

# Taller Minería de datos aplicados a la educación

*2ª parte*

*Presentación del software*

*PASW Modeler*

*27 de junio de 2011*

*Mercedes Torrado*

*Departamento Métodos de Investigación y  
Diagnóstico en Educación (MIDE)*

Este trabajo cuenta con licencia de Creative Commons:

Minería de datos aplicados a la educación: software PASW MODELER está sujeta a una  
licencia de Reconocimiento-NoComercial-SinObraDerivada 3.0 (CC BY-NC-ND 3.0)

Para citar la obra:

Torrado, M. (2011) *Minería de datos aplicados a la educación: software PASW MODELER*.  
Barcelona: Universidad de Barcelona. Deposito Digital <http://hdl.handle.net/2445/19862>

# *Presentación del software PASW Modeler*

- 1.- ¿Qué nos permite hacer este software?
- 2.- ¿Con qué terminología trabaja?
- 3.- ¿Qué preguntas nos debemos hacer?
- 4.- ¿Qué proceso debemos seguir?

# 1 ¿Qué nos permite hacer este software?

- **Definición del PASW Modeler**

*Es un conjunto de herramientas de minería de datos que permite desarrollar rápidamente **modelos** predictivos y utilizarlos para mejorar la toma de decisiones*

*Permite extraer información de grandes bases de datos*



- **Definición de MODELO**

*Es un conjunto de reglas, fórmulas o ecuaciones que puede utilizarse para pronosticar un resultado basándose en un conjunto de campos o variables de entrada*



# 1 ¿Qué nos permite hacer este software?

- **Objetivo del programa**

El objetivo principal del *PASW Modeler* es la capacidad de pronosticar un resultado a partir del análisis predictivo y la comprensión del proceso



# 1 ¿Qué nos permite hacer este software?

- **Por ejemplo**

*Identificar alumnos universitarios con mayor probabilidad de no persistir en la universidad en base a una serie de características personales y psicoeducativas*



## 2 ¿Con qué terminología trabaja?

- Terminología común • Terminología común
  - ESCALA DE MEDIDA
  - TÉRMINO DE PAPELES

Spss-win	PASW Modeler
Escala	Rango
N/a	Discreto
Nominal	Conjunto
Ordinal	Conjunto ordenado
N/a	Sin tipo
N/a	Por defecto

*Variable independiente o predictora*

**“entrada”**

*Variable dependiente*

**“salida” “objetivo”**

*Eliminación o anulación*

**“Ninguno”**



## 2 ¿Con qué terminología trabaja?

- **Terminología común**
  - Término de atributo, campo o variable
  - Término de registro, ejemplo o caso
  - “ruido” datos con errores





### 3.- ¿Qué preguntas nos debemos hacer?

- ¿Cuál es el problema?
- ¿Dónde tenemos los datos?
- ¿Qué datos son los más importantes?
- ¿Qué ruidos existen?
- ¿Qué técnica de modelado es la más adecuada?
- ¿El modelo es adecuado?

4.- ¿Qué proceso debemos seguir?



# TIPOLOGÍA DE MODELADOS

- *Modelos de clasificación*

*A partir del valor de uno o más campos de entrada (VI) predicen el valor de uno o más variables de salida (VD)*

- *Modelos de asociación*

*Permiten pronosticar varios resultados. Encuentran patrones de asociación entre datos*

- *Modelos de segmentación*

*Dividen los datos en segmentos o conglomerados con características similares y se desconoce el resultado*

- *Modelos de cribados*

*Permiten identificar anomalías / casos extraños o valores atípicos*



## 4.- ¿Qué proceso debemos seguir?

- **Creación de un modelo** ▶
  - Nodo de origen
  - Nodo de tipo
  - Nodo de modelado
  - Nodo de tabla y análisis
- **Exploración de un modelo** ▶
  - Nugget del modelo (examinar) visor
- **Evaluación de un modelo** ▶
  - Análisis

Glosario básico  
Ejemplo práctico



# Glosario básico

- Lienzo de ruta
- Administrador de rutas
- Proyectos
- Paleta de nodos
- Icono o nodo
- Rutas
- Flujo de datos
- Nuggets

