

PRG (ETS d'Enginyeria Informàtica) - Curs 2016-2017
Pràctica 3. Mesura empírica de la complexitat computacional
(2 sessions)

Departament de Sistemes Informàtics i Computació
Universitat Politècnica de València



Índex

1	Context i treball previ	1
2	Mesura del cost del mètode de Cerca Lineal	2
2.1	Definició del problema	2
2.2	Anàlisi de casos	2
2.3	Estructura d'un experiment de mesura	3
3	Representació gràfica: <i>Gnuplot</i>	4
4	Treball en el laboratori: anàlisi del cost de seleccio	8
5	Treball en el laboratori: anàlisi del cost d'insercio	10
6	Activitat addicional: anàlisi del cost de mergeSort	11
7	Avaluació	11

1 Context i treball previ

En el context acadèmic, esta pràctica correspon al “*Tema 2. Anàlisi d'algorismes. Eficiència. Ordenació.*”. Els objectius són els següents:

- Introduir l'anàlisi d'algorismes al laboratori, fent servir un entorn real de programació: anàlisi empírica.
- Representar gràficament el creixement dels recursos temporals emprats per confirmar els resultats teòrics.
- Inferir funcions aproximades per definir el comportament temporal d'un algorisme.
- Usar els resultats empírics per fer comparacions i prediccions.

Abans de la sessió de laboratori, l'alumne/a ha de llegir el butlletí de pràctiques i tractar de resoldre (tant com siga possible) els problemes proposats. La pràctica es realitzarà durant 2 sessions.

2 Mesura del cost del mètode de Cerca Lineal

En aquesta secció es presenta un problema complet d'anàlisi empírica. El problema és la cerca lineal, és a dir, la cerca d'un element en un array unidimensional que pot ser no estiga ordenat.

2.1 Definició del problema

El problema de la cerca lineal consisteix en, donat un array `a` d'elements d'un cert conjunt (per exemple, els números enters) i donat un element arbitrari `e` d'eixe conjunt (un número enter concret), tornar la primera posició de l'array que continga l'element `e`. Si `e` no es trobara en `a`, es tornaria un índex invàlid (per exemple, -1) com a resultat, indicant que l'element no s'ha trobat. El següent mètode resol aquest problema sobre un array d'int:

```
public static int cercaLineal(int[] a, int e) {
    int i = 0;
    while (i < a.length && a[i] != e) { i++; }
    if (i < a.length) { return i; }
    else { return -1; }
}
```

2.2 Anàlisi de casos

Quan s'enfrontem al problema d'analitzar un algorisme, el primer paràmetre que s'ha de definir és la *talla* del problema. En este cas, la talla del problema és clarament la grandària de l'array, ja que determina el nombre d'iteracions del seu bucle (per la condició `i < a.length`).

A més d'això, hem d'estudiar si l'algorisme presenta *instàncies significatives* o no. La cerca lineal presenta instàncies significatives, és a dir, *casos millor i pitjor*. Estos casos es defineixen per la segona condició del bucle (`a[i] != e`):

- **Cas millor:** quan l'element `e` es troba en la primera posició d'`a` (és a dir, `a[0] == e`), ja que en eixe cas el bucle no executaria cap de les seues iteracions previstes.
- **Cas pitjor:** quan l'element `e` no es troba en `a`, ja que per a verificar eixe fet s'ha d'explorar completament tot l'array.

Amb esta anàlisi prèvia, podem concloure que el cost en el cas millor és constant i que en el cas pitjor és lineal. Per tant, anomenant $n = a.length$, les fites asimptòtiques de l'algorisme són $T(n) \in \Omega(1), O(n)$.

El cas promedi és difícil de calcular. Es poden fer algunes assumpcions per simplificar els càlculs per este cas. Per exemple, podem assumir que l'element que es cerca sempre és dins de l'array, i que la probabilitat de trobar l'element en qualsevol posició és la mateixa. En este cas, la fita final del cas promedi ens diu que $T^\mu(n) \in \Theta(n)$.

