

2025

Proyecto



José Miguel Grande Abrajan	202123224
Orquídea Rubí Palacios Godínez	202246635
Mancilla Carbajal Miguel Antonio	202242736

1-12-2025

Resumen

El presente documento detalla la ejecución de un proyecto integral de Minería de Datos utilizando la suite de software Weka. El objetivo principal fue aplicar y evaluar técnicas de preprocesamiento, clasificación supervisada y agrupamiento no supervisado sobre la base de datos "optdigits", la cual contiene información numérica sobre el reconocimiento óptico de dígitos manuscritos.

La metodología empleada abarcó tres fases principales: primero, la reducción de dimensionalidad mediante algoritmos de extracción y selección de características (PrincipalComponents, CorrelationAttributeEval y GainRatioAttributeEval); segundo, la evaluación comparativa de tres clasificadores (J48, NaiveBayes e IBk) sobre las bases de datos original y reducidas; y tercero, la aplicación de técnicas de clustering ignorando la etiqueta de clase.

Para validar los resultados, se utilizaron las herramientas Explorer, Experimenter y KnowledgeFlow. Los resultados estadísticos (T-Test) demostraron que el algoritmo IBk (K-Vecinos Cercanos) ofrece el mejor rendimiento para este dominio, alcanzando una precisión superior al 98%. Asimismo, se concluyó que las técnicas de selección de atributos (Correlation y GainRatio) son efectivas para reducir el costo computacional sin sacrificar la capacidad predictiva de los modelos, a diferencia del análisis de componentes principales (PCA), que en este caso específico presentó una leve disminución en la precisión.

Introducción

En el ámbito de la Ciencia de Datos y el Aprendizaje Automático (Machine Learning), contar con plataformas que permitan preprocesar información y evaluar modelos de manera eficiente es crucial. Para el desarrollo de este proyecto, se ha utilizado Weka (Waikato Environment for Knowledge Analysis), una suite de software de código abierto desarrollada en la Universidad de Waikato, Nueva Zelanda. Escrito en Java, Weka integra una colección completa de algoritmos de aprendizaje automático para tareas de minería de datos, permitiendo al usuario aplicar métodos de preprocesamiento, clasificación, regresión, agrupamiento (clustering) y reglas de asociación directamente sobre sus conjuntos de datos.

Para cumplir con los objetivos de este análisis, hemos hecho uso de las tres interfaces principales que ofrece este entorno:

- **Weka Explorer (Explorador):** Es la interfaz principal y más utilizada para el análisis exploratorio de datos. Funciona como un entorno de "laboratorio" donde se cargan los datos crudos para procesarlos paso a paso. En este proyecto, el Explorer fue fundamental para la fase de Preprocesamiento (limpieza y reducción de dimensiones),

así como para realizar pruebas iniciales rápidas de los algoritmos de clasificación y agrupamiento. Nos permite visualizar la distribución de los atributos y evaluar el rendimiento de los modelos de forma individual e inmediata mediante matrices de confusión y estadísticas de error.

- **Weka Experimenter (Experimentador):** A diferencia del Explorer, que evalúa un modelo a la vez, el Experimenter es una herramienta diseñada para la automatización y la validación estadística rigurosa. Su función principal en este proyecto fue permitirnos comparar objetivamente el desempeño de múltiples algoritmos (como J48, NaiveBayes e IBk) sobre diferentes versiones de la base de datos (original y reducidas) de manera simultánea. Esta herramienta ejecuta pruebas de validación cruzada y aplica tests de significancia estadística (T-tests) para determinar si las diferencias en la precisión de los algoritmos son producto del azar o si un modelo es matemáticamente superior a otro.
- **Weka KnowledgeFlow (Flujo de Conocimiento):** Esta interfaz ofrece una perspectiva visual del proceso de minería de datos basada en el paradigma de "flujo de datos". En lugar de pestañas estáticas, utiliza un lienzo donde el usuario arrastra y conecta componentes (como cargadores de datos, filtros y evaluadores) para diseñar la tubería de procesamiento. En nuestra práctica, el KnowledgeFlow se utilizó para modelar gráficamente y replicar el procedimiento completo de agrupamiento (Clustering), demostrando cómo los datos fluyen desde su carga, pasan por la eliminación de la clase (filtro) y llegan al entrenamiento y evaluación del algoritmo, ofreciendo una representación clara de la arquitectura del proceso.

Descripción de las Técnicas de Minería de Datos

Para el análisis de la información contenida en la base de datos, se han implementado tres familias fundamentales de técnicas de minería de datos, cada una con un propósito específico en el flujo de descubrimiento de conocimiento:

1. Extracción y Selección de Características (Feature Extraction/Selection) Esta técnica, también conocida como reducción de dimensionalidad, consiste en el proceso de transformar o seleccionar un subconjunto de los atributos más relevantes de la base de datos original. En conjuntos de datos con un gran número de variables (como es el caso de optdigits), mucha información puede ser redundante, irrelevante o constituir "ruido" que confunde a los modelos. El objetivo de esta técnica es reducir la complejidad del espacio de datos sin perder la información esencial. Al aplicar algoritmos como el Análisis de Componentes Principales (PCA) o evaluaciones de correlación, logramos generar bases de datos "reducidas" que permiten un entrenamiento más rápido y, a menudo, modelos más precisos al eliminar variables confusas.

2. Clasificación (Classification) La clasificación es una técnica de Aprendizaje Supervisado. Su objetivo es construir un modelo matemático capaz de predecir la etiqueta de clase de una instancia nueva basándose en sus atributos. Para que esto funcione, el algoritmo "aprende" primero analizando un conjunto de datos de entrenamiento donde ya se conoce la respuesta

correcta (en nuestro caso, saber qué dígito representa cada imagen). Una vez entrenado, el modelo busca patrones generales que relacionan los atributos numéricos con la clase objetivo. En este proyecto, se aplicaron algoritmos de clasificación sobre las bases de datos original y reducidas para verificar qué tan bien podía el sistema distinguir entre los diferentes dígitos manuscritos.

3. Agrupamiento (Clustering) A diferencia de la clasificación, el agrupamiento es una técnica de Aprendizaje No Supervisado. Esto significa que el algoritmo no cuenta con una etiqueta de clase predefinida ni un "maestro" que le indique la respuesta correcta. El objetivo del clustering es explorar la estructura natural de los datos para encontrar grupos (clústers) de instancias que sean muy similares entre sí, pero muy diferentes a las instancias de otros grupos. En nuestra práctica, para aplicar correctamente esta técnica, fue necesario ignorar o remover el atributo de la clase, forzando al algoritmo a agrupar los datos basándose puramente en la similitud matemática de sus atributos.

Descripción de los Algoritmos Utilizados

A continuación, se describen los fundamentos teóricos de los algoritmos seleccionados para cada etapa del proceso de minería de datos en Weka.

1. Algoritmos para la Extracción y Selección de Características

Para reducir la dimensionalidad de la base de datos optdigits y optimizar el rendimiento computacional, se utilizaron tres evaluadores distintos acompañados del método de búsqueda **Ranker**:

- **PrincipalComponents (Análisis de Componentes Principales - PCA):**

Este algoritmo no simplemente selecciona atributos, sino que transforma el espacio de datos original. Realiza una transformación lineal ortogonal para convertir un conjunto de observaciones de variables posiblemente correlacionadas en un conjunto de valores de variables linealmente no correlacionadas llamadas "componentes principales". Su objetivo es proyectar los datos en un espacio de menor dimensión maximizando la varianza de los datos, conservando así la mayor cantidad de información posible con el menor número de atributos nuevos.

- **CorrelationAttributeEval:**

Este algoritmo evalúa la importancia de un atributo midiendo la correlación (generalmente de Pearson) entre este y la clase objetivo. La premisa es que los buenos atributos están altamente correlacionados con la clase, pero no correlacionados entre sí. Junto con el método de búsqueda Ranker, ordena los atributos desde el más correlacionado al menos correlacionado, permitiendo descartar aquellos que tienen una relación débil con el resultado.

- **GainRatioAttributeEval:**

Basado en la teoría de la información, este algoritmo evalúa el valor de un atributo calculando su "Ratio de Ganancia" (Gain Ratio) con respecto a la clase. Es una mejora sobre la medida de Information Gain (Ganancia de Información) que compensa el sesgo hacia atributos que tienen muchos valores distintos. Mide cuánto disminuye la entropía (incertidumbre) del sistema al conocer el valor del atributo.

2. Algoritmos de Clasificación

Para la tarea de predicción supervisada de los dígitos, se eligieron tres algoritmos representativos de diferentes paradigmas de aprendizaje:

- **J48 (Árboles de Decisión C4.5):**

J48 es la implementación en Java del famoso algoritmo C4.5. Es un clasificador basado en árboles de decisión. Funciona dividiendo recursivamente el conjunto de datos basándose en el atributo que proporciona la mayor ganancia de información normalizada en cada nodo. El resultado es una estructura en forma de árbol similar a un diagrama de flujo, donde cada nodo

interno representa una prueba en un atributo, cada rama representa el resultado de la prueba y cada hoja representa una etiqueta de clase. Es altamente interpretable por humanos.