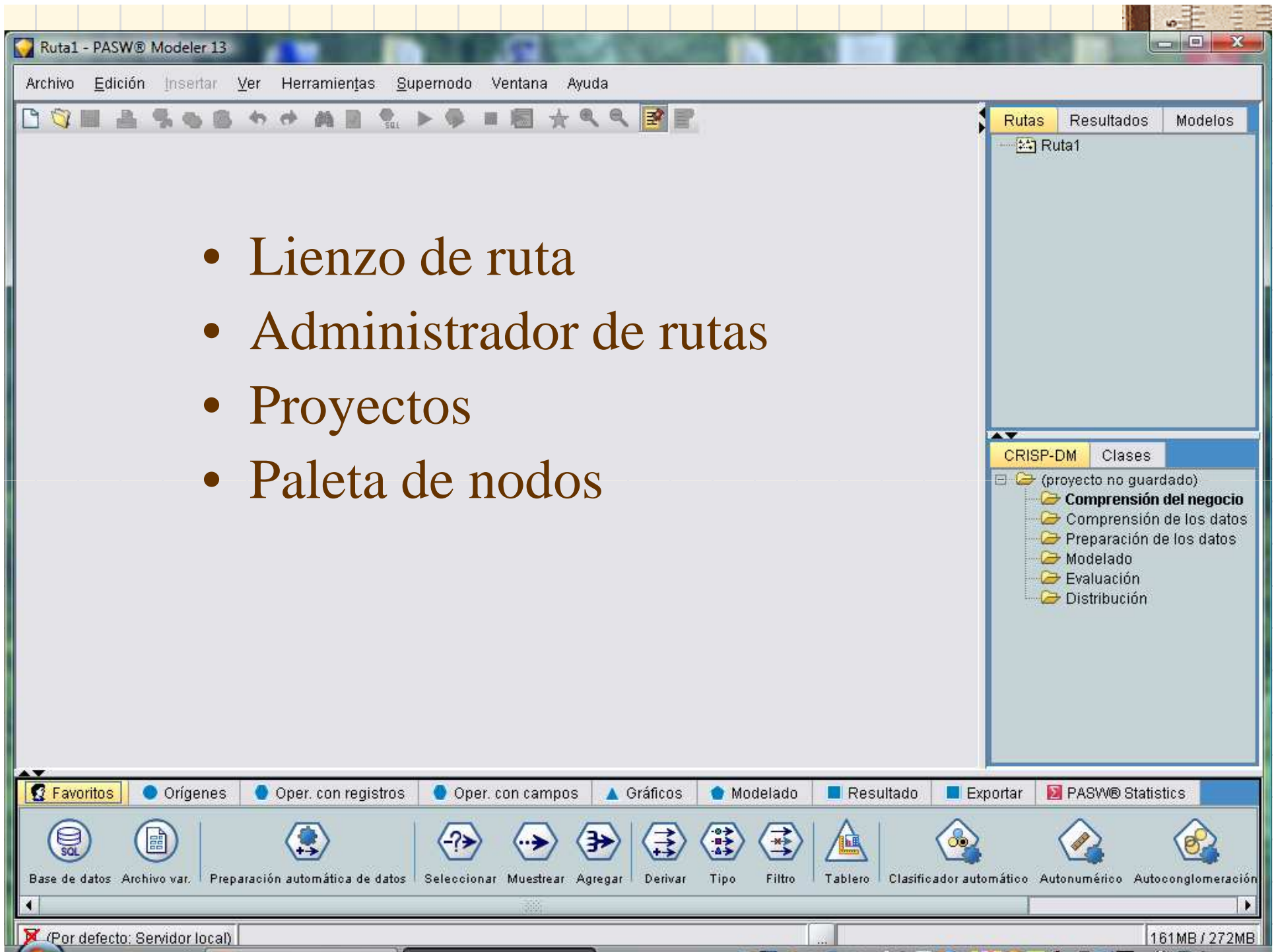


# Creación de un modelo

Para crear un modelo se necesitan mínimo tres elementos

- UN NODO DE ORIGEN que lea los datos
- UN NODO DE ORIGEN O NODO TIPO que especifique propiedades del campo, tipo de datos, el papel de cada campo (origen o predictor en el modelado)
- UN NODO DE MODELADO que genere un nugget de modelado cuando se ejecute la ruta
- UN NODO DE TABLA /ANÁLISIS para ver los resultados de puntuación después de crear el nugget de modelo y añadirlo a la ruta





