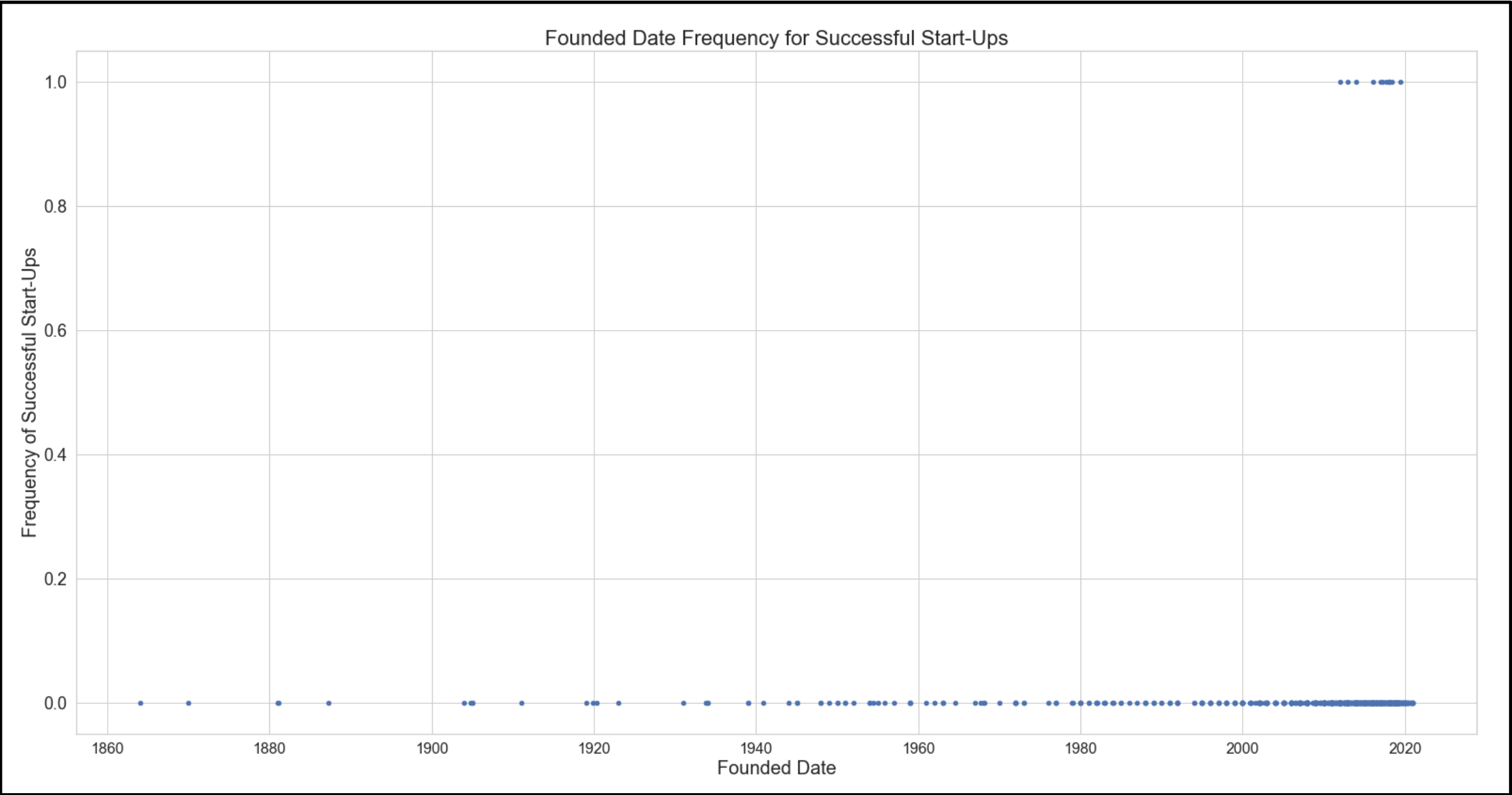
Si Importa Estimated Revenue Range para éxito de Startups.

Operating Status vs Successful Start-Ups



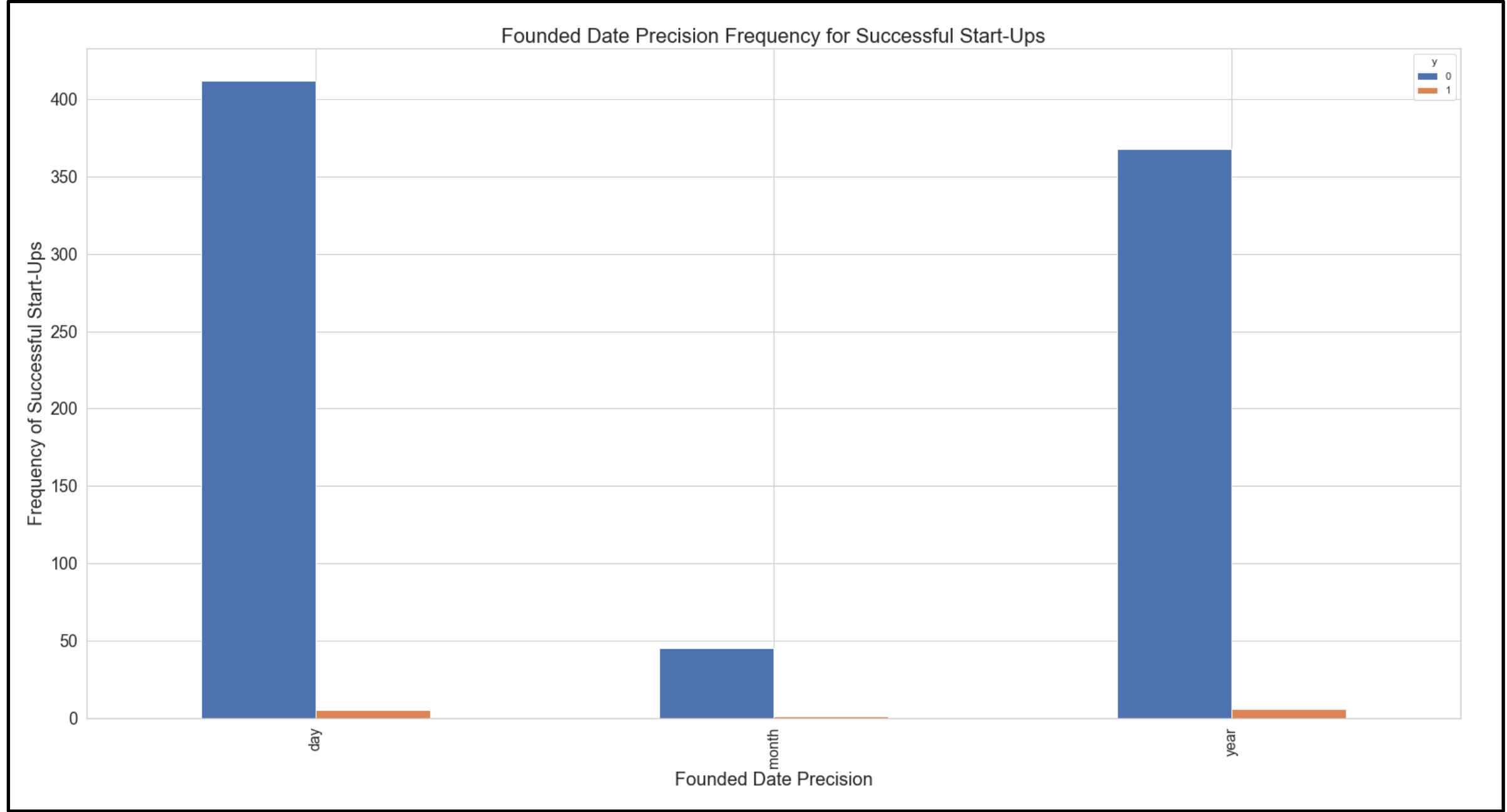
No importa Operating Status para éxito de Startups

