

# final

December 12, 2024

## 0.1 Abstract

## 0.2 Package import

```
[161]: import os

os.environ["KERAS_BACKEND"] = "tensorflow"

import ast
import numpy as np

from tensorflow import keras
```

```
[162]: #from tensorflow.keras import ops
from tensorflow.keras import layers
import pandas as pd

from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
from rdkit import Chem, RDLogger
from rdkit.Chem import BondType
from rdkit.Chem.Draw import MolsToGridImage
from rdkit.Chem import Draw
from rdkit import Chem
from rdkit.Chem import rdmolops, AllChem
from tensorflow.keras.regularizers import l1_l2
RDLogger.DisableLog("rdApp.*")
```

```
[163]: import tensorflow as tf
print("TensorFlow version:", tf.__version__)
print("GPU available:", tf.config.list_physical_devices('GPU'))
print("GPU in use:", tf.test.gpu_device_name())
```

TensorFlow version: 2.16.2

GPU available: []

GPU in use:

### 0.3 Database pharsing

```
[164]: '''  
        read the entire dataset  
        '''  
  
df = pd.read_csv('dataset1.csv')  
df.drop([0,1,2,3,4], inplace=True)  
df=df.rename(columns = {'PUBCHEM_EXT_DATASOURCE_SMILES':  
    ↳ 'SMILES', 'PUBCHEM_ACTIVITY_OUTCOME': 'Activity', 'PUBCHEM_ACTIVITY_SCORE':  
    ↳ 'Score'})  
columns_to_drop = [col for col in df.columns if col not in ['SMILES',  
    ↳ 'Activity', 'Score', 'Potency', 'Efficacy']]  
df = df.drop(columns = columns_to_drop)  
#df=df.drop(['Unnamed: 3', 'Unnamed: 4', 'Unnamed: 5'], axis=1)  
df = df.dropna(subset=['SMILES'])  
  
df=df.fillna(0)  
print(df.head())  
print(df.info())
```

	SMILES	Activity	Score \
5	CNCC1=NC2=C(C=C(C=C2)C1)C(=N1)C3=CC=CN3	Inactive	0.0
6	CCSC(=NC1=CC=C(C=C1)C(F)(F)F)N.C1	Inactive	0.0
7	CCN(CC1=CC(=CC=C1)S(=O)(=O)[O-])C2=CC=C(C=C2)C...	Inactive	0.0
8	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0
9	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0

	Potency	Efficacy
5	0.0	0.0
6	0.0	0.0
7	0.0	0.0
8	0.0	0.0
9	0.0	0.0

```
<class 'pandas.core.frame.DataFrame'>  
Index: 342051 entries, 5 to 342072  
Data columns (total 5 columns):  
#   Column      Non-Null Count  Dtype  
---  -  
0   SMILES      342051 non-null object  
1   Activity    342051 non-null object  
2   Score       342051 non-null float64  
3   Potency     342051 non-null float64  
4   Efficacy    342051 non-null float64  
dtypes: float64(3), object(2)  
memory usage: 15.7+ MB  
None
```

```
[165]: valid_indices = []
# Loop through each SMILES string in the DataFrame
for i in range(len(df)):
    smiles = df.iloc[i]['SMILES'] # Use iloc for positional indexing

    # Convert SMILES to molecule
    mol = Chem.MolFromSmiles(smiles)

    # Check if the molecule is valid and has <= 50 atoms
    if mol is not None and mol.GetNumAtoms() <= 50:
        valid_indices.append(i)
# Filter the DataFrame to include only valid molecules
df_50 = df.iloc[valid_indices]
```

```
[166]: df_50
```

```
[166]:
```

	SMILES	Activity	Score	\
5	CNCC1=NC2=C(C=C(C=C2)C1)C(=N1)C3=CC=CN3	Inactive	0.0	
6	CCSC(=NC1=CC=C(C=C1)C(F)(F)F)N.Cl	Inactive	0.0	
8	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0	
9	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0	
10	C1CN(CCN1C2=NC(=NC3=CC=CC=C32)C4=CC=CS4)S(=O)(...	Inactive	0.0	
...	...	...	...	
342068	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342069	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342070	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342071	CC(=O)NC1=CC=C(C=C1)C(=O)N(CC2=CC=CC=C2)CC3=CC...	Inactive	0.0	
342072	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	

	Potency	Efficacy
5	0.0	0.0
6	0.0	0.0
8	0.0	0.0
9	0.0	0.0
10	0.0	0.0
...	...	...
342068	0.0	0.0
342069	0.0	0.0
342070	0.0	0.0
342071	0.0	0.0
342072	0.0	0.0

[341260 rows x 5 columns]

```
[167]: def is_charged(smiles):
    mol = Chem.MolFromSmiles(smiles)
    if not mol:
```

```

    return False # Invalid SMILES
    return any(atom.GetFormalCharge() != 0 for atom in mol.GetAtoms())

# Test the function
print(is_charged("CC1=C(SC(=C1C#N)NC(=O)C2=CC(C=C2)OC) [N+] (=O)"))

```

True

```
[168]: df_50['Charged'] = df_50['SMILES'].apply(is_charged)
```

```

uncharged = df_50[df_50['Charged'] == False]
uncharged

```

/var/folders/jn/kkchdcr94t50xrmycsvkq2x80000gn/T/ipykernel\_85584/162626946.py:1:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df_50['Charged'] = df_50['SMILES'].apply(is_charged)
```

```
[168]:
```

	SMILES	Activity	Score	\
5	CNCC1=NC2=C(C=C(C=C2)C1)C(=N1)C3=CC=CN3	Inactive	0.0	
6	CCSC(=NC1=CC=C(C=C1)C(F)(F)F)N.C1	Inactive	0.0	
8	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0	
9	CC1=CC=C(C=C1)S(=O)(=O)N2CCN(CC2)C3=NC(=NC4=CC...	Inactive	0.0	
10	C1CN(CCN1C2=NC(=NC3=CC=CC=C32)C4=CC=CS4)S(=O)(...	Inactive	0.0	
...	...	...	...	
342068	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342069	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342070	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	
342071	CC(=O)NC1=CC=C(C=C1)C(=O)N(CC2=CC=CC=C2)CC3=CC...	Inactive	0.0	
342072	CC(=O)NC1=CC=C(C=C1)OCC2=C(C=CC(=C2)CN(CC3=CC=...	Inactive	0.0	

	Potency	Efficacy	Charged
5	0.0	0.0	False
6	0.0	0.0	False
8	0.0	0.0	False
9	0.0	0.0	False
10	0.0	0.0	False
...	...	...	...
342068	0.0	0.0	False
342069	0.0	0.0	False
342070	0.0	0.0	False
342071	0.0	0.0	False
342072	0.0	0.0	False

[322199 rows x 6 columns]

```
[169]: # Picking all "Active" molecules from the dataset
active_df = uncharged[uncharged['Activity'] == 'Active']
active_df.info()

# Picking all "Inactive" molecules from the dataset
inactive_df = uncharged[uncharged['Activity'] == 'Inactive']
inactive_df.info()

# Randomly sample from inactive_df to match the size of active_df
inactive_sampled = inactive_df.sample(n=len(active_df), random_state=42)

# Combine the active and sampled inactive molecules
balanced_df = pd.concat([active_df, inactive_sampled])

# Shuffle the combined dataset
balanced_df = balanced_df.sample(frac=1, random_state=42).reset_index(drop=True)

balanced_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 6273 entries, 13 to 341825
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   SMILES      6273 non-null   object
1   Activity    6273 non-null   object
2   Score       6273 non-null   float64
3   Potency     6273 non-null   float64
4   Efficacy    6273 non-null   float64
5   Charged     6273 non-null   bool
dtypes: bool(1), float64(3), object(2)
memory usage: 300.2+ KB
<class 'pandas.core.frame.DataFrame'>
Index: 304069 entries, 5 to 342072
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   SMILES      304069 non-null object
1   Activity    304069 non-null object
2   Score       304069 non-null float64
3   Potency     304069 non-null float64
4   Efficacy    304069 non-null float64
5   Charged     304069 non-null bool
dtypes: bool(1), float64(3), object(2)
memory usage: 14.2+ MB
<class 'pandas.core.frame.DataFrame'>
```

```

RangeIndex: 12546 entries, 0 to 12545
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   SMILES      12546 non-null  object
1   Activity    12546 non-null  object
2   Score       12546 non-null  float64
3   Potency     12546 non-null  float64
4   Efficacy    12546 non-null  float64
5   Charged     12546 non-null  bool
dtypes: bool(1), float64(3), object(2)
memory usage: 502.5+ KB

```

