Business Analytics & Data Mining Modeling Using R Dr. Gaurav Dixit Department of Management Studies Indian Institute of Technology, Roorkee

Lecture - 52 Logistic Regression-Part VII

Welcome to the course Business Analytics and Data Mining Modelling Using R. So, in previous few lectures we have been discussing different aspects of logistic regression so we will continue that discussion in this particular lecture as well.

So, in previous lecture we will be used apply details data set and then promotional offers data set as well. Some of the details regarding modelling exercise we could not cover.

(Refer Slide Time: 00:44)



So, we will do that and discuss. So, let us move to let us import this data set. So, we will let us import the library x plus x. So, we are again going to use this particular data set flight details.

(Refer Slide Time: 01:00)



So, let us import this. Now, let us remove any columns any rows, let us look at these are the observations.

(Refer Slide Time: 01:18)

er hår forte verw rhort verson kunt isetnig renter toos weip	
🐑 • 🔐 • 📄 🗐 🦀 (🕭 to to file/function) 📓 • Addies •	Project (Norse)
0) 10 logit 7.К.К. 0) 10 logit 3.К.К.	- Invironment History
COLD R Source on Save Q. Z . EL . Ration 1	🗣 🗣 Source + 🌫 🔐 🕞 🔐 Import Dataset + 🥑 📃 List + 🌀
1 library(xlsx)	· A clobal microment • Q
2	Data
<pre>3 # FlightDetails.xlsx 4 dfl=read vlsv(file choose() 1 header = T)</pre>	0 df1 107 obs of 13 variables
<pre>5 dfl:dfl, lapply(is.na(dfl), 2, all)] 6 dfl:dfl(apply(is.na(dfl), 1, all)] 7 head(dfl) 8 str(dfl)</pre>	
10 #fb=df1	
11 df15STD=strptime(format(df15STD, "%H:%M:%S"), "%H:%M:%S")	
<pre>12 dflsAtDistrptime(format(dflsAtD, RH:RM:RS), RH:RM:RS) 13 dfl5stA:strptime(format(dflssTA, "SH:SM:SS"), "SH:SM:SS")</pre>	Files Plots Packages Help Viewer
14 dflSATA=strptime(format(dflSATA, "%H:%M:%S"), "%H:%M:%S")	🔅 🔅 🖉 Zoom 🖓 Esport • 🍳 🧉
15 str(df1) 16 bad(df1)	
17	
101 (Top Level) 0	R Script 0
Consele C//Steudon 10/	11 28:00" 15:00" 11 12:00" 13:00" 14:1 15:2 32:30:18

Let us the structure for the data frame. So, we will follow some of the steps that we have gone through in previous lectures as well, so we will just go through them. (Refer Slide Time: 01:34)



Now, there is this particular exercise in the previous lecture we had used a separate grouping for departure time. So, again I have done certain changes into this grouping also, but this is not very important; however, let us run the model with a new grouping for departure time interval. So, this is the range we already familiar with.

(Refer Slide Time: 02:00)

Kinh		
HIP HAT CO	ter vorw Hote service kunst lantag menter lands weige	
Q. d	🚺 🔒 🔛 🖂 👘 Cato file/function	Project (None)
@ 10 lo	рг 7ж х 🔮 10 kupil 3 R х 👘 🗆	Environment History
0 19 20 21 22 23 24 25 26		Control of the second sec
27 28 29 30 31 32 33	<pre>levels(dfloay) levels(dfloay)=("Sunday","Nonday") dfl5FLTIME-as.difftime(as.character(dfl5FLTIME)) ktr(dfl) head(dfl) dfls_tst</pre>	breaks2 2017-08-28 12:00:00 pfrt chr (1:107) "0-12" "0-12" Fés Pob Packages Help Verent
34 35 311 console > brea > brea > brea > df1s > df1s > leve (1] "1 > leve > df1s > leve	dfb1of1 (Tpu[cv0] 5 (/Semine 10/ ⁽²⁾ (/Jrole-28 00.40.00 15T ^{**} 2017-08-28 20.00.00 15T ^{***} ks1sstrptime('00'00'00', "but'994'35') ks1sstrptime('10'00'00', "but'994'35') sife154(df154Tb-breaks2 d'df154Tb-breaks2 d'df154Tb-breaks2 d'df154Tb-breaks2 d'df154Tb-breaks2 d'df154Tb-breaks2 d'df154Tb-breaks1, "0-12", "12-24") cbind(df1, opr) 0E7Tas, factor(df150ry) 0E7Tas, factor(df156ry) 0E7Tas, factor(df156ry) 1s(df150ay) ************************************	

So, break is now 0 and 12, so these are the 2 breaks, 0 hours and 12 hours. Now this is how we are creating department and variable. So, if less than a breaks 2, breaks 2, and

breaks 1, then 0 to 12 so; that means, within if the timing is within these two you know hours 0 hours, and 12 hours.

Then first category that is 0 to 12 otherwise it is going to be the second category that is 12 to 14. So, let us create this variable let us append this to the data frame let us change it to a factor variable day variable as well let us cut the labels, change the labels flight time.

Let us also change it to appropriate format, now this is the structure that we have, now after taking backup we would not like to you know take forward some of the variables so let us get rid of them. Let us look at the structure again now these are the variables, now these are first few values for 6 values you can see everything is ok.

(Refer Slide Time: 03:04)

Kate					-	0 X
Q. • 🞯 • 🕞 🗐 🗁 🌧 (a to file,f	nction				B Project	(None)
0 10 logit 2.8 x 0 10 logit 3.8 x		-0	Environment	History		-0
25 offloer-as, factor(dfloe) 26 dfloer-as, factor(dfloe) 27 levels(dfloey) 28 levels(dfloey) 28 levels(dfloey) 29 dflst.ttmEas.difftime(a 30 31 str(dfl) 32 head(dfl) 33 dfl-dfl.r(1,3,5:8,10,1) 34 dfl-dfl 35 dfl-dfl,r(1,3,5:8,10,1) 35 str(dfl) 39 levels(dfl)flight.status 40 levels(dfl)flight.status 40 levels(dfl)flight.status 40 levels(dfl)flight.status	<pre>></pre>	_ ⊕ Rγ (b) (B) Source (≥	General Contractions of the second se	Import tatket * Import tatket * Import tatket * Import at the import of	C international control of the second	
41 * * * * * * * * * * * * * * * * * * *	variables: levels "Air India","Indigo", levels "Som","DEL","MAA", 11 levels "Gelayed", "Datus", 12 levels "Gelayed", "Datus", 12 me' atomic [1:107] 84 94 79 1 "mins" levels "O-12","12-24": 1 1 1 1	R50913				

Now, let us work on the outcome variable like we have been doing in previous lectures. So, let us change it to numeric code, let us change the reference category, now this is what it becomes now this is ok. (Refer Slide Time: 03:30)



Now, we can move ahead let us do our partitioning 2 partitions, training partition, testing partition and 90 percent for training partition and 10 of observation for test partition.

(Refer Slide Time: 03:40)

Wir Wir Wer Wer Vers weit Ding mee Des wei I to bop 7xx D bop 11x I to bop 11x I to bop 11x I to bop 11x <td< th=""></td<>
Q • • @ • initiation Q • Addime. Project Hance Q • Holey 72.8. Q • Dobget 22.8. Project Hance Project Hance Q • Holey 72.8. Q • Dobget 22.8. Project Hance Project Hance 39 Devels (df115/1)ght.status) Project Hance Project Hance 40 Term (df115/1)ght.status) Project Hance Project Hance 41 head (df115/1)ght.status) Project Hance Project Hance 42 df115/11ght.status) Project Hance Project Hance 44 head (df115/1)ght.status) Project Hance Project Hance 45 FP Arritioning: 10%: 10% Project Hance Project Hance 46 Hand(df115/1)ght.status) Project Hance Project Hance 47 partidx-sample(1.nrow(df1), 0.9*nrow(df1), replace = F) Project Hance Project Hance 48 df1train.mdf1[partidx.] Project Hance Project Hance Project Hance 50 partidx-sample(1.nrow(df1), 0.9*nrow(df1), replace = F) Project Hance Project Hance Project Hance 51 hod3:= pholes.status, family = binomial(1/ink = "logit"), data = df1train) Project + O (Project Hanc
Ibiogr 72%
30 If source name If evels (df1sF1)pit.status) 31 Tevels (df1sF1)pit.status) If evels (df1sF1)pit.status) 32 Tevels (df1sF1)pit.status) If evels (df1sF1)pit.status) 33 Tevels (df1sF1)pit.status) If evels (df1sF1)pit.status) 34 stridf1sF1)pit.status) If evels (df1sF1)pit.status) 35 for fifter status) If evels (df1sF1)pit.status) 36 add(df1sF1)pit.status) If evels (df1sF1)pit.status) 36 fifter ain off1 (apritids.) If evels (df1sF1)pit.status) 37 partidx-sample(1,invm(df1), 0.9*nrow(df1), replace = r) if df1 (afritids.) 30 fod3:=0[n(F1)pit.status, family = binomial(link = "logit"), data = df1train) If evels (df1sF1)pit.status, family = binomial(link = "logit"), data = df1train) 37 foptions (scipen:99) Epoint - 0 If evels (df1sF1)pit.status, family = binomial(link = "logit"), data = df1train)
<pre>39 levels(df1sf1pit.status) 40 levels(df1sf1pit.status)cr(1,0) 42 df1sf1pit.status)relevel(df1sf1pit.status, ref = "0") 43 str(df1sf1pit.status) 44 head(df1sf1pit.status) 45 marticlevels(df1sf1pit.status) 46 martitioning: 90%:10% 47 partitioning: 90%:10% 48 df1train.off1partidx,] 50 lod3:g1m(ff1ght.status - , family = binomial(link = "logit"), data = df1train) 53 soptions(scipen:99)</pre>
<pre>40 Tevels(df1)Fight.status)</pre> 40 Tevels(df1)Fight.status)40 Tevels(df1)Fight.status)40 Data 41 Tevels(df1)Fight.status) 42 df1)Fight.status) 43 str(df1)Fight.status) 44 Tevels(df1)Fight.status) 45 df1)Fight.status) 46 df1)Fight.status) 47 partitioning: 90%;10% 48 df1)Fight.status) 49 df1)Fight.status) 49 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 40 df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , family = binomial(link = "logit"), data = df1)Family 41 Df1)Fight.status - , famil
<pre>c dfl:Flight:Status:relevel(dfl:Flight.status, ref = "0") d dfl:Flight:Status:relevel(dfl:Flight.status) d head(dfl:Flight:status) d head(dfl:Flight:status) d flight:flight:status) d flight:status; d fligh</pre>
43 str(df)fijdt.status) 44 head(df)fijdt.status) 45 head(df)fijdt.status) 46 partidusample(1,nrow(df1), 0.9*nrow(df1), replace = r) 47 partidusample(1,nrow(df1), 0.9*nrow(df1), replace = r) 48 df1train(df)[apridk.] 49 df1text(sk.] 40 df1train(df)[apridk.] 50 hod3:g1m(Flight.status, family = binomial(link = "logit"), data = df1train) 51 podfies(scipens99)
44 head(df1)Flight.status) 55 # Partitioning: 90%:10% 67 partitioning: 90%:10% 68 dfltrain-df1[partidx,] 69 dfltrain-df1[partidx,] 60 dfltrain-df1[partidx,] 61 hed3ed[fliphi] 62 dfltrain-df1[partidx,] 63 hed3ed[fliphi] 64 dfltrain-df1[partidx,] 65 hed3ed[fliphi] 66 hed3ed[fliphi] 67 partidx,] 68 dfltrain-df1[partidx,] 69 hed3ed[fliphi] 60 hed3ed[fliphi] 60 hed3ed[fliphi] 60 hed3ed[fliphi] 61 hed3ed[fliphi] 62 hed3ed[fliphi] 63 hed1ed[fliphi] 64 hed1ed[fliphi] 65 hed1ed[fliphi] 66 hed2ed[fliphi] 67 hed2ed[fliphi] 68 hed2ed[fliphi] 69 hed2ed[fliphi] 60 hed2ed[fliphi] 60 hed2ed[fliphi] 60 </td
<pre>6 # Partitioning: 90%:10% 7 partidx=angle(1.nrow(df1), 0.9*nrow(df1), replace = F) 8 dfltrannoff[partidx,] 9 dfltest-df1[-partidx,] 5 1 hod3:glm(Flight.Status, family = binomial(link = "logit"), data = dfltrain) 5 2 # options(scipens99)</pre>
47 partidx:sample(l.nrow(dfl), 0.9*nrow(dfl), replace = F) 48 dfltrat.mdfl[partidx,] 9 dfltest-dfl[-partidx,] 50 hod3:glm(Flight.status = ., family = binomial(link = "logit"), data = dfltrain) 52 #options(scipens99)
48 dfiran-dfi[partidx,] 49 dfires-dfi-partidx,] 50 hod3-glm(flight.stuus, family = binomial(link = "logit"), data = dfitrain) 52 #options(scipans999)
<pre>50 billesteruir[sertinx]] 50 billestatus, family = binomial(link = "logit"), data = dfltrain) 51 bod3:glm(Flight.status, family = binomial(link = "logit"), data = dfltrain) 53 #options(scipen:999)</pre>
51 hod3:glm(flight.Status, family = binomial(link = "logit"), data = dfltrain) 52 #options(scipan:999)
52 53 #options(scipen=999)
33 #opcions(scipena33)
54 summary(mod3) -
55
511 (lop Level) 8 R Script 8
Console C// Acrolon 10/ 💬
(1) 0 0 1 0 1 0 -
Levels: 1 0
<pre>> dflFfight.status=relevel(dflFfight.status, ref = "0")</pre>
> strong is not in the second se
> head(df15F1ight.status)
(1) 0 0 1 0 1 0
Levels: 0 1
> diradiadi (artick)
> dfltest=dfl[-partidx,]

