

Theoriefragen zur Implementationsaufgabe

Küpper Joel, Ockenfels Malou, Schulz Daniel

January 23, 2017

1 5.3) Normalization and Tie-Breaking

1.1 a) Entscheidungsfunktion für getWinner

Problem: *Es kann vorkommen, dass mehr als eine Klasse am meisten Stimmen von den Nachbarn bekommen hat. Welche Klasse soll dann ausgewählt werden?*

Die Idee ist, dass falls mehr als eine Klasse am meisten Stimmen bekommen hat, jene Klasse ausgewählt wird, die generell im Data-set am Meisten vorkommt. Somit ist die Wahrscheinlichkeit, richtig zu liegen, am höchsten. Den Fall, dass die beiden *gewinnenden* Klassen auch gleichzeitig gleich häufig im Data-set vorkommen wird hier in der Implementierung nicht beachtet, da die Wahrscheinlichkeit doch sehr gering ist, vorallem bei den vorliegenden Data-sets.

Für diesen Fall könnte jedoch so vorgegangen werden, dass jene Klasse vorhergesagt wird, die am nächsten ($k=1, =2, \dots$) am Beispiel liegt. Dann gibt es aber Probleme bei zu hohem Rauschen in den Data-sets.

Eine weitere Möglichkeit wäre, dass man $k+=1$ inkrementiert und so einen weiteren Nachbarn (oder zwei, drei, ...) hinzuzieht, um eine eindeutigere Mehrheit zu finden. Dabei tritt aber das Problem der zu generellen Vorhersage auf. Da in der Aufgabe keine Angabe über jene Priorisierung gemacht wurde, werden die Denkansätze hier nur erwähnt und nicht implementiert.

Der folgende Pseudocode illustriert, wie das Zählen umgesetzt wurde.

Algorithm 1 getWinner Entscheidungsfunktion

```
1: function GETWINNER(data, votes)
2:    $max \leftarrow votes.getHighestVote()$ 
3:    $classesWithMaxVotes \leftarrow null$ 
4:   for each  $vote \in votes$  do
5:     if  $vote.numberOfVotes() == max$  then
6:        $classesWithMaxVotes.add(vote);$ 
7:   if  $classesWithMaxVotes.size == 1$  then
8:      $return \leftarrow classesWithMaxVotes.first();$ 
9:   else
10:     $counter.init(0);$ 
11:    for each  $instance \in data$  do
12:      if  $classesWithMaxVotes.contains(instance.getClass())$  then
13:         $counter.incrementClass(instance.getClass());$ 
14:     $return \leftarrow counter.getBiggest().getClass();$ 
```

1.2 c) Gleicher Abstand mehrerer Instanzen

Problem: *Es kann vorkommen, dass mehrere Instanzen den gleichen Abstand vom Beispiel haben. Wie soll dabei vorgegangen werden?*

Das ist prinzipiell nur ein Problem, wenn der k 'te Nachbar den gleichen Abstand hat wie $k+1$, $k+2$, ... Wenn die Nachbarn eins und zwei bspw. bei $k=4$ den gleichen Abstand haben, gibt es kein Problem.

Erst wenn Nachbar k und Nachbar $k+1$ den gleichen Abstand und unterschiedliche Klassen haben, gibt es Entscheidungsbedarf. Nun, hier ist die Idee, wie oben schon beschrieben, aber nun unausweichlich, k zu inkrementieren und alle Nachbarn, die den gleichen Abstand haben wie der k 'te Nachbar auch noch als nearest Neighbor zu bezeichnen und entsprechend zu berücksichtigen.

Auch hier der Pseudocode, wie es umgesetzt wurde:

Algorithm 2 getNearest Entscheidungsfunktion

```
1: function GETNEAREST(sortedListOfNeighbors, k)
2:   sortedNeighborsByDistance  $\leftarrow$  null;
3:   for each neighbor  $\in$  sortedListOfNeighbors do
4:     sortedNeighborsByDistance.addAt(neighbor.getDistance(), neighbor);
5:   nearests  $\leftarrow$  null;
6:   for each setOfNeighborsWithSameDistance  $\in$  sortedNeighborsByDistance do
7:     nearests.addAll(setOfNeighborsWithSameDistance);
8:     if nearest.size()  $\geq k$  then
9:       break;
10:  return  $\leftarrow$  nearests;
```
