

Part 1: Data from a mysterious celltype

Your collaborators performed ChIP-seq for two different factors in a mysterious new cell type.

They send you a file with DNA sequences and the matched profiles for a set of regions.

They believe there is an interaction between the two factors, but would like your help modeling and analyzing the data.

Try to answer their questions concisely.

```
In [1]: # Import dependencies
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import requests
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.model_selection import train_test_split
import io
import h5py
import tensorflow as tf
from tensorflow.keras.optimizers import SGD, RMSprop
from tensorflow.keras.layers import Conv1D, Dense, MaxPooling1D, Flatten,
from tensorflow.keras.models import Sequential
from tensorflow.keras.callbacks import EarlyStopping
import tensorflow.keras.backend as K
from tensorflow.keras.metrics import Accuracy
from tensorflow import argmax
from tensorflow import GradientTape
from tensorflow import Variable

# Below is required for saliency mapping
from tensorflow.python.framework.ops import disable_eager_execution
disable_eager_execution()

In [2]: import random
random.seed(123)

In [3]: def download_h5(dataset_url, filename = 'dataset.h5'):
    with requests.Session() as session:
        r = session.get(dataset_url, stream=True)
        r.raise_for_status()
        with open(filename, 'wb') as h5f:
            for chunk in r.iter_content(chunk_size=io.DEFAULT_BUFFER_SIZE):
                h5f.write(chunk)

In [4]: DATA_URL = "https://raw.githubusercontent.com/FuDenberg-Research-Group/20
download_h5(DATA_URL, filename = 'dataset1.h5')

In [5]: ## reading in the h5 data
with h5py.File('dataset1.h5', 'r') as h:
    sequences = h['segs'].asstr()[0:100]
    targets = h['targets'][0:100]

In [6]: ## To take a quick look at one of their profiles
plt.plot(targets[0,:], label='ChIP-protein1');
plt.legend()
plt.ylabel('signal')
plt.xlabel('position, bp')

Out[6]: Text(0.5, 0, 'position, bp')
```

1.0 Try to fit the data with a seq-to-profile neural network.

- Follow similar steps to those taken in the [Zou et al primer](#), e.g.:
 - one-hot encode the DNA,
 - set aside a test set,
 - define a convolutional architecture, perhaps with additional or modified layers
 - use some sort of loss function and activation to use.
 - Use [EarlyStopping](#)
 - Report your test set accuracy. Report an interpretable metric as well.
 - Where does the model do well or struggle?

Note: if the validation loss is not decreasing during training, the specified model may not be well-matched for the problem at hand. Consider how to modify the model (e.g. a model with sigmoid activation at the final layer would never be able to learn negative values).

```
In [7]: ## One-hot encode of DNA sequence
integer_encoder = LabelEncoder()
one_hot_encoder = OneHotEncoder(categories='auto')
input_features = []

for sequence in sequences:
    integer_encoded = integer_encoder.fit_transform(list(sequence))
    integer_encoded = np.array(integer_encoded).reshape(-1, 1)
    one_hot_encoded = one_hot_encoder.fit_transform(integer_encoded)
    input_features.append(one_hot_encoded.toarray())

input_features = np.stack(input_features)
input_labels = targets # signal values are labels
maxsignal = input_labels.max()
minsignal = input_labels.min()

# Divide the data: Create test and train data
from sklearn.model_selection import train_test_split
train_features, test_features, train_labels, test_labels = train_test_split(
    input_features, input_labels, test_size=0.25, random_state=42)

#print(train_features.shape) # (3750, 100, 4)
#print(train_labels.shape) # (3750, 100, 2)

In [8]: # Pearson for metric
from keras import backend as K
def pearson(r, y_true, y_pred):
    # pearson correlation coefficient
    # source: https://github.com/WenYanger/Keras_Metrics
    # provided after office hours
    epsilon = 1e-5
    x = y_true
    y = y_pred
    mx = K.mean(x)
    my = K.mean(y)
    xm = x - mx
    ym = y - my
    r_num = K.sum(xm * ym)
    x_square_sum = K.sum(xm * xm)
    y_square_sum = K.sum(ym * ym)
    r_den = K.sqrt(x_square_sum * y_square_sum)
    r = r_num / (r_den + epsilon)
    return K.mean(r)

In [9]: # Setup the model
model = Sequential() #sequential architecture
model.add(Conv1D(filters=32, kernel_size=16,
                  input_shape=(train_features.shape[1], 4), padding='same'
))
model.add(MaxPooling1D(pool_size=12, strides=1, padding='same'))
model.add(Dense(300, activation='relu'))
model.add(Dense(2))

model.compile(loss='mse', optimizer='adam', metrics=[pearson_r])
model.summary()
```

Layer (type)	Output Shape	Param #
conv1d (Conv1D)	(None, 100, 32)	2080
max_pooling1d (MaxPooling1D)	(None, 100, 32)	0
dense_1 (Dense)	(None, 100, 300)	9900
dense_1_1 (Dense)	(None, 100, 2)	602