Now, the same function glm that can be used to again model this. So, these are the results.

(Refer Slide Time: 03:53)

Klube	
Her Half Code sowe Posts weison must contag memor local weig	
🕗 • 📴 • 📄 💮 🧁 (a to file/function 👘 🔯 • Addex •	Indicat (None) •
0] 10 logit 7.8 x 0] 10 logit 3.8 x	Environment History
A G D francester Q Z + E +	
44 head(df1[s]ight_status)	
45	Global Invitorment • G
46 # Partitioning: 90%:10%	dtb 107 obs. of 13 variables
<pre>47 partidx=sample(1:nrow(df1), 0.9*nrow(df1), replace = F)</pre>	odfbl 107 obs. of 14 variables
48 dfltest-dfl[_partidx]	values
50	Obreaks1 2017-08-28
51 mod3=glm(Flight.Status, family = binomial(link = "logit"), data = dfltrain)	Obreaks2 2017-08-28 12:00:00
52	DEPT chr [1:107] "0-12" "0-12"
53 #options(scipen=999)	• mod3 List of 30
54 Summary(nods)	Fire Bale Backage Male Manue
56 # Measures of Goodness of fit	ries rios raciages nep viewer
<pre>57 gf=c(mod3Sdf.residual, mod3Sdeviance,</pre>	👷 🗇 🖉 Zoom 🖓 Export • 👰 🧃 🛞
<pre>58 100 table(dfltrainSFlight.Status)[[1]]/</pre>	
59 Tength(dfltrain)Flight.Status),mod3)iter,	
561 (RepLevel) 8 R Script 8	
Console G//Session 10/ 🕫	
Estimate Std. Error z value Pr(> z)	
(Intercept) -0.843187 0.861801 -0.978 0.32788	
Flight.CarrierIndigo -1.957116 0.747139 -2.619 0.00881 **	
Fight.carrierJet Alrways 0.094405 0.596515 0.156 0.8/46/	
SRCMAA -0.877616 0.708888 -1.238 0.21571	
DESTHYD -1.024817 0.617856 -1.659 0.09718 .	
DESTIXC -0.845506 0.673242 -1.256 0.20916	
FLTIME 0.014289 0.004909 2.911 0.00361 **	
DEPT12-24 1.0//030 1.234354 0.859 0.39034	
Signif_ codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	

Let us look at once again now you would see every time like we have been running this particular model on the same data set and every time. The significance levels have been changing that as I have been explaining smaller data set and therefore, subject to change in terms of as we as the observation that are part of training data set training partition this results will also slightly change and mainly with respect to significance level.

Now, you can see again flight carrier indigo has become significant to star level right and we can also see that destination has also become significant right. And then we can also see flight time has also you know modes you know higher level of significance.

However, more important thing is look at the p varies we can see that this one source for madras as well this is also smaller p values this is anyway significant flight carrier and destination anyway this is significant at ninety percent this is also significant. However, the new grouping that we have created out of department departure time intervals that also not comes out to be significant; however, p value is now smaller.

(Refer Slide Time: 05:14)

Kate	
He has code your that wasna mult listing theme look wap	
😢 🔹 🥶 🖬 🔛 🤮 🕼 Calto Ble Sunction 🔢 💹 🔸 Addies 🔹	B Project (None)
0 10 kopt 2 H x 0 10 kopt 3 R x =	Environment History
() () () () Source on Save () 2 · [] ·	🞯 🖶 🔐 Import Dataset + 🥑 📃 List + 🎯
<pre>44 head(df1flight.status) 45 # Partitioning: 90%:10% 47 partide:sample(limrow(df1), 0.9*nrow(df1), replace = F) 48 df1test.df1[-partidx] 49 df1test.df1[-partidx] 50 51 mod3=glm(Flight.status, family = binomial(link = "logit"), data = df1train) 52 # options(scipens-999) 53 # measures df Goodness of fit 54 measures df Goodness of fit 55 # measures df Goodness of fit 56 # measures df Goodness of fit 57 gftc(mod3)df.residual, mod3/deviane, 58 100*table(df1testinf)Flight.status)[[1]]/</pre>	• drbal inversionert • • • drbal inversionert • • • drbal 107 obs. of 13 variables • • drbal 107 obs. of 14 variables • • breaks1 2017-08-28 • • breaks2 2017-08-28 12:00:00 • • breaks1 1017-08-28 • • breaks1 2017-08-28 12:00:00 • • Breaks1
59 length(dfltrainSFlight.status),mod3Siter, 0 60	
Console G//Session 10/ P	
Flight.CarrierJet Airways 0.094405 0.598515 0.158 0.87467 SACDEL -0.121514 0.615738 -0.197 0.84356 SACMA -0.877616 0.70888 -1.238 0.21571 DESTINO -1.024817 0.617856 -1.659 0.09718 DESTINO -0.044817 0.617856 -1.659 0.09718 DESTINC -0.845506 0.673242 -1.256 0.20916 PLTDM 0.014289 0.004909 2.911 0.00361 ** DETT2-24 1.077556 1.254554 0.659 0.39034 Signif.codes: 0 **** 0.001 *** 0.01 ** 0.05 '.' 0.1 ' ' 1 (Dispersion parameter for binomial family taken to be 1) 1	

So, with this we will discuss the important aspect of logistic regression so that is the majors of goodness of fit.

(Refer Slide Time: 05:19)

Kinder	
Har half Color Youw Pools Sension Havills Carthag Protoc Loois Herp	
🔍 • 🞯 • 🔒 🔛 🖂 A cate file, function	B Project (None) •
0)10/007/00 0)10/00138 x	Environment History
Standarys Standarys	
(○) (Δ] [] IT Source on Save (G. Z. + []] + [] + [] + [] Source + ≥	😭 🔚 📑 Import Dataset * 🧃 🔄 I ist * 🕲
0 51 mad3_ala(flight factor, family - bigamia](ligh - "lagis") data - dflagaig)	* 🕼 Global i nvironment • Q,
51 mod3=gim(Fiight.Status ~ ., Tamiiy = Dinomial(link = logit), data = dritrain)	🔾 dfb 107 obs. of 13 variables 📃 -
St #antians(scinan=999)	O dfb1 107 obs. of 14 variables
54 summary(mod3)	values
55	Character 2017 08 38
56 # Measures of Goodness of fit	O breaks1 2017-08-28
57 gf=c(mod3idf.residual, mod3ideviance,	Obreaks2 2017-08-28 12:00:00
58 100*table(dfltrainSFlight.Status)[[1]]/	DEPT chr [1:107] "0-12" "0-12"
59 length(dfltrainSFlight.Status),mod3Siter,	O mod3 List of 30 -
60 1-(mod3Sdeviance/mod3Snull.deviance))	
61 gf=as.data.frame(gf, optional=T)	Files Plots Packages Help Viewer
63 ""K Success in training data" "#Ttarations used"	
64 "Multiple R-squared")	
65 of	R Litting Generalized Linear Models * Find in Topic
66	olm (state) B Documentation
59.42 (Top Level) 8 R Script 1	gin (auss) in cocomentation (
Console G//Session 10/ P	Fitting Generalized Linear
Elight Carrientet Airmans 0.094405 0.598515 0.158 0.87467	
SRCDEI -0.121514 0.615738 -0.197 0.84356	Models
SECMAA =0.877616 0.708888 =1.238 0.21571	
DESTRYD -1.024817 0.617856 -1.659 0.09718	Description
DESTIXC -0.845506 0.673242 -1.256 0.20916	Description
FLTIME 0.014289 0.004909 2.911 0.00361 **	
DEPT12-24 1.077656 1.254554 0.859 0.39034	gin is used to it generalized linear models, specified by
•••	a description of the error distribution
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	a description of the error distribution.
(Dispersion parameter for binomial family taken to be 1)	Usage
	1
and a second	

So, just like the linear regression multiple linear regression logistic regression also is a statistical technique primarily. And therefore, in a statistical modelling the main objective is to fit to data as we have talked about this aspect many times before. So, therefore, in multiple linear regression we have a particular metric called a multiple R square and adjusted R square which are used to assess how good the model is fitting to the data.

So, similarly since logistic regression is also a statistical technique so there also we are required to have matrix for good to measure the goodness of it how well model is fitting to the data. So, what are those matrix so because there is certain key differences in logistic regression and linear regression. So, we will talk about some of those matrix now.

So, you can see in the code that I have created a vector here gf first one is mod 3 the model that we have just computed and you can see degree of residual degree of freedom. So, this df dot residual is one of the value that is returned by and the glm function and gives us the residual degree of freedom then we have deviance.

So, this is again the returns the deviance value returned by the glm function right then few other things which are mainly for the descriptive purposes. For example, this table result which is for the outcome variable here and then divided by the full number observation that will give us percentage success in training data and then we have iterations.

So, as we talked about the particular estimation technique that is used in that is used in logistic regression is different from multiple linear regression we can look at for more details we can look at here. So, we talked about that Emily is used for typically for used for logistic regression; however, you can see that for example, we have been using glm function, and within that if we go we look up for the some of the arguments.

