

Sheet 11

Topic: Clustering, Gaussian Process Regression

Submission deadline: Tue 22.7.2008, 11:00 a.m. (before class)

Exercise 1: K-means Clustering

Use the k-means algorithm and the Euclidean distance to cluster the following 8 examples into 3 clusters: $p_1 = (2, 10)$, $p_2 = (2, 5)$, $p_3 = (8, 4)$, $p_4 = (5, 8)$, $p_5 = (7, 5)$, $p_6 = (6, 4)$, $p_7 = (1, 2)$, $p_8 = (4, 9)$

The distance matrix based on the Euclidean distance is given below:

	p_1	p_2	p_3	p_4	p_5	p_6	p_7	p_8
p_1	0	$\sqrt{25}$	$\sqrt{36}$	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{65}$	$\sqrt{5}$
p_2		0	$\sqrt{37}$	$\sqrt{18}$	$\sqrt{25}$	$\sqrt{17}$	$\sqrt{10}$	$\sqrt{20}$
p_3			0	$\sqrt{25}$	$\sqrt{2}$	$\sqrt{2}$	$\sqrt{53}$	$\sqrt{41}$
p_4				0	$\sqrt{13}$	$\sqrt{17}$	$\sqrt{52}$	$\sqrt{2}$
p_5					0	$\sqrt{2}$	$\sqrt{45}$	$\sqrt{25}$
p_6						0	$\sqrt{29}$	$\sqrt{29}$
p_7							0	$\sqrt{58}$
p_8								0

Suppose, the 3 initial seeds (centers of each cluster) are p_1, p_4 and p_7 . Run the k-means algorithm for one iteration only. At the end of this run show:

1. The new clusters (i.e. the examples belonging to each cluster)
2. The centers of the new clusters.

How many more iterations would be necessary for convergence? Draw the results for each iteration step.

Exercise 2: Gaussian Process Regression

- (a) Explain briefly, in your own words, the main idea of Gaussian process regression.

- (b) Standard GP regression makes the assumption that the noise is constant. Think up a real world example, where this assumption is violated.
- (c) Suppose you change the function value y_i of one of the targets (but not its input location x_i). How does this affect the variance of the predictive distribution at location x_i , and at any other location? Explain why this is counterintuitive.

Exercise 3: Programming Task (Gaussian Process Regression)

Download the source code stub from our website (you will need either Matlab or GNU Octave to run the scripts). The task is to use Gaussian process regression for estimating the values of a one-dimensional function f from which we have noisy observations $y_i = f(x_i) + \epsilon$ (stored in the vector `targets`) at input locations x_i (stored in the vector `inputs`).

- (a) Complete the function `covFunc` by implementing the one-dimensional squared exponential.
- (b) Complete the function `gpReg` by computing the covariance matrices \mathbf{K} and \mathbf{k}_* .
- (c) Now the `testGP` function should work and you will see a plot of the function f (black), the noisy observations (black), the predicted mean (blue) and the 3σ confidence interval of the prediction (red). Try different values for `sigmaNoise`, `sigma`, and `ell`, and describe how this affects the regression.
- (d) Implement a stochastic optimization method of your choice in the `trainGP` function and uncomment the two lines in the `testGP` file that call the optimization function. What are your optimized hyperparameters?