



31250 Introduction to Data Analytics

## **Assignment #3**

Data Mining in Action

**Florian Lubitz**

FEIT

13688799

[florian.lubitz@student.uts.edu.au](mailto:florian.lubitz@student.uts.edu.au)

## Contents

<b>1</b>	<b>The data mining problem</b>	<b>1</b>
<b>2</b>	<b>Preprocessing and transformation</b>	<b>1</b>
<b>3</b>	<b>Attacking the problem</b>	<b>3</b>
<b>4</b>	<b>Classification techniques</b>	<b>3</b>
<b>5</b>	<b>The best classifier</b>	<b>3</b>
	<b>List of Figures</b>	<b>4</b>
	<b>Listings</b>	<b>4</b>

## 1 The data mining problem

In the given scenario we have received a data set. This contains data points for various customers of an insurance company. The aim of data mining is to predict whether a customer will purchase an insurance.

The dataset contains multiple features, that have been explored in Assessment 2 of this subject. The list of features contains different type of features. Some of them are dichotomous, some are categorical and some are rational or ordinal features. This wide variety will require some preprocessing to build a good model. Most of the features contain characters as values, one of them contains a date in australian format.

The goal for this problem will be to categorize each datapoint. There are only two categories available: 1 for "will purchase insurance" and 0 for "will not purchase insurance".

## 2 Preprocessing and transformation

To provide a good starting point for the classifiers I preprocess the raw data. To solve the problem I use python with jupyter notebooks. The following explanations will be supported by short listings. The main sourcecode can be found in the appendix of this report.

After importing the data from csv, I start by converting the date into a unix timestamp and all ordinal features into numbers. I also convert dichotomous features that are coded with "Y" and "N" into machine readable "0" and "1", respectively.

```
1 def parse_data(df):
2     # Convert Date
3     df['Original_Quote_Date'] = df['Original_Quote_Date'].apply(
4         str_to_timestamp)
5     # Convert bool-values to int of 1 and 0
6     df['Field_info4'] = df['Field_info4'].apply(string_to_bool)
7     df['Personal_info1'] = df['Personal_info1'].apply(string_to_bool)
8     df['Property_info1'] = df['Property_info1'].apply(string_to_bool)
9     df['Sales_info4'] = df['Sales_info4'].apply(string_to_value)
10    df['Personal_info3'] = df['Personal_info3'].apply(string_to_value)
11    df['Property_info3'] = df['Property_info3'].apply(string_to_value)
12    # Convert special amount to int
```

```
12 df['Field_info3'] = df['Field_info3'].apply(format_amount)
```

Listing 1: An excerpt of the parse function

After parsing all data I convert categorical features into many flags. This makes it easier for following classifiers to work with these features. To do this I use a function of `pandas` called `get_dummies`. After adding those flags I delete the original feature as it would contain redundant information. The shown function can handle multiple feature conversions at once.

```
1 def categorical_to_many(df, columns, keep_columns=None):
2     # Change Categorical
3     if keep_columns is None:
4         keep_columns = []
5     dummies = dict()
6     for col in columns:
7         dummies[col] = pd.get_dummies(df[col]).add_prefix(col + '_')
8     for dum in dummies:
9         # Keep generated columns as they might include lots of empty(same)
10            values
11        keep_columns = keep_columns + list(dummies[dum].keys())
12        df.drop(columns=[dum], inplace=True)
13        df = pd.concat([df, dummies[dum]], axis=1)
14    return df, keep_columns
```

Listing 2: The function to convert one categorical feature into many dichotomous

After converting the features, the training set and test set could contain different amount of features (flags). To solve this, I populate the sets with all missing flags.

```
1 # Fill up train and test frame to have the same column length
2 for key in list(set(train_df.keys()) - set(test_df.keys())):
3     test_df.loc[:, key] = pd.Series(np.zeros(len(test_df['
4         Original_Quote_Date'])), index=test_df.index)
5 for key in list(set(test_df.keys()) - set(train_df.keys())):
6     train_df.loc[:, key] = pd.Series(np.zeros(len(train_df['
7         Original_Quote_Date'])), index=train_df.index)
```

Listing 3: Add missing flag features to both sets

### **3 Attacking the problem**

### **4 Classification techniques**

### **5 The best classifier**

Quote\_ID, Original\_Quote\_Date, QuoteConversion\_Flag, Field\_info1, Field\_info2, Field\_info3, Field\_info4, Coverage\_info1, Coverage\_info2, Coverage\_info3, Sales\_info1, Sales\_info2, Sales\_info3, Sales\_info4, Sales\_info5, Personal\_info1, Personal\_info2, Personal\_info3, Personal\_info4, Personal\_info5, Property\_info1, Property\_info2, Property\_info3, Property\_info4, Property\_info5, Geographic\_info1, Geographic\_info2, Geographic\_info3, Geographic\_info4, Geographic\_info5.

## List of Figures

### Listings

1	An excerpt of the parse function . . . . .	1
2	The function to convert one categorical feature into many dichotomous .	2
3	Add missing flag features to both sets . . . . .	2