



Universidad Politécnica de Madrid

ESCUELA TÉCNICA SUPERIOR DE INGENIEROS INDUSTRIALES

MÁSTER EN AUTOMÁTICA Y ROBÓTICA

APPLIED ARTIFICIAL INTELLIGENCE

Assignment 3.1: k-nn and Bayesian Classifier

Josep María BARBERÁ CIVERA (17048)

March 12, 2024

k-nn and Bayesian Classifier

Here are listed the main instructions for the task completion:

1. First load the synthetically created “data_D2_C2” data where the `p.value` dataset will be your training data and the `t.value` dataset will serve as the validation data.
2. Use the `nn-classifier` for computing the number of *miss-classifications* errors on `p.value` and on `t.value`, but remember to use in both cases `p.value` as the learning data.
3. Repeat again but now using the `Bayesian classifier`.
4. Compute the miss-classification error on `t.value` when using a Bayesian classifier with the following *a priori* probabilities: [1 0.1], [1 2] & [1 10].
5. Discuss the results.

1.1 Methodology

Thanks to the `fitcknn` and `fitcnb` functions, it is very easy to train the classifiers, parametric (`knn`) and non-parametric (`bayesian`) respectively.

On the other hand, once trained, the `predict` function, given a model and an input value, returns the estimate of the most probable class. Iterating over all the data of the training and validation sets, the behaviour of these algorithms has been evaluated and the results represented as confusion matrices (see next section (1.2)).

Finally, to better visualise the behaviour of the classifiers, a grid of points has been constructed that coincides with the domain of the data and the classification has been calculated for each node of this grid and the correct value of each point has been represented on it. It can then be seen where the classification fails and how the domains change depending on the initial value given to the `predict` function.

1.2 Discussion and Results

The figures obtained in the validation of the models are included below.

1.2.1 knn classifier

The k-nearest neighbour classifier is supervised and non-parametric, resulting in low generalisation with consequent noise learning. Figure 1.a shows the high specialisation of the model which for the data provided to learn the exact zones. While in Figure 1.b for the validation data, the classification fails in some points.

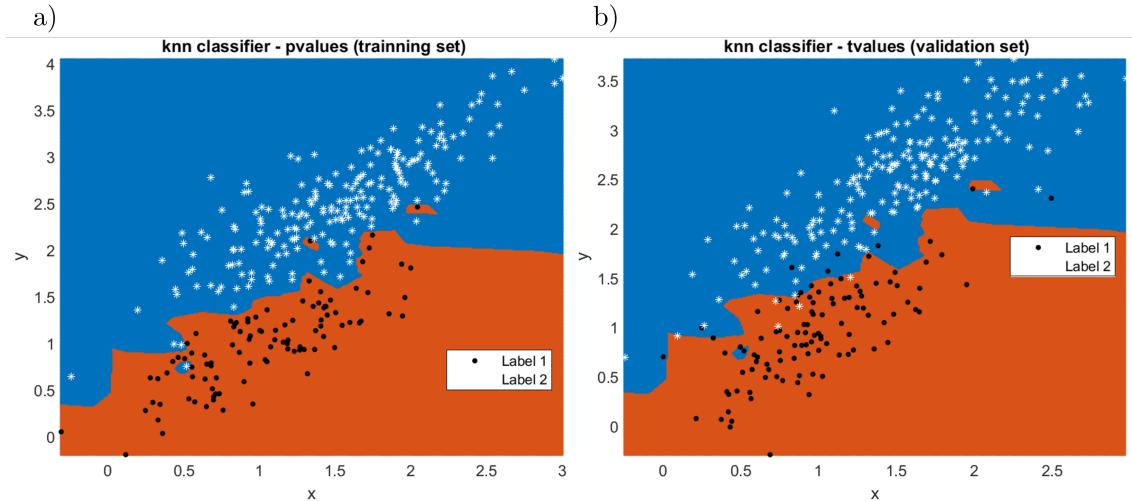


Figure 1: Map of predictions using the **knn** model together with the correctly labelled points. In a) the training data while in b) the validation data.

This means that the confusion matrix associated with the training data does not present any type of error, while for the validation data it does. It is seen in Figure 2.a for the `pvalues` data points and in Figure 2.b for the `tvalues` where there are about 15 erroneous points.

