

εллнпікн днмократіа Εдνικόν και Καποδιστριακόν Πανεπιστήμιον Αдηνών

Γραμμικά Μοντέλα

Σημειώσεις Εργαστηρίων

Διδάσκουσα: Λουκία Μελιγκοτσίδου

Τμήμα Μαθηματικών Εθνικό και Καποδιστριακό Πανεπιστήμιο Αθηνών

Contents

1	Εισαγωγή	3
	Cars	3
	Grades	3
2	Περιγραφική Στατιστική	5
	Είδη Μεταβλητών	5
	Ποιοτικές Μεταβλητές	5
	Ποσοτικές Μεταβλητές	8
	Έλεγχος Κανονικότητας	9
	Ασκήσεις	14
3	Συσχέτιση	15
	Διάγραμμα διασποράς	15
	Συντελεστής γραμμικής συσχέτισης Pearson	17
4	Απλή Γραμμική Παλινδρόμηση	20
	Εισαγωγή	20
	Εφαρμογή του μοντέλου	20
	Αποτελέσματα	23
	Πίνακας Model Summary	23
	Πίνακας ΑΝΟVΑ	23
	Πίνακας Coefficients	24
	Έλεγχος υποθέσεων γραμμικού μοντέλου	25
	Ασκήσεις	28
5	Πολλαπλή Γραμμική Παλινδρόμηση	29
	Εισαγωγή	29
	Εφαρμογή του μοντέλου	30
	Αποτελέσματα	34
	Πίνακας Model Summary	34
	Πίνακας ΑΝΟVΑ	34
	Πίνακας Coefficients	35

	Επιλογή Μοντέλου	36
	Αποτελέσματα για το Τελικό Μοντέλο	41
	Έλεγχος υποθέσεων γραμμικού μοντέλου	42
	Μετασχηματισμός Μοντέλου	44
	Ασκήσεις	49
6	Διαστήματα Εμπιστοσύνης & Έλεγχοι Υποθέσεων	50
	Διαστήματα Εμπιστοσύνης	50
	Έλεγχοι Υποθέσεων	51
	Παραδείγματα	53
	Διαστήματα Πρόβλεψης για την Εξαρτημένη Μεταβλητή	55
7	ANOVA & t-test	57
	One Sample t-test	57
	Independent Sample t-test	58
	One-Way ANOVA	63
	Two-Way ANOVA	67
8	Παράρτημα	71
	Μέτρα Θέσης και Διασποράς	71

1 Εισαγωγή

Οι σημειώσεις αυτές αποτελούν βοηθητικό υλικό για τα εργαστήρια του μαθήματος "Γραμμικά Μοντέλα". Περιλαμβάνουν τη στατιστική ανάλυση του σετ δεδομένων **cars.sav** σε SPSS και σε R. Παράλληλα, οι φοιτητές μπορούν να εξασκούνται στο αντίστοιχο σετ δεδομένων **grades.sav**.

Cars

Αυτό το σετ δεδομένων αφορά χαρακτηριστικά από 391 αυτοκίνητα. Συγκεκριμένα περιέχει τις μεταβλητές:

- Price: Η τιμή του αυτοκινήτου σε ευρώ.
- Engine: Το μέγεθος του κινητήρα σε cc (κυβικά εκατοστά).
- Horsepower: Η ιπποδύναμη του αυτοκινήτου σε ίππους.
- Weight: Η μάζα του αυτοκινήτου σε κιλά.
- Acceleration: Ο χρόνος που απαιτείται για να φτάσει το αυτοκίνητο ταχύτητα 100Km/hr, ξεκινώντας από την ηρεμία.
- Year: Το έτος κυκλοφορίας του αυτοκινήτου.
- Fuel: Η κατανάλωση καυσίμου του οχήματος, κωδικοποιημένη ως 1: Χαμηλή, 2: Μέτρια, 3: Υψηλή και 4: Πολύ υψηλή.
- Brand: Η μάρκα του αυτοκινήτου, κωδικοποιημένη ως 1: BMW, 2: Datsun και 3: Mercedes-Benz.
- Type: Ο τύπος του κινητήρα, κωδικοποιημένος ως 0: std, 1: turbo

Grades

Αυτό το σετ δεδομένων αφορά βαθμούς μαθητών δύο τάξεων γυμνασίου. Συγκεκριμένα περιέχει τις μεταβλητές:

- Grades: Ο γενικός μέσος όρος του μαθητή, στην κλίμακα του 100.
- Maths: Ο βαθμός του μαθητή στα μαθηματικά, στην κλίμακα του 100.
- Physics: Ο βαθμός του μαθητή στη Φυσική, στην κλίμακα του 100.

- Gymnastics: Ο βαθμός του μαθητή στη Γυμναστική, στην κλίμακα του 100.
- Hours: Οι ώρες καθημερινού διαβάσματος του μαθητή.
- GoingOut: Η συχνότητα εξόδων του μαθητή, κωδικοποιημένη ως 0: Καθόλου, 1: Μία φορά το μήνα,
 2: 2-3 φορές το μήνα, 3: Μία φορά τη βδομάδα, 4: 2-3 φορές τη βδομάδα, 5: Κάθε μέρα.
- Class: Το τμήμα του μαθητή, κωδικοποιημένο ως 1 ή 2.
- Sex: Το φύλο του μαθητή, κωδικοποιημένο ως 0: κορίτσι, 1: αγόρι.

2 Περιγραφική Στατιστική

Σύνοψη: Το πρώτο βήμα στην ανάλυση δεδομένων είναι η παρουσίαση και σύνοψη των πληροφοριών του δείγματος για τις μεταβλητές που περιλαμβάνονται σε αυτό. Στις ενότητες που ακολουθούν παραθέτουμε τη διαδικασία για τη συνοπτική παρουσίαση ποιοτικών και ποσοτικών δεδομένων.

Είδη Μεταβλητών

Οι μεταβλητές διακρίνονται σε ποσοτικές (scale) και ποιοτικές/κατηγορικές (categorical), οι οποίες με τη σειρά τους διακρίνονται σε διατάξιμες (ordinal) και ονοματικές (nominal). Ποσοτικές ονομάζονται οι μεταβλητές που μπορούν να μετρηθούν (έχουν δηλαδή αριθμητικές τιμές). Παραδείγματα ποσοτικών μεταβλητών είναι η ηλικία, το βάρος, το ύψος, η αξία μίας μετοχής κ.ά. Ποιοτικές ονομάζονται οι μεταβλητές που περιγράφουν χαρακτηριστικά του πληθυσμού που μεταβάλλονται κατά ποιότητα ή είδος. Τέτοιες μεταβλητές είναι το φύλο, το χρώμα των ματιών, η τάξη ενός μαθητή κ.ο.κ. Οι ποιοτικές μεταβλητές που έχουν τη δυνατότητα διάταξης ονομάζονται διατάξιμες (π.χ. η τάξη ενός μαθητή) ενώ οι άλλες ονοματικές (π.χ. το χρώμα των ματιών).

Παράδειγμα: Παρατηρούμε ότι στο αρχείο cars έχουμε 5 ποσοτικές μεταβλητές (Price, Engine, Horsepower, Weight, Acceleration) και 4 κατηγορικές (Year, Fuel, Brand, Type), από τις οποίες 2 είναι διατάξιμες (Year, Fuel) και 2 ονοματικές (Brand, Type).

Ποιοτικές Μεταβλητές

Η συνοπτική παρουσίαση των δεδομένων μίας ποιοτικής μεταβλητής με τον πίνακα συχνοτήτων και με τις γραφικές τους παραστάσεις, το ραβδόγραμμα (bar chart) και το κυκλικό διάγραμμα (pie chart). Ο πίνακας συχνοτήτων μιας ποιοτικής μεταβλητής προκύπτει από την απαρίθμηση και καταγραφή των δειγματικών τιμών στην αντίστοιχη κατηγορία. Ένας ολοκληρωμένος πίνακας συχνοτήτων μίας ποιοτικής μεταβλητής περιλαμβάνει τη στήλη των συχνοτήτων και τη στήλη των σχετικών συχνοτήτων. Επιπλέον, αν η μεταβλητή είναι διατάξιμη, μπορεί να συμπεριληφθεί η στήλη των αθροιστικών συχνοτήτων και των αθροιστικών σχετικών συχνοτήτων (ή των αθροιστικών ποσοστών). Ένας τρόπος άμεσης κατανόησης των χαρακτηριστικών της κατανομής των συχνοτήτων επιτυγχάνεται με μία ειδική γραφική παράσταση που ονομάζεται ραβδόγραμμα. Στον οριζόντιο άξονα ενός ραβδογράμματος συχνοτήτων σημειώνονται οι κατηγορίες στις οποίες τα μέλη του πληθυσμού κατατάσσονται, ενώ στον κατακόρυφο άξονα οι αντίστοιχες συχνότητες (εναλλακτικά οι αντίστοιχες σχετικές συχνότητες).

Παράδειγμα: Θα δημιουργήσουμε ένα πίνακα συχνοτήτων και ένα ραβδόγραμμα για την (διατάξιμη) μεταβλητή Year.

ta 👘								cars.sav [D	ataSet1]	- IBM S	PSS Statistic	s Data Edito	r					- 0	×				
<u>F</u> ile	<u>E</u> dit	View	<u>D</u> ata	<u>T</u> ransform	Analyze	Direct <u>M</u> arketing	<u>G</u> raphs	<u>U</u> tilities	Add- <u>o</u> ns	Window	<u>H</u> elp												
			i n		Report	ts	•	-1 ¥-						BC.									
					Descri	ptive Statistics	•	123 Frequer	ncies			14											
					Ta <u>b</u> les	3	•	Descrip	Descriptives		V												
		Price	э 🗌	Engine	Co <u>m</u> pa	are Means	•	A Evoloro		A Explore		A Explore		Jel	Brand	Туре	var	var	var	var	var	var	
1		34	4090	5571	Genera	al Linear Model	•	Connecto		1	1	1											
2		3	7887	6390	Genera	alized Linear Mode	ls 🕨	Tossia	108	1	1	1											
3		4	0960	7210	Mixed I	Models	•	P-P Plots	1	1	1												
4		4:	2170	7439	Correla	ate	*		ts	1	1	1											
5		3	0680	6554	Regres	ssion	•	🛃 Q-Q Plo	ts	1	1	1											
6		4	0478	7456	Logline	ear		D 2	017	1	1	1											
7		3	5023	6276	Neural	Networks) 2	016	1	1	1											
8		3	5780	7030	Classi	fv	•) 2	017	1	1	1											
9		4:	2794	7456	Dimen	sion Reduction) 2	016	1	1	1											
10	27204 4949		Scale	Sion requeston		1 2	016	1	1	1													
11		33146 5211		Nenne	rometrie Teste		1 2	016	1	1	1												
12		3	1601	5735	<u>In</u> onpa	rametric rests		2 2	017	1	1	1											
13		3	1300	4981	Foreca	isting		2 2	017	1	1	1											
14		2	5579	5030	Surviva	-		2 2	017	1	1	1											
15		24	4485	1982	Multiple	e Response	•	3 2	017	2	2	1											
16		3	9089	5211	🎎 Missing	g Value Anal <u>y</u> sis		4 2	018	1	1	1											
17		4	0799	5899	Multiple	e Imputation	•	4 2	018	1	1	1											
18		2	1037	1589	Comp	lex Samples	•	52	017	2	3	1											
19		2	1141	1753	🖶 Simula	ition		52	015	2	2	0											
20		2	2154	1851	<u>Q</u> uality	Control	•	52	016	2	3	0											
21		1	9962	3261	ROC C	urve		52	014	2	1	0											
22		3	9161	5030		1909		5 2	017	1	1	1											
23		2	0986	3261	9	7 1248	1	6 2	019	1	1	1											
24		2	3291	3244	9	1275	1	6 2	014	2	1	1											
25		1	9876	3277	8	1164	1	6 2	018	2	1	0							-				
		4													1949 P.								
Data V	/iew	Variable V	/iew																				
Freque	ncies.												IB	M SPSS Statis	tics Processor	is ready	Unicode:0	N					

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Descriptive Statistics \rightarrow Frequencies.

Μετακινούμε δεξιά την μεταβλητή Year, επιλέγουμε το Display Frequency Tables και στο μενού Charts επιλέγουμε το Bar charts για να δημιουργήσουμε το ραβδόγραμμα.

Frequencies ×	ta Frequencies: Charts ×
✓ Price Statistics ✓ Engine Year Charts ✓ Horsepower Charts Eormat ✓ Weight ✓ Style ✓ Fuel Style Bootstrap ✓ Type OK Paste Reset Cancel Help	Chart Type None Bar charts Pie charts Histograms: Show normal curve on histogram Chart Values Frequencies Percentages Continue Cancel Help

Το SPSS μας δίνει τον πίνακα συχνοτήτων (μαζί με τα αθροιστικά ποσοστά, αφού η Year είναι διατάξιμη μεταβλητή) και το ραβδόγραμμα. Κάνοντας κλικ πάνω στο γράφημα, ένα νέο παράθυρο ανοίγει και μας δίνεται η επιλογή να αλλάξουμε την εμφάνισή του γραφήματος (χρώμα, μέγεθος κ.ά.).



Ποσοτικές Μεταβλητές

Η συνοπτική παρουσίαση των δεδομένων ποσοτικών μεταβλητών γίνεται με τον υπολογισμό των τιμών διάφορων περιγραφικών μέτρων, όπως η μέση τιμή (mean), η τυπική απόκλιση (std deviation) κ.ά., και την κατασκευή του ιστογράμματος (histogram) και του θηκογράμματος (boxplot) της ποσοτικής μεταβλητής. Επιπλέον, γίνεται έλεγχος αν οι διαθέσιμες δειγματικές τιμές περιγράφονται ικανοποιητικά από την κανονική κατανομή.

Τα περιγραφικά μέτρα χωρίζονται σε μέτρα θέσης (που μας δείχνουν πού συγκεντρώνονται οι τιμές) και μέτρα διασποράς (που μας δείχνουν πόσο εξαπλώνονται οι τιμές). Τα σημαντικότερα μέτρα θέσης είναι η μέση τιμή (mean), η διάμεσος (median) και η επικρατούσα τιμή (mode). Σημαντικά μέτρα διασποράς είναι η τυπική απόκλιση (std deviation), το εύρος (range), η ελάχιστη (min) και μέγιστη (max) τιμή και τα ποσοστιαία σημεία (percentile values).

Είναι σημαντικό να μην υπάρχει σύγχυση ανάμεσα στη μέση τιμή (τον μέσο όρο των μετρήσεων), τη διάμεσο (την τιμή που χωρίζει τα δεδομένα σε δύο ίσα μέρη έτσι ώστε το πλήθος των μετρήσεων που βρίσκονται αριστερά της να είναι ίσο με το πλήθος των μετρήσεων που βρίσκεται δεξιά της) και την επικρατούσα τιμή (την τιμή με τη μεγαλύτερη συχνότητα).

Άλλα περιγραφικά μέτρα είναι ο συντελεστής ασυμμετρίας (skewness) και ο συντελεστής κύρτωσης (kyrtosis), που περιγράφουν την ασυμμετρία και την κύρτωση της κατανομής αντίστοιχα. Με βάση την ασυμμετρία, οι κατανομές διακρίνονται σε συμμετρικές όταν l = 0 (σε αυτές ανήκει η κανονική κατανομή), σε θετικά ασύμμετρες (ή λοξές δεξιά) όταν l > 0 και σε αρνητικά ασύμμετρες (ή λοξές αριστερά) όταν l < 0. Με βάση την κύρτωση οι κατανομές διακρίνονται σε λεπτόκυρτες όταν k > 3, σε μεσόκυρτες όταν k = 3 (σε αυτές ανήκει η κανονική κατανομή) και σε πλατύκυρτες όταν k < 3. Το S.P.S.S. υπολογίζει ως κύρτωση την τιμή k - 3, έτσι ώστε ο έλεγχος της κανονικότητας των δεδομένων να γίνεται με το μηδέν.

Τέλος, τα τεταρτημόρια έχουν την ιδιότητα να χωρίζουν το σύνολο των μετρήσεων σε τέσσερα ίσα μέρη (0%, 25%, 50%, 75%, 100%).

Για την γραφική αναπαράσταση των ποσοτικών μεταβλητών, χρησιμοποιούμε το ιστόγραμμα συχνοτήτων και το θηκόγραμμα. Το ιστόγραμμα παρουσιάζει ομαδοποιημένα τα δεδομένα με τη μορφή συνεχόμενων

8

ορθογωνίων, τα οποία έχουν ύψος ανάλογο με τη συχνότητα κάθε ομάδας. Οι τιμές της μεταβλητής τοποθετούνται στον οριζόντιο άξονα, ενώ οι συχνότητες στον κατακόρυφο άξονα.

Το θηκόγραμμα παριστάνει τα σημαντικότερα ποσοστημόρια των δεδομένων και έχει τη μορφή κουτιού με κάτω και άνω άκρα τα 25% και 75% ποσοστημόρια αντίστοιχα. Επιπλέον, διαθέτει δύο προεκτάσεις (whiskers) οι οποίες βοηθούν στον εντοπισμό ακραίων παρατηρήσεων. Η κάτω γραμμή φτάνει μέχρι το (1ο τεταρτημόριο – 1.5 * ενδοτεταρτημοριακό εύρος) ενώ η πάνω μέχρι (3ο τεταρτημόριο + 1.5 * ενδοτεταρτημοριακό εύρος).

Έλεγχος Κανονικότητας

Το μεγαλύτερο μέρος της Στατιστικής Συμπερασματολογίας, προϋποθέτει ότι τα δεδομένα προέρχονται από έναν πληθυσμό ο οποίος περιγράφεται ικανοποιητικά από την κανονική κατανομή. Η κανονικότητα ή μη των δεδομένων θα κρίνει τα μέτρα θέσης και διασποράς που τα περιγράφουν ορθότερα. Τα κανονικά δεδομένα περιγράφονται καλύτερα από την μέση τιμή και την τυπική απόκλιση, ενώ τα μη-κανονικά δεδομένα από τη διάμεσο και τα τεταρτημόρια. Ο έλεγχος της κανονικότητας γίνεται τόσο στατιστικά όσο και γραφικά.

