



LUNDS UNIVERSITET
Lunds Tekniska Högskola

EDAA35: R-programmering

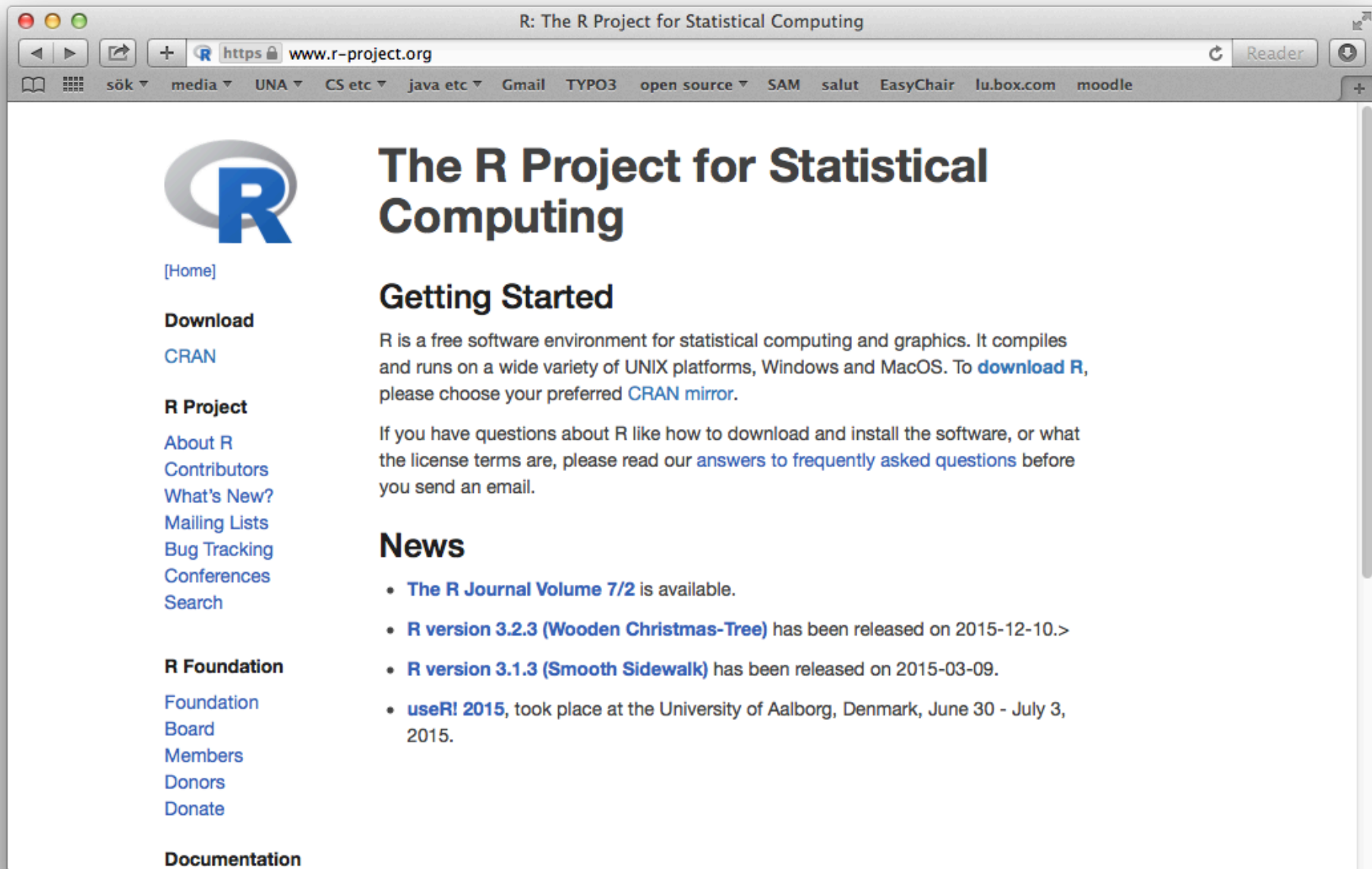
MARTIN HÖST



Idag

- Om R
- De vanligaste datatyperna
- Urval av data
- Analys av data
- Presentation av resultat: plot, utskrift, etc

www.r-project.org



The screenshot shows a browser window titled "R: The R Project for Statistical Computing" with the URL "https://www.r-project.org". The browser's address bar and tabs are visible at the top. The website content includes the R logo, a navigation menu on the left, and a main content area with sections for "Getting Started" and "News".

The R Project for Statistical Computing

[Home]

Download

[CRAN](#)

R Project

[About R](#)
[Contributors](#)
[What's New?](#)
[Mailing Lists](#)
[Bug Tracking](#)
[Conferences](#)
[Search](#)

R Foundation

[Foundation](#)
[Board](#)
[Members](#)
[Donors](#)
[Donate](#)

Documentation

Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To [download R](#), please choose your preferred [CRAN mirror](#).

