

pca_open_office

January 18, 2022

1 Principal Component Analysis

Below are the Principal Component Analysis on the first 7 features:

1. index: Index for each data point
2. sim: Index for each simulation
3. tem: Temperature of air supplied to the zone (°C)
4. mass: Total mass flow rate of air supply to the zone (kg/s)
5. n_occ: Number of occupants in zone
6. x: location of occupant of interest on x-axis (m)
7. y: location of occupant of interest on y-axis (m)

The last 3 features are treated as target predictions for evaluation purpose.

1. MRT: Mean radiant temperature of the occupant of interest's head and chest (°C)
2. T: Average temperature of air surrounding the occupant of interest (°C)
3. V: Average speed of air surrounding the occupant of interest occupant (°C)

```
[1]: data_loc = '../data/preprocessed/'
```

1.1 Open Office

1.1.1 1. Standardize the data

```
[2]: import pandas as pd
from sklearn.preprocessing import StandardScaler

open_office_data = pd.read_csv(data_loc+"open_office.csv")
open_office_data.isna().sum()
```

```
[2]: index      0
     sim        0
     tem        0
     mass        0
     n          0
     occ        0
     x          0
     y          0
     MRT        0
     T          0
```

```
V      8
dtype: int64
```

This dataset seems to have null values.

```
[3]: # Remove the index column
open_office_data.drop(columns=['index', 'sim'], inplace=True)

# Remove null data
open_office_data = open_office_data.dropna()

# Seperate the data into features and labels
X = open_office_data.drop(columns=['MRT', 'T', 'V'])
y_MRT = open_office_data['MRT']
y_T = open_office_data['T']
y_V = open_office_data['V']

# Standardize the data
scaler = StandardScaler()
scaler.fit(X)

X_std = pd.DataFrame(scaler.transform(X), columns=X.columns)
X_std.describe()
```

```
[3]:
```

	tem	mass	n	occ	x \
count	4.871000e+03	4.871000e+03	4.871000e+03	4.871000e+03	4.871000e+03
mean	7.243687e-16	-7.691104e-16	5.969813e-16	-2.552031e-15	1.479325e-15
std	1.000103e+00	1.000103e+00	1.000103e+00	1.000103e+00	1.000103e+00
min	-1.870995e+00	-1.405293e+00	-1.567282e+00	-1.414649e+00	-1.415396e+00
25%	-7.598156e-01	-8.609269e-01	-7.994826e-01	-1.414649e+00	-7.033713e-01
50%	7.356918e-02	-9.881505e-02	-3.168297e-02	7.068890e-01	6.777462e-02
75%	6.291590e-01	8.810431e-01	7.361167e-01	7.068890e-01	8.363501e-01
max	2.573724e+00	1.860901e+00	1.503916e+00	7.068890e-01	1.532952e+00

	y
count	4.871000e+03
mean	2.411744e-15
std	1.000103e+00
min	-1.117982e+00
25%	-1.117982e+00
50%	8.944685e-01
75%	8.944685e-01
max	8.944685e-01

1.1.2 2. Perform PCA

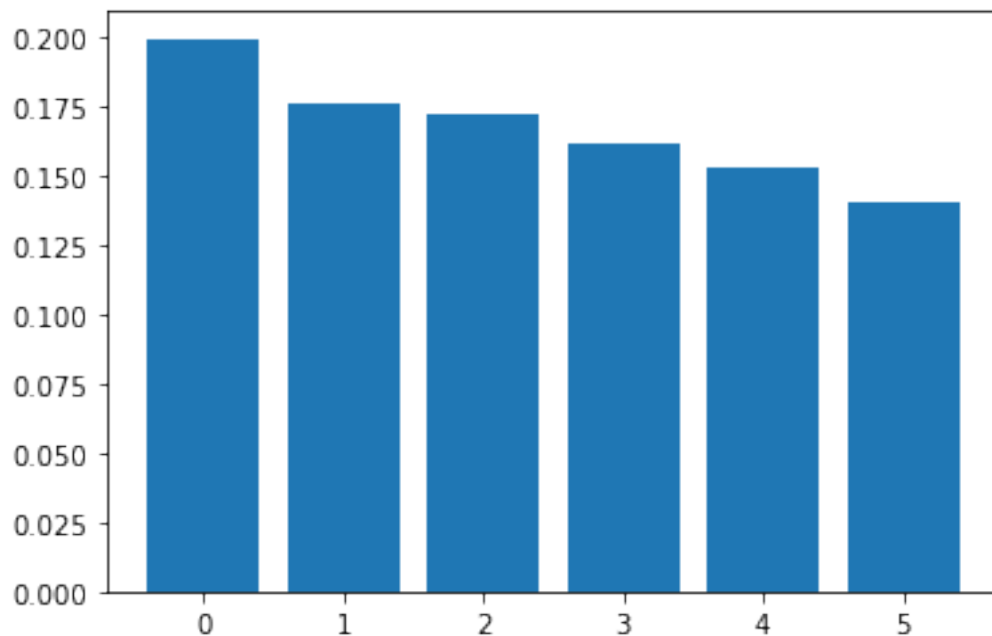
```
[4]: # Perform PCA
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import numpy as np
from numpy import savetxt

pca = PCA()
pca.fit(X_std)

e_vectors = pca.components_ # The eigenvectors
evr = pca.explained_variance_ratio_ # The variance explained by each eigenvector

# Save the eigenvectors to a text file
savetxt("../reports/data/pca_open_office_eigenvectors.csv", e_vectors,
        ↪delimiter=',')

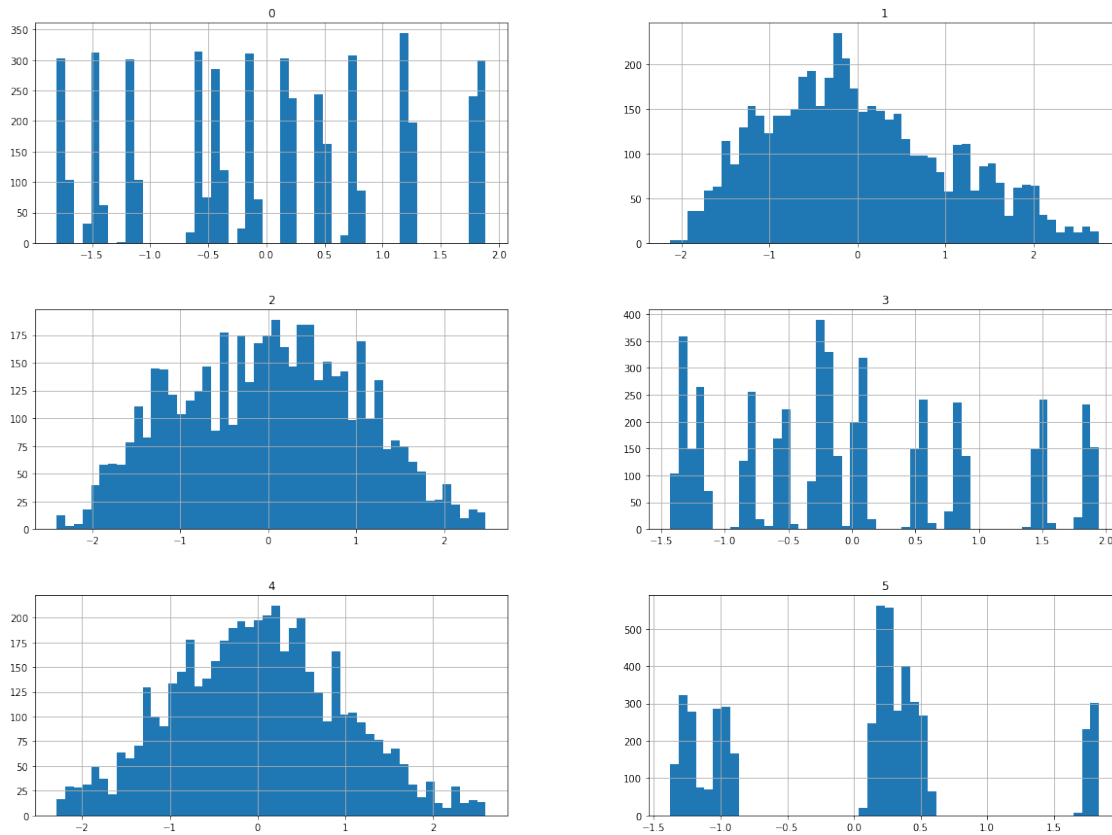
plt.bar(range(len(evr)), evr)
plt.savefig('../figures/evr_pca_open_office.png')
```



```
[5]: # Transform the data
X_pca = pd.DataFrame(pca.transform(X_std))
fig, ax = plt.subplots(figsize=(20,15))
X_pca.hist(bins=50, ax=ax)
fig.savefig('../figures/hist_pca_open_office.png')
```