Για τον στατιστικό έλεγχο κανονικότητας χρησιμοποιούμε το τεστ Shapiro-Wilk (ή το τεστ Kolmogorov-Smirnov για μικρά δείγματα, με n < 30). Οι έλεγχοι κανονικότητας έχουν τη μορφή:

 H_0 : τα δεδομένα προέρχονται από κανονική κατανομή

 H_1 : η κανονικότητα απορρίπτεται

Επομένως, αν ο έλεγχος δώσει αρκετά μικρό p-value, η υπόθεση της κανονικότητας απορρίπτεται.

Γραφικά, η κανονικότητα ελέγχεται με το Q-Q (quantile-quantile) γράφημα, το οποίο συγκρίνει τα ποσοστιαία σημεία (quantile) του δείγματος έναντι των πληθυσμιακών ποσοστιαίων σημείων της κανονικής κατανομής. Παρεκκλίσεις από την ευθεία δηλώνουν μη κανονικότητα. Παράδειγμα: Θα ελέγξουμε την κανονικότητα των (ποσοτικών) μεταβλητών Acceleration και Price.

ta							cars.sav [DataSet1] - IBM SP	SS Statisti	cs Data Edito	r					- 0	×
<u>F</u> ile <u>E</u>	dit <u>V</u> i	iew <u>D</u> ata	Transform	Analyze Direct	Marketing	<u>G</u> raphs	Utilities	Add- <u>o</u> ns	Window	<u>H</u> elp								
		<u>_</u> m.		Reports		*	*					ABC						
				Descriptive St	tatistics	•	123 Frequ	encies			14 🔍							
83 : Acce	leration	22,1	20057581275	7 Ta <u>b</u> les		•	Descr	iptives								Visit	ole: 9 of 9 Varia	ables
- 10		Price	Engine	Co <u>m</u> pare Mea	ans	*	A Explor	е	Jel	Brand	Туре	var	var	var	var	var	var	
77	_	15152	1555	<u>G</u> eneral Linea	ar Model	*	Cross	tahe	2		1 0							
78	_	16939	1589	Generalized L	inear Mode.	ls 🕨	Rotio		2		3 1							
79	_	16230	1573	Mixed Models		•			2		2 0							
80	_	15115	1982	Correlate		*	P-P PI	015	2		2 0							
81	_	16539	2294	Regression		*	<u>0</u> -Q P	lots	1		1 1							
82		14899	1966	L <u>o</u> glinear		*	þ	2015	2		2 1							
83		14437	1589	Neural Net <u>w</u> o	orks	*	2	2016	2		2 0							
84		30051	6554	Classify		*	þ	2018	1		1 1							
85		23107	5211	Dimension R	eduction	•	1	2017	1		1 1							
86		24622	5735	Sc <u>a</u> le		•	1	2015	1		1 1							
87		22284	5899	<u>N</u> onparametri	ic Tests	•	1	2018	1		1 1							
88		28594	7210	Forecasting		•	1	2018	1		1 1							
89		29043	7456	Survival		•	1	2016	1		1 1							
90		21811	4981	Multiple Resp	onse	*	2	2016	1		1 1							
91		27410	7030	🔛 Missing Value	Analysis		1	2016	1		1 1							
92		20677	6554	Multiple Imput	tation	*	2	2016	1		1 1							
93		22203	5211	Complex Sam	nples	•	3	2016	1		1 1							
94		26173	5735	Simulation			3	2018	1		1 1							
95		24646	6554	Quality Contro			4	2017	1		1 1							
96		20385	5735				3	2016	1		1 1							
97		20492	5735	KOC Culve			3	2017	1		1 0							
98		22537	5735	175	1845		14	2016	1		1 1							
99		20980	5751	158	1963		13	2019	1		1 1							
100		24725	5899	170	2094		13	2017	1		1 1							Ţ
	4													-				•
Data Vi	ew Va	riable View																
Explore												IBN	I SPSS Statist	ics Processor	is ready	Unicode	ON	

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Descriptive Statistics \rightarrow Explore.

Μετακινούμε δεξιά τις μεταβλητές Acceleration και Price, και στο μενού Charts επιλέγουμε Histogram και Normality plots with tests για να εκτελέσουμε τον έλεγχο κανονικότητας (ο οποίος θα παράγει αυτόματα το boxplot και το Q-Q Plot).

ta	Explore	×	Explore: Plots
 Engine Horsepower Weight Year Fuel Brand Type 	Dependent List: Image: Acceleration Image: Price Factor List: Image: Dependent List: Image: Label Cases by:	Statistics Plots Options Bootstrap	Boxplots □ Descriptive ● Factor levels together □ Stem-and-leaf ● Dependents together □ Histogram ● None □ Histogram ▼ Normality plots with tests Spread vs Level with Levene Test ● None □ Power estimation ● Transformed Power Natural log ▼
Both ◎ Statistics ◎	Plots		© <u>U</u> ntransformed
ок	Paste Reset Cancel Help		Continue Cancel Help

Ξεκινάμε την περιγραφική ανάλυση με τα ιστογράμματα των δύο μεταβλητών. Η Acceleration είναι συμμετρική, αλλά παρουσιάζει μία απότομη πτώση σε κεντρικές τιμές. Η Price εμφανίζει έντονη δεξιά ασυμμετρία (οι τιμές συγκεντρώνονται στα αριστερά).



Συνεχίζουμε με τα boxplot, τα οποία βοηθούν στον εντοπισμό ακραίων τιμών (outliers). Εδώ η Price εμφανίζει μεγάλο πλήθος ακραίων παρατηρήσεων (οι οποίες αντιστοιχούν στις τιμές του ιστογράμματος που προκαλούν την ασυμμετρία, κοντά στο 30000). Αντίθετα, η Acceleration εμφανίζει μόλις μία ακραία παρατήρηση.



Στα γραφήματα Q-Q Plot μπορούμε να δούμε πόσο καλά εφαρμόζουν τα σημεία της μεταβλητής Acceleration πάνω στα ποσοστημόρια της κανονικής κατανομής, σε αντίθεση με αυτά της Price που παρουσιάζουν σοβαρές αποκλίσεις. Υποψιαζόμαστε ότι ο έλεγχος θα απορρίψει την υπόθεση της κανονικότητας για την Price.

Έχουμε ένα δείγμα 391 αυτοκινήτων, επομένως κοιτάμε τον έλεγχο Shapiro-Wilk. Για την μεταβλητή Acceleration έχουμε p-value = 0.259(>0.05) επομένως δεν απορρίπτουμε την (H_0) υπόθεση της κανονικότητας σε επίπεδο σημαντικότητας 5%. Αντίθετα για την μεταβλητή Price έχουμε p-value < 0.05 (η τιμή δεν είναι πραγματικά 0. Κάνοντας κλικ πάνω στον έλεγχο βλέπουμε ότι p-value = $6, 2 \cdot 10^{-11}$). Επομένως απορρίπτουμε την υπόθεση της κανονικότητας για την μεταβλητή Price.



Επομένως, τα σωστά περιγραφικά μέτρα για την Acceleration είναι η μέση τιμή και η τυπική απόκλιση, ενώ για την Price η διάμεσος, το ελάχιστο, το μέγιστο και τα τεταρτημόρια. Για να τα επιλέξουμε, στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze — Descriptive Statistics — Frequencies. Επιλέγουμε τα επιθυμητά μέτρα για κάθε μεταβλητή (φυσικά, μπορούμε να ζητήσουμε περισσότερα, όπως η κύρτωση και η ασυμμετρία).

ta Frequencies: Sta	atistics ×	ta Frequencies: St	atistics ×
Percentile Values Quartiles Cut points for: 10 equal groups Percentile(s): Add Change Remove	Central Tendency ✓ Mean Median Mode Sum Values are group midpoints	Percentile Values Cut points for: 10 equal groups Percentile(s): Add Change Remove	Central Tendency ☐ Mean ☑ Median ☐ Mode ☐ Sum Values are group midpoints
Dispersion	Distribution	Dispersion	Distribution
Std. deviation 🔲 Minimum	Ske <u>w</u> ness	Std. deviation 🖌 Minimum	🔲 Ske <u>w</u> ness
🔲 Variance 📃 Ma <u>x</u> imum	Kurtosis	🔲 Variance 🛛 Maximum	🔲 <u>K</u> urtosis
🔲 Ra <u>n</u> ge 📃 S. <u>E</u> . mean		🗖 Ra <u>n</u> ge 📃 S. <u>E</u> . mean	
Continue	Help	Continue	Help

Τελικά, τα περιγραφικά μέτρα των δύο μεταβλητών εμφανίζονται στους παρακάτω πίνακες.

Acceleration	n	Price	
Μέτρο	Τιμή	Μέτρο	Τιμή
Μέση Τιμή	15.62	Ελάχιστο (0%)	11013
Τυπική Απόκλιση	2.70	1ο Τεταρτημόριο (25%)	15282.88
		Διάμεσος (50%)	17444.97
		3ο Τεταρτημόριο (75%)	20243.15
		Μέγιστο (100%)	30738

Ασκήσεις

Επαναλάβετε τις παραπάνω περιγραφικές αναλύσεις για τα δεδομένα του αρχείου grades.

- 1. Αξιολογήστε τις μεταβλητές ως ποιοτικές/ποσοτικές (και τις ποιοτικές σε διατάξιμες/ονοματικές).
- 2. Δημιουργήστε ένα πίνακα συχνοτήτων και ένα ραβδόγραμμα για τις ποιοτικές μεταβλητές Hours και Class.
- 3. Εφαρμόστε έλεγχο κανονικότητας για τις ποσοτικές μεταβλητές Grades και Physics. Τι εικόνα δίνουν τα γραφήματα; Γιατί δεν απορρίπτεται η υπόθεση της κανονικότητας;

3 Συσχέτιση

Σύνοψη: Το επόμενο βήμα στην ανάλυση των δεδομένων είναι να εξεταστεί κατά πόσο υπάρχει συσχέτιση μεταξύ των μεταβλητών. Αυτό επιτυγχάνεται τόσο γραφικά, με το διάγραμμα διασποράς, όσο και υπολογιστικά με το συντελεστή γραμμικής συσχέτισης Pearson.

Διάγραμμα διασποράς

Σε πολλές περιπτώσεις της ανάλυσης δεδομένων το ενδιαφέρον εστιάζεται στη μελέτη δύο ή περισσότερων χαρακτηριστικών του δείγματος. Για το λόγο αυτό, είναι λογικό να αναζητήσουμε μέτρα τα οποία μπορούν να εκφράσουν και να ποσοτικοποιήσουν την πιθανή συμμεταβολή-συσχέτιση των χαρακτηριστικών αυτών. Καθώς η οπτικοποίηση των δεδομένων αποτελεί πρωταρχικό στάδιο της ανάλυσης, το διάγραμμα το οποίο απεικονίζει τη συσχέτιση ματαξύ δύο μεταβλητών είναι το **διάγραμμα διασποράς (Scatter Diagram) ή σημειόγραμμα**

Η εφαρμογή θα γίνει στο σετ δεδομένων cars. Πιο συγκεκριμένα, θα εξετάσουμε γραφικά τη συσχέτιση μεταξύ των μεταβλητών Acceleration και Horsepower. Κατά πόσο, δηλαδή, ο χρόνος που απαιτείται για να φτάσει το αυτοκίνητο ταχύτητα 100Km/h ξεκινώντας από την ηρεμία, εξαρτάται από την ιπποδύναμή του.

🕼 cars.sav (D	cars.sav (DataSett) - IBM SPSS Statistics Data Editor – 🖸															×				
Eile Edit	<u>V</u> iew <u>D</u> ata	Transform	Analyze G	raphs <u>U</u> tilitie	es Extension	s <u>W</u> indow	Help													
2] 🗠	a 🖺	<u>C</u> hart Builder Graphboard	 Template Choo	ser	<mark>์ (</mark>													
				Lenary Dialo	ine		D nu										Vis	ible: 9 of 9 V	ariable	łs
	🖋 Price	🛷 Engine	Horsepow er	Weight	Accelerati on	d Year	<u>B</u> ar <u>111</u> <u>3</u> -D Bar		💰 Туре	var										
1	21041	5571	160	1624	9	2018	🛃 Line		1											8
2	27057	6390	190	1732	9	2018	🔼 Area		1											
3	29409	7210	215	1940	11	2017	Nig		1											1
4	29548	7439	220	1959	10	2018	High-Low		1											
5	24075	6554	150	1692	10	2016	Boxplot		1											
6	30506	7456	225	1388	10	2017	Bogpiot		1											
7	25868	6276	170	1603	9	2016	Error Bar		1											
8	26746	7030	198	1953	10	2017	Population	Pyramid	1											
9	29731	7456	225	1991	10	2016	Scatter/Do	t	1											
10	21840	4949	140	1552	10	2016	🚹 Histogram		1											
11	22121	5211	150	1546	11	2016	1	1	1											
12	24087	5735	165	1661	13	2017	1 1	1	1											
13	20134	4981	150	1544	13	2017	1	1	1											
14	20467	5030	130	1576	12	2017	1	1	1											
15	19741	1982	113	1005	13	2017	2	2	1											1
16	29638	5211	210	1971	14	2018	1	1	1											1
17	28743	5899	215	2076	14	2018	1	1	1											1
18	17382	1589	88	958	14	2017	2	3	1											1
19	17488	1753	90	1093	14	2015	5 2	2	0											1
20	17622	1851	95	1067	15	2016	5 2	3	0											1
21	14879	3261	90	1192	15	2014	2	! 1	0											1
22	26506	5030	200	1969	15	2017	1	1	1											1
23	19278	3261	97	1248	16	2019	1	1	1											1
24	15596	3244	95	1275	16	2014	2	1	1											1
25	15950	3277	85	1164	17	2018	2	! 1	0											1
26	19539	1704	95	1068	18	2017	2	2	1											1
27	15689	1802	87	1202	18	2017	·	2	0											Ē
_									***				-			-		_		4
Date Mound	Variable Menu																			

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Graphs \rightarrow Legacy Dialogs \rightarrow Scatter/Dot.

Έπειτα, επιλέγουμε το Simple Scatter. Στο συγκεκριμένο παράδειγμα επιλέγουμε να εξετάσουμε γραφικά τη συσχέτιση μεταξύ δύο μεταβλητών. Αν θέλαμε να εξετάσουμε τη συσχέτιση παραπάνω από δύο ποσοτικών μεταβλητών, θα επιλέγαμε το Matrix Scatter.

ta Scatter/Dot	×
Simple Matrix Scatter Scatter Dot	ole
Overlay Scatter	
Define Cancel Help	

Στο αναδυόμενο παράθυρο, μετακινούμε στο κελί Y axis τη μεταβλητή Acceleration και στο κελί X axis τη μεταβλητή Horsepower, όπως φαίνεται παρακάτω.

Simple Scatterplot		×
 Price Engine Weight Year Fuel Brand Type 	Y Axis: Acceleration X Axis: Horsepower Set Markers by: Label Cases by: Label Cases by: Rows: Nest variables (no empty rows) Columns: Nest variables (no empty columns)	<u></u> Options
Template	ns from:	
ОК	E Paste Reset Cancel Help	

Το διάγραμμα διασποράς φανερώνει μία γενικά αρνητική συσχέτιση μεταξύ των δύο μεταβλητών, καθώς όσο αυξάνεται η ιπποδύναμη του αυτοκινήτου τόσο ελαττώνεται ο χρόνος της επιτάχυνσής του.



Συντελεστής γραμμικής συσχέτισης Pearson

Από το διάγραμμα διασποράς φαίνεται ότι οι δύο μεταβλητές είναι αρνητικά συσχετισμένες, αλλά πόσο ισχυρή είναι αυτή η συσχέτιση; Δεδομένου του γεγονότος ότι και οι δύο μεταβλητές είναι ποσοτικές, χρησιμοποιείται ο δειγματικός συντελεστής συσχέτισης του Pearson

$$r_{XY} = \frac{s_{XY}}{s_X \cdot s_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \cdot \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Ο συντελεστής r δίνει ένα μέτρο του μεγέθους της γραμμικής συσχέτισης μεταξύ δύο μεταβλητών. Παίρνει τιμές στο [-1,1] και το πρόσημό του υποδηλώνει αν η συσχέτιση είναι θετική ή αρνητική. Πιο συγκεκριμένα, αν:

- $r = \pm 1$ τότε η γραμμική συσχέτιση είναι τέλεια.
- −0.3 ≤ 0.3 δεν υπάρχει γραμμική συσχέτιση, χωρίς αυτό να σημαίνει ότι δεν υπάρχει κάποιου άλλου είδους συσχέτιση (όπως καμπυλόγραμμη συσχέτιση).
- $0.5 < r \le -0.3$ ή $0.3 \le r < 0.5$ υπάρχει ασθενής γραμμική συσχέτιση.
- $-0.7 < r \le -0.5$ ή $0.5 \le r < 0.7$ υπάρχει μέση γραμμική συσχέτιση.
- $-0.8 < r \le -0.7$ ή $0.7 \le r < 0.8$ υπάρχει ισχυρή γραμμική συσχέτιση.
- 1 < r ≤ -0.8 ή $0.8 \leq r < 1$ υπάρχει πολύ ισχυρή γραμμική συσχέτιση.

Ο συντελεστής Pearson χρησιμοποιείται κυρίως για ποσοτικές μεταβλητές και έχει πιο ακριβή εφαρμογή, όταν αυτές είναι κανονικά κατανεμημένες.