Founded Date vs Successful Start-Ups



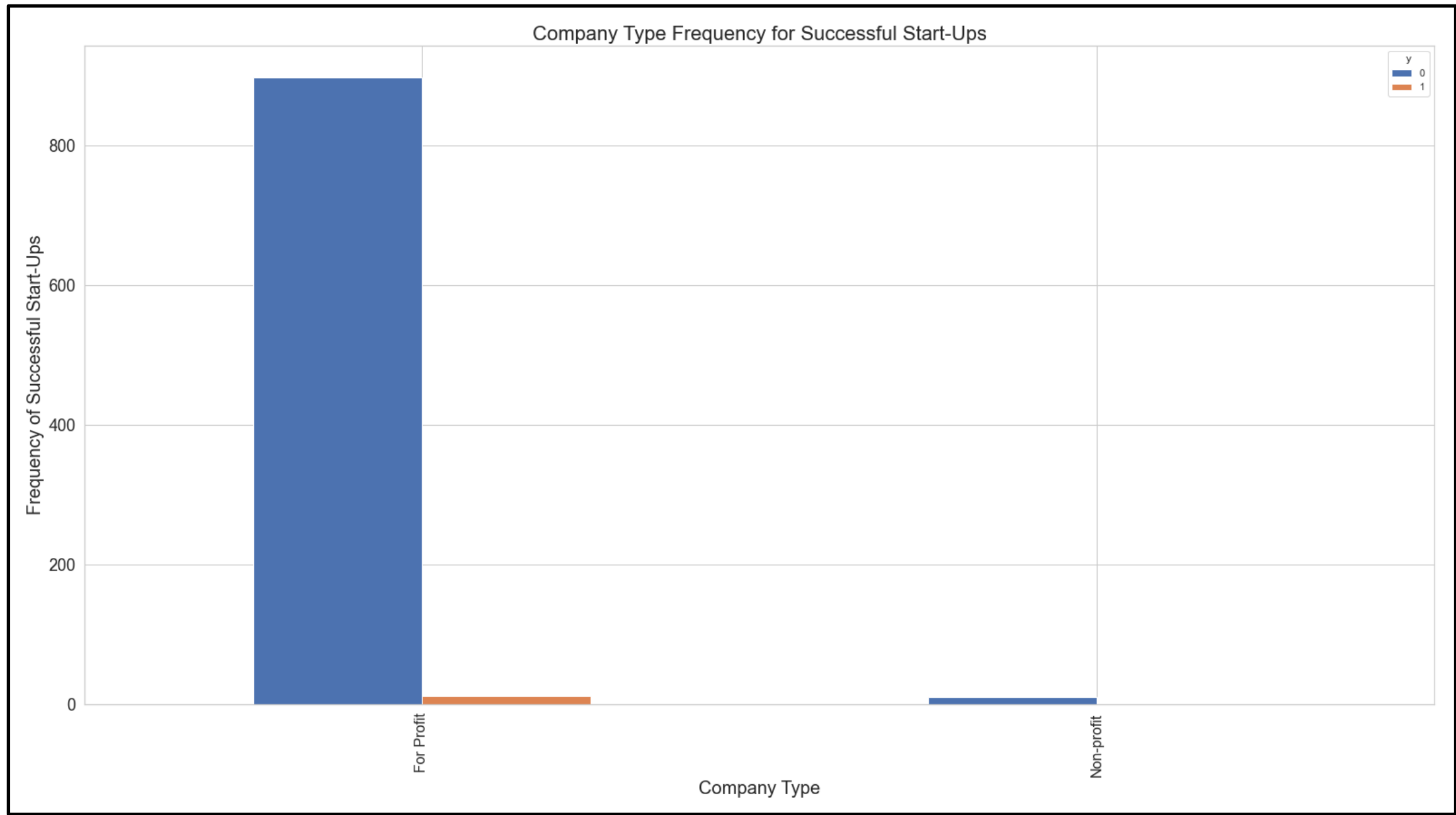
No importa Founded Date para éxito de Startups

Founded Date Precision vs Successful Start-Ups



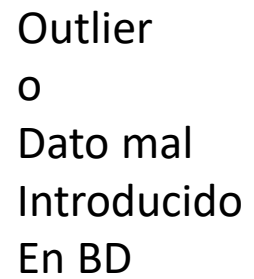
Parece importar Founded Date Precision para éxito de Startups, sin embargo no entiendo significado de la variable. La saco

Company Type vs Successful Start-Ups



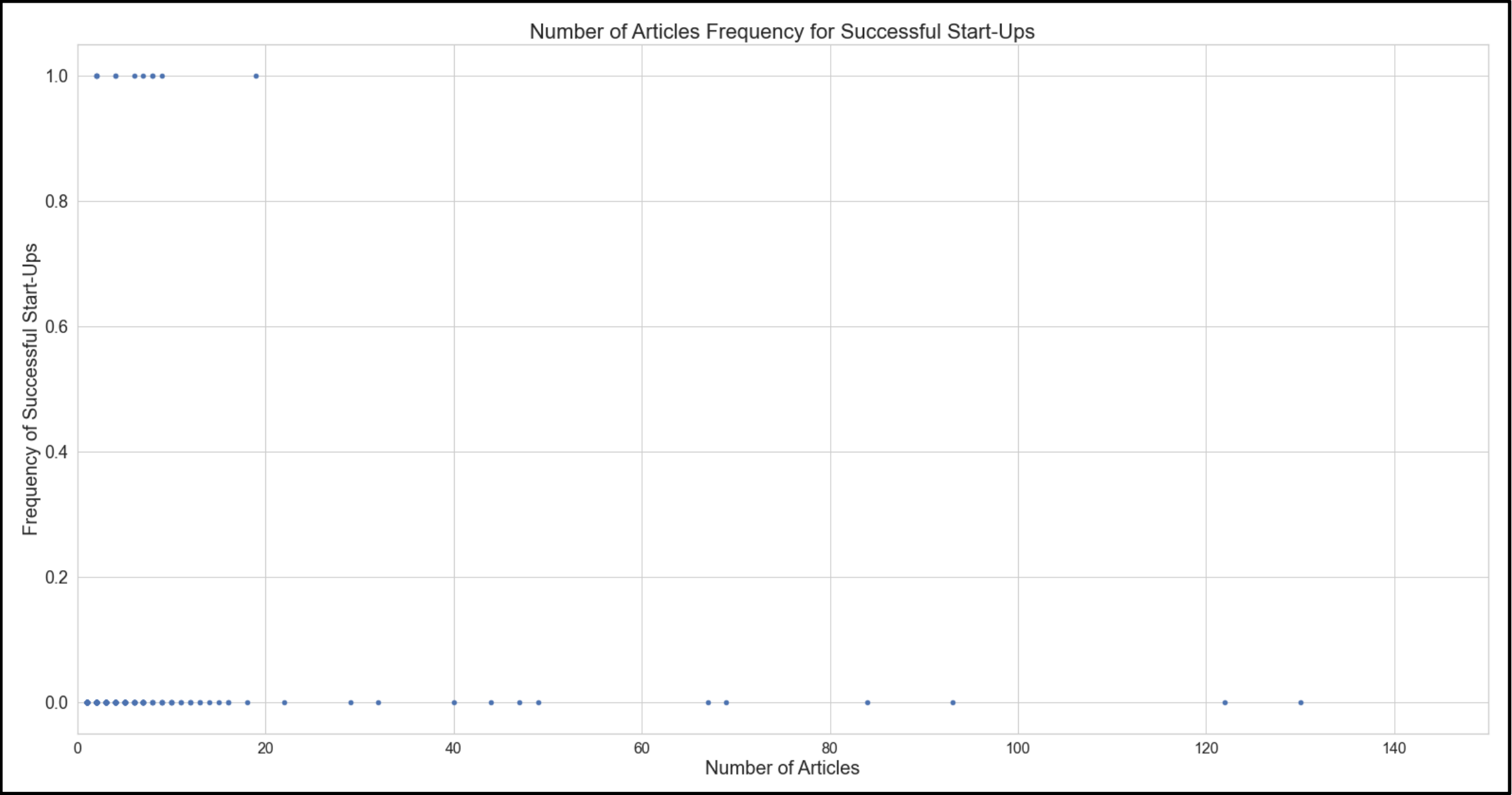
No importa Company Type para éxito de Startups

Si importa Number of Articles para éxito de Startups



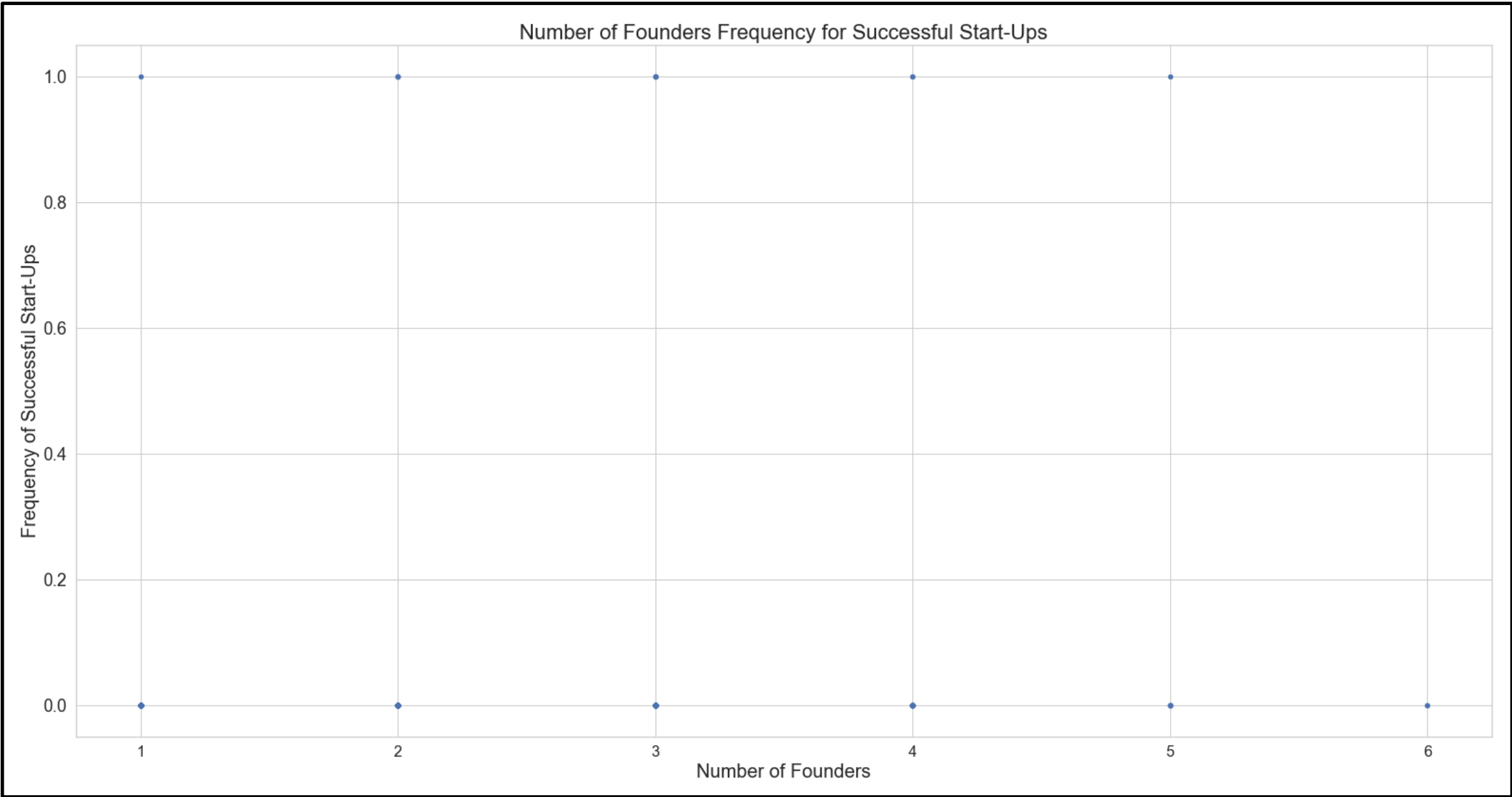
Number of Articles vs Successful Start-Ups

