Innlevering 3 IN2010 H24

# Kattungeoppgave

Kattunge-oppgaven er gjort og vedlagt som en java-fil. Den heter Kattunge.java

### **Sortering og Korrekthet**

Sorteringsalgoritmene som er implementert er BubbleSort, HeapSort, InsertionSort og MergeSort. Disse fire algoritmene er vedlagt som Python-filer. Jeg innså at det ville bli mer kode og litt mer komplisert hvis jeg skulle skrive algoritmene i java, og i tillegg så tenker jeg at det er enklere å implementere koden i python med tanke på syntaks og forståelse av algoritmene. Algoritmene kjører og utfører operasjonene riktig og hver av algoritmene får ut riktig og samme resultat fra outputfilen. Siden det ikke er gitt hvilke og hvor mange inputfiler vi skal teste for hver av disse algoritmene, så har jeg valgt å teste den mot example.txt. Etter å ha kjørt input på programmet, der hovedprogrammet kaller på alle algoritmene og skriver ut outputfilene, med riktig endelse.

Etter å ha kjørt programmet så har alle fått riktig output, som vil si at programmet har sortert tallene i riktig rekkefølge med den minste verdien på toppen/begynnelsen. Outputfilene har fått like svar.

Input fil:

|   | 80  |  |
|---|-----|--|
|   | 91  |  |
| 3 | 7   |  |
|   | 33  |  |
|   | 50  |  |
|   | 70  |  |
|   | 13  |  |
|   | 321 |  |
|   | 12  |  |
|   |     |  |

Output fil:

| atput iit. |   |     |  |
|------------|---|-----|--|
|            | 1 | 7   |  |
|            |   | 12  |  |
|            |   | 13  |  |
|            |   | 33  |  |
|            |   | 50  |  |
|            | 6 | 70  |  |
|            |   | 80  |  |
|            |   | 91  |  |
|            | 9 | 321 |  |
| _          |   |     |  |

Programmet ble også testet for større input-filer. Denne gangen testet vi med blant annet random\_100.txt. Her ble outputfilene resultatet og outputfilene helt identiske i tillegg til at de ble riktig sortert. Her var inputfilen helt random dvs. at algoritmene har litt mer arbeid å gjøre enn å velge den vanlige \_\_100.txt filen, som ikke var helt blandet sammen. Det blir litt mange tall å legge til en screenshot av input-/output filen, så derfor har jeg lagt til alle inputfilene og outputfilene jeg har testet med i programmet inne i en egen mappe som heter input. Her vil det ligge test-filen(input) sammen med alle .out-filene fra algoritmene.

# Sammenligninger, bytter og tid

Byttene som sorteringsalgoritmene ble målt med var ved hjelp av swaps fra filen «countswaps.py», mens sammenligninger «cmp» er målt i «countcompares.py» Ved hjelp av filene, blir funksjonene kalt når sorteringsalgoritmene kjøres noe som lar dem telle mens programmet kjører. Alle inputfilene som har blitt testet har lagt til en csv. -fil i input-mappen sammen med .out filene til hver sorteringsalgoritme.

## Eksperimenter

Nederst ligger diagrammer fra de genererte csv-filene

# 1. I hvilken grad stemmer kjøretiden overens med kjøretidsanalysene (store O) for de ulike algoritmene?

Ifølge grafene og store O sin representasjon, stemmer det overens med kjøretidskompleksiteten. De ulike sorteringsalgoritmene har disse kjøretidskompleksitetene (Store O):

BubbleSort =  $O(n^2)$  –

Kjøretiden og antall sammenligninger øker raskt når mengden av input øker. Den sammenligner hvert element med naboen og bytter dem hvis de er i feil rekkefølge. Algoritmen er derfor ineffektiv for store datasett. Ut ifra diagrammene ser vi at BubbleSort tar mye mer tid for større n, som stemmer med  $O(n^2)$ 

HeapSort = O(n log n) -

Dette er en algoritmen som bygger på maks-heap, som vi tidligere har sett, blitt sammenlignet med binære søketrær, men i «array-form». Her gjentar vi å bytte toppverdien med det siste elementet og omstrukturere heapen. Hjelpefunksjonen tar O (log n), og siden det utføres operasjoner på hvert element i listen på n elementer får vi totalt O (n log n). Her øker kjøretiden effektivt og balansert, uansett når inputstørrelsen øker, som betyr at den er litt bedre enn BubbleSort.

MergeSort = O(n log n) -

Her deler vi inputlisten rekursivt i mindre lister hvor vi sorterer dem og slår de sammen igjen. Siden listen deles i to i hvert steg har vi tiden O (n log n). På samme måte som HeapSort er denne veldig effektiv på store input og rask til å gjennomføre sorteringene. Selv om inputen øker, vil kjøretiden vokse sakte i forhold til logaritmiske økningen.

