

UNIVERSIDADE FEDERAL DA PARAÍBA 00
CENTRO DE CIÊNCIAS E TECNOLOGIA
CURSO DE MESTRADO EM ENGENHARIA CIVIL

MÉTODOS ESTATÍSTICOS MULTIVARIADOS
APLICADOS À ANÁLISE DE TRANSPORTES

por

WALTER SANTA CRUZ

CAMPINA GRANDE - PB

MAIO DE 1983

MÉTODOS ESTATÍSTICOS MULTIVARIADOS
APLICADOS À ANÁLISE DE TRANSPORTES



C957m Cruz, Walter Santa.
Métodos estatísticos multivariados aplicados à análise de transportes / Walter Santa Cruz. - Campina Grande, 1983. 92 f.

Dissertação (Mestrado em Ciências) - Universidade Federal da Paraíba, Centro de Ciências e Tecnologia, 1983. "Orientação : Prof. Dr. José Eugênio Leal".
Referências.

1. Transportes - Planejamento e Organização. 2. Transportes - Métodos Estatísticos. 3. Técnicas Estatísticas Multivariadas. 4. Dissertação - Ciências. I. Leal, José Eugênio. II. Universidade Federal da Paraíba - Campina Grande (PB). III. Título

CDU 656.07:519.2(043)

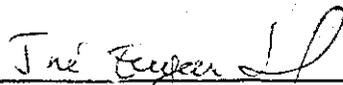
MÉTODOS ESTATÍSTICOS MULTIVARIADOS
APLICADOS À ANÁLISE DE TRANSPORTES

WALTER SANTA CRUZ
ENGENHEIRO CIVIL

DISSERTAÇÃO SUBMETIDA AO CORPO DOCENTE DA COORDENAÇÃO
DOS PROGRAMAS DE PÓS-GRADUAÇÃO E PESQUISA DO CENTRO DE CIÊN
CIAS E TECNOLOGIA DA UNIVERSIDADE FEDERAL DA PARAIBA COMO
PARTE DOS REQUISITOS NECESSÁRIOS PARA A OBTENÇÃO DO GRAU
DE MESTRE EM CIÊNCIAS (M.Sc.).

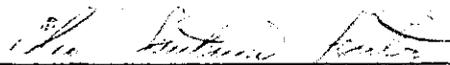
Aprovado por:

COMISSÃO



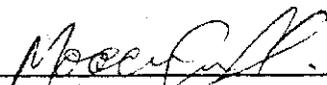
JOSE EUGENIO LEAL - Dr. Ing.

-Presidente-



ÉLIO SANTANA FONTES - M.Sc.

-Examinador Interno-



MOACIR GUILHERMINO DA SILVA - M.Sc.

-Examinador Externo-

CAMPINA GRANDE
ESTADO DA PARAÍBA-BRASIL
MAIO - 1983

À meus pais, irmãos, esposa e
filhas.

AGRADECIMENTOS

Ao orientador, Professor Dr. José Eugênio Leal pelo incentivo, assistência e orientação prestadas na elaboração deste trabalho.

Ao Professor Sebastião Guimarães Vieira, Pró-Reitor para Assuntos do Interior, pelo apoio e confiança depositada no autor.

Ao Professor Jean Pierre Demartinecourt pela valiosa ajuda na parte matemática deste trabalho.

Ao Professor Dr. Edmilson de Vasconcelos Pontes do CCEN/UFAL, pela orientação acadêmica prestada, durante a graduação do autor.

À equipe do Núcleo de Processamento de Dados do Centro de Ciências e Tecnologia - Campina Grande, pela ajuda prestada nos trabalhos de computação.

A todos os outros que, direta ou indiretamente, contribuíram para a realização deste trabalho.

MÉTODOS ESTATÍSTICOS MULTIVARIADOS
APLICADOS À ANÁLISE DE TRANSPORTES

DISSERTAÇÃO DE MESTRADO

POR

WALTER SANTA CRUZ

R E S U M O

Este trabalho tem como objetivo verificar a aplicabilidade das técnicas estatísticas multivariadas, Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante em estudos de análise de transportes.

São mostrados os enfoques matemáticos para cada uma dessas técnicas, e realizadas aplicações e discussões dos resultados em diversos problemas da análise de transportes.

MULTIVARIATE STATISTICAL METHODS
APPLIED TO TRANSPORTATION ANALYSIS

M.Sc. DISSERTATION

BY

WALTER SANTA CRUZ

A B S T R A C T

The objective of this work is to verify the applicability of the multivariate statistical techniques as Multiple Linear Regression Analysis, Factor Analysis and Discriminant Analysis in studies of transportation analysis.

It has been showed the mathematical approaches for each of these techniques as well as applications and discussions of the results for several problems of transportation analysis.

ÍNDICE

Capítulo	I - Introdução	
1.1	- Considerações Gerais	01
1.2	- Apresentação do Problema	01
1.3	- Objetivos	06
Capítulo	II - As Técnicas Estatísticas Multivariadas: Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante.	
2.1	- Introdução	07
2.2	- Análise de Regressão Linear Múltipla	07
2.2.1	- Considerações Iniciais Sobre a Análise de Regressão Linear Múltipla	07
2.2.2	- Enfoque Matemático da Análise de Regressão Linear Múltipla	09
2.2.3	- Interpretação dos Coeficientes de Regressão e da Parte Cons- tante de uma Equação de Regres- são	14
2.2.4	- Análise da Qualidade de uma Equação de Regressão Linear Múltipla	15

2.2.4.1 - Análise Numérica da Equação	16
2.2.4.2 - Análise Qualitativa da Equação	18
2.2.4.3 - Testes Estatísticos	19
2.3 - Análise Fatorial	21
2.3.1 - Considerações Iniciais Sobre a Análise Fatorial	21
2.3.2 - Enfoque Matemático da Análise Fatorial	22
2.3.3 - Rotação dos Fatores	29
2.4 - Análise Discriminante	31
2.4.1 - Considerações Iniciais Sobre a Análise Discriminante	31
2.4.2 - Enfoque Matemático da Análise Discriminante	32
2.4.3 - Interpretação Bayesiana da Análise Discriminante	37
2.4.3.1 - Análise Discriminante Para Dois Grupos. Uso da Razão de Verossimilhança	39
2.4.3.2 - Análise Discriminante Para Mais de Dois Grupos	41
Capítulo III - Aplicações e Resultados	
3.1 - Introdução	43

3.2 - Análise de Regressão Linear Múltipla.	
Aplicações e Resultados	43
3.2.1 - Aplicação 1	43
3.2.1.1 - Nível de Setor	46
3.2.1.2 - Nível de Zona	48
3.2.1.3 - Nível de Domicílio	49
3.2.1.4 - Considerações Sobre a Aplicação 1	50
3.2.2 - Aplicação 2	51
3.3 - Análise Fatorial - Aplicações e Resultados	56
3.3.1 - Aplicação 1	64
3.3.2 - Aplicação 2	64
3.4 - Análise Discriminante. Aplicações e Resultados	67
3.4.1 - Aplicação 1	67
3.4.2 - Aplicação 2	76
Capítulo IV - Conclusões e Sugestões para Pesquisas Futuras.	
4.1 - Conclusões	82
4.2 - Sugestões para Pesquisas Futuras	84

LISTA DE TABELAS

Tabela	Página
01 - Setores e suas zonas	44
02 - Matriz de correlação entre as variáveis a nível de setor	45
03 - Matriz de correlação entre as variáveis a nível de zona	45
04 - Matriz de correlação entre as variáveis a nível de domicílio	46
05 - Níveis de renda e número médio de automóveis por família em cada classe	52
06 - Número de automóveis e número de famílias, em cada nível de renda, por zona	54
07 - Matriz de correlação	53
08 - Numeração dos fatores, autovalores e porcentagem da variância do conjunto das variáveis explicada por cada fator	60
09 - Matriz dos carregamentos dos fatores após a rotação varimax	61
10 - Comunalidade das variáveis	61
11 - Matriz dos coeficientes para o cálculo dos escores dos fatores	62
12 - Fatores mais importantes na composição da variância das variáveis	63

Tabela	Página
13 - Fatores e suas definições em termos das variáveis originais dadas em valores relativos ...	63
14 - Numeração dos fatores, autovalores e porcentagem da variância do conjunto das variáveis explicada por cada fator	66
15 - Matriz dos carregamentos dos fatores após rotação varimax	67
16 - Comunalidade das variáveis	68
17 - Matriz dos coeficientes para o cálculo dos escores dos fatores	69
18 - Fatores mais importantes na composição da variância das variáveis	65
19 - Fatores e suas definições em termos das variáveis originais predominantes, em valores relativos.....	70
20 - Número de casos por grupo	72
21 - Lambda de Wilks e valor do F equivalente.....	72
22 - Valores de F e significância entre pares de grupos. Cada valor F tem 11 e 2451 graus de liberdade	73
23 - Variáveis discriminatórias	75
24 - Funções discriminantes	76
25 - Coeficientes na forma padronizada, das funções discriminantes	77

Tabela	Página
26 - Coeficientes das funções de classificação....	78
27 - Classificação das observações	80
28 - Resultados da classificação	81
29 - Número de casos por grupo	82
30 - Lambda de Wilks e $F_{\text{equivalente}}$	82
31 - Valores de F, para 8 e 229 graus de liberdade e significância para os pares de grupos	83
32 - Variáveis discriminatórias	84
33 - Coeficientes, na forma padronizada, das fun - ções discriminantes	85
34 - Classificação das observações	86
35 - Resultado da classificação	87

LISTA DE QUADROS

Quadro	Página
01 - Classificação do uso do solo	57
02 - Numeração e nomes dos bairros	58

CAPÍTULO I

INTRODUÇÃO

1.1 - Considerações Gerais

Na década de 50, o déficit da ferrovia brasileira e a criação da indústria automobilística, no Brasil, com perspectivas de desenvolvimento a curto prazo, contribuíram para o país relegar a segundo plano o sistema ferroviário e voltar a atenção para a construção de uma infra-estrutura viária capaz de absorver o tráfego de veículos automotores e permitir condições de fluxos, aceitáveis pelos usuários.

Com isso, todos os esforços se concentraram nas capacidades das rodovias e o planejamento dos transportes tinha, como um dos objetivos fundamentais, a determinação da demanda atual de tráfego, nas rodovias mais importantes, e o cálculo de coeficientes que permitissem estimar o tráfego futuro. O que se procurava, na realidade, era identificar o volume de tráfego por automóvel e, assim, toda estimativa futura era para prover as rodovias de condições adequadas ao uso do transporte individual.

1.2 - Apresentação do Problema

No Brasil, as cidades de médio e pequeno portes, ge-

ralmente, carecem de um planejamento sistêmico dos transportes. Ora os estudos realizados não são implantados devido a limitações de recursos, dos mais variados tipos, ora a deficiência na obtenção dos dados ou o uso de modelos inadequados produz resultados imprecisos.

No que concerne a parte de modelagem, tem-se, na fase de geração de viagens, os modelos à base do cálculo da correlação linear como os mais usados para explicar as viagens geradas pelas zonas de tráfego. Considera-se determinadas variáveis, em geral, agregadas a nível de zona, como explicativas das ocorrências de viagens zonais e busca-se, através do método dos mínimos quadrados, uma equação de regressão linear múltipla que associe as viagens geradas às variáveis do problema.

Encontrada a equação de regressão linear múltipla, o passo seguinte é decidir se essa equação se presta para estimar a demanda futura de viagens. Essa decisão deve ser tomada com base na análise do coeficiente de determinação múltipla, no erro padrão de estimativa e na significância estatística dos coeficientes de regressão parcial. É comum, no entanto, decidir-se sobre a qualidade de uma equação de regressão linear múltipla baseando-se, apenas, no valor do coeficiente de determinação múltipla. Isso deve ser evitado, pois, pode conduzir a erros.

Kassoff & Deutschman⁹ estudam a alteração sofrida por uma equação de regressão linear proveniente de dois aspectos das variáveis observadas: o primeiro, utiliza os valores totais das variáveis, para o nível de agregação considerado; o

segundo, utiliza os valores relativos das variáveis, isto é, os valores da razão entre os totais agregados de cada variável e o total agregado de uma variável específica. Eles mostram que as equações com as variáveis dadas por seus valores totais agregados são, apesar de possuírem o coeficiente de determinação múltipla maior do que o das equações com as variáveis dadas por seus valores relativos, qualitativamente, inferiores a estes últimos tipos de equações. Isso comprova que se pode incorrer em erro quando se julga uma equação de regressão linear múltipla com base, somente, no seu coeficiente de determinação.

Se a equação de regressão linear obtida possui um coeeficiente de determinação múltipla próximo da unidade, um pequeno erro padrão de estimativa e os coeficientes de regressão parcial, estatisticamente significantes a um nível de significância estabelecido, adota-se essa equação como equação de previsão de viagens.

De posse da equação para previsão de viagens, as viagens no ano meta são calculadas, simplesmente, atribuindo-se às variáveis independentes os seus respectivos valores futuros. Em assim procedendo, admite-se a hipótese de que os coeeficientes da equação são imutáveis, com o tempo, o que não parece razoável. Em outras palavras, significaria dizer que, dispondo dos valores reais das variáveis dependente e independentes, no ano meta, e calculando-se uma nova equação de regressão linear múltipla, essa equação não diferiria, substancialmente, da obtida no ano base. A inalteração dos coeeficientes de regressão parcial, para o ano meta, mostra a

insensibilidade, às mudanças que possam ocorrer em uma área de estudo, dos modelos de regressão linear para previsão de viagens.

Por outro lado, reportando-se novamente ao trabalho de Kassoﬀ & Deutschman ⁹, eles alertam para um outro tipo de problema: o da escolha das variáveis.

Não basta que as variáveis envolvidas numa equação de regressão linear múltipla satisfaçam as premissas básicas da técnica da regressão linear. É necessário que haja uma relação de causa-efeito entre cada variável independente e a variável dependente. Torna-se evidente, pois, que a escolha das variáveis, para explicar um dado fenômeno, é fundamental.

Não somente a escolha das variáveis, mas também, o nível de agregação em que elas se encontram, isto é, se os valores das variáveis são obtidos a nível de zona, domicílio, etc., afeta a qualidade de uma equação de regressão linear múltipla.

// McCarthy ⁷ mostra que, à medida que o nível de agregação das variáveis diminui, a qualidade de uma equação de regressão linear aumenta, porque a parcela correspondente à variância interna de cada variável vai se tornando menor e, conseqüentemente, maior é a contribuição da variância, entre as variáveis, para a variância total dos dados. Sabe-se que a variância total de um conjunto é igual a variância interna de cada elemento somada com a variância existente entre os elementos.

Portanto, pode-se questionar o uso da Análise de Regressão Linear Múltipla em estudos de demanda de transportes porque as variáveis observadas são, geralmente, agregadas a

nível de zona de tráfego e com os seus valores representando os totais agregados àquele nível de agregação.

Em um processo de planejamento dos transportes, composto das fases de geração, distribuição, divisão modal e alocação de viagens, em que a primeira fase fornece os modelos de regressão linear para a previsão de viagens, é de se esperar que os resultados obtidos careçam de precisão. A geração de viagens fornece os dados de entrada para os modelos, normalmente, usados na fase de distribuição, e esta, por sua vez, tem seus resultados como elementos de entrada nos modelos das fases posteriores. Assim, as falhas na geração de viagens se ampliam nas fases subsequentes.

Então, necessário se faz o estudo e aplicação de técnicas estatísticas mais refinadas a estudos de análise de transportes.

Além da Análise de Regressão Linear Múltipla, utiliza-se, neste trabalho, as técnicas multivariadas Análise Fatorial e Análise Discriminante. Estas duas últimas técnicas, quando comparadas com a primeira, apresentam vantagens, a saber:

- Abordam a complexidade do processo de viagem.
- Incorporam a variância produzida, tanto por fatores sociais (internos à zona), como por fatores locais (entre as zonas).
- Representam a viagem de um indivíduo não em termos de uma variável isolada (análise univariada), mas de uma série de variáveis (análise multivariada) representativas do comportamento do indivíduo.

- Permitem uma análise mais aprofundada do fenômeno dos transportes, em vista de analisar os fatores-causas da diferença de comportamento do indivíduo, em relação aos transportes.

1.3 - Objetivos

Os objetivos desse trabalho são:

- 1 - Apresentar as técnicas estatísticas multivariadas, Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante, com os seus respectivos desenvolvimentos matemáticos. Essa abordagem matemática objetiva dar uma visão dos métodos utilizados e permite uma melhor compreensão de cada método. Não se pretende, neste trabalho, esgotar a formulação matemática dos métodos tratados. O tratamento matemático formal pode ser visto na literatura especializada, citada na bibliografia.
- 2 - Propiciar uma visão crítica do uso da Análise de Regressão Linear Múltipla em estudos de demanda de viagens.
- 3 - Aplicar as técnicas estatísticas multivariadas, Análise Fatorial e Análise Discriminante, a estudos de análise de transportes. Esse constitui o principal objetivo desse trabalho.

