

Machine learning - random

Aheer Srabon

1 Section 1

```
from tensorflow import keras
from tensorflow.keras import layers

# Create a network with 1 linear unit
model = keras.Sequential([
    layers.Dense(units=1, input_shape=[3])
])
```

`keras.Sequential` creates a neural network as a stack of layers. The above example defines a model accepting three input features and producing a single output. The first argument `units` define how many outputs we want. The second argument `input_shape` tells keras the dimensions of the input. `input_shape` is the number of columns of the `DataFrame` except for the output column. It is a python list to permit use of more complex datasets.

Neural networks typically organize their neurons into layers. When we collect together linear units having a common set of inputs, we get a **dense** layer.

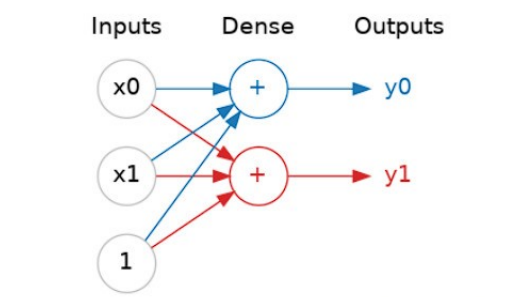


Figure 1: A dense layer of two linear units receiving two inputs and a bias.

Each layer in neural network performs some kind of relatively simple transformation. Through a deep stack of layers, a neural network can transform its inputs in more and more complex ways. In a well-trained neural network, each layer is a transformation getting us a little bit closer to a solution. A "layer" in Keras is a very general kind of thing. A layer can be, essentially, any kind of data transformation (like convolution, recurrence).

1.1 The activation function

It turns out that two dense layers with nothing in between are no better than a single dense layer by itself. Dense layers by themselves can never move us out of the world of lines and planes. An activation function brings non-linearity to the model. An **activation function** is simply some function that is applied to each of a layer's outputs (its activations). The most common is the rectifier function $\max(0, x)$. When we attach a rectifier to a linear unit, we get a rectified linear unit or **ReLU**. Applying a ReLU activation to a linear unit means the output becomes $\max(0, wx + b)$, which might be drawn like:

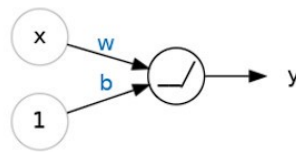


Figure 2: A rectified linear unit.

Stacking dense layers,

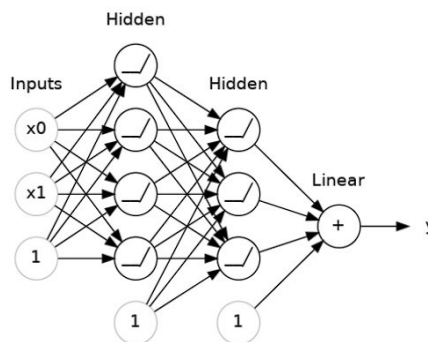


Figure 3: A stack of dense layers makes a "fully-connected" network.

The final (output) layer is a linear unit (with no activation function). That makes this network appropriate to a *regression* task. Other tasks (like classification) might require an activation function on the output.

1.2 Building sequential models

The Sequential model we've been using will connect together a list of layers in order from first to last: the first layer gets the input, the last layer produces the output. This creates the model in the figure above:

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
```

```

    # the hidden ReLU layers
    layers.Dense(units=4, activation='relu', input_shape=[2]),
    layers.Dense(units=3, activation='relu'),
    # the linear output layer
    layers.Dense(units=1),
])

```

Sequential takes a list of layers. Sometimes, some other layer are put between the Dense layer and its activation function. In these cases, the activation can be defined on its own Activation layer.

```

from tensorflow import keras
from tensorflow.keras import layers

```

```

model = keras.Sequential([
    # the hidden ReLU layers
    layers.Dense(units=4, activation='relu', input_shape=[2]),
    layers.Dense(units=3)
    # The activation is in its own layer. Between
    # the Dense layer and the Activation layer, some other layers can be
    # put
    layers.Activation('relu')
    # the linear output layer
    layers.Dense(units=1),
])

```

To train a neural network, the following things are needed,

- The model itself
- Training data
- A loss function (tells the network what problem to solve)
- An optimizer (tells the network how to solve the problem)

Some example of loss functions are,

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Huber loss

1.3 Stochastic gradient descent

Virtually all of the optimization algorithms used in deep learning belong to a family called **stochastic gradient descent**. They are iterative algorithms that train a network in steps.

1. Sample some training data and run it through the network to make predictions
2. Measure the loss between the predictions and the true values
3. Finally, adjust the weights in a direction that makes the loss smaller

Repeat the above process until the loss as small as necessary (or until it won't decrease any further). Each iteration's sample of training data is called a **minibatch** (or often just "batch"), while a complete round of the training data is called an **epoch**. A smaller learning rate means the network needs to see more **minibatches** before its weights converge to their best value. The number of epochs you train for is how many times the network will see each training example. *The learning rate and the size of the minibatches are the two parameters that have the largest effect on how the SGD training proceeds.*

Stochastic means "determined by chance". The training is stochastic because the minibatches are random samples from the dataset.

A common loss function for regression problems is **mean absolute error** or **MAE**.

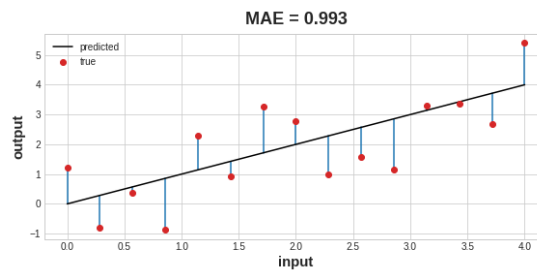


Figure 4: The mean absolute error is the average length between the fitted curve and the data points.

The loss function tells the network its objective. The optimizer is an algorithm that adjusts the weights of the model to minimize the loss. The learning rate and the size of the minibatches are the two parameters that have the largest effect on how the SGD training proceeds.

Adam is an SGD algorithm that has an adaptive learning rate that makes it suitable for most problems without any parameter tuning. Adam is a great general purpose optimizer. After defining a model, a loss function and optimizer with the models compile method can be defined,

