

Machine Learning Final Exam

Department of Computer Science, University of Copenhagen

Dhruv Chauhan

January 22, 2017

1 In a galaxy far, far away

1.1 Data preparation

The variance of the red-shifts in the spectroscopic training data was calculated to be:

$$0.0106$$

(where from now on, unless specified, values are shown to 3 significant places).

The MSE on the test SDSS predictions was calculated to be:

$$0.000812$$

This shows that the predictions were quite accurate.

1.2 Linear regression

The linear regression was done in Python, using the `sklearn` linear regression package. This performs an ordinary least squares linear regression. The error function is a Mean Squared Error.

The parameters of the model were (taken from the announcement):

$$\begin{aligned} &[0.0185134, 0.0479647, -0.0210943, -0.0274002, \\ &-0.0226798, 0.0064449, 0.0151842, 0.0120738, \\ &0.0103486, 0.00599684, -0.0294513, 0.069059, \\ &0.00630583, -0.00472042, -0.00873932, 0.00311043, \\ &0.0017252, 0.00435176] \end{aligned}$$

*** CHECK AND MAYBE PLUG THESE INTO THE MODEL ***

with bias term:

$$-0.801881$$

The error on the training data was calculated to be 0.00187, and on the test data was 0.00187 also. The errors normalised by the variance, σ_{red}^2 were equal to 0.176 for both the test and the training data.

This normalised error falling below one signifies that...

1.3 Non-Linear regression

For the non-linear regression, I chose to apply the K-nearest neighbours (KNN) algorithm. I chose this method for its simplicity (following Occam's razor), and therefore its intuitive understanding. The simplicity of the algorithm is also reflected in the single hyperparameter, k (if you consider a fixed distance metric), which means that there is less computation in tuning the hyperparameter.

I utilised the `neighbours` library from the `sklearn` package.

The KNN algorithm uses a distance metric to calculate the distance between a (set of) training point(s) and the other points. I used the Euclidian distance, given by $\|\mathbf{x} - \mathbf{x}'\|$, or $\sqrt{\mathbf{x}^T \mathbf{x}'}$. The algorithm works by calculating the distances from a test point to the other points, and then finding the nearest K points to that point. In a regression task, the test point is assigned the value of the mean of the nearest K neighbours.

My method involved using model selection methods such as cross-validation and grid search. Since this is model selection, cross validation was needed as we only use the training data during model selection. This gave us a better way to prevent overfitting of the training data. I used the `GridSearchCV` package from `sklearn`. The range of possible k values was given as the odd numbers between 1 and 29. This performed a 5-fold cross validation on each of the possible values of k , averaging out the resulting error (i.e. splitting the data into 5 equal chunks, using 4 as the training set, and 1 as the validation set, and then cycling through all possible 5 validation sets). The error function used in the algorithm was the mean squared error.

This method resulted in the optimum hyperparameter as $k = 7$, with a MSE of 0.00118 on the test data, and a MSE of 0.000870.

Clearly, the KNN Regressor worked better on the training data, which is to be expected due to the model's simplicity in using the k training points' average to return the regression results - therefore the training points are bound to have a low MSE. In comparison, the training data in the linear regressor performed the same as the test data. The difference in this would be due to the nature of a linear regressor, which would 'average out' the regression line over the points, thus leading to a small MSE on both the training and test data. On the test data, the KNN Regressor did perform a bit better than the linear regressor - perhaps due to the non-linear nature of KNN capturing the underlying nature of the data slightly more accurately. Overall, I believe the

KNN method worked well as we had a relatively low variance on the data - KNN can be skewed by big outliers in the data. The cross-validation helped to prevent overfitting on the data, which may have also caused the error on the test data to drop in comparison to the linear regressor.

2 Weed

2.1 Logistic Regression

The logistic regression was done in python, using the `LogisticRegression` package from `sklearn`. This implementation uses the logistic function. We are trying to classify testing points according to binary labels (0 - weeds, 1 - crops). The algorithm tries to minimise the following function:

$$\min_{w,c} \frac{1}{2} w^T w + C \sum_{i=1}^n \log e^{(-y_i(X_i^T w + c)) + 1}$$

It also uses a coordinate descent algorithm, which is a derivative-free optimization algorithm that performs a line search in one coordinate direction for the current point in each iteration. [1]

The parameters produced by the model were:

$$\begin{aligned} &[-0.0391283, 0.01549887, 0.00295803, 0.00033714, \\ &-0.00039954, 0.00348062, -0.00715483, 0.00473467, \\ &-0.02685346, -0.05592747, -0.04317431, 0.00579407, \\ &-0.00998736] \end{aligned}$$

with bias term:

$$-1.47771341e - 05$$

The zero-one loss on the training data was: 0.0200, and on the test data: 0.0348.