2.3 Estructura d'un experiment de mesura

L'anàlisi empírica s'ha de fer després de l'anàlisi teòrica. Al dissenyar una anàlisi empírica, s'han de prendre en consideració els següents punts:

- **La mesura de temps ha de fer-se per diverses talles:** l'objectiu és obtenir una funció de cost que té com paràmetre la talla del problema; per tant, han d'emprar-se diverses talles per obtenir el perfil de la funció.
- **Les instàncies significatives han de mesurar-se separatament:** els casos millor, pitjor i promedi presenten habitualment diferents taxes de creixement i, per tant, diferents funcions de cost; aleshores, han de mesurar-se a diferents parts del codi.
- **Per traure resultats significatius s'han de prendre diverses mesures:** una única mesura per cada talla no és significativa, ja que pot veure's afectada per condicions de l'entorn (per exemple, l'execució d'altres processos a l'ordinador); per tant, per garantir resultats correctes s'han de prendre diverses mesures, de les que s'haurà d'obtenir el seu valor mitjà; aquest valor en terme mitjà es podrà considerar com un resultat significatiu de la mesura.

La mesura de temps d'un algorisme pot presentar-se com el següent procés:

1. Prendre temps actual t_I (temps inicial).
2. Executar l'algorisme (mètode).
3. Prendre temps actual t_F (temps final).
4. La diferència entre t_F i t_I és el temps que ha emprat l'algorisme en resoldre el problema.

Este procés es pot fer usant un rellotge extern, però és més precís usar el rellotge intern. Java aporta el mètode `static long nanoTime()`, en `java.lang.System`, que retorna el valor actual del temporitzador més precís del sistema en nanosegons (encara que la resolució pot ser menor, però com a mínim és de mil·lisegons). Per tant, el codi Java per mesurar temps és semblant a:

```
long ti, tf, tt;
ti = System.nanoTime();
// Crida al mètode que es vol temporitzar
tf = System.nanoTime();
tt = tf - ti;
```

on a la variable `tt` s'enregistrarà el temps que el mètode ha invertit en resoldre el problema. Esta mesura s'ha de fer moltes voltes i es calcula el temps promedi. Tanmateix, per casos extremadament ràpids (per exemple, el cas millor de la cerca lineal), és habitual incloure el bucle de repeticions dins de la mesura de temps, considerant menyspreable la sobrecàrrega del bucle. O també es pot repetir la mesura un nombre de vegades bastant gran (per exemple, en el cas millor de la cerca lineal, un nombre de vegades molt més gran que per als casos pitjor i promedi).

Finalment, les mesures de temps s'han de representar apropiadament. La forma usual és utilitzar una taula que mostre en cada fila la talla del problema i els temps mesurats per a cadascuna de les instàncies. Aquests temps han de venir expressats en alguna unitat de mesura (microsegons, mil·lisegons, etc) que facilite la lectura dels valors, és aconsellable que la unitat aparega com un comentari a la taula. Cal notar, a més, que en tractar-se de valors mitjans, és raonable que vinguin expressats amb decimals. La forma típica d'aquesta taula és semblant a la següent:

```
# Cerca lineal. Temps en microsegons
# Talla      Pitjor      Millor      Promedi
#-----
10000        4.957        1.793        3.187
20000        8.306        1.790        4.964
30000       11.589        1.793        6.662
40000       15.002        1.793        8.353
50000       18.371        1.793       10.131
60000       21.803        1.793       11.856
.....
```

Activitat 1: creació del projecte *BlueJ* pract3

Obre *BlueJ* en el directori de treball de l'assignatura (prg) i crea un nou projecte **pract3** amb les classes `AlgorismesMesurables.java` que conté, entre altres, el mètode `cercaLineal(int[], int)` i `MesuraCercaLineal.java` que implementa el codi que realitza l'anàlisi per a les distintes instàncies significatives d'aquest algorisme. Tens disponibles aquestes classes en `Recursos/Laboratori/Pràctica 3`, dins de `poliformaT` de PRG.

