


```
In [27]:
import numpy as np
import pandas as pd
# this is also used in train validation splits through seed
def test_train_split(dataset: pd.DataFrame, test_size = 0.20, random_state = 0):
    # shuffle indices
    indices = np.random.permutation(len(dataset)).tolist()
    dataset_test_size = int(dataset.shape[0]*test_size)
    # separating test-train indices
    test_indices = indices[:dataset_test_size]
    train_indices = indices[dataset_test_size:]
    return (dataset.iloc[test_indices], dataset.iloc[train_indices])
```

Defining Standardization Function

```
In [28]:
def standardize(data: pd.DataFrame, skip: str=None) -> None:
    for col in data.columns:
        if skip != col:
            # standardizing data
            data[col] = (data[col] - data[col].mean())/data[col].std()
```

Defining Metric Functions

```
In [30]:
def MSE(Y_pred: np.ndarray, Y_test: np.ndarray) -> float:
    return np.sum(np.square(Y_pred - Y_test))/(Y_test.shape[0])

def rsquared_score(Y_pred: np.ndarray, Y_test: np.ndarray) -> float:
    data_var = np.sum((Y_test - np.mean(Y_test,axis=0))**2)
    model_var = np.sum((Y_pred - Y_test)**2)
    expl_var = data_var - model_var
    return float(expl_var/data_var)
```

2. ML Model 1

Support Vector Regression

Defining SVR function with SMO

Out of epsilon, nu, and least squares SVRs, we selected the epsilon SVR for this project. Epsilon SVR gives us direct control over the accuracy of the model, by controlling the allowed error. Although, this may lead to complex models (with lot of support vectors).

For faster training of our model we employ Sequential Minimal Optimization (SMO) algorithm, the referenced paper has been attached.

```
In [31]:
import pandas as pd
from typing import Callable
import math

def gaussian_kernel(x: np.ndarray, z: np.ndarray, sigma: float) -> np.ndarray:
    n = x.shape[0]
    m = z.shape[0]
    x2 = np.dot(np.sum(x**2, 1).reshape(n, 1), np.ones((1, m)))
    z2 = np.dot(np.sum(z**2, 1).reshape(m, 1), np.ones((1, n)))
    return np.exp(-(x2 + z2.T - 2 * np.dot(x, z.T)) / (2 * sigma**2))

def linear_kernel(x: np.ndarray, z: np.ndarray) -> np.ndarray:
    return np.matmul(x, z.T)

# E-toler is the amount of additional error that is allowed while checking the KKT conditions
def svr_train(X_train: np.ndarray, Y_train: np.ndarray, kernel: Callable[[np.ndarray, np.ndarray], np.ndarray], C=1.0,
              # variable initialization
              np.random.seed(random_seed),
              n_samples: n, features = X_train.shape
              kernel_matrix = kernel(X_train, X_train)
              is_changed = np.zeros(n_samples)
              alphas = np.zeros(n_samples)
              b = 0
              iter_num = 1
              alphas_pairs_changed = 0

# SMO training
while (alphas_pairs_changed > 0 or iter_num == 1) and iter_num <= max_iter:
    # After one pass through the training set, the outer Loop iterates over only those examples whose Lagrange multi
    seq_SV = list(range(n_samples))
    else:
        seq_SV = np.nonzero((alphas > -C)*(alphas < C))[0]
        for i in seq_SV:
            # Calculate error at index i
            Ei = alphas.dot(kernel_matrix[i,:].T) + b - Y_train[i]
            # Check KKT conditions
            condition_1 = alphas[i] == 0 and abs(Ei) < (epsilon - f_toler)
            condition_2 = alphas[i] != 0 and abs(alphas[i]) < C and
            abs(Ei) < (epsilon + f_toler) and abs(Ei) >= (epsilon - f_toler)
            condition_3 = abs(alphas[i]) == C and abs(Ei) > (epsilon - f_toler)
            if (condition_1 or condition_2 or condition_3):
                continue
            # Select a j (another data sample)
            = Ej =
            changed_seq = np.nonzero(is_changed)[0]
            if changed_seq.shape[0] == 0:
                random_seq = np.arange(n_samples)
                # remove i from the choice list
                random_seq = np.append(random_seq[:i], random_seq[i+1:])
                j = np.random.choice(random_seq)
            # Choose the new index j
            Ej = alphas.dot(kernel_matrix[j,:].T) + b - Y_train[j]
            else:
                temp_error = abs(temp_error - Ei)
                # take the i and Ej with max step (|Ej - Ei|)
                if temp_error >= max_step:
                    max_step = temp_error
                    j, Ej = c, temp_error
# all j of the support vectors now pass the KKT conditions
if j == -1:
    break
# Save the old i and j
alpha_i_old = alphas[i].copy()
alpha_j_old = alphas[j].copy()
# Calculate the lower and upper bound
lower_bound = max(-C, alpha_i_old + alpha_j_old - C)
upper_bound = min(C, alpha_i_old + alpha_j_old + C)
# This may not happen, but if the Lower-bound equals to upper_bound, continue
if lower_bound == upper_bound:
    continue
# Calculate eta
eta = kernel_matrix[i, j] + kernel_matrix[j, j] - 2.0 * kernel_matrix[i, j]
# May not happen
if eta <= 0:
    continue
# update j
is_updated = False
for i in range(1, n):
    temp_j = alpha_i_old + (Ei - Ej + epsilon)
    temp_j = np.sign((temp_j) - np.sign(temp_j)) == sign:
        is_updated = True
        break
if not is_updated:
    continue
alphas[j] = max(min(alphas[j], upper_bound), lower_bound)
if abs(alphas[j] - alpha_j_old) <= 1e-5:
    is_changed[j] = 1
    continue
alphas[j] += alpha_j_old - alphas[j]
# Calculate bi and bj
C = (Ei + (alphas[i] - alpha_i_old) * kernel_matrix[i, i]) \
    + (alphas[j] - alpha_j_old) * kernel_matrix[i, j]) + b
bj = -(Ei + (alphas[i] - alpha_i_old) * kernel_matrix[i, i]) \
    + (alphas[j] - alpha_j_old) * kernel_matrix[i, j]) + b
# Check if bi or bj is available (alpha is within bounds)
if abs(alphas[i]) < C:
    b = bi
    elif abs(alphas[j]) < C:
        b = bj
    else:
        b = (bi + bj) / 2
    is_changed[i] = 1
    is_changed[j] = 1
    alphas_pairs_changed += 1
iter_num += 1
return alphas, b

def svr_predict(x_train: np.ndarray, X_test: np.ndarray, alphas: np.ndarray, b: float, kernel: Callable[[np.ndarray, np
return alphas @ kernel(X_train, X_test) + b

"""Testing the model"""
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
standardize(data, skip = 'co2')
Y = data['co2']
X = data.drop(['co2'], axis = 1)

test, train = test_train_split(data)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

# hyper-params: c, gamma, epsilon
kernel = lambda x, y: gaussian_kernel(x, y, 1/len(X_train))
# Large c = more bias is overfitting = more time to train
alpha, b = svr_train(X_train, Y_train, kernel, 1, 0.01, 42, 1000, 0.01)
Y_pred = svr_predict(X_train, X_test, alphas, b, kernel)
print(alphas, b)
print(MSE(Y_pred, Y_test))
print(rsquared_score(Y_pred, Y_test))

7527.3136942675155
-1.152656178367875

3. ML Model 2
```

