

DSC 255R

1. ~~b.~~  $\|x - z\|_1 = \sum_{i=1}^m |x_i - z_i|$

mixed  
up the  
order

$$\sum_{i=1}^4 |x_i - z_i| = \begin{bmatrix} -2 \\ 0 \\ -2 \\ 0 \end{bmatrix} \quad \sum_{i=1}^4 |x_i - z_i| = \sqrt{4}$$

a.  $\sqrt{(-2)^2 + 0^2 + (-2)^2 + 0^2} = \sqrt{8} = 2\sqrt{2}$

c.  $\|x\|_\infty$

$\|x\|_\infty = \max_i |x_i - z_i|$   
 $= \max \begin{bmatrix} 2 \\ 0 \\ 2 \\ 0 \end{bmatrix} = 2$

2.  $x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

a.  $\|x\|_1 = 1 + 2 + 3 = 6$

b.  $\|x\|_2 = \sqrt{1^2 + 2^2 + 3^2} = \sqrt{14}$

c.  $\|x\|_\infty = \max_i |x_i| = 3$

3.  $d(x, y) \geq 0$  ~~is~~  $\checkmark$  all distances are nonnegative

$d(x, y) = 0$  iff  $x = y$   $\checkmark$  true for all entries, it's only equal to 0 when  $x$  and  $y$  are the same

$d(x, y) = d(y, x)$   $\checkmark$  symmetric for all  $x, y \in \{A, B, C, D\}$

$d(x, z) \leq d(x, y) + d(y, z)$

$\checkmark$   $x, d(A, D) \geq d(A, C) + d(C, D)$   
 $5 \geq 1 + 2$  ~~is~~

This distance function is not a metric because it violates the triangle inequality condition.

$$\begin{aligned}
 4. \quad d(p, q) &= \sum_{x \in X} p(x) \log_2 \frac{p(x)}{q(x)} \\
 &= \frac{1}{2} \log_2 \left( \frac{1/2}{1/4} \right) + \frac{1}{4} \log_2 \left( \frac{1/4}{1/4} \right) + \dots \\
 &\quad + \frac{1}{16} \log_2 \left( \frac{1/16}{1/16} \right) \\
 &= \frac{1}{2} + 0 + \frac{1}{8} \log_2 \frac{3}{4} + \frac{1}{16} \log_2 \frac{3}{8} + \frac{1}{16} \log_2 \frac{3}{8} \\
 &\approx 0.2712
 \end{aligned}$$

5. a. classification  
b. regression  
c. regression  
d. classification

6. a.  ~~$E(X)^2$~~   ~~$\neq E(X^2)$~~   $E(X^2) - E(X)^2$   
 $E(X) = 0$   $E(X^2) = 1$   
 $\text{Var}(X) = 1$

b.  $E(X) = 0$   $\text{Var}(X) = 0$

c.  $X \in \{0, 1\}$   $X=1$  w/  $p=1/4$   
 $E(X) = 1 \cdot \frac{1}{4} + 0 \cdot \frac{3}{4} = \frac{1}{4}$   $E(X^2) = \frac{1}{4}$   
 $\text{Var}(X) = \frac{1}{4} - \left( \frac{1}{4} \right)^2 = \frac{3}{16}$



$$7. a. \text{Cov}(X, Y) = E(X - E(X))(Y - E(Y)) \\ = E(XY) - E(X)E(Y)$$

$$E(X) = -1(0+0+\frac{1}{3}) + 0(\frac{1}{3}) + 1(\frac{1}{3}) = 0 \\ E(Y) = -1(\frac{1}{3} + \frac{1}{6}) + 0(\frac{1}{3}) + 1(\frac{1}{6} + \frac{1}{6}) \\ = -1(\frac{1}{3}) + 0 + 1 \cdot \frac{1}{3} = 0 \\ E(XY) = 1 \cdot -1 \cdot (\frac{1}{3}) + 1 \cdot -1 \cdot (\frac{1}{3}) = -2/3$$

$$\text{Cov}(X, Y) = -2/3 - 0 = \underline{-2/3}$$

$$b. \text{corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{std}(X)\text{std}(Y)} \\ \text{Var}(X) = \frac{1}{3} + \frac{1}{3} = \frac{2}{3} \quad \text{std}(X) = \sqrt{2/3} \\ \text{Var}(Y) = \frac{2}{3} \quad \text{std}(Y) = \sqrt{2/3} \\ \text{corr}(X, Y) = \frac{-2/3}{\sqrt{2/3} \cdot \sqrt{2/3}} = \underline{-1}$$

$$8. a. P(X=x, Y=y) = P(X=x) \cdot P(Y=y) \quad \forall x, y$$

$$P(X=-1, Y=-1) = 1/6$$

$$P(X=-1) \cdot P(Y=-1) = 1/3 \cdot 1/3 = 1/9$$

X and Y are not independent.