Η εφαρμογή στο SPSS γίνεται μεταξύ των μεταβλητών Acceleration και Horsepower.

tars.sav [Da	taSet1] - IBM SI	SS Statistics Da	ta Editor													-	- 0	×
Eile Edit	⊻iew Data	Transform	Analyze Graphs Utilities	Extensions	Window	Help												
a 🔳			Reports	•		A (
		• -	Descriptive Statistics			1												
			Bayesian Statistics													Vis	ible: 9 of 9 Var	riables
	🔗 Price	🛷 Engine	Ta <u>b</u> les		🛃 Year	🚮 Fuel	💰 Brand	💰 Туре	var									
	04044		Co <u>m</u> pare Means		0040													
1	21041	55/1	General Linear Model	•	2018	1	1	1										
2	27057	6390	Generalized Linear Models		2018		1	1										
3	29409	7210	Mixed Models		2017	1	1	1										
4	29540	7439	<u>C</u> orrelate	- F	Bivariate		1	1										
6	24075	7456	Regression		Partial		1	1										
7	25868	6276	L <u>og</u> linear		Distances		1	1										
8	26746	7030	Neural Networks		2017	1	1	1										
9	29731	7456	Classify		2016	1	1	1										
10	21840	4949	Dimension Reduction		2016	1	1	1										
11	22121	5211	Sc <u>a</u> le		2016	1	1	1										
12	24087	5735	Nonparametric Tests	•	2017	1	1	1										
13	20134	4981	Forecasţing		2017	1	1	1										
14	20467	5030	Survival	+	2017	1	1	1										
15	19741	1982	Multiple Response		2017	2	2	1										
16	29638	5211	💯 Missing Value Analysis		2018	1	1	1										
17	28743	5899	Multiple Imputation		2018	1	1	1										
18	17382	1589	Complex Samples		2017	2	3	1										
19	17488	1753	Simulation		2015	2	2	0										
20	17622	1851	Quality Control		2016	2	3	0										
21	14879	3261	ROC Curve		2014	2	1	0										
22	26506	5030	Spatial and Temporal Mode	ling 🕨	2017	1	1	1										
23	19278	3261	Direct Marketing		2019	1	1	1										
24	15596	3244			2014	2	1	1										_
25	15950	3277	85 1164	17	2018	2	1	0										
26	19539	1704	95 1068	18	2017	2	2	1										L L
27	15689	1802	87 1202	18	2017	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	2	0										
Data View	ariable View																	

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Correlate \rightarrow Bivariate.

Στο παράθυρο που εμφανίζεται παρακάτω μετακινούμε δεξιά τις μεταβλητές Acceleration και Horsepower, επιλέγουμε τον συντελεστή Pearson και τον αμφίπλευρο έλεγχο (two-tailed) στατιστικής σημντικότητας. Ο έλεγχος αυτός εξετάζει τη μηδενική υπόθεση ότι οι δύο μεταβλητές είναι ασυσχέτιστες έναντι της εναλλακτικής ότι υπάρχει στατιστικά σημαντική συσχέτιση μεταξύ τους.

 $H_0: \ \rho = 0$ $H_1: \ \rho \neq 0$

tai Bivariate Correlations		X							
 ✓ Price ✓ Engine ✓ Weight ✓ Year ✓ Fuel ✓ Brand ✓ Type 	Variables:	Options Style Bootstrap							
Correlation Coefficients -	tau-b 🥅 Spearman								
Test of Significance	Test of Significance								
✓ Elag significant correlations OK Paste Reset Cancel Help									

Ο πίνακας συσχέτισης που εξάγει το SPSS δείχνει μία μέση αρνητική γραμμική συσχέτιση μεταξύ των δύο μεταβλητών, καθώς ο συντελεστής Pearson είναι -0.684. Η συσχέτιση είναι στατιστικά σημαντική σε επίπεδο 1%, όπως φαίνεται και από την τιμή του p-value: Sig. << 0.01. Αυτό σημαίνει, ότι απορρίπτουμε την αρχική υπόθεση, ότι οι δύο μεταβλητές είναι ασυσχέτιστες.

	Correlations										
		Acceleration	Horsepower								
Acceleration	Pearson Correlation	1	- 684								
	Sig. (2-tailed)		,000								
	N	391	391								
Horsepower	Pearson Correlation	-,684	1								
	Sig. (2-tailed)	,000									
	Ν	391	391								

**. Correlation is significant at the 0.01 level (2-tailed).

4 Απλή Γραμμική Παλινδρόμηση

Σύνοψη: Αφού εξετάστηκε η συσχέτιση μεταξύ των μεταβλητών, ακολουθεί η εφαρμογή και ανάλυση της απλής γραμμικής παλινδρόμησης.

Εισαγωγή

Ο συντελεστής συσχέτισης μας δείχνει αν και κατά πόσο δύο μεταβλητές σχετίζονται, χωρίς όμως να μας παρέχει πληροφορίες για τον τρόπο με τον οποίο μεταβάλλονται οι τιμές της μίας σε σχέση με της άλλης. Ο στόχος της απλής γραμμικής παλινδρόμησης είναι να περιγράψει με μία μαθηματική σχέση τις παρατηρήσεις των δύο μεταβλητών.

Από το διάγραμμα διασποράς, προκύπτει ότι δεν υπάρχει κάποια καμπύλη να ενώνει όλα τα σημεία και για το λόγο αυτό αναζητείται εκείνη η καμπύλη που προσαρμόζεται καλύτερα στα δεδομένα. Η απλούστερη μορφή συσχέτισης δύο μεταβλητών είναι η γραμμική, η οποία εκφράζεται μαθηματικά από τη σχέση:

$$Y_i=\beta_0+\beta_1X_i+\epsilon_i,\quad i=1,..,n$$

Όπου:

- Υ είναι η εξαρτημένη μεταβλητή,
- Χ είναι η ανεξάρτητη μεταβλητή,
- β_0 είναι η τιμή της εξαρτημένης μεταβλητής για X = 0, χωρίς αυτή να έχει πάντα φυσική σημασία στην ανάλυση,
- β_1 είναι η κλίση της ευθείας της γραμμικής παλινδρόμησης και εκφράζει τη μεταβολή της Y για μοναδιαία μεταβολή της X και
- ϵ_i τα τυχαία σφάλματα, τα οποία στο απλό γραμμικό μοντέλο υποτίθενται ομοσκεδαστικά με $\epsilon_i \sim N(0, \sigma^2)$

Για να βρεθεί η ευθεία η οποία προσαρμόζεται καλύτερα στα δεδομένα πρέπει να εκτιμηθούν οι σταθερές β_0 και β_1 . Οι εκτιμημένες παράμετροι υπολογίζονται μέσω της μεθόδου ελαχίστων τετραγώνων (Μ.Ε.Τ.) των σφαλμάτων, όπως έχει ειπωθεί αναλυτικά και στη θεωρία.

Εφαρμογή του μοντέλου

Η εφαρμογή της απλής γραμμικής παλινδρόμησης πραγματοποιείται στο σετ δεδομένων cars, ως συνέχεια της προηγούμενης ενότητας. Θα θεωρήσουμε ως εξαρτημένη μεταβλητή (Y) την Acceleration και ως ανεξάρτητη (X) την Horsepower.

🍓 cars.sav [Da	carssav (Datšett) - IBM SPSS Statistics Data Editor – 🗗 X																			
<u>F</u> ile <u>E</u> dit	<u>V</u> iew <u>D</u> ata	Transform	Analyze Graphs	s <u>U</u> tilities E <u>x</u> l	ensions	Window	Help													
2			Reports Descriptive Sta	atistics	*		(0												
			Bayesian Stati	tistics														Visi	ble: 9 of 9 Var	riables
	🔗 Price	🖋 Engine	Tables			J Year	🚮 Fuel	💰 Bran	d	💰 Туре									1000	
	00077	0.554	Compare Mea	ans		0010					Var	VdI	Val	VdI	VdI	VdI	VdI	VdI	Val	
92	20677	6554	General Linea	ar Model		2016		1	1	1										-
93	22203	5211	Generalized Li	inear Models		2016		1	1	1										
94	26173	5/35	Mixed Models			2018		1	1	1										
95	24040	6726	Correlate			2017		1	1	1										
90	20305	5735	Regression		•	Automa	lic Linear Mod	laling	1	0										
98	20432	5735	Loglinear			Linear	uc Elliear mou	iennig	1	1										
99	20980	5751	Neural Networ	rks		Linear			1	1										
100	24725	5899	Classify			Curve E	stimation		1	1										
101	16950	1147	Dimension Re	eduction		📸 Partial L	.ea <u>s</u> t Squares	i	3	1										
102	18958	2556	Scale		*	🔛 Binary L	ogistic		3	0										
103	17356	2540	Nonparametri	ic Tests		Multinor	nial Logistic		1	1										
104	17064	1868	Forecasting			🔛 Or <u>d</u> inal.			2	1										
105	18578	1982	Survival			Probit			2	1										
106	23341	6554	Multiple Resp	onse		Nonline	ar		1	1										
107	21196	4949	Missing Value	Analysis		Weight I	Estimation		1	1										
108	23399	5211	Multiple Imput	tation		2.Stage	Loget Square		1	1										
109	19535	3801	Complex Sam	nles		E-Stage	Cease Square		1	1										
110	14663	1900	Rimulation	ipico.		Optimal	Scaling (CAT	REG)	2	0										
111	16911	1605	Ouelity Centre			2017	:	2	2	1										
112	18600	1982		л	,	2018		1	2	1										
113	17390	3244	ROC Curve			2017	1	2	1	1										
114	17830	3801	Spatial and Te	emporal Modeling.	. !	2017		1	1	0										
115	16909	1769	Direct Marketin	ng	•	2014		2	3	0										_
116	18296	4096	88	1359	17	2016		1	1	1										
117	18022	3687	105	1404	17	2016		1	1	1										
118	18224	4096	100	1475	18	2017		1	1	1								-		
D	(orighte View)									***										
Data View	anable view																			

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Regression \rightarrow Linear.

Στο παράθυρο που εμφανίζεται μετακινούμε δεξιά στο κελί Dependent την εξαρτημένη μεταβλητή Acceleration και στο κελί Independent(s) την ανεξάρτητη μεταβλητή Horsepower. Στο πεδίο Method μπορούμε να επιλέξουμε τη μέθοδο επιλογής του βέλτιστου μοντέλου, η οποία έχει νόημα μόνο όταν οι ανεξάρτητες μεταβλητές είναι παραπάνω από μία, οπότε αφήνουμε την προεπιλογή Enter ως έχει.

Στη συνέχεια, στην καρτέλα **Statistics** επιλέγουμε Estimates, Confidence Intervals και Model Fit.

ta Linear Regression	×	
Linear Regression	X Statistics Plots Save Qptions Style Bootstrap	Linear Regression: Statistics × Regression Coefficien Cigatimates Confidence intervals Level(%): 95 Cogariance matrix Residuals Dyrbin-Watson Casewise diagnostics @ Outliers outside: 3 standard deviations
WLS Weight		
OK Paste Reset Cancel Help		Cancel Help

Από την καρτέλα **Save** παρακάτω μας ενδιαφέρουν τα Predicted values και τα Residuals. Από τα Predicted values επιλέγουμε τα Standardized και τα Unstandardized, ενώ από τα Residuals επιλέγουμε τα Standardized. Στη συνέχεια πατάμε Continue και τέλος OK.

tinear Regression: Save	×								
Predicted Values Unstandardized Standardized Adjusted S.E. of mean predictions	Residuals Unstandardized Standardized Studentized Deleted Studentized deleted								
Distances Mahalanobis Cook's Leverage values Prediction Intervals Mean Individual Qonfidence Interval: 95 %	Influence Statistics DfBeta(s) Standardized DfBeta(s) DfFit Standardized DfFit Covariance ratio								
Coefficient statistics Create coefficient statistics Create a new dataset Dataset name: Write a new data file File									
Export model information to XML file									

Αποτελέσματα

Το ενδιαφέρον εστιάζεται στους πίνακες Model Summary, ANOVA και Coefficients.

Πίνακας Model Summary

Στον παρακάτω πίνακα η πρώτη τιμή αντιστοιχεί στο συντελεστή προσδιορισμού R του Pearson, ο οποίος αναλύθηκε στην προηγούμενη ενότητα. Καλύτερη φυσική ερμηνεία της συσχέτισης δύο μεταβλητών επιτυγχάνεται με το συντελεστή προσδιορισμού R^2 , ο οποίος παίρνει τιμές στο [0,1] και δίνεται από τον τύπο

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SSR}{SST}$$
όπου $SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$, $SSE = \sum_{i=1}^n (\hat{Y}_i - \hat{Y}_i)^2$ και $SST = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$

Η τιμή του συντελεστή προσδιορισμού στο συγκεκριμένο παράδειγμα είναι 0.468 ή περίπου 47%. Αυτό σημαίνει ότι περίπου το 47% της μεταβλητότητας της μεταβλητής Acceleration εξηγείται από την Horsepower και το 53% από τα τυχαία σφάλματα.

Model Summary ^b										
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate						
1	,684 ^a	,468	,467	1,969						

a. Predictors: (Constant), Horsepower

b. Dependent Variable: Acceleration

Πίνακας ΑΝΟVΑ

Η πρώτη στήλη του πίνακα ANOVA δείχνει τα τετραγωνικά αθροίσματα SSR = 1.328, 213, SSE = 1.508, 377 και SST = 2.836, 590. Οι βαθμοί ελευθερίας για το συγκεκριμένο μοντέλο είναι p = 1 για την παλινδρόμηση, n = p - 1 = 391 - 1 - 1 = 389 για τα κατάλοιπα και n - 1 = 390 για το συνολικό τετραγωνικό άθροισμα. Πραγματοποιείται ο παρακάτω έλεγχος (F-test) για την καλή προσαρμογή του μοντέλου στα δεδομένα:

$$H_0: \ \beta_1 = 0$$
$$H_1: \ \beta_1 \neq 0$$

Απορρίπτουμε τη μηδενική υπόθεση αν:

$$F = \frac{MSR}{MSE} > \mathcal{F}_{\nu_1,\nu_2}(\alpha)$$

όπου $\nu_1 = 1$, $\nu_2 = 389$ και επίπεδο στατιστικής σημαντικότητας $\alpha = 0.05$. Ωστόσο η τιμή της ελεγχοσυνάρτησης *F* δεν είναι πληροφοριακή από μόνη της, καθώς πρέπει να γνωρίζουμε και το άνω α - σημείο της κατανομής *F* με ν_1 και ν_2 βαθμούς ελευθερίας αντίστοιχα.

Για το λόγο αυτό καταφεύγουμε στην τελευταία στήλη του πίνακα ANOVA, όπου η μηδενική υπόθεση απορρίπτεται αν η τιμή του Sig. -που αντιστοιχεί στο p-value- είναι μικρότερη του 0.05. Στο παράδειγμά μας είναι Sig. << 0.05, οπότε συμπεραίνουμε ότι η παράμετρος β_1 είναι στατιστικά σημαντική και επομένως η ύπαρξη της μεταβλητής Horsepower έχει νόημα στο μοτέλο.

			ANOVAª			
Mode	el	Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1328,213	1	1328,213	342,537	,000 ^b
	Residual	1508,377	389	3,878		
	Total	2836,590	390			
a.	Dependent Varial	ole: Acceleration				
b.	Predictors: (Cons	tant). Horsepowe	r			

Πίνακας Coefficients

Ο πίνακας Coefficients είναι τελικά αυτός που δίνει δομή στην εξίσωση της απλής γραμμικής παλινδρόμησης, καθώς περίεχει τις εκτιμημένες παραμέτρους, τη στατιστική σημαντικότητά τους αλλά και τα διαστήματα εμπιστοσύνης.

				Coefficients ^a					
		Unstandardize	d Coefficients	Standardized Coefficients			95,0% Confidence Interval for B		
Mode	1	В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	20,647	,289		71,385	,000	20,078	21,215	
	Horsepower	-,048	,003	-,684	-18,508	,000	-,053	-,043	
a.	Dependent Variak	le: Acceleration							

Η πρώτη στήλη του πίνακα αφορά τις εκτιμήσεις των παραμέτρων.

Ερμηνεία β₀: Η τιμή του β₀ είναι 20, 647 και δηλώνει την τιμή της εξαρτημένης μεταβλητής όταν η ανεξάρτητη παίρνει την τιμή μηδέν. Ωστόσο, στη συγκεκριμένη περίπτωση δεν έχει φυσικό νόημα, αφού ένα αυτοκίνητο είναι αδύνατο να ξεκινήσει όταν δεν έχει ιπποδύναμη (!) και προφανώς δε θα έχει και επιτάχυνση.

Ερμηνεία β₁: Η τιμή του β₁ είναι -0.048 και είναι η αρνητική κλίση της ευθείας της γραμμικής παλινδρόμησης. Υποδηλώνει ότι για μοναδιαία αύξηση της ιπποδύναμης, ο χρόνος που απαιτείται για να φτάσει το αυτοκίνητο το 100km/h μειώνεται κατά 0.048 χρονικές μονάδες.

Ο πίνακας Coefficients περιλαμβάνει και τους ελέγχους στατιστικής σημαντικότας του μοντέλου για τις δύο παραμέτρους. Ο έλεγχος αυτός (t-test) είναι ισοδύναμος με τον έλεγχο F-test που αναφέρθηκε παραπάνω, αλλά μόνο στην απλή γραμμική παλινδρόμηση. Ο έλεγχος είναι ο εξής:

$$H_0: \ \beta_i = 0$$
$$H_1: \ \beta_i \neq 0$$

για i = 0, 1. Απορρίπτουμε την αρχική υπόθεση αν $|T| > t_{n-p-1,1-\frac{\alpha}{2}}$, όπου $T = \frac{\beta_i - 0}{s.e(\beta_i)}$ και $t_{n-p-1,1-\frac{\alpha}{2}}$ η κριτική τιμή της κατανομής Student-t με n - p - 1 β.ε. και επίπεδο σημαντικότητας α . Στο παράδειγμά μας έχουμε n = 391, p = 1 και $\alpha = 0.05$.

Ένας ευκολότερος τρόπος για να εξετάσουμε αν απορρίπτεται η μηδενική υπόθεση είναι να ελέγξουμε αν το p-value είναι μικρότερο από το επίπεδο σημαντικότητας (5%). Παρατηρούμε από την πέμπτη στήλη του πίνακα ότι *Sig.* << 0.05 και για τις δύο μεταβλητές, οπότε έχει νόημα η ύπαρξή τους στο μοντέλο.