Artificial Neural Network

```
In [32]:
# parameter of Leaky Rectified Linear Unit activation function (Leaky ReLU)
param = 0.01

# A normal layer of ANN (Linear)
class LinearLayer:
    def __init__(self, input_size: int, output_size: int):
        # weights of this layer
        self.weights = np.zeros((input_size, output_size))
        self.bias = np.zeros((1, output_size))

    # take 1-D input and produce 1-D output
    def forward_propagation(self, data: np.ndarray) -> np.ndarray:
        self.input = data
        # print(self.input.shape, self.weights.shape, self.bias.shape)
        return self.input @ self.weights + self.bias

    # computes dErr/dW (gradient), for a given dE/dY (output err), Returns dE/dX (input err)
    def backward_propagation(self, output_error: float, learning_rate: float):
        # update parameters
        self.weights = self.weights + learning_rate * self.input.T @ output_error
        self.bias = self.bias + learning_rate * output_error
        return output_error @ self.weights.T

# An activation layer of the ANN (non-linear)
class ActivationLayer:
    def __init__(self, activation_fn: Callable[[np.ndarray], np.ndarray], activation_fn_derivative: Callable[[np.ndarray]]
        self.afn = activation_fn
        self.afn_der = activation_fn_derivative

    # returns the activated input (non-linear transform using tanh)
    def forward_propagation(self, data):
        self.input = data
        return self.afn(data)

# This layer has no learnable parameters, so no need for any updation
def backward_propagation(self, output_error, learning_rate):
    return self.afn_der(self.input) * output_error

class ANN:
    def __init__(self, layers: list[LinearLayer| ActivationLayer]) -> None:
        # layers of ANN
        self.layers = layers

    def fit(self, X_train: np.ndarray, Y_train: np.ndarray, max_iters = 1000, learning_rate = 0.1):
        n_samples, n_features = X_train.shape
        for i in range(max_iters):
            error = 0
            for j in range(n_samples):
                # forward propagation (prediction) for one sample
                out = X_train[j].reshape((1, n_features))
                for layer in self.layers:
                    out = layer.forward_propagation(out)
                # backward propagation for the same sample
                error += 2*(out - Y_train[j])*(Y_train[j])
                for layer in reversed(self.layers):
                    error = layer.backward_propagation(error, learning_rate)

    def predict(self, X_test: np.ndarray):
        # sample dimension first
        samples = len(X_test)
        count = 1
        # run network over all samples
        for i in range(samples):
            # forward propagation
            output = X_test[i]
            for layer in self.layers:
                output = layer.forward_propagation(output)
            result.append(output)
        return np.array(result)

activation_fn = lambda x: np.maximum(param*x, x)
def activation_fn_derivative(x: np.ndarray):
    xc = np.full_like(x, param)
    xc[x >= 0] = 1
    return xc

"""Testing the model"""
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
standardize(data, skip = 'co2')
Y = data['co2']
X = data.drop(['co2'], axis = 1)

test, train = test_train_split(data, 0.5)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

# number of connections between two linear layers
n = 2
network = ANN([LinearLayer(X_train.shape[1], n), ActivationLayer(activation_fn, activation_fn_derivative), LinearLayer(n,
Y_pred = network.predict(X_test)

print(MSE(Y_pred, Y_test))
print(rsquared_score(Y_pred, Y_test))

12765297.36835541
-0.51212153204784
```