$$b. \text{Cov}(X, Y) \\ E(XY) = E(X)E(Y) \\ E(XY) = \sum_{x,y} xy \cdot P(X=x, Y=y) \\ = -1 \cdot -1 \cdot \frac{1}{6} + \dots + 1 \cdot 1 \cdot \frac{1}{6} \\ = 0$$

$$E(X) = -1 \cdot (\frac{1}{6} + \frac{1}{6}) + 0 + 1 \cdot (\frac{1}{6} + \frac{1}{6}) = 0$$

$$E(Y) = 0 \quad \text{Cov}(X, Y) = 0 \\ \text{X and Y are uncorrelated}$$

```
In [201... import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
import matplotlib.pyplot as plt
from sklearn.model_selection import LeaveOneOut
from sklearn.model_selection import KFold
```

```
In [93]: labels = ['NO', 'DH', 'SL']
data = np.loadtxt('hw2/spine-data.txt', converters={6: lambda s: labels.index(s) if s in labels else -1})
```

```
In [ ]: #data.view()
X = data[:, :5].reshape(310, 5)
Y = data[:, 6].reshape(310, 1)
```

```
In [107... X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 60/310)
#X_test.reshape(60,5)
#y_test.reshape(60,1)
```

```
In [112... neighbors_1 = KNeighborsClassifier(p=1).fit(X_train, y_train)
neighbors_2 = KNeighborsClassifier(p=2).fit(X_train, y_train)
```

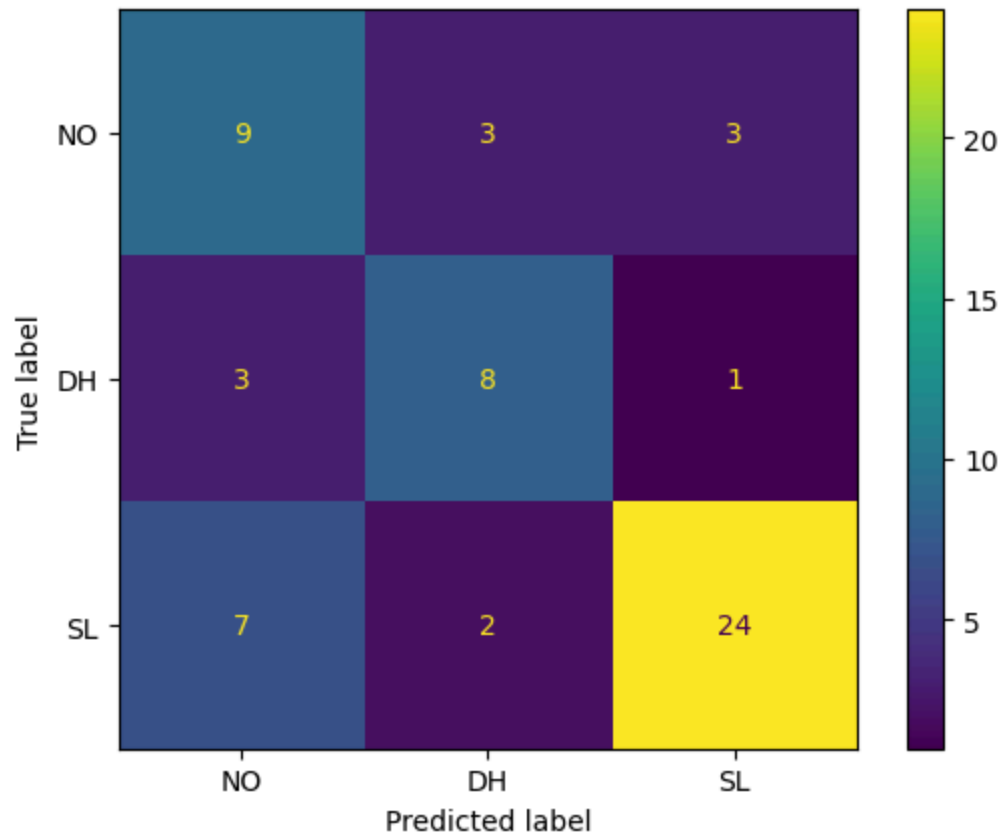
```
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/sklearn/neighbors/_classification.py:239: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
return self._fit(X, y)
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/sklearn/neighbors/_classification.py:239: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
return self._fit(X, y)
```

```
In [133... y_pred_1 = neighbors_1.predict(X_test)
y_pred_2 = neighbors_2.predict(X_test)
acc_1 = accuracy_score(y_test, y_pred_1)
acc_2 = accuracy_score(y_test, y_pred_2)
```

