

Miniprojekt: Programmering i R

January 16, 2015

Som en del av kursen i R programmering ska en analys av befintliga data göras med Rmarkdown genomföras. Miniprojektet är uppdelat i två delar. Den första delen handlar om att läsa in data från externa datakällor och beskriva dessa data.

I den andra delen av miniprojektet ska mer utförlig analys genomföras samt bearbeta och analysera denna data vidare.

För båda delarna gäller att:

- R-markdown ska användas.
- Rapporterna ska lämnas in som både **PDF** och **.Rmd**-fil.
- Samtliga material ska laddas in i R från webben som **externa datakällor**. Vill ni använda ett eget material får ni lägga upp det öppet på github, dropbox, google docs eller dylikt och läsa in det därifrån i R. Syftet är att rapporten ska vara helt reproducerbar och kunna skapas på godtycklig dator.
- **Rmd**-filen ska kunna köras och reproducera era resultat. D.v.s. den ska innehålla all er kod som behövs för analysen.
- **Namn**, **liu-id** och **gruppnummer** ska framgå i början av rapporten.

1 Deskriptiva data

Den första delen av miniprojektet är att samla in datamaterial och beskriva materialet kortfattat i en första del av rapporten.

Till den miniprojektet behöver ni **tre** datamaterial, två som innehåller kommunala data och ett material som innehåller en tidsserie.

Tänk på att välja material ni själva tycker är intressant!

Kommunala data Två datamaterial (data.frames) ska innehålla data på kommunnivå (d.s.v. för alla 290 kommuner). Ett exempel skulle kunna vara antal invånare i varje kommun. Dataseten ska ha minst **2 variabler** utöver kommunnamn. Ni väljer själv vilka variabler som ska ingå och vilken område data ska komma ifrån. Tanken är att i nästa del göra en enklare sambandsanalys mellan dessa variabler i nästa del av miniprojektet.

Tidsserie Hitta ett dataset som innehåller en **tidsserie**, det innebär att det finns en variabel som har observerats över tiden. Kravet är att data ska innehålla data på **månadsnivå** och innehålla data från minst 3 år (36 månader). Här ska ni alltså hitta en variabel som observerats under minst 36 tidpunkter, men fler går bra. Data ska alltså innehålla två kolumner, en med variabeln som vi är intresserade av och en med tidpunkterna.

1.1 Inlämning av del I

Den första inlämningsuppgiften handlar om att läsa in i R och beskriva de material ni valt med R-markdown. Ni ska beskriva era material i text samt sammanfatta de variabler ni valt med de beskrivande statistiska mått som ni själva finner lämpliga. Ta fram beskrivande statistik för **alla** variabler i data. Beroende på hur data ser ut så kan det vara medelvärden, frekvenstabeller mm. Ni kan göra relevanta transformationer av era variabler om ni vill, tex göra en numeriska variabel till en binär och räkna med andelar eller dela in kommunerna i stora, medelstora och små när det gäller befolkning.

Följande saker ska ni göra med data med basgrafiken i R:

1. Ni ska minst ha ett histogram eller barplot per variabel i kommunamaterialet
2. En tidsseriegraf/linjediagram för tidseriematerialet

Lämna in rapporten både som en fullt reproducerbar **Rmd**-fil och som **PDF** i LISAM.

- I denna del ska samtliga grafer vara skapade med basgrafiken i R.
- Tabeller ska vara "riktiga" tabeller (med ex. `kable()`), inte utskrifter av R-kod.

2 Analys

I den första delen av minprojektet har ni valt ut och beskrivit tre datamaterial. Nu ska vi fortsätta detta arbete med analyser av materialen. Ni som grupp kommer att ha en del frihet i hur ni utför datanalysen som beskrivs nedan. Det ni ska göra är att bearbeta data vidare, några enkla analyser, lite olika grafer i `ggplot2` och en linjär regression.