Activitat 2: obtenció de temps de cerca

Executa el mètode `mesuraCercaLineal()` de la classe `MesuraCercaLineal` per tal d'obtenir la taula de resultats de temps i guarda-la en un fitxer de nom `cercaLineal.out`; per a això pots utilitzar l'opció corresponent disponible a la finestra d'execució de *BlueJ* o redirigir l'eixida del programa executant des de la línia de comandaments:

```
java MesuraCercaLineal > cercaLineal.out
```

3 Representació gràfica: *Gnuplot*

Els resultats numèrics usualment s'interpreten millor amb la seua representació gràfica. En esta secció mostrem com emprar la ferramenta **Gnuplot** per obtindre representacions gràfiques dels resultats i per obtindre funcions de cost que s'aproximen als resultats empírics, les quals poden usar-se per comparar adequadament els algorismes i per obtindre prediccions. Per tal d'usar aquesta eina escriurem **gnuplot** en la línia d'ordres (terminal). *Gnuplot* accepta ordres amb modificadors; les ordres més importants són les següents:

- **plot**: dibuixa dades d'un fitxer o funcions predefinides; alguns modificadors són:

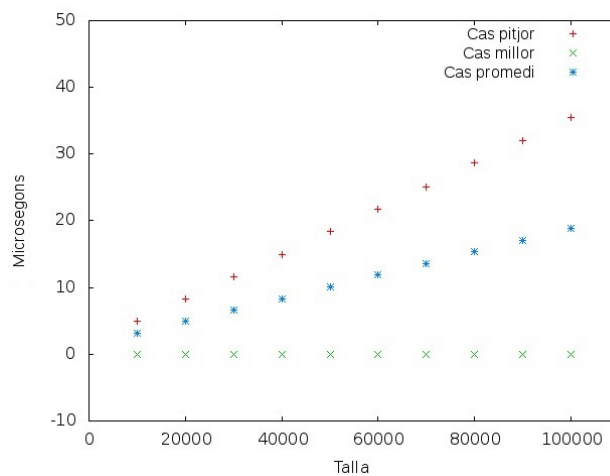
- *fitxer*: s'especifica entre cometes dobles i diu on és el fitxer amb les dades a dibuixar; les línies del fitxer que comencen amb el símbol # s'ignoren.
 - **title** *text*: especifica el títol que ha de donar-se als punts (llegenda).
 - **using** *i:j*: especifica les columnes del fitxer de dades que s'empraran (*i* per l'eix X, *j* per l'eix Y).
 - **with** *format*: especifica el format de dibuix (els usuals són **lines**, **points** i **linespoints**).
- **replot**: amb el mateix significat i modificadors que *plot*, però sense esborrar la finestra gràfica, i així permet veure diverses gràfiques alhora; *replot* només redibuixa la finestra gràfica.
 - **set xrange** [*inici:fi*], **set yrange** [*inici:fi*]: fixa el rang de valors de l'eix X (respectivament Y) entre *inici* i *fi*.
 - **set xtics** *interval*, **set ytics** *interval*: fixa l'interval entre marques a l'eix X (respectivament Y).
 - **set xlabel** *text*, **set ylabel** *text*: fixa l'etiqueta de l'eix X (respectivament Y).
 - **load** *fitxer*: carrega un fitxer de text amb ordres de *Gnuplot* que són executades per *Gnuplot*.
 - **fit** *funció fitxer* **using** *i:j* **via** *paràmetres*: permet ajustar una funció predefinida amb alguns paràmetres lliures a un conjunt de dades d'un fitxer. On:
 - *funció*: indica el nom de la funció a ajustar.
 - *fitxer*: indica el nom del fitxer amb les dades a ajustar (s'especifica entre cometes dobles).
 - **using** *i:j*: especifica les columnes del fitxer de dades que s'usaran (*i* per a l'eix X i *j* per a l'eix Y).
 - **via** *paràmetres*: especifica els paràmetres (separats per comes) de la funció a ajustar.

