

# BST 261: Data Science II

## Lecture 13

Heather Mattie

Department of Biostatistics  
Harvard T.H. Chan School of Public Health  
Harvard University

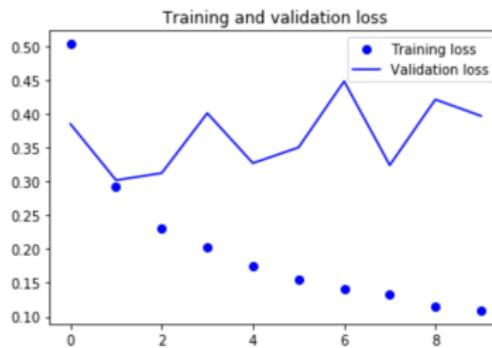
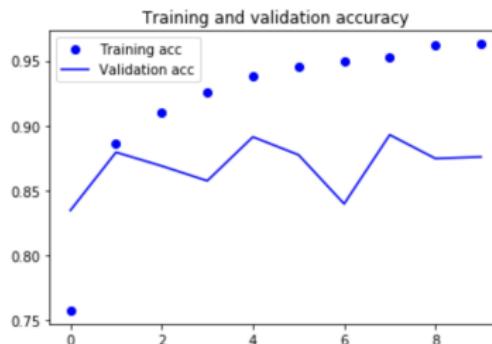
April 30, 2018

## LSTM in Keras

- Now that you have an idea of how LSTM works, let's implement it in Keras
- We set up a model using an LSTM layer and train it on the IMDB data
- The network is similar to the one with SimpleRNN that we discussed last week
- We only specify the output dimensionality of the LSTM layer, and leave every other argument (there are lots) to the Keras defaults

```
1 from keras.layers import LSTM
2
3 model = Sequential()
4 model.add(Embedding(max_features, 32))
5 model.add(LSTM(32)) (32 is circled)
6 model.add(Dense(1, activation='sigmoid'))
7
8 model.compile(optimizer='rmsprop',
9                 loss='binary_crossentropy',
10                metrics=['acc'])
11 history = model.fit(input_train, y_train,
12                      epochs=10,
13                      batch_size=128,
14                      validation_split=0.2)
```

- Now we get 89% accuracy, compared to 85% using the SimpleRNN layer
- This also performs better than when we used a feedforward network a few lectures ago, and here we are using less data



## Improving RNNs Performance and Generalization

- We will cover 3 techniques for improving RNNs:
  - **Recurrent dropout:** fights overfitting, different from the kind of dropout you are already familiar with
  - **Stacking recurrent layers:** increases generalizability, but comes with a higher computational cost
  - **Bidirectional recurrent layers:** increase accuracy and fight forgetting issues

- RNNs can be applied to any type of sequence data, not just text
- We will be using a weather timeseries data set recorded at the Weather Station at Max Planck Institute for Biochemistry in Jena, Germany
- 14 different variables were recorded **every 10 minutes** over several years, starting in 2003
  - Air temperature, atmospheric pressure, humidity, wind direction, etc.
- We will be using data from 2009-2016 to build a model that predicts air temperature 24 hours in the future using data from the last few days

6 recordings per hour

144 recordings per day

52,560 recordings per year

# Temperature Forecasting

- Load the data and take a look at the variables

```
1 cd ~/Downloads
2 mkdir jena_climate
3 cd jena_climate
4 wget https://s3.amazonaws.com/keras-datasets/jena_climate_2009_2016.csv.zip
5 unzip jena_climate_2009_2016.csv.zip
6
7 import os
8
9 data_dir = '/home/ubuntu/data/'
10 fname = os.path.join(data_dir, 'jena_climate_2009_2016.csv')
11
12 f = open(fname)
13 data = f.read()
14 f.close()
15
16 lines = data.split('\n') — each line is 1 recording
17 header = lines[0].split(',') — variable values are separated by
18 lines = lines[1:] — drop 1st row (header) Commas
19
20 print(header)
21 print(len(lines))
22
23 ['Date-Time', 'p (mbar)', 'T (degC)', 'Tpot (K)', 'Tdew (degC)', 'rh
   (%)', 'VPmax (mbar)', 'VPact (mbar)', 'VPdef (mbar)', 'sh (g/kg)', 'H2O
   C (mmol/mol)', 'rho (g/m***3)', 'wv (m/s)', 'max. wv (m/s)', 'wd (deg)']
24 420551 ≈ 8 years
```

