

Introductory Applied Machine Learning: Assignment 2

School of Informatics, University of Edinburgh

Instructors: Victor Lavrenko and Nigel Goddard
Assignment prepared by Sean J. Moran, revised by Boris Mitrovic

For due date and time, see the course web page.

Hard copy **and** electronic submission required. The time of the deadline will be strictly enforced. Ensure that your name does not appear on the document **AND** that your matriculation number does appear.

Remember that plagiarism is a university offence. Please read the policy at
<http://www.inf.ed.ac.uk/teaching/plagiarism.html> .

Marking Breakdown

70-100% results/answer correct plus extra achievement at understanding or analysis of results. Clear explanations, evidence of creative or deeper thought will contribute to a higher grade.

60-69% results/answer correct or nearly correct and well explained.

50-59% results/answer in right direction but significant errors.

40-49% some evidence that the student has gained some understanding, but not answered the questions properly.

0-39% serious error or slack work.

Mechanics

You should produce a word processed report in answer to this assignment (e.g. with $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$).

- **postscript or pdf** formats are acceptable for the report, other formats are not.
- you need to submit this report as a hard copy to the ITO **and** electronically as described below.

For the electronic submission place your report in a directory called `iamlans` and submit this using the `submit` command on a DICE machine. The format is

```
submit iaml 2 iamlans
```

You can check the status of your submissions with the `show-submissions` command.

NOTE: Your electronic submission will **not** count if you do not submit a hard copy of your report to the ITO.

Late submissions: The policy stated in the School of Informatics MSc Degree Guide is that normally you will not be allowed to submit coursework late. See

<http://www.inf.ed.ac.uk/teaching/years/msc/courseguide10.html#exam> for exceptions to this, e.g. in case of serious medical illness or serious personal problems.

Collaboration: You may discuss the assignment with your colleagues, provided that the writing that you submit is entirely your own. That is, you should **NOT** borrow actual text from other students. We ask that you list on the assignment sheet a list of the people who you've had discussions with (if any).

Important Instructions

- (a) In the following questions you are asked to run experiments using WEKA. The WEKA version installed on DICE is **Version 3.6.2**. If you are working on a machine other than DICE (e.g. your laptop), please make sure that you download and install the **same** version. This is important as your results need to be reproducible on DICE.
- (b) In many cases, the WEKA *Explorer* allows you to modify the random seed that will be used. Just using the default seed is fine. If you do change the seed you need to report the seed you have chosen.
- (c) You may find that WEKA crashes with an out-of-memory exception. If this should occur, refer to the instructions in IAML lab 1 in order to run WEKA with a larger memory allocation.
- (d) The .arff files that you will be using are available from the IAML website.
- (e) **IMPORTANT:** Keep your answers brief and concise. NOTE: you may **lose points** for a report longer than **900 words**.

Description of the dataset

This assignment is based on the automobile pricing dataset. Our goal will be to predict the price of automobiles based on various attributes. This data set consists of three types of entities: (a) the specification of an automobile in terms of various characteristics, (b) its assigned insurance risk rating, (c) its normalized losses in use as compared to other cars. The second rating corresponds to the degree to which the auto is more risky than its price indicates. Cars are initially assigned a risk factor symbol associated with its price. Then, if it is more risky (or less), this symbol is adjusted by moving it up (or down) the scale. Actuaries call this process "symboling". A value of +3 indicates that the auto is risky, -3 that it is probably pretty safe. The third factor is the relative average loss payment per insured vehicle year. This value is normalized for all autos within a particular size classification (two-door small, station wagons, sports/speciality, etc...), and represents the average loss per car per year.

To save you time and to make the problem manageable with limited computational resources, we preprocessed the original dataset. We removed any instances that had one or more missing values and randomized the data set. The resulting representation is much more compact and can be used directly to perform our experiments in WEKA.

1 Simple Linear Regression [30%]

We will begin by studying a simple Linear Regression model. Such a model will consider the relationship between a dependent (response) variable and only one independent (explanatory) variable. When applying machine learning in practice it can be prudent to start out simple in order to get a feeling for the dataset and for any potential difficulties that might warrant a more sophisticated model. In this Section we will consider one independent variable **engine-power** against the dependent variable **price**.

Important: *In order to reduce computation time you should use 5-fold cross validation unless explicitly specified otherwise (the default setting is 10).*

(a) Load the dataset **train_auto_partA.arff**. Have a look at the scatter plot of **price** against **engine-power** under the WEKA *Visualize* tab. Do you think that **engine-power** alone is sufficient for predicting the **price**? Please explain your answer in 2-3 sentences. Look at the distribution of the car prices. How would you preprocess it to improve the performance of linear regression? Don't do it at this stage, but instead in one sentence explain why you would do what you suggested.