Ένας άλλος τρόπος να ελέγξουμε αν οι παράμετροι είναι στατιστικά σημαντικές είναι να δούμε αν το Ο συμπεριλαμβάνεται στο 95% διάστημα επιστοσύνης. Και στις δύο περιπτώσεις το 0 δεν ανήκει στο διάστημα, οπότε καταλήγουμε στο ίδιο συμπέρασμα.

Τελικά, η εκτιμημένη εξίσωση της απλής γραμμικής παλινδρόμησης είναι:

$$(Acceleration) = 20.647 - 0.048 \cdot (Horsepower)$$

Έλεγχος υποθέσεων γραμμικού μοντέλου

Έλεγχος κανονικότητας: Για να ελέγξουμε αν τα κατάλοιπα του μοντέλου ακολουθούν κανονική κατανομή, ακολουθούμε την παρακάτω διαδικασία:

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Descriptive Statistics \rightarrow Explore.

Στο παράθυρο που εμφανίζεται στο κελί Dependent list μετακινούμε δεξιά τα Standardized Residuals τα οποία είχαμε αποθηκεύσει νωρίτερα. Στην καρτέλα Statistics επιλέγουμε Normality plots with tests.

Προκύπτει από τον παρακάτω πίνακα ότι η αρχική υπόθεση της κανονικότητας δεν μπορεί να απορριφθεί σε επίπεδο στατιστικής σημαντικότητας $\alpha = 5$ %, καθώς σύμφωνα με τους ελέγχους Kolmogorov-Smirnov και Shapiro-Wil Sig > 0.05.

Τέλος, από το Q-Q plot των Standardized residuals δεν φαίνεται γενικά να έχουμε πολύ μεγάλες αποκλίσεις από τα ποσοστημόρια της κανονικής κατανομής, αν εξαιρέσουμε κάποιες αποκλίσεις στο άνω δεξί τμήμα του γραφήματος.



Ομοσκεδαστικότητα και Ανεξαρτησία: Παρότι υπάρχουν στατιστικοί έλεγχοι για την ομοσκεδαστικότητα και την ανεξαρτησία των καταλοίπων, εδώ θα αρκεστούμε απλά σε ένα γραφικό έλεγχο. Συγκεκριμένα θα δημιουργήσουμε το σημειόγραμμα (scatterplot) των τυποποιημένων καταλοίπων (standardized residuals) ως προς τις τυποποιημένες προβλεπόμενες τιμές της γ (standardized predicted values).



Αρχικά, εξετάζουμε την ανεξαρτησία των καταλοίπων. Δεν παρατηρούμε τάσεις (όπως για παράδειγμα πολλά κατάλοιπα συγκεντρωμένα στην θετική μεριά, έπειτα πολλά συγκεντρωμένα στα αρνητικά κ.ο.κ.). Επομένως υποθέτουμε με μία σχετική βεβαιότητα ότι τα κατάλοιπα είναι ανεξάρτητα. Προχωρώντας όμως στην υπόθεση της ομοσκεδαστικότητας (σταθερής διασποράς), παρατηρούμε ότι το γράφημα δεν έχει τη μορφή ενός τυχαίου συννέφου γύρω από την ευθεία y = 0. Τα κατάλοιπα παρουσιάζουν ετεροσκεδαστικότητα, δείχνουν δηλαδή να μην κρατάνε σταθερή απόσταση από το 0 καθ' όλη τη διάρκεια του δείγματος, αλλά να αυξάνονται καθώς αυξάνονται και οι τιμές στον άξονα x. Τα παραπάνω υποδεικνύουν ότι πρέπει να δοκιμάσουμε εναλλακτικές μεθόδους μοντελοποίησης ώστε να έχουμε ομοσκεδαστικότητα.

Ασκήσεις

Κατασκευάστε απλά γραμμικά μοντέλα για τα δεδομένα του αρχείου grades.

- 1. Επιλέξτε ως εξαρτημένη μεταβλητή (y) την Grades και σαν ανεξάρτητη (x) την Maths.
- 2. Επιλέξτε ως εξαρτημένη μεταβλητή (y) την Grades και σαν ανεξάρτητη (x) την Physics.
- Ανταλλάξτε την εξαρτημένη μεταβλητή με την ανεξάρτητη. Επιλέξτε δηλαδή ως εξαρτημένη μεταβλητή
 (γ) την Maths και σαν ανεξάρτητη (x) την Grades. Συγκρίνετε τα μοντέλα (1) και (3). Θα μπορούσατε
 από το ένα να συμπεράνετε το άλλο;

5 Πολλαπλή Γραμμική Παλινδρόμηση

Σύνοψη: Αφού εξετάστηκε η συσχέτιση μεταξύ των μεταβλητών και κατασκευάστηκαν κατάλληλα μοντέλα απλής γραμμικής παλινδρόμησης, ακολουθεί η εφαρμογή και ανάλυση της πολλαπλής γραμμικής παλινδρόμησης.

Εισαγωγή

Συχνά οι τιμές μιας μεταβλητής ενδιαφέροντος εξαρτώνται από τις τιμές περισσότερων από μίας επεξηγηματικών μεταβλητών. Όταν μάλιστα αυτή η σχέση εξάρτησης είναι γραμμική μπορεί να εκφραστεί μέσα από ένα πολλαπλό γραμμικό μοντέλο. Ένα τέτοιο μοντέλο έχει τη μορφή:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \ldots + \beta_p X_{pi} + \epsilon_i, \quad i = 1, .., n$$

Όπου:

- Υ είναι η εξαρτημένη μεταβλητή,
- $X_j, j = 1, ..., p$ είναι η j-οστή ανεξάρτητη μεταβλητή,
- β₀ είναι η τιμή της εξαρτημένης μεταβλητής όταν όλες οι επεξηγηματικές μεταβλητές πάρουν την τιμή 0, χωρίς αυτή να έχει πάντα φυσική σημασία στην ανάλυση,
- $\beta_j, \ j = 1, ..., p$ εκφράζει τη μεταβολή της Y για μοναδιαία μεταβολή της X_j με την προϋπόθεση ότι οι τιμές των υπόλοιπων επεξηγηματικών μεταβλητών παραμένουν σταθερές και
- ϵ_i τα τυχαία σφάλματα, τα οποία υποθέτουμε πως είναι κανονικά κατανεμημένα και ομοσκεδαστικά με $\epsilon_i \sim N(0, \sigma^2)$

Σκοπός μας είναι η εκτίμηση των συντελεστών β_j , j = 0, 1, ..., p. Οι εκτιμήσεις των παραμέτρων υπολογίζονται μέσω της μεθόδου ελαχίστων τετραγώνων (Μ.Ε.Τ.), όπως έχει αναλυθεί στη θεωρία.

Εφαρμογή του Πλήρους Μοντέλου

Η εφαρμογή της πολλαπλής γραμμικής παλινδρόμησης πραγματοποιείται στο σετ δεδομένων cars, ως συνέχεια της προηγούμενης ενότητας. Θα θεωρήσουμε ως εξαρτημένη μεταβλητή (Y) την Acceleration και θα διερευνήσουμε το αν και πώς αυτή εξαρτάται από τις υπόλοιπες ποσοτικές μεταβλητές που έχουμε στα χέρια μας (Engine- X_1 , Price- X_2 , Weight- X_3 , Horsepower- X_4) συνολικά.

Αρχικά παίρνουμε μια πρώτη ιδέα για τις σχέσεις μεταξύ των μεταβλητών οπτικοποιώντας τα δεδομένα που διαθέτουμε όπως είδαμε στην ενότητα 3. Στη γραμμή εντολών ακολουθούμε τη διαδρομή Graphs → Legacy Dialogs → Scatter/Dot. Και αυτή τη φορά επιλέγουμε το Matrix Scatter πατώντας το κουμπί Define. Στο αναδυόμενο παράθυρο, μετακινούμε τις μεταβλητές που μας ενδιαφέρουν, εδώ Acceleration, Engine, Price, Weight και Horsepower, στο κελί Matrix Variables.

Catterplot Matrix		×
Fuel Brand Type	Matrix Variables: Acceleration Price Engine Weight Horsepower	Titles Options
	Set Markers by: Label Cases by: Panel by	
	Rows:	
	Nest variables (no empty columns)	
Template Use chart specifications from: File		

Κάνοντας διπλό κλικ στο γράφημα που προκύπτει μπορούμε να επεξεργαστούμε τα μεγέθη, τα χρώματα και ό,τι άλλο μας ενδιαφέρει. Το γράφημα που παίρνουμε είναι το ακόλουθο:

Acceleration					
Price					and the second second
Engine		ŵ.			in the second
Weight					J.
Horsepower		and the second s	and the second sec		
	Acceleration	Price	Engine	Weight	Horsepower

Στη συνέχεια προχωράμε με την εφαρμογή του μοντέλου. Ένα βασικό ζήτημα που καλούμαστε να αντιμετωπίσουμε είναι το ποιες επεξηγηματικές μεταβλητές θα συμπεριλάβουμε στην ανάλυσή μας. Αρχικά τις εισάγουμε όλες και θα συζητήσουμε στη συνέχεια περισσότερα.

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze \rightarrow Regression \rightarrow Linear όπως και στην απλή γραμμική παλινδρόμηση.

tars.sav [D	ataSet1] - IBM SI	PSS Statistics Da	ata Editor													-	- 0	×
<u>F</u> ile <u>E</u> dit	<u>V</u> iew <u>D</u> ata	Transform	Analyze Graphs Utilities	Extensions	Window	Help												
a b		, <u>e</u> 1	Reports Descriptive Statistics	+ +														
			Bayesian Statistics													Vis	ible: 9 of 9 Va	ariables
	🔗 Price	🛷 Engine	Ta <u>b</u> les		J Year	Fuel	💰 Brand	💰 Туре	var	Var	Var	var	var	var	var	var	var	
0.2	20677	6554	Compare Means		2016	1			, Tui	Yui	Yui	Yui	Yui	Yui	Tui	Yui	Yui	
93	22203	5211	General Linear Model		2010	1		1 .										- 1
0.4	26173	5735	Generalized Linear Models	s 🕨	2018	1		1 .										+
95	24646	6554	Mixed Models	*	2010	1		1 .										-1
96	20385	5735	<u>C</u> orrelate		2016	1		1 1										
97	20492	5735	Regression	•	Automat	tic Linear Model	ing	1 (
98	22537	5735	Loglinear		Linear			1 1										
99	20980	5751	Neural Networks		Cupia E	etimation		1 1										
100	24725	5899	Classify	*	Dertield	sunduon		1 1										
101	16950	1147	Dimension Reduction	•	Partial L	.east Squares		3 1										
102	18958	2556	Sc <u>a</u> le	•	📔 Binary L	.ogistic		3 (
103	17356	2540	Nonparametric Tests		Multinon	nial Logistic		1 1										
104	17064	1868	Forecasting		🔣 Orginal.			2 1										
105	18578	1982	Survival	*	Probit			2 1										
106	23341	6554	Multiple Response	*	Nonline	ar		1 1										
107	21196	4949	Missing Value Analysis		Weight 8	Estimation		1 1										
108	23399	5211	Multiple Imputation	•	2-Stane	Least Squares		1 1										
109	19535	3801	Complex Samples		Ontimal	Cooling (CATD		1 1										
110	14663	1900	Simulation		Opumai	Scaling (CATR	EG)	2 (
111	16911	1605	Quality Control		2017	2	:	2 1										
112	18600	1982		,	2018	1	:	2 1										
113	17390	3244		dellar b	2017	2		1 1										
114	17830	3801	Spatial and Temporal Mod	seling P	2017	1		1 (- 1
115	16909	1769	Direct Marketing	,	2014	2		3 (- 1
116	18296	4096	88 1359	17	2016	1		1 1										+
117	18022	3687	105 1404	17	2016	1		1 1										<u> </u>
118	18224	4096	100 1475	18	2017	1		1 1										

Data View	Variable View																	

Στο παράθυρο που εμφανίζεται μετακινούμε δεξιά, στο κελί Dependent, την εξαρτημένη μεταβλητή Acceleration και στο κελί Independent(s) τις ανεξάρτητες μεταβλητές Engine, Price, Weight και Horsepower. Στο πεδίο Method μπορούμε να επιλέξουμε τη μέθοδο επιλογής του βέλτιστου μοντέλου στην οποία αναφερθήκαμε παραπάνω. Αφήνοντας την προεπιλογή Enter δεν γίνεται επιλογή των μεταβλητών και συνεπώς εισάγονται όλες στο προς εκτίμηση μοντέλο. Στη συνέχεια, στην καρτέλα **Statistics** επιλέγουμε Estimates, Confidence Intervals και Model Fit.

🔚 Linear Regression		×		
 ✓ Price ✓ Engine ✓ Horsepower ✓ Weight ✓ Year Fuel ✓ Brand ✓ Type 	Dependent:	Statistics Plots Save Options Bootstrap	Linear Regression: Statistics	× ons
	Cancer Theip			

Από την καρτέλα **Save** παρακάτω μας ενδιαφέρουν τα Predicted values και τα Residuals. Από τα Predicted values επιλέγουμε τα Standardized και τα Unstandardized, ενώ από τα Residuals επιλέγουμε τα Standardized. Στη συνέχεια πατάμε Continue και τέλος ΟΚ.

tinear Regression: Save		\times					
Predicted Values Unstandardized Standardized Adjusted S.E. of mean predictions	Residuals Unstandardized Standardized Studentized Deleted Studentized deleted						
Distances Mahalanobis Cook's Leverage values Prediction Intervals Mean Individual Confidence Interval: 95 %	Influence Statistics DfBeta(s) Standardized DfBeta(s) DfFit Standardized DfFit Covariance ratio						
Coefficient statistics Create a new dataset Dataset name: Write a new data file File							
Export model information to XML file Browse Include the covariance matrix Continue Cancel Help							

Έτσι ζητάμε την εκτίμηση των συντελεστών, τα αντίστοιχα διαστήματα εμπιστοσύνης και άλλες ποσότητες που μας ενδιαφέρουν για το παρακάτω μοντέλο:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \epsilon_i, \quad i = 1, ..., 391$$

Αποτελέσματα για το Πλήρες Μοντέλο

Το ενδιαφέρον, και εδώ, εστιάζεται στους πίνακες Model Summary, ANOVA και Coefficients.

Πίνακας Model Summary

Στον παρακάτω πίνακα η πρώτη τιμή αντιστοιχεί στο συντελεστή προσδιορισμού R του Pearson, η δεύτερη στο συντελεστή προσδιορισμού R^2 , οι οποίοι αναλύθηκαν σε προηγούμενη ενότητα. Η τρίτη στήλη αντιστοιχεί στην τιμή του προσαρμοσμένου συντελεστή προσδιορισμού R^2_{adj} , ο οποίος παίρνει τιμές και αυτός στο [0,1] και λαμβάνει υπόψη του, επιπλέον, τον αριθμό των επεξηγηματικών μεταβλητών του εκάστοτε μοντέλου, δίνοντας μια ποινή για αύξηση του αριθμού αυτού. Αυτό έχει ως αποτέλεσμα σε περίπτωση που εισάγουμε, για παράδειγμα, στο μοντέλο μια μεταβλητή που δεν εξηγεί σημαντικό ποσοστό της εναπομένουσας μεταβλητότητας, να μειωθεί η τιμή του R^2_{adj} υπονοώντας ότι το μοντέλο χωρίς την τελευταία προσθήκη είναι καταλληλότερο. Η τελευταία ιδιότητα έρχεται σε αντίθεση με τη συμπεριφορά του συντελεστή R^2 , ο οποίος πάντα αυξάνεται με την εισαγωγή επιπλέον επεξηγηματικών μεταβλητών. Συνεπώς, ο προσαρμοσμένος συντελεστή προσδιορισμού είναι καταλληλότερος για τη σύγκριση μοντέλων με διαφορετικό αριθμό επεξηγηματικών μεταβλητών.

$$R_{adj}^2 = 1 - \frac{n-1}{n-p-1} \frac{SSE}{SST}$$

όπου $SSE = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2, \ SST = \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2$, η το πλήθος των παρατηρήσεων και ρ το πλήθος των επεξηγηματικών μεταβλητών.

Η τιμή του R^2_{adj} είναι 0.601 ή περίπου 60%, που σε σχέση με το απλό γραμμικό μοντέλο που είχαμε R^2_{adj} ίσο με 0.467 είναι μεγαλύτερη δηλώνοντας καλύτερη προσαρμογή.

······································									
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate					
1	,778 ^a	,606	,601	1,703					

Model Summary

a. Predictors: (Constant), Horsepower, Weight, Price, Engine

Πίνακας ΑΝΟVΑ

Στην πρώτη στήλη του πίνακα ANOVA μπορούμε να δούμε τα τετραγωνικά αθροίσματα SSR, SSE και SST. Οι βαθμοί ελευθερίας για το συγκεκριμένο μοντέλο είναι p = 4 για την παλινδρόμηση, n - p - p

1 = 391 - 4 - 1 = 386 για τα κατάλοιπα και n - 1 = 390 για το συνολικό τετραγωνικό άθροισμα. Πραγματοποιείται ο παρακάτω έλεγχος στατιστικής σημαντικότητας για το μοντέλο (F-test):

$$H_0: \ eta_1=eta_2=eta_3=eta_4=0$$

 $H_1: \ eta_j
eq 0$ για ένα τουλάχιστον $j, \ j=1,2,3,4$

Η H_0 αντιστοιχεί στη μη ύπαρξη εξάρτησης της εξαρτημένης μεταβλητής από καμία από τις ανεξάρτητες και επομένως δεν υπάρχει μοντέλο για να εκτιμήσουμε.

Απορρίπτουμε τη μηδενική υπόθεση αν το p-value είναι μικρότερο από το επιθυμητό επίπεδο στατιστικής σημαντικότητας, συνήθως 0.05. Όπως είδαμε και στην περίπτωση της απλής γραμμικής παλινδρόμησης, το p-value θα το βρούμε στον πίνακα ANOVA στην τελευταία στήλη με τίτλο Sig. Στην περίπτωσή μας είναι Sig<< 0.05, επομένως απορρίπτεται η μηδενική υπόθεση και έχει νήμα να προχωρήσουμε στην ανάλυσή μας.

