



Stockholms  
universitet

# Modeller för studieförframgång i Matematisk Analys IV

Filip Walldén

Kandidatuppsats 2015:9  
Matematisk statistik  
Juni 2015

[www.math.su.se](http://www.math.su.se)

Matematisk statistik  
Matematiska institutionen  
Stockholms universitet  
106 91 Stockholm

# Modeller för studieframgång i Matematisk Analys IV

Filip Walldén\*

Juni 2015

## Sammanfattning

Matematisk Analys IV är av många ansedd som en av de svåraste kurserna som Stockholms universitets matematiska institution erbjuder och läses redan andra terminen för kandidatprogrammet i matematik. I denna rapport kommer vi att arbeta med att ta fram modeller för att prediktera en students betyg utifrån faktorer som kön, ålder och examinator för att undersöka vad som kan tänkas påverka en elevs studieframgång. Vi kommer att finna två modeller, en enkel logistisk regressionsmodell för att få betygen C, B eller A som beror på faktorerna ålder, examinator, tidigare betyg och om man börjar på vår eller höst samt en ordinal logistisk regressionsmodell med faktorerna kön, ålder, examinator, tidigare betyg och om man börjar på höst eller vår. Den ordinala modellen kräver ytterligare arbete för att kunna verifieras men den logistiska modellen kommer att visa på att olika examinators skriver olika svåra eller lätta tentor, att tidigare betyg har stark påverkan, att åldern har en negativ inverkan samt att en vårstart har positiv inverkan på betyget.

---

\*Postadress: Matematisk statistik, Stockholms universitet, 106 91, Sverige.  
E-post: [filipwallden@hotmail.com](mailto:filipwallden@hotmail.com). Handledare: Martin Sköld och Jan-Olov Persson.

# Innehåll

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduktion</b>                             | <b>3</b>  |
| <b>2</b> | <b>Beskrivning av data</b>                      | <b>3</b>  |
| 2.1      | Respons . . . . .                               | 4         |
| 2.1.1    | Betyg i Analys IV . . . . .                     | 4         |
| 2.2      | Prediktor . . . . .                             | 5         |
| 2.2.1    | Kön . . . . .                                   | 5         |
| 2.2.2    | Ålder . . . . .                                 | 5         |
| 2.2.3    | Program . . . . .                               | 5         |
| 2.2.4    | Höst eller vår . . . . .                        | 5         |
| 2.2.5    | Krävda terminer . . . . .                       | 5         |
| 2.2.6    | Examinator . . . . .                            | 6         |
| 2.2.7    | Tidigare betyg . . . . .                        | 6         |
| <b>3</b> | <b>Teori</b>                                    | <b>6</b>  |
| 3.1      | Regressionsmodeller . . . . .                   | 6         |
| 3.1.1    | Odds och oddskvot . . . . .                     | 6         |
| 3.1.2    | Logistisk regression . . . . .                  | 7         |
| 3.1.3    | Ordinal logistisk regression . . . . .          | 7         |
| 3.2      | Selektion och verifiering av modeller . . . . . | 8         |
| 3.2.1    | Akaike informationskriterium . . . . .          | 8         |
| 3.2.2    | Variabelselektion . . . . .                     | 8         |
| 3.2.3    | Goodness of fit . . . . .                       | 9         |
| <b>4</b> | <b>Analys</b>                                   | <b>9</b>  |
| 4.1      | Program . . . . .                               | 10        |
| 4.2      | Logistisk regression . . . . .                  | 10        |
| 4.2.1    | Exempel . . . . .                               | 12        |
| 4.3      | Ordinal logistisk regression . . . . .          | 12        |
| 4.3.1    | Exempel . . . . .                               | 14        |
| 4.4      | Verifiering av modell 2 . . . . .               | 14        |
| <b>5</b> | <b>Diskussion</b>                               | <b>15</b> |
| 5.1      | Resultat . . . . .                              | 16        |
| 5.1.1    | Ålder . . . . .                                 | 16        |
| 5.1.2    | Tidigare Betyg . . . . .                        | 16        |
| 5.1.3    | Start på vår eller höst . . . . .               | 16        |
| 5.1.4    | Examinator . . . . .                            | 17        |
| 5.2      | Förslag till förbättringar . . . . .            | 17        |
| <b>6</b> | <b>Sista ord</b>                                | <b>19</b> |
| <b>7</b> | <b>Appendix</b>                                 | <b>20</b> |

## 1 Introduktion

Matematisk Analys IV är värd 7.5 hp och läses både höst- och vårtermin och behandlar sådant som generaliserade integraler, potensserier, kurvintegraler, trippelintegraler, Gauss och Stokes satser samt Greens formel. Kursen går till stor del ut på att översätta sina analytiska kunskaper från det tvådimensionella planet till det tredimensionella rummet, och sedan bygga på med nya koncept, termer, formler och bevis. Analys IV är en svår kurs och många får betyg som inte lever upp till deras standard eller förhoppningar. Stockholms universitets studenter är inte heller ensamma om detta och det kan höras liknande besvikelser från KTH gällande deras flervariabelsanalys. Men det finns också de som visar stolthet, som kämpat hårt och klarat kursen bra, de som vet hur svår kursen är och fortfarande får betyg bättre än förväntat.

I denna rapport kommer vi att arbeta för att finna logistiska modeller, både binär och ordinal, för att prediktera en individs betyg utifrån faktorer såsom kön, ålder och examinatorer med syftet att undersöka vilka faktorer som är relevanta och hur dessa påverkar studieframgången, både i styrka och riktning.

I kapitel 2 kommer vi att presentera det data vi arbetar med och alla dess tänkbara variabler. Vi kommer sedan, i kapitel 3, att presentera koncept som odds och oddskvot och hur dessa leder till den binära logistiska regressionsmodellen samt till den ordinala logistiska regressionsmodellen. I detta avsnitt kommer vi också att presentera de mått, statistikor samt algoritmer som använts för att finna och verifiera tänkbara modeller. I kapitel 4 presenteras analysen, det är här vi utför all analys och där de resulterande modellerna kommer att presenteras. Vi kommer också kortfattat att gå igenom hur modellerna bör tolkas. I kapitel 5 kommer modellerna att diskuteras och vi kommer ge förslag på hur modellerna kan förbättras för fortsatta studier inom liknande områden.

## 2 Beskrivning av data

Datamaterialet innefattar två utskrifter från Ladok, en för kursen Analys III och en för Analys IV. Varje datafil innehåller alla studenter som fått ett godkänt betyg i kursen sedan december 2007 till augusti 2014 och deras personnummer, program, starttermin, datum för avklarad kurs samt betyg. Arkivet över tidigare tentor användes sedan för att finna examinatorerna för de relevanta tentorna i Analys IV. En del av studenterna saknar betyg i någon av kurserna och ett par av tentorna är ej arkiverade, både tidigare betyg och examinator kommer att visa sig högst relevanta för en students framgång och av denna anledning kommer datasetet att reduceras sådant att alla studenter har fullständig information. Utifrån detta dataset extraherades sedan så många tänkbara förklarande variabler som möjligt och de resulterande variablerna är redovisade i tabellen nedan.

Det reducerade datasetet innehåller 315 unika individer och det finns en hel del kombinationer av variabelnivåer som saknar utfall såsom att datasetet inte innehåller en enda student som fått betyget E i Analys III och sedan A i Analys IV. Jag kommer i fortsättningen att hänföra till detta som att datamaterialet innehåller tomma celler.

