



Chapter Outline

- 1) Overview
- 2) Basic Concept
- 3) Relation to Regression and ANOVA
- 4) Discriminant Analysis Model
- 5) Statistics Associated with Discriminant Analysis
- 6) Multiple Discriminant Analysis
- 7) Stepwise Discriminant Analysis
- 8) Logit Regression

Simila Regre	rities and ession, and	Differences b I Discriminan	etween ANOVA, t Analysis]
	ANOVA	REGRESSION	DISCRIMINANT (Logit) ANALYSIS	
Similarities Number of dependent variables	One	One	One	
Number of independent variables	Multiple	Multiple	Multiple	
<u>Differences</u> Nature of the dependent variables	Metric	Metric	Categorical	
Nature of the independent variables	Categorical	Metric	Metric	



























	Inform	ation o	on Res	ort Visits	: Analy	sis San	nple
	-					Table 18.	2
No.	Resort Visit	Annual Family Income (\$000)	Attitude Toward Travel	Importance Attached to Family Vacation	Household Size	Age of Head of Household	Amount Spent on Family Vacation
1	1	50.2	5	8	3	43	M (2)
2	1	70.3	6	7	4	61	H (3)
3	1	62.9	7	5	6	52	H (3)
4	1	48.5	7	5	5	36	L (1)
5	1	52.7	6	6	4	55	H (3)
6	1	75.0	8	7	5	68	H (3)
7	1	46.2	5	3	3	62	M (2)
8	1	57.0	2	4	6	51	M (2)
9	1	64.1	/	5	4	57	H (3)
10	1	08.1 72.4	6	6	5	45	H (3)
11	1	73.4	5	/	2	44 64	п (3) Ч (3)
12	1	71.9 56.2	5	0	4	54	FI (3) M (2)
14	1	40.3	4	2	3	56	H (3)
15	1	62.0	5	6	2	58	H (3)

[1	Info	rmatio Ho	on on F Idout S	Resort ample	Visits: ; Table	18.3
No.	Resort Visit	Annual Family Income (\$000)	Attitude Toward Travel	Importance Attached to Family Vacation	Household Size	Age of Head of Household	Amount Spent on Family Vacation
1 2 3 4 5 6 7 8 9 10 11 12	1 1 1 1 2 2 2 2 2 2 2 2 2 2	50.8 63.6 54.0 45.0 68.0 62.1 35.0 49.6 39.4 37.0 54.5 38.2	4 7 6 5 6 5 4 5 6 2 7 2	7 4 6 6 3 3 5 6 3 2	3 7 4 3 6 3 4 5 3 5 3 3 3	45 55 58 60 46 56 54 39 44 51 37 49	M(2) H (3) M(2) H (3) H (3) L (1) L (1) H (3) L (1) M(2) L (1)





[Results	; of Two-C	Group Discrimi	nant Analysis Table 18.4 cont.
	CANONICAL	DISCRIMINANT	FUNCTIONS	
Function 1*	Eigenvalue 1.7862	% of Cum Variance %	n Canonical After V Correlation Function : 0 0.00 0.8007 :	Nilks' λ Chi-squaredfSignificance0.358926.13050.0001
* marks t	he 1 canonica	l discriminant fur	nctions remaining in the a	inalysis.
	Standard Ca	nonical Discrimin	nant Function Coefficients	
	F	UNC 1		
INCON TRAVE VACAT HSIZE AGE	ME O EL O TION O O O	.74301 .09611 .23329 .46911 .20922		
Pooled wi (variables	Structure Ma thin-groups co ordered by si	atrix: orrelations betwe ize of correlation	een discriminating variabl within function)	les & canonical discriminant functions
	F	UNC 1		
INCO HSIZE VACAT TRAVE AGE	ME 0 TION 0 EL 0	.82202 .54096 .34607 .21337 .16354		