```
[170]: filtered_df = balanced_df
filtered_df
```

```
[170]:
```

	SMILES	Activity	Score	\
0	<chem>CC1=C(C=CC=C1Br)NC(=O)C2=C(C=CS2)N3C=CC=C3</chem>	Active	82.0	
1	<chem>CCCCC(C(C)CC(=O)NC1CCCC1)C(=O)O</chem>	Active	43.0	
2	<chem>CC1=CC=C(C=C1)S(=O)(=O)NC2=NN3C(C=C(NC3=N2)C)C...</chem>	Active	41.0	
3	<chem>CC1=CC(=O)OC2=C1C=C(C=C2)OCC(=O)NC3=CC=CC(=C3)...</chem>	Inactive	0.0	
4	<chem>CC1=CC(=C(N1C)C)C(=O)COC(=O)C23CC4CC(C2)CC(C4)...</chem>	Inactive	0.0	
...	...	...	...	
12541	<chem>C1CN(CCN1C(=O)C2=CC=CC=C2CC3=CC=CC=C3)S(=O)(=O)...</chem>	Inactive	0.0	
12542	<chem>C1=CC=C(C=C1)OC2=NC=NC(=C2)N3C=NC=N3</chem>	Active	64.0	
12543	<chem>CC1=C(C(=CC=C1)N2CCN(CC2)C3=NC4=CC=CC=C4C(=O)N...</chem>	Active	42.0	
12544	<chem>CCC(C)NC(=O)CSC1=NC2=CC=CC=C2C3=NC(C(=O)N31)C4...</chem>	Active	42.0	
12545	<chem>CC(C)C1=CC=C(C=C1)S(=O)(=O)NC2CCCC2</chem>	Inactive	0.0	
...	...	...	...	
12541				Potency
0				8.9125
1				12.5893
2				22.3872
3				0.0000
4				0.0000
...				...
12541				0.0000
12542				2.8184
12543				17.7828
12544				15.8489
12545				0.0000
				Efficacy
0				140.7280
1				136.6590
2				166.6580
3				0.0000
4				0.0000
...				...
12541				0.0000
12542				74.9734
12543				126.5240
12544				139.3040
12545				0.0000
				Charged
0				False
1				False
2				False
3				False
4				False
...				...
12541				False
12542				False
12543				False
12544				False
12545				False

[12546 rows x 6 columns]

## 0.4 Parameter setting

```
[171]: '''  
scan through all the molecules to obtain unique atom types  
'''  
  
smiles = filtered_df['SMILES'].tolist()  
search_elements=[]  
for smile in smiles:  
    mol = Chem.MolFromSmiles(smile)  
    atoms = list(set([atom.GetSymbol() for atom in mol.GetAtoms()]))  
    search_elements += atoms  
    search_elements = list(set(search_elements))  
search_elements.append("H")  
print(search_elements)
```

```
['C', 'F', 'N', 'I', 'O', 'P', 'Br', 'B', 'S', 'Cl', 'As', 'H']
```

```
[172]: '''  
Setting up the atom mapping and bond mapping.  
Code adopted from https://keras.io/examples/generative/molecule\_generation/  
'''  
  
SMILE_CHARSET = str(search_elements)  
bond_mapping = {"SINGLE": 0, "DOUBLE": 1, "TRIPLE": 2, "AROMATIC": 3}  
bond_mapping.update(  
    {0: BondType.SINGLE, 1: BondType.DOUBLE, 2: BondType.TRIPLE, 3: BondType.  
    ↪AROMATIC}  
)  
SMILE_CHARSET = ast.literal_eval(SMILE_CHARSET)  
  
MAX_MOLSIZE = max(filtered_df['SMILES'].str.len())  
SMILE_to_index = dict((c, i) for i, c in enumerate(SMILE_CHARSET))  
index_to_SMILE = dict((i, c) for i, c in enumerate(SMILE_CHARSET))  
atom_mapping = dict(SMILE_to_index)  
atom_mapping.update(index_to_SMILE)  
print(atom_mapping)  
print("Max molecule size: {}".format(MAX_MOLSIZE))  
print("Character set Length: {}".format(len(SMILE_CHARSET)))
```

```
{'C': 0, 'F': 1, 'N': 2, 'I': 3, 'O': 4, 'P': 5, 'Br': 6, 'B': 7, 'S': 8, 'Cl':  
9, 'As': 10, 'H': 11, 0: 'C', 1: 'F', 2: 'N', 3: 'I', 4: 'O', 5: 'P', 6: 'Br',  
7: 'B', 8: 'S', 9: 'Cl', 10: 'As', 11: 'H'}
```

```
Max molecule size: 117
```

```
Character set Length: 12
```

## 0.5 Hyperparameters

```
[173]: '''  
Defining the Hyperparameters of the model  
'''  
  
NUM_ATOMS = 50 #Max number of atoms  
ATOM_DIM = len(SMILE_CHARSET) # Number of atom types  
BOND_DIM = 5 # Number of bond types
```

## 0.6 Molecule featurization

```
[174]: '''  
Defining functions to convert smiles string into node graph and recover  
↳ molecule structure from it.  
Code referenced from: https://keras.io/examples/generative/molecule\_generation/  
'''  
  
def smiles_to_graph(smiles):  
    '''  
    Reference: https://keras.io/examples/generative/wgan-graphs/  
    '''  
    # Converts SMILES to molecule object  
    molecule = Chem.MolFromSmiles(smiles)  
    #molecule = Chem.AddHs(molecule)  
    # Initialize adjacency and feature tensor  
    adjacency = np.zeros((BOND_DIM, NUM_ATOMS, NUM_ATOMS), "float32")  
    features = np.zeros((NUM_ATOMS, ATOM_DIM), "float32")  
  
    # loop over each atom in molecule  
    for atom in molecule.GetAtoms():  
        i = atom.GetIdx()  
        atom_type = atom_mapping[atom.GetSymbol()]  
        features[i] = np.eye(ATOM_DIM)[atom_type]  
        # loop over one-hop neighbors  
        for neighbor in atom.GetNeighbors():  
            j = neighbor.GetIdx()  
            bond = molecule.GetBondBetweenAtoms(i, j)  
            bond_type_idx = bond_mapping[bond.GetBondType().name]  
            adjacency[bond_type_idx, [i, j], [j, i]] = 1  
  
        # Where no bond, add 1 to last channel (indicating "non-bond")  
        # Notice: channels-first  
        adjacency[-1, np.sum(adjacency, axis=0) == 0] = 1  
  
    # Where no atom, add 1 to last column (indicating "non-atom")
```



```

features[np.where(np.sum(features, axis=1) == 0)[0], -1] = 1

return adjacency, features

def graph_to_molecule(adjacency, features):
    # RWMol is a molecule object intended to be edited
    molecule = Chem.RWMol()
    # Remove "no atoms" & atoms with no bonds
    keep_idx = np.where(
        (np.argmax(features, axis=1) != ATOM_DIM - 1)
        & (np.sum(adjacency[:-1], axis=(0, 1)) > 0))[0]

    features = features[keep_idx]
    adjacency = adjacency[:, keep_idx][:, :, keep_idx]

    # Add atoms to molecule
    for atom_type_idx in np.argmax(features, axis=1):
        atom = Chem.Atom(atom_mapping[atom_type_idx])
        _ = molecule.AddAtom(atom)

    added_bonds = set()
    (bonds_ij, atoms_i, atoms_j) = np.where(np.triu(adjacency) == 1)
    for (bond_ij, atom_i, atom_j) in zip(bonds_ij, atoms_i, atoms_j):
        if atom_i == atom_j or bond_ij == BOND_DIM - 1:
            continue
        bond_type = bond_mapping.get(bond_ij, None)
        if (atom_i, atom_j) in added_bonds or (atom_j, atom_i) in added_bonds:
            continue
        molecule.AddBond(int(atom_i), int(atom_j), bond_type)
        added_bonds.add((atom_i, atom_j))

    # Sanitize without Kekulization
    try:
        Chem.SanitizeMol(molecule, sanitizeOps=Chem.SanitizeFlags.SANITIZE_ALL_
    ↪ Chem.SanitizeFlags.SANITIZE_KEKULIZE)
    except Exception as e:
        print(f"Sanitization failed: {e}")
        return None