(Refer Slide Time: 07:48)

MT (AV WW RVN VISIAN Raid 12hug RVNV IAVIS MID	1-100
• 🞯 • 🕞 😥 🖂 🔿 Calts flatturition 🔡 • 🛛 Addex •	Noject (Nor
2) 10 log# 2 И ж 2) 10 log# 3 Я н = =	Environment History
Sol 21 R T Source on Save Q / - E - Brun 59 De Source -	2 🔐 🔒 🔐 Import 134aset + 🧹 📃 1 ist +
50	· Cabal I missement • Q
51 mod3=glm(Flight.Status ~ ., family = binomial(link = "logit"), data = dfltrain)	O dfb 107 obs. of 13 variables
53 Fontions(scipen=999)	O dfb1 107 obs. of 14 variables
54 summary(mod3)	values
55	0 breaks1 2017-08-28
56 # Measures of Goodness of fit	obreaks2 2017-08-28 12:00:00
58 100°table(df1train)Flight, Status)[[1]]/	DEPT chr [1:107] "0-12" "0-12"
59 length(dfltrainSFlight.Status),mod3Siter,	0 mod3 List of 30
60 1-(mod3Sdeviance/mod3Snull.deviance))	and the second second second
62 rownames(of)=r("Residual df" "Std Dev Estimate"	Files Plots Packages Help Viewer
63 "% success in training data", "#Iterations used",	🖕 🤿 🏠 🚍 🖉 🔍 q. g/m 🛛 🔍
64 "Multiple R-squared")	It Litting Generalized Linear Models + Find in Topic
65 gf	
29.42 (Replayed) # R Scr	upio na.action, start NULL, etastart, t
	control = list(), model = TRUE, r
onsole 6//Session 10/ 💬	A FALSE, Y TRUE, CONTRASTS NUL
light.CarrierJet Airways 0.094405 0.598515 0.158 0.87467	glm.fit(x, y, weights rep(1, nobs),
RCDEL -U.121514 U.615/38 -U.19/ U.84356 RCMAA -D.877616 0.708888 -1.238 0.21571	start = NULL, clustart = NULL, a
ESTRYD -1.024817 0.617856 -1.659 0.09718	offset rep(0, nobs), family
ESTIXC -0.845506 0.673242 -1.256 0.20916	control - anothy incorpt - in
TIME 0.014289 0.004909 2.911 0.00361 **	## S3 method for class 'gim'
LP112-24 1.0//030 1.234334 0.839 0.39034	weights(object, type c("prior", "Work
ignif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	Arguments
Dispansion parameter for binomial family taken to be 1)	
pispersion parameter for binomial family taken to be 1)	formula an object of class "formula" (or one
0.1.2 0.1.2 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1	

Specifically for this purpose the estimation technique purpose will get more detail. So, we will see that glm for dot fit.

(Refer Slide Time: 08:00)

r HM CAAP York Velant Kult Ishug Petror lank wep	
🛛 • 🞯 • 📄 🗐 🔮 🏘 (La to file/function	Project (Nor
0 10 log# 7.8 × 0 10 log# 3.8 ×	Environment History
Sol 2 Run 59 Source on Save Q 2 . E . Brun 59 De Source . 2	🞯 🔒 📑 import i Jataset + 🧃 📃 i ist +
<pre>50 51 mod3=glm(Flight.Status, family = binomial(link = "logit"), data = dfltrain) 53 grotions(scipane999) 54 summary(mod3) 55 6# Measures of Goodness of fit 57 gf=c(mod3)df.residual, mod3ideviance, 51 100'table(dfltrainf)ight.Status)(l1]// 59 length(dfltrainf)ight.Status)(l1]// 50 length(dfltrainf)ight.Status)(l1]// 51 length(dfltrainf)ight.Status)(l1]// 52 gf=as data.frame(gf.optional=1) 52 gf=as data.frame(gf.optional=1) 53 gf=as data.frame(gf.optional=1) 54 monames(gf)=("Residual df", "Std. Dev. Estimate", 55 monames(gf)=("Residual df", "Std. Dev. Estimate", 55 monames(gf)=("Residual df", "Std. Dev. Estimate", 56 monames(gf)=("Residual df", "Std. Dev. Estimate", 57 monames(gf)=("Residual df", "Std. Dev. Estimate", 57 monames(gf)=("Residual df", "Std. Dev. Estimate", 58 monames(gf)=("Residual df", "Std. Dev. Estimate", 59 monames(gf)=("Residual df", "Std. Dev. Estimate", 59 monames(gf)=("Residual df", "Std. Dev. Estimate", 51 monames(gf)=("Residual df", "Std. Dev. Estimate", "Residual df", "Std. Dev. Estimate", "Residu</pre>	a. Jabai Insurement - Q Q dfb 107 obs. of 13 variables Q dfb1 107 obs. of 14 variables Q breaks1 2017-08-28 Q breaks1 2017-08-28 12:200:00 DEFT cref [1:107] "0-12" "0-12" 0 mod3 List of 30 Here Pols Prockages Help Q g/m Q g/m
65 gf 2742 Obplaced B RSarget B Camaba C2//Shadan 10/ P ⁰ m□ Robit C2//Shadan 10/ P ⁰ 0.051238 Robit C2//Shadan 10/ P ⁰ 0.09128. Robit C2//Shadan 10/ P ⁰ 0.09128. Robit C2//Shadan 10/ P ⁰ 0.0015128 Robit C2//Shadan 10/ P ⁰ 0.011111 Dispersion parameter for binonial family taken to be 1) m 0	uses teratively reweighted least squares (WLS) the attentive "model". Li cume: "Future Stems the model frame and does no fitting. Use-supplied fitting functions can be supplied either as a function or a character string tamong a function, with a function which lakes the same arguments as n the "Li is looked up from within the stats namespace. x. y For julk logical values indicating within the stats namespace.

We will see that glm dot fit method. So, this we can see iteratively re weighted least square is used the particular function that we are using glm iteratively re weighted least square function is used which is quite similar in approach with respect to Emily estimation technique that we talked about.

(Refer Slide Time: 08:41)



So, MLM techniques, sorry MLM maximum likelihood method that we talked about in our discussion as we can look in these slides as well maximum likelihood method MLM method that we talked about. So, this is quite similar to what this the discussion that we have.

So, in particular R's implementation glm function that we are using it is the iteratively re weighted least square that estimation technique that is used to estimate the coefficients that we have been doing quite similar to a MLM and as we talked about that number of iterations have to be performed to reach to estimate these parameters so that we can get the best model which is the model best model which is fitting the data.

So, number of iterations actually indicate that and then we have one matrix which is quite similar to what we have multiple R square and linear regression right. So, this is actually computed where using the deviance value. So, null deviance and the standard deviation estimate or deviance that we have can look at the returned values here.

So, you can see deviance is the one of the return value so this is one and then we also null deviance is also a return. So, this is null deviance is come with respect to the Knave rule. So, 1 minus this deviance divided by null deviance gives us a value which is quite similar to what we have in multiple linear regression multiple R square.

So, this particular value will be will give us a metric which can be used to understand the goodness of fit follows degradation model and as on its own deviance also can be used it is quite similar to what we have there in as I see some of squares error.

So, this is quite similar deviance is quite similar to that and then we can have one metric as I talked about similar to multiple R square. So, these metrics can be used to assess; the assess the fitness or model goodness of fitness goodness of fit of alloys degradation model.

So, let us compute some of these things so, residual degree of freedom deviance which is similar to SSC in linear regression, then we have this proportions percentage in successive training data. Then we have number of iteration and then we have a multiple R square kind of metric let us compute this.

(Refer Slide Time: 11:34)

Date	-0-
HAT LODE YOR SHOULD REAL DATE THE HAS HER	
Addex •	Project (None)
() 10 kg#78 x () 10 kg#38 x	Environment History
S S D Source on Save Q Z + E + Brun 59 D Source + 2	🞯 🔒 🔐 Import Dataset + 🧃 📃 List + 🔘
<pre>St summary(mod3) # Measures of Goodness of fit for gfrc(mod3)df.residual, mod3ideviance, for gfrc(mod3)df.residual, mod3ideviance, for gfras.data.fram(gfr, optional=1) for grownames(gfr, optional=1) for grownames(gfr, optional=1) for grownames(gfr, optional=1) for gfr = for enduction in deviance (g-value from chi-sq distr.) for pchisq(mod3imull.deviance.mod3ideviance, length(mod3iscoefficients)=1, lower tail for a classify observations using a cutoff value of 0.5 for the classify observations using a cutoff value of 0.5 for th</pre>	dabal namonent • • d fb1 10/ 0bs. of 13 variables d fb1 10/ obs. of 14 variables of db1 10/ obs. of 14 variables of db1 10/ obs. of 14 variables of brasks1 2017-08-28 obreaks2 2017-08-28 obreaks2 2017-08-28 obreaks2 2017-08-28 miles filted • obreaks2 00/0-08-28 obreaks2 00/0-08-28 obreaks2 00/0-08-28 obreaks2 00/0-08 per chr (1:107) familes filted by quasi- likelhood the value is in.
Console C(/Neusion 10) ←	Inclinite Viewer The deviance The deviance The deviance for the null model comparable with durations. The null model will indicate the offset, and an intercept if there is one in the model. Note that its will be incorrect if the link function depends on the data other than through the find mean. specify a zero offset to force a correct calculation. It ev the number of iterations of MVLS used

Let us create a data frame and row names we have given some and these are the values. So, you can see residual df df is 87 here, so as you can see let us again have a look df 1 training partition we have 96 observations.

(Refer Slide Time: 11:56)

(Sale	
e Hat Lode wew most sesses must comp mote tools wep	
🛛 • 📴 • 🕞 🗐 🗁 🖗 Calto Bladurdice 🔤 • Addies •	Project (Nors)
0] 10 log# 7.8 x 0] 10 log# 3.8 x = 0	Environment History
0 2 B T Source on Save Q Z + E + BRun B D Source + 3	🛿 🔐 🕞 🔐 Import 13staset + 🧃 🗐 Ilitt + 🎯
54 summary(mod3) 55 56 # Measures of Goodenss of fit 7 gfrc(mod3)df residual, mod3ideviance, 7 gfrc(mod3)df residual, mod3ideviance, 7 genth(dftrainFilght, status)[1]]/ 9 = length(dftrainFilght, status), mod3itter, 60 = 1-(mod3)deviance/mod3imul] deviance)) 10 = gfra.st.ar.frame(gf, optional=7) 20 = rownames(gf)	dfl 107 obs. of 6 variables dfl 107 obs. of 6 variables dfltest 11 obs. of 6 variables dfltest 96 obs. of 6 variables dfb 107 obs. of 13 variables dfb 107 obs. of 14 variables off 107 obs. of 14 variables
of f Multiple R-squared") f pr Check significance for reduction in deviance (p-value from chi-sq distr.) f pchisq(mod3snull.deviance-mod3ideviance, length(mod3icoefficients)-1, lower.tail f classify observations using a cutoff value of 0.5 d classify observations using a cut	File Pots Packages Help Verent
Cansele Gr/Areadon 10/ ∞ Coefficients: Estimate Std. Error z value Pr(>[2]) (Intercept) -0.84187 0.861801 -0.978 0.32288 Flight.CarrierIndigo -1.95718 0.74718 -2.619 0.00881 ** Flight.CarrierIndigo -1.95718 0.74718 -2.619 0.00881 ** Flight.CarrierIndigo -1.95718 0.74718 -2.619 0.00881 ** Flight.CarrierIndigo -1.95718 0.7456 SROWA -0.877616 0.708888 -1.238 0.21571 DESTMC -0.4845506 0.673342 -1.236 0.20918 DESTMC -0.4845506 0.673422 -1.256 0.0918 FLITME 0.014289 0.00490 2.911 0.00516 ** DEFTL224 1.077656 1.254554 0.859 0.39034	null deviance for the null model, comparable with deviance. The null model will include the offset, and an intercept if there is one in the model. Note that his will be incorrect faith offset and one offset offset and other the intercept if the init function depends on the data other that https://www.initercept.com/ the initercept.com/ the initercept.com/ the initercept.com/ the initercept.com/ the initercept.com/ the initercept.com/ the initercept.com/ WLS used.