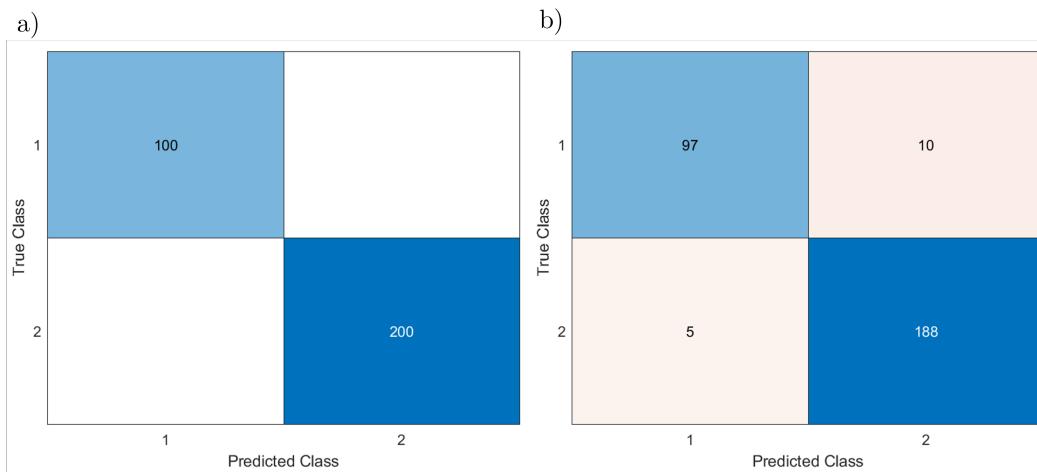


Figure 2: Confusion matrix for the **knn** model. In a) for the same training values while in b) for the validation values.

The accuracy associated with the confusion matrices are the following:

$$\text{ACC}_1 = 100\% ; \quad \text{ACC}_2 = 95\%$$

1.2.2 Bayesian classifier

The **Bayesian classifier**, on the other hand, has an inherently generalising capacity as opposed to the data storage that the **knn classifier** used to perform. On the other hand, this model assumes that the data follow a normal distribution, which in our case is not true, leading to poor behaviour.

First, the results obtained for the Bayesian model without any prior estimation are shown. Thus, Figure 3 shows the frontier obtained against the points with their correct classification. It contrasts the smoothness of the model that presents little noise, but fails catastrophically in the classification. This occurs both for the training values themselves (FIG. 3.a) and for the validation data (FIG. 3.b).

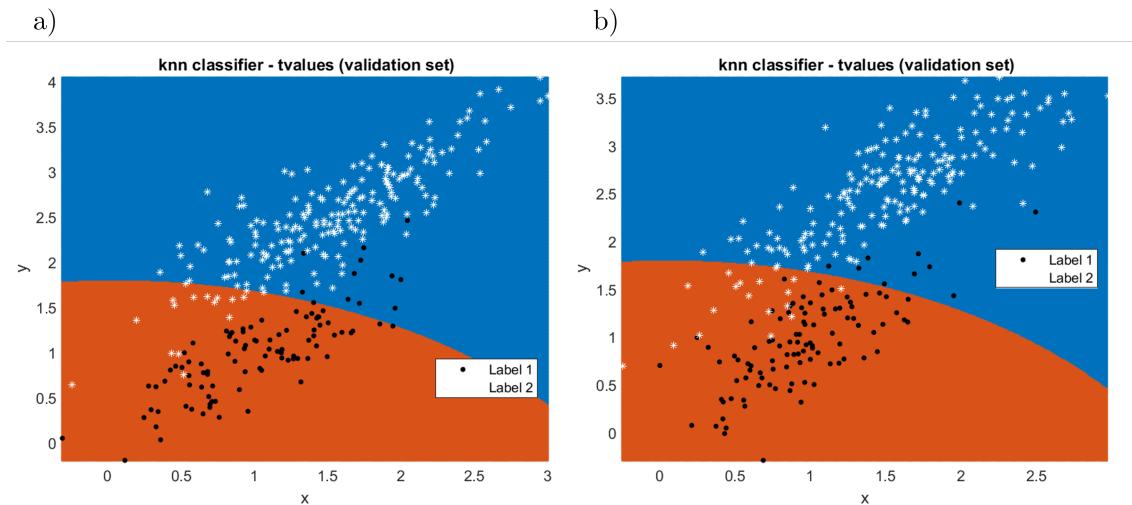


Figure 3: Map of predictions using the **bayesian** model together with the correctly labelled points. In a) the training data while in b) the validation data.

