

Tarea 4

5271

2 de abril de 2019

1. Algoritmos generadores de grafos.

Gracias a la versatilidad de Los grafos como modelo de representación de datos, los procesos aleatorios de generación de grafos también son relevantes en aplicaciones que van desde la física, biología a la sociología [8] .

Para realizar la tarea se seleccionaron tres algoritmos generadores de grafos de la biblioteca de networkx, con el objetivo de crear los grafos con pesos con distribución normal (como se muestra en la figura 1 de la página 2) a los que se le aplicaran los algoritmos de flujo máximo para realizar los experimentos requeridos. Los tres algoritmos escogidos son de generación aleatorias.

Los algoritmos escogidos fueron los siguientes:

- *Erdos renyi graph* (en el modelo $G(n, p)$, el grafo se construye conectando los n vértices al azar. Cada arista se incluye en el grafo con probabilidad p independiente de cualquier otro borde) [3].
- *Fast gnp random graph* (este algoritmo recibe como parámetros n números de vértices y probabilidad p de ocurrencia de aristas, el mismo devuelve un grafo aleatorio) [4].
- *Binomial graph* (este algoritmo recibe como parámetros n números de vértices y probabilidad p de ocurrencia de aristas, el mismo devuelve un grafo aleatorio, para grafos dispersos (para valores pequeños de p), *Fast gnp random graph* es un algoritmo más rápido) [5].

En la figura 2 de la 3 se muestra ejemplos de grafos generados por los algoritmos antes mencionados.

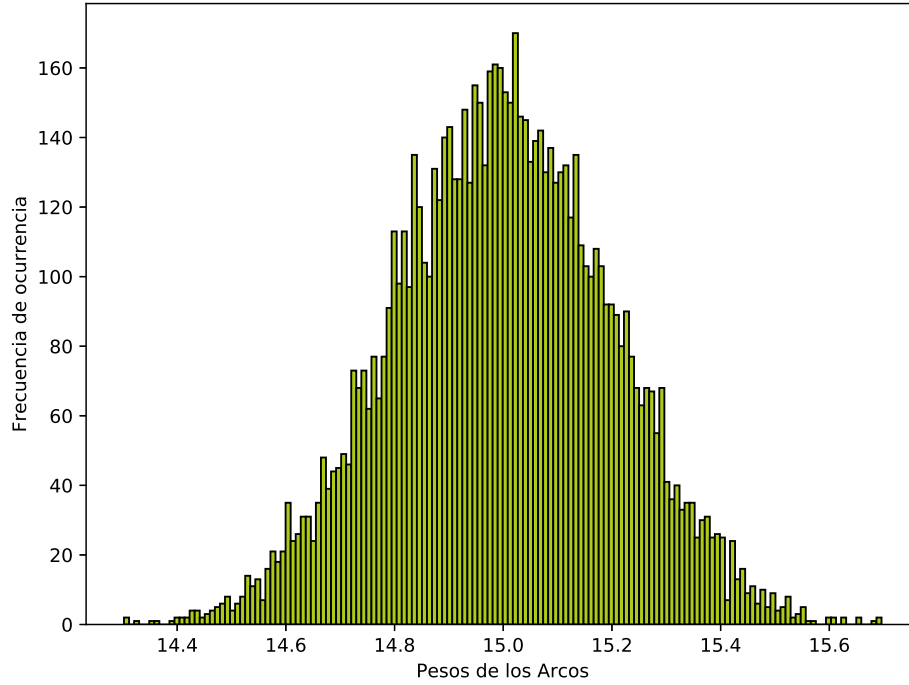


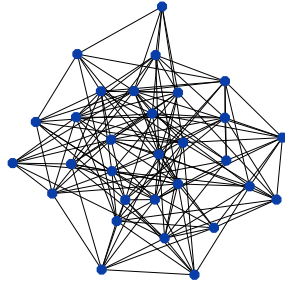
Figura 1: Histograma de distribución de los pesos.

2. Algoritmos de flujo máximo.

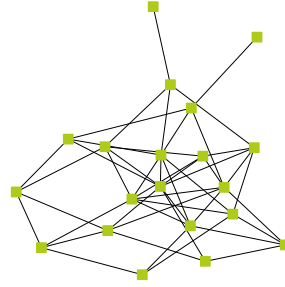
En los problemas de flujo en redes, las aristas representan vías por las que puede circular elementos: datos, agua, corriente eléctrica, entre otras. Los pesos de las aristas representan la capacidad máxima de una vía: velocidad de una conexión, volumen máximo de agua, voltaje de una línea eléctrica, entre otras; aunque es posible que la cantidad real de flujo sea menor.