Regresión\_conglomerados\* - PASW® Modeler 13

Archivo Edición Insertar Ver Herramientas Supernodo Ventana Ayuda

68 Campos Auditar datos de tabla... Satisfacción\_académi..

matriz\_modeler.sav Tipo Satisfacción\_académi..

NotaPAU v. Edad Autoconglomeración K-medias K-medias

- Icono o nodo
- Rutas
- Flujo de datos
- Nuggets

Rutas Resultados Modelos

Satisfacción...

K-medias

P-DM Clases

(proyecto no guardado)

- Comprensión del negocio
- Comprensión de los datos
- Preparación de los datos
- Modelado
- Evaluación
- Distribución

Favoritos Orígenes Oper. con registros Oper. con campos Gráficos Modelado Resultado Exportar PASW® Statistics

Base de datos Archivo var. Preparación automática de datos Seleccionar Muestrear Agregar Derivar Tipo Filtro Tablero Clasificador automático Autonumérico Autoconglomeración Tabla Archivo plano Base de datos

Servidor: Servidor local 173MB / 286MB

ES 21:19 26/06/2011

Presentación preliminar

Leer valores | Borrar valores | Borrar todos los valores

Campo	Tipo	Valores	Perdidos	Comprobar	Dirección
EDUCATE	Conjunto ordenado	5,6,7,8,9,10,11,1...		Ninguna	Entrada
GENDER	Conjunto	0,1		Ninguna	Entrada
AGE	Rango	[18,89]		Ninguna	Entrada
TVDAY	Rango	[0,12]		Ninguna	Entrada
ORGS	Conjunto ordenado	0,1,2,3,4,5,6,7,8		Ninguna	Entrada
CHILDS	Conjunto ordenado	0,1,2,3,4,5,6,7,8		Ninguna	Entrada
INC	Conjunto ordenado	1,2,3,4,5,6		Ninguna	Entrada
NEWSCHAN	Marca	1/0		Ninguna	Salida

Ver campos actuales | Ver configuración de campos no utilizados

Tipos | Formato | Anotaciones

Aceptar | Cancelar

Presentación preliminar desde nodo Type (8 campos, 10 registros)

	EDUCATE	GENDER	AGE	TVDAY	ORGS	CHILDS	INC	NEWSCHAN
1	20	0	35	1	0	1	4	1
2	12	1	25	5	0	0	1	0
3	14	1	64	2	1	2	5	1
4	9	0	72	2	2	0	3	1
5	12	1	67	4	0	5	\$n...	1
6	15	0	33	2	0	0	6	1
7	14	0	23	4	0	1	3	0

Tabla | Anotaciones

Aceptar



Ruta1\* - PASW® Modeler 13

Archivo Edición Insertar Ver Herramientas Supernodo Ventana Ayuda

matriz\_modeler.sav → Tipo → NotaPAU v. Edad

22 Campos

Auditar datos de [22 campos]

Campo	Gráfico	Tipo	Mín	Máx	Media	Desv. típica	Asimetría	Únicos	Válidos
NotaPAU		Rango	5.070	9.162	6.564	0.750	0.588	--	
Nota_Admissió		Rango	5.000	12.930	8.527	1.233	0.400	--	
Ordre_Eleccio		Conjunto	1.000	7.000	--	--	--	7	
rordre_eleccio		Rango	1.000	4.000	1.352	0.843	2.337	--	
rtipo_residencia		Rango	1.000	5.000	4.618	0.983	-2.650	--	
rvia		Rango	1.000	5.000	1.781	1.199	1.110	--	
rerevia		Rango	1.000	5.000	1.781	1.199	1.110	--	
titulación_M1		Conjunto	1.000	2.000	--	--	--	2	
Edad		Rango	17.000	57.000	19.860	3.778	4.123	--	