This is corroborated by the confusion matrices in Figure 4, since the model does not change its boundary either, both confusion matrices present similar values of accuracy, being in this case equal, although this need not always be the case.

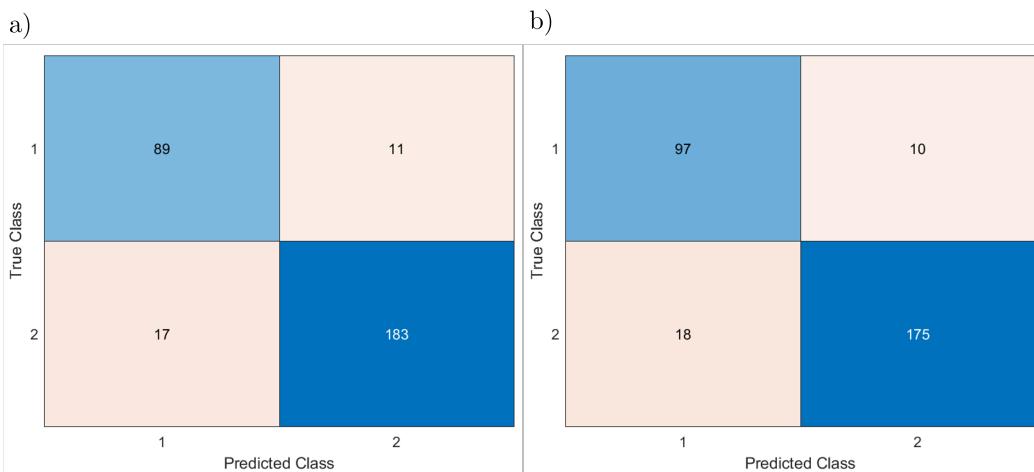


Figure 4: Confusion matrix for the **bayesian** model. In a) for the same training values while in b) for the validation values.

The accuracy associated with the confusion matrices are the following:

$$ACC_1 = 90.67\%; \quad ACC_2 = 90.67\%$$

The same process has been repeated for three different values of **initial probabilities** or **initial values** in the training of the model. These values are the next: [1 0.1], [1 2] & [1 10].

Figure 5 shows the evolution of the frontier for the three possible cases and the training data (right) or the validation data (left).