    # Add explicit hydrogens
    molecule_with_h = Chem.AddHs(molecule)

    # Fix aromaticity in aromatic rings
    for atom in molecule_with_h.GetAtoms():
        if atom.GetIsAromatic():

```

```

        atom.SetIsAromatic(False) # Clear aromaticity if needed

# Force Kekulization to alternate bond orders in aromatic rings
try:
    Chem.Kekulize(molecule_with_h, clearAromaticFlags=True)
except Chem.KekulizeException as e:
    print(f"Kekulization failed: {e}")
    return molecule_with_h # Return molecule without Kekulé bonds

return molecule_with_h

```

## 0.7 Building model

```

[175]: '''
        Defining GCN
        Reference: https://keras.io/examples/generative/wgan-graphs/
        The Encoder takes as input a molecule's graph adjacency matrix and feature_
        ↪matrix.
    '''
    class RelationalGraphConvLayer(keras.layers.Layer):
        def __init__(
            self,
            units=128,
            activation="relu",
            use_bias=False,
            kernel_initializer="glorot_uniform",
            bias_initializer="zeros",
            kernel_regularizer=None,
            bias_regularizer=None,
            **kwargs
        ):
            super().__init__(**kwargs)

            self.units = units
            self.activation = keras.activations.get(activation)
            self.use_bias = use_bias
            self.kernel_initializer = keras.initializers.get(kernel_initializer)
            self.bias_initializer = keras.initializers.get(bias_initializer)
            self.kernel_regularizer = keras.regularizers.get(kernel_regularizer)
            self.bias_regularizer = keras.regularizers.get(bias_regularizer)

        def build(self, input_shape):
            bond_dim = input_shape[0][1]
            atom_dim = input_shape[1][2]

            self.kernel = self.add_weight(
                shape=(bond_dim, atom_dim, self.units),

```

```

        initializer=self.kernel_initializer,
        regularizer=self.kernel_regularizer,
        trainable=True,
        name="W",
        dtype=tf.float32,
    )

    if self.use_bias:
        self.bias = self.add_weight(
            shape=(bond_dim, 1, self.units),
            initializer=self.bias_initializer,
            regularizer=self.bias_regularizer,
            trainable=True,
            name="b",
            dtype=tf.float32,
        )

    self.built = True

    def call(self, inputs, training=False):
        adjacency, features = inputs
        # Aggregate information from neighbors
        x = tf.matmul(adjacency, features[:, None, :, :])
        # Apply linear transformation
        x = tf.matmul(x, self.kernel)
        if self.use_bias:
            x += self.bias
        # Reduce bond types dim
        x_reduced = tf.reduce_sum(x, axis=1)
        # Apply non-linear transformation
        return self.activation(x_reduced)

```

## 0.8 Build the Encoder and Decoder

```

[176]: '''
    defining function to build encoder and decoder.
    Code adopted and modified from https://keras.io/examples/generative/
    ↪ molecule_generation/
    '''

    def get_encoder(gconv_units, latent_dim, adjacency_shape, feature_shape,
    ↪ dense_units, dropout_rate, regularizer=None):
        adjacency = keras.layers.Input(shape=adjacency_shape,
        ↪ name="adjacency_input")
        features = keras.layers.Input(shape=feature_shape, name="feature_input")
        scores = keras.layers.Input(shape=(1,), name="score_input") # Conditional
        ↪ input (scalar)

```

```

# Graph convolution layers
features_transformed = features
for units in gconv_units:
    features_transformed = RelationalGraphConvLayer(units)(
        adjacency, features_transformed
    )

# Reduce 2D representation to 1D
x = keras.layers.GlobalAveragePooling1D()(features_transformed)

# Concatenate the score (condition) to the reduced graph representation
x = keras.layers.Concatenate()([x, scores])

# Fully connected layers
for units in dense_units:
    x = layers.Dense(units, activation="relu",
        ↪kernel_regularizer=regularizer)(x)
    x = layers.Dropout(dropout_rate)(x)

# Latent space
z_mean = layers.Dense(latent_dim, name="z_mean")(x)
z_log_var = layers.Dense(latent_dim, name="z_log_var")(x)

# Create encoder model
encoder = keras.Model(inputs=[adjacency, features, scores],
    ↪outputs=[z_mean, z_log_var], name="encoder")
encoder.summary()
return encoder

class SymmetrizeLayer(layers.Layer):
    def call(self, x):
        return (x + tf.transpose(x, (0, 1, 3, 2))) / 2

def get_decoder(dense_units, latent_dim, adjacency_shape, feature_shape,
    ↪dropout_rate, regularizer=None):
    latent_input = keras.Input(shape=(latent_dim,), name="latent_input")
    scores = keras.Input(shape=(1,), name="score_input") # Conditional input
    ↪(scalar)

    # Concatenate latent input with the conditional score
    x = keras.layers.Concatenate()([latent_input, scores])

    # Dense layers
    for units in dense_units:

```

```

        x = keras.layers.Dense(units, activation="tanh",
kernel_regularizer=regularizer)(x)
        x = keras.layers.Dropout(dropout_rate)(x)

        # Adjacency reconstruction
        adj_output = keras.layers.Dense(tf.math.reduce_prod(adjacency_shape).
numpy().astype(int))(x)
        adj_output = keras.layers.Reshape(adjacency_shape)(adj_output)
        adj_output = SymmetrizeLayer()(adj_output)
        adj_output = keras.layers.Softmax(axis=1)(adj_output)

        # Feature reconstruction
        feat_output = keras.layers.Dense(tf.math.reduce_prod(feature_shape).numpy().
astype(int))(x)
        feat_output = keras.layers.Reshape(feature_shape)(feat_output)
        feat_output = keras.layers.Softmax(axis=2)(feat_output)

        # Create decoder model
        decoder = keras.Model(inputs=[latent_input, scores], outputs=[adj_output,
feat_output], name="decoder")
        decoder.summary()
        return decoder

```

## 0.9 Build the VAE

```

[177]: '''
defining the VAE
Code adopted and modified from https://keras.io/examples/generative/
molecule_generation/
'''

class VAE(keras.Model):
    def __init__(self, encoder, decoder, beta=1.0, **kwargs):
        super(VAE, self).__init__(**kwargs)
        self.encoder = encoder
        self.decoder = decoder
        self.beta = beta

    def call(self, inputs):
        adjacency, features, scores = inputs
        z_mean, z_log_var = self.encoder([adjacency, features, scores])
        z = self.reparameterize(z_mean, z_log_var)
        return self.decoder([z, scores])

    def sampling(self, args):
        """
        Reparameterization trick: Sample from a Gaussian distribution using

```

```

        z = z_mean + epsilon * exp(z_log_var / 2), where epsilon is sampled
        ↪from N(0, 1).
        """
        z_mean, z_log_var = args
        batch = tf.shape(z_mean)[0]
        dim = tf.shape(z_mean)[1]
        epsilon = tf.keras.backend.random_normal(shape=(batch, dim)) #
        ↪Standard normal noise
        return z_mean + tf.exp(0.5 * z_log_var) * epsilon

```

## 0.10 Model training

```

[178]: '''
        splitting the dataset into training and testing
        '''

train, test = train_test_split(filtered_df, test_size=0.2, random_state=42)
train_df, val_df = train_test_split(train, test_size=0.2, random_state=42)
train_df.reset_index(drop=True, inplace=True)
val_df.reset_index(drop=True, inplace=True)
test.reset_index(drop=True, inplace=True)

adj_train, fea_train, score_train = [], [], []
adj_val, fea_val, score_val = [], [], []

for idx in range(len(train_df)):
    adjacency, features = smiles_to_graph(train_df.loc[idx]["SMILES"])
    score = train_df.loc[idx]["Score"]
    adj_train.append(adjacency)
    fea_train.append(features)
    score_train.append(score)

for idx in range(len(val_df)):
    adjacency, features = smiles_to_graph(val_df.loc[idx]["SMILES"])
    score = val_df.loc[idx]["Score"]
    adj_val.append(adjacency)
    fea_val.append(features)
    score_val.append(score)

adj_train = np.array(adj_train)
fea_train = np.array(fea_train)
score_train_ = np.array(score_train).reshape(-1,1)

adj_val = np.array(adj_val)
fea_val = np.array(fea_val)
score_val_ = np.array(score_val).reshape(-1,1)

```

```
[179]: from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
score_train_n = scaler.fit_transform(score_train_)
```

```
score_val_n = scaler.transform(score_val_)
```

```
[180]: print(adj_train.shape)
print(fea_train.shape)
print(score_train_.shape)
print(adj_val.shape)
print(fea_val.shape)
print(score_val_.shape)
```

