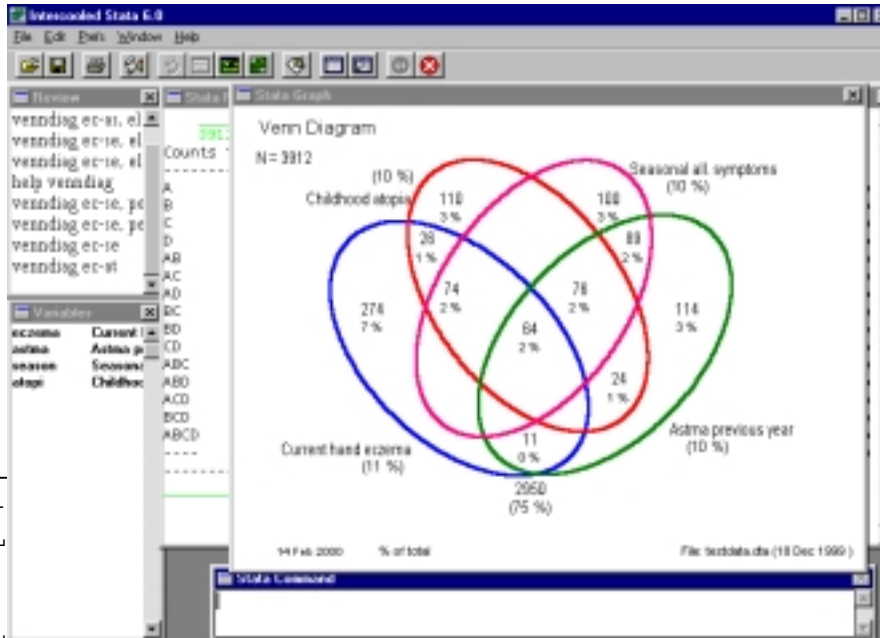


Stata-Introduktion



```

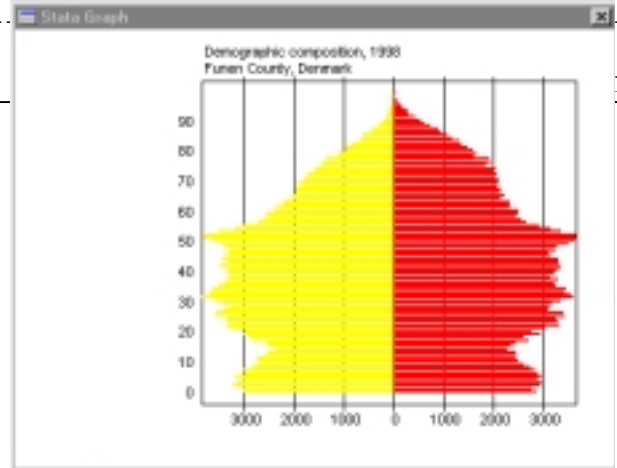
of obs = 240
chi2(1) = 14.89
chi2 = 0.0001
R2 = 0.0453
  
```

test	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
sr5a	.3386243	.0980107	-3.741	0.000	.1920208 .5971564

```

-> xi: logistic test sr5b if sportn == 1 , cluster(hold)
Logit estimates
Log likelihood = -126.28648
Number of obs = 192
Wald chi2(1) = 4.27
Prob > chi2 = 0.0387
Pseudo R2 = 0.0420
(standard errors adjusted for clustering on hold)
  
```

test	Odds Ratio	Robust Std. Err.
------	------------	------------------



Jens M. Lauritsen ©

- Forslag til arbejdsmåde
- Dokumentation
- Oparbejdning og afledte variable
- De mest anvendte analysefunktioner
- Nyttige hjælpeprogrammer
- Stata Technical Bulletin
- Tips til opsætning

- 3. udgave 2000 (31 marts 2000)

Stata - Introduktion .

© Jens M. Lauritsen

3. udgave, ultimo marts 2000 (kontrollér dato på forsiden)

ISBN-87-987843-0-7 (trykt udgave)

ISBN-87-987843-1-5 (Elektronisk distribueret udgave)

Den trykte udgave distribueres via Studenterbogladsen ved Syddansk Universitet.

E- Mail: studenter@boghandel.sdu.dk.

Materialet inklusive øvelsesfiler må ikke kopieres, udgives eller anvendes til undervisning uden forudgående aftale med forfatteren. Materialet kan anvendes til personligt brug.

