

## Problem 2 (30 Points)

Data-driven field prediction models can be used as a substitute for performing expensive calculations/simulations in design loops. For example, after being trained on finite element solutions for many parts, they can be used to predict nodal von Mises stress for a new part by taking in a mesh representation of the part geometry.

Consider the plane-strain compression problem shown in "data/plane-strain.png".

In this problem you are given node features for 100 parts. These node features have been extracted by processing each part shape using a neural network. You will perform feature selection to determine which of these features are most relevant using feature selection tools in sklearn.

*You are welcome to use any of the code provided in the lecture activities.*

Summary of deliverables:

SciKit-Learn Models: Print Train and Test MSE

- `LinearRegression()` with all features
- `DecisionTreeRegressor()` with all features
- `LinearRegression()` with features selected by `RFE()`
- `DecisionTreeRegressor()` with features selected by `RFE()`

Feature Importance/Coefficient Visualizations

- Feature importance plot for Decision Tree using all features
- Feature coefficient plot for Linear Regression using all features
- Feature importance plot for DT showing which features RFE selected
- Feature coefficient plot for LR showing which features RFE selected

Stress Field Visualizations: Ground Truth vs. Prediction

- Test dataset shape index 8 for decision tree and linear regression with all features
- Test dataset shape index 16 for decision tree and linear regression with RFE features

Questions

- Respond to the 5 prompts at the end

Imports and Utility Functions:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.tree import DecisionTreeRegressor
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.feature_selection import RFE
```

```

def plot_shape(dataset, index, model=None, lims=None):
    x = dataset["coordinates"][index][:,0]
    y = dataset["coordinates"][index][:,1]

    if model is None:
        c = dataset["stress"][index]
    else:
        c = model.predict(dataset["features"][index])

    if lims is None:
        lims = [min(c), max(c)]

    plt.scatter(x, y, s=5, c=c, cmap="jet", vmin=lims[0], vmax=lims[1])
    plt.colorbar(orientation="horizontal", shrink=.75,
pad=0, ticks=lims)
    plt.axis("off")
    plt.axis("equal")

def plot_shape_comparison(dataset, index, model, title=""):
    plt.figure(figsize=[6,3.2], dpi=120)
    plt.subplot(1,2,1)
    plot_shape(dataset, index)
    plt.title("Ground Truth", fontsize=9, y=.96)
    plt.subplot(1,2,2)
    c = dataset["stress"][index]
    plot_shape(dataset, index, model, lims = [min(c), max(c)])
    plt.title("Prediction", fontsize=9, y=.96)
    plt.suptitle(title)
    plt.show()

def load_dataset(path):
    dataset = np.load(path)
    coordinates = []
    features = []
    stress = []
    N = np.max(dataset[:,0].astype(int)) + 1
    split = int(N*.8)
    for i in range(N):
        idx = dataset[:,0].astype(int) == i
        data = dataset[idx,:]
        coordinates.append(data[:,1:3])
        features.append(data[:,3:-1])
        stress.append(data[:, -1])
    dataset_train = dict(coordinates=coordinates[:split],
features=features[:split], stress=stress[:split])
    dataset_test = dict(coordinates=coordinates[split:],
features=features[split:], stress=stress[split:])
    X_train, X_test = np.concatenate(features[:split], axis=0),
np.concatenate(features[split:], axis=0)

```

```

    y_train, y_test = np.concatenate(stress[:split], axis=0),
np.concatenate(stress[split:], axis=0)
    return dataset_train, dataset_test, X_train, X_test, y_train,
y_test

def get_shape(dataset,index):
    X = dataset["features"][index]
    y = dataset["stress"][index]
    return X, y

def plot_importances(model, selected = None, coef=False, title=""):
    plt.figure(figsize=(6,2),dpi=150)
    y = model.coef_ if coef else model.feature_importances_
    N = 1+len(y)
    x = np.arange(1,N)

    plt.bar(x,y)

    if selected is not None:
        plt.bar(x[selected],y[selected],color="red",label="Selected
Features")
        plt.legend()

    plt.xlabel("Feature")

    plt.ylabel("Coefficient" if coef else "Importance")
    plt.xlim(0,N)
    plt.title(title)
    plt.show()

```

## Loading the data

First, complete the code below to load the data and plot the von Mises stress fields for a few shapes.

You'll need to input the path of the data file, the rest is done for you.

All training node features and outputs are in `X_train` and `y_train`, respectively. Testing nodes are in `X_test`, `y_test`.

`dataset_train` and `dataset_test` contain more detailed information such as node coordinates, and they are separated by shape.

Get features and outputs for a shape by calling `get_shape(dataset,index)`. `N_train` and `N_test` are the number of training and testing shapes in each of these datasets.