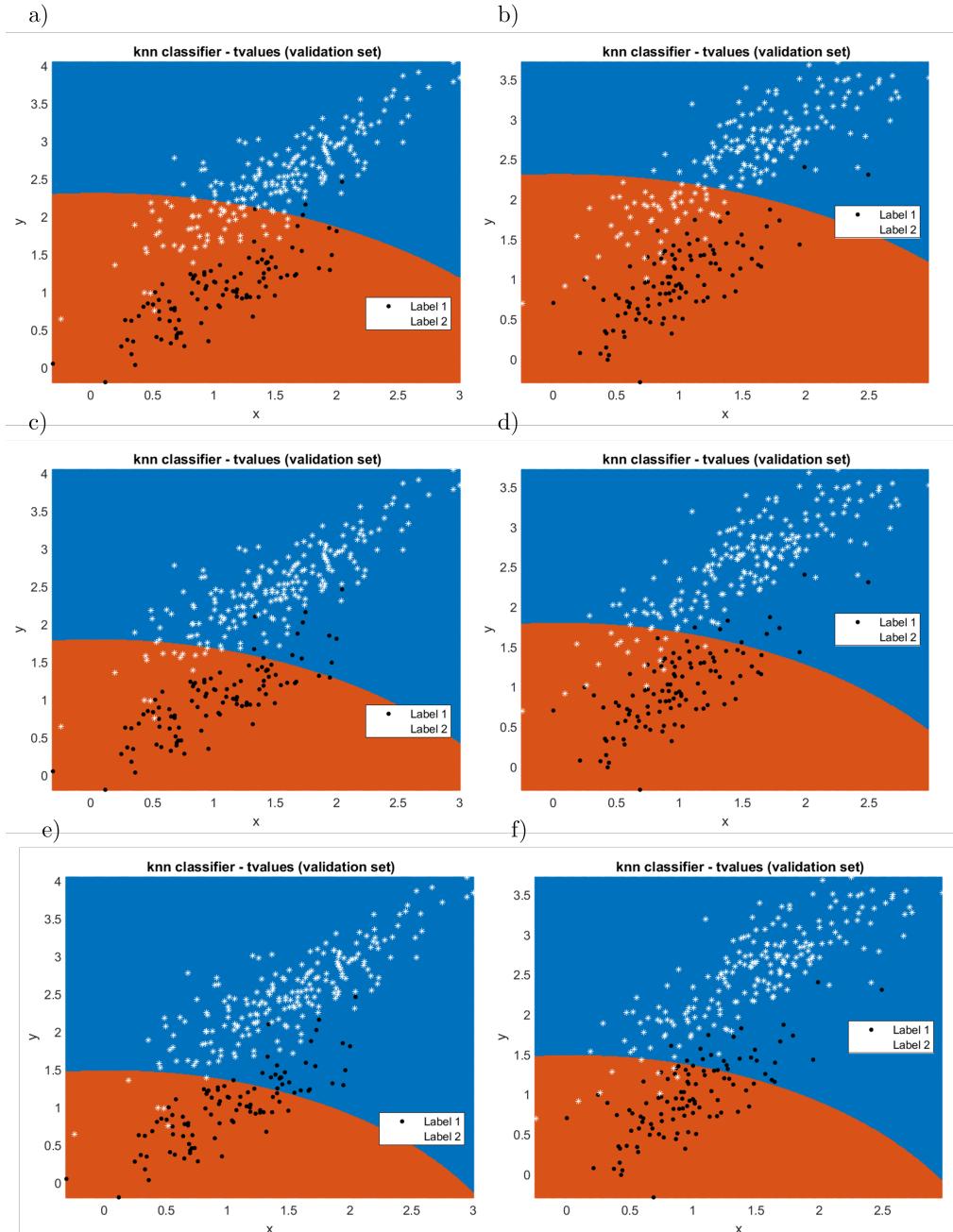


Figure 5: Map of predictions using the **bayesian** model along with some prior probabilities values, together with the correctly labelled points. In a), c) and e) the training data while in b), d) and f) the validation data. In a) and b) the initial guess is [1 0.1], in c) and d) is [1 2], while in e) and f) is [1 10].

The associated confusion matrices are included in Figure 6.