[Re	esults of T	wo-Group Di	scrimina	Int Anal	ysis
				Table 1	8.4 cont. 📕
INCOME TRAVEL VACATIOI HSIZE AGE (constant	Unstandardized FUNG 0.49 0.49 0.42 0.42 0.42 0.42 0.42 0.42 0.42 0.42	Canonical Discriminan C 1 54710E-01 54455E-01 22813 54380E-01 5476 minant functions evalu C 1 9118 9118 sults for cases selected	t Function Coe ated at group I for use in ana	fficients means (group lysis	centroids)
	Actual Group	Predicted No. of Cases	Group Me 1	embership 2	
				_	
Group	1	15	12 80.0%	3 20.0%	
Group Group	1 2	15 15	12 80.0% 0 0.0%	3 20.0% 15 100.0%	









		- F			aryono	
(Table 1	8.5	
6 m m						
AMOUNT	INCOME	TRAVEL V	ACATION H	SIZE AG	E	
1	38 57000	4 50000	4 70000	3 10000	50 30000	
2	50.11000	4.00000	4.20000	3.40000	49.50000	
3	64.97000	6.10000	5.90000	4.20000	56.00000	
Total	51.21667	4.86667	4.93333	3.56667	51.93333	
Group	Standard Devi	ations				
1	5.29718	1.71594	1.88856	1.19722	8.09732	
2	6.00231	2.35702	2.48551	1.50555	9.25263	
3	8.61434	1.19722	1.66333	1.13529	7.60117	
Total	12.79523	1.97804	2.09981	1.33089	8.57395	
Poole	d Within-Group	s Correlatio	n Matrix			
	INCOME	TRAVEL	VACATION	HSIZE	AGE	
INCOME	1.00000					
TRAVEL	0.05120	1.00000				
VACATION	0.30681	0.03588	1.00000			
HSIZE	0.38050	0.00474	0.22080	1.00000		
AGE	-0.20939	-0.34022	-0.01326	-0.02512	1.00000	

ſ	Results	of Thr	ree-G	roup [Discrim	ninant _{Tabl}	Analys e 18.5 c	sis ont.]
	Wilks' (U-s	tatistic) and	d univaria	ite <i>F</i> ratio v	vith 2 and 2	27 degrees	of freedom	ı.	
	Variable	Wilks' Lam	ıbda	F	Significan	ce			
	INCOME TRAVEL VACATION HSIZE AGE	0.26215 0.78790 0.88060 0.87411 0.88214		38.00 3.634 1.830 1.944 1.804	0.0000 0.0400 0.1797 0.1626 0.1840				
	CANONICAL	DISCRIMI	NANT FU	NCTIONS					
Function	Eigenvalue	% of Variance	Cum (%	Canonical Correlatio	After N n Function : 0	Wilks' λ Chi 0.1664	-square df 44.831	Signif 10	icance 0.00
1*	3.8190	93.93	93.93	0.8902	: 1	0.8020	5.517	4	0.24
2*	0.2469	6.07	100.00	0.4450	:				
* marks t	he two canon	ical discrim	inant fun	ctions rem	aining in th	e analysis.			
	Standardize	d Canonical	l Discrimi	nant Funct	ion Coeffici	ents			
		F	UNC 1	FUNC	2				
			1.04740	-0.4	2076 6851				
	VACATION	-(0.14198	0.7	3354				
	HSIZE	-(0.16317	0.1	2932				



Resi	uits o	t Inree-C	PLOUD DIS	scrimina	ant Ar	iaiysis
				Та	ble 18.	5 cont.
C	Classific	ation Results	51	Dradistad	Crown M	omborchin
		Actual Group	No. of Cases	1	2	<u>3</u>
	Group	1	10	9 90.0%	1	0
	Group	2	10	1 10.0%	9 90.0%	0 0.0%
	Group	3	10	0	2	8
P	Percent of	f grouped cases	correctly classif	fied: 86.67%	20.0 /0	00.070
C	Classifi	cation resu	lts for hold	out samp	le	
		Actual Group	No. of Cases	Predicted 1	Group M 2	embership <u>3</u>
	Group	1	4	3 75.0%	1 25.0%	0 0.0%
	Group	2	4	0 0.0%	3 75.0%	1 25.0%
	Group	3	4	1 25.0%	0 0.0%	3 75.0%