```

# YOUR CODE GOES HERE
# Define data_path

```

```

dataset_train, dataset_test, X_train, X_test, y_train, y_test =
load_dataset("data/stress_nodal_features.npy")

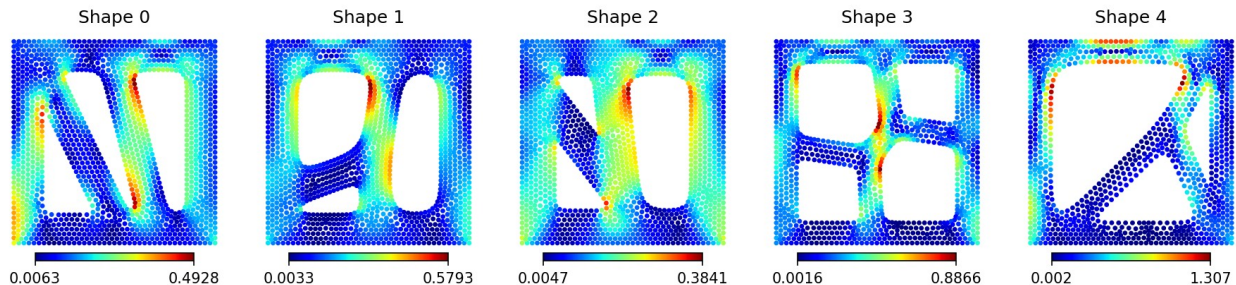
```

```

N_train = len(dataset_train["stress"])
N_test = len(dataset_test["stress"])

plt.figure(figsize=[15,3.2], dpi=150)
for i in range(5):
    plt.subplot(1,5,i+1)
    plot_shape(dataset_train,i)
    plt.title(f"Shape {i}")
plt.show()

```



## Fitting models with all features

Create two models to fit the training data `X_train`, `y_train`:

1. A `LinearRegression()` model
2. A `DecisionTreeRegressor()` model with a `max_depth` of 20

Print the training and testing MSE for each.

```

# YOUR CODE GOES HERE
model_names = ["Linear Regression", "Decision Tree Regressor"]
# linear regression model
LSR = LinearRegression()
# decision tree regressor model with max depth 20
DTR = DecisionTreeRegressor(max_depth=20, random_state=0)

# train the models and print train & test MSE
models = [LSR, DTR]
for model in models:
    model.fit(X_train, y_train)
    train_pred = model.predict(X_train)
    test_pred = model.predict(X_test)
    print(f"{model_names[models.index(model)]} Train MSE: {mean_squared_error(y_train, train_pred)}")
    print(f"{model_names[models.index(model)]} Test MSE: {mean_squared_error(y_test, test_pred)}")
    print()

Linear Regression Train MSE: 0.008110600523650646
Linear Regression Test MSE: 0.009779523126780987

```