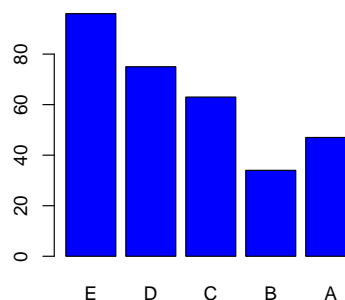
Tabell 1: Variabler

| Förklarande variabel     | Beskrivning  |
|--------------------------|--|
| <i>Kön</i>               | Man eller Kvinna   |
| <i>ÅlderAnalys4</i>      | Individens ålder vid avklarandet av Analys IV                                |
| <i>Program</i>           | Programmet personen läser <sup>1</sup>                                       |
| <i>HöstEllerVår</i>      | Om eleven påbörjade kursen på höst eller vår                                 |
| <i>KrävdaTerminer</i>    | Antal extra terminer personen krävde för att uppnå godkänt betyg i Analys IV |
| <i>ExaminatorAnalys4</i> | Examinatoren för studentens godkända tentamen i Analys IV                    |
| <i>BetygAnalys3</i>      | Betyg E-A som individen fick i Analys III                                    |
| Responsvariabel          | Beskrivning  |
| <i>BetygAnalys4</i>      | Betyg E-A som individen fick i Analys IV                                     |
| <i>ÖverB</i>             | Ja eller Nej för om eleven fick A  |
| <i>ÖverC</i>             | Ja eller Nej för om eleven fick B eller bättre                               |
| <i>ÖverD</i>             | Ja eller Nej för om eleven fick C eller bättre                               |
| <i>ÖverE</i>             | Ja eller Nej för om eleven fick D eller bättre                               |

## 2.1 Respons

### 2.1.1 Betyg i Analys IV

Bland alla godkända studenter fick 96 stycken E, 75 stycken D, 63 stycken C, 34 stycken B och 47 stycken A som illustrerats i figur 1. Utseendet på figuren verkar tyda på att det är svårt att få bra betyg då majoriteten av studenterna antingen fick E eller D och betygsfördelningen i allmänhet verkar vara avtagande. Undantaget är att det förekommer fler studenter med betyget A i Analys IV än studenter med betyget B som verkar bero på att minst en examinator skriver relativt enkla tentor där många studenter får A.



Figur 1: Frekvens för betyg i Analys IV

<sup>1</sup>Programkoder avkodas exempelvis på [sisu.it.su.se](http://sisu.it.su.se)

## 2.2 Prediktor

### 2.2.1 Kön

Datasetet innehåller 105 kvinnor och 210 män. Medianbetyget för båda könen är betyget D men män verkar få högre betyg i större utsträckning än kvinnor.

Tabell 2: Frekvens för betyg i Analys IV per kön

|         | E      | D      | C      | B      | A      |
|---------|--------|--------|--------|--------|--------|
| Män     | 27.14% | 26.67% | 18.10% | 11.90% | 16.19% |
| Kvinnor | 37.14% | 18.10% | 23.81% | 8.57%  | 12.38% |

### 2.2.2 Ålder

Genomsnittsåldern för avklarandet av Analys IV är 23.84 år och medianen är 22. Den yngsta som klarade kursen var 19 år och den äldsta var 48, 25% och 75% kvantilen är 20 respektive 25 år.

### 2.2.3 Program

Datasetet inkluderar 26 stycken olika program och majoriteten av alla studenter läser Analys IV fristående (77 studenter), eller går kandidatprogrammet i matematik (64 studenter) eller kandidatprogrammet i fysik (45 studenter) och 18 program har under 10 stycken studenter. Det skulle därför vara lämpligt att försöka slå ihop ett antal av de mindre programmen. Bland programmen finns 3 stycken matematik och ekonomiprogram och 2 stycken lärarprogram som kan slås ihop. De resterande programmen är dock inte lika naturliga att slå ihop och kommer undersökas analytiskt senare i rapporten.

Matematik och ekonomiprogrammen, och lärarprogrammen slås ihop till SMA och LÄR.

### 2.2.4 Höst eller vår

196 av de godkända studenterna påbörjade kursen på hösten och 119 på våren. Medianbetyget med vårstart är betyget C och medianbetyget med höststart är D.

### 2.2.5 Krävda terminer

Denna variabel beskriver hur många extra terminer studenten krävde för att uppnå godkänt betyg, dvs. en elev som klarar kursen första terminen kommer att ha variabeln `KrävdaTerminer` = 0 och en student som klarade kursen ett år efter kursstart kommer att ha `KrävdaTerminer` = 2. En termin är programmerad att löpa från mars till september och oktober till februrari för att ge möjlighet för en omtentamen utan att behöva ta ut en extra termin.

Bland de 315 studenterna uppnådde 240 stycken (76.19%) godkänt betyg första terminen. Medelvärdet är 0.7619 terminer och det största värdet i datamaterialet är 13 extra terminer.

### 2.2.6 Examinator

I det ursprungliga datasetet fanns det 27 stycken skrivningsdatum för Analys IV och 23 av dessa är arkiverade, och bland dessa 23 tentamer finns det 5 stycken olika examinatorer. Majoriteten av tentamerna är skrivna av examinator 1 och examinator 5 och endast ett fåtal tentamer är skrivna av de övriga examinatorerna.

Tabell 3: Godkända studenter per examinator

| Exam 1 | Exam 2 | Exam 3 | Exam 4 | Exam 5 |
|--------|--------|--------|--------|--------|
| 121    | 14     | 34     | 30     | 116    |

### 2.2.7 Tidigare betyg

Den absoluta majoriteten av studenter har godkänt betyg i Analys III, men ett antal studenter har inte det. Detta kan bero på att studenten läst motsvarande kurs på annan ort eller att studenten inte skrivit ett godkänt betyg i Analys III än men har gjort det i Analys IV, som kan vara möjligt om eleven t.ex. går kandidatprogrammet i matematik som läser Analys III och Analys IV samma termin.

## 3 Teori

### 3.1 Regressionsmodeller

#### 3.1.1 Odds och oddskvot

En stor del av statistiken innefattar kategoriserad data, dvs. data som ej är kontinuerlig och som kan kategoriseras in i ett antal grupper, såsom Ja och Nej, Blå och Röd. Ett centralt begrepp när man handskas med kategorisk data är odds och oddskvoten för olika händelser. Oddset att utfall  $y$  inträffar definieras som:

$$\text{Odds}(Y = y) = \frac{P(Y = y)}{1 - P(Y = y)}.$$

Och oddskvoten definieras som kvoten mellan två stycken odds med olika nivåer på en underliggande variabel. Detta innebär att om vi låter variabeln  $X$  skifta nivå från  $x_0$  till  $x_1$  så definieras oddskvoten för  $y$  givet  $x_1$  som:

$$\text{Oddskvot}(Y = y|X = x_1) = \frac{\text{Odds}(Y = y|X = x_1)}{\text{Odds}(Y = y|X = x_0)}.$$

Oddset är ett mått på proportionen mellan att lyckas mot att misslyckas och oddskvoten beskriver hur mycket denna proportionen ökar eller minskar när en underliggande



variabel skiftar nivå, i detta fall från nivå  $x_0$  till nivå  $x_1$ .

När man arbetar med kategorisk data sätter man en av variabelns nivåer till basnivå som sedan alla andra odds för denna variabel jämförs mot (i ovanstående formel är  $x_0$  basnivå) som ger möjlighet att jämföra oddskvoter mot varandra för att bilda en uppfattning om hur en variabels nivåer påverkar responsen, både i riktning och storlek.

Om variabeln är kontinuerlig ger oddskvoten ett mått på hur mycket som oddset ökar eller minskar när variabeln ökar en enhet, t.ex. 1 år eller 1 kilogram.

