

MODELO DE DETECCIÓN DE FRAUDES EN SEGUROS DE VEHÍCULOS UTILIZANDO COMPONENTES PRINCIPALES Y ANÁLISIS PRIDIT

Roa Heydi¹, Sandoya Fernando²

Resumen: El presente trabajo tiene como propósito proporcionar una metodología dirigida a cuantificar la probabilidad de fraude en las declaraciones de siniestros vehiculares, denominada Análisis PRIDIT, específicamente cuando se tienen variables cualitativas involucradas. Cabe notar que en el Ecuador aun no se han introducido este tipo de técnicas y mas bien se lo realiza de manera rudimentaria. La metodología PRIDIT es una técnica no-paramétrica simple y fácil de entender e implementar, que puede satisfacer necesidades gerenciales debido a que sus aplicaciones pueden extenderse a clasificaciones mucho más finas que la prueba binaria de detección fraude/no fraude, y provee valor adicional en su capacidad de probar la consistencia del modelo de puntuaciones con los patrones de las variables de entrada. Específicamente, los pesos y las puntuaciones obtenidos de la metodología PRIDIT son representativos de las variables de entrada y pueden ser probados a través de la correlación con otros modelos de puntuaciones existentes.

Palabras Claves: Análisis multivariado, Análisis PRIDIT, Análisis de Componentes Principales.

1. INTRODUCCIÓN

Un término desafortunadamente muy utilizado en nuestra sociedad es el de “fraude”. Entendemos por fraude a cualquier actividad en la que para derivar un beneficio económico, se crean situaciones ficticias o se exageran daños.

El fraude está considerado como una de las industrias criminales más grandes en la sociedad, y según estudios de varios investigadores aumentan en épocas en las cuales la gente necesita dinero tales como la Navidad, Fin de año, etc.

A pesar de que el fraude se da en casi todas las actividades económicas, los fraudes en los seguros se han convertido en una práctica común. El mercado asegurador considera el fraude como un factor ineludible de riesgo y, hoy en día las entidades aseguradoras luchan por desarrollar un frente de acción hacia el mismo.

En este estudio se desarrolla un modelo de detección de fraude a una cartera de seguros de automóviles en base a las declaraciones de siniestros, para clasificar y cuantificar el nivel de fraude de estas declaraciones realizadas por los asegurados.

Se obtienen resultados eficientes si logramos:

- Reducir la incertidumbre e incrementar las oportunidades de clasificar las demandas correcta y eficientemente a cada grupo (fraudulentas / no fraudulentas) sin importar el tipo de variables que intervengan.

- Transformar respuestas categóricas en un conjunto de valores numéricos que estén dentro del intervalo $[-1,1]$, que refleje la relativa anomalía de una respuesta en particular.
- Determinar una ponderación de fraude para cada variable involucrada en el análisis.
- Encontrar una medida de poder discriminatorio que permita clasificar las demandas en fraudulentas y no fraudulentas.

2. CARACTERÍSTICAS GENERALES DE LOS FRAUDES VEHICULARES.

La influencia de las acciones deshonestas de los asegurados se nota tanto en el número de siniestros declarados como en la cuantía de los mismos. Si consideramos el peso que ello puede tener a la hora de justificar la aparición de resultados técnicos negativos durante los últimos años en el seguro de vehículos y el incremento del valor de las primas por la contratación de estos mismos seguros, queda justificada la necesidad de diseñar herramientas que ayuden a las entidades de seguros a detectar y luchar contra el fraude.

En el Ecuador no se han hecho estudios profundos respecto al porcentaje de fraudes, sin embargo en otros países este tipo de estudios están muy adelantados, como por ejemplo: En España, según la Investigación Cooperativa entre Entidades Aseguradoras y Fondos de Pensión (ICEA), el sector del automóvil es el que más fraudes registra, puesto que de los 46.228 casos detectados en el 2001, el 90 por ciento correspondía a esta rama. En lo que respecta al año (2004) según datos de la ICEA, más del 75 por ciento de los casos de fraude detectados corresponde a la rama de vehículos.