```
Decision Tree Regressor Train MSE: 0.0004944875978805109
Decision Tree Regressor Test MSE: 0.00838549943577
```

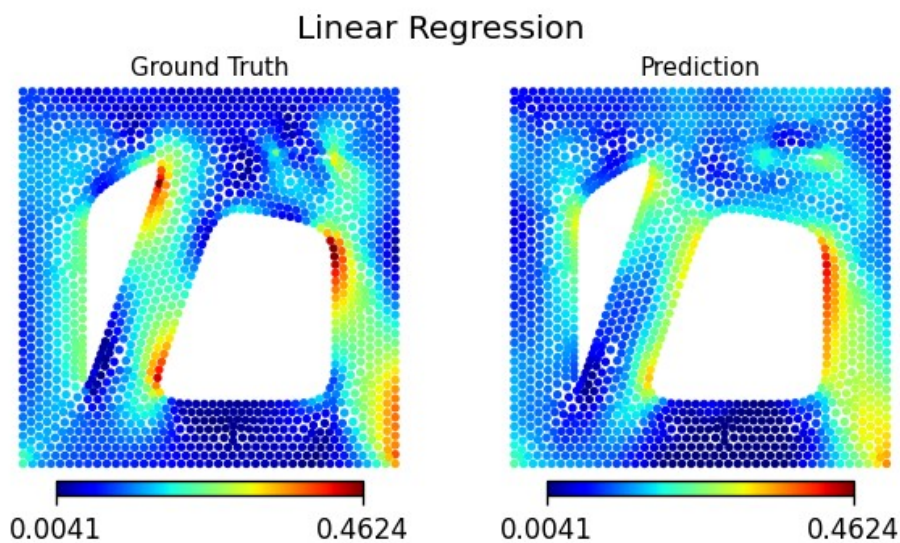
## Visualization

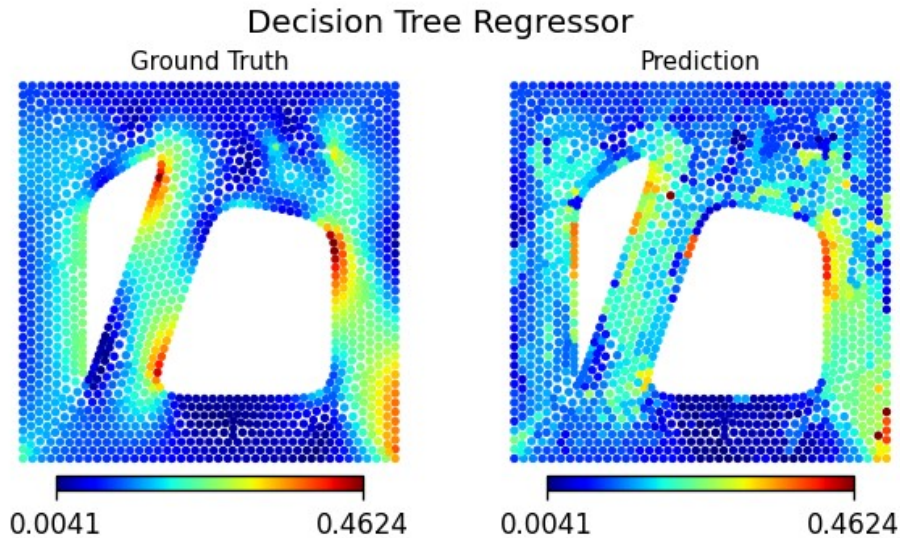
Use the `plot_shape_comparison()` function to plot the index 8 shape results in `dataset_test` for each model.

Include titles to indicate which plot is which, using the `title` argument.