- **NaiveBayes (Bayes Ingenuo):**

Es un clasificador probabilístico fundamentado en el Teorema de Bayes. Se denomina "ingenuo" (Naive) porque asume una fuerte independencia entre los atributos; es decir, asume que la presencia de una característica particular en una clase no está relacionada con la presencia de ninguna otra característica. A pesar de esta simplificación teórica, es extremadamente rápido y muy efectivo en problemas de alta dimensión como el reconocimiento de patrones.

- **IBk:**

Conocido comúnmente como K-Vecinos Más Cercanos (KNN). Es un algoritmo de aprendizaje "vago" (lazy learning), lo que significa que no construye un modelo interno explícito durante el entrenamiento. Simplemente almacena las instancias de entrenamiento y, para clasificar una nueva instancia, busca las "k" instancias más cercanas (usando una métrica de distancia, como la Euclidiana) en el espacio de características y asigna la clase más común entre esos vecinos.

3. Algoritmos de Agrupamiento (Clustering)

Para el análisis no supervisado, se utilizaron algoritmos particionales para encontrar estructuras subyacentes en los datos sin usar la etiqueta de clase:

- **SimpleKMeans:**

Es uno de los algoritmos de clustering más utilizados. Su objetivo es particionar las n observaciones en k grupos en el que cada observación pertenece al grupo cuyo valor medio (centroide) es el más cercano. El algoritmo funciona iterativamente: primero asigna centroides aleatorios, luego asigna cada punto a su centroide más cercano y finalmente recalcula los centroides basándose en el promedio de los puntos asignados, repitiendo el proceso hasta que los centroides se estabilizan.

- **FarthestFirst:**

Es una variante heurística del algoritmo K-Means diseñada para la velocidad. En lugar de seleccionar centroides aleatorios iniciales y reajustarlos iterativamente, coloca el primer centroide aleatoriamente y selecciona cada centroide subsiguiente como el punto que está más alejado de los centroides ya existentes. Es extremadamente rápido, ideal para bases de datos grandes o como paso preliminar.

- **Canopy:**

Es un algoritmo de agrupamiento a menudo utilizado como un "preprocesador" para K-Means. Utiliza dos umbrales de distancia (T_1 y T_2) para crear grupos superpuestos llamados

"canopies". Es muy eficiente computacionalmente y ayuda a estimar el número óptimo de clústers o a proporcionar puntos de inicio más inteligentes para algoritmos más costosos como K-Means.

Descripción de la Base de Datos Utilizada

Para la realización de este proyecto experimental, se seleccionó la base de datos "Optical Recognition of Handwritten Digits" (optdigits), un conjunto de datos clásico y ampliamente utilizado en la comunidad de aprendizaje automático, disponible originalmente en el repositorio de la UCI (University of California, Irvine).

Características Técnicas:

- **Volumen de Datos:** La base de datos contiene un total de **5,620** instancias (muestras). Este tamaño fue seleccionado estratégicamente para el proyecto, ya que proporciona suficientes datos para validar estadísticamente los resultados sin comprometer el tiempo de procesamiento limitado para la ejecución de múltiples algoritmos.
- **Dimensionalidad:** Cuenta con 64 atributos de entrada (predictivos) y 1 atributo de clase.
- **Tipo de Atributos:** Todos los atributos de entrada son **numéricos (enteros)**.

Atributo de Clase (Target):

El objetivo de la base de datos es identificar el dígito escrito. Por lo tanto, el atributo de clase es nominal y posee 10 valores posibles: {0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

Justificación de la Elección:

Se eligió esta base de datos por dos razones principales:

1. **Naturaleza Numérica:** Al ser una base de datos puramente numérica, facilita la aplicación directa de algoritmos matemáticos como *PrincipalComponents* y *K-Means* sin necesidad de transformaciones complejas previas.
2. **Alta Dimensionalidad:** Al tener 64 atributos, es un candidato ideal para demostrar la efectividad de las técnicas de **reducción de características** (paso clave de este proyecto), permitiendo observar claramente cómo se puede comprimir la información manteniendo la capacidad predictiva.

El desarrollo del proyecto se llevó a cabo siguiendo una metodología estructurada en cinco fases secuenciales, utilizando las diferentes herramientas proporcionadas por el entorno Weka. A continuación, se detalla el procedimiento técnico ejecutado en cada etapa:

Fase 1: Preprocesamiento y Reducción de Dimensionalidad (Weka Explorer)

Como paso inicial, se cargó la base de datos original optdigits.arff en el módulo **Explorer**. Con el objetivo de optimizar el conjunto de datos y eliminar redundancia entre los 64 atributos de entrada, se aplicaron tres técnicas distintas de selección y extracción de características:

1. **Reducción por PCA:** Se utilizó el evaluador PrincipalComponents junto con el método de búsqueda Ranker. Esto transformó el espacio de atributos original en nuevos componentes ortogonales. Se guardó el conjunto resultante como una nueva base de datos reducida.
2. **Reducción por Correlación:** Se aplicó el evaluador CorrelationAttributeEval para identificar los atributos con mayor correlación lineal respecto a la clase. Se generó un segundo archivo reducido conservando únicamente los atributos mejor rankeados.
3. **Reducción por Ratio de Ganancia:** Se empleó GainRatioAttributeEval para medir la ganancia de información de cada píxel. Se generó un tercer archivo reducido con esta selección.

Como resultado de esta fase, se obtuvieron cuatro versiones del conjunto de datos: la Original y tres versiones Reducidas (PCA, Correlation, GainRatio), listas para ser comparadas.

Fase 2: Ejecución de Modelos de Clasificación (Weka Explorer)

En esta etapa se evaluó el rendimiento predictivo. Para garantizar la validez de los resultados, se utilizó el esquema de evaluación Cross-Validation (Validación Cruzada) de 10 pliegues (folds), que asegura que el modelo sea probado en datos que no vio durante el entrenamiento.

Se ejecutaron tres algoritmos de clasificación de familias distintas: J48 (Árboles), NaiveBayes (Probabilístico) y IBk (Instancias). Este procedimiento se repitió sistemáticamente para cada una de las 4 bases de datos generadas en la fase anterior, registrando métricas de precisión (Accuracy) y matrices de confusión para su posterior análisis.

Fase 3: Validación Estadística Automatizada (Weka Experimenter)

Para dar rigor científico a la comparación, se trasladó el flujo de trabajo al entorno **Experimenter**.

- **Configuración:** Se definió un experimento cargando las 4 bases de datos (datasets) y los 3 algoritmos seleccionados. Se configuró la salida de resultados a un archivo .arff para su registro permanente.

- **Ejecución:** El sistema ejecutó automáticamente 10 repeticiones de la validación cruzada de 10 pliegues para cada combinación de algoritmo-dataset, generando una muestra estadística robusta.
- **Análisis:** Se aplicó una prueba de significancia estadística (Paired T-Test con un nivel de confianza del 95% o 0.05) utilizando como base de comparación el dataset original o el algoritmo estándar, permitiendo identificar matemáticamente qué combinaciones ofrecían una mejora real ("v") o un deterioro significativo ("*") en el rendimiento.

Fase 4: Análisis de Agrupamiento (Clustering)

Se procedió a realizar un análisis no supervisado utilizando la pestaña Clúster del Explorer.

- **Preparación:** Dado que el clustering busca patrones naturales sin etiquetas predefinidas, se configuró el entorno para ignorar el atributo de clase ("Ignore attributes") durante el entrenamiento.
- **Algoritmos:** Se aplicaron los algoritmos **SimpleKMeans**, **FarthestFirst** y **Canopy** sobre la base de datos original y las reducidas.
- **Evaluación:** Se utilizó el modo "Classes to clústers evaluation" para contrastar cómo los grupos formados matemáticamente coincidían con los dígitos reales de la base de datos.

Fase 5: Modelado de Flujo de Datos (KnowledgeFlow)

Finalmente, se diseñó un pipeline visual para replicar y automatizar el proceso de agrupamiento. Se construyó un diagrama de flujo conectando los siguientes componentes:

1. **ArffLoader:** Para la ingesta de datos.
2. **ClassAssigner / Remove Filter:** Configurado específicamente para eliminar el atributo de la clase (índice last), simulando un escenario real de clustering no supervisado.
3. **TrainingSetMaker:** Para generar los lotes de datos de entrenamiento.
4. **Clusterer (SimpleKMeans):** Para ejecutar el agrupamiento.
5. **ClustererPerformanceEvaluator:** Para medir la calidad de los clústers.
6. **TextViewer:** Para visualizar los resultados finales y los centroides generados.

Este flujo permitió verificar la reproducibilidad del experimento de manera gráfica y modular.

Algoritmo	Base Original (%)	Base Reducida PCA (%)	Base Reducida Corr (%)	Base Reducida Gain (%)
J48	90.694 %	87.3665 %	90.694 %	90.694 %
NaiveBayes	91.3345 %	87.4911 %	91.3345 %	91.3345 %
IBk (KNN)	98.6121 %	96.0498 %	98.6121 %	98.6121 %

Dataset	(1) trees.J4 (2) lazy. (3) bayes			
optdigits	(100)	90.52	98.70 v	91.39
'optdigits-weka.filters.u(100)		90.52	98.70 v	91.39
'optdigits-weka.filters.u(100)		90.52	98.70 v	91.39
'optdigits-weka.filters.s(100)		87.41	96.12 v	87.42

		(v/ /*)	(4/0/0)	(0/4/0)

Comparación Estadística de los Modelos de Clasificación

Para validar rigurosamente el desempeño de los modelos, se ejecutó una prueba de significancia estadística (Paired T-Test) con un nivel de confianza del 95% utilizando el módulo Experimenter de Weka.

