

SOS3003
**Anvendt statistisk
dataanalyse i
samfunnsvitenskap**
Forelesingsnotat, vår 2003

Erling Berge
Institutt for sosiologi og statsvitenskap
NTNU

Forelesing VII

- Logistisk regresjon I
Hamilton Kap 7 s217-223
- Manglende data
Allison 2002 "Missing Data"

LOGIT REGRESJON eller LOGISTISK REGRESJON

- **Skal nyttast når avhengig variabel er på nominalnivå**
- Føreset at Y har verdiane 0 eller 1
- Modellen av den betinga forventninga til Y , $E[Y | X]$, nyttar den logistiske funksjonen
- Men
Kvifor kan ikkje $E[Y | X]$ vere ein lineær funksjon også her?

Den lineære sannsynsmodellen: LPM

- Den lineære sannsynsmodellen (LPM) brukt på Y_i når Y_i berre kan ta to verdier (0,1) føreset at vi kan tolke $E[Y_i | \mathbf{X}]$ som eit sannsyn
- $E[Y_i | \mathbf{X}] = b_0 + \sum_j b_j x_{ji} = \Pr[Y_i = 1]$
- Dette fører til problem

Er føresetnadene rette i LPM?

- Ein føresetnad i LPM er at residualen e_i stettar krava til OLS
- Residualen er anten $e_i = 1 - (b_0 + \sum_j b_j x_{ji})$ eller $e_i = 0 - (b_0 + \sum_j b_j x_{ji})$
- Dette tyder heteroskedastisitet (residualen varierer med storleiken på x-variablane)
- Det finst estimeringsmetodar som kan komme rundt dette problemet (2-steps vektta minste kvadrats metode til dømes)

LPM er feil modell

- Eit anna problem, at ein for rimelege verdiar av x-ane kan får ein verdi av predikert y der $E[Y_i | \mathbf{X}] > 1$ eller $E[Y_i | \mathbf{X}] < 0$, kan ein ikkje gjere noko med
- LPM er substansielt sett feil modell
- Det trengst ein modell der ein alltid har $0 < E[Y_i | \mathbf{X}] < 1$
- Eit eksempel på LPM:

Regresjon av dikotom avhengig variabel på variabelen "år budd i byen"

ANOVA tabell	Sum of Squares	df	Mean Square	F	Sig.
Regression	3,111	1	3,111	13,648	,000(a)
Residual	34,418	151	,228		
Total	37,529	152			

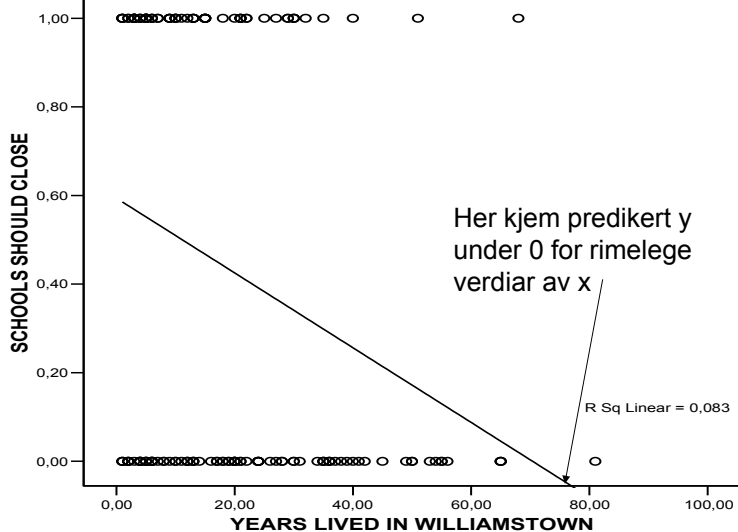
Dependent Variable: SCHOOLS SHOULD CLOSE	B	Std. Error	t	Sig.
(Constant)	,594	,059	10,147	,000
YEARS LIVED IN TOWN	-,008	,002	-3,694	,000

Regresjonen ser heilt OK ut i desse tabellane.