2.2 Binary classification using support vector machines

I used the `sklearn` python package, specifically the `svm.SVC` and `GridSearchCV` packages - in addition to `numpy`. Their model selection allows cross validation and grid-searching across values. It does this by splitting the training data into the training and validation sets across 5 folds and then running with the specified grid combination, for each option, finally returning the option with the lowest classification error. The error was calculated as the zero-one loss.

I used a kernel of the form:

$$k(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$$

Using the formula given, I calculated:

$$\begin{aligned}\sigma_{\text{Jaakkola}} &= 609 \\ \gamma_{\text{Jaakkola}} &= 1.35e - 06\end{aligned}$$

From that, the grid search looked over values of C and γ as given in the instructions, with $b = 10$.

From the grid search, the optimum hyperparameters were found to be:

$$\begin{aligned}C &= 1000 \\ \gamma &= 1.35e - 8\end{aligned}$$

The accuracy for the training and test sets is shown below:

$$\begin{aligned}\text{accuracy}_{\text{training}} &= 0.978 \\ \text{accuracy}_{\text{test}} &= 0.969\end{aligned}$$

From the accuracy results above, on both the training and test sets, the SVMs performed very well overall. This demonstrates the effectiveness of SVMs, especially in comparison to logistic regression.

2.3 Normalisation

The normalisation was performed by calculating the mean and the standard deviation of the training data. A function, f_{norm} was formed as below that would result in the training data having mean = 0 and variance = 1.

$$f_{\text{norm}} = \frac{\mathbf{x} - \mu}{\sigma}$$

This function was used to transform the training data to have the above mean and variance. The function was then also used to encode the test data, however with the *training* mean and variance (as we do not know the test mean and variance in most circumstances).

This normalised data was then used on the SVMs, resulting in:

$$\begin{aligned}\sigma_{\text{Jaakkola}} &= 1.38 \\ \gamma_{\text{Jaakkola}} &= 0.0300\end{aligned}$$

with optimum hyperparameters:

$$C = 100$$

$$\gamma = 0.0261$$

and test and training accuracy:

$$\begin{array}{ll} \text{accuracy}_{\text{training}} & = 0.983 \\ \text{accuracy}_{\text{test}} & = 0.969 \end{array}$$

For the logistic regression, the results on the normalised data were:

The parameters of the model were:

$$\begin{aligned} &[-0.15301671, 0.41886541, 1.20183536, 0.73123418, \\ &0.47797006, 1.60814454, -0.9043635, -1.69624452, \\ &-3.2546309, -2.55718762, -1.45943251, -0.27085139, \\ &-0.80955427] \end{aligned}$$

with bias:

$$-3.10513732$$

The zero-one loss on the training data was: 0.0300, and on the test data: 0.0383.

For SVMs, the accuracy improved marginally on the training data, but stayed constant for the test data. Since the choice of kernel uses the squared Euclidian distance ($\|x - x'\|^2$), this places a different level of importance on different features. In the training data, the ranges of the different features spanned from 112 for feature 1, to 7205 for feature 4. This squared Euclidian distance measure means that a difference of 1 in the 4th feature (a tiny amount of difference, relative to the range is equivalent to, about 0.01%) would be equivalent to a difference of approximately 0.90% in the 1st feature. This shows that small differences are exaggerated in feature 4 from the Euclidian distance calculation. By normalising the data, we even out these differences in ranges to adjust for this problem. In the transformed data, the range of feature 1 is 10.78, and feature 4 is 4.80, which means distance in both is much more relative. The results also converged faster than without normalising the data first.

However, for Logistic Regression, the loss increased for the training data samples, but decreased (by a very, very marginal amount, 0.001) for the test data. As we can see, the normalisation did not have as much of an effect as with the SVMs. This is due to the fact that algorithm looks at the proportional relationships between the coefficients, which doesn't change with normalisation, so it would be expected to perform similarly

to the non-normalised data.

With random forests, normalisation does not affect the performance. This is as in the algorithm, the magnitude of one feature is never directly compared to another feature - only one feature is split at each stage. Convergence and precision don't affect the performance of the random forest algorithm. The formal way of describing this is saying that random forests are invariant to monotonic transformations of individual features. [2]

2.4 Principal Component Analysis

The PCA was performed using a combination of the `sklearn.decomposition.PCA` package in python, and previous code I had written for an assignment. The package uses Singular Value Decomposition to project the data to a lower dimensionality.

