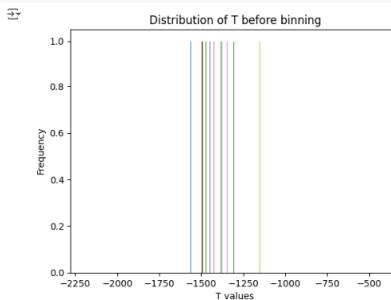


```
[1] all_features = np.hstack((X_flattened, R_flattened, Z_encoded, P_reshape))
print(f"All features shape after trimming: {all_features.shape}")
```

↪ All features shape after trimming: (7165, 998)

Convert to csv file

```
[1] plt.hist(T, bins=50)
plt.title("Distribution of T before binning")
plt.xlabel('T values')
plt.ylabel('Frequency')
plt.show()
```



```
[1] df = pd.DataFrame(data)

# Save to CSV
df.to_csv('Features.csv', index=False)
```

```
[1] df2=pd.read_csv('Features.csv')
df2.head()
```

```
↪
```

	0	1	2	3	4	5	6	7	8	9	...	981	982	983	984	985	986	987	988	989	Label
0	36.858105	2.907633	2.907612	2.907564	2.905349	0.000000	0.000000	0.000000	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	5578.0	-417.96
1	36.858105	12.599944	2.902000	1.473118	1.473101	2.901973	2.901886	1.473102	0.000000	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2504.0	-712.42
2	36.858105	14.261827	1.503703	2.924997	2.924732	1.503680	0.000000	0.000000	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	4851.0	-564.21
3	36.858105	15.871878	2.979434	1.401225	0.000000	0.000000	0.000000	0.000000	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	3130.0	-404.88
4	73.516693	17.885317	10.561490	4.355064	2.062530	2.069614	1.581991	1.590780	1.261482	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	6670.0	-808.87

5 rows × 991 columns

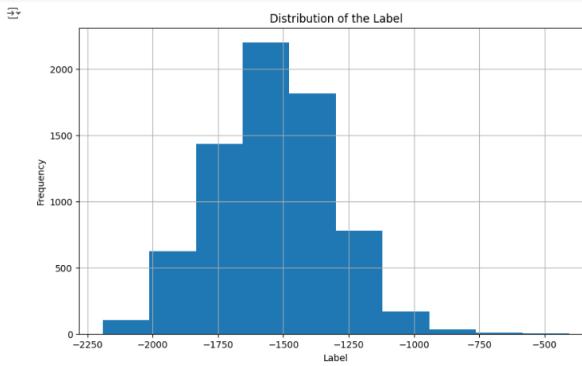
```
[1] df2.describe()
```

```
↪
```

	0	1	2	3	4	5	6	7	8	9	...	981	982	983	984	985	986	987	988	989	Label
count	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	7165.000000	...	7165.0	7165.0	7165.0	7165.0	7165.0	7165.0	7165.0	7165.0	7165.000000	
mean	78.447518	13.247392	13.651229	12.761609	11.313269	9.741235	9.741219	9.992688	1.709067	1.595967	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	3562.000000	
std	65.448339	6.134428	6.168350	5.406752	4.984658	4.785995	4.318243	0.912114	0.721890	0.698753	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	223.934546	
min	36.858105	2.907633	1.503703	1.172837	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	
25%	53.358707	9.065171	9.293125	9.019773	7.737596	6.671409	4.668823	1.421404	1.297994	1.192658	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1791.000000	
50%	73.516693	12.434812	12.495984	10.668070	10.247300	8.685217	6.699063	1.717358	1.562530	1.470488	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	3582.000000	
75%	73.516693	15.262316	17.618086	17.504143	14.998941	10.677979	9.189721	2.049483	1.957705	1.774468	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	5373.000000	
max	388.023438	47.143398	47.179157	47.015327	29.794344	29.790239	29.817144	4.379876	4.373618	4.364484	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	7164.000000	

8 rows × 991 columns

```
[1] plt.figure(figsize=(10, 6))
df2['Label'].hist()
plt.title("Distribution of the Label")
plt.xlabel('Label')
plt.ylabel('Frequency')
plt.show()
```



Run on Machine Learning Models

```
[1] # Check the shape of the concatenated data
print(df2.shape)

# Define the feature matrix X and target vector y
X = df2.loc[:, :-1].values
y = df2.loc[:, -1].values

# Check the distribution of y
print(np.unique(y, return_counts=True))

X=np.array(X)
y=np.array(y)

(7165, 991)
(array([1, 2, 3, 4, 5, 6, 7, 8, 9]), array([1, 1, 1, 1, 1, 1, 1, 1, 1]))
```

```
[1] # Standardize the features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