La configuración del experimento estableció al algoritmo **J48 (Árboles de Decisión)** como la **base de comparación** (Columna 1), contrastando su rendimiento contra **IBk** (Columna 2) y **NaiveBayes** (Columna 3).

Interpretación de los Resultados:

La tabla de resultados arrojó los siguientes hallazgos estadísticos (donde v indica una mejora significativa y el espacio vacío indica similitud estadística):

1. Análisis Comparativo de Algoritmos:

- J48 (Base) vs. IBk (K-Vecinos Cercanos):

El algoritmo IBk demostró una superioridad estadística contundente en todos los escenarios (marcado con el símbolo v). Mientras que J48 obtuvo una precisión promedio de ~90.5% en la base original, IBk alcanzó un 98.70%.

- *Interpretación:* Esto confirma que, para la base de datos optdigits (imágenes pixeladas), un enfoque basado en instancias y distancia (como IBk) es mucho más efectivo que un enfoque basado en reglas de decisión (árboles), logrando clasificar los dígitos casi a la perfección.

- J48 (Base) vs. NaiveBayes:

Al comparar J48 contra NaiveBayes, la ausencia de símbolos (v o *) indica que no existe una diferencia estadísticamente significativa entre ellos. Aunque NaiveBayes muestra números absolutos ligeramente superiores (91.39% frente a 90.52%), matemáticamente se considera un "empate técnico". Ambos algoritmos tienen un desempeño competente, pero inferior a IBk.

2. Impacto de las Técnicas de Reducción:

- Efectividad de Correlation y GainRatio:

Al observar las filas correspondientes a CorrelationAttributeEval y GainRatioAttributeEval, los resultados numéricos son idénticos a los de la base Original para los tres algoritmos.

- *Conclusión:* Esto indica una selección de características ideal. Se logró reducir la dimensionalidad de la base de datos eliminando atributos irrelevantes sin perder ni una fracción de capacidad predictiva.

- Desempeño con PrincipalComponents (PCA):

La técnica PCA (última fila) provocó una caída generalizada en el rendimiento.

- El algoritmo base **J48** cayó de 90.52% a **87.41%**.
- El algoritmo **IBk** cayó de 98.70% a **96.12%**.
- *Conclusión:* La transformación geométrica de los datos realizada por PCA, en este caso específico, dificultó ligeramente la tarea de clasificación en comparación con trabajar con los píxeles originales o filtrados.

NAIVEBAYES

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""**

Test options:
☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds **10**
☐ Percentage split % **66**
More options...

(Nom) class: **Start** **Stop**

Result list (right-click for options):
13:42:29 - treesJ48
13:42:47 - **bayes.NaiveBayes**
13:42:58 - lazyIBk

Classifier output

```
Incorrectly classified instances 50 / 5620 %
Kappa statistic 0.9037
Mean absolute error 0.0174
Root mean squared error 0.1265
Relative absolute error 9.6604 %
Root relative squared error 42.1705 %
Total Number of Instances 5620
```

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.977	0.000	0.998	0.977	0.987	0.986	0.999	0.995	0
	0.825	0.010	0.904	0.825	0.863	0.849	0.982	0.884	1
	0.908	0.008	0.928	0.908	0.918	0.910	0.986	0.950	2
	0.888	0.009	0.922	0.888	0.905	0.894	0.985	0.945	3
	0.923	0.008	0.926	0.923	0.924	0.916	0.977	0.946	4
	0.892	0.003	0.967	0.892	0.928	0.922	0.982	0.961	5
	0.970	0.002	0.978	0.970	0.974	0.971	0.995	0.987	6
	0.966	0.010	0.913	0.966	0.939	0.932	0.995	0.965	7
	0.919	0.026	0.795	0.919	0.853	0.838	0.987	0.923	8
	0.868	0.020	0.830	0.868	0.849	0.832	0.975	0.860	9
Weighted Avg.	0.913	0.010	0.916	0.913	0.914	0.905	0.986	0.942	

=== Confusion Matrix ===

```
a b c d e f g h i j <-- classified as
541 0 0 0 9 1 1 0 2 0 | a = 0
0 471 21 0 2 0 5 1 41 30 | b = 1
0 11 506 2 0 2 0 0 31 5 | c = 2
0 1 6 508 1 8 0 7 25 16 | d = 3
1 1 1 0 524 0 3 22 8 8 | e = 4
0 0 2 13 3 498 1 5 7 29 | f = 5
0 6 0 0 7 2 541 0 2 0 | g = 6
0 0 4 2 5 0 0 547 4 4 | h = 7
0 23 4 1 3 2 2 2 509 8 | i = 8
0 8 1 25 12 2 0 15 11 488 | j = 9
```

IBk

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""**

Test options:
☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds **10**
☐ Percentage split % **66**
More options...

(Nom) class: **Start** **Stop**

Result list (right-click for options):
13:42:29 - treesJ48
13:42:47 - bayes.NaiveBayes
13:42:58 - **lazyIBk**

Classifier output

```
Incorrectly classified instances 70 / 5620 %
Kappa statistic 0.9846
Mean absolute error 0.0031
Root mean squared error 0.0526
Relative absolute error 1.7363 %
Root relative squared error 17.5458 %
Total Number of Instances 5620
```

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.998	0.000	1.000	0.998	0.999	0.999	0.999	0.998	0
	0.993	0.005	0.959	0.993	0.976	0.973	0.994	0.954	1
	0.998	0.000	0.996	0.998	0.997	0.997	0.999	0.995	2
	0.983	0.002	0.983	0.983	0.983	0.981	0.990	0.967	3
	0.989	0.001	0.991	0.989	0.990	0.989	0.994	0.982	4
	0.984	0.001	0.989	0.984	0.987	0.985	0.991	0.975	5
	0.996	0.001	0.993	0.996	0.995	0.994	0.998	0.990	6
	0.993	0.001	0.988	0.993	0.990	0.989	0.996	0.981	7
	0.962	0.001	0.991	0.962	0.976	0.974	0.981	0.957	8
	0.964	0.003	0.973	0.964	0.969	0.965	0.981	0.942	9
Weighted Avg.	0.986	0.002	0.986	0.986	0.986	0.985	0.992	0.974	

=== Confusion Matrix ===

```
a b c d e f g h i j <-- classified as
553 0 0 0 0 0 0 1 0 0 0 | a = 0
0 567 1 0 0 0 0 1 1 1 | b = 1
0 1 556 0 0 0 0 0 0 0 | c = 2
0 1 1 562 0 3 0 2 1 2 | d = 3
0 2 0 0 562 0 1 1 0 2 | e = 4
0 0 0 2 0 549 1 0 1 5 | f = 5
0 2 0 0 0 0 0 556 0 0 | g = 6
0 0 0 1 1 0 0 562 0 2 | h = 7
0 16 0 1 0 0 1 0 533 3 | i = 8
0 2 0 6 4 3 0 3 2 542 | j = 9
```

optdigits_CorrelationAttributeEval.arff

treesJ48

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier: Choose IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

Test options

☐ Use training set

☐ Supplied test set Set...

☒ Cross-validation Folds 10

☐ Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)

- 13:42:29 - treesJ48
- 13:42:47 - bayes.NaiveBayes
- 13:42:58 - lazy.IBk
- 13:45:29 - treesJ48
- 13:45:42 - bayes.NaiveBayes
- 13:45:48 - lazy.IBk

Classifier output

Kappa statistic 0.8966

Mean absolute error 0.0204

Root mean squared error 0.131

Relative absolute error 11.3403 %

Root relative squared error 43.6503 %

Total Number of Instances 5620

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.966	0.005	0.959	0.966	0.962	0.958	0.983	0.950	0
	0.902	0.013	0.886	0.902	0.894	0.882	0.949	0.830	1
	0.923	0.009	0.915	0.923	0.919	0.910	0.960	0.866	2
	0.890	0.014	0.876	0.890	0.883	0.870	0.943	0.833	3
	0.912	0.011	0.904	0.912	0.908	0.898	0.958	0.845	4
	0.928	0.008	0.927	0.928	0.927	0.919	0.969	0.906	5
	0.950	0.005	0.958	0.950	0.954	0.949	0.976	0.924	6
	0.928	0.006	0.948	0.928	0.938	0.931	0.964	0.896	7
	0.843	0.017	0.846	0.843	0.844	0.828	0.930	0.755	8
	0.829	0.016	0.852	0.829	0.840	0.823	0.910	0.711	9
Weighted Avg.	0.907	0.010	0.907	0.907	0.907	0.897	0.954	0.852	

=== Confusion Matrix ===

	a	b	c	d	e	f	g	h	i	j	<-- classified as
535	0	0	3	5	3	2	0	5	1		a = 0
2	515	5	10	6	2	5	4	15	7		b = 1
3	5	514	9	0	1	3	0	14	8		c = 2
2	3	7	509	3	11	1	3	14	19		d = 3
3	8	3	2	518	3	4	7	8	12		e = 4
2	6	1	7	3	518	2	0	4	15		f = 5
2	7	5	0	7	3	530	1	3	0		g = 6
1	5	1	6	6	1	1	525	11	9		h = 7
6	25	17	9	6	7	5	2	467	10		i = 8
2	7	9	26	19	10	0	12	11	466		j = 9

NaiveBayes

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier: Choose IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

Test options

☐ Use training set

☐ Supplied test set Set...