And if we go back to our summary results right if we go back to our summary results we can see that how many variables we have here 1, 2, 3, 4, 5, 6, 7, 8; 8 8 variables. And we can see that 87 is the residual.

(Refer Slide Time: 12:06)



So, 87 plus 7 that makes it 94 that is n minus 1, so that is the computation that is how the degrees of freedom have been computed, so this is a correct value here. And then we have deviance value which is also called as standard deviation estimated by some software's some statistical commercial statistical software's.

And then so this is all similar to what we have SSE sum of squares error in multiple linear regression then we have number of iteration that have been used to arrive at the particular model and that we have that is not 3 in this case. Then we have a value similar to multiple are squares we can see 20 percent, of the variability in the outcome variable has been explained by this model. So, all we talk about that this is being computed by 1 minus mod 3 deviance, divided by null deviance.

Now, whether so in terms of on further in terms of deviance the null deviance represents the knave rule value. So, we have to see how much our model has been able to how much reduction in deviance has been done by our model and whether that that is significant or not so that can also be that can also be performed using this chi square test that we can do.

So, we have one function p chi square. So, there we can actually use these two values or we can take a difference of null deviance and deviance so that would be the reduction in deviance from a knave rule and you know how much deduction that our model has done. And we can look at the number of predictors as degrees of freedom I could be used the number predictors that we use have used could be used degrees of freedom because these are freedom degree of freedom that had been used to reduce the deviance as we talked about from 95 minus 1 available degrees of freedom two we have reduced up to 87 that is residual degree of freedom so 7 predictors have been used.

So, that information can be use to perform chi square test and to find out the significance of whether the reduction has been significant or not. So, the third argument is lowered tail that is specified as false and we can compute this chi square value so we can see that this is a small value. So, therefore, it seems that this redis reduction is reduction deviance is significant which are also clear by the difference between the deviance and null deviance values.

So, we can also compute that we can see this is null deviance; this is null deviance. So, let us look at the value and we can have the deviance. Let us look at the value so you can see.

(Refer Slide Time: 15:17)

(Krabe)		-	0
HIR HAR LOAD YOW PURE SESSOR KURE LINEAR PHONE LANS HUTP			
🝳 • 🞯 • 🔛 🔛 🗁 🕐 Co to file/function		I Project	t (Nonc)
• 10 log# 7.8 ж • • 10 log# 3.8 ж	Environment	History	-0
Source on Save Q Z . E . BRun De Dource . 2	C 8 8	nport Dataset * 🧃 📃 Lie	. @
<pre>57 gf-c(mod3idf.residual, mod3ideviance, 100'table(df1rtain)Filght.status)[[1]]/ 9 length(df1rtain)Filght.status)[[1]]/ 10 gf-as.dcat.frame(gf, optional=1) 10 gf-as.dcat.frame(gf, optional=1) 10 gf-as.dcat.frame(gf, optional=1) 10 gf-as.dcat.frame(gf, optional=1) 11 gf-as.dcat.frame(gf) 11 gf-as.dcat.frame(gf) 12 rommames(gf)=c("Residual df", "std. Dev. Estimate", 13 "% Success in training data", "#Iterations used", 14 "Multiple R-squared") 15 gf 16 gf c(beck significance for reduction in deviance (p-value from chi-sq distr.) 17 pchisq(mod3inul1.deviance-mod3ideviance, length(mod3icoefficients)-1, lower.tail 16 gf classify observations using a cutoff value of 0.5 17 mod3trainc:ifalse(mod3ifited values>0.5, 1, 0) 17 mod3trainc:ifalse(mod3if</pre>	G (Jobbil Invice G (Jobbil In	107 obs. of 6 variables 11 obs. of 6 variables 10 obs. of 6 variables 96 obs. of 1 variables 107 obs. of 13 variables 107 obs. of 14 variable Packages Hete Vener Q. gm	
72 table("Actual value"-dfltrain5Flight.status, "Predicted value"-mod3trainc) - 73 ************************************	null.devia	families fitted by quasi- likelihood the value is in in. The deviance for the null model, comparable with deviance. The null model include the offset, and an intercept if there is one in model. Note that this will be incorrect if the link function depends on the data often than through the fitted me specify a zero offset to fo correct calculation.	el will the e n r an: rce a
> md336ev1ance [1] 104.2421 	iter	the number of iterations of IWLS used.	1

So, the difference is so there is good enough Reduction and deviance and that is why it also came as significant you know difference. So, these are some of the matrix that can actually be used to understand to a measure the goodness of it to asses goodness of it for a model as we talked about I talked about that in a statistical set setting. So, these are some of the matrix which would be more useful. So, in a statistical modelling we stopped at when we build the model using the training partition. So, typically all the observations are used that are present and then we look asses the model with respect to some of these some of these matrix. Now, let us move forward to our next discussion point in logistic regression so that is let us move forward so that is this particular point whether linear regression can be used for a categorical outcome variable right.

So, there are there are some situations where linear regression can be used as a categorical outcome variable which we will discuss later; however, right now we are discussing some of the more important points with respect to overall general applicability of linear regression for a categorical outcome variable.

So, can be done technically it can be applied so we can treat the outcome variable as continuous. So, the categorical outcome variable can be treated as continuous variable. So, we can essentially do the numeric coding and keep it as a numeric variable and technically we can apply we will get the results; however, the results would be meaningful or not that we need to understand.

So, technically it can be applied we can read the categorical outcome variable as continuous variable we can code it numerically so that can be done. However, there are going to be anomalies that would lead to spurious modelling so what could be some of these things.

So, number one predict predictions can take any value not just any values so for example, that binary logistic regression model that we have been performing for on some of the data sets so they are the our outcome variable it is it typically has two classes class 1, and class 0.

And so, the values the remaining variable we will take is to 0 for class 1 and 0 and 1 for class 1; however, when we apply a linear regression model to a categorical outcome variable the prediction can take any values any real value can be taken and not just the dummy values 0 and 1 so that is one challenge.

How do we map some of these predators values which can which can be any real value to the actual values of the outcome variable 0 and 1. Now, outcome variable or residuals do not follow normal distribution. So, as we have discussed during a linear regression that this is one of the important assumption that dependent variable that is outcome variable all residuals should follow normal distribution, but that is not the case as we can understand that categorical outcome variable will have just 2 values 0 and 1.

So, it is actually it actually follows a binomial distribution so this particular assumption would also be violated; however, we talked about that because we are in data mining modelling context. So, where as we talked about even if for you know prediction purposes even if this particular you know assumption is violated in terms of prediction it might not much matter much because generally check performance on validation partition and test partitions; however, this case is different.

The deviation from normal distribution is much higher it is actually different distribution binomial distribution so that is one problem. Now, the another assumptions that we talked about in multiple linear regression is homoscedasticity; however, if we apply linear regression to an a to a categorical outcome variable this particular assumption would also be violated. The variance of outcome variable that we that we expect to be constant across all the records that is the homoscedasticity property so we for to apply multiple interrogation we want our outcome variable to follow this to have this property.

So, variance should be constant across all time; however, if we look at the variance for our categorical outcome variable it is going to be this particular value n times p into 1 minus p and as you can see because this is dependent on the value of p. So, therefore, the variance will change as the value of p changes.

So, when the value of a probability value is actually close to 0, then the variance would be on the as lower side and when the value of p is approaches 1 then the variance would be on the higher side. So, therefore, the variance will be will not be constant and it will be it will increase as the probability value you know in case is from 0 to 1.

So, so some of these are some of the problems that we can see in that we can directly understand and why a linear regression and cannot be applied to category outcome variable in a general sense and the problems that we can see here. So, what we will do we will do an exercise in R to understand the same thing to understand this particular aspect. So, what we will do we will apply a linear regression model on a logistic partition and see the see its applicability and see some of the anomalies or violations that are that could be there. So, for this purpose we are going to use as you can see multiple here the comment is multiple linear regression model for a categorical response. So, promotion offers is the data set that we are going to use for this particular exercise. So, let us import the data set let us edit load.

(Refer Slide Time: 22:16)

States MI CAN	P Yow York Vision Kurk Dobug Polity Lance was		-	0
•	• 🔒 🚳 🖂 🕐 Controllingtion 👘 🛛 🔯 • 🕴 Addees •		🖲 Project	t (Non
10 log	# 7.8 x @ 10 logil 3.8 x	Environment	History	
130 131 132 133 134 135 136 137 138 139 140 141 142 143 144	<pre>Sourcensaw Q_Z +</pre>	Charles to the second s	mportazior • 107 obs. of 6 variables 11 obs. of 6 variables 10 obs. of 10 variables 107 obs. of 13 variables 107 obs. of 14 variables 107 obs. of 14 variables Pedage Holy Veeer Pedage (g dm (g dm (g)))	
insole insole td. De Succe iterat altipl pchis 1] 0.0 mod3 1] 131 mod3 1] 104 df2=	Atriat/ (Dpicwi): C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 1: C/Jestain 2: C/Jestain	null.devia	families filted by quasi- likelihood the value is NA- model, comparable with deviance for the null include the offset, and an intercept if there is one in model. Note that this will be incorrect if the link function depends on the data other than through the filted me specify azer offset to be correct actuation. the number of farations o VVC sued	al will the e n r an: rce a

So, once the observations have been loaded into environment we as we will see in the environment section we will do some of these steps I think it is has been loaded yes. So, df 2 we can see 5000 observation, 9 variables so let us remove any columns, or any rows if there are any let us look at the structure so this is the data set. So, we are already familiar with this.

(Refer Slide Time: 22:41)



So, let us go through some of the some of the competitions some of the transformations that we have been doing you know in previous lectures as well take a backup let us get rid of let us select just these 4 variables that is income promotional offer of our outcome variable then we have I think family size and online activity. So, let us select them. So, these are the variable selected for this particular exercise income promotion offer and family size and online.

(Refer Slide Time: 23:12)

Kale	
Mit care www. Ports session multi-landing methic cares werp	
💽 🔹 😭 🗐 🗁 🛛 🚁 Cante fürstundion 👘 🖾 🔹 Addies. •	S Project (N
0] 10 kogit 2.8 x 0] 10 kogit 3.8 x =	Environment History
Sol 21 R T Source on Save Q 7 . E . Bun 59 De Source - 3	🗧 🔐 🕞 🔐 import 13staset + 🥑 📃 1 ist +
144 df2=df2[, c(1,3,7,9)]	· Chalai (minoment • Q.
145 str(df2)	O of 5 obs. of 1 variable
146 #df25Promoffer=as.factor(df25Promoffer)	values
147 #df25dmine=as.tactor(df25dmine) 148 #str(df2)	Obreaks1 2017-08-28
149	Obreaks1 2017-08-28 13:00:00
<pre>150 # Partitioning: Tr:Te->60%:40%</pre>	of peaks2 2017-00-20 12:00:00
<pre>151 partidx=sample(1:nrow(df2), 0.6*nrow(df2), replace = F)</pre>	0 mod2 ist of 20
152 df2train=df2(partidx,)	O modd Lange la (12 alamante 201
154 mod4=lm(Promoffer, df2train)	Umod4 Large Im (12 elements, 801
155	Files Plots Packages Help Viewer
156 #op=options()	
157 #options(scipen=999) 158 J	AAR BO (IM A)
159 mod4summ=summary(mod4); mod4summ	ic riming Generalized Linear Models * Find in Topic
160	x likelihood the value is no
1581 (Replayed) 8 R Script	
for the deside that (2)	null.deviance The deviance for the null
AFb2=dF2	model, comparable with
df2=df2[, c(1,3,7,9)]	include the offset, and an
str(df2)	intercept if there is one in the
data.frame': 5000 obs. of 4 variables:	model. Note that this will be
S Income : num 49 35 10 101 45 31 71 23 80 182	incorrect if the link function
Fromorrer : num 00000001	than through the fitted mean
Sonline : num 0000011010	specify a zero offset to force
<pre>partidx=sample(1:nrow(df2), 0.6*nrow(df2), replace = F)</pre>	correct calculation.
df2train=df2[partidx,]	then umber of iterations of
mod4=Im(Promotter ~ ., df2train)	IWLS used.
	· ·

So, the variable that we are going to use for our outcome variable is going to be the promotional offer as you can see that we have commented out the lines of code, which we used earlier to convert these numeric variables going to convert numeric variables into the categorical variable. So, promotional offer an online so they are actually categorical variable factor variable, but we are not converting them into factor variable because we are going to apply a linear regression modelling. So, we will give them as numeric and we will apply.