var Weight	les;	Length?	Length3 F	un; Meight Widt	h.
run.	Dengeni	Lengenz	Lengens I	leigne hiue	,
-	Class	Level Infor	mation		
	Variable				
Species	Name	Frequency	Weight	Proportion	
Ducan	Deserv	24	24 0000	0 215100	
Bream	Bream Darkki	34	11 0000	0.215190	
Perch	Perch	56	56.0000	0.354430	
Pike	Pike	17	17.0000	0.107595	
Roach	Roach	20	20.0000	0.126582	
Smelt	Smelt	14	14.0000	0.088608	
Whitefish	Whitefish	6	6.0000	0.037975	
Mult	ivariate Sta S=6	M=-0.5	F Approximat	ions	
Statistic		Value F	7alue Num	DF Den DF	Pr >
Wilks' Lambda	0.000	36325	90.71	36 643.89	<.000
Pillai's Trace	3.104	65132	26.99	36 906	<.000
Hotelling-Lawley Trace	e 52.057	99676 2	09.24	36 413.64	<.000
	20 124	00000	04 00	c	

pr ru	coc cl va un;	candisc da ass Specie r Weight L	uta=fish nc es; wength1 Len	an=3 out=ou gth2 Length	itcan; n3 Height Width;	
		The	ANDISC Proce	dure		
			Adjusted	Approximate	Squared	
		Canonical	Canonical	Standard	Canonical	
		Correlation	Correlation	Error	Correlation	
;	1	0.987463	0.986671	0.001989	0.975084	
	2	0.952349	0.950095	0.007425	0.906969	
	3	0.838637	0.832518	0.023678	0.703313	
	4	0.633094	0.623649	0.047821	0.400809	
	5	0.344157	0.334170	0.070356	0.118444	
	6	0.005701		0.079806	0.000033	
			Eigenvalues	of Inv(E) *H		
			= CanRsq/	(1-CanRsq)		
		Eigenvalue	Difference	Proportion	Cumulative	
	1	39.1350	29.3859	0.7518	0.7518	
	2	9.7491	7.3786	0.1873	0.9390	
	3	2.3706	1.7016	0.0455	0.9846	
	4	0.6689	0.5346	0.0128	0.9974	
	5	0.1344	0.1343	0.0026	1.0000	
	6	0.0000		0.0000	1.0000	

proc c cla var	andisc data ss Species; Weight Len	=fish ncan= gth1 Length	3 out=ou 2 Lengtl	ıtcan; n3 Heigh	t Width;	
run •	5	5 5	2	5		
ran,						
	Test of HU:	The canonical c	orrelation	s in the		
	current ro	The canonical c w and all that	orrelation follow are	s in the zero		
	current ro	The canonical c w and all that	orrelation follow are	s in the zero		
	test of HU: current ro Likelihood	The canonical c w and all that Approximate	orrelation follow are	s in the zero		
	Test of HU: current ro Likelihood Ratio	The canonical c w and all that Approximate F Value	orrelation follow are Num DF	zero Den DF	Pr > F	
	Test of HO: current ro Likelihood Ratio	The canonical c w and all that Approximate F Value	orrelation follow are Num DF	Den DF	Pr > F	
1	Test of HO: current ro Likelihood Ratio 0.00036325	The canonical c w and all that Approximate F Value 90.71	orrelation follow are Num DF 36	Den DF	Pr > F <.0001	
1 2	Likelihood Ratio 0.00036325 0.01457896	The canonical c w and all that Approximate F Value 90.71 46.46	orrelation follow are Num DF 36 25	Den DF 643.89 547.58	Pr > F <.0001 <.0001	
1 2 3	Test of H0: current ro Likelihood Ratio 0.00036325 0.01457896 0.15671134	The canonical c w and all that Approximate F Value 90.71 46.46 23.61	orrelation follow are Num DF 36 25 16	Den DF 643.89 547.58 452.79	Pr > F <.0001 <.0001 <.0001	
1 2 3 4	Test of H0: current ro Likelihood Ratio 0.00036325 0.01457896 0.15671134 0.52820347	The canonical c w and all that Approximate F Value 90.71 46.46 23.61 12.09	orrelation follow are Num DF 36 25 16 9	Den DF 643.89 547.58 452.79 362.78	Pr > F <.0001 <.0001 <.0001 <.0001	
1 2 3 4 5	Test of H0: current ro Likelihood Ratio 0.00036325 0.01457896 0.1567134 0.52820347 0.88152702	The canonical co w and all that Approximate F Value 90.71 46.46 23.61 12.09 4.88	orrelation follow are Num DF 36 25 16 9 4	Den DF 643.89 547.58 452.79 362.78 300	<pre>Pr > F <.0001 <.0001 <.0001 <.0001 0.0008</pre>	



