

Homework 3: Ensemble Methods and Recurrent Neural Networks

Large-Scale Data Analysis 2020

Shahriyar Mahdi Robbani (XQR418)

May 25, 2020

1 AdaBoost (?? pts.)

Let $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \in (\mathbb{R}^d \times \{-1, 1\})^n$ be the training data, $t = 1, \dots, T$ denote the boosting rounds. At round t , $w_i^{(t)}$ is the weight of training pattern i , $h^{(t)}$ is the base classifier from round t , $\alpha^{(t)}$ is the importance of this classifier, $\epsilon^{(t)} = \sum_{j=1}^n w_j \mathbb{I}(h^{(t)}(\mathbf{x}_j) \neq y_j)$, and the aggregated classifier is $f^{(t)} = \sum_{i=1}^t \alpha^{(i)} h^{(i)}$.

1.1 Step 1

The goal is to prove for each i at any round t

$$w_i^{(t+1)} = w_i^{(t)} \frac{\exp(-\alpha^{(t)} h^{(t)}(\mathbf{x}_i) y_i)}{2\sqrt{\epsilon^{(t)}(1 - \epsilon^{(t)})}}. \quad (1)$$

We start from the update of the weights as defined in the lecture:

$$w_i^{(t+1)} = \begin{cases} w_i^{(t)} / (2\epsilon^{(t)}) & \text{if } h^{(t)}(\mathbf{x}_j) \neq y_j \\ w_i^{(t)} / (2(1 - \epsilon^{(t)})) & \text{otherwise} \end{cases} \quad (2)$$

And importance $\alpha^{(t)}$ of each weight defined in the lecture as:

$$\alpha^{(t)} = \frac{1}{2} \ln \left(\frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right) \quad (3)$$

The prediction $h^{(t)}(\mathbf{x}_j)$ and class y_j can only take the values 1 and -1. Therefore when $h^{(t)}(\mathbf{x}_j) = y_j$, we can rewrite equation (1) as:

$$w_i^{(t+1)} = w_i^{(t)} \frac{\exp(-\alpha^{(t)})}{2\sqrt{\epsilon^{(t)}(1 - \epsilon^{(t)})}} \quad (4)$$

Substituting equation (3) into equation (4) we get:

$$\begin{aligned} w_i^{(t+1)} &= w_i^{(t)} \frac{\exp\left(-\frac{1}{2} \ln \left(\frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right)\right)}{2\sqrt{\epsilon^{(t)}(1 - \epsilon^{(t)})}} \\ &= w_i^{(t)} \frac{\left(\frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right)^{-\frac{1}{2}}}{2\sqrt{\epsilon^{(t)}(1 - \epsilon^{(t)})}} \\ &= w_i^{(t)} \frac{\sqrt{1 - \epsilon^{(t)}}}{\sqrt{\epsilon^{(t)}}} \frac{1}{2\sqrt{\epsilon^{(t)}} \sqrt{1 - \epsilon^{(t)}}} \\ &= w_i^{(t)} / (2\epsilon^{(t)}) \end{aligned} \quad (5)$$

Similarly, when $h^{(t)}(\mathbf{x}_j) \neq y_j$, we can rewrite equation (1) as:

$$w_i^{(t+1)} = w_i^{(t)} \frac{\exp(\alpha^{(t)})}{2\sqrt{\varepsilon^{(t)}(1 - \varepsilon^{(t)})}} \quad (6)$$

Substituting equation (3) into equation (6) we get:

$$\begin{aligned} w_i^{(t+1)} &= w_i^{(t)} \frac{\exp\left(\frac{1}{2} \ln\left(\frac{1-\varepsilon^{(t)}}{\varepsilon^{(t)}}\right)\right)}{2\sqrt{\varepsilon^{(t)}(1 - \varepsilon^{(t)})}} \\ &= w_i^{(t)} \frac{\left(\frac{1-\varepsilon^{(t)}}{\varepsilon^{(t)}}\right)^{\frac{1}{2}}}{2\sqrt{\varepsilon^{(t)}(1 - \varepsilon^{(t)})}} \\ &= w_i^{(t)} \frac{\sqrt{\varepsilon^{(t)}}}{\sqrt{1 - \varepsilon^{(t)}}} \frac{1}{2\sqrt{\varepsilon^{(t)}}\sqrt{1 - \varepsilon^{(t)}}} \\ &= w_i^{(t)} / (2(1 - \varepsilon^{(t)})) \end{aligned} \quad (7)$$

We can see equations (5) and (7) are the two cases of equation (2), we have proved that equation (1) equals equation (2).