### 3.1.2 Logistisk regression

Resultaten från föregående avsnitt ger möjlighet att definiera en ny typ av regressionsmodell, nämligen logistiska regressionsmodellen. Den logistiska regressionsmodellen definieras som

$$\log [\text{Odds}(Y = y | \mathbf{X} = \mathbf{x})] = \alpha + \sum_{i=1}^p \beta_{x_i}^{X_i} \quad (1)$$

där  $\mathbf{X} = (X_1, X_2, \dots, X_p)$ ,  $\mathbf{x} = (x_1, x_2, \dots, x_p)$  och  $\beta_{x_i}^{X_i}$  är parametern för variabel  $X_i$  när  $X_i$  antar nivå  $x_i$  (när  $X_i$  är kontinuerlig är  $\beta_{x_i}^{X_i} = \beta^{X_i} * x_i$ ). Denna modell lämpar sig väl för binära responsvariabler men är bristfällig för att förklara responsvariabler med fler än 2 nivåer på grund av att oddset endast kan beskriva proportionen mellan två utfall, att den specifika händelsen  $Y = y$  inträffar och att händelsen inte inträffar, och det kan därför vara fördelaktigt att välja en annan modell som kan förklara alla nivåer hos responsvariabeln. Av denna anledning kommer modellen ibland refereras till som den binära logistiska regressionsmodellen. Modellen är också praktisk eftersom oddskvoterna för alla förklarande variabler är enkla att beräkna från parametrarna med  $\text{Oddskvot}(Y = y | X_i = x_i) = e^{\beta_{x_i}^{X_i}}$ . Parametrarna för den logistiska regressionsmodellen skattas vanligtvis med iterativt reweighted least square metoden för att finna maximum likelihood skattningarna.

### 3.1.3 Ordinal logistisk regression

Betyget i en kurs är ett exempel på en responsvariabel som kan anta fler än 2 nivåer och modeller som kan förklara fler än 2 nivåer är därför att föredra för att beskriva betyget i Analys IV. Den ordinala logistiska regressionsmodellen är en sådan modell och kräver att responsvariabeln är ordnad, som betyder att det finns en tydlig ordning bland variabelns nivåer, och transformerar responsvariabeln för att skapa ett antal nya binära variabler. Med betyg som responsvariabeln utförs 4 stycken transformationer, en för om betyget är bättre än E, en för om betyget är bättre än D, en för om betyget är bättre än C och en för om betyget är bättre än B. Denna transformation ger 4 stycken binära responsvariabler som tillsammans fullständigt beskriver den ursprungliga responsvariabeln med 5 nivåer och ger möjlighet att använda den logistiska regressionsmodellen genom att sätta upp 4 stycken logistiska modeller för oddset för de olika betygsgränserna:

$$\begin{aligned}
\log [Odds (\text{betyg} > B)] &= \alpha_1 + \sum_{i=1}^p \beta_{x_i}^{X_i} \\
\log [Odds (\text{betyg} > C)] &= \alpha_2 + \sum_{i=1}^p \beta_{x_i}^{X_i} \\
\log [Odds (\text{betyg} > D)] &= \alpha_3 + \sum_{i=1}^p \beta_{x_i}^{X_i} \\
\log [Odds (\text{betyg} > E)] &= \alpha_4 + \sum_{i=1}^p \beta_{x_i}^{X_i}
\end{aligned} \tag{2}$$

Det är viktigt att notera att de 4 stycken logistiska regressionsmodellerna inte skattas oberoende eftersom de 4 modellerna har samma  $\beta_{x_i}^{X_i}$  parametrar och det enda som skiljer de 4 modellerna åt är deras intercept  $\alpha_i$ . Av denna anledning kräver modellen att alla förklarande variabelers nivåer har ungefär samma inverkan på alla betygsgränser, detta kallas för proportional odds assumption och behöver verifieras för alla ordinala logistiska modeller. Om en faktor skattas ha signifikant negativ inverkan vid en av betygsgränserna och sedan signifikant positiv inverkan på en annan så finns det misstankar att tro att antagandet inte håller. Det finns också teststatistikor som kan användas för att undersöka antagandet, men de flesta kräver att datamaterialet inte innehåller tomma celler eller kontinuerliga variabler och det har visat sig vara svårt att hitta ett test som kan appliceras på detta datamaterial. Av denna anledningen har jag inte gått igenom någon teststatistika för att testa antagandet och kommer diskutera alternativa möjligheter för att verifiera eller förkasta modeller i analysdelen av denna rapport. Modellen skattas, precis som den logistiska modellen, vanligtvis med iterativt reweighted least square metoden för att finna maximum likelihood skattningarna.

## 3.2 Selektion och verifiering av modeller

### 3.2.1 Akaike informationskriterium

Akaike informationskriterium (AIC) är ett mått på hur väl en modell förhåller sig till det ursprungliga datat som tar hänsyn till antalet förklarande variabler i modellen och bestraffar modeller med många variabler. AIC definieras enligt Agresti (2013, s. 212) som  $AIC = -2(\text{maximerade loglikelihood} - \text{antalet parametrar i modellen})$ . Man vill att detta mått ska vara så litet som möjligt och det används ofta för att jämföra modeller mot varandra.

### 3.2.2 Variabelselektion

När man sedan har en regressionsmodell är det dags att undersöka vilka förklarande variabler som kan tänkas ha inflytande på responsen. Det finns ett par alternativ för att hitta relevanta variabler och i denna rapport kommer vi att gå igenom forward selection.

Forward selection startar med den minsta tänkbara modellen, modellen som innehåller interceptet men som helt saknar förklarande variabler, och testas sedan att utöka

modellen med en variabel och noterar dess AIC. Sedan tas denna variabel bort och modellen testas utöka med en annan variabel och antecknar dess AIC, detta upprepas tills man testat att utöka modellen med alla möjliga förklarande variabler och forward selection väljer den modellen som har lägst AIC. Forward selection testas sedan att utöka denna modell med alla variabler precis som i tidigare steg och väljer återigen den med lägst AIC. Dessa steg upprepas tills ingen av de utökade modeller har lägre AIC än den man lägger till variablerna till och forward selection väljer då den icke utökade modellen. Ett exempel på denna typ av variabelselektion är bifogat i appendix.

Forward selection kan också välja modell baserat på variablernas signifikans i modellen och algoritmen är mestadels densamma med skillnaden att modellen noterar p-värdet för variabelns signifikans i modellen istället för AIC och väljer att lägga till den variabel med lägst p-värde. Forward selection väljer i detta fall att stanna när ingen av de testade variablerna är signifikant nollskild i modellen, där gränsen för vad som är signifikant vanligtvis är 5% eller 10%.

### 3.2.3 Goodness of fit

För att undersöka om en modell kan tänkas prediktera det ursprungliga datamaterialet krävs det ett goodness of fit test. Det finns många olika test för att undersöka detta för binära logistiska regressionsmodeller, så som Pearsons  $\chi^2$ -test eller Deviance  $G^2$ . Dessa konvergerar dock inte mot en  $\chi^2$ -fördelning när datamaterialet innehåller ogrupperad data eller kontinuerliga eller nästan kontinuerliga variabler (Agresti, 2013, s. 172) och av denna anledning används ett Hosmer-Lemeshow test för goodness of fit som är något mer komplicerad men som har en konvergerande fördelning. Testet börjar med att skatta sannolikheten att lyckas, under modellen som testas, för alla utdrag i datamaterialet. Datamaterialet delas sedan in i ett antal mindre grupper, ofta 10 stycken, efter storleksordning hos de skattade sannolikheterna. De skattade sannolikheterna jämförs sedan mot det observerade utfallet inom dessa grupper och Hosmer-Lemeshows test statistika definieras, enligt Agresti (2013, s. 173), som:

$$H = \sum_{i=1}^g \frac{(\sum_j y_{ij} - \sum_j \hat{\pi}_{ij})^2}{(\sum_j \hat{\pi}_{ij})[1 - (\sum_j \hat{\pi}_{ij})]}$$

Där  $y_{ij}$  är det observerade utfallet, 1 eller 0, och  $\hat{\pi}_{ij}$  är den skattade sannolikheten för individ  $j$  i grupp  $i$  och  $g$  är antalet grupper som specificerats.