٦

The first canonical variable, Can1, shows that the linear combination of the centered variables Can1= $-0.0006 \times \text{Weight} - 0.33 \times \text{Length1} - 2.49 \times \text{Length2} + 2.60 \times \text{Length3} + 1.12 \times \text{Height} - 1.45 \times \text{Width separates the species most effectively (see Figure 21.5).}$

Can	Can2	Canl	Variable
-0.005596193	-0.005231659	-0.000648508	Weight
-2.93432410	-0.626598051	-0.329435762	Length1
4.04503889	-0.690253987	-2.486133674	Length2
-1.13926491	1.803175454	2.595648437	Length3
0.28320255	-0.714749340	1.121983854	Height
0.74148668	-0.907025481	-1.446386704	Width

PROC CANDISC computes the means of the canonical variables for each class. The first canonical variable is the linear combination of the variables Weight, Length1, Length2, Length3, Height, and Width that provides the greatest difference (in terms of a univariate F-test) between the class means. The second canonical variable provides the greatest difference between class means while being uncorrelated with the first canonical variable.

Fish Measurement Data

The CANDISC Procedure

Class Means on Canonical Variables

Species	Canl	Can2	Can3
Bream	10.94142464	0.52078394	0.23496708
Parkki	2.58903743	-2.54722416	-0.49326158
Perch	-4.47181389	-1.70822715	1.29281314
Pike	-4.89689441	8.22140791	-0.16469132
Roach	-0.35837149	0.08733611	-1.10056438
Smelt	-4.09136653	-2.35805841	-4.03836098
Whitefish	-0.39541755	-0.42071778	1.06459242





В

proc discrim data=fish; class Species; run;

The coefficients of the linear discriminant function are displayed (in Figure 31.4) with the default options METHOD=NORMAL and POOL=YES.

Figure 31.4 Linear Discriminant Function

		Linear I	Discriminant	Function for	r Species		
Variable	Bream	Parkki	Perch	Pike	Roach	Smelt	Whitefish
Constant	-185.91682	-64.92517	-48.68009	-148.06402	-62.65963	-19.70401	-67.44603
Weight	-0.10912	-0.09031	-0.09418	-0.13805	-0.09901	-0.05778	-0.09948
Length1	-23.02273	-13.64180	-19.45368	-20.92442	-14.63635	-4.09257	-22.57117
Length2	-26.70692	-5.38195	17.33061	6.19887	-7.47195	-3.63996	3.83450
Length3	50.55780	20.89531	5.25993	22.94989	25.00702	10.60171	21.12638
Height	13.91638	8.44567	-1.42833	-8.99687	-0.26083	-1.84569	0.64957
Width	-23.71895	-13.38592	1.32749	-9.13410	-3.74542	-3,43630	-2.52442

A summary of how the discriminant function classifies the data used to develop the function is displayed last. In Figure 31.5, you see that only three of the observations are misclassified. The error-count estimates give the proportion of misclassified observations in each group. Since you

