

Name: Atharva Chundurwar ASU ID: 1233380696

Importing Pandas and Reading the Dataset

```
import pandas as pd

df = pd.read_csv(r"titanic/train.csv")

df.head(), df.isna().sum()

(   PassengerId  Survived  Pclass \
0              1         0      3
1              2         1      1
2              3         1      3
3              4         1      1
4              5         0      3

                                                Name     Sex   Age
SibSp \
0                               Braund, Mr. Owen Harris    male  22.0
1  Cumings, Mrs. John Bradley (Florence Briggs Th...  female  38.0
2                               Heikkinen, Miss. Laina  female  26.0
0  Futrelle, Mrs. Jacques Heath (Lily May Peel)  female  35.0
1                               Allen, Mr. William Henry    male  35.0
0

   Parch      Ticket     Fare Cabin Embarked
0     0        A/5 21171  7.2500   NaN      S
1     0          PC 17599  71.2833  C85      C
2     0      STON/O2. 3101282  7.9250   NaN      S
3     0        113803  53.1000  C123      S
4     0        373450  8.0500   NaN      S  ,
PassengerId      0
Survived         0
Pclass           0
Name             0
Sex              0
Age            177
SibSp           0
Parch           0
Ticket          0
Fare            0
Cabin          687
Embarked        2
dtype: int64)
```

Building a Preprocessing Pipeline

```
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder
from sklearn.tree import DecisionTreeClassifier, plot_tree,
export_text
from sklearn.model_selection import GridSearchCV
import matplotlib.pyplot as plt

# Defining the target and features
target = "Survived"
features = ["Pclass", "Sex", "Age", "SibSp", "Parch", "Fare",
"Embarked"]

X = df[features].copy()
y = df[target].copy()

# Grouping the columns into Numerical and Categorical
num_cols = ["Age", "SibSp", "Parch", "Fare"]
cat_cols = ["Sex", "Embarked", "Pclass"]

# Numerical Columns Transformer

num_tf = Pipeline(steps=[
    ("imputer", SimpleImputer(strategy="median")),
])

# Categorical Columns Transformer

cat_tf = Pipeline(steps=[
    ("imputer", SimpleImputer(strategy="most_frequent")),
    ("ohe", OneHotEncoder(handle_unknown="ignore",
sparse_output=False))
])

#

preprocess = ColumnTransformer(
    transformers=[
        ("num", num_tf, num_cols),
        ("cat", cat_tf, cat_cols),
    ],
    remainder="drop"
)

# Decision Tree Pipeline

pipe_dt = Pipeline(steps=[
    ("prep", preprocess),
```

```
    ("clf", DecisionTreeClassifier(random_state=42))
])

# Defining the parameter grid

param_grid = {
    "clf_criterion": ["gini", "entropy", "log_loss"],
    "clf_max_depth": [3, 4, 5, 6, None],
    "clf_min_samples_split": [2, 5, 10],
    "clf_min_samples_leaf": [1, 2, 5]
}

# Defining the Grid Search CV

grid_dt = GridSearchCV(
    estimator=pipe_dt,
    param_grid=param_grid,
    cv=5,
    scoring="accuracy",
    n_jobs=-1
)
grid_dt.fit(X, y)

GridSearchCV(cv=5,
             estimator=Pipeline(steps=[('prep',
                                         ColumnTransformer(transformers=[('num',
                                         Pipeline(steps=[('imputer',
                                         SimpleImputer(strategy='median'))]),
                                         ['Age'],
                                         ['SibSp'],
                                         ['Parch'],
                                         ['Fare']),
                                         ('cat',
                                         Pipeline(steps=[('imputer',
                                         SimpleImputer(strategy='most_frequent')),
                                         ('ohe',
                                         OneHotEncoder(handle_unknown='ignore'),
```

```

sparse_output=False))),

['Sex',
'Embarked',
'Pclass'))),
('clf',
DecisionTreeClassifier(random_state=42))),
n_jobs=-1,
param_grid={'clf_criterion': ['gini', 'entropy',
'log_loss'],
'clf_max_depth': [3, 4, 5, 6, None],
'clf_min_samples_leaf': [1, 2, 5],
'clf_min_samples_split': [2, 5, 10]},
scoring='accuracy')

# Taking the best estimator from the grid search

best_dt = grid_dt.best_estimator_
best_dt.get_params()

{'memory': None,
'steps': [('prep', ColumnTransformer(transformers=[('num',
Pipeline(steps=[('imputer',
SimpleImputer(strategy='median'))),
['Age', 'SibSp', 'Parch',
'Fare']),
('cat',
Pipeline(steps=[('imputer',
SimpleImputer(strategy='most_frequent')),
('ohe',
OneHotEncoder(handle_unknown='ignore',
sparse_output=False))]),
['Sex', 'Embarked', 'Pclass'))),
('clf',
DecisionTreeClassifier(criterion='entropy', max_depth=5,
min_samples_leaf=2,
min_samples_split=10, random_state=42))),
'transform_input': None,
'verbose': False,
'prep': ColumnTransformer(transformers=[('num',
Pipeline(steps=[('imputer',