/tmp/ipykernel_21455/4076575879.py:4: UserWarning: To output multiple subplots, the figure containing the passed axes is being cleared

```
X_pca.hist(bins=50, ax=ax)
```



After performing PCA on the features, the distribution graph shows more normal distribution features (8 features).

1.1.3 3. Compare models' results with and without PCA on predicting MRT

3.1 Split the data into training and testing

```
[6]: # Split the data into training and testing
from sklearn.model_selection import train_test_split

# Without PCA
X_train, X_test, y_train, y_test = train_test_split(X_std, y_MRT, test_size=0.
↪2, random_state=42)

# With PCA
X_train_pca, X_test_pca, y_train_pca, y_test_pca = train_test_split(X_pca,
↪y_MRT, test_size=0.2, random_state=42)
```

3.2 Train and evaluate the Linear Regression model

3.2.1 Without PCA

```
[7]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import pickle
import time

# Train linear regression model on training data
model = LinearRegression()
start = time.time()
model.fit(X_train, y_train)
stop = time.time()
lr_train_time = stop - start
print(f"Training time: {lr_train_time}s")

# Save the model to a pickle file
filename = '../reports/models/lr_MRT_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

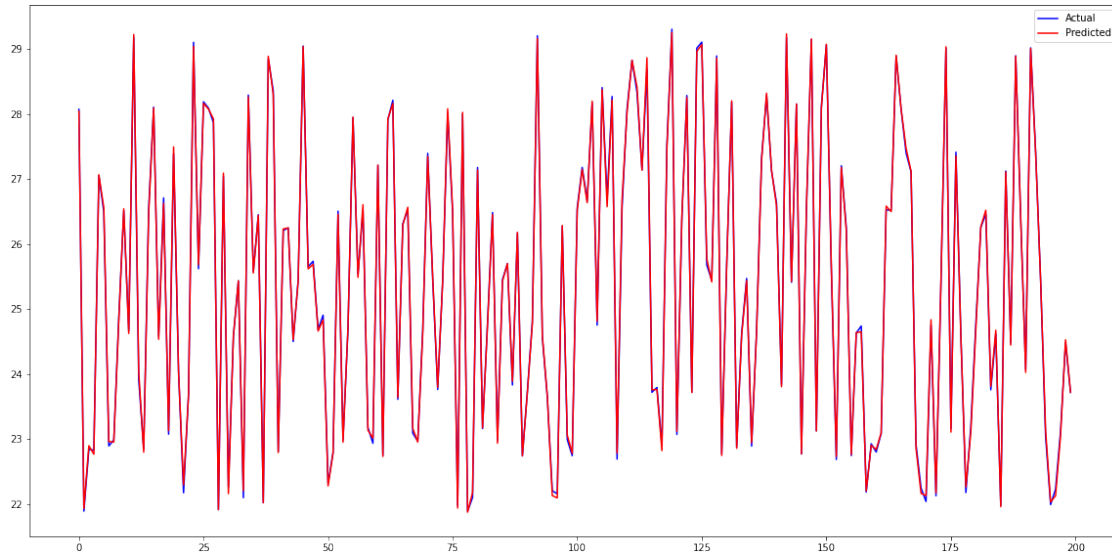
# Predict on test data
y_pred = model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

Training time: 0.005346536636352539s

Mean squared error: 0.002502969797609541



3.2.2 With PCA

```
[8]: # Train linear regression model on training data
start = time.time()
model.fit(X_train_pca, y_train_pca)
stop = time.time()
lr_pca_train_time = stop - start
print(f"Training time: {lr_pca_train_time}s")