So, a partitioning is the same 60 percent, 40 percent in this case so we can see that let us do partitioning. So, df 2 train you can see 3000 observations, 4 variables. Now the same align function is going to be used now the promotional offer is going to be request against all the predictors that are present in this particular dataset df 2 train. So, let us run this.

(Refer Slide Time: 24:30)

Kinde		- 0
e has code view more senses word tomag methor toos weep		
၊ 🞯 • 🕞 😭 🦀 🕼 to to file/function 👘 🔯 • 🛛 Addient •		Project (Non
0 10 logit 2 H x 0 10 logit 3 R x	- C Environment	History
5 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	Source . 2 C B	Import Dataset + 🧹 📃 List +
148 #str(df2)	· Cabal Inv	inneret - Q
149	values	
150 # Partitioning: Tr:Te->60%:40%	O breaks1	2017-08-28
151 partick=sample(1:hrow(dr2), 0.0-hrow(dr2), replace = F) 152 df2train=df2[partidx]	hreaks?	2017-08-28 12:00:00
153	OFET	che [1:107] "0-12" "0-12"
154 mod4=lm(Promoffer ~ ., df2train)	0 mod3	List of 20
155	0 mous	List of 50
156 #op=options() 157 #op=options()	0 mod4	Large Im (12 elements, 801
157 #options(scipen=999) 158	0 mod4sum	n List of 11
159 mod4summ=summary(mod4); mod4summ	Files Plots	Packages Help Viewer
160 mod4ava=anova(mod4); mod4ava		
161	6 9 X	
162 DF=C(mod4summ)fstatistic[[numor]], 163 mod4summ(fstatistic[["deodf"]]]	R: Litting Gener	alized Linear Models • Find in Topic
164		families fitted by quasi-
1601 (Rep Level) 0	R Script 0	likelihood the value is NA.
Console 6// Sension 10/ 🕫 📼		ance The deviance for the null
		model, comparable with
Min 10 Median 30 Max	1	include the offset and an
0.381/2 -0.13/44 -0.0303/ 0.03823 0.32044		intercept if there is one in the
oefficients:		model. Note that this will be
Estimate Std. Error t value Pr(> t)		incorrect if the link function
Intercept) -0.2388135 0.0148069 -16.129 <2e-16 🎬*		depends on the data other
ncome 0.0033865 0.0001005 33.706 <2e=16 ***		than through the fitted mean.
1119.5126 0.03/2885 0.0040443 9.220 <2e-16 ***		correct calculation
-0.0030494 0.0033213 -0.027 0.33		
ignif. codes: 0 '===' 0.001 '==' 0.01 '=' 0.05 '.' 0.1 ' ' 1	iter	the number of iterations of
		INCO USED.
the set and the set		

Now, what we will do we will look at the summary table. Let us look at the results. So, as we can see that income is intercepted significant income is significant and family size is significant we can look at the different estimates for example, income quite a small value from this family size is also 0.03, online is not significant it has not only found to significant as we can see.

(Refer Slide Time: 24:55)



And we can have a look at the other values we can see adjusted R square R square and multiple R square. So, we can see if we you know 27 percent multiple R square value we look at the p well it is quite as small. So, the model is significant ah; however, as we talked about certain problems as we have discussed could be there certain anomalies could be there.

(Refer Slide Time: 25:21)

Kodu	-0
e hat case sow not sessa nuit isolag neme laos wap	
🕽 • 🤠 • 🔒 🔝 🔮 🏘 (alto Blafunction 🔤 🛯 🔯 • Addies •	B Project (No
0 10 logir 7.8 × 0 10 logil 3.R ×	Environment History
O D B T Source on Save Q Z • E • → Run 59 → Source • ≥	🔐 🕞 🖙 Import Dataset + 🥑 📃 Liet +
158	Global I mirconment Q,
159 mod4summ=summary(mod4); mod4summ	O mod4ava 4 obs. of 5 variables
160 mod4ava=anova(mod4); mod4ava	values
162 DF=c(mod4summifstatistic[["numdf"]].	O breaks1 2017-08-28
<pre>163 mod4summSfstatistic[["dendf"]],</pre>	O breaks2 2017-08-28 12:00:00
<pre>164 mod4summSfstatistic[["numdf"]]+mod4summSfstatistic[["dendf"]]) 165</pre>	DEPT chr [1:107] "0-12" "0-12"
103 166 SS=c(sum(head(mod4ava[,"Sum So"],=1)),	O mod3 List of 30
167 mod4ava["Residuals","sum sq"],	0 mod4 Large lm (12 elements, 801
168 sum(mod4ava[,"sum sq"]))	Files Plots Packages Help Viewer
109 170 MS:c(mean(head(moddaya["Mean So"] -1))	
171 mod4ava["Residuals","Mean Sq"],"")	
172	R Litting Generalized Linear Models * Find in Tupic
1/3 Estatistic or (mod4summifstatistic) ["value"]]	families fitted by quasi-
16218 (Rop Level) 0 R Script 8	
Family (1) Decision 10/ D	null.deviance The deviance for the null
mod4ava=anova(mod4); mod4ava	deviance. The null model will
malysis of Variance Table	include the offset, and an
	intercept if there is one in the
response: Promotter	incorrect if the link function
Income 1 67.073 67.073 1064.8362 <2e-16 ***	depends on the data other
Family.size 1 5.345 5.345 84.8536 <2e-16 ***	than through the fitted mean:
Dn11ne 1 0.025 0.025 0.3938 0.5304	specify a zero offset to force a
112 100013 1370 1001110 01003	Correct Calculation.
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1	iter the number of iterations of
	* *

So, what we will do? We will look at look at computing some of those things to check whether those anomalies are there in this particular case. So, what we will do we are running and over to extract some of the parameters here. So, we can see so sum of square errors mean of square errors and F value of a statistic and the probability values for these predictors income family size is given there in the ANOVA table.

So, what we are going to do is we are going to compute these particular values in a in a format that can be used for tabular presentation later on. So, mod 4 so on; we have a F statistics value. So, that is written as part of the summary function of models.

(Refer Slide Time: 26:19)



So, when we apply summery on the model object so you can see this is nothing, but F statistic value for the model and then we have degree of freedom and here. So, we can see the degree of freedom residuals degree of freedom and the residuals degree of freedom here so that is going to be stored in this so in this data frame.

First we have the regression degree of freedom, then we have the residual degree of freedom and then we have the total so this data frame is about degree of freedom as is sustained by its named DF. Then we will compute some other square error so first we can see that in the ANOVA table, from the ANOVA table we are trying to extract out this these values.

(Refer Slide Time: 27:21)

Kinder	
e hat code were their sesson must clean preter seen map	
🕐 😋 • 🔛 🔝 🖂 🍙 Cata Bir/function 👘 🔯 • Addies •	B Project (Non)
	I manual street
0 10 log# 7.8 x 0 10 logit 3.8 x	Environment History
♦ Ø B Source on Save Q Z • E • B Run D D Source • 2	🞯 🔒 📑 Import Dataset + 🧃 📃 List + 🗍
<pre>162 DF=c(mod4summifstatistic[["numdf"]]),</pre>	Gobal i nvironment • Q
<pre>163 mod4summSfstatistic[["dendf"]],</pre>	Obseaks2 2017-08-28 12:00:00
<pre>164 mod4summSfstatistic[["numdf"]]+mod4summSfstatistic[["dendf"]])</pre>	Obreaks2 2017-08-28 12:00:00
165	DEPT Chr [1:107] 0-12 0-12
100 SS=C(Sum(head(mod4ava[, Sum Sq"],-1)), 167 mod4ava["Besiduals" "Sum Sq"]	DF num [1:3] 3 2996 2999
168 sum(moddava[residuals , sum sq],	Fstatistic chr [1:3] "383.36120294222
160 sum(mou+ava[, sum sq]/)	Q mod3 List of 30
170 MS=c(mean(head(mod4ava[,"Mean So"],-1)),	O mod4 Large lm (12 elements, 801.
171 mod4ava["Residuals", "Mean Sq"], "")	O mod4summ List of 11
172	
<pre>173 Fstatistic=c(mod4summSfstatistic[["value"]],"","")</pre>	Files Plots Packages Help Viewer
<pre>174 P=pf(mod4summSfstatistic[[1]], mod4summSfstatistic[[2]], mod4summSfstatistic[[3]]</pre>	
175 lower.tail = F)	
176 pvalue=c(P, ,)	R: The EDittribution * Find in Topic
177	
17421 (RepLayed) # R Second #	PDist (stats) R Documentation
Console G://Session 10/ 🛱	The F Distribution
value numdf dendf	. The Plothoutern
383.3612 3.0000 2996.0000	
<pre>DF=c(mod4summ\$fstatistic[["numdf"]],</pre>	Description
<pre>mod4summSfstatistic[["dendf"]],</pre>	
<pre>mod4summSfstatistic[["numdf"]]+mod4summSfstatistic[["dendf"]])</pre>	Density, distribution function, quantile function and
<pre>SS=c(sum(head(mod4ava[, "Sum sq"],-1)),</pre>	random generation for the F distribution with dF1 and
modyaval Kestouals', Sum Sq"],	carameter nam)
Sum(moovava(, Sum Sq))) MSac(man(haad(moddava("Maan So") -1))	parameter mapp
moddava["Baciduals" "Maan Co"] "")	lleane
Estatisticsc(mod4sumSfstatistic[["value"]],"","")	e suge
	dity dil di7 non in = Faters
1	

So, first some other square for the you know regressors regression and then for the residuals, and then total. So, this would be recorded in this particular variable essence. And then we look at the mean or mean of square errors, so that is also being extracted from the ANOVA table results, you can see this particular column mean square and these values are being instructed.

So, first 3 values so this head function as you can see now this role of this head function is quite different and both these computations we have used head function.

So, you can see the second argument is minus 1 so what it actually does is it gives us all the values except the last value in the vector. So, for example, mean square or some other square. So, there are 4 values and these 2 columns, and these 2 vectors.

So, except the last values that is corresponding to residuals the first 3 values are going to be written; that means, last value will be left out and the remaining and remaining n minus 1 values are going to be written. So, that is what we want. So, that will give us the corresponding sum of squares or mean sum of squares are for the deviation.

So, first that then residuals so let us compute this. Then let us also extract the F statistics, from the this particular vector that we have already seen. So, let us do that then what we are trying to do we are trying to compute the corresponding probability value.

So, probability value corresponding to the F for test so that that is how it is being committed pf is the function that can be used for more detail you can go into the help section and find out more information about pf. So, you can see this is F distribution.