Hosmer-Lemeshows test är approximativt  $\chi_{g-2}^2$  fördelat.

## 4 Analys

Vi kommer i försetningen att se kvinnor som basnivå för variabeln **Kön**, start på hösten som basnivå för variabeln **HöstEllerVår**, examinator 5 som basnivå för **ExaminatorAnalys4**, fristående som basnivå för **Program** och betyget E som basnivå för **BetygAnalys3**. Detta innebär att parameterskattningarna för dessa variabelnivåer kommer att fixeras vara lika med 0 i de kommande modellerna och parameterskattningarna för de övriga nivåerna kommer bli ett mått på hur mycket denna variabel

skiljer sig från basnivån.

## 4.1 Program

Tidigare i rapporten upptäckte vi att ett par program nästan är helt tomma och lyckades slå ihop 5 program till 2 nya, ett matematik och ekonomiprogram samt ett lärarprogram. Det verkade dock inte som det fanns några fler program som hade en naturlig anledning att slås ihop och kunde inte reducera antalet program ytterligare. Vi har nu introducerat den logistiska regressionsmodellen som ger oss en ny möjlighet att undersöka programmen. Vi undersöker programmen genom att sätta upp ett antal binära logistiska regressionsmodeller med betygen i Analys IV som respons och alla förklarande variabler och undersöker sedan vilka program som är signifikant skilda från de som läser Analys IV fristående. Det visar sig dock att endast ett par program är signifikant skilda från fristående, och detta gäller endast vid högst ett betygssteg. Det är också endast program med få studenter som visar sig vara signifikanta, såsom SMAEM, ett av matematik och ekonomiprogrammen, med 6 studenter (som ej förblir signifikant efter ekonomiprogrammen slagits ihop). Vi har alltså inte lyckats hitta något program som verkar vara signifikant skild från de som läser kursen fristående och modellen verkar tyda på att det finns få skillnader mellan programmen. Detta gör det svårt att slå ihop program eftersom vi inte kunnat hitta ett enda program som inte kan slås ihop med de som läser kursen fristående och av denna anledning har vi valt att gå vidare utan att kombinera några program. Parameterskattningarna för modellen med responsen  $\text{Betyg} \geq C$  har bifogats i appendix tillsammans med ett anova typ test på modellen, som också visar på att det inte finns några signifikanta skillnader mellan programmen.

## 4.2 Logistisk regression

Datasetet inkluderar många celler som helt saknar utfall, såsom att ingen elev har fått E i Analys III och sedan fått A i Analys IV, detta innebär ett problem för den ordinala logistiska regressionsmodellen som ofta kräver att datamaterialet är stort och inte innehåller tomma celler (Agresti, 2010, s. 58). Av denna anledning inleder vi med att modellera och verifiera en enkel logistisk regressionsmodell för att få ett bra betyg innan vi försöker utöka modellen för att beskriva alla 5 betyg med en ordinal logistisk regressionsmodell.

Vad som är ett bra betyg skiljer från student till student, men i denna studie har vi valt att se betygen C eller bättre som bra betyg då denna gräns närmast delar datamaterialet på mitten och låter därför responsvariabeln vara om studenten fick betyget C eller bättre. Variabeln för tidigare betyg kan ses både som kontinuerlig där E motsvarar 1, D motsvarar 2, C motsvarar 3, B motsvarar 4 och A motsvarar 5 eller som kategorisk med nivåerna A, B, C, D och E. De har båda sina fördelar och nackdelar, och den första tolkningen har fördel i att det endast krävs en skattning för hela variabeln och den andra tolkningen predikterar i allmänhet bättre men kräver 4 skattningar. Vi vet inte på förhand vilket av dessa alternativ som passar bäst för detta datamaterial och av denna anledning låter vi utföra forward selection, för att minimera AIC, 2 gånger med funktionen `step`, första gången där vi låter tidigare betyg vara en kategorisk variabel och sedan en gång när vi transformerar tidigare betyg till en kontinuerlig variabel. Detta resulterar i 2 modeller som har samma för-

klarande variabler men där AIC blir något mindre när betyget ses som kontinuerligt, 316 mot 321, och baserat på AIC väljer vi därför modellen med tidigare betyg sedd som kontinuerlig.

Tabell 4: Parameterskattningarna för att få betyg  $\geq C$

|                  | Estimate | Std. Error | Pr(> z ) |
|------------------|----------|------------|----------|
| $\alpha$         | -1.62    | 0.93       | -        |
| Ålder            | -0.09    | 0.03       | 0.007    |
| Start på hösten  | 0        | -          | -        |
| Start på våren   | 0.83     | 0.36       | 0.022    |
| Betyg Analys III | 1.01     | 0.13       | 0.000    |
| Examinator 1     | -1.05    | 0.38       | 0.006    |
| Examinator 2     | -1.88    | 0.74       | 0.011    |
| Examinator 3     | 1.79     | 0.54       | 0.001    |
| Examinator 4     | 0.08     | 0.56       | 0.880    |
| Examinator 5     | 0        | -          | -        |

I tabell 4 motsvarar den första raden  $\alpha$  skattningen och de resterande raderna är  $\beta$  skattningarna för de förklarande variabler och deras nivåer i modell formel 1. Modellen behöver nu verifieras, och detta gör vi med ett Hosmer-Lemeshows goodness of fit test med funktionen `hoslem.test()` från paketet `ResourceSelection` (2014). Testet ger ett p-värde på 85% som betyder att det inte finns några bevis som tyder på att modellen inte håller och modellen accepteras. I tabell 4 finner vi att examinator 4 inte är signifikant skild från examinator 5 men att alla andra variabler är signifikant skilda från sina basnivåer och vi går vidare med att beräkna oddskvoterna för modellen för att tolka variabelernas inflytande på responsen. Modellen visar

Tabell 5: Oddskvoter för den logistiska modellen

|                  | Oddskvot |
|------------------|----------|
| Ålder            | 0.92     |
| Start på hösten  | 1        |
| Start på våren   | 2.30     |
| Betyg Analys III | 2.75     |
| Examinator 1     | 0.35     |
| Examinator 2     | 0.15     |
| Examinator 3     | 6.01     |
| Examinator 4     | 1.09     |
| Examinator 5     | 1        |

att åldern har ett negativt inflytande på en students chanser att få bra betyg, detta ser vi då oddset att få betyget C eller bättre reduceras med 8% per år. Vårstart och tidigare betyg är båda starkt positiva, en start på våren ökar oddset för betyget C eller bättre med 130% och varje ökat betygssteg i Analys III ökar oddsen att få bra betyg i Analys IV med 175%.