¹ Roa Heydi; Ingeniera en Estadística Informática, (e-mail: heydi_roa@yahoo.com)

² Sandoya Fernando, M.Sc., Profesor Agregado de la Escuela Superior Politécnica del Litoral (ESPOL); (e-mail: fernandoss@yahoo.com)

Para el caso de las compañías de seguros de nuestro país les es más práctico pagar los siniestros reclamados por los asegurados que entrar en trámites legales, ya que en el país no se cuenta con una unidad investigativa a la que puedan recurrir las aseguradoras para verificar la honestidad de los reportes de siniestros. Esto podría corregirse si se cuenta con métodos que permitan detectar y clasificar una demanda como fraudulenta o no, con cierta probabilidad de acierto.

3. EL SEGURO VEHICULAR EN EL ECUADOR DURANTE LA ÚLTIMA DÉCADA

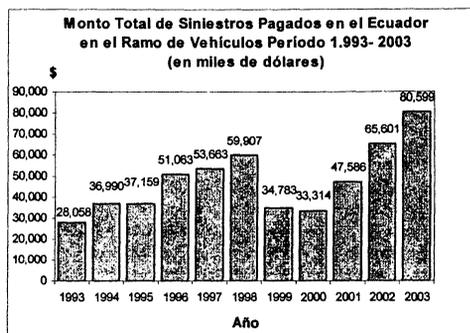
Durante la década (1993 – 2003) en el Ecuador, muchas de las aseguradoras que ofrecen seguros para automóviles, detectaron que el monto de indemnización pagado a sus asegurados por causa de un siniestro era muy alto. El mercado asegurador del Ecuador durante esta última década ha visto en el ramo de seguros de vehículos, el ramo con mayor monto pagado por siniestros en lo que respecta a los demás ramos de seguros.

Según la Superintendencia de Bancos y Seguros del Ecuador los montos totales de siniestros pagados de todas las aseguradoras en el ramo de vehículos se muestran en la Tabla I y la Figura 1.

TABLA. I
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Predit
Monto total de indemnización líquida por siniestros vehiculares 1993 – 2003.

	Indemnización Pagada (en miles de dólares)
Dec-93	28,058
Dec-94	36,990
Dec-95	37,159
Dec-96	51,063
Dec-97	53,663
Dec-98	59,907
Dec-99	34,783
Dec-00	33,314
Dec-01	47,586
Dec-02	65,601
Dec-03	80,599

FIG. 1
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Predit
Montos de indemnización pagadas por siniestros vehicular



En estos datos se puede visualizar la tendencia creciente de los pagos de siniestros a los asegurados en los últimos años, notándose claramente la reducción de los pagos en el año 1999 y 2000 por la crisis económica que había en ese período; sin embargo a partir del 2001 vuelve la tendencia creciente, una vez adoptado el dólar como moneda oficial.

Para reflejar con precisión la evolución de los montos de siniestros pagados por las aseguradoras se puede observar los valores utilizando índices con año base 1993 = 100. La conversión de los datos a índices facilita la estimación de la tendencia en una serie compuesta por números muy grandes como los que se está manejando (miles de dólares). Ver Tabla II, Figura 2

TABLA. II
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Predit
Monto total de indemnización pagada e índices con año base 1993 = 100

	Indemnización Pagada (en miles de dólares)	Ajustado a 100 con Año Base 1.993
Dec-93	28,058	100
Dec-94	36,990	132
Dec-95	37,159	132
Dec-96	51,063	182
Dec-97	53,663	191
Dec-98	59,907	214
Dec-99	34,783	124
Dec-00	33,314	119
Dec-01	47,586	170
Dec-02	65,601	234
Dec-03	80,599	287

FIG 2

Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Pridit
Índices del Monto Total de Sinistros Liquidados con Año Base 1993 = 100



Respecto a las empresas aseguradoras, actualmente en nuestro país, existen 40 compañías dedicadas al mercado asegurador, de las cuales 27 ofrecen pólizas de seguros para vehículos. Esta situación general se refleja también a nivel particular para cada empresa aseguradora. Así según datos proporcionados por una compañía aseguradora del país involucrada con el seguro de vehículos en la última década, los montos pagados por siniestros vehiculares en cada una de las sucursales desde el 2000 hasta el 2003, se muestran en la Tabla III.

La información siguiente coincide con la de la Tabla I, lo que confirma que los montos pagados por siniestros de vehículos aumentan año a año.