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1717,662	4	429,416	148,137	,000 ^b
	Residual	1118,928	386	2,899		
	Total	2836,590	390			

a. Dependent Variable: Acceleration

b. Predictors: (Constant), Horsepower, Weight, Price, Engine

Πίνακας Coefficients

Ο πίνακας Coefficients είναι τελικά αυτός που δίνει δομή στην εξίσωση της παλινδρόμησης, καθώς περιέχει τις εκτιμημένες παραμέτρους, τη στατιστική σημαντικότητά τους αλλά και τα αντίστοιχα διαστήματα εμπιστοσύνης.

Πριν ασχοληθούμε όμως, με τις εκτιμήσεις των παραμέτρων πρέπει να ελέγξουμε τη στατιστική σημαντικότητα για κάθε έναν από τους συντελεστές του μοντέλου. Αυτό γίνεται με τη πραγματοποίηση τεσσάρων διαφορετικών t-tests (για *i* = 1, ..., 4), ελέγχους της μορφής:
$$H_0: \ \beta_i = 0$$
$$H_1: \ \beta_i \neq 0$$

Την τιμή της κάθε ελεγχοσυνάρτησης θα τη βρούμε στη στήλη του πίνακα Coefficients, με τίτλο t, ενώ το αντίστοιχο p-value, που είναι και αυτό που μας ενδιαφέρει για τη λήψη απόφασης, θα το βρούμε στη διπλανή στήλη με τίτλο Sig. Όπως μπορούμε να δούμε όλοι οι συντελεστές προκύπτουν στατιστικά σημαντικοί, εκτός από εκείνον τη μεταβλητής Price.

				Standardized				
		Unstandardized Coefficients		Coefficients			95,0% Confiden	ce Interval for B
Model		В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound
1	(Constant)	17,288	,761		22,707	,000	15,791	18,785
	Price	-7,015E-006	,000	-,010	-,097	,922	,000	,000,
	Engine	-,001	,000	-,341	-3,273	,001	-,001	,000,
	Weight	,006	,001	,915	10,116	,000	,005	,008
	Horsepower	-,082	,009	-1,158	-9,504	,000	-,098	-,065

Coefficients ^a

a. Dependent Variable: Acceleration

Επιλογή Μοντέλου

Σε αυτό το σημείο έφτασε η ώρα να επανέλθουμε στο θέμα της επιλογής επεξηγηματικών μεταβλητών που αναφέρθηκε στην αρχή της ενότητας. Ένας τρόπος να προχωρήσουμε, είναι να αφαιρέσουμε "χειροκίνητα" τη μεταβλητή με το μεγαλύτερο p-value και να εκτιμήσουμε από τη αρχή ένα νέο μοντέλο, αυτή τη φορά με τρεις επεξηγηματικές μεταβλητές. Προσοχή, ακόμα και αν έχουμε εικόνα για το ποιες μεταβλητές είναι στατιστικά σημαντικές από το πλήρες μοντέλο, θα πρέπει να ξαναγίνουν εκ νέου οι έλεγχοι για το νέο, μειωμένο μοντέλο.

Συνεχίζουμε, λοιπόν, επιλέγοντας στη γραμμή εντολών τη διαδρομή Analyze — Regression — Linear όπως και στην προηγούμενη περίπτωση, με τη διαφορά πως αυτή τη φορά δεν προσθέτουμε τη μεταβλητή Price στο κελί Independent(s), όπως φαίνεται παρακάτω:

🕼 Linear Regression		×
 ✓ Price ✓ Engine ✓ Horsepower ✓ Weight ✓ Year ✓ Fuel ✓ Brand ✓ Type 	Dependent:	Statistics Plots Save Options Bootstrap
	Selection Variable: Case Labels: WLS Weight: Paste Reset Cancel Help	

Στο παράθυρο των αποτελεσμάτων, αρχικά εξετάζουμε τον πίνακα Model Summary ο οποίος για την ώρα δεν μας παρέχει κάποια ιδιαίτερα χρήσιμη πληροφορία καθώς διαφέρει ελάχιστα από τον αντίστοιχο για το πλήρες μοντέλο.

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,778 ^a	,606	,602	1,700

a. Predictors: (Constant), Horsepower, Weight, Engine

Ο πίνακας που μας ενδιαφέρει περισσότερο για τη συνέχεια της ανάλυσής μας είναι ο πίνακας Coefficients από τον οποίο μπορούμε να διαπιστώσουμε, κοιτώντας τη στήλη Sig με τα p-value, πως όλοι οι συντελεστές είναι στατιστικά σημαντικοί. Αυτό αποτελεί ένδειξη πως αν αφαιρέσουμε κάποια από τις υπάρχουσες επεξηγηματικές μεταβλητές θα οδηγηθούμε σε ένα λιγότερο ερμηνευτικό μοντέλο. Καλείστε να το ελέγξετε συγκρίνοντας τους προσαρμοσμένους συντελεστές προσδιορισμού για τα διάφορα μοντέλα.

Καταλήγουμε λοιπόν σε ένα μοντέλο με τρεις επεξηγηματικές μεταβλητές: Engine- X_1 , Horsepower- X_2 και Weight- X_3 για την ανεξάρτητη μεταβλητή Acceleration. Το οποίο είναι και αυτό που θα εκτιμήσουμε

στη συνέχεια.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \epsilon_i, \quad i = 1, ..., 391$$

Η μέθοδος επιλογής μοντέλου που περιγράφηκε αντιστοιχεί σε μια εκτέλεση της μεθόδου backward, που είναι και μία από τις μεθόδους που μπορεί να εκτελέσει το SPSS αυτόματα όπως θα δούμε στη συνέχεια.

Επειδή συχνά οι διαθέσιμες επεξηγηματικές μεταβλητές είναι πολύ περισσότερες από ότι στο παράδειγμα που εξετάζουμε, όπως μπορείτε να φανταστείτε, η εκτέλεση της παραπάνω διαδικασίας μπορεί να γίνει πολύ πιο χρονοβόρα. Όπως ήδη αναφέρθηκε στο SPSS υπάρχουν διαθέσιμες κάποιες αυτοματοποιημένες μέθοδοι επιλογής. Εν συντομία αυτές είναι οι ακόλουθες:

Forward: Η μέθοδος αυτή ξεκινάει από το κενό/μηδενικό μοντέλο (null model), χωρίς επεξηγηματικές μεταβλητές και εισάγει στο μοντέλο τη μεταβλητή την περισσότερο στατιστικά σημαντική (με το μικρότερο p-value). Η διαδικασία επαναλαμβάνεται έως ότου δεν υπάρχει μεταβλητή που βελτιώνει στατιστικά σημαντικά το μοντέλο.

Backward: Η μέθοδος αυτή ξεκινά με το μοντέλο που περιλαμβάνει όλες τις υποψήφιες επεξηγηματικές μεταβλητές και συνεχίζει αφαιρώντας κάθε φορά τη λιγότερο στατιστικά σημαντική (με το μεγαλύτερο p-value). Η διαδικασία επαναλαμβάνεται μέχρι όλες οι επεξηγηματικές μεταβλητές που περιλαμβάνονται στο μοντέλο να είναι στατιστικά σημαντικές.

Stepwise: Η παρούσα μέθοδος αποτελεί ένα συνδυασμό των δύο παραπάνω. Είναι μια τροποποίηση της μεθόδου Forward, όπου σε κάθε βήμα, μετά την εισαγωγή της εκάστοτε μεταβλητής, ελέγχεται αν πρέπει να αφαιρεθεί κάποια από τις ήδη υπάρχουσες με βάση κάποιο κριτήριο επιλογής. Η μέθοδος σταματά όταν δεν μπορεί να εισαχθεί ή να αφαιρεθεί κάποια μεταβλητή.

Remove: Η συγκεκριμένη μέθοδος είναι το αντίθετο της Enter. Δηλώνουμε τις μεταβλητές που θέλουμε να μην μπουν στο μοντέλο μας. Σπάνια την χρησιμοποιούμε.

Σημειώνεται ότι διαφορετικές μέθοδοι επιλογής μπορεί να οδηγήσουν σε διαφορετικά μοντέλα.

Ας δούμε τώρα πως μπορούν να εφαρμοστούν τα όσα αναλύσαμε στο SPSS. Εφόσον πήραμε μια εικόνα για

την Bacward μέθοδος επιλογής στη "χειροκίνητη" εκδοχή της ας δούμε πως μπορούμε να εφαρμόσουμε τη Forward μέθοδο. Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze — Regression — Linear ακριβώς όπως πριν. Στο παράθυρο που εμφανίζεται μετακινούμε δεξιά, στο κελί Dependent, την εξαρτημένη μεταβλητή Acceleration και στο κελί Independent(s) όλες τις υποψήφιες επεξηγηματικές μεταβλητές Engine, Price, Weight και Horsepower. Η διαφορά αφορά στο πεδίο Method, από το οποίο, όπως αναφέρθηκε, μπορούμε να επιλέξουμε τη μέθοδο επιλογής του βέλτιστου μοντέλου. Διαλέγουμε λοιπόν την επιλογή Forward. Στη συνέχεια, στην καρτέλα **Statistics** επιλέγουμε Estimates, Confidence Intervals και Model Fit, όπως κάναμε και πριν.



Εφαρμόζοντας τα παραπάνω, στα αποτελέσματα θα δούμε πέντε πίνακες αυτή τη φορά. Αρχικά, εστιάζουμε τους πίνακες Model Summary, ANOVA και Coefficients, που είδαμε και παραπάνω, με τη διαφορά πως αυτή τη φορά μας παρέχουν πληροφορίες για όλα τα μοντέλα που εξετάστηκαν και εκτιμήθηκαν μέχρι να προκύψει το τελικό. Κάθε γραμμή αντιστοιχεί σε διαφορετικό μοντέλο και μπορούμε να την ερμηνεύσουμε όπως αναλύθηκε παραπάνω.

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,684ª	,468	,467	1,969
2	,771 ^b	,595	,592	1,722
3	,778°	,606	,602	1,700

a. Predictors: (Constant), Horsepower

b. Predictors: (Constant), Horsepower, Weight

c. Predictors: (Constant), Horsepower, Weight, Engine

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1328,213	1	1328,213	342,537	,000 ^b
	Residual	1508,377	389	3,878		
	Total	2836,590	390			
2	Regression	1686,587	2	843,294	284,519	°000,
	Residual	1150,003	388	2,964		
	Total	2836,590	390			
3	Regression	1717,635	3	572,545	198,019	^b 000,
	Residual	1118,955	387	2,891		
	Total	2836,590	390			

a. Dependent Variable: Acceleration

b. Predictors: (Constant), Horsepower

c. Predictors: (Constant), Horsepower, Weight

d. Predictors: (Constant), Horsepower, Weight, Engine

Co	effi	cier	ntsa
~~			

		Unstandardized Coefficients		Standardized Coefficients		
Model		В	Std. Error	Beta	t	Sig.
1	(Constant)	20,647	,289		71,385	,000
	Horsepower	-,048	,003	-,684	-18,508	,000
2	(Constant)	18,439	,323		57,104	,000
	Horsepower	-,091	,005	-1,291	-20,186	,000
	Weight	,005	,000	,703	10,996	,000
3	(Constant)	17,231	,487		35,355	,000
	Horsepower	-,082	,005	-1,168	-15,872	,000
	Weight	,006	,001	,915	10,129	,000
	Engine	-,001	,000	-,341	-3,277	,001

a. Dependent Variable: Acceleration

Οι δύο επιπλέον πίνακες που εμφανίζονται, με τίτλους Variables Entered/Removed και Variables Excluded μας πληροφορούν για το ποιές μεταβλητές εισήχθησαν ή αφαιρέθηκαν από το μοντέλο και με ποιά κριτήρια,

ο πρώτος, και για το ποιές από τις υποψήφιες μεταβλητές εξαιρέθηκαν και δεν μπήκαν στο μοντέλο σε κάθε επανάληψη της διαδικασίας, ο δεύτερος.

Model	Variables Entered	Variables Removed	Method
1	Horsepower		Forward (Criterion: Probability-of- F-to-enter <= , 050)
2	Weight		Forward (Criterion: Probability-of- F-to-enter <= , 050)
3	Engine		Forward (Criterion: Probability-of- F-to-enter <= , 050)

Variables Entered/Removed^a

a. Dependent Variable: Acceleration

Collinearity Statistics Partial Correlation Tolerance Beta In t Sig. Model Price ,009^b 1 ,074 ,941 ,004 ,097 Engine ,412^b 5,043 ,000, ,248 ,193 Weight ,703^b 10,996 .000 ,487 ,255 2 Price -,008° ,097 -,072 ,942 -,004 Engine -,341° -3,277 .001 -,164 ,094 3 Price -.010^d -.097 .922 -,005 ,097

Excluded Variables^a

a. Dependent Variable: Acceleration

b. Predictors in the Model: (Constant), Horsepower

c. Predictors in the Model: (Constant), Horsepower, Weight

d. Predictors in the Model: (Constant), Horsepower, Weight, Engine

Μπορούμε, λοιπόν, να δούμε από την πρώτη γραμμή του πίνακα Variables Entered/Removed, πως η πρώτη μεταβλητή που εισήχθη στο μοντέλο είναι η Horsepower, ενώ στην πρώτη γραμμή του πίνακα Variables Excluded μπορούμε να δούμε τις μεταβλητές που εξαιρέθηκαν από την εισαγωγή, συμπληρωματικά. Παρατηρούμε, επομένως, ότι η σειρά εισαγωγής των μεταβλητών είναι Horsepower, Weight και Engine, ενώ η μεταβλητή Price δεν επιλέγεται ποτέ.

Αποτελέσματα για το Τελικό Μοντέλο

Το μοντέλο στο οποίο καταλήγουν και οι δύο μέθοδοι επιλογής που εφαρμόσαμε είναι αυτό με τρεις επεξηγηματικές μεταβλητές, τις Engine- X_1 , Horsepower- X_2 και Weight- X_3 .

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \epsilon_i, \quad i = 1, ..., 391$$

Στον πίνακα Coefficients παρακάτω στη στήλη B, μπορούμε να δούμε τις εκτιμήσεις των συντελεστών του γραμμικού μοντέλου, ενώ από τις δύο τελευταίες με τίτλους Lower Bound και Upper Bound παίρνουμε τα αντίστοιχα διαστήματα εμπιστοσύνης. Έτσι, για παράδειγμα η εκτίμηση του συντελεστή της μεταβλητής X_2 , Horsepower, είναι $\hat{\beta}_2 = -0.082$, με 95% διάστημα εμπιστοσύνης (-0.092, -0.07200). Τελικά, η εκτιμημένη εξίσωση της πολλαπλής γραμμικής παλινδρόμησης είναι:

$$Y_i = 17.231 - 0.001X_{1i} - 0.082X_{2i} + 0.006X_{3i}, \quad i = 1, ..., 391$$

Comolence									
		Unstandardized Coefficients		Standardized Coefficients			95,0% Confiden	ce Interval for B	
Model		В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	17,231	,487		35,355	,000	16,273	18,189	
	Engine	-,001	,000	-,341	-3,277	,001	-,001	,000	
	Horsepower	-,082	,005	-1,168	-15,872	,000	-,092	-,072	
	Weight	.006	.001	.915	10.129	.000	.005	.008	

Coefficients	Co	eff	ici	eı	nts	a
--------------	----	-----	-----	----	-----	---

a. Dependent Variable: Acceleration

Έλεγχος υποθέσεων γραμμικού μοντέλου

Έλεγχος κανονικότητας: Για να ελέγξουμε αν τα κατάλοιπα του μοντέλου ακολουθούν κανονική κατανομή, ακολουθούμε τη διαδικασία που περιγράφηκε στο τέλος της ενότητας 4.

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Analyze ightarrow Descriptive Statistics ightarrow Explore.

Στο παράθυρο που εμφανίζεται στο κελί Dependent list μετακινούμε δεξιά τα Standardized Residuals τα οποία είχαμε αποθηκεύσει νωρίτερα. Στην καρτέλα Plots επιλέγουμε Normality plots with tests.

Βλέπουμε πως ο έλεγχος Shapiro-Wilk απορρίπτει την κανονικότητα των καταλοίπων σε επίπεδο στατιστικής σημαντικότητας 5% (p-value=2%). Από τα τρία γραφήματα είναι αρκετά εμφανής η δεξιά ουρά των καταλοίπων, η οποία προκαλεί πολλές ακραίες παρατηρήσεις.



Δεν μπορούμε να σταματήσουμε εδώ την ανάλυσή μας, αφού η υπόθεση της κανονικότητας δεν πληρείται.

Μετασχηματισμός Μοντέλου

Ας ξανακοιτάξουμε το scatterplot της acceleration με την horsepower. Παρατηρούμε ότι η σχέση τους δεν είναι ακριβώς γραμμική. Το γράφημα σχηματίζει μία γωνία η οποία θυμίζει πιο πολύ κυρτή συνάρτηση παρά γραμμική.



Θα επιχειρήσουμε να μετασχηματίσουμε τη μεταβλητή Horsepower ώστε να αναγκάσουμε τη σχέση να γίνει γραμμική. Αφού το γράφημα εμφανίζει μία κυρτή καμπύλη, επιλέγουμε μία κοίλη συνάρτηση για τον μετασχηματισμό μας. Δημιουργούμε λοιπόν την μεταβλητή lhorse=ln(Horsepower).