```

```
SimpleImputer(strategy='median'))],
                           ['Age', 'SibSp', 'Parch', 'Fare']),
('cat',
 Pipeline(steps=[('imputer',
 SimpleImputer(strategy='most_frequent')),
 ('ohe',
 OneHotEncoder(handle_unknown='ignore',
 sparse_output=False))]),
 ['Sex', 'Embarked', 'Pclass'))),
 'clf': DecisionTreeClassifier(criterion='entropy', max_depth=5,
 min_samples_leaf=2,
 min_samples_split=10, random_state=42),
 'prep_force_int_remainder_cols': 'deprecated',
 'prep_n_jobs': None,
 'prep_remainder': 'drop',
 'prep_sparse_threshold': 0.3,
 'prep_transformer_weights': None,
 'prep_transformers': [(['num',
 Pipeline(steps=[('imputer', SimpleImputer(strategy='median'))]),
 ['Age', 'SibSp', 'Parch', 'Fare']),
 ('cat',
 Pipeline(steps=[('imputer',
 SimpleImputer(strategy='most_frequent')),
 ('ohe',
 OneHotEncoder(handle_unknown='ignore',
 sparse_output=False))]),
 ['Sex', 'Embarked', 'Pclass'])],
 'prep_verbose': False,
 'prep_verbose_feature_names_out': True,
 'prep_num': Pipeline(steps=[('imputer',
 SimpleImputer(strategy='median'))]),
 'prep_cat': Pipeline(steps=[('imputer',
 SimpleImputer(strategy='most_frequent')),
 ('ohe',
 OneHotEncoder(handle_unknown='ignore',
 sparse_output=False))]),
 'prep_num_memory': None,
 'prep_num_steps': [('imputer', SimpleImputer(strategy='median'))],
 'prep_num_transform_input': None,
 'prep_num_verbose': False,
 'prep_num_imputer': SimpleImputer(strategy='median'),
 'prep_num_imputer_add_indicator': False,
 'prep_num_imputer_copy': True,
 'prep_num_imputer_fill_value': None,
 'prep_num_imputer_keep_empty_features': False,
 'prep_num_imputer_missing_values': nan,
 'prep_num_imputer_strategy': 'median',