2.1 Inlämning del II

Den fulla rapporten ska lämnas in som en fullt reproducerbar **Rmd**-fil och som ett **PDF**-dokument i LISAM. Nedan framgår exakt vilka analyser som ska genomföras.

- I denna del ska samtliga grafer vara skapade med `ggplot2`
- Tabeller ska vara "riktiga" tabeller (med ex. `kable()`), inte utskrifter i R-kod.

2.1.1 Dataanalys av kommundata

Slå samman de två dataseten med kommundata så det blir ett dataset som innehåller variablerna från båda dataseten. Om ni gör rätt här så ska ni få ett dataset med en variabel över kommun och minst 4 andra variabler. Detta kan göras på olika sätt, ett är att använda funktionen `merge()`. Här finns en video för hur ni kan använda `merge()`.

Följande saker ska ni göra med data:

1. Ta fram minst en scatterplot mellan två variabler.
2. Ta fram minst en scatterplot/histogram/stapeldiagram som är grupperat på en annan variabel i minst två grupper.
3. Göra minst ett hypotestest, där ni ställer upp en nollhypotes och sen testar om ni kan förkasta den. Beroende på hur er data ser ut så kan det vara ett t-test, test av andelar eller ett chitvå-test/fishers test. I rapporten ska ni skriva upp både nollhypotesen och mothypotesen, ange även p-värdet.
4. Om data är numerisk (ex. frekvenser) så ska ni beräkna korrelationer mellan alla sådana variabler.
5. Beräkna två olika linjära regressionsmodeller. Ni väljer en responsvariabel (y -variabel) ni tror kan bero/ha ett samband med de andra variablerna. Välj sedan ut två olika variabler som ni kan ha som förklarande variabler (x -variabler). Skatta sedan två regressionsmodeller med de två olika förklarande variablerna i var sin regressionsmodell, dvs om era förklarande variabler heter X_1 och X_2 , så blir modellerna: $y = \beta_0 + \beta_1 * X_1$ och $y = \beta_0 + \beta_1 * X_2$. Det är okej att använda frekvensdata i regressionsmodellerna. Följande saker ska vara med från modellerna:

- (a) Scatterplot mellan förklarande variabel och responsvariabel, tillsammans med skattad regressionslinje.
 - (b) Histogram över residualerna (felen)
 - (c) Hurvida lutningen på regressionslinjen är signifikant skild från noll, dvs om p-värdet för β_1 i de olika modellerna är mindre än 0.05.
6. Ni ska ha med minst en Sverigekarta där ni plottar någon beskrivande statistik för en variabel över kommunerna. Ni behöver ladda ner "Kommun_SCB.zip" och packa upp den för att få tillgång till en karta på kommunnivå.

2.1.2 Dataanalys av tidseriedata

Låt X vara en variabel i tidsseriematerialet. Utför nu följande:

1. Gör en linjeplot mellan X och en tidsvariabel. Skalan på x-axeln ska vara en lämplig tidsskala.
2. Beräkna månadsmedelvärden, spara dessa i `monthMeans`.
3. Använd funktionen `summary()` för att fram beskrivande statistik för varje år (det ska vara minst tre år i data)
4. Subtrahera månadsmedelvärden från X , så ni tar bort säsongsvariationen i data. Månadsmedelvärdet för januari ska subtraheras från alla januarivärden i data, och likadant för de andra månaderna. Spara den nya tidserien som `newX`.
5. Gör en linjeplot mellan `newX` och tid.
6. Skatta en linjär regressionsmodell mellan `newX` och dess tidsindex. Rapportera följande:
 - (a) Scatterplot (punkter) mellan förklarande variabel och responsvariabel, tillsammans med skattad regressionslinje.
 - (b) Histogram över residualerna (felen)
 - (c) Hurvida lutningen på regressionslinjen är signifikant skild från noll, dvs om p-värdet för β_1 i de olika modellerna är mindre än 0.05.
7. Verkar det finnas någon trend i data? Dvs ökar/minskar data med tiden, eller är data konstant över tid.