4. ML Model 3

Random Forest Regressor

Defining Trees of the Random Forest

This decision tree uses mean absolute error as the impurity measure. These trees only perform binary split wherever applicable

```
In [33]:
import math
import numpy as np

class DecisionNode:
    def __init__(self, feature=None, threshold=None, left=None, right=None, value=None):
        self.feature = feature
        self.threshold = threshold
        self.left = left
        self.right = right
        self.value = value

class DecisionTree:
    def __init__(self, max_depth=20, min_samples=10, random_forest=False):
        self.max_depth = max_depth
        self.min_samples = min_samples
        self.tree = []
        self.random = random_forest

    def fit(self, X_train: np.ndarray, Y_train: np.ndarray):
        self.tree = self.grow_tree(X_train, Y_train)

    def predict(self, X_test: np.ndarray):
        return np.array([self.traverse_tree(row, self.tree) for row in X_test])

    def best_split(self, X: np.ndarray, y: np.ndarray):
        best_feature: int = None
        best_threshold: float = None
        mse_gain = -1

        # go through a random subset of features (and all thresholds) for random forest, else go through all and find the s
        for b in features:
            thresholds = np.unique(X[:, b])
            for threshold in thresholds:
                gain = self.split_gain(X[:, b], y, threshold)
                if gain > best_gain:
                    best_gain = gain
                    best_feature = b
                    best_threshold = threshold
                return best_feature, best_threshold

    def split_gain(self, X: column, y: threshold):
        # define a impurity finding function
        def impurity(samples: np.ndarray) -> float:
            s_avg = np.mean(samples)
            # consider the mean absolute error as the impurity measure
            return np.mean(np.abs(samples - s_avg))

        n = len(y)
        parent_imp = impurity(y)

        left_indexes = np.argwhere(X[:, b] <= threshold).flatten()
        right_indexes = np.argwhere(X[:, b] > threshold).flatten()
        child_imp = 0

        # take child's impurity's weighted average across the two way split
        if len(left_indexes) != 0:
            child_imp += (len(left_indexes) / n) * impurity(y[left_indexes])
        if len(right_indexes) != 0:
            child_imp += (len(right_indexes) / n) * impurity(y[right_indexes])

        return parent_imp - child_imp

    def grow_tree(self, X, y, depth=0):
        n_samples, n_features = X.shape
        if n_samples <= self.min_samples or depth >= self.max_depth:
            return TreeNode(value = np.mean(y))

        best_feature, best_threshold = self.best_split(X, y)

        left_indexes = np.argwhere(X[:, best_feature] <= best_threshold).flatten()
        right_indexes = np.argwhere(X[:, best_feature] >= best_threshold).flatten()

        if len(left_indexes) == 0 or len(right_indexes) == 0:
            return TreeNode(value=np.average(y))

        # build the tree on both children with remaining samples
        left_data = self.grow_tree(X[left_indexes, :], y[left_indexes], depth+1)
        right_data = self.grow_tree(X[right_indexes, :], y[right_indexes], depth+1)

        # return the current node
        return TreeNode(best_feature, best_threshold, left, right)

    def traverse_tree(self, x: np.ndarray, tree: list[TreeNode]):
        # return the value of the node if we reach a leaf
        if tree.value is not None:
            return tree.value
        # else continue travers
        if left_feature <= tree.threshold:
            return self.traverse_tree(x, tree.left)
        return self.traverse_tree(x, tree.right)

"""Testing the model"""
# get the test train split
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
standardize(data, skip = 'co2')
Y = data['co2']
X = data.drop(['co2'], axis = 1)

test, train = test_train_split(data)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

# max depth and min samples per leaf are hyperparams
cl = DecisionTree()
cl.fit(X_train, Y_train)
Y_pred = cl.predict(X_test)
print(MSE(Y_pred, Y_test))
print(rsquared_score(Y_pred, Y_test))

0.95788644392173
0.9918023447525804
```