Vår 2004

© Erling Berge 2004

7



Spreingsplott med regresjonslinje. Figur 7.1 Hamilton

Vår 2004

© Erling Berge 2004

8

Den logistiske funksjonen

Den generelle logistiske funksjonen er

- $Y_i = \alpha / (1 + \gamma \cdot \exp[-\beta X_i]) + \varepsilon_i$

Set ein $\alpha = 1$ (øvre grense) og $\gamma = 1$

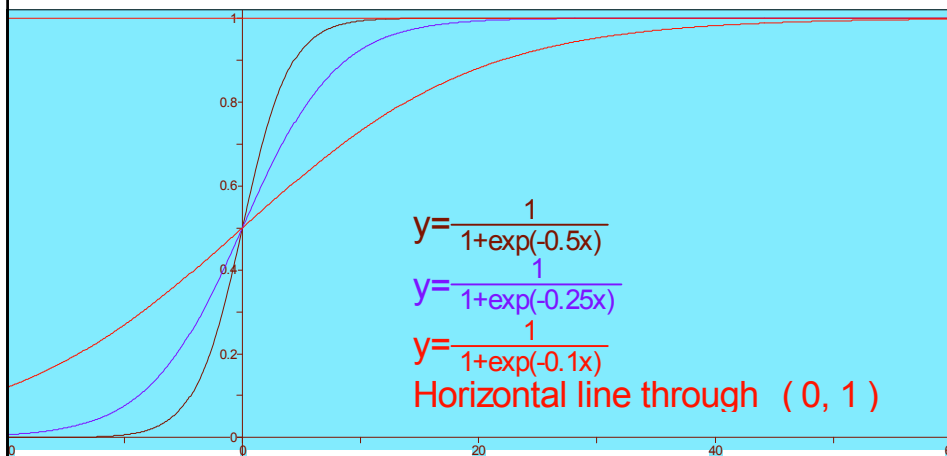
(γ gir plasseringa av det horisontale punkt for rask vekst)

Vil ein alltid ha

- $0 < 1 / (1 + \exp[-\beta X_i]) < 1$

Den logistiske funksjonen vil for alle verdiar av x liggje mellom 0 og 1

Logistiske kurver for ulike β



LOGISTISK REGRESJON: MODELL (1)

Definisjonar

- Sannsynet for at person i skal ha verdien 1 på variabelen Y skriv vi $\Pr(Y_i=1)$. Da er $\Pr(Y_i \text{ ulik } 1) = 1 - \Pr(Y_i=1)$
- Oddsen for at person i skal ha verdien 1 på variabelen Y_i , her kalla O_i , er tilhøvet mellom to sannsyn:

$$O_i(y_i = 1) = \frac{\Pr(y_i = 1)}{1 - \Pr(y_i = 1)} = \frac{p_i}{1 - p_i}$$

LOGISTISK REGRESJON: MODELL (2)

Definisjonar:

- LOGITEN, L_i , er den naturlege logaritmen til oddsen, O_i , for person i :

$$L_i = \ln(O_i)$$

- Modellen føreset at L_i er ein lineær funksjon av forklaringsvariablane x_j ,
- dvs:
- $L_i = \beta_0 + \sum_j \beta_j x_{ji}$, der $j=1, \dots, K-1$, og $i=1, \dots, n$

LOGISTISK REGRESJON: MODELL (3)

- Sett \mathbf{X} = (samlinga av alle x_j), da er sannsynet for at $Y_i = 1$ for person nr i

$$\Pr(y_i = 1) = E[y_i | \mathbf{x}] = \frac{1}{1 + \exp(-L_i)} = \frac{\exp(L_i)}{1 + \exp(L_i)}$$

$$\text{der } L_i = \beta_0 + \sum_{j=1}^{K-1} \beta_j X_{ji}$$

Grafen til dette sambandet er nyttig for tolkinga av kva ei endring i x tyder

LOGISTISK REGRESJON: MODELL (4)

I modellen $Y_i = E[Y_i | \mathbf{X}] + \varepsilon_i$ er feilen enten

- $\varepsilon_i = 1 - E[Y_i | \mathbf{X}]$ med sannsyn $E[Y_i | \mathbf{X}]$
(sidan $\Pr(Y_i = 1) = E[Y_i | \mathbf{X}]$),
- eller feilen er
- $\varepsilon_i = - E[Y_i | \mathbf{X}]$ med sannsyn $1 - E[Y_i | \mathbf{X}]$
- **mao** feilen har ei fordeling kjent som binomialfordelinga med $p_i = E[Y_i | \mathbf{X}]$