```
test_idx = 8

# YOUR CODE GOES HERE
for model in models:
    plot_shape_comparison(dataset_test, test_idx, model,
                          model_names[models.index(model)])
```





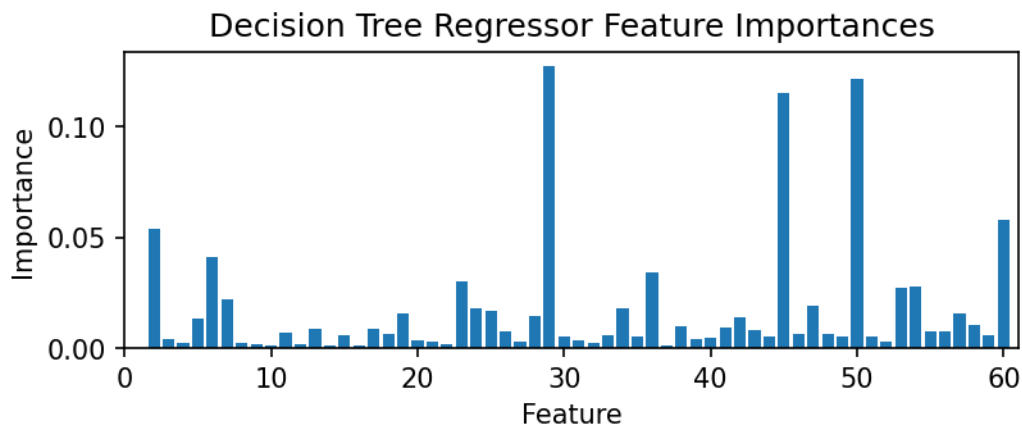
## Feature importance

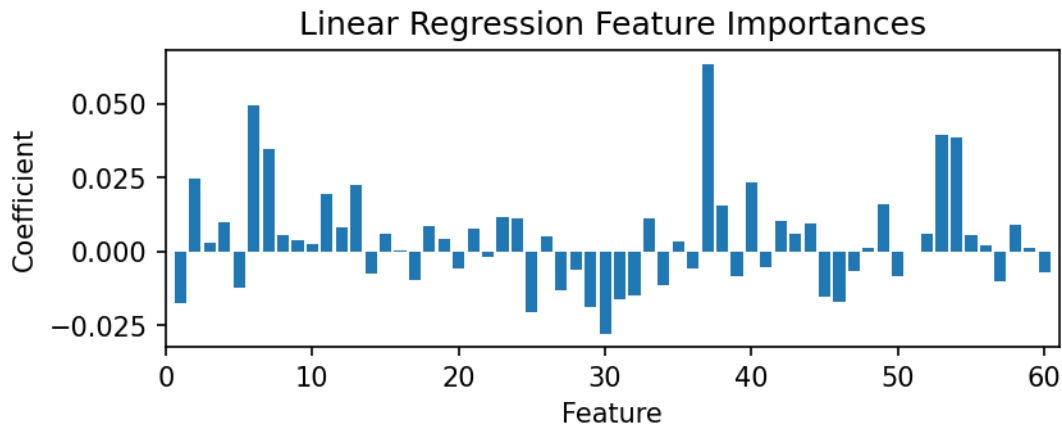
For a tree methods, "feature importance" can be computed, which can be done for an sklearn model using `.feature_importances_`.

Use the provided function `plot_importances()` to visualize which features are most important to the final decision tree prediction.

Then create another plot using the same function to visualize the linear regression coefficients by setting the "coef" argument to `True`.

```
# YOUR CODE GOES HERE
# plot the feature importances for the decision tree regressor
plot_importances(DTR, title="Decision Tree Regressor Feature
Importances")
# plot the feature importances for the linear regression model
plot_importances(LSR, coef=True, title="Linear Regression Feature
Importances")
```





## Feature Selection by Recursive Feature Elimination

Using `RFE()` in sklearn, you can iteratively select a subset of only the most important features.

For both linear regression and decision tree (depth 20) models:

1. Create a new model.
2. Create an instance of `RFE()` with `n_features_to_select` set to 30.
3. Fit the RFE model as you would a normal sklearn model.
4. Report the train and test MSE.

Note that the decision tree RFE model may take a few minutes to train.

Visit [https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_selection.RFE.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.RFE.html) for more information.

```
# YOUR CODE GOES HERE
# RFE with n_features_to_select=30

model_names = ["Linear Regression RFE", "Decision Tree Regressor RFE"]
# RFE with n_features_to_select=30
LSR_RFE = RFE(LSR, n_features_to_select=30)
DTR_RFE = RFE(DTR, n_features_to_select=30)

# train the RFE models and print train & test MSE
models = [LSR_RFE, DTR_RFE]
for model in models:
    model.fit(X_train, y_train)
    train_pred = model.predict(X_train)
    test_pred = model.predict(X_test)
    print(f"{model_names[models.index(model)]} Train MSE:
{mean_squared_error(y_train, train_pred)}")
    print(f"{model_names[models.index(model)]} Test MSE:
{mean_squared_error(y_test, test_pred)}")
    print()
```



```
Linear Regression RFE Train MSE: 0.008508718572556973
Linear Regression RFE Test MSE: 0.010150418616831303
```