```
(8028, 5, 50, 50)
```

```
(8028, 50, 12)
```

```
(8028, 1)
```

```
(2008, 5, 50, 50)
```

```
(2008, 50, 12)
```

```
(2008, 1)
```

```
[181]: print(np.max(score_train_n))
```

```
0.9999999999999999
```

```
[182]: #Hyperparameters
```

```
BATCH_SIZE = 64
```

```
EPOCHS = 20
```

```
VAE_LR = 3e-4 # changed to 1e-3
```

```
LATENT_DIM = 256 # Size of the latent space
```

```
[183]: '''
compiling the VAE
'''
```

```
encoder = get_encoder(
    gconv_units=[16],
    adjacency_shape=(BOND_DIM, NUM_ATOMS, NUM_ATOMS),
    feature_shape=(NUM_ATOMS, ATOM_DIM),
    latent_dim=LATENT_DIM,
    dense_units=[128, 256, 512],
    dropout_rate=0,
    regularizer=l1_l2(l1=1e-6, l2=1e-3)
)
decoder = get_decoder(
    dense_units=[128, 256, 512],
    dropout_rate=0.3,
    latent_dim=LATENT_DIM,
```

```

adjacency_shape=(BOND_DIM, NUM_ATOMS, NUM_ATOMS),
feature_shape=(NUM_ATOMS, ATOM_DIM),
regularizer=l1_l2(l1=1e-4, l2=1e-2)
)
vae = VAE(encoder, decoder)

vae.compile(optimizer=keras.optimizers.Adam(learning_rate=VAE_LR))

```

Model: "encoder"

Layer (type)	Output Shape	Param #	Connected to
adjacency_input ( <a href="#">InputLayer</a> )	(None, 5, 50, 50)	0	-
feature_input ( <a href="#">InputLayer</a> )	(None, 50, 12)	0	-
relational_graph_c... ( <a href="#">RelationalGraphCo...</a> )	(None, 50, 16)	960	adjacency_input[... feature_input[0]...
global_average_poo... ( <a href="#">GlobalAveragePool...</a> )	(None, 16)	0	relational_graph...
score_input ( <a href="#">InputLayer</a> )	(None, 1)	0	-
concatenate_8 ( <a href="#">Concatenate</a> )	(None, 17)	0	global_average_p... score_input[0][0]
dense_28 ( <a href="#">Dense</a> )	(None, 128)	2,304	concatenate_8[0]...
dropout_20 ( <a href="#">Dropout</a> )	(None, 128)	0	dense_28[0][0]
dense_29 ( <a href="#">Dense</a> )	(None, 256)	33,024	dropout_20[0][0]
dropout_21 ( <a href="#">Dropout</a> )	(None, 256)	0	dense_29[0][0]
dense_30 ( <a href="#">Dense</a> )	(None, 512)	131,584	dropout_21[0][0]
dropout_22 ( <a href="#">Dropout</a> )	(None, 512)	0	dense_30[0][0]
z_mean ( <a href="#">Dense</a> )	(None, 256)	131,328	dropout_22[0][0]



z\_log\_var ([Dense](#)) (None, 256) 131,328 dropout\_22[0][0]

Total params: 430,528 (1.64 MB)

Trainable params: 430,528 (1.64 MB)

Non-trainable params: 0 (0.00 B)

Model: "decoder"

Layer (type)	Output Shape	Param #	Connected to
latent_input ( <a href="#">InputLayer</a> )	(None, 256)	0	-
score_input ( <a href="#">InputLayer</a> )	(None, 1)	0	-
concatenate_9 ( <a href="#">Concatenate</a> )	(None, 257)	0	latent_input[0][... score_input[0][0]
dense_31 ( <a href="#">Dense</a> )	(None, 128)	33,024	concatenate_9[0]...
dropout_23 ( <a href="#">Dropout</a> )	(None, 128)	0	dense_31[0][0]
dense_32 ( <a href="#">Dense</a> )	(None, 256)	33,024	dropout_23[0][0]
dropout_24 ( <a href="#">Dropout</a> )	(None, 256)	0	dense_32[0][0]
dense_33 ( <a href="#">Dense</a> )	(None, 512)	131,584	dropout_24[0][0]
dropout_25 ( <a href="#">Dropout</a> )	(None, 512)	0	dense_33[0][0]
dense_34 ( <a href="#">Dense</a> )	(None, 12500)	6,412,500	dropout_25[0][0]
reshape_8 ( <a href="#">Reshape</a> )	(None, 5, 50, 50)	0	dense_34[0][0]
dense_35 ( <a href="#">Dense</a> )	(None, 600)	307,800	dropout_25[0][0]

```

symmetrize_layer_4      (None, 5, 50, 50)          0  reshape_8[0][0]
(SymmetrizeLayer)

reshape_9 (Reshape)      (None, 50, 12)            0  dense_35[0][0]

softmax_8 (Softmax)      (None, 5, 50, 50)          0  symmetrize_layer...

softmax_9 (Softmax)      (None, 50, 12)            0  reshape_9[0][0]

```

Total params: 6,917,932 (26.39 MB)

Trainable params: 6,917,932 (26.39 MB)

Non-trainable params: 0 (0.00 B)

```

[184]: val_loss_list = []
       train_loss_list = []
       kl_theshold = 1.0

[185]: train_dataset = tf.data.Dataset.from_tensor_slices((adj_train, fea_train, ↵
           ↵score_train)).batch(BATCH_SIZE)
       val_dataset = tf.data.Dataset.from_tensor_slices((adj_val, fea_val, ↵
           ↵score_val)).batch(BATCH_SIZE)

[186]: for epoch in range(EPOCHS):
       print(f"Epoch {epoch + 1}/{EPOCHS}")
       if epoch < 10:
           beta = 0.05
       else:
           beta = epoch*0.01
       # Training Loop
       train_loss = 0
       for (adjacency, features, scores) in train_dataset:
           with tf.GradientTape() as tape:
               # Forward pass
               z_mean, z_log_var = vae.encoder([adjacency, features, scores])
               z = vae.sampling([z_mean, z_log_var])
               adj_reconstruction, feature_reconstruction = vae.decoder([z, ↵
           ↵scores])

               # Compute losses
               adj_loss = tf.reduce_mean(
                   tf.reduce_sum(keras.losses.binary_crossentropy(adjacency, ↵
           ↵adj_reconstruction), axis=(1, 2))

```

```

    )
    feat_loss = tf.reduce_mean(
        tf.reduce_sum(keras.losses.categorical_crossentropy(features,
↪feature_reconstruction), axis=1)
    )
    reconstruction_loss = adj_loss + feat_loss
    kl_loss = -0.5 * tf.reduce_mean(
        tf.reduce_sum(1 + z_log_var - tf.square(z_mean) - tf.
↪exp(z_log_var), axis=1)
    )
    total_loss = reconstruction_loss + beta * kl_loss

    # Backpropagation
    grads = tape.gradient(total_loss, vae.trainable_weights)
    vae.optimizer.apply_gradients(zip(grads, vae.trainable_weights))

    train_loss += total_loss

train_loss /= len(train_dataset)
train_loss_list.append(train_loss)

    print(f"Train Loss: {train_loss.numpy()}, KL Loss: {kl_loss.numpy()},
↪Reconstruction Loss: {reconstruction_loss.numpy()}")

    # Validation Loop
    val_loss = 0
    for (val_adjacency, val_features, val_scores) in val_dataset:
        # Forward pass
        z_mean, z_log_var = vae.encoder([val_adjacency, val_features,
↪val_scores])
        z = vae.sampling([z_mean, z_log_var])
        val_adj_reconstruction, val_feat_reconstruction = vae.decoder([z,
↪val_scores])

        # Compute losses
        val_adj_loss = tf.reduce_mean(
            tf.reduce_sum(keras.losses.binary_crossentropy(val_adjacency,
↪val_adj_reconstruction), axis=(1, 2))
        )
        val_feat_loss = tf.reduce_mean(
            tf.reduce_sum(keras.losses.categorical_crossentropy(val_features,
↪val_feat_reconstruction), axis=1)
        )
        val_reconstruction_loss = val_adj_loss + val_feat_loss
        val_kl_loss = -0.5 * tf.reduce_mean(

```

```

        tf.reduce_sum(1 + z_log_var - tf.square(z_mean) - tf.
↪exp(z_log_var), axis=1)
    )
    val_total_loss = val_reconstruction_loss + beta * val_kl_loss

    val_loss += val_total_loss

val_loss /= len(val_dataset)
val_loss_list.append(val_loss)