Der er udarbejdet et sæt øvelsesfiler til denne note. Kan downloades fra <http://www.bola.suite.dk> Se også side 13.

Forbehold: Da det er et meget stort arbejde at udarbejde noter, holde styr på figurer og undersøge om øvelserne fungerer må læserne bære over med uoverensstemmelser mellem de præsenterede skærbilleder og de billeder der kommer frem ved at udføre øvelserne. Fejl og mangler vil jeg dog gerne have oplysninger om. Den elektroniske form tilpasses efterhånden, mens der kan være – forhåbentlig mindre – uoverensstemmelser med indholdet heri.

Kommentarer til noten sendes til: JM.Lauritsen@dadlnet.dk.

Angiv venligst hvilken udgave af noten, som kommentaren knytter sig til (Se forsiden).

Nogle nyttige internet sider om Stata, biostatistik eller hjælpeprogrammer:

Statistik rutiner og uddybende forklaringer af forskellige analysetyper i Stata:

<http://www.oac.ucla.edu/training/stata/> og <http://www.stata.com/links/resources1.html>

Instruktionsrutiner og forklaring af analyse af forskellige datatyper, inklusive øvelsesdata. Er baseret på EpiInfo, men principperne er generelt gyldige uanset statistikprogram:

<http://www.sjsu.edu/faculty/gerstman/EpiInfo>

Enkle manualer

http://mkn.co.uk/help/extra/people/Brixton_Books

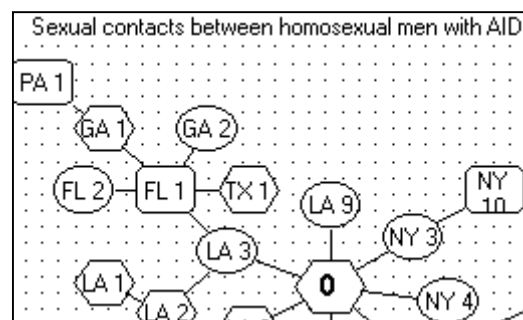
Gratis software til særlige formål:

EpiData til indtastning af data på en nem og pålidelige måde win95/98/NT/2000. Send en mail til adressen <mailto:epidata@dadlnet.dk> for nærmere oplysning. Programmet skriver datafiler, der direkte kan analyseres i Stata. (+Dbase, Excel og EpiInfo)

Epicalc 2000 (Simple beregninger-nærmest lommeregner) og et simpelt program til at tegne flowcharts og diagrammer (**EpiGram**):



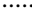
<http://www.myatt.demon.co.uk/>

Uddrag af diagram tegnet med EpiGram. Ved endelige brug af diagrammet kopieres prikkerne ikke.