Activitat 3: representació i anàlisi dels resultats empírics

Per tal de representar els temps de la taula emmagatzemada al fitxer `cercaLineal.out`, arranca *Gnuplot* i escriu les següents ordres:

```
gnuplot> set xrange [0:110000]
gnuplot> set yrange [-10:50]
gnuplot> set xtics 20000
gnuplot> set ytics 10
gnuplot> set xlabel "Talla"
gnuplot> set ylabel "Microsegons"
gnuplot> plot "cercaLineal.out" using 1:2 title "Cas pitjor"
gnuplot> replot "cercaLineal.out" using 1:3 title "Cas millor"
gnuplot> replot "cercaLineal.out" using 1:4 title "Cas promedi"
```

La imatge mostrada ha de ser semblant a la presentada en la següent figura:



Per a guardar l'eixida gràfica en un fitxer .jpg, has d'executar les ordres següents:

```
gnuplot> set term jpeg
gnuplot> set output "cercaLineal.jpg"
gnuplot> replot
```

El fitxer `cercaLineal.jpg` s'emmagatzemarà en el directori actual amb els continguts de la vista gràfica. Per a redirigir l'eixida de nou al terminal s'ha d'executar:

```
gnuplot> set term x11
gnuplot> set output
```

Activitat 4: aproximació de funcions als resultats empírics

Gnuplot permet ajustar els valors de les columnes d'un fitxer de dades a una funció definida prèviament. Per exemple, si es sospita que determinats valors segueixen una distribució quadràtica es pot, mitjançant *Gnuplot*, ajustar aquests valors a una paràbola. El resultat del procés d'ajust, seran els coeficients de la paràbola que millor s'aproxime a les dades de l'arxiu.

Per fer aquest tipus d'ajust s'utilitza l'ordre `fit` que, com s'ha dit, permet obtenir funcions aproximades que mostren el comportament d'un algorisme. Per exemple, com el cas pitjor i promedi de la cerca lineal presenten un cost lineal, és possible conèixer la diferència entre tots dos, ajustant prèviament cadascuna de les columnes de dades corresponents a una funció d'aquest tipus (lineal) per, després, comparar les funcions obtingudes.

Executa la següent seqüència de comandaments *Gnuplot* per realitzar l'ajust de les dades del cas pitjor mitjançant una funció lineal:

```
gnuplot> f(x)=a*x+b
gnuplot> fit f(x) "cercaLineal.out" using 1:2 via a,b
```

El resultat serà similar al que segueix:

```
...
Final set of parameters          Asymptotic Standard Error
=====
a                                = 0.000339198      +/- 1.051e-06    (0.3098%)
b                                = 1.47013        +/- 0.0652      (4.435%)
...
```

Ara fes el mateix però ajustant les dades del cas promedi mitjançant una altra funció lineal:

```
gnuplot> g(x)=c*x+d
gnuplot> fit g(x) "cercaLineal.out" using 1:4 via c,d
```

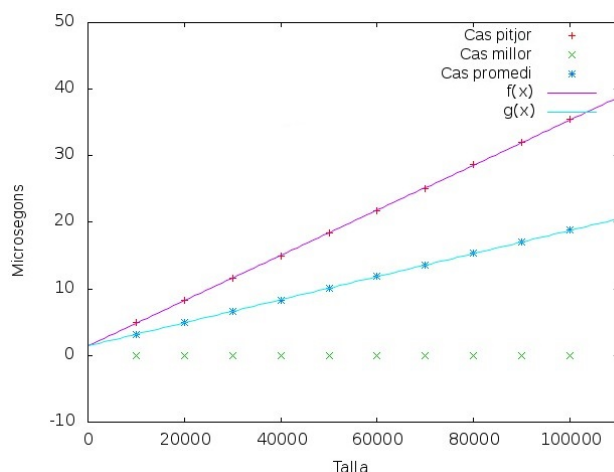
```
...
Final set of parameters          Asymptotic Standard Error
=====
c                                = 0.000173287      +/- 3.002e-07    (0.1733%)
d                                = 1.4602         +/- 0.01863      (1.276%)
...
```