Det verkar finnas starka skillnader bland examinatorerna som också kan verifieras av anova typ testet från föregående avsnitt som ger p-värdet att alla examinatorer har

samma sannolikhet att ge betyget C eller bättre till mindre än 0.001. Examinator 1 har 65% reducerad odds att ge betyg C eller bättre i jämförelse med examinator 5, examinator 2 har 85% reducerad odds, och examinator 3 har 501% ökad odds. Examinator 4 har 9% ökad odds att ge betyget C eller bättre men kan inte signifikant skiljas från examinator 5.

#### 4.2.1 Exempel

Låt oss ta ett exempel för att illustrera hur modellen används. Jag är 23 år och låt oss säga att jag påbörjade kursen föregående höst, skrev under examinator 5 och hade betyget C i Analys III, vad är mina sannolikheter att få betyget C eller bättre och vad är mina sannolikheter att få E eller D?

För detta använder vi våra parameterskattningar från ovanstående modell och modellformel 1 och får

$$\begin{aligned} \log [\text{Odds} (\text{Betyg} \geq C)] &= \alpha + \beta_x^{\mathbf{X}} \\ &= \alpha + \beta^{age} * age + \beta_{h/v}^{termin} + \beta^{betyg} * betyg + \beta_{exam}^{examinator} \\ &= -1.62 - 0.09 * 23 + 0 + 1.01 * 3 + 0 \\ &= -0.66 \end{aligned}$$

som sedan kan användas för att finna sannolikheten att få betyget C eller bättre.

$$P(\text{Betyg} \geq C) = \frac{e^{-0.66}}{1 + e^{-0.66}} = 0.34$$

Fullständiga uträkningar för att beräkna en sannolikheten från en oddskvot är bifogat i appendix. Resultatet säger oss att jag har ungefär 34% chans att få betyget C eller bättre under dessa omständigheter, eller 36% med ej avrundade parameterskattningar. Om man testar att byta examinator, t.ex. till examinator 1, finner vi samma sannolikhet till ungefär 16%. Oddskvoten mellan dessa två händelser är  $\frac{0.16/(1-0.16)}{0.36/(1-0.36)} = 0.34$  som stämmer bra överens med oddskvoten för examinator 1 som står i tabellen och beräknades enligt  $e^{-1.05} = 0.35$ . Vid ej avrundade parameterskattningar och sannolikheter gäller exakt likhet mellan dessa oddskvoterna.

### 4.3 Ordinal logistisk regression

Nu när den logistiska regressionsmodellen är accepterad går vi vidare med att undersöka om det är möjligt att hitta och verifiera en ordinal logistisk regressionsmodell. Vi börjar med att finna den modell som minimerar AIC och detta gör vi med forward selektion, och låter utföra algoritmen en gång där tidigare betyg ses som en kategorisk variabel och en gång som kontinuerlig, precis som vid den logistiska regressionen. Vi finner återigen att tidigare betyg bör ses som kontinuerlig eftersom vi finner AIC till 844 respektive 838, men att den ordinala regressionen tar med kön i modellen, något som inte den logistiska modellen gjorde. I tabell 6 visas alla parameterskattningar för den ordinala logistiska regressionsmodell som minimerar AIC. Tabellen tolkas till stor del som den logistiska modellen, de första 8 raderna är  $\beta$  skattningarna för de förklarande variablerna och deras nivåer och de resterande 4 raderna är  $\alpha$  skattningarna i formeln för modell 2. Vi observerar att de flesta faktorer är signifikant skilda från sina basnivåer men att examinator 2 och 4 inte är signifikant skilda

Tabell 6: Parameterskattningarna för den ordinala modellen

|                  | Estimate | Std. Error | Pr(> z ) |
|------------------|----------|------------|----------|
| Kvinna           | 0        | -          | -        |
| Man              | 0.34     | 0.23       | 0.132    |
| Ålder            | -0.06    | 0.02       | 0.008    |
| Start på höst    | 0        | -          | -        |
| Start på vår     | 0.62     | 0.28       | 0.023    |
| Betyg Analys III | 0.90     | 0.10       | 0.000    |
| Examinator 1     | -0.92    | 0.28       | 0.001    |
| Examinator 2     | -1.06    | 0.56       | 0.059    |
| Examinator 3     | 1.54     | 0.39       | 0.000    |
| Examinator 4     | 0.42     | 0.40       | 0.286    |
| Examinator 5     | 0        | -          | -        |
| $\alpha_1$       | -4.23    | 0.70       | -        |
| $\alpha_2$       | -3.29    | 0.69       | -        |
| $\alpha_3$       | -1.96    | 0.67       | -        |
| $\alpha_4$       | -0.52    | 0.67       | -        |

från examinator 5 och att män inte är signifikant skilda från kvinnor på 5% nivån. Vi beräknar nu oddskvoterna för att försöka tolka hur de förklarande variablerna påverkar betyget. Oddskvoterna för den ordinala logistiska regressionsmodellen är

Tabell 7: Oddskvoter för den ordinala

|                  | Oddskvot |
|------------------|----------|
| Kvinna           | 1        |
| Man              | 1.41     |
| Ålder            | 0.94     |
| Start på höst    | 1        |
| Start på vår     | 1.87     |
| Betyg Analys III | 2.45     |
| Examinator 1     | 0.40     |
| Examinator 2     | 0.35     |
| Examinator 3     | 4.65     |
| Examinator 4     | 1.53     |
| Examinator 5     | 1        |

exakt samma för alla fyra betygsgränser som betyder att om oddset att få betyget D eller bättre ökar med en viss procent så ökar oddset att få C eller bättre, oddset att få B eller bättre och oddset att få A alla med precis samma procent. Av denna anledning kommer jag i detta avsnitt endast redovisa oddset för bra betyg, och när jag säger att en variabelnivå ökar oddset att få bra betyg med  $x\%$  innebär detta att oddset att få betygen D eller bättre, oddset att få C eller bättre, oddset att få B eller bättre och oddset att få A alla ökar med  $x\%$ .

Den ordinala regressionen visar på att åldern har negativt inflytande på betyget då oddset för bra betyg reduceras med 6% per år. En start på våren och tidigare

betyg påverkar båda positivt på betyget i Analys IV, en vårstart ökar oddset att få ett bra betyg med 87% och varje ökat betygssteg i Analys III ökar oddset att få ett bra betyg i Analys IV med 145%. Examinator 1 reducerar oddset för bra betyg med 60% och examinator 3 ökar oddset för bra betyg med 365%. Män har 41% icke signifikant ökad odds att skriva bra betyg, examinator 2 har 65% icke signifikant reducerad odds och examinator 4 har 53% icke signifikant ökad odds.

### 4.3.1 Exempel

Med hjälp av den ordinala logistiska regressionsmodellen kan vi nu beräkna sannolikheten för alla betygssteg för en student, vi väljer att fortsätta från föregående exempel där jag är 23 år, man, startade på hösten, skrev under examinator 5 och fick betyget C i Analys III. Beräkandet blir till stor del analog med föregående exempel och behöver utföras 4 stycken gånger, en för varje logistiska modell, för att finna 4 stycken sannolikheter,  $P(\text{Betyg} > B)$ ,  $P(\text{Betyg} > C)$ ,  $P(\text{Betyg} > D)$  och  $P(\text{Betyg} > E)$ .

I detta fall så blir de kumulativa sannolikheterna

$$P(\text{Betyg} > B) = 0.07$$

$$P(\text{Betyg} > C) = 0.16$$

$$P(\text{Betyg} > D) = 0.42$$

$$P(\text{Betyg} > E) = 0.75$$

Vi ser att denna modell skiljer sig lite från den logistiska regressionsmodellen eftersom vi fann sannolikheten för betyget C eller bättre till 36% i föregående exempel och här finner vi sannolikheten för samma händelse till 42%.