CAPÍTULO II

AS TÉCNICAS ESTATÍSTICAS MULTIVARIADAS: ANÁLISE DE REGRESSÃO LINEAR MÚLTIPLA, ANÁLISE FATORIAL E ANÁLISE DISCRIMINANTE

2.1 - Introdução

Análise Multivariada compreende um conjunto de procedimentos estatísticos, dentre estes a Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante, que consistem em analisar propriedades distintas, observadas simultaneamente, acerca de cada elemento de uma amostra ou de uma população.

Este capítulo trata do desenvolvimento matemático das técnicas estatísticas multivariadas, Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante.

Na seção referente às considerações iniciais sobre cada técnica, dá-se uma visão geral da técnica apresentada. Em seguida, apresenta-se o enfoque matemático e os testes estatísticos, quando houver, para cada técnica.

2.2 - Análise de Regressão Linear Múltipla

2.2.1 - Considerações iniciais sobre a Análise de Regressão Linear Múltipla.

Os estudos de demanda para o planejamento dos transportes são, usualmente, baseados em modelos compostos de quatro etapas principais, a saber: geração, distribuição, divisão modal e alocação de viagens.

Tem-se, na fase de geração de viagens, os modelos à base do cálculo da correlação linear como os mais usados para explicar as viagens geradas pelas zonas de tráfego. Considera-se determinadas variáveis, em geral, agregadas a nível de zona, como explicativas das ocorrências de viagens zonais e busca-se, através do método dos mínimos quadrados, uma equação de regressão linear múltipla que associe as viagens geradas às variáveis do problema.

Assim, o que se obtém é uma tentativa de descrição das viagens atuais por uma expressão do tipo $Y = a + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$ onde, Y é a variável dependente que representa o número de viagens atuais, a é a parcela de Y não explicada pela equação de regressão, b_i é o i -ésimo coeficiente de regressão parcial e X_i é a i -ésima variável explicativa ou independente, $i = 1, 2, \dots, p$.

Se $Y = a + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$ possui um coeficiente de determinação, R^2 , próximo da unidade, um pequeno erro padrão de estimativa, S_e , e os coeficientes de regressão parcial, b_i , estatisticamente significantes, a um nível de significância estabelecido, adota-se essa equação como equação de previsão de viagens e, dessa forma, as viagens no ano meta são calculadas, pela equação de regressão, simplesmente, atribuindo-se às variáveis independentes os seus respectivos valores futuros.

Assim, os coeficientes de regressão parcial são supostamente imutáveis, com o tempo, o que equivale a se supor não haver modificações no comportamento espacial da população, no ano meta, que possa alterar os coeficientes de regressão.

Análise de Regressão Linear Múltipla, também, tem aplicações em outras fases do planejamento dos transportes, como na divisão modal, na calibração do modelo gravitacional, etc.

2.2.2 - Enfoque Matemático da Análise de Regressão Linear Múltipla.

Análise de regressão linear múltipla é uma técnica estatística cujo principal objetivo é analisar a relação entre uma variável e um conjunto de variáveis. A primeira é chamada variável dependente e as últimas, independentes.

Uma equação de regressão linear múltipla é uma expressão da forma:

$$\hat{X}_{kn} = b_0 + b_1 X_{k1} + b_2 X_{k2} + \dots + b_{n-1} X_{k,n-1}$$

$$k = 1, 2, \dots, m$$

onde: m = número de observações de cada variável

n = número total de variáveis

\hat{X}_{kn} = valor estimado de cada observação da variável dependente X_n .

X_{ki} = k-ésima observação da variável independente X_i , $i = 1, 2, \dots, n-1$.

b_0 = constante não explicada pelas variáveis inde

pendentes X_i .

b_i = coeficiente de regressão parcial da variável independente X_i .

Utilizando-se o método dos mínimos quadrados, procura-se encontrar o conjunto de valores b_1, \dots, b_{n-1} tal que a soma dos quadrados dos desvios de X_n em relação ao valor estimado de X_n (\hat{X}_n) seja mínima. Em outras palavras, quer-se achar o $\text{Min } \sum (X_{kn} - \hat{X}_{kn})^2$.

$$\text{Seja } S = \sum_{k=1}^m (X_{kn} - \hat{X}_{kn})^2$$

O valor de S será mínimo quando as derivadas parciais de S em relação a b_1, \dots, b_{n-1} forem nulas. Essas expressões derivadas e arrumadas, convenientemente, fornecem o seguinte sistema de equações.

$$\left. \begin{aligned} b_1 SQ_{11} + b_2 SP_{12} + \dots + b_{n-1} SP_{1\ n-1} &= SP_{1\ n} \\ b_1 SP_{21} + b_2 SQ_{22} + \dots + b_{n-1} SP_{2\ n-1} &= SP_{2\ n} \\ \dots & \\ b_1 SP_{n-1\ 1} + b_2 SP_{n-1\ 2} + \dots + b_{n-1} SQ_{n-1\ n-1} &= SP_{n-1\ n} \end{aligned} \right\} \quad (1)$$

onde: $SQ_i = \sum_{k=1}^m (X_{ki} - \bar{X}_i)^2$ = soma dos quadrados dos

desvios de cada variável independente X_i em relação à sua média \bar{X}_i .

$SP_{ij} = \sum_{k=1}^m (X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j)$ = soma dos produtos cruzados entre as variáveis X_i e X_j .

A resolução do sistema de equações (1) conduz aos va-

lores de b_1, b_2, \dots, b_{n-1} .

O valor de b_0 é dado por:

$$b_0 = \bar{X}_n - b_1 \bar{X}_1 - \dots - b_{n-1} \bar{X}_{n-1}$$

$$\text{sendo } \bar{X}_i = \frac{\sum_{k=1}^m X_{ki}}{m}, \quad i = 1, 2, 3, \dots, n$$

O sistema (1) em termos matriciais, torna-se:

$$\begin{pmatrix} SQ_1 & SP_{12} & \dots & SP_{1 \ n-1} \\ SP_{21} & SQ_2 & \dots & SP_{2 \ n-1} \\ \dots & \dots & \dots & \dots \\ SP_{n-11} & SP_{n-12} & \dots & SQ_{n-1} \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{pmatrix} = \begin{pmatrix} SP_{1 \ n} \\ SP_{2 \ n} \\ \dots \\ SP_{n-1 \ n} \end{pmatrix}$$

$$SQ \cdot b = SP_n$$

$$b = SQ^{-1} \cdot SP_n$$

onde: b = vetor coluna $(n-1 \times 1)$ dos coeficiente parciais de regressão

SQ = matriz $(n-1 \times n-1)$ da soma dos quadrados e produtos cruzados

SP_n = vetor coluna $(n-1 \times 1)$ da soma dos produtos cruzados entre cada variável independente e a variável dependente.

Transformando-se as variáveis X_i nos seus valores padronizados, Z_i , $i=1, 2, \dots, n$, a equação de regressão linear múltipla torna-se:

$$Z_{kn} = B_1 Z_{k1} + B_2 Z_{k2} + \dots + B_{n-1} Z_{kn-1}$$

sendo:

$$Z_{ki} = \frac{X_{ki} - \bar{X}_i}{\sigma_{X_i}}, \quad i = 1, 2, \dots, n$$

σ_{X_i} = desvio padrão da variável X_i .

A função a ser minimizada passa a ser:

$$S_1 = \sum_{k=1}^m (Z_{kn} - \hat{Z}_{kn})^2$$

Derivando-se parcialmente S_1 em relação a B_1, B_2, \dots, B_{n-1} e igualando essas derivadas a zero, obtem-se um sistema de $(n-1)$ equações normais com $(n-1)$ incógnitas. Estas equações normais são da forma

$$\left\{ \begin{array}{l} B_1 + r_{12} B_2 + r_{13} B_3 + \dots + r_{1n-1} B_{n-1} = r_{1n} \\ r_{21} B_1 + B_2 + r_{23} B_3 + \dots + r_{2n-1} B_{n-1} = r_{2n} \\ \dots \\ r_{n-11} B_1 + r_{n-12} B_2 + r_{n-13} B_3 + \dots + B_{n-1} = r_{n-1n} \end{array} \right. \quad (2)$$

onde: r_{ij} = coeficiente de correlação entre as variáveis Z_i e Z_j

B_i = coeficiente de regressão padronizado da variável Z_i , $i = 1, 2, \dots, n-1$

Seja $R =$

$$\begin{array}{cccc|c} r_{11} & r_{12} & \dots & r_{1n-1} & r_{1n} \\ r_{21} & r_{22} & \dots & r_{2n-1} & r_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ r_{n-11} & r_{n-12} & \dots & r_{n-1n-1} & r_{n-1n} \\ \hline r_{n1} & r_{n2} & \dots & r_{nn-1} & r_{nn} \end{array}$$

dividida em blocos

$$\text{Sejam } R_{11} = \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1n-1} \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ r_{n-11} & r_{n-12} & \cdots & r_{n-1n-1} \end{pmatrix}, \quad R_{12} = \begin{pmatrix} r_{1n} \\ r_{2n} \\ \cdot \\ \cdot \\ r_{n-1n} \end{pmatrix}$$

$$R_{21} = (r_{n1}, r_{n2}, \dots, r_{nn-1}) \text{ e } R_{22} = r_{nn}. \text{ Ent\~{a}o}$$

$$R = \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \quad \text{e o sistema de equa\~{c}o\~{e}s} \quad (2)$$

pode ser escrito como

$$R_{11} \cdot B = R_{12}$$

$$B = R_{11}^{-1} \cdot R_{12}$$

onde: B = vetor coluna $(n-1 \times 1)$ dos coeficientes de regress\~{a}o padronizados

R_{11}^{-1} = inversa da matriz de correla\~{c}o\~{e} entre as vari\~{a}veis independentes.

R_{12} = vetor coluna $(n-1 \times 1)$ dos coeficientes de correla\~{c}o\~{e} entre cada vari\~{a}vel independente e a vari\~{a}vel dependente.

Pode se obter os coeficientes de regress\~{a}o n\~{a}o padronizados b_1, b_2, \dots, b_{n-1} a partir dos coeficientes de regress\~{a}o padronizados, calculando-se

$$b = AB$$

onde:

b = vetor coluna ($(n-1 \times 1)$) dos coeficientes de regressão não padronizados

A = matriz diagonal em que

$a_{ii} = \frac{\sigma_{X_n}}{\sigma_{X_i}}$, $i = 1, 2, \dots, n-1$, σ_{X_n} é o desvio padrão da variável dependente X_n e σ_{X_i} é o desvio padrão da variável independente X_i . As variáveis X_j , $j=1, \dots, n$ não estão na forma padronizada.

B = vetor coluna ($(n-1 \times 1)$) dos coeficientes de regressão padronizados.

O valor de b_0 é dado por:

$$b_0 = \bar{X}_n - (b^t \cdot M)$$

onde:

\bar{X}_n = média aritmética da variável dependente

b^t = vetor transposto de b

M = vetor coluna ($(n-1 \times 1)$) das médias das variáveis independentes.

2.2.3 - Interpretação dos coeficientes de regressão e da parte constante de uma equação de regressão linear múltipla

Os coeficientes de regressão parcial b_i , associados a cada variável independente X_i , $i = 1, 2, \dots, n-1$, representam a variação ocorrida no valor observado da variável dependente X_n quando a variável X_i varia de uma unidade, des

de que as demais variáveis independentes permaneçam constantes.

Para a obtenção dos coeficientes de regressão padronizados, transforma-se as variáveis X_j nas variáveis Z_j , $j=1, 2, \dots, n$, com média zero e desvio padrão igual a 1. Isso elimina o efeito das dimensões diferenciadas entre as variáveis X_j e propicia uma análise mais realista da contribuição de cada variável independente para o valor estimado da variável dependente.

Nesse caso, o coeficiente de regressão na forma padronizada B_i indica a variação sofrida por X_n , quando ocorre uma variação de um desvio padrão na variável independente X_i , mantendo-se constante as outras variáveis independentes. Isso possibilita a identificação daquelas variáveis independentes que mais contribuem para a variação da variável dependente.

A parte constante de uma equação de regressão linear múltipla reflete a incapacidade das variáveis independentes para explicar a variável dependente.

Para se ter idéia da influência da magnitude da constante na qualidade de uma equação de regressão linear, geralmente, compara-se o seu valor com a média da variável dependente. Em termos práticos, pode-se considerar que a constante tem um valor razoável quando a porcentagem que ela representa da média da variável dependente for inferior a 20%.

2.2.4 - Análise da qualidade de uma equação de regressão linear múltipla.

Obtida a equação de regressão linear múltipla, $\hat{X}_{kn} = b_0 + b_1 X_{k1} + \dots + b_{n-1} X_{kn-1}$, a etapa seguinte é analisar a sua validade, como modelo para previsão de viagens.

Uma equação de regressão linear múltipla pode ser avaliada através dos seguintes critérios:

- Análise numérica da equação;
- Análise qualitativa da equação;
- Testes estatísticos.

2.2.4.1 - Análise numérica da equação

A análise numérica de uma equação de regressão linear múltipla é feita, normalmente, com base na análise de dois parâmetros: o erro padrão de estimativa e o coeficiente de determinação.

O erro padrão de estimativa S_e e o coeficiente de determinação R^2 , podem ser calculados, respectivamente, por:

$$S_e = \sqrt{\frac{\sum_{i=1}^N (X_{in} - \hat{X}_{in})^2}{N - K - 1}}$$

$$R^2 = \frac{\sum_{i=1}^N (\tilde{X}_{in} - \bar{X}_n)^2}{\sum_{i=1}^N (X_{in} - \bar{X}_n)^2}$$

onde:

X_{in} = i-ésima observação da variável dependente X_n

\hat{X}_{in} = i-ésima estimativa da variável dependente X_n

\bar{X}_n = média aritmética da variável dependente X_n

N = número de observações

K = número de variáveis independentes

Em termos matriciais R^2 é dado por:

$$R^2 = B^t R_{12} \text{ sendo } B^t \text{ o vetor transposto de } B.$$

Quanto mais próximas forem as observações da variável dependente dos seus valores estimados, menores serão os desvios ou resíduos entre esses valores e, conseqüentemente, menor será o valor de S_e e mais próximo de 1 é o valor de R^2 .

O valor de R^2 , expresso em porcentagem, indica a quantidade da variância da variável dependente explicada pelas variáveis independentes.

Em uma equação de regressão, a soma dos quadrados dos desvios das observações da variável dependente X_n , em relação a sua média, \bar{X}_n , é igual a soma dos quadrados dos desvios dos valores estimados da variável dependente, \hat{X}_n , em relação a média \bar{X}_n , mais a soma dos quadrados dos desvios das observações da variável dependente, em relação aos valores estimados, \hat{X}_n .

Em termos matemáticos,

$$\sum_{i=1}^N (X_{in} - \bar{X}_n)^2 = \sum_{i=1}^N (\hat{X}_{in} - \bar{X}_n)^2 + \sum_{i=1}^N (X_{in} - \hat{X}_{in})^2$$

ou

$$SS_{X_n} = SS_{\hat{X}_n} + SS_{\text{resíduo}}$$

$$R^2 = \frac{SS_{\hat{X}_n}}{SS_{X_n}} = \frac{SS_{X_n} - SS_{\text{resíduo}}}{SS_{X_n}}, \text{ o que permite ver o in}$$

tervalo de variação de R^2 . Quando o resíduo é uma parcela

substancial de SS_{X_n} , $SS_{X_n} - SS_{\text{resíduo}}$ tende para zero e R^2 , também. Se $SS_{\text{resíduo}}$ é, aproximadamente, igual a zero, $SS_{X_n} - SS_{\text{resíduo}} = SS_{X_n}$ e R^2 é igual a 1. Portanto,

$$0 \leq R^2 \leq 1.$$

O coeficiente de correlação múltipla $R = \sqrt{R^2}$ pode ser visualizado como o coeficiente de correlação simples r entre X_n e \hat{X}_n , porque \hat{X}_n pode ser considerada como uma variável independente construída a partir da equação de regressão ².

Uma boa equação de regressão deve ter o menor valor possível para S_e e um valor de R^2 próximo da unidade.

2.2.4.2 - Análise qualitativa da equação.

A qualidade de uma equação de regressão linear está, intimamente, ligada à qualidade das variáveis independentes. Se as variáveis independentes possuem uma relação lógica de causa-efeito com a variável dependente, a equação obtida passa a ter um sentido lógico definido e, dessa forma, pode-se analisar os coeficientes de regressão quanto a sua coerência, em termos de magnitude e sinal.

É de se esperar que uma equação de regressão linear construída sem uma relação linear bem definida entre cada variável independente e a dependente não produza bons resultados.

Um outro aspecto importante, é que entre as variáveis independentes não deve haver relação linear. Caso contrário, haverá o aparecimento de multicolinearidade entre as

variáveis independentes, e isso pode implicar em distorção na magnitude e no sinal dos coeficientes de regressão, afetando a interpretação lógica da equação.

Portanto, a análise qualitativa de uma equação de regressão linear deve abranger os seguintes passos:

- Análise das relações lineares entre as variáveis independentes e entre cada variável independente e a dependente.
- Análise dos coeficientes de regressão na tentativa de detectar alguma incoerência, quanto aos seus valores e sinais.
- Análise, do ponto de vista lógico, da relação de causa e efeito entre cada variável independente e a variável dependente bem como análise lógica da equação como um todo.