LOGISTISK REGRESJON: Estimering

- Metoden brukt for å estimere parametrene i modellen heiter Maximum Likelihood
- ML-metoden gir oss dei parametrene som maksimerer sannsynet (Likelihood) for å finne dei observasjonane vi faktisk har
- Dette sannsynet skal vi kalle \mathcal{L}

Maximum Likelihood (1)

- Likelihooden er lik produktet av sannsynet for kvar einiskild observasjon. For ein dikotom variabel der $\Pr(Y_i = 1) = P_i$ kan dette skrivast

$$\mathcal{L} = \prod_{i=1}^n \left\{ P_i^{Y_i} (1 - P_i)^{(1-Y_i)} \right\}$$

Maximum Likelihood (2)

- For lettere å kunne maksimere sannsynet \mathcal{L} tar ein den naturlege logaritmen til \mathcal{L} :

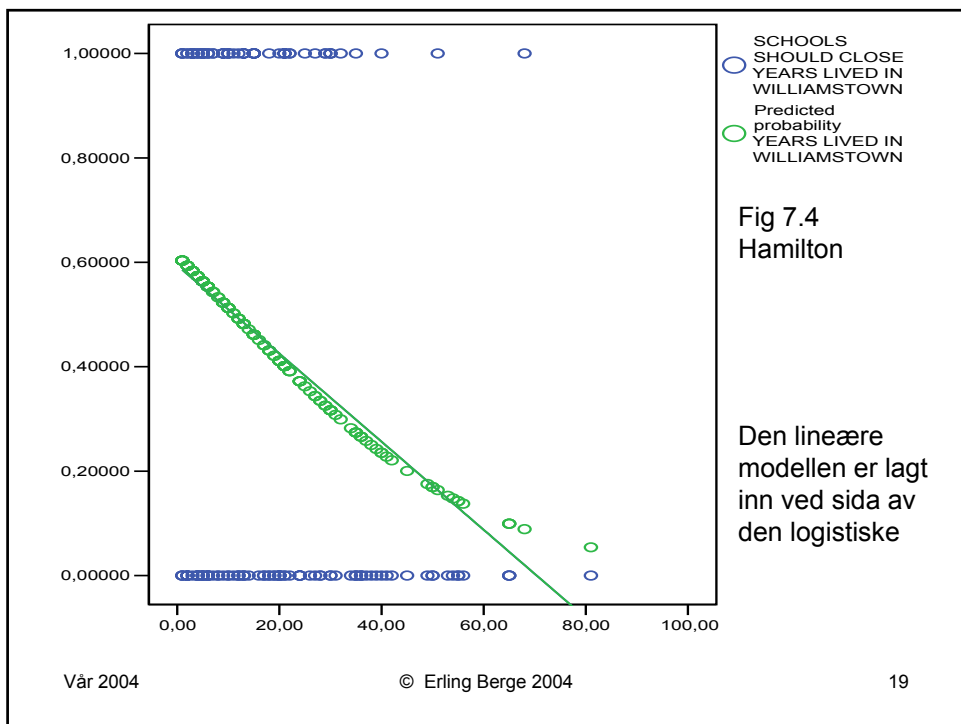
$$\ln(\mathcal{L}) = \sum_{i=1}^n \{ y_i \ln P_i + (1 - y_i) \ln(1 - P_i) \}$$

- Den naturlege logaritmen til \mathcal{L} kallar vi LogLikelihooden, Vi kan kalle den \mathcal{LL} .
- \mathcal{LL} har ei sentral rolle i logistisk regresjon.

Logistisk modell i staden for LPM

Iteration	-2 Log Likelihood	Coefficients	
		Constant	Lived in town
Step 0	209,212	-,275	0
1	195,684	,376	-,034
2	195,269	,455	-,041
3	195,267	,460	-,041
4	195,267	,460	-,041

Dependent: Schools should close	B	S.E.	Wald	df	Sig.	Exp(B)
Lived in town	-,041	,012	11,399	1	,001	,960
Constant	,460	,263	3,069	1	,080	1,584



Tolkning (1)

- Skilnaden mellom den lineære modellen og den logistiske er stor i nærleiken av 0 og 1
- LPM er lett å tolke: $Y_i = \beta_0$ når $x_{1i} = 0$, og når x_{1i} veks med ei eining veks Y_i med β_1 einigar
- Logitmodellen er vanskelegare å tolke. Den er ikkje-lineær både i høve til oddsen og sannsynet.

