

**SGN-41007 Pattern Recognition and Machine Learning**  
**Exam 11.4.2017**  
**Heikki Huttunen**

- ▷ Use of calculator is allowed.
  - ▷ Use of other materials is not allowed.
  - ▷ The exam questions need not be returned after the exam.
  - ▷ You may answer in English or Finnish.
1. Are the following statements true or false? No need to justify your answer, just T or F. Correct answer: 1 pts, wrong answer:  $-\frac{1}{2}$  pts, no answer 0 pts.
- (a) Maximum likelihood estimators are unbiased.
  - (b) The Receiver Operating Characteristics curve plots the probability of detection versus the probability of false alarm for all thresholds.
  - (c) The number of support vectors of a support vector machine equals the total number of samples.
  - (d) Logistic regression classifier has linear decision boundary between classes.
  - (e) Dropout regularization improves the generalization of an LDA classifier.
  - (f) Stratified cross-validation resamples the training data such that all classes have equal number of samples.
2. The *Rayleigh distribution* is a probability distribution used *e.g.*, when modeling magnitude of a vector field. The density is defined as

$$p(x; \sigma) = \frac{x}{\sigma^2} \exp\left(-\frac{x^2}{2\sigma^2}\right) \quad \text{for } x > 0$$

We measure  $N$  samples:  $x_0, x_1, \dots, x_{N-1}$  and assume they are Rayleigh distributed and independent of each other.

- (a) Compute the probability  $p(\mathbf{x}; \sigma)$  of observing the samples  $\mathbf{x} = (x_0, x_1, \dots, x_{N-1})$ . (2p)
  - (b) Compute the logarithm of  $p(\mathbf{x}; \sigma)$  and differentiate the result with respect to  $\sigma$ . (2p)
  - (c) Find the maximum of the function, *i.e.*, the value where  $\frac{\partial}{\partial \sigma} \log p(\mathbf{x}; \sigma) = 0$ . (2p)
3. A dataset consists of two classes, containing four samples each. The samples are shown in Figure 1. The classes are linearly separable, and there are many linear decision boundaries that classify the training set with 100 % accuracy.
- (a) Find one such linear classifier. You can use whatever method you want (except the LDA), but justify your answer. Present the decision rule for sample  $\mathbf{x} \in \mathbb{R}^2$  in the following format:

$$\text{Class}(\mathbf{x}) = \begin{cases} 1, & \text{if } \boxed{\text{something}} \\ 2, & \text{otherwise} \end{cases}$$

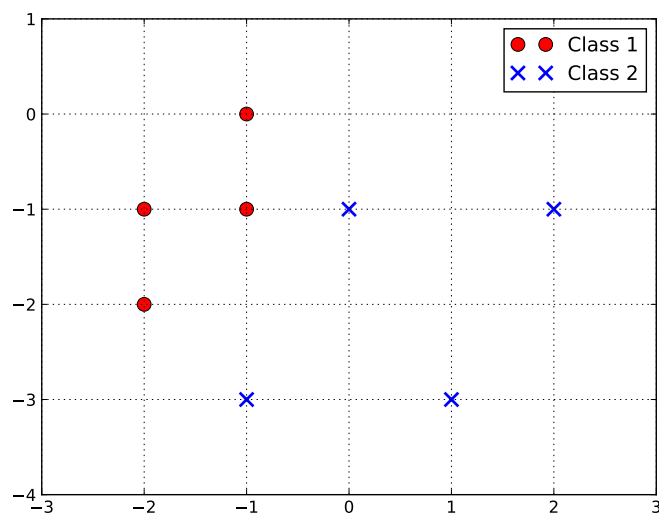


Figure 1: Training sample of question 3.

	Prediction	True label
Sample 1	0.8	1
Sample 2	0.5	1
Sample 3	0.6	0
Sample 5	0.4	0
Sample 4	0.2	0

Table 1: Results on test data for question 5a.

- (b) Find the Linear Discriminant Analysis (LDA) classifier for this data. You can choose the threshold arbitrarily, but the projection vector has to be the LDA projection. Present the decision rule in the above format in this case as well.

The covariances of the classes are

$$\mathbf{C}_1 = \frac{1}{3} \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \quad \mathbf{C}_2 = \frac{1}{3} \begin{pmatrix} 5 & 2 \\ 2 & 4 \end{pmatrix}.$$

4. (6 pts) Consider the Keras model defined in Listing 1. Inputs are  $128 \times 128$  color images from 10 categories.
- Draw a diagram of the network.
  - Compute the number of parameters for each layer, and their total number over all layers.
5. (a) (3p) A random forest classifier is trained on training data set and the `predict_proba` method is applied on the test data of five samples. The predictions and true labels are in Table 1. Draw the receiver operating characteristic curve. What is the Area Under Curve (AUC) score?

