

Application Background

Computer vision and image analysis are becoming more and more important in sports analytics, the science of analyzing and modeling processes underlying sporting events. Sports with a high media coverage create a demand for systematic review and objective evaluation of the performance of individual athletes as well as of teams. Across almost all sports, management and coaches make use of statistics and categorized video material to support their strategies.

In this assignment, we consider classification of the protagonists in soccer. A pattern detection algorithm has extracted regions of interest (ROIs) from a live video recording. Now each ROI has to be classified into one out of seven main classes. We have the obvious seven classes: outfielder (two classes, one for each team), goalkeeper (two classes) and referees. Further, there is a category that groups ROIs that contain at least one outfielder from each team. Finally, there is a class for irrelevant objects (i.e., false detections by the ROI detection algorithm). Sample ROIs are shown in Figure 1, the classes are summarized with the corresponding class index (label) in Table 1.

Because teams and referees can be identified based on the color of the clothing, color histogram features are extracted from the ROIs. The color model is HSV (hue, saturation, value), considering 3, 3, and 2 bits per channel, respectively.

The Exercises

Question 1

(data understanding and preprocessing). Download and extract the data.

Consider the training data `trainInput.csv` and the corresponding labels `trainTarget.csv`.

Plot a histogram showing the probability distribution of classes in the training data.

The i th row in `trainInput.csv` are the features of the i th training pattern. The class label of the i th pattern is given in the i th row of `trainTarget.csv`.

Deliverables: description of software used; single histogram plot

Question 2

(principal component analysis). Perform a principal component analysis of the training data `trainInput.csv`. Plot the eigenspectrum (see Figure 12.4 in [1] for an example). How many components are necessary to "explain 90% of the variance"? Visualize the data by a scatter plot of the data projected on the first two principal components.

Use `di`

different colors for the d_i

different classes in the plot (see Figure 12.8 in [1] for an example, which, however, lacks proper axes labels).

Deliverables: description of software used; plot of the eigenspectrum; indicate number of components necessary to explain 90% of variance; scatter plot of the data projected on the first two principal components with d_i

different colors indicating the 7 d_i

different classes

Question 3

(clustering). Perform 7-means clustering of `trainInput.csv` (feel free to play around with the number of clusters). After that, project the cluster centers to the first two principal components of the training data. Then visualize the clusters by adding the cluster centers to the plot from the previous exercise. Briefly

discuss the results: Did you

get meaningful clusters?

Deliverables: description of software used; one plot with cluster centers and data points;

short discussion of results

Question 4

(overtting). John Langford, who is \Doctor of Learning at Yahoo Research", maintains a very interesting blog (web log). Read the very true blog entry: \Clever methods of overtting," <http://hunch.net/?p=22>, 2005. Discuss if and how the di

erent types of overtting can occur when applying machine learning techniques to our sample sports

4

analytics application. Ignore the last type of overtting. You need not discuss issues related to reviewing of scientific papers (still, it is good to keep them in mind).

Deliverables: short discussion addressing the rst 10 \methods of overtting" listed in the blog entry

Question 5

(multi-class classication). Use a linear and a non-linear classication method (picking from the methods presented in the course) for classifying the 7 image classes, for example k-nearest neighbor and linear discriminant analysis (LDA). Use trainInput.csv and trainTarget.csv for training. After you trained a model, use the test data in testInput.csv and testTarget.csv to evaluate it. Report the classication error on both training and test set.

Deliverables: description of software used; arguments for your choice of classication methods; a short description of how you proceeded and what training and test results you achieved

Question 6

(binary classication using support vector machines). Now we consider binary classication using support vector machines (SVMs). To this end, we reduce the problem to distinguishing between referees and regions not containing persons of interest. Please use the data les trainInputBinary.csv, trainTargetBinary.csv, testInputBinary.csv, and testTargetBinary.csv. These are real-world data. The splitting into training and testing data has been done by the people providing the images. The class frequencies in training and testing data set di

er considerably. This could be bad luck, but could

also indicate a sampling bias (i.e., a violation of the i.i.d. assumption). The di

erences between training and test set make the problem dicult for some learning algorithms.

Please ignore this phenomenon in the assignment, that is, do not use asymmetric loss functions, class-dependent regularization parameters (i.e., di

erent\C-values" depending

on the class), etc.¹

For this exercise, use standard C-SVMs as introduced in the lecture. Employ radial Gaussian kernels of the form

$$k(x; z) = \exp(-$$

$$kx \dots zk^2) :$$

Here

$\gamma > 0$ is a bandwidth parameter that has to be chosen in the model selection process.

Note that instead of

often the parameter $\gamma =$

$\frac{1}{2}$

$\frac{1}{2}$

) is considered.