☒ Cross-validation Folds 10

☐ Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)

- 13:42:29 - treesJ48
- 13:42:47 - bayes.NaiveBayes
- 13:42:58 - lazy.IBk
- 13:45:29 - treesJ48
- 13:45:42 - bayes.NaiveBayes
- 13:45:48 - lazy.IBk

Classifier output

Kappa statistic 0.9037

Mean absolute error 0.0174

Root mean squared error 0.1265

Relative absolute error 9.6604 %

Root relative squared error 42.1705 %

Total Number of Instances 5620

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.977	0.000	0.998	0.977	0.987	0.986	0.999	0.995	0
	0.825	0.010	0.904	0.825	0.863	0.849	0.982	0.884	1
	0.908	0.008	0.928	0.908	0.918	0.910	0.986	0.950	2
	0.888	0.009	0.922	0.888	0.905	0.894	0.985	0.949	3
	0.923	0.008	0.926	0.923	0.924	0.916	0.977	0.946	4
	0.892	0.003	0.967	0.892	0.928	0.922	0.982	0.961	5
	0.970	0.002	0.978	0.970	0.974	0.971	0.995	0.987	6
	0.966	0.010	0.913	0.966	0.939	0.932	0.995	0.965	7
	0.919	0.026	0.795	0.919	0.853	0.838	0.987	0.923	8
	0.868	0.020	0.830	0.868	0.849	0.832	0.975	0.860	9
Weighted Avg.	0.913	0.010	0.916	0.913	0.914	0.905	0.986	0.942	

=== Confusion Matrix ===

	a	b	c	d	e	f	g	h	i	j	<-- classified as
541	0	0	0	9	1	1	0	2	0		a = 0
0	471	21	0	2	0	5	1	41	30		b = 1
0	11	506	2	0	2	0	0	31	5		c = 2
0	1	6	508	1	8	0	7	25	16		d = 3
1	1	1	0	524	0	3	22	8	8		e = 4
0	0	2	13	3	498	1	5	7	29		f = 5
0	6	0	0	7	2	541	0	2	0		g = 6
0	0	4	2	5	0	0	547	4	4		h = 7
0	23	4	1	3	2	2	2	509	8		i = 8
0	8	1	25	12	2	0	15	11	488		j = 9

IBk

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""**

Test options

☐ Use training set

☐ Supplied test set Set...

☒ Cross-validation Folds **10**

☐ Percentage split % **66**

More options...

(Nom) class

Start Stop

Result list (right-click for options)

- 13:42:29 - trees.J48
- 13:42:47 - bayes.NaiveBayes
- 13:42:58 - lazy.IBk
- 13:45:29 - trees.J48
- 13:45:42 - bayes.NaiveBayes
- 13:45:48 - lazy.IBk**

Classifier output

Kappa statistic 0.9846

Mean absolute error 0.0031

Root mean squared error 0.0526

Relative absolute error 1.7363 %

Root relative squared error 17.5458 %

Total Number of Instances 5620

=== Detailed Accuracy By Class ===

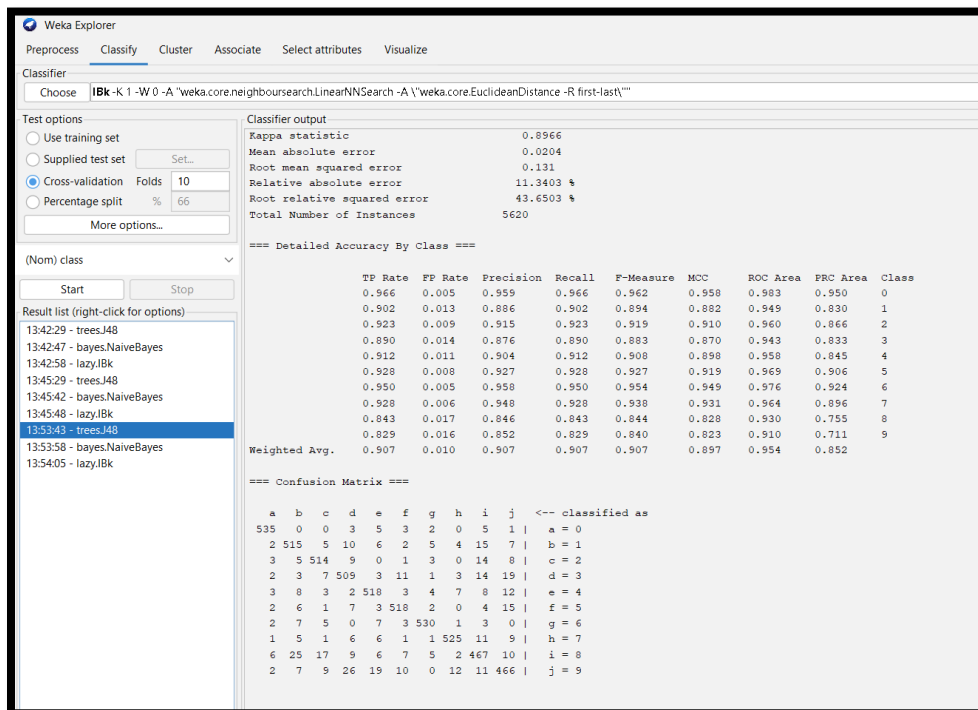
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.998	0.000	1.000	0.998	0.999	0.999	0.999	0.998	0
	0.993	0.005	0.959	0.993	0.976	0.973	0.994	0.954	1
	0.998	0.000	0.996	0.998	0.997	0.997	0.999	0.995	2
	0.983	0.002	0.983	0.983	0.983	0.981	0.990	0.967	3
	0.989	0.001	0.991	0.989	0.990	0.989	0.994	0.982	4
	0.984	0.001	0.989	0.984	0.987	0.985	0.991	0.975	5
	0.996	0.001	0.993	0.996	0.995	0.994	0.998	0.990	6
	0.993	0.001	0.988	0.993	0.990	0.989	0.996	0.981	7
	0.962	0.001	0.991	0.962	0.976	0.974	0.981	0.957	8
	0.964	0.003	0.973	0.964	0.969	0.965	0.981	0.942	9
Weighted Avg.	0.986	0.002	0.986	0.986	0.986	0.985	0.992	0.974	

=== Confusion Matrix ===

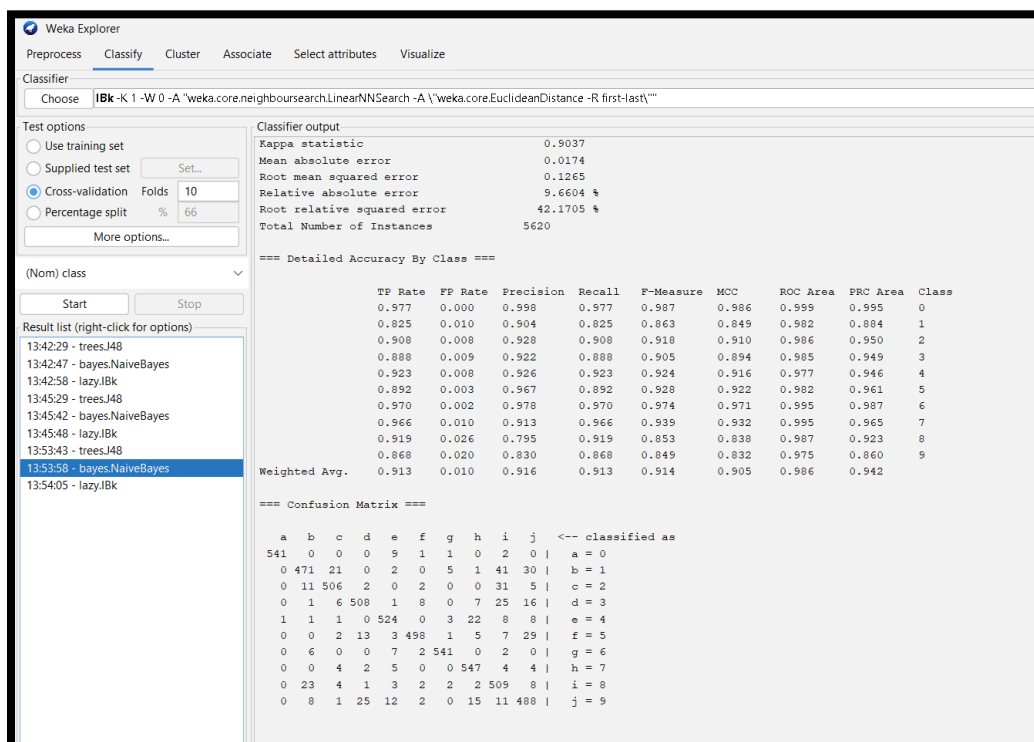
	a	b	c	d	e	f	g	h	i	j	<-- classified as
553	0	0	0	0	0	0	1	0	0	0	a = 0
0	567	1	0	0	0	0	0	1	1	1	b = 1
0	1	556	0	0	0	0	0	0	0	0	c = 2
0	1	1	562	0	3	0	2	1	2	1	d = 3
0	2	0	0	562	0	1	1	0	2	1	e = 4
0	0	0	2	0	549	1	0	1	5	1	f = 5
0	2	0	0	0	0	556	0	0	0	0	g = 6
0	0	0	1	1	0	0	562	0	2	1	h = 7
0	16	0	1	0	0	1	0	533	3	1	i = 8
0	2	0	6	4	3	0	3	2	542	1	j = 9

optdigits_GainRatioAttributeEval.arff

treesJ48



NaiveBayes



IBk

Weka Explorer

Preprocess

Classify

Cluster

Associate

Select attributes

Visualize

Classifier

Choose IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

Test options

Use training set

Supplied test set

Cross-validation Folds 10

Percentage split % 66

More options...