Στη γραμμή εντολών ακολουθούμε τη διαδρομή Transform \rightarrow Compute Variable. Μπορούμε να ορίσουμε τη νέα μεταβλητή είτε χειροκίνητα, είτε από τις επιλογές στο παράθυρο που μας εμφανίζεται.

tars.sav [DataSet1] - IBM Si	PSS Statistics Data I	Editor						đ	🗎 Compute Variable			×
Constant of the second	View Data View Data Price 21041 27057 29409 29548 24075 30506 25868 26746 25731 21840 21212 24087 20134 20467 19741 29638 28743 17382 28733	55 Statistics Data II Transform & Compute Value Sitt Values Count Value Sitt Values Count Value Sitt Values Count Value Sitt Values Recode Into Recode Into Sitt Replace Mis Replace Mis Recode Into Sitt Replace Mis Sitt Repla	Editor Inalyze Diri ariable es within Cas suthin Cas suthin Cas suthin Cas suthin Cas suthin Cas suthin Cas suthin Cas suthin Cas suthin Cas ing ning Different Varial bothere Va	ect <u>Marketing</u> les hables ng ators 1576 1971 2076 958	Graphh tri+G	Ution Image: Constraint of the second s	Vear Add_on 2018 2018 2018 2017 2018 2016 2017 2016 2016 2017 2016 2017 2016 2017 2016 2017 2016 2017 2016 2017 2017 2016 2017 2017 2017 2017 2017 2017 2018 2018 2018 2017	s <u>Wind</u>	•	Target Variable: Ihorse Type & LabeL Price Engine Horsepower Hors] =	Numgic Expression: LN(Horsepower) • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < > 7 • < < > 7 • < < > 7 • < < > 7 • < < > 7 • < < > 7 • < < > 7 • < < > 7 • < < > 7 • < < < 7 • < < < < 7 • < < < < < < 1 • < < < < < < < < < < < < < < < < < < <	Function group: All Antimetic COF & Noncentral CDF Convertion Outrent Date/Imme Date Antimetic Date Creation The Date Creation The Convertion The C

Συγκρίνοντας τώρα τα scatterplot των δύο μεταβλητών, βλέπουμε βελτίωση στη γραμμικότητα της σχέσης.



Είμαστε σε θέση να κατασκευάζουμε ξανά το γραμμικό μας μοντέλο, αυτή τη φορά με την lhorse στην θέση της Horsepower. Παρακάτω βλέπουμε τον πίνακα των συντελεστών του μοντέλου. Η εκτιμημένη εξίσωση της πολλαπλής γραμμικής παλινδρόμησης είναι:

$$Y_i = 52.076 - 0.001X_{1i} - 9.771X_{2i} + 0.008X_{3i}, \quad i = 1, ..., 391$$

		Unstandardize	d Coefficients	Standardized Coefficients		
Model		В	Std. Error	Beta	t	Sig.
1	(Constant)	52.076	2.080		25.042	.000
	Engine	001	.000	512	-5.580	.000
	lhorse	-9.771	.535	-1.240	-18.263	.000
	Weight	.008	.001	1.152	13.141	.000

Coefficients^a

a. Dependent Variable: Acceleration

Ας συγκρίνουμε τα γραφήματα των καταλοίπων για το αρχικό μοντέλο και το μετασχηματισμένο. Στο ιστόγραμμα παρατηρούμε βελτίωση στην συμμετρία, αν και κάποιες ακραίες παρατηρήσεις φαίνεται να υπάρχουν ακόμα.



Στο boxplot βλέπουμε πως όντως κάποιες ακραίες παρατηρήσεις επιμένουν, αλλά είναι λιγότερες και επιπλέον εμφανίζονται και στις δύο ουρές (κάτι που υποδηλώνει συμμετρία).



Στο Q-Q Plot βλέπουμε εμφανή βελτίωση στην δεξιά ουρά.



Τέλος, κάνουμε τον έλεγχο κανονικότητας Shapiro-Wilk, ο οποίος μας δίνει p-value=0.18. Επομένως, η κανονικότητα δεν απορρίπτεται σε επίπεδο στατιστικής σημαντικότητας 5%.

	Kolm	ogorov-Smir	nov ^a	ę	Shapiro-Wilk	
	Statistic df				df	Sig.
Standardized Residual	.033	391	.200	.995	391	.180

Tests of Normality

*. This is a lower bound of the true significance.

a. Lilliefors Significance Correction

Ομοσκεδαστικότητα και Ανεξαρτησία: Παρότι υπάρχουν στατιστικοί έλεγχοι για την ομοσκεδαστικότητα και την ανεξαρτησία των καταλοίπων, εδώ θα αρκεστούμε απλά σε ένα γραφικό έλεγχο. Συγκεκριμένα θα δημιουργήσουμε το σημειόγραμμα (scatterplot) των τυποποιημένων καταλοίπων (standardized residuals) ως προς τις τυποποιημένες προβλεπόμενες τιμές της γ (standardized predicted values).



Παρατηρούμε πώς για τιμές μικρότερες εκατέρωθεν του x=1, τα κατάλοιπα παρουσιάζουν πολύ μεγάλες διαφορές. Για μικρές τιμές του x, τα κατάλοιπα εμφανίζουν μεγάλη συγκέντρωση στα θετικά και μία πολύ ακραία παρατήρηση. Επομένως τα κατάλοιπα δεν φαίνονται να είναι ανεξάρτητα. Επιπλέον, τα κατάλοιπα μέχρι το σημείο x=1 φαίνεται να έχουν μία φθίνουσα πορεία, επομένως παρουσιάζουν και ετεροσκεδαστικότητα. Τα παραπάνω υποδεικνύουν ότι πρέπει να δοκιμάσουμε εναλλακτικές μεθόδους μοντελοποίησης ώστε να έχουμε τόσο ανεξαρτησία, όσο και ομοσκεδαστικότητα των καταλοίπων.

Ασκήσεις

Εξελίξτε τα απλά γραμμικά μοντέλα για τα δεδομένα του αρχείου grades που κατασκευάσατε στην προηγούμενη ενότητα. Επιλέξτε ως εξαρτημένη μεταβλητή την Grades και σαν ανεξάρτητες τις Maths, Physics και Gymnastics.

6 Διαστήματα Εμπιστοσύνης & Έλεγχοι Υποθέσεων

Σύνοψη: Σε αυτή την ενότητα θα εμβαθύνουμε στην κατασκευή διαστημάτων εμπιστοσύνης και τους ελέγχους υποθέσεων που μπορούμε να εκτελέσουμε για τους συντελεστές (β) ενός γραμμικού μοντέλου, καθώς και στις μελλοντικές προβλέψεις (y).

Διαστήματα Εμπιστοσύνης

Αρχικά, θα κατασκευάσουμε ξανά ένα γραμμικό μοντέλο. Ως εξαρτημένη μεταβλητή επιλέγουμε την acceleration, ενώ ως ανεξάρτητες τις price, engine, horsepower και weight.



Στο μενού Statistics επιλέγουμε επιπλέον το confidence intervals που προς το παρόν αφήνουμε στο 95% που έχει ως προεπιλογή το SPSS. Στο μενού Save επιλέγουμε Mean και Individual, ξανά με την προεπιλογή 95%.

Regression Coefficients Estimates Confidence intervals Level(%): 95 Covariance matrix	Model fit R squared change Descriptives Part and partial correlations Collinearity diagnostics
Residuals	
Durbin-Watson	
Casewise diagnostics	3
Outliers outside:	3 standard deviations

ta	Linear Regress	ion: Save ×
Predicted Values Vunstandardized Standardized Adjusted S.E. of mean gred	ictions	Residuals Unstandardized Standardized Studentized Dejeted Studentized deleted
Distances Mahalanobis Cooks Leverage values Prediction Intervals Mean Individua Confidence Interval:	ıl 95 %	Influence Statistics
Coefficient statistics Create coefficient Create a new dat Dataset nar Write a new data	statistics aset ne:	
Export model informa	tion to XML file	Browse

Τρέχουμε το γραμμικό μοντέλο και παρατηρούμε τον πίνακα Coefficients. Μπορούμε να δούμε ότι το SPSS κατασκεύασε 95% διαστήματα εμπιστοσύνης για κάθε συντελεστή $\hat{\beta}_i$. Τα διαστήματα που κατασκευάζει το SPSS είναι συμμετρικά, αφήνουν δηλαδή ένα σφάλμα 2.5% προς κάθε κατεύθυνση. Μπορούμε να αλλάξουμε το επίπεδο σημαντικότητας ώστε να πάρουμε τα διαστήματα της αρεσκείας μας. Με τα διαστήματα αυτά μπορούμε να πάρουμε μία καλύτερη εικόνα για το "πού είναι πιθανό να βρίσκεται η πραγματική τιμή β_i ".

				Coefficients ^a					
		Unstandardize	d Coefficients	Standardized Coefficients			95,0% Confider	ice Interval for B	
Mode	1	В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	17,288	,761		22,707	,000	15,791	18,785	
	Price	-7,015E-6	,000	-,010	-,097	,922	,000	,000	(-0.000148, 0.000134)
	Engine	-,001	,000	-,341	-3,273	,001	-,001	,000	(-0.000858,-0.000214)
	Horsepower	-,082	,009	-1,158	-9,504	,000	-,098	-,065	
	Weight	,006	,001	,915	10,116	,000	,005	,008	
a De	ependent Variable	Acceleration							

Σε αυτό το σημείο πρέπει να δώσουμε ιδιαίτερη προσοχή: Επαναλαμβάνοντας πολλές φορές το πείραμά μας (παίρνοντας διαφορετικά δείγματα και δημιουργώντας τα γραμμικά μοντέλα), θα παίρνουμε ελαφρώς διαφορετικά αποτελέσματα. Αυτό που αλλάζει από επανάληψη σε επανάληψη είναι η εκτίμησή μας $\hat{\beta}_i$, όχι το β_i το οποίο πραγματικό, σταθερό, αλλά άγνωστο σε εμάς. Επομένως είναι λάθος να πούμε "το 95% των φορών που επαναλαμβάνουμε το πείραμα το β_i θα πέφτει μέσα στο διάστημα εμπιστοσύνης" γιατί αυτό υποδηλώνει πώς το διάστημά μας είναι σταθερό ενώ το β_i αλλάζει τιμές, ενώ στην πραγματικότητα συμβαίνει το αντίθετο. Μια σωστή περιγραφή θα ήταν "το 95% των φορών που επαναλαμβάνουμε το β_i ".

Έλεγχοι Υποθέσεων

Πέρα από τα διαστήματα εμπιστοσύνης θα θέλαμε να εκτελέσουμε και κάποιους ελέγχους υποθέσεων για τα β_i . Αυτοί οι έλεγχοι μπορεί να είναι τριών διαφορετικών μορφών:

$$\begin{array}{ll} H_0: \ \beta_i = 0 & H_0: \ \beta_i = 0 & H_0: \ \beta_i = 0 \\ \\ H_1: \ \beta_i \neq 0 & H_1: \ \beta_i > 0 & H_1: \ \beta_i < 0 \end{array}$$

Στο SPSS δεν μπορούμε να εκτελέσουμε απ'ευθείας τους ελέγχους για τα β_i . Μπορούμε όμως να χρησιμοποιήσουμε την άλλη όψη του ίδιου νομίσματος, τα διαστήματα εμπιστοσύνης.

Έστω ότι θέλουμε να εκτελέσουμε έναν έλεγχο της πρώτης μορφής σε επίπεδο σημαντικότητας 5%. Αυτό ισοδυναμεί με το να κατασκευάσουμε ένα διάστημα εμπιστοσύνης όπως παραπάνω, όπου επιτρέπουμε ένα σφάλμα 2.5% προς την κατεύθυνση $\beta_i > 0$ και ένα σφάλμα 2.5% προς την κατεύθυνση $\beta_i < 0$. Αν το Ο εμπεριέχεται στο διάστημά μας, τότε μένουμε στην H_0 , αν όχι, τότε την απορρίπτουμε και δεχόμαστε την H_1 .

Αν θέλουμε να εκτελέσουμε έναν έλεγχο της δεύτερης μορφής σε επίπεδο σημαντικότητας 5%, τότε όλο το σφάλμα 5% συγκεντρώνεται στην κατεύθυνση $\beta_i > 0$. Για να κατασκευάσουμε το αντίστοιχο διάστημα εμπιστοσύνης, ζητάμε από το SPSS να κατασκευάσει 90% διαστήματα, ώστε το διάστημα (a, b) να αφήνει 5% σφάλμα προς κάθε κατεύθυνση. Το διάστημα που μας ενδιαφέρει τότε είναι το $(a, + \inf)$. Αν το 0 εμπεριέχεται στο διάστημά μας, τότε μένουμε στην H_0 , αν όχι, τότε την απορρίπτουμε και δεχόμαστε την H_1 . Αντίστοιχα για ελέγχους της τρίτης μορφής τα διαστήματά μας έχουν τη μορφή $(-\inf, b)$.

Οι έλεγχοι της μορφής $\beta_i = 0$ μας δείχνουν την σημαντικότητα των μεταβλητών μας, αν δηλαδή έχει νόημα να έχουμε τη αντίστοιχη x_i στο μοντέλο μας. Αν θέλουμε να εκτελέσουμε ελέγχους της μορφής $\beta_i = c$ τότε ενεργούμε όπως παραπάνω και απλά ελέγχουμε αν το c εμπεριέχεται στα διαστήματα (αντί για το 0). Ας δούμε ένα βοηθητικό γράφημα και μερικά παραδείγματα:



Παραδείγματα

Παράδειγμα 1: Για το μοντέλο που κατασκευάσαμε, θέλουμε να εκτελέσουμε τον παρακάτω δίπλευρο έλεγχο σε επίπεδο σημαντικότητας 5%:

$$H_0: \ \beta_{weight} = 0$$
$$H_1: \ \beta_{weight} \neq 0$$

Ζητάμε από το SPSS να μας κατασκευάσει 95% διαστήματα εμπιστοσύνης (όπως παραπάνω):

				Coefficients ^a					
		Unstandardize	d Coefficients	Standardized Coefficients			95,0% Confider	nce Interval for B	
Model		В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	17,288	,761		22,707	,000	15,791	18,785	
	Price	-7,015E-6	,000	-,010	-,097	,922	,000	,000	(-0.000148, 0.000134)
	Engine	-,001	,000	-,341	-3,273	,001	-,001	,000	(-0.000858,-0.000214)
	Horsepower	-,082	,009	-1,158	-9,504	,000	-,098	-,065	
	Weight	,006	,001	,915	10,116	,000	,005	,008	

a. Dependent Variable: Acceleration

Παρατηρούμε πως έχουμε $b_{weight} = 0.006$. Το διάστημα εμπιστοσύνης που κατασκευάσαμε είναι το (0.005, 0.008). Βλέπουμε πως δεν περιέχει το 0, επομένως απορρίπτουμε την μηδενική υπόθεση.

Παράδειγμα 2: Για το μοντέλο που κατασκευάσαμε, θέλουμε να εκτελέσουμε τον παρακάτω μονόπλευρο έλεγχο σε επίπεδο σημαντικότητας 2.5%:

$$\begin{aligned} H_0: & \beta_{price} = 0 \\ H_1: & \beta_{price} > 0 \end{aligned}$$

Το διάστημα εμπιστοσύνης που θέλουμε είναι της μορφής $(a, + \inf)$. Το SPSS μας κατασκευάζει μόνο συμμετρικά διαστήματα (a, b), επομένως θα του ζητήσουμε να έχει 2.5% σφάλμα σε κάθε μεριά, δηλαδή να κατασκευάσει ένα 95% διάστημα εμπιστοσύνης (ξανά, όπως παραπάνω).

				Coefficients ^a					
		Unstandardize	d Coefficients	Standardized Coefficients			95,0% Confider	ice Interval for B	
Model		В	Std. Error	Beta	l t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	17,288	,761		22,707	,000	15,791	18,785	
	Price	-7,015E-6	,000	-,010	-,097	,922	,000	,000	(-0.000148, 0.000134)
	Engine	-,001	,000	-,341	-3,273	,001	-,001	,000	(-0.000858,-0.000214)
	Horsepower	-,082	,009	-1,158	-9,504	,000	-,098	-,065	
	Weight	,006	,001	,915	10,116	,000	,005	,008	

a. Dependent Variable: Acceleration

Παρατηρούμε πως έχουμε $b_{price} = -0.000007015$. Το διάστημα που παίρνουμε είναι $(-0.000148, + \inf)$. Βλέπουμε πως περιέχει το 0, επομένως δεν απορρίπτουμε την μηδενική υπόθεση. Αυτό μας δείχνει πως η μεταβλητή price πρέπει να βγει από το μοντέλο.

Παράδειγμα 3: Για το μοντέλο που κατασκευάσαμε, θέλουμε να εκτελέσουμε τον παρακάτω μονόπλευρο έλεγχο σε επίπεδο σημαντικότητας 5%:

$$H_0: \ \beta_0 = 16$$

 $H_1: \ \beta_0 < 16$

Το διάστημα εμπιστοσύνης που θέλουμε είναι της μορφής $(-\inf, b)$. Το SPSS μας κατασκευάζει μόνο συμμετρικά διαστήματα (a, b), επομένως θα του ζητήσουμε να έχει 5% σφάλμα σε κάθε μεριά, δηλαδή να κατασκευάσει ένα 90% διάστημα εμπιστοσύνης.

				coemeienta				
		Unstandardize	d Coefficients	Standardized Coefficients			90,0% Confider	ice Interval for B
Model		В	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound
1	(Constant)	17,288	,761		22,707	,000	16,033	18,543
	Price	-7,015E-6	,000	-,010	-,097	,922	,000	,000
	Engine	-,001	,000	-,341	-3,273	,001	-,001	,000
	Horsepower	-,082	,009	-1,158	-9,504	,000	-,096	-,067
	Weight	,006	,001	,915	10,116	,000	,005	,008

Coefficients^a

a. Dependent Variable: Acceleration

Παρατηρούμε πως έχουμε $\beta_0 = 17.288$. Το διάστημα που παίρνουμε είναι (- inf, 18.543). Βλέπουμε πως περιέχει το 16, επομένως δεν απορρίπτουμε την μηδενική υπόθεση. Αν στη θέση του 16 είχαμε το 19, θα απορρίπταμε την μηδενική υπόθεση.