PLOT ZOOM



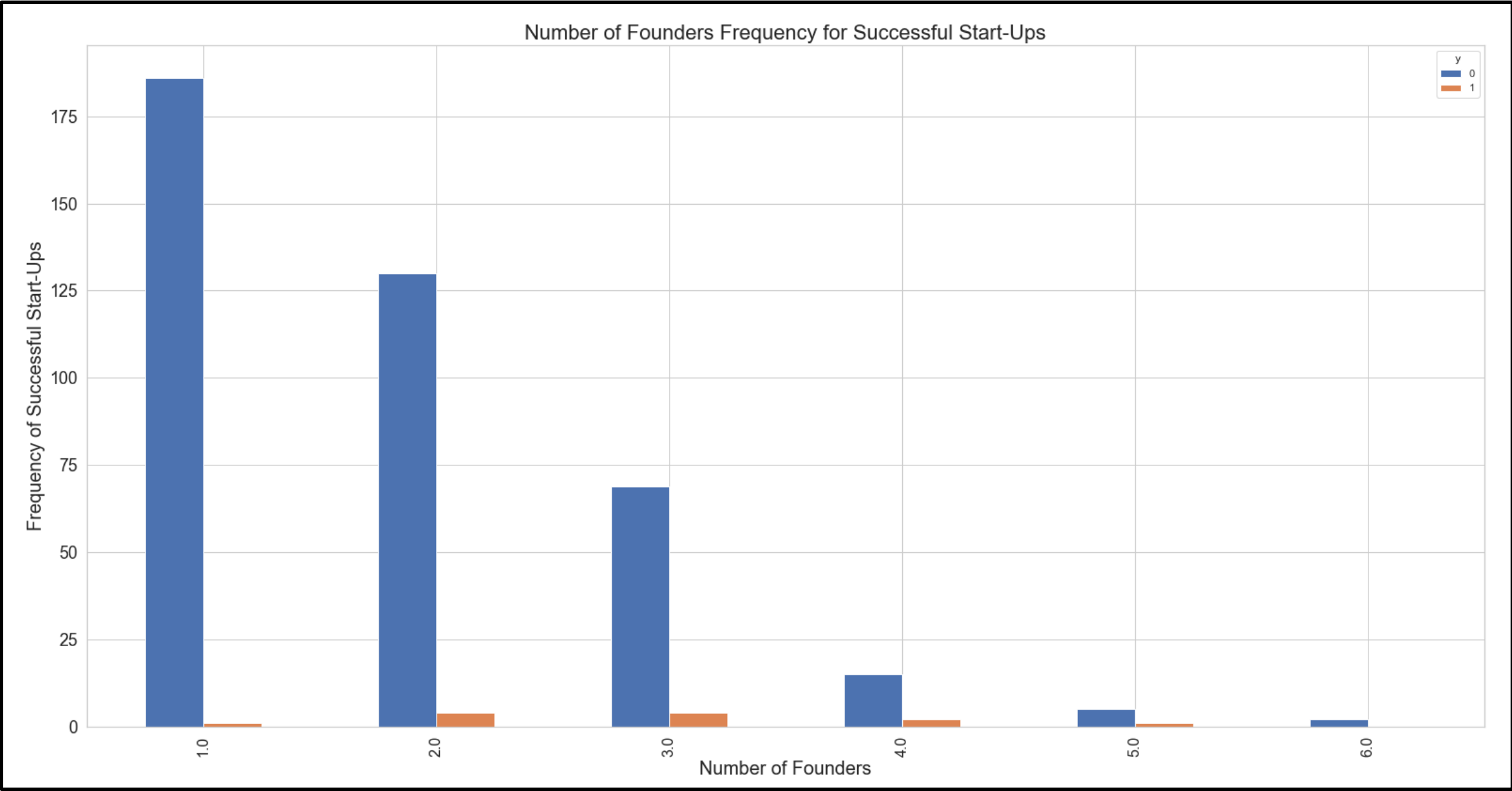
Si importa Number of Articles para éxito de Startups

Number of Founders vs Successful Start-Ups (Variable N merica)



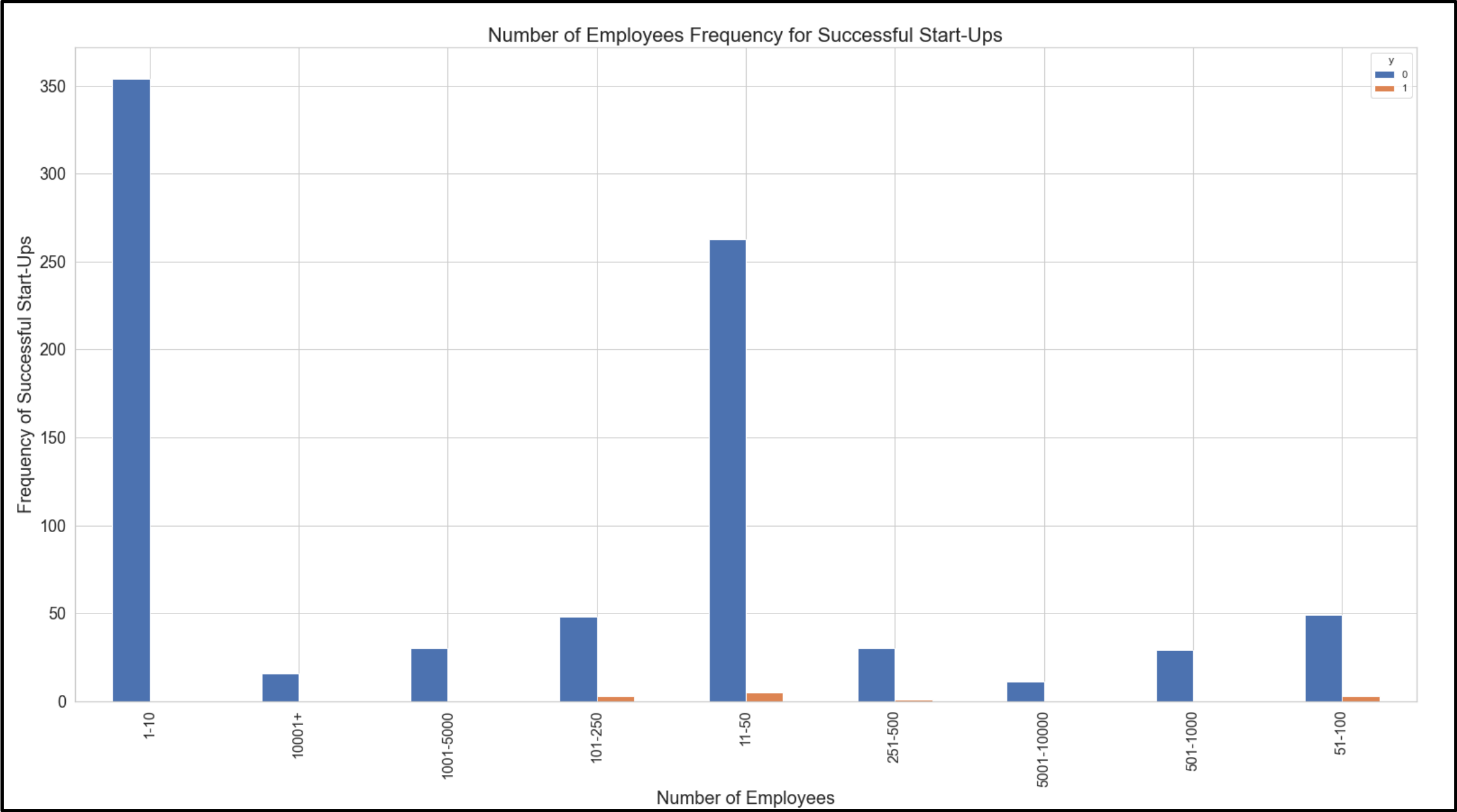
Si importa Number of Founders para  xito de Startups

Number of Founders vs Successful Start-Ups (Variable Categórica)