El problema del flujo máximo consiste en lo siguiente: dado un grafo con pesos, $G = (V, A, W)$, que representa las capacidades máximas de los canales, un vértice fuente f y otro sumidero s en V , encontrar la cantidad máxima de flujo que puede circular desde f hasta s .

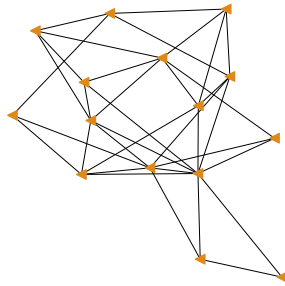
En la librería de networkx encontramos varios algoritmos con los que podemos atacar los problemas de flujo máximo, para la realización de esta tarea se escogerán tres algoritmos de dicha librería.



(a) *Binomial graph*, con 30 vértices



(b) *Erdos renyi graph*, con 20 vértices



(c) *Fast gnp random graph*, con 15 vértices

Figura 2: Ejemplo de grafos generados con los algoritmos seleccionados.

Los algoritmos escogidos fueron los siguientes:

- *Shortest augmenting path* (es uno de los enfoques más clásicos para la máxima coincidencia y los problemas de flujo máximo. Sorprendentemente, aunque esta idea es una de las técnicas más básicas, está lejos de ser completamente entendida. Es más fácil hablar de ello introduciendo el problema de emparejamiento bipartito en línea [1]. Este algoritmo encuentra el flujo máximo de un solo producto utilizando la ruta de aumento más corto y devuelve la red residual resultante después de calcular el flujo máximo).
- *maximum flow* (Encuentra la ruta por la cual pasa la máxima cantidad de flujo, recibe como parámetros un grafo G , una fuente f , un sumidero s y además una capacidad que de no tenerla, se considera que el borde tiene una capacidad infinita. Se puede aplicar en grafos tanto dirigidos como no dirigidos.) [6].

- *Preflow push* (encuentra un flujo máximo de un solo producto utilizando el algoritmo de empuje previo al flujo de la etiqueta más alta. Esta función devuelve la red residual resultante después de calcular el flujo máximo. Este algoritmo tiene un tiempo de ejecución de $O(n^2\sqrt{m})$ para n vértices y m aristas.) [7].

3. Generación de datos.

Con el objetivo de realizar las mediciones de los tiempos de ejecución de los algoritmos de flujo máximo seleccionados se desarrolló el siguiente código.

En primer lugar, se crea una función (*Algoritmo_{FM}*()) que recibe como parámetros el algoritmo de flujo máximo (*algoritmo_F*) que se va a utilizar, el grafo al que se le aplicara el algoritmo (*Graf*), la fuente (*funte*) y sumidero (*sumidero*). Al llamar esta función se genera con cada ejecución los datos que son guardados en un *data frame* del cual se muestra un fragmento en el cuadro 1 de la página 5.

```

1 def Algoritmo_FM(algoritmoF, Graf, funte, sumidero):
2     tiempos_ejecucion = []
3     for medicion in range(1, mediciones + 1):
4         t_inicio = dt.datetime.now()
5         obj = algoritmos_flujo[algoritmoF](Graf, funte, sumidero,
6         capacity="capacity")
7         t_fin = dt.datetime.now()
8         tiempo_consumido_segundos = (t_fin - t_inicio).
9         total_seconds()
10        tiempos_ejecucion.append(tiempo_consumido_segundos)
11        media = stats.mean(tiempos_ejecucion)
12
13        t_csv["grafo"].append("vertices" + str(inst_g_n) + "aristas" +
14        str(aristas))
15        t_csv["generador_grafo"].append(generador_grafo)
16        t_csv["vertices"].append(inst_g_n)
17        t_csv["densidad"].append(nx.density(Grafo))
18        t_csv["aristas"].append(aristas)
19        t_csv["f"].append(f)
20        t_csv["s"].append(s)
21        t_csv["algoritmo_fm"].append(algoritmo)
22        t_csv["media"].append(round(media, 5))
23        t_csv["mediana"].append(round(stats.median(tiempos_ejecucion),
24        5))
25        t_csv["varianza"].append(round(stats.pvariance(
26        tiempos_ejecucion, mu=media), 5))
27        t_csv["desv"].append(round(stats.pstdev(tiempos_ejecucion, mu=
28        media), 5))
29    return t_csv

```

Generar_datos.py

Cuadro 1: Fracmento de *data frame* generado.