Stata - introduktion.....	7
Forord.....	7
Hvorfor Stata ?	7
Noten er skrevet ud fra følgende forudsætninger:	8
Hjælpeprogrammer.....	8
Manualer, Stata moduler og øvrige indkøb.	8
Yderligere kurser om Stata.....	8
Krav til dokumentation af data og forskning.	9
A. Mulighed for at finde tilbage til originalmateriale.....	9
B. Kvalitetskontrol, dokumentation og arkivering.....	9
Indsamling, oparbejdning og analyse af data.....	9
Øvelsesdata – og hjælperutiner.....	11
Datasæt 1 – Bromaraton.	11
Datasæt 2 - Euroqol.	11
Datasæt 3 – Befolkningen i Fyns Amt pr 1.1.1998.....	12
Datasæt 4 – Simulerede data over forekomst af astma, høfeber og eksem.	12
Datasæt 5 – Data over varighed for brug af hoftebeskytter hos 42 ældre.....	12
Hent øvelsesfiler og læg dem ind på PC'en.....	13
Øvelse 1 Opret en mappe med navnet c:\data (eller d:\data)	13
Øvelse 2 Hent øvelsesfilerne på internet	13
Øvelse 3 Pak filerne ud på PC'en.....	13
Øvelse 4 Installér Stata.....	13
Øvelse 5-s Opsæt PC'en på en god måde	13
Første start af Stata (Delvis med QUEST hjælpe/menusystemet).	14
Øvelse 6 Start Stata	14
Øvelse 7 Læs en datafil	14
Øvelse 8 Aftegn en graf for de to variable "alder" og "dectime"	15
Øvelse 9 Start menusystemet quest.....	15
Øvelse 10 Regnemaskine.....	16
Stata - hjælp.....	17
Øvelse 11 Opslag i hjælpesystemet - 1	17
Øvelse 12-s Opslag i hjælpesystemet - 2.....	17
Øvelse 13-s Opslag i hjælpesystemet - 3.....	17
Øvelse 14 Opsætning af Stata skærm.	18
Øvelse 15 Gem kommandoer i en "do" fil	18
Anden start af Stata – kopi af resultater.....	19
Øvelse 16 Start en logfil ud fra menu systemet.	19
Øvelse 17 Frekvenstabel af køn og graf af løbstid og alder.....	19
Øvelse 18 Gem kommandoer i en "do" fil	20
Indhold i "do" filer : Opbygning af kommandoer i Stata.....	20
Øvelse 19 Kopiér grafen fra Stata til tekstbehandling.	21
Øvelse 20 Kopi'er en tabel direkte fra Stata til tekstbehandling.	21
Øvelse 21 Print logfil	22
Øvelse 22 Afslut logfil	22
Øvelse 23 Overfør log fil og graf til tekstbehandling	22
Øvelse 24 Installering af de rutiner der indgår i kursus materialet.	23
Øvelse 25-s Opslag På internet direkte fra Stata 6	23
Tredie start af Stata – rettelse i "do" filer	24
Øvelse 26 Gentag analysen ud fra filen <i>first.do</i>	24
Øvelse 27 Kontrollér indholdet af <i>first.do</i> i "Do-file" editor.....	24
Øvelse 28 Udfør kommandoer direkte fra "Do-file" editor.....	24
Øvelse 29 Bestem selv over logfiler.....	25

Øvelse 30-s gentag analysen og få kopi til logfil automatisk – logfil overskrives.....	26
Øvelse 31-s Se på indholdet af first.log og sammenlign med den tidligere logfil.	26
Øvelse 32-s Gentag analysen og få tilføjet til tidligere logfil.....	26
Fjerde start af Stata – labels og beregning af nye variable.....	27
Øvelse 33 Beregn alder og aldersgrupper	27
Øvelse 34 Omsæt tiden fra kryptisk tid til decimaltid.	27
Øvelse 35 Afprøv den udregnede decimaltid.	27
Øvelse 36 Materiale beskrivelse for bromaraton.....	27
Øvelse 37 Association alder-løbstid.....	27
Øvelse 38 Strategi 1 : Labels ud fra menusystem.....	28
Øvelse 39 Strategi 2. Tildel labels med kommandoer.....	29
Øvelse 40 Alle labels og kodning af varighed for brug af hoftebeskytter.	29
Øvelse 41 Materiale beskrivelse for hoftebeskytter.	29
Andre nyttige muligheder med STATA.....	30
Øvelse 42 Tutorials – selvstudieprogrammer	30
Øvelse 43 Befolkningspyramider	30
Øvelse 45 Venn Diagram	31
Øvelse 46 Stata som regnemaskine. Statistik ved indtastning	31
Øvelse 47-s Meta analyser	33
Tablet fra Stata til tekstbehandling.....	34
Graf fra Stata til tekstbehandling	35
Konvertering af data.....	37
Øvelse 48-s Oversæt data til Stata format.	37
Samlet oparbejdning af Bromaraton data med Stata	38
Øvelse 49 Kig på tallene.	38
Opgave 1: Hvor mange personer og omfang af uoplyst i de enkelte variable ?	39
Øvelse 50 Start Stata og gør klar til dagens opgave.....	40
Øvelse 51 Indlæs data fra originalfilen	40
Øvelse 52 Gem kommandoer i en ”do” fil	40
Øvelse 53-s Gem kommandoer i en ”do” fil (alternativ strategi)	40
Øvelse 54 Kontrollér indholdet af oparbejd.do i ”Do-file” editor.....	40
Strategi 2 ”Ekstern Editor” Alternativ for erfarne brugere.....	41
Øvelse 55-s Kontrollér indholdet af oparbejd.do med brug af Win-Commander.....	41
Øvelse 56-s Se på resultatet fra logfilen.....	41
Øvelse 57 Insufficente data for enkeltindivider ?.....	43
Beregn afledte variable.	44
Øvelse 58 Kontrol af variabelen pnr	44
Øvelse 59 Stata som regnemaskine. Afprøv kommandoen ”display”.....	46
Øvelse 60 Rækkefølge af linier med rekodning !.....	47
Øvelse 61 Beregn afstand og aftegn en frekvenstabel over grupperet afstand.	48
Øvelse 62 Beregn alder, aldersgrupper og decimaltid.....	49
Øvelse 63 Dokumentation, anvendte filer og beslutninger samles	51
Bortfaldsanalyse.....	51
Øvelse 64 Udfør en bortfaldsanalyse	51
Endelig besvarelse af hypoteser fra Bromaraton løbet.....	51
Øvelse 65 Belysning af de opstillede hypoteser.....	51
Nyhedsgruppe på internet for Stata og www.stata.com.....	52
Øvelse 66 Opslag på Stata’s internetadresse (www.Stata.com).	52
Tilmelding til nyhedsgruppen.....	52
Installering af Stata, Stat/Transfer og Quest	53
Øvelse 67 Installér Stata og Stat/Transfer.....	53
Øvelse 68-s Speciel opstart af Stata (ændring af menuen).	54
Opdatering og tilføjelser til STATA.....	55