Defining RandomForest Class

```
In [34]:
class RandomForest():
    def __init__(self, n_estimators, sample_size=None, min_leaf=3, max_depth=10):
        self.n_estimators = n_estimators
        self.sample_size = sample_size
        self.min_leaf = min_leaf
        self.max_depth = max_depth
        self.trees = None

    def fit(self, X_train: np.ndarray, Y_train: np.ndarray, random_seed=0):
        self.trees = []
        np.random.seed(random_seed)
        if self.sample_size is None:
            self.sample_size = len(X_train)
        for _ in range(self.n_estimators):
            indexes = np.random.randint(len(Y_train), size=self.sample_size)
            dt = DecisionTree(self.max_depth, self.min_leaf, True)
            dt.fit(X_train[indexes, :], Y_train[indexes])
            self.trees.append(dt)

    def predict(self, X):
        tree_predictions = np.array([tree.predict(x) for tree in self.trees])
        prediction = np.mean(tree_predictions, axis=0)
        return prediction

"""Testing the model"""
# max depth, min samples per leaf, sample size per tree and no of estimators are hyperparams
cl = RandomForest(n_estimators=50, sample_size=len(X_train), min_leaf=3, max_depth=10)
cl.fit(X_train, Y_train)
Y_pred = cl.predict(X_test)
print(MSE(Y_pred, Y_test))
print(rsquared_score(Y_pred, Y_test))

20.80739526231540
0.9915285518715435
```

5. ML Model 4

Elastic Net Regression

```
In [50]:
import numpy as np
import pandas as pd

# assumes dummy ones in X-test
def elastic_net_predict(X_test: np.ndarray, W: np.ndarray) -> np.ndarray:
    return np.dot(X_test, W)

def elastic_net_train(X_train: np.ndarray, Y_train: np.ndarray, learning_rate = 0.01, l1_penalty = 50, l2_penalty = 3,
                    n_samples, n_features = X_train.shape):
    W = np.zeros(n_features + 1)
    X_train = np.concatenate((np.ones((n_samples, 1)), X_train), axis = 1)
    for _ in range(max_iters):
        Y_pred = elastic_net_predict(X_train, W)
        print(MSE(Y_pred, Y_test))
        W = (learning_rate/n_samples)*grad
    return W

"""Testing the model"""
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
standardize(data, skip = 'co2')
Y = data['co2']
X = data.drop(['co2'], axis = 1)

test, train = test_train_split(data)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()
W = elastic_net_train(X_train, Y_train, 0.01, 0.01, 0.001, 10000)
Y_pred = elastic_net_predict(np.concatenate((np.ones((len(X_test), 1)), X_test), axis = 1), W)
print(MSE(Y_pred, Y_test))
print(rsquared_score(Y_pred, Y_test))

250.81338919 2.20645808 24.29668318 13.35625935
20.83147469 -6.6131026 -30.23899655 -14.12379727 -13.91612893
34.4365688877235
0.990811489593665

6. Comparison of insights drawn from the models
```

Identifying the Best/Safest Car Model

Based on Exploratory Data Analysis, we observe that Fuel Consumption, Engine Size, Cylinders affect CO2 Emissions the most. The best car model will be the one which simultaneously minimizes all three since all of them show strong positive correlation, minimizing one leads to a decrease in other two. Out of categorical data, Fuel Type and Make affect the Emissions a little. Cars using Diesel or Regular Gasoline should be preferred. Usage of Natural Gas is included, but it is inconclusive, due to the lack of samples in that category.

Comparing the effects of considering fuel consumption separately

Checking for Elastic Search

```
In [40]:
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
# standardize all columns except the target
standardize(data, skip = 'co2')

mse1 = r21 - mse2 = r22 - mse3 = r23 = 0
# checking avg value across multiple runs
for _ in range(5):
    # dropping all fuel consumption combo columns
    data2 = data.drop(['fuel_cons_combo', 'fuel_cons_combo_mpg'], axis = 1)
    test, train = test_train_split(data2, random_state=1)
    Y_train = train['co2'].to_numpy()
    X_train = train.drop(['co2'], axis = 1).to_numpy()
    Y_test = test['co2'].to_numpy()
    X_test = test.drop(['co2'], axis = 1).to_numpy()

    W = elastic_net_train(X_train, Y_train, 0.01, 0.01, 0.001, 10000)
    Y_pred = elastic_net_predict(np.concatenate((np.ones((len(X_test), 1)), X_test), axis = 1), W)
    print(MSE(Y_pred, Y_test))
    r21 += rsquared_score(Y_pred, Y_test)

# leaving only fuel consumption combo (in l/m)
data2 = data.drop(['fuel_cons_city', 'fuel_cons_hwy', 'fuel_cons_combo_mpg'], axis = 1)
test, train = test_train_split(data2, random_state=1)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

W = elastic_net_train(X_train, Y_train, 0.01, 0.01, 0.001, 10000)
Y_pred = elastic_net_predict(np.concatenate((np.ones((len(X_test), 1)), X_test), axis = 1), W)
print(MSE(Y_pred, Y_test))
r22 += rsquared_score(Y_pred, Y_test)

# leaving only fuel consumption combo (in mpg)
data2 = data.drop(['fuel_cons_city', 'fuel_cons_hwy', 'fuel_cons_combo'], axis = 1)
test, train = test_train_split(data2, random_state=1)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

W = elastic_net_train(X_train, Y_train, 0.01, 0.01, 0.001, 10000)
Y_pred = elastic_net_predict(np.concatenate((np.ones((len(X_test), 1)), X_test), axis = 1), W)
print(MSE(Y_pred, Y_test))
r23 += rsquared_score(Y_pred, Y_test)

print("Fuel consumption considered separately")
print(mse1/5)
print(r21/5)
print("Fuel consumption combo in 1/100km")
print(r22/5)
print(mse2/5)
print("Fuel consumption combo in mpg")
print(r23/5)
print(mse3/5)

Fuel consumption considered separately
17.154143067895557
0.99593358082643
Fuel consumption combo in 1/100km
34.4365688877235
0.9894949819595663
Fuel consumption combo in mpg
247.5744834624199
0.992778683783999

Inference: We see in both models, Fuel consumption combo in 1/100km or separately yields better results on average than using Fuel Consumption Combo in mpg. Thus to minimise the loss and increase efficiency, it is best to only use fuel consumption combo in 1/100km.
```