\* Indica un resultado de varios modos    \* Indica un resultado muestreado

Auditar Calidad Anotaciones

Servidor: Servidor local

167MB / 293MB

ES 20:37 26/06/2011



# Explorar el modelo

Cuando se finaliza la ejecución del **modelo el nugget** se añade a la paleta de modelos. Para ver los detalles se pulsa el botón derecho del ratón y se selecciona “examinar”. La ficha muestra los detalles del modelo (variables, valores, etc.) La ficha visor muestra también un gráfico de la importancia de la variable en la estimación del modelo.



Ruta1\* - PASW® Modeler 13

Archivo Edición Insertar Ver Herramientas Supernodo Ventana Ayuda

The main window displays a workflow diagram with the following components and connections:

- matriz\_modeler.sav** (File icon) → **Tipo** (Hexagon icon)
- 68 Campos** (Table icon) → **Tipo**
- Tipo** → **Auditar datos titula..** (Table icon)
- Tipo** → **Satisfacción\_académi..** (Table icon)
- Tipo** → **Satisfacción\_académi..** (Table icon)
- Tipo** → **NotaPAU v. Edad** (Line graph icon)

The right-hand window, titled "Satisfacción\_académica", displays a bar chart titled "Importancia de variable". The chart shows the relative importance of various variables. The x-axis ranges from 0.0 to 0.4. The y-axis lists the variables. A legend at the bottom identifies the variables: "NotaPAU" (light blue) and "Mi motivación ACTUAL por los estudios es ..." (dark blue).

Variable	Importancia (aproximada)
Mi motivación ACTUAL por los estudios es ...	0.35
Expectativas de autoeficacia (EA)	0.20
Apoyo académico (AA)	0.18
Apoyo social y familiar (ASF)	0.12
Mi adaptación academica INICIAL fue ...	0.11
Nota_Admissió	0.05
Cuantas veces has percibido dificultades	0.01
Edad	0.01
NotaPAU	0.01

Ver: Importancia de variable

Modelo Resumen Avanzado Anotaciones

Aceptar

Servidor: Servidor local

211MB / 355MB

ES 20:58 26/06/2011

K-medias

Archivo Generar Ver

Conglomerado	conglomerado-1	conglomerado-2	conglomerado-4	conglomerado-3	conglomerado-5
Etiqueta					
Descripción					
Tamaño	24,2% (200)	23,2% (192)	23,1% (191)	18,9% (156)	10,6% (88)
Características	Apoyo académico (AA) -0,07	Apoyo académico (AA) 0,33	Apoyo académico (AA) 0,14	Apoyo académico (AA) -0,73	Apoyo académico (AA) 0,44
	Apoyo social y familiar (ASF) 0,25	Apoyo social y familiar (ASF) 0,26	Apoyo social y familiar (ASF) 0,00	Apoyo social y familiar (ASF) -0,83	Apoyo social y familiar (ASF) 0,36
	Confianza: combinar estudios con otras actividades 3,41	Confianza: combinar estudios con otras actividades 4,22	Confianza: combinar estudios con otras actividades 3,99	Confianza: combinar estudios con otras actividades 3,24	Confianza: combinar estudios con otras actividades 3,91
	Confianza: completar todos los requisitos académicos 3,63	Confianza: completar todos los requisitos académicos 4,32	Confianza: completar todos los requisitos académicos 3,88	Confianza: completar todos los requisitos académicos 3,25	Confianza: completar todos los requisitos académicos 3,98
	Confianza: comprender los contenidos de las asig....	Confianza: comprender los contenidos de las asig....	Confianza: comprender los contenidos de las asig....	Confianza: comprender los contenidos de las asig....	Confianza: comprender los contenidos de las asig....