(Nom) class

Start Stop

Result list (right-click for options)

13:42:29 - trees.J48

13:42:47 - bayes.NaiveBayes

13:42:58 - lazy.IBk

13:45:29 - trees.J48

13:45:42 - bayes.NaiveBayes

13:45:48 - lazy.IBk

13:53:43 - trees.J48

13:53:58 - bayes.NaiveBayes

13:54:05 - lazy.IBk

Classifier output

Kappa statistic 0.9846

Mean absolute error 0.0031

Root mean squared error 0.0526

Relative absolute error 1.7363 %

Root relative squared error 17.5458 %

Total Number of Instances 5620

=== Detailed Accuracy By Class ===

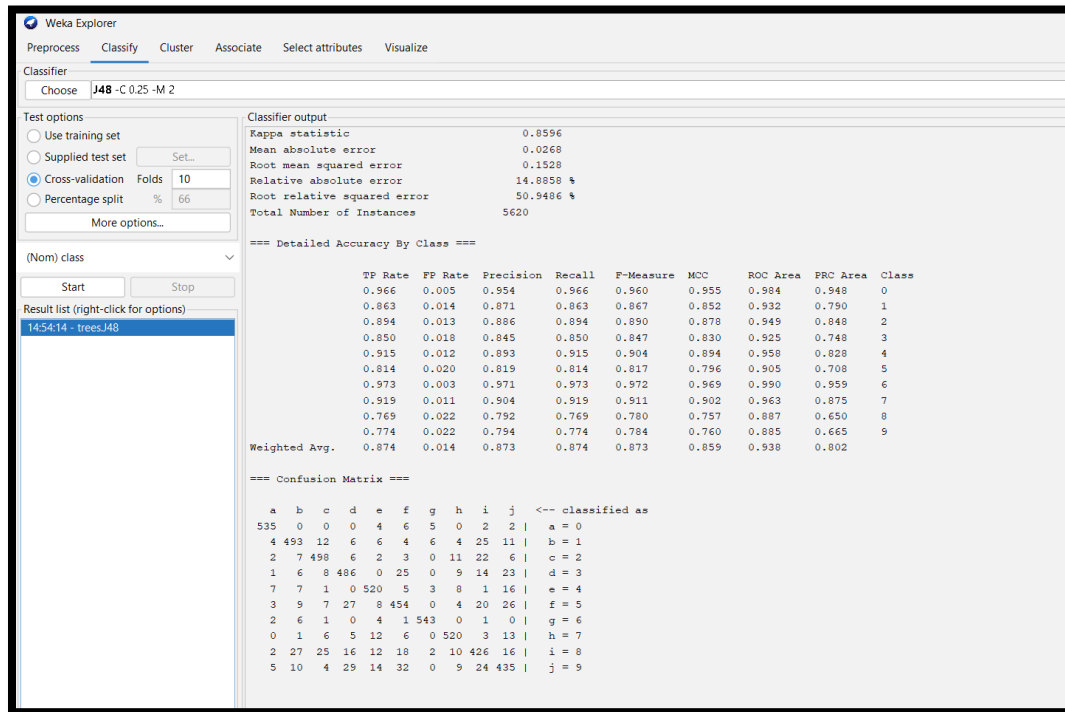
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.998	0.000	1.000	0.998	0.999	0.999	0.999	0.998	0
	0.993	0.005	0.959	0.993	0.976	0.973	0.994	0.954	1
	0.998	0.000	0.996	0.998	0.997	0.997	0.999	0.995	2
	0.983	0.002	0.983	0.983	0.983	0.981	0.990	0.967	3
	0.989	0.001	0.991	0.989	0.990	0.989	0.994	0.982	4
	0.984	0.001	0.989	0.984	0.987	0.985	0.991	0.975	5
	0.996	0.001	0.993	0.996	0.995	0.994	0.998	0.990	6
	0.993	0.001	0.988	0.993	0.990	0.989	0.996	0.981	7
	0.962	0.001	0.991	0.962	0.976	0.974	0.981	0.957	8
	0.964	0.003	0.973	0.964	0.969	0.965	0.981	0.942	9
Weighted Avg.	0.986	0.002	0.986	0.986	0.986	0.985	0.992	0.974	

=== Confusion Matrix ===

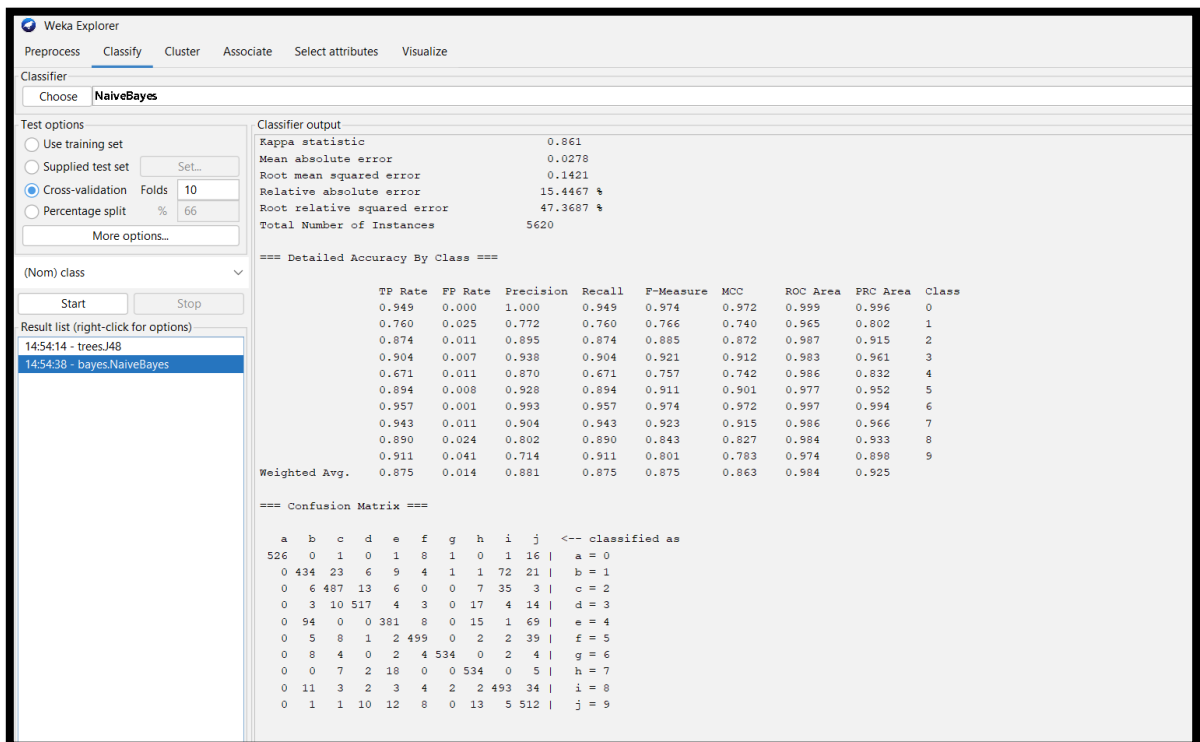
	a	b	c	d	e	f	g	h	i	j	<-- classified as
553	0	0	0	0	0	0	1	0	0	0	a = 0
0	567	1	0	0	0	0	1	1	1	1	b = 1
0	1	556	0	0	0	0	0	0	0	0	c = 2
0	1	1	562	0	3	0	2	1	2	1	d = 3
0	2	0	0	562	0	1	1	0	2	1	e = 4
0	0	0	2	0	549	1	0	1	5	1	f = 5
0	2	0	0	0	0	556	0	0	0	1	g = 6
0	0	0	1	1	0	0	562	0	2	1	h = 7
0	16	0	1	0	0	1	0	533	3	1	i = 8
0	2	0	6	4	3	0	3	2	542	1	j = 9

optdigits_Principal Components.arff

treesJ48



NaiveBayes



IBk

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Classifier: Choose **IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""**

Test options:

- ☐ Use training set
- ☐ Supplied test set Set...
- ☒ Cross-validation Folds **10**
- ☐ Percentage split % **66**

More options...