1.2 Step 2

The goal is to prove for each i at any round t

$$w_i^{(t)} \frac{\exp(-\alpha^{(t)} h^{(t)}(\mathbf{x}_i) y_i)}{2\sqrt{\varepsilon^{(t)}(1 - \varepsilon^{(t)})}} = \frac{\exp(-f^{(t)}(\mathbf{x}_i) y_i)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1 - \varepsilon^{(\tau)})}}. \quad (8)$$

Since each weight depends on the previous weight, we can calculate the weights recursively. So at round t we have a weight:

$$w_i^{(t)} = w_i^{(t-1)} \frac{\exp(-\alpha^{(t-1)} h^{(t-1)}(\mathbf{x}_i) y_i)}{2\sqrt{\varepsilon^{(t-1)}(1 - \varepsilon^{(t-1)})}} \quad (9)$$

Which is the same as:

$$w_i^{(t)} = w_i^{(t-2)} \frac{\exp(-\alpha^{(t-2)} h^{(t-2)}(\mathbf{x}_i) y_i)}{2\sqrt{\varepsilon^{(t-2)}(1 - \varepsilon^{(t-2)})}} \frac{\exp(-\alpha^{(t-1)} h^{(t-1)}(\mathbf{x}_i) y_i)}{2\sqrt{\varepsilon^{(t-1)}(1 - \varepsilon^{(t-1)})}} \quad (10)$$

Going all the way down to the base case, where $w_1 = \frac{1}{n}$, this becomes:

$$\begin{aligned} w_i^{(t)} &= \frac{1}{n} \frac{\prod_{\tau=1}^t \exp(-\alpha^{(\tau)} h^{(\tau)}(\mathbf{x}_i) y_i)}{\prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1 - \varepsilon^{(\tau)})}} \\ &= \frac{\exp(\sum_{\tau=1}^t -\alpha^{(\tau)} h^{(\tau)}(\mathbf{x}_i) y_i)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1 - \varepsilon^{(\tau)})}} \end{aligned} \quad (11)$$

The decision function at round t is defined as:

$$f^{(t)} = \sum_{i=1}^t \alpha^{(i)} h^{(i)} \quad (12)$$

Substituting this into equation (11) we get:

$$w_i^{(t)} = \frac{\exp(-f^{(t)}(\mathbf{x}_i) y_i)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1 - \varepsilon^{(\tau)})}} \quad (13)$$

1.3 Step 3

The goal is to prove for each i at any round t

$$\frac{\exp(-f^{(t)}(\mathbf{x}_i)y_i)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})}} = \frac{\exp(-f^{(t)}(\mathbf{x}_i)y_i)}{\sum_{l=1}^n \exp(-f^{(t)}(\mathbf{x}_l)y_l)} . \quad (14)$$

The weights are normalized, that is:

$$\begin{aligned} \sum_{i=1}^n w_i^{(t)} &= 1 \\ \sum_{i=1}^n \frac{\exp(-f^{(t)}(\mathbf{x}_i)y_i)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})}} &= 1 \\ \sum_{i=1}^n \frac{1}{n} \sum_{i=1}^n \exp(-f^{(t)}(\mathbf{x}_i)y_i) \sum_{i=1}^n \frac{1}{\prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})}} &= 1 \\ \sum_{i=1}^n \exp(-f^{(t)}(\mathbf{x}_i)y_i) &= \sum_{i=1}^n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})} \\ \sum_{i=1}^n \exp(-f^{(t)}(\mathbf{x}_i)y_i) &= n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})} \quad (15) \end{aligned}$$

Substituting equation (15) into equation (13) we get:

$$w_i^{(t)} = \frac{\exp(-f^{(t)}(\mathbf{x}_i)y_i)}{\sum_{l=1}^n \exp(-f^{(t)}(\mathbf{x}_l)y_l)} . \quad (16)$$

2 Gradient Boosting

3 Recurrent Neural Networks

Replace the simple recurrent layer in `LSDA2020_RNN1.ipynb` notebook with a LSTM layer. Follow the equations from `LSDA2020_RNN2.ipynb` to implement the LSTM. You will need to add more variables and extend the step function. Also remember that the LSTM not only transfers the hidden state h_t to the next time step $t+1$, but also the cell state c_t .