### Tamaños de conglomerados

Legend: Conglomerada (light blue), conglomerera (red), conglomerera (dark blue), conglomerera (green), conglomerera (yellow)

Tamaño de conglomerado más pequeño	88 (10,6%)
Tamaño de conglomerado más grande	200 (24,2%)
Cociente de tamaños: Conglomerado mayor / Conglomerado menor	2,27

Ver: Conglomerados Casillas: Centros de conglomerados Representación... Restablecer

Clasificar características por: Importancia global Clasificar conglomerados por: Tamaño

Ver: Tamaños de conglomerados

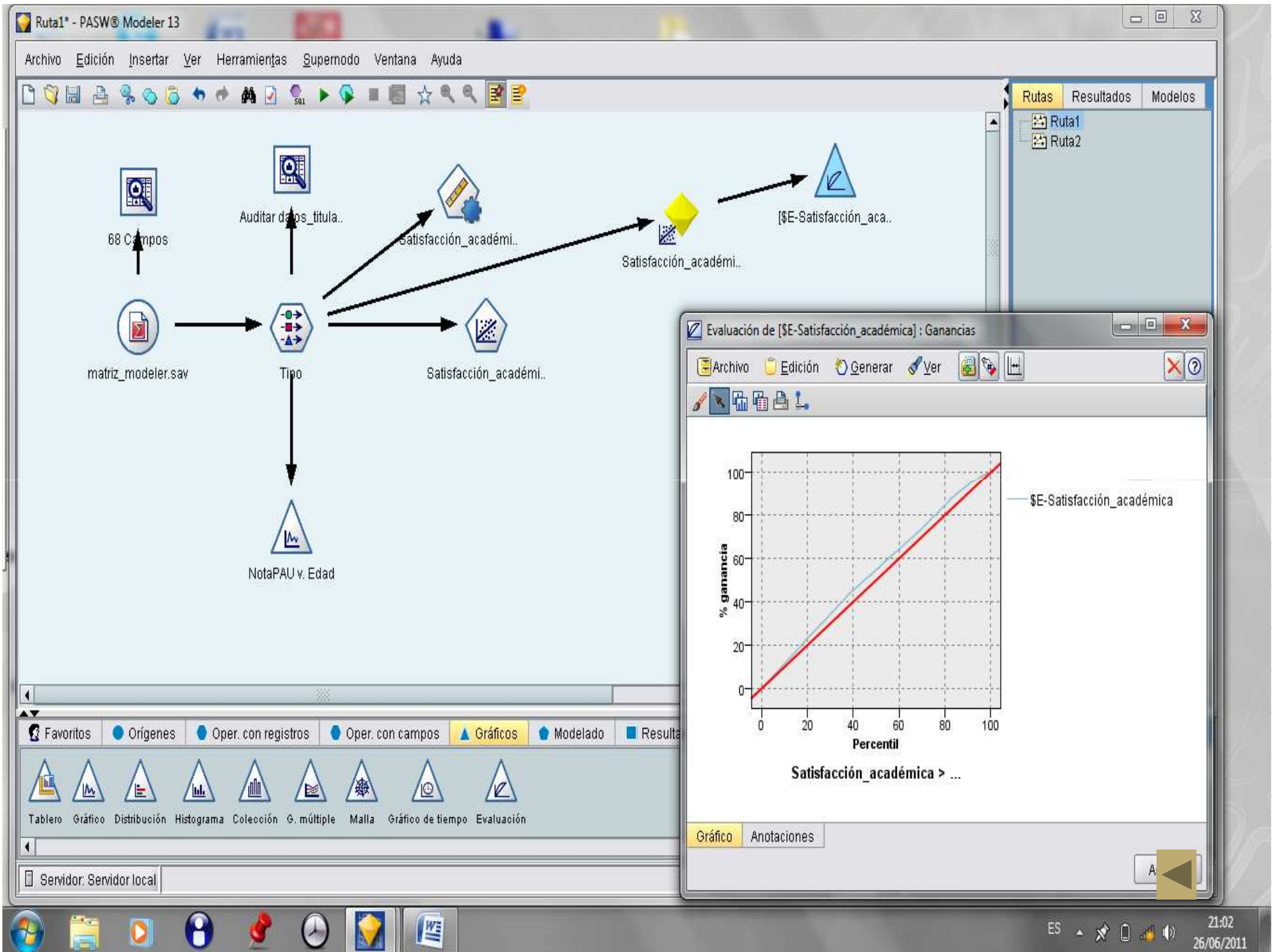
Modelo Resumen Anotaciones

21:09 26/06/2011

# Evaluación del modelo

Para evaluar con que precisión trabaja en modelo, debe puntuar varios registros y comparar las respuestas pronosticadas por el modelo con los resultados reales. Se debe pulsar “añadir ruta”. Se adjunta el modelo al nodo Tipo. Para ver las puntuaciones o pronósticos se adjunta el nodo Tabla y se ejecuta. Se añadirá al final de la tabla el valor pronosticado. Para descubrir cuántos pronósticos son correctos puede leer la tabla o puede adjuntar el nodo análisis que lo hace automáticamente.





# Ejemplo práctico

Investigación sobre *“La persistencia y el abandono en el primer año de universidad en ciencias sociales: bases para la mejora de la retención”* (Ref: EDU2009-10351) que se lleva a cabo por el grupo de investigación TRALS de la Facultad de Pedagogía de la UB