# Adjust beta if KL loss is very low
if kl_loss < kl_theshold:
    beta = 0.05
print(f"Validation Loss: {val_loss.numpy()}, KL Loss: {val_kl_loss.
↪numpy()}, Reconstruction Loss: {val_reconstruction_loss.numpy()}")
print('BETA is: ', beta)

```

Epoch 1/20

2024-12-12 18:30:03.743240: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 68.51957702636719, KL Loss: 6.51322603225708, Reconstruction Loss:  
36.936180114746094

2024-12-12 18:30:04.408546: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 38.39116668701172, KL Loss: 7.85448694229126, Reconstruction  
Loss: 35.157981872558594

BETA is: 0.05

Epoch 2/20

2024-12-12 18:30:12.789587: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 38.01901626586914, KL Loss: 7.983117580413818, Reconstruction Loss:  
35.29029846191406

2024-12-12 18:30:13.312265: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 36.62051010131836, KL Loss: 9.628422737121582, Reconstruction  
Loss: 33.35990524291992

BETA is: 0.05

Epoch 3/20

2024-12-12 18:30:21.823862: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 33.5779914855957, KL Loss: 12.101990699768066, Reconstruction Loss:  
31.358749389648438

2024-12-12 18:30:22.342605: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 32.068939208984375, KL Loss: 13.179827690124512, Reconstruction  
Loss: 29.782800674438477  
BETA is: 0.05  
Epoch 4/20

2024-12-12 18:30:30.727471: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 31.387109756469727, KL Loss: 13.011499404907227, Reconstruction  
Loss: 29.32843589782715

2024-12-12 18:30:31.229339: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.902626037597656, KL Loss: 15.3453369140625, Reconstruction  
Loss: 27.747791290283203  
BETA is: 0.05  
Epoch 5/20

2024-12-12 18:30:40.645041: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.642066955566406, KL Loss: 15.282196044921875, Reconstruction  
Loss: 29.32878875732422

2024-12-12 18:30:41.190827: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.674419403076172, KL Loss: 17.057682037353516, Reconstruction  
Loss: 27.417327880859375  
BETA is: 0.05  
Epoch 6/20

2024-12-12 18:30:49.656011: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.340274810791016, KL Loss: 15.912251472473145, Reconstruction  
Loss: 28.87799644470215

2024-12-12 18:30:50.170870: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.2147159576416, KL Loss: 18.694887161254883, Reconstruction  
Loss: 27.021533966064453  
BETA is: 0.05  
Epoch 7/20

2024-12-12 18:30:58.656666: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.166858673095703, KL Loss: 16.057947158813477, Reconstruction  
Loss: 28.732450485229492

2024-12-12 18:30:59.186481: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.019132614135742, KL Loss: 17.847206115722656, Reconstruction  
Loss: 26.88463592529297

BETA is: 0.05

Epoch 8/20

2024-12-12 18:31:07.878724: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 29.954687118530273, KL Loss: 15.391514778137207, Reconstruction  
Loss: 28.946693420410156

2024-12-12 18:31:08.409757: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 29.946964263916016, KL Loss: 16.681718826293945, Reconstruction  
Loss: 26.971174240112305

BETA is: 0.05

Epoch 9/20

2024-12-12 18:31:16.945733: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 29.84228515625, KL Loss: 16.165109634399414, Reconstruction Loss:  
29.082962036132812

2024-12-12 18:31:17.472882: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.09977912902832, KL Loss: 17.07674217224121, Reconstruction  
Loss: 26.848041534423828

BETA is: 0.05

Epoch 10/20

2024-12-12 18:31:25.910867: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 29.79227638244629, KL Loss: 16.24569320678711, Reconstruction Loss:  
28.732189178466797

2024-12-12 18:31:26.418563: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 29.83408546447754, KL Loss: 17.454748153686523, Reconstruction  
Loss: 26.80592918395996

BETA is: 0.05

Epoch 11/20

2024-12-12 18:31:34.882094: W tensorflow/core/framework/local\_rendevvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.638765335083008, KL Loss: 11.345499992370605, Reconstruction  
Loss: 28.695213317871094

2024-12-12 18:31:35.423871: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.51316261291504, KL Loss: 13.088162422180176, Reconstruction  
Loss: 26.347801208496094

BETA is: 0.1

Epoch 12/20

2024-12-12 18:31:44.083496: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 31.29293441772461, KL Loss: 12.0257568359375, Reconstruction Loss:  
28.44754981994629

2024-12-12 18:31:44.627711: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.589923858642578, KL Loss: 12.816065788269043, Reconstruction  
Loss: 26.69342613220215

BETA is: 0.11

Epoch 13/20

2024-12-12 18:31:53.110507: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.741384506225586, KL Loss: 10.480046272277832, Reconstruction  
Loss: 29.36916732788086

2024-12-12 18:31:53.634911: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.553064346313477, KL Loss: 11.884610176086426, Reconstruction  
Loss: 27.17499542236328

BETA is: 0.12

Epoch 14/20

2024-12-12 18:32:02.139204: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.516401290893555, KL Loss: 10.94199275970459, Reconstruction Loss:  
28.466699600219727

2024-12-12 18:32:02.662070: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.66318130493164, KL Loss: 11.293120384216309, Reconstruction  
Loss: 26.394275665283203

BETA is: 0.13

Epoch 15/20

2024-12-12 18:32:11.226403: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.85893440246582, KL Loss: 13.422137260437012, Reconstruction Loss:  
28.727540969848633

2024-12-12 18:32:11.763504: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 31.54447364807129, KL Loss: 11.864672660827637, Reconstruction  
Loss: 28.32206153869629

BETA is: 0.14

Epoch 16/20

2024-12-12 18:32:20.163615: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.816301345825195, KL Loss: 8.611360549926758, Reconstruction Loss:  
28.985681533813477

2024-12-12 18:32:20.679793: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.67209243774414, KL Loss: 9.73166561126709, Reconstruction  
Loss: 27.155027389526367

BETA is: 0.15

Epoch 17/20

2024-12-12 18:32:29.136181: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.708045959472656, KL Loss: 10.172674179077148, Reconstruction  
Loss: 28.612720489501953

2024-12-12 18:32:29.659116: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.850566864013672, KL Loss: 11.591927528381348, Reconstruction  
Loss: 26.911741256713867

BETA is: 0.16

Epoch 18/20

2024-12-12 18:32:38.161630: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.703624725341797, KL Loss: 8.000271797180176, Reconstruction Loss:  
28.36505126953125

2024-12-12 18:32:38.686689: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.541950225830078, KL Loss: 8.589282989501953, Reconstruction  
Loss: 26.045917510986328

BETA is: 0.17

Epoch 19/20

2024-12-12 18:32:47.286038: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.978403091430664, KL Loss: 7.179697513580322, Reconstruction Loss:  
28.462646484375



2024-12-12 18:32:47.823747: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Validation Loss: 30.83968734741211, KL Loss: 7.570359706878662, Reconstruction  
Loss: 26.57512855529785

BETA is: 0.18

Epoch 20/20

2024-12-12 18:32:56.395421: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

Train Loss: 30.546899795532227, KL Loss: 7.951568603515625, Reconstruction Loss:  
28.570207595825195

Validation Loss: 30.662322998046875, KL Loss: 9.32971477508545, Reconstruction  
Loss: 27.113357543945312

BETA is: 0.19

2024-12-12 18:32:56.946421: W tensorflow/core/framework/local\_rendezvous.cc:404]  
Local rendezvous is aborting with status: OUT\_OF\_RANGE: End of sequence

```
[187]: '''  
Checking the model's ability to reconstruct a molecule from the training dataset  
'''
```