Jaakkola's heuristic provides a reasonable initial guess for the bandwidth parameter or

of a Gaussian kernel [3]. To estimate a good value for γ , consider all pairs consisting of a training input vector from the positive class and a training input vector from the negative class. Compute the difference in input space between all pairs. The median of these distances can be used as a measure of scale and therefore as a guess for γ . More formally, compute

$$G = \frac{1}{|S^+ \times S^-|} \sum_{(x_i, y_i) \in S^+ \times S^-} \sum_{(x_j, y_j) \in S^+ \times S^-} \|x_i - x_j\|$$

Be careful, because, for instance, the SVM from the Matlab Bioinformatics Toolbox may by default use different regularization parameters depending on the class and the class frequency.

5
 based on your training data S . Then set γ equal to the median of the values in G :
 $\gamma = \text{median}(G)$
 Compute the bandwidth parameter σ from γ using the identity given above.
 Use grid-search to determine appropriate SVM hyperparameters γ and C . Look at all combinations of $C \in \{0.1, 1, 10\}$ and $\gamma \in \{0.01, 0.1, 1, 10\}$.

2

$\gamma \in \{0.01, 0.1, 1, 10\}$

where the base b can be chosen to be either 2, the base e of the natural logarithm (Euler's number), or 10. Feel free to vary this grid. For each pair, estimate the performance of the SVM using 5-fold cross validation (see section 1.3 in [1]). Pick the hyperparameter pair with the lowest average 0-1 loss (classification error) and use it for training an SVM with the complete training dataset. Only use `trainInputBinary.csv` and `trainTargetBinary.csv` in the model selection and training process.

Report the values for C and γ you found in the model selection process. Compute the classification accuracy based on the 0-1 loss on the training data as well as on the test data `testInputBinary.csv` and `testTargetBinary.csv`. An accuracy on the test set significantly larger than 75% can be expected.
 Deliverables: description of software used; a short description of how you proceeded; initial

or value suggested by Jaakkola's heuristic; optimal C and γ found by grid search; classification accuracy on training and test data

Question 7

(linear regression vs. classification I.). 2 In this exercise, we explore the relation between regression and linear binary classification. Let the input space be $X = \mathbb{R}^p$ and the output space be Y . We are given n training patterns $S = \{(x_1; y_1), \dots, (x_n; y_n)\}$. Let n_1 of these patterns belong to the first class and n_2 to the second class, respectively (i.e. $n = n_1 + n_2$).

Show that the LDA classification rule as introduced in the lecture classes (see also

section 4.3 in [2]) to the second class if

$$x^T (\mu_2 - \mu_1) > \frac{1}{2} (\mu_2 - \mu_1)^T (\mu_2 + \mu_1) + \ln \frac{\sigma_1}{\sigma_2}$$

and to the first class otherwise. Here Σ is the empirical covariance matrix and μ_1 and μ_2 are the means of the training data from the first and second class, respectively. Assume that someone uses least squares regression in order to solve a binary classification problem. Further, assume that the two classes are encoded by $Y = \mu_1 \mathbf{1} + \mu_2 \mathbf{g}$. This exercise has been taken from [2].

6 Consider the regression objective function

$$\sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2;$$

with regression coefficients θ_0 and θ_1 and scalar offset/bias/intercept parameter θ_0 . Show that the optimal solution minimizing this least squares criterion is in the direction

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^{-1} (y_i - \bar{y})$$

(/ means proportional) :

Therefore the least squares regression coefficient is identical to the LDA coefficient, up to a scalar multiple.3 Hints (feel free to ignore): There are several ways to prove the statements. For the second part, you may find it convenient to get rid of θ_0 as in (3.3) in [1]. You may want to consult [4] (in particular section 2.4.2) to get the derivatives of vectors and matrices right. Deliverables: proofs including intermediate steps

Question 8

(linear regression vs. classification II.). As the previous exercise shows, abusing regression for binary classification is not a completely crazy thing to do. Does this also

hold for multi-class classification, in which the same least squares regression criterion is used and the class indices serve as the targets for regression (e.g., $Y = f_0; 1; \dots; 6g$)? You may conduct an experiment to support your answer, but you have to provide arguments beyond your empirical findings.

Deliverables: reason(s) why this is a reasonable approach or not.

References

[1] C. Bishop. Pattern Recognition and Machine Learning. Springer, 1 edition, 2006.

[2] T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning. Springer, 2009.

[3] T. Jaakkola, M. Diekhaus, and D. Haussler. Using the Fisher kernel method to detect remote protein homologies. In T. Lengauer, R. Schneider, P. Bork, D. Brutlad, J. Glasgow, H.-W. Mewes, and R. Zimmer, editors, Proceedings of the Seventh International Conference on Intelligent Systems for Molecular Biology, pages 149{158. AAAI Press, 1999.

[4] K. B. Petersen and M. S. Pedersen. The matrix cookbook. Technical University of Denmark, 2008.

3Actually, this result holds for any distinct coding of the two classes. However, in general the o

-

set/intercept/bias parameter di

ers depending on whether LDA or least squares regression is used to

determine it. For balanced class frequencies the two approaches coincide.

7