Øvelse 69 Opdatering og tilføjelser til Stata version 6 direkte fra internet.	55
Øvelse 70 Hent opdateringsfiler til Stata på internet.	56
Øvelse 71 Hent Quest systemet.	56
Øvelse 72 Installér den nye Wstata.bin, som er hentet fra internet.	56
Øvelse 73 Installér ny udgave af ADO systemet, som er hentet fra internet.	57
Øvelse 74 Installér <i>quest.zip</i> , som er hentet fra internet.	57
Øvelse 75 Hent STB udvidelser som zip filer og installér fra harddisken.	58
Tilføjelse af rutiner direkte fra internet.	58
Øvelse 76 Tilføjelse af rutiner direkte fra internet.	58
Eksempel på svar fra FAQ listen på www.Stata.com	59
What are some of the problems with stepwise regression?.....	59
Hjælpeprogrammer.....	60
Installér Hjælpeprogrammer	61
Øvelse 77-s Hent win-commander fra internettet:	61
Øvelse 78-s Installér Win-Commander	61
Øvelse 79-s Afprøv Win-Commander	62
Øvelse 80-s Opsætning af Win-Commander – 1 (grundlæggende).	62
Øvelse 81-s Hent editor filerne fra internet:	63
Øvelse 82-s Installering af "Programmers file editor".....	63
Øvelse 83-s Opsætning af Win-Commander – 2 (til brug af Pfe).....	63
Øvelse 84-s Afprøv Programmers File Editor - Pfe.....	63
Øvelse 85-s Opsætning Pfe – 1 (stata filtyper vises)	63
Øvelse 86-s Opsætning Pfe – 3 (Standard tilpasning)	64
Øvelse 87-s Opsætning Pfe – 2 (tastatur)	64
Øvelse 88-s Opsætning Pfe – 4 (Avanceret tilpasning)	64
Backup og kopi af data	65
Øvelse 89-s Pak en kopi af "do" og datafiler til et arkiv på harddisken	65
Øvelse 90-s Pak en kopi af "do" og datafiler til en diskette	65
Øvelse 91-s Kopiér fra zip filen tilbage til c:\data.....	66
Vejledning for udformning af forsøgsprotokoller og forsøgsrapporter, datadokumentation og opbevaring af data inden for sundhedsvidenskabelig basalforskning	67
Vejledning for udformning af undersøgelsesplaner, datadokumentation og opbevaring af data inden for klinisk og klinisk-epidemiologisk forskning	68
Euroqol 5-d	69
God opsætning af start for Stata og Stat/Transfer.....	71
Øvelse 92-s God opbygning af startmenuen (tast ).	71
Øvelse 93-s Opbygning af STATA mappe på startmenuen (tast ).	72
Øvelse 94-s Hjælpeprogrammer i startmenuen. ()	76
Hyppegt anvendte tastaturkombinationer i windows.	77
General Commands in Stata	78

Stata - introduktion.