```
[1] models = [
    'KernelRidge': KernelRidge(alpha=1e-3),
    'SVR': SVR(),
    'LinearRegression': LinearRegression(),
    'XGBRegressor': xgb.XGBRegressor(n_estimators=200)
]

kf = KFold(n_splits=10, shuffle=True, random_state=1)

Kernels = ['poly', 'rbf', 'sigmoid']
```

```
[1] def evaluate_model(model, X, y, cv):
    ...
```

▼ Visualize Results

```
[ ] plt.figure(figsize=(14, 10))

# MSE CV
plt.subplot(3, 2, 1)
sns.barplot(x=results_df.index, y=results_df['MSE (CV)'])
plt.title('Mean Squared Error (CV)')
plt.xlabel('Model')
plt.ylabel('MSE (CV)')
plt.yscale('log')
plt.xticks(rotation=45)

# R2 Score (CV)
plt.subplot(3, 2, 2)
sns.barplot(x=results_df.index, y=results_df['R2 Score (CV)'])
plt.title('R2 Score (CV)')
plt.xlabel('Model')
plt.ylabel('R2 Score (CV)')
plt.ylim(-1, 1)
plt.xticks(rotation=45)

# MSE std
plt.subplot(3, 2, 3)
sns.barplot(x=results_df.index, y=results_df['MSE (Test)'])
plt.title('Mean Squared Error standard')
plt.xlabel('Model')
plt.ylabel('MSE standard')
plt.yscale('log')
plt.xticks(rotation=45)

# R2 std
plt.subplot(3, 2, 4)

```

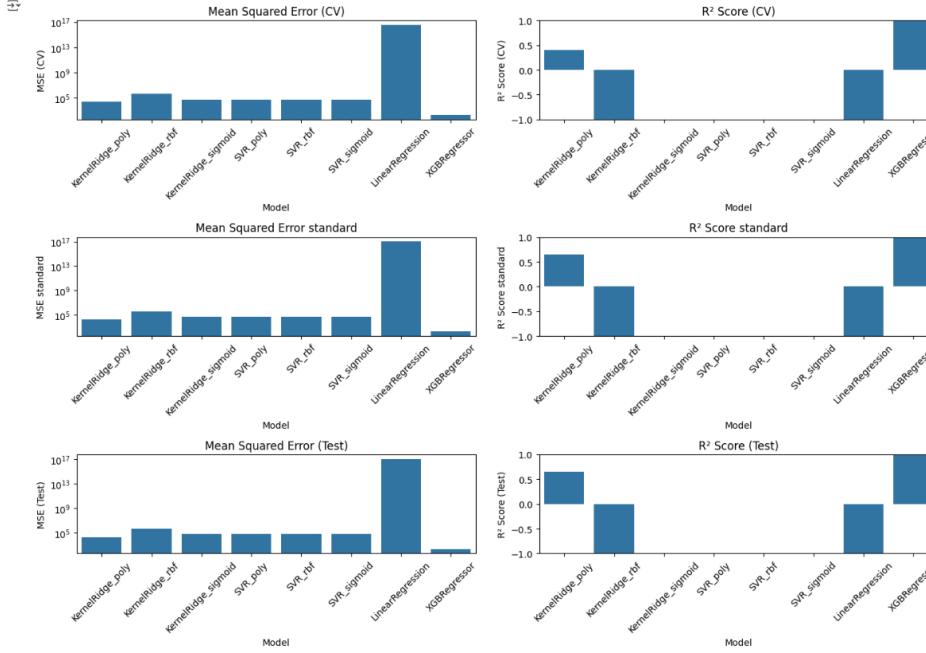
```

sns.barplot(x=results_df.index, y=results_df['R2 Score (Test)'])
plt.title('R2 Score standard')
plt.xlabel('Model')
plt.ylabel('R2 Score standard')
plt.ylim(-1, 1)
plt.xticks(rotation=45)

# MSE (Test)
plt.subplot(3, 2, 5)
sns.barplot(x=results_df.index, y=results_df['MSE (Test)'])
plt.title('Mean Squared Error (Test)')
plt.xlabel('Model')
plt.ylabel('MSE (Test)')
plt.yscale('log')
plt.ylim(-1, 1)
plt.xticks(rotation=45)

# R2 Score (Test)
plt.subplot(3, 2, 6)
sns.barplot(x=results_df.index, y=results_df['R2 Score (Test)'])
plt.title('R2 Score (Test)')
plt.xlabel('Model')
plt.ylabel('R2 Score (Test)')
plt.ylim(-1, 1)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

```



GNN