Ur dessa kumulativa sannolikheter går det också att beräkna sannolikheterna för alla betyg som

$$\begin{aligned} P(\text{Betyg} = A) &= P(\text{Betyg} > B) \\ &= 0.07 \end{aligned}$$

$$\begin{aligned} P(\text{Betyg} = B) &= P(\text{Betyg} > C) - P(\text{Betyg} > B) = 0.16 - 0.07 \\ &= 0.09 \end{aligned}$$

$$\begin{aligned} P(\text{Betyg} = C) &= P(\text{Betyg} > D) - P(\text{Betyg} > C) = 0.42 - 0.16 \\ &= 0.26 \end{aligned}$$

$$\begin{aligned} P(\text{Betyg} = D) &= P(\text{Betyg} > E) - P(\text{Betyg} > D) = 0.75 - 0.42 \\ &= 0.33 \end{aligned}$$

$$\begin{aligned} P(\text{Betyg} = E) &= 1 - P(\text{Betyg} > E) = 1 - 0.75 \\ &= 0.25 \end{aligned}$$

## 4.4 Verifiering av modell 2

Den ordinala logistiska regressionen kan beskriva proportionen mellan alla godkända betyg i jämförelse med den binära logistiska regressionsmodellen som endast kan beskriva förhållandet mellan bra och mindre bra betyg, och den ordinala är därför att föredra om den kan undersökas och accepteras. För att kunna acceptera modellen måste proportionella odds villkoret undersökas, något som är betydligt svårare att



göra för ett datamaterial som innehåller tomma celler eller kontinuerliga variabler. Agresti (2010, s. 68) skriver att det, vid publiceringen av hans bok, finns 2 möjliga test men att dessa inte är förprogrammerade i något av de nuvarande statistiska programspråken och det har visat sig vara svårt att hitta programpaket som kan utföra något av dessa test och vi använder ett alternativt tillvägagångssätt publicerat av Brant (1990).

Metoden går ut på att sätta upp en ny ordinal logistiska regressionsmodell med samma förklarande variabler utan att fixera  $\beta$ -skattningarna och sedan jämföra  $\beta$ -skattningarna, både inom modellen och mot den ordinala med fixerat  $\beta$ . Detta utförs enklast genom att sätta upp 4 stycken logistiska regressionsmodeller, en för varje betygssteg, med de förklarande variablerna kön, ålder, start på höst- eller vårtermin, tidigare betyg sett som kontinuerlig samt examinator.

Vi har 8 stycken variabelnivåer som vi vill undersöka och varje variabelnivå har 5 stycken skattningar, 4 stycken från de logistiska och 1 för den ordinala, och bifogar tabellen för  $\beta$  skattningarna för examinator 1 här nedan och de resterande skattningarna i appendix.

Tabell 8: Examinator 1

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | -0.92    | 0.28       |
| Log[Betyg > B] | -2.93    | 0.81       |
| Log[Betyg > C] | -2.05    | 0.50       |
| Log[Betyg > D] | -1.05    | 0.38       |
| Log[Betyg > E] | -0.39    | 0.34       |

Skattningarna skiljer sig kraftigt och är ett tecken på att antagandet om proportionella odds inte är uppfyllt.

Skattningarna för många av de resterande variablerna är betydligt bättre och verkar kunna uppfylla antagandet, men för att den ordinala logistiska regressionsmodellen ska hålla så måste alla förklarande variabler ha proportionella odds. Det finns därför väldigt starka misstankar om att den ordinala logistiska regressionsmodellen inte håller eftersom vi funnit att examinator 1 inte verkar uppfylla antagandet för proportionella odds.

## 5 Diskussion

Vi har funnit två modeller, en binär logistisk och en ordinal logistisk, och har kunnat verifiera den logistiska modellen men har starka misstankar om att den ordinala modellen inte håller. Av denna anledning kommer resultaten i denna diskussionsdel vara baserade på den binära logistiska modellen och den ordinala logistiska regressionen blir ett möjligt område att utöka studien.

I resultatdelen har ett 2-sidigt 95% konfidensintervallen beräknats för alla oddskvoter och står inom parentes för alla signifikanta prediktorer.

## 5.1 Resultat

Tabell 9: Den logistiska modellens oddskvoter med 95% konfidensintervall

|                  | Oddskvot | 2.5% kvantil | 97.5% kvantil |
|------------------|----------|--------------|---------------|
| Ålder            | 0.92     | 0.86         | 0.98          |
| Start på hösten  | 1        | 1            | 1             |
| Start på våren   | 2.30     | 1.13         | 4.70          |
| Betyg Analys III | 2.75     | 2.12         | 3.56          |
| Examinator 1     | 0.35     | 0.16         | 0.74          |
| Examinator 2     | 0.15     | 0.04         | 0.66          |
| Examinator 3     | 6.01     | 2.07         | 17.44         |
| Examinator 4     | 1.09     | 0.37         | 3.23          |
| Examinator 5     | 1        | 1            | 1             |

### 5.1.1 Ålder

Oddskvoten för åldern visar på att oddset att få bra betyget signifikant reduceras med 8% (2, 14) per år som kan förklaras på många sätt, men jag tror att den huvudsakliga anledningen är brist på tid. Jag tror att ju äldre man blir, desto mindre fritid har man som kan använda för att fokusera på sina studier. Någon som kommer direkt från gymnasium kan ha ett jobb på sidan men studerar annars på heltid, men när man blir äldre kan man ha familj och jobb som kommer i första hand. Det blir därför väldigt naturligt att man låter studierna släpa efter när arbete eller familj kräver ens uppmärksamhet, och det är inte konstigt om betygen lider. Det är kanske ett något deprimerande resultat att komma fram till, men det var tyvärr inte alldeles för oväntat heller.

### 5.1.2 Tidigare Betyg

Varje ökat betygssteg i Analys III ökar signifikant oddset att få bra betyg i Analys IV med 175% (112, 256). Att tidigare kunskaper har stor betydelse för slutbetyget i Analys IV är knappast någon överraskning och kommer därför inte att diskuteras i allt för stor utsträckning i sig själv men ger möjlighet att jämföra olika variabelers inflytande på ett mått som kan vara enklare att förstå för någon som inte arbetat med odds och oddskvoter tidigare.

### 5.1.3 Start på vår eller höst

Den logistiska modellen visar på att en start på våren signifikant ökar oddset att få bra betyg med 130% (13, 370) jämfört med oddset att få bra betyg om man börjar på hösten. Det var spekulerat att en vårstart skulle vara fördelaktigt men inte till denna omfattning, modellen visar på att en vårstart istället för en höststart påverkar betyget i Analys IV nästan lika mycket som ett helt betygssteg i Analys III. Det skulle kunna röra som om att det är olika föreläsare på höst- och vårtermin och jag undersökte detta genom att titta på gamla kurshemsidor. Jag kunde inte finna något mönster bland föreläsarna och termin med blotta ögat men frågan kan undersökas

genom att lägga till en ny prediktor för föreläsare och undersöka om några föreläsare bidrar signifikant negativt eller positivt till studenternas studieframgång och sedan undersöka om dessa är korrelerade med variabeln **HöstEllerVår**. Detta har inte jag gjort i denna rapport och accepterar att detta kan vara en möjlig förklarande faktor jag inte tagit hänsyn till.

Det finns också andra faktorer som kan förklara detta resultat.