Tolkning (2)

- Når alle x er lik 0 er $L_i = \beta_0$ Det tyder at oddsen for at $y_i = 1$ i det høvet er $\exp\{\beta_0\}$
- Dersom ein held alle x -ane fast (set dei lik ein konstant) medan x_1 aukar med 1 vil oddsen for at $y_i = 1$ verte multiplisert med $\exp\{\beta_1\}$ Det tyder at den vil endre seg med $100(\exp\{\beta_1\} - 1) \%$
- Sannsynet $\Pr\{y_i = 1\}$ vil endre seg med ein faktor som er påverka av alle elementa i logiten

LOGISTISK REGRESJON: FØRESETNADER

- Modellen er korrekt spesifisert
 - logiten er lineær i parametrene
 - alle relevante variablar er med
 - ingen irrelevante er med
- x -variablane er målt utan feil
- Observasjonane er uavhengige
- Ikkje perfekt multikollinearitet
- Ikkje perfekt diskriminering
- Stort nok utval

Logistisk regresjon

- Dersom føresetnadene er korrekte vil logistisk regresjon gi oss normalfordelte, forventningsrette og variansminimale estimat av parametrene

Manglande Data

ref.: Allison, Paul 2002 "Missing data"

Data manglar av mange grunnar

- Personar nektar å svar
- Personar gløymer eller overser nokre spørsmål
- Personar veit ikkje noko svar
- Spørsmålet er irrelevant
- I administrative register kan somme dokument ha gått tapt
- I forskingsdesign for vanskeleg målbare variablar

Manglane data fører til problem

- Det er eit praktisk problem sidan alle statistiske prosedyrar føreset fullstendige datamatriser
- Det er eit analytisk problem sidan manglande data som regel gir skeive estimat av parametranne
- Det er eit viktig skilje mellom data som manglar av tilfeldige årsaker og dei som manglar av systematiske årsaker

Den enkle løysinga: fjern alle case med manglande data

- Listwise/ casewise fjerning av manglande data tyder at ein fjernar alle case som manglar data på ein eller fleire variablar inkludert i modellen
- Metoden har gode eigenskapar, men kan i somme høve ta ut av analysen mesteparten av casa
- Vanlege alternativ, som parvis ("pairwise") fjerning, har vist seg å vere dårlegare
- Nyare metodar som "maximum likelihood" og "multiple imputation" har betre eigenskapar men er krevjande
- Det løner seg å gjere god arbeid i datainnsamlinga

Typar av tilfeldig missing

- **MCAR: missing completely at random**
 - Tyder at mangel på data for ein person i på variabelen y ikkje er korrelert med verdien på y eller med verdien på nokon anna variabel i datasettet (dette hindrar ikkje at missing i seg sjølv kan korrelere internt case for case)
- **MAR: missing at random**
 - Tyder at mangel på data for ein person i på variabelen y ikkje er korrelert med verdien på y når ein kontrollerer for dei andre variablane i modellen
 - Meir formelt: $\Pr(Y=\text{missing} \mid Y, X) = \Pr(Y=\text{missing} \mid X)$

Prosesen som gir missing

- Kan ignorerast (ignorable)
 - Prosesen kan ignorerast dersom resultatet er MAR og parametrane som styrer missing prosessen ikkje er relatert til dei som skal estimerast
- Kan ikkje ignorerast (non-ignorable)
 - Prosesen kan ikkje ignorerast dersom resultatet ikkje er MAR. Modellestimering krev da ein eigen modell for missing-prosessen (sjå Breen 1996 "Regression Models: Censored, Sample Selected, or Truncated Data", QASS Paper 111, London, Sage,
- I det følgjande er det situasjonen med MAR data som vert drøfta

Konvensjonelle metodar

Vanlege metodar ved MAR data:

- Listewis utelating (Listwise deletion)
- Parwis utelating (Pairwise deletion)
- Dummy variabel korreksjon
- Innsetjing av verdi (Imputation)

Ingen av dei vanleg brukte metodane er tydeleg betre enn listewis utelating

Listevis utelating (1)