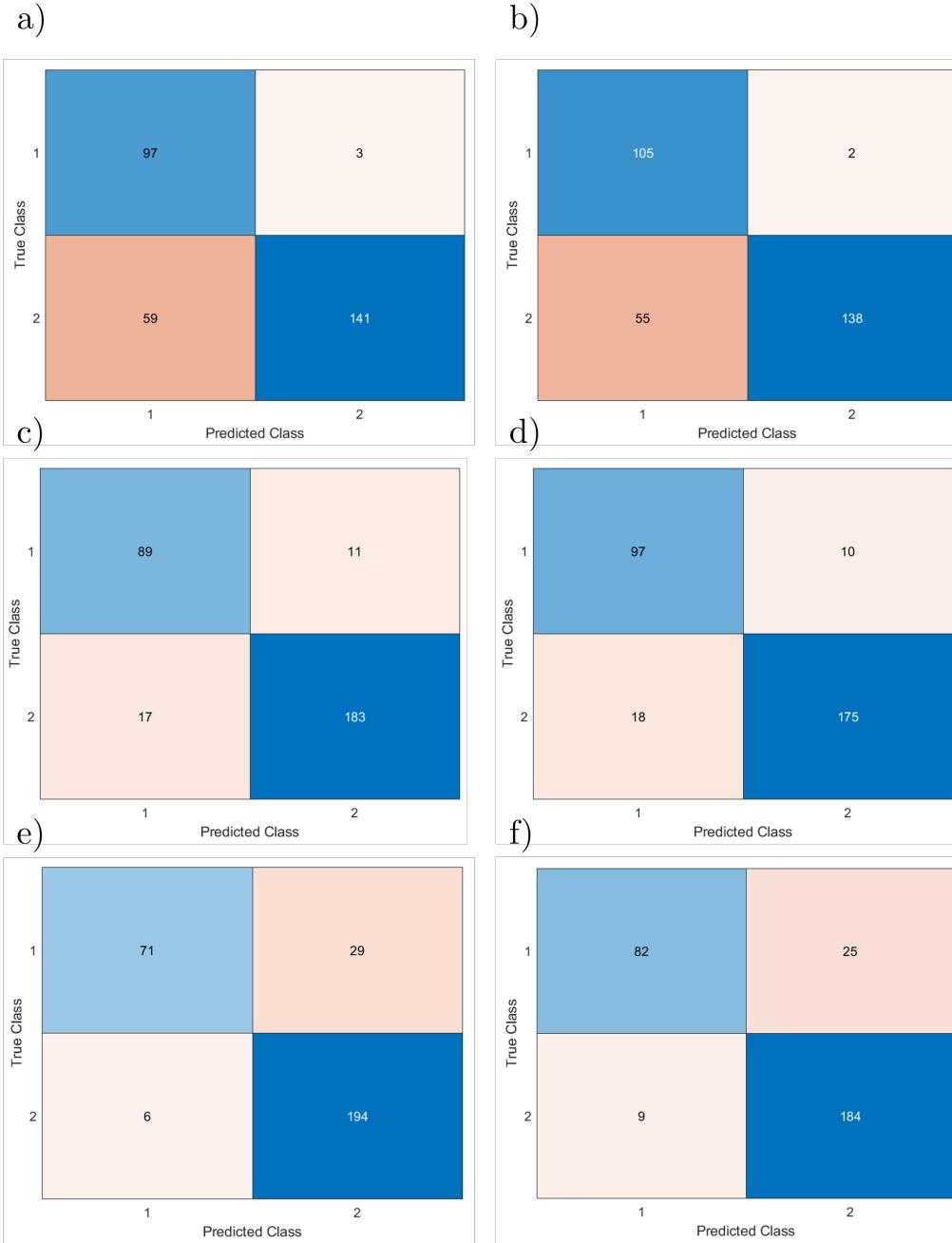


Figure 6: Confusion matrix for the **bayesian** model with some prior probabilities values. In a), c) and e) the training data while in b), d) and f) the validation data. In a) and b) the initial guess is [1 0.1], in c) and d) is [1 2], while in e) and f) is [1 10].

The respective precision values are as follows:

$$\begin{aligned}
 ACC_a &= 79.33\%; & ACC_b &= 81\% \\
 ACC_c &= 90.67\%; & ACC_d &= 90.67\% \\
 ACC_e &= 88.33\%; & ACC_f &= 88.67\%
 \end{aligned}$$

1.3 Relevant Code

The code prepared for this assignment is shown in the next pages in a wider style.

```

1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %                               Master in Robotics
3 %                               Applied Artificial Intelligence
4 %
5 % Assingment 3.1: k-nn and Bayesian Classifier
6 % Student: Josep Barbera Civera
7 % ID: 17048
8 % Date: 11/03/2024
9 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
10
11 % 0. Load data_D2_C2: p.value and t.value
12 % 1. Use the nn-classifier for computing the number of
13 % missclassifications errors on p.value and on t.value,
14 % using in both cases p.value as the learning data.
15 % 2. Same as before using the Bayesian classifier
16 % 3. Compute the miss-classification error on t.value when
17 % using a Bayessian classifier with following a priori
18 % probabilities: [1 0.1], [1 2] & [1 10]. Discuss the results.
19
20 load data_D2_C2.mat
21
22 %% Accesing Data
23 pvalues = p.value;
24 plabels = p.class;
25 tvalues = t.value;
26 tlabels = t.class;
27
28 [Mp, Np] = size(pvalues);
29 [Mt, Nt] = size(tvalues);
30
31 %% Models Trainning
32 % knn = fitcknn(pvalues', plabels');
33 % bayes = fitcnb(pvalues', plabels');
34 % prior = [1 0.1];
35 % prior = [1 2];
36 prior = [1 10];
37 bayes= fitcnb(pvalues', plabels', 'Prior', prior);
38
39 %% Model Prediction: training set
40 for i=1:Np
41     % predKnn(:,i) = predict(knn, pvalues(:,i)');
42     predBayes(:,i) = predict(bayes, pvalues(:,i)');
43 end
44 pred = predBayes;
45 %% Confusion matrix plotting
46 C = confusionmat(plabels',pred');
47 confusionchart(C);
48 set(gcf, 'Color', 'white')
49 saveas(gcf, 'confusion_chart_pvalues_bayes_prior_3.png'); % do not
   forget to change the name!
50
51 %% Model Prediction: test set
52 for i=1:Nt