Checking for Random Forest

```
In [42]:
data = pd.read_csv("Team16_Preprocessed.csv", index_col=0, header=0)
standardize(data, skip = 'co2')

mse1 = r21 - mse2 = r22 - mse3 = r23 = 0
# defining the class
cl = RandomForest(n_estimators=20, sample_size=len(data), min_leaf=3, max_depth=8)
# checking avg value across multiple runs
for _ in range(5):
    # dropping all fuel consumption combo columns
    data2 = data.drop(['fuel_cons_city', 'fuel_cons_hwy', 'fuel_cons_combo_mpg'], axis = 1)
    test, train = test_train_split(data2, random_state=1)
    Y_train = train['co2'].to_numpy()
    X_train = train.drop(['co2'], axis = 1).to_numpy()
    Y_test = test['co2'].to_numpy()
    X_test = test.drop(['co2'], axis = 1).to_numpy()

    cl.fit(X_train, Y_train, random_seed=1)
    Y_pred = cl.predict(X_test)
    mse1 += MSE(Y_pred, Y_test)
    r21 += rsquared_score(Y_pred, Y_test)

# leaving only fuel consumption combo (in l/m)
data2 = data.drop(['fuel_cons_city', 'fuel_cons_hwy', 'fuel_cons_combo_mpg'], axis = 1)
test, train = test_train_split(data2, random_state=1)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

cl.fit(X_train, Y_train, random_seed=1)
Y_pred = cl.predict(X_test)
mse2 += MSE(Y_pred, Y_test)
r22 += rsquared_score(Y_pred, Y_test)

# leaving only fuel consumption combo (in mpg)
data2 = data.drop(['fuel_cons_city', 'fuel_cons_hwy', 'fuel_cons_combo'], axis = 1)
test, train = test_train_split(data2, random_state=1)
Y_train = train['co2'].to_numpy()
X_train = train.drop(['co2'], axis = 1).to_numpy()
Y_test = test['co2'].to_numpy()
X_test = test.drop(['co2'], axis = 1).to_numpy()

cl.fit(X_train, Y_train, random_seed=1)
Y_pred = cl.predict(X_test)
mse3 += MSE(Y_pred, Y_test)
r23 += rsquared_score(Y_pred, Y_test)

print("Fuel consumption considered separately")
print(mse1/5)
print(r21/5)
print("Fuel consumption combo in 1/100km")
print(r22/5)
print(mse2/5)
print("Fuel consumption combo in mpg")
print(r23/5)
print(mse3/5)

Fuel consumption considered separately
17.154143067895557
0.99593358082643
Fuel consumption combo in 1/100km
34.4365688877235
0.9894949819595663
Fuel consumption combo in mpg
247.5744834624199
0.992778683783999

Inference: We see in both models, Fuel consumption combo in mpg does not give good results, while the other two give close results. Thus we only use fuel consumption combo in 1/100km for best results and efficiency.
```