```
i=10  
adjacency_check, features_check = smiles_to_graph(train_df.loc[i]["SMILES"])  
score_check = [train_df.loc[i]["Score"]]  
molobj = Chem.MolFromSmiles(train_df.loc[i]["SMILES"])  
adj0 = np.expand_dims(adjacency_check,axis=0)  
feature0 = np.expand_dims(features_check,axis=0)  
score0 = np.expand_dims(score_check,axis=0)  
print(adj0.shape)  
print(feature0.shape)  
print(score0.shape)
```

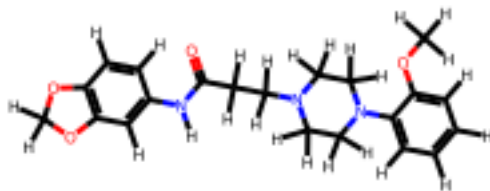
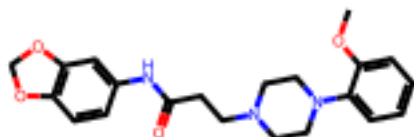
(1, 5, 50, 50)

(1, 50, 12)

(1, 1)

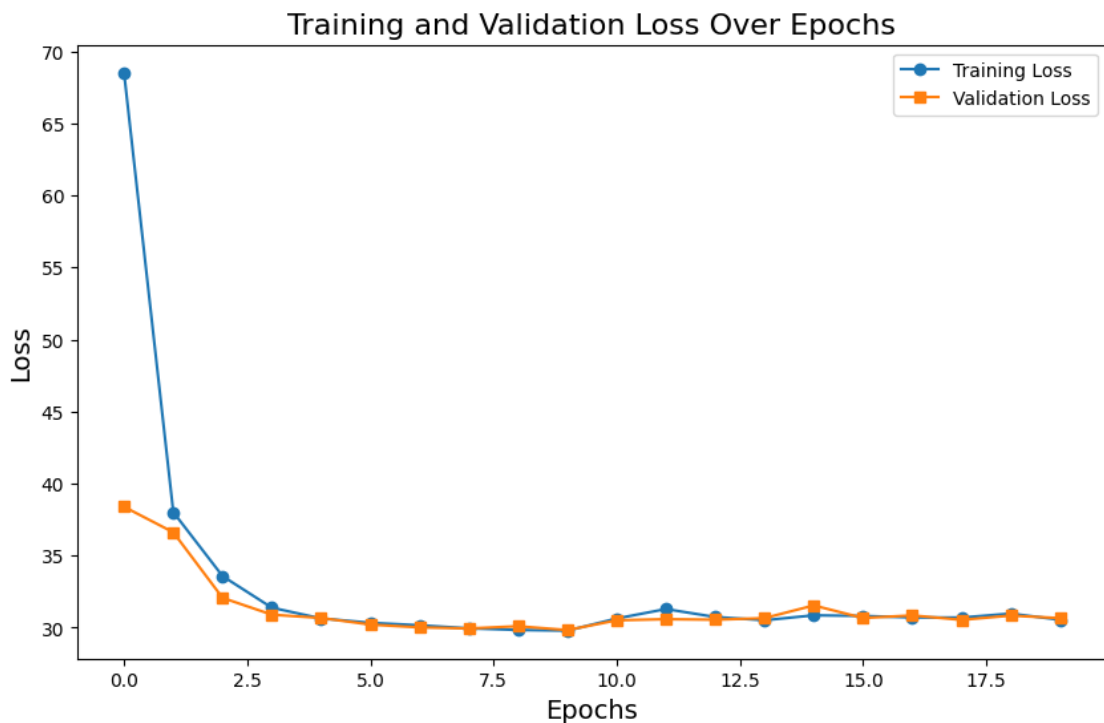
```
[188]: mole_pred = graph_to_molecule(adj0[0], feature0[0])  
Draw.MolsToGridImage([molobj,mole_pred], molsPerRow=2,)
```

```
[188]:
```



```
[189]: plt.figure(figsize=(10, 6))
plt.plot(range(EPOCHS), train_loss_list, label='Training Loss', marker='o')
plt.plot(range(EPOCHS), val_loss_list, label='Validation Loss', marker='s')

# Add title and labels
plt.title('Training and Validation Loss Over Epochs', fontsize=16)
plt.xlabel('Epochs', fontsize=14)
plt.ylabel('Loss', fontsize=14)
plt.legend()
plt.show()
```



## 0.11 Visualize latent space

```
[190]: adj_test, fea_test, score_test = [], [], []

for idx in range(len(test)):
    adjacency, features = smiles_to_graph(test.loc[idx]["SMILES"])
    score = test.loc[idx]["Score"]
    adj_test.append(adjacency)
    fea_test.append(features)
    score_test.append(score)

adj_test = np.array(adj_test)
fea_test = np.array(fea_test)
score_test_ = np.array(score_test).reshape(-1,1)

score_test_n = scaler.transform(score_test_)

[191]: ls_train = vae.encoder.predict([adj_train, fea_train, score_train_])
ls_test = vae.encoder.predict([adj_test, fea_test, score_test_])
```

```
251/251          0s 986us/step
79/79           0s 912us/step
```

```
[192]: ls_train_ = np.array(ls_train)
ls_test_ = np.array(ls_test)
```

```
[193]: z_mean, _ = vae.encoder.predict([adj_test, fea_test, score_test_])
```

79/79                      0s 901us/step

```
[194]: latent_noise = np.random.normal(scale=0.1, size=z_mean.shape) # Adjust scale
↳as needed
adj_pred, feature_pred = vae.decoder.predict([z_mean, score_test_])
print("Shape of adj_pred:", adj_pred.shape)
print("Shape of feature_pred:", feature_pred.shape)

# Reconstruct molecules
gen_molecules = [
    graph_to_molecule(adj_pred[i], feature_pred[i])
    for i in range(adj_pred.shape[0])
]
```

79/79                      0s 4ms/step  
Shape of adj\_pred: (2510, 5, 50, 50)  
Shape of feature\_pred: (2510, 50, 12)