If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

News

- [The R Journal Volume 7/2](#) is available.
- [R version 3.2.3 \(Wooden Christmas-Tree\)](#) has been released on 2015-12-10.>
- [R version 3.1.3 \(Smooth Sidewalk\)](#) has been released on 2015-03-09.
- [useR! 2015](#), took place at the University of Aalborg, Denmark, June 30 - July 3, 2015.

Vad är R?

- En variant av S
 - Högnivåspråk för dataanalys och grafik utvecklat av AT&T Bell Laboratories. Fick "ACM Software System Award" 1998
- Open source version (GNU GPL), mycket väl använd
 - "free"
 - Mycket hjälp finns t ex på nätet
- "Statistikverktyg" och programspråk

Typisk användning

- Inläsning av data
- Sortering, urval av data, etc
- Statistisk analys
 - Deskriptiv statistik
 - Mer avancerade funktioner
- Presentation av resultat
 - Skärm, fil, grafik,...
 - Spara data
- Många paket för olika funktioner

Library Highlights (Pandas)

- A fast and efficient **DataFrame** object for data manipulation with integrated indexing;
- Tools for **reading and writing data** between in-memory data structures and different formats: CSV and text files, Microsoft Excel, SQL databases, and the fast HDF5 format;
- Intelligent **data alignment** and integrated handling of **missing data**: gain automatic label-based alignment in computations and easily manipulate messy data into an orderly form;
- Flexible **reshaping** and pivoting of data sets;
- Intelligent label-based **slicing, fancy indexing, and subsetting** of large data sets;
- Columns can be inserted and deleted from data structures for **size mutability**;
- Aggregating or transforming data with a powerful **group by** engine allowing split-apply-combine operations on data sets;
- High performance **merging and joining** of data sets;
- **Hierarchical axis indexing** provides an intuitive way of working with high-dimensional data in a lower-dimensional data structure;
- **Time series**-functionality: date range generation and frequency conversion, moving window statistics, moving window linear regressions, date shifting and lagging. Even create domain-specific time offsets and join time series without losing data;
- Highly **optimized for performance**, with critical code paths written in [Cython](#) or C.
- Python with *pandas* is in use in a wide variety of **academic and commercial** domains, including Finance, Neuroscience, Economics, Statistics, Advertising, Web Analytics, and more.

Programspråket

- Dynamiskt typat
- Script-språk, körs från terminalen utan kompilering

Lab 1

- Görs som quiz i moodle
- Man måste klara alla frågorna
- Jobba i par på en inloggning i moodle. Godkännande inte i moodle (i "SAM")
- Ibland möjligt att "gissa" svar... Gör inte det, ni har nytta av att förstå på nästa lab...
- "enrolment key": EDAA3520
 - OBS! Välj rätt kursinstans
(Utvärdering av programvarusystem 2020)

Några alternativ till R

- Matlab – används av många inom elektro, signalbehandling, etc. Mycket bra på t ex matriser
- Python – generellt språk. Med "numpy" och "pandas" får man ungefär samma funktioner som i R (t ex "dataramar")
- ...men har man förstått hur man jobbar i ett av språken kan man enkelt lära sig ett annat

-
- Visa i terminalen

Operationer

<code>:: :::</code>	åtkomst av variabler i 'namespace' *
<code>\$ @</code>	komponent, @*
<code>[[[</code>	indexering
<code>^</code>	exponent
<code>- +</code>	unär plus och minus
<code>:</code>	sekvens
<code>%any%</code>	special-operator *
<code>* /</code>	multiplikation, division
<code>+ -</code>	binär addition och subtraktion
<code>< > <= >= == !=</code>	jämförelser
<code>!</code>	negation
<code>& &&</code>	och
<code> </code>	eller
<code>~</code>	för att beskriva modeller *
<code>-> ->></code>	tilldelning, tilldelning i 'enclosing environment'*
<code><- <<-</code>	tilldelning, tilldelning i 'enclosing environment'*
<code>=</code>	tilldelning (motsvarar '<-')
<code>?</code>	hjälp

Några funktioner för beräkningar

<code>sum(..., na.rm = FALSE)</code>	summa
<code>prod(..., na.rm = FALSE)</code>	produkt
<code>cumsum(x), cumprod(x)</code>	kumulativ summa, produkt
<code>mean(x, trim = 0, na.rm = FALSE, ...)</code>	medelvärde
<code>median(x, na.rm = FALSE)</code>	median
<code>var(x, na.rm = FALSE)</code>	varians
<code>sd(x, na.rm = FALSE)</code>	standardavvikelse
<code>exp(x)</code>	e^x
<code>log(x, base = exp(1)), log10(x), log2(x)</code>	logaritmer
<code>max(..., na.rm = FALSE)</code>	max
<code>min(..., na.rm = FALSE)</code>	minimum
<code>range(..., na.rm = FALSE)</code>	'range', dvs min och max
<code>which.max(x), which.min(x)</code>	vilket värde som är max, min
<code>sqrt(x), abs(x)</code>	\sqrt{x} , $ x $
<code>sin(x), cos(x), tan(x)</code>	<i>sin()</i> , <i>cos()</i> , <i>tan()</i> (radianer)
<code>round(x, digits = 0)</code>	avrundning (digits decimaler)
<code>floor(x), ceiling(x)</code>	avrundning nedåt, uppåt
<code>factorial(x), choose(n, k)</code>	$x!$, $n!/((n - k)!k!)$