Keele	
e hat cade vow more vesson must coming more soon was	
🛛 • 🔄 • 🔒 📵 🗁 (🖈 cate flayfunction) 🔛 • Addex •	Project (None)
0 10 logi 2 # × 0 10 logi 3 # ×	Environment History
00 2 0 1 Source on Save Q 7 . 0	🞯 🕞 🔐 Import I Jata set + 🥑 📃 I ist + 🙆
167 mod4ava("Residuals","Sum Sq"), 07	Gabal I mircoment Q
168 sum(mod4ava[,"Sum Sq"]))	Data
170 MS=c(mean(head(mod4ava[,"Mean So"],-1)).	anovadf 3 obs. of 5 variables
171 mod4ava["Residuals","Mean Sq"],"")	0 df1 107 obs. of 6 variables
	Odfitest 11 obs. of 6 variables
<pre>1/3 FStatistic=c(mod4summ)fstatistic[[]] mod4summ(fstatistic[[]]] mod4summ(fstatistic[[]]]</pre>	Odfltrain 96 obs. of 6 variables
175 lower.tail = F)	0 df2 5000 obs of 4 variables
176 pvalue=c(P,"","")	Odf2train 3000 obs of 4 variables
1// 178 annuadé-data énama/DE CC MC Estatistic nualua)	The Red Balance Made Manual
179 rownames(anovadf)=c("Regression", "Error", "Total")	Files Plots Packages Help Viewer
180 anovadf	
181	R: the I Distribution . Find in Topic
182 # Prediction for a new observation:	random generation for the F distribution with dill and
1791 (Replayed) 8 R Script 1	af2 degrees of freedom (and optional non-centrality
	parameter ncp).
Console G//Section 10/ P/ multi 11. moddsum (feratistic [["deadf"]])	lleage
ss=c(sum(head(mod4ava[,"sum so"],-1)).	Coage
mod4ava["Residuals", "Sum Sq"],	df(x, dfl, df2, ncp, log FALSE)
sum(mod4ava[,"sum sq"]))	pl(q, dil, di2, ncp, lower.tail = TRUE, 1
<pre>MS=c(mean(head(mod4ava[,"Mean Sq"],-1)), mod4ava[["mec6duals" "Mean Sq"],"])</pre>	<pre>qf(p, df1, df2, ncp, lower.tail TRUE, 1 rl(n, dl1, dl2, ncp)</pre>
Fstatisticuc(mod4summ\$fstatistic[["value"]].""."")	and the search and story i
<pre>P=pf(mod4summSfstatistic[[1]], mod4summSfstatistic[[2]], mod4summSfstatistic[[3]],</pre>	Arguments
lower.tail = F)	
<pre>> pvalue=c(P, ,)</pre>	x, q vector of quantiles.
> anovaor=oata.rrametur, ss, ms, rstatistic, pvalue)	p vector of probabilities.
	•

(Refer Slide Time: 29:09)

So, within F distribution we have F function which can be used to compute the corresponding p value. So, what we need is F statistics and the first argument then we need degree of freedom as second argument that is responding to a regression here and the second is then last one is residuals degree of freedom third argument and then we have a specified lower tail as false.

So, we will get the p value corresponding p value. So, let us record it in this format once this is done we can create this table data frame. And let us assign a names for this and let us have a look at this particular table. So, we can see that now once computed.

(Refer Slide Time: 29:50)



So, first column we have degrees of freedom. So, for regression there are 3 and then because as remember that we have used just 4 variable, so one being outcome variables, so 3 predictors. So, therefore, degree of freedom degree of freedom is here is 3 for regression then residuals it is the remaining that is n minus 1 is 2999, total number of observation and is 3000.

So, residual degree of freedom 2996 sum of a square for regression and for residual is also there. So, you can see that residual sum of square is much higher so that also is we can we saw that the variance was also on the lower side right earlier we saw that the variance that we had computed was on the lower side so that is also indicated here, mean square of errors is also there and then we have F statis F statistic and then we have p value is small value.

(Refer Slide Time: 31:00)



So, this gives us some information about the model and now what we will do is we will use this model to the score of a particular observation. So, let us do a prediction for a new observation let us say this is our new observation is annual income of rupees 5 lakhs with 2 family members who is not active online.

So, this is information about the particular customers customer whether is going to accept the promotional offer or not. So, annual income is 5 lakh family size is 2 and not acting online. So, you can use the predict function first argument is as usual the model object mod 4 then in a data frame we are trying to pass on the values operate predictors for example, income 5.

So, it should be in the same unit as was used for the modelling exercise in the training partition. So, you can see 5 and family size is 2 and then online is because this particular customer is not active so 0. So, once we do this we will be able to predict this particular observation. So, you can see this comes out to be a negative value. Now that was one of the first point that we discussed in the slide.

(Refer Slide Time: 32:11)



Let us go back you can see here that the predictions can take any value not just dummy values. So, we can see that a negative value has been taken here and so that is in one anomaly that we can clearly see. And set of values for them was no offer that we already know 0 and 1, and the value is that comes out to be predicted value comes out to be minus 0.14 now.

So, one difficulty is how do we do our classification in this case? Now, let us look at the residuals so second anomaly that I discussed that outcome variable of residuals do not follow normal distribution. So, let us look at that let us plot a histogram and find out whether this is being followed.

(Refer Slide Time: 33:07)



You can see residuals when we plot residuals when we get histogram you can see this is clearly not being followed normal distribution is not being followed one grouping here and other grouping where this is of this lower values, lower frequency and this is higher frequency. So, this is clearly looks binomial or different groups kind of thing. So, this definitely not following so normal distribution and distortions due to real binomial distribution can be seen here.

So, the typically what exercises and the discussion that we have been doing was pertaining 2 classification tasks and as we talked about that this is a statistical technique logistic regression being a statistical technique is also applied in so used in statistical modelling. So, the kind of task that we generally do in a statistical modelling are quite similar to what we can call profiling tasks.

So, as we talked about in the you know starting lecture of logistic regression and that it is about understanding similarities and differences between two groups. So, logistic regression can also be used in can also use to understand what are the variables which you know which can bring out some of the similarities or differences between group, so let us discuss that aspect as well.

So, in profiling tasks when we talk about profiling tasks the situation is slightly different. So, in the classification task we typically build our model and look at the performance of that classifier that particular model using the classification matrix using the overall accuracy or overall error matrix and you know some deviations when we have a class of interest.

So, those are the things that we typically do we also typically look at left chart especially when we have a class of interest to in left chart and decide charge to see whether it is still despite you know higher error whether it still the model is useful in class of interest with respect to knave you know knave rule or a or a average case.

So, some of those things that we do in classification tasks; however, in profiling tasks what we do in classification we follow that. So, apart from model performance on validation partition we also asses models fit to data on training partition right because as I taught what in a statistical technique typically the whole sample is used; however, since we are using training partition for model building.

So, the models fit to data is assessed on training partition; however, model performance is assessed on validation partition for profiling tasks and models fit to data assessed on training partition. So, some of these things we talked about when we talk about goodness of fit measures we talked about the deviance, we talked about 1 minus deviance divided by null deviance that is equivalent of multiple R square some of those things you can see here as well.

So, models fit to data is assessed on training partition ah; however, we still focus on avoiding over fitting because as we talked about when we have matrix which look for and when we do modelling to achieve the you know goodness of fit then it can it can lead to over fitting. So, it still we would like to avoid over fitting and still be able to you know check the performance still be able to have good classification performance as well.

So, usefulness of usefulness of predictors is also examined in this particular case profiling so because it is about understanding similarities and differences between 2 groups. So, therefore, which predictor is more helpful in terms of bringing out those differences those similarities so when we build our model. We also look at the significance levels of some of these predictors we look at which predictors are significant whether the not significant predictors can be dropped from the model.

And this has also should be looked at from the perspective of model performance because as we have taught about data mining model we would like to keep the insignificant variables also in the order if they provide some practical importance in terms of scoring new observation.

How and logistic statistical modelling we just drop the insignificant variables because we are just interested in understanding the phenomena, understanding the underlying relationship. So, profiling is quite similar to you know that that approach; however, because we are doing data mining modelling. So, we have to balance between these two, we have to balance between performance and also the main profiling tasks.

So, we have to really see whether the predictors can be dropped just like in statistical training or whether they have to be kept in the model because they also provide some practical significance for scoring new observation so that balance has to be achieved.

So, we have to avoid over fitting, we have to look for using usefulness or predictors in both the context data mining context prediction point of view and also statistical context in terms of understanding you know finding understanding variables which differentiate the groups, which bring out similarity or differences between groups.

So, this kind of exercise is done in profiling tasks and goodness of fit matrix that through an exercise in R that we have already understood over phila. So, we look to understand first we look to understand the overall fit of the model, so if the overall fit of the model is good only then we go ahead and look at the individual variables.

So, first step typically is in profiling on a statistical modelling we look at the overall fit of the model. So, in this particular case logistic regression as we talked about the deviance is the metric we taught one previous lecture that could be used and we also said that this is equal to SSE in linear regression and 1 minus deviance is divided by null devian that is equivalent to multiple R square in linear regression.

So, these are two matrix that could be used to assess the overall fit of the model and then once this is done then we look at the single predictors. We look at whether they are significant or not as I talked about and whether we can strike a balance in terms of prediction performance, classification performance, versus the profiling that is and also statistical modelling context.

(Refer Slide Time: 40:04)



So, with this we move to our next discussion that is about so till now the exercises that we have been doing they were mainly focused on binary classification, and binary logistic regression model we just had 2 classes class 0, and class 1. Can logistic regression we extended to a scenario where we are dealing with more than 2 classes where we are dealing with you know m classes. So, yes it is possible so we will discuss some of those things.

So, first one is multinomial logistic regression so multinomial logistic regression. So, the categories the classes that we have they are you know so the categorical variable is nominal, so, in that case we can apply multinomial logistic regression. So, what happens in multinomial logistic regression first out of those m classes that we have we have to select one we have to pick one as the reference category, and for the remaining m minus 1 classes we create separate binary logistic regression model.

So, for each of the m minus 1 classes apart from the reference category class. So, we will have n minus 1 classes so for each of the n minus 1 classes will create separate binary logistic regression model; that means, for a class 1 we will have the scenario where the observation probability of belonging to class 1 and probability of long not belonging to class 1 so that kind of binary scenario we will have and that for each of the n minus 1 classes.

So, we will be dealing with m minus 1 binary logistic regression equations and so using that we can compute all those probabilities values with the help of the predictors and then the remaining reference class of that probability for that can always be computed by the m minus 1 probabilities values for the m minus 1 classes.

So, we can just subtract that from one for you know and then we can get for the probability value for the different class. And once we have the probabilities values for all m classes then we can apply our most probable class method routine where the class having highest probability value would be assigned to the new observation.

So, this is how we can go about applying logistic regression to an outcome variable with m classes. So, this is called this is a multinomial logistic regression and this is applicable mainly to nominal categorical nominal variable right, sos the categorical variable having nominal classes.

The second scenario second scenario is about when we have a categorical variable with ordinal classes we have an ordinal variable. So, in those cases we can apply ordinal logistic regression. So, so again within this ordinal logistic regression we can have 2 scenarios. So, as we have understood in some of the initial lectures and supplementary lectures that ordinal variables they have order among different labels is also important that is also meaningful.

So, less than or equal to or greater than or equal to operations they are also applicable in this case, so the first scenario is large number of ordinal classes. So, if our outcome variable which is categorical variable with ordinal classes if that variable have is having large number of ordinal classes, then one solution is treat that ordinal variable as continuous variable and apply multiple linear regression right.

(Refer Slide Time: 44:10)



So, when we have a categorical variable with ordinal classes so we have a categorical variable. Categorical outcome variable when m ordinal classes, and m is large then I as I talked about we can treat this particular variable as a continuous variable and apply multiple linear regression.

So, multiple linear regression can be applied and reason is one reason one justification for this is as we talked about earlier that in binary situation we have just two values for a categorical variable and the predicted values up using MLR can range anywhere any real value so, that was the one main problem here.

But when we have m is large; that means, set of values could be you know many more. So, it could be you know if there are 50 groups let us so in this fashion we can go on up to this. So, the number of values that can be taken by this particular this particular ordinal variable are many more.

