

# basic\_regression

April 24, 2019

**Copyright 2018 The TensorFlow Authors.**

```
In [0]: #@title Licensed under the Apache License, Version 2.0 (the "License");  
# you may not use this file except in compliance with the License.  
# You may obtain a copy of the License at  
#  
# https://www.apache.org/licenses/LICENSE-2.0  
#  
# Unless required by applicable law or agreed to in writing, software  
# distributed under the License is distributed on an "AS IS" BASIS,  
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.  
# See the License for the specific language governing permissions and  
# limitations under the License.
```

```
In [0]: #@title MIT License  
#  
# Copyright (c) 2017 François Chollet  
#  
# Permission is hereby granted, free of charge, to any person obtaining a  
# copy of this software and associated documentation files (the "Software"),  
# to deal in the Software without restriction, including without limitation  
# the rights to use, copy, modify, merge, publish, distribute, sublicense,  
# and/or sell copies of the Software, and to permit persons to whom the  
# Software is furnished to do so, subject to the following conditions:  
#  
# The above copyright notice and this permission notice shall be included in  
# all copies or substantial portions of the Software.  
#  
# THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR  
# IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY,  
# FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL  
# THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER  
# LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING  
# FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER  
# DEALINGS IN THE SOFTWARE.
```

# 1 Regression: predict fuel efficiency

[<a target="\\_blank" href="https://www.tensorflow.org/tutorials/keras/basic\\_regression">](https://www.tensorflow.org/tutorials/keras/basic_regression)

[<a target="\\_blank" href="https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/basic\\_regression.ipynb">](https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/basic_regression.ipynb)

[<a target="\\_blank" href="https://github.com/tensorflow/docs/blob/master/site/en/tutorials/keras/basic\\_regression.ipynb">](https://github.com/tensorflow/docs/blob/master/site/en/tutorials/keras/basic_regression.ipynb)

In a *regression* problem, we aim to predict the output of a continuous value, like a price or a probability. Contrast this with a *classification* problem, where we aim to select a class from a list of classes (for example, where a picture contains an apple or an orange, recognizing which fruit is in the picture).

This notebook uses the classic [Auto MPG](#) Dataset and builds a model to predict the fuel efficiency of late-1970s and early 1980s automobiles. To do this, we'll provide the model with a description of many automobiles from that time period. This description includes attributes like: cylinders, displacement, horsepower, and weight.

This example uses the `tf.keras` API, see [this guide](#) for details.

```
In [0]: # Use seaborn for pairplot
        !pip install seaborn
```

```
In [0]: from __future__ import absolute_import, division, print_function, unicode_literals

        import pathlib

        import matplotlib.pyplot as plt
        import pandas as pd
        import seaborn as sns

        import tensorflow as tf
        from tensorflow import keras
        from tensorflow.keras import layers

        print(tf.__version__)
```

## 1.1 The Auto MPG dataset

The dataset is available from the [UCI Machine Learning Repository](#).

### 1.1.1 Get the data

First download the dataset.

```
In [0]: dataset_path = keras.utils.get_file("auto-mpg.data", "http://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data")
        dataset_path
```

Import it using pandas

```
In [0]: column_names = ['MPG', 'Cylinders', 'Displacement', 'Horsepower', 'Weight',
                        'Acceleration', 'Model Year', 'Origin']
raw_dataset = pd.read_csv(dataset_path, names=column_names,
                           na_values = "?", comment='\t',
                           sep=" ", skipinitialspace=True)

dataset = raw_dataset.copy()
dataset.tail()
```

### 1.1.2 Clean the data

The dataset contains a few unknown values.

```
In [0]: dataset.isna().sum()
```

To keep this initial tutorial simple drop those rows.

```
In [0]: dataset = dataset.dropna()
```

The "Origin" column is really categorical, not numeric. So convert that to a one-hot:

```
In [0]: origin = dataset.pop('Origin')

In [0]: dataset['USA'] = (origin == 1)*1.0
dataset['Europe'] = (origin == 2)*1.0
dataset['Japan'] = (origin == 3)*1.0
dataset.tail()
```

### 1.1.3 Split the data into train and test

Now split the dataset into a training set and a test set.

We will use the test set in the final evaluation of our model.

```
In [0]: train_dataset = dataset.sample(frac=0.8, random_state=0)
test_dataset = dataset.drop(train_dataset.index)
```

### 1.1.4 Inspect the data

Have a quick look at the joint distribution of a few pairs of columns from the training set.

```
In [0]: sns.pairplot(train_dataset[["MPG", "Cylinders", "Displacement", "Weight"]], diag_kind=
```

Also look at the overall statistics:

```
In [0]: train_stats = train_dataset.describe()
train_stats.pop("MPG")
train_stats = train_stats.transpose()
train_stats
```

### 1.1.5 Split features from labels

Separate the target value, or "label", from the features. This label is the value that you will train the model to predict.

```
In [0]: train_labels = train_dataset.pop('MPG')
        test_labels = test_dataset.pop('MPG')
```

### 1.1.6 Normalize the data

Look again at the `train_stats` block above and note how different the ranges of each feature are.

It is good practice to normalize features that use different scales and ranges. Although the model *might* converge without feature normalization, it makes training more difficult, and it makes the resulting model dependent on the choice of units used in the input.

Note: Although we intentionally generate these statistics from only the training dataset, these statistics will also be used to normalize the test dataset. We need to do that to project the test dataset into the same distribution that the model has been trained on.

```
In [0]: def norm(x):
        return (x - train_stats['mean']) / train_stats['std']
        normed_train_data = norm(train_dataset)
        normed_test_data = norm(test_dataset)
```

This normalized data is what we will use to train the model.