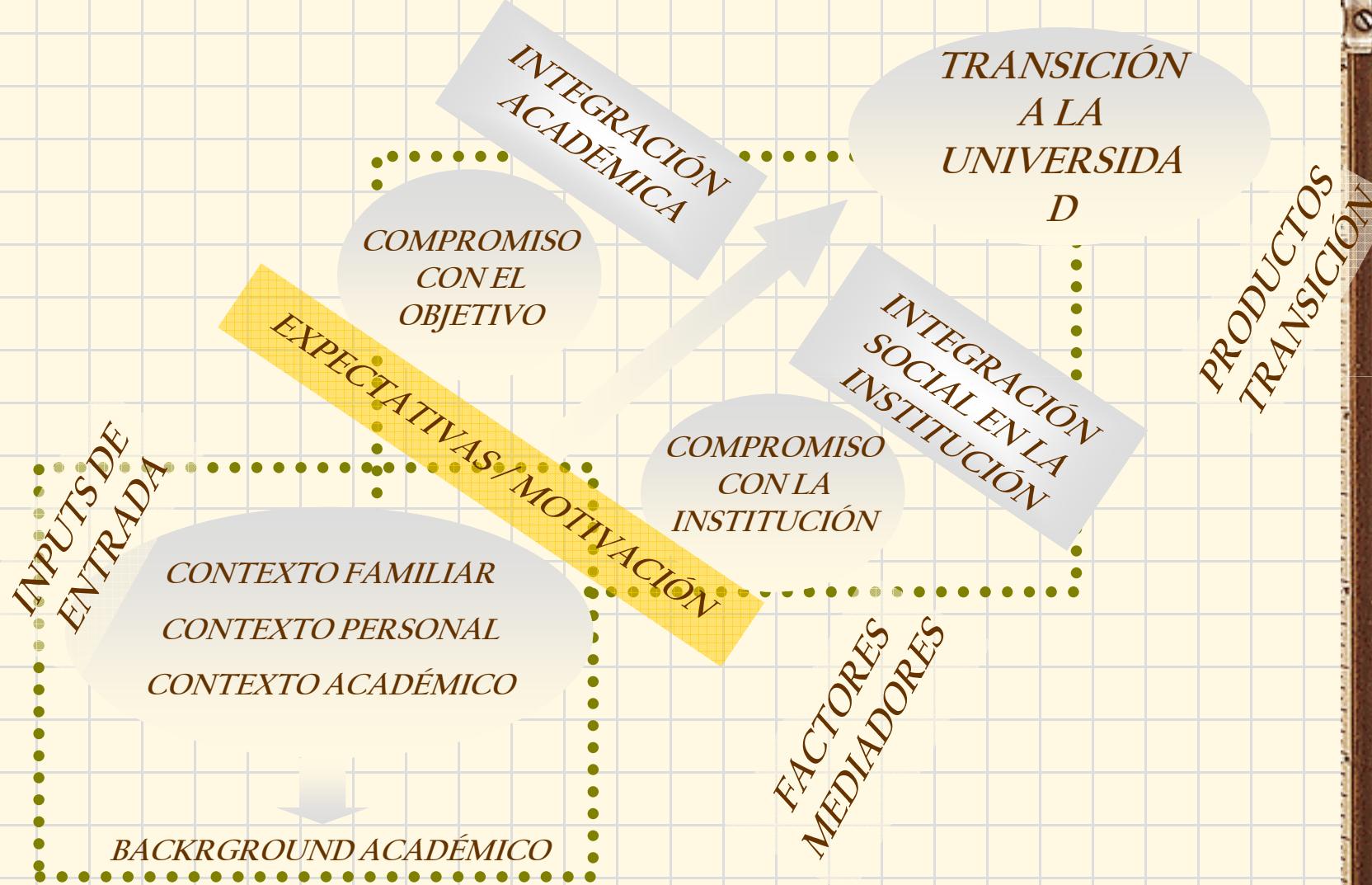


# Objetivo de la investigación

- Analizar los factores personales e institucionales de persistencia académica al finalizar el primer año de universidad.
- Plantea un proceso de investigación longitudinal de carácter descriptivo y comprensivo, basado en la metodología multiestratégica y sistémica en dos niveles de análisis: Macro (conjunto de la promoción) y Micro (análisis de tipologías).
- La muestra está constituida por el total de una cohorte de estudiantes de la promoción 2010 de las titulaciones de Pedagogía y Administración de Empresas.



# MODELO DE TRANSICIÓN DE LA INVESTIGACIÓN



# 1ª fase. Objetivo

- Analizar el proceso de integración inicial y la percepción de los estudiantes en relación a la primera etapa de transición académica (primer semestre de curso).
- Se parte del modelo de integración de Robert W. Lent, que analiza los factores sociocognitivos relacionados con la adaptación inicial a los estudios.

# Qué se quiere demostrar

1.- Una persona se siente satisfecha en el ámbito académico de la universidad cuando se siente competente en lo que hace (expectativas de autoeficacia), sabe que va a obtener beneficios haciéndolo (expectativas de resultados), puede implicarse y participar en el desarrollo de la tarea (actividad hacia la meta) y cuenta con recursos y apoyos necesarios que le ayudaran a alcanzar sus objetivos (soporte).