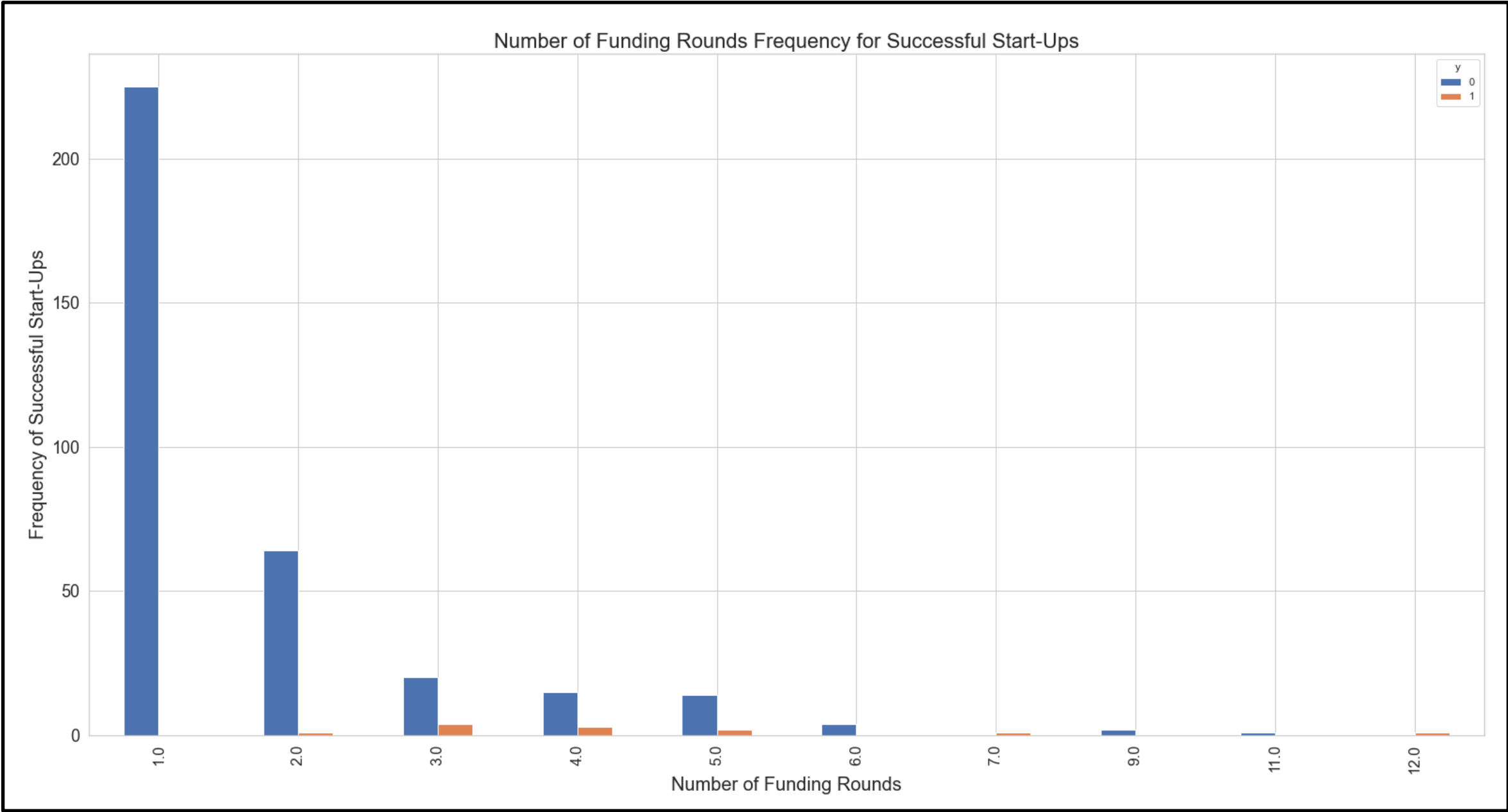
Si importa Number of Founders para éxito de Startups

Number of Employees vs Successful Start-Ups



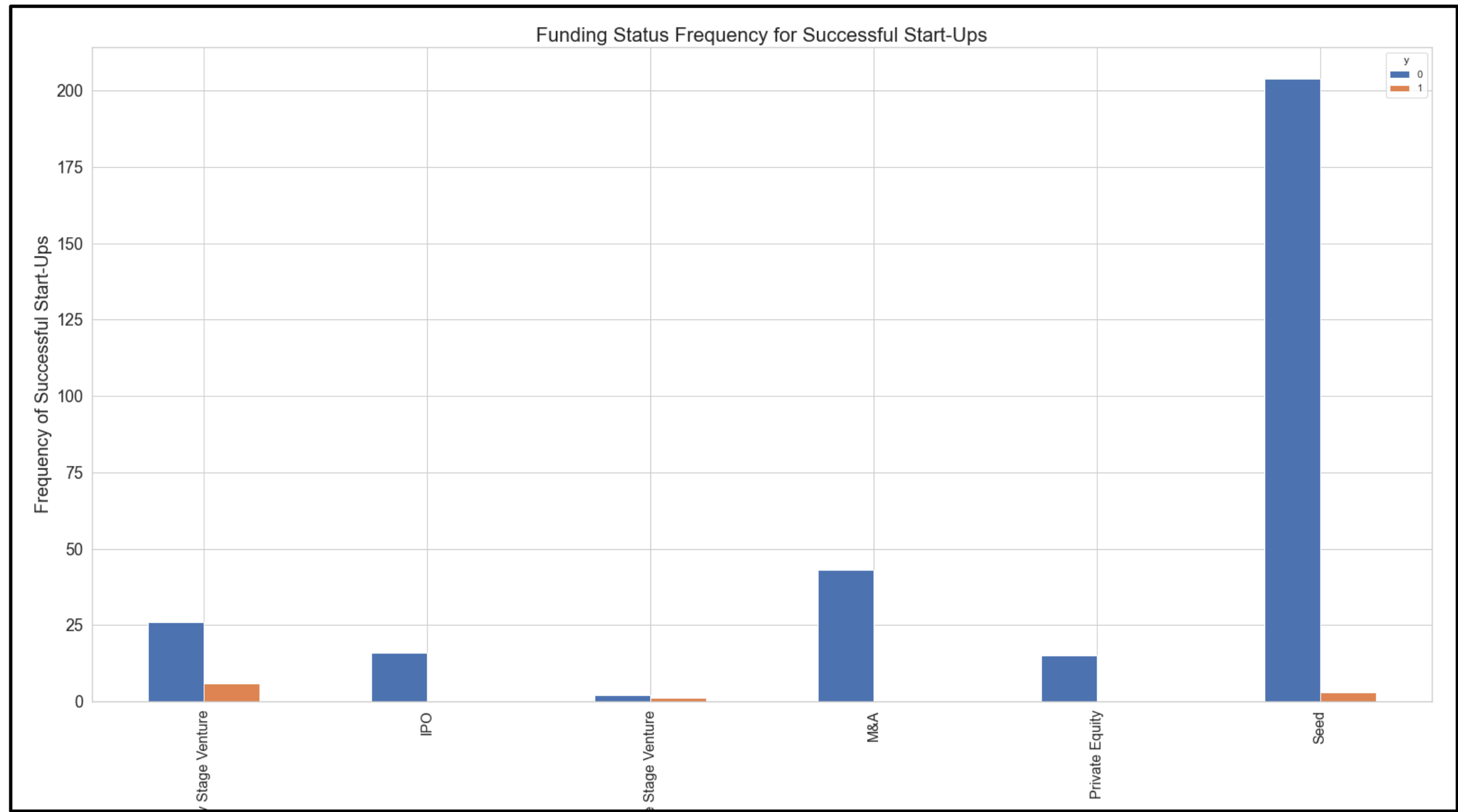
Si importa Number of Employees para éxito de Startups