TABLA III

Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Pridit
Montos pagados por Siniestros Vehiculares por seguros equinoccial (Periodo 2000 – 2003)

		VALOR PAGADO (en dólares)			
Sucursal	Causas	2000	2001	2002	2003
CUENCA	DAÑO PARCIAL	51,048	58,949	129,634	87,864
	DAÑO TOTAL	20,954	27,404	42,029	81,016
	ROBO PARCIAL	5,189	3,016	6,668	13,443
	ROBO TOTAL	20,580	9,941	8,730	32,899
Total CUENCA		97,770	99,310	187,061	215,223
GUAYAQUIL	DAÑO PARCIAL	109,411	161,996	215,893	539,112
	DAÑO TOTAL	51,080	75,260	119,330	152,875
	ROBO PARCIAL	40,947	23,149	25,525	73,717
	ROBO TOTAL	123,591	67,015	40,452	43,817
Total GUAYAQUIL		325,030	327,420	401,200	809,522
QUITO	DAÑO PARCIAL	615,191	971,582	1,531,313	1,997,830
	DAÑO TOTAL	163,257	425,845	1,105,290	1,093,183
	ROBO PARCIAL	180,119	247,309	292,380	422,735
	ROBO TOTAL	258,075	312,623	333,322	244,809
Total QUITO		1,216,643	1,957,359	3,262,305	3,758,557
TOTAL		1,639,443	2,384,088	3,850,565	4,783,301

4. LA METODOLOGÍA PRIDIT

La metodología PRIDIT constituye una nueva técnica no-paramétrica simple y fácil de implementar, difiere de otras pruebas estadísticas no-paramétricas, como la Chi Cuadrado, que asume un orden natural de los datos, en que esta técnica estadística puede emplearse con variables categóricas que pueden estar en escalas subjetivas (por ejemplo: severo, moderado, menor) que pueden tomar una forma numérica en la cual el sistema de medición confía fuertemente en el método experimental o en la habilidad técnica del

científico involucradas en la medición de la variable en cuestión.

4.1. TIPO DE DATOS

La mayoría de los métodos estadísticos utilizados para este tipo de estudios requieren que los datos sean de tipo intervalo y en algunos casos requieren que los datos sean continuos y estén normalmente distribuidos. Pero la metodología PRIDIT no pone ninguna restricción en lo que respecta a los datos, esta metodología trabaja con variables ordinales discretas o con variables

categorías. Para el caso de seguros de automóviles la mayoría de las variables son dicotómicas. Para el desarrollo de esta investigación se obtuvo información de 100

demandas de siniestros vehiculares, y se determinó como variables importantes para la calificación de la demandas las que se muestran en la Tabla IV.

TABLA IV
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Prídít
Montos pagados por Siniestros Vehiculares por seguros equinoccial (Período 2000 – 2003)

COBERT	El asegurado tiene cobertura de daños a terceros.
FRAQCIA	El asegurado tiene franquicia en la póliza.
ACCESOR	El asegurado tiene cobertura de accesorios.
Características del Vehículo	
VEHUSO	El vehículo es de uso privado del asegurado.
Características del Accidente	
ACULPA	El asegurado acepta la culpa del siniestro.
ZNURB	El siniestro ha ocurrido en una zona no urbana.
ACCNOCHE	El siniestro ocurrió en la noche.
ACCFINS	El siniestro ocurrió en el fin de semana.
TESTIGOS	Existen testigos.
REPPOLIC	Existe un reporte policial.
ZONA1	El siniestro ocurrió en una zona de elevada siniestralidad.
ZONA3	El siniestro ocurrió en una zona de baja siniestralidad.
REPSOSP	Existe presencia de relatos sospechosos.
PARENTEZ	Existe coincidencia de apellidos entre las partes.
RETRASO	El siniestro no fue reportado a la compañía aseguradora dentro del período establecido.

4.2. TRANSFORMACIÓN DE LOS DATOS CATEGÓRICOS EN PUNTUACIONES RIDIT

Un método de puntuación de mejor desempeño para variables categóricas, es el método de Puntuaciones RIDIT. El objetivo de este método de puntuación es cuantificar el nivel de sospecha de fraude producida por la representación categórica de cada variable indicadora de fraude en una declaración de siniestro.