- Convert all lines of data into a numpy array

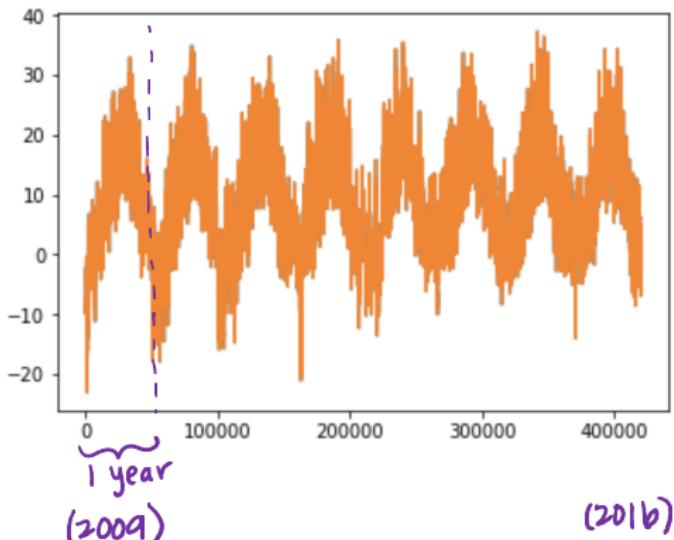
```
1 import numpy as np      420,551
2                                     [ ]   14
3 float_data = np.zeros((len(lines), len(header) - 1))
4 for i, line in enumerate(lines):
5     values = [float(x) for x in line.split(',')][1:]
6     float_data[i, :] = values
```

- One line of 14 values is saved as 1 row of 14 values in a 2D tensor
- Data shape: (420,551 , 14)

# Temperature Forecasting

- Plot the temperature over time (years)

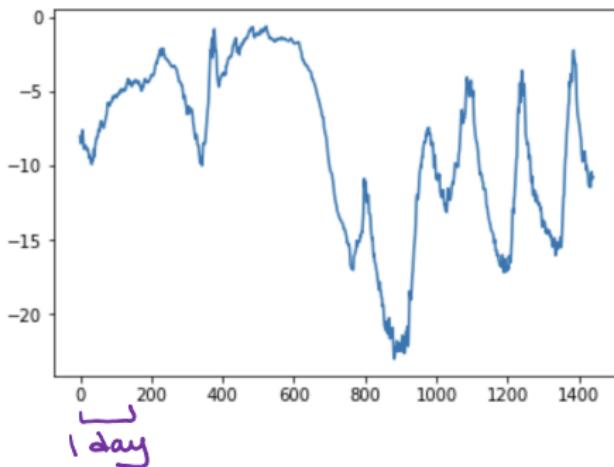
```
1 from matplotlib import pyplot as plt
2
3 temp = float_data[:, 1] # temperature (in degrees Celsius)
4 plt.plot(range(len(temp)), temp)
5 plt.show()
```



# Temperature Forecasting

- Plot the temperature over time (a few days)
- Notice that there is periodicity present, but that it isn't as consistent as the last plot
  - this will make predicting the weather in the next 24 hours using data from a few days beforehand more challenging

```
1 # Plot the temperature for the last 10 days
2 plt.plot(range(1440), temp[:1440])
3 plt.show()
```



- Task: given data going as far back as lookback timesteps (here a timestep is 10 minutes) and sampled every steps timesteps, can you predict the temperature in delay timesteps?
- lookback = 1440 - we will go back 10 days
- steps = 6 - observations will be sampled at one data point per hour
- delay = 144 - targets will be 24 hours in the future
- Process the data:
  - Normalize all variables  $\rightarrow \text{mean} = 0, \text{sd} = 1$

```
1 mean = float_data[:200000].mean(axis=0)
2 float_data -= mean
3 std = float_data[:200000].std(axis=0)
4 float_data /= std
```



$\text{axis}=0$  : collapse by rows  
 $\rightarrow$  find  $\bar{x}$  for each column

- Generate samples
  - data: The original array of floating point data, which we just normalized in the code on the last slide
  - lookback: How many timesteps back should our input data go
  - delay: How many timesteps in the future should our target be
  - min\_index and max\_index: Indices in the data array that delimit which timesteps to draw from. This is useful for keeping a segment of the data for validation and another one for testing.
  - shuffle: Whether to shuffle our samples or draw them in chronological order
  - batch\_size: The number of samples per batch
  - step: The period, in timesteps, at which we sample data. We will set it 6 in order to draw one data point every hour.

1440

144

6

# Temperature Forecasting

```
1 def generator(data, lookback, delay, min_index, max_index, shuffle=False,
2     batch_size=128, step=6):
3     if max_index is None:
4         max_index = len(data) - delay - 1
5     i = min_index + lookback
6     while 1:
7         if shuffle:
8             rows = np.random.randint(
9                 min_index + lookback, max_index, size=batch_size)
10        else:
11            if i + batch_size >= max_index:
12                i = min_index + lookback
13                rows = np.arange(i, min(i + batch_size, max_index))
14                i += len(rows)
15
16        samples = np.zeros((len(rows),
17                            lookback // step,
18                            data.shape[-1]))
19        targets = np.zeros((len(rows),))
20        for j, row in enumerate(rows):
21            indices = range(rows[j] - lookback, rows[j], step)
22            samples[j] = data[indices]
23            targets[j] = data[rows[j] + delay][1]
24
25        yield samples, targets
```