Tuning ANN

```
In [49]:
# Defining a function for finding error in svr across random train and validation sets
def svr_train(kernel: Callable[[np.ndarray, np.ndarray], np.ndarray], C=10.0, epsilon=0.1, random_seed=0, max_iter=match
n_train = len(X_train)
n_validation = len(X_validation)
for i in range(count):
    validate, train = test_train_split(train_validate, random_state=random_seed)
    X_train = train.drop(['co2'], axis = 1).to_numpy()
    Y_validate = validate.drop(['co2'], axis = 1).to_numpy()
    X_validation = validate.drop(['co2'], axis = 1).to_numpy()
    Y_pred = svr_predict(X_train, X_validation, alphas, b, kernel)
    mse_validate += MSE(Y_pred, Y_validate)
    r2_validate += rsquared_score(Y_pred, Y_validate)
    mse_train += MSE(Y_pred, Y_train)
    r2_train += rsquared_score(Y_pred, Y_train)
    print("validation MSE=", mse_validate/count)
    print("train MSE=", mse_train/count)
    print("train r2=", r2_train/count)

linear = lambda x, y: linear_kernel(x, y)
# computing
for c in [0.1, 1, 10, 100]:
    print("\nC=", c)
    svr = for_linear(C=c)

print()

# comparing epsilon
for e in [0.01, 0.1, 1, 10, 100]:
    print("\nepsilon=", e)
    svr = for_linear(epsilon=e)

print()

# comparing kernels
kernel = lambda x, y: gaussian_kernel(x, y, 1/len(X_train))
svr = for_kernel(kernel)

C=0.1
validation MSE= 4386.992232872809
validation r2= 0.2045943718757805
train MSE= 4149.165847937683
train r2= -0.1829541086172784

C=1
validation MSE= 4206.799897298887
validation r2= -0.15323458789595252
train MSE= 3961.6259504871430545

C=10
validation MSE= 3287.981171851158
validation r2= 0.49481639367812695
train MSE= 3257.88144487578
train r2= 0.0717866347876896

C=100
validation MSE= 3645808.3765848973
validation r2= -994.7529736807712
train MSE= 1637322.1562332167
train r2= -1037.4483857666437

epsillon=0.01
validation MSE= 3287.981171851158
validation r2= 0.49481639367812695
train MSE= 3257.88144487578
train r2= 0.0717866347876896

epsillon=0.1
validation MSE= 3287.981171851158
validation r2= 0.49481639367812695
train MSE= 3257.88144487578
train r2= 0.0717866347876896

epsillon=1
validation MSE= 6505.9823691366467
validation r2= -0.8978721275982899
train MSE= 6394.96036916174
train r2= -0.8175426595862622

epsillon=10
validation MSE= 9669.82675860198
validation r2= -1.071329192941349
train MSE= 9665.007361224563
train r2= -1.7513426340152563

epsillon=100
validation MSE= 38839.99898254254
validation r2= -9.886815748383924
train MSE= 38208.01627989037
train r2= -10.126343243758888

Gaussian kernel
validation MSE= 4468.035494988804
validation r2= 0.21839647284145395
train MSE= 4770.77114819591
train r2= 0.025988728932193

Inference: SVR is not a good model for the given problem. We look towards the next models for better results.
```

Improvement Suggestions: Due to lack of computational resources, we are unable to train complex svr models and apply standard algorithms for svr like quadratic programming. Also more kernels, like polynomial kernel, can be tested than the ones mentioned here for better results.

Tuning ANN


```
In [65]:
# Defining a function for finding error in ann across random train and validation sets
def ein_error(layers: List[LinearLayer|ActivationLayer]):
    mse_train = mse_validate = r2_validate = 0
    count = 1
    for i in range(count):
        network = ANN(layers)
        validate, train = test_train_split(train_validate, random_state=i)
        Y_train = train['co2'].to_numpy()
        X_train = train.drop(['co2'], axis = 1).to_numpy()
        Y_validate = validate['co2'].to_numpy()
        X_validate = validate.drop(['co2'], axis = 1).to_numpy()
        network.fit(X_train, Y_train)
        mse_train += MSE(Y_pred, Y_validate)
        r2_validate += rsquared_score(Y_pred, Y_validate)
        Y_pred = network.predict(X_validate)
        mse_train += MSE(Y_pred, Y_train)
        r2_train += rsquared_score(Y_pred, Y_train)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)
    print("train MSE=", mse_train/count)
    print("train r2=", r2_train/count)

print("\nSingle Layer with ReLU activation function")
ann_error([LinearLayer(train_validate.shape[1]-1, 1), ActivationLayer(activation_fn, activation_fn_derivative)])

Single Layer with ReLU activation function
validation MSE= 53864918.721634544
validation r2= 15113.369728508851
train MSE= 214080456.02578485
train r2= -68665.45645551516
```

Inference: Training ANN is a resource intensive task, thus we were unable to test it thoroughly. According to the testing we have done:

1. Using linear layers alone is like using linear regression without any penalty. This leads to very large errors which the system is unable to handle (nan values).

Improvement Suggestions: Exploring other activation functions like tanh will be helpful. More layers can be added.