InsertionSort =  $O(n^2)$  –

Denne må ta imot en liste og sortere den ved å flytte elementer til riktig posisjon for ett og ett element om gangen. Denne passer godt for små datasett og nesten sorterte liste, som vi ser fra diagrammene, i tillegg til den lineære veksten i kjøretid. Denne har også mindre sammenligninger enn HeapSort, når vi har en nesten sortert liste.

# 2. Hvordan er antall sammenligninger og antall bytter korrelert med kjøretiden?

Algoritmene HeapSort og MergeSort har mer avanserte sammenligningsteknikker, som vil si at de bruker mindre bytter og sammenligninger som igjen gjør det til raskere kjøretid. De kan uansett få mange sammenligninger, men kjøretiden til å finne plasseringen for elementet er raskere enn f.eks. BubbleSort og insertion.

De andre to: BubbleSort og InsertionSort har ekstremt mange sammenligninger og bytter for randomized lister, som fører til lengre kjøretid, eller stopp på algoritmen siden de har brukt mer tid enn ønsket. Vi kan se på diagrammet til random 100 at disse to algoritmene stoppet og fikk en høy «spike» når de nådde over halvparten av elementene i listen, som mest sannsynlig betyr at de fikk et element som trengte å søkes gjennom hele listen som tok ekstremt my tid.

# 3. Hvilke sorteringsalgoritmer utmerker seg positivt når n er veldig liten? Og når n er veldig stor?

Når n er liten er InsertionSort det beste valget, siden tiden er lineært og i disse tilfellene er det få elementer, som utmerker seg når elementene er nesten sortert. Den bruker også lite minne i henhold til MergeSort, og har en enkel implementasjon. Grunnen til at det ikke er HeapSort er fordi at den trenger med vedlikehold og høyrer overhead som krever at vi bygger den etter hvert element. For små elementer vil den bruke unødvendig mer tid enn den trenger. På samme måte er BubbleSort også ikke helt ideelt, siden den gjør flere sammenligninger enn insertion (se diagram under for example.txt testen).

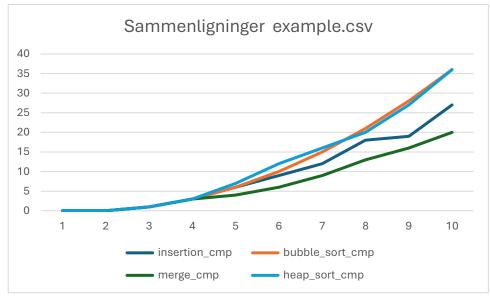
Når n er veldig stor passer MergeSort og HeapSort veldig god på grunn av kjøretidskompleksiteten, som gjør den effektiv mot store datasett. HeapSort bruker den vanlige heap-strukturen til å oppnå effektiv sortering, og den har en fordel om å bruke lite minne i forhold til MergeSort, som krever ekstra minne for å slå sammen listene. På samme måte er MergeSort veldig effektiv siden den deler listene i mindre deler som blir utført rekursivt, som sorterer i en systematisk måte.

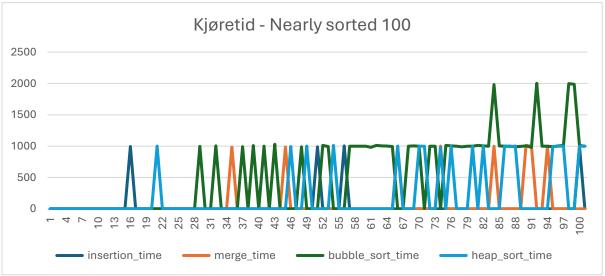
# 4. Hvilke sorteringsalgoritmer utmerker seg positivt for de ulike inputfilene? Når vi testet nearly sorted lister: av alle sorteringsalgoritmer var det InsertionSort som fungerte veldig godt blant disse. Den var raskest og sorterte listene med kortest tid. Algoritmen trenger bare å gjøre noen få bytter, og når elementene er nesten sortert, blir antallet cmp og bytter mye mindre.

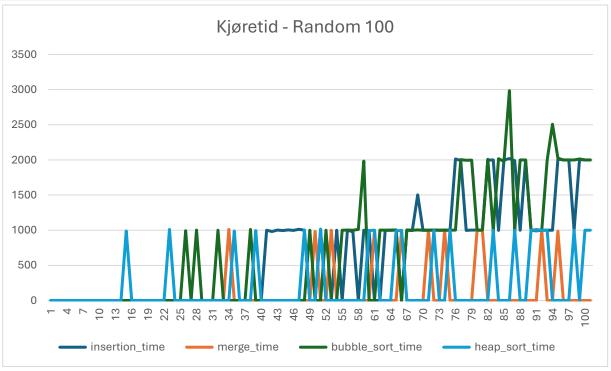
På random data, det vil si elementer som er helt blandet i listen, så var det MergeSort og HeapSort som dominerte hele kjøretiden. De to andre algoritmene ble stoppet siden de brukte ekstra mye tid, og var mye mer tidskrevende i sorteringsfasen. Siden MergeSort alltid deler listene i mindre lister så blir den ikke påvirket av graden av sortering eller usortering i listene. Vi kan også nesten si det samme for HeapSort. Her er den også veldig rask men ikke like rask som MergeSort i praksis. Hvis vi ser på random 10000-diagrammet ser vi at HeapSort stopper opp. Dette var veldig overraskende, siden jeg trodde den også ville sortere ferdig hele listen. En av grunnene kan være at algoritmen har nådd den maksimale minnet som programmet kan bruke eller at den har brukt mer tid enn det som maksimalt er avsatt.