2.2.4.3 - Testes estatísticos

Para a verificação das qualidades estatísticas de uma equação de regressão linear múltipla deve-se analisar do ponto de vista estatístico o coeficiente de correlação múltipla e os coeficientes de regressão parcial.

Para o teste do coeficiente de correlação múltipla, adota-se a hipótese nula de que o coeficiente de correlação é nulo e testa-se essa hipótese contra a hipótese alternativa de que ele não é nulo.

Computa-se, então,

$$F_{\text{calculado}} = \frac{R^2/K}{(1-R^2)/(N-K-1)} \quad \text{e compara-se esse valor}$$

com o $F_{\text{crítico}}$, de uma tabela de distribuição F, para K e (N-K-1) graus de liberdade, e um nível de significância α .

Se $F_{\text{calculado}} > F_{\text{crítico}}$, rejeita-se a hipótese nula e, em consequência, aceita-se a hipótese alternativa, ao nível de significância α . Esse teste para R é equivalente a se testar a nulidade de todos os coeficientes de regressão parcial. A rejeição da hipótese nula assegura, apenas, que um ou alguns coeficientes de regressão parcial, mas não necessariamente todos, são diferentes de zero. Esse fato mostra a necessidade de se testar individualmente cada coeficiente de regressão parcial. A hipótese nula, agora, é a de que cada coeficiente de regressão parcial é nulo e a alternativa é a de que cada um deles é diferente de zero.

Calcula-se F através da expressão

$$F_{\text{calculado}} = \frac{\text{incremento de } R^2 \text{ devido à entrada da variável ind. } X_i \text{ na eq./1}}{(1-R^2)/(N-K-1)}$$

Se $F_{\text{calculado}} > F_{\text{crítico}}$, para 1 e (N-K-1) graus de liberdade e um nível de significância α , rejeita-se a hipótese nula referente ao coeficiente de X_i . Nesse caso diz-se que esse coeficiente é estatisticamente significativo ao nível de significância α . Em outras palavras, isto significa que a probabilidade de o coeficiente analisado ser nulo é menor do que α .

Se $F_{\text{calculado}} < F_{\text{crítico}}$ aceita-se a hipótese nula e, portanto, aquele coeficiente de regressão parcial não é estatisticamente significativo ao nível de significância α .

A questão que se apresenta agora é a de considerar como válida ou não uma equação onde um ou vários coeficientes

não são estatisticamente significantes.

Um critério geral a adotar é o de rejeitar toda equação que tenha os coeficientes não significantes. Esse critério pode, no entanto, ocasionalmente, contradizer o critério da qualidade de uma equação baseado no valor de R^2 e no erro padrão de estimativa.

Pode ocorrer, por exemplo, uma situação em que uma equação de regressão linear múltipla possua um coeficiente de determinação, aproximadamente, igual a 1, um pequeno erro padrão de estimativa, e, no entanto, um dos coeficientes de regressão parcial, associado a uma variável com uma parcela de contribuição, relativamente, pequena para a variável dependente, não seja estatisticamente significativa, ao nível de significância adotado. Nesse caso, deveria se construir uma nova equação excluindo aquela variável não importante e refazer toda a análise sobre a qualidade dessa nova equação.

Em suma, o critério recomendado neste trabalho para avaliar uma equação de regressão linear múltipla como modelo de previsão é que a equação tenha um menor erro padrão de estimativa possível, um coeficiente de determinação próximo de 1 e todos os coeficientes de regressão sejam estatisticamente significantes ao nível de significância adotado.

2.3 - Análise Fatorial

2.3.1 - Considerações iniciais sobre a Análise Fatorial.

Segundo Cooley ⁴, Análise Fatorial vem se tornando um

termo genérico para uma variedade de procedimentos desenvolvidos com a finalidade de analisar as interrelações dentro de um conjunto de variáveis.

A característica mais importante da Análise Fatorial é sua capacidade de transformar o conjunto de dados em um outro conjunto menor de fatores que, em geral, são não correlacionados, entre si, e são extraídos do conjunto de variáveis originais de tal forma que o primeiro fator explique a maior quantidade de variância existente dentro do conjunto de dados, o segundo fator explique a maior quantidade da variância remanescente, isto é, variância não explicada pelo primeiro, o terceiro fator explique parte da variância não explicada pelos dois fatores anteriores, e, assim, sucessivamente.

O método de extração dos fatores, explicado resumidamente acima, é chamado método dos componentes principais e faz parte do programa de Análise Fatorial, da bateria de programas SPSS (Statistical Package for the Social Sciences), utilizado neste trabalho.

A Análise Fatorial tem um vasto campo de aplicações. Além do seu amplo emprego em áreas concernentes à Psicologia e Sociologia, essa técnica da Análise Multivariada, também, é bastante utilizada na Geografia e nos diversos tipos de planejamento, como planejamento do uso do solo, dos transportes, regional, etc.

2.3.2 - Enfoque Matemático da Análise Fatorial

Seja X uma matriz do tipo $(N \times m)$ onde cada coluna cor

N-dimensional, na forma padronizada, no qual se quer projetar X.

Como em geral, X é de posto m, pode-se encontrar um vetor g ($m \times 1$) tal que

$$Xg = y \quad (1)$$

Além disso,

$$\frac{y'y}{N} = g' \frac{X'X}{N} g = g'Rg \quad (2)$$

onde:

$\frac{X'X}{N} = R =$ matriz de correlação entre os m vetores coluna de X do tipo ($m \times m$) e y' , g' e X' são os transpostos de y , g e X , respectivamente.

Por outro lado, $\frac{y'y}{N}$ é a variância de y e, portanto, pelo fato de y ser um vetor na forma padronizada

$$\frac{y'y}{N} = 1 \quad (3)$$

Comparando-se as expressões (2) e (3), obtém-se:

$$g'Rg = 1 \quad (4)$$

A projeção de X sobre y é dada pelo produto escalar de X' pelo vetor unitário da direção y , $\frac{y}{\sqrt{N}}$. Entretanto, trabalhar-se com os vetores coluna da matriz $\frac{X}{\sqrt{N}}$, matriz em que seus vetores coluna são de comprimento unitário, equivale, simplesmente, a transformar a escala do sistema de referência, e as projeções dos vetores de $\frac{X}{\sqrt{N}}$ sobre o vetor y isto é, $\frac{X'}{\sqrt{N}} \cdot \frac{y}{\sqrt{N}}$, fornecem as projeções da variância unitária

das variáveis originais. Isso vai ao encontro da estratégia da Análise Fatorial de buscar os vetores y ortogonais, entre si, sobre cada um dos quais a projeção da variância é máxima. Como os vetores $\frac{X^j}{\sqrt{N}}$ e $\frac{Y}{\sqrt{N}}$ são unitários a projeção se torna igual ao cosseno do ângulo formado por cada vetor coluna de X e o vetor y , e esse cosseno é o coeficiente de correlação linear entre os vetores considerados.

Seja f a projeção procurada.

$$\rightarrow f = \frac{X'Y}{N}$$

De (1), $y = Xg$. Então:

$$f = \frac{X'Xg}{N} = Rg \quad (5)$$

O objetivo, agora, é maximizar a soma dos quadrados das projeções f . Em termos vetoriais, a soma dos quadrados dessas projeções é dada por $f'f$, onde f' é o vetor transposto de f .

$$f'f = g'R'Rg = g'R^2g \quad (6)$$

Quer-se, portanto, maximizar $f'f$ sujeita à condição de $g'Rg = 1$.

Usando o método dos multiplicadores indeterminados de Lagrange, obtém-se:

$$Z = g'R^2g - \lambda (g'Rg - 1)$$

Derivando-se Z em relação a g' e igualando-se o resultado a zero:

$$\frac{\partial Z}{\partial g'} = 2R^2g - 2\lambda Rg = 0 \quad \therefore R^2g = \lambda Rg \quad (7)$$

Ainda,

$$R^2 g = R \cdot Rg = \lambda Rg$$

De (5), $Rg = f$. Tem-se, pois:

$$Rf = \lambda f \quad Rf - \lambda f = 0 \quad (8)$$

A expressão (8) mostra f como um autovetor correspondente ao autovalor λ , de R . Se R possui m autovalores reais não nulos e distintos, a matriz F , formada pelos autovetores f_i , dispostos em colunas, correspondentes aos autovalores λ_i de R , é de posto m . Generalizando a expressão (8), tem-se:

$$RF = F\Lambda \quad (9)$$

onde:

$F = (f_{ij})$ = matriz dos carregamentos dos fatores

Λ = matriz diagonal formada pelos autovalores de R , em ordem decrescente, isto é, $\lambda_{11} > \lambda_{22} > \dots > \lambda_{mm}$.

Multiplicando-se ambos os membros de $R^2 g = \lambda Rg$, à esquerda por g' , vem:

$$g'R^2 g = g'\lambda Rg = \lambda g'Rg$$

De (4) e (6), tem-se, respectivamente,

$$g'Rg = 1 \text{ e } g'R^2 g = f'f. \text{ Então,}$$

$$f'f = \lambda \quad (10)$$

Portanto, para se obter a maior soma dos quadrados das projeções f , basta tomar o maior autovalor λ , de R .

Generalizando a expressão (10), vem:

$$F'F = \Lambda \quad (11)$$

Substituindo a expressão (11) em (9), resulta:

$$RF = FF'F$$

Como a matriz F é inversível,

$$RFF^{-1} = FF'FF^{-1} \quad \dots$$

$$R = FF' \quad (12)$$

A expressão (7) pode, também, ser escrita como:

$$R^2g - \lambda Rg = 0$$

$$R.Rg - R\lambda g = 0$$

Como $R = \frac{X'X}{N}$ tem posto m , R^{-1} existe. Assim:

$$R^{-1}R.Rg - R^{-1}R\lambda g = 0$$

$$Rg = \lambda g$$

De (5), $Rg = f$. Então:

$$f = \lambda g \quad \dots$$

$$g = f\lambda^{-1}$$

Generalizando o vetor g para uma matriz G , tem-se:

$$G = F\Lambda^{-1} \quad (13)$$

A matriz Y , formada pelos vetores y é dada por:

$$Y = XG = XF\Lambda^{-1} \quad (14)$$

$$\begin{aligned} \frac{Y'Y}{N} &= \frac{G'X'XG}{N} = G'RG = \Lambda^{-1} F' (RF) \Lambda^{-1} = \\ &= \Lambda^{-1} F' F \Lambda^{-1} = \Lambda^{-1} \Lambda \Lambda^{-1} = I \cdot I = I \end{aligned}$$

Portanto, Y é ortogonal

De $R = FF'$,

$$R = f_1 f_1' + f_2 f_2' + \dots + f_m f_m'$$

A parcela $f_j f_j'$ representa a contribuição do fator j para a matriz de correlação.

De (14) pode-se escrever:

$$X F = Y \Lambda$$

$$X F - Y \Lambda = 0$$

De (11), $\Lambda = F'F$. Daí,

$$X F - Y F'F = 0$$

Como a matriz F possui inversa:

$$X F F^{-1} - Y F' F F^{-1} = 0$$

$$X = Y F'$$

Assim, para qualquer vetor coluna da matriz X tem-se:

$$X_i = Y F_i' \quad (15)$$

onde F_i' é o transposto do i -ésimo vetor linha F_i da matriz F .

Multiplicando ambos os membros de (15) por $\frac{X_i'}{N}$, vem:

$$\frac{X_i' X_i}{N} = \frac{F_i Y_i' Y_i F_i'}{N} \quad (16)$$

Como X possui os seus vetores coluna na forma padronizada, $\frac{X_i' X_i}{N}$ representa a variância do vetor coluna X_i e, portanto, é igual a 1. Ainda, conforme mostrado anteriormente, $\frac{Y'Y}{N}$ é igual a matriz iden-

tidade I. Dessa forma:

$$1 = F_i F_i' = f_{i1}^2 + \dots + f_{im}^2$$

Portanto, os termos f_{ik}^2 , $k=1,2,\dots,m$, são os componentes da variância unitária de X_i .

Se a matriz F final é de dimensão m,

$$\sum_{k=1}^m f_{ik}^2 = 1$$

Se F final é de dimensão $p < m$,

$$\sum_{k=1}^p f_{ik}^2 < 1, \text{ e esse somatório é chamado } \underline{\text{comunalidade}}$$

da variável X_i .

$$f'f = \lambda \text{ equivale a } \sum_{i=1}^N f_{ij}^2 = \lambda_j \text{ e representa a con-}$$

tribuição do fator j para a variância total das variáveis.

Como a variância total é igual ao número de variáveis m, $\frac{\lambda_j}{m} \times 100$ fornece a porcentagem da variância total dos dados explicada pelo fator j.

2.3.3 - Rotação dos fatores

Partindo-se da matriz dos fatores iniciais quer-se fazer uma mudança nos carregamentos das variáveis de modo que os novos valores dos carregamentos sejam as coordenadas das variáveis em um novo sistema de eixos ortogonais, obtidos de uma rotação dos eixos antigos.

Essa rotação possibilita uma interpretação lógica dos fatores. Em geral, é difícil a interpretação dos fatores como eles se apresentam na matriz dos fatores iniciais.

Em termos matriciais, se $F = (f_{ij})$ é a matriz dos fa-

tores iniciais $F^r = (f_{ij}^r)$, matriz dos fatores após a rotação, é definida por

$$F^r = F \cdot T$$

onde:

T = matriz de rotação final

A matriz T é o produto de sucessivas rotações que podem ser organizadas em uma ordem conveniente de forma que cada eixo seja rotacionado com cada um dos outros eixos apenas uma vez.

O processo de rotação apresentado pode ser generalizado para um espaço de dimensão p . Nesse caso, a matriz de rotação T torna-se:

$$T = \begin{pmatrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1p} \\ \alpha_{21} & \alpha_{22} & & \alpha_{2p} \\ \dots & \dots & \dots & \dots \\ \alpha_{p1} & \alpha_{p2} & & \alpha_{pp} \end{pmatrix}$$

As colunas da matriz acima são os cossenos diretores dos eixos $F_1^r, F_2^r, \dots, F_p^r$ em relação a $F_1, F_2, F_3, \dots, F_p$.

Como a rotação é ortogonal, os cossenos diretores devem satisfazer a condição

$$\sum_{r=1}^p \alpha_{rk} \alpha_{rl} = \delta_{kl} \quad (k, l = 1, 2, \dots, p; k < l)$$

onde:

δ_{kl} é o delta de Kronecker, isto é,

$$\delta_{kl} = 0, \text{ se } k \neq l$$

$$\delta_{kl} = 1, \text{ se } k = l$$

Nesse trabalho é usado o método de rotação ortogonal VARIMAX. Esse método tem como principal finalidade rotacionar os eixos de modo a maximizar a variância dos quadrados dos carregamentos de cada coluna da matriz F.

2.4 - Análise Discriminante

2.4.1 - Considerações iniciais sobre a Análise Discriminante.

As N observações de um conjunto de m variáveis formam uma matriz $X(N \times m)$.

Se essas N linhas podem ser arrumadas de forma a constituírem grupos de observações segundo um critério estabelecido, isto é, se há uma classificação prévia dos elementos de X, essa matriz fica dividida em matrizes de dimensões menores formada pelos grupos de observações, onde a soma do número de linhas de todas essas matrizes é igual a N.

Feita essa classificação inicial, tem-se a matriz X preparada para a aplicação de uma Análise Discriminante.

Análise Discriminante é uma técnica estatística multivariada que tem como objetivo principal construir combinações lineares das variáveis originais que maximizem a separação entre os grupos. Essa técnica, também, pode ser usada para testar se os grupos são, realmente, distintos, determinar o poder discriminatório das variáveis iniciais e identificar aquelas que mais contribuem para a separação dos grupos, avaliar a dispersão interna dos grupos, identificar elementos indevidamente classificados, alocar elementos a um dos grupos existentes e melhorar a qualidade da classificação ini

cial por meio de iterações que tem como finalidade minimizar a variância interna dos grupos e maximizar a variância entre os grupos.

Nesse trabalho será usado o programa computacional de Análise Discriminante, da bateria de programas SPSS (Statistical Package for the Social Sciences).

2.4.2 - Enfoque matemático da Análise Discriminante.

Seja X uma matriz do tipo $(N \times m)$, formada por grupos de observações, em que cada vetor coluna representa uma variável específica com seus valores dados em termos de desvios da sua média.

Então, cada observação de X pode ser escrita como:

$$x_{ijk} = m_{ij} + C_{ijk} \quad (17)$$

onde:

x_{ijk} = k -ésima observação da variável X_i no grupo j .

m_{ij} = média de X_i no grupo j .

C_{ijk} = k -ésimo componente do erro da variável X_i no grupo j .