- Kan alltid nyttast
- Dersom data er MCAR gir det eit enkelt tilfeldig utval av det opphavslege utvalet
- Mindre n gir sjølvsagt større variansestimater
- Også når data er MAR og missing på x-variablar er uavhengig av verdien på y vil listevis utelating gi forventingsrette estimat

Listevis utelating (2)

- I logistisk regresjon er listevis utelating problematisk berre dersom missing er relatert både til avhengig og uavhengige variablar
- Når missing berre er avhengig av den uavhengige variabelen sine egne verdier er listevis betre enn maximum likelihood og multiple imputation

Parvis utelating

- Tyder at alle utrekningar baserer seg på alt tilgjengeleg materiale sett parvis for alle par av variablar som inngår i analysen
- Dette fører til at ulike parametrar er rekna ut på grunnlag av ulike utval (variasjon i n frå observator til observator)
- Da er alle variansestimater og vanlege testobservatorar skeivt estimert
- Bruk ikkje parvis utelating!

Dummy variabel korreksjon

Dersom data manglar på den uavhengige variabelen x

- Sett $x^* = x$ dersom x ikkje er missing og $x^* = c$ (ein vilkårleg konstant) når x er missing
- Definer $D=1$ hvis x er missing, 0 elles
- Bruk x^* og D i regresjonen i staden for x
- I nominalskalavariabel kan missing få sin eigen dummy

Studiar viser at sjølv med MCAR data er parameterestimata skeive

Bruk ikkje dummy-variabel korreksjon!

Innsetjing av verdi (imputasjon)

- Målet her er å erstatte missing verdier med rimelege gjettingar på kva verdien kunne vere før ein gjennomfører analysen som om dette var verkelege verdier, t.d.
 - Gjennomsnitt av valide verdier
 - Regresjonsestimat basert på mange variablar og case med gyldige observasjonar
- Parameterestimata er konsistente, men variansestimata er skeive (systematisk for små) og testobservatorar er for store
- Unngå om mogeleg å nytte enkel imputasjon

Oppsummering om konvensjonelle metodar for manglande data

- Vanlege metodar for korreksjon av manglande data gjer problema verre
- Ver nøye med datainnsamlinga slik at det er eit minimum av manglande data
- Prøv å samle inn data som kan hjelpe til med å modellere prosessen som fører til missing
- Der data manglar **bruk listevis utelating** dersom ikkje maximum likelihood eller multiple imputasjon er tilgjengeleg

Nye metodar for ignorerbare manglande data (MAR data): Maximum Likelihood

- Konklusjonar
 - Baserer seg på sannsynet for å observere nett dei variabelverdiane vi har funne i utvalet
 - ML gir optimale parameterestimater i store utval når data er MAR
 - Men ML krev ein modell for den felles fordelinga av alle variablane i utvalet som manglar data, og den er vanskeleg å bruke for mange typar modellar

ML-metoden: eksempel (1)

- Observerer y og x for 200 case
- 150 er fordelt som vist
- For 19 case med $Y=1$ er x missing og for 31 case med $Y=2$ er x missing
- Vi ønskjer å finne sannsyna p_{ij} i populasjonen

	Y=1	Y=2
X=1	52	21
X=2	34	43

	Y=1	Y=2
X=1	p_{11}	p_{12}
X=2	p_{21}	p_{22}

ML-metoden: eksempel (2)

- I ein tabell med I rekkjer og J kolonner, fullstendig informasjon om alle case og med n_{ij} case i celle ij er Likelihooden

$$\mathcal{L} = \prod_{i, j} \left(p_{ij} \right)^{n_{ij}}$$

Dvs produktet av alle sannsyn for kvar tabellcelle opphøgd med celfrekvensen som potens

ML-metoden: eksempel (3)

For ein firefeltstabell vert Likelihooden:

$$\mathcal{L} = \left(p_{11} \right)^{n_{11}} \left(p_{12} \right)^{n_{12}} \left(p_{21} \right)^{n_{21}} \left(p_{22} \right)^{n_{22}}$$

For dei 150 casa i tabellen ovanfor der vi har alle observasjonane vert den

$$\mathcal{L} = \left(p_{11} \right)^{52} \left(p_{12} \right)^{21} \left(p_{21} \right)^{34} \left(p_{22} \right)^{43}$$

ML-metoden: eksempel (4)