one batch of input data      corresponding target temperatures

training set  
randomly choose points in time

Choose patches of timesteps in chronological order

validation & test sets

- Use the generator function to instantiate three generators, one for training, one for validation and one for testing. Each will look at different temporal segments of the original data: the training generator looks at the first 200,000 timesteps, the validation generator looks at the following 100,000, and the test generator looks at the remainder.

```
1  lookback = 1440
2  step = 6
3  delay = 144
4  batch_size = 128
5
6  train_gen = generator(float_data,
7                  lookback=lookback,
8                  delay=delay,
9                  min_index=0,
10                 max_index=200000,
11                 shuffle=True,
12                 step=step,
13                 batch_size=batch_size)
14 val_gen = generator(float_data,
15                  lookback=lookback,
16                  delay=delay,
17                  min_index=200001,
18                  max_index=300000,
19                  step=step,
20                  batch_size=batch_size)
```

only shuffle the training data  
- does not shuffle the timesteps,  
it randomly chooses points in the  
training data to use as training  
examples

# Temperature Forecasting

```
1 test_gen = generator(float_data,
                      lookback=lookback,
                      delay=delay,
                      min_index=300001,
                      max_index=None,
                      step=step,
                      batch_size=batch_size)
2
3 # This is how many steps to draw from 'val_gen' in order to see the whole
4     validation set:
5 val_steps = (300000 - 200001 - lookback) // batch_size
6
7 # This is how many steps to draw from 'test_gen' in order to see the whole test
8     set:
9 test_steps = (len(float_data) - 300001 - lookback) // batch_size
```



floor division

# Temperature Forecasting Baseline

- We need to come up with a baseline benchmark to beat
- Common-sense approach: always predict that the temperature 24 hours from now will be equal to the temperature now
- We'll use mean absolute error (MAE) to measure loss

```
1 # Define MAE
2 np.mean(np.abs(preds - targets))
3
4 def evaluate_naive_method():
5     batch_maes = []
6     for step in range(val_steps):
7         samples, targets = next(val_gen)
8         preds = samples[:, -1, 1]
9         mae = np.mean(np.abs(preds - targets))
10        batch_maes.append(mae)
11    print(np.mean(batch_maes))
12
13 evaluate_naive_method()
14 0.289735972991
```

$$\frac{1}{m} \sum |y_i - \hat{y}_i|$$

- We get MAE = 0.29. Since our temperature data has been normalized to be centered at 0 and have a standard deviation of 1, this number is not immediately interpretable. It translates to an average absolute error of  $0.29 * \text{temperature\_std}$  degrees Celsius, i.e.  $2.57^\circ\text{C}$ . That's a fairly large average absolute error – now the task is to leverage our knowledge of deep learning to do better.

# Temperature Forecasting Simple Model

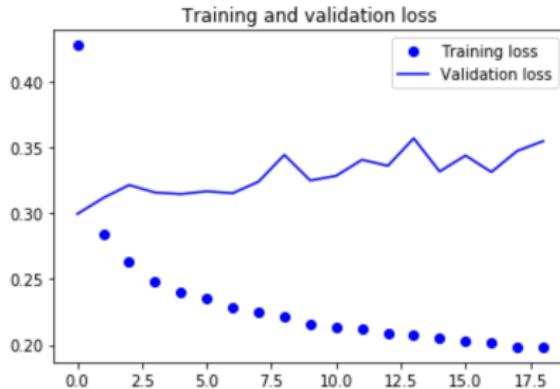
- Let's first try a simple model before developing a more complex one
- In general it's best to start with a basic model and then work your way up in complexity

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.Flatten(input_shape=(lookback // step, float_data.shape[-1])))
7 model.add(layers.Dense(32, activation='relu'))
8 model.add(layers.Dense(1))
9
10
11 model.compile(optimizer=RMSprop(), loss='mae')
12 history = model.fit_generator(train_gen,
13                               steps_per_epoch=500,
14                               epochs=20,
15                               validation_data=val_gen,
16                               validation_steps=val_steps)
```

$1440 // 6 = 240$       14

# Temperature Forecasting Simple Model

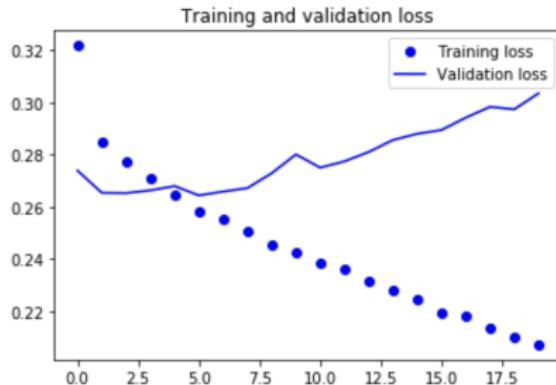
- The loss hovers around the loss value of the common-sense baseline, proving that having a baseline to beat is useful
- Why isn't the loss better than the baseline?
  - It doesn't take time into account
  - The hypothesis space of possible models is large



- Now let's fit a model with a GRU layer

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.GRU(32, input_shape=(None, float_data.shape[-1])))
7 model.add(layers.Dense(1))
8
9 model.compile(optimizer=RMSprop(), loss='mae')
10 history = model.fit_generator(train_gen,
11                               steps_per_epoch=500,
12                               epochs=20,
13                               validation_data=val_gen,
14                               validation_steps=val_steps)
```

- This model is better than the previous simple model and the common-sense baseline
- Evidence of overfitting, so let's try dropout next



# Recurrent Dropout

- It turns out that the classic technique of dropout we saw in earlier lectures can't be applied in the same way for recurrent layers
- Applying dropout before a recurrent layer impedes learning rather than helping to implement regularization
- The proper way to apply dropout with a recurrent network was discovered in 2015
  - Yarin Gal, "Uncertainty in Deep Learning (PhD Thesis),"  
[http://mlg.eng.cam.ac.uk/yarin/blog\\_2248.html](http://mlg.eng.cam.ac.uk/yarin/blog_2248.html)
  - The **same pattern of dropped units** should be applied at every timestep
  - This allows the network to properly propagate its learning error rate through time - a temporally random dropout pattern would disrupt the error signal and hinder the learning process
- Yarin's mechanism has been built into Keras
  - Every recurrent layer has 2 dropout-related arguments:
    - `dropout`: a float number specifying the dropout rate for input units of the layer
    - `recurrent_dropout`: a float number specifying the dropout rate of the recurrent units

- Gal and Ghahramani, (2015) “A Theoretically Grounded Application of Dropout in Recurrent Neural Networks”
- <https://pdfs.semanticscholar.org/58ab/4ad843ce2a3494578ca62a70d189c10d7e3b.pdf>

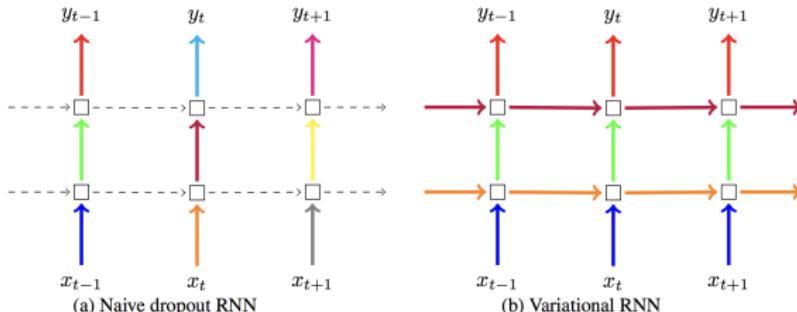


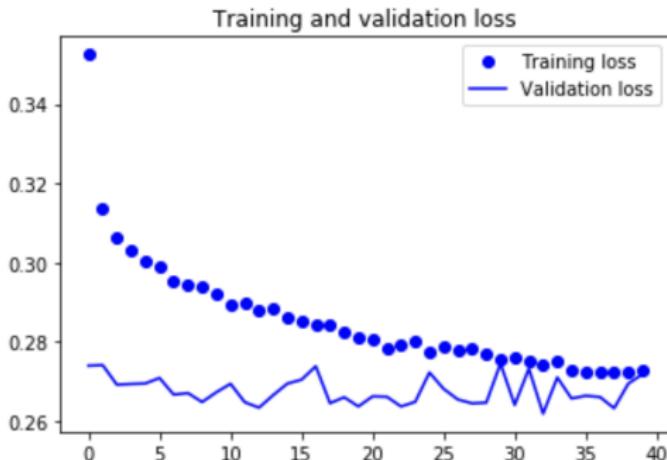
Figure 1: **Depiction of the dropout technique following our Bayesian interpretation (right) compared to the standard technique in the field (left).** Each square represents an RNN unit, with horizontal arrows representing time dependence (recurrent connections). Vertical arrows represent the input and output to each RNN unit. Coloured connections represent dropped-out inputs, with different colours corresponding to different dropout masks. Dashed lines correspond to standard connections with no dropout. Current techniques (naive dropout, left) use different masks at different time steps, with no dropout on the recurrent layers. The proposed technique (Variational RNN, right) uses the same dropout mask at each time step, including the recurrent layers.

