

# Implementierung kNN



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

---

## Projekt der Vorlesung Data Mining und Maschinelles Lernen, WS 15/16

---

### 1 Über dieses Projekt

---

Im Rahmen dieses Projektes werden Sie einen k-Nearest Neighbor Klassifizierer implementieren, wie er auch in der Vorlesung und Übung vorgestellt wird. Hierzu wird ein Framework bereitgestellt, welches Kompatibilität mit WEKA gewährleistet. Bei korrekter Implementation sind die Resultate des Klassifizierers identisch mit der WEKA-Version von IBK. Neben dem Framework werden auch JUnit TestCases zur Verfügung gestellt mit welchen die korrekte Funktionsweise überprüft werden kann.

Die Aufgaben 5.1 und 5.2 sind fast ausschließlich mit Informationen aus dem Vorlesungsskript zu lösen, wohingegen Aufgabe 5.3 etwas Transferleistung und eigene Überlegungen erfordert.

### 2 Installation

---

Das Projekt ist in Java geschrieben und es wird empfohlen, Eclipse<sup>1</sup> als IDE zu verwenden. Das Projektarchiv enthält entsprechende Eclipse Projektsignaturen. Zudem wird JUnit 4 benötigt, welches aber bereits Bestandteil der Eclipse Java Development Tools<sup>2</sup> ist. Die Anbindung an WEKA<sup>3</sup> kann entweder über das Hinzufügen der Datei weka.jar zum classpath erfolgen oder indem man den Quellcode als Projekt einbindet. Das Projekt wurde für WEKA 3.7.11 geschrieben, ist aber auch mit WEKA 3.7.10 folgend kompatibel (getestet bis 3.7.12). Die entsprechende Datei finden Sie unter Developer Version auf der WEKA Download Seite.

### 3 Abgabe

---

Die Abgabe erfolgt zusammen mit der letzten Abgabe des theoretischen Projektes. Die Fragen des Implementationsprojekts sollen in einem getrennten Dokument beantwortet werden, nicht direkt im Code. Die Abgabe des Quellcodes besteht nur aus der NearestNeighbor.java. Weitere Dateien oder Links finden keine Berücksichtigung. Die Verwendung von weiteren Klassen oder Libraries ist also nicht möglich. Im Code sollen zudem keine Annahmen über die Reihenfolge des Ausführenden der Methoden gemacht werden. Z.b. wird zu Testzwecken `determineManhattanDistance` ausgewertet, ohne dass vorher `learnModel` aufgerufen wird.

### 4 Das Framework

---

Zu implementieren ist die Klasse `tud.ke.ml.project.classifier.NearestNeighbor`. Diese erbt von `tud.ke.ml.project.framework.classifier.ANearestNeighbor`, welche ein lightweight Framework definiert so wie das Handling der Konfiguration. Die Klasse `weka.classifiers.lazy.keNN` definiert die Schnittstelle zu WEKA und darf genauso wie die abstrakte Basisklasse nicht modifiziert werden.

Instanzen sind im Framework als Liste von Objekten gespeichert. Ein jedes Objekt ist entweder ein `String`, und repräsentiert in diesem Falle ein nominales Attribut, oder ein `Double`, welches für ein numerisches Attribut steht. Das Klassenattribut ist immer nominal. Der Index des Klassenattributes entspricht dem Rückgabewert von `getClassAttribute` und muss nicht immer dem letzten Attribut entsprechen.

Gestartet wird das Projekt über `tud.ke.ml.project.main.RunWEKA`. Der neue Classifier ist nun in allen Teilen des Programmes unter `lazy.keNN` zu finden. Die Klasse `main.SimpleRun` ist ein einfaches Beispiel, welches während der Entwicklung genutzt werden kann, falls Sie nicht immer die WEKA GUI verwenden wollen. Im Package `tud.ke.ml.project.junit` befinden sich zudem die erwähnten JUnit-Tests. `junit.SimpleValidation` und `junit.AdvancedValidation` testen verschiedene Konfigurationen des Klassifizierers. Sollten diese Tests alle erfolgreich durchlaufen, ist dieser sehr wahrscheinlich korrekt

---

<sup>1</sup> <https://www.eclipse.org/>

<sup>2</sup> <http://www.eclipse.org/jdt/>

<sup>3</sup> <http://www.cs.waikato.ac.nz/ml/weka/>

---

implementiert.