- For tabellar er ML estimatoren for $p_{ij} = n_{ij}/n$
- Dette gir oss gode estimat i den tabellen der vi ikkje har manglande data (listevis utelating av case)
- Korleis kan ein ta omsyn til det vi veit om y for dei 50 som manglar data på x ?
- Sidan vi har MAR må dei 50 ekstra casa med kjent Y følgje marginalfordelinga til y
- $\Pr(Y=1) = (p_{11} + p_{21})$ og $\Pr(Y=2) = (p_{12} + p_{22})$

ML-metoden: eksempel (5)

- Når vi tar omsyn til alt vi veit om dei 200 casa blir Likelihooden

$$\mathcal{L} = (p_{11})^{52} (p_{12})^{21} (p_{21})^{34} (p_{22})^{43} (p_{11} + p_{21})^{19} (p_{11} + p_{21})^{31}$$

- ML-estimatorane vil no vere

$$\hat{p}_{ij} = \hat{p}(x = i | y = j) \hat{p}(y = j)$$

ML-metoden: eksempel (6)

- Tar vi omsyn til informasjonen vi har om case med manglende data får vi andre estimat av parametrene

Estimat av	Missing utelatt	Missing med
p_{11}	0.346	0.317
p_{21}	0.227	0.208
P_{12}	0.140	0.156
p_{22}	0.287	0.319

ML-metoden

- I det generelle tilfellet av manglende data finst det to tilnærmingar
 - EM metoden, ein tostegsmetode der ein startar med ein forventa verdi på dei manglende data som vert nytta til å estimere parametarar som igjen vert nytta til å gi betre gjetting på forventa verdi som igjen
 - (metoden gir skeive estimat av standardfeil)
 - Direkte ML estimat er betre (men er tilgjengeleg berre for lineære og log-lineære modellar)

Nye metodar for ignorerbare manglande data (MAR data): Multippel Imputasjon

- Konklusjonar
 - Baserer seg på ein tilfeldig komponent som vert lagt til estimat av dei einsskilte manglande opplysningane
 - Har like gode eigenskapar som ML og er enklare å implementere for alle slags modellar.
 - Men den gir ulike resultat for kvar gong den blir brukt

Multiple Imputasjon (1)

- MI har dei same optimale eigenskapane som ML, kan brukast på alle slags data og med alle slags modellar, og kan i prinsippet utførast med vanleg analyseverktøy
- Bruken av MI kan vere temmeleg krokete slik at det er lett å gjere feil. Og sjølv om det vert gjort rett vil ein aldri få same resultat to gonger på grunn av bruken av ein tilfeldig komponent i gjettinga (imputasjonen)

Multiple Imputasjon (2)

- Bruk av data frå enkel imputasjon (med eller utan ein tilfeldig komponent) vil underestimere variansane til parametrane. Konvensjonelle teknikkar klarer ikkje å justere for at data faktisk er generert ved imputasjon
- Løysinga for imputasjon med tilfeldig komponent er å gjenta prosessen mange gonger og bruke den observerte variasjonen i parameterestimat til å justere estimata av variansane
- Allison, side 30-31 forklarar korleis dette kan gjerast

Multiple Imputasjon (3)

- MI krev ein modell som kan nyttast til å gjette på manglande data. Som regel er det føresetnad om normalfordelte variablar og lineære samband. Men modellar kan lagast særskilt for kvart problem
- MI kan ikkje handtere interaksjon
- MI modellen bør ha med alle variablane i analysemodellen (også avhengig variabel)
- MI fungerer berre for måleskalavariabel. Tar ein med nominalskalavariabel trengst spesiell programvare
- Testing av fleire koeffesientar under eitt er meir komplisert

Data som manglar systematisk

- Krev som regel ein modell av korleis fråfallet oppstår
- ML og MI tilnærmingane kan framleis nyttast, men med mye strengare restriksjonar og resultatata er svært sensitive for brot på føresetnadene

Oppsummering

- Dersom nok data vert igjen er listevis utelating den enklaste løysinga
- Dersom listevis utelating ikkje fungerer bør ein freiste med multippel imputasjon
- Dersom ein har mistanke om at data ikkje er MAR må ein lage ein modell for prosessen som skaper missing. Denne kan eventuelt nyttast saman med ML eller MI. Gode resultat krev at modellen for missing er korrekt