Number of Funding Rounds vs Successful Start-Ups



Si importa Number of Funding Rounds para éxito de Startups

Funding Status vs Successful Start-Ups

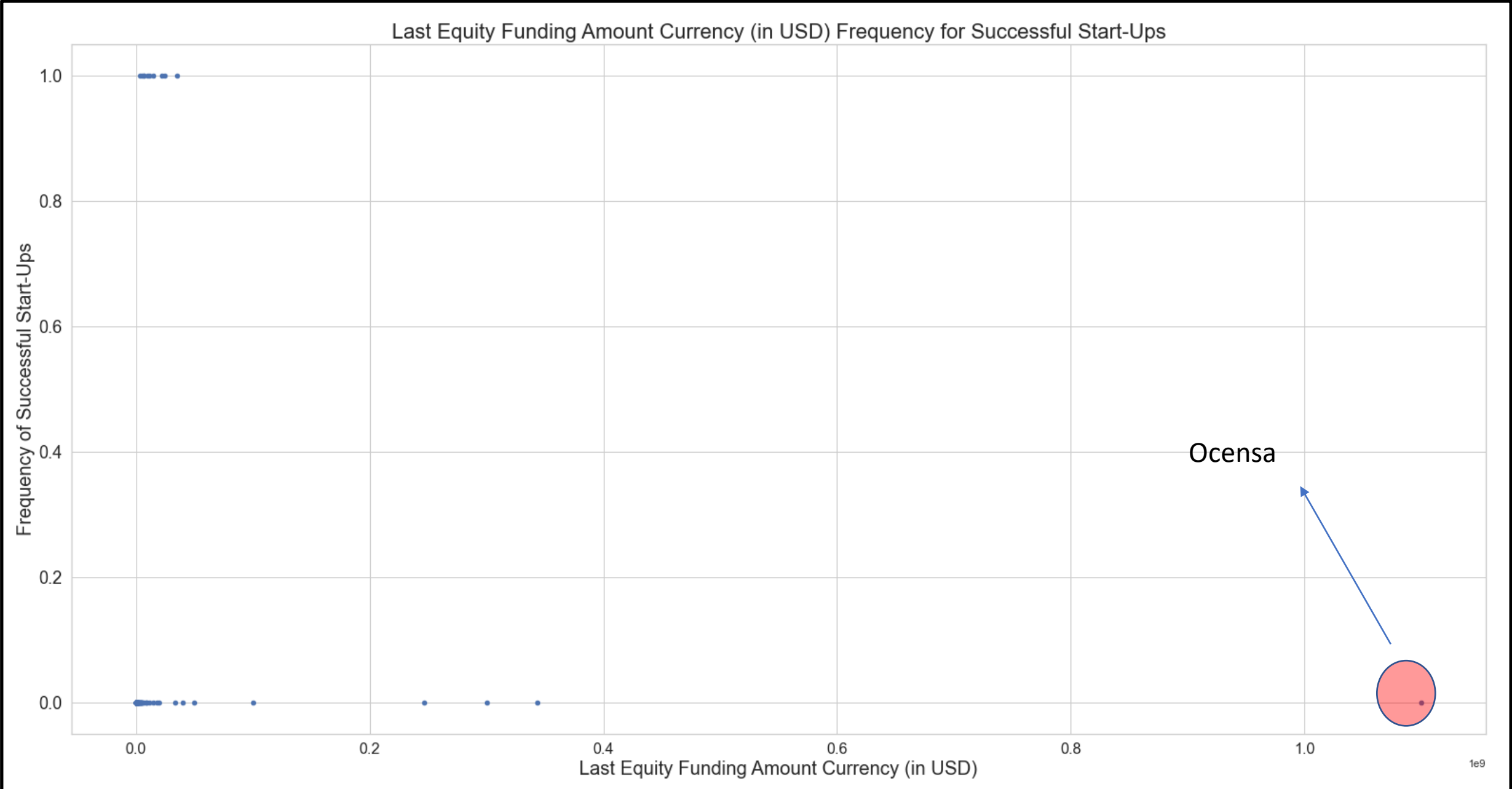


Si importa Funding Status para éxito de Startups

[illegible]

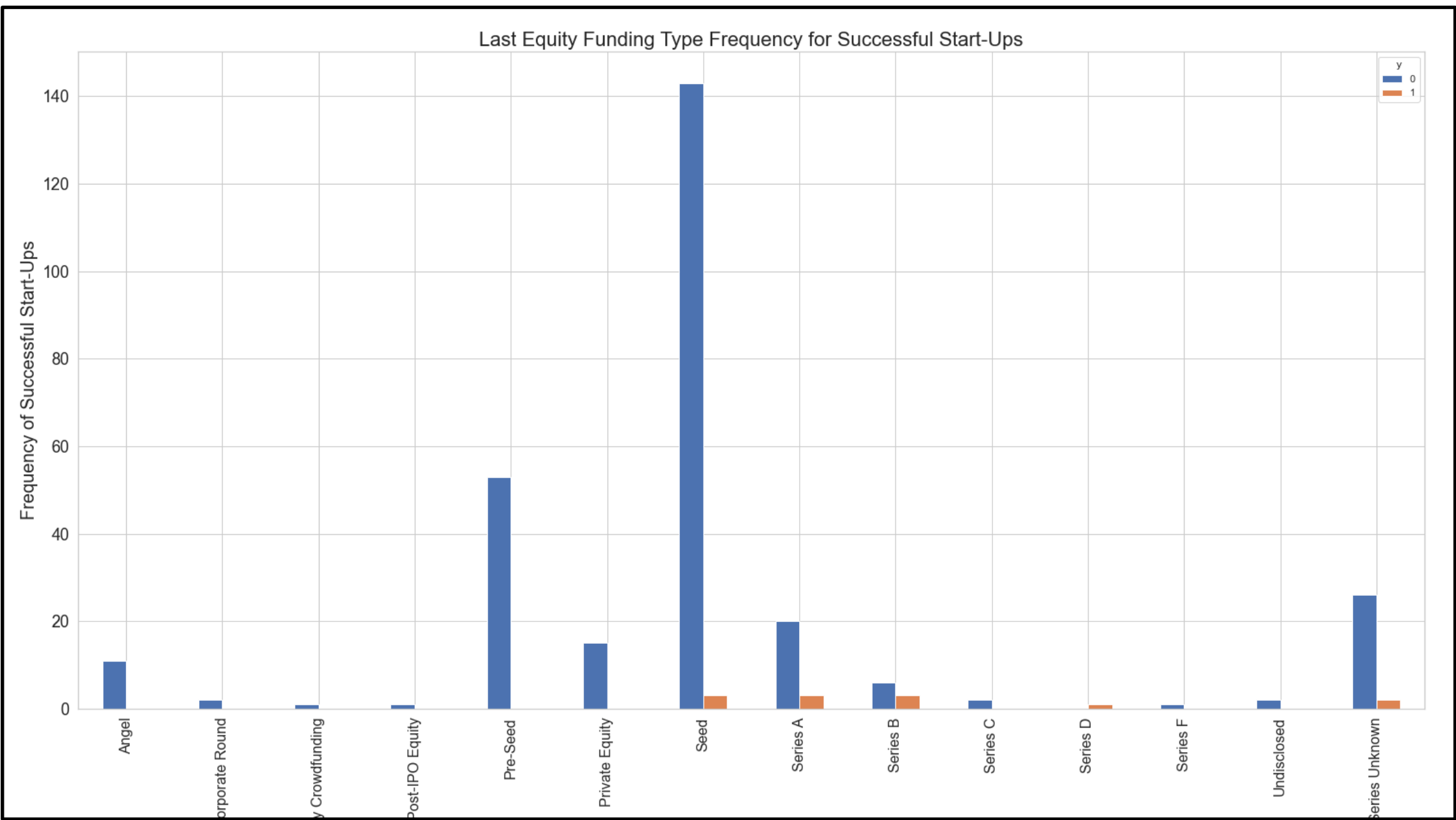
Si importa Last Funding Amount Currency (in USD) para éxito de Startups

Last Equity Funding Amount Currency (in USD) vs Successful Start-Ups



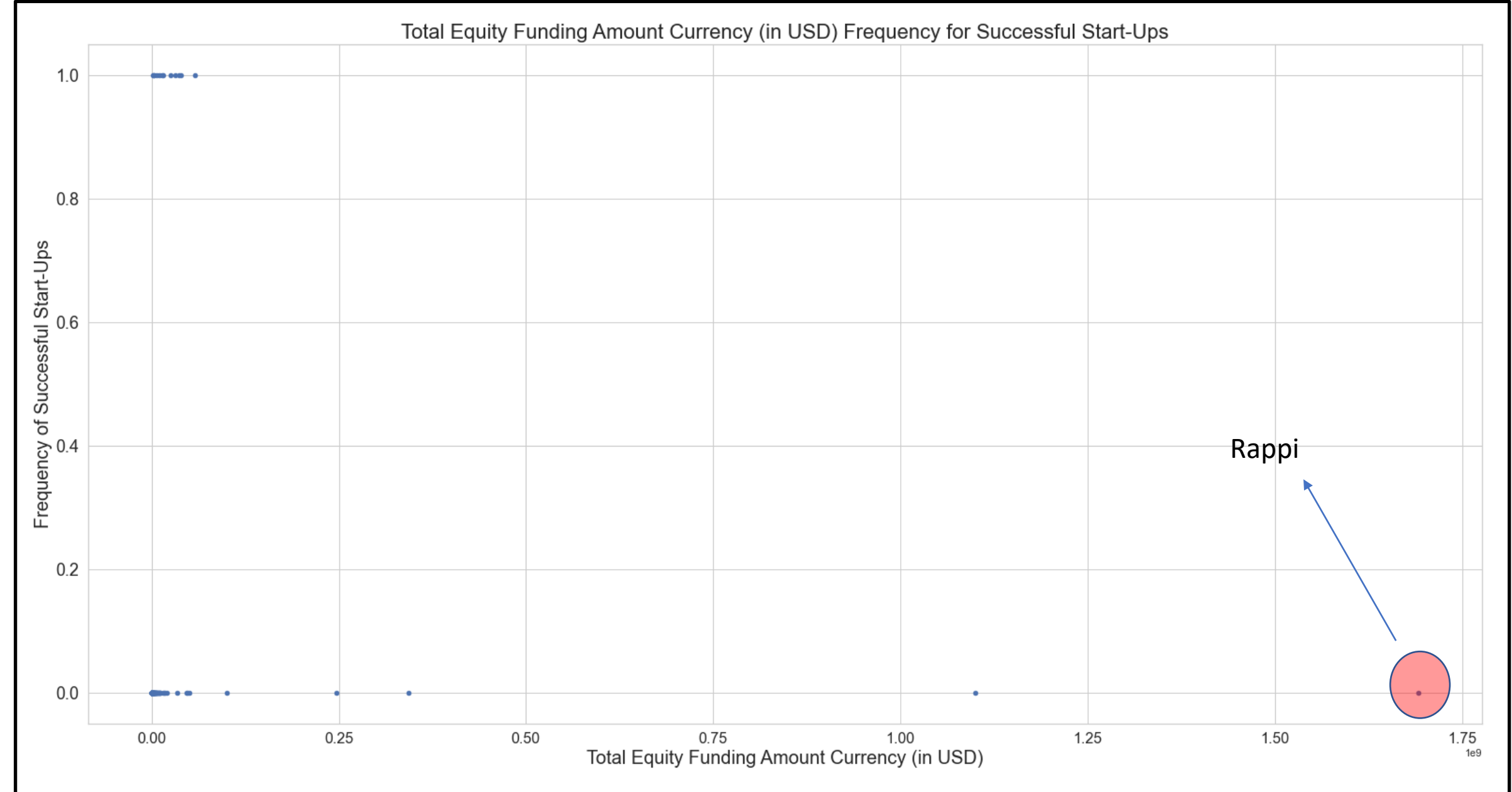
Si importa Last Equity Funding Amount Currency (in USD) para éxito de Startups