per tant, s'hi pot veure que la relació de creixement entre els casos pitjor i promedi és la relació entre les pendents de les funcions lineals corresponents, és a dir, $0.000339198/0.000173287 \approx 2$ per a les dades de la nostra taula. Per tant, es pot inferir que el cas pitjor és aproximadament el doble de lent que el cas promedi.

Les funcions d'ajust recentment estimades es poden graficar junt als anteriors resultats:

```
gnuplot> replot f(x),g(x)
```

apreciant així que les dades experimentals segueixen amb gran fidelitat les funcions calculades:



L'estimació d'aquestes funcions aproximades es pot utilitzar per a fer prediccions. Per exemple, en el cas promedi, per a un array de mil milions d'enters (10^9), el temps requerit per a la seua execució segons aquestos càlculs s'estimaria de la següent manera:

$$g(x) = c * x + d = 0.000173287 * x + 1.4602$$

i substituint x pel valor de la talla del qual volem fer la predicció, 10^9 ,

$$g(10^9) = 0.000173287 * 10^9 + 1.4602 \approx 173288$$

microsegons, es a dir, prop d'un terç de segon, aproximadament.

4 Treball en el laboratori: anàlisi del cost de seleccio

L'estratègia Selecció Directa ordena un array amb una complexitat temporal $\Theta(n^2)$, sent n el nombre d'elements de l'array, i no presenta instàncies significatives per al cost. Per això, n'hi ha prou en realitzar l'estudi del cost independentment del contingut de l'array, considerant aleshores que les dades del mateix es generen aleatòriament.

Activitat 5: mètode per a generar un array amb valors aleatoris

Afegeix al teu projecte `pract3` la classe `MesuraOrdenacio.java` disponible en **Recursos: Laboratorio: Práctica 3** de la **PoliformaT** de PRG. Abans de començar amb l'anàlisi del cost del mètode `seleccio(int[])` que implementa aquesta estratègia d'ordenació, definit en la classe `AlgorismesMesurables`, has d'escriure en la classe `MesuraOrdenacio` un mètode per a crear un array d'enters generats de forma aleatòria d'una grandària determinada d'acord al següent perfil:

```
/* Genera un array d'int de talla t amb valors compresos entre 0 i t.
 * @param t int, la talla.
 * @result int[], l'array generat.
 */
private static int[] crearArrayAleatori(int t) { ... }
```

Activitat 6: temps d'execució d'una única crida al mètode

Completa el mètode `mesuraSeleccio()` de la classe `MesuraOrdenacio` amb les instruccions necessàries per tal de:

1. Crear un array aleatori `a` de talla 100 utilitzant el mètode `crearArrayAleatori(int)`.
2. Cridar al mètode `System.nanoTime()` per a obtenir en una variable `ti` (de tipus `long`) el valor del rellotge (en nanosegons) abans de cridar al mètode que volem temporitzar.
3. Cridar al mètode `seleccio(int[])` de la classe `AlgorismesMesurables` per a ordenar l'array `a`.
4. Tornar a cridar al mètode `System.nanoTime()` per a obtenir en la variable `tf` el valor del rellotge una vegada acabada l'ordenació.
5. Calcular la diferència de temps (`tf - ti`) per a saber el temps que ha requerit el mètode `seleccio(int[])` per a ordenar l'array `a`.
6. Mostrar per pantalla una filera de dades amb la talla de l'array i el temps en microsegons.