```
[195]: from scipy.stats import pearsonr

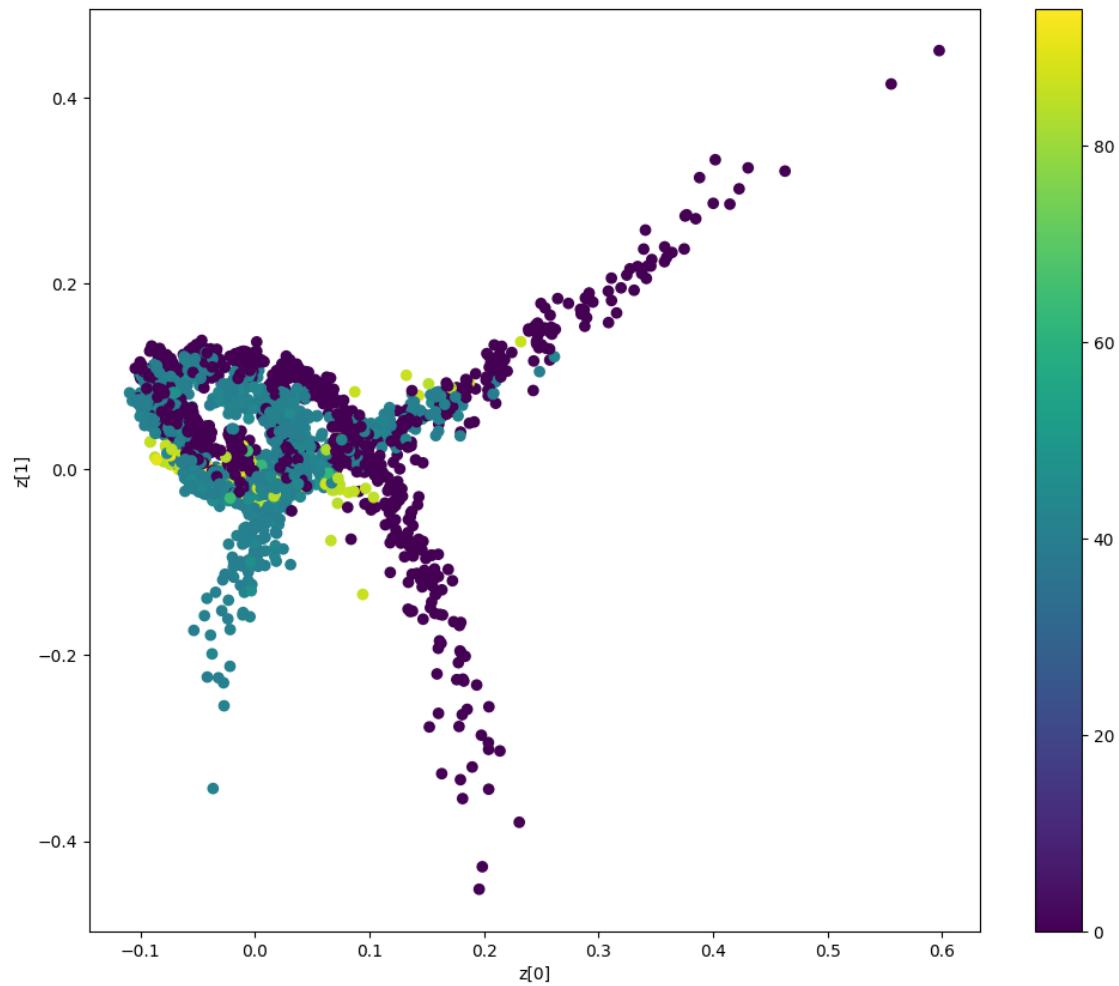
# Correlate latent dimensions with molecular scores
correlations = [pearsonr(z_mean[:, i], score_test_.flatten())[0] for i in
↳range(z_mean.shape[1])]
print("Correlations between latent dimensions and scores:", correlations)
```

Correlations between latent dimensions and scores: [-0.2525738640061273,  
-0.22531715906543592, 0.6434128304765752, 0.16416502511708225,  
0.377734084860448, -0.6168508355229856, 0.04529268882718179,  
0.31648487567027583, -0.6071813838126771, 0.37298063040288965,  
0.6475608991398247, 0.05062614560130395, -0.6725152629300453,  
-0.7275736364998957, -0.6777652913953159, -0.398432134044331,  
-0.011174325095537462, 0.7886020869785735, -0.35100706637056933,  
-0.8342214030132228, -0.2706327792927208, -0.30185522100914935,  
-0.520738661131469, 0.6585300802525533, -0.3836164444661537,  
-0.1637607510695243, 0.06494133135039279, -0.5515315568597173,  
-0.18564734524934356, 0.14752081652144017, -0.6355185396288621,  
0.05645750308494606, 0.33142908865011345, 0.5565947368238634,  
-0.4566253367777652, -0.3628239920225188, 0.355072562760383,  
-0.2875118825419305, -0.3211540940831799, 0.006743348127974842,  
0.20924810626746998, -0.658763989585168, -0.6916150260840441,  
0.2641736232403855, 0.15182369997423772, -0.009232126321090951,  
0.05137397439650175, -0.24857355385896981, 0.4514361945973758,  
0.5375556772263842, 0.6405663404471982, 0.5521444022262929,  
-0.004015038591551433, 0.5553090691210543, 0.4838979112957975,

-0.7345957436231475, 0.2122121551302846, -0.5298476123275657,  
0.1322087486428239, 0.33464027012245867, -0.7099346838880389,  
0.4592721806704094, -0.3728833075752813, 0.691684919798807, 0.1907818598525636,  
-0.8359908755562853, -0.7444948254027949, -0.49184831267100343,  
-0.7819237406066634, -0.10335499408641885, -0.06461806497364828,  
0.14749428802442874, 0.6689118520445754, 0.42666746158705693,  
-0.19340164990442515, 0.27567080777455844, 0.5992467335900125,  
-0.49977760525246295, 0.7815087695435072, -0.22993017916586694,  
0.4607956722684273, -0.3109664395621272, 0.031053118436473733,  
-0.6028690958676335, -0.47916385586989596, -0.2561115620657417,  
0.5434334135981906, 0.4191332235102636, 0.37378931381040026,  
-0.4693325963529912, 0.16942316613040104, 0.6648893438653742,  
-0.4460740970055913, -0.007520387449553773, -0.7854434980164499,  
-0.42797813815258673, 0.36034241396720523, 0.7998182942879877,  
0.5210640336736438, 0.6816254272015225, 0.3530143643756275, 0.32159789955775686,  
0.39048285611848776, 0.23458481617100957, 0.2482032561485154,  
-0.5385767602647706, 0.22370395960453993, 0.08592278749584163, 0.71447781693762,  
-0.28000258155327373, 0.34542786100248973, 0.17071778267461135,  
-0.17273338703063587, -0.2658114884633452, -0.014503051164727911,  
-0.39477271106691636, 0.5068105594138588, -0.38395930688475033,  
0.15392800898189193, 0.3057418357480992, 0.44467777216472915,  
-0.09349526691902746, -0.04827008117032223, 0.6796960964376653,  
0.5581831940441111, -0.7900245543693032, 0.5389386836793071,  
-0.2591188025002161, 0.5517966500781457, -0.6523759094381925,  
-0.6612989909382467, 0.15982879213911733, -0.33585624584953594,  
0.04832415479206585, 0.2300808165104499, 0.49485417772022394,  
0.7808751549101858, 0.3867100385602841, -0.7256171110926146, 0.6809956325553594,  
0.3276018761272377, 0.29865478464129674, -0.475563757626004,  
-0.23316142236788273, 0.6833894452474063, 0.37758884351597494,  
-0.25861548991262295, 0.09147501486360442, -0.4012246527056428,  
-0.11599972507707307, -0.15825190016678842, -0.07392071488504987,  
0.09652715928120834, -0.3178283654923457, 0.1362099584225469,  
-0.5441344381519206, -0.6910120110541981, -0.5975899201424688,  
0.8626472796968624, -0.590474066244329, 0.4924486757461424, -0.6391937659558888,  
0.3148718085740265, -0.6269102329518971, -0.34152805536457154,  
-0.5161443324819008, -0.06281111751291149, -0.4336747328813675,  
-0.63563253480028, -0.7140911746807455, -0.5515945828532701, 0.547254646306786,  
-0.15408154502863025, 0.4211045758927283, 0.4005631724681573,  
-0.6895329162132362, -0.6460775996701611, -0.17293419641709418,  
0.6018582815238452, 0.7622105727574835, -0.12425198656903118,  
-0.34968843072030914, 0.08524790640069263, -0.6590360911204893,  
0.7838266324055991, 0.7369548821162654, -0.6568361128918575,  
-0.5886696215010763, -0.07647461636492404, 0.4473084245368717,  
-0.49641130394444294, -0.38275206013725926, -0.7840294934266233,  
0.38055192069034266, 0.44086636796503625, 0.4775440241812461,  
0.25886502155040625, 0.0032617371804553597, -0.7279186786889933,  
-0.39374294211159555, 0.333566625754558, -0.7588185875410652,  
-0.2704265873880967, -0.3685473374422552, 0.7328656124576152,

```
0.5368845116407457, -0.5374222859536435, 0.3006389010924029, 0.7862764680410097,  
-0.06320008978731523, 0.27496766530997885, -0.8046648587551491,  
-0.028700429288091838, 0.5173942487901868, 0.6052213374753702,  
-0.6145979565011078, 0.835499663422969, -0.45009974281443343,  
-0.20556434659629924, -0.4416224950854324, 0.5064963214016164,  
-0.4044944555228373, -0.6607587389402322, -0.33120116680865447,  
-0.3153047500068516, 0.13309230648063122, -0.2253790824550045,  
0.48873888061161697, 0.6715725223915265, 0.3463351743299379,  
0.22961006091144226, 0.036280414899933375, -0.10284186498638838,  
-0.18438030951528828, 0.5951504391672071, 0.05686813535302522,  
0.7999676372960672, 0.5918693845758821, -0.13289414832821403,  
0.6770356807538347, -0.5055634447035431, -0.4973159964794378,  
0.07102434740251909, 0.6699824501284019, -0.37219353602707017,  
-0.8485310221655771, -0.7725100938821036, 0.4363570438522917,  
0.5525741430687774, 0.7068287539318006, 0.18081485362827196,  
0.29837231341281834, -0.5196994438070892, -0.6922818949750787,  
0.25646382558614267, -0.45431420932303174]
```

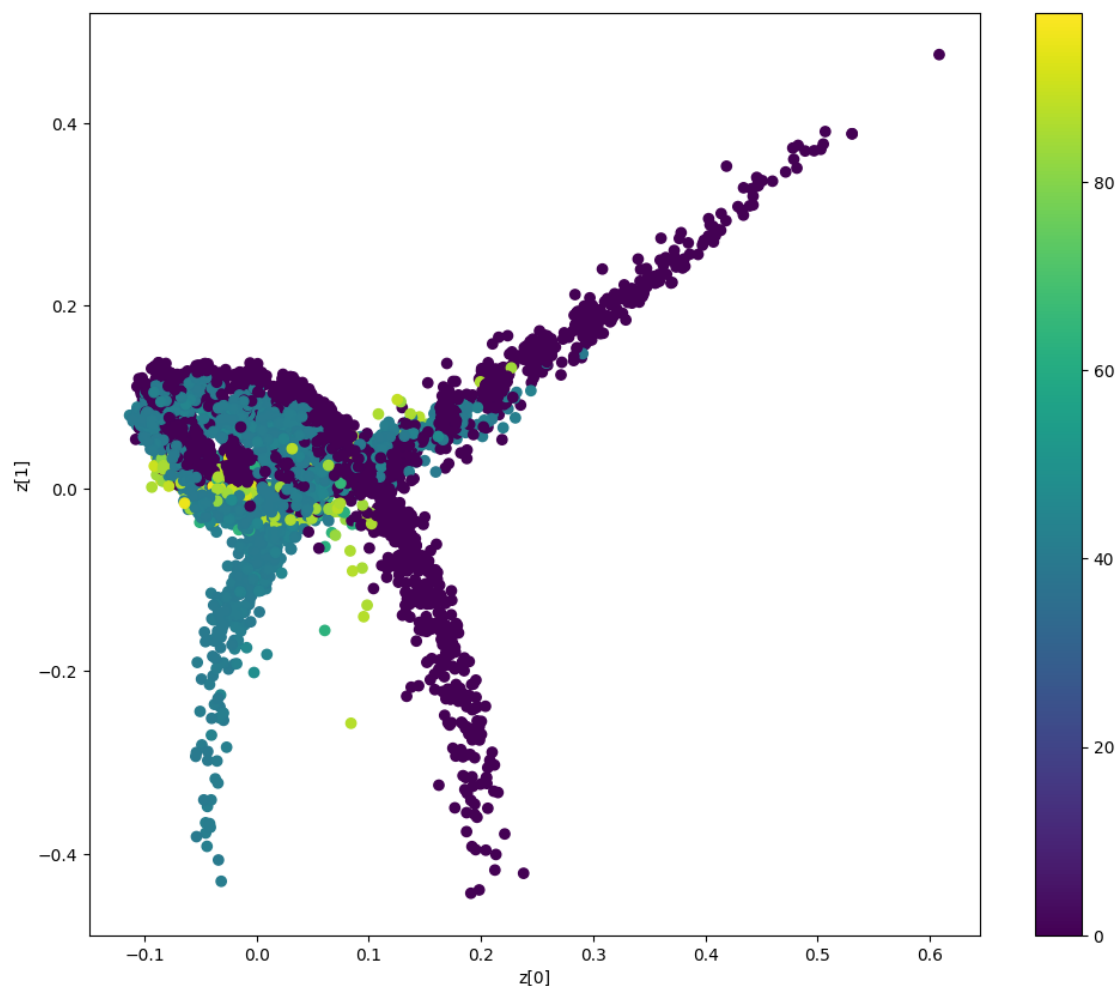
```
[196]: plt.figure(figsize=(12, 10))  
plt.scatter(z_mean[:, 0], z_mean[:, 1], c=score_test_)  
plt.colorbar()  
plt.xlabel("z[0]")  
plt.ylabel("z[1]")  
plt.show()
```