```
Decision Tree Regressor RFE Train MSE: 0.0005351821126843501
Decision Tree Regressor RFE Test MSE: 0.009041167760149434
```

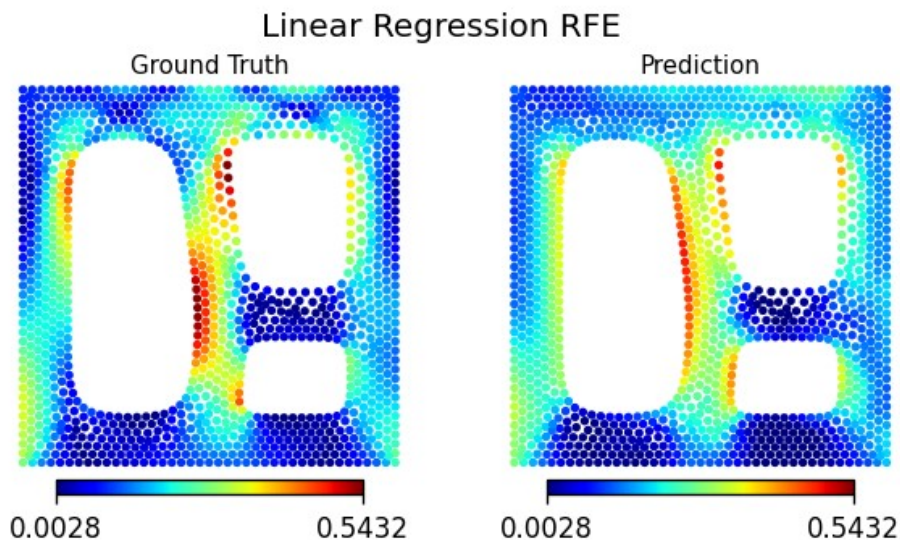
## Visualization

Use the `plot_shape_comparison()` function to plot the index 16 shape results in `dataset_test` for each model.

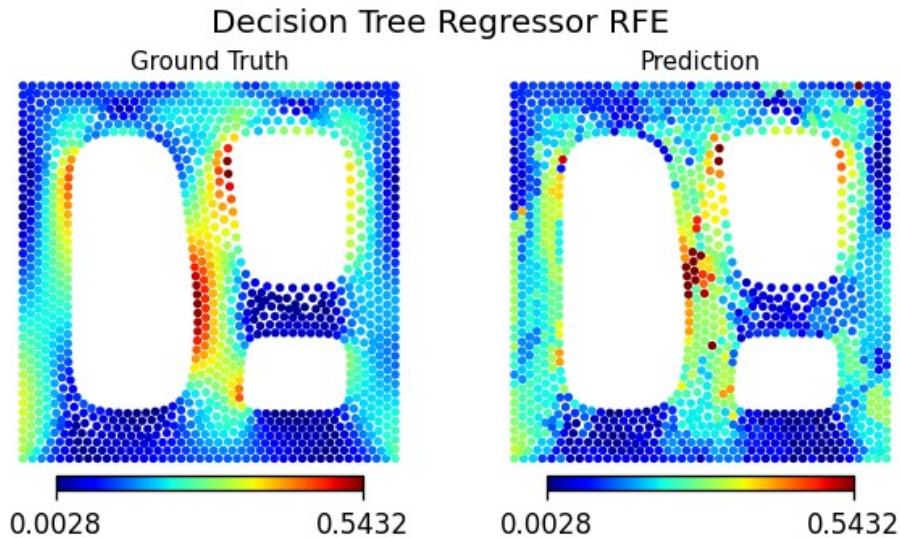
As before, include titles to indicate which plot is which, using the `title` argument.