	Numbe	er of Obse	rvations a	nd Percent	Classifie	d into Spe	cies	
From								
Species	Bream	Parkki	Perch	Pike	Roach	Smelt	Whitefish	Total
Bream	34	0	0	0	0	0	0	34
	100.00	0.00	0.00	0.00	0.00	0.00	0.00	100.00
Parkki	0	11	0	0	0	0	0	11
	0.00	100.00	0.00	0.00	0.00	0.00	0.00	100.00
Perch	0	0	53	0	0	3	0	56
	0.00	0.00	94.64	0.00	0.00	5.36	0.00	100.00
Pike	0	0	0	17	0	0	0	17
	0.00	0.00	0.00	100.00	0.00	0.00	0.00	100.00
Roach	0	0	0	0	20	0	0	20
	0.00	0.00	0.00	0.00	100.00	0.00	0.00	100.00
Smelt	0	0	0	0	0	14	0	14
	0.00	0.00	0.00	0.00	0.00	100.00	0.00	100.00
Whitefish	0	0	0	0	0	0	6	6
	0.00	0.00	0.00	0.00	0.00	0.00	100.00	100.00
Total	34	11	53	17	20	17	6	158
	21.52	6.96	33.54	10.76	12.66	10.76	3.80	100.00
Paul and								

С	•	proc c run;	c ste <u>r</u> class ;	disc Speci	data=: es;	fish;				
L	-			The	STEPDISC	Procedur	e			
				stepw	ise Selec	tion Sum	mary			
N Step 1 2 3 4 5 6	umber In 2 3 4 5 6	Entered Height Length2 Length3 Width Weight Length1	Removed	Partial R-Square 0.7553 0.9229 0.8826 0.5775 0.4461 0.2987	F Value 77.69 299.31 186.77 33.72 19.73 10.36	Pr > F <.0001 <.0001 <.0001 <.0001 <.0001	Wilks' Lambda 0.24466983 0.01886065 0.00221342 0.00093510 0.00051794 0.00036325	Pr < Lambda <.0001 <.0001 <.0001 <.0001 <.0001	Average Squared Canonical Correlation 0.12588836 0.25905822 0.38427100 0.45200732 0.49488458 0.51744189	Pr > ASCC <.0001 <.0001 <.0001 <.0001 <.0001
Fiç All	Figure 67.5. Step Summary All the variables in the data set are found to have potential discriminatory power.						ſ.			

596	PART III • DATA COLLECTION, PREPARATION, ANALYSIS, AND REPORTING
	Edition. The Multivariate>Discriminant Analysis task offers both two-group and multiple discriminant analysis. Both two-group and multiple discriminant analysis can be performed using the Discriminant Analysis task within the SAS Learning Edition. To select this task click:
	The steps for running three-group discriminant analysis are similar to these steps. To run logit analysis or logistic regression using the SAS Learning Edition, click:
	Analyze>Regression>Logistic
	The following are the detailed steps for running logit analysis with brand loyalty as the dependent variable and attitude toward the brand, attitude toward the product category, and attitude toward shopping as the independent variables using the data of Table 18.6.
	 Select ANALYZE from the SAS Learning Edition menu bar. Click REGRESSION and then LOGISTIC. Move LOYALTY to the Dependent variable task role. Move BRAND, PRODUCT, and SHOPPING to the Quantitative variables task role. Select MODEL EFFECTS. Choose BRAND, PRODUCT, and SHOPPING as Main Effects. Select MODEL OPTIONS. Check SHOW CLASSIFICATION TABLE and enter 0.5 as the critical probability value. Click RUN.