(Nom) class v

Start Stop

Result list (right-click for options)

- 14:54:14 - trees.J48
- 14:54:38 - bayes.NaiveBayes
- 14:54:53 - lazy.IBk**

Classifier output

Kappa statistic 0.9561
Mean absolute error 0.0082
Root mean squared error 0.0888
Relative absolute error 4.5778 %
Root relative squared error 29.5996 %
Total Number of Instances 5620

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.991	0.000	0.998	0.991	0.995	0.994	0.995	0.990	0
	0.982	0.007	0.938	0.982	0.960	0.955	0.988	0.923	1
	0.982	0.002	0.986	0.982	0.984	0.982	0.990	0.970	2
	0.951	0.009	0.925	0.951	0.938	0.931	0.971	0.885	3
	0.967	0.002	0.980	0.967	0.973	0.970	0.982	0.951	4
	0.948	0.004	0.964	0.948	0.956	0.951	0.972	0.919	5
	0.995	0.002	0.984	0.995	0.989	0.988	0.996	0.979	6
	0.981	0.002	0.981	0.981	0.981	0.978	0.989	0.963	7
	0.894	0.006	0.946	0.894	0.919	0.911	0.944	0.856	8
	0.915	0.010	0.907	0.915	0.911	0.901	0.952	0.838	9
Weighted Avg.	0.960	0.004	0.961	0.960	0.960	0.956	0.978	0.927	

=== Confusion Matrix ===

	a	b	c	d	e	f	g	h	i	j	<-- classified as
549	0	0	1	1	0	1	0	1	1	1	a = 0
0	561	2	1	1	0	1	1	3	1	1	b = 1
0	6	547	0	0	0	0	0	4	0	1	c = 2
0	3	2	544	0	1	0	1	9	12	1	d = 3
0	3	0	0	549	0	3	3	0	10	1	e = 4
1	0	0	10	2	529	0	0	3	13	1	f = 5
0	2	0	0	0	0	555	0	1	0	1	g = 6
0	0	0	1	2	0	0	555	0	8	1	h = 7
0	20	4	12	1	10	3	1	495	8	1	i = 8
0	3	0	19	4	9	1	5	7	514	1	j = 9

EXPERIMENTAL

Weka Experiment Environment

SetupRunAnalyse

Experiment Configuration ModeSimple

Open...Save...New

Results Destination

ARFF fileFilename:Browse...

Experiment Type

Cross-validation

Number of folds:10

☒ Classification☐ Regression

Iteration Control

Number of repetitions:10

☒ Data sets first☐ Algorithms first

Datasets

Add new...Edit selected...Delete selected

☐ Use relative paths

C:\Users\miguel\Desktop\ProyectoDeMineria\optdigits.arff
C:\Users\miguel\Desktop\ProyectoDeMineria\optdigits_CorrelationAt
C:\Users\miguel\Desktop\ProyectoDeMineria\optdigits_GainRatioAttr
C:\Users\miguel\Desktop\ProyectoDeMineria\optdigits_Principal Corr

UpDown

Algorithms

Add new...Edit selected...Delete selected

J48 -C 0.25 -M 2
NaiveBayes
IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"we

Load options...Save options...UpDown

Notes

Weka Experiment Environment

SetupRunAnalyse

Source

Got 1200 results

File...Data...

Actions

Perform testSave outputOpen Explorer...

Configure test

Testing withPaired T-Tester (corrected)

Select rows and colsRowsColsSwap

Comparison fieldPercent_correct

Significance0.05

Sorting (asc) by<default>

Test baseSelect

Displayed ColumnsSelect

Show std. deviations

Output FormatSelect

Result list

150631 - Available resultsets

150652 - Percent_correct - treesJ48 -C 0.25 -M 2 -217733168

Test output

st.PairedCorrectedTTester -G 4,5,6 -D 1 -R 2 -S 0.05 -result-matrix "weka.experiment.ResultMatrixPlainText" -mean-prec 2 -stddev-prec 2 -col-n

st

ted)

\$

(1) trees.J4 | (2) lazy. (3) bayes

(100) 90.52 | 98.70 v 91.39

(100) 90.52 | 98.70 v 91.39

(100) 90.52 | 98.70 v 91.39

(100) 87.41 | 96.12 v 87.42

(v/ /*) | (4/0/0) (0/4/0)

2* -217733168393644444

\\weka.core.neighboursearch.LinearNNSearch -A \\weka.core.EuclideanDistance -R first-last\\\\" -3080186098777067172

995231201785697655

Clustering

optdigits.arff

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose **FarthestFirst -N 2 -S 1**

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst

Clusterer output

```
input63      2.0936      2.1894      1.8741
input64      0.2541      0.3267      0.0878

Time taken to build model (full training data) : 0.28 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      3912 ( 70%)
1      1708 ( 30%)

Class attribute: class
Classes to Clusters:

0  1 <-- assigned to cluster
2 552 | 0
547 24 | 1
547 10 | 2
572  0 | 3
73 495 | 4
511 47 | 5
2 556 | 6
563  3 | 7
541 13 | 8
554  8 | 9

Cluster 0 <-- 3
Cluster 1 <-- 6

Incorrectly clustered instances :      4492.0      79.9288 %
```

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose **FarthestFirst -N 2 -S 1**

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst

Clusterer output

```
0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 <-- assigned to cluster
0  0  0  0  6  67  0  481  0  0  0  0  0  0  0  0  0  0  0  0  0
228 39 0 18 1 37 0 3 125 111 0 0 0 2 0 0 0 0 0 0 0 0 1
75 204 1 1 0 147 0 4 0 0 2 0 0 0 0 0 0 0 0 13 0 110 1 2
57 36 458 0 0 1 0 4 0 3 2 0 1 0 0 0 4 4 2 0 0 13
4 0 0 240 5 0 0 1 70 5 5 135 0 0 0 102 0 0 0 1 0 14
21 2 139 1 6 0 6 0 2 3 0 0 163 139 48 0 2 26 0 0 0 15
2 6 0 40 509 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 16
10 0 0 0 0 0 0 0 375 2 3 107 0 0 0 0 0 43 0 0 26 0 17
413 28 71 1 8 1 3 2 8 5 8 0 1 3 0 0 1 1 0 0 0 18
20 9 313 0 1 0 13 13 5 89 5 0 0 5 0 9 12 56 12 0 0 19

Cluster 0 <-- 8
Cluster 1 <-- 2
Cluster 2 <-- 3
Cluster 3 <-- 4
Cluster 4 <-- 6
Cluster 5 <-- No class
Cluster 6 <-- 0
Cluster 7 <-- 7
Cluster 8 <-- 1
Cluster 9 <-- 9
Cluster 10 <-- No class
Cluster 11 <-- No class
Cluster 12 <-- 5
Cluster 13 <-- No class
Cluster 14 <-- No class
Cluster 15 <-- No class
Cluster 16 <-- No class
Cluster 17 <-- No class
Cluster 18 <-- No class
Cluster 19 <-- No class
Cluster 20 <-- No class

Incorrectly clustered instances :      2563.0      45.605 %
```

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Clusterer

ChooseFarthestFirst - N 2 - S 1

Cluster mode

☐ Use training set

☐ Supplied test set

☐ Percentage split

☒ Classes to clusters evaluation

☒ Store clusters for visualization

Set...

%66

(Nom) class

Ignore attributes

StartStop

Result list (right-click for options)

152255 - SimpleKMeans

152301 - Canopy

152809 - FarthestFirst

Clusterer output

0.0 0.0 5.0 16.0 16.0 19.0 1.0 0.0 0.0 0.0 10.0 9.0 10.0 16.0 2.0 0.0 0.0 0.0 0.0 1.0 12.0 10.0 0.0 0.0 0.0 0.0 16.0 8.0 0.0 0.0 0.0 0.0 2.0

Cluster 1

0.0 0.0 0.0 3.0 14.0 7.0 0.0 0.0 0.0 0.0 12.0 12.0 0.0 3.0 8.0 0.0 0.0 9.0 14.0 0.0 2.0 16.0 6.0 0.0 4.0 16.0 2.0 0.0 10.0 14.0 0.0 1.0 15.0 8.0 3.0

Time taken to build model (full training data) : 0.1 seconds

=== Model and evaluation on training set ===

Clustered Instances

05110 (51%)

1510 (5%)

Class attribute: class

Classes to Clusters:

01 <-- assigned to cluster

5513 | 0

5710 | 1

5570 | 2

5720 | 3

147421 | 4

5580 | 5

47583 | 6

5642 | 7

5540 | 8

5611 | 9

Cluster 0 <-- 3

Cluster 1 <-- 4

Incorrectly clustered instances :4627.082.331 %

optdigits_CorrelationAttributeEval.arff

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer: Choose **FarthestFirst -N 2 -S 1**

Cluster mode:

- ☐ Use training set
- ☐ Supplied test set
- ☐ Percentage split %
- ☒ Classes to clusters evaluation (Nom) class
- ☒ Store clusters for visualization

Ignore attributes

Start

Result list (right-click for options):

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans**
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst

Cluster output:

```

input63      2.0936      2.1894      1.8741
input64      0.2541      0.3267      0.0878

Time taken to build model (full training data) : 0.17 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      3912 ( 70%)
1      1708 ( 30%)

Class attribute: class
Classes to Clusters:

0 1 <-- assigned to cluster
2 552 | 0
547 24 | 1
547 10 | 2
572 0 | 3
73 495 | 4
511 47 | 5
2 556 | 6
563 3 | 7
541 13 | 8
554 8 | 9

Cluster 0 <-- 3
Cluster 1 <-- 6

Incorrectly clustered instances :      4492.0      79.9288 %

```

Status

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer: Choose **FarthestFirst -N 2 -S 1**

Cluster mode:

- ☐ Use training set
- ☐ Supplied test set
- ☐ Percentage split %
- ☒ Classes to clusters evaluation (Nom) class
- ☒ Store clusters for visualization

Ignore attributes

Start

Result list (right-click for options):

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans**
- 17:32:01 - Canopy**
- 17:35:36 - FarthestFirst

Cluster output:

```

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 <-- assigned to cluster
0 0 0 6 67 0 481 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | 0
228 39 0 18 1 37 0 3 125 111 0 0 0 2 0 0 0 0 7 0 0 | 1
75 204 1 1 0 147 0 4 0 0 2 0 0 0 0 0 0 0 13 0 110 | 2
57 36 458 0 0 1 0 4 0 3 2 0 1 0 0 0 4 4 2 0 0 | 3
4 0 0 240 5 0 0 1 70 5 5 135 0 0 0 102 0 0 0 1 0 | 4
21 2 139 1 6 0 6 0 2 3 0 0 163 139 48 0 2 26 0 0 | 5
2 6 0 40 509 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | 6
10 0 0 0 0 0 0 0 375 2 3 107 0 0 0 0 0 43 0 0 26 | 7
413 28 71 1 8 1 3 2 8 5 8 0 1 3 0 0 1 1 0 0 0 | 8
20 9 313 0 1 0 13 13 5 89 5 0 0 5 0 9 12 56 12 0 0 | 9

Cluster 0 <-- 8
Cluster 1 <-- 2
Cluster 2 <-- 3
Cluster 3 <-- 4
Cluster 4 <-- 6
Cluster 5 <-- No class
Cluster 6 <-- 0
Cluster 7 <-- 7
Cluster 8 <-- 1
Cluster 9 <-- 9
Cluster 10 <-- No class
Cluster 11 <-- No class
Cluster 12 <-- 5
Cluster 13 <-- No class
Cluster 14 <-- No class
Cluster 15 <-- No class
Cluster 16 <-- No class
Cluster 17 <-- No class
Cluster 18 <-- No class
Cluster 19 <-- No class
Cluster 20 <-- No class

Incorrectly clustered instances :      2563.0      45.605 %

```


Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose **FarthestFirst -N 2 -S 1**

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation

(Nom) class ▾

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst**

Clusterer output

```
0.0 0.0 5.0 16.0 16.0 15.0 1.0 0.0 0.0 0.0 10.0 9.0 10.0 16.0 2.0 0.0 0.0 0.0 1.0
Cluster 1
0.0 0.0 0.0 3.0 14.0 7.0 0.0 0.0 0.0 0.0 0.0 12.0 12.0 0.0 3.0 8.0 0.0 0.0 14.0

Time taken to build model (full training data) : 0.05 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      5110 ( 91%)
1       510 (  9%)

Class attribute: class
Classes to Clusters:

0 1 <-- assigned to cluster
551 3 | 0
571 0 | 1
557 0 | 2
572 0 | 3
147 421 | 4
558 0 | 5
475 83 | 6
564 2 | 7
554 0 | 8
561 1 | 9

Cluster 0 <-- 3
Cluster 1 <-- 4

Incorrectly clustered instances :      4627.0   82.331  %
```

optdigits_GainRatioAttributeEval.arff

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer
Choose **FarthestFirst -N 2 -S 1**

Cluster mode
☐ Use training set
☐ Supplied test set Set...
☐ Percentage split % 66
☒ **Classes to clusters evaluation**
 (Nom) class
☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst
- 17:38:54 - SimpleKMeans**
- 17:39:04 - Canopy
- 17:42:36 - FarthestFirst

Cluster output

input63	2.0936	2.1894	1.8741
input64	0.2541	0.3267	0.0878

Time taken to build model (full training data) : 0.17 seconds

=== Model and evaluation on training set ===

Clustered Instances

0	3912 (70%)
1	1708 (30%)

Class attribute: class

Classes to Clusters:

0	1	<-- assigned to cluster
2	552	0
547	24	1
547	10	2
572	0	3
73	495	4
511	47	5
2	556	6
563	3	7
541	13	8
554	8	9

Cluster 0 <-- 3
Cluster 1 <-- 6

Incorrectly clustered instances : 4492.0 79.9288 %

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer
Choose **FarthestFirst -N 2 -S 1**

Cluster mode
☐ Use training set
☐ Supplied test set Set...
☐ Percentage split % 66
☒ **Classes to clusters evaluation**
 (Nom) class
☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst
- 17:38:54 - SimpleKMeans
- 17:39:04 - Canopy**
- 17:42:36 - FarthestFirst

Cluster output

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	<-- assign
0	0	0	6	67	0	481	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
228	39	0	18	1	37	0	3	125	111	0	0	0	0	2	0	0	0	7	0	0	1
75	204	1	1	0	147	0	4	0	0	2	0	0	0	0	0	0	0	13	0	110	12
57	36	458	0	0	1	0	4	0	3	2	0	1	0	0	0	4	4	2	0	0	13
4	0	0	240	5	0	0	1	70	5	5	135	0	0	0	102	0	0	0	1	0	4
21	2	139	1	6	0	6	0	2	3	0	0	163	139	48	0	2	26	0	0	0	15
2	6	0	40	509	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	16
10	0	0	0	0	0	0	375	2	3	107	0	0	0	0	0	43	0	0	26	0	17
413	28	71	1	8	1	3	2	8	5	8	0	1	3	0	0	1	1	0	0	0	18
20	9	313	0	1	0	13	13	5	89	5	0	0	5	0	9	12	56	12	0	0	19

Cluster 0 <-- 8
Cluster 1 <-- 2
Cluster 2 <-- 3
Cluster 3 <-- 4
Cluster 4 <-- 6
Cluster 5 <-- No class
Cluster 6 <-- 0
Cluster 7 <-- 7
Cluster 8 <-- 1
Cluster 9 <-- 9
Cluster 10 <-- No class
Cluster 11 <-- No class
Cluster 12 <-- 5
Cluster 13 <-- No class
Cluster 14 <-- No class
Cluster 15 <-- No class
Cluster 16 <-- No class
Cluster 17 <-- No class
Cluster 18 <-- No class
Cluster 19 <-- No class
Cluster 20 <-- No class

Incorrectly clustered instances : 2563.0 45.605 %

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Clusterer

ChooseFarthestFirst -N 2 -S 1

Cluster mode

☐ Use training set

☐ Supplied test set

Set...

☐ Percentage split

%66

☒ Classes to clusters evaluation

(Nom) class

☒ Store clusters for visualization

Ignore attributes

Start

Stop

Result list (right-click for options)

15:22:55 - SimpleKMeans

15:23:01 - Canopy

15:28:09 - FarthestFirst

17:31:54 - SimpleKMeans

17:32:01 - Canopy

17:35:36 - FarthestFirst

17:38:54 - SimpleKMeans

17:39:04 - Canopy

17:42:36 - FarthestFirst

Clusterer output

0.0 0.0 5.0 16.0 16.0 15.0 1.0 0.0 0.0 0.0 10.0 9.0 10.0

Cluster 1

0.0 0.0 0.0 3.0 14.0 7.0 0.0 0.0 0.0 0.0 12.0 12.0

Time taken to build model (full training data) : 0.03 seconds

=== Model and evaluation on training set ===

Clustered Instances

05110 (91%)

1510 (9%)

Class attribute: class

Classes to Clusters:

01<-- assigned to cluster

5513|0

5710|1

5570|2

5720|3

147421|4

5580|5

47583|6

5642|7

5540|8

5611|9

Cluster 0 <-- 3

Cluster 1 <-- 4

Incorrectly clustered instances :4627.082.331%

optdigits_Principal Components_v2.arff

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance" -R first

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst
- 17:38:54 - SimpleKMeans
- 17:39:04 - Canopy
- 17:42:36 - FarthestFirst
- 17:43:49 - SimpleKMeans**

Clusterer output

0.416input51+0.337input31+0.317input37+0.276input7-0.246input23...
-0.31input41-0.289input4-0.281input54+0.249input60+0.216input61...

Time taken to build model (full training data) : 0.11 seconds

=== Model and evaluation on training set ===

Clustered Instances

0	3564	(63%)
1	2056	(37%)

Class attribute: class

Classes to Clusters:

0	1	<-- assigned to cluster
6	548	0
521	50	1
531	26	2
115	457	3
557	11	4
295	263	5
497	61	6
559	7	7
412	142	8
71	491	9

Cluster 0 <-- 7
Cluster 1 <-- 0

Incorrectly clustered instances : 4513.0 80.3025 %

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer

Choose **Canopy** -N 1 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t2 -1.0 -t1 -1.25 -S 1

Cluster mode

☐ Use training set

☐ Supplied test set Set...