```

```

'prep_cat_memory': None,
'prep_cat_steps': [('imputer',
SimpleImputer(strategy='most_frequent')),
 ('ohe', OneHotEncoder(handle_unknown='ignore',
sparse_output=False))],
'prep_cat_transform_input': None,
'prep_cat_verbose': False,
'prep_cat_imputer': SimpleImputer(strategy='most_frequent'),
'prep_cat_ohe': OneHotEncoder(handle_unknown='ignore',
sparse_output=False),
'prep_cat_imputer_add_indicator': False,
'prep_cat_imputer_copy': True,
'prep_cat_imputer_fill_value': None,
'prep_cat_imputer_keep_empty_features': False,
'prep_cat_imputer_missing_values': nan,
'prep_cat_imputer_strategy': 'most_frequent',
'prep_cat_ohe_categories': 'auto',
'prep_cat_ohe_drop': None,
'prep_cat_ohe_dtype': numpy.float64,
'prep_cat_ohe_feature_name_combiner': 'concat',
'prep_cat_ohe_handle_unknown': 'ignore',
'prep_cat_ohe_max_categories': None,
'prep_cat_ohe_min_frequency': None,
'prep_cat_ohe_sparse_output': False,
'clf_ccp_alpha': 0.0,
'clf_class_weight': None,
'clf_criterion': 'entropy',
'clf_max_depth': 5,
'clf_max_features': None,
'clf_max_leaf_nodes': None,
'clf_min_impurity_decrease': 0.0,
'clf_min_samples_leaf': 2,
'clf_min_samples_split': 10,
'clf_min_weight_fraction_leaf': 0.0,
'clf_monotonic_cst': None,
'clf_random_state': 42,
'clf_splitter': 'best'}

# Extract feature names after preprocessing (for readable plot)
ohe =
best_dt.named_steps["prep"].named_transformers_["cat"].named_steps["oh
e"]
cat_names = ohe.get_feature_names_out(cat_cols)
feature_names = list(num_cols) + list(cat_names)

dt_model = best_dt.named_steps["clf"]

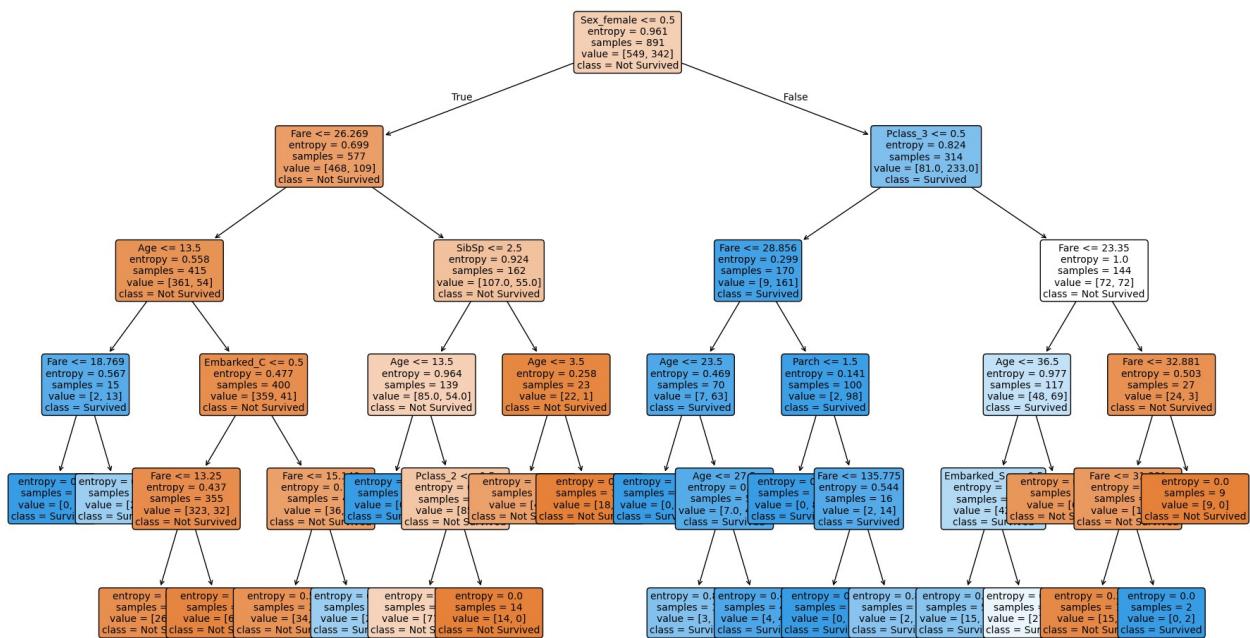
plt.figure(figsize=(18, 10))
plot_tree(
    dt_model,