grafo	generador	algoritmo_fm	vertices	densidad	aristas	f	s	mediana	varianza	desv
v256a8140	fast_gnp	shortest_a_path	256	0.24939	8140	33	101	0.06247	0.00004	0.00618
v256a8141	fast_gnp	edmonds_karp	256	0.24939	8140	33	101	0.04692	0.00001	0.00318
v256a8142	fast_gnp	preflow_push	256	0.24939	8140	33	101	0.08225	0.00025	0.01596
v256a8143	fast_gnp	shortest_a_path	256	0.24939	8140	33	101	0.06251	0.00004	0.00638
v256a8144	fast_gnp	edmonds_karp	256	0.24939	8140	33	101	0.05291	0.00007	0.00834
v256a8145	fast_gnp	preflow_push	256	0.24939	8140	33	101	0.08253	0.00017	0.01296
v256a8146	fast_gnp	shortest_a_path	256	0.24939	8140	33	101	0.06903	0.00020	0.01398
v256a8147	fast_gnp	edmonds_karp	256	0.24939	8140	33	101	0.05296	0.00003	0.00566

En este otro fragmentó del código es donde se generan los grafos y se llama la función (*Algoritmo_FM()*).

```

1 numero_intancias = 10
2 mediciones = 5
3 archivo_CSV = "Datost4.csv"
4 control_iteraciones = 0
5
6 generadores_grafos = {
7     "fast_gnp_random_graph": nx.fast_gnp_random_graph,
8     "binomial_graph": nx.binomial_graph,
9     "erdos_renyi_graph": nx.erdos_renyi_graph
10 }
11
12 algoritmos_flujo = {
13     "shortest_augmenting_path": shortest_augmenting_path,
14     "maximum_flow": maximum_flow,
15     "preflow_push": preflow_push
16 }
17
18 t_csv = {
19     "grafo": [], "generador_grafo": [],
20     "algoritmo_fm": [], "vertices": [],
21     "densidad": [], "aristas": [],
22     "f": [], "s": [], "media": [],
23     "mediana": [], "varianza": [],
24     "desv": []
25 }
26 for generador_grafo in generadores_grafos:
27     for inst_g_n in [round(pow(2, value + 1))
28                     for value in
29                         range(7, 11 )]:
30
31         for grafo in range(1, numero_intancias + 1):
32             f = np.random.randint(1, high=(inst_g_n - 1), dtype="
33 int")
34             s = np.random.randint(1, high=(inst_g_n - 1), dtype="
35 int")
36             while s == f:

```

```

37         f = np.random.randint(1, high=(inst_g_n - 1), dtype
    = "int")
38         s = np.random.randint(1, high=(inst_g_n - 1), dtype
    = "int")
39         Grafo = generadores_grafos[generador_grafo](inst_g_n,
    0.25, seed=None)
40         aristas = Grafo.number_of_edges()
41         p_n_distribuidos = np.random.normal(15, 0.2, aristas)
42         incremento = 0
43
44         for (u, v) in Grafo.edges():
45             Grafo.edges[u, v]["capacity"] = p_n_distribuidos[
    incremento]
46             incremento += 1
47             pesos = []
48             for (u, v) in Grafo.edges():
49                 t = Grafo.edges[u, v]["capacity"]
50                 pesos.append(t)
51             plt.figure(figsize=(8, 6))
52             n = plt.hist(pesos, bins=70, color="#932525", alpha=1,
    edgecolor='black', linewidth=1)
53             plt.xlabel('Pesos de los Arcos')
54             plt.ylabel('Frecuencia de ocurrencia')
55
56             plt.savefig("histograf.png")
57             plt.savefig("histograf.eps")
58             for instancia_grafo in range(1, 6):
59                 for algoritmo in algoritmos_flujo:
60                     d = Algoritmo_FM(algoritmo, Grafo, f, s,
    control_iteraciones)
61                     control_iteraciones=d
62             ds = pd.DataFrame(t_csv)
63             ds.to_csv(archivo_CSV, encoding="utf-8", index=None)

```

Generar_datos.py

4. Resultados del Análisis de los datos.

Con los datos recopilado de las ejecuciones de los tres algoritmos de flujo máximo con los ciento veinte grafos y sus cinco combinaciones de fuentes y sumidero, se realizó una serie de diagramas y pruebas estadísticas. A continuación, se muestra un el fragmento de código donde se lleva a cabo lo antes mencionado.

```

1 df = pd.read_csv("Datost4.csv", index_col=None, usecols=[1,2,3,4,9],
    dtype={'generador_grafo': 'category',
2
    'algoritmo_fm': 'category', 'vertices': 'category', '
    densidad': np.float64, 'mediana': np.float64} )
3 modelo = ols('mediana_log ~ generador_grafo+algoritmo_fm+vertices+
    densidad+generador_grafo*algoritmo_fm+algoritmo_fm*vertices+
    vertices*densidad+generador_grafo*vertices+generador_grafo*
    densidad+algoritmo_fm*densidad', data=df).fit()