# Predict on test data
y_pred_pca = model.predict(X_test_pca)

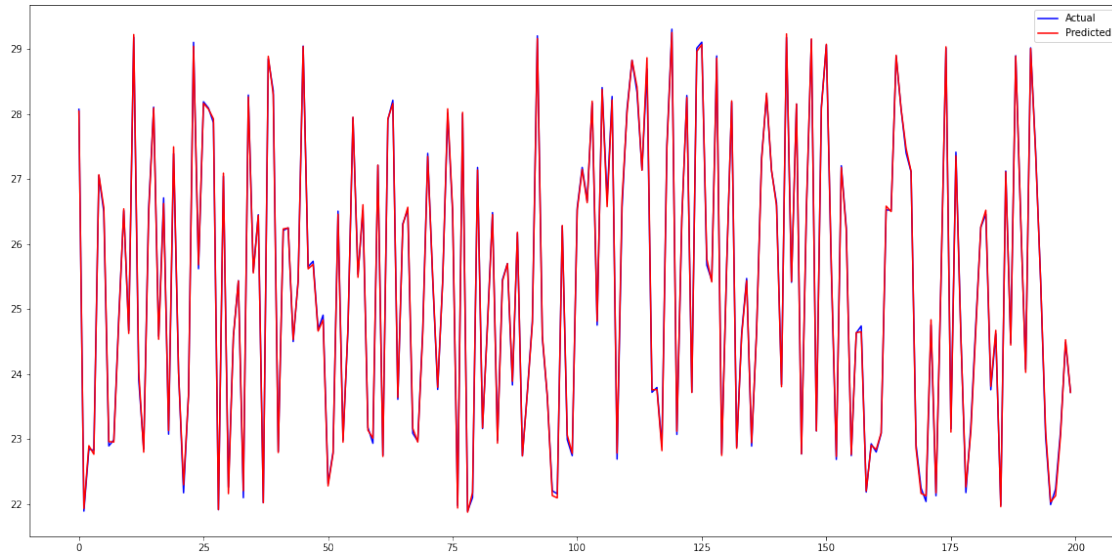
# Save the model to a pickle file
filename = '../reports/models/lr_MRT_pca_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

# Evaluate the model
mse = mean_squared_error(y_test_pca, y_pred_pca)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

Training time: 0.007262229919433594s

Mean squared error: 0.0025029697976095343



3.3 Train and evaluate the Feedforward Neural Network model

3.2.1 Without PCA

```
[9]: # Create neural network model
from keras.models import Sequential
from keras.layers import Dense

N_NEURONS = 1024
N_LAYERS = 4

model = Sequential()
model.add(Dense(units=N_NEURONS, input_dim=X.shape[1], activation='relu'))
for i in range(N_LAYERS-1):
    model.add(Dense(units=N_NEURONS, activation='relu'))
model.add(Dense(units=1, activation='linear')) # Output layer
model.compile(loss='mean_squared_error', optimizer='adam')

# Train the model
start = time.time()
model.fit(X_train, y_train, epochs=100, verbose=0)
stop = time.time()
nn_train_time = stop - start
print(f"Training time: {nn_train_time}s")

# Predict on test data
y_pred_nn = model.predict(X_test)
model.save('../reports/models/nn_MRT_open_office_model.pkl')
```

```

# Evaluate the model
mse = mean_squared_error(y_test, y_pred_nn)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

```

2022-01-18 18:10:58.452571: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcudart.so.10.1
2022-01-18 18:11:00.235584: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcuda.so.1
2022-01-18 18:11:00.277612: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.278207: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1716] Found device 0 with
properties:
pciBusID: 0000:04:00.0 name: GeForce GTX 1060 6GB computeCapability: 6.1
coreClock: 1.7845GHz coreCount: 10 deviceMemorySize: 5.93GiB
deviceMemoryBandwidth: 178.99GiB/s
2022-01-18 18:11:00.278261: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcudart.so.10.1
2022-01-18 18:11:00.280817: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcublas.so.10
2022-01-18 18:11:00.282432: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcufft.so.10
2022-01-18 18:11:00.282790: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcurand.so.10
2022-01-18 18:11:00.286776: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcusolver.so.10
2022-01-18 18:11:00.288329: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcusparsparse.so.10
2022-01-18 18:11:00.293957: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully

```



```

opened dynamic library libcudnn.so.7
2022-01-18 18:11:00.294147: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.294814: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.295351: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1858] Adding visible gpu
devices: 0
2022-01-18 18:11:00.309030: I
tensorflow/core/platform/profile_utils/cpu_utils.cc:104] CPU Frequency:
2194710000 Hz
2022-01-18 18:11:00.311605: I tensorflow/compiler/xla/service/service.cc:168]
XLA service 0x563e63c88e40 initialized for platform Host (this does not
guarantee that XLA will be used). Devices:
2022-01-18 18:11:00.311652: I tensorflow/compiler/xla/service/service.cc:176]
StreamExecutor device (0): Host, Default Version
2022-01-18 18:11:00.557467: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.558124: I tensorflow/compiler/xla/service/service.cc:168]
XLA service 0x563e63cf48f0 initialized for platform CUDA (this does not
guarantee that XLA will be used). Devices:
2022-01-18 18:11:00.558147: I tensorflow/compiler/xla/service/service.cc:176]
StreamExecutor device (0): GeForce GTX 1060 6GB, Compute Capability 6.1
2022-01-18 18:11:00.558401: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.558958: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1716] Found device 0 with
properties:
pciBusID: 0000:04:00.0 name: GeForce GTX 1060 6GB computeCapability: 6.1
coreClock: 1.7845GHz coreCount: 10 deviceMemorySize: 5.93GiB
deviceMemoryBandwidth: 178.99GiB/s
2022-01-18 18:11:00.559028: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcudart.so.10.1
2022-01-18 18:11:00.559071: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcublas.so.10
2022-01-18 18:11:00.559110: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcufft.so.10

```