Kursen går på andra perioden på både höst och vår som ger studenterna en halv termin att ställa om sig från sin sommarledighet, och förklaringen att hösten bidrar negativt på grund av att terminen följer sommarlovet verkar därför inte allt för trolig men ger större möjlighet för en annan faktor, nämligen årstiderna i sig. Jag tror att det är fullt möjligt att årstidernas mörker och allmänna stämning kan påverka studieframgången hos studenter, men för att undersöka signifikans och styrka behöver vi eliminera alla andra faktorer som är korrelerade med höst- och vårtermin. En annan förklaring är att olika program möjligtvis läser Analys IV på olika terminer, men att programmen inte slagits ihop tillräckligt för att de ska visas vara signifikanta vid något test och att program variabeln innehåller alldeles för många nivåer för att den ska väljas baserat på AIC. För att ta reda på om detta är fallet så måste programmen studeras närmare, något som vi inte gjort i denna rapport, och kan därför inte svara på om det är primärt en faktor som bidrar positivt eller negativt eller om det är kombination av de ovanstående.

#### 5.1.4 Examinator

Modellen visar på att det finns signifikanta skillnader mellan examinatorer, och att de är kraftiga. Examinator 4 är ej är signifikant skild från examinator 5 och kommer därför inte diskuteras. Att skriva för examinator 1 och examinator 2 reducerar signifikant oddset för bra betyg med 65% (26, 84) respektive 85% (34, 96) jämfört med oddset att skriva bra betyg för examinator 5. Detta är ekvivalent med att skriva 1 respektive 2 betygssteg lägre på tentamen i Analys III och är en enorm skillnad mot examinator 3 som signifikant ökar oddset för bra betyg med 501% (107, 1644), som är ekvivalent med att skriva strax under 2 betygssteg bättre i Analys III. De skiljer nästan 3 respektive 4 betygssteg mellan examinator 1 och 3 samt examinator 2 och 3.

Man bör dock notera att examinator 2, 3 och 4 endast skrivit ett fåtal tentor under denna period, men både examinator 1 och 5 har skrivit många tentor och det verkar därför inte vara en engångångföreteelse att examinator 1 och 5 skriver olika svåra tentor. Det är dock inte min uppgift att avgöra om examinator 1 skriver för svåra tentamina eller om examinator 5 skriver för enkla.

## 5.2 Förslag till förbättringar

Den ordinala logistiska regressionsmodellen kräver stora datamaterial eftersom denna typ av modell har svårigheter att behandla tomma celler (Agresti, 2010, s. 58-68) och den största förbättringen av modellen är att fylla dessa celler. Detta kan åstadkommas genom att studera programmen noggrannare, både analytiskt och bakgrundsmässigt, för att dela in programmen i färre grupper för att reducera antalet celler, dvs. antalet möjliga kombinationer av variabelnivåer. Om antalet totala celler reduceras behöver inte heller datamaterialet vara lika stort för att fylla alla celler och därmed

eliminera tomma cellerna. Modellen har också svårigheter med kontinuerliga variabler och det är därför att föredra om man valde att se tidigare betyg som kategorisk och delade in åldern i ett antal kategoriserade nivåer, möjligtvis binär för om man är över eller under median ålder. Detta skulle möjliggöra de vanliga teststatistikorna som används för att verifiera ordinala modeller, så som Pearsons  $\chi^2$  och Deviance  $G^2$ .

Detta datamaterial kan dock inte korrigeras för att få icke tomma celler utan att ta bort tidigare betyg som en prediktor, något som inte rekommenderas eftersom variabeln förklarar mycket av responsen, och är därför inte ett lämpligt datamaterial att bygga den ordinala logistiska regressionen på. Jag föreslår därför att fortsatta studier med ordinal logistiska regression bör utföras på ett större datamaterial där materialet har korrigerats likt ovanstående. Större datamaterial kan åstakommas genom att exempelvis undersöka kurser med fler studenter eller genom att utöka studien till att innehålla flera universitet och högskolor och byta examineringsvariabeln mot en platsvariabel.

## 6 Sista ord

Jag vill tacka både Martin Sköld och Jan-Olov Persson som varit mina handledare för detta arbete. Martin hjälpte mig hämta data och visade mig den ordinala logistiska regressionsmodellen som jag tillägnat mestadels av mitt arbete att undersöka, jag har sedan fått arbeta helt fritt men de har funnits att prata och diskutera med när jag haft mina funderingar. Tack vare Jan-Olov och Martin så har jag lärt mig mycket och hade det inte varit för dem så hade jag inte kunnat presentera mitt arbete vid detta tillfälle och för det är jag oerhört tacksam.

## 7 Appendix

### Ett exempel på forward selection

```
## Start: AIC=436
## ÖverD ~ 1
##
##           Df Deviance AIC
## + as.numeric(BetygAnalys3) 1     358 362
## + ExaminatorAnalys4       4     408 418
## + KrävdaTerminer         1     417 421
## + ÅlderAnalys4           1     419 423
## + HöstEllerVår          1     419 423
## <none>                   434 436
## + Kön                    1     434 438
## + Program                22     400 446
##
## Step: AIC=362
## ÖverD ~ as.numeric(BetygAnalys3)
##
##           Df Deviance AIC
## + ExaminatorAnalys4  4     314 326
## + ÅlderAnalys4      1     348 354
## + HöstEllerVår     1     350 356
## + KrävdaTerminer   1     353 359
## <none>             358 362
## + Kön              1     358 364
## + Program          22     328 376
##
## Step: AIC=326
## ÖverD ~ as.numeric(BetygAnalys3) + ExaminatorAnalys4
##
##           Df Deviance AIC
## + ÅlderAnalys4    1     305 319
## + HöstEllerVår   1     308 322
## + KrävdaTerminer 1     310 324
## <none>           314 326
## + Kön            1     314 328
## + Program        22     285 341
##
## Step: AIC=319
## ÖverD ~ as.numeric(BetygAnalys3) + ExaminatorAnalys4 + ÅlderAnalys4
##
##           Df Deviance AIC
## + HöstEllerVår   1     300 316
## + KrävdaTerminer 1     303 319
## <none>           305 319
```

```

## + Kön          1      305 321
## + Program      22      275 333
##
## Step:  AIC=316
## ÜberD ~ as.numeric(BetygAnalys3) + ExaminerAnalys4 + ÅlderAnalys4 +
##       HöstEllerVår
##
##              Df Deviance AIC
## <none>              300 316
## + KrävdaTerminer  1      299 317
## + Kön              1      300 318
## + Program          22      273 333

```

Vi ser att modellen med endast intercept har ett AIC på 436 och R listar sedan upp alla variabler som vi skulle kunna utöka modellen med och dess resulterande AIC. Vi väljer sedan modellen med lägst AIC som är tidigare betyg med ett AIC på 362. R listar sedan återigen upp alla variabler som modellen kan utökas med och deras AIC, vi finner att examinerarna har lägst AIC på 326 och vi väljer den modellen. Modellen har nu både tidigare betyg och examinerer som förklarande variabler med ett AIC på 326. Vi fortsätter att utöka modellen, först med ålder och om man börjar på höst eller vår för att få ner AIC till 316. Vi testar nu att utöka modellen med antalet extra terminer och får ett AIC på 317, sedan med kön och får ett AIC på 318 och sedan med program och får ett AIC på 333. Ingen av dessa tre modeller har lägre AIC än den icke utökade modellen med betyg, examiner, ålder och höst/vårstart och vi väljer att inte utöka modeller något mera. Modellen blir:  $\text{ÜberD} = \text{BetygAnalys4} + \text{ExaminerAnalys4} + \text{ÅlderAnalys4} + \text{HöstEllerVår}$