Sea k_t el número de categorías disponibles y ordenadas para la variable indicadora t , y sea $\hat{p}_t = (\hat{p}_{t1}, \dots, \hat{p}_{tk_t})$ el vector de las proporciones de respuesta observadas en el conjunto entero de demandas. Se supone que las categorías están ordenadas en relación a la probabilidad de sospecha de fraude, en forma decreciente; por lo

que una respuesta categórica más alta indica menos sospecha de fraude. Para la opción categórica i de la variable t , se asigna el valor numérico o puntuación siguiente:

$$B_{ti} = \sum_{j < i} \hat{p}_{tj} - \sum_{j > i} \hat{p}_{tj} \quad i = 1, 2, \dots, k_t$$

Así, este procedimiento transforma cualquier conjunto de respuestas categóricas en un conjunto de valores numéricos dentro de un intervalo $[-1, 1]$, lo cual refleja la relativa "anormalidad" de una respuesta particular.

En la Tabla V se muestran las puntuaciones para valores simulados de las respuestas.

TABLA V
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Pridit
Cálculo de puntuaciones Ridit

Variables	% Si	B _{t1} ("si")	B _{t2} ("no")
El asegurado tiene cobertura de daños a terceros	96%	-0.04	0.96
El asegurado tiene franquicia en la póliza	90%	-0.10	0.90
El asegurado tiene cobertura de accesorios	93%	-0.07	0.93
El vehículo es de uso privado del asegurado	55%	-0.45	0.55
El asegurado acepta su culpa	6%	-0.94	0.06
Accidente ocurrido en zona no urbana	10%	-0.90	0.10
Accidente ocurrido en la noche	57%	-0.43	0.57
Accidente ocurrido durante un fin de semana	38%	-0.62	0.38
Existen testigos	12%	-0.88	0.12
Existe un reporte policial	95%	-0.05	0.95
Siniestro ocurrido en zona de elevada siniestralidad	81%	-0.19	0.81
Siniestro ocurrido en zona de baja siniestralidad	21%	-0.79	0.21
Existe presencia de relatos sospechosos	3%	-0.97	0.03
Existe coincidencia de apellidos entre las partes	1%	-0.99	0.01
El siniestro no fue reportado a la compañía aseguradora dentro del periodo establecido.	14%	-0.86	0.14

4.3 OBTENCIÓN DE PONDERACIONES PRIDIT Y PUNTUACIÓN GLOBAL PARA UNA DEMANDA ENTERA

Sea $F = (f_{it})$ la notación de la matriz de las puntuaciones individuales RIDIT para cada una de las $t = 1, 2, \dots, m$ variables y para cada una de las $i = 1, 2, \dots, N$ demandas. Esto es, $f_{it} = B_{ik}$ si la demanda i contiene k niveles de respuestas categóricas para la variable t . Se obtiene una puntuación de sospecha global para cada demanda simplemente sumando las puntuaciones de las variables individuales respectivas.

Sea $W^{(0)} = (1, 1, \dots, 1)$, entonces el vector de la suma global de las puntuaciones de sospecha de fraude obtenidas para cada demanda denotada en la matriz es: $S^{(0)} = FW^{(0)}$. Ahora, se tiene una medida de consistencia de la variable indicadora t con las puntuaciones globales de sospecha de fraude para las demandas.

A través del siguiente Teorema, se garantiza que el valor de W converge y la fijación del peso de la variable $\hat{W}^{(\infty)}$ es la primera componente principal de $F'F$.

Teorema: Las sucesiones de los pesos de las variables predictoras $\{W^{(n)}\}$ y la sumatoria global de las puntuaciones de sospecha de las demandas $\{S^{(n)}\}$ convergen. De modo que, la fijación de peso de la variable predictorora $\hat{W}^{(\infty)}$ es la primera componente principal de $F'F$, la cual es una estimación consistente de la primera componente principal $W^{(\infty)}$ de $E(F'F)$, la t -ésima componente principal se puede obtener explícitamente como:

$$W_t^{(\infty)} = \frac{A_t}{(\mu_1 - U_{tt}) \sqrt{\sum_{s=1}^m A_s^2 / (\mu_1 - U_{ss})^2}}$$