```

2022-01-18 18:11:00.559149: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcurand.so.10
2022-01-18 18:11:00.559191: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcusolver.so.10
2022-01-18 18:11:00.559230: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcusparsesparse.so.10
2022-01-18 18:11:00.559270: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcudnn.so.7
2022-01-18 18:11:00.559379: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.559975: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:00.560485: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1858] Adding visible gpu
devices: 0
2022-01-18 18:11:00.560546: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully
opened dynamic library libcudart.so.10.1
2022-01-18 18:11:01.236847: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1257] Device interconnect
StreamExecutor with strength 1 edge matrix:
2022-01-18 18:11:01.236888: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1263]          0
2022-01-18 18:11:01.236898: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1276] 0:      N
2022-01-18 18:11:01.237188: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:01.237710: I
tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:982] successful NUMA node
read from SysFS had negative value (-1), but there must be at least one NUMA
node, so returning NUMA node zero
2022-01-18 18:11:01.238129: I
tensorflow/core/common_runtime/gpu/gpu_device.cc:1402] Created TensorFlow device
(/job:localhost/replica:0/task:0/device:GPU:0 with 4679 MB memory) -> physical
GPU (device: 0, name: GeForce GTX 1060 6GB, pci bus id: 0000:04:00.0, compute
capability: 6.1)
2022-01-18 18:11:01.918703: I
tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully

```

opened dynamic library libcublas.so.10

Training time: 43.761592864990234s

WARNING:tensorflow:From /home/khiem/anaconda3/lib/python3.8/site-packages/tensorflow/python/training/tracking/tracking.py:111:

Model.state_updates (from tensorflow.python.keras.engine.training) is deprecated and will be removed in a future version.

Instructions for updating:

This property should not be used in TensorFlow 2.0, as updates are applied automatically.

WARNING:tensorflow:From /home/khiem/anaconda3/lib/python3.8/site-

packages/tensorflow/python/training/tracking/tracking.py:111: Layer.updates (from tensorflow.python.keras.engine.base_layer) is deprecated and will be removed in a future version.

Instructions for updating:

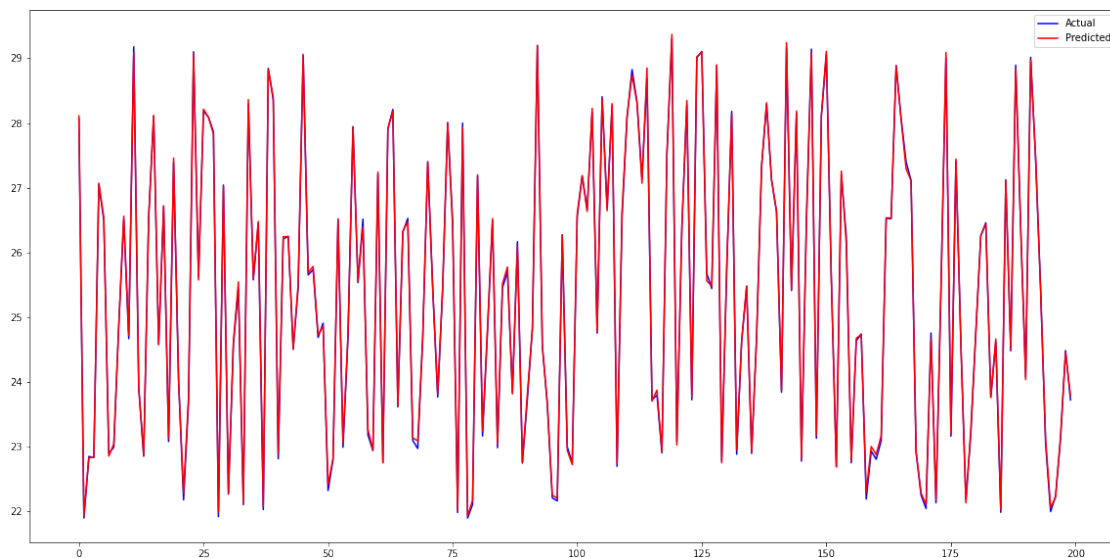
This property should not be used in TensorFlow 2.0, as updates are applied automatically.

2022-01-18 18:11:45.662754: W tensorflow/python/util/util.cc:348] Sets are not currently considered sequences, but this may change in the future, so consider avoiding using them.

INFO:tensorflow:Assets written to:

../reports/models/nn_MRT_open_office_model.pkl/assets

Mean squared error: 0.002490786140369075



3.3.2 With PCA

```
[10]: # Train the model
start = time.time()
model.fit(X_train_pca, y_train_pca, epochs=100, verbose=0)
```

```

stop = time.time()
nn_pca_train_time = stop - start
print(f"Training time: {nn_pca_train_time}s")

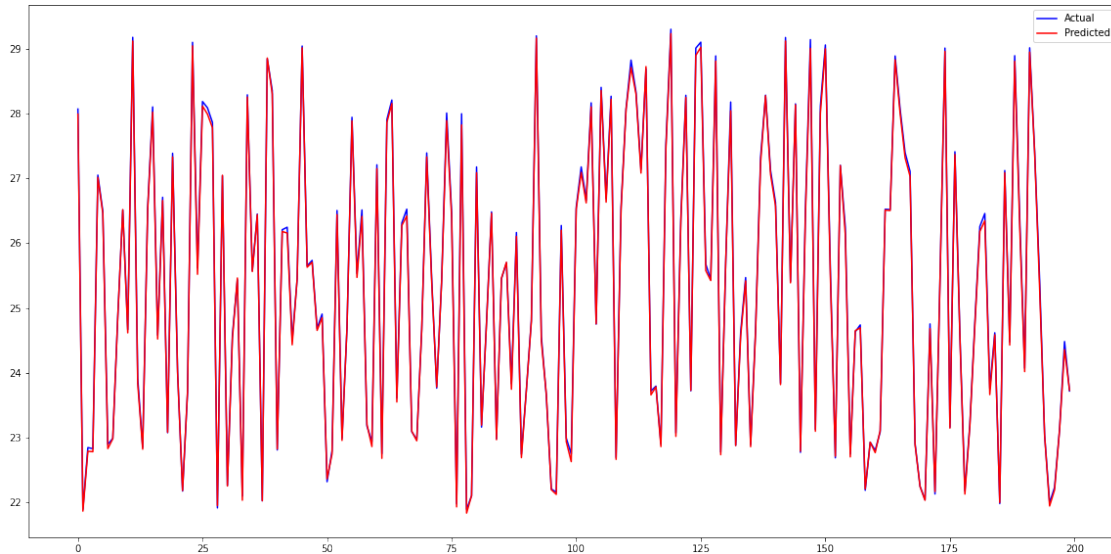
# Predict on test data
y_pred_nn_pca = model.predict(X_test_pca)
model.save(' ../reports/models/nn_MRT_pca_open_office_model.pkl')

# Evaluate the model
mse = mean_squared_error(y_test_pca, y_pred_nn_pca)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

