



Blatt 3

Eines der aktuellen Forschungsgebiete in der Biologie ist die Analyse von Genexpressionen. Ausgehend von dem Wissen, dass verschiedene Krebsarten mit der Bildung verschiedener Proteine zusammenhängen - einerseits solche, die für den Krebs verantwortlich sind, andererseits solche, die ihn bekämpfen, ist man daran interessiert herauszufinden, welche Proteine verstärkt/vermindert (die Genexpression ist also erhöht oder verringert) gebildet werden.

Dazu wird präparierte und fluoreszierend markierte mRNA aus Gewebeproben auf ein sogenanntes Microarray aufgebracht, das komplementäre Basenfolgen der zu untersuchenden Proteine enthält. Die Proteine der Gewebeprobe binden sich an die zugehörigen Basen - je nach Anzahl der gebundenen Proteine erscheint dieser Ausschnitt auf dem Microarray durch die Fluoreszenz heller oder dunkler. Der Chip wird eingescannt und der Helligkeitswert jedes Proteins wird in einen Zahlenwert umgewandelt.

In diesem Falle wurden mehrere Knochenmark-Proben mit je etwa 35403 Genen (jede Probe ein Chip) von Patienten mit Leukämie (ALL) oder Ewing (Knochenmark-Krebs), sowie gesundem Knochenmark (KM) eingescannt.

Nun sollen auf die Daten dieser drei Arten von Gewebeproben verschiedene Klassifikationsverfahren angewendet, also möglichst ähnliche Proben in ein Cluster gebracht werden. Dabei ist zu beachten, dass man gerade bei der Verwendung von zufälligen Anfangspunkten völlig verschiedene Ergebnisse erhalten kann. Für Tests ist also auch interessant, zwischen welchen Werten der Quantisierungsfehler schwankt. Zur Vereinfachung wird nur eine Auswahl von Genen betrachtet.

Die Daten sind, neben den Übungsblättern, im Web zu finden - in GenDaten.txt die bereits klassifizierten Proben, und in ProbenXYZ.txt 2 noch unbekanntes Proben.

Aufgabe 3.1 Es soll eine Auswahl der gegebenen Genexpressionen als Merkmale genutzt werden. Welche Kriterien sind bei der Auswahl der Gene sinnvollerweise zu beachten?

Ermitteln Sie für die Unterscheidung zwischen den Probensorten nützliche Gene. Bedenken Sie dabei, dass die tatsächliche Anzahl der Gene mit 35403 erheblich höher ist, als die Zusammenstellung von Genen für diese Aufgaben - eine Auslese per Hand ist also (i.A.) nicht möglich. Schreiben Sie ein Programm, das solche Merkmale findet und aussortiert.

Welches Ergebnis liefert Ihr Programm?

Aufgabe 3.2 Programmieren Sie den k-means-Algorithmus mit mindestens 2 verschiedenen Distanzfunktionen. Testen Sie Ihr Programm sowohl mit allen vorgegebenen Genen als auch nur mit den aussortierten. Welche Cluster werden gebildet und welche Proben werden korrekt bzw. falsch zugeordnet? Zu welchen Clustern werden die zwei unbekanntes Proben am ehesten zugeordnet? Geben Sie jeweils die Quantisierungsfehler an.

Aufgabe 3.3 Vergleichen Sie die Klassifikationsergebnisse (alle/ausgewählte Gene, NN-Klassifikator, mit/ohne k-means, falls mögl. Hyperquaderkl.).