

Applied Data Mining Midterm Practice

Xun Zhao, xz2827

Problem 1

1.

No, the Bayes-optimal classifier is trying to minimize the risk $R(f)$, whose value depends on the loss function.

2.

No, it is likely to overfitting for small k value. For example, assuming $k = 1$, the decision boundary can be highly non-linear.

3.

The data points have labels, such as discrete ys (classification) or continuous ys (regression).

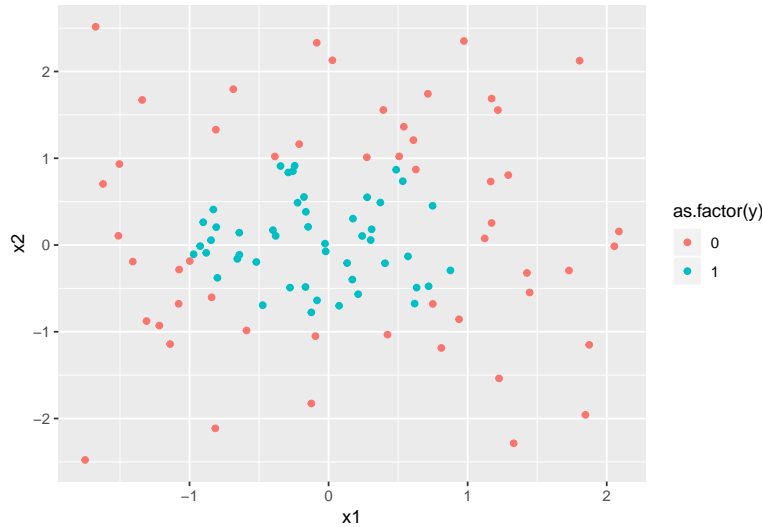
Problem 2

Problem 3

When the decision boundary is not linear, the maximum margin classifier is not optimal.

For example,

```
x1 = rnorm(100, 0, 1)
x2 = rnorm(100, 0, 1)
y = ifelse(x1 ^ 2 + x2 ^ 2 < 1, 1, 0)
d = data.frame(x1, x2, y)
library(ggplot2)
ggplot(d, aes(x = x1, y = x2, color = as.factor(y))) + geom_point()
```



In this case, it is hard to define margin if the classifier use a linear function, although using *kernel function* can solve it.

Problem 4

Fig 1.

Would: Linear classifier with gradient descent learning

Why: It is because that the linear classifier can finally give a straight line as the boundary. With the gradient descent and squared-error, we can learn the slope and intercept to minimize the loss.

Would not: Perceptron classifier

Why: It is because the perceptron cannot separate non-linearly separatable data. No matter what boundary it learnt, there would always be some points that are misclassified, which change the slope and intercept again.

Fig 2.

Would: Naive Bayes

Why: The boundary is very smooth curve, so the boundary may be well defined by a function, which can be derived from the $P(class1) = P(class2)$ equation. Plus, the distribution model of data points have to be complex enough to draw such a complex boundary.

Would not: Linear classifier

Why: the boundary is not linear.

Fig 3.

Would: K-nearest classifier

Why: the boundary is not smooth, but is connected by many line segments, which can be seen as the edges of K-nearest polygons of those data points.

Would not: Linear classifier

Why: the boundary is not linear.

Problem 5

1.

The sample space is \mathbb{R}^{k+m} , whose first k elements are from \vec{a} and last m elements are from \vec{b} . The labels are in two classes: 1 (match) and 0 (not match).

2.

The loss function's penalty is too large, making the classifier is highly overfitting the "match" data. In this case, only new data that fits the training data can get "match".

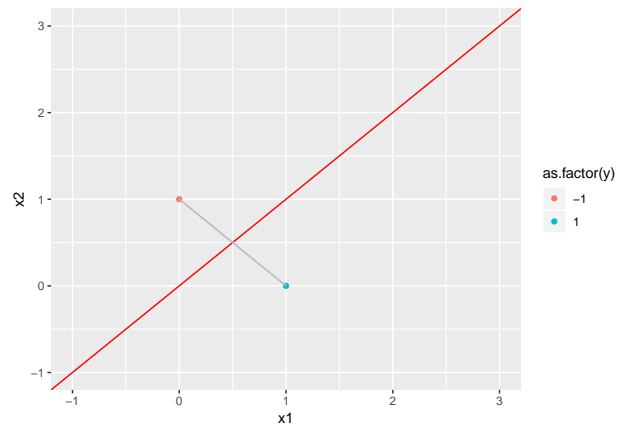
3.

Divide the training data into two parts, which are not overlapping with each other. Then use the first part, namely, training set to training the classifier as usual. After that, apply classifier on second part, namely, test set to calculate the loss.

Repeat the process several times for different training and test sets. Then choose the classifier that has the lowest loss on test set.

This process means to generalize the classifier to fits the new data. In the learning process, the test set can be seen as new data to classifier trained from training set. If the classifier is overfitting with training set, the loss on test will be large. Otherwise, the loss is low and the classifier can be well generalized to new data.

Problem 6



Problem 7

1.

$$\begin{aligned} \langle \vec{v}_H, \vec{x}_1 \rangle - c &= \left(\frac{1}{\sqrt{2}} \times (-3) + \frac{1}{\sqrt{2}} \times 0 \right) - \frac{1}{2\sqrt{2}} \\ &= -\frac{7}{2\sqrt{2}} < 0 \end{aligned}$$

$$\begin{aligned} \langle \vec{v}_H, \vec{x}_2 \rangle - c &= \left(\frac{1}{\sqrt{2}} \times \frac{1}{2} + \frac{1}{\sqrt{2}} \times \frac{1}{2} \right) - \frac{1}{2\sqrt{2}} \\ &= \frac{1}{2\sqrt{2}} > 0 \end{aligned}$$

So \vec{x}_1 is classified to class -1 , and \vec{x}_2 is classified to class $+1$.

2.

No, SVM is the improvement of linear classifier. It maximize the margin which makes sure that all the data points are away from the boundary as far as possible. In this case, SVM can better fits the new data than simple linear classifier.

So, same as linear classifier, SVM also calculates the $\langle \vec{v}_H, \vec{x} \rangle - c$ and use the sign as the class label.

3.

Sigmoid function is the approximate of 0 – 1 loss function. We approximate it because the sigmoid function is continuous, and we can calculate its derivative and use the gradient descent learning.