Διαστήματα Πρόβλεψης για την Εξαρτημένη Μεταβλητή

Έχοντας δει τα αυτοκίνητα του σετ δεδομένων μας, θα θέλαμε να προβλέψουμε την επιτάχυνση ενός αυτοκινήτου γνωρίζοντας μόνο τα άλλα χαρακτηριστικά του. Θα θέλαμε δηλαδή να κατασκευάσουμε διαστήματα εμπιστοσύνης και για την εξαρτημένη μεταβλητή ενός μοντέλου. Υπάρχουν δύο είδη τέτοιων διαστημάτων. Ας τα γνωρίσουμε μέσα από δύο σενάρια:

Σενάριο 1: Βρισκόμαστε σε μία έκθεση αυτοκινήτων, όπου μπορούμε να δούμε και να μελετήσουμε τα 391 αυτοκίνητα. Γνωρίζουμε πως σύντομα θα καταφτάσει στην έκθεση ένα ακόμα αυτοκίνητο, το No. 392. Οι οργανωτές της έκθεσης μας πληροφορούν για τα χαρακτηριστικά του (*price* = 20000, *engine* = 6000, *horsepower* = 200, *weight* = 1800) αλλά όχι για την επιτάχυνσή του (*acceleration*). Θα θέλαμε να προβλέψουμε την επιτάχυνση του συγκεκριμένου αυτοκινήτου και να δημιουργήσουμε ένα διάστημα πρόβλεψης για αυτήν.

Σενάριο 2: Βρισκόμαστε σε μία έκθεση αυτοκινήτων, όπου μπορούμε να δούμε και να μελετήσουμε τα 391 αυτοκίνητα. Σύντομα αναχωρούμε για το εργοστάσιο παραγωγής αυτοκινήτων με συγκεκριμένα χαρακτηριστικά τα οποία γνωρίζουμε (price = 20000, engine = 6000, horsepower = 200, weight = 1800). Θα θέλαμε να προβλέψουμε τη μέση επιτάχυνσή τους και να δημιουργήσουμε ένα διάστημα πρόβλεψης για αυτήν.

Παρότι τα δύο παραπάνω σενάρια δείχνουν ίδια με μια πρώτη ανάγνωση, αν τα εξετάσουμε αναλυτικά, θα δούμε πως διαφέρουν. Πράγματι, η σημειακή πρόβλεψη που θα κάνουμε για την επιτάχυνση θα ταυτίζεται στα δύο σενάρια, αφού το μόνο που κάνουμε είναι να αντικαταστήσουμε τα συγκεκριμένα *x* που μας δίνονται στο γραμμικό μοντέλο και να υπολογίσουμε την *y*. Τα διαστήματα που θα κατασκευάσουμε όμως διαφέρουν. Όταν το αυτοκίνητο Νο. 392 φτάσει στην έκθεση, είμαστε αρκετά αβέβαιοι για το πόσο κοντά θα πέσουμε στην πρόβλεψή μας, ενώ όταν πάμε στο εργοστάσιο παραγωγής γνωρίζουμε πως η μέση επιτάχυνση όλων των αυτοκινήτων και η πρόβλεψή μας θα είναι σχετικά κοντά. Τα δύο διαστήματα δηλαδή κατασκευάζονται χρησιμοποιώντας διαφορετική διασπορά.

Για να κατασκευάσουμε τα εν λόγω διαστήματα, δημιουργούμε μία ακόμα σειρά στο σετ δεδομένων μας στην οποία βάζουμε τα χαρακτηριστικά που θέλουμε, αφήνοντας την επιτάχυνση κενή. Με αυτόν τον τρόπο, το SPSS δεν χρησιμοποιεί την παρατήρηση αυτή στο γραμμικό μοντέλο. Στην παλινδρόμηση,

55

επιλέγουμε Mean και Individual Prediction Intervals όπως δείξαμε στην αρχή. Το SPSS δημιουργεί καινούριες μεταβλητές που δείχνουν το κάτω και το άνω άκρο των προβλέψεων μας, τόσο για αυτοκίνητα με χαρακτηριστικά σαν τα 391 που είχαμε στο σετ δεδομένων μας εξ αρχής, όσο και για την τελευταία παρατήρηση που βάλαμε.

t	*cars.sav [DataSet1] - IBM SPSS Statistics Data Editor 🛛 – 🗇 🗙																	
<u>F</u> ile <u>E</u> dit	<u>V</u> iew D	ata	<u>T</u> ransform	<u>A</u> nalyze D	Direct <u>M</u> arketing	<u>G</u> raphs	Utilities Add	ons <u>W</u>	indow	<u>H</u> elp								
😂 H		.		1		r A	*,		4									
392 : Price		2000	0													Visi	ble: 15 of 15	5 Variables
510	Price		Engine	Horsepower	Weight	Acceleration	Year	Fuel	Brand	Туре	PRE_1	ZRE_1	LMCI_1	UMCI_1	LICI_1	UICI_1	var	V.
377	18-	153	2294	92	1289	17	2016	2	1	1	1 16,78	-,10	16,45	17,11	13,42	20,14		-
378	18	462	2966	110	1325	17	2016	2	1	1	1 15,18	1,00	14,98	15,39	11,83	18,54		
379	13	171	1769	70	1010	17	2015	3	3	3	0 17,08	,05	16,78	17,39	13,72	20,44		
380	15	374	4293	85	1356	17	2013	3	1	1	0 16,73	,06	16,21	17,26	13,35	20,12		
381	14	575	1605	70	956	16	2016	3	1	1	0 16,81	-,38	16,57	17,06	13,46	20,17		
382	16	086	2474	90	1327	17	2016	2	1	1	1 17,11	-,21	16,84	17,38	13,75	20,47		
383	12	518	1491	68	886	17	2016	3	3	3	0 16,60	,39	16,19	17,00	13,22	19,97		
384	13	650	1835	88	1077	18	2016	3	1	1	1 16,01	1,34	15,53	16,49	12,63	19,39		
385	16	141	2474	90	1230	19	2016	2	1	1	1 16,48	1,29	16,27	16,69	13,12	19,83		
386	14:	278	1491	68	911	19	2016	3	:	3	0 16,75	1,47	16,48	17,01	13,39	20,10		
387	16	213	1966	79	1181	19	2017	2	1	1	0 17,33	,81	17,08	17,58	13,97	20,69		
388	19033 1835 88 1188 19 2016 2 1 1 16,69 1,38 15,25 17,14 13,32 20,07																	
389	15323 1950 52 1224 18 2016 3 1 0 17,38 ,36 17,09 17,67 14,02 20,74 16,00 16,01 16,02 18,01 16,00 17,01 17,67 14,02 20,74 16,00 16,01 16,111 16,11 16,11 16,11 16,11 16,11 16,1																	
390	16	025	1835	88	1172	19	2017	2	1		1 16,61	1,63	16,33	16,89	13,25	19,97		_
391	12	538	1589	52	958	20	2013	4	2	2	0 18,32	1,10	17,99	18,64	14,95	21,68		_
392	20	000	6000	200	1800						. 9,28		8,11	10,45	5,74	12,83		
393																		
394																		
200																		
397																		+
398																		
399																		
400																		
100																		

Data View	Variable Vi	ew																
												IBM S	PSS Statistic	s Processo	r is ready	Unicod	le:ON	

Βλέπουμε πώς η σημειακή πρόβλεψη είναι 9.28 για την επιτάχυνση (και για τα δύο σενάρια). Το Mean Confidence Interval (MCI) που προβλέπουμε για το εργοστάσιο είναι το (8.11, 10.45) ενώ το Individual Confidence Interval (ICI) που προβλέπουμε για το αυτοκίνητο που θα φτάσει στην έκθεση είναι (5.74, 12.83), αρκετά μεγαλύτερο του MCI.

7 ANOVA & t-test

Σύνοψη: Σε αυτή την ενότητα θα εμβαθύνουμε σε ελέγχους για τον μέσο μίας μεταβλητής και τους μέσους διαφορετικών υποκατηγοριών μιας μεταβλητής. Συγκεκριμένα θα ασχοληθούμε με τα t-test και τα μοντέλα ANOVA.

One Sample t-test

Ας υποθέσουμε πως θέλουμε να κάνουμε έναν έλεγχο υποθέσεων για την μέση επιτάχυνση των αυτοκινήτων. Αν τρέξουμε κάποια περιγραφικά στατιστικά, θα δούμε πώς η μέση επιτάχυνση των αυτοκινήτων των δεδομένων μας είναι 15.62 δευτερόλεπτα. Θέλουμε λοιπόν να εφαρμόσουμε τον παρακάτω έλεγχο:

$$H_0: \ \mu_{accel} = 15$$
$$H_1: \ \mu_{accel} \neq 15$$

Ο κατάλληλος έλεγχος για να το κάνουμε αυτό ονομάζεται t-test και συγκεκριμένα one-sample t-test (θα δούμε και άλλα είδη t-test παρακάτω). Απαραίτητη προϋπόθεση για το one sample t-test είναι η μεταβλητή να προέρχεται από κανονική κατανομή, κάτι που έχουμε ήδη ελέγξει για την Acceleration σε προηγούμενη ενότητα. Για να εκτελέσουμε τον έλεγχο στο SPSS, ακολουθούμε τη διαδρομή Analyze — Compare Means — One-Sample T Test και συμπληρώνουμε το παράθυρο που εμφανίζεται όπως παρακάτω:



Τρέχουμε το t-test και το SPSS μας εμφανίζει δύο πίνακες, έναν με περιγραφικά στατιστικά και έναν όπου εκτελείται ο έλεγχος. Παίρνουμε p - value = 0.000007 < 0.05 επομένως απορρίπτουμε την μηδενική υπόθεση σε επίπεδο στατιστικής σημαντικότητας 5%. Ακόμα, μας δίνεται η διαφορά ανάμεσα στο δειγματικό μέσο και την τιμή που θέσαμε ως πραγματικό μέσο (0.62) και ένα 95% διάστημα εμπιστοσύνης για αυτή τη διαφορά (0.35, 0.89).

T-Test

One-Sample	Statistics
one-sumple	Statistics

	Ν	Mean	Std. Deviation	Std. Error Mean
Acceleration	391	15.62	2.697	.136

One-Sample Test

	Test Value = 15								
				Mean	95% Confidence Differ	e Interval of the ence			
	t	df	Sig. (2-tailed)	Difference	Lower	Upper			
Acceleration	4.553	390	.000	.621	.35	.89			

Με το παραπάνω t-test κάνουμε έλεγχο για τον μέσο μιας μεταβλητής ως προς μια συγκεκριμένη τιμή που θέτουμε εμείς. Προχωράμε σε ένα t-test που αφορά την σύγκριση των μέσων δύο υποκατηγοριών της μεταβλητής.

Independent Sample t-test

Τα αυτοκίνητα στο σετ δεδομένων μας έχουν δύο είδη κινητήρα, std και turbo. Αυτό μπορούμε να το δούμε από την μεταβλητή Type, η οποία έχει κωδικοποιημένα τα δύο είδη με 0 και 1 για std και turbo αντίστοιχα. Στο SPSS, πάμε στην καρτέλα Variable View, στην στήλη Values, όπως φαίνεται παρακάτω:

🤖 car	s.sav [D	ataSet1] - IBM SP	SS Statistics Da	ita Editor													-	٥	\times
<u>F</u> ile	Edit	<u>V</u> iew <u>D</u> ata	Transform	Analyze	Direct <u>M</u> arke	ting <u>G</u> ra	aphs <u>U</u>	tilities Add- <u>o</u> r	ns <u>W</u> indow	Help									
			, 🛌 -	∽ 📱		با			- 42		A (14		ABG						
		Name	Туре	Width	Decimal	i La	ibel	Values	Missing	Columns		Align	Measure	Role					
1		Price	Numeric	8	0			None	None	8	遭 R	ight	🛷 Scale	💊 Input					4
2	2	Engine	Numeric	8	0			None	None	8	ĩ目R	ight	🛷 Scale	💊 Input					
3	;	Horsepower	Numeric	8	0			None	None	8	彊 R	ight	🛷 Scale	🦒 Input					
4	ļ į	Weight	Numeric	8	0			None	None	8	ĩ目R	ight	🛷 Scale	💊 Input					
Ę	i i	Acceleration	Numeric	8	0			None	None	8	彊 R	ight	🛷 Scale	🦒 Input					
6	i	Year	Numeric	8	0			None	None	8	泪 R	ight	📲 Ordinal	🦒 Input					
ī	'	Fuel	Numeric	8	0			{1, Low Con	None	8	彊 R	ight	drdinal	S Input					
8	;	Brand	Numeric	8	0			{1, BMW}	None	8	遭R	ight	🗞 Nominal	S Input					
9		Туре	Numeric	8	0			{0, std}	None	8	彊 R	ight	🗞 Nominal	S Input					
1	0			-		_													
1	1				🔄 🍓 Value	Labels					×								
1	2																		
1	3				Value	Labels													
1	4				Value					Spelling)								
1	5			_	Label														
1	6			_		0	- "etd"												
1	(Add 1	= "turbo"												
1	8																		
1	9					nange													
2	4			_	- 6	emove													
2	2																		
2	2				-			Concel	Halp										
2	3			-	-			Calicer	Help										
2	+																		-
		1				_	_				_	_						_	
Data	View	Variable View																	
													IB	M SPSS Statistics	Processor is ready	Unicod	e'ON		

Ας υποθέσουμε λοιπόν ότι θέλουμε να συγκρίνουμε τα δύο είδη αυτοκινήτων και συγκεκριμένα την επιτάχυνσή τους. Σημαντική προϋπόθεση για αυτό το t-test είναι οι δύο υποκατηγορίες να προέρχονται από κανονική κατανομή. Πρέπει λοιπόν να κάνουμε έλεγχο κανονικότητας για τα std και turbo ξεχωριστά. Για να "σπάσουμε" το δείγμα μας σε κομμάτια σύμφωνα με την μεταβλητή Type, στο μενού Explore που ανοίγουμε για τους ελέγχους κανονικότητας, βάζουμε την Type στο κουτί Factor List, όπως φαίνεται παρακάτω:



Από τον έλεγχο προκύπτει πως δεχόμαστε την κανονικότητα για την κατηγορία turbo, αλλά όχι για την std. Αυτό σημαίνει πώς το t-test που θα κάνουμε θα ισχύει προσεγγιστικά (και δεν θα μπορούσα να το εμπιστευτώ για μικρά δείγματα).

Αρχικά, θα κατασκευάσουμε ένα error bar για να έχουμε μια καλύτερη εικόνα των δύο κατηγοριών. Στη γραμμή εντολών ακολουθούμε τη διαδρομή Graphs → Legacy Dialogs → Error Bar και επιλέγουμε το Simple. Στο μενού που εμφανίζεται επιλέγουμε ως Variable την Acceleration και ως Category Axis την Type. To SPSS κατασκευάζει το γράφημα που του ζητήσαμε. Παρατηρούμε πως ο άξονας x έχει τις δύο κατηγορίες της μεταβλητής Type. Για κάθε μία από τις δύο κατηγορίες, μπορούμε να δούμε τον δειγματικό μέσο του acceleration (μαύρη κουκκίδα) και ένα 95% διάστημα γύρω από αυτόν με πλάτος ανάλογο με τη δειγματική τυπική απόκλιση (γραμμές και φράχτες). Ο μέσος για τα αυτοκίνητα τύπου std βρίσκεται πολύ υψηλότερα από τον μέσο των αυτοκινήτων turbo. Αυτό σημαίνει πώς τα std χρειάζονται περισσότερο χρόνο για να επιταχύνουν (είναι δηλαδή πιο αργά). Αν έπρεπε να βγάλουμε ένα συμπέρασμα μόνο από τους μέσους έχουμε κατασκευάσει και τα διαστήματα εμπιστοσύνης τα οποία δεν τέμνονται καν.



Προσοχή: Το box plot (που κατασκευάσαμε σε προηγούμενα εργαστήρια) και το error bar είναι διαφορετικά γραφήματα. Το box plot έχει σαν κέντρο την διάμεσο και διάστημα βασισμένο σε ποσοστημόρια, ενώ το error bar έχει σαν κέντρο τον δειγματικό μέσο και διάστημα βασισμένο στην τυπική απόκλιση. Το box plot δηλαδή είναι ανεξάρτητο κατανομής, ενώ το error bar στηρίζεται στην κανονική κατανομή της μεταβλητής y (εδώ η Acceleration).

Θέλουμε να ελέγξουμε όσα παρατηρήσαμε στο γράφημα με έναν στατιστικό έλεγχο. Θέλουμε δηλαδή να ελέγξουμε αν οι δειγματικοί μέσοι των δύο κατηγοριών είναι ίσοι ή οχι. Αυτός ο έλεγχος θα μπορούσε να περιγραφεί ως:

$$H_0: \quad \mu_{std} = \mu_{turbo}$$
$$H_1: \quad \mu_{std} \neq \mu_{turbo}$$

Ο έλεγχος αυτός ονομάζεται independent sample t-test. Το independent sample δηλώνει πως τα δύο υποσετ αυτοκινήτων είναι ανεξάρτητα μεταξύ τους, αφού δεν έχουμε λόγο να υποστηρίξουμε πως ένα αυτοκίνητο (std) επηρεάζει κάπως ένα άλλο αυτοκίνητο (turbo). Αν για παράδειγμα εξετάζαμε τους βαθμούς ενός μαθητή στα μαθηματικά και τη φυσική, δε θα μπορούσαμε να υποστηρίξουμε πώς είναι ανεξάρτητοι, αφού καλή γνώση μαθηματικών γενικά ενισχύει την επίδοση στη φυσική. Για να εκτελέσουμε τον έλεγχο στο SPSS, ακολουθούμε τη διαδρομή Analyze → Compare Means → Independent-Samples T Test. Στο μενού που εμφανίζεται τοποθετούμε τις μεταβλητές Acceleration και Type όπως παρακάτω. Για την μεταβλητή Type, πρέπει να δώσουμε στο SPSS τις τιμές 0 και 1 τις οποίες θέλουμε να συγκρίνουμε. Αυτό πρέπει να το κάνουμε γιατί θα μπορούσε η Type να είχε πολλές κατηγορίες, ενώ το t-test μπορεί να συγκρίνει μόλις δύο.