```

```

        filled=True,
        feature_names=feature_names,
        class_names=["Not Survived", "Survived"],
        rounded=True,
        fontsize=10
    )
plt.tight_layout()
plt.show()

```



```

from sklearn.model_selection import StratifiedKFold, cross_val_score

cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

scores = cross_val_score(
    estimator=best_dt,
    X=X, y=y,
    cv=cv,
    scoring="accuracy",
    n_jobs=-1
)

print("Fold accuracies:", [round(s, 4) for s in scores])
print(f"Average accuracy: {scores.mean():.4f}")
print(f"Std dev: {scores.std():.4f}")

Fold accuracies: [np.float64(0.838), np.float64(0.8258),
np.float64(0.7978), np.float64(0.8202), np.float64(0.8371)]

```

```
Average accuracy: 0.8238
Std dev: 0.0146

from sklearn.ensemble import RandomForestClassifier
# Creating Random Forest Pipeline

pipe_rf = Pipeline(steps=[
    ("prep", preprocess),
    ("clf", RandomForestClassifier(random_state=42, n_jobs=-1))
])
# Defining parameter grid for Random Forest

param_grid_rf = {
    "clf__n_estimators": [200, 400, 800],
    "clf__max_depth": [None, 6, 10, 14],
    "clf__min_samples_split": [2, 5, 10],
    "clf__min_samples_leaf": [1, 2, 4],
    "clf__max_features": ["sqrt", "log2", 0.5] # 0.5 = half the features
}
# Applying Grid Search CV for Random Forest

grid_rf = GridSearchCV(
    estimator=pipe_rf,
    param_grid=param_grid_rf,
    cv=5,
    scoring="accuracy",
    n_jobs=-1
)
grid_rf.fit(X, y)

GridSearchCV(cv=5,
            estimator=Pipeline(steps=[('prep',
ColumnTransformer(transformers=[('num',
Pipeline(steps=[('imputer',
SimpleImputer(strategy='median'))]),
['Age'],
['SibSp'],
['Parch'],
['Fare']),
```

```

('cat',
Pipeline(steps=[('imputer',
SimpleImputer(strategy='most_frequent')),
('ohe',
OneHotEncoder(handle_unknown='ignore',
sparse_output=False))]),
['Sex',
'Embarked',
'Pclass'))),
('clf',
RandomForestClassifier(n_jobs=-1,
random_state=42))),
n_jobs=-1,
param_grid={'clf__max_depth': [None, 6, 10, 14],
            'clf__max_features': ['sqrt', 'log2', 0.5],
            'clf__min_samples_leaf': [1, 2, 4],
            'clf__min_samples_split': [2, 5, 10],
            'clf__n_estimators': [200, 400, 800]},
scoring='accuracy')

best_rf = grid_rf.best_estimator_
print("Best RF params:", grid_rf.best_params_)

Best RF params: {'clf__max_depth': None, 'clf__max_features': 0.5,
'clf__min_samples_leaf': 2, 'clf__min_samples_split': 2,
'clf__n_estimators': 200}

# Applying Stratified K-Fold CV to the best Random Forest model

cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
rf_scores = cross_val_score(best_rf, X, y, cv=cv, scoring="accuracy",
n_jobs=-1)

print("RF fold accuracies:", [round(s, 4) for s in rf_scores])
print(f"RF average accuracy: {rf_scores.mean():.4f}")
print(f"RF std dev: {rf_scores.std():.4f}")

RF fold accuracies: [np.float64(0.8603), np.float64(0.8371),
np.float64(0.8202), np.float64(0.8258), np.float64(0.8427)]
RF average accuracy: 0.8372
RF std dev: 0.0140

```

For this task, we can see that the Random Forest has achieved an average accuracy of 83.72% and Decision Tree has received an average accuracy of 82.38%, this tells us that after the similar preprocessing steps we took to clean the data a bit, Random Forest performs a bit better than Decision Tree model because Random Forest tries to reduce the overfitting issue using techniques like bagging and feature sub-sampling.

Decision tree model tends to capture the dataset patterns but overfits but Random forest provides more stable and generalised results.

Therefore while the average accuracy is not that different, we can still conclude that Random Forest is a better algorithm than Decision Tree for this task using the preprocessing pipeline we created.

①

ASU ID - 1233380696

Name - Atharva Chundurwar

DM - HW 2

Task 2

- a) The training error simply tells us about the rate of misclassified samples.

$$\text{Training error} = \frac{\text{total misclassified samples}}{\text{total samples}}$$

For this example, the prediction of the leaf is whichever class count is higher. The other class count is misclassified.

\therefore We have 6 leaf nodes : $D_0, D_1, E_0, E_1, C_0, C_1$.

Nodes	Misclassification #
D_0	5
D_1	6
E_0	2
E_1	6
C_0	5
C_1	5

$$\text{Total misclassifications} = 29.$$

$$\begin{aligned} \text{Total samples} &= (14+5)+(6+7)+(2+10)+(8+6)+(22+20) \\ &= \underline{100}. \end{aligned}$$

$$\therefore \text{Training error rate} = \frac{29}{100} = \underline{\underline{29\%}}.$$

②

- b) For the test instance : { A=0, B=1, C=1, D=1, E=0 } .

We start at the root node and traverse the tree according to the labels for each nodes.

Root (A) $\rightarrow 0 \rightarrow$ go left side to B.

Node (B) $\rightarrow 1 \rightarrow$ go right side to E.

Node (E) $\rightarrow 0 \rightarrow$ go left to leaf node.

We get (+ : 2, - : 10).

\therefore the decision tree will assign this test instance T to "-" class.

#

Task 3

- Q1) We are given 10 samples, out of which we have 4 positive samples and 6 negative ones.

$$\text{Gini} = 1 - \sum (h_i)^2$$

where h_i is the probability of each class label.

$$(h_+) = 4/10 = 0.4 \quad (h_-) = 6/10 = 0.6$$

$$\therefore \text{Gini} = 1 - [(0.4)^2 + (0.6)^2] = 1 - 0.16 - 0.36$$

$\therefore \underline{\text{Gini}} = 0.48$ before splitting.