Zu allen relevanten Methoden ist eine Javadoc-Dokumentation vorhanden.

---

#### 4.1 JUnit

---

JUnit-Tests können ausgeführt werden, indem die entsprechende Library dem Projekt hinzugefügt wird. In Eclipse ist diese bereits vorhanden und muss daher nicht mehr zusätzlich hinzugefügt werden. Zum Zesten muss eine JUnit-Java-Datei ausgewählt werden und per Kontextmenü *Run As*→*JUnit Test* ausgeführt werden. Rote Kreuze bedeuten, dass der Test eine Exception ausgelöst hat. Blaue Kreuze repräsentieren einen fehlgeschlagenen JUnit-Test. Für eine genauere Beschreibung von JUnit siehe <http://www.vogella.com/tutorials/JUnit/article.html>.

---

### 5 Die Aufgabe (12 Punkte)

---

Im Folgenden ist zu beachten, dass zum Testen alle Methoden implementiert sein müssen. Es ist aber natürlich möglich, hierfür erst einmal nur Dummy-Methoden zu verwenden, die einfach einen fixen Wert zurückgeben. Definieren sie zuerst `getGroupNumber`, da diese Funktion für die finale Bewertung verwendet wird! (**Achtung:** In einer vorherigen Version wurde noch gefordert, die Matrikelnummern zurückzuliefern.)

---

#### 5.1 Der Basis-Klassifizierer

---

- Implementieren Sie die Methode `learnModel`, welche die als Argument übergebenen Trainingsdaten intern speichert. Des Weiteren soll die Methode `vote` implementiert werden, indem Sie `getUnweightedVotes` aufrufen. Die daraus resultierende Map muss an `getWinner` übergeben werden um die aus den Votes resultierende Klasse zu bestimmen.
- `getNearest` muss für alle Instanzen des gespeicherten Modells die Manhattan-Distanz zur übergebenen Test-Instanz berechnen. Dies soll mit Hilfe der zu implementierenden Methode `determineManhattanDistance` erfolgen. Die Liste, die zurückgegeben wird, muss auf die nächsten `getkNearest` Elemente beschränkt sein. Die Liste beinhaltet `Pair` Objekte, bestehend aus der Instanz sowie dem Abstand.
- Wenn Sie nun noch die Methode `getWinner` implementieren, sollten Sie eine erste, lauffähige Implementation haben. Die genannte Methode soll den Name/Wert der Klasse mit den meisten Votes zurückgeben.

---

#### 5.2 Inverse Distance Weighting und Euclidean Distance

---

- Implementieren Sie nun `getWeightedVotes`, welches die Stimmen als Summe der inversen Distanzen berechnet. `vote` muss diese Methode in Abhängigkeit von `isInverseWeighting` aufrufen.
- Genauso muss die Methode `getNearest` in Abhängigkeit von `getMetric` die euklidische Distanz als Entfernungsmaß nutzen. Hierzu soll die Methode `determineEuclideanDistance` genutzt und implementiert werden.

---

#### 5.3 Normalization and Tie-Breaking

---

- Es ist meistens sinnvoll, die Attributwerte zu normalisieren. Dafür soll die Methode `normalizationScaling` die nötigen Skalierungs- und Translationsfaktoren zurückliefern, abhängig von `isNormalizing`. In der Methode `getNearest` muss nun das resultierende Array aufgeteilt und als `protected double[] scaling` für die Skalierung, bzw. `protected double[] translation` für die Verschiebung gespeichert werden, wobei für jedes Attribut (auch das potentielle Klassenattribut) ein Eintrag vorhanden sein muss. Diese Arrays **müssen** nun genutzt werden, um die Attributwerte entsprechend zu normalisieren, d.h. das Resultat von `normalizationScaling` darf nicht direkt übernommen werden!
- In der Methode `getWinner` kann es vorkommen, dass mehr als eine Klasse die meisten Stimmen bekommen hat. Implementieren Sie hierfür eine Entscheidungsfunktion und begründen Sie diese. (1 Punkt)
- In der Methode `getNearest` kann es passieren, dass mehrere Instanzen den gleichen Abstand haben. Wie sollte man hierbei vorgehen? Begründen und implementieren Sie ihre Methode. (1 Punkt)