Caution: The statistics used to normalize the inputs here (mean and standard deviation) need to be applied to any other data that is fed to the model, along with the one-hot encoding that we did earlier. That includes the test set as well as live data when the model is used in production.

## 1.2 The model

### 1.2.1 Build the model

Let's build our model. Here, we'll use a Sequential model with two densely connected hidden layers, and an output layer that returns a single, continuous value. The model building steps are wrapped in a function, `build_model`, since we'll create a second model, later on.

```
In [0]: def build_model():
        model = keras.Sequential([
            layers.Dense(64, activation=tf.nn.relu, input_shape=[len(train_dataset.keys())]),
            layers.Dense(64, activation=tf.nn.relu),
            layers.Dense(1)
        ])

        optimizer = tf.keras.optimizers.RMSprop(0.001)

        model.compile(loss='mean_squared_error',
                      optimizer=optimizer,
                      metrics=['mean_absolute_error', 'mean_squared_error'])
        return model

In [0]: model = build_model()
```

### 1.2.2 Inspect the model

Use the `.summary` method to print a simple description of the model

```
In [0]: model.summary()
```

Now try out the model. Take a batch of 10 examples from the training data and call `model.predict` on it.

```
In [0]: example_batch = normed_train_data[:10]
        example_result = model.predict(example_batch)
        example_result
```

It seems to be working, and it produces a result of the expected shape and type.

### 1.2.3 Train the model

Train the model for 1000 epochs, and record the training and validation accuracy in the history object.

```
In [0]: # Display training progress by printing a single dot for each completed epoch
        class PrintDot(keras.callbacks.Callback):
            def on_epoch_end(self, epoch, logs):
                if epoch % 100 == 0: print('')
                print('.', end='')

        EPOCHS = 1000

        history = model.fit(
            normed_train_data, train_labels,
            epochs=EPOCHS, validation_split = 0.2, verbose=0,
            callbacks=[PrintDot()])
```

Visualize the model's training progress using the stats stored in the history object.

```
In [0]: hist = pd.DataFrame(history.history)
        hist['epoch'] = history.epoch
        hist.tail()

In [0]: def plot_history(history):
        hist = pd.DataFrame(history.history)
        hist['epoch'] = history.epoch

        plt.figure()
        plt.xlabel('Epoch')
        plt.ylabel('Mean Abs Error [MPG]')
        plt.plot(hist['epoch'], hist['mean_absolute_error'],
                  label='Train Error')
        plt.plot(hist['epoch'], hist['val_mean_absolute_error'],
                  label = 'Val Error')
```

```
plt.ylim([0,5])
plt.legend()

plt.figure()
plt.xlabel('Epoch')
plt.ylabel('Mean Square Error [$MPG^2$]')
plt.plot(hist['epoch'], hist['mean_squared_error'],
         label='Train Error')
plt.plot(hist['epoch'], hist['val_mean_squared_error'],
         label = 'Val Error')
plt.ylim([0,20])
plt.legend()
plt.show()
```

```
plot_history(history)
```

This graph shows little improvement, or even degradation in the validation error after about 100 epochs. Let's update the `model.fit` call to automatically stop training when the validation score doesn't improve. We'll use an *EarlyStopping callback* that tests a training condition for every epoch. If a set amount of epochs elapses without showing improvement, then automatically stop the training.

You can learn more about this callback [here](#).

```
In [0]: model = build_model()
```

```
# The patience parameter is the amount of epochs to check for improvement
early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', patience=10)

history = model.fit(normed_train_data, train_labels, epochs=EPOCHS,
                    validation_split = 0.2, verbose=0, callbacks=[early_stop, PrintDot])

plot_history(history)
```

The graph shows that on the validation set, the average error is usually around +/- 2 MPG. Is this good? We'll leave that decision up to you.

Let's see how well the model generalizes by using the **test** set, which we did not use when training the model. This tells us how well we can expect the model to predict when we use it in the real world.

```
In [0]: loss, mae, mse = model.evaluate(normed_test_data, test_labels, verbose=0)

print("Testing set Mean Abs Error: {:.2f} MPG".format(mae))
```

## 1.2.4 Make predictions

Finally, predict MPG values using data in the testing set:

```
In [0]: test_predictions = model.predict(normed_test_data).flatten()

plt.scatter(test_labels, test_predictions)
plt.xlabel('True Values [MPG]')
plt.ylabel('Predictions [MPG]')
plt.axis('equal')
plt.axis('square')
plt.xlim([0,plt.xlim()[1]])
plt.ylim([0,plt.ylim()[1]])
_ = plt.plot([-100, 100], [-100, 100])
```

It looks like our model predicts reasonably well. Let's take a look at the error distribution.

```
In [0]: error = test_predictions - test_labels
plt.hist(error, bins = 25)
plt.xlabel("Prediction Error [MPG]")
_ = plt.ylabel("Count")
```

It's not quite gaussian, but we might expect that because the number of samples is very small.

### 1.3 Conclusion

This notebook introduced a few techniques to handle a regression problem.

- Mean Squared Error (MSE) is a common loss function used for regression problems (different loss functions are used for classification problems).
- Similarly, evaluation metrics used for regression differ from classification. A common regression metric is Mean Absolute Error (MAE).
- When numeric input data features have values with different ranges, each feature should be scaled independently to the same range.
- If there is not much training data, one technique is to prefer a small network with few hidden layers to avoid overfitting.
- Early stopping is a useful technique to prevent overfitting.