Tuning Random Forest

```
In [63]:
# Defining a function for finding error in random forest across random train and validation sets
def rf_error(n_estimators=20, min_leaf=3, max_depth=8):
    mse_train = r2_train = mse_validate = r2_validate = 0
    count = 3
    rf = RandomForest(n_estimators=n_estimators, min_leaf=min_leaf, max_depth=max_depth)
    for i in range(count):
        validate, train = test_train_split(train_validate, random_state=i)
        Y_train = train['co2'].to_numpy()
        X_train = train.drop(['co2'], axis = 1).to_numpy()
        Y_validate = validate['co2'].to_numpy()
        X_validate = validate.drop(['co2'], axis = 1).to_numpy()
        rf.fit(X_train, Y_train, 1)
        Y_pred = rf.predict(X_validate)
        mse_validate += MSE(Y_pred, Y_validate)
        r2_validate += rsquared_score(Y_pred, Y_validate)
        Y_pred = rf.predict(X_train)
        mse_train += MSE(Y_pred, Y_train)
        r2_train += rsquared_score(Y_pred, Y_train)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)
    print("train MSE=", mse_train/count)
    print("train r2=", r2_train/count)

print("\nFinding the best number of estimators")
for e in [10, 20, 30, 40, 50, 60]:
    print("number of trees=", e)
    rf = rf_error(n_estimators=e)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)

for d in range(5, 15):
    print("depth=", d)
    rf = rf_error(max_depth=d)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)

for l in range(1, 3):
    print("Samples per leaf=", l)
    rf = rf_error(min_leaf=l)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)

Finding the best number of estimators
number of trees= 10
validation MSE= 42.610870613223993533
validation r2= 0.988285823257035
train MSE= 35.6228772857085
train r2= 0.989493263748126
number of trees= 20
validation MSE= 40.6943890347693
validation r2= 0.988152223993533
train MSE= 34.18186448412649
train r2= 0.9902593803838695
number of trees= 30
validation MSE= 39.788091926115974
validation r2= 0.9890621748031427
train MSE= 34.22031569937146
train r2= 0.9904544848412649
number of trees= 40
validation MSE= 39.24865238481688
validation r2= 0.989214223293078
train MSE= 32.95588828721234
train r2= 0.990690831329366
number of trees= 50
validation MSE= 40.26679457766181
validation r2= 0.988935933348887
train MSE= 33.879612866346465
train r2= 0.990345179894998
number of trees= 60
validation MSE= 39.77114425506524
validation r2= 0.9890664457594
train MSE= 34.16084029138152
train r2= 0.9902671058384045

Finding the best depth of each tree
depth= 5
validation MSE= 141.94882878392534
validation r2= 0.9610801355443937
train MSE= 140.4452579211163
train r2= 0.959978871143817
depth= 6
validation MSE= 91.86867795925444
validation r2= 0.974783968798172
train MSE= 87.50297430386225
train r2= 0.9794571882172862
depth= 7
validation MSE= 65.18295142590854
validation r2= 0.982073943043493
train MSE= 59.12932232480334
train r2= 0.9831675804244929
depth= 8
validation MSE= 40.6943890347693
validation r2= 0.988815223293533
train MSE= 34.18186448412649
train r2= 0.9902593803838695
depth= 9
validation MSE= 32.49854476989028
validation r2= 0.991800117186507
train MSE= 26.47922512721999
train r2= 0.99239876353595
depth= 10
validation MSE= 28.271913289675293
validation r2= 0.99082220593016573
train MSE= 21.707627952805887
train r2= 0.9938149984588585
depth= 11
validation MSE= 30.8145888961521
validation r2= 0.9917535477911613
train MSE= 24.9854142542483
train r2= 0.992877924534197
depth= 12
validation MSE= 39.38218436741496
validation r2= 0.9891403437645535
train MSE= 38.867969287528677
train r2= 0.9912896657386048
depth= 13
validation MSE= 27.864636839021376
validation r2= 0.9925611152838036
train MSE= 28.93871225916165
train r2= 0.9948326245635358
depth= 14
validation MSE= 26.350714072436272
validation r2= 0.99075916524352
train MSE= 19.079861631821824
train r2= 0.9945645538469462

Finding the best number of leaves
Samples per leaf= 1
validation MSE= 41.12379945069733
validation r2= 0.986845789366629
train MSE= 35.748382651084956
train r2= 0.9898162460441072
Samples per leaf= 2
validation MSE= 43.210283719636219
validation r2= 0.9881855672665947
train MSE= 39.226403684685884
train r2= 0.988879245347124
```

Inference: Random Forest is found to be a very good model for the given data, giving over 99% accuracy with the correct parameters.

1. Higher number of estimators (trees) give better results, but take more time to train, same can be said about the depth.
2. number of samples per leaf make little difference in the final results (but fewer mean more time to train, but better results).

Improvement Suggetions: Some other impurity measures can be tested out for training each tree