Forord

Hensigten med denne note er, at introducere brugen af statistik programmet Stata version 6. Noten giver anvisninger til den konkrete måde at arbejde på og kan **ikke** erstatte håndbøger, manualer og lærebøger i statistik og forskningsmetode. Tidligere udgaver af noten har angivet flere alternative metoder til arbejdet, hvilket i et vist omfang har forvirret "ikke erfarne" edb/statistik brugere. Denne udgave søger at fokusere på de muligheder der findes ved at bruge Stata, som programmet er udviklet. Dermed bliver noten også mere uafhængig af styresystem. idet Stata fungerer stort set ens uanset om du arbejder ud fra Windows, Linux eller Macintosh.

Forslag til hjælpeprogrammer er samlet i et særligt afsnit, herunder hvordan backup og lignende kan foretages.

Hvor meget edb erfaring er nødvendigt for at bruge Stata effektivt? Ikke så meget. Det er nødvendigt at mestre grundlæggende tekstbehandling, at kunne kopiere filer, oprette mapper og være i stand til at hente filer fra internet. Desuden bør du være i stand til at arbejde med pakkede filer, især hvis din PC ikke har adgang til internet.

Filen kan pakkes på forskellige måder, fx såkaldt zip, arj eller uar format. Når en fil "pakkes" sker der det, at unødigt tom plads fjernes. Fx kan en halv tom side i et dokument jo ændres til kommandoen "hop til næste side". Når en sådan fil pakkes ud igen indsættes den halve tomme side blot igen. Størrelsen på pakkede filer varierer fra 10-90 % af originalfilen. Det er nyttigt at arbejde med pakkede filer, hvis man skal sende over internet fordi transmissionstiden så nedsættes, men det forudsætter selvfølgelig at modtager kan pakke ud igen.

Der findes ingen universel **rigtig** måde at arbejde på. Find din egen og sørg for at blive god til at arbejde på den måde. Det er vigtigt, at du arbejder **på en reproducerbar måde**. Minimér brugen af "mus", minimér brug af menu'er og maksimer reproducerbarhed og dokumentation af analyser og arbejdsformer. Du sparer dig selv for megen statisk muskelbelastning ved at bruge tastaturkombinationer hyppigt, se side 77 og 78

Ved nyindkøb af computer foreslås følgende prioritering: 1. En bedre kvalitet af maskinen med langsommere processor fremfor en hurtigere processor på en dårligere generel kvalitet (pris er et godt udgangspunkt for kvalitet, 3 års garanti også udtryk for kvalitet). 2. Mindst 64 Mb Ram. 3. Skærnkvalitet (køb en dyrere 17" skærm, der har flere punkter pr enhed). 4. Et bedre skærnkort der har en hurtigere "refresh rate". 5. Lavest muligt støjniveau. 6. Hurtigere og bedre harddisk, evt. 4 Gb. 6. Optimalt backup udstyr. 7. Bedre stol, bord og lampe. 8. Arbejdstid til at blive god til at bruge systemet, fx kurser. 9. Hurtigere maskine, Mere Ram end 64Mb. osv..

Hvorfor Stata ?

Principielt er det ligegyldigt hvilket statistikprogram der bruges, hvis blot det kan udføre de analyser der er brug for, på en måde der er til at finde ud af og at programmet kan købes til en rimelig pris. Det kan ikke betale sig at gøre en masse ud af hvorfor det ene program er bedre end det andet. Nogle af begrundelserne er udtryk for vane, smag og behag. For mig er der følgende gode grunde til at anvende Stata som statistikprogram fremfor andre (spss, sas, systat, ...).

- Det understøtter den allernyeste udvikling indenfor biostatistik. Det omhandler stort set alle analysetyper der er relevante i biomedicin. Dvs ét program er tilstrækkeligt.
- Det anvendes af en række gode biostatistikere og epidemiologer, som har skrevet uddybende rutiner på flere områder. (Bl.a. David Clayton, Michael Hills, Hosmer, Lemeshaw)¹ Det er veldokumenteret og der er særdeles fyldige og gode eksempler på de fleste biomedicinske analysetyper i manualerne. Nye rutiner dokumenteres og skrives af avancerede brugere og distribueres gratis. Der er muligt at foretage omfattende matrix beregninger. Det håndterer datoer på en effektiv måde.
- Det arbejder hurtigt og fylder lidt (Kan være på tre disketter). En komplet installation svarende til noten fylder kun 6-10 Mb afhængig af antal ekstra rutiner der anvendes.
- Det er særdeles velegnet til analyse og beskrivelse af follow-up data, som fx overlevelsesdata. Både "simple" data med kun ind- og udgangstid og data med gentagne outcome, fx gentagne sygdomsepisoder.
- De fleste regressionsmetoder har indbygget mulighed for at tilføje en klynge (cluster) effekt. Samt anvende "robust methods".
- Når man har købt en licens er denne "uendelig". Dvs. at en typisk bruger har programmet til ejendom. I modsætning til en årlig afgift. Prisen er overkommelig (< 5000 kr) og specielt hvis man arbejder i forskningsmiljø (< 1000 kr). For køb se www.stata.com. Evt. www.sdu.dk/dou eller www.metrika.se.