# Qué se quiere demostrar

2.- Los resultados sobre el proceso de integración inicial es diferente según el contexto académico.



# UNA APROXIMACIÓN AL PROGRAMA DE MINERÍA DE DATOS

*PASW Modeler*



# Modelo de clasificación/relación

- **Árboles de decisión QUEST**

Método de clasificación **binario** para generar árboles de decisión.

VI – continua    VD categórica

- **Árbol de decisión CHAID**

Método de clasificación utilizando estadísticos del Chi-cuadrado para identificar las divisiones óptimas.

VI – continua/categórica    VD continua/categórica



# Modelo de clasificación/relación

- **Regresión Lineal**

Permite pronosticar valores en las VD a partir de la VI

Ajusta en una superficie o línea recta las discrepancias existentes entre los valores de salida reales y los pronosticados

VI – continua    VD continua

- **Regresión logística**

Es similar a la regresión lineal pero parte de valores categóricos en la VD





# Modelo de clasificación/relación

- **Análisis discriminante**

Permite identificar el peso de discriminación de las variables independientes

VI – continua    VD categórica

- **Análisis factorial**

Permite reducir datos para reducir la complejidad de los mismos. Busca combinaciones lineales



# Modelo de cribado

- **Detección de anomalías**

Identifica valores extraños o atípicos que no se ajustan a los patrones de los demás datos