Generalizando a expressão (17), vem:

$$X = M + E, \quad (18)$$

onde as linhas da matriz M , em cada grupo, são iguais.

$$\text{Seja } T = X'X \quad (19)$$

T é chamada matriz da soma dos quadrados e produtos cruzados.

$$T = X'X = (M+E)'(M+E) = M'M + E'E + M'E + E'M \quad (20)$$

Adotando a hipótese de não haver correlação entre as matrizes dos erros e as das médias, $M'E = E'M = 0$ e \forall

$$T = M'M + E'E = B + W \quad (21)$$

As matrizes B, W e T podem, também, ser definidas por:

$$B = (b_{ij}), \quad b_{ij} = \sum_{k=1}^G Ng (\bar{X}_{ik} - \bar{X}_i) (\bar{X}_{jk} - \bar{X}_j)$$

$$W = (w_{ij}), \quad w_{ij} = \sum_{k=1}^G \sum_{n=1}^{Ng} \left[(X_{ikn} - \bar{X}_{ik}) (X_{jkn} - \bar{X}_{jk}) \right]$$

$$T = (t_{ij}), \quad t_{ij} = \sum_{p=1}^N (X_{ip} - \bar{X}_i) (X_{jp} - \bar{X}_j)$$

onde:

N = número total de observações

G = número de grupos

Ng = número de elementos no grupo g *(Linhas)*

m = número de variáveis originais

B = matriz de dispersão entre grupos

W = matriz de dispersão interna dos grupos

T = B + W = matriz de dispersão total,

\bar{X}_{ik} = média aritmética da variável X_i no grupo k

$i = 1, 2, \dots, m$ e $k = 1, 2, \dots, G$

\bar{X}_i = média aritmética da variável X_i

X_{ikn} = n-ésima observação da variável X_i no grupo k

$n = 1, 2, \dots, Ng$

X_{ip} = p-ésima observação da variável X_i , $p = 1, 2, \dots, N$

Com a finalidade de maximizar a razão da dispersão em

tre os grupos para a dispersão interna dos grupos, pode-se projetar a matriz X em um vetor k , cujas componentes são os seus cossenos diretores.

Seja $y = XK$ a projeção de X sobre K .

$$y'y = K'X'XK = K'TK = K'(B+W)K = K'BK + K'WK \quad (22)$$

onde:

$y'y$ = soma dos quadrados das projeções de X sobre K .

$K'BK$ = parcela de $y'y$ correspondente à dispersão entre os grupos. \gg

$K'WK$ = parcela de $y'y$ correspondente à dispersão interna dos grupos. \ll

Quer-se, portanto, maximizar a razão $\frac{K'BK}{K'WK}$.

$$\text{Seja } \lambda = \frac{K'BK}{K'WK} \quad (23)$$

$$\log \lambda = \log(K'BK) - \log(K'WK)$$

Derivando-se $\log \lambda$ em relação a K e igulando-se o resultado a zero, tem-se:

$$\frac{\partial \log \lambda}{\partial K} = \frac{2BK}{K'BK} - \frac{2WK}{K'WK} = 0 \quad (24)$$

Multiplicando-se os termos da expressão (24) por $K'BK$,

vem:

$$2BK - 2WK (K'BK)/K'WK = 0$$

De (23), $(K'BK)/K'WK = \lambda$. Então:

$$BK - \lambda WK = 0 \quad (25)$$

Multiplicando-se a expressão (25) à esquerda por W^{-1} ,

inversa de W , resulta.

$$W^{-1}BK - \lambda K = 0 \quad \text{ou}$$

$$(W^{-1}B - \lambda I)K = 0 \quad (26)$$

Dessa forma, K é um autovetor associado ao autovalor λ de $W^{-1}B$.

Para se obter K deve-se encontrar as soluções não nulas do sistema de equações homogêneo (26).

A solução não trivial do sistema de equações homogêneo $(W^{-1}B - \lambda I)K = 0$ requer:

$$\left| W^{-1}B - \lambda I \right| = 0$$

onde:

$$\left| W^{-1}B - \lambda I \right| = \text{determinante da matriz } (W^{-1}B - \lambda I)$$

Para cada autovalor de $W^{-1}B$ obtém-se um autovetor correspondente, K , em que suas componentes são os coeficientes de uma combinação linear das variáveis originais e essas combinações lineares são chamadas funções discriminantes.

Como o número de autovalores reais é igual ao mínimo de $\left[(Ng - 1), m \right]$, isso implica em que se tenha esse mesmo número de autovetores de $W^{-1}B$ e, conseqüentemente, esse mesmo número de funções discriminantes.

A forma geral de uma função discriminante é:

$$D_i = d_{i1} z_1 + d_{i2} z_2 + \dots + d_{ip} z_p, \quad p \leq m$$

$$D_i = \text{i-ésima função discriminante } i=1, \dots, \min \left[(Ng - 1), m \right]$$

Z_j = variável discriminante j , na forma padronizada,

$$j = 1, \dots, p$$

d_{ij} = coeficiente da variável discriminatória j , na i -ésima função discriminante.

m = número de variáveis originais

p = número de variáveis discriminatórias.

As funções discriminantes são obtidas em ordem decrescente de importância. A primeira função explica o máximo da variância entre os grupos; a segunda função, ortogonal à primeira, o máximo da variância remanescente, e assim, sucessivamente, até se esgotar a variância existente na matriz $W^{-1}B$.

Um dos critérios, normalmente usados para se avaliar, qualitativamente uma função discriminante é a relação entre o seu autovalor e o somatório dos autovalores para todas as funções, expresso em porcentagem, $\frac{\lambda_j}{\sum_{i=1}^j \lambda_i} \times 100$,

onde:

$$S = \min \left[(Ng - 1), m \right]$$

λ_i = i -ésimo autovalor de $W^{-1}B$

O poder discriminatório do conjunto das variáveis é geralmente, verificado através do lambda de Wilks.

$$\Lambda = \frac{|W|}{|T|}, \quad 0 \leq \Lambda \leq 1$$

Quanto maior o poder discriminatório das variáveis menor o valor de Λ .

Se, no entanto, se pretende analisar o poder discriminatório de uma variável, tomada isoladamente, pode-se exami

nar o seu coeficiente na função discriminante em que a variável aparece, pela primeira vez, e baseando-se na importância da função, tem-se uma idéia do poder discriminatório dessa variável.

Uma outra alternativa para se analisar o poder discriminatório de cada variável é realizar o teste F de significância para cada variável.

2.4.3 - Interpretação Bayesiana da Análise Discriminante.

O objetivo principal da interpretação Bayesiana da Análise Discriminante é, conhecida uma observação x , alocar essa observação para um dos grupos anteriormente definidos. Essa alocação é feita para o grupo ao qual x tem a maior probabilidade de pertencer.

Sejam H_1, H_2, \dots, H_r , r hipóteses mutuamente exclusivas e $p(H_1), p(H_2), \dots, p(H_r)$ as probabilidades associadas a essas hipóteses.

$p(H_g)$, $g = 1, 2, \dots, r$, pode ser interpretada como a probabilidade de ser correta a sentença: H_g é verdadeira para a observação x .

$$p(H_g/x) = \frac{p(H_g) \cdot p(x/H_g)}{p(x)} \quad (27)$$

onde:

$p(H_g/x)$ = probabilidade de H_g ser verdadeira se x é observado

$p(x/H_g)$ = probabilidade condicional de x ser observado se H_g é verdadeira.

$p(x)$ = probabilidade associada à observação x .

Então, se o conjunto das hipóteses contém todas as alternativas possíveis, uma delas deve ser verdadeira. Tem-se, pois:

$$\sum_i p(H_i/x) = 1 \quad (28)$$

Portanto, o somatório de (27) para todo i variando de 1 a r , torna-se:

$$\sum_i p(H_i/x) = \frac{\sum_i [p(H_i) \cdot p(x/H_i)]}{p(x)}$$

De (28), $\sum_i p(H_i/x) = 1$. Resulta, então:

$$p(x) = \sum_i [p(H_i) \cdot p(x/H_i)] \quad (29)$$

Substituindo-se (29) em (27), vem:

$$p(H_g/x) = \frac{p(H_g) \cdot p(x/H_g)}{\sum_i [p(H_i) \cdot p(x/H_i)]} \quad (30)$$

Similarmente, para uma outra hipótese H_f , tem-se:

$$p(H_f/x) = \frac{p(H_f) \cdot p(x/H_f)}{\sum_i [p(H_i) \cdot p(x/H_i)]} \quad (31)$$

Dividindo-se (30) por (31), chega-se a:

$$\frac{p(H_g/x)}{p(H_f/x)} = \frac{p(H_g)}{p(H_f)} \cdot \frac{p(x/H_g)}{p(x/H_f)} = \frac{p(H_g)}{p(H_f)} \cdot L \quad (32)$$

onde:

$$\frac{p(H_g/x)}{p(H_f/x)} = \text{razão das probabilidades a posterior}$$

$$\frac{p(H_g)}{p(H_f)} = \text{razão das probabilidades a priori}$$

$$\frac{p(x/H_g)}{p(x/H_f)} = L = \text{razão de verossimilhança}$$

Uma estimativa da razão das probabilidades a priori, $p(H_g)/p(H_f)$, pode ser obtida de N_1/N_2 .

onde:

N_1 = número de observações do grupo 1

N_2 = número de observações do grupo 2

A expressão (32) mostra L como um fator de correção da razão das probabilidades a priori para se obter a razão das probabilidades a posteriori. Isso permite, dada uma observação, inferir que hipótese é a mais provável.

2.4.3.1 - Análise discriminante para dois grupos. Uso da razão de verossimilhança.

Sejam duas hipóteses H_1 e H_2 para as quais se conheça as distribuições de probabilidade de uma observação x , $f_1(x)$ e $f_2(x)$.

Assim,

$$\frac{f_1(x)}{f_2(x)} = L \quad (33)$$

Em geral, torna-se mais simples, em vez de se considerar L, considerar-se o logaritmo de L, isto é, $\log L$, principalmente se $f_1(x)$ e $f_2(x)$ são distribuições normais ou multinormais.

No caso multinormal,

$$\log \left[\frac{f_1(x)}{f_2(x)} \right] = -\frac{1}{2} m \log(2\pi) - \frac{1}{2} \log |V_1| - \frac{1}{2} (x - \bar{x}_1)' V_1^{-1} (x - \bar{x}_1)$$

$$\log [f_2(x)] = -\frac{1}{2} m \log (2^m) - \frac{1}{2} \log |V_2| - \frac{1}{2} (x - \bar{x}_2)' V_2^{-1} (x - \bar{x}_2)$$

onde:

$|V_1|$ e $|V_2|$ = determinantes das matrizes de variância-covariância interna dos grupos 1 e 2, respectivamente.

x = vetor cujas componentes são os valores das variáveis correspondentes a uma observação considerada.

\bar{x}_1 e \bar{x}_2 = vetores cujas componentes são os valores médios das observações das variáveis para os grupos 1 e 2, respectivamente.

$$\log \left[\frac{f_1(x)}{f_2(x)} \right] = \log L = -\frac{1}{2} \log |V_1| + \frac{1}{2} \log |V_2| - \frac{1}{2} (x - \bar{x}_1)' V_1^{-1} (x - \bar{x}_1) + \frac{1}{2} (x - \bar{x}_2)' V_2^{-1} (x - \bar{x}_2) \quad (34)$$

Se $V_1 = V_2 = V$, (34) torna-se:

$$\begin{aligned} \log L &= -\frac{1}{2} (x - \bar{x}_1)' V^{-1} (x - \bar{x}_1) + \frac{1}{2} (x - \bar{x}_2)' V^{-1} (x - \bar{x}_2) \\ \log L &= x' V^{-1} (\bar{x}_1 - \bar{x}_2) - \frac{1}{2} (\bar{x}_1' V^{-1} \bar{x}_1 - \bar{x}_2' V^{-1} \bar{x}_2) \quad (35) \end{aligned}$$

A matriz V de variância-covariância interna dos grupos pode ser estimada através de ¹⁰:

$$V = \frac{1}{N_1 + N_2 - 2} \left[\sum_i (x_i - \bar{x}_1)(x_i - \bar{x}_1)' + \sum_i (x_i - \bar{x}_2)(x_i - \bar{x}_2)' \right]$$

Quando $\log L = 0$, $L = 1$ e com base em (33), $f_1(x) = f_2(x)$. Tem-se, portanto, da expressão (32),

$$\frac{p(H_g/x)}{p(H_f/x)} = \frac{p(H_g)}{p(H_f)} \quad (36)$$

Nesse caso, (36) indica que a razão das probabilida-

des a priori é igual a razão das probabilidades a posteriori, implicando em que não se pode estabelecer se a observação x provém do primeiro grupo ou do segundo.

A curva associada a $\log L$ é chamada curva de verossimilhança. Se $\log L = 0$, ou seja, $L = 1$, obtém-se a curva discriminante, isto é, aquela que representa o lugar geométrico das observações com a mesma distribuição de probabilidade para as hipóteses H_1 e H_2 .

2.4.3.2 - Análise Discriminante para mais de dois grupos. Uso da razão de verossimilhança.

A abordagem da Análise Discriminante para mais de dois grupos é, apenas, uma generalização do caso de dois grupos. Para três grupos, por exemplo, adota-se as hipóteses H_1 , H_2 e H_3 e calcula-se as curvas de verossimilhança para os grupos, tomados dois a dois, como segue:

$$\begin{aligned} \log L_{12} &= \log \left[\frac{f_1(x)}{f_2(x)} \right] = \log [f_1(x)] - \log [f_2(x)] \\ \log L_{23} &= \log \left[\frac{f_2(x)}{f_3(x)} \right] = \log [f_2(x)] - \log [f_3(x)] \\ \log L_{31} &= \log \left[\frac{f_3(x)}{f_1(x)} \right] = \log [f_3(x)] - \log [f_1(x)] \end{aligned} \quad (37)$$

Quando $\log L_{12} = \log L_{31} = 0$, o espaço das observações fica dividido em regiões delimitadas pelos hiperplanos $\log L_{23} = 0$ e $\log L_{31} = 0$, e pode-se verificar que a terceira equação do sistema (37) é o resultado da soma das duas primeiras. Assim, despreza-se a última equação e as probabilidades a posteriori, dada uma observação de x , para as hipó-

teses H_1 , H_2 e H_3 podem ser dadas por:

$$\frac{p(H_1/x)}{p(H_2/x)} = \frac{p(H_1)}{p(H_2)} \cdot L_{12}$$

$$\frac{p(H_2/x)}{p(H_3/x)} = \frac{p(H_2)}{p(H_3)} \cdot L_{23}$$

Nesse caso, a observação x será alocada à hipótese que tenha a maior probabilidade de ocorrer.

CAPÍTULO III

APLICAÇÕES E RESULTADOS

3.1 - Introdução

Neste capítulo são apresentadas aplicações dos métodos estatísticos multivariados: Análise de Regressão Linear Múltipla, Análise Fatorial e Análise Discriminante.

Os dados utilizados nas aplicações para Campina Grande provêm das seguintes fontes: de uma pesquisa domiciliar, realizada em 1978, pelo GEIPOT e de um trabalho realizado pelo professor Masayuki Doi¹¹, sob o título "The land use and zoning study in relation to systematic transportation study".

Os resultados das aplicações foram obtidos através do uso do pacote computacional SPSS² no computador IBM-370, do Núcleo de Processamento de Dados da Universidade Federal da Paraíba, Campus de Campina Grande.

3.2 - Análise de Regressão Linear Múltipla. Aplicações e resultados.

3.2.1 - Aplicação 1

Esta aplicação tem por finalidade verificar se o nível de agregação das variáveis envolvidas numa equação de regressão linear tem influência na qualidade dessa equação.

Mc Carthy ⁷ construiu equações de regressão linear com as variáveis agregadas a nível de distrito, setor, zona e domicílio e verificou que essas equações tornavam-se melhores à medida que o nível de agregação diminuía.

Para esta aplicação, utilizou-se os dados da pesquisa domiciliar em Campina Grande realizada pelo Geipot, em 1978, e considerou-se essa cidade dividida a nível de setor, de zona e de domicílio.

Os setores foram definidos como o agrupamento de determinadas zonas e a Tabela 1 mostra as zonas que compõem cada setor. O zoneamento adotado foi o utilizado pelo Geipot na pesquisa domiciliar em 1978, em que a cidade foi dividida em 23 zonas de tráfego.

Tabela 1 - Setores e suas zonas

Setor	Zonas
1	1, 2, 3, 4
2	5, 6, 7
3	8, 9, 10, 11, 12
4	13, 14, 15
5	16, 17, 18, 19, 20
6	21, 22, 23

Para cada nível computou-se as seguintes variáveis:

X_1 = viagens diárias com base domiciliar por domicílio.

X_2 = renda por domicílio

X_3 = tamanho da família por domicílio

X_4 = número de veículos por domicílio

X_5 = pessoas maiores ou iguais a 5 anos por domicílio.

Essas variáveis foram as mesmas empregadas por Mc Carthy⁷.

Com base no raciocínio de se considerar as variáveis em diversos níveis de agregação, empregado por Mc Carthy⁷, e nos coeficientes de correlação linear entre as variáveis, mostrados nas Tabelas 2, 3 e 4, construiu-se as equações de regressão linear de X_1 em X_4 e X_2 e de X_1 em X_4 , tanto para os níveis de setor como para o de zona, e as equações de X_1 em X_5 e X_3 e de X_1 em X_5 , para o nível de domicílio.