Tuning ELNET

```
In [63]:
# Defining a function for finding error in elnet across random train and validation sets
def eln_error(l1_penalty=50, l2_penalty=3, learning_rate=0.01):
    mse_train = r2_train = mse_validate = r2_validate = 0
    count = 3
    for i in range(count):
        validate, train = test_train_split(train_validate, random_state=i)
        Y_train = train['co2'].to_numpy()
        X_train = train.drop(['co2'], axis = 1).to_numpy()
        Y_validate = validate['co2'].to_numpy()
        X_validate = validate.drop(['co2'], axis = 1).to_numpy()
        W = elastic_net_predict(np.concatenate((np.ones((len(X_validate), 1)), X_validate), axis=1), W)
        mse_validate += MSE(Y_pred, Y_validate)
        r2_validate += rsquared_score(Y_pred, Y_validate)
        Y_pred = elastic_net_predict(np.concatenate((np.ones((len(X_train), 1)), X_train), axis=1), W)
        mse_train += MSE(Y_pred, Y_train)
        r2_train += rsquared_score(Y_pred, Y_train)
    print("validation MSE=", mse_validate/count)
    print("validation r2=", r2_validate/count)
    print("train MSE=", mse_train/count)
    print("train r2=", r2_train/count)

print("\nFinding a good l1_penalty")
for l1 in [0.01, 0.1, 1, 5, 10, 50, 100]:
    print("\nl1=", l1)
    eln_error(l1_penalty=l1)

print("\nFinding a good l2_penalty")
for l2 in [0.01, 0.1, 1, 5, 10, 50, 100]:
    print("\nl2=", l2)
    eln_error(l2_penalty=l2)

print("\nFinding a good learning_rate")
for lr in [0.01, 0.1, 1, 5, 10, 50, 100]:
    print("\nlr=", lr)
    eln_error(learning_rate=lr)

Finding a good l1_penalty

l1= 0.01
validation MSE= 33.6201644247844
validation r2= 0.9907671222369493
train MSE= 31.716471855757867
train r2= 0.990962716443866

l1= 0.1
validation MSE= 33.620232082629204
validation r2= 0.9907671037941753
train MSE= 31.7165721621597
train r2= 0.9909627024471291

l1= 1
validation MSE= 33.620911535344845
validation r2= 0.9907670184212525
train MSE= 31.718181350115758
train r2= 0.990962261761995

l1= 5
validation MSE= 33.62400791805499
validation r2= 0.9907649691809922
train MSE= 31.733943983667633
train r2= 0.990977515375805

l1= 10
validation MSE= 33.62805368843892
validation r2= 0.9907649691809922
train MSE= 31.733943983667633
train r2= 0.990977515375805

l1= 50
validation MSE= 33.66743399825922
validation r2= 0.990754206235399
train MSE= 31.812273903817644
train r2= 0.9909354279142724

l1= 100
validation MSE= 33.73420581483528
validation r2= 0.990735928628021
train MSE= 31.928558881817545
train r2= 0.9909822897771549

Finding a good l2_penalty

l2= 0.01
validation MSE= 33.443409180144984
validation r2= 0.990815958287257
train MSE= 31.43956188907629
train r2= 0.9910416471351399

l2= 0.1
validation MSE= 33.44294925655269
validation r2= 0.9908160801374729
train MSE= 31.443651365691665
train r2= 0.9910404816218547

l2= 1
validation MSE= 33.46326339634354
validation r2= 0.9908160408743222
train MSE= 31.50921441970553
train r2= 0.9910217964527589

l2= 5
validation MSE= 34.08419889496558
validation r2= 0.990639507663854
train MSE= 31.735400214372163
train r2= 0.9907891684756859

l2= 10
validation MSE= 36.0073081525224
validation r2= 0.990104389710899
train MSE= 34.478264815665575
train r2= 0.990175613546603

l2= 50
validation MSE= 88.206755316510081
validation r2= 0.975753346120085
train MSE= 87.85587241466884
train r2= 0.9749640553124743

l2= 100
validation MSE= 223.4528194529535
validation r2= 0.9385517214698342
train MSE= 223.13787987828075
train r2= 0.936411504672553

Finding a good learning_rate

lr= 0.01
validation MSE= 33.66743399825922
validation r2= 0.990754206235399
train MSE= 31.812273903817644
train r2= 0.9909354279142724

lr= 0.1
validation MSE= 33.62495393095919
validation r2= 0.9907657805626961
train MSE= 31.72347341727844
train r2= 0.99096073861806716

lr= 1
validation MSE= nan
validation r2= nan
train MSE= nan
train r2= nan

lr= 10
validation MSE= nan
validation r2= nan
train MSE= nan
train r2= nan

lr= 100
validation MSE= nan
validation r2= nan
train MSE= nan
train r2= nan
```

Inference: Elastic Net is also found to be a good model for the given dataset, with over 99% accuracy.

1. l1_penalty makes little difference in the model, as irrelevant features have already been eliminated
2. High l2_penalty leads to the shrinking of relevant features, leading to increase in error.
3. lower learning rate leads to slower learning but better accuracy, higher learning rate may even lead to nan values (errors too high for the system)

Improvement Suggestions: A basis function can be used to transform the points for better performance.

7. References

1. Libraries * **numpy**: <https://numpy.org/doc/stable/user/index.html> * **pandas**: https://pandas.pydata.org/docs/getting_started/index.html#getting-started * **matplotlib**: <https://matplotlib.org/stable/users/index> * **seaborn**: <https://seaborn.pydata.org/tutorial.html#user-guide-and-tutorial>
2. Support Vector Regression with Sequential Minimal Optimization
 - * **MSMD** Paper***: <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/smo-book.pdf>
 - * **SVNR** Theory***: <https://jin.mathworks.com/help/stats/understanding-support-vector-machine-regression.html>
2. Artificial Neural Network * **ANN Theory**: <https://towardsdatascience.com/math-neural-network-from-scratch-in-python-dda09f23ce65>
3. Random Forest Regression * None
4. Elastic Net Regression * **ELNET Theory**: <https://datawisdom.ca/paper/2013-JSMProceedings-ElasticNet.pdf>

NOTE: Accuracy refers to 100% if a model wherever referred to in this notebook.