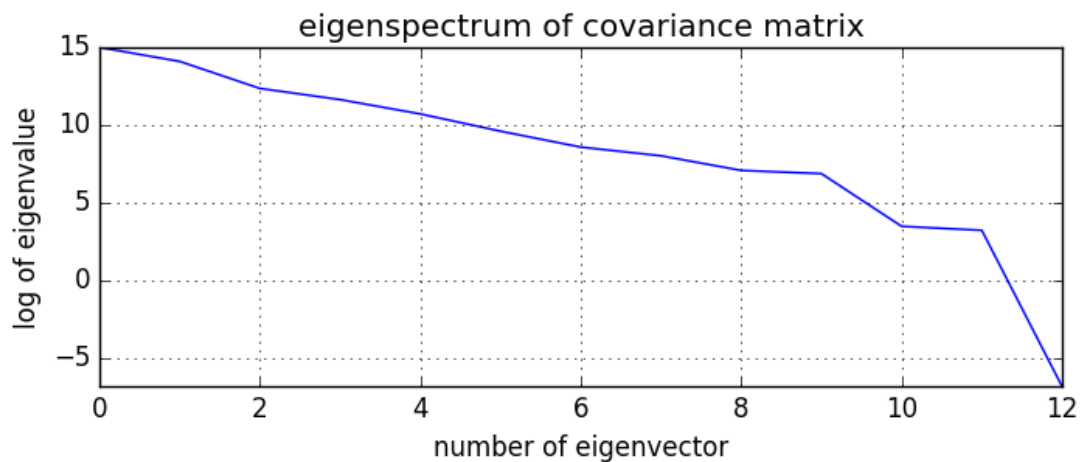


Figure 1: Eigenspectrum of the data

The PCA gave 13 principal components, of which only 2 were necessary to explain 90% of the variance. This is verified by looking at the eigenspectrum plot in Figure 1. This plots the log of the eigenvalues against the principal component number.

The scatter plot in Figure 2 shows the data projected onto the first two principal components of the resulting 13. The legend can be used to identify the class of the points.

Looking at the plot, one can see how important these first 2 principal components are, as visually the clusters are mainly separated by their colours. In comparison, a plot

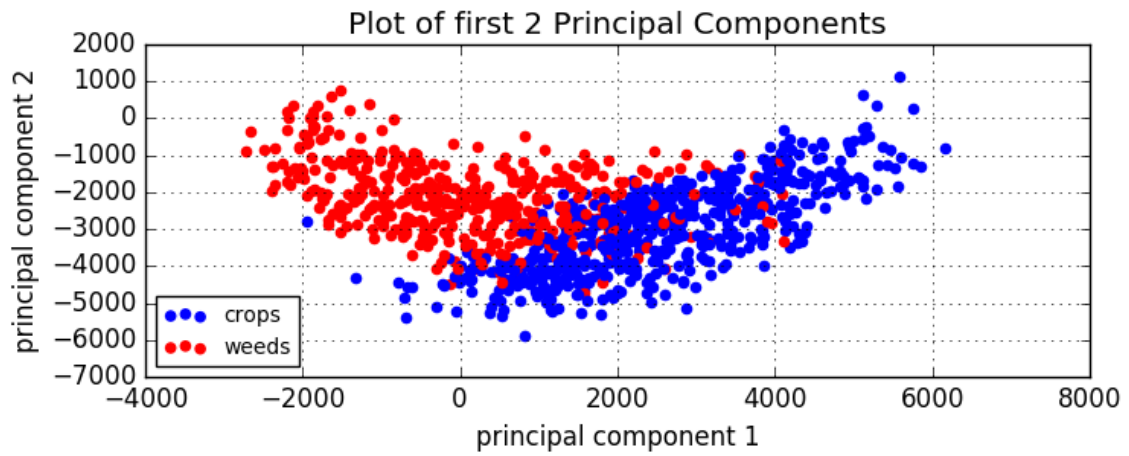


Figure 2: Plot of the data projected onto the first two principal components

of the data projected against the final 2 principal components is mostly concentrated around the center.

2.5 Clustering

References

- [1] Wikipedia, *Coordinate descent - wikipedia the free encyclopedia*, [Online; accessed 3-November-2016], 2016. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Coordinate_descent&oldid=747699222.
- [2] Ansari, *Stack overflow - normalisation svm, random forests*, 2017. [Online]. Available: <http://stats.stackexchange.com/questions/57010/is-it-essential-to-do-normalization-for-svm-and-random-forest>.