(b) Now we will build a **SimpleLinearRegression** model. Under the *Classify* tab, choose *Functions* > *SimpleLinearRegression*. Keep the default settings for the model. Train the model and observe the results.

Examine the learnt Linear Regression model in the WEKA results buffer. What happens to the price as one more unit of `engine-power` is added? By examining the magnitude of the regression coefficient is it possible to tell whether or not `engine-power` is an important influential variable on `price`? Explain your answer in 1-2 sentences.

(c) Record the Root Mean Squared Error (RMSE), Mean Absolute Error (MAE) and Correlation Coefficient (CC) from the `SimpleLinearRegression` model you build in (b). What do these metrics intend to measure? Relate the values of CC, MAE and RMSE to the observations you made in question (a).

(d) Load the new dataset `train_auto_partA_base.arff`. Build a `SimpleLinearRegression` model on this dataset. Report the RMSE, MAE and CC metrics for this model. What is the *simplest* baseline model for the purposes of regression? Relate your answer to the regression model you have just built as part of this question.

2 Multivariate Linear Regression [70%]

In this Section we will fit a Multivariate Linear Regression model (`LinearRegression`) to the dataset. In contrast to Question 1, we will now train a model with multiple explanatory variables and ascertain how they affect our ability to predict the retail price of a car. One of our foremost concerns will be to determine exactly which attributes to include in the model and which may be left out.

Important: *In order to reduce computation time you should use 5-fold cross validation unless explicitly specified otherwise (the default setting is 10). When training a `LinearRegression` model be sure to set `attributeSelectionMethod` to `No attribute selection` and the `eliminateCovariateAttributes` option to `False`. Leave all other parameters on default unless otherwise indicated in the questions that follow.*

Important: *When performing various transformations to the data you must ensure that the “class” variable is set to `price` in the WEKA drop down list (this drop down list is located sandwiched between `More Options` and the `Start` button in the `Classify` tab) before running your model. If you do not do this your results will be vastly optimistic and entirely incorrect.*

(a) Load the dataset `train_auto_partB_numeric.arff` into WEKA. This version of the dataset contains the numeric attributes only. Go to the Visualize tab and examine whether or not any of the other attributes are particularly good at predicting the price. Can you find any? Do any attributes appear useless at predicting the price? Do any attributes exhibit significant correlations? As you answer these questions list two attributes for each question but do not modify the dataset at this stage. Of the attributes you listed, which ones could you safely remove?

(b) We will now make a first attempt at building a Multivariate `LinearRegression` model using the numeric attributes only. Under the `Classify` tab, choose `Functions > LinearRegression`. Now run the regression and note down the Root Mean Squared Error (RMSE), Mean Absolute Error (MAE) and Correlation Coefficient (CC). Comment on each in comparison to what you have obtained for the `SimpleLinearRegression` model in 1(c).

(c) Examine the histogram for the `engine-size` attribute. Is the distribution expected to cause a problem for regression? Explain your answer in 2-3 sentences. Transform this attribute using an appropriate simple technique from the lecture 3. Re-build a `LinearRegression` model on the transformed dataset and report the RMSE, MAE and CC metrics. How has the performance of your model changed? In 1-2 sents explain why.

(d) Reload the `train_auto_partB_numeric.arff` dataset. Regression performance can occasionally be enhanced by adding quadratic (interaction) terms to the model. For example, multiplying the attributes `Width` and `Height` together (`Width*Height`) would constitute an interaction term. What does adding interaction terms attempt to capture? We will experiment with interaction terms that involve `engine_size` and one other attribute. Only multiplication should be used as the means of combination. Add various interaction

terms of this type and report the results. Can you find any interaction terms that improved the regression performance of your model? *Hint:* The interaction terms can be added by using the `AddExpression` filter.

(e) So far we have performed regression with numeric attributes. We will now attempt to integrate nominal (categorical) attributes into our regression model. Load the dataset `train_auto_partB_full.arff` into WEKA. This dataset contains a mixture of numeric and nominal attributes. Why can we not use the nominal attributes in their current form for the purposes of regression? Convert the *nominal* attributes into a form that is more suitable for regression. In 2-3 sentences explain what you have done and why. Save this modified dataset as `train_auto_partB_edit.arff` as you will need it in question (f).

(f) In this question we will build a Multivariate `LinearRegression` model on the transformed dataset `train_auto_partB_edit.arff`. Under the *Classify* tab, choose *Functions > LinearRegression*. Now run the regression and observe the results. Record the RMSE, MAE, and CC. How does this more complex model perform with respect to your best performing model from either question (b), (c) or (d)? List one advantage and one disadvantage of using the more complex model.