```

```

4 print(modelo.summary())
5 modelo.csv = open("Anova.Mult.csv", 'w')
6 aov_table = sm.stats.anova_lm(modelo, typ=2)
7 df1=pd.DataFrame(aov_table)
8 df1.to_csv("modelo.csv")
9 for column in range(0, df["densidad"].count()):
10     pass
11     if df.iat[column, 3] >=0.2061 and df.iat[column, 3] <
12     0.20854:
13         df.iat[column, 3] = 1
14     elif df.iat[column, 3] >=0.20854 and df.iat[column, 3] <
15     0.21098:
16         df.iat[column, 3] = 2
17     else:
18         df.iat[column, 3] = 3
19 print(df["densidad"])
20 df['densidad'].replace({1:"baja", 2: 'media', 3:'alta' }, inplace=
21 True)
22 print(df["densidad"])
23 logX = np.log1p(df['mediana'])
24 df = df.assign(mediana_log=logX.values)
25 df.drop(['mediana'], axis= 1, inplace= True)
26 factores=["vertices", "generador-grafo", "densidad", "algoritmo-fm"]
27 plt.figure(figsize=(8, 6))
28 for i in factores:
29     print(rp.summary_cont(df['mediana_log'].groupby(df[i])))
30
31     anova = pg.anova (dv='mediana_log', between=i, data=df,
32     detailed=True ,)
33     pg._export_table (anova, ("ANOVAs"+i+".csv"))
34
35     ax=sns.boxplot(x=df["mediana_log"], y=df[i], data=df, palette="
36     cubehelix")
37
38     plt.savefig("boxplot_" + i + ".eps", bbox_inches='tight')
39     tukey = pairwise_tukeyhsd(endog = df["mediana_log"], groups= df
40     [i], alpha=0.05)
41
42     tukey.plot_simultaneous(xlabel='Tiempo', ylabel=i)
43     plt.vlines(x=49.57,ymin=-0.5,ymax=4.5, color="red")
44     plt.savefig("simultaneous_tukey" + i + ".eps", bbox_inches='
45     tight')
46
47     print(tukey.summary())
48     t_csv = open("Tukey"+i+".csv", 'w')
49     with t_csv:
50         writer = csv.writer(t_csv)
51         writer.writerow(tukey.summary())
52     plt.show()

```

Analisis_datos.py

4.1. Análisis de varianza (ANOVA).

El análisis de varianza (ANOVA) es la técnica central en el análisis de datos experimentales. La idea general de esta técnica es separar la variación total en las partes con las que contribuye cada fuente de variación en el experimento. En el caso de los diseños completamente al azar se separan la variabilidad debida a los tratamientos y la debida al error. Cuando la primera predomina sobre la segunda, es cuando se concluye que las medias son diferentes. Cuando los tratamientos no dominan contribuyen igual o menos que el error, por lo que se concluye que las medias son iguales [2].

Para analizar si los diferentes factores (algoritmo generador de grafo, algoritmo de flujo máximo, número de vértices y densidad del grafo) influían en la variable dependiente *tiempo de ejecución* se realizó un ANOVA de un factor para cada caso.

4.1.1. Influencia del algoritmo generador de grafos en el tiempo de ejecución.

El siguiente cuadro muestra el resultado de la aplicación del ANOVA.

Cuadro 2: ANOVA, relación del algoritmo generador con el *tiempo de ejecución*

Factor	SS	DF	MS	F	p-unc	np2
generador_grafo	0.280	2	0.140	0.354	0.702	0
<i>Within</i>	712.085	1797	0.396	-	-	-

En el cuadro 2 se muestra que no existen diferencia entre las medianas de los grupos de factores ya que el $p - unc$ es mayor que 0,05 por lo que se acepta la hipótesis de que el tipo de generador no influye en el *tiempo de ejecución*. Esto se puede observar en la figura 3 de la página 9.