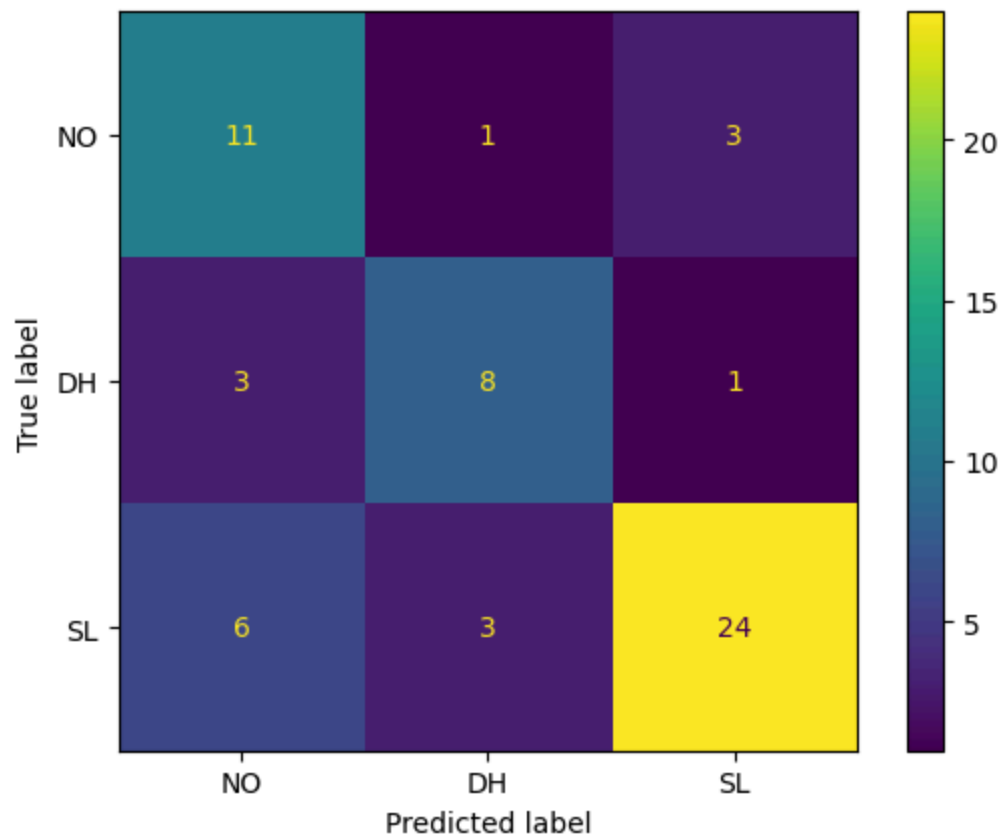
```
In [134... cm_1 = confusion_matrix(y_test, y_pred_1)
cm_2 = confusion_matrix(y_test, y_pred_2)
print("model with l1 norm error rate: ", 1 - acc_1)
print("model with l2 norm error rate: ", 1 - acc_2)
```

```
model with l1 norm error rate: 0.31666666666666665
model with l2 norm error rate: 0.28333333333333333
```

```
In [129... cm_1_disp = ConfusionMatrixDisplay(cm_1, display_labels=['NO', 'DH', 'SL'])
cm_1_disp.plot()
plt.show()
```



```
In [130... cm_2_disp = ConfusionMatrixDisplay(cm_2, display_labels=['NO', 'DH', 'SL'])  
cm_2_disp.plot()  
plt.show()
```