Τρέχουμε το t-test και το SPSS μας εμφανίζει δύο πίνακες. Ο πρώτος έχει περιγραφικά στατιστικά για τις δύο κατηγορίες. Μπορούμε να δούμε ότι έχουμε 181 αυτοκίνητα με κινητήρα std και 210 με turbo. Ο δεύτερος πίνακας εκτελεί το t-test. Βασική προϋπόθεση είναι να γνωρίζουμε αν οι δύο κατηγορίες δείχνουν να έχουν ίδια διασπορά ή όχι. Για να το ελέγξουμε αυτό, το SPSS τρέχει το Levene's Test (πρώτες δύο στήλες του πίνακα). Μπορούμε να δούμε ότι έχουμε p - value = 0.137 > 0.05 επομένως μένουμε στην αρχική υπόθεση ότι οι διασπορές είναι ίσες. Επομένως θα κοιτάξουμε την πρώτη γραμμή του πίνακα (Equal Variances Assumed). Αν είχαμε απορρίψει την ισότητα, θα κοιτούσαμε την δεύτερη γραμμή (αν συγκρίνουμε τις δύο γραμμές, μπορούμε να δούμε πως τα δύο t-test είναι ελαφρώς διαφορετικά).

Βλέπουμε ότι το t-test δίνει $p - value = 1.6 \cdot 10^{-11}$ επομένως απορρίπτουμε την ισότητα των δύο μέσων, όπως αναμέναμε. Επιπλέον, μας δίνει την διαφορά των μέσων (1.794), καθώς και ένα 95% διάστημα εμπιστοσύνης για αυτή την διαφορά (από 1.286 έως 2.302).

T-Test

	Group Statistics										
	Туре	N	Mean	Std. Deviation	Std. Error Mean						
Acceleration	std	181	16.58	2.375	.177						
	turbo	210	14.79	2.687	.185						

ndenendent Samnles Test					
	ndep	endent	Sam	ples	Test

Levene's Test for Equality of Variances						t-test for Equality	of Means			
					Mean		Std. Error	95% Confidence Interval of the Difference		
		F	Sig.	t	df	Sig. (2-tailed)	Difference	Difference	Lower	Upper
Acceleration	Equal variances assumed	2.216	.137	6.945	389	.000	1.794	.258	1.286	2.302
	Equal variances not assumed			7.008	388.745	.000	1.794	.256	1.291	2.297

Προσοχή: Είναι σημαντικό να καταλάβουμε πως για τον παραπάνω έλεγχο χρησιμοποιήσαμε την ανεξαρτησία των δύο κατηγοριών. Αν είχαμε εξαρτημένες παρατηρήσεις, ο κατάλληλος έλεγχος θα ήταν το Paired-Samples T test.

One-Way ANOVA

Το t-test μπορεί να συγκρίνει μόνο δύο υποκατηγορίες μίας μεταβλητής. Αν θέλουμε να συγκρίνουμε την επιτάχυνση με βάση την μάρκα των αυτοκίνητων, για την οποία έχουμε 3 κατηγορίες (BMW, Datsun, Mercedes-Benz) πρέπει να χρησιμοποιήσουμε μία γενίκευση του t-test, την One-Way ANOVA (κατά ένα παράγοντα). Όπως και με το t-test, έτσι και με την ANOVA, προϋπόθεση αποτελεί η κανονικότητα των υποκατηγοριών. Στην πραγματικότητα, η ANOVA είναι ένα γραμμικό μοντέλο με τις ποιοτικές ανεξάρτητες μεταβλητές. Επομένως, μπορούμε να εκτελέσουμε έλεγχο κανονικότητας στα κατάλοιπα.

Αρχικά, θα κατασκευάσουμε ένα error bar όπως παραπάνω.



Μπορούμε να δούμε πώς τα αυτοκίνητα μάρκας BMW επιταχύνουν πολύ πιο γρήγορα από τις άλλες δύο μάρκες. Επομένως αναμένουμε πως ο μέσος των BMW θα διαφέρει στατιστικά σημαντικά από τους άλλους δύο μέσους. Για την σύγκριση Datsun και Mercedes-Benz, δεν μπορούμε να είμαστε σίγουροι,

αφού τα δύο διαστήματα εμφανίζουν μεγάλη αλληλεπικάλυψη, με τον μέσο της μίας μάρκας να βρίσκεται μέσα στο διάστημα της άλλης. Επιθυμούμε λοιπόν να κατασκευάσουμε το μοντέλο ANOVA για να δούμε αν ισχύουν όσα παρατηρήσαμε.

Για να τρέξουμε ένα μοντέλο ANOVA στο SPSS, ακολουθούμε τη διαδρομή Analyze → General Linear Model → Univariate και μιμούμαστε τις επιλογές που φαίνονται παρακάτω:

🔄 Univariate X	🕼 Univariate: Post Hoc Multiple Comparisons for Observed Means 🛛 🗙
 ✓ Price ✓ Engine ✓ Horsepower ✓ Weight ✓ Year ✓ Fuel ✓ True 	Factor(s): Brand Fequal Variances Assumed Post Hoc Tests for: Brand Fequal Variances Assumed
Covariate(s): WLS Weight: OK Paste Reset Cancel Help	LSD S-N-K Waller-Duncan Bonferroni Tukey Type I/Type II Error Ratio: 100 Sidak Tukey's-b Dunnatt Scheffe Duncan Control Category: R-E-G-W-F Hochberg's GT2 Test R-E-G-W-Q Gabriel @ 2-sided < Control
	Continue Cancel Help
CEstimated Marginal Means	🛱 Univariate: Save 🛛 🗙
Factor(s) and Factor Interactions: Display Means for: (OVERALL) (OVERALL) Brand Image: Compare main effects Compare main effects Confidence interval adjustment. LSD(none) Total	Predicted Values Unstandardized Unstandardized Unstandardized Weighted Standard error Diagnostics Studentized Cook's distance Deleted Leverage values Coefficient Statistics
Display Descriptive statistics Estimates of effect size Observed power Parameter estimates Lack of fit Contrast coefficient matrix General estimable function Significance level: O5 Continue Cancel Hein	Create coefficient statistics Create a new dataset Dataset name: Write a new data file File Continue Cancel Help

Παρατήρηση: Για την One-Way ANOVA θα μπορούσαμε να ακολουθήσουμε την διαδρομή Analyze → Compare Means → One-Way ANOVA. Αυτό το μενού όμως δεν μας επιτρέπει να αποθηκεύσουμε τα κατάλοιπα, τα οποία χρειαζόμαστε για τον έλεγχο κανονικότητας. To SPSS μας εμφανίζει δύο πινακάκια με γενικά περιγραφικά στατιστικά για την επιτάχυνση κάθε μάρκας ξεχωριστά. Ο πίνακας που μας ενδιαφέρει είναι ο τρίτος, (Tests of Between-Subjects Effects):

Dependent Variable: Acceleration									
Source	Type III Sum of Squares	df	Mean Square	F	Sig.				
Corrected Model	183.984 ^a	2	91.992	13.456	.000				
Intercept	73555.299	1	73555.299	10759.026	.000				
Brand	183.984	2	91.992	13.456	.000				
Error	2652.606	388	6.837						
Total	98246.722	391							
Corrected Total	2836.590	390							

Tests of Between-Subjects Effects

a. R Squared = .065 (Adjusted R Squared = .060)

Από την πρώτη γραμμή του πίνακα βλέπουμε πως το μοντέλο είναι στατιστικά σημαντικό ($p - value = 2 \cdot 10^{-7}$), δηλαδή το μοντέλο είναι στατιστικά σημαντικό. Κοιτώντας την γραμμή Brand μπορούμε να δούμε ότι έχει το ίδιο p-value και επομένως η μέση επιτάχυνση αυτοκινήτων διαφορετικής μάρκας διαφέρουν (χωρίς να ξέρουμε λεπτομέρειες). Στην υποσημείωση (a) βλέπουμε πως το μοντέλο έχει $R^2 = 0.065$, $R^2_{adj} = 0.060$, δηλαδή αρκετά ασθενές.

To SPSS έπειται εμφανίζει δύο ακόμα πίνακες με τον δειγματικό μέσο και διαστήματα εμπιστοσύνης για την επιτάχυνση κάθε κατηγορίας. Αυτοί οι πίνακες έχουν την ίδια λειτουργία με τα error bar. Ο δεύτερος πίνακας που μας ενδιαφέρει είναι ο "Multiple Comparisons", που συγκρίνει τις κατηγορίες ανά δύο, ώστε να εντοπίσουμε αυτές που διαφέρουν μεταξύ τους. Από τους αστερίσκους του πίνακα μπορούμε να βρούμε τις στατιστικά σημαντικές διαφορές. Εδώ παρατηρούμε ότι η BMW διαφέρει και από τις δύο άλλες μάρκες, οι οποίες μεταξύ τους δεν διαφέρουν, όπως ακριβώς αναμέναμε.

Multiple Comparisons

Dependent Variable: Acceleration Tukey HSD

		Mean Difference (I-			95% Confide	ance Interval
(I) Brand	(J) Brand	J)	Std. Error	Sig.	Lower Bound	Upper Bound
BMW	Datsun	-1.632	.359	.000	-2.48	79
	Mercedes-Benz	-1.176	.338	.002	-1.97	38
Datsun	BMW	1.632	.359	.000	.79	2.48
	Mercedes-Benz	.456	.433	.543	56	1.47
Mercedes-Benz	BMW	1.176	.338	.002	.38	1.97
	Datsun	456	.433	.543	-1.47	.56

*. The mean difference is significant at the 0.05 level.

Ο πίνακας Homogeneous Subsets είναι προέκταση του προηγούμενου πίνακα και χωρίζει τις κατηγορίες με διαφορετικούς μέσους. Εδώ βλέπουμε πως έχουμε σε μία ομάδα την BMW (μόνη της) και σε μία άλλη τις Mercedes-Benz και Datsun.

Homogeneous Subsets

Acceleration

Tukey HSD ^{a,b}									
		Subset for alpha = 0.05							
Brand	Ν	1	2						
BMW	244	15.10							
Mercedes-Benz	79		16.28						
Datsun	68		16.73						
Sig.		1.000	.452						

Means for groups in homogeneous subsets are displayed.

a. Uses Harmonic Mean Sample Size = 95.352.

b. The group sizes are unequal. The harmonic mean of the group sizes is used. Type I error levels are not guaranteed.

Τέλος, πρέπει να κάνουμε τον έλεγχο κανονικότητας των καταλοίπων (κατά τα γνωστά). Ο έλεγχος Shapiro-Wilk βγάζει p - value = 0.04. Επομένως, η κανονικότητα των καταλοίπων απορρίπτεται σε επίπεδο σημαντικότητας 5%, αλλά όχι σε επίπεδο 1%. Σε κάθε περίπτωση, τα αποτελέσματά μας ισχύουν προσεγγιστικά, λόγω του μεγάλου μέγεθους του δείγματος.

Two-Way ANOVA

Ας υποθέσουμε τώρα ότι θέλουμε να εξετάσουμε την διαφορά της μέσης επιτάχυνσης για τις διάφορες μάρκες και τους διάφορους τύπους κινητήρων ταυτόχρονα (δηλαδή για τις $3 \cdot 2 = 6$ κατηγορίες). Θα κατασκευάσουμε μία Two-Way ANOVA (κατά δύο παράγοντες). Βασική προϋπόθεση είναι η κανονικότητα, την οποία θα ελέγξουμε μέσα από τα κατάλοιπα.

Για να την τρέξουμε στο SPSS, ακολουθούμε τη διαδρομή Analyze \rightarrow General Linear Model \rightarrow Univariate και μιμούμαστε τις επιλογές που φαίνονται παρακάτω:



Αρχικά, κοιτάμε το γράφημα, όπου μπορούμε να δούμε τους μέσους κάθε μίας από τις 6 κατηγορίες. Με μπλε χρώμα συμβολίζονται τα αυτοκίνητα με κινητήρα std και με πράσινο τα turbo. Μπορούμε να δούμε σημαντικές διαφορές τόσο στις τρεις μάρκες, όσο και στα δύο χρώματα. Το ερώτημα που θέλουμε να θέσουμε σε αυτό το σημείο είναι αν αυτοκίνητα ίδιας μάρκας και διαφορετικού κινητήρα έχουν σταθερή διαφορά για τα όλα τα είδη κινητήρα, καθώς και αν αυτοκίνητα ίδιου κινητήρα και διαφορετικής μάρκας έχουν σταθερή διαφορά για όλες τις μάρκες.



Ο πίνακας που εκτελεί τους ελέγχους που μας ενδιαφέρουν είναι ο Tests of Between-Subjects Effects:

Dependent Variable: Acceleration									
Source	Type III Sum of Squares	df	Mean Square	F	Sig.				
Corrected Model	369.959 ^a	5	73.992	11.549	.000				
Intercept	56805.782	1	56805.782	8866.433	.000				
Brand	52.627	2	26.313	4.107	.017				
Туре	129.303	1	129.303	20.182	.000				
Brand * Type	8.579	2	4.289	.670	.513				
Error	2466.632	385	6.407						
Total	98246.722	391							
Corrected Total	2836.590	390							

Tests of Between-Subjects Effects

a. R Squared = .130 (Adjusted R Squared = .119)

Από την πρώτη γραμμή του πίνακα βλέπουμε πως το μοντέλο είναι στατιστικά σημαντικό ($p - value = 2 \cdot 10^{-10}$), δηλαδή τουλάχιστον δύο κατηγορίες διαφέρουν μεταξύ τους. Κοιτώντας την γραμμή Brand μπορούμε να δούμε ότι και αυτή έχει πολύ μικρό p-value και επομένως η μέση επιτάχυνση αυτοκινήτων διαφορετικής μάρκας διαφέρουν (κάτι που είδαμε και στην One-Way ANOVA). Ομοίως για την γραμμή Type. Στην υποσημείωση (a) βλέπουμε πως το μοντέλο έχει $R^2 = 0.130$, $R^2_{adj} = 0.119$, που είναι καλύτερο από το One-Way ANOVA, αλλά παραμένει αρκετά ασθενές.

Η γραμμή που μας ενδιαφέρει όμως είναι η Brand*Type, που μας δείχνει την αλληλεπίδραση των δύο μεταβλητών. Παρατηρούμε ότι έχουμε p - value = 0.513 επομένως ο όρος αυτός δεν είναι στατιστικά σημαντικός και μπορεί να αφαιρεθεί από το μοντέλο. Αυτό σημαίνει πώς η διαφορά στην επιτάχυνση δύο αυτοκινήτων ίδιας μάρκας με διαφορετικό κινητήρα παραμένει σταθερή, από όποια μάρκα κι αν προέρχονται τα δύο αυτοκίνητα. Ομοίως, η διαφορά στην επιτάχυνση δύο αυτοκινήτων ίδιου κινητήρα

Οι επιλογές που κάναμε στο μενού Post Hoc μας εμφανίζουν ακριβώς τον ίδιο πίνακα Multiple Comparisons όπως στην One-Way ANOVA. Οι υπόλοιποι πίνακες αφορούν περιγραφικά στατιστικά και διαστήματα εμπιστοσύνης για τους μέσους.

Σε αυτό το σημείο, επαναλαμβάνουμε το Two-Way ANOVA, αυτή τη φορά χωρίς την αλληλεπίδραση των δύο παραγόντων (interaction term). Ανοίγουμε ξανά το παράθυρο για την Two-Way ANOVA, αλλά πέραν των άλλων επιλογών, ρυθμίζουμε και την επιλογή Model όπως φαίνεται παρακάτω:



Τρέχουμε το μοντέλο και παρατηρούμε πως και η Type και η Brand είναι στατιστικά σημαντικές, ενώ το R^2 έχει μια πολύ μικρή μείωση. Αυτό είναι το τελικό μας μοντέλο.

Dependent Variable: Acceleration									
Source	Type III Sum of Squares	df	Mean Square	F	Sig.				
Corrected Model	361.380 ^a	3	120.460	18.834	.000				
Intercept	70314.632	1	70314.632	10993.716	.000				
Brand	48.479	2	24.239	3.790	.023				
Туре	177.395	1	177.395	27.736	.000				
Error	2475.211	387	6.396						
Total	98246.722	391							
Corrected Total	2836.590	390							

Tests of Between-Subjects Effects

a. R Squared = .127 (Adjusted R Squared = .121)

Τέλος, κάνουμε τον έλεγχο κανονικότητας των καταλοίπων. Ο έλεγχος Shapiro-Wilk δίνει p - value = 0.005, επομένως η κανονικότητα απορρίπτεται σε επίπεδο σημαντικότητας 5% ή και 1%. Τα αποτελέσματά μας ισχύουν προσεγγιστικά λόγω του μεγάλου μέγεθους του δείγματος.

8 Παράρτημα

Μέτρα Θέσης και Διασποράς

Παρακάτω δίνονται τα σημαντικότερα περιγραφικά μέτρα, τα δειγματικά τους ανάλογα και οι τύποι υπολογισμού τους. Ιδιαίτερη προσοχή στη σύγχυση παραμέτρου (π.χ. μέση τιμή μ) και εκτιμήτριας (π.χ. δειγματικός μέσος \bar{X}).

Μέτρο	Εκτιμήτρια	Τύπος
Μέση Τιμή	Δειγματικός Μέσος	$\bar{X} = \frac{\sum_{i=1}^{n} X_i}{n}$
Διασπορά	Δειγματική Διασπορά	$S^2 = \frac{\sum_{i=1}^{n} (X_i - \bar{X})^2}{n-1}$
Τυπική Απόκλιση	Τυπικό Σφάλμα	$S = \sqrt{S^2}$
Ασυμμετρία	Δειγματική Ασυμμετρία	$l = \frac{\sum_{i=1}^{n} (X_i - \bar{X})^3}{nS^3}$
Κύρτωση	Δειγματική Κύρτωση	$k = \frac{\sum_{i=1}^{n} (X_i - \bar{X})^4}{nS^4}$

Table 1: Μέτρα Θέσης και Διασποράς