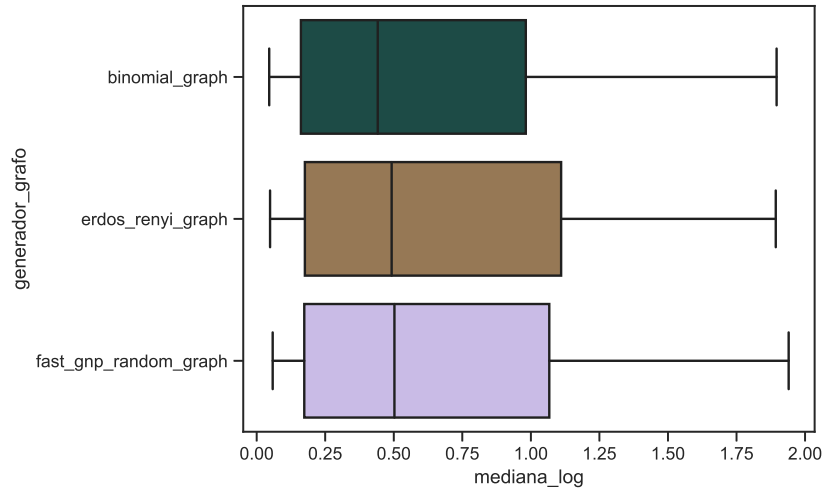


Figura 3: Diagrama de caja y bigotes que relaciona los tiempos de ejecución con los algoritmos generadores de grafos.

4.1.2. Influencia del algoritmo de flujo máximo en el tiempo de ejecución.

El siguiente cuadro muestra el resultado de la aplicación del ANOVA.

Cuadro 3: ANOVA, relación del algoritmo de flujo máximo con el *tiempo de ejecución*

Factor	SS	DF	MS	F	p-unc	np2
algoritmo_fm	1.499	2	0.749	1.895	0.151	0.002
<i>Within</i>	710.867	1797	0.396	-	-	-

En el cuadro 3 se muestra que no existen diferencia entre las medianas de los grupos de factores ya que el $p - unc$ es mayor que 0,05 por lo que se acepta la hipótesis de que el tipo de algoritmo no influye en el *tiempo de ejecución*. Esto se puede observar en la figura 4 de la página 10.

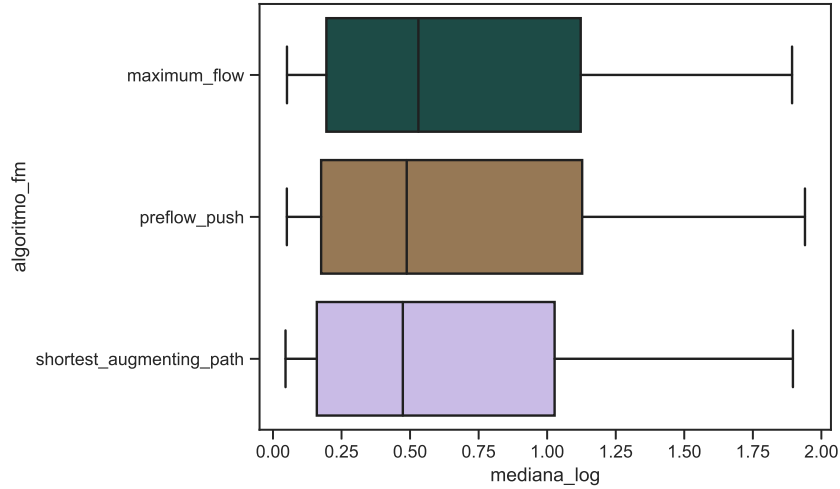


Figura 4: Diagrama de caja y bigotes que relaciona los tiempos de ejecución con los algoritmos de flujo máximo.

4.1.3. Influencia del número de vértices en el tiempo de ejecución.

El siguiente cuadro muestra el resultado de la aplicación del ANOVA.

Cuadro 4: ANOVA, relación del número de vértices con el *tiempo de ejecución*

Factor	SS	DF	MS	F	p-unc	np2
vértices	706.976	3	235.659	78536.187	0	0.992
<i>Within</i>	5.389	1796	0.003	-	-	-

En el cuadro 4 se muestra que existen grandes diferencia entre las medianas de los grupos de factores ya que el $p - unc$ es menor que 0,05 por lo que se rechaza la hipótesis de que la cantidad de vértices no influye en el *tiempo de ejecución*. Esto se puede observar en la figura 5 de la página 11. Por tal motivo se realiza la prueba de Tukey mostrara las diferencias entre las medianas de los factores, lo que se evidencia en el cuadro 5 de la página 11.

En el cuadro 5 se puede observar que en todos los casos se rechaza la hipótesis. En la figura 6 de la página 12 muestra claramente este hecho.