In [144... `!pip install pandas`

Requirement already satisfied: pandas in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (2.2.3)  
 Requirement already satisfied: numpy>=1.26.0 in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (from pandas) (2.2.1)  
 Requirement already satisfied: python-dateutil>=2.8.2 in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (from pandas) (2.9.0.post0)  
 Requirement already satisfied: pytz>=2020.1 in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (from pandas) (2024.2)  
 Requirement already satisfied: tzdata>=2022.7 in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (from pandas) (2024.2)  
 Requirement already satisfied: six>=1.5 in /Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)

In [191... `wine_data = np.loadtxt('hw2/wine.csv', delimiter=',')`  
`x_wine = wine_data[:, 1:]`  
`y_wine = wine_data[:, :1]`

In [192... `all_train_model = KNeighborsClassifier(n_neighbors=1)`  
`all_train_model.fit(x_wine, y_wine)`

/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/sklearn/neighbors/\_classification.py:239: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples,), for example using ravel().  
 return self.\_fit(X, y)

Out[192... `KNeighborsClassifier`  
`KNeighborsClassifier(n_neighbors=1)`

```
In [ ]: predictions = []
actual_values = []
loov = LeaveOneOut()
loov.get_n_splits(x_wine)

for i, (train_index, test_index) in enumerate(loov.split(x_wine)):
    #print(f"Fold {i}:")
    #print(f"  Train: index={train_index}")
    #print(f"  Test: index={test_index}")
    X_train, X_test = x_wine[train_index], x_wine[test_index]
    y_train, y_test = y_wine[train_index], y_wine[test_index]

    all_train_model.fit(X_train, y_train)

    y_pred = all_train_model.predict(X_test)

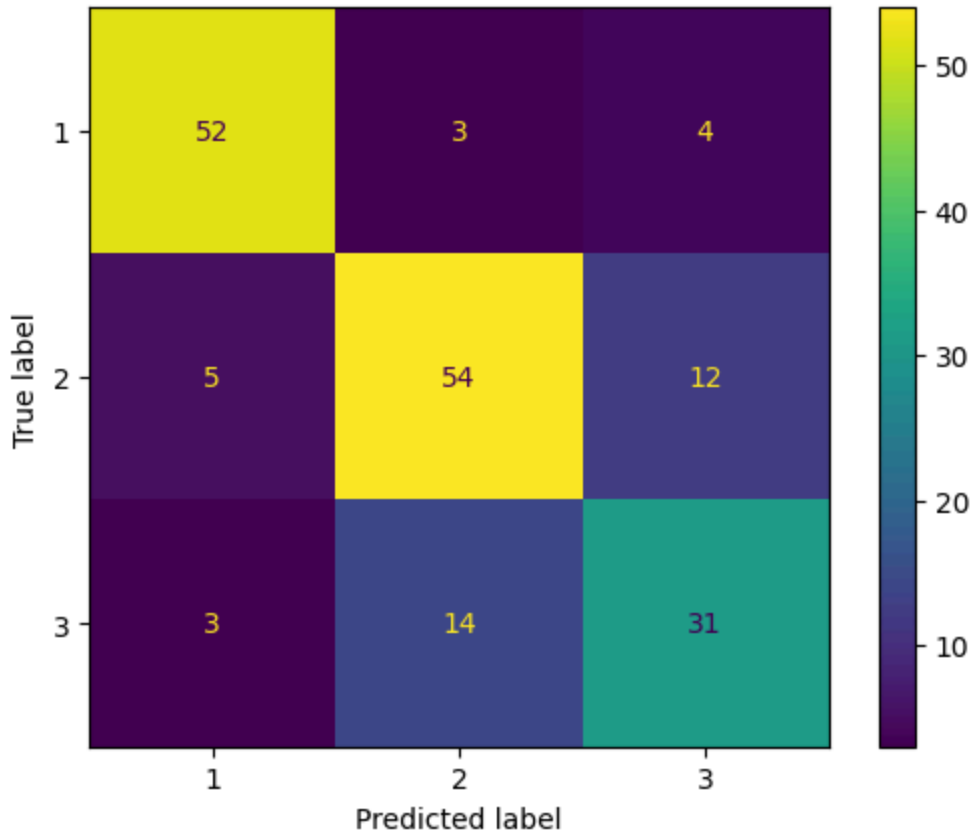
    predictions.append(y_pred[0])
    actual_values.append(y_test[0])
```

In [208... `print("L00CV accuracy: ", accuracy_score(actual_values, predictions))`