Activitat 7: temps d'execució per a una única talla donada

Com s'ha comentat anteriorment, prendre una única mesura per estimar el cost d'un mètode per a una determinada talla no és un procediment adequat, ja que aquesta mesura pot veure's afectada per condicions de l'entorn. Així, per garantir resultats correctes s'ha de repetir la mesura i després promediar sobre el nombre de mesures que s'hagen pres.

Defineix la constant `REPETICIONS` (amb valor 200) a la classe `MesuraOrdenacio` i completa el mètode `mesuraSeleccio()` d'aquesta classe amb un bucle per a repetir aquest nombre de vegades el bloc d'instruccions que ja tenim escrit de l'activitat anterior i calcula el temps promedi. Adona't que per a minimitzar la dependència del nostre experiment respecte d'una instància particular de l'array generat de forma aleatòria, resulta convenient que **en cada repetició es genere un nou array aleatori diferent**. De no ser així, fixa't que a més a més l'array a ordenar ja estaria ordenat a priori a partir de la primera repetició, fet que no invalidaria l'experiment formalment però estadísticament no seria significatiu. Calcula el temps promedi per repetició (en μs) i mostra'l per pantalla junt a la talla de l'array.

Activitat 8: temps d'execució per a diferents talles

Defineix en la classe `MesuraOrdenacio` les constants `INITALLA`, `MAXTALLA` i `INCRRTALLA` per a representar respectivament el valor de la talla més menut a considerar (1000), el valor més gran (10000) i l'increment de talla (1000). Completa el mètode `mesuraSeleccio()` d'aquesta classe afegint un bucle per a repetir el càlcul del temps d'execució per a talles des de `INITALLA` fins `MAXTALLA` amb increments de `INCRRTALLA`; és a dir, per a talles 1000, 2000, 3000, ..., 10000. El mètode ha de mostrar per pantalla una taula com la que segueix en la que el temps es done en microsegons:

```
# Selecció directa. Temps en microsegons
# Talla      Promedi
# -----
# 1000      403.346
# 2000     1348.255
# 3000     3061.948
# ...
```

Activitat 9: representació gràfica dels resultats

- Executa el mètode `mesuraSeleccio()`, guarda la taula resultat en un fitxer i, utilitzant *Gnuplot*, mostra els resultats en una gràfica en la que l'eix X represente la talla i l'eix Y el temps d'ordenació en microsegons.
- Ajusta els resultats obtinguts a una funció quadràtica ($f(x)=a*x*x+b*x+c$), observant els valors dels paràmetres d'ajust.
- Torna a construir la gràfica mostrant, a més dels punts experimentals, la funció d'ajust. Etiqueta adequadament els eixos, els punts, i la funció d'ajust, afegeix un títol a la gràfica i guarda-la en un fitxer .jpg.
- Utilitza la funció d'ajust per a predir quin seria el temps necessari per a ordenar amb el mètode `seleccio(int[])` un array de 800000 enters.

5 Treball en el laboratori: anàlisi del cost d'insercio

L'estratègia Inserció Directa ordena un array amb una complexitat temporal $\Omega(n)$ i $O(n^2)$, sent n el nombre d'elements de l'array, presentant instàncies significatives per al cost: el cas millor quan l'array ja està ordenat (de forma creixent), i el cas pitjor quan l'array també està ordenat, però a l'inrevés, es a dir, de manera decreixent. Per això, és necessari realitzar l'estudi del comportament del mètode en el cas millor (amb arrays ja ordenats), en el cas pitjor (amb arrays ordenats de forma decreixent) i en el cas promedi (amb arrays generats aleatòriament).