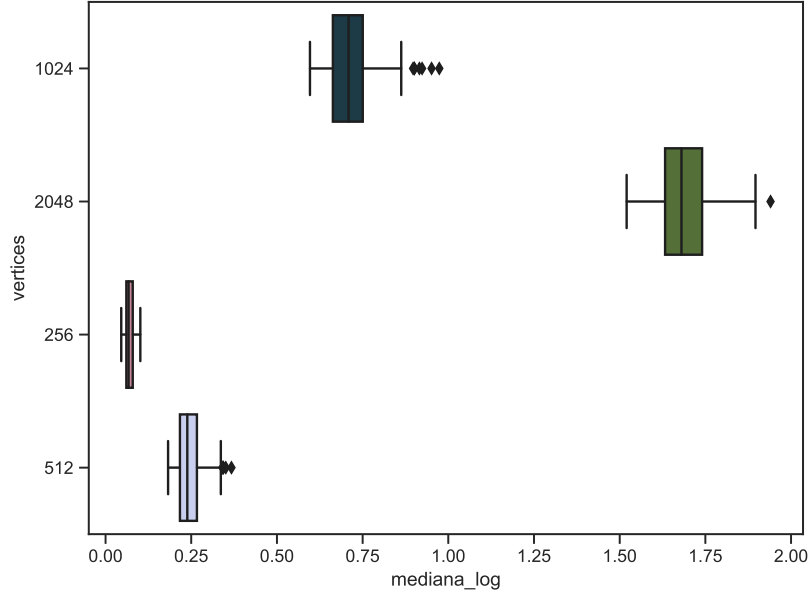


Figura 5: Diagrama de caja y bigotes que relaciona los tiempos de ejecución con el número de vértices.

Cuadro 5: Tukey, influencia del número de vértices en el tiempo de ejecución.

<i>group1</i>	<i>group2</i>	<i>meandiff</i>	<i>lower</i>	<i>upper</i>	<i>reject</i>
1024	2048	0.97	0.9606	0.9794	<i>True</i>
1024	256	-0.6445	-0.6539	-0.6352	<i>True</i>
1024	512	-0.4689	-0.4782	-0.4595	<i>True</i>
2048	256	-1.6146	-1.624	-1.6052	<i>True</i>
2048	512	-1.4389	-1.4483	-1.4295	<i>True</i>
256	512	0.1757	0.1663	0.1851	<i>True</i>

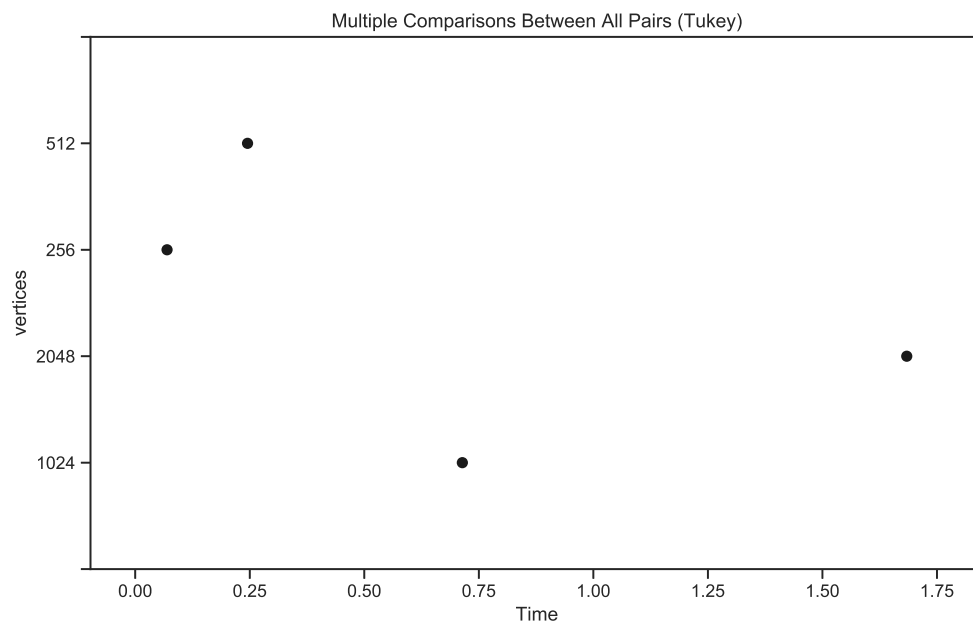


Figura 6: Diagrama simultaneo que relaciona los tiempos de ejecución con los grupos del factor número de vértices.

4.1.4. Influencia de la densidad de los grafos en el tiempo de ejecución.

En el caso del factor densidad se hizo una categorización por rangos para hacer cómoda la visualización de su relación con la variable dependiente estos rangos se obtuvieron a través de los contenedores que devolvió el histograma que se muestra en la figura 7 de la página 13.