```

[1] # Convert train and test data into sparse matrix
sparse_matrix_train = csr_matrix(X_train)
sparse_matrix_test = csr_matrix(X_test)

# Convert sparse matrix into dense matrix
node_features_train = torch.tensor(sparse_matrix_train.toarray(), dtype=torch.float)
node_features_test = torch.tensor(sparse_matrix_test.toarray(), dtype=torch.float)

# Convert dense matrix into sparse representation
edge_index_train, edge_attr_train = dense_to_sparse(node_features_train)
edge_index_test, edge_attr_test = dense_to_sparse(node_features_test)

# Convert labels into tensor
labels_train = torch.tensor(y_train, dtype=torch.float)
labels_test = torch.tensor(y_test, dtype=torch.float)

print("Train Edge Index:\n" + str(edge_index_train))
print("Train Node Feature Shape: " + str(node_features_train.shape))
print("Train Labels Shape: " + str(labels_train.shape))

[2] Train Edge Index:
tensor([[ 0,  0,  0, ..., 5731, 5731, 5731],
        [ 0,  1,  2, ..., 955, 972, 989]])
Train Node Features Shape: torch.Size([5732, 990])
Train Labels Shape: torch.Size([5732])

[1] class GCNRegression(torch.nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super(GCNRegression, self).__init__()
        self.conv1 = GCONConv(input_dim, hidden_dim)
        self.conv2 = GCONConv(hidden_dim, output_dim)

    def forward(self, node_features, edge_index):
        x = self.conv1(node_features, edge_index)
        x = F.relu(x)
        x = self.conv2(x, edge_index)
        return x

[1] model = GCNRegression(input_dim=node_features_train.shape[1], hidden_dim=64, output_dim=1)

# Use Adam optimizer and Mean Squared Error loss
optimizer = torch.optim.Adam(model.parameters(), lr=0.01)
criterion = torch.nn.MSELoss()

[1] # List to store loss values
train_losses = []

# Training loop with loss storage
model.train()
for epoch in range(100):
    optimizer.zero_grad()
    out = model(node_features_train, edge_index_train).squeeze()
    loss = criterion(out, labels_train)
    loss.backward()
    optimizer.step()

    train_losses.append(loss.item())

    if epoch % 5 == 0:
        print(f"Epoch {epoch}, Loss: {loss.item()}")

[2] Epoch 0, Loss: 456272.0
Epoch 5, Loss: 5468117.0
Epoch 10, Loss: 3582781.25
Epoch 15, Loss: 3096413.0
Epoch 20, Loss: 2744736.0
Epoch 25, Loss: 2452165.5
Epoch 30, Loss: 2252135.5
Epoch 35, Loss: 2162911.0
Epoch 40, Loss: 2133829.5
Epoch 45, Loss: 2125043.25
Epoch 50, Loss: 2122889.0
Epoch 55, Loss: 2121201.0

```

```

Epoch 69, Loss: 2115464.5
Epoch 65, Loss: 2117798.5
Epoch 70, Loss: 2116322.0
Epoch 75, Loss: 2115862.75
Epoch 80, Loss: 2113924.75
Epoch 85, Loss: 2112798.75
Epoch 90, Loss: 2111627.75
Epoch 95, Loss: 2110426.5

[ ] # After training, evaluate and visualize the predictions
model.eval()
with torch.no_grad():
    output_test = model(music_features_test, music_index_test).squeeze(0)
    mse_test = criterion(output_test, labels_test).item()
    print(f'Test Mean Squared Error: {mse_test}')

# Convert tensors to numpy for plotting
output_test_np = output_test.cpu().numpy()
labels_test_np = labels_test.cpu().numpy()

⇒ Test Mean Squared Error: 1697099.625

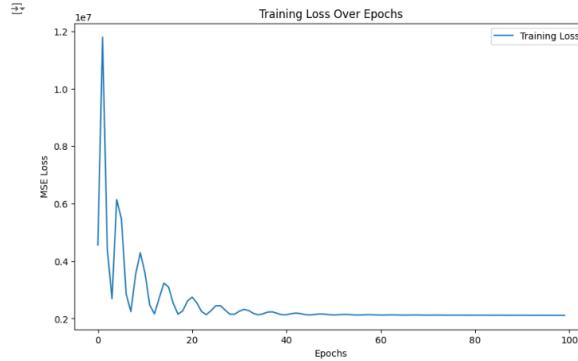
```

Visualize Result

```

[ ] # Plot the training loss over epochs
plt.figure(figsize=(10, 6))
plt.plot(train_losses, label='Training Loss')
plt.xlabel('Epochs')
plt.ylabel('MSE Loss')
plt.title('Training Loss Over Epochs')
plt.legend()
plt.show()