Tabela 2 - Matriz de correlação entre as variáveis a nível de setor

	X_1	X_2	X_3	X_4	X_5
X_1	1,0000	0,2369	-0,5536	0,5666	-0,0861
X_2		1,0000	0,1375	-0,0956	0,1252
X_3			1,0000	-0,6325	0,8503
X_4				1,0000	-0,3722
X_5					1,0000

Tabela 3 - Matriz de correlação entre as variáveis a nível de zona

	X_1	X_2	X_3	X_4	X_5
X_1	1,0000	0,5464	-0,5117	0,6911	-0,5393
X_2		1,0000	-0,3778	0,7822	-0,2409
X_3			1,0000	-0,5861	0,9480
X_4				1,0000	-0,4991
X_5					1,0000

Tabela 4 - Matriz de correlação entre as variáveis a nível de domicílio

	X ₁	X ₂	X ₃	X ₄	X ₅
X ₁	1,0000	0,2304	0,5667	0,2561	0,6140
X ₂		1,0000	0,0730	0,3820	0,0917
X ₃			1,0000	0,1021	0,9535
X ₄				1,0000	0,1150
X ₅					1,0000

As equações com os seus respectivos parâmetros são:

3.2.1.1 - Nível de setor

Número de observações = 6

$$X_1 = 4,5046 X_4 + \frac{0,1709}{10^4} X_2 + 5,7191 \quad (38)$$

$$B_{X_4} = 0,5947; \quad B_{X_2} = 0,2937$$

$$R^2 = 0,4066; \quad Se = 1,0225; \quad F = 1,0277; \quad F_3^2 (5\%) = 9,55$$

$$F_{X_4} = 1,7772; \quad F_{X_2} = 0,432; \quad F_3^1 (5\%) = 10,13$$

onde:

B_{X_4} e B_{X_2} = coeficientes de regressão padronizados de X_4 e X_2 , respectivamente

R^2 = coeficiente de determinação múltipla

Se = erro padrão de estimativa

F = F calculado para a equação

F_{X_4} e F_{X_2} = F calculado para, respectivamente, os coeficientes de X_1 e X_2

$F_{GLD}^{GLN(u)}$ = $F_{crítico}$, isto é, tirado de uma tabela de distribuição F, para GLN graus de liberdade do numerador, GLD graus de liberdade do denominador e um nível de significância α .

Pelo critério recomendado neste trabalho de se adotar como equação de previsão aquelas equações possuidoras de um coeficiente de determinação R^2 próximo de 1, um erro padrão de estimativa o menor possível, e os coeficientes de regressão parcial estatisticamente significantes a um nível de significância, previamente, adotado, a equação (38) não poderia ser usada como equação de previsão de demanda de viagens, pois, não apresenta o coeficiente de correlação múltipla R e os coeficientes de regressão parcial, estatisticamente significantes ao nível de significância de 5% porque $F < F_3^2(5\%)$ e F_{X_4} e F_{X_2} ambos são menores do que o $F_3^1(5\%)$.

A outra equação de regressão obtida, ainda, para o nível de setor foi:

$$X_1 = 4,2919 X_4 + 6,0381 \quad (39)$$

$$B_{X_4} = 0,5666$$

$$R^2 = 0,3211; S_e = 0,9472; F = 1,8918; F_4^1(5\%) = 7,71$$

Como a equação (39) é uma equação de regressão linear simples, o $F_{calculado}$ para o coeficiente de regressão é igual ao $F_{calculado}$ para a equação e, portanto, desnecessário escrevê-lo.

Nota-se que a equação (39) também não possui o seu coeficiente de regressão estatisticamente significativo a 5% de

nível de significância, além de um pequeno valor para R^2 .

Portanto, essa equação, também, não satisfaz às condições exigidas, nesse trabalho, para uma equação de previsão.

3.2.1.2 - Nível de zona

Número de observações = 19

Apesar de Campina Grande estar dividida em 23 zonas de tráfego, as zonas 2, 4, 6 e 17 foram eliminadas porque não possuem domicílios. Portanto, restaram 19 zonas que correspondem ao número de observações utilizadas nos cálculos das equações de regressão linear, a esse nível de agregação das variáveis.

O processo de obtenção das equações foi o "Stepwise" onde a entrada das variáveis independentes na equação se dá em ordem decrescente de sua capacidade de explicar a variância da variável dependente. Pode ocorrer, pois, que uma ou mais variáveis independentes fiquem fora da equação de regressão final, em virtude de essas variáveis não darem contribuição adicional para a explicação da variância da variável dependente.

Para o nível de zona, na equação de regressão linear de X_1 em X_4 e X_2 , a variável X_2 foi eliminada; com isso se obteve, somente, a equação de regressão linear, dada por:

$$X_1 = 9,1482 X_4 + 3,8220 \quad (40)$$

$$B_{X_4} = 0,6911$$

$$R^2 = 0,4776; \text{ Se} = 1,8031; F = 15,5449; F_{17}^1 (5\%) = 4,45$$

A equação (40) possui o valor de R^2 pequeno e, embora

o coeficiente de regressão seja estatisticamente significativo a 5% de nível de significância, o valor do coeficiente de regressão parece excessivo, visto que ele indica um acréscimo de 9 viagens, aproximadamente, para um acréscimo de um automóvel por domicílio. Essa equação, também, não deveria ser usada como modelo de previsão de demanda de viagens.

3.2.1.3 - Nível de domicílio

Número de observações = 1940

As equações obtidas a nível de domicílio foram:

$$X_1 = 2,1455 X_5 - 0,5258 X_3 + 0,6341 \quad (41)$$

$$B_{X_5} = 0,8112; B_{X_3} = -0,2068$$

$$R^2 = 0,3809; S_e = 4,8466; F = 595,8394; F_{1937}^2(5\%) = 4,61$$

$$F_{X_5} = 187,007; F_{X_3} = 12,159; F_{1937}^1(5\%) = 6,63$$

A equação (41) possui um valor de R^2 , ainda, pequeno porém a equação e os coeficientes de regressão são estatisticamente significantes ao nível de significância de 5%.

Vale notar, por exemplo, que o coeficiente de X_3 não possui o sinal coerente com a definição dessa variável. Tem-se, nesse caso, um decréscimo de 0,52 viagens para um acréscimo de uma pessoa no tamanho da família, mantendo-se constante a variável X_5 .

Inicialmente, era de se esperar que a equação (41) produzisse melhores resultados. Isto porque os coeficientes de correlação linear entre X_1 e X_5 e entre X_1 e X_3 , apesar de não serem altos, em termos absolutos, podem ser aceitos como, relativamente, altos, quando comparados com os demais

coeficientes de correlação de X_1 com as outras variáveis, mostrados na Tabela 4. Entretanto, a matriz de correlação entre as variáveis envolvidas em uma equação de regressão linear pode antecipar alguma indicação sobre um possível resultado. Esse é o caso da Tabela 4 que mostra o coeficiente linear entre X_5 e X_3 igual a 0,9535, portanto, muito alto. Surge, então, o problema da multicolinearidade e, podia-se ter suspeitado, desde antes de se obter a equação (41), da qualidade da equação de regressão linear de X_1 em X_5 e X_3 .

A outra equação para o nível de domicílio é dada por:

$$X_1 = 1,6239 X_5 + 0,4068 \quad (42)$$

$$B_{X_5} = 0,61401$$

$$R^2 = 0,3770; S_e = 4,8605; F = 1172,7676; F_{1938}^1(5\%) = 6,64$$

Conforme mostra a equação (42), o valor do coeficiente de determinação múltipla, R^2 , é pequeno. Portanto, a equação (42) não possui todas as qualidades, para ser usada como modelo de previsão de viagens, recomendadas neste trabalho.

3.2.1.4 - Considerações sobre a aplicação 1.

Comparando-se os resultados das equações obtidas para os três níveis de agregação das variáveis, pode-se ser levado a supor que a equação obtida para o nível de zona é melhor do que as obtidas para o nível de domicílio, porque apresentam um R^2 ligeiramente maior e um erro padrão de estimativa menor.

No entanto, Mc Carty⁷ alerta para o significado do

coeficiente de determinação. O coeficiente de determinação, R^2 , mede a variância da variável dependente que é explicada pela combinação linear das variáveis independentes. Assim, quanto maior o valor de R^2 maior é a relação de linearidade entre as variáveis dependente e independentes, não significando, porém, uma maior explicação da relação de causa-efeito, por parte da equação.

Mc Carthy⁷, através da análise de variância mostrou que somente os fatores entre as unidades de agregação são levados em conta nos modelos de regressão linear, deixando-se de lado os fatores internos. Isso mostra a limitação dos modelos à base da correlação linear e a necessidade de se considerar modelos que incorporem esses fatores internos.

É evidente, porém, que a falta de dados que permitam um estudo desagregado do problema da demanda de viagens pode-se utilizar modelos à base da correlação linear embora deva-se estar consciente das limitações dos modelos empregados.

Nesta aplicação, foi seguido o raciocínio empregado por Mc Carthy⁷ e os resultados encontrados mostraram uma tendência similar à verificada pelos resultados dele.

3.2.2 - Aplicação 2

A finalidade dessa aplicação é obter uma equação de regressão linear múltipla que relacione o número de automóveis em cada zona de tráfego com o número de famílias em cada nível de renda, por zona.

Inicialmente, considerou-se a renda familiar das pes

soas residentes em Campina Grande, em 1978, dividida em 20 classes de amplitude igual a Cr\$ 1.000,00, e determinou-se, para cada classe, o quociente automóveis/família, através da razão $NA(J)/NF(J)$,

onde:

$NA(J)$ = número total de automóveis pertencentes às famílias de classe de renda J , $J=1,2,\dots,20$.

$NF(J)$ = número total de famílias de classe de renda J .

Isso possibilitou a verificação e agrupamento daquelas classes que apresentaram o quociente automóveis/família, aproximadamente, igual e tornou-se possível reduzir, de 20 para 4, o número de classes. A tabela 5 mostra os quatro níveis de renda considerados, as classes de renda correspondentes, e o número de automóveis/família, para cada uma dessas classes.

Tabela 5 - Níveis de renda e número médio de automóveis por família em cada classe.

Níveis de renda	Classes	Automóveis família
1	0 - 6.999	0,174
2	7.000 - 13.999	0,701
3	14.000 - 20.999	1,010
4	> 21.000	1,232

Com os níveis de renda definidos, computou-se, por zona, o número de automóveis e o número de famílias, em cada

nível de renda. A Tabela 6 mostra esse resultado.

Com os dados dessa Tabela, exceto os referentes às zonas 2,4,6,11 e 17, calculou-se a equação de regressão linear múltipla $Y = f(X_1, X_2, X_3, X_4)$,

onde:

Y = número de automóveis por zona

X_1 = número de famílias de nível de renda 1, por zona

X_2 = número de famílias de nível de renda 2, por zona

X_3 = número de famílias de nível de renda 3, por zona

X_4 = número de famílias de nível de renda 4, por zona

A Tabela 7 mostra a matriz de correlação entre essas variáveis.

Tabela 7 - Matriz de correlação

	Y	X_1	X_2	X_3	X_4
Y	1,0000	0,4942	0,9411	0,8248	0,7952
X_1		1,0000	0,4430	0,1348	0,0732
X_2			1,0000	0,8073	0,7171
X_3				1,0000	0,8864
X_4					1,0000

A equação obtida foi:

$$Y = 0,1206X_1 + 1,3866X_2 + 0,0729X_3 + 2,6269X_4 - 6,0240 \quad (43)$$

$$B_{X_1} = 0,2043; B_{X_2} = 0,5927; B_{X_3} = 0,0181; B_{X_4} = 0,3392$$

$$R^2 = 0,9434; S_e = 7,9591; F = 54,1583; F_{14}^4(5\%) = 3,34$$

$$F_{X_1} = 6,237; F_{X_2} = 18,804; F_{X_3} = 0,011; F_{X_4} = 5,569; F_{14}^1(5\%) = 4,60$$

A equação (43) possui um valor alto para R^2 e o coefi

Tabela 6 - Número de automóveis e número de famílias, em cada nível de renda, por zona.

Zona	Autos	Número de famílias em cada nível de renda.			
		1	2	3	4
1	51	43	26	14	11
2	0	0	0	0	0
3	96	91	50	24	9
4	0	0	0	0	0
5	89	97	28	12	10
6	0	0	0	0	0
7	30	70	11	7	3
8	28	164	12	2	2
9	18	23	6	1	3
10	73	121	22	20	11
11	0	1	0	0	0
12	12	65	4	1	1
13	7	19	5	0	4
14	59	124	25	5	5
15	20	101	9	2	1
16	40	170	18	1	2
17	0	0	0	0	0
18	3	35	5	1	1
19	12	69	5	4	3
20	6	23	3	0	0
21	12	38	4	1	1
22	58	160	26	2	2
23	19	92	9	0	0

ciente de correlação múltipla, R , e os coeficientes de regressão, exceto o da variável X_3 , são estatisticamente significantes ao nível de significância de 5%.

Entretanto, se os coeficientes de regressão da equação (43) representassem o número médio de automóveis por família, para cada nível de renda, como era de se esperar, a equação (43) se transformaria em:

$$Y = 0,174X_1 + 0,701X_2 + 1,010X_3 + 1,232X_4 \quad (44)$$

Nota-se, portanto, grandes discrepâncias entre os coeficientes das equações (43) e (44). A ocorrência dessas discrepâncias pode ser, principalmente, devida ao efeito da multicolinearidade, pois, as variáveis X_2 , X_3 e X_4 , tomadas duas a duas, possuem o coeficiente de correlação linear, relativamente, alto, conforme mostrado na Tabela 7.

Na equação (43) o coeficiente da variável X_3 é, aproximadamente, nulo. Então, eliminou-se essa variável e calculou-se uma equação de regressão linear múltipla de Y em X_1 , X_2 , e X_4 .

A equação obtida foi:

$$Y = 0,1195X_1 + 1,4049X_2 + 2,7087X_4 - 6,1294 \quad (45)$$

$$B_{X_1} = 0,2025; B_{X_2} = 0,6005; B_{X_4} = 0,3497; F = 54,1583; F_{15}^3 (5\%) = 3,34$$

$$R^2 = 0,9434; S_e = 7,9591$$

$$F_{X_1} = 6,897; F_{X_2} = 29,627; F_{X_4} = 12,435; F_{15}^1 (5\%) = 4,60$$

Analisando-se a equação (44) constata-se que ela possui qualidades que permitem o seu uso como equação de previsão.

Por outro lado, comparando-se os resultados das equações (43) e (45) verifica-se que elas possuem o mesmo valor para R^2 , o mesmo S_e e os valores dos coeficientes de regressão são, praticamente, iguais. Portanto, pode-se usar, também, a equação (43) como equação de previsão, sem deixar, no entanto, de se reconhecer as suas limitações para esse fim.

As equações (43), (44) e (45) podem ser usadas, por exemplo, para se estimar o número de automóveis em uma zona, desde que se conheça o número de famílias em cada nível de renda, por zona de tráfego. Essa aplicação pode ser feita tanto para novas zonas de tráfego como para as zonas já existentes, em situações futuras.

3.3 - Análise Fatorial. Aplicações e Resultados

3.3.1 - Aplicação 1

Essa aplicação utiliza variáveis que descrevem o uso do solo dos bairros de Campina Grande e tem como objetivo reduzir o número dessas variáveis, sem perda significativa de informações. Para isso, empregou-se a Análise Fatorial do tipo R.

As variáveis originais utilizadas representam as frações da área total, de cada bairro, destinadas a cada tipo de uso do solo.

De posse da divisão em 54 bairros feita pela Companhia de Desenvolvimento de Campina Grande - COMDECA, Masayuki Doi¹¹ fez a classificação do uso do solo de Campina Grande conforme mostra o Quadro 1.

Quadro 1 - Classificação do uso do solo

1 - Residencial

1.1 - Residencial unifamiliar

RU₁ - multicômodo isolado

RU₂ - multicômodo conjugado

RU₃ - unicômodo isolado, geminado, conjugado

1.2 - Residencial multifamiliar

RM₁ - edifício multicômodo

RM₂ - edifício unicômodo

2 - Comercial

2.1 - Região comercial (atende a população de todos os bairros)

CR₁ - compras grandes (movelaria, implementos agrícolas, etc.)

CR₂ - compras pequenas (camisas, sapatos, etc.)

CR₃ - serviços (restaurante, postos de gasolina, cinema, etc.)

2.2 - Comercial de bairro (atende, principalmente, a população do bairro)

CB - comércio de bairro (bares, etc.)

2.3 - Comercial de manufatura

CM - comércio de manufatura (oficinas, atacadistas, etc.)

3 - Industrial

I - Indústria

4 - Agricultura

A - agricultura, fazendas, etc.

Continuação do Quadro 1

5 - Público

P - igrejas, escolas, prefeituras, etc.

6 - Espaço Aberto

EA - Campos de futebol, parques, etc.

7 - Utilidade (Equipamentos de utilidade pública de grande consumo de solo)

U - Cagepa, aeroporto, estação ferroviária, etc.

8 - Não Desenvolvido

ND - lotes desocupados.

O Quadro 2 mostra a denominação dos bairros de Campina Grande.

Quadro 2 - Numeração e nomes dos bairros.