- (b) (3p) In the lectures we saw that the kernel trick  $\kappa(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^2$  for  $\mathbf{x} = (x_1, x_2)$  and  $\mathbf{y} = (y_1, y_2)$  corresponds to the mapping

$$\begin{pmatrix} u \\ v \end{pmatrix} \mapsto \begin{pmatrix} u^2 \\ v^2 \\ \sqrt{2}uv \end{pmatrix}$$

Find the explicit mapping corresponding to the inhomogeneous kernel  $\kappa(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^2$  with  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ .

Listing 1: A CNN model defined in Keras

```
model = Sequential()

w, h = 3, 3
sh = (3, 128, 128)

model.add(Convolution2D(32, w, h, input_shape=sh, border_mode='same'))
model.add(MaxPooling2D(pool_size=(4, 4)))
model.add(Activation('relu'))

model.add(Convolution2D(48, w, h, border_mode='same'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Activation('relu'))

model.add(Convolution2D(64, w, h, border_mode='same'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Activation('relu'))

model.add(Flatten())
model.add(Dense(100))
model.add(Activation('relu'))

model.add(Dense(10, activation = 'softmax'))
```

## Related Wikipedia pages

### Inversion of 2 × 2 matrices [\[ edit \]](#)

The *cofactor equation* listed above yields the following result for 2 × 2 matrices. Inversion of these matrices can be done as follows:<sup>[6]</sup>

$$\mathbf{A}^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det \mathbf{A}} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

### ROC space [\[ edit \]](#)

The contingency table can derive several evaluation "metrics" (see infobox). To draw a ROC curve, only the true positive rate (TPR) and false positive rate (FPR) are needed (as functions of some classifier parameter). The TPR defines how many correct positive results occur among all positive samples available during the test. FPR, on the other hand, defines how many incorrect positive results occur among all negative samples available during the test.

A ROC space is defined by FPR and TPR as *x* and *y* axes respectively, which depicts relative trade-offs between true positive (benefits) and false positive (costs). Since TPR is equivalent to sensitivity and FPR is equal to 1 – specificity, the ROC graph is sometimes called the sensitivity vs (1 – specificity) plot. Each prediction result or instance of a confusion matrix represents one point in the ROC space.

For degree-*d* polynomials, the polynomial kernel is defined as<sup>[2]</sup>

$$K(x, y) = (x^T y + c)^d$$

where *x* and *y* are vectors in the *input space*, i.e. vectors of features computed from training or test samples and *c* ≥ 0 is a free parameter trading off the influence of higher-order versus lower-order terms in the polynomial. When *c* = 0, the kernel is called homogeneous.<sup>[3]</sup> (A further generalized polykernel divides *x*<sup>T</sup>*y* by a user-specified scalar parameter *a*.<sup>[4]</sup>)

As a kernel, *K* corresponds to an inner product in a feature space based on some mapping *φ*:

$$K(x, y) = \langle \varphi(x), \varphi(y) \rangle$$

The nature of *φ* can be seen from an example. Let *d* = 2, so we get the special case of the quadratic kernel. After using the [multinomial theorem](#) (twice—the outermost application is the [binomial theorem](#)) and regrouping,

$$K(x, y) = \left( \sum_{i=1}^n x_i y_i + c \right)^2 = \sum_{i=1}^n (x_i^2) (y_i^2) + \sum_{i=2}^n \sum_{j=1}^{i-1} (\sqrt{2} x_i x_j) (\sqrt{2} y_i y_j) + \sum_{i=1}^n (\sqrt{2c} x_i) (\sqrt{2c} y_i) + c^2$$

From this it follows that the feature map is given by:

$$\varphi(x) = \langle x_n^2, \dots, x_1^2, \sqrt{2} x_n x_{n-1}, \dots, \sqrt{2} x_n x_1, \sqrt{2} x_{n-1} x_{n-2}, \dots, \sqrt{2} x_{n-1} x_1, \dots, \sqrt{2} x_2 x_1, \sqrt{2c} x_n, \dots, \sqrt{2c} x_1, c \rangle$$

feature space by an SVM is an ellipse in the input space.