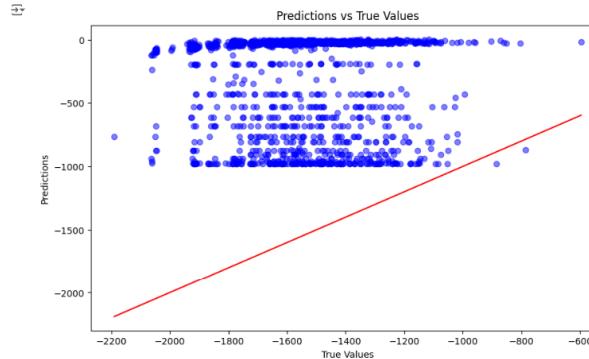
```



```

[ ] # Plot predictions vs true values
plt.figure(figsize=(10, 6))
plt.scatter(labels_test_np, output_test_np, color='blue', alpha=0.5)
plt.plot([min(labels_test_np), max(labels_test_np)], [min(labels_test_np), max(labels_test_np)], color='red')
plt.xlabel('True Values')
plt.ylabel('Predictions')
plt.title('Predictions vs True Values')
plt.show()

```



Features Extraction

```

[ ] all_features = np.hstack((X.flatten(), R.flatten(), P_reshape))
print(f'All features shape after trimming: {all_features.shape}')

⇒ All features shape after trimming: (7165, 599)

[ ] # Check the shape of the concatenated data
print(all_features.shape)

# Define the feature matrix X and target vector y
X = all_features[:, :-1]
y = all_features[:, -1]

# Check the distribution of y
print(np.unique(y, return_counts=True))

X=np.array(X)
y=np.array(y)

⇒ (7165, 599)
(array([0.000e+00, 1.000e+00, 2.000e+00, ..., 7.162e+03, 7.163e+03,
       7.164e+03]), array([1, 1, 1, ..., 1, 1, 1]))

[ ] X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Standardizing the features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

[ ] models = [
    'KernelRidge': KernelRidge(alpha=1e-3),
    'SVR': SVR(),
    'LinearRegression': LinearRegression(),
    'XGBRegressor': xgb.XGBRegressor(n_estimators=200)
]

kf = KFold(n_splits=10, shuffle=True, random_state=1)
kernels = ['poly', 'rbf', 'sigmoid']

[ ] def evaluate_model(model, X, y, cv):
    # Tính toán MSE và lưu lại tất cả các điểm số
    mse_scores = -cross_val_score(model, X, y, scoring='neg_mean_squared_error', cv=cv)
    mse_mean = mse_scores.mean()
    mse_std = mse_scores.std()

    # Tính toán R2 và lưu lại tất cả các điểm số
    r2_scores = cross_val_score(model, X, y, scoring='r2', cv=cv)
    r2_mean = r2_scores.mean()

```

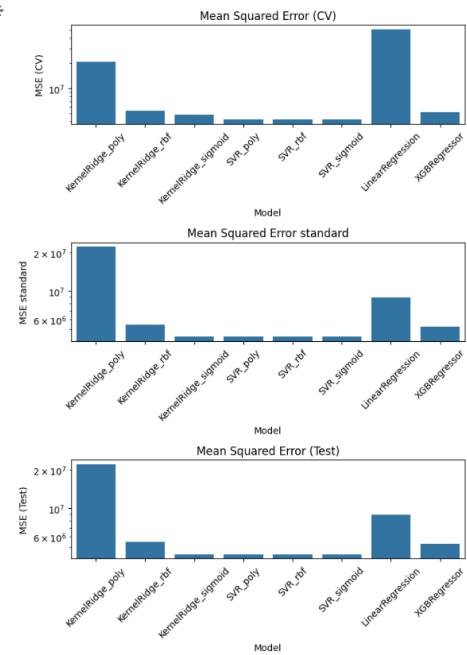


```

plt.yscale('log')
plt.xticks(rotation=45)

# R2 Score (Test)
sns.barplot(x=results_df.index, y=results_df['R2 Score (Test)'])
plt.title('R2 Score (Test)')
plt.xlabel('Model')
plt.ylabel('R2 Score (Test)')
plt.ylim(-1, 1)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

```



GNN For Features Extraction

```

[1] # Convert train and test data into sparse matrix
sparse_matrix_train = csr_matrix(X_train)
sparse_matrix_test = csr_matrix(X_test)