So, therefore, the predicted values this can be easily mapped to some of these values could be close to some of these values, and probably multiple linear regression can be still applied so this is one way. When we have a ordinal variable with many number of classes with large number of classes, when m is large then is still you know instead of logistic regression we can apply multiple linear regression so that is first scenario.

The second one is whatever we have a small number of ordinal classes. So, if m is in this case m ordinal classes that we have if this is a small m is small, then probably we will we will run into the same problem like for binary classification. So, similar problem would be there we will have only few values 0, 1, 2, 3 let us say these many.

So, again small number of classes small number of ordinal classes, so we will have I will run into same problem. So, so we what we do is we use a different version of logistic regression called proportional odds or cumulative logic method as indicated in the slide, so small number of ordinal classes so we would like to use proportional odds, or cumulative cumulative, or logit method.

So, what we do here in this particular method is we create separate binary logistic regression model for m minus 1 cumulative probabilities, so, we talked about that when we so, when we discussed multinomial logistic regression; So, for all m minus 1 classes will have separate binary logistic regression model.

(Refer Slide Time: 47:50)



However, if you see that here will have separate binary logistic regression model not for m minus 1 classes not for presence of that class or absence of that class, but m minus 1 cumulative probabilities.

So, let us understand what we mean by that so let us take an example for a 3 class case as written in this slide for a 3 class case C 1, C 2, and C 3. Let us say these are the 3 classes

ordinal classes C 1, C 2, C 3 and a single predictor x 1 that is being used. So, our logit equations could be something like this logit for C 1 it could be alpha 0 plus beta 1 x 1 and logit for C 1, or C 2 so; that means, from C 1 first logit equation is just for you know observation belonging to C 1. The second is observation belonging to C 1, or C 2 so that gives us the cumulative sense. So, as we talked about that ordinal the order is important so that that means, you know different classes can be compared.

So, therefore, C 1, or C 2 is a you know is a meaningful here in the sense that if we look at the rights part of the equation beta 0 plus beta 1, x 1 you can see that beta 1 is same and both these equation x 1 so you can see because the this is the comparison can be done.

So, the coefficient so the intercept R only difference because the comparison so one is when we talk about ordinal, so one particular class this ordering this ordering is you know this is meaningful. So, when the ordering of classes is meaningful; that means, one can be you know said you know less than one particular class the or higher than one particular class just like the categories that we might have high, you know low medium, upper medium, medium and then low.

So, this kind of this kind of classes we might have all we might have you know a strongly agree to strongly disagree. So, those kind of ordinal classes where the order is meaningful, where the order is meaningful then in those regression we can have something like this cumulative probabilities values, and beta coefficient is can be used the same beta coefficient can be used for both these logistic regression.

And from this we can compute the cumulative probabilities values and once these cumulative probabilities value values for two of the classes for these classes have been computed the actual probabilities value for C 1, C 2, C 3 can be derived from using the probability value that formulation that we have so, it is from logit, we can derive the probabilities values.

So, once these cumulative so from there once we have the probabilities values for all these classes then we can again apply the most probable class method and also assign the class based on the based on the probability value and again this so in this fashion in this fashion we can apply ordinary logistic regression to a scenario with fewer number of ordinal classes. So, what we will do we will go through R Studio and do an exercise for this when we have classes more than m is greater than 2, so a logistic regression modelling for classes greater than 2. So, what we will do we will create a hypothetical data set here, so number of observations are 100 in this case, so as you can see and it is 1000.

So, let us create this now what we are going to do is we are going to create a data frame having that data frame where we have 2 variables x 1 and x 2, x 1 as you can see we are using run if for x 1 and x 2 both so the observations. So, the wave values would lie between 0 and 100 and n number of values would be created that is 1000 values would be created lying between 0 and 100.

(Refer Slide Time: 52:23)



So, let us create this data frame, let us look at some of the observation. For 6 observations you can see all the values they are lying for both the variables x 1 and x 2; they are between 0 and 100.

Now, what we will do we will create a categorical outcome variable with 3 classes. So, this particular data set that we are trying to create we are going to use it for both the scenarios multinomial scenario, and ordinal scenario right, so this is just for illustration purpose so we are not specifying whether the variable is ordinal or not s, we are just going to use it for both the scenarios.

So, what we are going to do is we are using transform function to create our outcome variable, categorical outcome variables. So, you can see that we are trying to compute y as 1 plus if else and if the value is less than 0, this particular value 100 minus x 1 minus x 2 plus again a certain value is being taken from normal distribution no standard deviation 10000 values.

So, if it is less than 0, then you know we are using this information and that is information based on x 1, x 2, and certain additional computation we are trying to assign it a class 0 or then the second if that is not true then second computation so it will get a class 1 or class 2. So, let us compute this once this is done let us look at the observations you can see another variable y that is categorical variable has been created having you can see 1 and 3 2 values are being taken. Let us look at the structure of this particular data frame.

(Refer Slide Time: 54:00)

er hat can' vow more when well long more land more	
🛃 • 😹 • 🔛 🚱 🕼 (🛦 to file/function 👘 🖾 • 🛛 Addiex •	B Project (Nore)
0 10 logi 28 × 0 To logi 18 ×	Environment History
0 0 0 0 0 5 Source on Save 9 2 - 0 - Pun 59 0 Source - 2	🞯 🔒 📑 Import I Jatavet + 🥑 📃 i ist + 🌀
<pre>1 # sysothetical dataset # no. of observations=1000 n=1000 5 # create predictors: v1 and v2 6 df.dts.frame(x1erunif(n,0,100), x2erunif(n,0,100)) 7 head(df) 8 & Create categorical outcome variable with three classes 9 df.rtmsform(df, y=1+ifelse(100-x1-x2ernorm(n,sd=10)<0, 0, 10 ifelse(100-2*x2ernorm(n,sd=10)<0, 1, 2))) 11 head(df) 12 str(df) 13 str(df) 14 palette(15 palette(gray(0:3/3)))</pre>	Abbdi Innovenet • • 0 df 1000 obs. of 3 variables 0 df1 107 obs. of 6 variables 0 df1rest 11 obs. of 6 variables 0 df1rest 11 obs. of 6 variables 0 df2rain 1000 obs. of 4 variables 0 df2rain 1000 obs. of 4 variables 0 df1rest 11 obs. of 6 variables 0 df2rain 1000 obs. of 4 variables 0 df1rest 1000 obs. of 4 variables 0 df2rain 1000 obs. of 1 variables 0 df1rest 1000 obs. of 3 variables 0 df2 rain 1000 obs. of 3 variables 0 df2 rain 1000 obs. of 3 variables 0 df1rest 1000 obs. of 3 variables 0 df2 rain 1000 obs. of 3 variables 0 df1rest 1000 obs. of 3 variables
17 rance (dfs1) 141 (Bip Lewith 2) 130.04808 58.066454 1 2 16:30433 19:30933 3 35:40761 0.506333 4 25:1094 5:005623 3 25:23010 6:628991 > > str(df) 5 x2: mum 30:63 35:52.51 3:x2: mum 36:63 35:52.51 3:x2: mum 35:61 39:16 4:x: 5 x2: mum 33:11 3:12	6 6 10 -05 00 05 10 Residuals

So, these are 3 variables x 1, x 2 values between 0 and 100, and y is taken value 3 values 1, 2, and 3. So, let us plot this particular data set. So, this is our default palette so; however, I would like to use this palette agree 3 sets. So, in this fashion we can have sets of no number of sets that we require. So, let us look at the range of x 1, y, and x 2 which is already.

(Refer Slide Time: 54:25)



Because we have just now created these variables studies actually clearly understood as well so limits 0 to 100, and 0 to 100 and then colouring is with using this particular factor y. So, you can see as dot factor we are using this particular variable and let us create this plot.

(Refer Slide Time: 54:58)



So, you can see so this is our plotting. So, a one group is here, the second group is having some medium level gray set here and the third group is here, this is lighter gray code

colour. So, these are the values that we are going to use x 1, x 2 this is plot between x 1, and x 2 and the outcome variable has been color coded, so, 3 categories.

(Refer Slide Time: 55:25)

HAR CADE Your Ports Session Ruid Linning Protect Lands Help	
• 📴 • 🔄 🗐 🔐 🏟 La to file/function 🔄 🔯 • 🛛 Addient •	Project (Nork
2 10 logit 28 × 0 Tra logit 18 ×	Environment History
S D D R T Source on Save Q / • E • BRun De De Source • 3	🛿 🔐 🕞 🔐 import Dataset + 🥑 📃 List + 🖉
1/ range(dT)X1)	Gabal minoment Q
18 range(dfsy) 19 range(dfsy)	o df 1000 obs of 3 variables
20 $plot(df_x)$, df_x^2 , $xlim = c(0, 100)$, $vlim = c(0, 100)$.	odfi 107 abs of 6 variables
<pre>21 xlab="x1", ylab="x2", col=as.factor(dfSy),</pre>	Odti 10/ 005. of 6 variables
<pre>22 pch=19, cex=0.8, panel.first = grid())</pre>	- Odfitest 11 obs. of 6 variables
23	Odfltrain 96 obs. of 6 variables
24 F Multinomial Logistic Regression	Odf2 5000 obs. of 4 variables
26 mod_multinom(v - data = df)	Odf2train 3000 obs. of 4 variables
27	Odfh 107 obs of 13 variables
28 summary(mod)	Eiles Plots Parkages Help Viewer
29	
30 head(modSfitted.values)	😳 💿 🖉 Zoom 🌁 Export • 👰 🔮 •
32 modtrain-product(mod_df(-3)_type = "probs")	
33 head(modtrain)	
N1 (Top Level) 8 R Script	12
Console G//Session 10/ 💬 😁	B. (B. (4) 21 ((1) 21 (1) (1)
L] "black" "red" "green3" "blue" "cyan" "magenta" "yellow" "gray"	
palette(gray(0:3/3))	N - 355 (0.2 (0.2 (0.2 (0.2 (0.2 (0.2 (0.2 (0.2
range (drax1)	× o +
range(df\$y)	TARA SARANSI S
1] 1 3	35537205200F8
range(df\$x2)	
1] 0.004542479 99.993072520	0 20 40 60 80 100
plot(df3x1, df3x2, xlim = c(0, 100), ylim = c(0, 100),	20 40 00 100
xiab= x1 , yiab= x2 , col=as.factor(ofSy),	
ochely cevel X papel first = ocid(1)	V4

Now, what we will do is we will use the multinomial logistic regression. So, for this we need this package the library, and net package this package actually for a neural network, but it provides us offers us this function which can be used for multinomial logistic regression.

So, a multi norm is the function so what we are now going to do is the outcome variable y we are going to regress it with the remaining variables that is predictors in this particular data. So, all the observations we are going to use here, so let us run, this let us look at the summary.

(Refer Slide Time: 56:00)



So, we can see coefficients x 1 and x 2 right. We can also see the error values residual deviance and the AIC values, are also indicated here. Now, fitted values also we can look at first 6 observation here.