```
[197]: z_mean2, _2 = vae.encoder.predict([adj_train, fea_train, score_train_])
```

251/251                      0s 902us/step

```
[198]: plt.figure(figsize=(12, 10))
plt.scatter(z_mean2[:, 0], z_mean2[:, 1], c=score_train_)
plt.colorbar()
plt.xlabel("z[0]")
plt.ylabel("z[1]")
plt.show()
```



[ ]:

## 0.12 Model Inferencing

We would be inferring our model to predict over random latent space and try to generate 100 new valid molecules.

### 0.12.1 Generate unique Molecules with the model

```
[199]: def inference(model=vae, batch_size=1000, dim = LATENT_DIM, activity=10):
        z = np.random.normal(size=(batch_size, dim))
        activityarray = (np.zeros(batch_size) + activity).reshape(-1,1)

        reconstruction_adjacency, reconstruction_features = model.decoder.
        ↪predict([z,activityarray])
        # obtain one-hot encoded adjacency tensor
```



```

adjacency = tf.argmax(reconstruction_adjacency, axis=1)
adjacency = tf.one_hot(adjacency, depth=BOND_DIM, axis=1)
# Remove potential self-loops from adjacency
adjacency = tf.linalg.set_diag(adjacency, tf.zeros(tf.shape(adjacency)[:
↪-1]))
# obtain one-hot encoded feature tensor
features = tf.argmax(reconstruction_features, axis=2)
features = tf.one_hot(features, depth=ATOM_DIM, axis=2)

return [
    graph_to_molecule(adjacency[i].numpy(), features[i].numpy())
    for i in range(batch_size)
]

```

```

[200]: gen_mols = inference(batch_size=1000,activity=10)
MolsToGridImage([m for m in gen_mols if m is not None][:1000], molsPerRow=5,
↪subImgSize=(260, 160))

```

32/32                      0s 4ms/step

Sanitization failed: Explicit valence for atom # 1 O, 23, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 Cl, 48, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 Cl, 112, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 P, 61, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 P, 61, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 O, 91, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 Br, 42, is greater than permitted

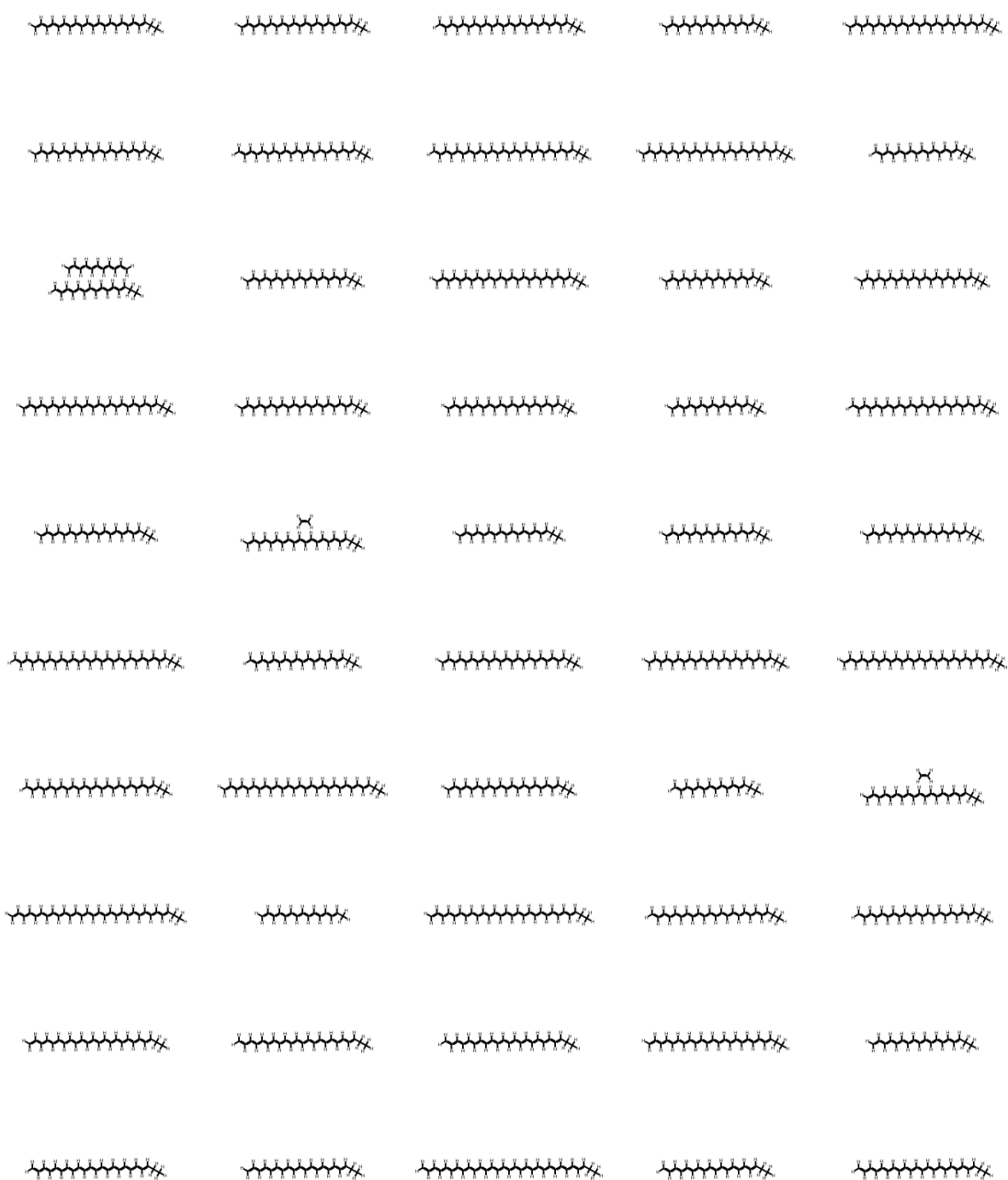
Sanitization failed: Explicit valence for atom # 0 Cl, 62, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 Cl, 82, is greater than permitted

Sanitization failed: Explicit valence for atom # 0 O, 93, is greater than permitted

/Users/thinh/Library/Python/3.12/lib/python/site-packages/rdkit/Chem/Draw/IPythonConsole.py:261: UserWarning: Truncating the list of molecules to be displayed to 50. Change the maxMols value to display more.  
warnings.warn(

[200]:



[ ]: