

Analítica Social de los Datos y Herramientas para Analítica de Datos

Jose Aguilar
CEMISID, Escuela de Sistemas
Facultad de Ingeniería
Universidad de Los Andes
Mérida, Venezuela

El análisis sociales de datos es un estilo de análisis en el que es considerado las personas trabajan en un contexto social, de colaboración, para darle sentido a los datos.

El análisis de datos sociales se compone de dos partes:

- Los datos generados a partir de los sitios de redes sociales (o a través de aplicaciones sociales), y
- Análisis de los datos, en tiempo real (o casi en tiempo real), en los cuales se incluyen medidas para entender, y apropiadamente pesar, factores como la influencia, alcance y relevancia del contexto de los datos, y se incluye el horizonte de tiempo.

La analítica social de datos implica el análisis de las redes sociales, con el fin de comprender la percepción y actividad social que se inserta dentro de los datos.

Por ejemplo, analizar

- Interacciones directas con los demás (por ejemplo, mensajería, etc.),
- Uso de plataformas (por ejemplo, búsquedas, etiquetados).

En un sistema de análisis social de datos:

- los usuarios almacenan conjuntos de datos y crean representaciones visuales.
- Los conjuntos de datos y representaciones son accesibles para otros usuarios de la red o sitio web.
- Los usuarios pueden crear nuevas e interesantes representaciones, así como comentarios asociados a las.
- Se pueden usar un blogs y wikis para conducir esta inteligencia social.

Por lo general, podemos recuperar los datos sociales de una variedad de redes sociales, como Twitter, Facebook, nos sentimos bien, Wikipedia, etc.

- Dado que la mayoría de las redes sociales nos proporcionan la API, no es difícil para nosotros para recuperar los datos.
- El uso de API para obtener los datos es como enviar una solicitud a la página web
 y luego el sitio web envía los datos solicitados en forma de XML o en forma de
 JSON.

Métodos de análisis

- En la mayoría de los casos, queremos averiguar las relaciones entre los datos sociales y otro evento, o queremos predecir algunos eventos.
- Con el fin de lograr estos objetivos, necesitamos los métodos adecuados para hacer los análisis: métodos estadísticos, de aprendizaje de máquinas o de minería de datos.

Cuando se habla de análisis de datos sociales, hay una serie de factores que es importante tener en cuenta:

- Análisis de datos sofisticados: El análisis de datos sociales debe tomar en consideración una serie de factores (contexto, el contenido, el sentimiento) para proporcionar información adicional.
- La consideración del tiempo: Lo más relevante de un día (o incluso una hora) puede no ser en la siguiente. Ser capaz de ejecutar con rapidez el análisis es imperativo.
- Análisis de la influencia: la comprensión del impacto potencial de individuos específicos puede ser clave en la comprensión de cómo los mensajes podrían estar resonando. No se trata sólo de la cantidad, también tiene mucho que ver con la calidad.
- Análisis de las Redes: los datos social migran, crece n(o mueren) en base a cómo los datos se propagan a través de la red. Es como una actividad viral, que se inicia y se propaga.

Por ejemplo, Analítica Social del aprendizaje (SLA) centra la atención en los elementos de aprendizaje que son relevantes cuando se está aprendiendo en una cultura participativa.

- Ella busca hacer visible, comportamientos y patrones en el ambiente de aprendizaje.
- Los alumnos no son solitarios, y hacen actividades sociales, ya sea interactuando con los demás, o dejando trazas en las plataformas de sus actividades
- Los logros individuales, se transmiten a través de la interacción y la colaboración.

Los aprendices construyen el conocimiento, influenciados por sus objetivos, los sentimientos y las relaciones, las cuales cambian según el contexto.

- El éxito es una combinación de los conocimientos, habilidades individuales, el medio ambiente, el uso de herramientas, y la capacidad de trabajar juntos.
- La comprensión de aprendizaje obliga a prestar atención a los procesos del grupo en la construcción del conocimiento.
- La atención debe centrarse no sólo en los alumnos, sino también en las herramientas y contextos.

El análisis de aprendizaje desde una perspectiva social destaca un tipo de analítica que es empleada para darle sentido a la actividad de aprendizaje en un entorno social.

- Se deben identificar los comportamientos sociales y los patrones que significan un proceso eficaz, en entornos de aprendizaje emergente.
- Debe ser aplicable a diferentes escalas: redes nacionales e internacionales, pequeños grupos y alumnos individuales.



Algunos horizontes nuevos de análisis que permite SLA.

Necesitamos desarrollar nuevos conjuntos de análisis que se pueden utilizar para apoyar el aprendizaje y la enseñanza, considerando:

- Factores tecnológicos
- El cambio a "libre" y "abierto"
- La demanda de habilidades para la era del conocimiento
- La innovación requerida en el aprendizaje social
- Los nuevos desafíos de las instituciones educativas.

Factores tecnológicos

Los cambios en la tecnología no necesariamente implican cambios en la pedagogía.

Aquellos que ven la educación como la transferencia de información, utilizará medios interactivos para el almacenamiento, y la transmisión de la información;

Los que buscan el cambio conceptual explotan sus cualidades interactivas.

- Buscan con las tecnologías conocimientos que se basen en conjuntos de datos, para comprender el movimiento hacia el aprendizaje social en línea y analizar los procesos asociados.
- Conceptualmente, es introducir la idea de conocimiento distribuido, colaboración o innovación.

El cambio a lo libre y abierto

El Internet hace posibles modelos completamente nuevos de generación de ingresos

- Compañías en línea son capaces de ofrecer servicios de calidad de forma gratuita, produciendo una enorme variedad de herramientas y fuentes de contenido libre hospedado 'en la nube'.
- Los Recursos Abiertos Educativos (REA) son un vehículo de gran alcance, que hacen materiales de aprendizaje de alta calidad disponibles, no sólo de forma gratuita, sino también en formatos que promueven la re-mezcla,

Esto es amplificado por los esfuerzos para hacer que los datos sean abiertos al procesamiento de la máquina, así como a la interpretación humana.

- Ejemplos son las comunidades de Linked Data y de Web Semántic, el Gobierno Abierto, Ciencia 2.0 y 2.0 de la Salud
- La provisión de contenidos, las herramientas, etc., pueden esperarse que sean de forma gratuita, mientras que los estudiantes paguen por otros servicios: tutorías personalizados, orientación profesional, acreditación

La demanda de habilidades en la era del conocimiento

- Las habilidades de la era del conocimiento están relacionadas no sólo a un imperativo económico, sino a un deseo y al derecho a saber, una extensión de las oportunidades educativas,
- Además, es un aprendizaje que debe desarrollar habilidades y competencias.
- Esto implica la necesidad de apoyar el desarrollo de la creatividad y la curiosidad, habilidades de colaboración

La innovación requiere aprendizaje social

El aprendizaje social es la única manera en que las organizaciones pueden hacer frente en el cambiante mundo de hoy.

Nuestra comprensión de los contenidos se construye socialmente a través de conversaciones sobre ese contenido ya través de interacciones

Una característica importante del paradigma de la Web 2.0 es el grado de personalización que los usuarios finales esperan ahora.

Centrado en el aprendiz, no en el usuario

- Lo que quiero no es necesariamente lo que necesito, porque mi comprensión de la materia, es incompleta
- El trabajo es enseñar a la gente a pensar, y que requiere un aprendizaje más profundo, dejando al lado la seguridad cognitiva y emocional donde las hipótesis son meramente reforzadas.
- Esto implica un desafío para llevar a los alumnos fuera de su zona de confort, lo que subraya la importancia de la afirmación y el aliento al alumno para darle seguridad de salir.

Dos analíticas inherentemente sociales, y tres de análisis socializados:

- Inherentemente sociales tienen sentido en un contexto colectivo:
 - Analítica de Redes Sociales- relaciones interpersonales en plataformas sociales.
 - Analítica del Discurso- el lenguaje es una herramienta fundamental para la construcción del conocimiento.
- Análisis socializados: son analíticas personales, pertinentes socialmente en un contexto colectivo:
 - Analítica del contenido -es una de las características definitorias de la Web 2.0
 - Analítica de la Disposición- Motivación para aprender (centro de la innovación)
 - Analítica del Contexto- clave en entornos móviles, aprendizaje informal, etc.

Análisis de redes sociales

El aprendizaje en red implica el uso de las TIC para promover conexiones entre un alumno y otros alumnos, entre alumnos y tutores, y entre las comunidades de aprendizaje y recursos de aprendizaje.

- Estas redes se componen de actores (tanto de personas como de recursos) y las relaciones entre ellos.
- Actores con una relación entre ellos se dice que están vinculados, y estos lazos pueden ser clasificados como fuertes o débiles, dependiendo de su frecuencia, calidad o importancia.
- El análisis de redes sociales investiga los procesos en la red, las propiedades de las relaciones, los roles y la formación de la red, para entender cómo la gente desarrolla y mantiene estas relaciones para apoyar el aprendizaje

Ejemplo implementación de análisis de redes sociales de aprendizaje es SNAPP (Social Networks Adapting Pedagogical Practice), una herramienta de visualización de la red que incluye:

- Identificar a los estudiantes desconectados
- Identificar a los agentes de información claves dentro de una clase
- Indicar el grado en que una comunidad de aprendizaie se está desarrollando.

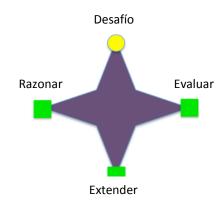
La analítica del discurso aprendizaje social

El análisis del discurso es el término colectivo para el análisis de la serie de eventos comunicativos.

- Algunos de estos enfoques se centran en la cara-a-cara o interacción oral.
- Otros proporcionan nuevas formas de entender las grandes cantidades de texto generados en los cursos en línea y conferencias.
- Otros lo usan para comprender la relación entre las dimensiones interactivas, cognitivas, la interacción en línea, y el seguimiento de los intercambios (IRF).
- En las discusiones asíncronas, muestra cómo los grupos crean y mantienen la comunidad y la coherencia a través del uso de dispositivos discursivos.

Dos ejemplos de análisis de discurso son:

- El diálogo exploratorio, que incluye descubrir en los foros desafíos ('pero si', 'no creo'), evaluaciones ("Buen punto", "importante",), el razonamiento ('quiere decir eso', 'mi entendimiento es') y las extensiones ('Siguiente paso', 'se refiere a'), por las palabras y frases claves, y
- El desarrollo Mapas Conceptuales para la deliberación y argumentación.



Analítica de la disposición para el aprendizaje

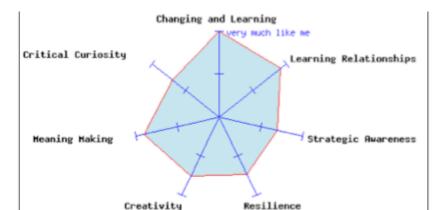
- Es la única que se origina del campo de la investigación educativa
- Es definida por un modelo de siete dimensiones
- Mezcla la experiencia, la motivación y la capacidad de una persona para toda la vida aprender e identificar a las oportunidades de aprendizaje
- Disposición para el aprendizaje no es "estilo de aprendizaje": las disposiciones de aprendizaje varían según el contexto,

El compromiso en el aprendizaje es una combinación compleja de las identidades de los alumnos, disposiciones, valores, actitudes y habilidades.

Cuando éstos son frágiles, los estudiantes luchan por alcanzar su potencial en las evaluaciones convencionales, y no están preparados para la novedad, la complejidad de los retos, y las muchas otras esferas de la vida que requieren cualidades como la capacidad de recuperación, el pensamiento crítico y habilidades de colaboración.

Las disposiciones de aprendizaje pueden ser modelados como un constructo multidimensional, llamado **Aprendizaje de energía**, actualmente evaluados a través de un cuestionario web llamado ELLI

ELLI genera una visualización analítica que se utiliza para apovar la auto-reflexión y cambio



Análisis de contenidos de aprendizaje social

- Originario de los campos relacionados con los sistemas de recomendación y recuperación de la información
- Es la variedad de métodos automatizados que se pueden utilizar para examinar, indizar y filtrar los recursos multimedia en línea, con la intención de guiar a los estudiantes a través del océano de los recursos potenciales disponibles para ellos.
- En combinación con el análisis de contexto de aprendizaje o con los términos de búsqueda definidos, análisis de contenidos se pueden utilizar para proporcionar recomendaciones de recursos que se adaptan bien a las necesidades de un individuo o de las necesidades de un grupo de alumnos.
- La recuperación de información representa las técnicas para la indexación automática y el filtrado de contenido, ya sea textual o multimedia (por ejemplo, imágenes, vídeo o música).
- En conjunto, estos elementos se pueden utilizar para proporcionar nuevos métodos para sugerir, de navegación o búsqueda de medios educativos.

Análisis de contenidos de aprendizaje social

Los componentes básicos de análisis de contenidos de aprendizaje social, ofrecen la opción de ver las palabras clave, enlaces directos o imágenes relacionadas con la página web abierta. Esta información sobre las imágenes se puede combinar con la búsqueda de similitud visual, para identificar y recomendar otros recursos que hacen uso de estas

imágenes,



Análisis del contexto de aprendizaje

- El aprendizaje se pueden aplicar a una amplia variedad de contextos que se extiende mucho más allá de los sistemas institucionales.
- Pueden ser utilizados en contextos formales como escuelas, colegios y universidades, en contextos informales en las que los estudiantes eligen tanto el proceso como el objetivo de su aprendizaje.
- En algunos casos, los estudiantes están en entornos síncronos, sobre la base de que los participantes son co-presente en el tiempo, o se encuentran en entornos asíncronos, donde se supone que van a participar en diferentes momentos.
- Puede ocurrir en una red, en un grupo de afinidad, en las comunidades de investigación, comunidades de interés, etc.

Bajo el título análisis del "contexto" se agrupan las diversas herramientas analíticas que exponen, hacen uso de, o tratan de comprender estos contextos en relación con el aprendizaje.

Análisis del contexto de aprendizaje

Se puede caracterizar el contexto de una entidad (por ejemplo, un estudiante) en función de cinco categorías distintas:

- contexto Individualidad, incluye información sobre la entidad dentro del contexto. En el caso de los estudiantes, esto podría incluir su idioma, su comportamiento, sus preferencias y sus objetivos
- contexto del Tiempo, que incluye puntos en los tiempos, y oscila historias, por lo puede tener en cuenta el flujo de trabajo, cursos de larga duración y las historias de interacción
- contexto de Localización, puede incluir la ubicación absoluta, ubicación en relación con personas o recursos, o la ubicación virtual (dirección IP)
- contexto de la Actividad, se refiere a los objetivos, tareas y acciones
- contexto de la Relaciones, capta las relaciones de la entidad con otras entidades, por ejemplo, entre los alumnos, los profesores y los recursos.

- Si las herramientas y los datos de SLA se colocan en las manos de los estudiantes, el equilibrio de poder cambia significativamente.
- Si la analítica están llamando la atención a los alumnos para su desarrollo como, estudiantes, intrínsecamente los motiva conscientes de sí mismos, moviendo en la dirección opuesta a depender pasivamente en la institución o plataforma para decirles cómo lo están haciendo y qué hacer a continuación.
- Si la analítica se centran en proporcionar retroalimentación formativa para mejorar el proceso de aprendizaje, en lugar de hacer juicios automatizados sobre el dominio en un tema determinado, puede haber un menor número de preocupaciones en torno a la eliminación de los tutores humanos del circuito de retroalimentación.

"Libre y abierto" es una expectativa fundamental y dinámico dentro de aprendizaje social en línea.

- Muchas herramientas de SLA están disponibles en las versiones de código abierto.
- Se convierte en normal que los patrones de SLA y los datos sean abiertos.
- Alumnos están dispuestos a pagar por características más potentes, una vez que las herramientas de mayor éxito se han ganado su derecho a cobrar.

Las instituciones que carecen de la infraestructura necesaria para SLA pedirá los servicios de SLA en la computación "en la nube",

- Los alumnos individuales o comunidades que necesitan este tipo de servicios también utilizaran estos servicios.
- instituciones educativas explotarán su experiencia pedagógica para proporcionar servicios de consultoría de SLA.

Aspiraciones de todas las culturas han ido cambiando hacia un creciente deseo de participación y libre expresión. La web social es una expresión de este cambio,

 Las herramientas SLA se convierten en una parte importante del sentido de la identidad de los individuos, y su capacidad para evidenciar sus habilidades. Por ejemplo, podríamos ver Placas tales como: "yo puedo llevar debates complejos ", "yo puedo guiar a los estudiantes en la construcción de su creatividad."

La innovación en entornos complejos y turbulentos requiere infraestructuras de creación de conocimiento y negociación sociales construidas sobre las relaciones de calidad y conversaciones, en orden a trabajar juntos para resolver "problemas perversos".

SLA se convierten en una parte integral del juego de herramientas,

El papel de las instituciones educativas está cambiando. Se están moviendo cada vez más para proporcionar apoyo personalizado para aprender a pensar profundamente, y aprender cómo ser un miembro efectivo de las comunidades.

- Las instituciones educativas ya no son la única opción para evidenciar el aprendizaje avanzado.
- Analítica se convierten en una nueva forma de pruebas de confianza, que se genera a partir de conjuntos de datos verificables, públicas, privadas, que no podrían haber sido fabricados.

Weka



Weka (Data Mining Tool)

- Colección de algoritmos open source de aprendizaje automático
 - pre-procesamiento
 - clasificadores
 - clustering
 - Reglas de asociación



- Herramientas de pre-procesamiento, algoritmos de aprendizaje y métodos de evaluación
- Interfaz Grafica (incl. visualización de datos)
- Ambiente para comparer algoritmos de aprendizaje



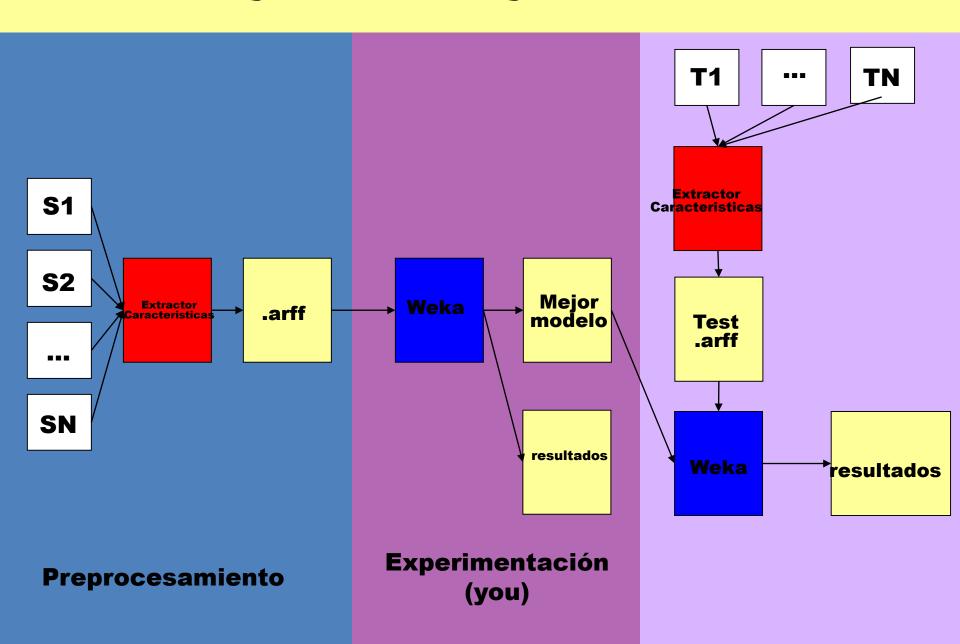
Weka (Data Mining Tool)

- Weka es una herramienta de minería de datos de código abierto desarrollado en Java.
- Se utiliza para la investigación, la educación, y las aplicaciones.
- Basado en Java

http://www.cs.waikato.ac.nz/ml/weka/

- Se puede ejecutar en Windows, Linux y Mac.
- 3 modos of operation
 - GUI: aplicar directamente a un conjunto de datos (con interfaz gráfica de usuario)
 - Linea de comando
 - Java API: llamados desde su propio código Java (usando biblioteca Weka de Java).

Flujo de trabajo con Weka



Weka

- Entrada de datos a Weka (Input)
 - Usa archivos planos
 - Formatos ".arff", C4.5, CSV,
 - Datos pueden ser leidos desde URL o BD SQL (usando JDBC)

Weka para Data Mining

Salida desde Weka (Output)

Entrada de datos a Weka (Input)

 El más popular formato is "arff" ("arff" es la extensión del nombre del archivo).

FILE FORMAT @relation RELATION_NAME @attribute ATTRIBUTE NAME ATTRIBUTE TYPR @data DATAROW1 DATAROW2 DATAROW3

Entrada de datos a Weka (Input)

archivo "arff"

```
@relation heart-disease-simplified

Atributo numeric

@attribute age numeric

@attribute sex { female, male}

@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angina}

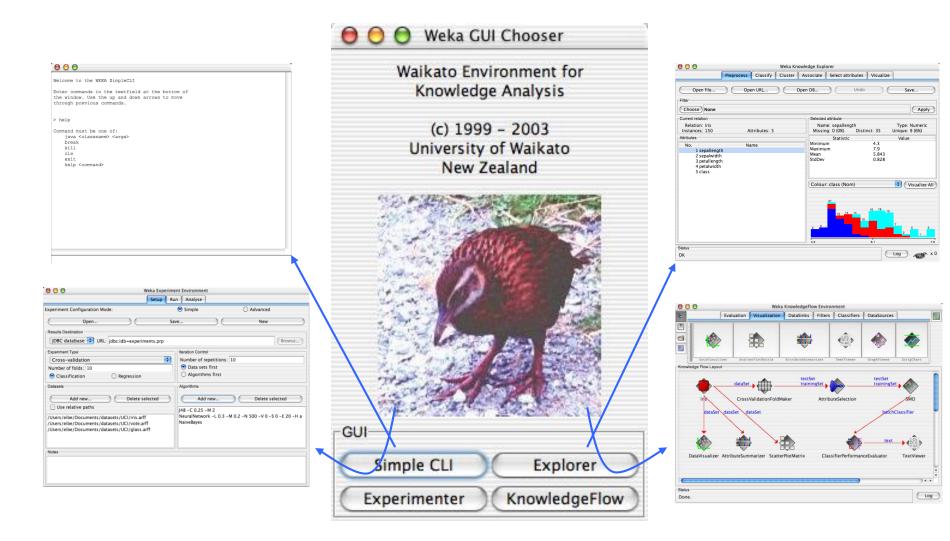
@attribute cholesterol numeric

@attribute exercise_induced_angina { no, yes}

@attribute class { present, not_present}
```

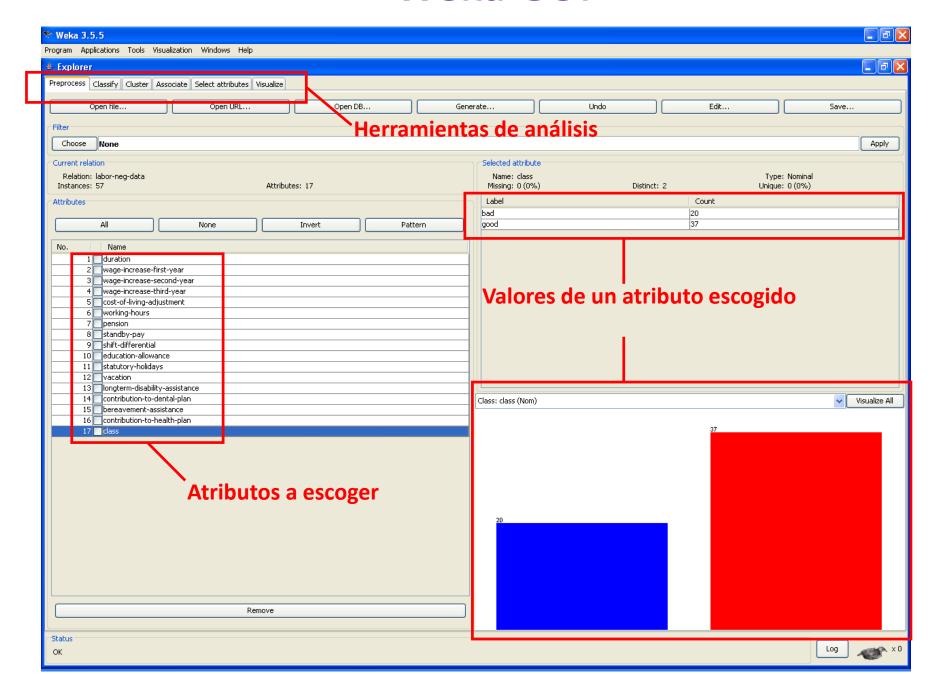
@data
63,male,typ_angina,233,no,not_present
67,male,asympt,286,yes,present
67,male,asympt,229,yes,present
38,female,non_anginal,?,no,not_present

Con interfaz gráfica de usuario

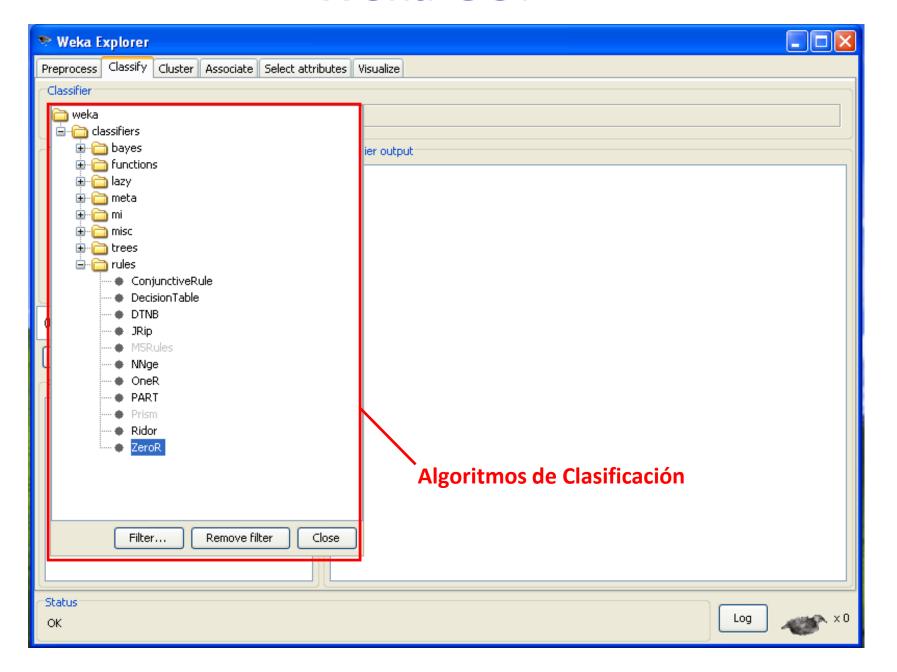


o usando biblioteca Weka de Java

Weka GUI



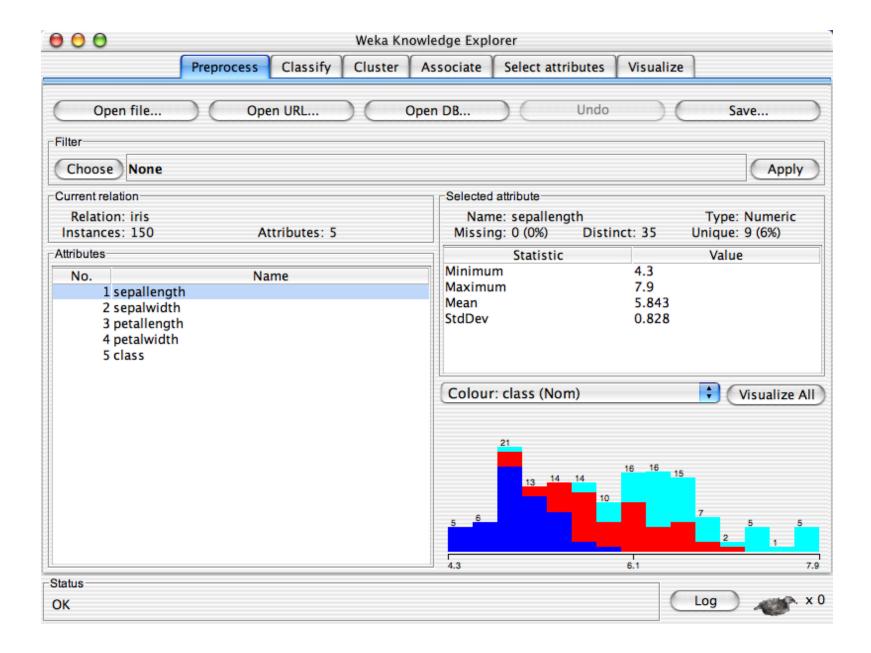
Weka GUI

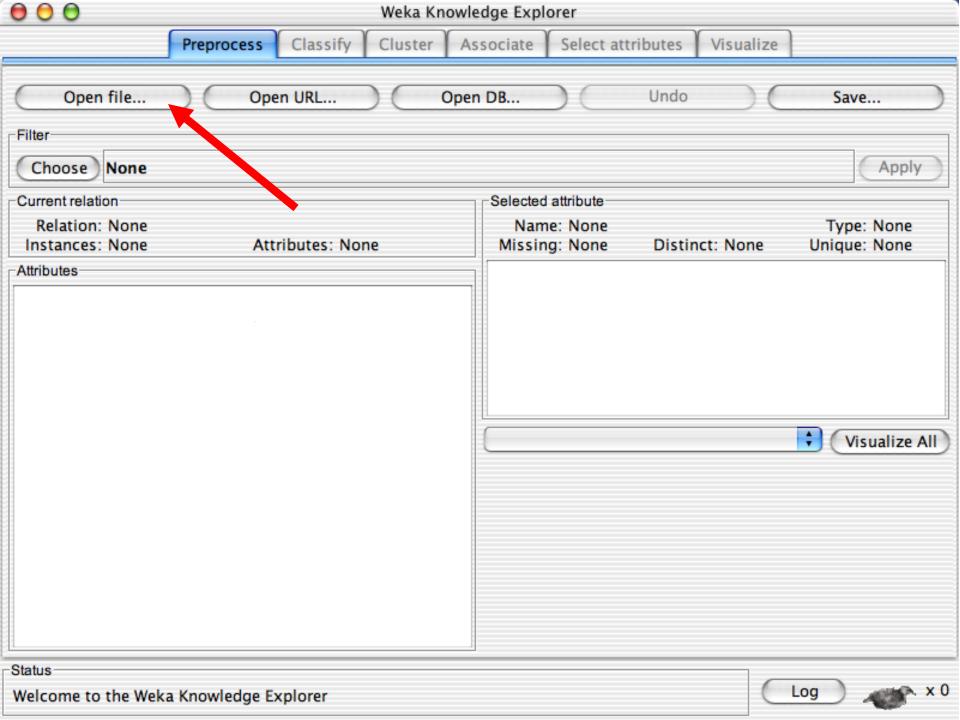


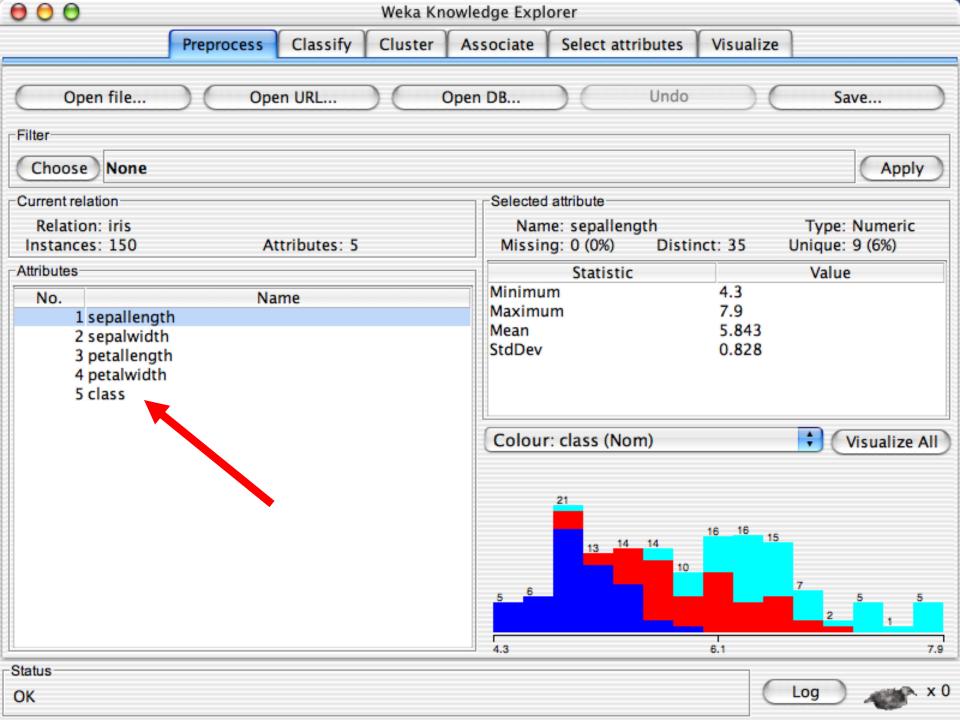
Explorar: Preprocesamiento

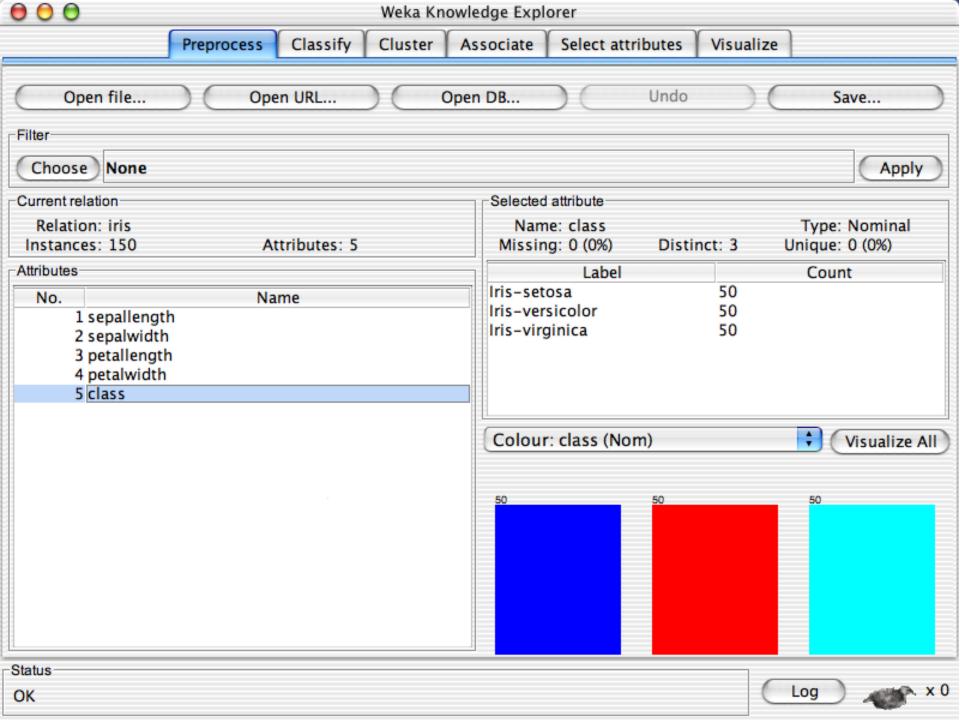
 Las herramientas de pre-procesamiento en WEKA son llamados "filtros"

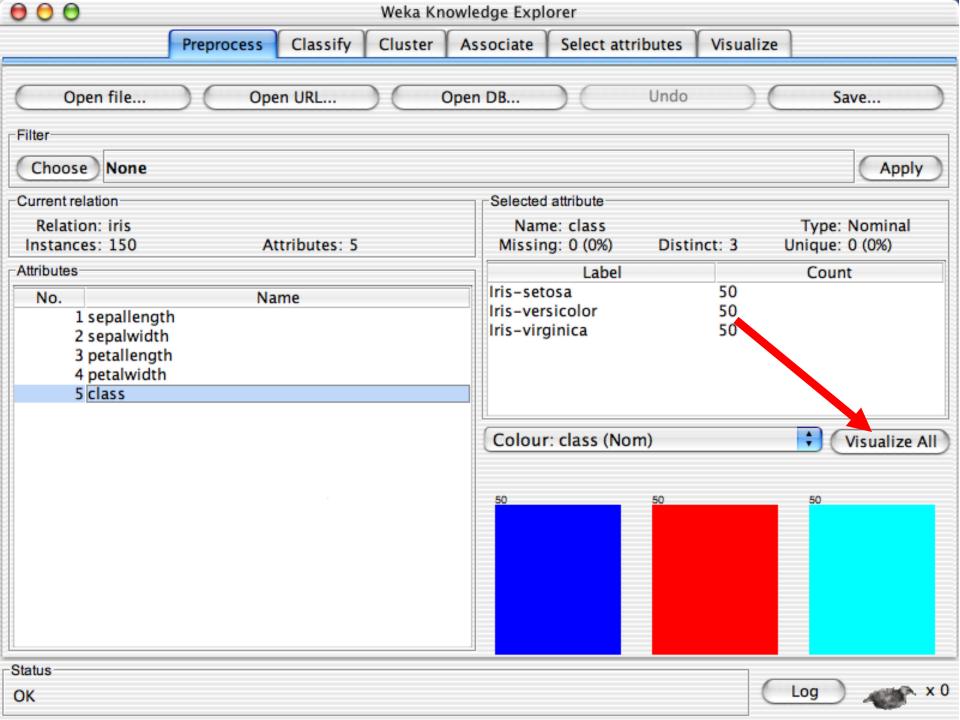
- Los filtros son:
 - Discretizar, normalizar, selección de atributos, transformar, combinar atributos, etc.

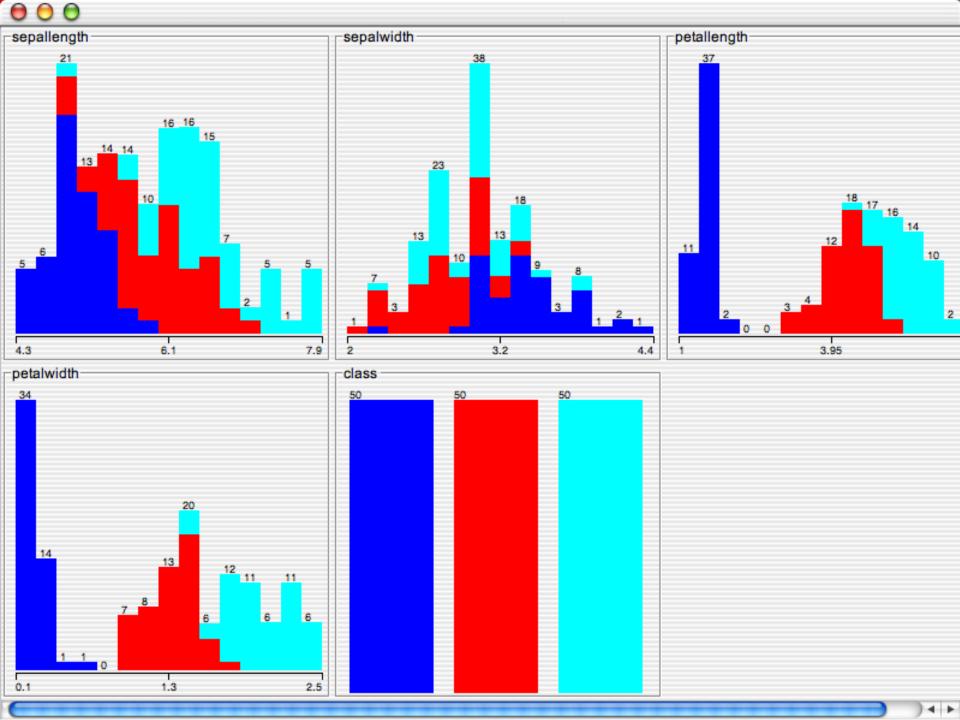


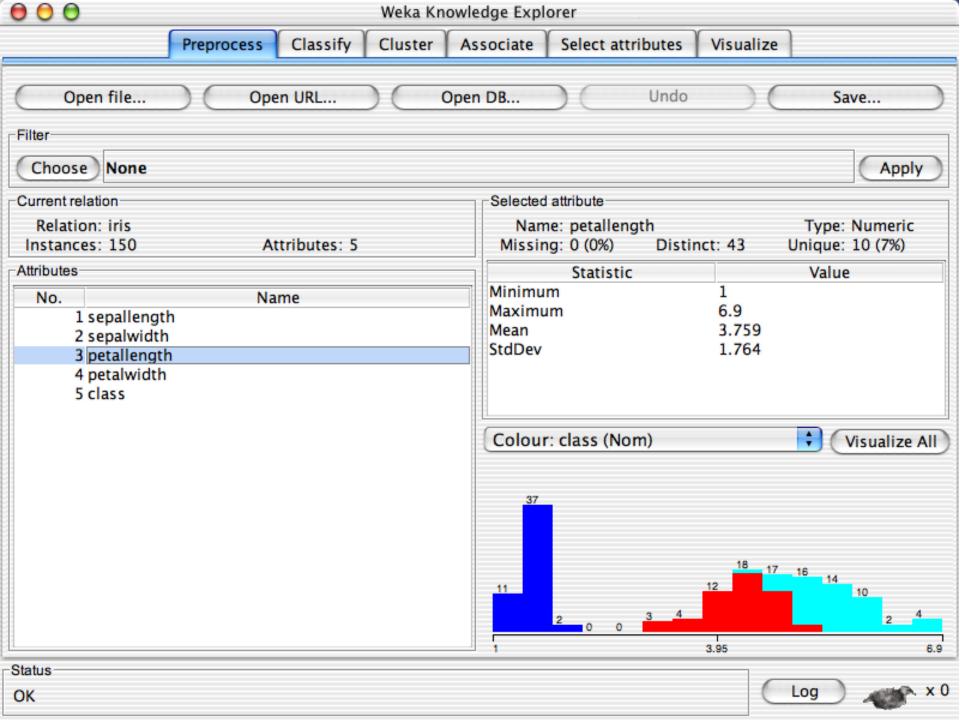


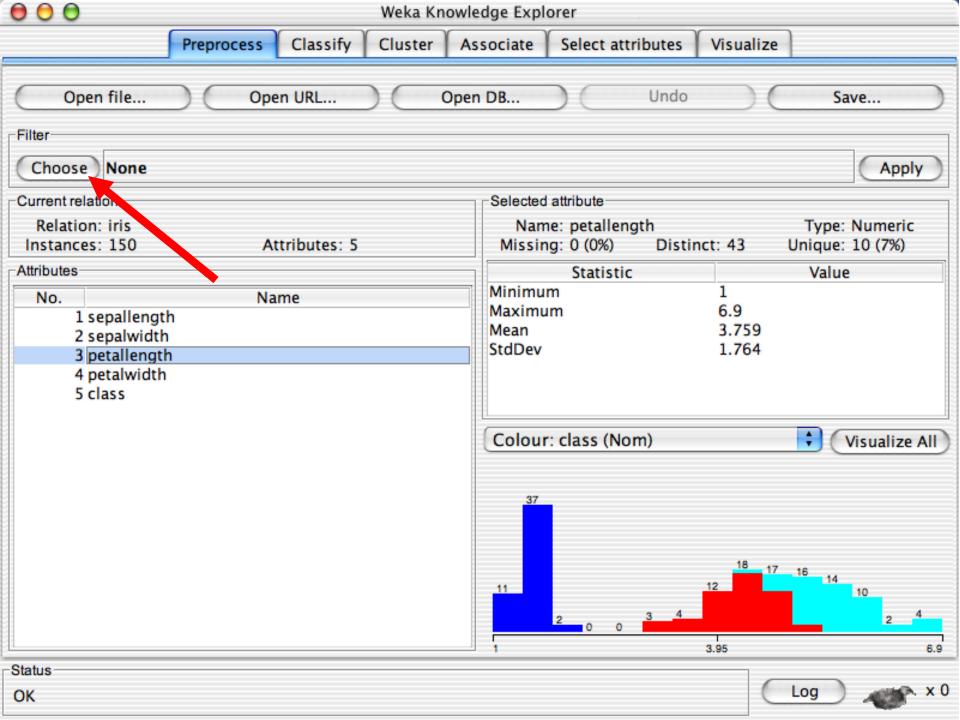


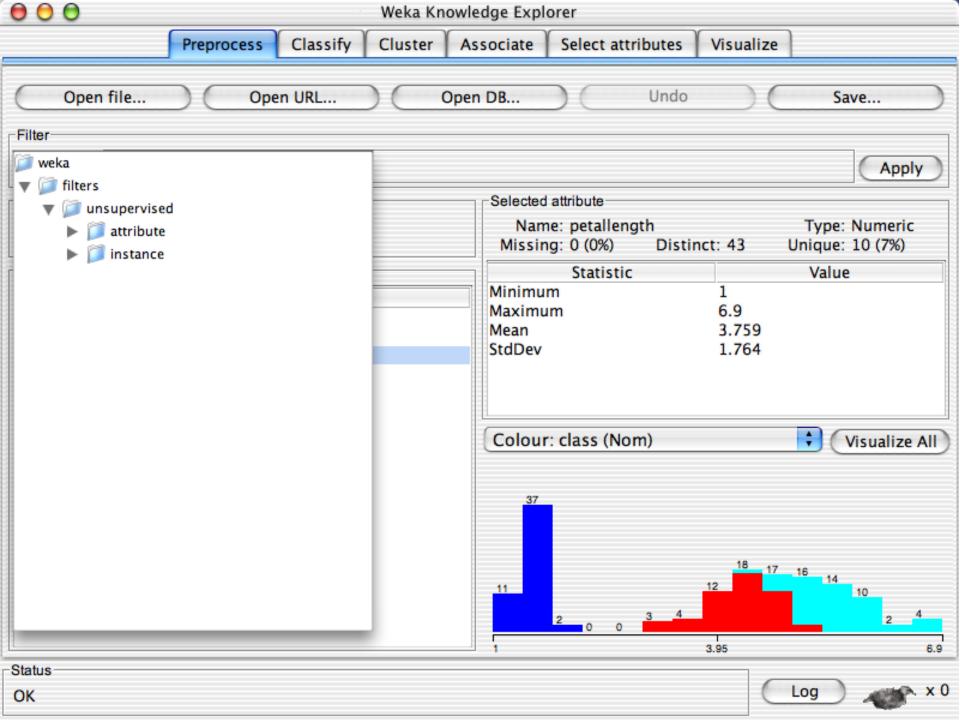


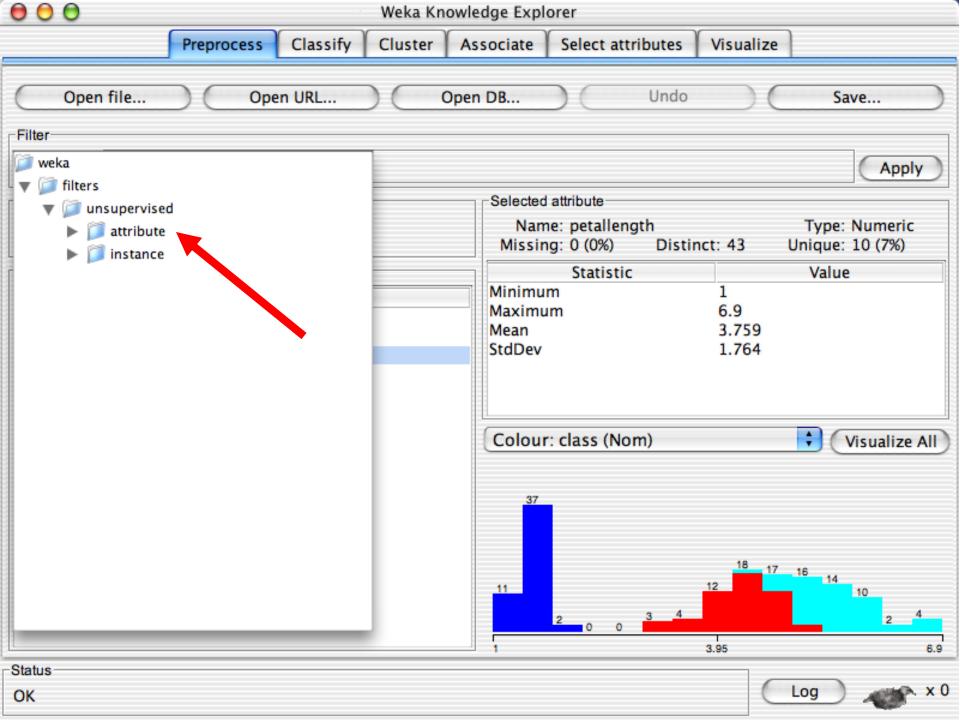


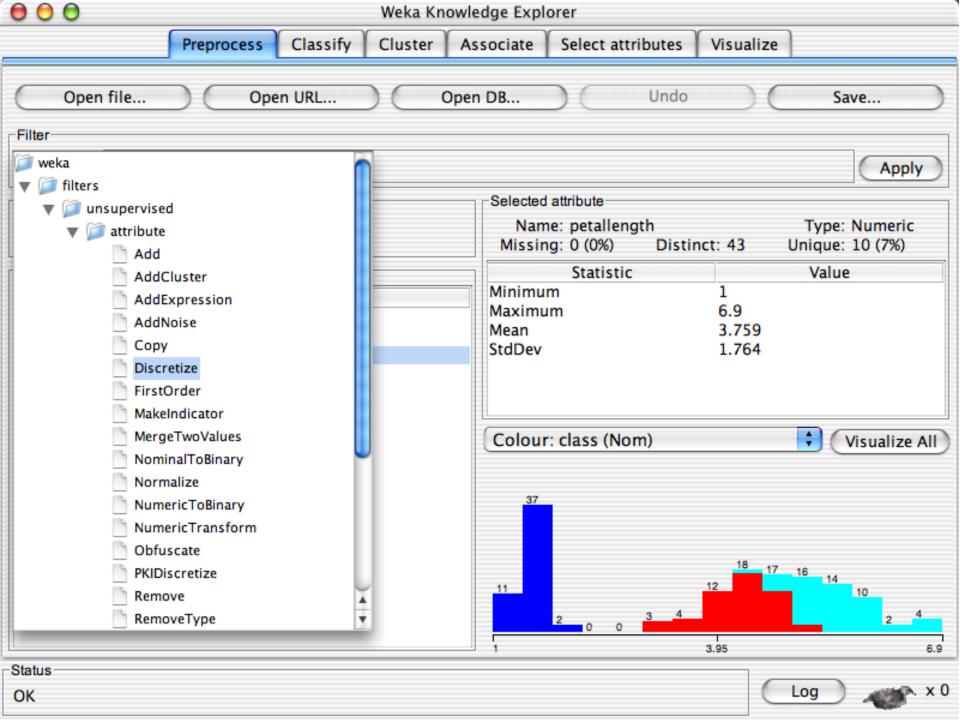


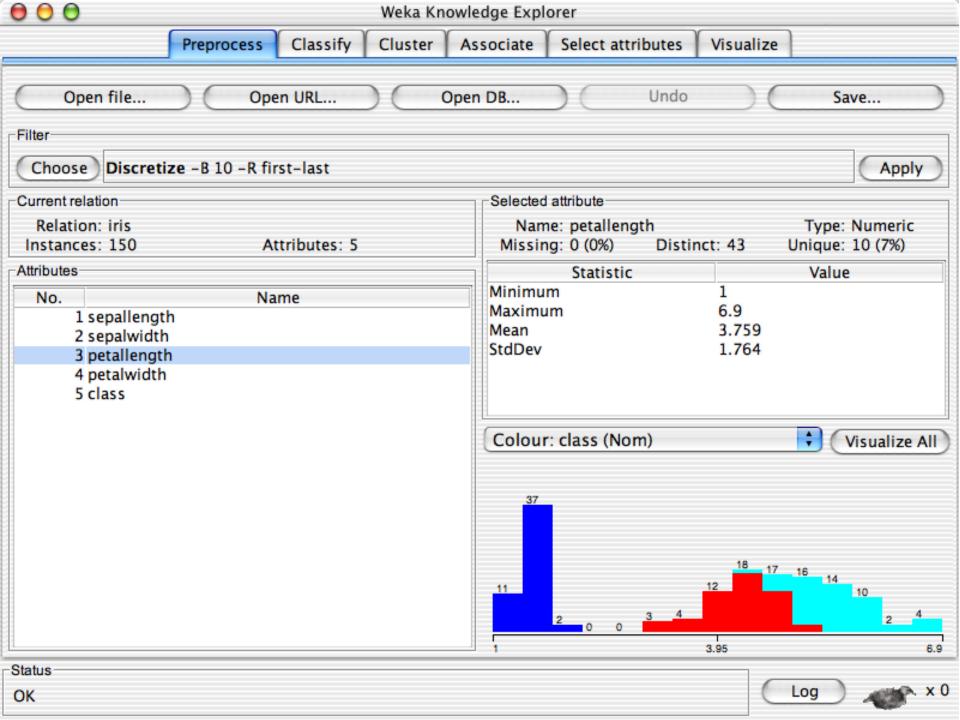


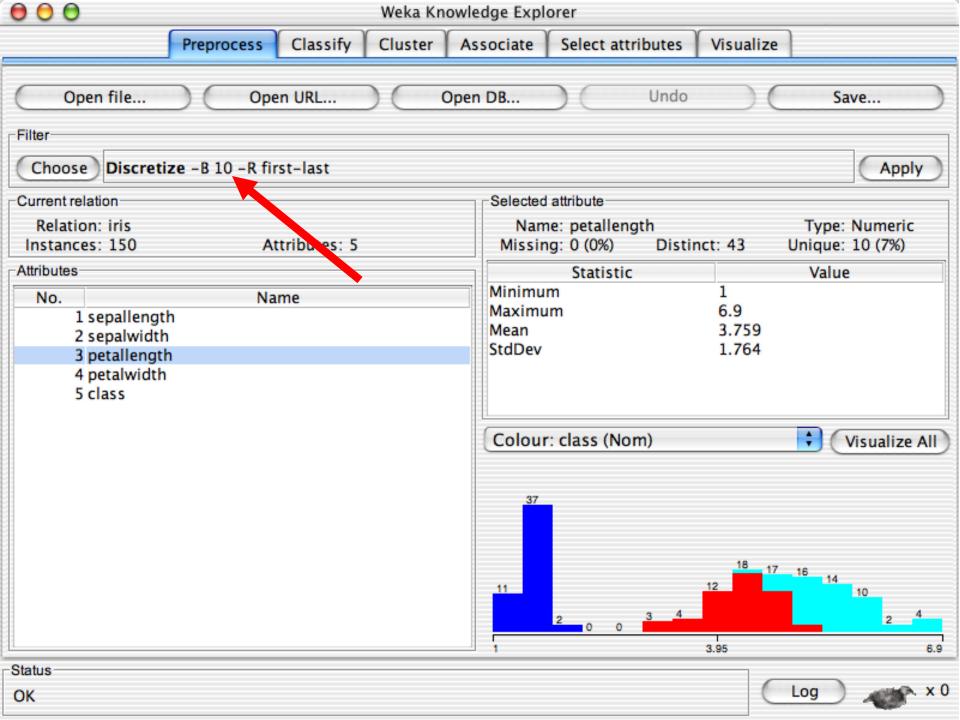


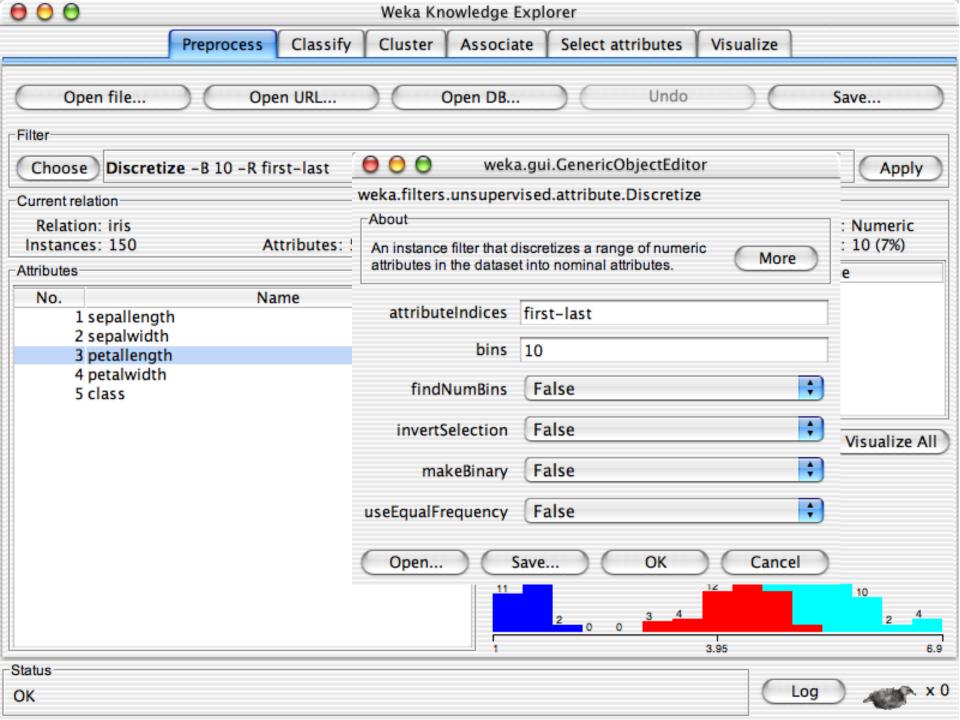


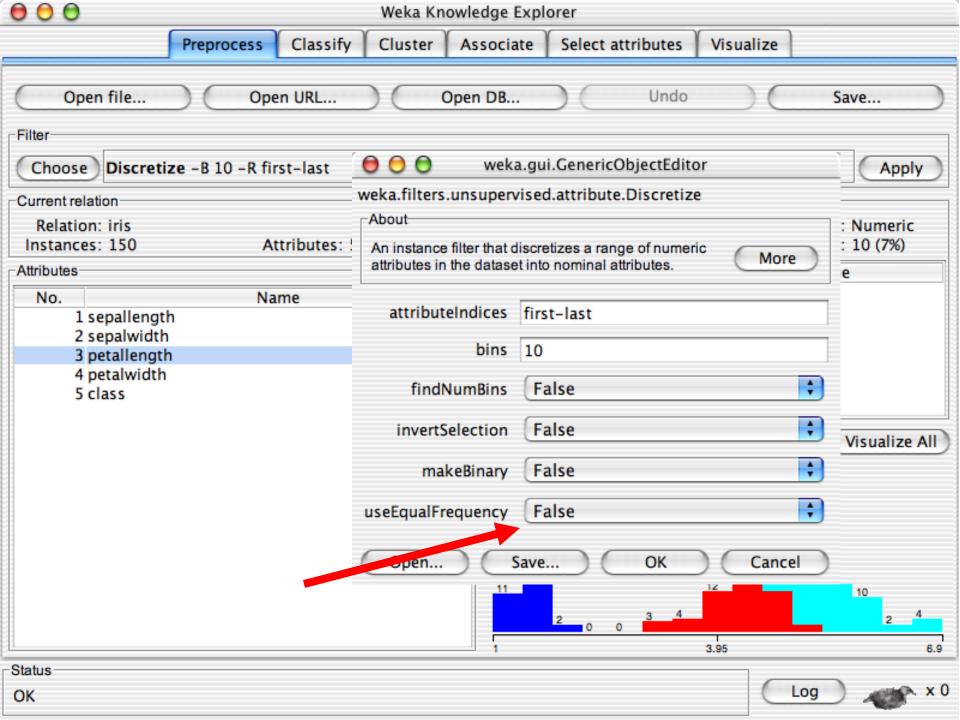


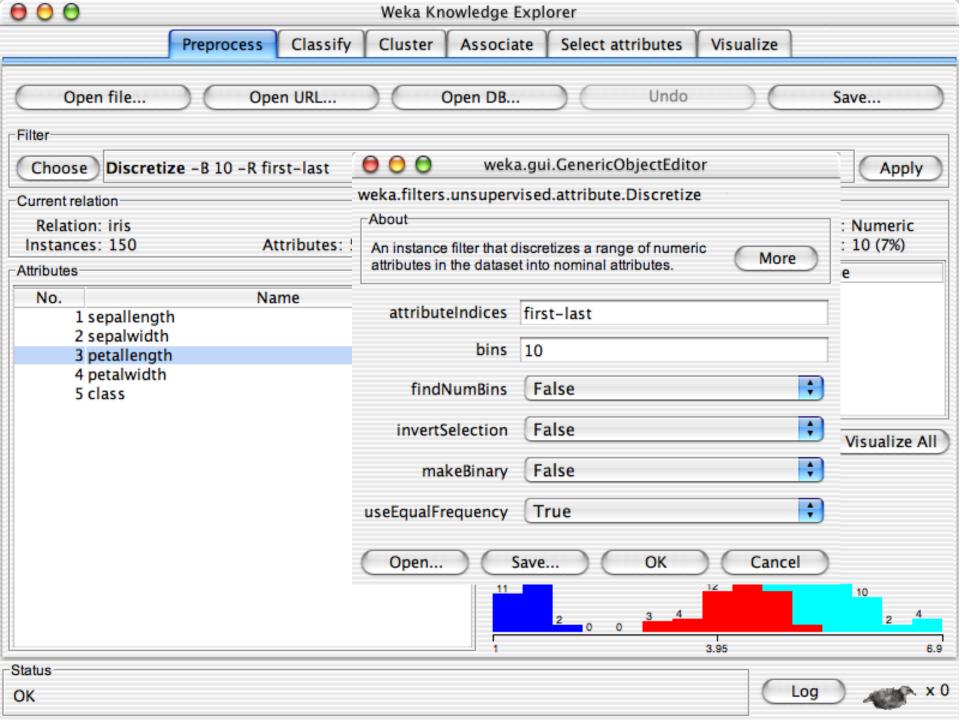


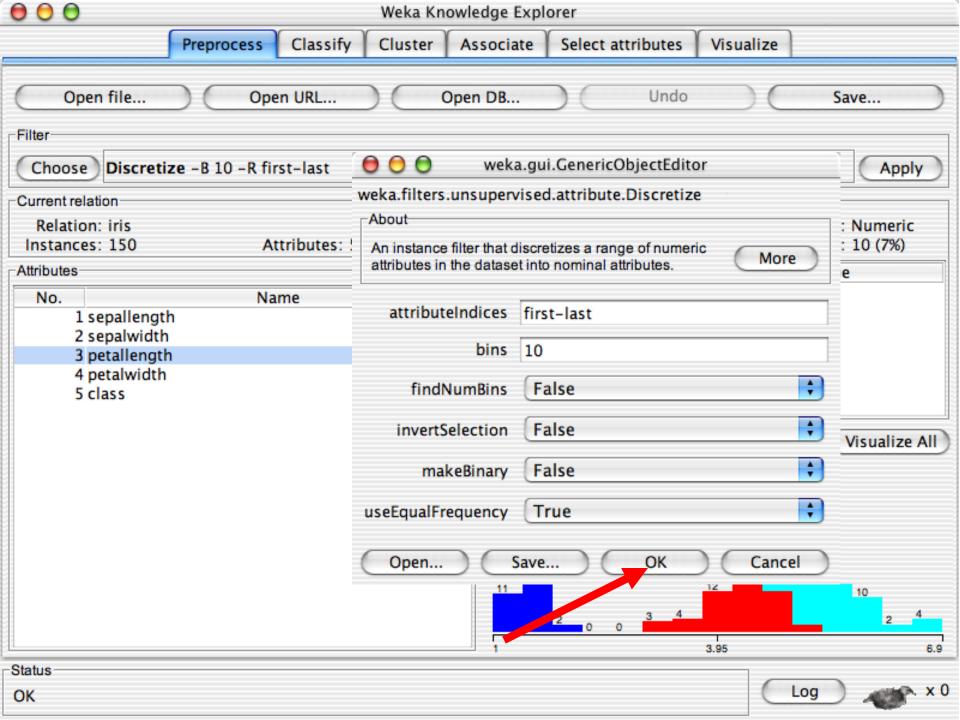


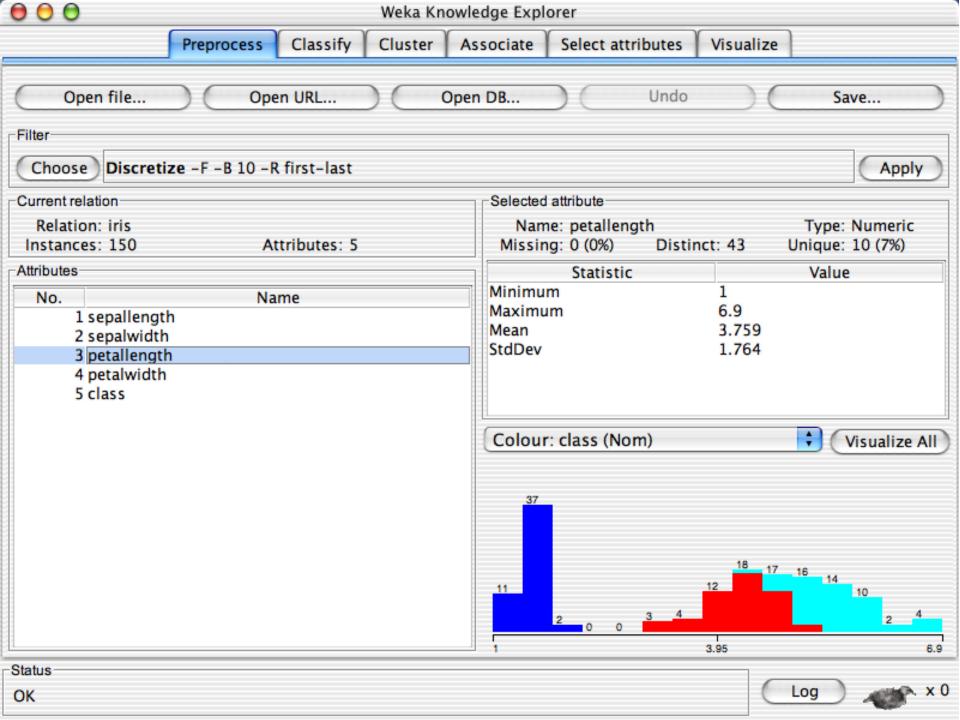


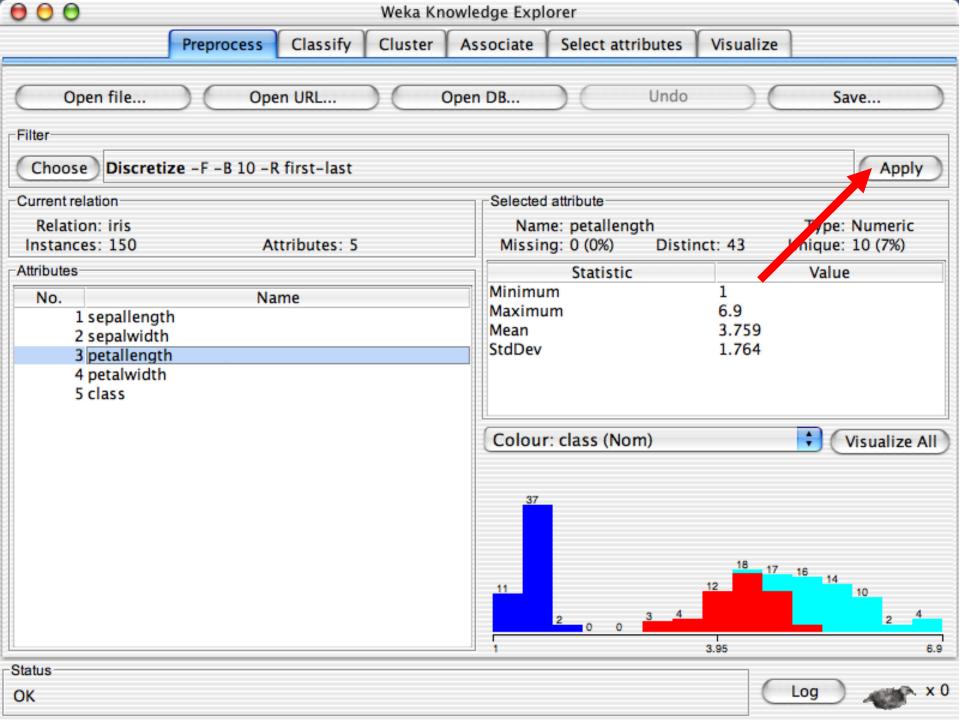


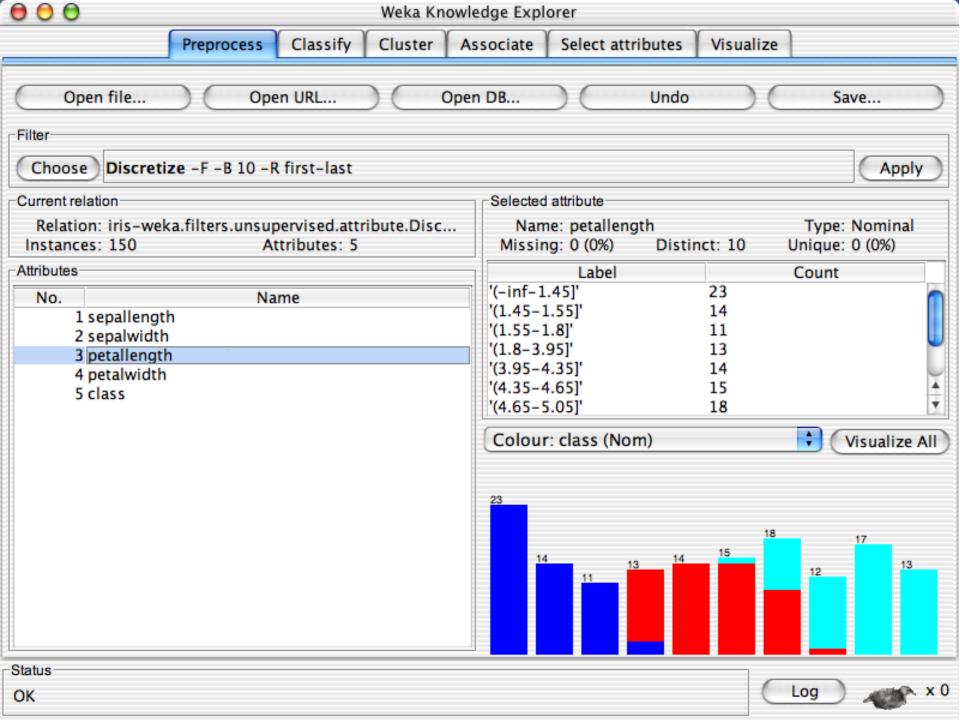






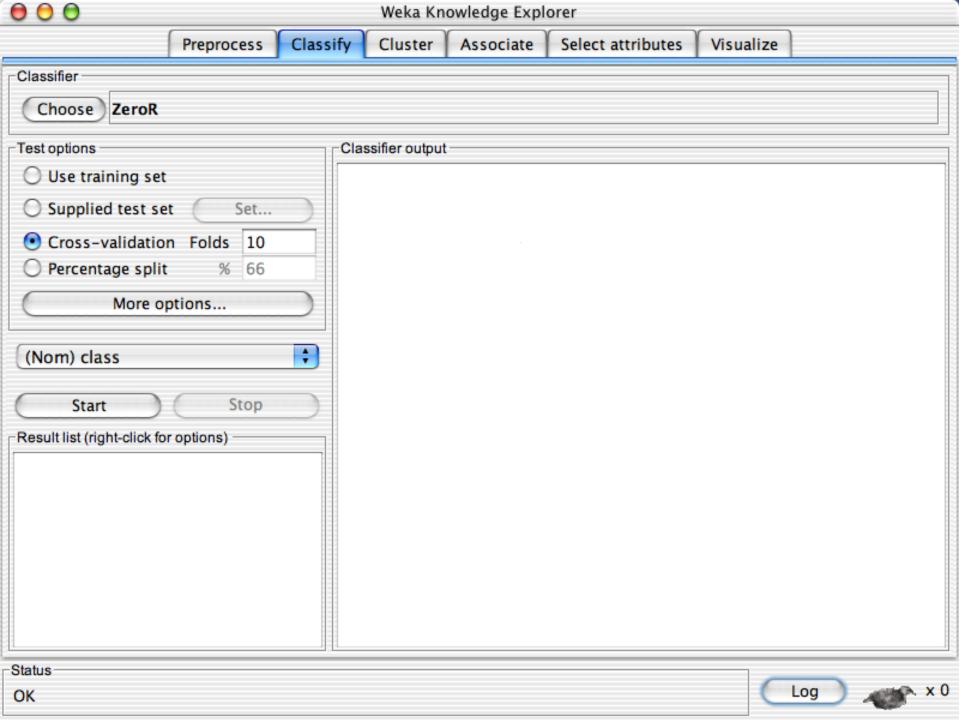


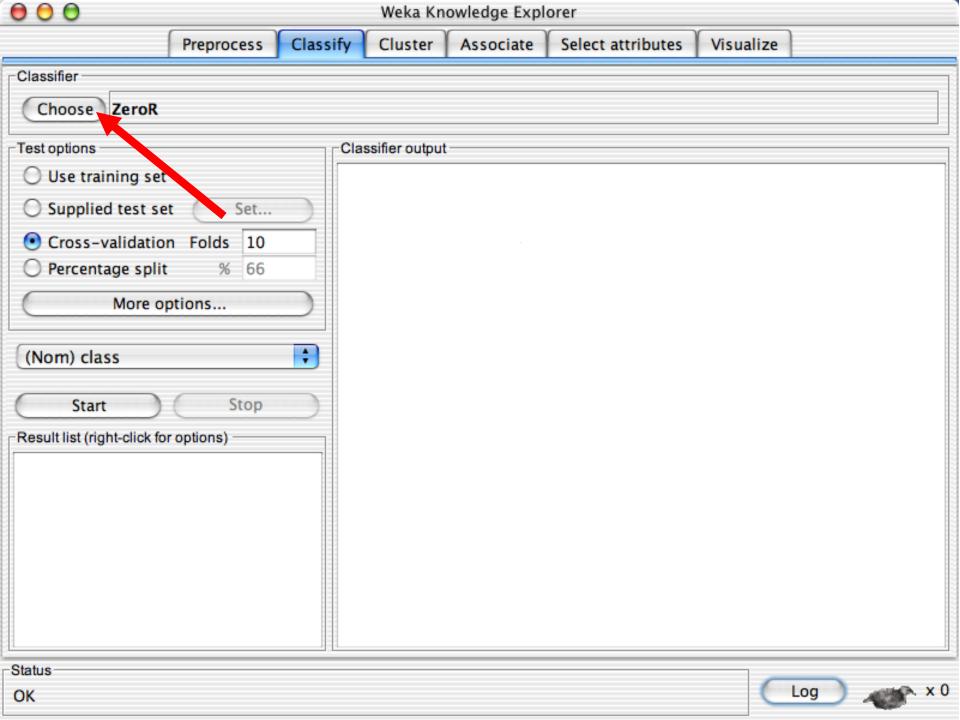


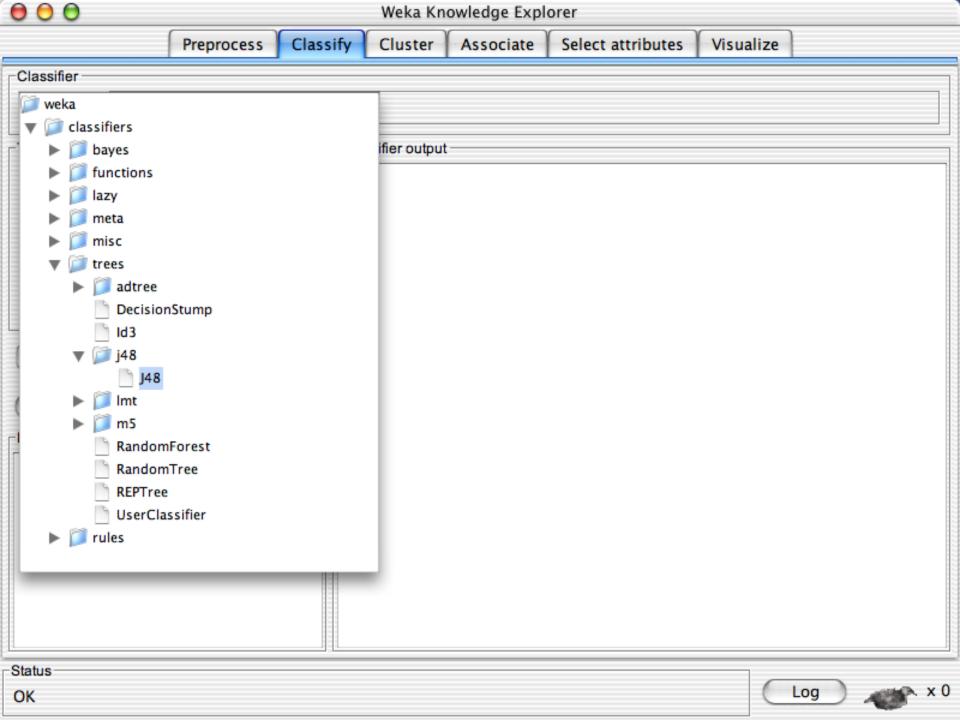


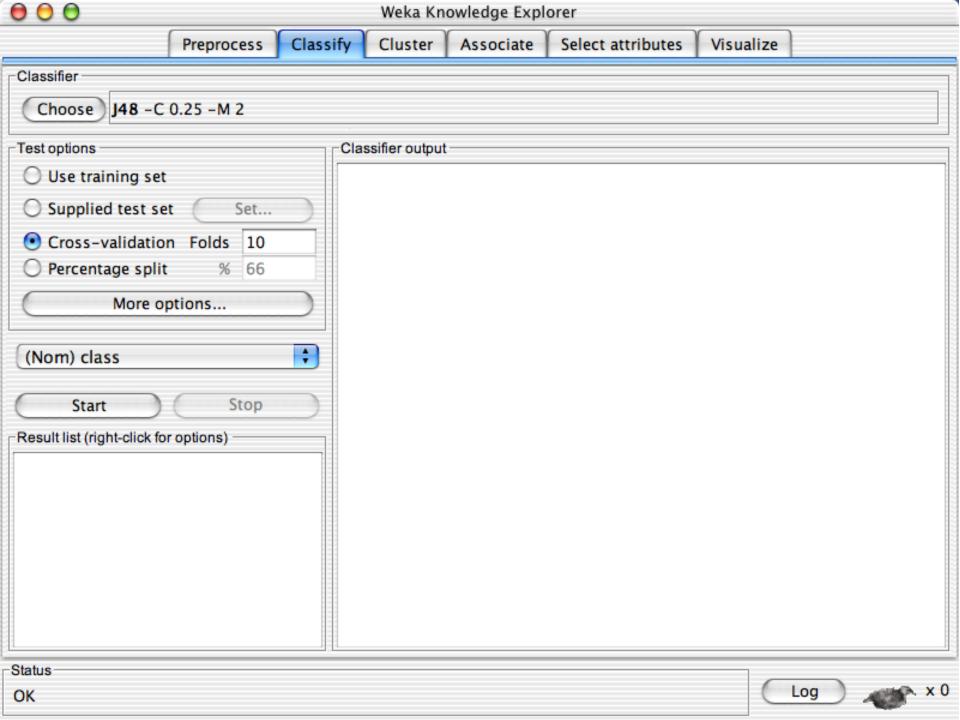
Explorar: construir "clasificadores"

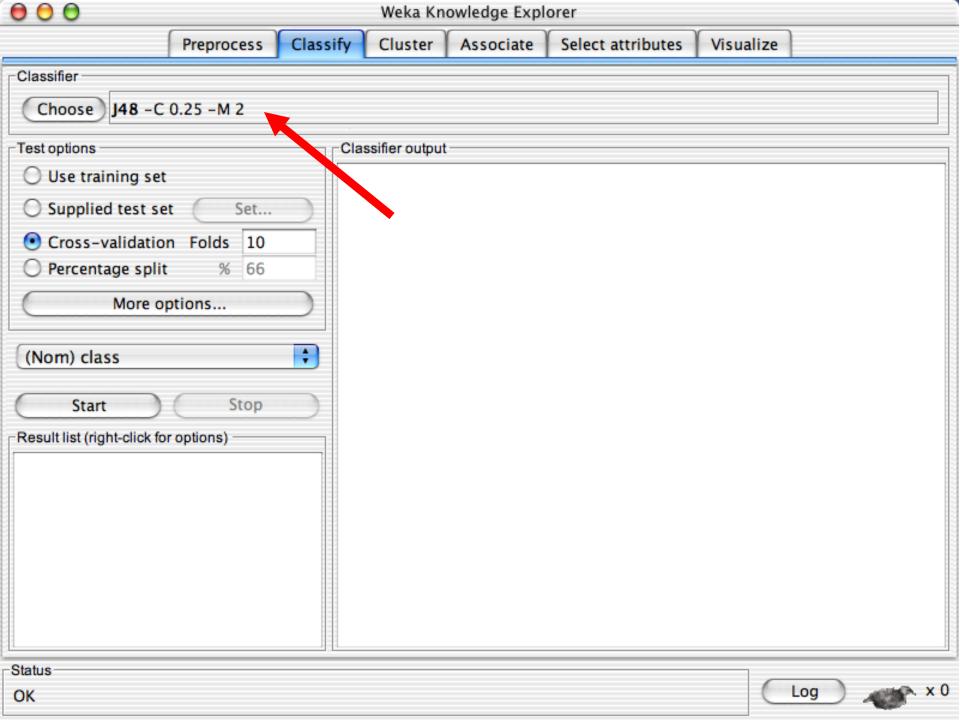
- Son modelos para predecir cantidades nominales o numericas
- Ejemplo son:
 - Árboles de decisión
 - Support vector machines,
 - Perceptron Multi-capa,
 - Regresión lógica,
 - Red de Bayes, ...

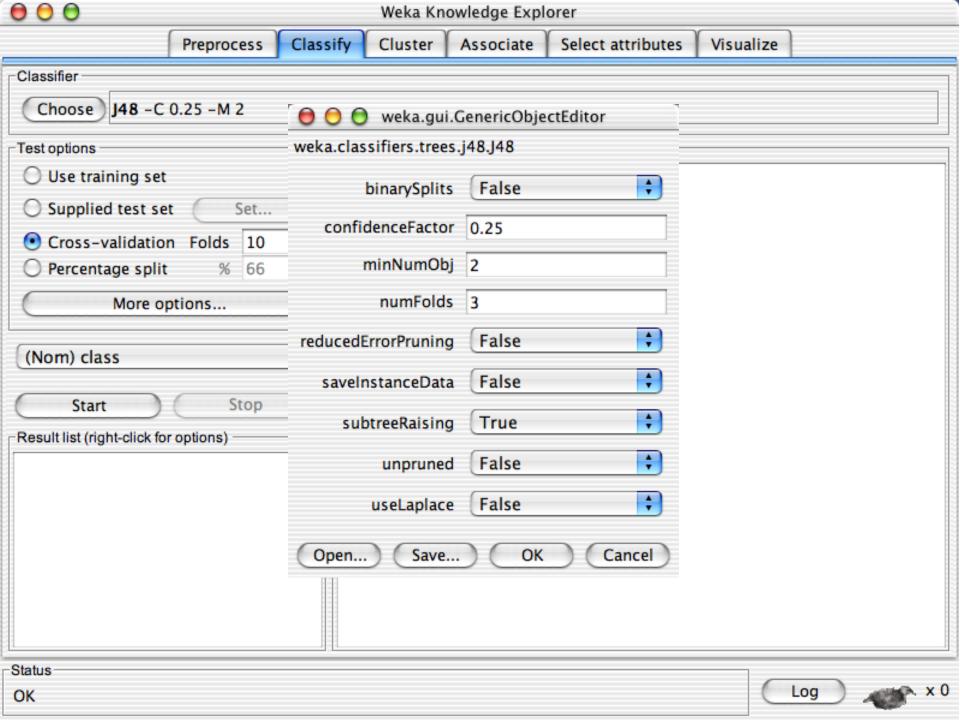


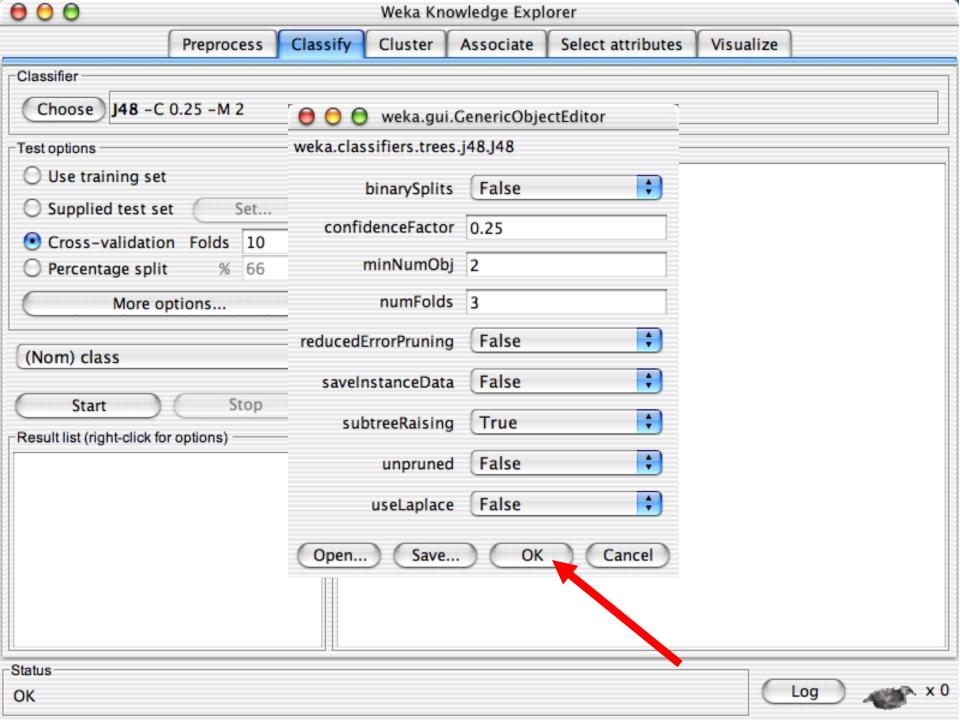


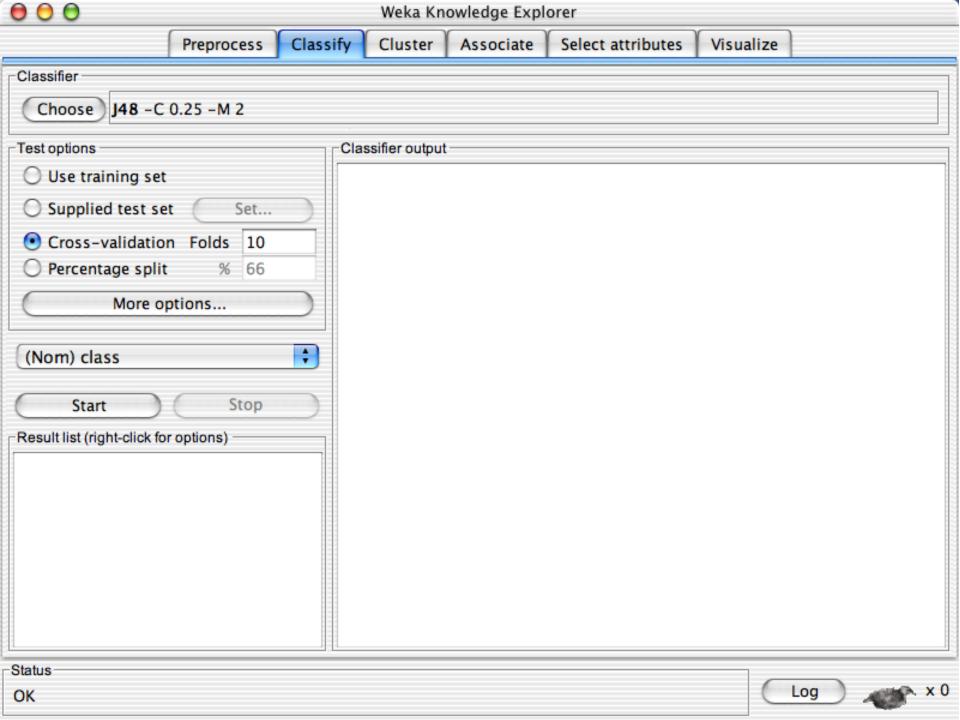


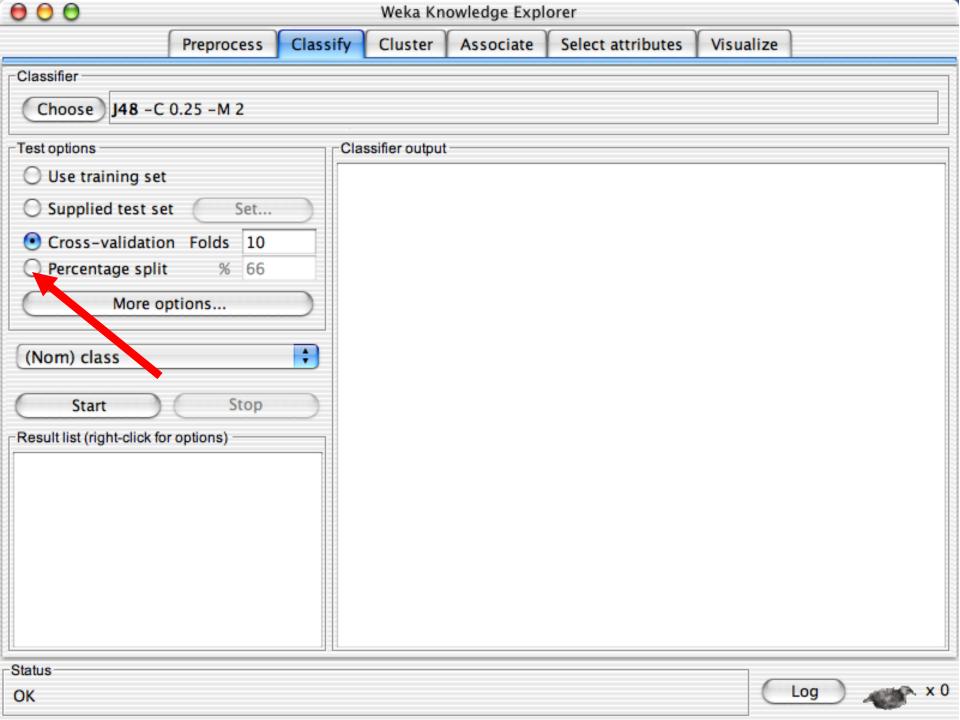


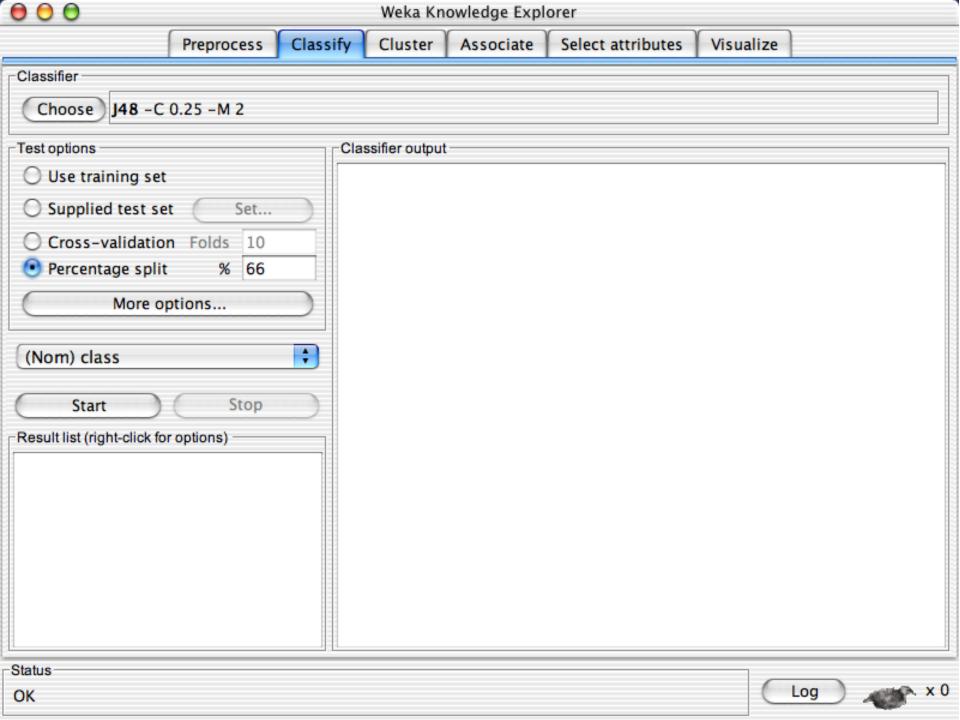


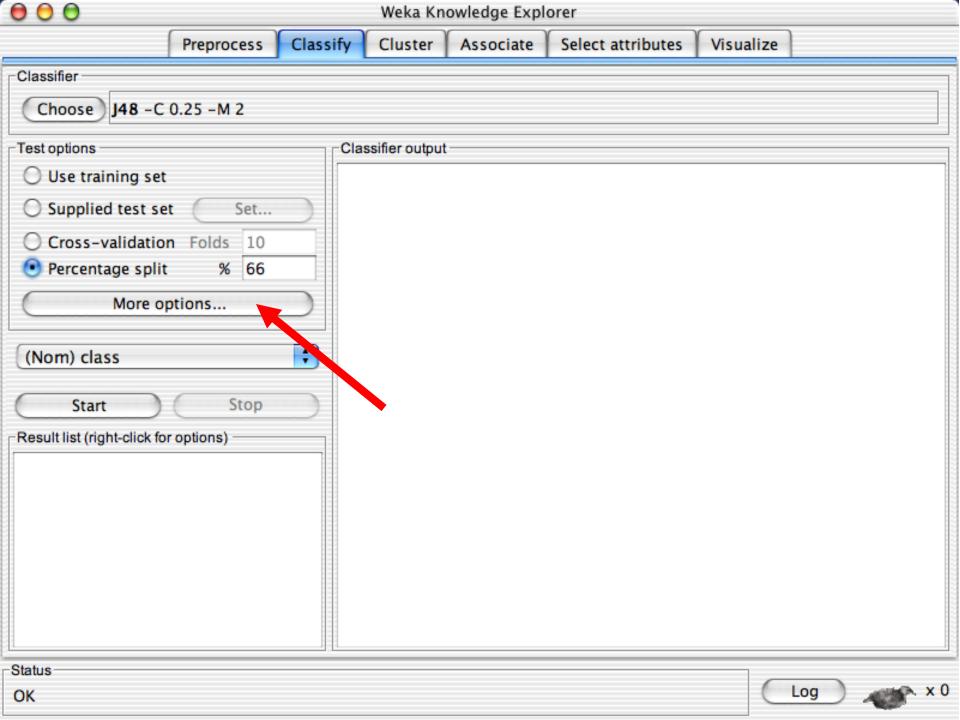


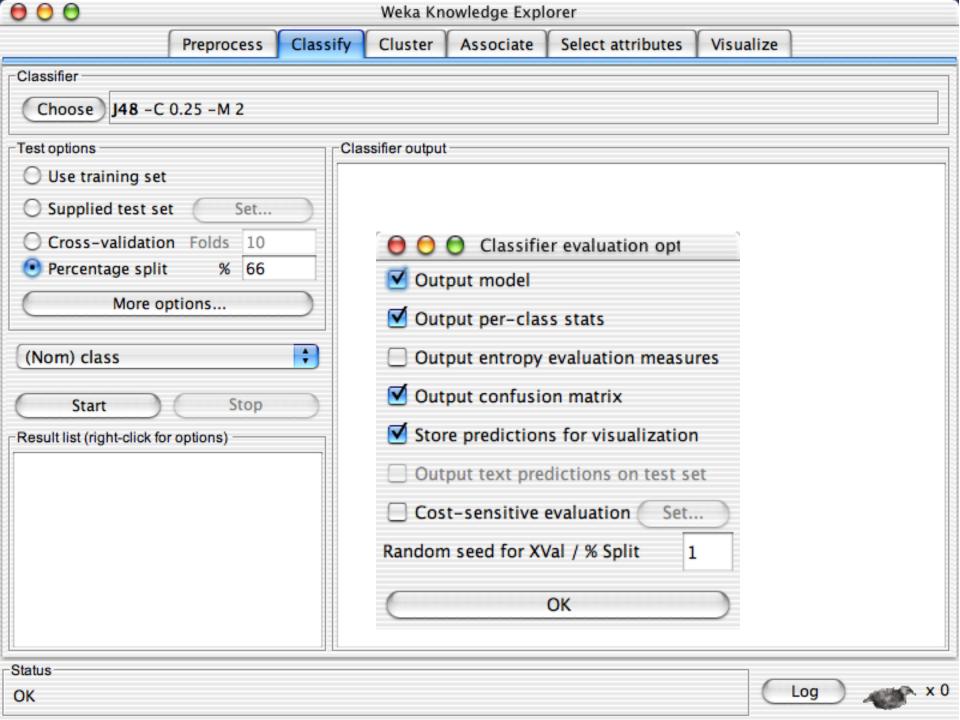


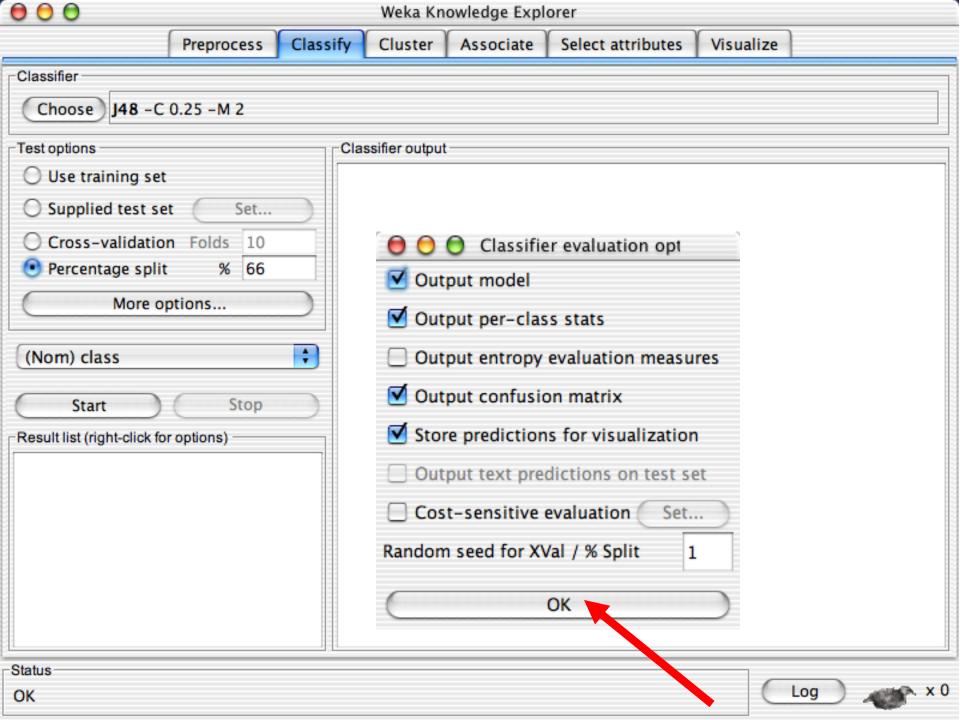


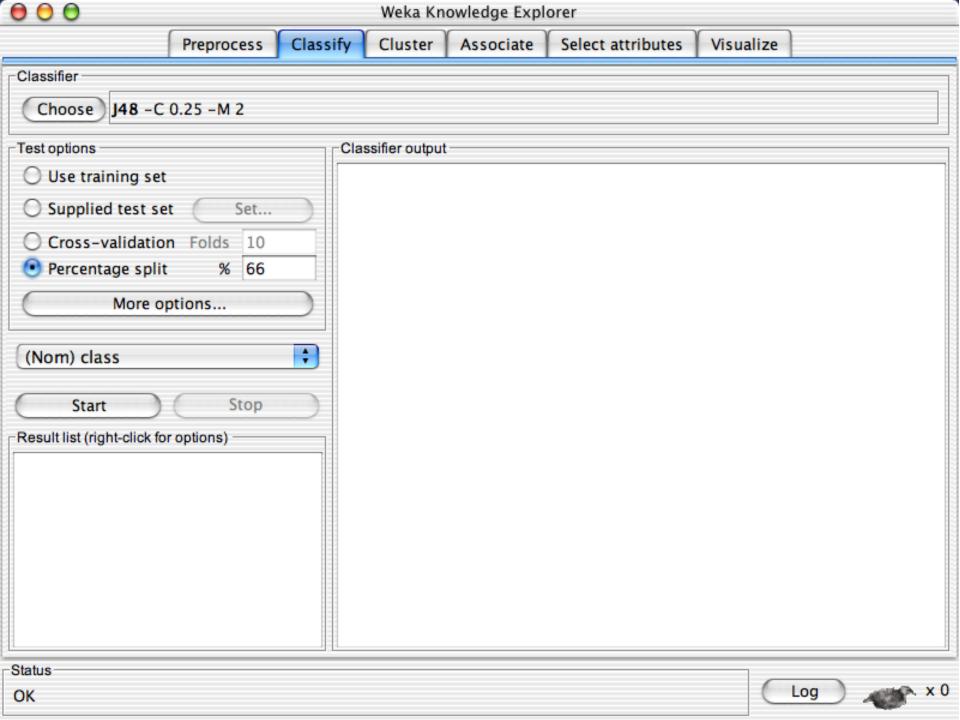


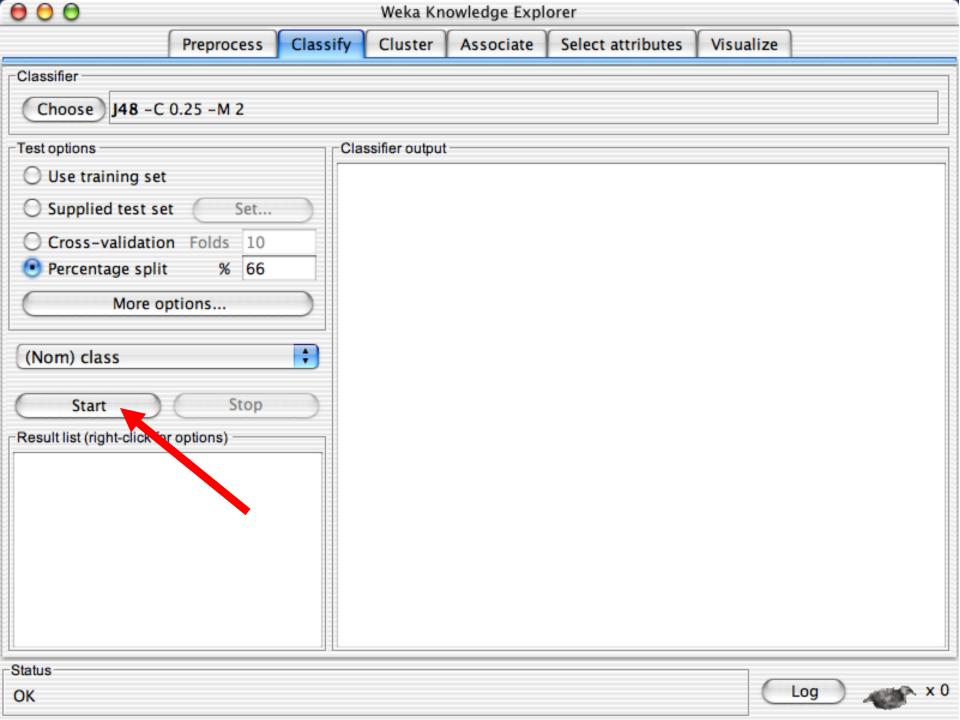


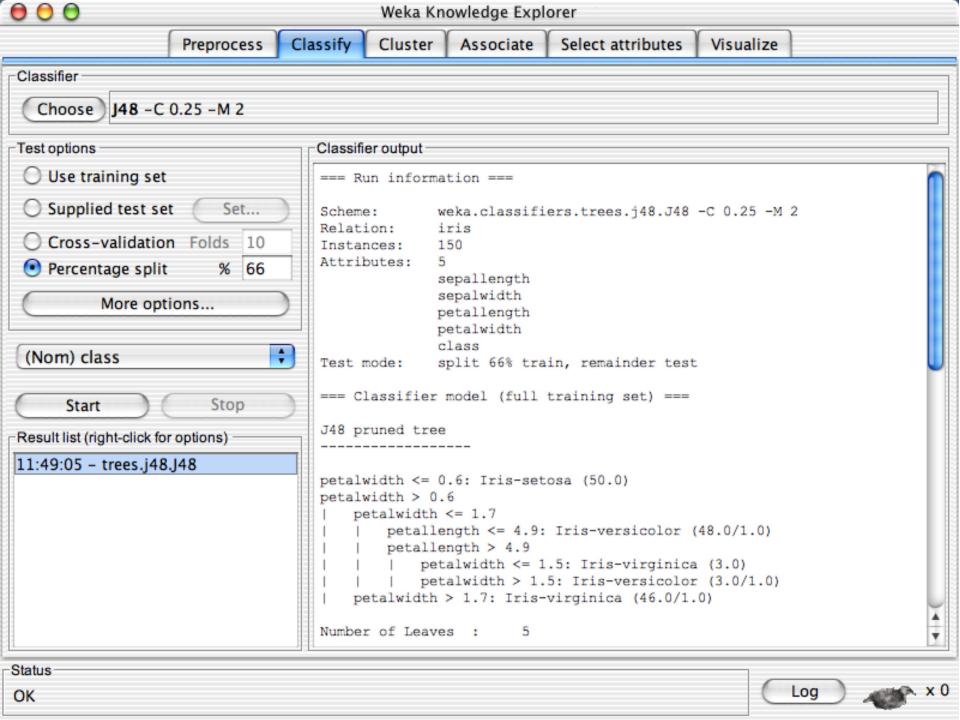


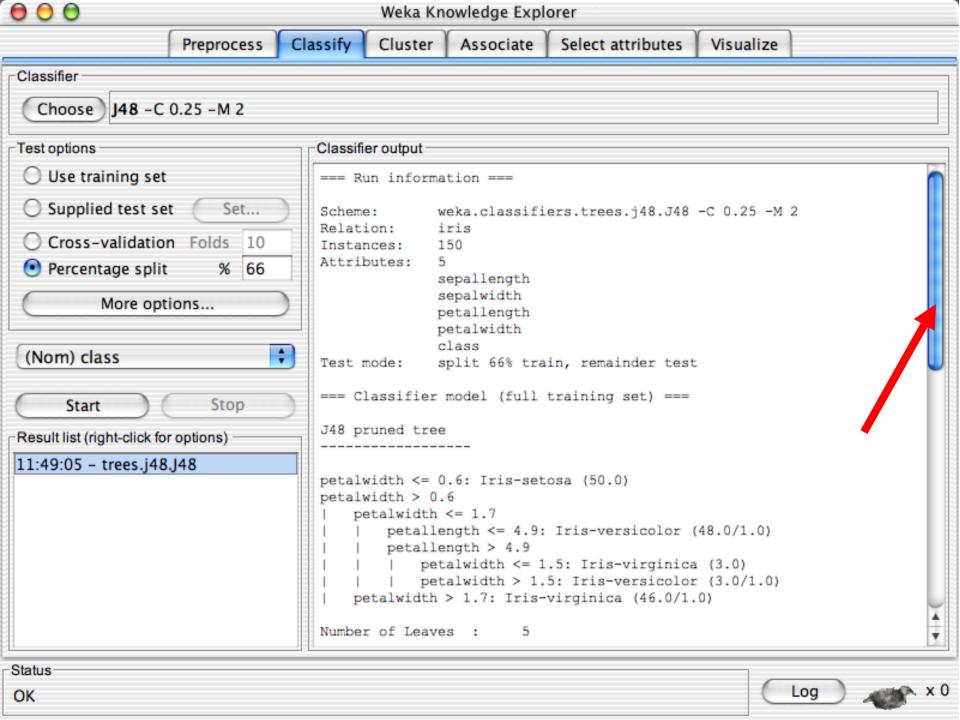


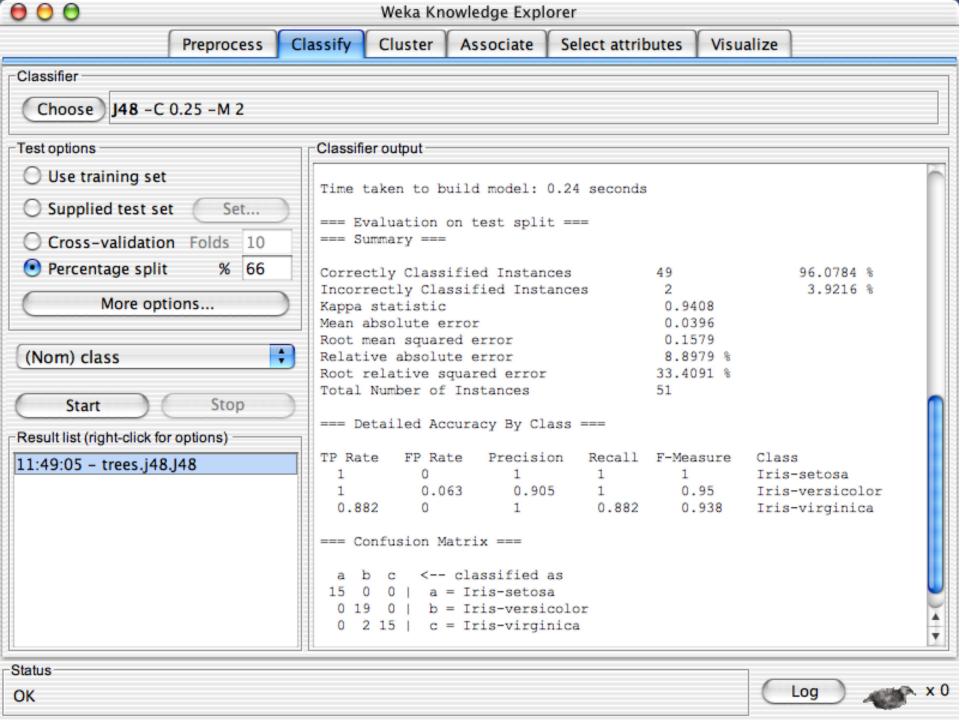


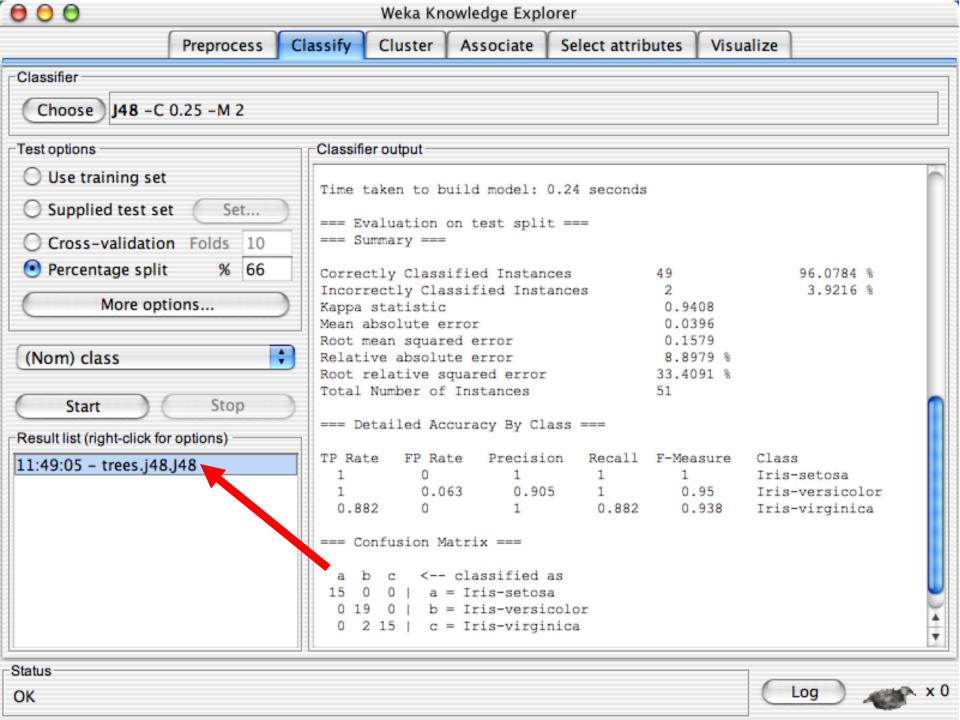


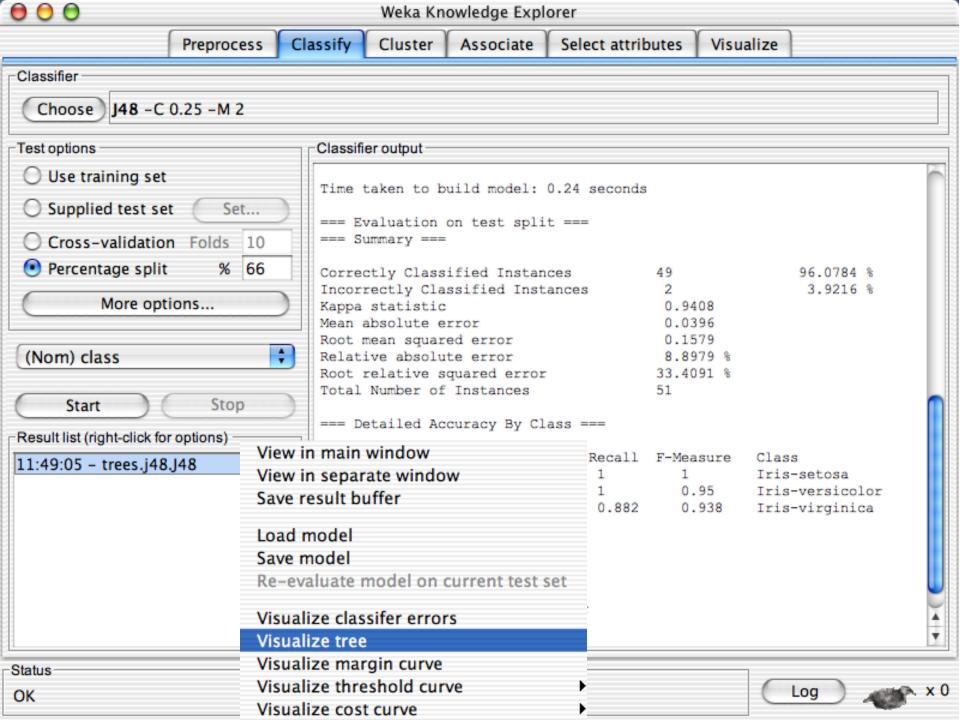


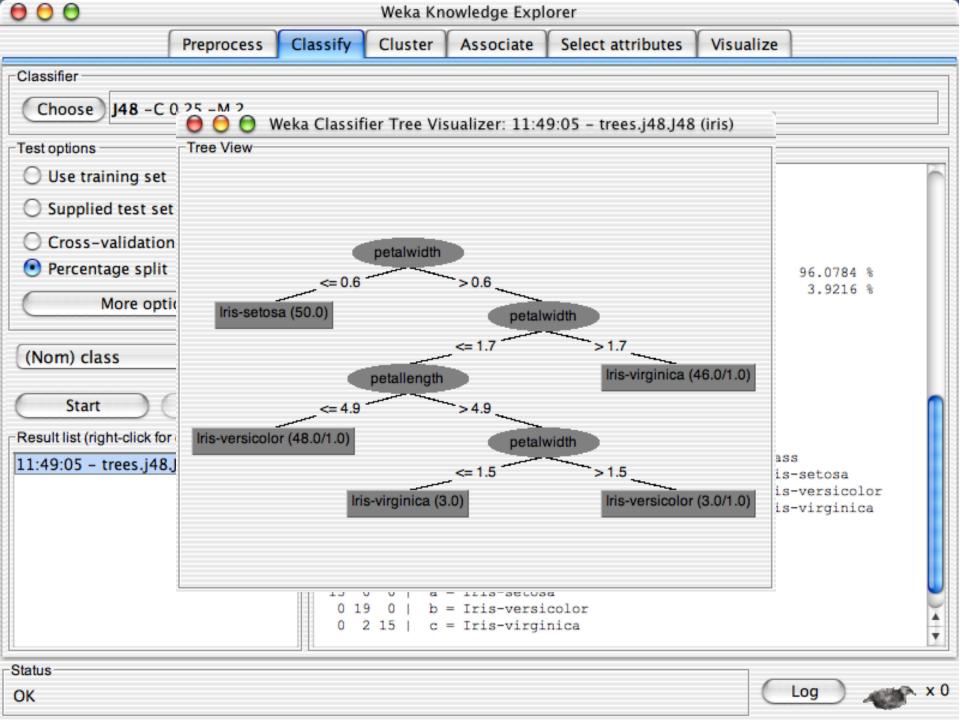


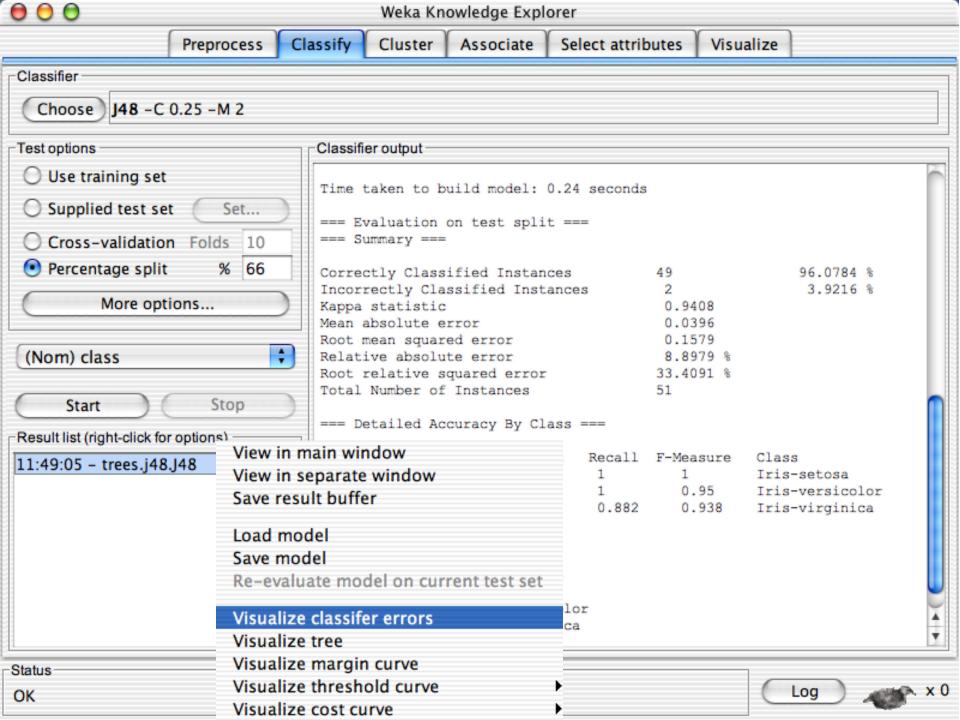


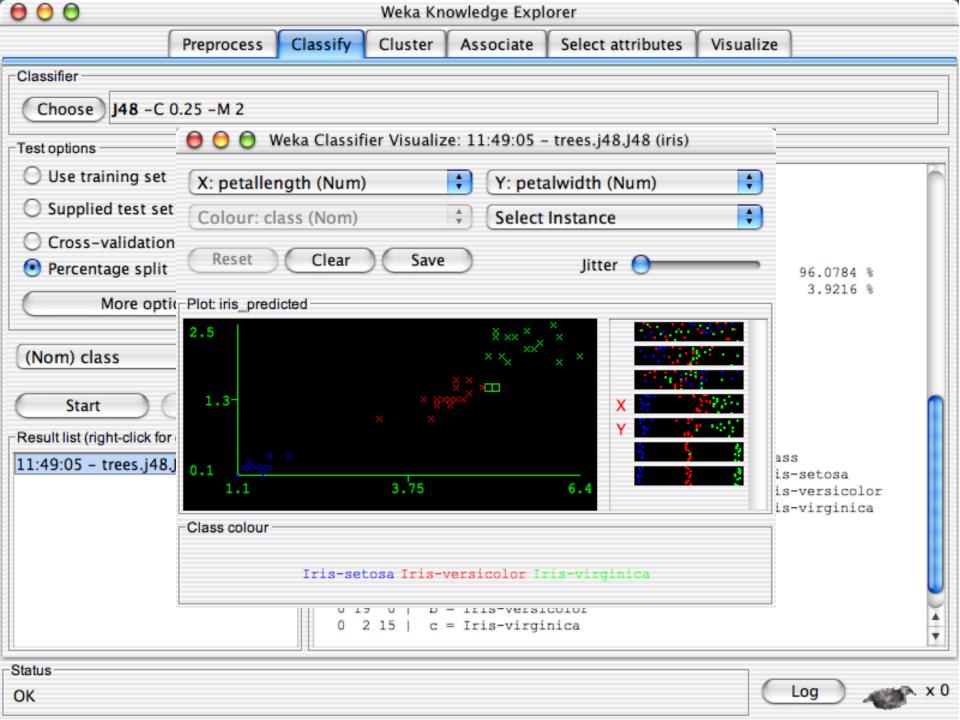


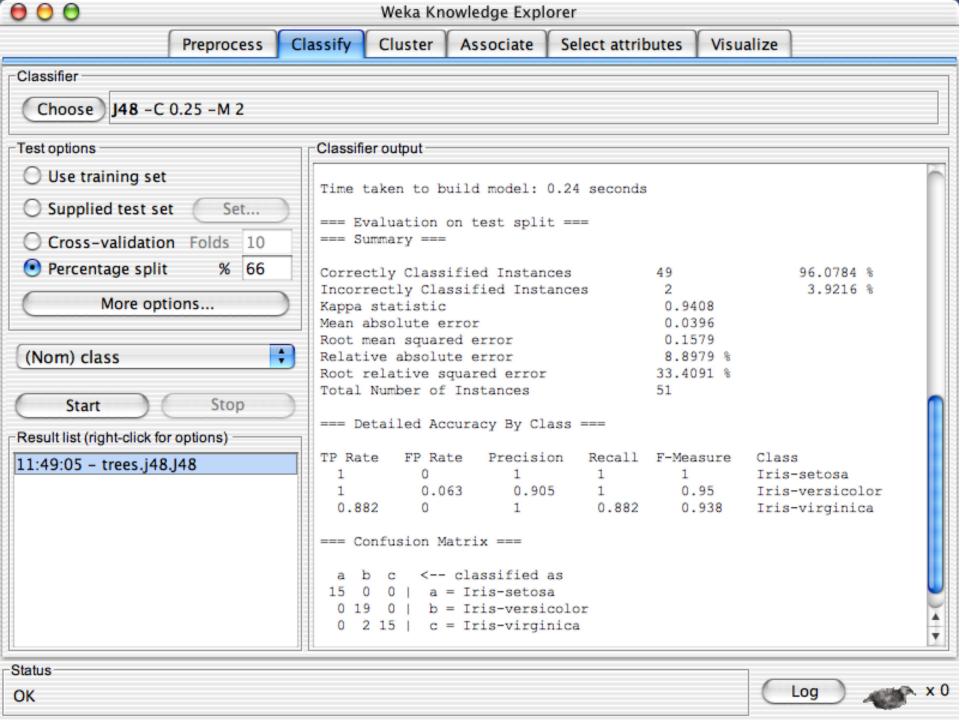


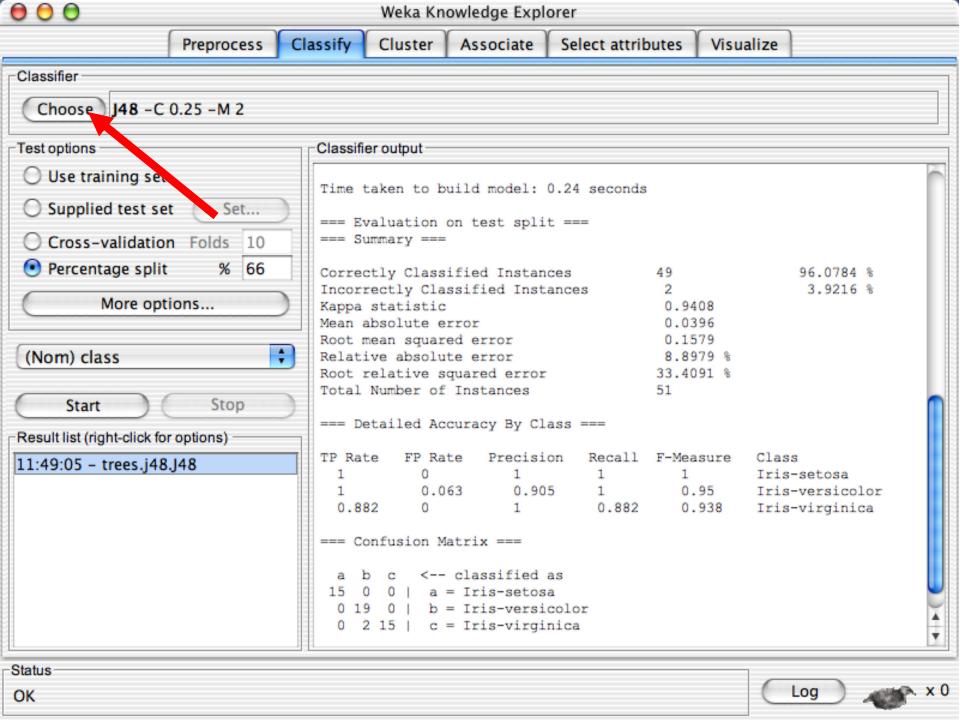


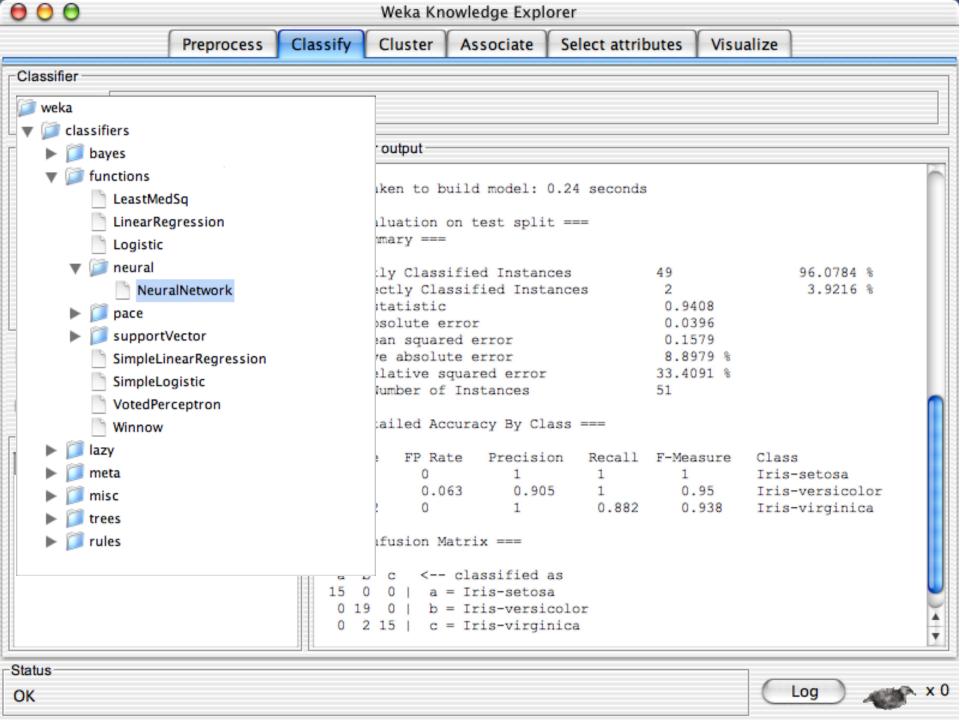


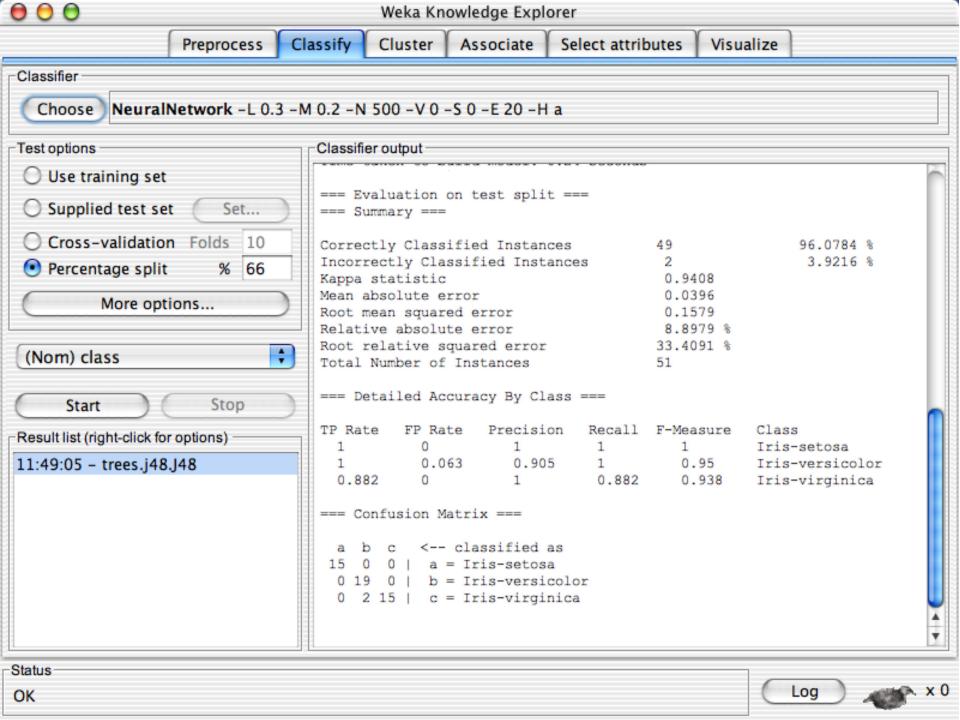


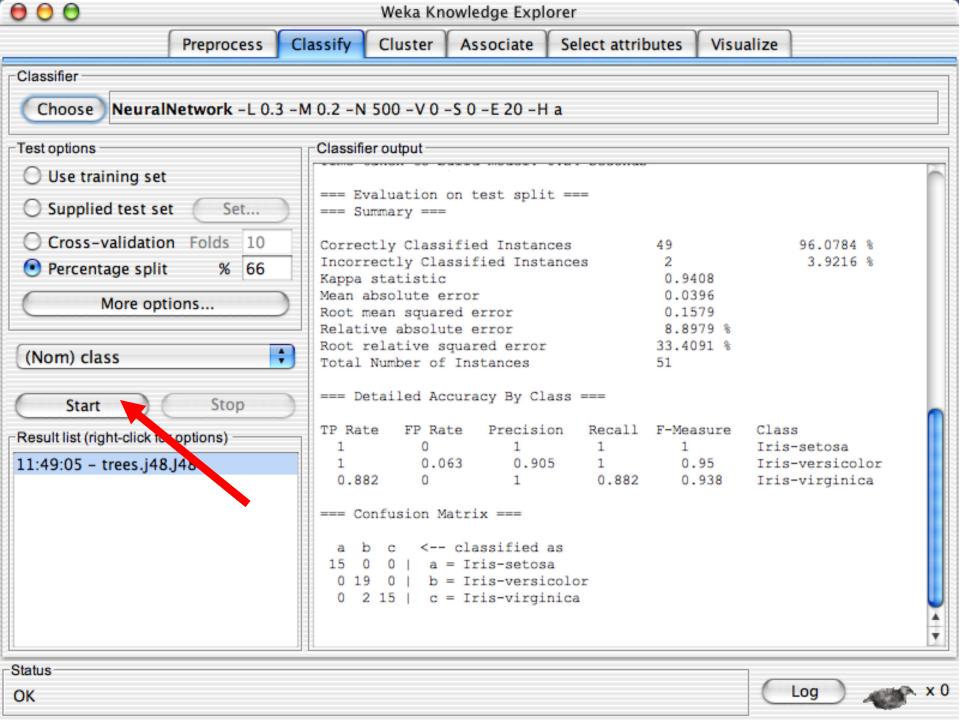


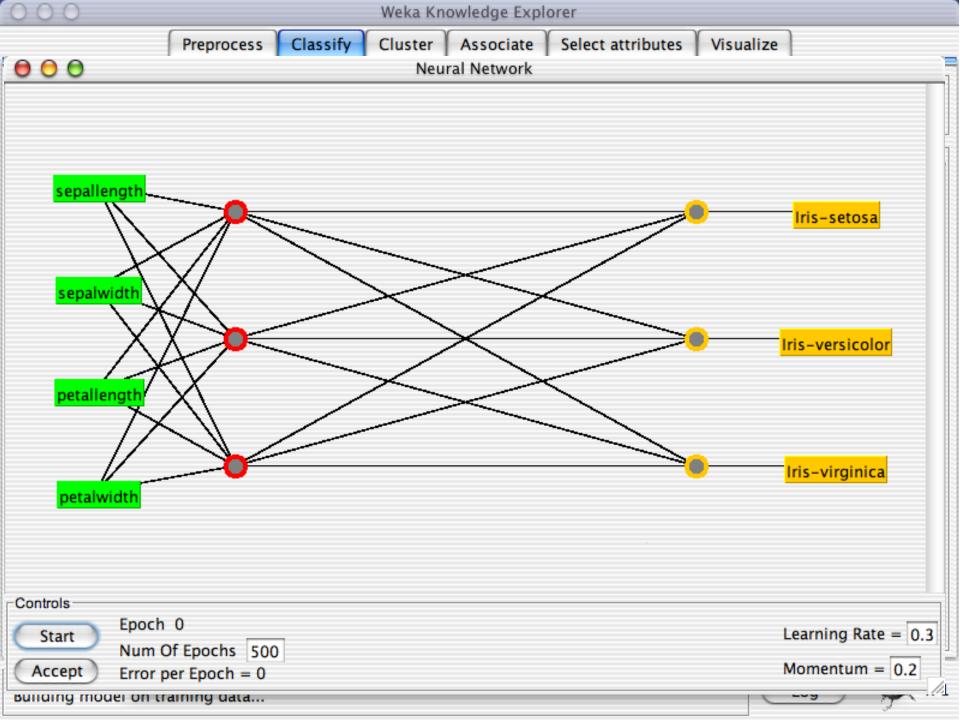


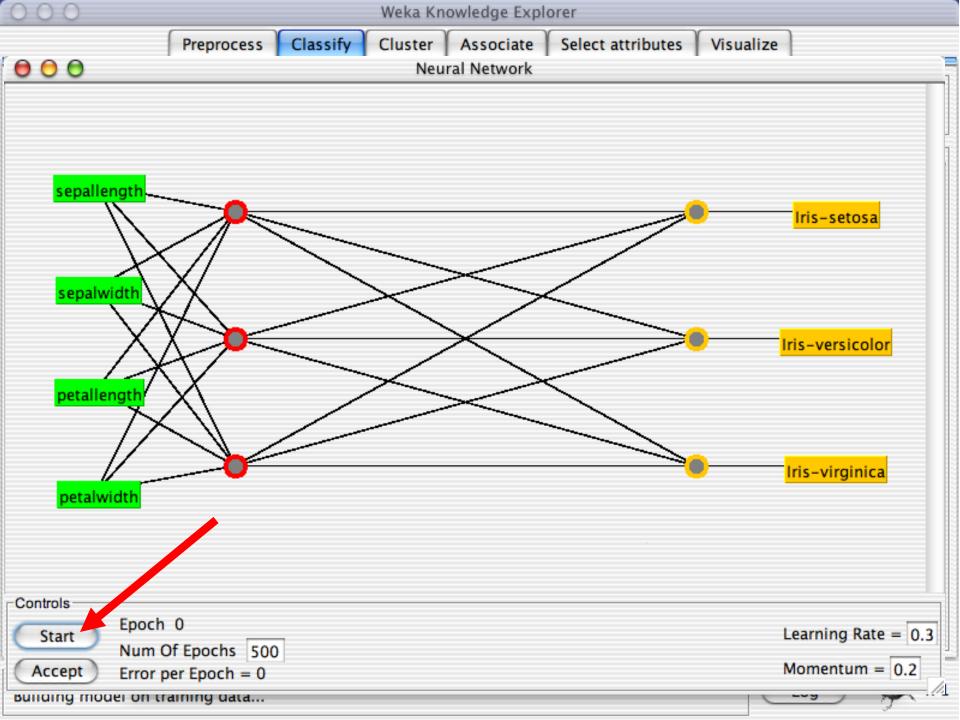


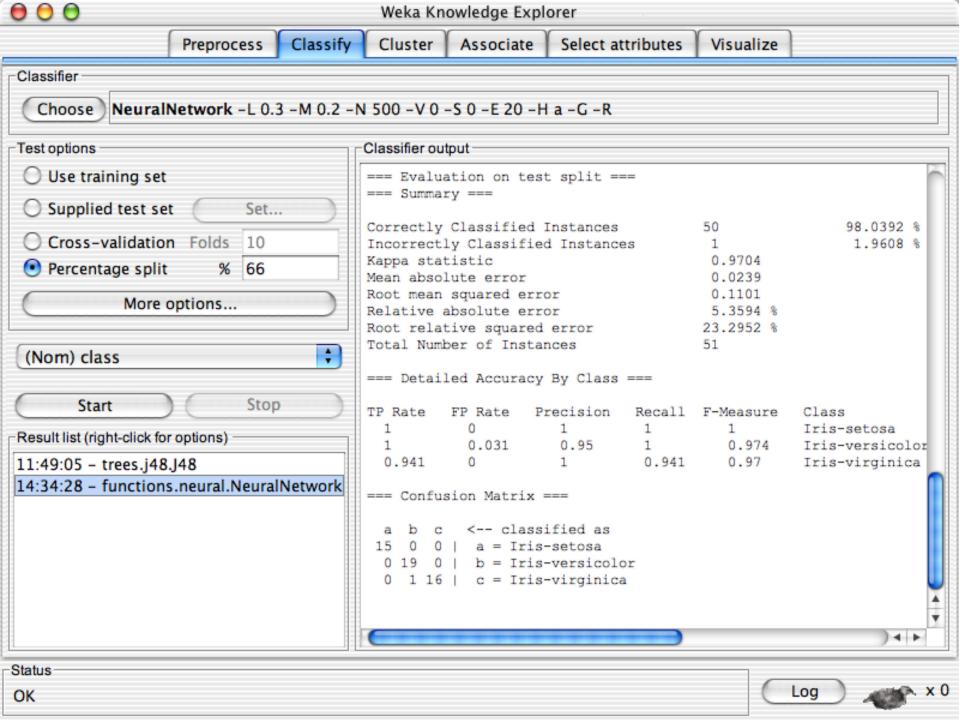


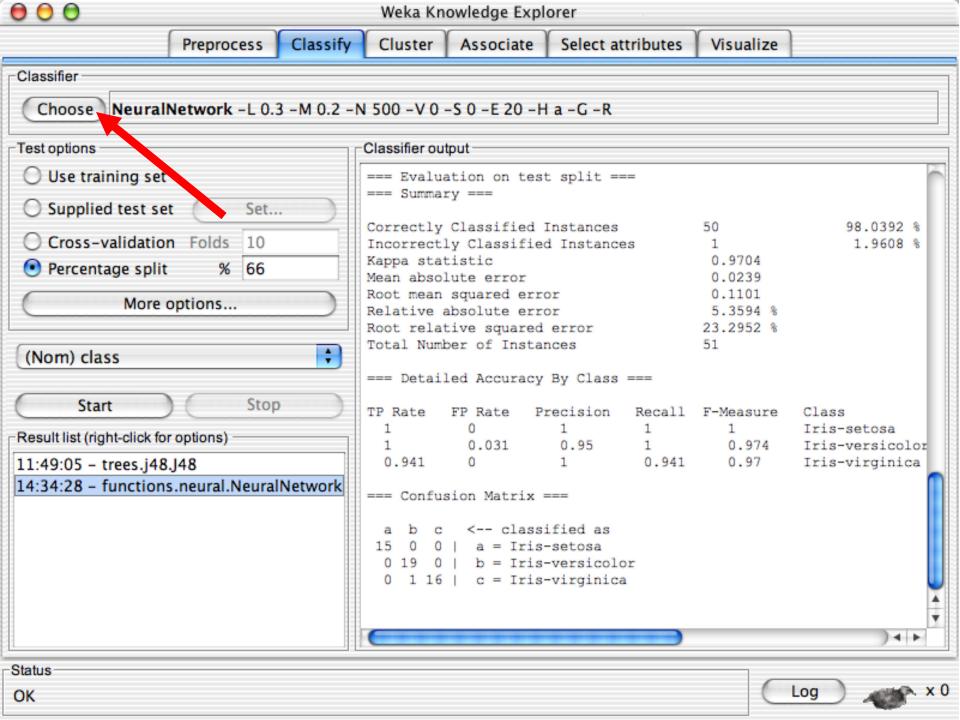


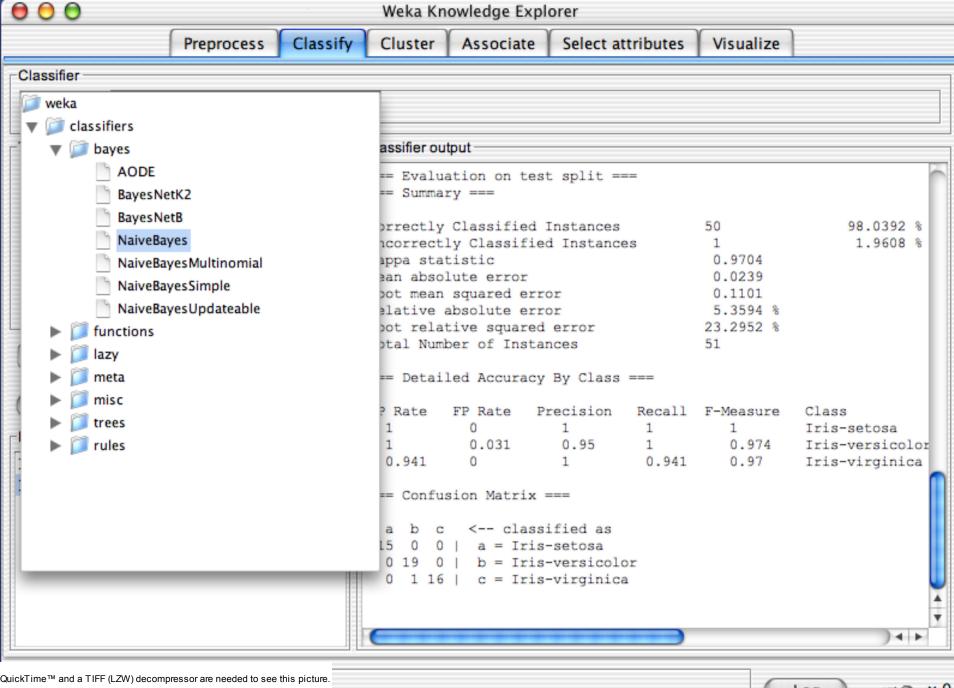


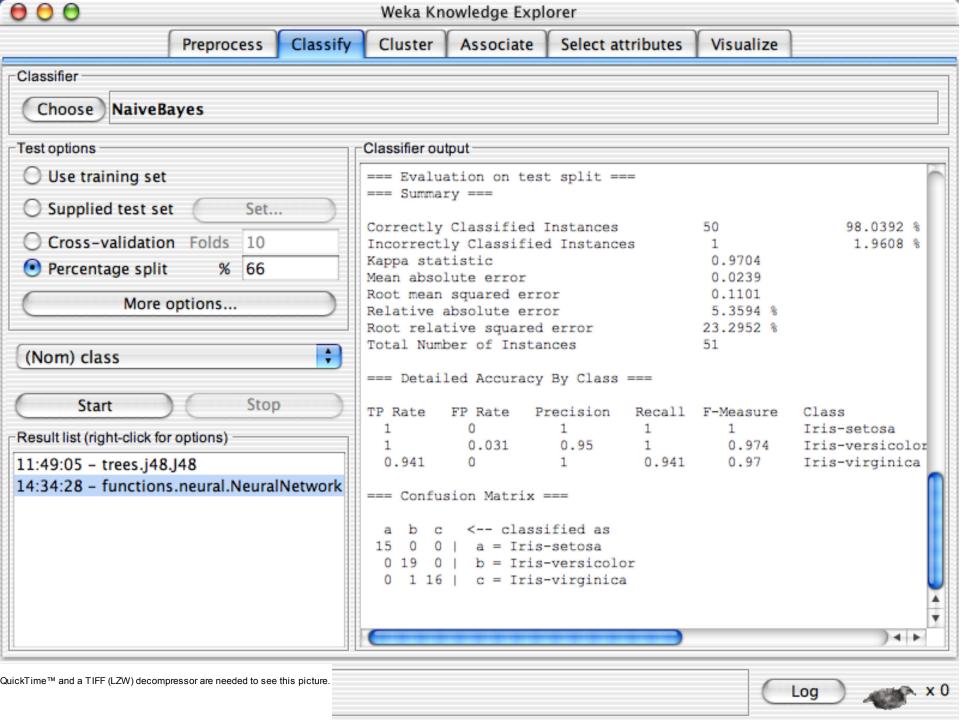


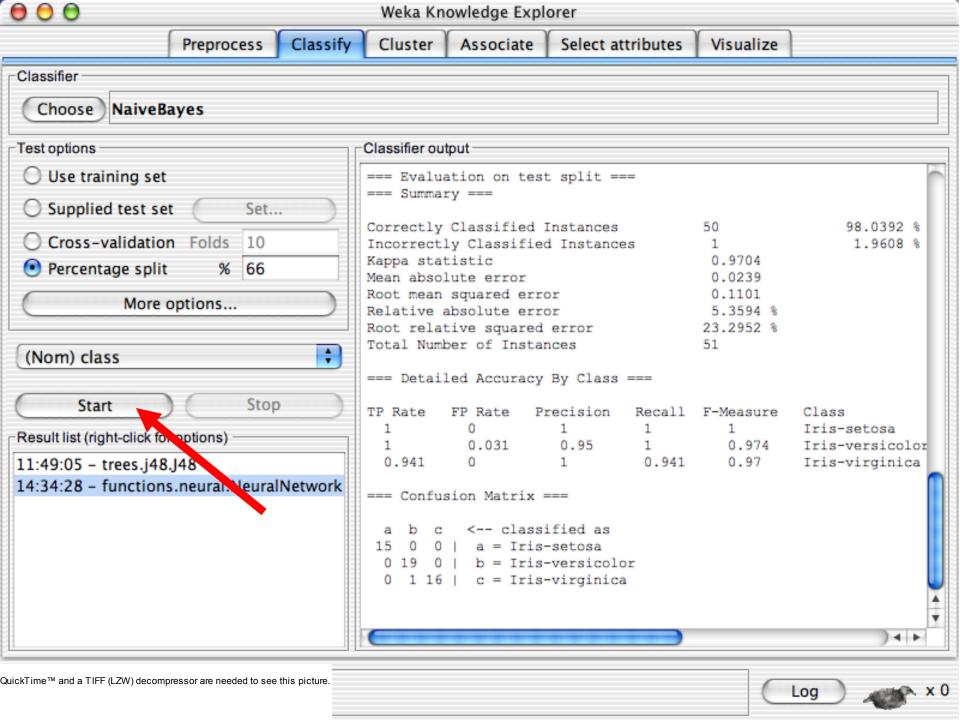


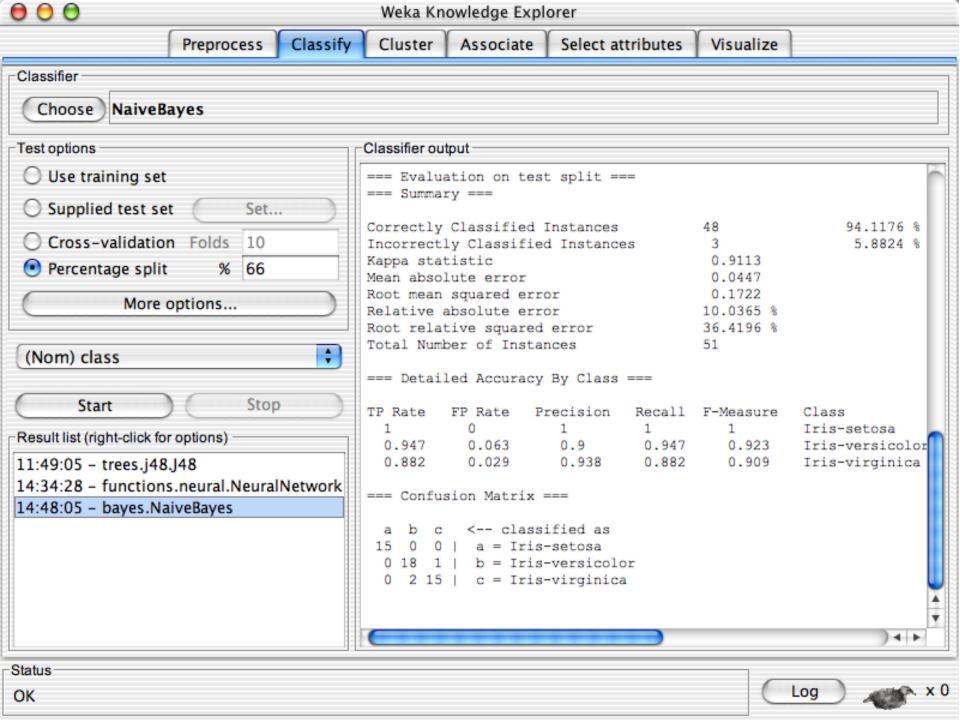


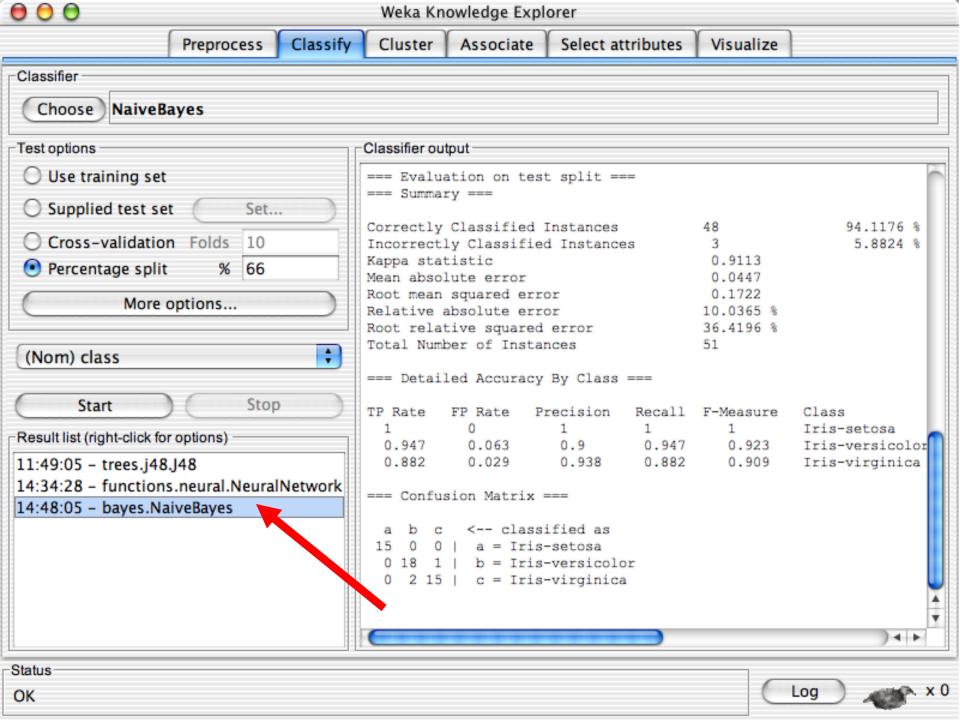


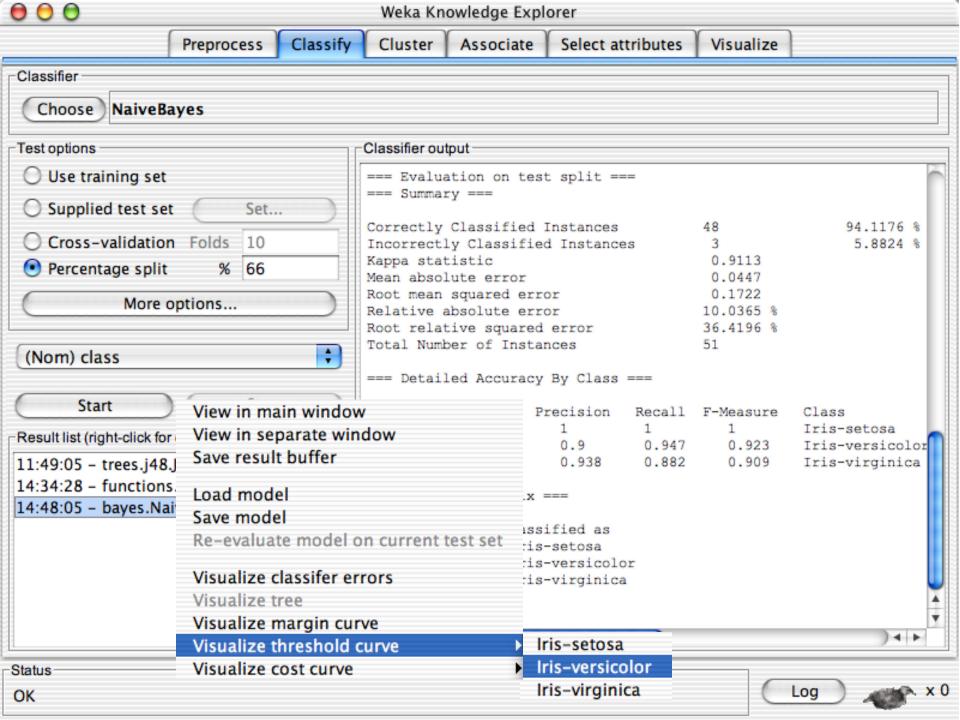


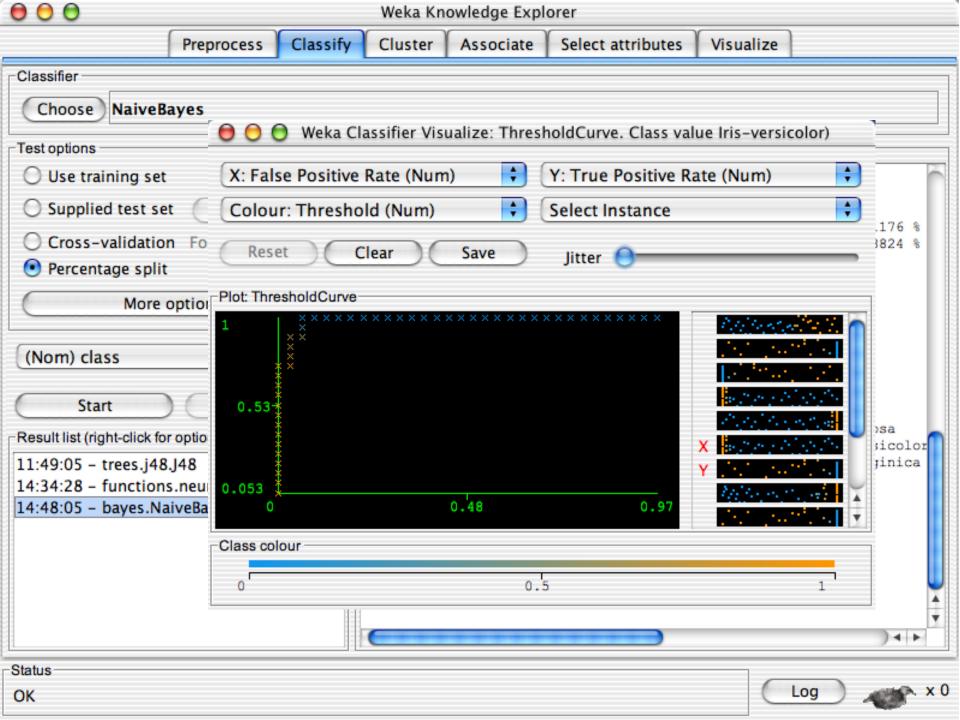


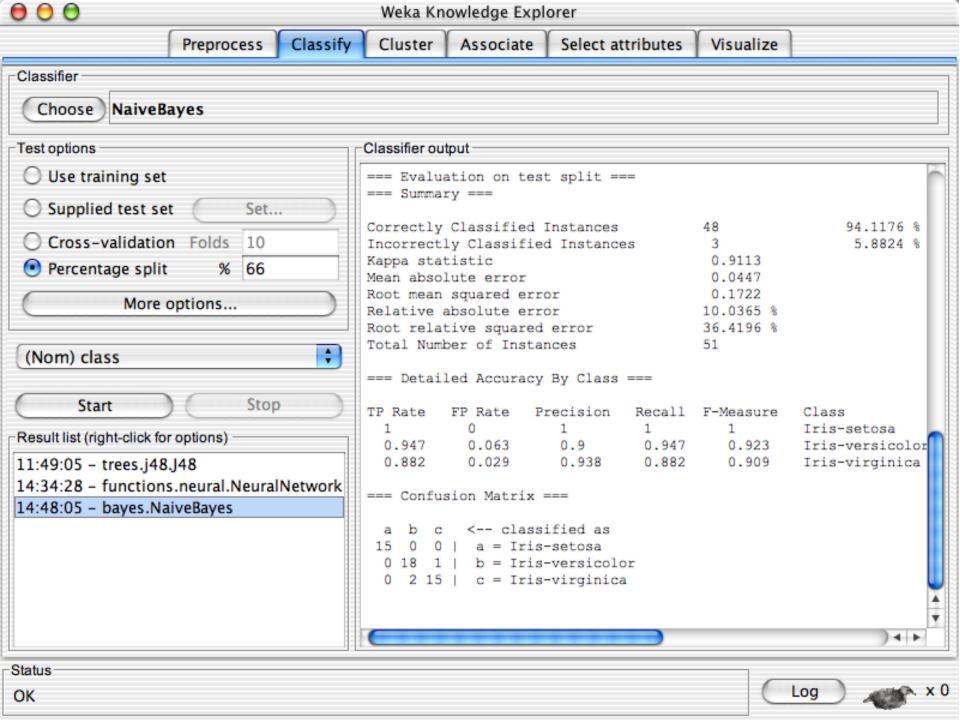


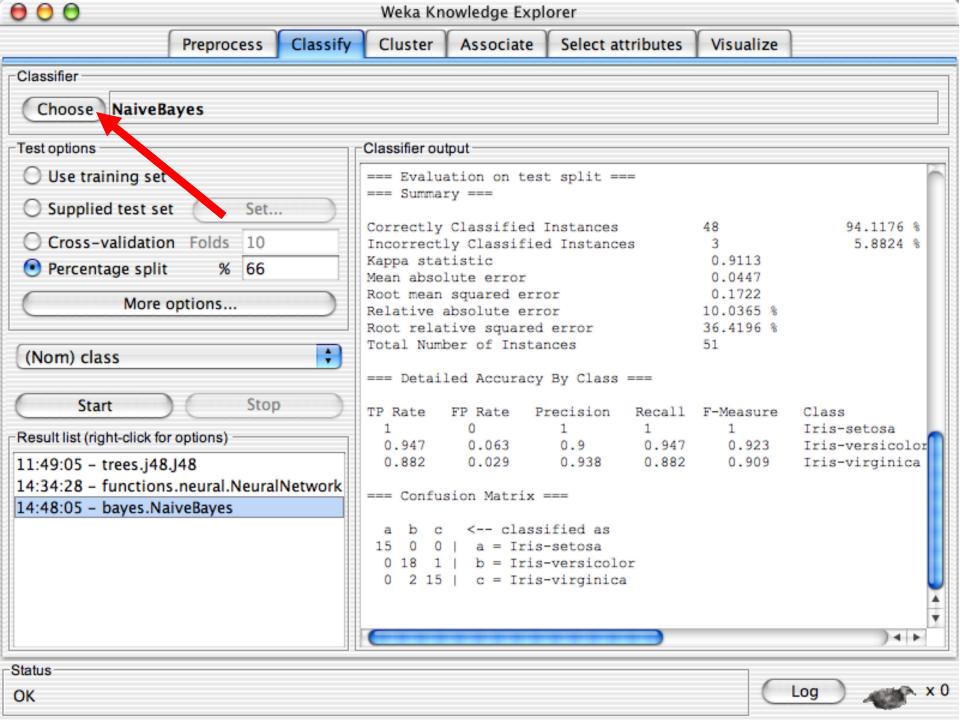


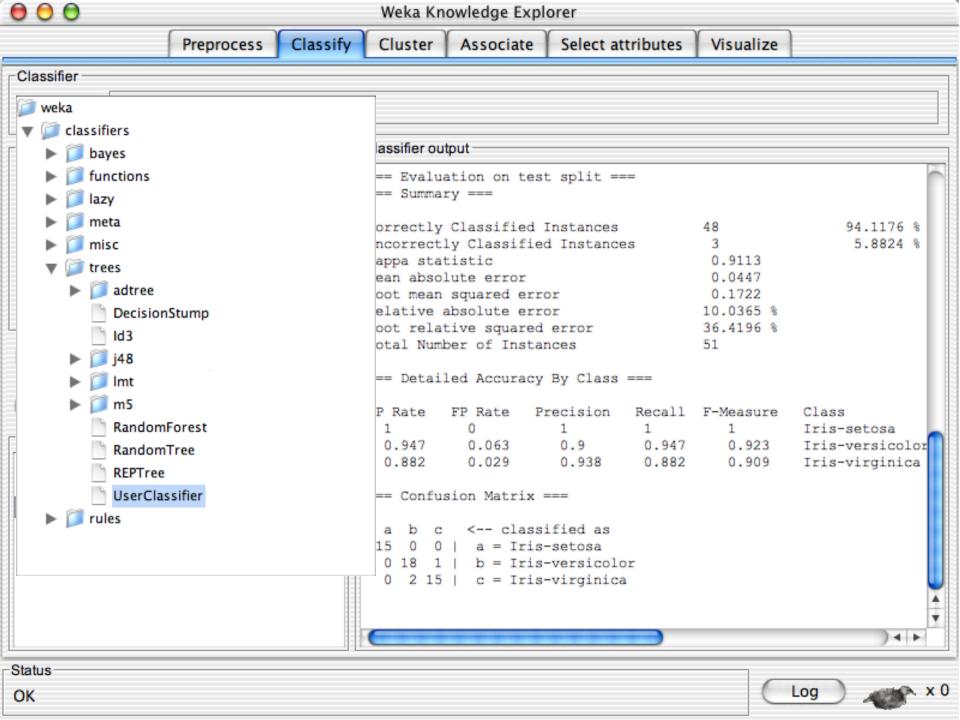