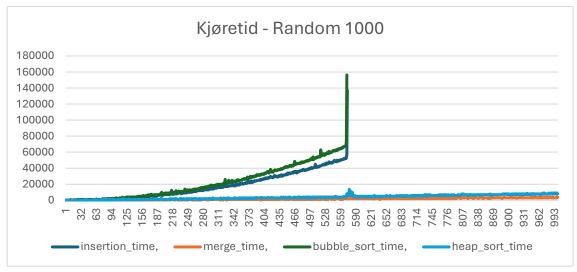
# 5. Har du noen overraskende funn å rapportere?

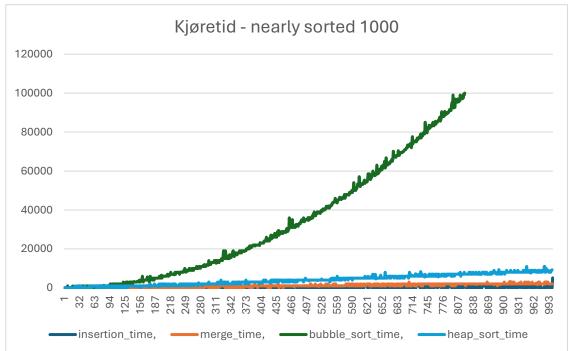
Jeg ble litt overasket over hvor godt BubbleSort fungerte. Jeg trodde den ville stoppe helt og ikke kjøre når jeg testet den med nearly\_sorted 1000, men den kom et stykke fram. Grunnen til at jeg tenkte det slik var at kjøretiden allerede var høy når den ble testet mot nearly og random 100, og antall sammenligninger var også veldig høy og den hadde høyest kjøretid i y-aksen enn de andre.

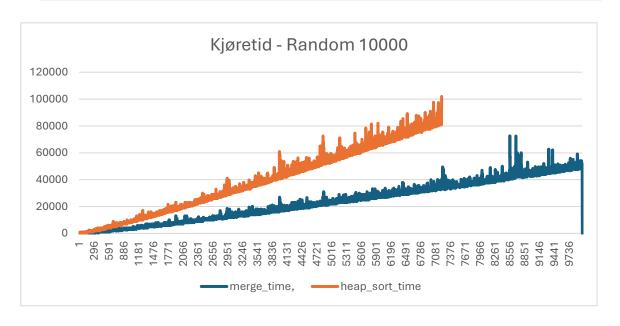


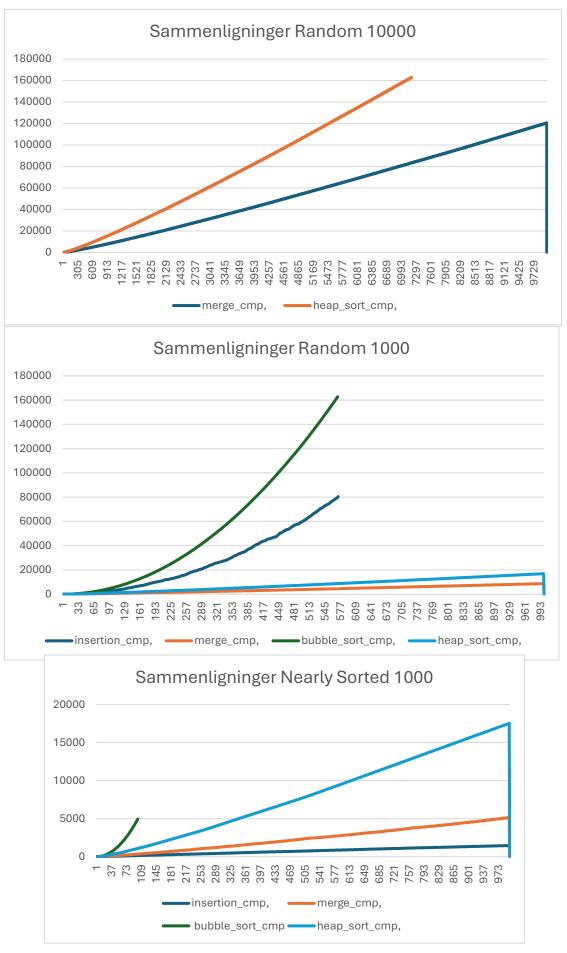












Side 7 av 7