Last Equity Funding Type vs Successful Start-Ups



Si importa Last Equity Funding Type para éxito de Startups

Total Equity Funding Amount Currency (in USD) vs Successful Start-Ups



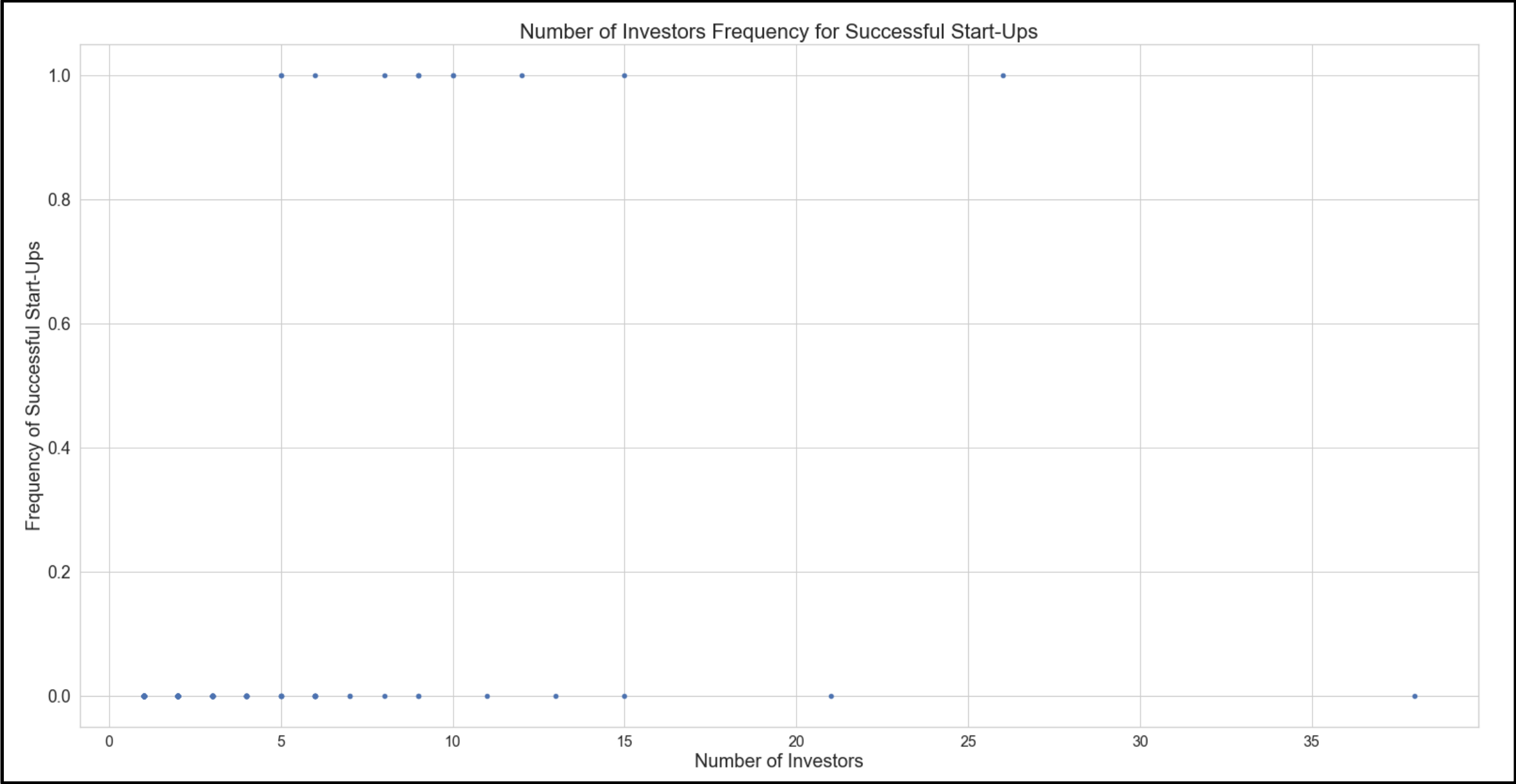
Si importa Total Equity Funding Amount Currency (in USD) para éxito de Startups

The scatter plot displays the frequency of successful start-ups across different total funding amounts in USD. The x-axis represents the total funding amount in billions of USD, ranging from 0.00 to 2.00. The y-axis represents the frequency of successful start-ups, ranging from 0.0 to 1.0. The data shows a high concentration of start-ups with very low funding amounts (below 0.25 billion USD), where the frequency is either 0.0 or 1.0. A notable outlier is Avianca Holdings, which received approximately 2.00 billion USD in funding and has a frequency of 0.0. This point is highlighted with a large red circle and a blue arrow.

Total Funding Amount (USD)	Frequency of Successful Start-Ups
0.00	1.00
0.01	1.00
0.02	1.00
0.03	1.00
0.04	1.00
0.05	1.00
0.06	1.00
0.07	1.00
0.08	1.00
0.09	1.00
0.10	1.00
0.11	1.00
0.12	1.00
0.13	1.00
0.14	1.00
0.15	1.00
0.16	1.00
0.17	1.00
0.18	1.00
0.19	1.00
0.20	1.00
0.21	1.00
0.22	1.00
0.23	1.00
0.24	1.00
0.25	1.00
0.26	1.00
0.27	1.00
0.28	1.00
0.29	1.00
0.30	1.00
0.31	1.00
0.32	1.00
0.33	1.00
0.34	1.00
0.35	1.00
0.36	1.00
0.37	1.00
0.38	1.00
0.39	1.00
0.40	1.00
0.41	1.00
0.42	1.00
0.43	1.00
0.44	1.00
0.45	1.00
0.46	1.00
0.47	1.00
0.48	1.00
0.49	1.00
0.50	1.00
0.51	1.00
0.52	1.00
0.53	1.00
0.54	1.00
0.55	1.00
0.56	1.00
0.57	1.00
0.58	1.00
0.59	1.00
0.60	1.00
0.61	1.00
0.62	1.00
0.63	1.00
0.64	1.00
0.65	1.00
0.66	1.00
0.67	1.00
0.68	1.00
0.69	1.00
0.70	1.00
0.71	1.00
0.72	1.00
0.73	1.00
0.74	1.00
0.75	1.00
0.76	1.00
0.77	1.00
0.78	1.00
0.79	1.00
0.80	1.00
0.81	1.00
0.82	1.00
0.83	1.00
0.84	1.00
0.85	1.00
0.86	1.00
0.87	1.00
0.88	1.00
0.89	1.00
0.90	1.00
0.91	1.00
0.92	1.00
0.93	1.00
0.94	1.00
0.95	1.00
0.96	1.00
0.97	1.00
0.98	1.00
0.99	1.00
1.00	1.00
1.01	1.00
1.02	1.00
1.03	1.00
1.04	1.00
1.05	1.00
1.06	1.00
1.07	1.00
1.08	1.00
1.09	1.00
1.10	1.00
1.11	1.00
1.12	1.00
1.13	1.00
1.14	1.00
1.15	1.00
1.16	1.00
1.17	1.00
1.18	1.00
1.19	1.00
1.20	1.00
1.21	1.00
1.22	1.00
1.23	1.00
1.24	1.00
1.25	1.00
1.26	1.00
1.27	1.00
1.28	1.00
1.29	1.00
1.30	1.00
1.31	1.00
1.32	1.00
1.33	1.00
1.34	1.00
1.35	1.00
1.36	1.00
1.37	1.00
1.38	1.00
1.39	1.00
1.40	1.00
1.41	1.00
1.42	1.00
1.43	1.00
1.44	1.00
1.45	1.00
1.46	1.00
1.47	1.00
1.48	1.00
1.49	1.00
1.50	1.00
1.51	1.00
1.52	1.00
1.53	1.00
1.54	1.00
1.55	1.00
1.56	1.00
1.57	1.00
1.58	1.00
1.59	1.00
1.60	1.00
1.61	1.00
1.62	