L00CV accuracy: 0.7696629213483146

```
In [209... all_model_cm = confusion_matrix(actual_values, predictions)
cm_disp = ConfusionMatrixDisplay(all_model_cm, display_labels=['1', '2', '3'])
cm_disp.plot()
```

```
Out[209... <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x136feb1d0>
```



```
In [ ]: kfold = list(range(2, 100, 5))
kfold_acc = []
for i in range(2, 100, 5):
    kf = KFold(n_splits=i)
    for i, (train_index, test_index) in enumerate(kf.split(x_wine)):
        #print(f"Fold {i}:")
        #print(f"  Train: index={train_index}")
        #print(f"  Test: index={test_index}")
        X_train, X_test = X[train_index], X[test_index]
        y_train, y_test = y[train_index], y[test_index]

        all_train_model.fit(X_train, y_train)

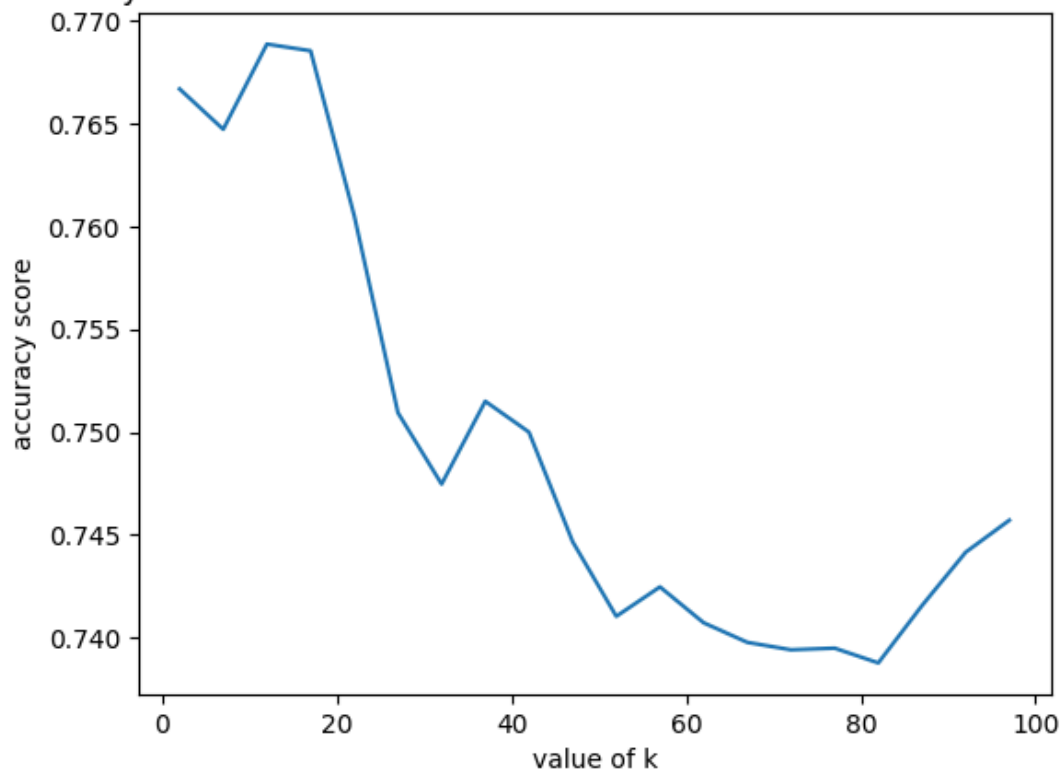
        y_pred = all_train_model.predict(X_test)

        predictions.append(y_pred[0])
        actual_values.append(y_test[0])

    kfold_acc.append(accuracy_score(actual_values, predictions))
```

```
In [214... plt.plot(kfold, kfold_acc)
plt.xlabel('value of k')
plt.ylabel('accuracy score')
plt.title('accuracy score of k-fold cross-validation at each value of k in i
plt.show()
```