```
model.compile(
    optimizer="adam",
    loss="mae",
)
```

Splitting the data into training data and validation data,

```
import pandas as pd
from IPython.display import display

red_wine = pd.read_csv('../input/dl-course-data/red-wine.csv')

# Create training and validation splits
df_train = red_wine.sample(frac=0.7, random_state=0)
df_valid = red_wine.drop(df_train.index)
display(df_train.head(4))

# Scale to [0, 1]
max_ = df_train.max(axis=0)
min_ = df_train.min(axis=0)
```

```
df_train = (df_train - min_) / (max_ - min_)
df_valid = (df_valid - min_) / (max_ - min_)
```

```
# Split features and target
X_train = df_train.drop('quality', axis=1)
X_valid = df_valid.drop('quality', axis=1)
y_train = df_train['quality']
y_valid = df_valid['quality']
```

```
_, input_shape = X_train.shape
```

Choose a three-layer network with over 1500 neurons.

```
from tensorflow import keras
from tensorflow.keras import layers
```

```
model = keras.Sequential([
    layers.Dense(512, activation='relu', input_shape=[11]),
    layers.Dense(512, activation='relu'),
    layers.Dense(512, activation='relu'),
    layers.Dense(1),
])
```

Deciding the architecture of the model should be part of the process. Start simple and use the validation loss as guide. Compile the optimizer and loss function,

```
model.compile(
    optimizer='adam',
    loss='mae'
)
```

The following code tells Keras to feed the optimizer 256 rows of the training data at a time (the batch.size) and to do that 10 times all the way through the dataset (the epochs),

```
history = model.fit(
    X_train, y_train,
    validation_data = (X_valid, y_valid),
    batch_size=256,
    epochs=10,
)
```

```
# convert the training history to a dataframe
history_df = pd.DataFrame(history.history)
history_df['loss'].plot()
```

1.4 Overfitting and Underfitting

Information of the training data as being two kinds,

- *Signal* - is the part that generalizes, the part that can help the model to make predictions from new data. Useful.

- *Noise* - is the part that is only true for the training data. It is all of the random fluctuation that comes from data in the real-world or all of the incidental, non-informative patterns that can't actually help the model make predictions. Not useful.

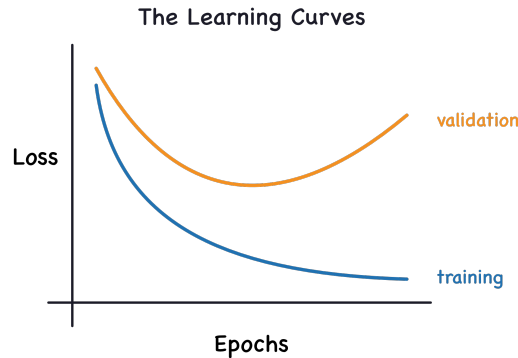


Figure 5: The validation loss gives an estimate of the expected error on unseen data.

The training loss will go down either when the model learns signal or when it learns noise. But the validation loss will go down only when the model learns signal (whatever noise the model learned from the training set won't generalize to new data). So, when a model learns signal, both curves go down, but when it learns noise, a gap is created in the curves. The size of the gap tells us how much noise the model has learned.

Underfitting vs overfitting,

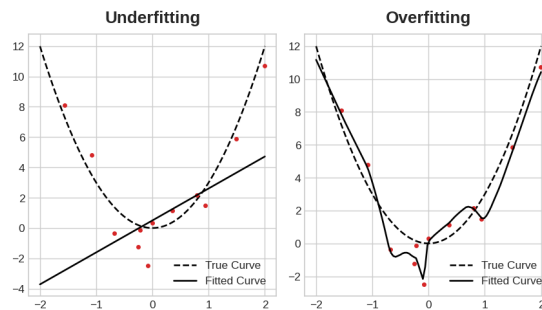


Figure 6: Underfitting vs overfitting

- **Underfitting** the training set is when the loss is not as low as it could be because the model hasn't learned enough *signal*.
- **Overfitting** the training set is when the loss is not as low as it could be because the model learned too much noise.

A model's **capacity** refers to the size and complexity of the patterns it is able to learn,

- Wider neural network – more units to existing layers. Learns linear relationships easily
- Deeper neural network – adding more layers. Learns non-linear relationships easily

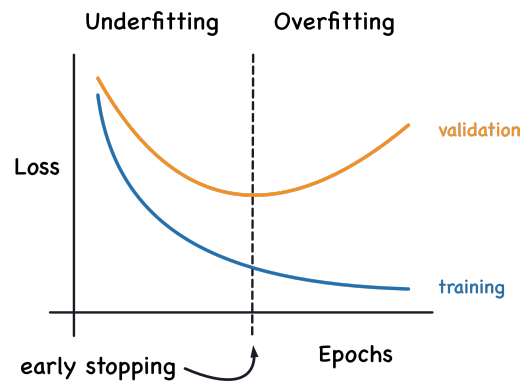


Figure 7: We keep the model where the validation loss is at a minimum - early stopping