Número	Bairro	Número	Bairro
1	São José	28	Conceição
2	Açude Novo	29	Louzeiro
3	Açude Velho	30	Pirineus
4	Centro	31	Bela Vista
5	Ariús	32	Monte Santo
6	Acauã	33	Geremias
7	Vila Cabral	34	Redentorista
8	Sandra Cavalcanti	35	Araxá
9	Catolé	36	Cidade Universitária
10	Prado	37	Centenário
11	Sandra	38	Casa de Pedra
12	Provisão	39	Dona Merquinha

Continuação do Quadro 2

Número	Bairro	Número	Bairro
13	Passa Tempo	40	Produção Mineral
14	Tambor	41	Quarenta
15	Santíssima	42	Santana
16	Cachoeira	43	Santa Rosa
17	Monte Castelo	44	Jardim Nordeste
18	Nova Brasília	45	Moita
19	Vila Castelo Branco	46	Cruzeiro
20	Santo Antonio	47	Liberdade
21	Tavares	48	Graças
22	Lauritzen	49	Jardim Paulistano
23	Alto Branco	50	Três Irmãs
24	José Pinheiro	51	Distrito Industrial
25	Palmeira	52	Adrianópolis
26	Areias	53	Bodocongô
27	Prata	54	Zoobotânico.

Masayuki Doi¹¹ contém os valores das variáveis para cada uso do solo mencionado no Quadro 1, em m².

Definiu-se, para cada bairro, as variáveis:

$$XX1 = (RU_1 + RU_2 + RU_3 + RM_1 + RM_2) / \text{TOTAL}$$

$$XX2 = CR_1 / \text{TOTAL}$$

$$XX3 = (CR_2 + CB) / \text{TOTAL}$$

$$XX4 = CR_3 / \text{TOTAL}$$

$$XX5 = CM / \text{TOTAL}$$

$$XX6 = I / \text{TOTAL}$$

$$XX7 = A / \text{TOTAL}$$

XX 8 = P/TOTAL

XX 9 = EA/TOTAL

XX10 = U/TOTAL

XX11 = ND/TOTAL

Construiu-se, então a matriz X do tipo (54x11) e fez-se uma Análise Fatorial tipo R. A Tabela 8 é um extrato da listagem de computador para essa Análise Fatorial.

Tabela 8 - Numeração dos fatores, autovalores e porcentagem da variância do conjunto das variáveis explicada por cada fator

Fator Fi	Autovalor λ_i	% da variância	% Acumulada da Variância
1	3,14453	28,6	28,6
2	1,67022	15,2	43,8
3	1,22993	12,1	55,9
4	1,21159	11,0	66,9
5	1,03236	9,4	76,2
6	0,92752	8,4	84,7
7	0,84164	7,7	92,3
8	0,66189	6,0	98,3
9	0,16462	1,5	99,8
10	0,01337	0,1	100,0
11	0,01357	0,0	100,0

Conforme visto no enfoque matemático, $\frac{\lambda_i}{m} \times 100$ fornece a porcentagem da variância do conjunto de dados explicada pelo fator i, Fi. Portanto, o primeiro elemento da terceira coluna, 28,6, é igual a $\frac{3,14453}{11} \times 100$. A quarta coluna é obtida somando-se cada elemento da terceira coluna à soma dos valores da terceira coluna, anteriores a esse elemento.

Como se quer reduzir o número de variáveis, o número de fatores da solução final deve ser menor do que 11. Para isso, adota-se um autovalor como mínimo e despreza-se todo

Com base na Tabela 9, pode-se escrever qualquer variável como uma combinação linear dos fatores. Tem-se, por exemplo,

$$XX1 = 0,17549F_1 + 0,05183F_2 + 0,86013F_3 + 0,05274F_4 - 0,21618F_5 \quad (46)$$

Desde que se conheça os componentes de cada fator pode-se calcular os correspondentes valores da variável original $XX1$, dada por (46). As componentes ou escores de um fator podem ser obtidas através da Tabela 11.

Tabela 11 - Matriz dos coeficientes para o cálculo dos escores dos fatores.

	Fator 1	Fator 2	Fator 3	Fator 4	Fator 5
XX1	-0.02061	-0.12757	0.57440	0.05882	-0.17401
XX2	0.45385	-0.07627	-0.09835	-0.07508	-0.00309
XX3	0.43617	-0.15742	0.01215	-0.04920	-0.00293
XX4	0.18454	0.36471	-0.13835	0.02260	0.03849
XX5	0.01644	0.41382	-0.03538	0.12210	-0.00747
XX6	-0.08155	0.18342	-0.06983	0.04638	0.48682
XX7	-0.00747	-0.10301	-0.31164	0.48930	0.06015
XX8	-0.08045	-0.05900	0.43667	0.09898	0.28220
XX9	-0.16048	0.49716	-0.05978	-0.12202	-0.06567
XX10	-0.07955	0.20973	-0.08250	0.04598	-0.72613
XX11	0.06922	-0.01824	-0.19821	-0.63844	0.04127

Os escores de F_1 podem ser dados por:

$$F_1 = -0,02061XX1 + 0,45385XX2 + 0,43617XX3 + \dots + 0,06922XX11 \quad (47)$$

Entrando-se com os valores padronizados de cada observação das variáveis originais em (47), pode-se calcular os escores de F_1 . De forma semelhante, calcula-se os escores dos outros fatores e, construindo-se uma matriz onde os vetores coluna sejam os escores de cada fator, a correlação linear entre cada par de vetores-coluna é zero, visto que os

fator correspondente a um autovalor menor do que o especificado como mínimo.

A menos que se diga o contrário, o autovalor adotado como mínimo, nas aplicações de Análise Fatorial neste trabalho, é 1.

Assim, apenas, os cinco primeiros fatores da Tabela 9 permanecerão na solução final, passando-se de 11 variáveis originais para 5 variáveis novas, chamadas fatores.

A matriz dos carregamentos dos fatores após a rotação é:

Tabela 09 - Matriz de carregamentos dos fatores após a rotação varimax.

	Fator 1	Fator 2	Fator 3	Fator 4	Fator 5
XY1	0.17540	0.05185	-0.06013	0.05274	0.21418
XY2	0.07520	0.07425	0.02619	0.02619	0.01191
XY3	0.04471	0.01530	0.19017	0.04720	-0.02340
XY4	0.00437	0.07601	0.05714	0.12875	0.09513
XY5	0.01088	0.07167	0.16461	0.22540	0.14711
XY6	-0.11945	0.31390	-0.05728	0.05051	-0.52115
XY7	-0.02246	-0.29232	-0.58494	-0.71662	0.04030
XY8	-0.00121	0.12494	-0.63784	0.09052	0.01353
XY9	-0.15201	-0.07032	0.06931	-0.17806	-0.00191
XY10	-0.09279	0.17121	-0.06836	0.08417	-0.00226
XY11	-0.13091	-0.14919	-0.27789	0.01350	0.06005

A Tabela 10 fornece as comunalidades das variáveis

Tabela 10 - Comunalidade das Variáveis

Variável	Comunalidade
XY1	0.02000
XY2	0.0500
XY3	0.02247
XY4	0.03010
XY5	0.03570
XY6	0.24170
XY7	0.04450
XY8	0.52800
XY9	0.56471
XY10	0.09332
XY11	0.05714

4
chevada!

fatores são ortogonais.

Analisando-se a Tabela 9, pode-se destacar em que fator cada variável tem projetada a maior quantidade de sua variância. Isso é mostrado na Tabela 12.

Tabela 12 - Fatores mais importantes na composição da variância das variáveis

Fator 1	Fator 2	Fator 3	Fator 4	Fator 5
XX2	XX5	XX1	XX11	XX10
XX3	XX9	XX8	XX 7	XX 6
	XX4			

Com base na tabela 12, pode-se, então, considerar as 5 variáveis da solução final como uma junção das variáveis originais, como mostra a Tabela 13.

Tabela 13 - Fatores e suas definições em termos das variáveis originais dadas em valores relativos

Fatores	Variáveis Originais	Característica dos Fatores
F ₁	$(CR_1 + CR_2 + CB)$	Variável que caracteriza usos do solo destinados a comércio.
F ₂	$CM + EA + CR_3$	Variável que caracteriza usos do solo destinados a serviços e lazer
F ₃	$RU_1 + RU_2 + RU_3 + RU_4 + RU_5 + P$	Variável que caracteriza usos do solo destinados a residências e a equipamentos de uso público
F ₄	$A + ND$	Variável que caracteriza usos do solo, destinados a agricultura e solos não ocupados.
F ₅	$I + U$	Variável que caracteriza usos do solo destinados a indústria e a equipamentos de utilidade pública.

As grandezas TOTAL 1, TOTAL 2, ..., TOTAL 54 significam o valor da área total de cada bairro.

Os resultados dessa Análise Fatorial são mostrados nas Tabelas 14 a 17.

A Tabela 14 mostra que com os onze fatores iniciais já foi possível explicar 100% da variância dos dados originais. Entretanto, somente os seis fatores iniciais tomarão parte na solução final, em virtude de corresponderem a autovalores maiores que o adotado como mínimo. A Tabela 14 mostra, ainda, que os fatores da solução final explicam 95,6% da variância total dos dados.

Após a rotação varimax, tem-se os carregamentos dos fatores, conforme a Tabela 15.

As comunalidades das variáveis são mostradas na Tabela 16.

Para o cálculo dos escores dos fatores usa-se os elementos da Tabela 17.

Partindo-se da Tabela 15 pode-se construir a Tabela 18.

Tabela 18 - Fatores mais importantes na composição da variância das variáveis.

F ₁	BB3, BB5, BB8, BB9, BB11, BB13, BB15, BB16, BB17, BB18, BB19, BB21, BB23, BB29, BB30, BB31, BB33, BB39, BB40, BB41, BB42, BB43, BB 44, BB48, BB49, BB51, BB53.
F ₂	BB6, BB7, BB12, BB14, BB34, BB35, BB36, BB38, BB45, BB46, BB50, BB52, BB54.
F ₃	BB1, BB2, BB4, BB20, BB22, BB25, BB26, BB27.
F ₄	BB24, BB28, BB32, BB37.
F ₅	BB10.

Tabela 14 - Numeração dos fatores, autovalores e porcentagem da variância do conjunto das variáveis explicada por cada fator.

Fator	Autovalor	%da Variância	%Acumulada da Variância
1	27.11455	50.2	50.2
2	11.77557	21.9	72.1
3	6.63826	12.3	84.3
4	3.12819	5.8	90.1
5	1.97277	3.7	93.8
6	1.03546	1.9	95.6
7	0.78839	1.4	97.1
8	0.71896	1.3	98.4
9	0.41943	0.8	99.2
10	0.30569	0.6	99.7
11	0.17642	0.3	100.0
12	0.27758	0.5	100.0
13	0.00132	0.0	100.0
14	0.00503	0.0	100.0
15	0.00073	0.0	100.0
16	0.00032	0.0	100.0
17	0.00031	0.0	100.0
18	0.00031	0.0	100.0
19	0.00030	0.0	100.0
20	0.00009	0.0	100.0
21	0.00009	0.0	100.0
22	0.00009	0.0	100.0
23	0.00009	0.0	100.0
24	0.00009	0.0	100.0
25	0.00009	0.0	100.0
26	0.00009	0.0	100.0
27	0.00009	0.0	100.0
28	0.00009	0.0	100.0
29	0.00009	0.0	100.0
30	0.00009	0.0	100.0
31	-0.00009	-0.0	100.0
32	-0.00009	-0.0	100.0
33	-0.00009	-0.0	100.0
34	-0.00009	-0.0	100.0
35	-0.00009	-0.0	100.0
36	-0.00009	-0.0	100.0
37	-0.00009	-0.0	100.0
38	-0.00009	-0.0	100.0
39	-0.00009	-0.0	100.0
40	-0.00009	-0.0	100.0
41	-0.00009	-0.0	100.0
42	-0.00009	-0.0	100.0
43	-0.00009	-0.0	100.0
44	-0.00009	-0.0	100.0
45	-0.00009	-0.0	100.0
46	-0.00009	-0.0	100.0
47	-0.00009	-0.0	100.0
48	-0.00009	-0.0	100.0
49	-0.00009	-0.0	100.0
50	-0.00009	-0.0	100.0
51	-0.00009	-0.0	100.0
52	-0.00001	-0.0	100.0
53	-0.00001	-0.0	100.0
54	-0.00003	-0.0	100.0

Tabela 15 - Matriz dos carregamentos dos fatores após rotação varimax.

Variável	Fator 1	Fator 2	Fator 3	Fator 4	Fator 5	Fator 6
B31	-0.03042	0.11603	0.27030	-0.05344	0.03903	0.39394
B32	0.26271	-0.07136	0.03944	0.12906	-0.11862	-0.13190
B33	-0.62034	-0.01093	0.20217	0.45405	0.54492	0.17713
B34	-0.11934	-0.13034	0.31936	0.33470	-0.00542	0.03574
B35	-0.02156	-0.00365	0.16634	0.24504	0.06547	0.15932
B36	0.57959	0.58241	0.01228	0.15126	0.53559	0.04824
B37	-0.02532	0.97517	-0.06258	0.12214	-0.02512	0.00305
B38	0.88004	0.44519	0.02011	0.05150	0.03765	-0.01542
B39	0.93744	-0.00079	0.06690	0.11826	0.05682	-0.01137
B310	0.07046	0.06086	0.15187	0.08917	0.92583	-0.01955
B311	0.65041	-0.02333	0.05956	0.33036	0.37762	-0.07597
B312	0.52902	0.81686	0.16414	-0.08981	-0.02016	-0.00815
B313	0.90594	-0.00284	0.02742	0.14660	-0.04020	-0.15475
B314	0.56601	0.81795	-0.03020	0.00273	0.03169	0.00871
B315	0.72370	-0.07021	0.15760	0.41266	0.43287	0.09791
B316	0.98365	0.02950	-0.01833	0.10163	-0.02099	0.00271
B317	0.66058	-0.01419	0.35031	0.64165	-0.07479	0.02214
B318	0.94309	0.01269	0.09112	0.23153	0.11646	-0.05862
B319	0.90504	0.02324	0.03784	0.04910	0.02184	0.00652
B320	0.33057	-0.01688	0.01310	0.10204	-0.00195	-0.07405
B321	0.90708	0.06033	0.01476	0.01462	0.00081	0.00512
B322	0.46711	-0.01901	0.80255	0.32776	0.04525	-0.04657
B323	0.86676	0.08718	0.32212	0.34645	0.03318	-0.02693
B324	0.23098	-0.06363	0.18760	0.93271	0.00304	0.07263
B325	0.26737	-0.06359	0.89703	0.23876	0.10285	-0.05311
B326	-0.12555	-0.16678	0.75100	0.09373	0.30129	-0.05250
B327	-0.01723	-0.05519	0.96831	0.15036	0.01201	0.13797
B328	0.03246	-0.08010	0.34010	0.67919	0.09115	0.57611
B329	0.07477	0.43614	0.12836	0.14285	0.03725	0.02047
B330	0.88208	0.08020	0.43501	0.08174	0.01283	-0.02635
B331	0.65399	0.05345	0.42320	0.55642	0.07161	0.16896
B332	0.37512	0.11319	0.23936	0.87747	0.12269	0.04416
B333	0.01308	-0.00331	0.09283	0.54755	-0.01740	-0.06999
B334	-0.05242	0.99136	-0.06964	-0.03924	0.02545	0.07175
B335	-0.10685	0.99045	-0.07461	-0.02566	0.01343	0.00936
B336	0.45771	0.88265	-0.04784	-0.04091	0.04116	0.02049
B337	0.40621	-0.02907	0.27013	0.38468	0.03045	-0.03965
B338	0.47452	0.87553	-0.00972	-0.02441	0.03337	-0.00789
B339	0.99654	0.03717	0.02167	-0.00655	0.01986	-0.00691
B340	0.91264	0.00992	-0.04990	0.13193	-0.03583	0.53489
B341	0.97648	0.15084	0.00911	0.11837	0.01366	0.01574
B342	0.07832	0.02665	0.18819	0.08321	-0.00951	0.02568
B343	0.72139	0.51924	0.01620	0.41261	-0.02214	0.02322
B344	0.06515	0.06413	0.07307	0.22342	0.02751	-0.03905
B345	0.00622	0.99044	-0.09444	0.02760	-0.01072	0.00903
B346	-0.10654	0.98941	-0.07933	-0.02603	0.01123	0.00712
B347	-0.01610	-0.03607	0.44511	0.80740	0.14953	-0.16094
B348	0.02260	-0.01547	0.12563	0.22903	0.08541	0.11547
B349	0.09491	0.05327	0.00052	0.01764	0.03010	0.01141
B350	0.00920	0.09135	-0.05938	0.05256	0.00741	0.00330
B351	0.32724	-0.00128	-0.05771	-0.07789	0.53217	0.03452
B352	-0.11451	0.78041	-0.18045	-0.08352	-0.11746	-0.21925
B353	0.92237	0.33152	0.05714	0.07422	0.14420	0.06373
B354	-0.04738	0.98403	-0.08452	-0.05457	0.13561	0.02013

Tabela 16 - Comunalidade das Variáveis.