Der er selvfølgelig også nogle begrænsninger i Stata. Der kunne godt være flere eksakte test, fx i visse non-parametriske test.

¹ Clayton D, Hills M. Statistical Models in Epidemiology. Oxford, OUP, 1993. Analyser fra bogen er uddybet i: Clayton D, Hills M. Analysis of follow-up studies with Stata 5.0 STB Reprints Vol 7: 253-268. Clayton D, Hills M. Analysis of case-control and prevalence studies STB Reprints Vol 5: 227-233. Levy & Lemeshow Sampling of Populations: Methods and Applications. Wiley, 1999. Hosmer & Lemeshaw: "Applied Survival Analysis" 1998.

General Commands in Stata

clear	Clear data from memory
display	Display values from functions
do	Execute commands
exit	Exit Stata
help	Find help
log	Echo output into a file
save	Save data for future use
use	Read a copy of a Stata file with data into memory
now	Start do file with automatic log (part of <i>kursus.zip</i>)
webseek	Find material on internet

Data manipulation and Management

collapse	aggregate data into a table
count	count number of observations
describe	which data are in memory
destring	change string variable to numerical
drop	eliminate variables or observations
encode	create numeric variable from string
expand	duplicate observations
format	specify display format
generate	create new variable
infile	read fixed format data
input	add data via keyboard
label	descriptive text (values, vars., files)
list	list some variables
listjl	list id variable (also by group)
merge	combine files (check also mmerge)
move	change sequence of variables in file
mvdecode	recode to missing (e.g. from 9 to .)
mvencode	recode from “.” to value
notes	add documentation to data + variable
order	change sequence of variables
outfile	write data in fixed format
pattern	indicate patterns of missing
query	show set options
quietly	do a command without showing output
recode	recode numeric
rename	change name of variable
replace	change contents of variable
set	set general options
sort	sort observations by variables.

Other useful

venndiag	Venn diagrams and creation of coupled variables
hbar	Horizontal bar charts. E.g. Population Pyramids
cdf	Cumulative plot of continuous var.
outreg	Regression output to publication
qnorm	Q normal probability plot

Descriptive Statistics

list	list values of variables
graph	graph Y ₁ Y ₂ Y ₃ X, xlabel ylabel xscale t1(“”) b1(“”) title(“”) Saving(“”)
hist	histogram
summarize	display summary statistics
table	multiway tables
tabulate	one and two way frequency
etab	One variable against many
codebook	Simple standard tables

General Statistics

anova	Analysis of Variance
correlate	Correlation
oneway	One-way analysis of variance
ranksum	Wilcoxon ranksum test
tabulate	crosstables incl. tests of homogeneity, gamma, exact r*c test
table	Summary statistics in tables
ttest	Mean comparison for small samples

Cohort / survival analysis

stmh	Mantel-Haenzel rates
poisson	poisson regression
strate	table of rates
stset	define survival time rates
staalen	Aalen cumulative hazards
stcox	Cox regression (+stptest)
stsplot	expand data according to a lexis diagram
sts graph	Kaplan Meyer survival functions in graph
sts test	Logrank test

Case Control / Cross Sectional analysis

mhodds	Mantel Haenzel Odds Ratios
partgam	Partial Gamma Coefficient
tabodds	Table of Odds

Regression Models

clomit	conditional logistic regression
logit	logistic regression (output on log scale)
logistic	logistic regression (output on antilog scale)
lrtest	likelihood ratio test
poisson	poisson regression
predict	obtain prediction and indicators of model fit
regress	linear regression
Categorical variables- “dummy” indicator: xi: regress Y i.varname	
sort X /* stratify on X */ by X: command variables, options for var variables: command X ...	