	Standardized canonical coefficients	F	value	Prob > F	
ART	-0.0517		1.08	0.3014	
CGS/S	-0.2782		0.15	0.6958	
APT	-0.1720		0.66	0.4184	
INVT	0.1710		3.74	0.0552*	
R&D/S	-0.2515	-	23.85	$< 0.0001^{***}$	
SG&A/S	-0.2504	1	37.89	< 0.0001***	
Dep/S	-0.3525		2.01	0.1587	
FAT	-0.1317		0.14		
Tax/S	0.2282		19.35		
Eigenvalue	Canonical correlation	Likelihood ratio	F value	Prob > F	
1.1121	0.725628	0.47346	15.82	< 0.0001	
Classification results	used for cross-validation ^a				
Groups	Competitive adv	antage Competitive disadvantage		Tota	
a	e 64		13	77	
Competitive advantag			46	61	







Health Care Manage Sci (2008) 11:353–358 DOI 10.1007/s10729-008-9054-y

Detecting hospital fraud and claim abuse through diabetic outpatient services

Fen-May Liou • Ying-Chan Tang • Jean-Yi Chen

Received: 3 July 2007 / Accepted: 8 January 2008 / Published online: 19 January 2008 Springer Science + Business Media, LLC 2008

Abstract Hospitals and health care providers tend to get involved in exaggerated and fraudulent medical claims initiated by national insurance schemes. The present study applies data mining techniques to detect fraudulent or **1** Introduction

Healthcare fraud and abuse are of major concern in many countries, in some cases costing public and private financial

			Health Care N	fanage Sci (2008) 11:353–3
e 1 Descriptive statistics ormal and fraudulent	Variable (per case)	Normal Hospital		Fraudulent Hospital	
ospitals		Mean	SD	Mean	SD
	Average days of drug dispense	7.72	5.60	7.39	1.50
and To	Average drug cost	221.63	274.06	208.25	88.1
	Average consultation and treatment fees	358.71	176.88	259.58	113.
	Average diagnosis fees	265.42	42.93	265.00	43.4
	Average dispensing service fees	24.48	8.13	30.01	11.2
100	Average medical expenditure	548.04	408.33	584.75	145.
	Average amount claimed	487.99	394.69	511.81	131.
	Average drug cost per day	28.81	27.79	33.82	10.0
	Average medical expenditure per day	134.37	92.29	173.13	73.5

claim was associated with the value "0" if regular, and "1" if irregular.

Г

Stepwise logistic regression was performed on each variable individually to identify the most effective factors

weighted sum of the input variables, and transform that sum to an output signal using some kind of threshold function (typically a step function or sigmoid) [30, 31]. The output layer (often a single node) receives a weighted sum of the

Г	THE WALL STREET JOURNAL.
_	2013/11/19 07:43:43
	正文 評論(1) □ 投稿 ● 打印 ● 韓登 ● ● ● ● ● ● ● ● ● ● ● ● ● ● ●
	查潔林 ➡, R上有些人具備獨特的洞察力,可以從腐朽中發現神奇,物理學家音萊斯(Derek 」John de Solla Price)就是其中一位。20世紀60年代初,普萊斯的議角累積了一堆老 審論文集,他在隨便翻閱中發現了一条列有關論文的數字規律,例如在所發表的論文中, 約35%的論文從來未被引用,49%的論文被引用一次,9%被引用兩次,3%被引用三次, 2%被引用四次,而被引用次數大於導於五次的論文,共佔2%。他進而對這個現象進行了 条統地分析並發表了裡程碑性質的論文:科學的網絡。他發現的規律後來被很多人用不同 的論文集合驗証,數字號然略有不同,但規律不會變,因為這反映了自組織社會的過性, 被現在的學者們稱之為"無尺度網絡"。
	無獨有偶,美國結構語言學博士尤金·加菲爾德(Eugene Garfield)在普萊斯之前也發現 論文的引文關系可以提供很多重要的信息,例如論文的被引數可以反映論文的質量。因 此,他成立了一家公司,專門收集和發布學術論文家引的統計數據,也即現在學術界流行 的SC和SSCI文獻家引數據庫,如今這兩個統計系統已經成為國際上各類學術期刊、論文 和學者排行榜的重要參照,論文被引數也被認可為較好地反映了論文的原創性和權威性。