## Full logistisk modell för att få betyget C eller bättre

Tabell 10: Parameterskattningar för betyg  $\geq C$

|                          | Estimate | Std. Error | Pr(> z ) |
|--------------------------|----------|------------|----------|
| (Intercept)              | -1.36    | 1.20       | 0.257    |
| KönMan                   | -0.03    | 0.35       | 0.939    |
| as.numeric(BetygAnalys3) | 1.10     | 0.15       | 0.000    |
| ProgramKOMBM             | -2.51    | 1.39       | 0.072    |
| ProgramKP002             | 0.40     | 1.99       | 0.842    |
| ProgramNASTK             | -1.64    | 0.91       | 0.070    |
| ProgramNASTM             | 2.40     | 1.58       | 0.129    |
| ProgramNBERK             | 15.66    | 2399.54    | 0.995    |
| ProgramNBFFK             | 0.89     | 1.99       | 0.656    |
| ProgramNBIBK             | 14.90    | 1695.41    | 0.993    |
| ProgramNBIMA             | -1.31    | 1.93       | 0.497    |
| ProgramNBIMK             | 1.39     | 1.24       | 0.265    |
| ProgramNDATK             | -1.79    | 1.72       | 0.297    |
| ProgramNDAVK             | 15.10    | 2399.54    | 0.995    |
| ProgramNFYSK             | -0.35    | 0.54       | 0.516    |
| ProgramNFYSM             | 1.13     | 0.80       | 0.159    |
| ProgramNMATK             | 0.16     | 0.50       | 0.743    |
| ProgramNMETE             | -0.72    | 1.05       | 0.492    |
| ProgramNMETK             | 0.73     | 0.85       | 0.390    |
| ProgramNMFJK             | 0.54     | 1.13       | 0.630    |
| ProgramNSFKY             | 0.02     | 1.00       | 0.981    |
| ProgramNSFPY             | -0.80    | 0.86       | 0.350    |
| ProgramNSUFY             | 0.24     | 1.04       | 0.814    |
| ProgramLÄR               | -16.01   | 1054.50    | 0.988    |
| ProgramSMA               | -1.45    | 0.98       | 0.138    |
| ÅlderAnalys4             | -0.09    | 0.04       | 0.012    |
| ExaminatorAnalys4Exam1   | -1.41    | 0.45       | 0.002    |
| ExaminatorAnalys4Exam2   | -2.05    | 0.81       | 0.011    |
| ExaminatorAnalys4Exam3   | 1.77     | 0.64       | 0.006    |
| ExaminatorAnalys4Exam4   | -0.20    | 0.62       | 0.747    |
| KrävdaTerminer           | -0.14    | 0.13       | 0.285    |
| HöstEllerVårV            | 0.67     | 0.51       | 0.188    |



## Anova typ test för den fulla logistisk modellen

Tabell 11: Anova typ test för betyg  $\geq C$

|                          | Df | Deviance | Resid. Df | Resid. Dev | Pr(>Chi) |
|--------------------------|----|----------|-----------|------------|----------|
| NULL                     |    |          | 314       | 434.37     |          |
| Kön                      | 1  | 0.06     | 313       | 434.31     | 0.810    |
| as.numeric(BetygAnalys3) | 1  | 76.28    | 312       | 358.03     | 0.000    |
| Program                  | 25 | 32.63    | 287       | 325.41     | 0.141    |
| ÅlderAnalys4             | 1  | 13.10    | 286       | 312.30     | 0.000    |
| ExaminatorAnalys4        | 4  | 39.35    | 282       | 272.95     | 0.000    |
| KrävdaTerminer           | 1  | 1.77     | 281       | 271.19     | 0.184    |
| HöstEllerVår             | 1  | 1.26     | 280       | 269.93     | 0.261    |

## Fortsättning från exempel i avsnittet Logistisk regression

Vi låter  $P = P(\text{Betyg} \geq C)$  och får

$$\text{Odds}(\text{Betyg} \geq C) = e^{-0.66}$$

$$\frac{P}{1-P} = e^{-0.66}$$

$$P = e^{-0.66} * (1 - P)$$

$$P = e^{-0.66} - e^{-0.66} * P$$

$$P + e^{-0.66} * P = e^{-0.66}$$

$$P * (1 + e^{-0.66}) = e^{-0.66}$$

$$P = \frac{e^{-0.66}}{1 + e^{-0.66}} = 0.34$$

## Parameterskattningarna för den ordinala modellen samt för de logistiska

Tabell 12: Man

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | 0.34     | 0.23       |
| Log[Betyg > B] | 0.26     | 0.42       |
| Log[Betyg > C] | 0.44     | 0.35       |
| Log[Betyg > D] | -0.06    | 0.31       |
| Log[Betyg > E] | 0.51     | 0.29       |

Tabell 13: Ålder

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | -0.06    | 0.02       |
| Log[Betyg > B] | -0.04    | 0.05       |
| Log[Betyg > C] | -0.08    | 0.04       |
| Log[Betyg > D] | -0.09    | 0.03       |
| Log[Betyg > E] | -0.04    | 0.03       |

Tabell 14: Start på vår

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | 0.62     | 0.28       |
| Log[Betyg > B] | 0.77     | 0.49       |
| Log[Betyg > C] | 0.68     | 0.41       |
| Log[Betyg > D] | 0.84     | 0.36       |
| Log[Betyg > E] | 0.23     | 0.35       |

Tabell 15: Tidigare betyg

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | 0.90     | 0.10       |
| Log[Betyg > B] | 1.07     | 0.21       |
| Log[Betyg > C] | 1.07     | 0.17       |
| Log[Betyg > D] | 1.01     | 0.13       |
| Log[Betyg > E] | 0.73     | 0.12       |

Tabell 16: Examinator 1

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | -0.92    | 0.28       |
| Log[Betyg > B] | -2.93    | 0.81       |
| Log[Betyg > C] | -2.05    | 0.50       |
| Log[Betyg > D] | -1.05    | 0.38       |
| Log[Betyg > E] | -0.39    | 0.34       |

Tabell 17: Examinator 2

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | -1.06    | 0.56       |
| Log[Betyg > B] | -1.48    | 0.88       |
| Log[Betyg > C] | -1.68    | 0.80       |
| Log[Betyg > D] | -1.88    | 0.74       |
| Log[Betyg > E] | -0.33    | 0.71       |

Tabell 18: Examinator 3

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | 1.54     | 0.39       |
| Log[Betyg > B] | 1.00     | 0.64       |
| Log[Betyg > C] | 1.70     | 0.58       |
| Log[Betyg > D] | 1.79     | 0.54       |
| Log[Betyg > E] | 1.57     | 0.59       |

Tabell 19: Examinator 4

|                | Estimate | Std. Error |
|----------------|----------|------------|
| Ordinal        | 0.42     | 0.40       |
| Log[Betyg > B] | 0.02     | 0.52       |
| Log[Betyg > C] | -0.02    | 0.52       |
| Log[Betyg > D] | 0.07     | 0.56       |
| Log[Betyg > E] | 0.74     | 0.63       |

## Referenser

- Agresti, Alan (2010). *Analysis of Ordinal Categorical Data*. Wiley.
- Agresti, Alan (2013). *Categorical Data Analysis*. Wiley.
- Brant, Rollin (1990). Assessing Proportionality in the Proportional Odds Model for Ordinal Logistic Regression. *Biometrics* **46**, 1171-1178
- Subhash R. Lele, Jonah L. Keim and Peter Solymos (2014). ResourceSelection: Resource Selection (Probability) Functions for Use-Availability Data. R package version 0.2-4. <http://CRAN.R-project.org/package=ResourceSelection>