# Convert sparse matrix into dense matrix
node_features_train = torch.tensor(sparse_matrix_train.toarray(), dtype=torch.float)
node_features_test = torch.tensor(sparse_matrix_test.toarray(), dtype=torch.float)

# Convert dense matrix into sparse representation
edge_index_train, edge_attr_train = dense_to_sparse(node_features_train)
edge_index_test, edge_attr_test = dense_to_sparse(node_features_test)

# Convert labels into tensor
labels_train = torch.tensor(y_train, dtype=torch.float)
labels_test = torch.tensor(y_test, dtype=torch.float)

print("Train Edge Index:\n", edge_index_train)
print("Train Node Features Shape:", node_features_train.shape)
print("Train Labels Shape:", labels_train.shape)

[2] Train Edge Index:
tensor([[ 0,  0, ..., 5731, 5731, 5731],
       [ 0,  1, ..., 568, 569, 570]])
Train Node Features Shape: torch.Size([5732, 598])
Train Labels Shape: torch.Size([5732])

[3] class GCNRegression(torch.nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super(GCNRegression, self).__init__()
        self.conv1 = GCNConv(input_dim, hidden_dim)
        self.conv2 = GCNConv(hidden_dim, output_dim)

    def forward(self, node_features, edge_index):
        x = self.conv1(node_features, edge_index)
        x = F.relu(x)
        x = self.conv2(x, edge_index)
        return x

[4] model = GCNRegression(input_dim=node_features_train.shape[1], hidden_dim=64, output_dim=1)

# Use Adam optimizer and Mean Squared Error loss
optimizer = torch.optim.Adam(model.parameters(), lr=0.01)
criterion = torch.nn.MSELoss()

[5] # List to store loss values
train_losses = []

# Training loop with loss storage
model.train()
for epoch in range(100):
    optimizer.zero_grad()
    out = model(node_features_train, edge_index_train).squeeze()
    loss = criterion(out, labels_train)
    loss.backward()
    optimizer.step()

    train_losses.append(loss.item())

    if epoch % 5 == 0:
        print(f'Epoch {epoch}, Loss: {loss.item()}')

[6] Epoch 0, Loss: 16677674.0
Epoch 5, Loss: 15684064.0
Epoch 10, Loss: 15664284.0
Epoch 15, Loss: 15643793.0
Epoch 20, Loss: 15643379.0
Epoch 25, Loss: 15687472.0
Epoch 30, Loss: 15588737.0
Epoch 35, Loss: 15577462.0
Epoch 40, Loss: 15551173.0
Epoch 45, Loss: 15551147.0
Epoch 50, Loss: 15539169.0
Epoch 55, Loss: 15528964.0
Epoch 60, Loss: 15520058.0
Epoch 65, Loss: 15504469.0
Epoch 70, Loss: 15497388.0
Epoch 75, Loss: 15490698.0
Epoch 80, Loss: 15494376.0
Epoch 85, Loss: 15488455.0
Epoch 90, Loss: 15478455.0
Epoch 95, Loss: 15472950.0

[7] # After training, evaluate and visualize the predictions
model.eval()
with torch.no_grad():
    output_fact = model(node_features_test, edge_index_test).cpu().numpy()

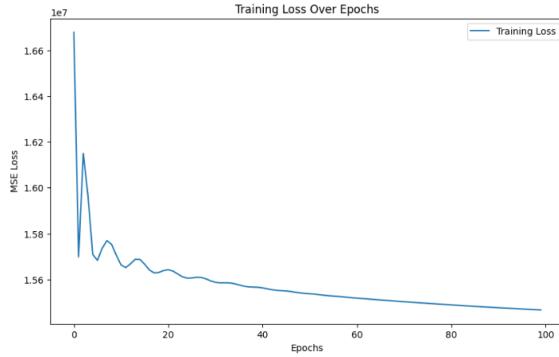
```

```
mse_test = criterion(output_test, labels_test).item()
print(f'Test Mean Squared Error: {mse_test}')

# Convert tensors to numpy for plotting
output_test_np = output_test.cpu().numpy()
labels_test_np = labels_test.cpu().numpy()
```

```
Test Mean Squared Error: 13136575.0
```

```
[1] # Plot the training loss over epochs
plt.figure(figsize=(10, 6))
plt.plot(train_losses, label='Training Loss')
plt.xlabel('Epochs')
plt.ylabel('MSE Loss')
plt.title('Training Loss Over Epochs')
plt.legend()
plt.show()
```



Các sản phẩm có tính phí của Colab - Huỷ hợp đồng tại đây