```

```

53 % predKnn(:, i) = predict(knn, tvalues(:, i)');
54 predBayes(:, i) = predict(bayes, tvalues(:, i)');
55 end
56 pred = predBayes;
57 %% Confusion matrix plotting
58 C = confusionmat(tlabels', pred');
59 confusionchart(C);
60 set(gcf, 'Color', 'white');
61 saveas(gcf, 'confusion_chart_tvalues_bayes_prior_3.png'); % do not
   forget to change the name!
62
63 %% pvalues prediction map vs predicted values
64 % model = knn;
65 model = bayes;
66 x1 = min(pvalues(1,:)):0.01:max(pvalues(1,:));
67 x2 = min(pvalues(2,:)):0.01:max(pvalues(2,:));
68 [x1G, x2G] = meshgrid(x1, x2);
69 XGrid = [x1G(:, ), x2G(:, )];
70 allPred = predict(model, XGrid);
71
72 figure
73 gscatter(XGrid(:, 1), XGrid(:, 2), allPred, [0.8500 0.3250 0.0980; 0 0.4470
   0.7410])
74 hold on
75
76 X = pvalues';
77 y = plabels';
78
79 h1 = plot(X(y == 1, 1), X(y == 1, 2), '.', 'MarkerSize', 10, '
   MarkerEdgeColor', 'k', 'MarkerFaceColor', "w");
80 hold on;
81 h2 = plot(X(y == 2, 1), X(y == 2, 2), '*', 'MarkerSize', 5, '
   MarkerEdgeColor', 'w', 'MarkerFaceColor', "w");
82
83 xlabel('x')
84 ylabel('y')
85 title('{\bf knn classifier - tvalues (validation set)}')
86 legend([h1, h2], {'Label 1', 'Label 2'}, 'Location', 'best');
87 legend_box = findobj(gcf, 'Type', 'legend');
88 set(legend_box, 'Color', [0.8 0.8 0.8]); % Gray background color
89 axis tight
90 hold off
91 set(gcf, 'Color', 'white');
92 saveas(gcf, 'prediction_map_pvalues_bayes_prior_3.png'); % do not forget
   to change the name!
93
94
95 %% tvalues prediction map vs predicted values
96 % model = knn;
97 model = bayes;
98 x1 = min(tvalues(1,:)):0.01:max(tvalues(1,:));
99 x2 = min(tvalues(2,:)):0.01:max(tvalues(2,:));
100 [x1G, x2G] = meshgrid(x1, x2);
101 XGrid = [x1G(:, ), x2G(:, )];

```

```

102 allPred = predict(model,XGrid);
103
104 figure
105 gscatter(XGrid(:,1),XGrid(:,2),allPred,[0.8500 0.3250 0.0980;0 0.4470
106     0.7410])
107 hold on
108 X = tvalues';
109 y = tlabels';
110
111 h1 = plot(X(y == 1, 1), X(y == 1, 2), '.', 'MarkerSize', 10, '
112     MarkerEdgeColor', 'k', 'MarkerFaceColor', "w");
113 hold on;
114 h2 = plot(X(y == 2, 1), X(y == 2, 2), '*', 'MarkerSize', 5, '
115     MarkerEdgeColor', 'w', 'MarkerFaceColor', "w");
116 xlabel('x')
117 ylabel('y')
118 title('{\bf knn classifier - tvalues (validation set)}')
119 legend([h1, h2], {'Label 1', 'Label 2'}, 'Location', 'best');
120 legend_box = findobj(gcf, 'Type', 'legend');
121 set(legend_box, 'Color', [0.8 0.8 0.8]); % Gray background color
122 axis tight
123 hold off
124 set(gcf, 'Color', 'white');
125 saveas(gcf, 'prediction_map_tvalues_bayes_prior_3.png'); % do not forget
126     to change the name!

```