En la tabla V se observan los componentes principales para los valores obtenidos en la Tabla IV

TABLA VI
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Prídít
 Ponderaciones Prídít para las variables indicadoras de fraude

Matriz de Componentes

	Componente1
	W
COBERT	.468
FRAQCIA	-.460
ACCESOR	.354
VEHUSO	.507
ACULPA	-.560
ZNURB	.397
ACCNOCHE	-.105
ACCFINS	-.252
TESTIGOS	.177
REPPOLIC	-.341
ZONA1	.468
ZONA3	-.328
REPSOSP	.317
PARENTEZ	-.104
RETRASO	-.362

La Tabla VI muestra que las variables indicadoras de fraude **VEHUSO** y **ACULPA** llevan la mayor ponderación PRIDIT de las 100 demandas analizadas, seguidas de **COBERT**, **ZONA1** y **FRAQCIA**, es decir que éstas son las variables que tienen mayor peso a la hora de determinar un nivel de sospecha de fraude.

TABLA VII
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Prídít
 Puntuación Global para una demanda completa

	S
COBERT	0.808
FRAQCIA	-0.797
ACCESOR	-1
VEHUSO	0.856
ACULPA	-0.913
ZNURB	0.709
ACCNOCHE	-0.203
ACCFINS	-0.474
TESTIGOS	0.339
REPPOLIC	-0.624
ZONA1	0.868
ZONA3	-0.603
REPSOSP	0.585
PARENTEZ	-0.201
RETRASO	-0.657

Los valores A_t de la Tabla VII son una medida del poder discriminatorio de las t variables indicadoras de fraude para una demanda entera.

4.4 CLASIFICACIÓN DE LAS DEMANDAS POR MEDIO DE LAS PONDERACIONES PRIDIT

Con respecto a la clasificación, se consideran dos casos para la proporción de demandas del grupo 1, Cuando θ es conocido y cuando θ es desconocido.

Cuando θ es conocido, se ordenan las N demandas por medio de sus puntuaciones unidimensionales $S = \sum_{t=1}^m W_t^{(\theta)} X_t$, y luego se clasifican las primeras $N\theta$ demandas en el nivel alto del grupo 1 de sospecha de fraude. Aquí X_t es la puntuación calculada de la demanda obtenida para la variable t , $X_t = \sum_{i=1}^k B_{it} I_A$ donde I_A es el indicador del conjunto A .

Si θ es desconocido (como es este caso), se separan los dos grupos de acuerdo a las puntuaciones globales positivas o negativas y se clasifican las demandas dentro del grupo de bajo nivel de sospecha de fraude.

Si la puntuación global del nivel de sospecha de fraude es positivo, se ubican en el grupo de las demandas no fraudulentas y si la puntuación de valor negativo se ubican dentro del grupo de demandas fraudulentas.

En la TABLA VII se demuestra que cada una de las demandas tiene una puntuación que permite clasificar cada demanda dentro del grupo correspondiente.

Una puntuación positiva significa que la demanda pertenece al Grupo 2, es decir es una demanda no fraudulenta mientras que un score negativo significa que la demanda debe ser discriminada al grupo 1, que es el grupo de las demandas con una alta probabilidad de ser consideradas fraudulentas..

TABLA VII
Modelo de detección de fraudes en seguros de vehículos utilizando componentes principales y análisis Pridit
Clasificación de las demandas