# Temperature Forecasting Simple RNN Model with Dropout

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.GRU(32,
7                     dropout=0.2,
8                     recurrent_dropout=0.2,
9                     input_shape=(None, float_data.shape[-1])))
10 model.add(layers.Dense(1))
11
12 model.compile(optimizer=RMSprop(), loss='mae')
13 history = model.fit_generator(train_gen,
14                               steps_per_epoch=500,
15                               epochs=40,
16                               validation_data=val_gen,
17                               validation_steps=val_steps)
```

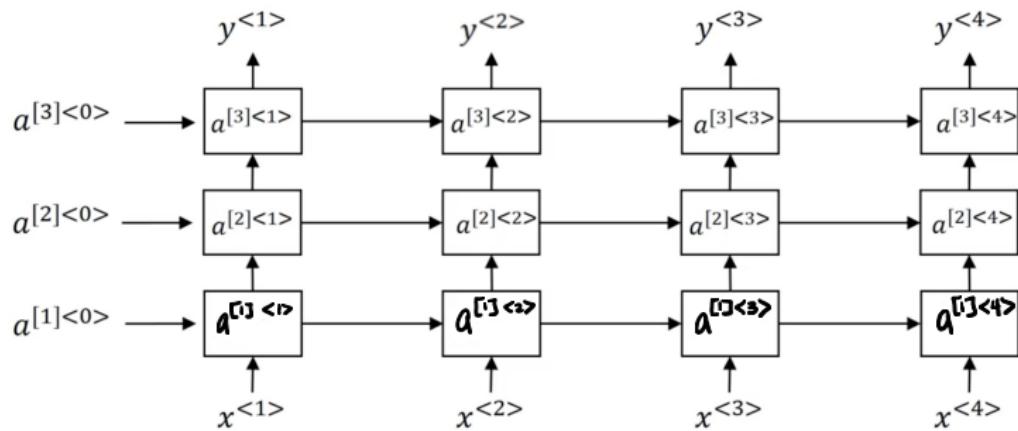
# Temperature Forecasting Simple RNN Model with Dropout

- Success: we are no longer overfitting during the first 30 epochs
- We have more stable evaluation scores, but our best scores are not much lower than they were previously



## Stacking Recurrent Layers

# Stacking Recurrent Layers



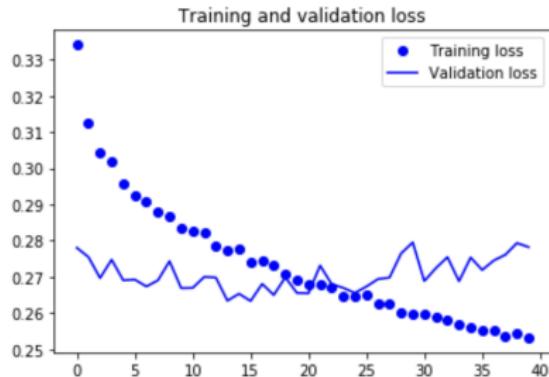
# Temperature Forecasting Stacking Recurrent Layers

- We have dealt with overfitting, but our model still doesn't perform much better than the common-sense baseline
- Let's add a recurrent layer to increase the complexity of the model

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.GRU(32,
7                     dropout=0.1,
8                     recurrent_dropout=0.5,
9                     return_sequences=True,
10                    input_shape=(None, float_data.shape[-1])))
11 model.add(layers.GRU(64, activation='relu',
12                     dropout=0.1,
13                     recurrent_dropout=0.5))
14 model.add(layers.Dense(1))
15
16 model.compile(optimizer=RMSprop(), loss='mae')
17 history = model.fit_generator(train_gen,
18                               steps_per_epoch=500,
19                               epochs=40,
20                               validation_data=val_gen,
21                               validation_steps=val_steps)
```

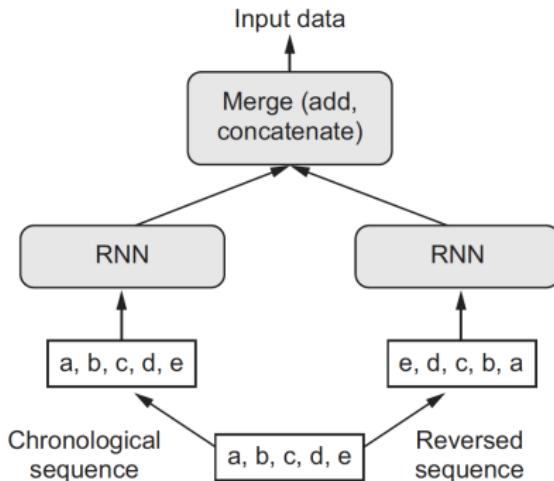
# Temperature Forecasting Stacking Recurrent Layers

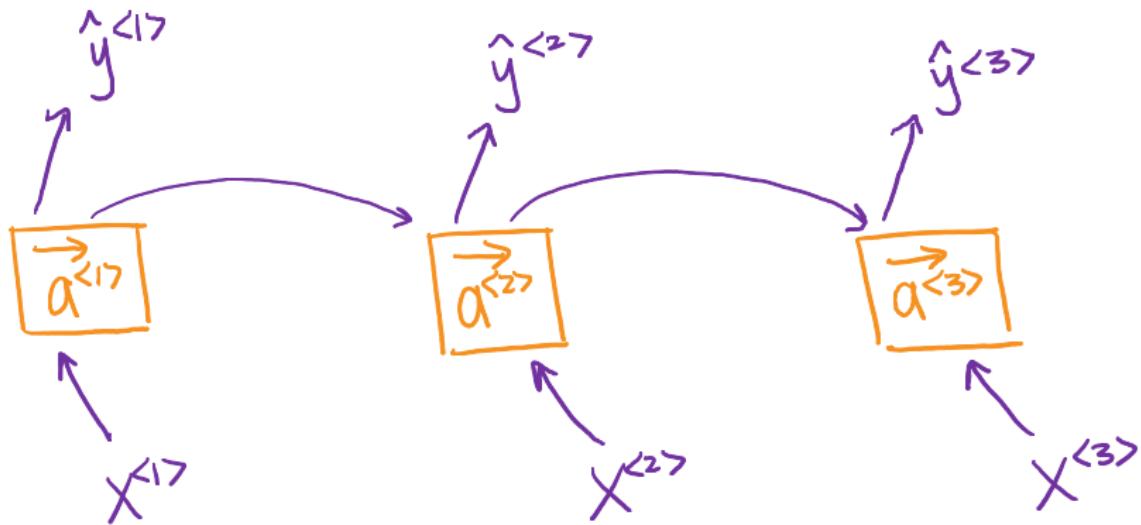
- The added layer does improve the results, but not significantly
- We could add more layers, but this will be more computationally expensive
- Let's now try another approach - a bidirectional RNN



## Bidirectional RNNs

- A **bidirectional RNN** (BRNN) can offer greater performance on certain tasks
- Frequently used in natural-language processing (NLP)
- BRNNs exploit the order sensitivity of RNNs
  - Uses 2 regular RNNs, each of which processes the input sequence in one direction (chronologically and antichronologically), and then merges their representations
  - Catches patterns that may be overlooked by a regular RNN





# Temperature Forecasting with Antichronological RNN

```
1 def reverse_order_generator(data, lookback, delay, min_index, max_index,
2                             shuffle=False, batch_size=128, step=6):
3     if max_index is None:
4         max_index = len(data) - delay - 1
5     i = min_index + lookback
6     while 1:
7         if shuffle:
8             rows = np.random.randint(
9                 min_index + lookback, max_index, size=batch_size)
10        else:
11            if i + batch_size >= max_index:
12                i = min_index + lookback
13            rows = np.arange(i, min(i + batch_size, max_index))
14            i += len(rows)
15
16        samples = np.zeros((len(rows),
17                            lookback // step,
18                            data.shape[-1]))
19        targets = np.zeros((len(rows),))
20        for j, row in enumerate(rows):
21            indices = range(rows[j] - lookback, rows[j], step)
22            samples[j] = data[indices]
23            targets[j] = data[rows[j] + delay][1]
24        yield samples[:, ::-1, :], targets
```

only this line changes - reverses the order

# Temperature Forecasting with Antichronological RNN

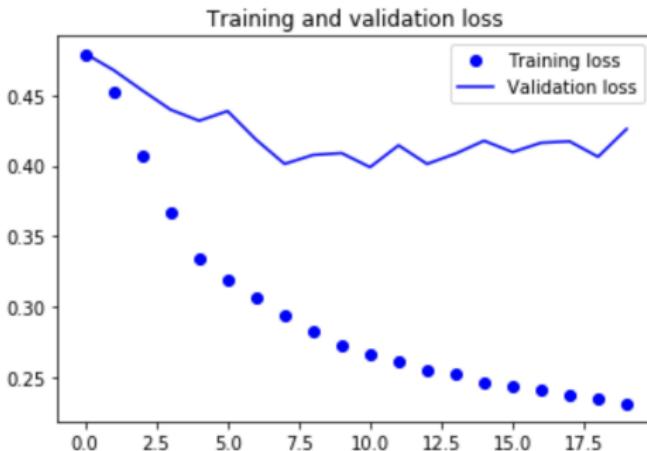
```
1 train_gen_reverse = reverse_order_generator(
2     float_data,
3     lookback=lookback,
4     delay=delay,
5     min_index=0,
6     max_index=200000,
7     shuffle=True,
8     step=step,
9     batch_size=batch_size)
10 val_gen_reverse = reverse_order_generator(
11     float_data,
12     lookback=lookback,
13     delay=delay,
14     min_index=200001,
15     max_index=300000,
16     step=step,
17     batch_size=batch_size)
```

# Temperature Forecasting with Antichronological RNN

```
1 model = Sequential()
2 model.add(layers.GRU(32, input_shape=(None, float_data.shape[-1])))
3 model.add(layers.Dense(1))
4
5 model.compile(optimizer=RMSprop(), loss='mae')
6 history = model.fit_generator(train_gen_reverse,
7                               steps_per_epoch=500,
8                               epochs=20,
9                               validation_data=val_gen_reverse,
10                             validation_steps=val_steps)
```

# Temperature Forecasting with Antichronological RNN

- The antichronological RNN significantly underperforms the original, chronological RNN as well as the common-sense baseline
- Why?
  - The GRU layer will be better at remembering the recent past than the distant past
  - More recent weather data points are more predictive than older data points
- Note that this isn't true for other types of problems
  - Order typically doesn't matter as much in NLP



# Temperature Forecasting with Bidirectional RNN

- To instantiate a bidirectional RNN in Keras, use the `Bidirectional` layer, which takes as first argument a recurrent layer instance
- `Bidirectional` will create a second, separate instance of this recurrent layer, and will use one instance for processing the input sequences in chronological order and the other instance for processing the input sequences in reversed order

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.Bidirectional(
7     layers.GRU(32), input_shape=(None, float_data.shape[-1])))
8 model.add(layers.Dense(1))
9
10 model.compile(optimizer=RMSprop(), loss='mae')
11 history = model.fit_generator(train_gen,
12                               steps_per_epoch=500,
13                               epochs=40,
14                               validation_data=val_gen,
15                               validation_steps=val_steps)
```

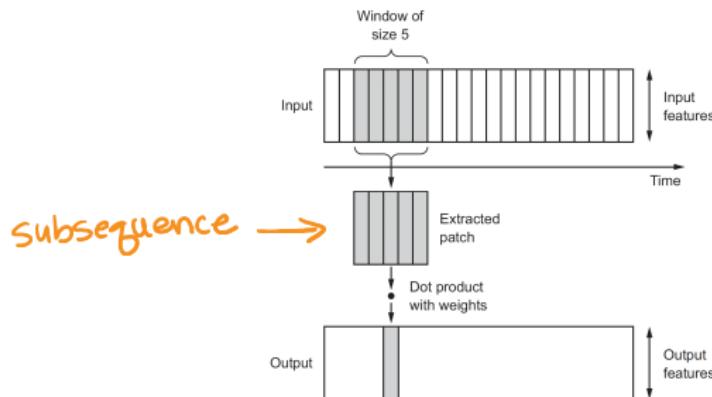
- This model performs about as well as the regular GRU layer (Loss  $\approx 0.3$ )
- All of the predictive capacity comes from the chronological half of the network

- There are several other things you can try to improve performance
  - Change the number of units in each recurrent layer
  - Try using LSTM layers instead of GRU layers
  - Change the learning rate used by the RMSprop optimizer
  - Try a bigger densely connected classifier on top of the recurrent layers

## 1D Convolution for Sequence Data

# 1D Convolution for Sequence Data

- Recall that CNNs can extract features from local input patches and then recognize them anywhere
- If we think of time as a spatial dimension, we can use 1D CNNs for sequence data
- Great for audio generation and machine translation
- Faster to run than RNNs
- 1D CNNs extract subsequences (patches) from sequences and perform the same transformation on each subsequence
  - A pattern learned at a specific position of a sequence can be recognized at a different position (translation invariant)
- Pooling in this case is similar to what we have seen before - output the maximum or average value of a subsequence



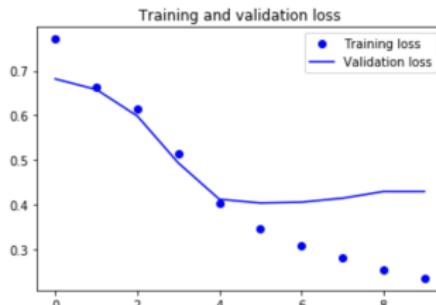
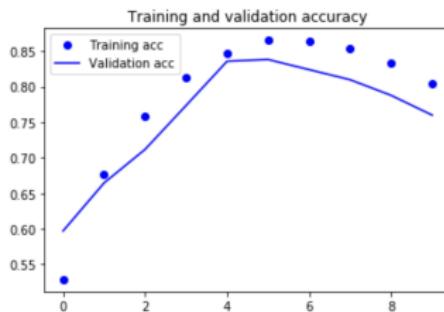
- In Keras, use the Conv1D layer
  - Takes as input (samples, time, features)
  - Returns 3D tensors
- Combine Conv1D layer with a MaxPooling1D layer
- End with a GlobalMaxPooling or Flatten layer
- Can use larger windows for 1D CNNs - typically windows of size 7 or 9
  - In a 2D convolution layer, a  $3 \times 3$  filter contains 9 feature vectors
  - A 1D convolution layer with a window of size 3 contains only 3 vectors

# 1D CNN IMDB Example

```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.Embedding(max_features, 128, input_length=max_len))
7 model.add(layers.Conv1D(32, 7, activation='relu'))
8 model.add(layers.MaxPooling1D(5))
9 model.add(layers.Conv1D(32, 7, activation='relu'))
10 model.add(layers.GlobalMaxPooling1D())
11 model.add(layers.Dense(1))
12
13 model.summary()
14
15 model.compile(optimizer=RMSprop(lr=1e-4),
16                 loss='binary_crossentropy',
17                 metrics=['acc'])
18 history = model.fit(x_train, y_train,
19                       epochs=10,
20                       batch_size=128,
21                       validation_split=0.2)
```

# 1D CNN IMDB Example

- The validation accuracy is a little lower than when using an RNN with an LSTM layer, but the running time for this model is much lower

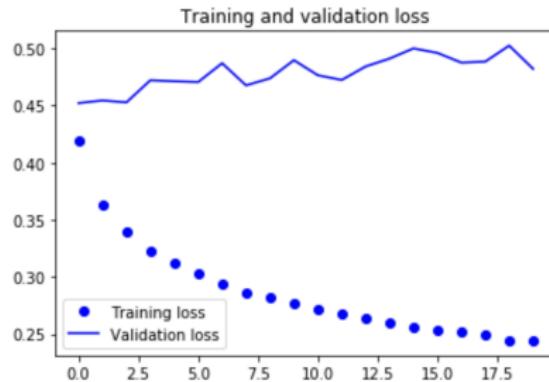


# Stacking 1D CNN Layers

- 1D CNNs process subsequences independently and aren't sensitive to the order of the timesteps - thus, they don't perform well when faced with long sequences
- Could try stacking 1D CNN layers, but still doesn't induce order sensitivity

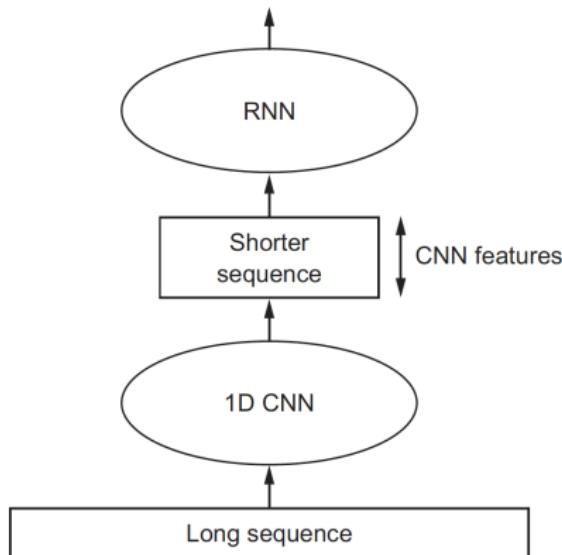
```
1 from keras.models import Sequential
2 from keras import layers
3 from keras.optimizers import RMSprop
4
5 model = Sequential()
6 model.add(layers.Conv1D(32, 5, activation='relu',
7                      input_shape=(None, float_data.shape[-1])))
8 model.add(layers.MaxPooling1D(3))
9 model.add(layers.Conv1D(32, 5, activation='relu'))
10 model.add(layers.MaxPooling1D(3))
11 model.add(layers.Conv1D(32, 5, activation='relu'))
12 model.add(layers.GlobalMaxPooling1D())
13 model.add(layers.Dense(1))
14
15 model.compile(optimizer=RMSprop(), loss='mae')
16 history = model.fit_generator(train_gen,
17                               steps_per_epoch=500,
18                               epochs=20,
19                               validation_data=val_gen,
20                               validation_steps=val_steps)
```

- Doesn't beat the common-sense baseline
  - Due to the fact that 1D CNNs aren't time sensitive
- Let's try combining CNNs and RNNs



# Combining CNNs and RNNs

- To combine the speed of 1D CNNs with the time sensitivity of RNNs, could do the following:
  - Use the 1D CNN as a preprocessing step
  - Feed this output into an RNN
- The CNN turns long sequences into much shorter sequences



# Temperature Forecasting with CNN and RNN

- With the combination of a CNN and RNN we can now either look at data from longer ago (increase the `lookback` parameter), or look at high-resolution timeseries (decrease the `step` parameter)
- Let's decrease the `step` parameter by half - this gives us sequences that are twice as long
- Temperature data is now sampled at a rate of 1 point per 30 minutes

```
1 # This was previously set to 6 (one point per hour).
2 # Now 3 (one point per 30 min).
3 step = 3
4 lookback = 720 # Unchanged
5 delay = 144 # Unchanged
6
7 train_gen = generator(float_data,
8                     lookback=lookback,
9                     delay=delay,
10                    min_index=0,
11                    max_index=200000,
12                    shuffle=True,
13                    step=step)
14 val_gen = generator(float_data,
15                     lookback=lookback,
16                     delay=delay,
17                     min_index=200001,
18                     max_index=300000,
19                     step=step)
```

# Temperature Forecasting with CNN and RNN

```
1 test_gen = generator(float_data,
2                      lookback=lookback,
3                      delay=delay,
4                      min_index=300001,
5                      max_index=None,
6                      step=step)
7 val_steps = (300000 - 200001 - lookback) // 128
8 test_steps = (len(float_data) - 300001 - lookback) // 128
9
10 model = Sequential()
11 model.add(layers.Conv1D(32, 5, activation='relu',
12                       input_shape=(None, float_data.shape[-1])))
13 model.add(layers.MaxPooling1D(3))
14 model.add(layers.Conv1D(32, 5, activation='relu'))
15 model.add(layers.GRU(32, dropout=0.1, recurrent_dropout=0.5))
16 model.add(layers.Dense(1))
17
18 model.summary()
19
20 model.compile(optimizer=RMSprop(), loss='mae')
21 history = model.fit_generator(train_gen,
22                               steps_per_epoch=500,
23                               epochs=20,
24                               validation_data=val_gen,
25                               validation_steps=val_steps)
```

# Temperature Forecasting with CNN and RNN

- This model doesn't perform as well as the regularized GRU model, but it does run a lot faster
- This model also looks at twice as much data as the other model, but that doesn't seem to help with accuracy in this case - it could for other data sets