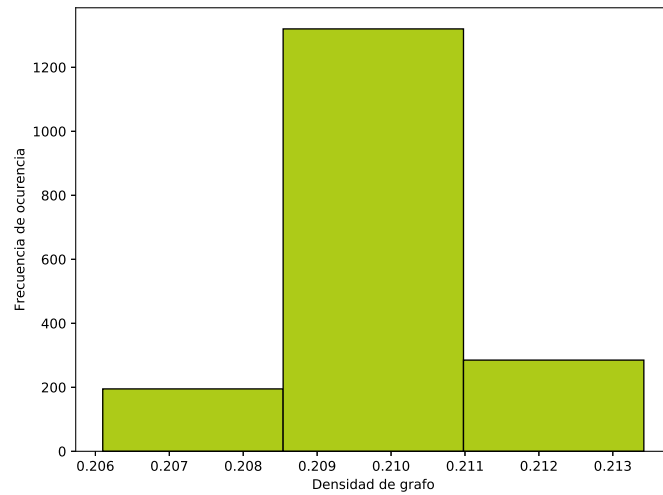


Figura 7: Histograma de densidad de los grafos.

El siguiente cuadro muestra el resultado de la aplicación del ANOVA.

Cuadro 6: ANOVA, influencia de la densidad de los grafos en el tiempo de ejecución.

Factor	SS	DF	MS	F	p-unc	np2
densidad	170.777	2	85.388	283.320	0.000	0.24
<i>Within</i>	541.589	1797	0.301	-	-	-

En el cuadro 7 se muestra que como en el cuadro 4 existen grandes diferencia entre las medianas de los grupos de factores ya que el $p-unc$ es menor que 0,05 por lo que se rechaza la hipótesis de que la densidad de los grafos no influye en el *tiempo de ejecución*. Esto se puede observar en la figura 8 de la página 14.

Por tal motivo se realiza la prueba de Tukey mostrara las diferencias entre las medianas de los factores, lo que se evidencia en el cuadro 7 de la página 14.

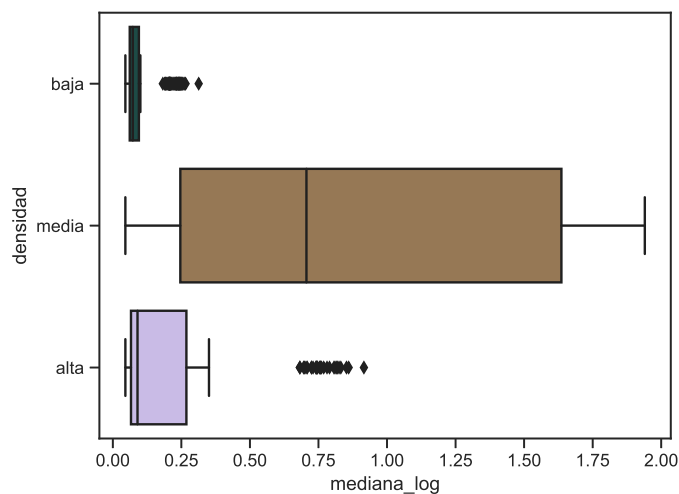


Figura 8: Diagrama de caja y bigotes que relaciona los tiempos de ejecución con la densidad de los grafos.

Cuadro 7: Tukey, influencia de la densidad de los grafos en el tiempo de ejecución.

<i>group1</i>	<i>group2</i>	<i>meandiff</i>	<i>lower</i>	<i>upper</i>	<i>reject</i>
alta	baja	-0.1085	-0.2281	0.0112	<i>False</i>
alta	media	0.6497	0.5656	0.7338	<i>True</i>
baja	media	0.7582	0.6594	0.8569	<i>True</i>

En el cuadro 7 se puede observar que en dos de los tres casos se rechaza la hipótesis y en el pareo de densidad alta y baja no existen diferencias estadísticas en cuanto a la mediana. En la figura 9 de la página 15 muestra claramente este hecho.