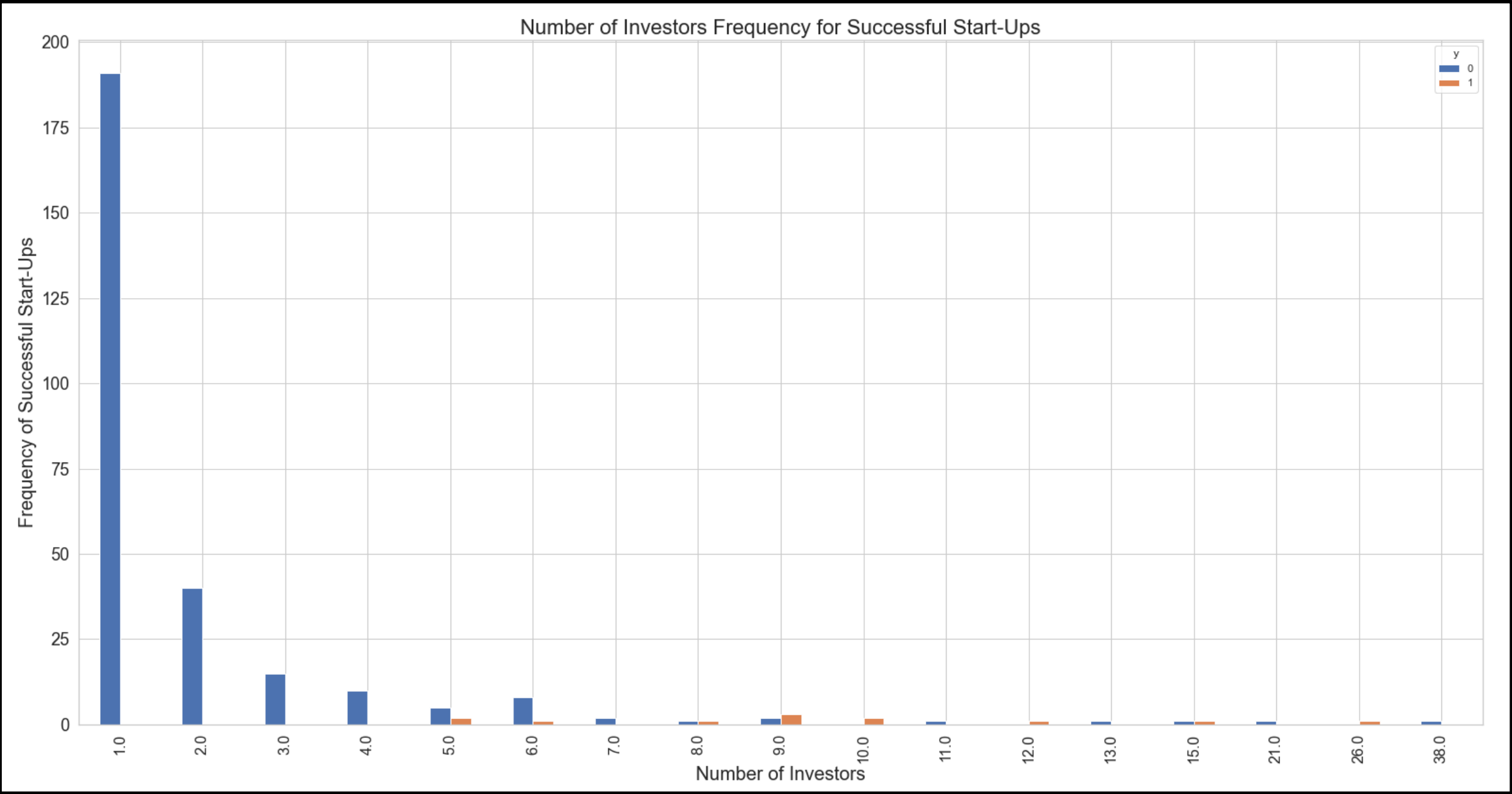
Si importa Total Funding Amount Currency (in USD) para éxito de Startups

Number of Investors vs Successful Start-Ups (Númerica)



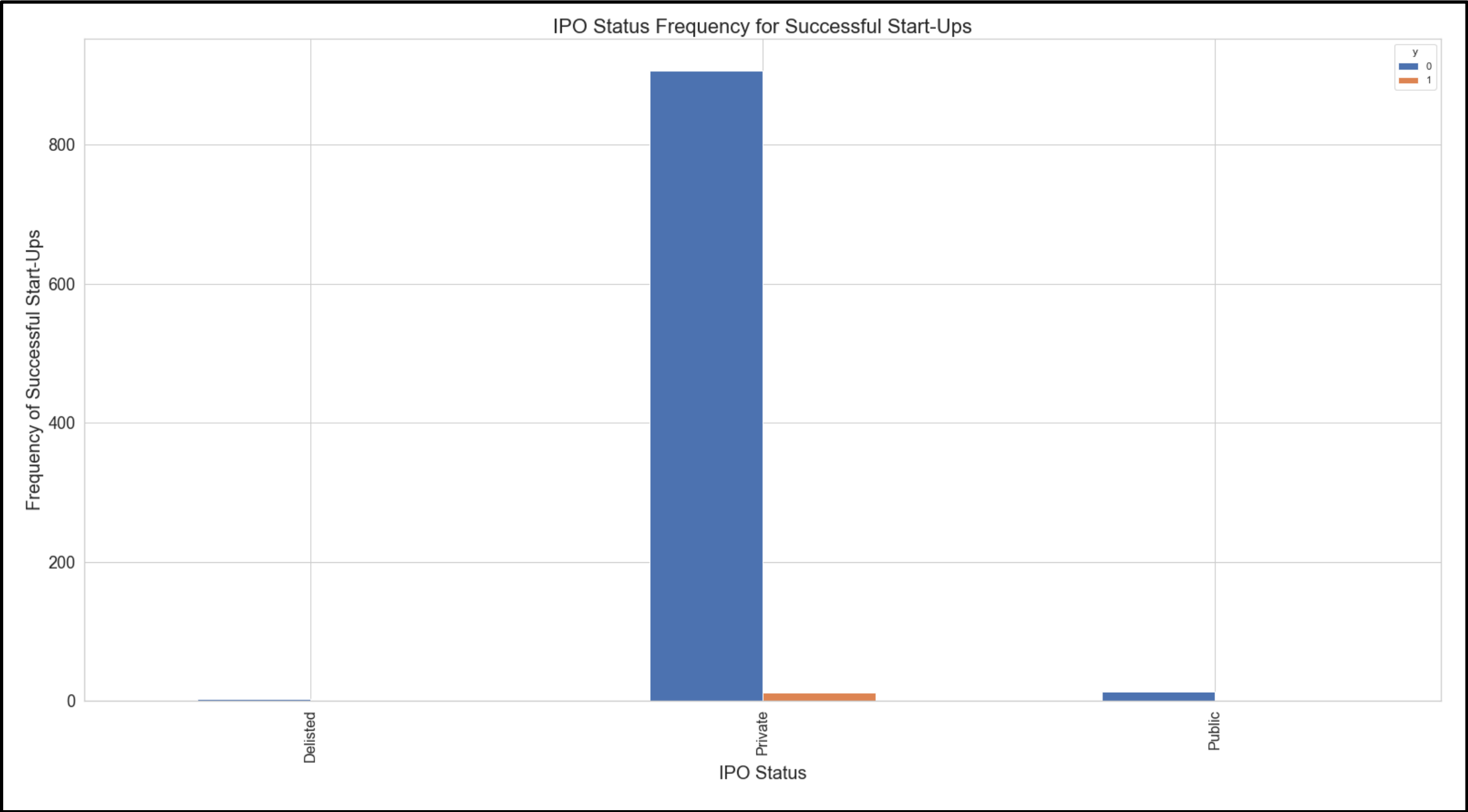
Si importa Number of Investors para éxito de Startups

Number of Investors vs Successful Start-Ups (Categórica)



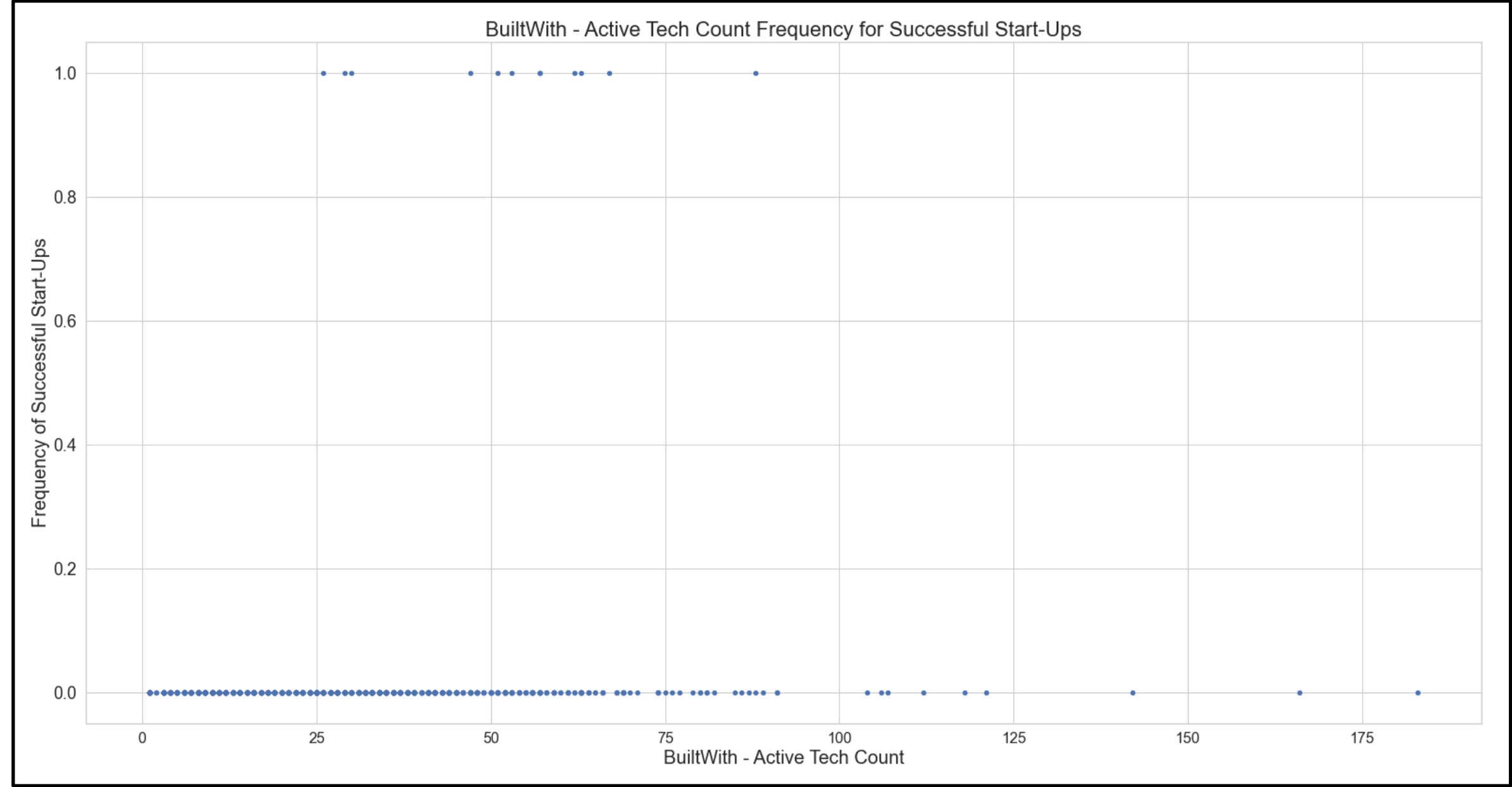
Si importa Number of Investors para éxito de Startups