Arbeta med dataobjekt

<code>nrow(a)</code>	antal rader i <code>a</code>
<code>ncol(a)</code>	antal kolumner i <code>a</code>
<code>rbind(a, b)</code>	sammanslagning där raderna i <code>b</code> kommer efter raderna i <code>a</code> , dvs en radvis sammanslagning.
<code>cbind(a, b)</code>	sammanslagning där kolumnerna i <code>b</code> kommer efter kolumnerna i <code>a</code> , dvs en kolumnvis sammanslagning.
<code>summary(a)</code>	sammanfattar innehållet i <code>a</code> , t ex medelvärde och median
<code>str(a)</code>	sammanfattar <code>a</code> :s struktur
<code>head(a)</code>	de första raderna i <code>a</code>
<code>tail(a)</code>	de sista raderna i <code>a</code>
<code>sort(a)</code>	sortera <code>a</code>
<code>order(a)</code>	index som sorterar vektor <code>a</code> , dvs <code>a[order(a)]</code> motsvarar <code>sort(a)</code>
<code>table(a)</code>	räknar hur många element i <code>a</code> det finns av varje värde
<code>unlist(x)</code>	omvandlar listan <code>x</code> till en vektor
<code>append(x, values, after = length(x))</code>	lägger in vektorn <code>values</code> i vektorn <code>x</code> efter position <code>after</code> (om <code>after==0</code> så hamnar <code>values</code> först)

Läsa data från fil

- `read.csv(file, header = TRUE, sep = ",", quote = "\"", dec = ".", fill = TRUE, comment.char = "", ...)`
 - För komma-separerade data-filer
- `read.table()`
 - Mer generell

+ många fler paket och funktioner för att importera data

Instruktioner i script

- Körs med
`> source("filnamn")`

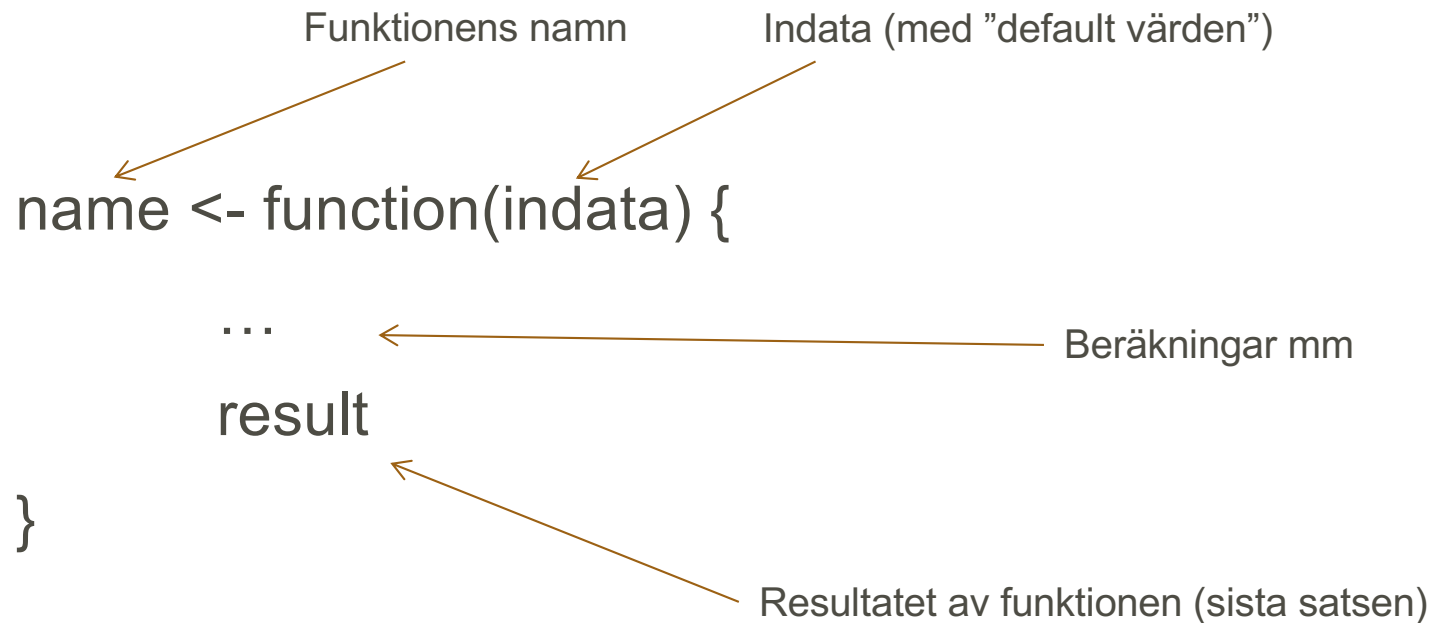
- T ex:

```
> source("script.R")
```

```
[1] "hej"
```

```
a <- "hej"  
print(a)
```

Funktioner, definition



(Kan vara anonyma)

-
- Visa i terminalen