Activitat 10: creació d'arrays ordenats

En la classe `AlgorismesMesurables` està definit el mètode `insercio(int[])` que implementa aquesta estratègia d'ordenació. Per a poder analitzar aquest mètode s'han d'escriure en aquesta mateixa classe dos mètodes per a crear arrays d'enters de diferents grandàries, de manera que el seu contingut estiga ordenat, respectivament, de forma creixent i decreixent; els seus perfils serien:

```
/* Genera un array d'int de talla t ordenat de forma creixent.
 * @param t int, la talla.
 * @result int[], l'array generat.
 */
private static int[] crearArrayOrdCreixent(int t) { ... }

/* Genera un array d'int de talla t ordenat de forma decreixent.
 * @param t int, la talla.
 * @result int[], l'array generat.
 */
private static int[] crearArrayOrdDecreixent(int t) { ... }
```

Activitat 11: anàlisi empírica del cost del mètode insercio

Completa el mètode `mesuraInsercio()` de la classe `MesuraOrdenacio` per a estudiar el comportament del mètode `insercio(int[])` per als casos pitjor, millor, i promedi, talles des de `INITALLA` fins `MAXTALLA` amb increments de `INCR TALLA`; és a dir, per a talles 1000, 2000, 3000, ..., 10000. El mètode ha de mostrar per pantalla una taula com la que segueix:

```
# Inserció directa. Temps en microsegons
# Talla Millor Pitjor Promedi
# -----
# 1000 0.025 422.532 134.647
# 2000 0.029 848.167 405.849
# 3000 0.040 1904.827 919.622
# ...
```

Fixa't que per al algorisme `insercio`, la recomanació que us vam fer a l'activitat 7, sobre que **en cada repetició s'havia de regenerar un nou array aleatori diferent**, en aquest cas resulta absolutament obligatòria. Si no, els resultats en els casos pitjor i promedi no serien vàlids ja que l'array ja estaria ordenat (sent el cas millor) en les restants repeticions.

Activitat 12: representació gràfica dels resultats

- Executa el mètode `mesuraInsercio()`, guarda la taula resultat en un fitxer i, utilitzant *Gnuplot*, mostra els resultats en una gràfica en la que l'eix X represente la talla i l'eix Y el temps d'ordenació en microsegons. Has de dibuixar els punts experimentals obtinguts per als tres casos.
- Ajusta els resultats obtinguts a les funcions previstes de l'anàlisi teòrica (cas millor a funció lineal, casos pitjor i promedi a funcions quadràtiques) observant els valors dels paràmetres d'ajust.
- Torna a construir la gràfica mostrant, a més dels punts experimentals, les tres funcions d'ajust. Etiqueta adequadament els eixos, els punts, i les funcions d'ajust, afegeix un títol a la gràfica i guarda-la en un fitxer .jpg.
- Utilitza les funcions d'ajust per a predir quin seria el temps necessari per a ordenar un array de 800000 enters mitjançant el mètode `insercio(int[])` si: a) l'array ja està ordenat creixentment; b) si l'array també està ordenat, però en sentit decreixent; i c) per a un array amb valors aleatoris.

6 Activitat addicional: anàlisi del cost de mergeSort

Completa el mètode `mesuraMergeSort()` de la classe `MesuraOrdenacio` per a estudiar el comportament del mètode `mergeSort(int[], int, int)` definit en la classe `MesuraOrdenacio`. En aquest exercici, et suggerim que facis servir talles que siguin potències de 2. Per a això, pots definir en la classe `MesuraOrdenacio` les constants `INITALLA_MERGE` i `MAXTALLA_MERGE` per representar respectivament el valor de la talla més menuda a considerar (2^{10}) i el valor més gran (2^{19}), i que l'increment de talla siga: `talla *= 2`. El nombre de crides al mètode d'ordenació per obtenir valors significatius, pot ser el mateix que en l'anàlisi dels altres dos algorismes d'ordenació (200).

7 Avaluació

Aquesta pràctica forma part del primer bloc de pràctiques de l'assignatura que serà avaluat en el primer parcial d'aquesta. El valor d'eixe bloc és d'un 40% respecte al total de les pràctiques. El valor percentual de les pràctiques en l'assignatura és d'un 20% de la nota final.