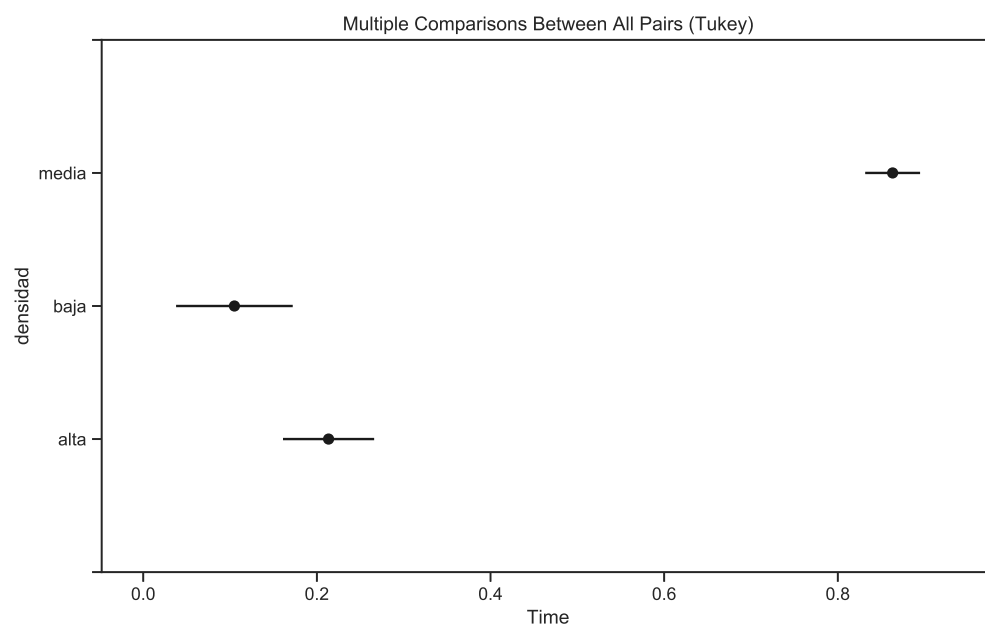


Figura 9: Diagrama simultaneo que relaciona los tiempos de ejecución con los grupos del factor densidad de los grafos.

4.1.5. Influencia de los cuatro factores (algoritmo generador de grafo, algoritmo de flujo máximo, numero de vértices y densidad del grafo) en el tiempo de ejecución.

Para analizar este caso se realizo ANOVA multifactorial dando como resultado el cuadro 8 de la página 16.

Cuadro 8: ANOVA multifactor, influencia de los cuatro factores en el tiempo de ejecución.

	sum_sq	df	F	PR(>F)
generador_grafo	0.0291	2	33.0499	0.0000
algoritmo_fm	-0.0001	2	-0.1439	1.0000
vertices	46.4767	3	35163.6974	0.0000
densidad	0.0000	2	0.0000	0.9999
generador_grafo:algoritmo_fm	0.0008	4	0.4714	0.7567
algoritmo_fm:vertices	0.0001	6	0.0248	0.8748
vertices:densidad	0.0001	6	0.0425	0.9584
generador_grafo:vertices	0.0651	6	24.6313	0.0000
generador_grafo:densidad	0.0000	4	0.0181	0.8930
algoritmo_fm:densidad	0.0001	4	0.0425	0.9584
Residual	0.7798	1770		

En el cuadro 8 se puede observar que los factores que más influyen en el tiempo de ejecución son el número de vértices y la densidad de los grafos, que además existe una relación entre el número de nodos y la densidad de los grafos.

Referencias

- [1] Bartłomiej Bosek, Dariusz Leniowski, Piotr Sankowski, and Anna Zych-Pawlewicz. Shortest augmenting paths for online matchings on trees. *Theory of Computing Systems*, 62(2):337–348, Feb 2018.
- [2] Humberto Gutiérrez and Román de la Vara. *Análisis y diseño de experimentos*. The McGraw-Hill Companies, Inc., segunda edición edition, 2008. 60–74.
- [3] Desarrolladores NetworkX. https://networkx.github.io/documentation/stable/reference/generated/networkx.generators.random_graphs.erdos_renyi_graph.html#networkx.generators.random_graphs.erdos_renyi_graph. Accessed: 01-04-2019.
- [4] Desarrolladores NetworkX. https://networkx.github.io/documentation/stable/reference/generated/networkx.generators.random_graphs.fast_gnp_random_graph.html#networkx.generators.random_graphs.fast_gnp_random_graph. Accessed: 01-04-2019.
- [5] Desarrolladores NetworkX. https://networkx.github.io/documentation/stable/reference/generated/networkx.generators.random_graphs.binomial_graph.html#networkx.generators.random_graphs.binomial_graph. Accessed: 01-04-2019.
- [6] Desarrolladores NetworkX. https://networkx.github.io/documentation/stable/reference/algorithms/generated/networkx.algorithms.flow.maximum_flow.html#networkx.algorithms.flow.maximum_flow. Accessed: 01-04-2019.
- [7] Desarrolladores NetworkX. https://networkx.github.io/documentation/networkx-1.10/reference/generated/networkx.algorithms.coloring.greedy_color.html. Accessed: 18-03-2019.
- [8] Sadegh Nobari, Xuesong Lu, Panagiotis Karras, and Stéphane Bressan. Fast random graph generation. In *Proceedings of the 14th International Conference on Extending Database Technology, EDBT/ICDT '11*, pages 331–342, New York, NY, USA, 2011. ACM.