(Refer Slide Time: 56:19)

har code were then sensor must sensor more soon were	
🔸 🤠 • 🔄 🔛 🔮 🕼 Calto Blattunction 🔤 🛛 🔤 • 🛛 Addies •	B Project (Nors)
10 logit 2.R × 9 10 logit 1.H ×	Environment History
o a B = 5 Source on Save Q Z + E + → → Bun 😁 → Sou	nte • ≷ 🔐 📴 import Dataset • 🧃 📃 tiet • 🕻
24 # Multinomial Logistic Repression	* 🚯 Global Environment • Q,
25 library(nnet)	values
<pre>26 mod=multinom(y ~ ., data = df)</pre>	0 breaks1 2017-08-28
27	Obseals2 2017 08 28 12:00:00
28 summary(mod)	orex che [1:107] "0.12" "0.12"
30 head(mod(fitted values)	DEP1 CHP [1:10/] 0-12 0-12
31	Fernigric che [1:3] "203 26120204222
<pre>32 modtrain=predict(mod, df[,-3], type = "probs")</pre>	Pstatistic tir [1.5] 505.5020054222
33 head(modtrain)	a whod List of 20
34 # Classify observations	Files Plots Packages Help Viewer
<pre>35 modtrainc=ifelse(modtrain[,1]>modtrain[,2] & modtrain[,1]>modtrain[,3], 1, 35</pre>	
30 ITelse(modtrain[,2]>modtrain[,3], 2, 3)) 37	👳 👳 💆 Zoom 🎿 Export • 😢 🖉 🧐 •
<pre>38 table("actual Value"=dfSv "Dredicted Value"=modtrainc)</pre>	
JU LEUTE, ACCUST VALUE -DIJY, Predicted Value -modelame)	
39 #classfication accuracy	
240 (hptech) t	R Script 1
and cancer accurate any produced value inductancy 39 classification accuracy 249 (dipliced) i annole C//tenion 10/ P	R Script 2
ander cytension ng / Angel and angel and angel	R Souge 2
y classification accuracy y classification accuracy 248 (Replayed) 5 anote Cr/sension 10 ∞ psidual Deviance: 478,5526	Richard 1
All Carlos Contractor accuracy All Safety and Carlos Contractor and Carlos Contractor All Safety and Carlos Contractor	
and address of the action act	
246 (apple to accuracy 246 (figured) 3 246 (apple to accuracy 246 (figured) 3	
30 Cash (*** Construction accuracy) 36 Cash (***) 2245 (Replaced) 3 Cash (***) 3 Cash (***) 4 2 3 Cash (***) 1 2 3 Cash (***) 3 Cash (***) 3 Cash (***) 4 2 3 Cash (***) 1 3 1 1336438-01 1 3502428-01 1 1334348-01 3 Cash (***)	
30 Last for a court y 30 Last for a court y 24 (Replaced a 25 (Replaced a 2 108104+015 2 3 2.108104+015 5956144+011 1.1182000+02 1.38438+018 2.8181324+010 3.5902424+011 3.581324+02 3.584394+002 3.581324+02 3.584394+002 3.581324+02 3.584394+002 3.584394+02 5.502744+011	
246 Charles Los and Lange Strategy 246 Charles Los accuracy 247 Charles Los accuracy 248 Charles Los accuracy 249 Charles Los accuracy 249 Charles Los accuracy 240 Charles Los accuracy 241 Charles Los accuracy 242 Charles Los accuracy 243 Charles Los accuracy 244 Charles Los accuracy 245 Charles Los accuracy 246 Charles Los accuracy 247 Charles Los accuracy 248 Charles Los accuracy 249 Charles Los accuracy 240 Charles Los accuracy 241 Charles Los accuracy 25 Charles Los accuracy 26 Charles Los accuracy 27 Charles Los accuracy	
	R Sorget 1 C C C C C C C C C C C C C

So, we will get these fitted values which are nothing, but estimated probabilities values for these 3 classes 1, 2, and 3 so what this is what we have. Now, if we had another partition validation or test partition you can use predict function to the score those partitions and have probabilities values, estimated probability values, for those new observation.

However for the demonstration purpose, we are applying predict function on the training partition that is the full data set itself. So, you can see another argument that we can see is type which is props which is proper probabilities values. So, let us run this and you would see we will have the same probabilities values which were fitted. So, we can see here that same values you can see first row it is same, same values. So, the fitted values and all we have got the same one using predict function.

Now, let us classify these observations. So, this is how we can classify if the probability value for 1 is greater than the probability value for class 2 and also greater than the probability value over class 3, then of course, this observation has to be classified to class 1.

Otherwise we will again compare the probability value for class 2, with probability value for class 3, if again it is greater and I know class 3 value divided then the class 2 is assigned, otherwise class 3.

So, in this fashion we can assign all the observation into appropriate classes so this is implementation of most probable class method. So, once this is done you can see mod train see for all 1000 variables have been created and observations have been assigned to go to appropriate classes.

(Refer Slide Time: 58:13)



Now, let us look at the classification matrix. So, we can see actual value, predicted values, now till now the classification matrices that we have been observing that we have been creating they had 2 values 0 and 1. And we had 2 by 2 classification matrix now this time we are seeing 3 by 3 classification matrix.

So, 3 possibilities for actual values 1, 2, and 3 and the predicted values and we have the corresponding numbers here. So, you can see that again the diagonal element; that means, these 3 elements they represent the values records which have been correctly classified, and the remaining observation off diagonal elements you know they have the records, they have the counts of record, which have been incorrectly classified right.

So, from this we can compute the classification accuracy that is 91 percent in this case, and error that is the remaining 9 percent. So, in this fashion we can apply multinomial logistic regression to a particular data set. Now, let us move to next part that is ordinal logistic regression.

So, in this product case ordinal logistic regression as I talked about 2 scenarios, so 1 where m is large; So, in the in that scenario we can apply multiple regression one exercise that we had done in previous lecture, we had applied multiple linear regression, but that was for the class which had just two you know for a variable category which just had 2 classes 0 and 1.

In the same fashion for the first scenario the multiple linear regression, can be applied. So, there is not much much different in terms of applicability. So, well what we will do we will do an exercise for this scenario where we have to apply this ordinary ordinal logistic regression the cumulative probabilities method, cumulative logit method so for this we need this package mass.

(Refer Slide Time: 60:33)

Mada	
He HAT CARE Were Plans Session Ruild Carbing Planter Local Marp	
💽 • 🔐 • 🔒 🗐 🦀 🕼 to file/function	B Project (None)
© 10 hgs 12 K x © 10 hgs 11 K x	Environment History
<pre>36</pre>	Omod3 List of 30 0 mod3 List of 30 0 mod4 Large Im (12 elements, 801. 0 mod4summ List of 11 modtrainc Named num (1:1000) 2 3 3 3. MS chr (1:3) "24.14772431505. n 1000 P 1.00014894386617e-210
47 48 summary(mod1) 49 50 head(mod1)fitted.values) 51 40 (mod1)fitted.values) 52 mod1=sisuenedite/mod1_dff_=11_*uma = "menbe") 40 fill (mod1) = 12 fill	File Plots Packages Help Viewer Image: Condema to point Image: Condema to po
Conder C//sendon N/ © table("Actual Value"=dfsy, "Predicted Value"=modtrainc) Predicted Value Actual Value 1 2 3 1 440 14 23 2 20 97 11 > mean(modtrainc=dfsy) [1] 0.91 > mean(modtrainc=dfsy) [1] 0.92 > library(MASS) >	Ordered Logistic or Probit Regression Description Fits a logistic or probit regression model to an ordered factor regionse. The default logistic case is proportional doci logistic regression, after which the function is named.

So, let us load this particular library and we have polr is the function for this. So, this is actually for profit however, it can also be used for it can also be used for the cumulative logit method. So, we will just see in the in the in the help section. So, you can see polr this is ordered logistic, or probit regression right.

So, what we are interested in ordered logistic regression, so so here ordered logistic regression. So, you can see in the method argument the first one is draw logistic this is actually ordered logistic method. This is also called proportional odds logistic regression which we have discussed right.

So, let us go and build this model so polr is the function. So, first as you can see y variable I have converted into a factor variable and then the request against all other variables that which are predictors data is full, then we need this argument as well as which is mainly if we want to apply for summary function later on the model object, so, this is also true.

So, now let us apply summary, so we will get the results we can you can see the coefficient values here x 1, and x 2, the value and the error and the T values are also there. And we have the residual deviance, and AIC value as well for this particular model and we are interested in finding the fitted values so that is also returned by the model. So, let us look at some of these values.

(Refer Slide Time: 62:00)



So, you can see for each of these classes class 1, 2, and 3, so these values are actually probability estimated probabilities values. Now using these values we can again apply the most probable class method so first we need to compute the first we need to compute and do the assignment as per the most probable class method like we did in previous exercise.

So, as I talked about predict function can again be used to score new data. So, in this case we are scoring the training partition itself again. So, we expect to get the same values as you can see last row, you can see here, here also the same values are there. So, it is scoring off for the training partition itself. So, we will get the same observation.

Now, what we will do classify observation. So, as we did in previous exercise mod 1 train for class 1 and probability value for class 1 greater than value class 2, and greater than class 3. Then assign it to class 1 otherwise again we look on more comparison the probability value for class 2 is greater than probability value class 3, then assign it to

class 2, otherwise class 3. So, in this way we will have the appropriate classification scores.

(Refer Slide Time: 63:17)

Ruik.	=.0.
e har cone sow more sense wate contag more cone wep	
🔹 🔐 🔹 🔗 🧀 in the function	Project (No
е) 10 юря 28 × е) 10 юря 18 × —	Environment History -
O O O O R IT Source on Save Q / - E - → Run De De Source - ≥	😭 📄 🔐 Import Datavet + 🧹 📃 List +
48 summary(mod1) .	Cablal Invironment • Q
49 50 head(mod1(fitted values)	Obreaks2 2017-08-28 12:00:00
51	DEPT chr [1:107] "0-12" "0-12"
<pre>52 modltrain=predict(mod1, df[,-3], type = "probs")</pre>	DF num [1:3] 3 2996 2999
53 head(modifrain)	Fstatistic chr [1:3] "383.36120294222
<pre>55 modltrainc=ifelse(modltrain[.1]>modltrain[.2] & modltrain[.1]>modltrain[.3]. 1.</pre>	0 mod List of 26
<pre>56 ifelse(modltrain[,2]>modltrain[,3], 2, 3))</pre>	0 mod1 List of 18
57	modltrainc Named num [1:1000] 2 3 2 3
56 table("Actual Value"=df3y, "Predicted Value"=moditrainc) 59 #classification accuracy	
60 mean(modltrainc==dfSy)	Files Plots Packages Help Viewer
61 #misclassification error	🖕 🤿 🏠 🗁 🖉 🔍 Q, polr 🛛 🛛
62 mean(modltrainc!=df5y)	R: Ordered Logistic or Probit Regression + Find in Topic
03	
d1 (Top Level) 8 R Script 8	Ordered Lesistic er Drehit
	Ordered Logistic or Probit
onsole C://Session 10/ 🕫 📼	Regression
<pre>modltrain=predict(modl, df[,-3], type = "probs") head(modltain=)</pre>	
	Description
0.3188775524 0.4415837004 0.2395387471	
0.0515084371 0.2176331067 0.7308584562	Fits a logistic or probit regression model to an ordered
0.2585359098 0.4442402184 0.2972238718	factor response. The default logistic case is
0.9697389188 0.0256803680 0.0045807132	proportional odds logistic regression, after which the function is named
0.9964164603 0.0030534632 0.0005300764	
<pre>modltrainc=ifelse(modltrain[,1]>modltrain[,2] & modltrain[,1]>modltrain[,3], 1,</pre>	Usage
<pre>ifelse(moditrain[,2]>moditrain[,3], 2, 3))</pre>	

Now, let us generate the classification matrix here. so you can see 3 by 3 matrix we have 3 classes actual values 3 possibilities, predicted values, 3 predicted classes, so again diagonal elements they represent the correct classification values and off diagonal elements they represent the incorrect classifications.

So, what we can do is so let us compute a classification accuracy, and you can see eighty two percent and the remaining is error.

(Refer Slide Time: 63:49)



So, with this we have completed our discussion on logistic regression and so today we have been also able to cover the scenarios where more than 2 classes are present in our categorical variable. What happens when the classes are nominal and how we can apply logistic regression when classes are ordinal, so we have seen that we have also done an exercise in R.

So, we stop here and we will continue our discussions in next structure for a new technique.

Thank you.