Variável	Comunalidade
801	0.93290
802	0.84431
803	1.06070
804	0.91883
805	0.97164
806	0.96733
807	0.97574
808	0.99157
809	0.90744
810	0.90006
811	0.92461
812	0.98763
813	0.97022
814	0.99372
815	0.91607
816	0.97006
817	0.95745
818	0.96504
819	0.99552
820	0.96465
821	0.99035
822	0.97657
823	0.98452
824	0.97204
825	0.97500
826	0.70990
827	0.90373
828	0.92423
829	0.99582
830	0.93141
831	0.95951
832	0.99753
833	0.97620
834	0.99713
835	0.99951
836	0.99466
837	0.99899
838	0.99423
839	0.99542
840	0.96777
841	0.99077
842	0.99405
843	0.96132
844	0.99613
845	0.99389
846	0.99753
847	0.99663
848	0.54048
849	0.99773
850	0.99760
851	0.97900
852	0.70110
853	0.99439
854	0.99594

Tabela 17 - Matriz dos coeficientes para o cálculo dos escores dos fatores.

VARIÁVEL	FATOR 1	FATOR 2	FATOR 3	FATOR 4	FATOR 5	FATOR 6
B01	-0.00789	0.01934	0.15983	-0.12156	-0.01210	0.36701
B02	0.00311	0.00442	0.14781	-0.02132	-0.00740	-0.15685
B03	-0.00459	-0.01131	-0.02137	0.04309	-0.01955	0.11106
B04	-0.02676	0.00764	0.12231	0.02136	-0.00839	-0.00188
B05	0.04060	-0.01191	-0.00204	-0.00436	-0.02123	0.13661
B06	-0.00085	0.03903	-0.01624	-0.00983	0.23151	0.01210
B07	-0.02371	0.09168	-0.00736	0.04787	-0.02760	0.00821
B08	0.04341	0.02572	0.00109	-0.02766	-0.01966	-0.02530
B09	0.00293	-0.01619	-0.00547	-0.02108	-0.00993	-0.03155
B010	-0.03434	-0.00242	0.00194	-0.02948	0.45386	-0.07238
B011	0.02519	-0.02034	-0.03035	0.03083	0.14847	-0.12204
B012	0.02032	0.06967	0.05029	-0.05802	-0.01845	-0.00543
B013	0.05245	-0.01640	-0.00976	0.00443	-0.05113	-0.16362
B014	0.01024	0.06577	0.00515	-0.02132	-0.01426	0.00920
B015	0.00003	-0.01920	-0.02387	0.03919	0.18766	0.03075
B016	0.05717	-0.01467	-0.01815	-0.01647	-0.05209	-0.01933
B017	0.01222	-0.00747	-0.04546	0.13510	-0.00943	-0.00936
B018	0.04127	-0.01477	-0.01480	0.02051	0.01451	-0.00985
B019	0.05792	-0.01436	-0.00439	-0.03794	-0.02724	-0.00787
B020	0.01173	0.01011	0.16056	-0.06240	-0.03327	-0.10305
B021	0.06001	-0.01222	-0.00487	-0.04449	-0.03486	-0.00172
B022	0.02721	0.06615	0.11865	-0.00197	-0.00134	-0.10927
B023	0.03293	-0.00002	0.02584	0.02403	-0.03477	-0.06298
B024	-0.02276	-0.00064	-0.04909	0.21410	-0.05225	0.01079
B025	-0.00562	0.00756	0.13906	-0.01433	0.01376	-0.00832
B026	-0.02460	-0.00246	0.12465	-0.04385	0.13948	-0.00846
B027	-0.01603	0.01812	0.16401	-0.06083	-0.02546	0.10504
B028	-0.03320	0.00586	-0.01073	0.11167	-0.02277	0.01628
B029	0.03707	0.02814	0.01222	-0.01230	-0.02740	0.00219
B030	0.04630	-0.00003	0.06985	-0.05524	-0.00385	-0.05091
B031	0.02675	0.00529	0.02265	0.06748	-0.02325	0.11571
B032	-0.02624	0.01241	-0.03430	0.13855	0.00313	-0.02077
B033	0.02576	-0.01065	-0.03341	0.10311	-0.05828	-0.10975
B034	-0.02073	0.00324	0.00082	-0.00186	0.00122	0.03256
B035	-0.02369	0.00866	0.00800	0.00433	-0.00025	0.02159
B036	0.01202	0.07328	0.00720	-0.02585	-0.00506	0.02494
B037	-0.01728	-0.00058	-0.02451	0.13686	-0.03833	-0.10045
B038	0.01359	0.07304	0.01046	-0.02362	-0.00961	-0.00467
B039	0.06093	-0.01476	-0.00226	-0.05044	-0.02502	-0.01745
B040	0.04056	-0.00896	-0.03735	-0.02816	-0.00749	0.01444
B041	0.05137	-0.00292	-0.01542	-0.01577	-0.03488	-0.00049
B042	0.05800	-0.01140	0.02789	-0.06197	-0.04464	0.01292
B043	0.01657	0.03874	-0.02910	0.07715	-0.06126	-0.00241
B044	0.04702	-0.00778	-0.01055	0.00011	-0.02982	-0.06310
B045	-0.01819	0.00181	-0.00852	0.02077	-0.01801	0.01860
B046	-0.02362	0.00349	-0.00724	0.00983	-0.00111	0.01962
B047	-0.04500	0.00700	0.01411	-0.17861	0.04040	-0.22218
B048	0.04137	-0.01471	-0.00754	-0.03360	-0.00540	0.03439
B049	0.05938	-0.01227	0.00272	-0.04871	-0.02280	-0.00159
B050	-0.01490	0.00002	0.00309	0.02006	-0.01337	0.00916
B051	0.03736	-0.00640	-0.02337	-0.07856	0.24090	0.00556
B052	-0.01058	0.02261	0.01196	0.01392	-0.00377	-0.10530
B053	0.04184	0.01454	-0.00746	-0.03617	0.02042	0.04412
B054	-0.02278	0.00009	0.00306	-0.00594	0.00984	0.02756

Uma análise detalhada das características dos fatores da solução final revela as variáveis originais que predominam na formação dos fatores. A Tabela 19 apresenta os fatores e suas características.

Tabela 19 - Fatores e suas definições em termos das variáveis originais predominantes, em valores relativos.

Fatores	Variáveis Originais Predominantes	Características
F ₁	ND	bairros já loteados com baixo índice de uso do solo
F ₂	ND + A	bairros com baixo índice de uso do solo, parcialmente loteado e com uso do solo para agricultura
F ₃	R + CR + P	bairros residencial e/ou comercial e/ou de uso público
F ₄	R + (P+EA) + ND	bairro residencial com área para lazer e com loteamento
F ₅	I	bairro industrial

A Tabela 19 mostra o resultado do agrupamento dos bairros com padrões semelhantes de uso do solo.

Na Tabela 18, tem-se os bairros pertencentes a cada um dos grupos e, identificando-se aqueles bairros, especialmente vizinhos, pertencentes a um mesmo grupo, pode-se juntá-los para formar novas unidades de agregação, de mesmo padrão de uso do solo. Então, pode-se usar esse resultado para, por exemplo, estudar o processo de geração de viagens dessas novas unidades de agregação.

3.4. - Análise Discriminante. Aplicações e resultados

3.4.1. - Aplicação 1

O objetivo dessa aplicação é, dados três grupos de observações formados por pessoas que exercem atividade trabalho no setor secundário, no terciário e em ambos, no secundário e terciário, verificar se esses grupos apresentam com portamentos distintos, usando a Análise Discriminante.

Para essa aplicação, observou-se em que setor da economia era exercida a atividade trabalho das pessoas ocupadas, residentes em Campina Grande, e agrupou-se, separadamente, os indivíduos com atividade trabalho no setor secundário, no terciário e aqueles com atividade trabalho nos setores secundário e terciário.

Em seguida, dividiu-se o período de 24 horas em cinco intervalos de tempo, 5hs às 7hs, 7hs às 9hs, 9hs às 13hs, 13hs às 15hs e 15hs às 24 hs, e considerou-se os seguintes modos de transporte: ônibus, carro-motorista, carro-acompanhante, taxi, a pé e bicicleta. Os modos a pé e bicicleta foram agrupados passando a formar um único modo.

Para cada indivíduo, construiu-se um perfil do seu comportamento, em relação aos transportes. Esse perfil era representado por um vetor-linha de 30 componentes, isto é, 30 variáveis, onde as primeiras 25 variáveis correspondiam aos cinco intervalos de tempo para cada modo de transporte, e as 5 últimas à duração da atividade trabalho do indivíduo, em cada intervalo. Inicialmente, o vetor perfil era zerado e, para cada viagem casa-trabalho realizada, alocava-se o

número 1 à posição correspondente, simultaneamente, ao modo de transporte utilizado e ao intervalo de tempo no qual a viagem teve início. Ao intervalo de tempo no qual o indivíduo começou a sua atividade, alocou-se o valor da duração daquela atividade. Dessa forma obteve-se os três grupos de observação.

A Tabela 20 mostra o número de casos em cada grupo.

Tabela 20 - Número de casos por grupo

Grupo	Número de casos
1	249
2	2208
→ 3	7 ← ?
Total	2464

Com base no número de casos apresentados na Tabela 20, pode-se questionar a representatividade do grupo 3, quanto ao comportamento das pessoas com atividade trabalho, nos setores secundário e terciário, visto que esse grupo contém um número pequeno de casos.

Ao término do processo da Análise Discriminante, verificou-se, com base na Tabela 21, que a atividade trabalho é um fator significativa para a distinção dos grupos.

Tabela 21 - Lambda de Wilks e valor do $F_{\text{equivalente}}$

		Graus de Liberdade	significância
Lambda de Wilks	0,9673821		
$F_{\text{equivalente}}$	3,725320	20 e 4904	0.0000

O lambda de Wilks ou o F equivalente calculado a partir desse vetor dá indicativos sobre o poder discriminatório do conjunto de variáveis incluídas na análise. Para saber se o valor de F é significativo, testa-se com o F crítico para os números de graus de liberdade fornecidos na tabela 21 e para um nível de significância estabelecido.

Para esta aplicação, foi encontrado que o valor de F é estatisticamente significativa a 5% de nível de significância, em virtude de o F equivalente ser maior do que o F crítico para 20 e 4904 graus de liberdade. Isso, em outras palavras, equivale a dizer que o critério usado na construção dos grupos, realmente, produz diferença significativa nos grupos.

Pode-se, também, verificar a significância estatística da diferença entre pares de grupos, e isso é mostrado na Tabela 22.

Tabela 22 - Valores de F e significância entre pares de grupos. Cada valor F tem 11 e 2451 graus de liberdade.

Grupos	1	2
2	6,2743*	
	0,0000**	
3	1,5003	1,1880
	0,1242	0,2895

* - Valor de F

** - Significância entre pares de grupos

O valor de F apresentado na Tabela 22 é o F calculado

para cada par de grupos. Admite-se as hipóteses:

H_0 : Os grupos não são diferentes

H_1 : Os grupos são diferentes

Compara-se o $F_{\text{calculado}}$ com o $F_{\text{crítico}}$ tirado de uma tabela F para 11 e 2451 graus de liberdade e um nível de significância desejado; nesse trabalho está sendo considerado um nível de significância de 5%. Se $F_{\text{calculado}} > F_{\text{crítico}}$, rejeita-se H_0 e tem-se que os grupos são diferentes; se $F_{\text{calculado}} < F_{\text{crítico}}$, aceita-se H_0 e, nesse caso, os grupos não são significativamente diferentes.

Nessa aplicação tem-se $F_{\text{crítico}} = 1,79$ e, portanto, apenas os grupos 1 e 2 são significativamente diferentes, entre si.

A significância entre pares de grupos na Tabela 22 indica o nível de significância a partir do qual se tem $F_{\text{calculado}} > F_{\text{crítico}}$.

A Análise Discriminante selecionou como variáveis discriminatórias as variáveis originais mostradas na Tabela 23.

Na Tabela 23 não aparecem variáveis que representam o uso de modos de transporte no 2º intervalo de tempo. No entanto, esse tipo de variável é importante no comportamento dos grupos porque se trata de valores para o intervalo entre 7hs e 9hs, horário em que há um grande número de viagens para o trabalho, principalmente, para o setor terciário.

O fato de a Análise Discriminante Stepwise trabalhar com o objetivo de encontrar a máxima discriminação entre os grupos, excluindo as variáveis que não contribuem significativamente para tal objetivo, fez com que fosse excluída essa

Tabela 23 - Variáveis discriminatórias

Variável	Modo	Intervalo de Tempo
X ₃	Ônibus	30
X ₆	Carro-Motorista	10
X ₈	Carro-Motorista	30
X ₁₁	Carro-Motorista	10
X ₁₆	Taxi	10
X ₂₁	Pé+Bicicleta	10
X ₂₃	Pé+Bicicleta	30
X ₂₆	Duração	10
X ₂₇	Duração	20
X ₂₉	Duração	40
X ₃₀	Duração	50

variável, importante no comportamento dos grupos. Essa é uma limitação do uso da análise discriminante Stepwise nesse tipo de aplicação. A garantia da inclusão dessas variáveis na Análise Discriminante pode ser obtida, pela realização de uma Análise Discriminante, Método Direto, forçando-se a entrada de todas as variáveis na solução final.

A Análise Discriminante realizada nessa aplicação deve fornecer, no máximo, duas funções discriminantes. Essas funções são combinações lineares das variáveis discriminatórias e o valor numérico ou score de cada observação, em termos de função discriminante, é obtido entrando-se com os valores observados das variáveis discriminatórias, na forma padronizada, nas funções discriminantes. A média aritmética

dos escores das observações de um grupo, em uma função, é chamada média do grupo relativo aquela função. Para um mesmo grupo, a média dos escores em todas as funções é chamada centróide do grupo.

A Tabela 24 dá algumas informações sobre as funções discriminantes obtidas, e a Tabela 25 mostra os coeficientes das variáveis discriminatórias, em cada função discriminante.

Tabela 24 - Funções discriminantes

	Função	
	1	2
Autovalor	0.02842	0.00515
% da Variância	84.66	13.34
% Acumulada da Variância	84.66	100.00
Lambda de Wilks	0.9673821	0.0048769
Qui-quadrado	81.445	12.615
Graus de Liberdade	22	10
Significância	0.0000	0.2460

Tabela 25 - Coeficientes, na forma padronizada, das funções discriminantes.

Variável	Função 1	Função 2
X ₃	-0.14137	0.24622
X ₆	-0.49900	0.34869
X ₈	-0.09623	-0.58622
X ₁₁	-0.18422	0.10923
X ₁₆	-0.21263	0.05284
X ₂₁	-0.33142	0.11730
X ₂₃	0.46099	0.12840
X ₂₆	0.98837	0.06180
X ₂₇	-0.25150	0.28344
X ₂₉	0.06199	0.41729
X ₃₀	0.06315	-0.45406

Da Tabela 25, tem-se a função discriminante 1 como:

$$D_1 = -0,14137X_3 - 0,49900X_6 - 0,09623X_8 - \dots + 0,06315X_{30} \quad (48)$$

A função discriminante 2 é obtida de maneira semelhante à (48).

Através dos coeficientes de (48) pode-se avaliar a importância de cada variável discriminatória como expressão da diferença de comportamento entre os grupos. Aquela variável de maior coeficiente em (48) dá maior contribuição para a discriminação dos grupos, isto é, expressa melhor a diferença de comportamento entre os grupos.

Conforme mostrado no enfoque matemático da Análise Discriminante, a interpretação Bayesiana permite alocar uma dada observação a um dos grupos já definidos. A alocação se dá para o grupo no qual a observação tem a maior probabilidade de pertencer.

Uma outra forma de se fazer a classificação de uma observação é com o uso das funções de classificação, também chamadas de funções discriminantes lineares de Fisher.

A Tabela 26 mostra os coeficientes das funções de classificação, para cada grupo.

Tabela 26 - Coeficientes das funções de classificação.

X ₁	1,2513	0,0000	0,0000
X ₂	-1,3961	0,0000	0,0000
X ₃	1,4584	0,0000	0,0000
X ₄	0,0000	0,0000	0,0000
X ₅	0,0000	0,0000	0,0000
X ₆	0,0000	0,0000	0,0000
X ₇	0,0000	0,0000	0,0000
X ₈	0,0000	0,0000	0,0000
X ₉	0,0000	0,0000	0,0000
X ₁₀	0,0000	0,0000	0,0000
X ₁₁	0,0000	0,0000	0,0000
X ₁₂	0,0000	0,0000	0,0000
X ₁₃	0,0000	0,0000	0,0000
X ₁₄	0,0000	0,0000	0,0000
X ₁₅	0,0000	0,0000	0,0000
X ₁₆	0,0000	0,0000	0,0000
X ₁₇	0,0000	0,0000	0,0000
X ₁₈	0,0000	0,0000	0,0000
X ₁₉	0,0000	0,0000	0,0000
X ₂₀	0,0000	0,0000	0,0000
(Constante)	2,0178	1,1540	0,0000

A função de classificação para o grupo 1 é dada por:

$$FC_1 = 1,2513X_3 - 1,3961X_6 + 1,4584X_8 - \dots - 2,0178 \quad (49)$$

As demais funções de classificação FC_2 e FC_3 são construídas de maneira similar à (49)

Entrando-se com os valores das variáveis discriminatórias de uma observação nas funções FC_1 , FC_2 e FC_3 obtém-se

os escores daquela observação nessas funções. A observação, então, será alocada para o grupo correspondente à função de mais alto escore.