```
test_idx = 16

# YOUR CODE GOES HERE
for model in models:
    plot_shape_comparison(dataset_test, test_idx, model,
                          model_names[models.index(model)])
```





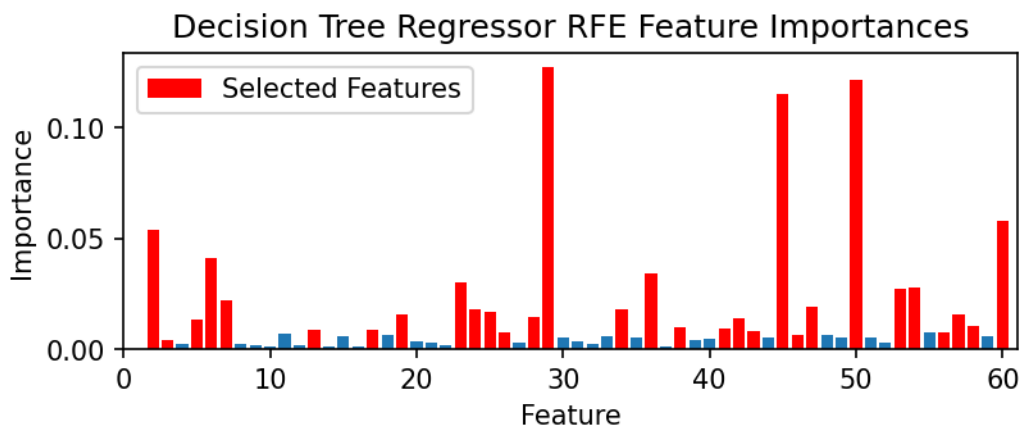


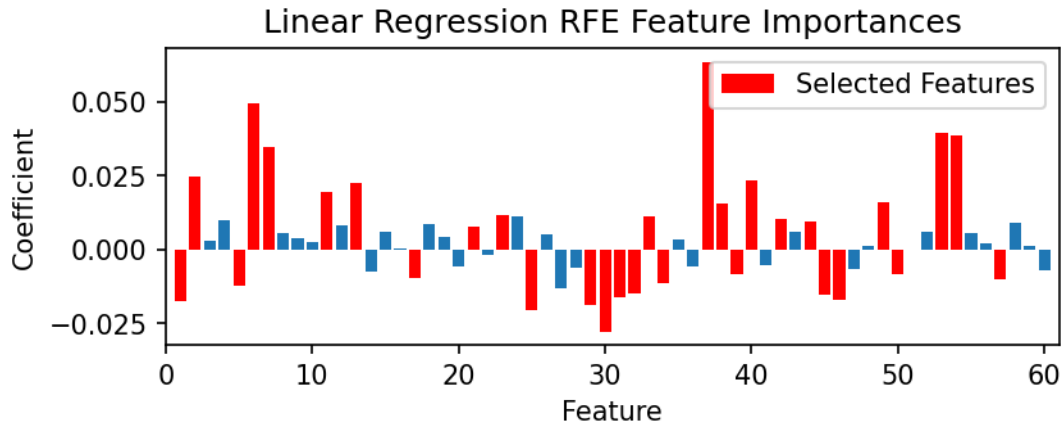
## Feature importance with RFE

Recreate the 2 feature importance/coefficient plots from earlier, but this time highlight which features were ultimately selected after performing RFE by coloring those features red. You can do this by setting the `selected` argument equal to an array of selected indices.

For an RFE model `rfe`, the selected feature indices can be obtained via `rfe.get_support(indices=True)`.

```
# YOUR CODE GOES HERE
# plot the feature importances for the decision tree regressor RFE
plot_importances(DTR, selected=DTR_RFE.get_support(indices=True),
title="Decision Tree Regressor RFE Feature Importances")
# plot the feature importances for the linear regression model RFE
plot_importances(LSR, selected=LSR_RFE.get_support(indices=True),
coef=True, title="Linear Regression RFE Feature Importances")
```





## Questions

1. Did the MSE increase or decrease on test data for the Linear Regression model after performing RFE?

The MSE increased on test data for the linear regression model after performing RFE.

2. Did the MSE increase or decrease on test data for the Decision Tree model after performing RFE?

The MSE increased on test data for the decision tree model after performing RFE.

3. Describe the qualitative differences between the Linear Regression and the Decision Tree predictions.

The decision tree has a lower MSE compared to the linear regression model. For the linear regression model, the predicted nodal von Mises stress has a smoother mesh representation compared to the decision tree model.

4. Describe how the importance of features that were selected by RFE compare to that of features that were eliminated (for the decision tree).

In the decision tree, the coefficients that were selected by RFE are of higher importance and influence on the model outcome compared to the features that were eliminated. From the above figure, we can observe that the selected features (red) have a higher importance score compared to the eliminated features (blue) that have a lower importance score.

5. Describe how the coefficients that were selected by RFE compare to that of features that were eliminated (for linear regression).

In linear regression, the coefficients that were selected (red) by RFE are expected to be larger positive or negative coefficients that are non-zero, which has the most influence on the model outcome. The features that are eliminated (blue) are features

that are close to zero or zero, which have little to no influence on the model outcome.