Training time: 44.0262987613678s
 INFO:tensorflow:Assets written to:
 ../reports/models/nn_MRT_pca_open_office_model.pkl/assets
 Mean squared error: 0.0027803035776146155



```

[11]: # Save training times to a csv file
import csv

with open(' ../reports/data/training_times_open_office_MRT.csv', 'w') as csvfile:

```

```

writer = csv.writer(csvfile)
writer.writerow(['Model', 'Training time (s)'])
writer.writerow(['Linear regression', lr_train_time])
writer.writerow(['Linear regression with PCA', lr_pca_train_time])
writer.writerow(['Neural network', nn_train_time])
writer.writerow(['Neural network with PCA', nn_pca_train_time])

```

1.1.4 4. Compare models' results with and without PCA on predicting T

4.1 Split the data into training and testing

```

[12]: # Split the data into training and testing
from sklearn.model_selection import train_test_split

# Without PCA
X_train, X_test, y_train, y_test = train_test_split(X_std, y_T, test_size=0.2,
    random_state=42)

# With PCA
X_train_pca, X_test_pca, y_train_pca, y_test_pca = train_test_split(X_pca, y_T,
    test_size=0.2, random_state=42)

```

4.2 Train and evaluate the Linear Regression model

4.2.1 Without PCA

```

[13]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
import pickle
import time

# Train linear regression model on training data
model = LinearRegression()
start = time.time()
model.fit(X_train, y_train)
stop = time.time()
lr_train_time = stop - start
print(f"Training time: {lr_train_time}s")

# Save the model to a pickle file
filename = '../reports/models/lr_T_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

# Predict on test data
y_pred = model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)

```

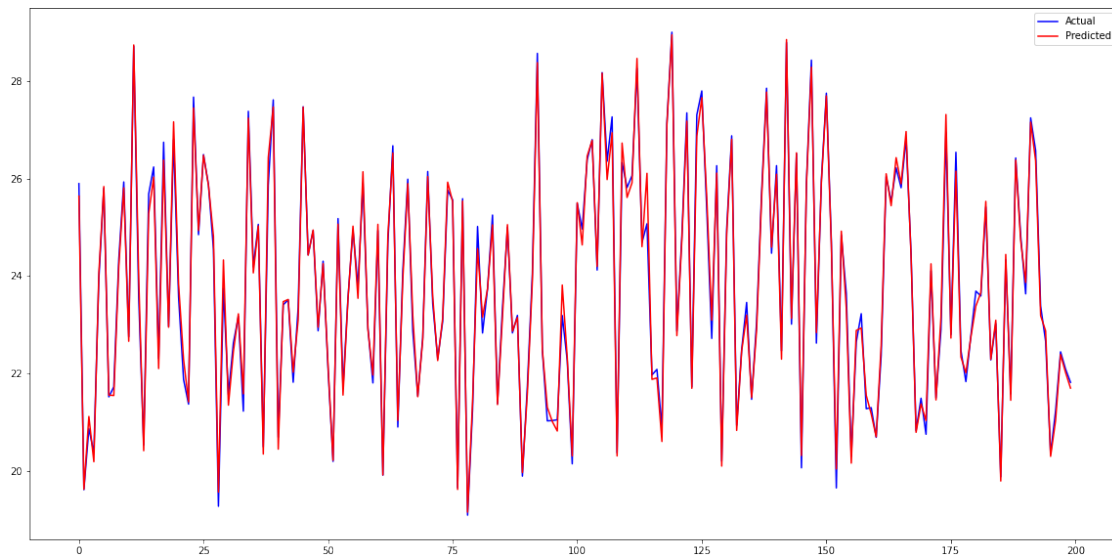
```

print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

Training time: 0.006529092788696289s
Mean squared error: 0.04119175262435056



4.2.2 With PCA

```

[14]: # Train linear regression model on training data
start = time.time()
model.fit(X_train_pca, y_train_pca)
stop = time.time()
lr_pca_train_time = stop - start
print(f"Training time: {lr_pca_train_time}s")

# Predict on test data
y_pred_pca = model.predict(X_test_pca)

# Save the model to a pickle file
filename = '../reports/models/lr_T_pca_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

# Evaluate the model

```

```

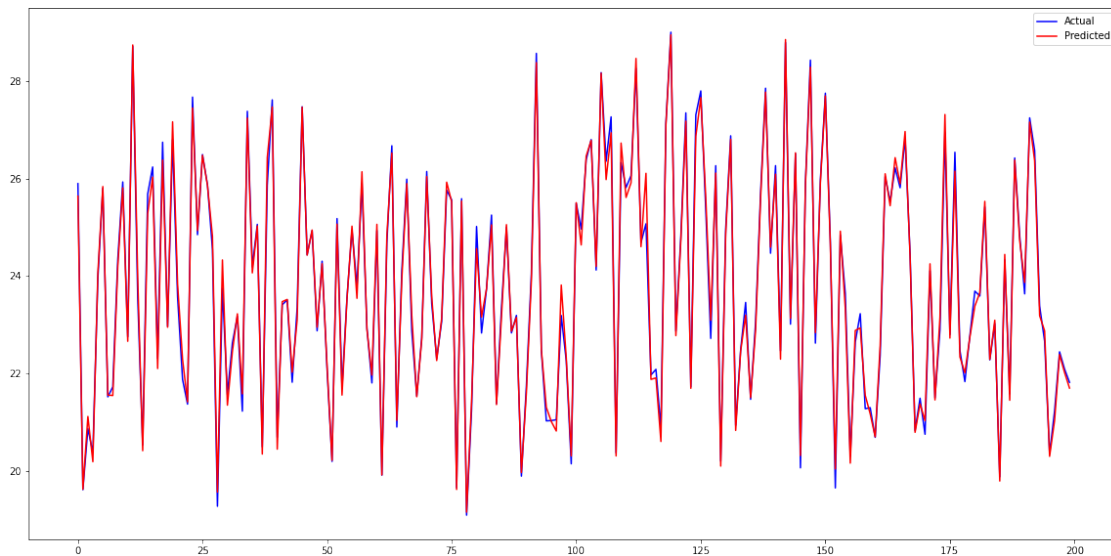
mse = mean_squared_error(y_test_pca, y_pred_pca)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