Total params: 12,582
Trainable params: 12,582
Non-trainable params: 0

```
In [10]: # Define callback function to prevent overfitting
callback = EarlyStopping(monitor='loss', patience=3)

# Fit the model
model.fit(train_features, train_labels, epochs=100, verbose=0,
          validation_data=(test_features, test_labels), callbacks=[callback])

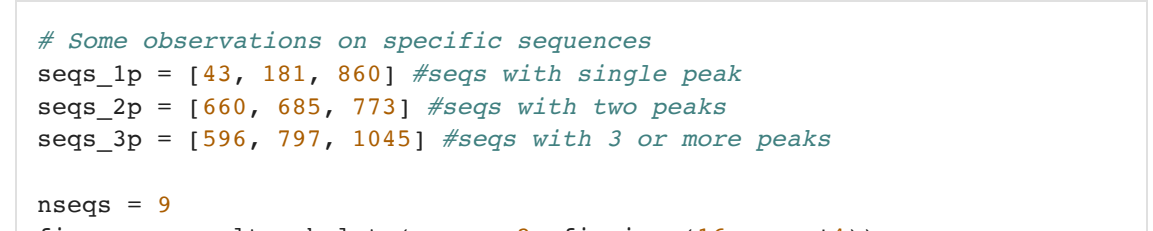
updates = self.state_updates
```

```
In [11]: # Figure 1-A: Epoch vs Loss
fig1, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))

ax1.plot(history.history['loss'])
ax1.plot(history.history['val_loss'])
ax1.set_title('Mean Squared Error')
ax1.set_ylabel('Loss (mse)', xlabel='Epoch')
ax1.legend(['Train Set', 'Validation Set'])

# Figure 1-B: Epoch vs Loss
ax2.plot(history.history['pearson_r'])
ax2.plot(history.history['val_pearson_r'])
ax2.set_title('Pearson Correlation')
ax2.set_ylabel('r', xlabel='epoch')
ax2.legend(['Train Set', 'Validation Set'])

fig.tight_layout(pad=2.0)
fig1.show()
```



```
In [12]: # Evaluate the models accuracy
score_train = model.evaluate(train_features, train_labels, verbose = True)
score_test = model.evaluate(test_features, test_labels, verbose = True)
```

```
In [13]: # Use test set for predictions
test_prediction = model.predict(test_features, verbose = True)
```

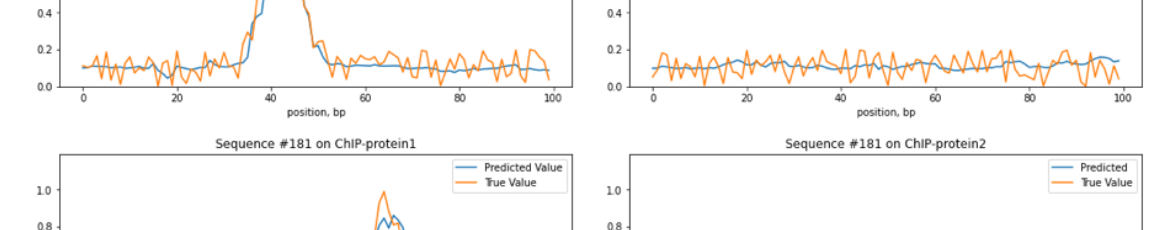
/usr/local/lib/python3.8/dist-packages/keras/engine/training_v1.py:2045: UserWarning: Model.state_updates will be removed in a future version. This property should not be used in TensorFlow 2.0, as 'updates' are applied automatically.
updates=self.state_updates

```
In [14]: # Investigation of the prediction on randomly selected 2 sequences
nsegs = 2
fig, axes = plt.subplots(nsegs, 2, figsize=(16,nsegs*4))
for i in range(nsegs):
    seq = random.randint(0,1250)

    axes[i][0].plot(test_prediction[seq,:])
    axes[i][0].plot(test_labels[seq,:])
    axes[i][0].set_ylabel('signal', xlabel='position, bp')
    axes[i][0].set_title('ChIP-protein1 vs. Sequence #' + str(seq))
    axes[i][0].legend(['Predicted Value', 'True Value'])

    axes[i][1].plot(test_prediction[seq,:])
    axes[i][1].plot(test_labels[seq,:])
    axes[i][1].set_ylabel('signal', xlabel='position, bp')
    axes[i][1].set_title('ChIP-protein2 vs. Sequence #' + str(seq))
    axes[i][1].set_title('signal', xlabel='position, bp')
    axes[i][1].legend(['Predicted', 'True Value'])

fig.tight_layout(pad=2.0)
fig.show()
```



```
In [15]: # Some observations on specific sequences
seqs_1p = [43, 181, 860] #segs with single peak
seqs_2p = [560, 685, 773] #segs with two peaks
seqs_3p = [596, 797, 1045] #segs with 3 or more peaks

nsegs = 9
fig, axes = plt.subplots(nsegs, 2, figsize=(16,nsegs*4))
for i in range(nsegs):
    if i in range(0,3): seq = seqs_1p[i]
    if i in range(3,6): seq = seqs_2p[i-3]
    if i in range(6,9): seq = seqs_3p[i-6]

    axes[i][0].plot(test_prediction[seq,:])
    axes[i][0].set_ylabel('signal', xlabel='position, bp')
    axes[i][0].set_title('ChIP-protein1 vs. Sequence #' + str(seq))
    axes[i][0].set_title('signal', xlabel='position, bp')
    axes[i][0].legend(['Predicted Value', 'True Value'])

    axes[i][1].plot(test_prediction[seq,:])
    axes[i][1].plot(test_labels[seq,:])
    axes[i][1].set_ylabel('signal', xlabel='position, bp')
    axes[i][1].set_title('ChIP-protein2 vs. Sequence #' + str(seq))
    axes[i][1].set_title('signal', xlabel='position, bp')
    axes[i][1].legend(['Predicted', 'True Value'])

fig.tight_layout(pad=2.0)
fig.show()
```