Sob a hipótese de uma distribuição normal multivariada, os escores de classificação podem ser convertidos em probabilidades de a observação pertencer a um grupo. Assim, alocar a observação para o grupo de mais alto escore de classificação é equivalente a se alocar a observação para o grupo no qual ela tem a maior probabilidade de pertencer. Isso vai ao encontro da abordagem Bayesiana na Análise Discriminante. É sempre conveniente usar a abordagem Bayesiana na classificação das observações quando os prejuízos causados por uma classificação errada são muito grandes, ou quando os grupos são de tamanhos muito diferentes ou, ainda, quando se deseja tirar proveito do conhecimento das probabilidades a priori de uma observação pertencer a um grupo².

Para essa aplicação é conveniente usar a abordagem Bayesiana na classificação dos grupos devido a diferença acentuada nos tamanhos dos grupos.

A Tabela 27 mostra a classificação Bayesiana para uma parte das observações.

Como resultado da classificação de todas as observações, tem-se:

Tabela 27 - Classificação das observações

Arquivo	Nº do caso	Grupo real	Maior Probabilidade		Segunda maior Probabilidade		Escores	Discriminante	
			Grupo	P (X/G)	Grupo	P (G/X)			
CM4T81.0	241	1 ***	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	242	1	1	0.9939	0.4779	2	0.3956	0.5437	0.0026
CM4T81.0	243	1	1	0.2951	0.4741	2	0.2444	2.0334	0.3208
CM4T81.0	244	1 ***	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	245	1 ***	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	246	1	1	0.5617	0.6282	2	0.3317	1.3993	-0.5700
CM4T81.0	247	1 ***	2	0.7069	0.5492	1	0.3493	-0.5610	-0.6852
CM4T81.0	248	1 ***	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	249	1 ***	3	0.7350	0.7564	2	0.1260	0.0347	2.0421
CM4T81.0	250	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	251	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	252	2	2	0.6033	0.5543	1	0.3754	-0.4373	-0.9274
CM4T81.0	253	2	2	0.6589	0.5483	1	0.2595	-0.4453	-0.7803
CM4T81.0	254	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	255	2	2	0.5974	0.5607	1	0.3546	-0.4934	-0.9192
CM4T81.0	256	2	2	0.6909	0.5457	1	0.3619	-0.4916	-0.7476
CM4T81.0	257	2	2	0.4029	0.5633	1	0.3953	-0.2550	-1.3210
CM4T81.0	258	2	2	0.5216	0.5715	1	0.3650	-0.5360	-1.0411
CM4T81.0	259	2	2	0.5637	0.5630	1	0.3596	-0.5561	-0.9434
CM4T81.0	260	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	261	2	2	0.4733	0.3587	1	0.3469	-0.5921	-0.7143
CM4T81.0	262	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	263	2 ***	1	0.4733	0.5593	2	0.3598	0.9430	-0.7507
CM4T81.0	264	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	265	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	266	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	267	2	2	0.7448	0.5375	1	0.3684	-0.4254	-0.6794
CM4T81.0	268	2 ***	3	0.0226	0.8740	1	0.0140	0.3146	4.0462
CM4T81.0	269	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	270	2	2	0.3663	0.4644	1	0.2975	-0.5831	0.0865
CM4T81.0	271	2	2	0.3925	0.4447	1	0.3006	-0.4905	0.1385
CM4T81.0	272	2	2	0.5849	0.5657	1	0.3595	-0.5507	-0.9159
CM4T81.0	273	2	2	0.9262	0.4817	1	0.4072	-0.0612	-0.3979
CM4T81.0	274	2	2	0.1931	0.5804	1	0.4008	-0.3550	-1.9204
CM4T81.0	275	2	2	0.7993	0.5295	1	0.3609	-0.4401	-0.5541
CM4T81.0	276	2	2	0.8496	0.4732	1	0.2959	-0.6243	0.0411
CM4T81.0	277	2	2	0.4753	0.3570	1	0.3023	-0.9537	-0.6212
CM4T81.0	278	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	279	2	2	0.5213	0.5629	1	0.3528	-0.3799	-0.3281
CM4T81.0	280	2	2	0.1933	0.5124	1	0.4796	0.1525	-1.7932
CM4T81.0	281	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	282	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	283	2	2	0.5935	0.5574	1	0.3740	-0.4515	-0.9430
CM4T81.0	284	2 ***	3	0.0073	0.9911	2	0.0052	-0.5809	4.4285
CM4T81.0	285	2	2	0.8857	0.4675	1	0.4149	-0.0441	-0.4939
CM4T81.0	286	2	2	0.8097	0.4901	1	0.2925	-0.7039	-0.0455
CM4T81.0	287	2	2	0.3796	0.3911	1	0.3062	-0.2390	0.4658
CM4T81.0	288	2	2	0.9673	0.4781	1	0.3731	-0.2141	-0.2066

Tabela 28 - Resultados da classificação

Grupo real	Número de casos	Grupo estimado		
		1	2	3
1	249	89	144	16
		35,7%	57,8%	6,4%
2	2208	421	1650	137
		19,1%	74,7%	6,2%
3	7	0	5	2
		0%	71,4%	28,6%

Com esse resultado verificou-se que existem pessoas de um grupo que se comportam como pertencentes a outro grupo.

Para essa aplicação o número total de casos corretamente classificados foi de 70,66%.

3.4.2 - Aplicação 2

O objetivo dessa aplicação é, para as pessoas com atividade trabalho no setor secundário e utilizando-se as mesmas variáveis da aplicação 3.4.1, fazer uma Análise Discriminante para verificar se o nível de renda é um fator da distinção significativa dos grupos.

Nesse caso, considerou-se quatro níveis de renda, em valores de 1978, e os grupos ficaram constituídos por:

Grupo 1 - NR_1 : Renda < 1111

Grupo 2 - NR_2 : 1111 < Renda < 2222

Grupo 3 - NR_3 : 2222 < Renda < 3333

Grupo 4 - NR_4 : Renda > 3333

onde NR_1 , NR_2 , NR_3 e NR_4 são os níveis de renda de 1 a 4, respectivamente.

O número de casos em cada grupo é mostrado na Tabela 29.

Tabela 29 - Número de casos por grupo

Nível de renda	Número de casos
1	26
2	69
3	48
4	97
TOTAL	240

Conforme mostrado na Tabela 29, não há grandes diferenças nos tamanhos dos grupos, tornando-se, praticamente, impossível avaliar-se inicialmente, se os grupos podem, ou não, ser representativos do comportamento das pessoas de cada nível de renda.

As informações acerca do critério nível de renda como fator da distinção dos grupos são dadas na Tabela 30.

Tabela 30 - Lambda de Wilks e $F_{\text{equivalente}}$.

		Graus de Liberdade	Significância
Lambda Wilks	0,8018216		
$F_{\text{equivalente}}$	2,191764	24 e 664,8	0,0009

A coluna referente à significância indica que F equivalente é maior do que qualquer valor de F com 24 e 664,8 graus de liberdade a um nível de significância maior do que 0.0009%. Portanto, como se está usando um nível de significância de 5%, concluí-se que o fator nível de renda distingue, significativamente, os grupos.

A tabela 31 mostra os valores de F e as significâncias entre pares de grupos.

Tabela 31 - Valores de F, para 8 e 229 graus de liberdade, e significância para os pares de grupos.

Grupo	1	2	3
2	1.6331 0.1164		
3	1.5783 0.1322	0.90482 0.5132	
4	2.0147 0.0457	4.0744 0.0001	2.0749 0.0392

Baseando-se na tabela 31, verifica-se que o grupo 4 difere dos demais grupos a 5% de nível de significância. Porém, para esse nível de significância, os outros pares de grupos não apresentam diferenças significantes. Esse fato poderia ser usado para se tentar redefinir os grupos, juntando-se aqueles com diferenças menos significantes. Por exemplo, poderia se juntar os grupos 2 e 3 dessa aplicação e fazer Análise Discriminante com os grupos 1, (2+3) e 4, isto

é, com três grupos.

As variáveis discriminatórias são mostrada na tabela 32.

Tabela 32 - Variáveis discriminatórias

Variável	Modo	Intervalo
X ₁	Ônibus	1º
X ₃	Ônibus	3º
X ₆	Carro-Motorista	1º
X ₇	Carro-Motorista	2º
X ₁₂	Carro-Acompanhante	2º
X ₂₁	Pé+Bicicleta	1º
X ₂₂	Pé+Bicicleta	2º
X ₂₇	Duração	2º

A tabela 32 mostra que as variáveis que discriminam o comportamento dos indivíduos, segundo o nível de renda, são os modos de transporte, exceto taxi, usados, principalmente, no período entre 5 e 9hs, bem como a duração da atividade no 2º intervalo. Esse resultado é coerente, pois, é de se supor que o modo de transporte utilizado pelas pessoas ocupadas no secundário, em Campina Grande, seja um indicador do nível de renda dessas pessoas.

A tabela 33 fornece os coeficientes das funções discriminantes.

Tabela 33 - Coeficientes, na forma padronizada, das funções discriminantes.

Variáveis	Função 1	Função 2	Função 3
X ₁	0.14008	0.37929	0.83615
X ₃	0.38398	0.08126	0.14618
X ₆	0.30927	0.11965	0.14654
X ₇	0.37472	0.27236	0.25502
X ₁₂	0.41472	0.05356	0.28903
X ₂₁	0.50813	0.58468	0.04371
X ₂₂	0.15318	1.02007	0.38854
X ₂₇	0.26781	0.96556	0.24579

Então, a função discriminante 1 é dada por:

$$D_1 = - 0,14008X_1 + 0,38398X_3 + \dots + 0,26781X_{27} \quad (48)$$

A tabela 34 mostra a classificação Bayesiana para uma parte das observações.

Os resultados dessa classificação são mostrados na Tabela 35.

Tabela 34 - Classificação das observações

Arquivo	Nº do Grupo caso		Maior Probabilidade Grupo P (X/G) P (G/X)		Segunda maior Probabilidade Grupo P (G/X)	Escores Discriminantes		
MIWEN.D	1	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	2	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	3	1	1	0.3149 0.6364	3 0.1328	0.7646	-2.3650	-0.3381
MIWEN.D	4	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	5	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	6	1	1	0.3149 0.6364	3 0.1328	0.7646	-2.3650	-0.3381
MIWEN.D	7	1 ***	3	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	8	1	1	0.3263 0.6304	3 0.1355	0.7719	-2.3384	-0.3157
MIWEN.D	9	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	10	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	11	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	12	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	13	1 ***	2	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	14	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	15	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	16	1	1	0.3263 0.6304	3 0.1355	0.7719	-2.3384	-0.3157
MIWEN.D	17	1 ***	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	18	1 ***	4	0.3407 0.5344	3 0.2034	2.1756	0.3215	0.0906
MIWEN.D	19	1	1	0.2795 0.6567	3 0.1245	0.7417	-2.4522	-0.4082
MIWEN.D	20	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	21	1 ***	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	22	1	1	0.3092 0.6237	3 0.1325	0.7301	-2.3073	-0.2995
MIWEN.D	23	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	24	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	25	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	26	1	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	27	2 ***	3	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	28	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	29	2 ***	3	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	30	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	31	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	32	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	33	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	34	2 ***	1	0.3858 0.5584	3 0.1501	0.3104	-2.1536	-0.1930
MIWEN.D	35	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	36	2 ***	3	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	37	2 ***	3	0.3021 0.4854	2 0.1937	0.4417	-0.4919	2.1500
MIWEN.D	38	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	39	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	40	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	41	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	42	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	43	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	44	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	45	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	46	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969
MIWEN.D	47	2	2	0.6094 0.4500	3 0.3170	-1.5293	0.5634	0.6395
MIWEN.D	48	2 ***	1	0.9258 0.2876	2 0.2519	-0.1991	-0.1623	-0.6969

Tabela 35 - Resultado da classificação

Grupo Real	Número de Casos	Grupo estimado			
		1	2	3	4
1	26	21	2	2	1
		80.8%	7.7%	7.7%	3.8%
2	69	39	22	7	1
		56.5%	31.9%	10.1%	1.4%
3	48	20	12	11	5
		41.7%	25.0%	22.9%	10.4%
4	97	53	10	5	26
		54.6%	10.3%	8.2%	26.8%

A tabela 35 mostra que a maioria das observações dos grupos 2, 3 e 4 estão alocadas inadequadamente. Daí a porcentagem dos casos corretamente classificados ter sido 33,33%, isto é, muito baixa.

Portanto, em situações semelhantes a essa recomenda-se realizar uma Análise Discriminante para uma nova definição dos grupos, através da junção daqueles grupos que apresentem diferenças menos significantes, quando testados dois a dois.

CAPÍTULO IV

CONCLUSÕES E SUGESTÕES PARA PESQUISAS FUTURAS

4.1 - Conclusões

As conclusões tiradas com o desenvolvimento deste trabalho podem ser enumeradas por:

- 1 - O conhecimento das abordagens matemáticas das técnicas estatísticas multivariadas, usadas neste trabalho, propicia um mais rápido entendimento das aplicações práticas de cada técnica. No entanto, o analista de transportes não deve superestimar o enfoque matemático em detrimento das relações lógicas aparentes de cada caso. Deve, isto sim, tirar proveito do seu entendimento teórico, acerca da técnica usada, para facilitar a análise dos resultados e a elaboração de conclusões mais detalhadas.
- 2 - No tocante ao emprego da Análise de Regressão Linear Múltipla em estudos de demanda de viagens, comprova-se que o processo comumente usado, maior R^2 e menor Se , para a escolha de uma equação de regressão linear múltipla como modelo de previsão, não é suficiente. Deve-se, além disso, analisar a qualidade da equação com base nas relações lineares entre as variáveis mostradas na matriz

de correlação, na coerência quanto à magnitude e sinal dos coeficientes de regressão, no sentido lógico da equação, como um todo, e nos testes estatísticos tanto para a equação como para os coeficientes de regressão. A decisão sobre que equação usar como modelo de previsão, fica a critério do analista.

- 3 - O emprego de técnicas estatísticas multivariadas, como a Análise Fatorial e a Análise Discriminante em análise de transportes, proporcionam um alto grau de detalhamento na análise dos resultados, visto que, permitem, por exemplo, identificar indivíduos de comportamento semelhantes, em relação aos transportes, e agrupá-los segundo um comportamento padrão. Assim, é mais confiável tirar conclusões acerca de grupos, construídos a partir de indivíduos de comportamentos semelhantes, do que se tirar conclusões a respeito de cada indivíduo, partindo-se de grupos não homogêneos, como é feito nos modelos agregados que considera a média da zona como representativa do comportamento de todos os indivíduos da zona. Portanto, os modelos comportamentais ou desagregados para a análise de transportes abordam a complexidade do processo de viagem, pois, a análise das viagens de um indivíduo é feita em termos de uma série de variáveis, representativas do comportamento desse indivíduo, e não de uma só variável como se faz tradicionalmente. Além disso, permite incorporar as variâncias da amostra produzidas por fatores sociais ou locais da unidade de agregação considerada.

4.2 - Sugestões para pesquisas futuras

- 1 - Construir um modelo comportamental de geração-repartição modal integrado, para Campina Grande.
- 2 - Utilizar a Análise Fatorial e/ou Análise Discriminante para analisar os perfis de atividade das pessoas residentes em Campina Grande.
- 3 - Construir e avaliar modelos de escolha modal baseados na Análise de Regressão Linear Múltipla, na Análise Discriminante e na abordagem Multinomial Logit.

BIBLIOGRAFIA

01. Draper N.R. e H. Smith - Applied Regression Analysis.
02. Nie, Norman H. - Statistical Package for the Social Sciences. SPSS - Editora McGraw Hill - 1975.
03. Faissol, Esperidião - Tendências Atuais na Geografia Urbana/Regional - 1978.
04. Cooley, Willian W. e Paul R. Lohnes - Multivariate Procedures for the Behavioral Sciences.
05. Van De Geer - Introduction to Multivariate Analysis for the Social Sciences. Freemar e Company - 1971.
06. Rummel, R. J. - Applied Factor Analysis - 1970.
07. McCarthy, Gerald M. - Multiple - Regression Analysis of Household Trip Generation - A Critique.
08. Harman, H. H. - Modern Factor Analysis - 1967.
09. Kassoff e Deutschman - Trip Generation: A Critical Appraisal
10. Leal, José E. - Notas de aula do Curso de Modelos de Planejamento - 1981 (Não publicada)
11. Doi, Masayuki - The land Use and Zoning Study in Relation to Systematic Transportation Study in Campina Grande - Parte do projeto de planejamento de Transportes - 1978 - (Não publicado).

12. Cruz, Walter S. - A Necessidade de Novas Abordagens à Demanda de Viagens - Anais do Simpósio sobre Modelos Urbanos, Regionais e de Transportes - IPT - São Paulo - Março/1983 - pp 203 - 209.
13. Leal, José E. - Bases para Um Modelo Desagregado de Geração-Repartição - Anais do Simpósio sobre Modelos Urbanos, Regionais e de Transportes - IPT - São Paulo Março/1983.