Training time: 0.005934476852416992s

Mean squared error: 0.041191752624350535



4.3 Train and evaluate the Feedforward Neural Network model

4.2.1 Without PCA

```

[15]: # Create neural network model
from keras.models import Sequential
from keras.layers import Dense

N_NEURONS = 1024
N_LAYERS = 4

model = Sequential()
model.add(Dense(units=N_NEURONS, input_dim=X.shape[1], activation='relu'))
for i in range(N_LAYERS-1):
    model.add(Dense(units=N_NEURONS, activation='relu'))
model.add(Dense(units=1, activation='linear')) # Output layer

```

```

model.compile(loss='mean_squared_error', optimizer='adam')

# Train the model
start = time.time()
model.fit(X_train, y_train, epochs=100, verbose=0)
stop = time.time()
nn_train_time = stop - start
print(f"Training time: {nn_train_time}s")

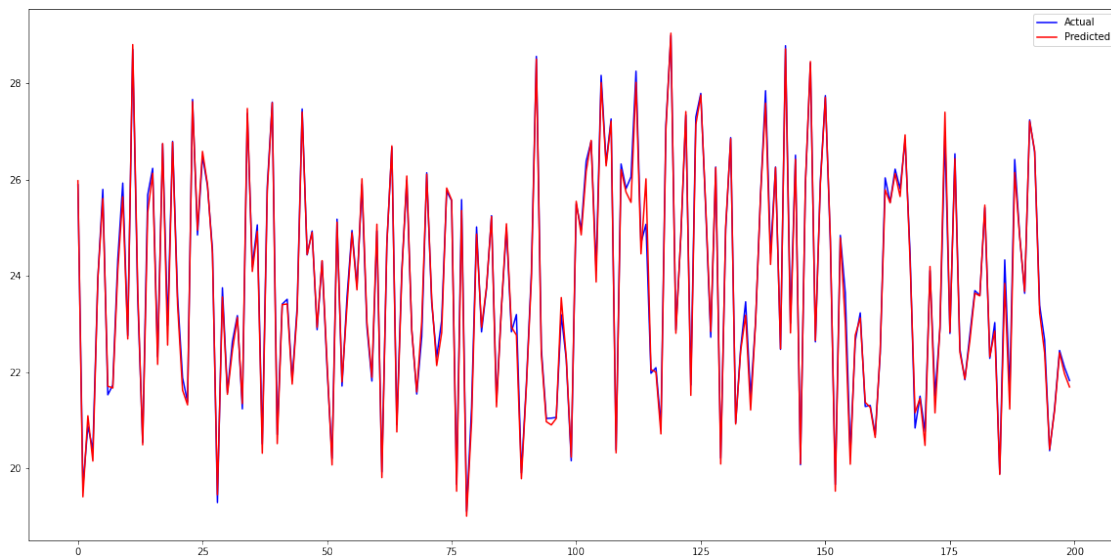
# Predict on test data
y_pred_nn = model.predict(X_test)
model.save('../reports/models/nn_T_open_office_model.pkl')

# Evaluate the model
mse = mean_squared_error(y_test, y_pred_nn)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

Training time: 45.3539400100708s
 INFO:tensorflow:Assets written to:
 ../reports/models/nn_T_open_office_model.pkl/assets
 Mean squared error: 0.026644152618731248



4.3.2 With PCA

```
[16]: # Train the model
start = time.time()
model.fit(X_train_pca, y_train_pca, epochs=100, verbose=0)
stop = time.time()
nn_pca_train_time = stop - start
print(f"Training time: {nn_pca_train_time}s")

# Predict on test data
y_pred_nn_pca = model.predict(X_test_pca)
model.save(' ../reports/models/nn_T_pca_open_office_model.pkl')

# Evaluate the model
mse = mean_squared_error(y_test_pca, y_pred_nn_pca)
print(f'Mean squared error: {mse}')

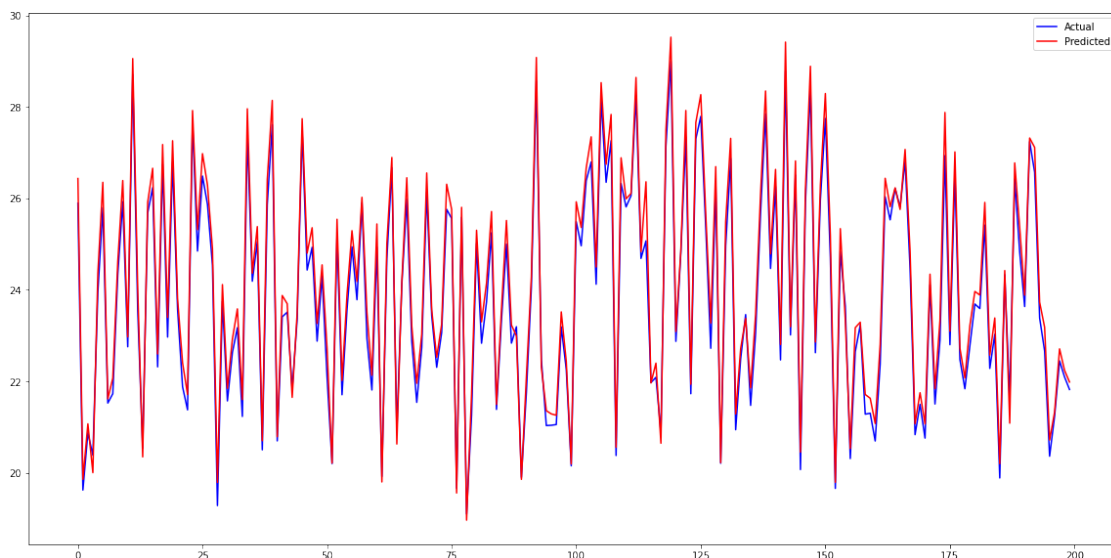
# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

Training time: 46.28824257850647s

INFO:tensorflow:Assets written to:

../reports/models/nn_T_pca_open_office_model.pkl/assets

Mean squared error: 0.13827998860915247



```
[17]: # Save training times to a csv file
import csv

with open('../reports/data/training_times_open_office_T.csv', 'w') as csvfile:
    writer = csv.writer(csvfile)
    writer.writerow(['Model', 'Training time (s)'])
    writer.writerow(['Linear regression', lr_train_time])
    writer.writerow(['Linear regression with PCA', lr_pca_train_time])
    writer.writerow(['Neural network', nn_train_time])
    writer.writerow(['Neural network with PCA', nn_pca_train_time])
```

1.1.5 5. Compare models' results with and without PCA on predicting V

5.1 Split the data into training and testing

```
[18]: # Split the data into training and testing
from sklearn.model_selection import train_test_split

# Without PCA
X_train, X_test, y_train, y_test = train_test_split(X_std, y_V, test_size=0.2,
    ↪random_state=42)

# With PCA
X_train_pca, X_test_pca, y_train_pca, y_test_pca = train_test_split(X_pca, y_V,
    ↪test_size=0.2, random_state=42)
```

5.2 Train and evaluate the Linear Regression model

5.2.1 Without PCA

```
[19]: X_train.isna().sum()
```

```
[19]: tem      0
      mass    0
      n       0
      occ     0
      x       0
      y       0
      dtype: int64
```

```
[20]: y_V.isna().sum()
```

```
[20]: 0
```

```
[21]: from sklearn.linear_model import LinearRegression
      from sklearn.metrics import mean_squared_error
      import pickle
      import time
```

```

# Train linear regression model on training data
model = LinearRegression()
start = time.time()
model.fit(X_train, y_train)
stop = time.time()
lr_train_time = stop - start
print(f"Training time: {lr_train_time}s")

# Save the model to a pickle file
filename = '../reports/models/lr_V_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

# Predict on test data
y_pred = model.predict(X_test)

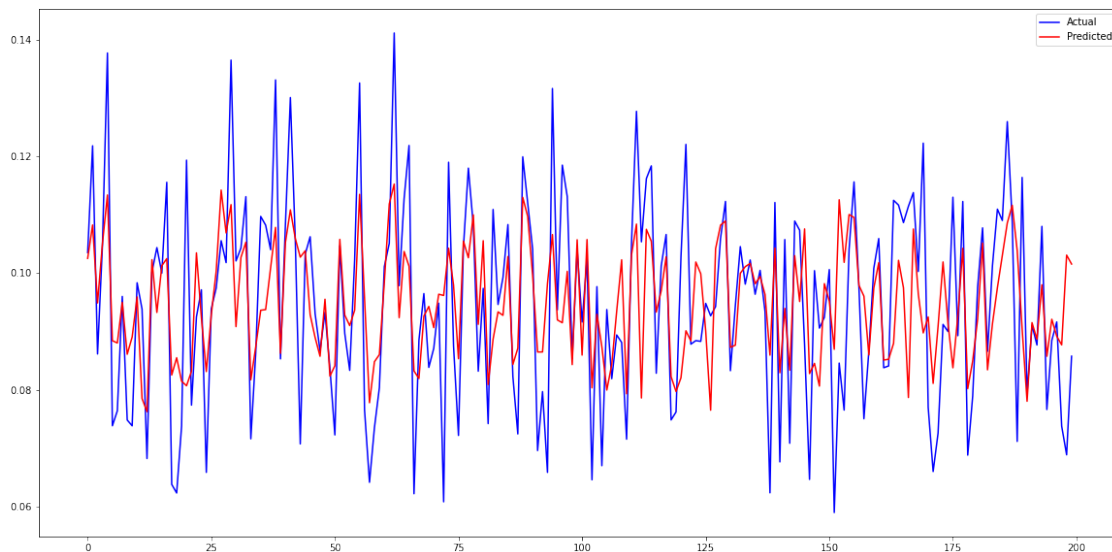
# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred[:200], color='red', label='Predicted')
plt.legend()
plt.show()

```

Training time: 0.007966041564941406s

Mean squared error: 0.0002148180985995211



5.2.2 With PCA

```
[22]: # Train linear regression model on training data
start = time.time()
model.fit(X_train_pca, y_train_pca)
stop = time.time()
lr_pca_train_time = stop - start
print(f"Training time: {lr_pca_train_time}s")

# Predict on test data
y_pred_pca = model.predict(X_test_pca)

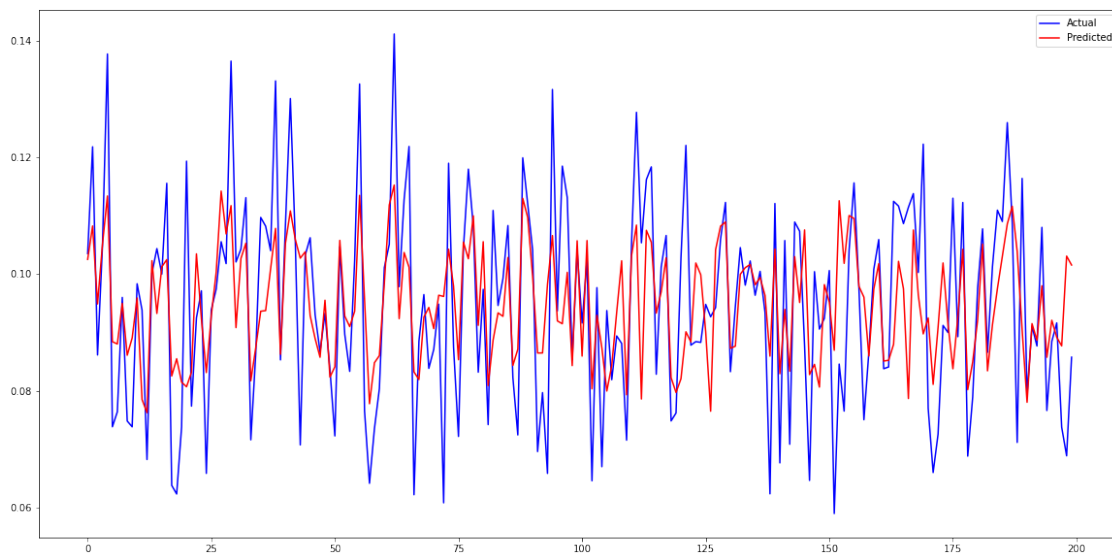
# Save the model to a pickle file
filename = '../reports/models/lr_V_pca_open_office_model.pkl'
pickle.dump(model, open(filename, 'wb'))

# Evaluate the model
mse = mean_squared_error(y_test_pca, y_pred_pca)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

Training time: 0.006435394287109375s

Mean squared error: 0.0002148180985995211



5.3 Train and evaluate the Feedforward Neural Network model

5.2.1 Without PCA

```
[23]: # Create neural network model
from keras.models import Sequential
from keras.layers import Dense