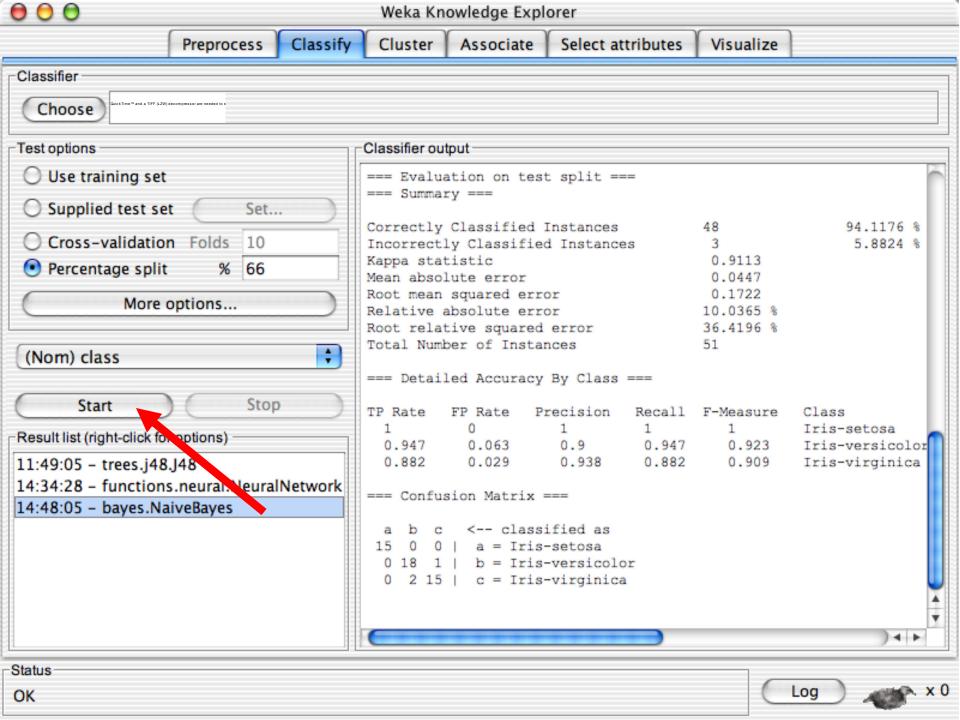


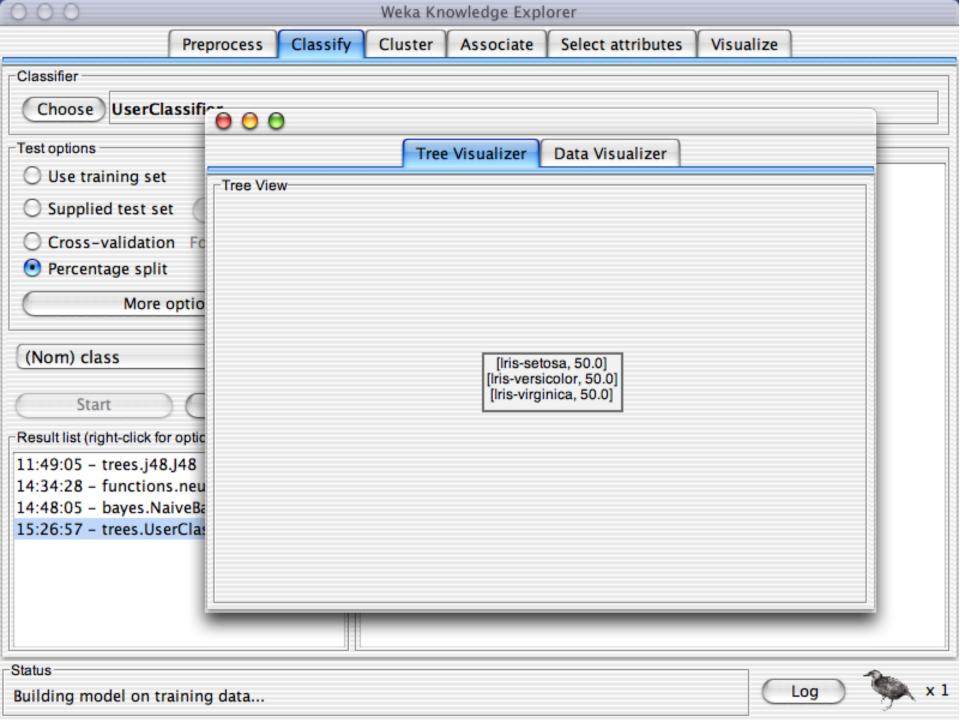


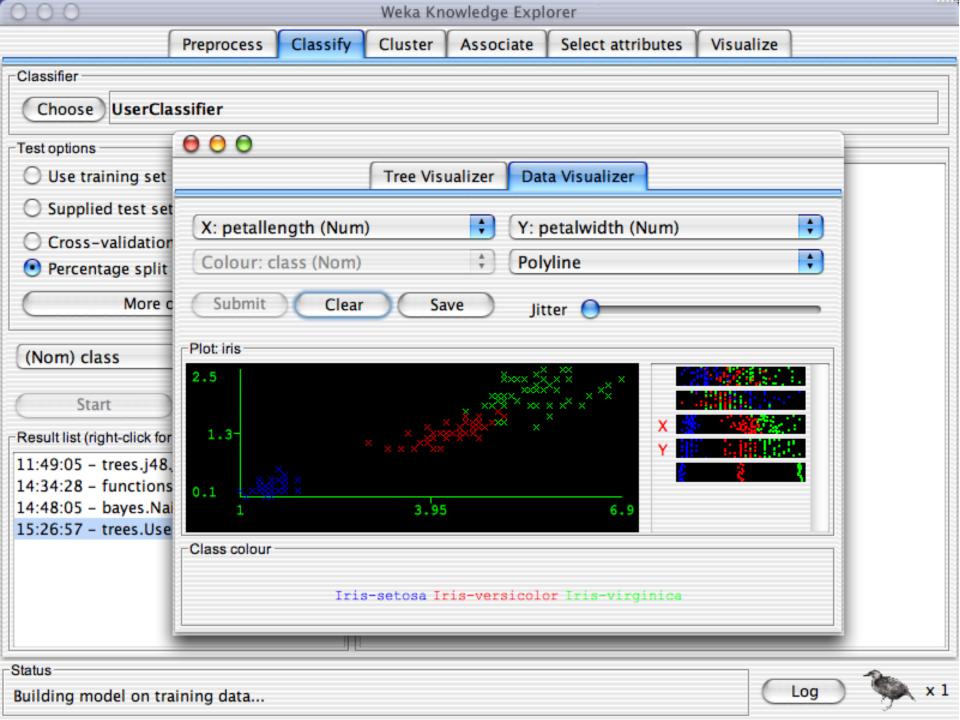


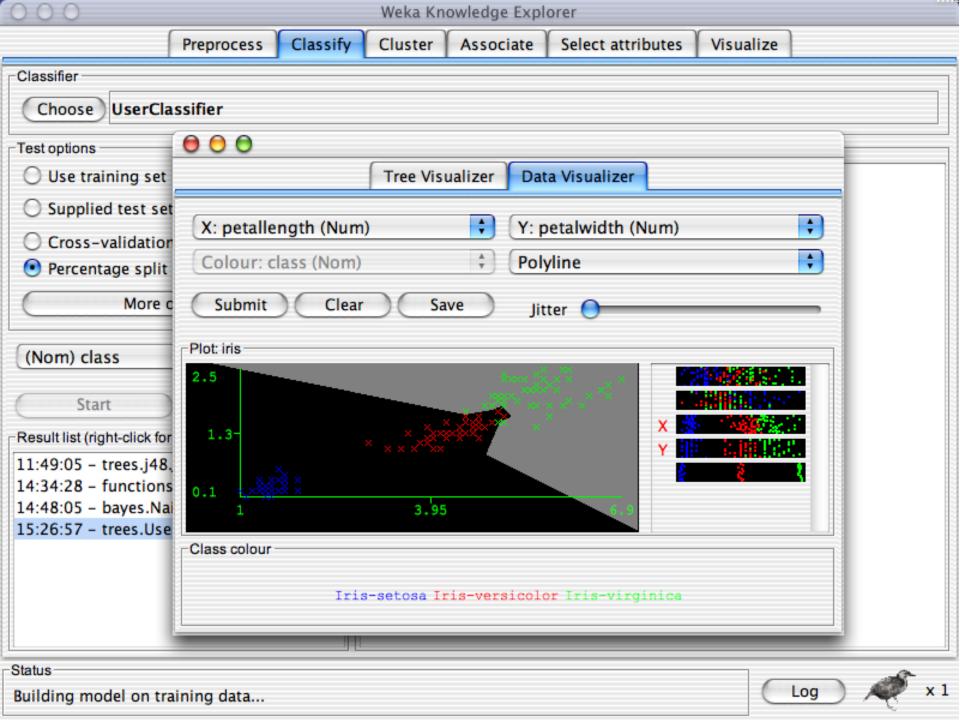


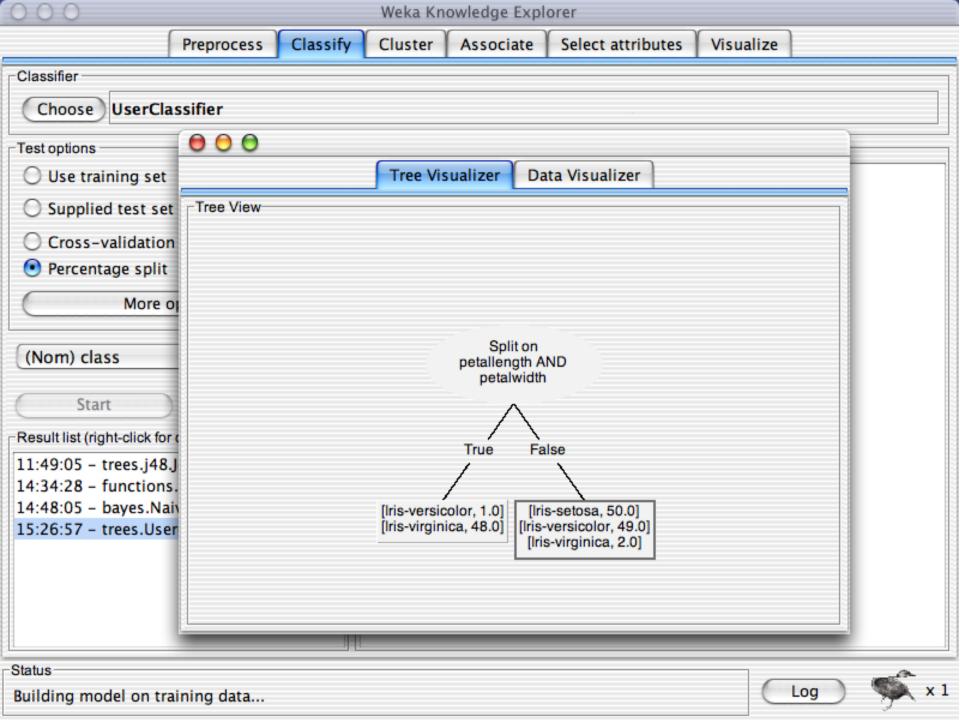


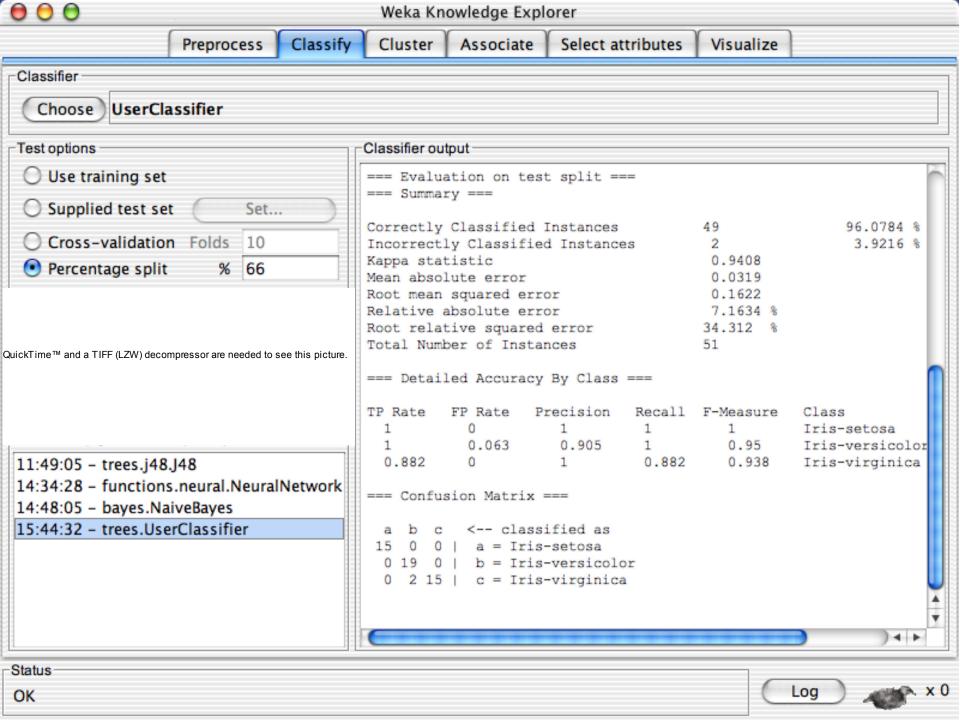


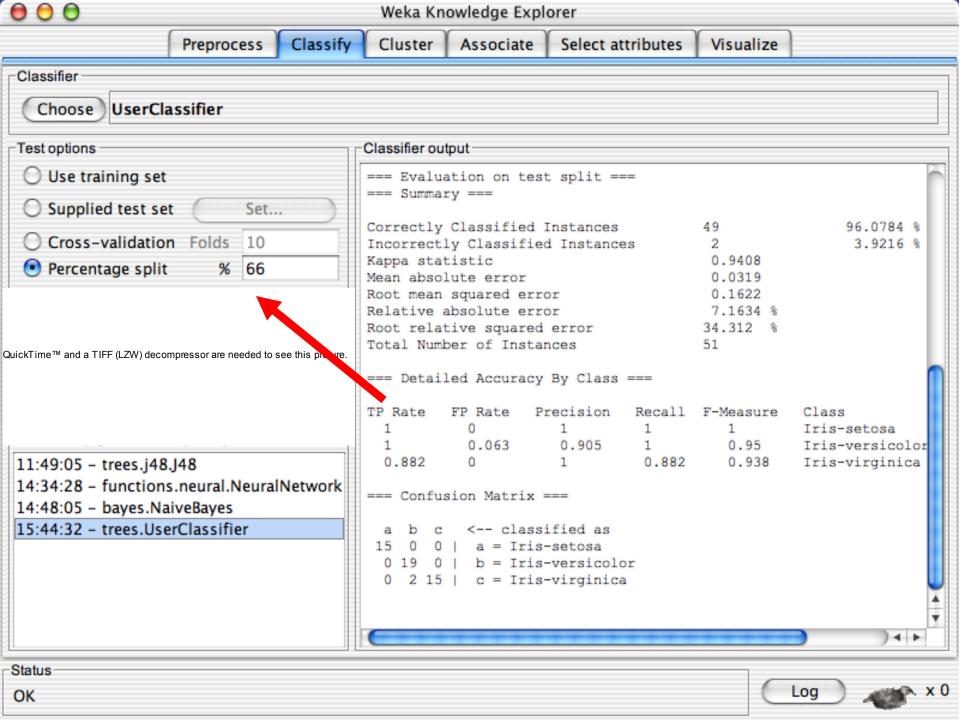


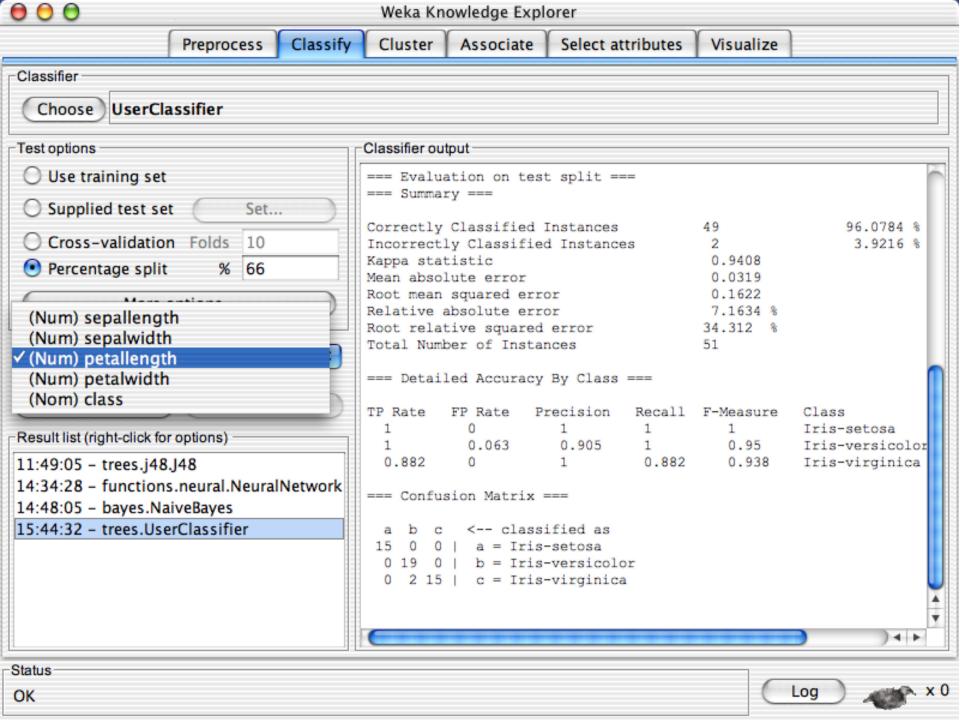






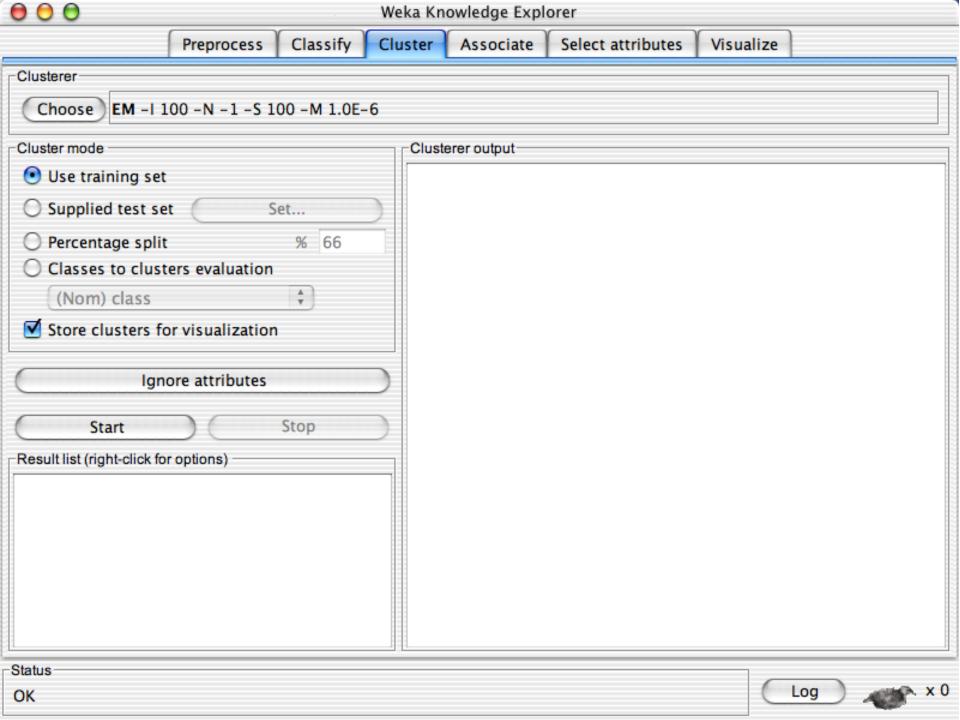


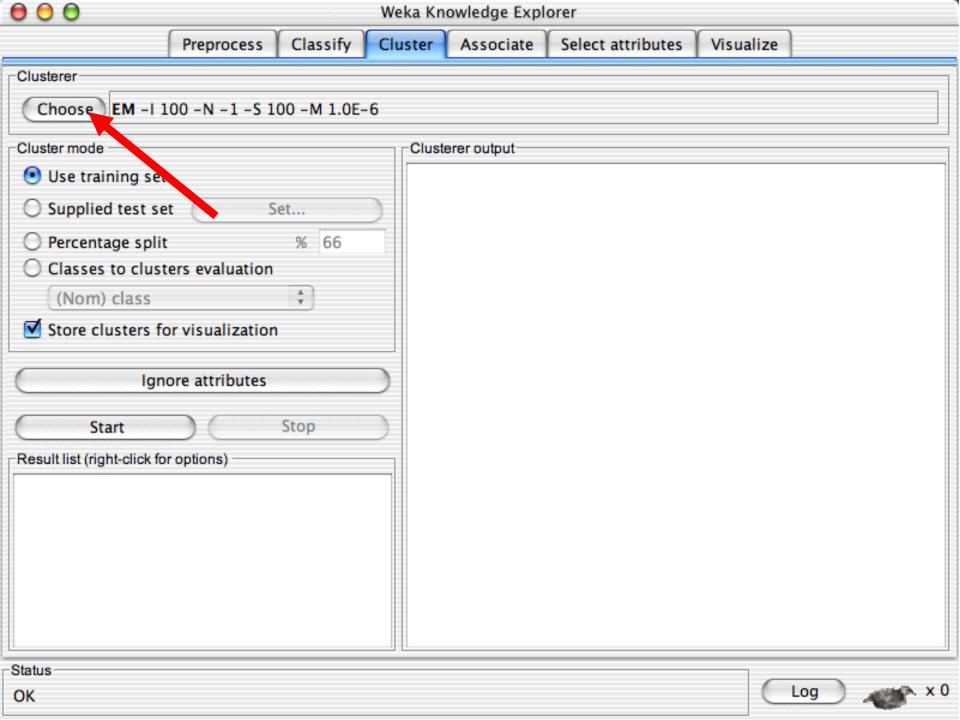


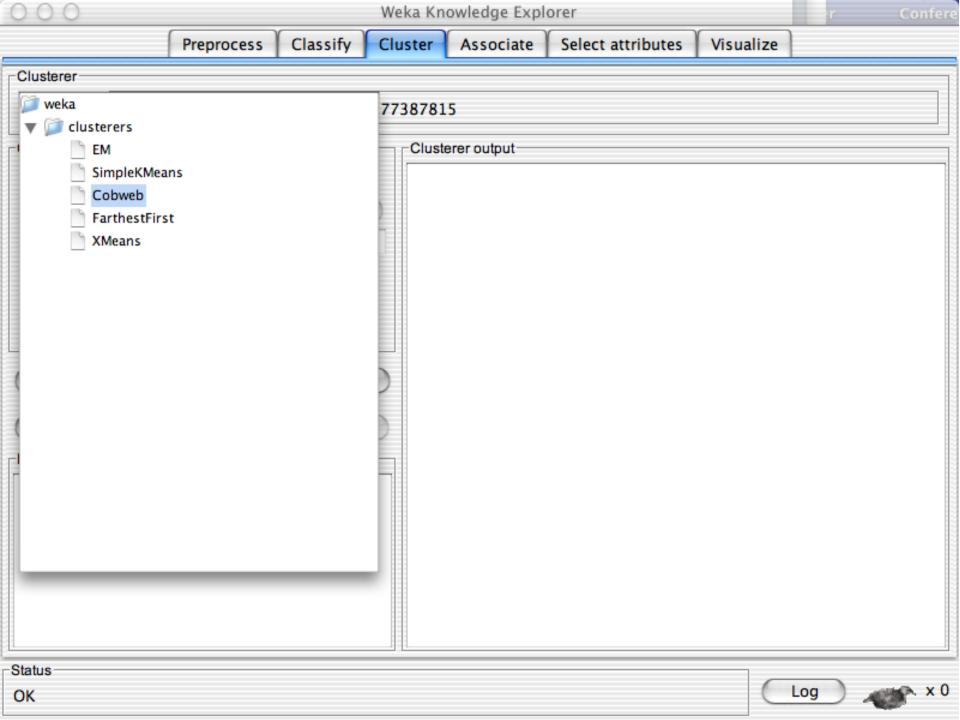


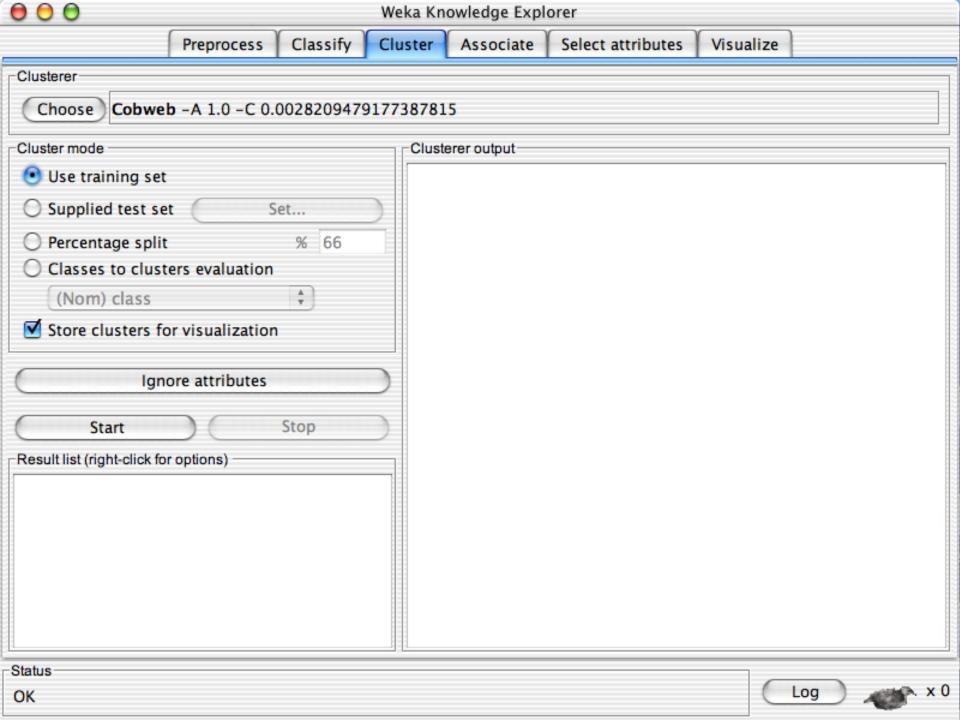
Explorar: clustering

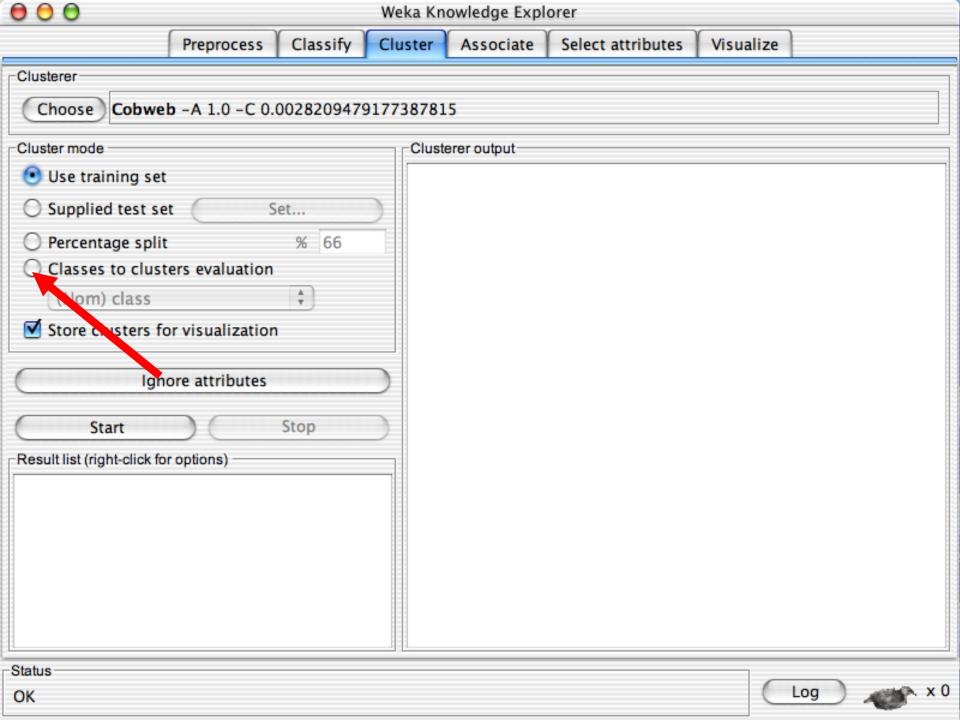
- WEKA contiene "clusterers" para buscar grupos simulares
- Algunos algoritmos:
 - k-Means, EM, Cobweb, X-means, FarthestFirst
- Clusters pueden ser visualizados y comparados con el correcto cluster (si es dado)

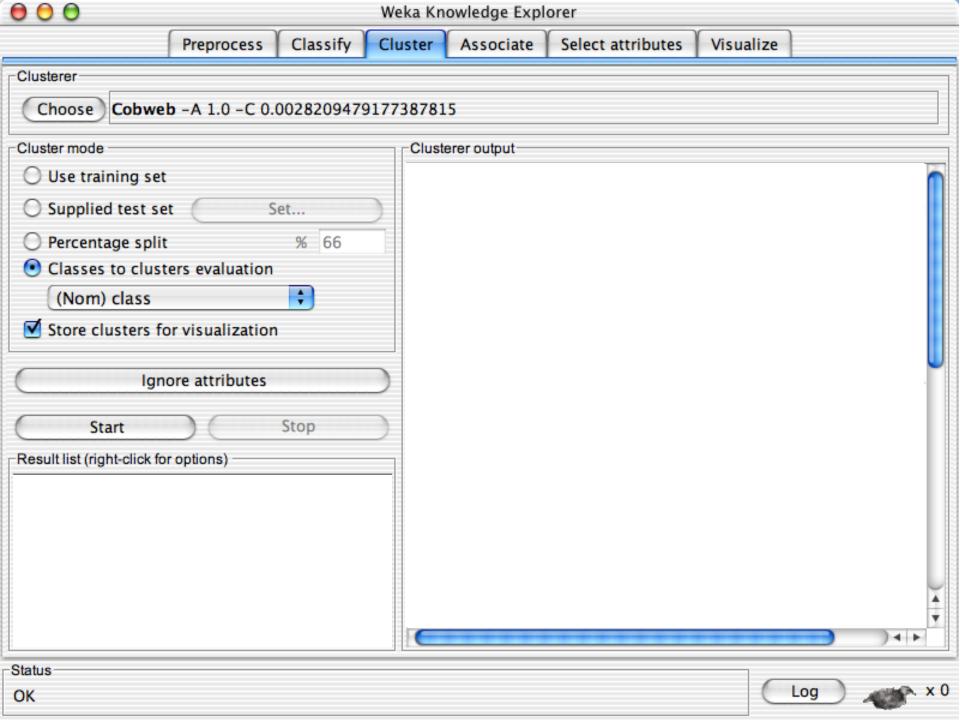


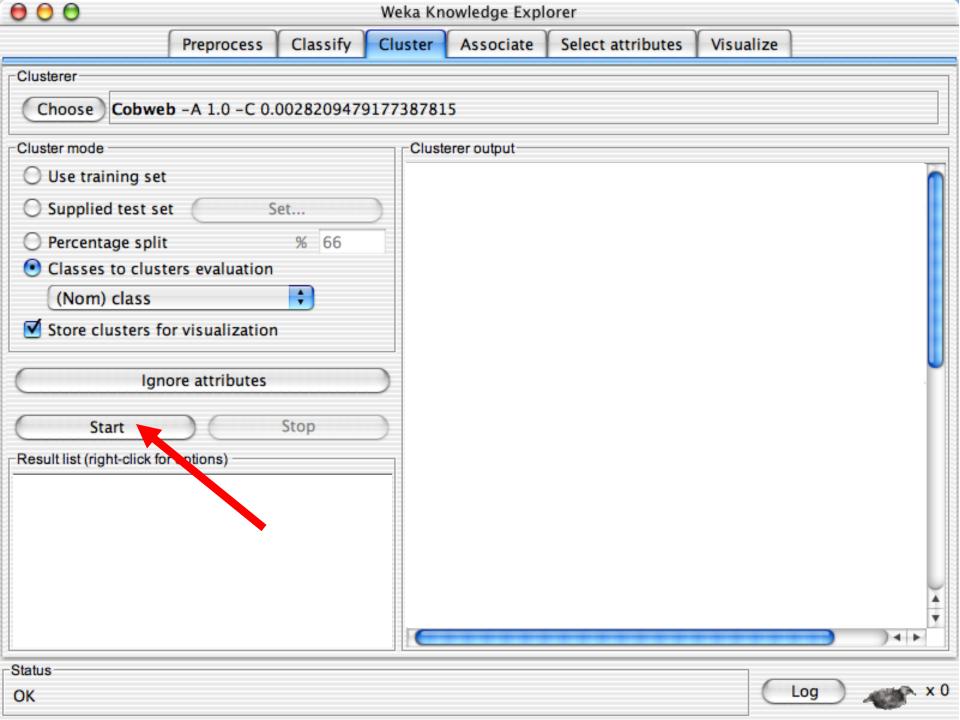


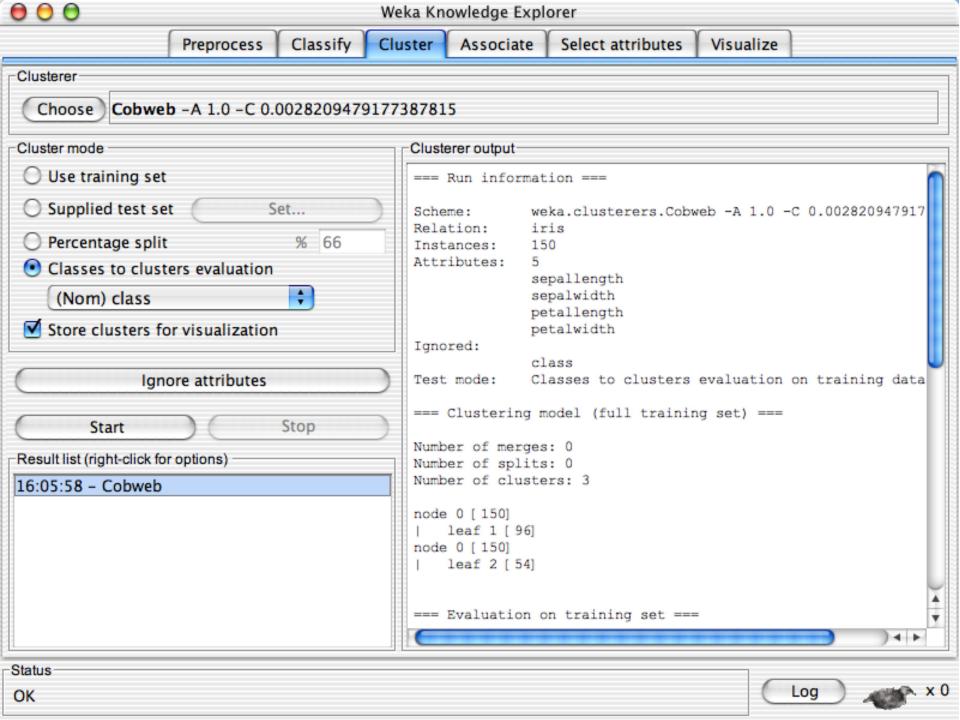


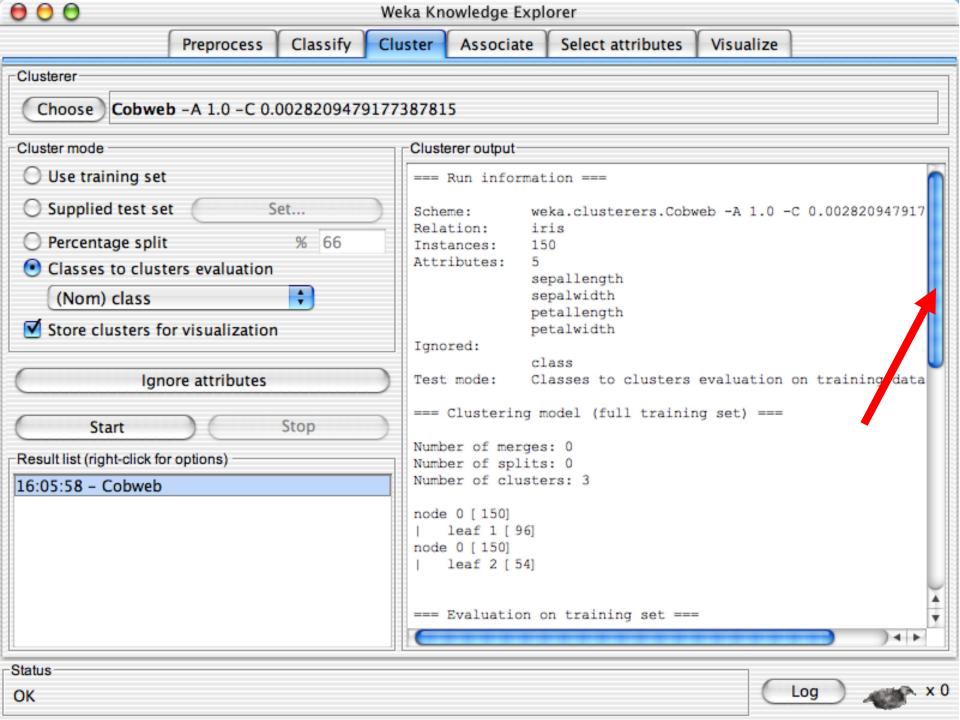


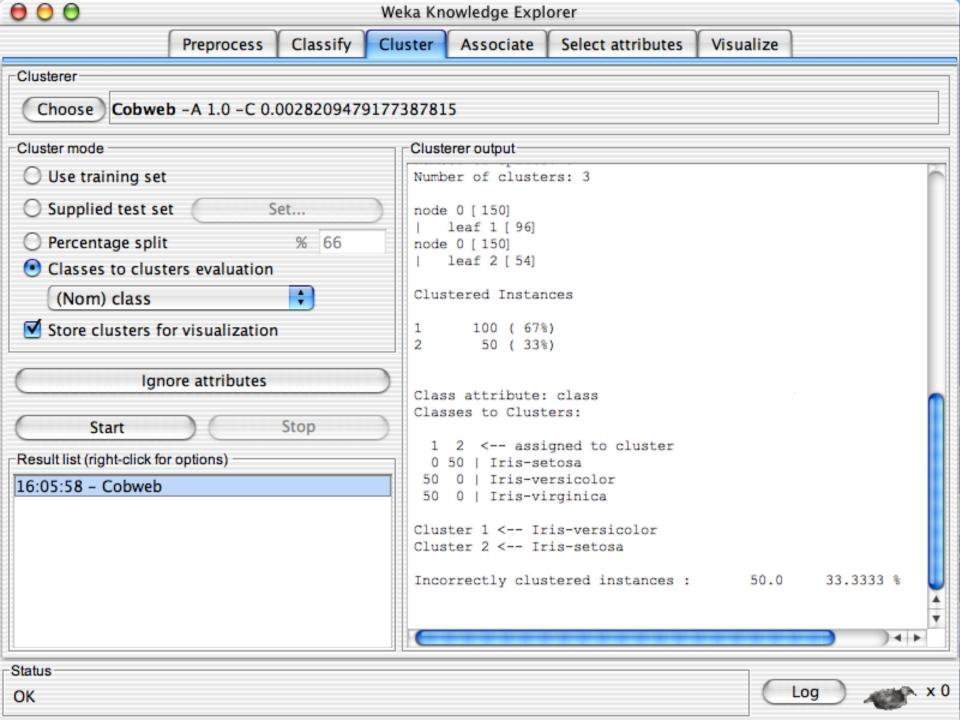


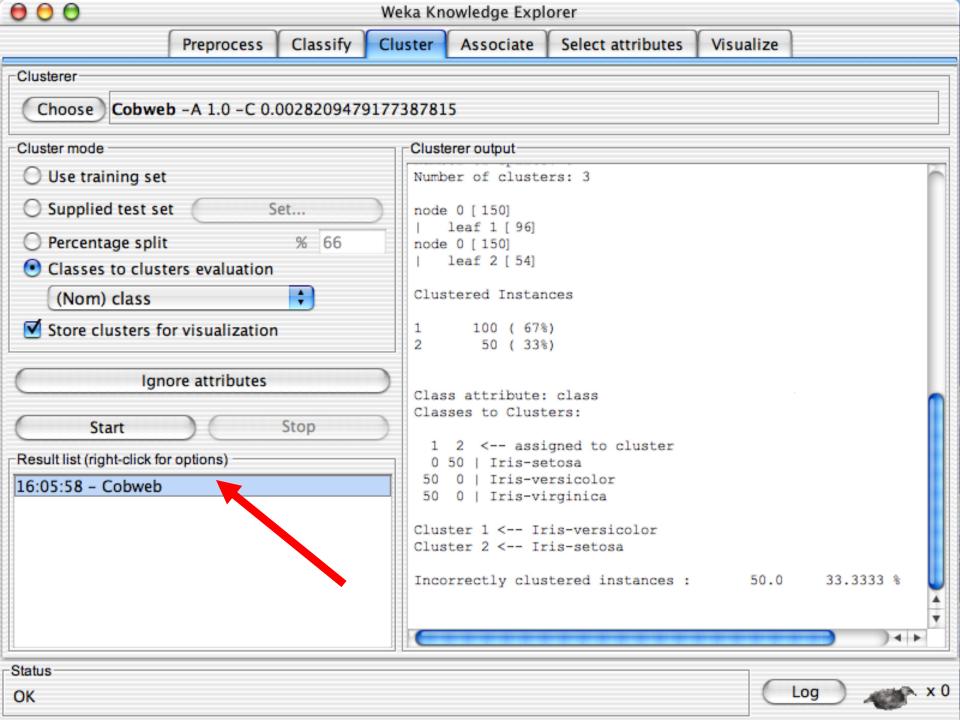


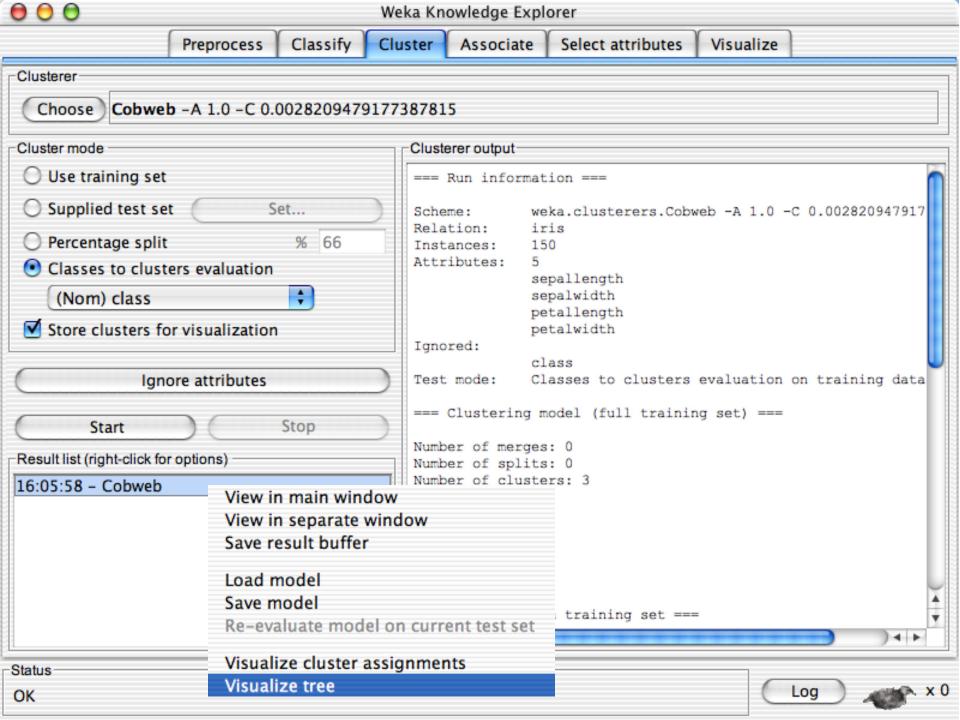


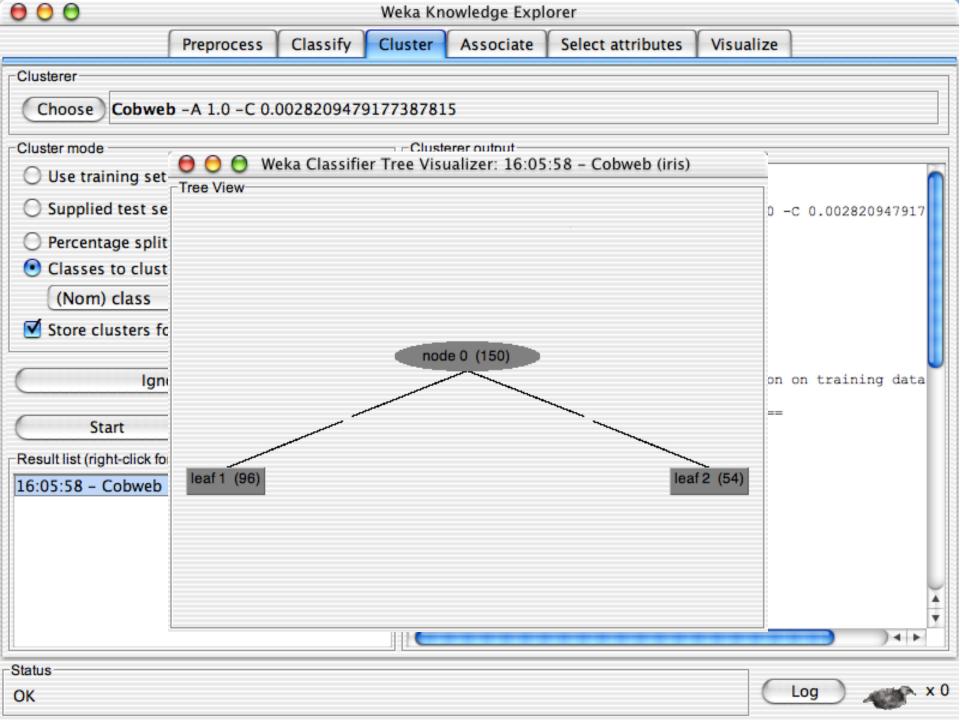


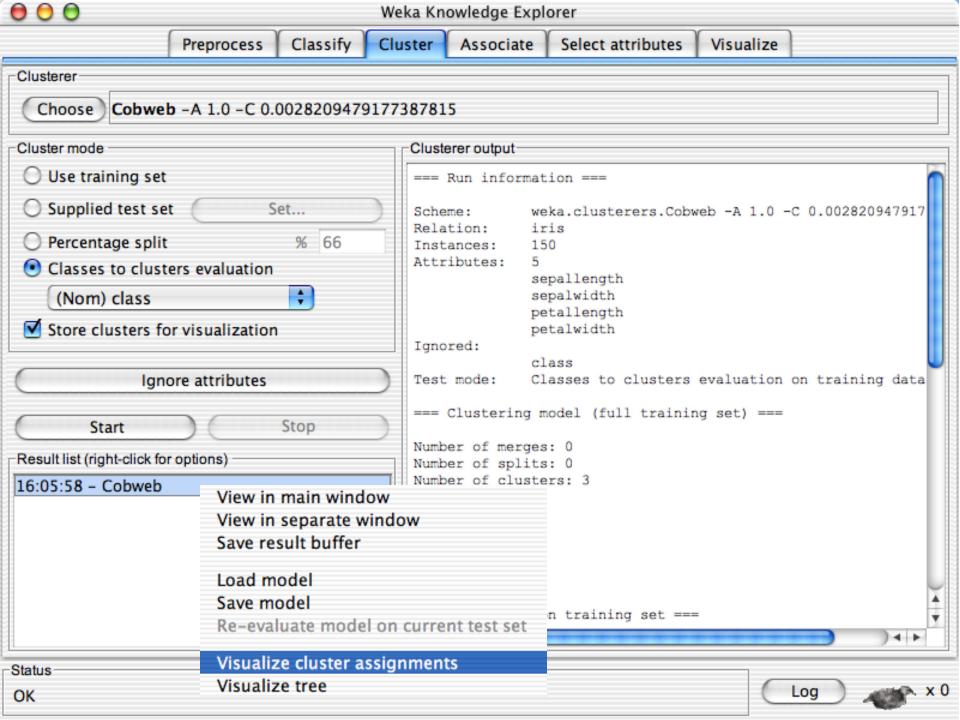


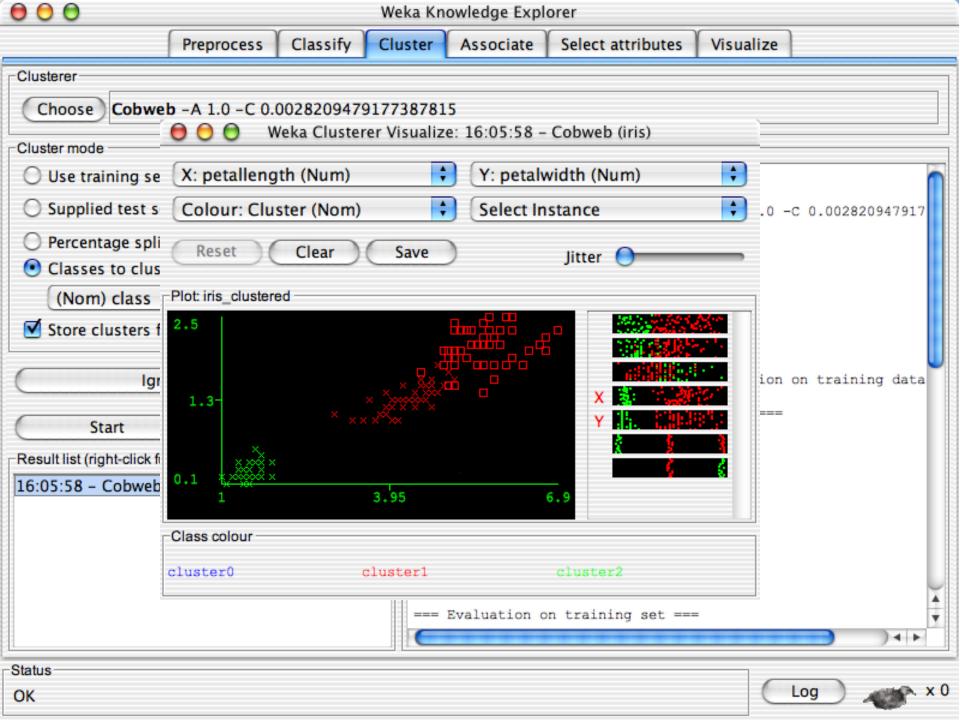






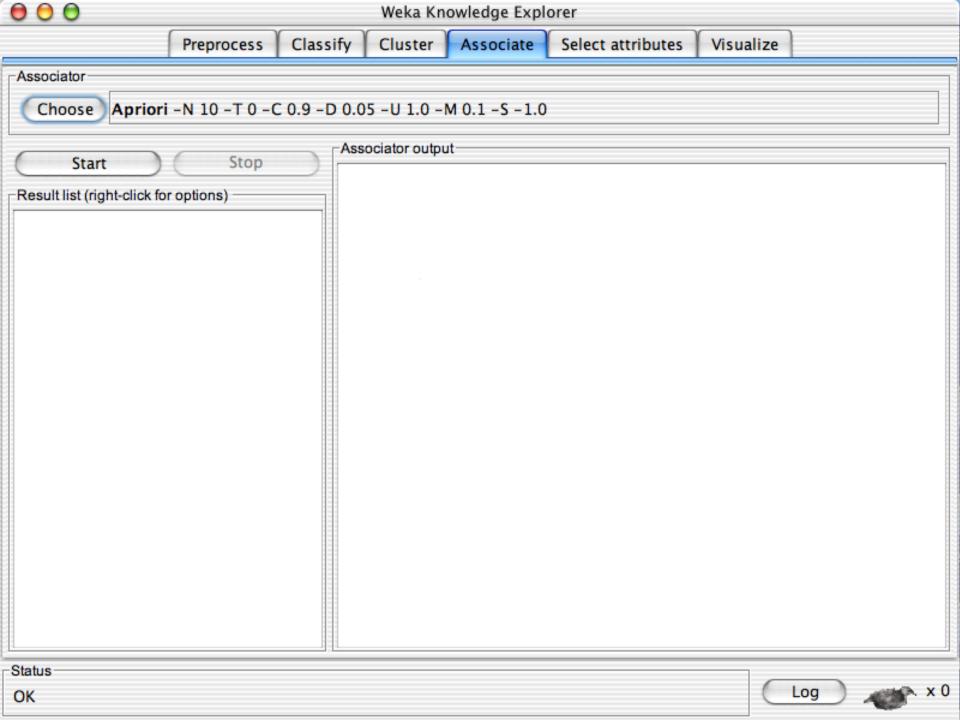


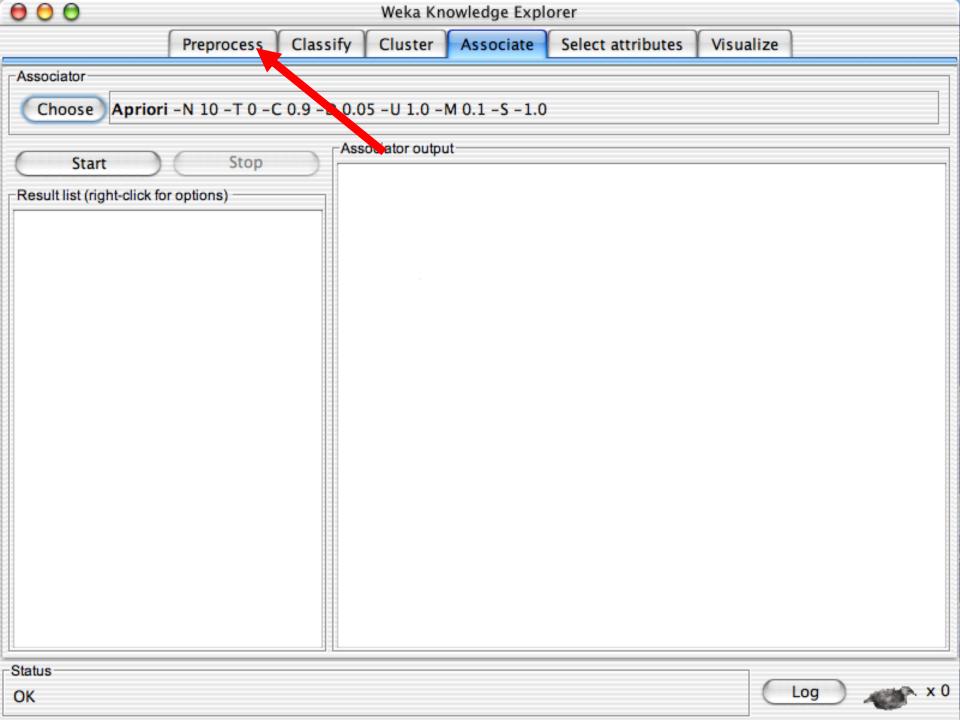


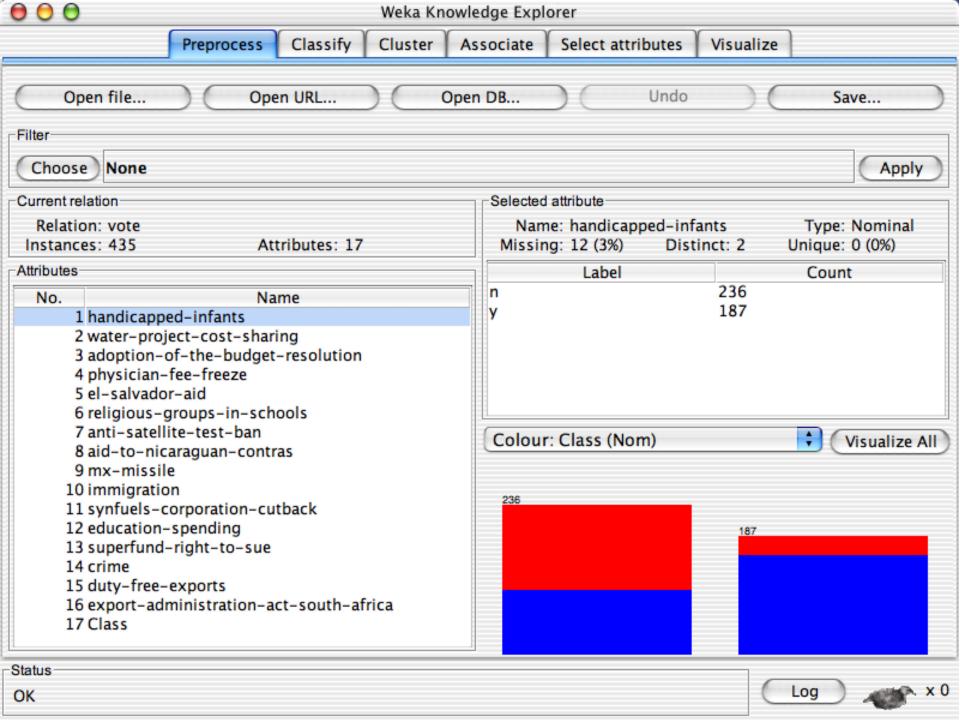


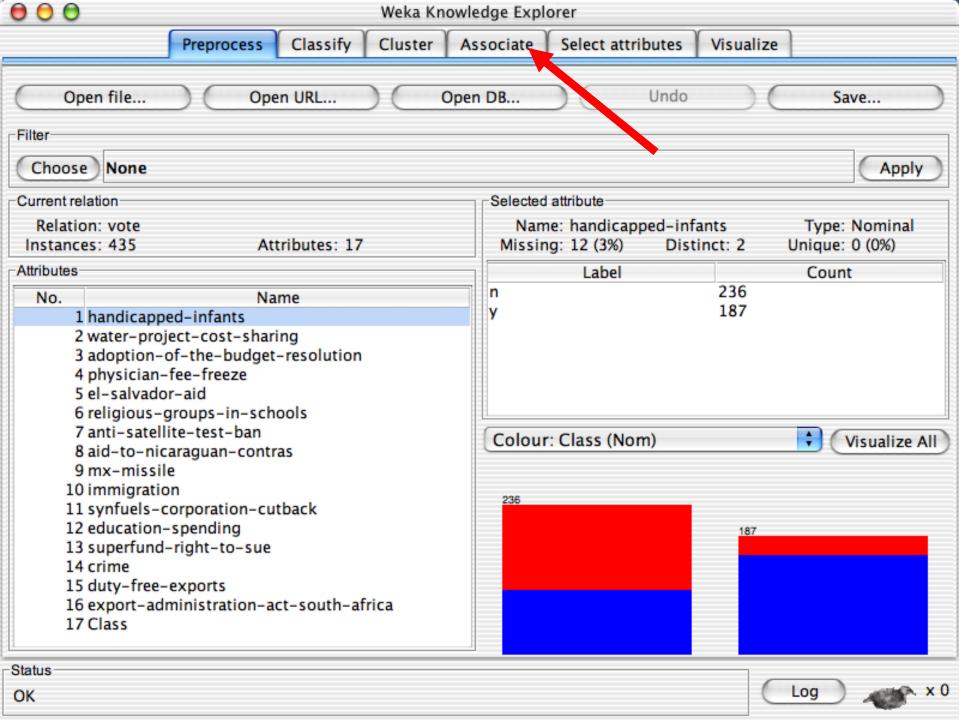
Explorar: asociaciones

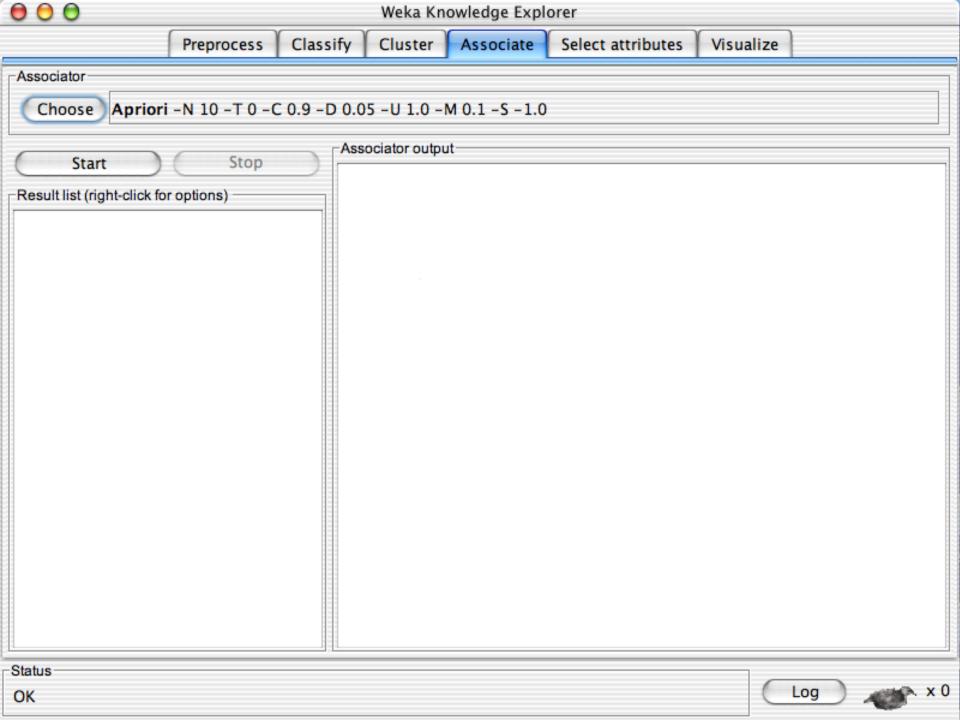
- WEKA implementa el algoritmo Apriori para aprender reglas de asociación
 - Solo con datos discretos
- Puede identificar dependencias estadísticas entre grupos de atributos:
 - Leche, matequilla ⇒ pan, huevos(con confianza de 0.9 y soporte de 2000)
- Apriori puede calcular reglas con mínimo soporte y que excedan una dada confianza

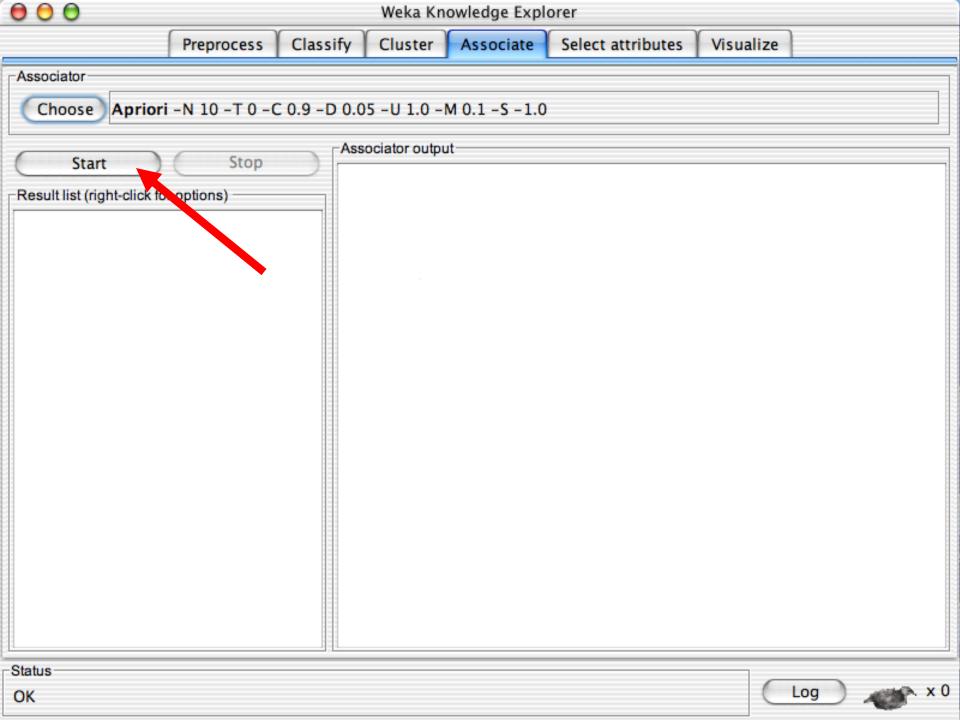


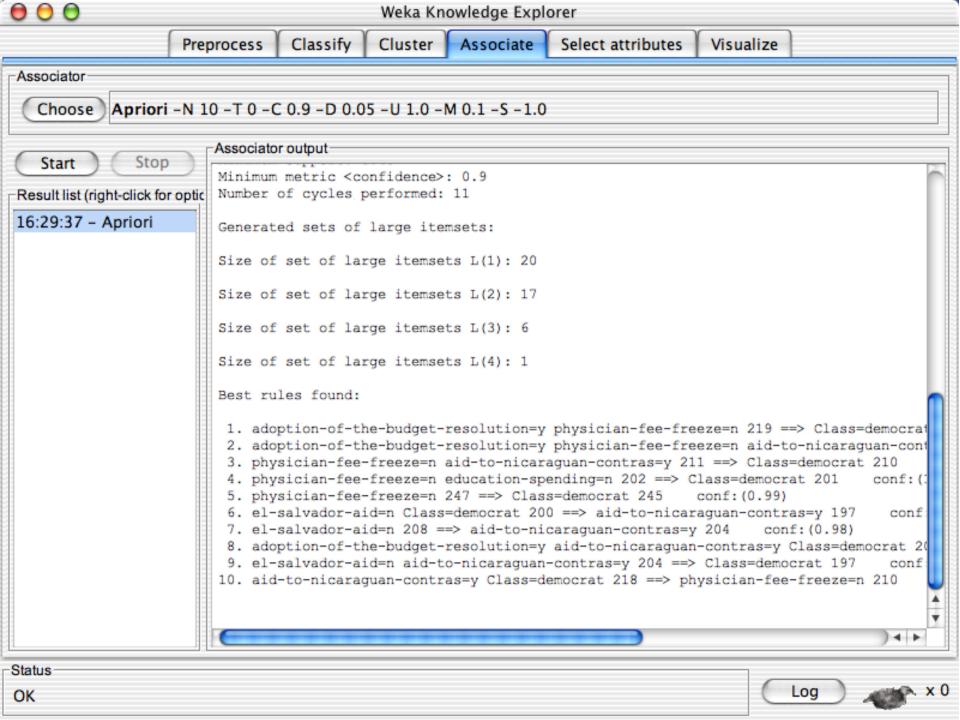












Explorar: selección de atributos

- Panel para investigar los atributos más predictivos
- Métodos de selección:
 - Métodos de búsqueda: best-first, aleatorio, exhaustivo, algoritmos genéticos, rankeo
 - Métodos de evaluación: basado en correlación, wrapper, ganancia de información, chi-squared, ...
- Se pueden combinar

Explorar: visualizar datos

- Visualizar es muy útil y práctico: p.e. ayuda a determinar dificultades en los procesos de aprendizaje
- WEKA puede visualizer un atributo simple (1d) y pares de atributos (2-d)

Weka Java library

Clases para carga de datos

Clases para los clasificadores

Clases para la evaluación

Clases para carga de datos

- weka.core.Instances
- weka.core.Attribute

Cada DataRow -> Instance,

```
# Load a file as Instances
FileReader reader;
reader = new FileReader(path);
Instances instances = new Instances(reader);
```

– ¿Cómo recuperar un valor de una instancia?

```
# Get Instance
Instance instance = instances.instance(index);
# Get Instance Count
int count = instances.numInstances();
```

– ¿Cómo recuperar un atributo?

```
# Get Attribute Name
Attribute attribute = instances.attribute(index);
# Get Attribute Count
int count = instances.numAttributes();
```

Clases para carga de datos

– ¿Cómo recuperar el valor del atributo para cada Instancia?

```
# Get value instance.value(index); or instance.value(attrName);
```

Indice de Clase

Clases para los clasificadores

- Clases Weka para C4.5, Naïve Bayes, and SVM
 - Clasificadores:
 - C4.5: weka.classifier.trees.J48
 - NaiveBayes: weka.classifiers.bayes.NaiveBayes
 - SVM: weka.classifiers.functions.SMO
- Cómo construir un clasificador?

```
# Build a C4.5 Classifier

Classifier c = new weka.classifier.trees.J48();

c.buildClassifier(trainingInstances);

Build a SVM Classifier

Classifier e = weka.classifiers.functions.SMO();

e.buildClassifier(trainingInstances);
```

Clases para la evaluación

- weka.classifiers.CostMatrix
- weka.classifiers.Evaluation

Como usarlas?

```
# Use Classifier To Do Classification
CostMatrix costMatrix = null;
Evaluation eval = new Evaluation(testingInstances, costMatrix);

for (int i = 0; i < testingInstances.numInstances(); i++){
    eval.evaluateModelOnceAndRecordPrediction(c,testingInstances.instance(i));
    System.out.println(eval.toSummaryString(false));
    System.out.println(eval.toClassDetailsString());
    System.out.println(eval.toMatrixString());
}
```

Clases para la evaluación

Validación Cruzada

 Divide un único conjunto de datos en N partes iguales.

 Toma la N-1 como un conjunto de datos de entrenamiento, el resto se utilizará como prueba de conjunto de datos.

Clases para la evaluación

Obtención conjunto de entrenamiento



Tres pasos para el uso

- Asignar el archivo de datos primero
- Seleccionar funcionalidad
- Ejecutar Función utilizando rápido Minero

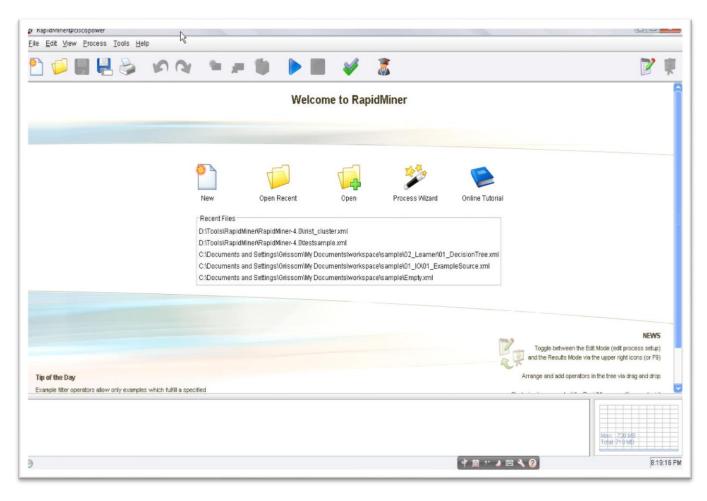
Los conjuntos de datos pueden utilizar formato ARFF, pero otros formatos támbien

RapidMiner

- Software open-source software que permite
 - Analisis inteligente de datos
 - Mineria de datos
 - Descubrir Conocimiento
 - Aprendizaje automático
 - Analitica predictivAa
 - Inteligencia de Negocio.
- Esta en Java y es GPL licencia
- http://rapid-i.com

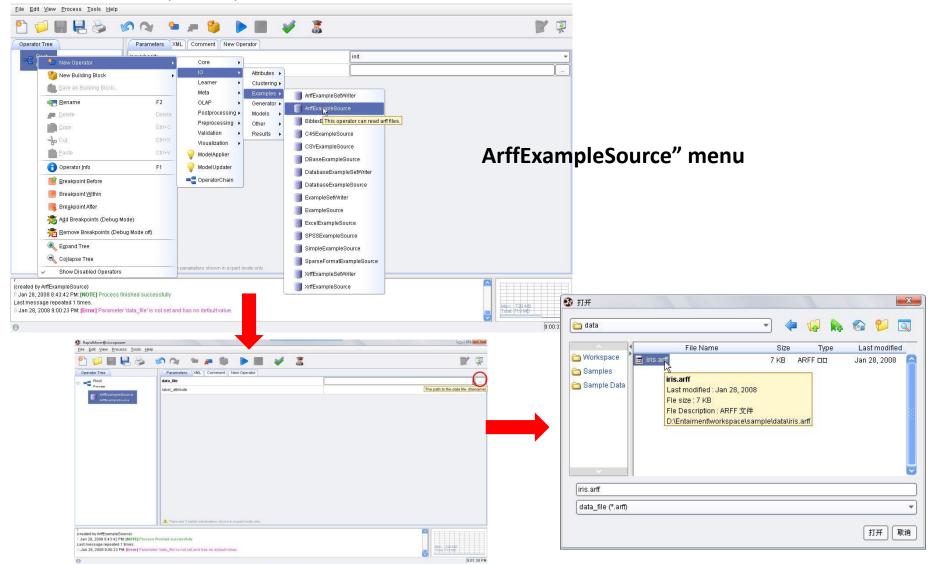
http://rapidi.com/content/view/130/82/

➤ Crear un proyecto Rapid Miner (opción "new")



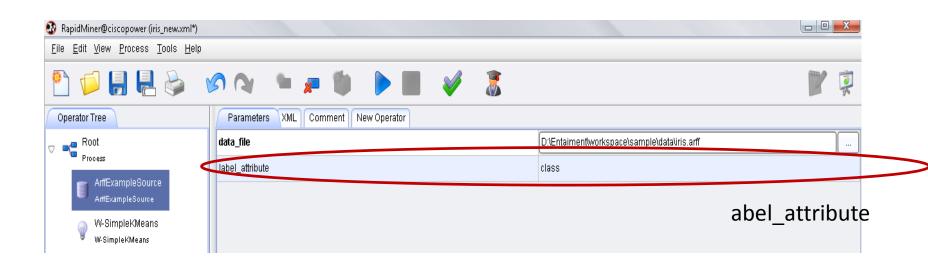
Establecer una fuente de datos para el proyecto

➤ Arff, excel , etc



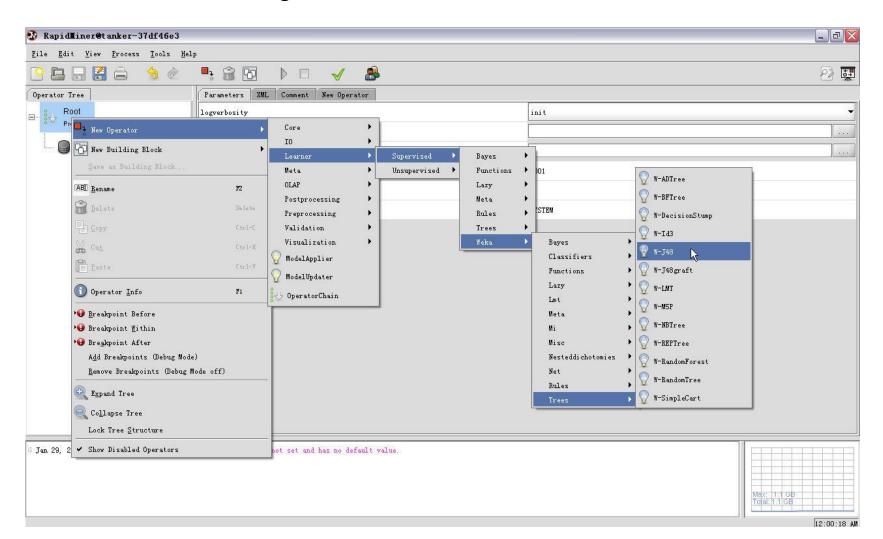
Escoger etiqueta atributo

>qué atributo en el origen de datos es el atributo de clase.

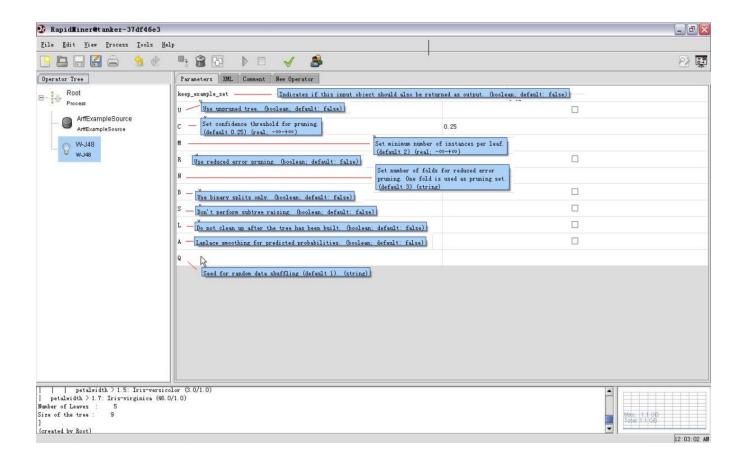


Selecccionar funcionalidad

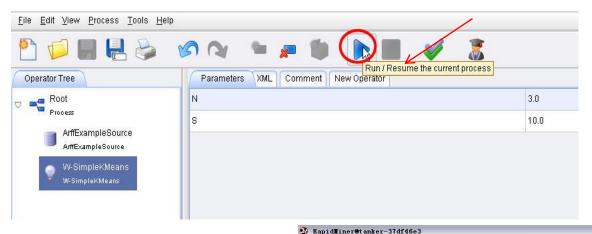
clustering, classification, decision tree, etc.



Escoger parámetros



Ejecutar



resultados

```
File Edit View Process Tools Help
 Model
● Text Kiew O Graph View
      Changes to a textual view of this model.
 W-J48
 J48 pruned tree
 petalwidth <= 0.6: Iris-setosa (50.0)
 petalwidth > 0.6
 petalwidth <= 1.7
   petallength <= 4.9: Iris-versicolor (48.0/1.0)
 petallength > 4.9
 petalwidth <= 1.5: Iris-virginica (3.0)
 petalwidth > 1.5: Iris-versicolor (3.0/1.0)
 | petalwidth > 1.7: Iris-virginica (46.0/1.0)
 Number of Leaves : 5
 Size of the tree: 9
   petalwidth > 1.5: Iris-versicolor (3.0/1.0)
petalwidth > 1.7: Iris-virginica (46.0/1.0)
Number of Leaves :
Size of the tree :
(created by Root)
```



Qué (y cómo) es R

- R es un entorno libre de análisis estadístico de datos y de creación de gráficos estadísticos
- Se basa en una interfaz de usuario de líneas de comandos
- Lenguaje Interpretado
- Guarda, análiza y gráfica datos
- Software Open Source

Razones para utilizar (y no) utilizar

Pros

- Es libre
- Flexibilidad
- Procedimientos disponibles
- Se aprende estadística
- Estado del arte en estadistica
- 2do detras de MATLAB para graficos.
- Comunidad de usuarios
- Lo oblige a pensar al usuarios sobre el analisis
- Interfaz con BD (SQL)

□Contras

- Al principio suele ser durp
- Cuesta un poco si se está habituado a trabajar con otro tipo de programas
- No amigable: curva de aprendizaje
- Fácil cometer errores

Definiciones Básicas

- Objetos
 - Funciones
 - nombre.de.la.función(argumento/s=, opción/es=)
- Espacio o área de trabajo
- Directorio de trabajo
- Paquetes
- Archivos de comandos o scripts

Como trabajar en R...

- Leer datos desde otras fuentes
- Usar paquetes, librerias, y funciones
- Escribir funciones si es necesario
- Conducir Análisis estadístico de Datos Guardar salidas en archivos
- Guardar espacio de trabajo de R si es necesario

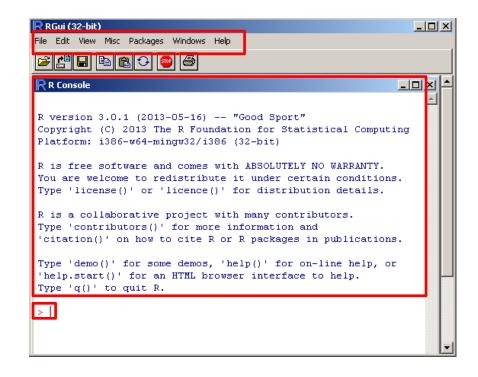
Corriendo R

- Sistemas de ventanas(Console)
- Usando Editores
 - Notepad, WinEdt, Tinn-R: Windows
 - Xemacs, ESS (Emacs speaks Statistics)
- En el Editor:
 - -source("filename.R")
 - Salida puede ser direccionada con
 - sink("filename.Rout")

Descripción de la interfaz gráfica de usuario en R

 Ejecutar el programa o el acceso directo del escritorio

 Apertura de la consola de comandos

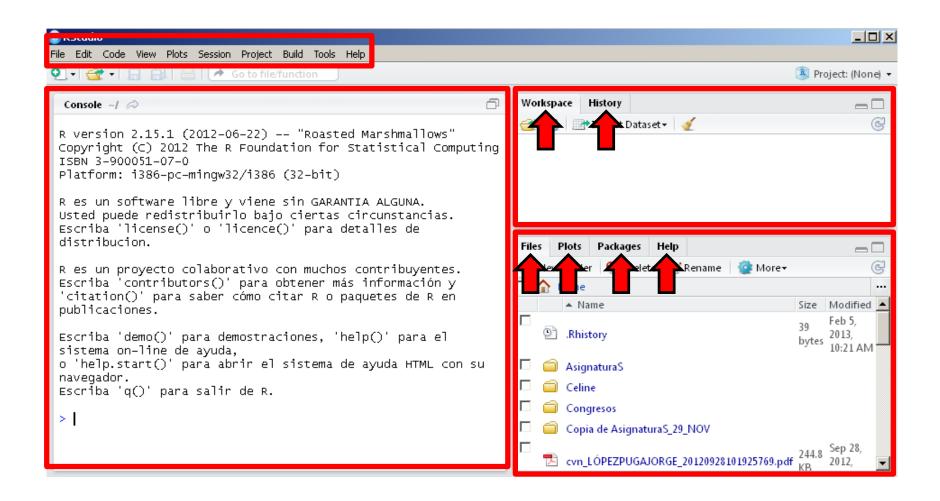


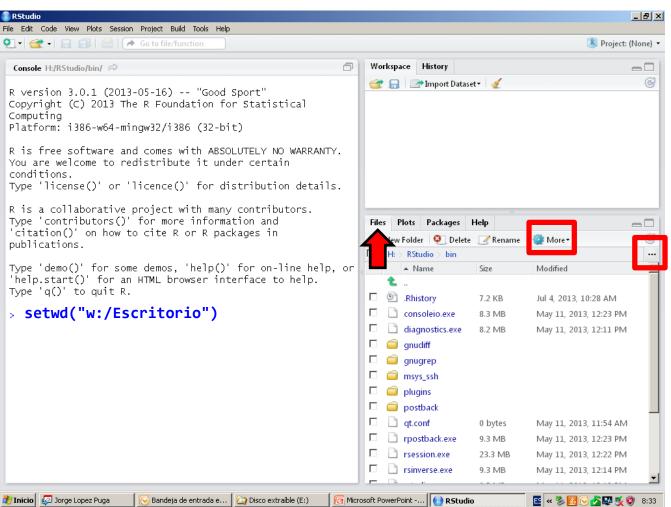
Otras Interfaces Interesantes

Hay variedad de interfaces para

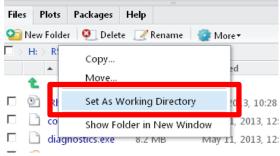
- R Commander es un paquete de R
- RKWard http://rkward.sourceforge.net
- Tinn-R http://www.sciviews.org/Tinn-R
- Emacs http://www.gnu.org/software/emacs
- RStudio http://www.rstudio.com

Interfaz gráfica de RStudio









Comenzar con R

- Basico: asignaciones y operaciones.
- Operaciones Aritmeticas:

- Aritmetica de Matrices.
 - * is element wise multiplication
 - %*% es una multiplicación de matrices
- Asignación

Descripción de paquetes

- Función packageDescription()
- Función library(help="nombre.del.paquete")

Cargar paquetes

• La función library()

Instalar paquetes

• La función install.packages("nombre.del.paquete")

Desactivar paquetes

La función

```
detach("package:nombre.del.paquete")
```

Desinstalar paquetes

La función

```
remove.packages("nombre.del.paquete")
```

Actualizar paquetes

La función update.packages()

Los objetos de R

- Uno objeto es una estructura de datos con la que R puede trabajar
- Por ejemplo:
 - Un conjunto de datos base de datos
 - El resultado de un análisis estadístico
 - Una tabla de datos
 - Una función

Tipos de objetos

Soporta cualquier tipo de datos

- Vector (numéricos, caracteres, lógicos)
- Matrices (y arrays)
- Listas (pueden contener información de diferente tipo)
- Tablas (de frecuencia o de contingencia)
- Data frame o base de datos

variables

```
> a = 49
> sqrt(a)
[1] 7
```

numeric

```
> a = "The dog ate my homework"
> sub("dog","cat",a)
[1] "The cat ate my homework"
```

character string

logical

Datos

	Lineal	Rectangular
Del mismo Tipo	VECTORS	MATRIX*
Mezcla	LIST	DATA FRAME

Tareas específicas

- search(): ver cuales directories y datos existen
- ls(): ver objetos guardados
- attach(NameOfTheDataset, expression): incluir un conjunto de datos en el camino para proceso de análisis
- detach(NameOfTheDataset): quitar un conjunto de datos del camino del proceso de análisis despues de culminarlo

Creación de objetos

- La creación de objetos se lleva a cabo realizando una asignación
- Para ello utilizamos los símbolos <-, -> o =
 - Los códigos parecidos a flecha funcionan en ambos sentidos
 - El símbolo igual sólo en un sentido

R como calculador

> log2(32)

[1] 5

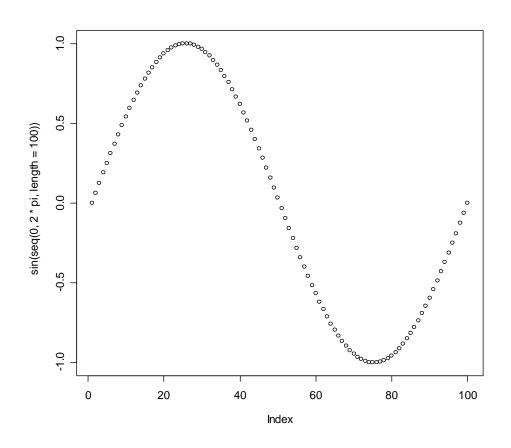
> sqrt(2)

[1] 1.414214

> seq(0, 5, length=6)

[1] 0 1 2 3 4 5

> plot(sin(seq(0, 2*pi, length=100)))



Archivos de código

- R no genera archivos de gráficos propios
- Por tanto, es conveniente guardarlos como scripts
- Recordar la extensión .R
- Comentar líneas para aclaraciones
- El dispositivo gráfico
- Abrir un dispositivo gráfico
 - Al ejecutar una función de alto nivel – consola de R
 - Con la función
 windows() → dibujar
 gráfico

Guardar como
Copiar para el área de transferencia
Copiar para el área de transferencia
Como un Bitmap CTRL+C
como un Metafile CTRL+W

CETRL+P

CETRL+P

Archivo.

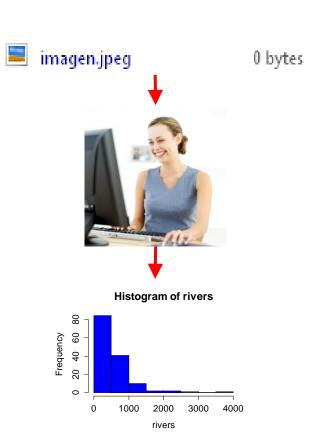
P Graphics: Device 3 (ACTIVE)

Histórico Redimensionar

Utilizando código

- Habría que dar tres pasos
 - 1. Inicializar generar el archivo
 - 2. Generar el gráfico
 - 3. Cerrar el dispositivo gráfico
- Ejemplo:

```
jpeg("imagen.jpeg")
hist(rivers, col="blue")
dev.off()
```



Creación de vectores

- La función c()
- La función assingn()
- Los dos puntos :
- La función seq()
- Distribuciones estadísticas
- Ejemplos

Modos de crear matrices y arrays

- A partir de un vector con la función dim()
- Con la función matrix()
- Con la función array()

vectores, matrices

vector: mismo tipo

$$> a = c(1,2,3)$$

 $> a*2$
[1] 2 4 6

- En R, un numero es un especial caso de un vector con 1 elemento
- Other vector types: character strings, logical

matriz: tabla rectangular de datos del mismo tipo

• ejemplo: bipsias descritas por 10000 genes y 30 tejidos: una matriz de 10000 filas y 30 columnas.

Listas

- vector: acceso.
- > a = c(7,5,1)
- > a[2]
- [1] 5

lista: colección de datos de tipo arbitrario.

- > doe = list(nombre="jose",edad=28,casado=F)
- > doe\$nombre
- [1] "jose"
- > doe\$edad
- [1] 28
- Elementos vectores/matrices se accesan por indices (un entero): x[1], y[1,] y[,1],
- Elementos listas/frames por sus nombres (cadena caracter): myframe\$age.
- Ahora, ambos tipos soportam ambos métodos.

Creación de tablas y su utilidad

- Creación (entre otras)
 - Con la función table()
 - Con la función as.table()

- Utilidad:
 - Para realizar análisis de frecuencias
 - Para crear gráficos (por ejemplo, de barras o de sectores)

Creación de tablas y su utilidad

- Función read.table()
 - Lee datos directamente
 - La primera línea del archivo tiene nombre de cada variable en la tabla
 - Cada línea adicional comienza con una etiqueta de fila y despues los valores de cada variable.

l	Precio	Piso	Area	Hab	Edad	Calent
01	52.00	11	830	5	6.2	no
02	54.75	8	710	5	7.5	no
03	57.50	1	1000	5	4.2	no
04	57.50	13	690	6	8.8	no
05	59.75	9	900	5	1.9	yes

. . .

¿Qué es una base de datos?

- En R se le llama data frame
- Una especie de matriz bi-dimensional
 - Columnas representan variables
 - Filas representan personas, registros o casos
- Contiene diferentes tipos de datos (es un tipo de lista)
- Ejemplo: iris

Creación de bases de datos

- Existe la función data.frame()
 - Hay que crear previamente un conjunto de vectores
- A partir de otros objetos as.data.frame()
- También podemos importar datos desde archivos externos (Excel, SPSS, .csv, .txt)

Data frames

Tabla rectangular de filas y columnas; y datos en cada columna son del mismo tipo (p.e. numero, texto, lógico), pero diferentes columnas pueden tener diferentes tipos.

> a			
	localización	tamatumor	progreso
XX348	proximo	6.3	FALSE
XX234	distan	8.0	TRUE
XX987	proximo	10.0	FALSE

Gráficos

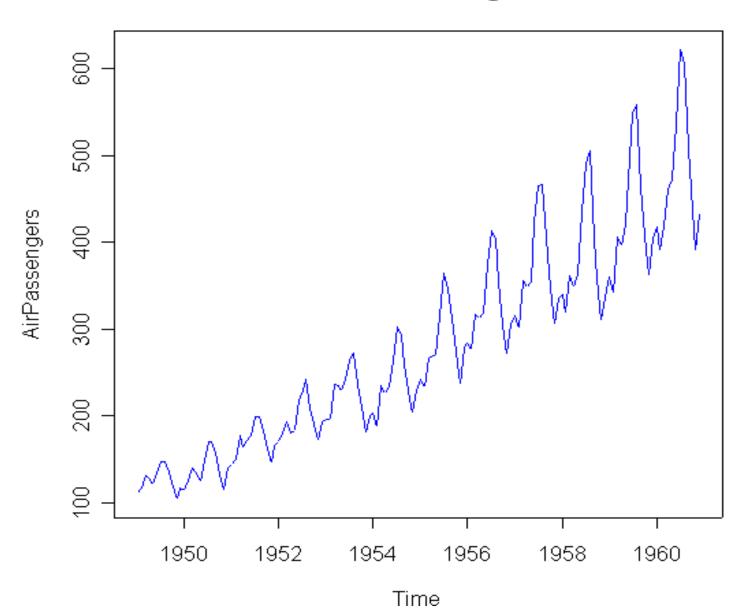
Dibujar un objecto: plot(num.vec)

- Dibujo enviar a un dispositivo gráfico
 - Debe decirse dispositivo
 - postscript, gif, jpeg, etc...
 - Se pueden apagar o encender: dev.off()

- 2 tipos
 - Alto nivel: dibujo gráfico
 - Bajo nivel: agregar información a uno existente

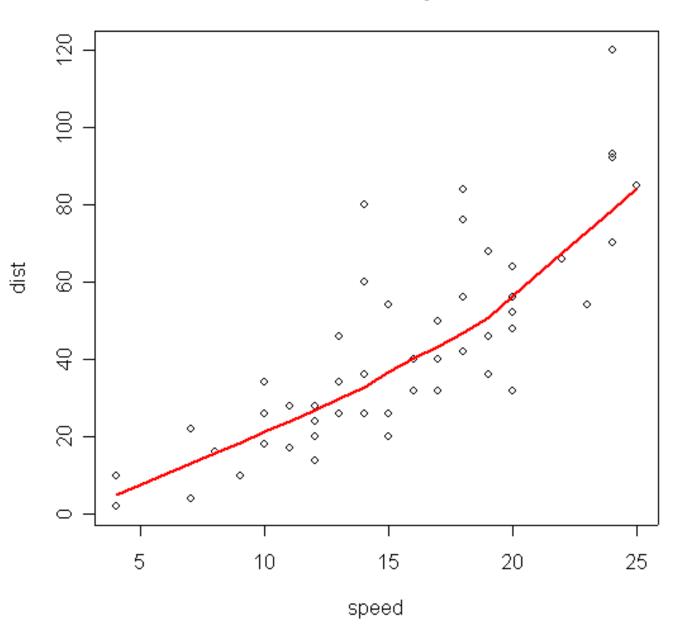
Alto nivel: generado con plot()

Number of Airline Passengers over time



Bajo nivel

distance vs speed



Programando con R

- Funciones & Operadores
- Expresiones definidas en {}
- Código separado por nueva linea, ";" no es necesario
- Repetición
- Condicional

Condicional

```
if (logical expression) {
  statements
} else {
  alternative statements
}
```

else branch is optional

Lazos

```
for(i in 1:10) {
    print(i*i)
}

i=1
while(i<=10) {
    print(i*i)
    i=i+sqrt(i)
}</pre>
```

- repeat expr statement.
- instrucción break usado para terminar un lazo.
- instrucción next usado para saltar un ciclo y pasar al siguiente.

lapply, apply

• Si una tarea debe ejecutarse sobre los elementos de una lista o matriz

lapply(li, function)

```
> li = list("klaus","martin","georg")
lapply(li, toupper)
> [[1]]
> [1] "KLAUS"
> [[2]]
> [1] "MARTIN"
> [[3]]
> [1] "GEORG"
```

 Aplicar la función fct en una dimension de una matriz arr, y regresarlo en un vector o matriz de apropiada tamaño.

```
apply( arr, margin, fct )
```

```
> x

[,1] [,2] [,3]

[1,] 5 7 0

[2,] 7 9 8

[3,] 4 6 7

[4,] 6 3 5

> apply(x, 1, sum)

[1] 12 24 17 14

> apply(x, 2, sum)

[1] 22 25 20
```

Guardar datos

- > save(x, file="x.Rdata")
- > load("x.Rdata")

Importar y exportar datos

```
> x = read.delim("filename.txt")
also: read.table, read.csv
```

> write.table(x, file="x.txt", sep="\t")

Funciones estadisticas en R

Descriptivas

- mean, median, first, third quartiles,
- stem() gives stem-leaf plots
- table() gives tabulation of categorical variables

Modelado

- Regresión: Lineal y Lógica
- Series de tiempo
- Más de 400 funciones
 - lm, glm, aov, ts
- Muchas librerias & paquetes
 - survival, coxph, tree (recursive trees), nls, ...

Multivariables

Descriptivas

- Tendencia central
 - Media
 - Mediana
 - Segmentación por factores
- Opciones interesantes

Tendencia central

- Para describir variables de manera general
- La función median()
- Por ejemplo:

```
median(rivers) [1] 425
```

- La media aritmética
- Función mean()
- Por ejemplo:

```
mean(austres) [1] 15273.45
```

Opciones interesantes

- Aprovechando que acabamos de hablar de la media...
- El argumento trim
 - Sirve para estimar medias recortadas valor de 0 a 0.5
- Por ejemplo:

Descriptivas

- Dispersión:
 - La función para la varianza: var()
 - La función para la desviación típica: sd()

- Forma
 - La función para el sesgo o asimetría: skewness()
 - Función para el apuntamiento o curtosis: kurtosis()

Dispersión

- Está referida a la heterogeneidadhomogeneidad
- La función var() para estimar la varianza
- Puede funcionar
 - Con vectores o variables
 - Con matrices
- Por ejemplo:

```
var(rivers)
[1] 243908.4
```

Dispersión

- Cuando el argumento es una matriz
 - Matriz de varianzas-covarianzas

```
data(iris)
iris$Species <- NULL
var(iris)</pre>
```

```
Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length
                0.6856935
                           -0.0424340
                                          1.2743154
                                                      0.5162707
Sepal.Width
               -0.0424340 0.1899794
                                         -0.3296564
                                                     -0.1216394
Petal.Length
                                          3.1162779
                1.2743154
                            -0.3296564
                                                      1.2956094
Petal Width
                0.5162707
                            -0.1216394
                                          1.2956094
                                                      0.5810063
```

Dispersión

- La función sd() para estimar la desviación típica
- El argumento (x) ha de ser un vector
- Acepta el argumento na.rm para casos perdidos
- Por ejemplo:

```
sd(rivers)

[1] 493.8708

vp <- c(45, 75, 32, NA, 22, 14, 85)
sd(vp, na.rm = TRUE)

[1] 28.83574</pre>
```

Como modelar

Específicar modelo:

$$y \sim x_i + c_i$$

- -y = variable resultado, x_i = variables de entrada, c_i = constantes, + = operador
- Operadores
 - + = agrega, : = interactua, / = anida, etc...

```
carReg <- Im(speed~dist, data=cars)
carReg = es un objecto</pre>
```

Ejemplo de Modelo de Regresión

plot

- La función plot()
- Gráficos de dispersión
- Gráficos de caja
- Gráficos de series temporales
- Gráficos de barras
- Gráficos de Mosaico
- Gráficos de densidad

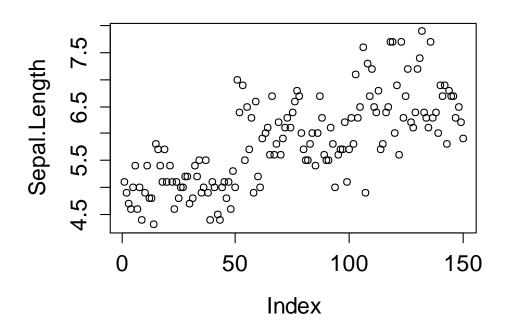
- Es una función genérica para dibujar gráficos x-y
- Sus dos argumentos básicos son x e y
- Otro argumento interesante que puede tomar es type
 o superpuestos

o s – salto de "escalera"

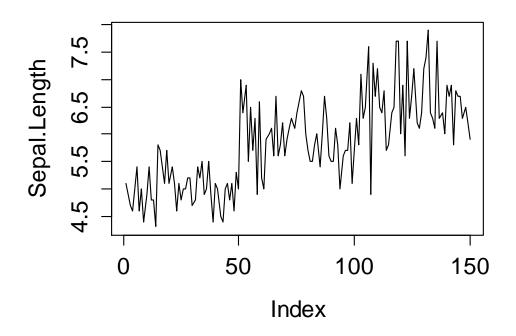
- p puntos (por defecto)^{h histograma}
- 1 línea
- b ambos

- Ejemplos
- Usemos la variable Sepal. Length contenida en la base de datos iris
- Recuerda utilizar attach()
- Es decir:
 - attach(iris)

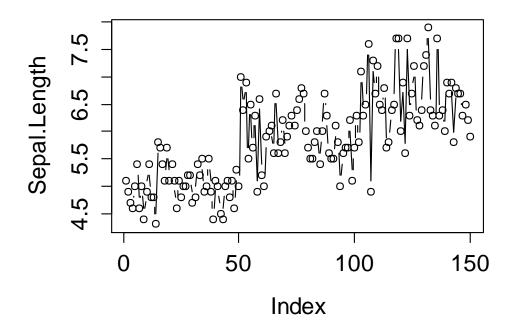
• plot(Sepal.Length, type="p") =
 plot(Sepal.Length)



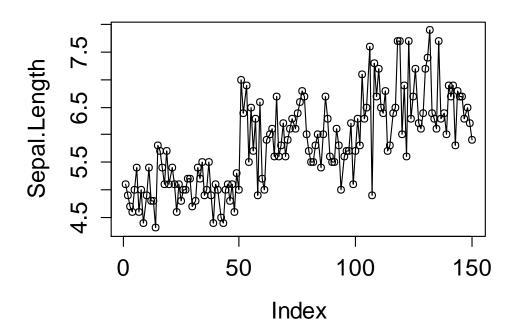
plot(Sepal.Length, type="1")



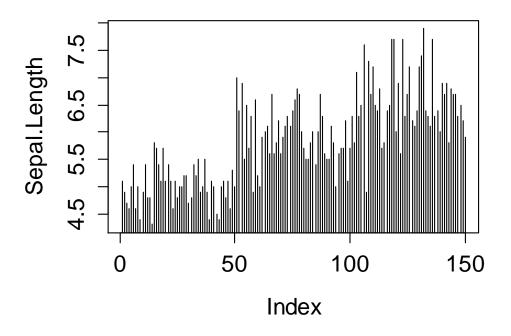
plot(Sepal.Length, type="b")



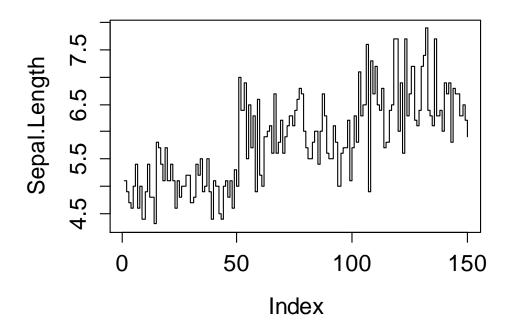
plot(Sepal.Length, type="o")



plot(Sepal.Length, type="h")

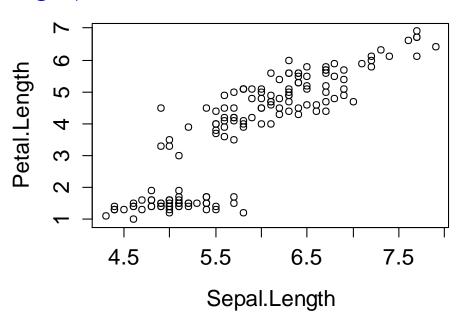


plot(Sepal.Length, type="s")



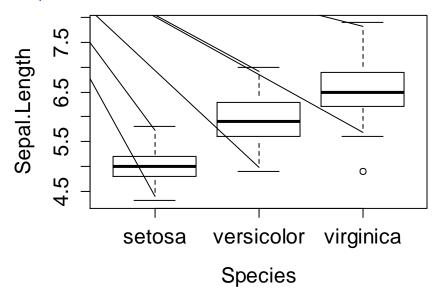
Gráficos de Dispersión

- Necesitamos dos variables (x e y)
- Como coordenadas (x, y)
 - plot(Sepal.Length, Petal.Length)
- Con interfaz tipo fórmula
 - plot(Petal.Length ~ Sepal.Length)



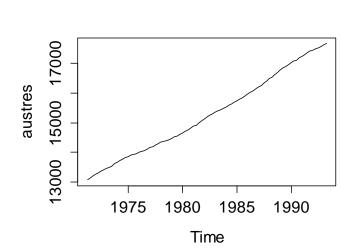
Gráficos de Caja

- Combinando un factor y una variable cuantitativa
- Aludiendo a lo representado en los ejes
 - plot(Species, Sepal.Length)
- Interfaz tipo fórmula
 - plot(Sepal.Length ~ Species)

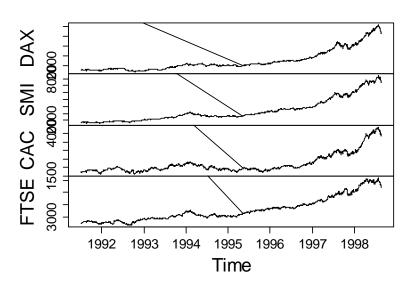


Series Temporales

- Cuando se utiliza como objeto una serie temporal
 - plot(austres)
 - plot(EuStockMarkets)



EuStockMarkets



Gráficos de Barras

- Cuando el argumento es una tabla
 - plot(table(Sepal.Length))

Frecuencia absoluta

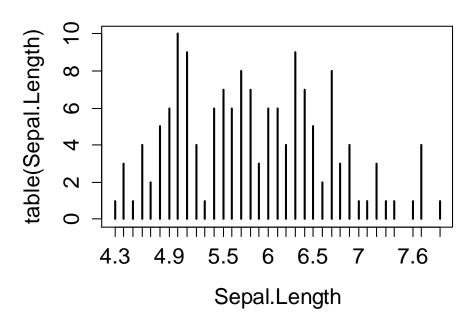
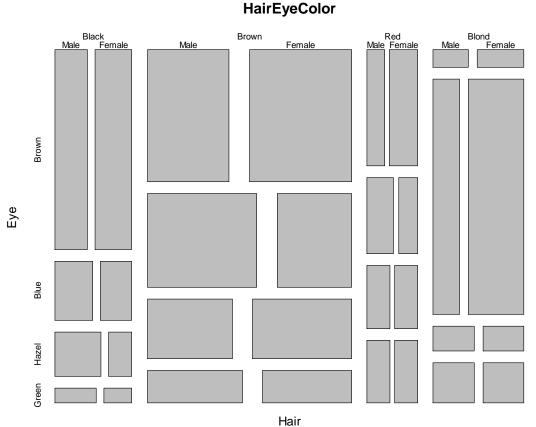


Gráfico de Mosaico

- Cuando el argumento es una tabla bi o multi dimensional
- Por ejemplo
 - HairEyeColor
 - plot(HairEyeColor)

Variables:

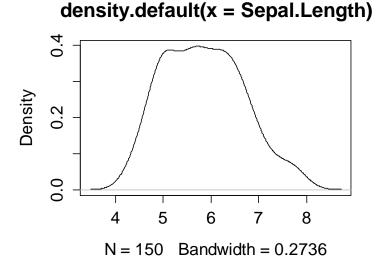
Color del Pelo Color de los ojos Sexo



Gráficos de Densidad

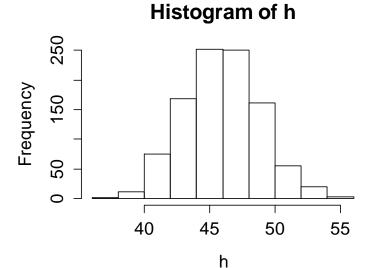
- Podemos obtener la función de densidad de una variable utilizando la función density()
- Podríamos trazar el gráfico de densidad incrustando la función density() dentro de plot()

Consideremos a la variable
Sepal.Length
plot(density(Sepal.Length))

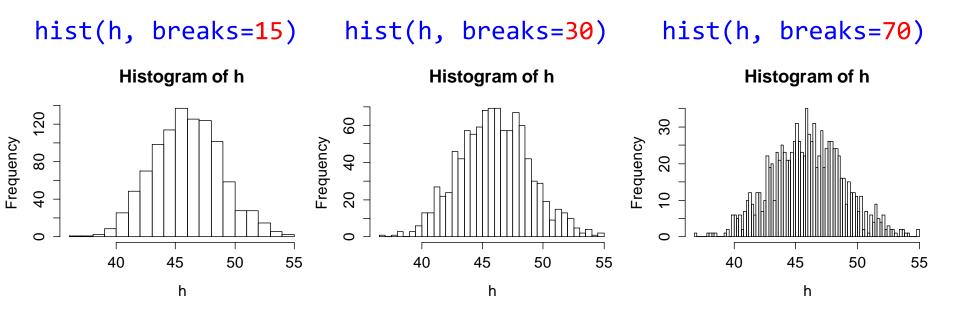


La función hist()

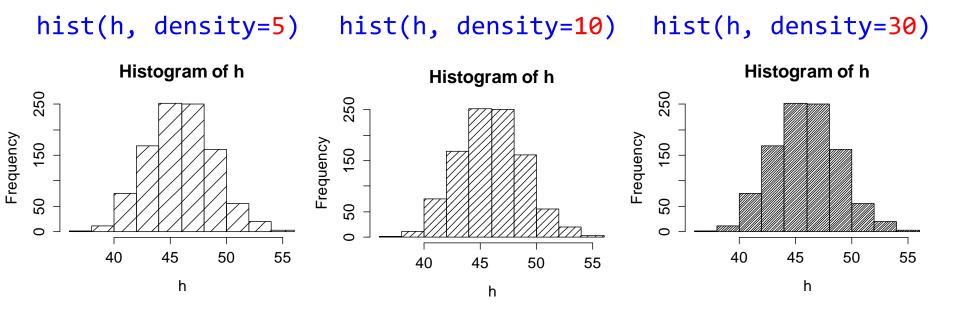
- Es una función genérica para dibujar histogramas
- Toma como argumento un vector
- Por ejemplo:
 - 1. h = rnorm(1000, 46, 3)
 - 2. hist(h)



La opción breaks

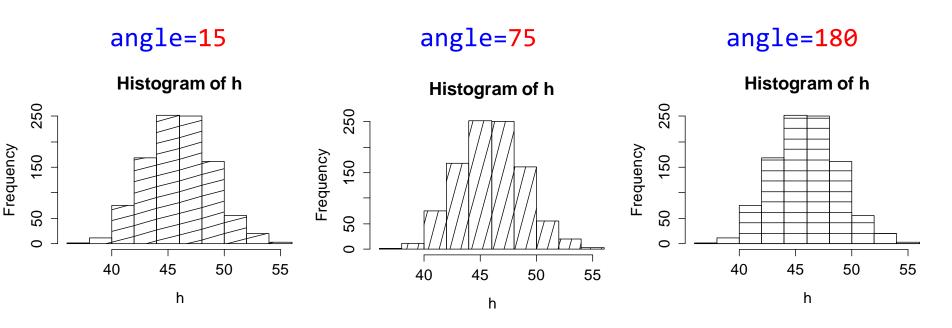


La opción density



La opción angle

density=5

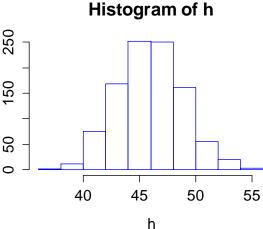


La opción border

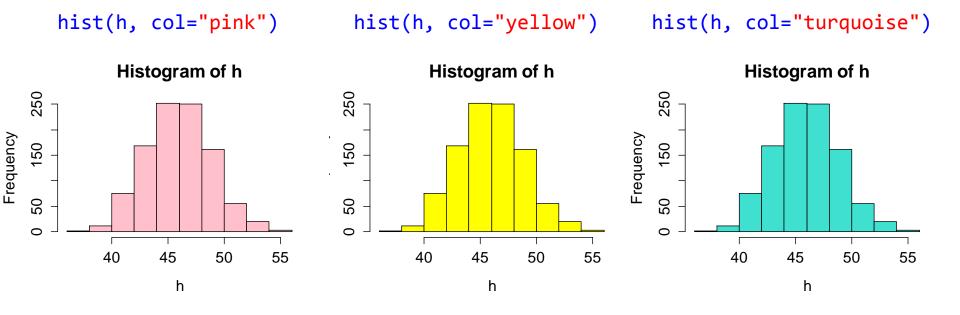
colors()

hist(h, border="red") hist(h, border="green") Histogram of h Histogram of h 250 250 250 Frequency Frequency 150 150 150 50 50 50 45 50 55 40 45 50 40 55 40 h h

hist(h, border="blue")



La opción col



Varias opciones en un mismo histograma

 Las opciones que hemos visto (y más) se pueden combinar

Ejemplo:

```
Histograma

Lecnencia

40 45 50 55

h
```

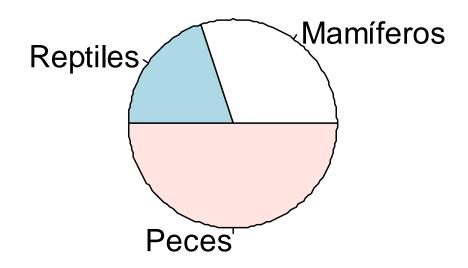
```
hist(h,
          col = "limegreen",
          border = "blue",
          density = 15,
          angle = 180,
          main = "Histograma",
          ylab = "Frecuencia"
          )
```

La función pie()

- Sirve para crear gráficos de sectores
- El argumento principal es x
 - Vector de cantidades no negativas áreas del gráfico
- Otro argumento interesante es labels
- Ejemplo:

```
pie(x=c(3,2,5), labels=c("Mamíferos", "Reptiles", "Peces"))
```

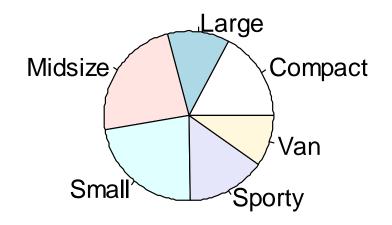
La función pie()



```
pie(x=c(3,2,5), labels=c("Mamíferos", "Reptiles", "Peces"))
```

La función pie()

- ¿Cómo funciona con bases de datos?
- Tenemos que crear una tabla de frecuencias
- Ejemplo:

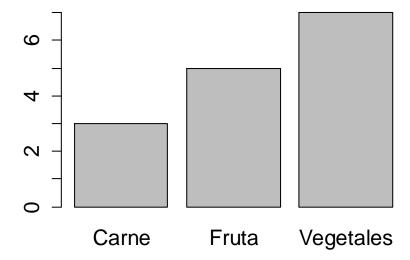


```
data(Cars93, package="MASS")
attach(Cars93)
pie(table(Type))
```

La función barplot()

- Sirve para generar gráficos de barras
- El argumento principal son alturas de barras height
- El argumento que controla las etiquetas de las barras es names.arg
- Ejemplo:

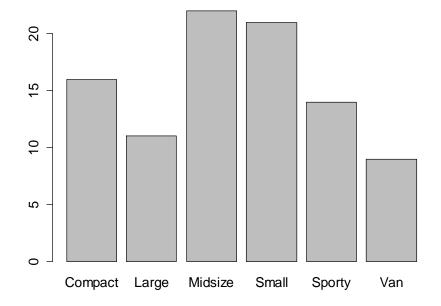
barplot(height=c(3,5,7),names.arg=c("Carne","Fruta","Vegetales"))



La función barplot()

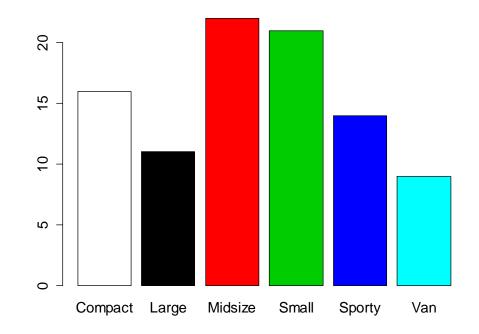
- ¿Cómo lo haríamos desde una base de datos?
- Tendríamos que generar una tabla de frecuencias
- Ejemplo:
- Utilizando la base de datos Cars93 y la variable Type

barplot(table(Type))



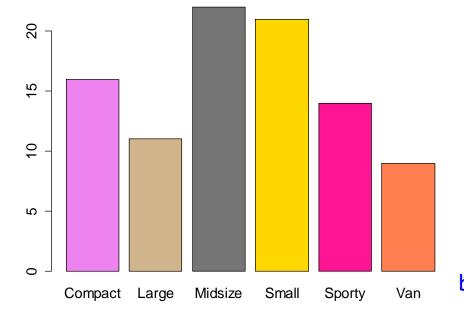
Dar color

- Con números del 0 al 8
- Con su nombre lista en colors()
- Paletas específicas
 - rainbow
 - topo.colors
 - cm.colors

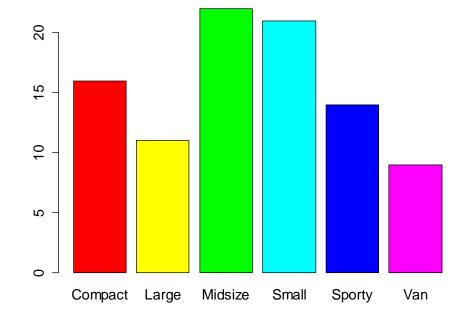


barplot(table(Type), col=0:5)

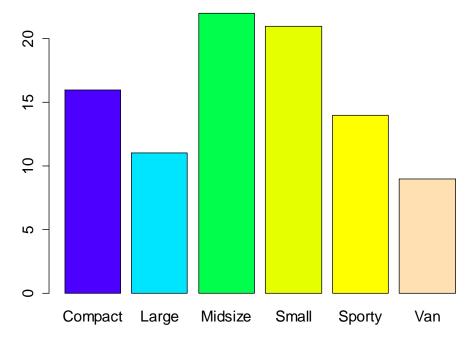
barplot(table(Type), col=c("violet","tan","gray46","gold","deeppink","coral"))



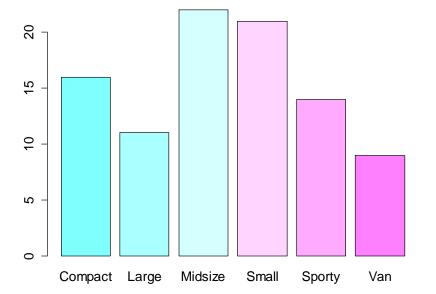
barplot(table(Type), col=rainbow(6))



barplot(table(Type), col=topo.colors(6))



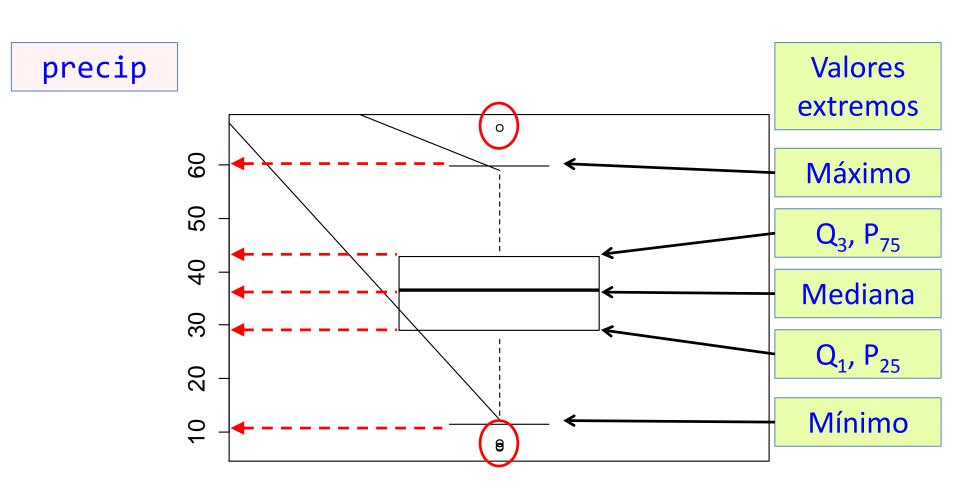
barplot(table(Type), col=cm.colors(6))



¿Qué es un diagrama de caja?

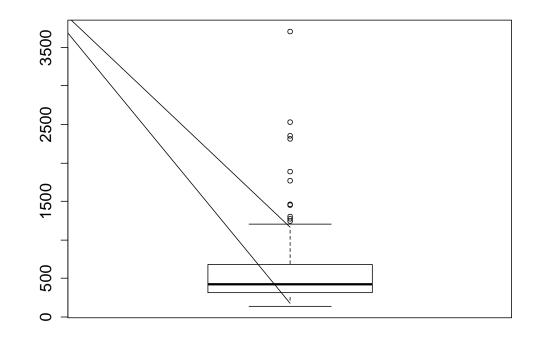
- Una caja que con sus bordes delimita ciertos cuantiles
- En concreto señala dónde están
 - El mínimo bajo ciertas condiciones
 - El máximo bajo ciertas condiciones
 - Los cuartiles Percentiles 25, 50 y 75
 - Datos extremos, alejados o atípicos
- Ejemplo:

¿Qué es un diagrama de caja?



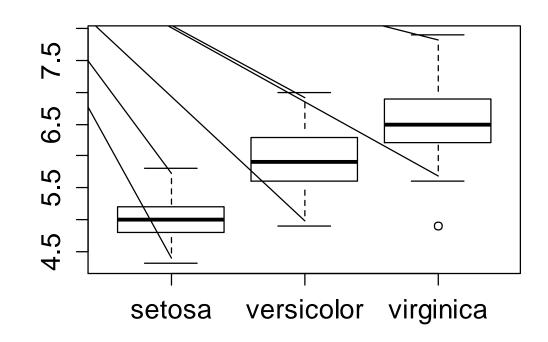
La función boxplot()

- Sirve para dibujar gráficos de caja
- Su principal argumento es x
- Ejemplo: boxplot(rivers)



La función boxplot()

- Puede tomar una interfaz de tipo fórmula
- Ejemplo:

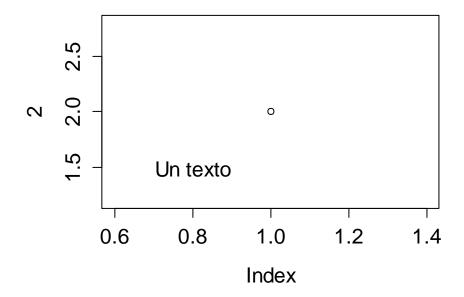


boxplot(Sepal.Length ~ Species)

La función text()

Sirve para añadir texto en coordenadas concretas

text(0.8, 1.5, "Un texto")



La función points()

Para añadir puntos en coordenadas concretas

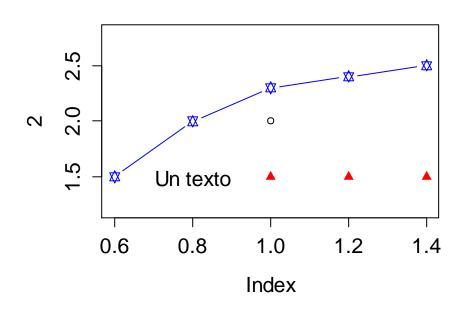
```
points(
  x = c(1,1.2,1.4),

y = c(rep(1.5,3)),
  pch = 17,
  col = "red"
                                           Un texto
                                       0.6
                                             8.0
                                                   1.0
                                                                1.4
                                                  Index
```

La función lines()

Sirve para añadir líneas

```
lines(
    x = c(0.6,0.8,1,1.2,1.4),
    y = c(1.5,2,2.3,2.4,2.5),
    type = "b",
    col = 4,
    pch = 11
)
```



La función abline()

Añade una línea con intercepto y pendiente arbitrarios

```
abline(

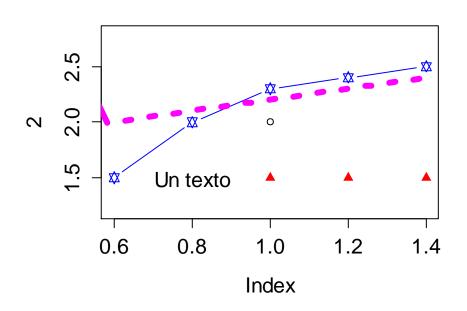
a = 1.7,

b = 0.5,

lwd= 6,

lty = 3,

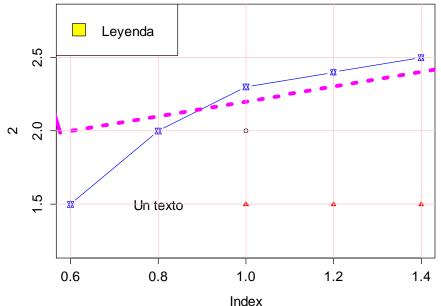
col = 6
```



La función legend()

Sirve para añadir leyendas

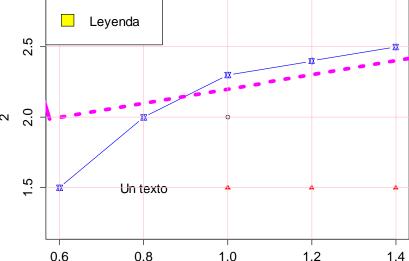
```
legend(
  x = "topleft",
  legend = "Leyenda",
  fill = "yellow"
)
```



La función title()

Sirve para añadir o modificar títulos

```
title(
  main = "Es una prueba",
  sub = "Base",
  col.main = "green"
)
```



Index

Es una prueba