# Modelo de segmentación

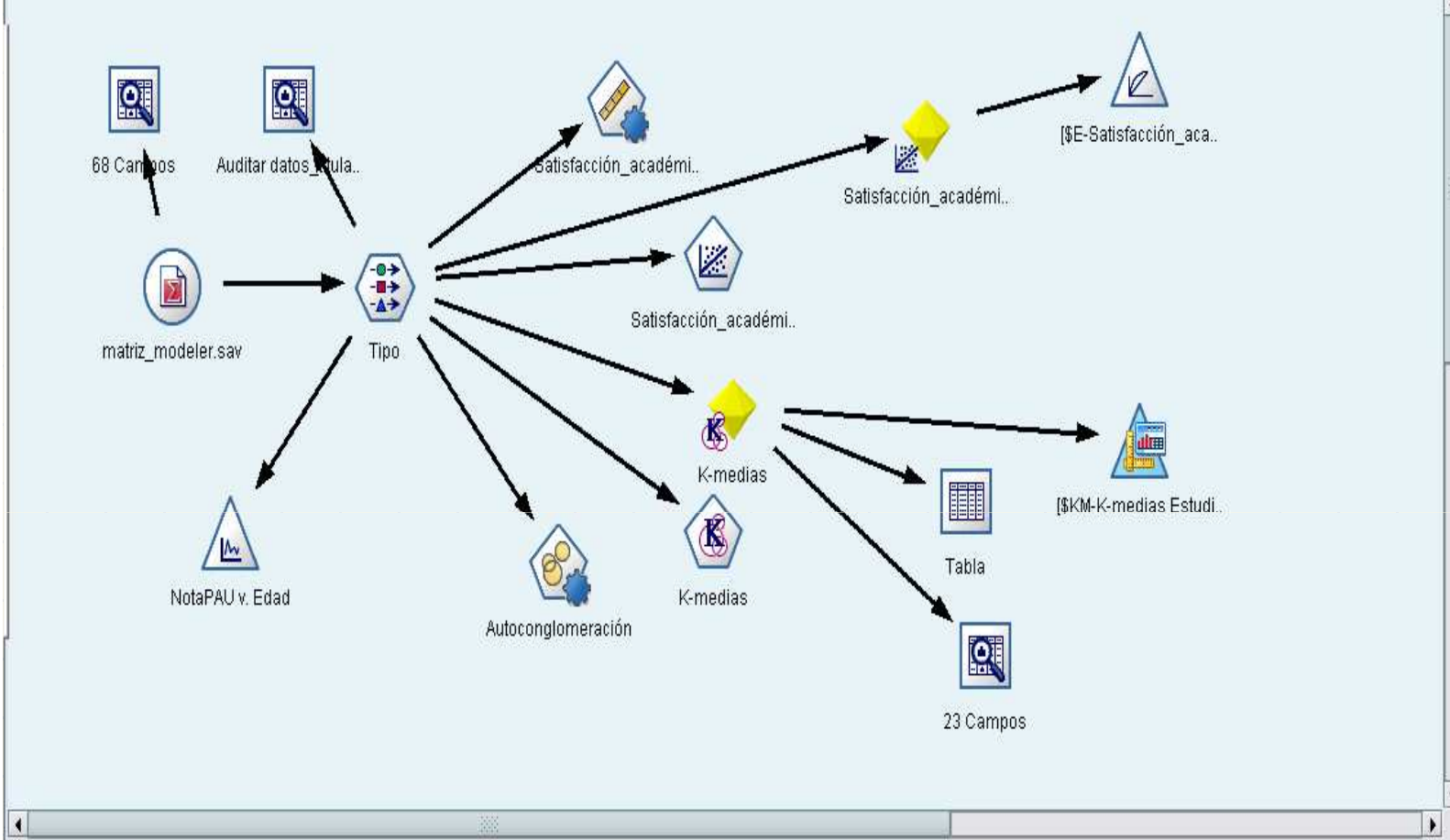
- **K-medias**

Agrupar conjunto de datos en conglomerados. Define un número fijo de grupos

- **Bietápico**

Es un método de dos pasos. Estima automáticamente el número óptimo de conglomerados para los datos





Modelos

Rutas Resultados

- Regr...
- Ruta2

CRISP-DM Clases

- (proyecto no guardado)
- Comprensión del negocio
  - Comprensión de los datos
  - Preparación de los datos
  - Modelado
  - Evaluación
  - Distribución