accuracy score of k-fold cross-validation at each value of k in increments of 5



```
In [226... x_data = wine_data[:, 1:]
normalized_wine = (x_data-np.min(x_data))/(np.max(x_data)-np.min(x_data))
print(normalized_wine)
X = normalized_wine
y = wine_data[:, :1]
normalized_model = KNeighborsClassifier(n_neighbors=1)
#normalized_model.fit(X, y)
```

```
[8.39350664e-03 9.40548971e-04 1.36915357e-03 ... 5.41708585e-04
 2.25612696e-03 6.33900242e-01]
[7.78036396e-03 9.82218862e-04 1.19652116e-03 ... 5.47661426e-04
 1.94657920e-03 6.24970980e-01]
[7.75655259e-03 1.32748367e-03 1.51202176e-03 ... 5.35755743e-04
 1.80966384e-03 7.05334341e-01]
...
[7.82203385e-03 2.47042926e-03 1.26795526e-03 ... 2.73830713e-04
 8.51256347e-04 4.96984886e-01]
[7.76250543e-03 1.46439903e-03 1.33343652e-03 ... 2.79783555e-04
 8.86973397e-04 4.99961307e-01]
[8.33397822e-03 2.36327811e-03 1.55369165e-03 ... 2.85736396e-04
 8.75067714e-04 3.33281742e-01]]
```

```
In [ ]: predictions = []
actual_values = []
```



```

loov = LeaveOneOut()
loov.get_n_splits(X)

for i, (train_index, test_index) in enumerate(loov.split(X)):
    #print(f"Fold {i}:")
    #print(f"  Train: index={train_index}")
    #print(f"  Test:  index={test_index}")
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    normalized_model.fit(X_train, y_train)

    y_pred = normalized_model.predict(X_test)

    predictions.append(y_pred[0])
    actual_values.append(y_test[0])

```

In [228... accuracy\_score(actual\_values, predictions)

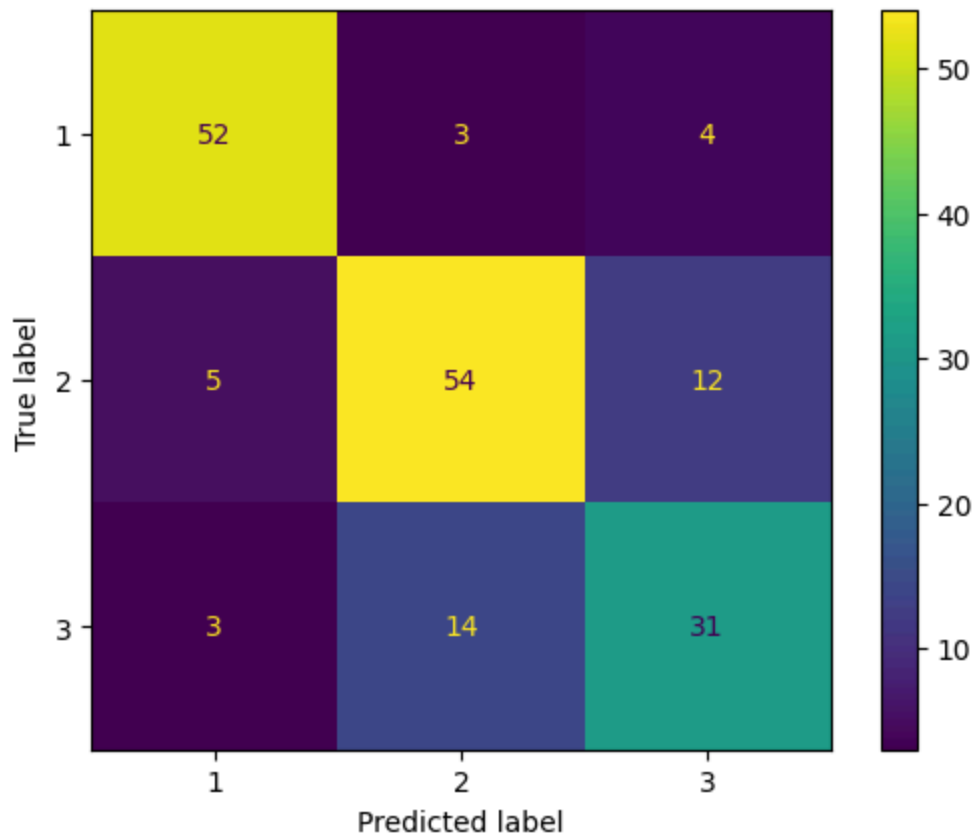
Out[228... 0.7696629213483146

```

In [230... normalized_cm = confusion_matrix(actual_values, predictions)
cm_disp = ConfusionMatrixDisplay(normalized_cm, display_labels=['1', '2', '3'])
cm_disp.plot()

```

Out[230... <sklearn.metrics.\_plot.confusion\_matrix.ConfusionMatrixDisplay at 0x1580313d0>



The accuracy score seems to be about the same with the normalization.