☐ Percentage split % 66

☒ Classes to clusters evaluation (Nom) class

☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 15:22:55 - SimpleKMeans
- 15:23:01 - Canopy
- 15:28:09 - FarthestFirst
- 17:31:54 - SimpleKMeans
- 17:32:01 - Canopy
- 17:35:36 - FarthestFirst
- 17:38:54 - SimpleKMeans
- 17:39:04 - Canopy
- 17:42:36 - FarthestFirst
- 17:43:49 - SimpleKMeans
- 17:44:09 - Canopy**

Clusterer output

4	25	(0%)
5	7	(0%)
6	9	(0%)
7	5	(0%)
8	2	(0%)
9	5	(0%)

Class attribute: class

Classes to Clusters:

0	1	2	3	4	5	6	7	8	9	<-- assigned to cluster
554	0	0	0	0	0	0	0	0	0	0
563	0	0	0	0	0	0	4	0	4	1
555	0	0	0	0	0	0	0	2	0	2
572	0	0	0	0	0	0	0	0	0	3
527	0	5	4	15	7	9	0	0	1	4
555	0	0	2	0	0	0	1	0	0	5
547	11	0	0	0	0	0	0	0	0	6
539	0	0	22	5	0	0	0	0	0	7
554	0	0	0	0	0	0	0	0	0	8
555	1	0	1	5	0	0	0	0	0	9

Cluster 0 <-- 3
Cluster 1 <-- 6
Cluster 2 <-- No class
Cluster 3 <-- 7
Cluster 4 <-- 4
Cluster 5 <-- No class
Cluster 6 <-- No class
Cluster 7 <-- 5
Cluster 8 <-- 2
Cluster 9 <-- 1

Incorrectly clustered instances : 4993.0 88.8434 %

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Clusterer

ChooseFarthestFirst -N 2 -S 1

Cluster mode

☐ Use training set

☐ Supplied test set

Set...

☐ Percentage split

%66

☒ Classes to clusters evaluation

(Nom) class

☒ Store clusters for visualization

Ignore attributes

StartStop

Result list (right-click for options)

15:22:55 - SimpleKMeans

15:23:01 - Canopy

15:28:09 - FarthestFirst

17:31:54 - SimpleKMeans

17:32:01 - Canopy

17:35:36 - FarthestFirst

17:38:54 - SimpleKMeans

17:39:04 - Canopy

17:42:36 - FarthestFirst

17:43:49 - SimpleKMeans

17:44:09 - Canopy

17:44:39 - FarthestFirst

Clusterer output

-0.727197 2.404908 1.358906 2.475858 0.469974 0.179989 -1.4008

Cluster 1

5.702828 -3.994317 -1.787413 -1.226637 13.604697 12.928888 6.7

Time taken to build model (full training data) : 0.02 seconds

=== Model and evaluation on training set ===

Clustered Instances

05617 (100%)

13 (0%)

Class attribute: class

Classes to Clusters:

01 <-- assigned to cluster

5540 | 0

5710 | 1

5570 | 2

5720 | 3

5653 | 4

5580 | 5

5580 | 6

5660 | 7

5540 | 8

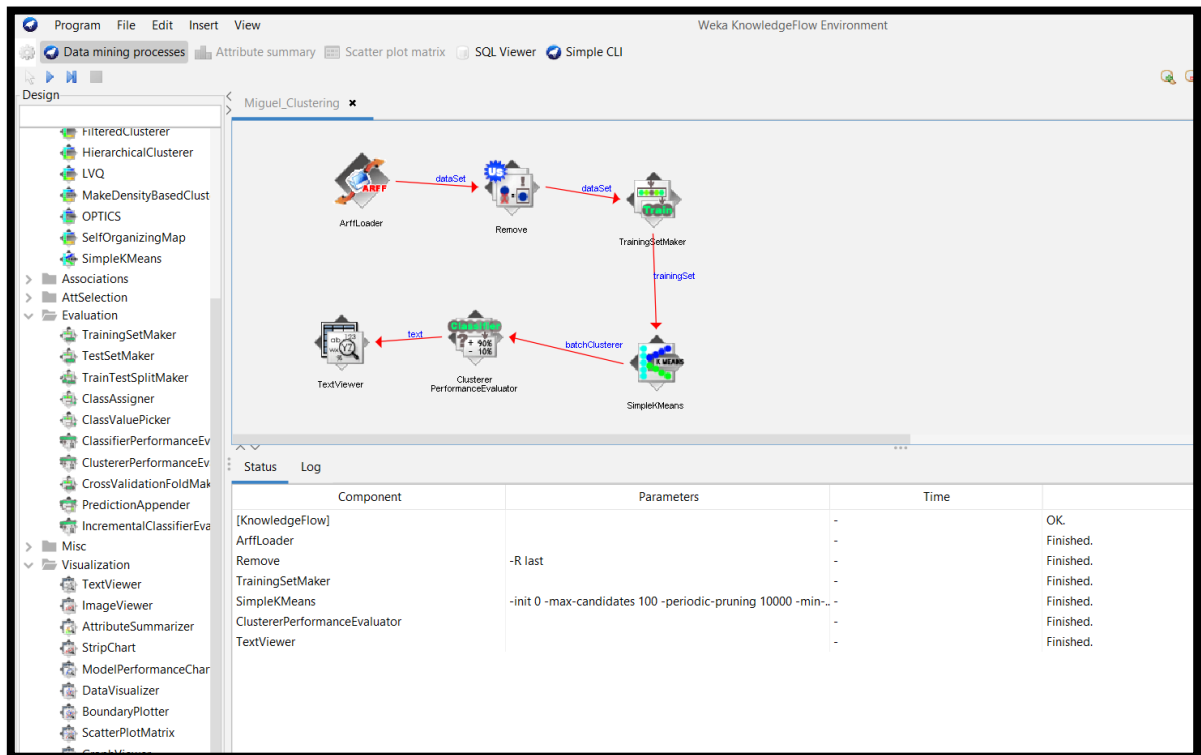
5620 | 9

Cluster 0 <-- 3

Cluster 1 <-- 4

Incorrectly clustered instances :5045.089.7687 %

knowledgeFlow.



Conclusiones

La realización de este proyecto permitió consolidar el conocimiento práctico sobre el flujo de descubrimiento de conocimiento en bases de datos, destacando la importancia de la validación estadística en la toma de decisiones.

A partir de la experimentación con la base de datos optdigits, se derivan las siguientes conclusiones técnicas:

- Efectividad del Modelo IBk:** El análisis comparativo confirmó que el algoritmo basado en instancias (IBk) es significativamente superior (con una precisión de ~98.7%) a los enfoques basados en árboles (J48) o probabilidades (NaiveBayes) para problemas de reconocimiento de patrones visuales digitalizados.
- Selección vs. Transformación:** Se evidenció que las técnicas de *selección* de atributos (Correlation y GainRatio) fueron más exitosas que la *transformación* (PCA). Las primeras lograron mantener la alta precisión del modelo eliminando datos redundantes, mientras que PCA ocasionó una pérdida de información espacial relevante para la clasificación de los dígitos.

3. **Importancia de la Validación Estadística:** El uso del módulo Experimenter fue crucial para determinar que las diferencias de rendimiento observadas no eran aleatorias. Esto valida que la elección del algoritmo correcto tiene un impacto matemático real en el resultado final.
4. **Agrupamiento Autónomo:** Finalmente, mediante el flujo de datos en KnowledgeFlow, se verificó la capacidad de los algoritmos de clustering (como K-Means) para identificar correctamente la estructura de los 10 dígitos manuscritos sin necesidad de supervisión humana, validando la consistencia interna de la base de datos.

En conclusión, se determinó que la combinación óptima para este problema consiste en utilizar una técnica de selección de atributos (como CorrelationAttributeEval) junto con el clasificador IBk, logrando un equilibrio ideal entre precisión y eficiencia computacional.

Bibliografía

1. **Witten, I. H., Frank, E., & Hall, M. A. (2011).** *Data Mining: Practical Machine Learning Tools and Techniques* (3ra ed.). Morgan Kaufmann. (Este es el libro "biblia" de Weka, cubre el funcionamiento del software, el Explorer, el Experimenter y algoritmos como NaiveBayes y J48).
2. **Dua, D. & Graff, C. (2019).** *UCI Machine Learning Repository*. Irvine, CA: University of California, School of Information and Computer Science. Recuperado de [<http://archive.ics.uci.edu/ml>]. (Referencia obligatoria por haber usado la base de datos optdigits que proviene de este repositorio).
3. **Han, J., Kamber, M., & Pei, J. (2011).** *Data Mining: Concepts and Techniques* (3ra ed.). Elsevier. (Referencia teórica fundamental para definir los conceptos de Minería de Datos, Clasificación, Clustering y reducción de dimensionalidad incluidos en la introducción).
4. **Quinlan, J. R. (1993).** *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers. (Referencia específica para el algoritmo J48, que es la implementación en Java del algoritmo C4.5 descrito en este libro).
5. **The University of Waikato. (2024).** *Weka 3: Data Mining Software in Java (Documentation)*. Recuperado de [<https://www.cs.waikato.ac.nz/ml/weka/documentation.html>] (Fuente oficial para el uso de las interfaces gráficas KnowledgeFlow y los filtros de preprocesamiento).