Demanda	cobertur	fracqia	accesor	usovehi	aculpa	znurb	accnoche	accfins	testigos	reppolic	zona1	zona3	repsosp	parentez	retraso	Score	Clase
1	-0.04	-0.1	-0.07	0.55	0.06	0.1	-0.43	0.38	-0.88	-0.05	-0.19	-0.79	0.03	0.01	0.14	0.10	2
2	-0.04	-0.1	-0.07	-0.45	0.06	0.1	0.57	0.38	0.12	0.95	0.81	0.21	0.03	0.01	-0.86	-0.54	1
3	-0.04	-0.1	-0.07	-0.45	0.06	0.1	0.57	-0.62	-0.88	-0.05	-0.19	-0.79	0.03	0.01	0.14	0.39	2
4	-0.04	-0.1	-0.07	0.55	0.06	0.1	-0.43	0.38	0.12	-0.05	-0.19	0.21	0.03	0.01	0.14	0.09	2
5	-0.04	-0.1	-0.07	0.55	0.06	0.1	-0.43	0.38	0.12	-0.05	-0.19	0.21	-0.97	0.01	0.14	0.45	2
6	-0.04	0.9	-0.07	-0.45	0.06	0.1	-0.43	0.38	-0.88	-0.05	-0.19	0.21	0.03	0.01	0.14	0.06	2
7	-0.04	-0.1	-0.07	0.55	0.06	0.1	-0.43	0.38	0.12	-0.05	-0.19	0.21	0.03	0.01	0.14	-0.29	1
8	-0.04	-0.1	-0.07	-0.45	0.06	0.1	-0.43	0.38	0.12	-0.05	-0.19	0.21	0.03	0.01	0.14	0.33	2
9	-0.04	-0.1	0.93	0.55	-0.94	0.1	-0.43	-0.62	0.12	-0.05	0.81	-0.79	0.03	0.01	0.14	-0.54	1
10	-0.04	-0.1	-0.07	-0.45	0.06	0.1	-0.43	0.38	0.12	-0.05	-0.19	0.21	0.03	0.01	0.14	0.48	2

5. CONCLUSIONES

1. La utilización de la técnica PRIDIT en la detección de fraudes dentro del campo de seguros de vehículos, es más eficiente que las técnicas estadísticas tradicionales cuando dentro de las variables involucradas existen variables con respuestas del tipo categóricas, puesto que la técnica PRIDIT transforma el conjunto de respuestas categóricas en un conjunto de valores numéricos dentro de un intervalo [-1,1], lo cual refleja la relativa anormalidad de una respuesta en particular.
2. La técnica PRIDIT provee una medida que permite determinar qué variables son más consistentes, dando ponderaciones más altas a las variables indicadoras de fraude. En esta investigación las variables con mayor ponderación son: ACULPA, VEHUSO, COBERT, ZONA1, FRAQCIA, que se observan en la tabla IV. Estas variables son las que indican mayormente fraude.
3. Otra de las ventajas de la utilización de la metodología PRIDIT que se comprobó en esta investigación, es que provee una medida del poder discriminatorio de las variables indicadoras de fraude. Las variables con mayor poder discriminatorio en esta investigación fueron: ACCESOR, ACULPA, ZONA1, VEHUSO, COBERT, FRAQCIA, las cuales se observan e la tabla IV.
4. La medida cuantitativa del poder discriminatorio que resulta de la técnica PRIDIT provee además la capacidad de determinar correlaciones con otras medidas cuantitativas tales como: edad del conductor, número de accidentes anteriores del asegurado, número de años que el asegurado tiene licencia, etcétera. Este procedimiento puede dejarse para un estudio posterior puesto que no es el objetivo de la actual investigación. Además que no se cuenta con datos reales que permitan determinar la verdadera correlación entre las variables.

REFERENCIAS BIBLIOGRÁFICAS Y ELECTRÓNICAS

- [1] **ROA, H.** , (2004) "*Modelo de Detección de Fraudes en los Seguros de Vehículos utilizando Componentes Principales y Análisis RIDIT*". (Tesis Grado, Instituto de Ciencias Matemáticas, Escuela Superior Politécnica del Litoral.
- [2] **JOHNSON, D.** , (2000) "*Métodos Multivariados aplicados al análisis de datos*", Internacional Thompson Editores, México DF - México.
- [3] **FERRAN, A.** , (2001) "*SPSS para Windows: Análisis Estadístico*", Ed. McGraw-Will, Madrid - España.
- [4] **BROCKETT, P., DERRIG, R., GOLDEN, L., LEVINE, A., MARK, A., THE JOURNAL OF RISK AND INSURANCE,** (2002) "*Fraud Classification Using Principal Component Analysis of RIDIT's,*" Vol.69, No. 3.
- [5] **SUPERINTENDENCIA DE BANCOS Y SEGUROS DEL ECUADOR,** (2002) "*Monto total de indemnización liquidada por siniestros vehiculares, periodo (1993 – 2003)*", http://www.superban.gov.ec/pages/seguros_privados.htm
Ultima visita: Octubre del 2004, Guayaquil - Ecuador