IPO Status vs Successful Start-Ups



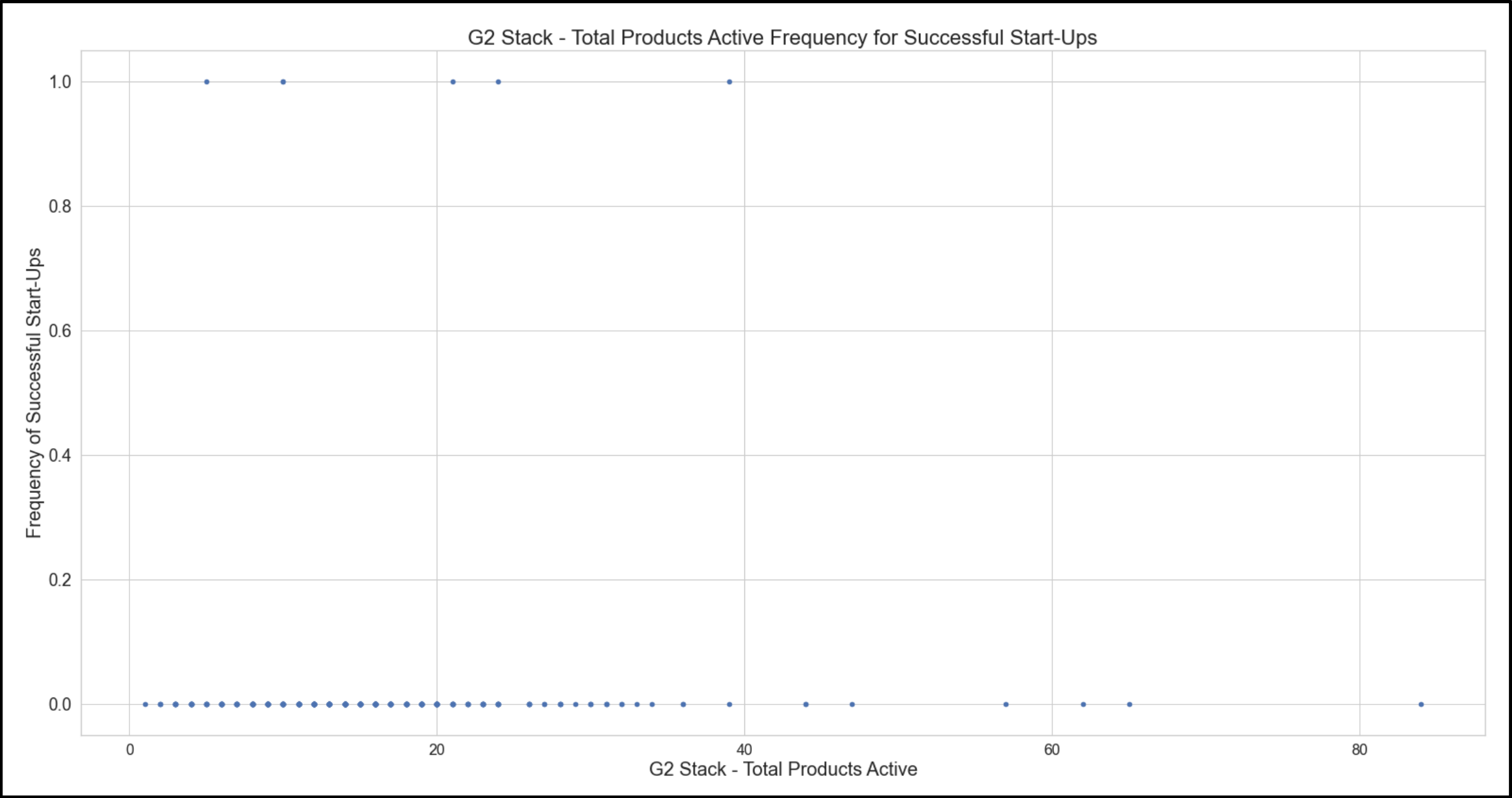
No importa IPO Status para éxito de Startups

BuiltWith - Active Tech Count vs Successful Start-Ups (Categórica)



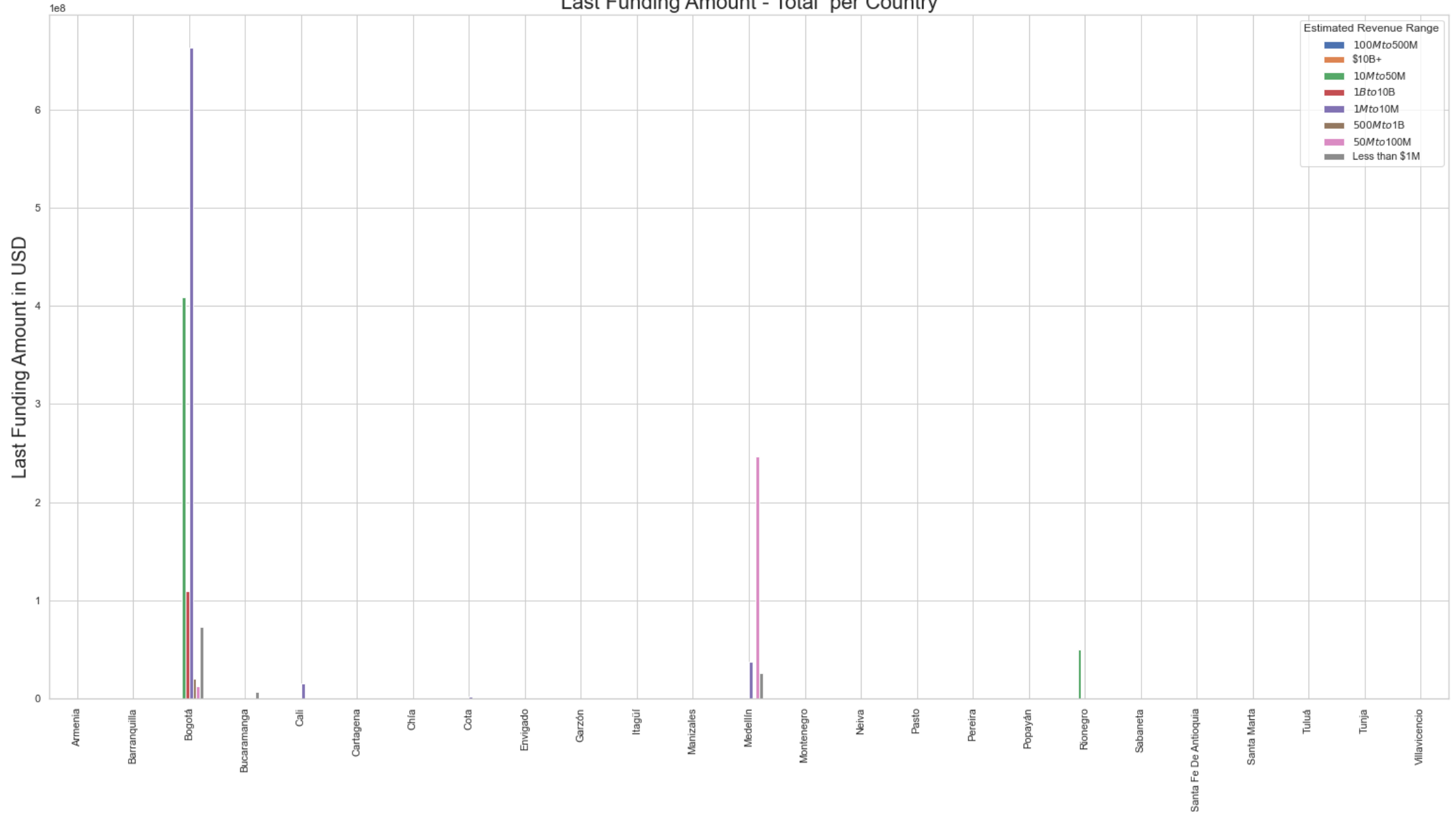
Si importa BuiltWith - Active Tech Count para éxito de Startups

G2 Stack - Total Products Active vs Successful Start-Ups



Si importa G2 Stack - Total Products Active para éxito de Startups

Last Funding Amount - Total per Country



Set of variables

<u>CBRank</u>	Númerica
Estimated Revenue Range	Categórica
Number of Articles	Númerica
Number of Founders	Númerica
Number of Employees	Categórica
Number of Funding Rounds	Númerica
Funding Status	Categórica
Last Funding Amount Currency (in USD)	Númerica
Last Equity Funding Amount Currency (in USD)	Númerica
Last Equity Funding Type	Categórica
Total Equity Funding Amount Currency (in USD)	Númerica
Total Funding Amount Currency (in USD)	Númerica
Number of Investors	Númerica
BuiltWith - Active Tech Count	Númerica
G2 Stack - Total Products Active	Númerica