N_NEURONS = 1024
N_LAYERS = 4

model = Sequential()
model.add(Dense(units=N_NEURONS, input_dim=X.shape[1], activation='relu'))
for i in range(N_LAYERS-1):
    model.add(Dense(units=N_NEURONS, activation='relu'))
model.add(Dense(units=1, activation='linear')) # Output layer
model.compile(loss='mean_squared_error', optimizer='adam')

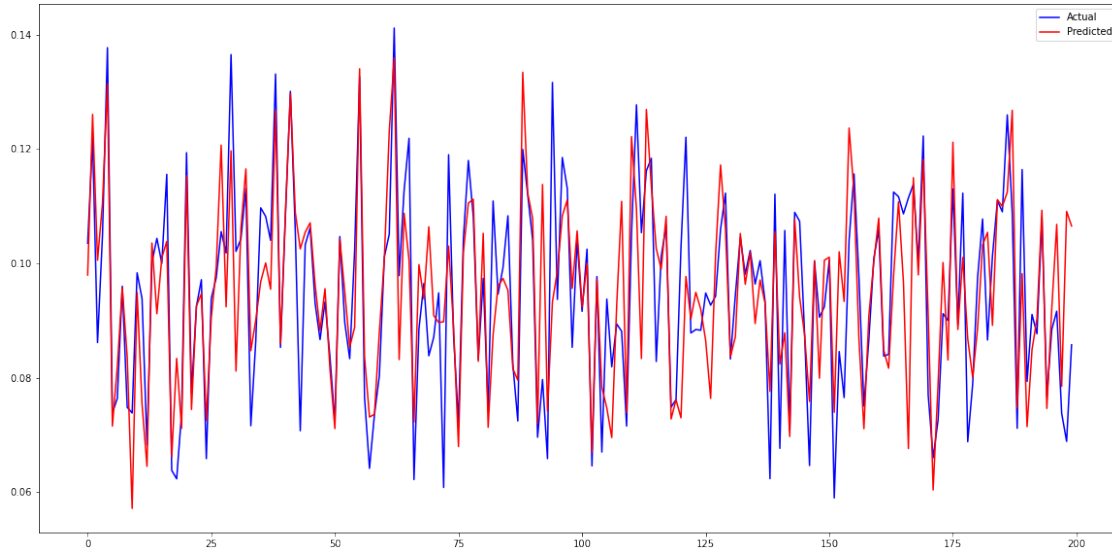
# Train the model
start = time.time()
model.fit(X_train, y_train, epochs=100, verbose=0)
stop = time.time()
nn_train_time = stop - start
print(f"Training time: {nn_train_time}s")

# Predict on test data
y_pred_nn = model.predict(X_test)
model.save('../reports/models/nn_V_open_office_model.pkl')

# Evaluate the model
mse = mean_squared_error(y_test, y_pred_nn)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

```
Training time: 47.35318875312805s
INFO:tensorflow:Assets written to:
../reports/models/nn_V_open_office_model.pkl/assets
Mean squared error: 0.0001326037106000784
```



5.3.2 With PCA

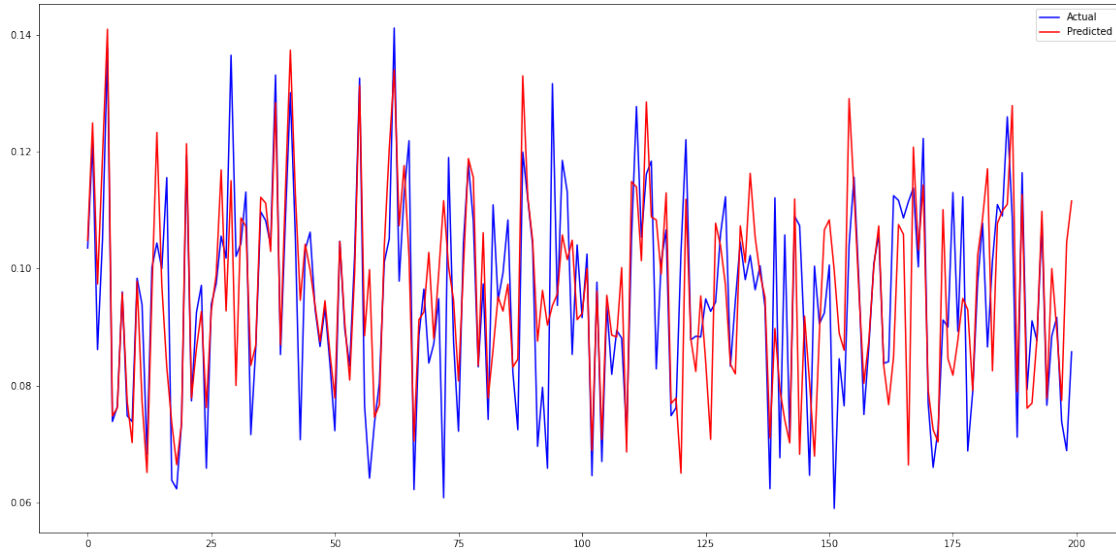
```
[24]: # Train the model
start = time.time()
model.fit(X_train_pca, y_train_pca, epochs=100, verbose=0)
stop = time.time()
nn_pca_train_time = stop - start
print(f"Training time: {nn_pca_train_time}s")

# Predict on test data
y_pred_nn_pca = model.predict(X_test_pca)
model.save('../reports/models/nn_V_pca_open_office_model.pkl')

# Evaluate the model
mse = mean_squared_error(y_test_pca, y_pred_nn_pca)
print(f'Mean squared error: {mse}')

# Plot the predictions and actual values
plt.figure(figsize=(20,10))
plt.plot(range(200), y_test_pca[:200], color='blue', label='Actual')
plt.plot(range(200), y_pred_nn_pca[:200], color='red', label='Predicted')
plt.legend()
plt.show()
```

```
Training time: 44.344775438308716s
INFO:tensorflow:Assets written to:
../reports/models/nn_V_pca_open_office_model.pkl/assets
Mean squared error: 0.00018198284814540812
```



```
[25]: # Save training times to a csv file
import csv

with open('../reports/data/training_times_open_office_V.csv', 'w') as csvfile:
    writer = csv.writer(csvfile)
    writer.writerow(['Model', 'Training time (s)'])
    writer.writerow(['Linear regression', lr_train_time])
    writer.writerow(['Linear regression with PCA', lr_pca_train_time])
    writer.writerow(['Neural network', nn_train_time])
    writer.writerow(['Neural network with PCA', nn_pca_train_time])
```