4 Recurrent Neural Networks in Keras

Add different components to the RNN from `LSDA2020_RNN2.ipynb` and report the results on the validation set (the changed parts in the code).

1. Add bidirectional sequence processing by utilizing `tf.keras.layers.Bidirectional`.
2. Stack 2 LSTM layers. What is the difference to bidirectional processing? (you may need to use the `return_sequences` parameter)
3. Add a 1-d convolution layer (`tf.keras.layers.Conv1D`) before the recurrent part. You will need to reshape the data, you can use the `tf.keras.layers.Reshape` layer for that.
4. Gradient clipping can be a helpful to train recurrent networks. Keras offers to clip gradients directly through the optimizer. Try this with clip values of 0.1, 1, and 10.

4.1 Bidirectional

```
def get_model_bidirectional(name="Bidirectional", shape=INPUT_SHAPE):
    inp = Input(shape)
    x = Bidirectional(LSTM(64))(inp)
    x = Dense(1)(x)

    model = Model(inp, x, name=name)
    model.summary()
    model.compile(sgd, loss='mae')
    return model
```

Model: "Bidirectional"

Layer (type)	Output Shape	Param #
input_11 (InputLayer)	[(None, 19, 22)]	0
bidirectional_1 (Bidirection	(None, 128)	44544
dense_10 (Dense)	(None, 1)	129

Total params: 44,673
Trainable params: 44,673
Non-trainable params: 0

Model: Bidirectional
Mean Error: 1.9972226952917502

4.2 Stacked LSTM

```
def get_model_stacked(name="stacked", shape=INPUT_SHAPE):
    inp = Input(shape)
    x = LSTM(64, return_sequences=True)(inp)
    x = LSTM(64, go_backwards=True)(x)
    x = Dense(1)(x)

    model = Model(inp, x, name=name)
    model.summary()
    model.compile(sgd, loss='mae')
    return model
```

Model: "stacked"

Layer (type)	Output Shape	Param #
input_12 (InputLayer)	[(None, 19, 22)]	0
lstm_11 (LSTM)	(None, 19, 64)	22272
lstm_12 (LSTM)	(None, 64)	33024
dense_11 (Dense)	(None, 1)	65

Total params: 55,361
Trainable params: 55,361

Non-trainable params: 0

Model: stacked

Mean Error: 1.9909247734882982

The bidirectional layers outputs two concatenated LSTM layers, while the stacked layers output one single LSTM layer. In the stacked layer, the input of the backward layer is affected by the forward layer since it passes through the forward layer, but the two layers are independent in the bidirectional layer.

4.3 Convolutional

```
def get_model_convolutional(name="convolutional", shape=INPUT_SHAPE):
    inp = Input(shape)
    x = Conv1D(filters=32, kernel_size=3, strides=1)(inp)
    x = LSTM(64, return_sequences=False)(x)
    x = Dense(1)(x)

    model = Model(inp, x, name=name)
    model.summary()
    model.compile(sgd, loss='mae')
    return model
```

Model: "convolutional"

Layer (type)	Output Shape	Param #
input_13 (InputLayer)	[(None, 19, 22)]	0
conv1d (Conv1D)	(None, 17, 32)	2144
lstm_13 (LSTM)	(None, 64)	24832
dense_12 (Dense)	(None, 1)	65

=====

Total params: 27,041

Trainable params: 27,041

Non-trainable params: 0

Model: convolutional

Mean Error: 2.0492551517846094

A recurrent neural network uses every piece of the input data to make a prediction. A convolutional neural network extracts features from the input data a subset of the input data and uses the features to make a prediction. Not all input data may be necessary to make a good prediction therefore convolutional networks can be less computationally expensive.

4.4 Gradient Clipping

```
optimizer = tf.keras.optimizers.Adam(learning_rate=lr, clipvalue=1)
get_model_convolutional().compile(optimizer, loss='mae')
```