③

Q2) We have 7 true samples and 3 false samples for attribute A.

Now we will calculate Gini_T & Gini_F .

For $A = \text{true}$ (7 samples) + : 4, - : 3
 $\therefore h_+ = 4/7, h_- = 3/7$.

$$\begin{aligned}\therefore \text{Gini}_T &= 1 - \left[\left(\frac{4}{7}\right)^2 + \left(\frac{3}{7}\right)^2 \right] = \frac{49}{49} - \frac{25}{49} \\ &= 24/49\end{aligned}$$

For $A = \text{false}$ (3 samples) + : 3, - : 0
 $\therefore h_+ = 1, h_- = 0$

$$\therefore \text{Gini}_F = 1 - [1^2 + 0^2] = 0$$

Now we will calculate weighted gini after split.

$$\therefore \text{Gini}_{\text{split}} = \frac{7}{10} (\text{Gini}_T) + \frac{3}{10} (\text{Gini}_F)$$

$$\therefore \text{Gini}_{\text{split}} = \frac{7}{10} \cdot \frac{24}{49} = \underline{\underline{0.34}}$$

Now, $\text{Gain} = \text{Gini}_{\text{before}} - \text{Gini}_{\text{after split}}$

$$\therefore \text{Gain} = 0.48 - 0.34 = \underline{\underline{0.14}}$$

\therefore Gini gain after splitting on attribute A is $\underline{\underline{0.14}}$.

(4)

Q3) For attribute B, true \rightarrow 4 samples
 false \rightarrow 6 samples

For B = true, + : 3, - : 1

$$\therefore \text{Gini}_T = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 1 - \frac{10}{16} \\ = \frac{6}{16} = 0.375$$

For B = false, + : 1, - : 5

$$\therefore \text{Gini}_F = 1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right] \\ = 0.2778$$

Weighted gini after split

$$\begin{aligned} \text{Gini}_{\text{split}} &= \frac{4}{10} (\text{Gini}_T) + \frac{6}{10} (\text{Gini}_F) \\ &= \frac{4}{10} (0.375) + \frac{6}{10} (0.2778) \\ &= 0.3167 \end{aligned}$$

$$\begin{aligned} \text{Gini gain} &= \text{Gini before} - \text{Gini after split} \\ &= 0.48 - 0.3167 = 0.1633 \end{aligned}$$

\therefore Gini gain after splitting on B
 is 0.1633.

⑤

Q4) We pick the splitting attribute which has larger gini gain.

As we calculated, $\text{gini gain}(A) = 0.137$
 $\& \text{gini gain}(B) = 0.163$

\therefore We will choose or decision tree will choose B as the root split as it will give purer children and lower impurity than A.

Task 4

Q1) No decision tree classifiers are not linear classifiers. A linear classifier is the one which tries to separate the data using a straight line. On the other hand, decision tree tries to split the data in step by step manner using different conditions. Therefore, because of these multiple cuts, decision trees can have non-linear patterns in the data and the final boundary is not a straight line.

Q2) Gini index is a better splitting criterion than misclassification error.

The misclassification error criterion only checks how many samples are wrongly classified but it does not show how pure a node is.

But gini index gives importance to this, it becomes smaller as the node get purer.

Thus gini index helps the decision tree find better and cleaner splits.

⑦

#

Task 5

Bagging Weaknesses:

Bagging creates many decision trees using random samples of data and then combines all the results. But a big weakness is that all the trees use same features when splitting.

Because of this, trees may become correlated and result in similar decisions. This reduces the benefit of combining them.

Difference between Bagging & Random Forest

Random forest is improved version of bagging. In random forest, at each split in tree, the algorithm chooses a random subset of features instead of using all features, this makes every tree slightly different.

Why this helps:

The random feature selection of trees in random forest generates less correlated and more diverse trees. As a result the final prediction is more stable and accurate.

Task 7

circle equation with (a, b) as center and radius r : $(x_1 - a)^2 + (x_2 - b)^2 - r^2 = 0$

$$x_1^2 - 2ax_1 + a^2 + x_2^2 - 2bx_2 + b^2 - r^2 = 0$$

$$\therefore (x_1^2 + x_2^2) + (-2a)x_1 + (-2b)x_2 + (a^2 + b^2 - r^2) = 0$$

Given the feature map $\phi(x) = (x_1, x_2, x_1^2, x_2^2)$ in this feature space, the expanded circle is a linear equation of the form $w \cdot \phi(x) + b = 0$.

where, $w = (-2a, -2b, 1, 1)$ & $b = (a^2 + b^2 - r^2)$

→ The points present inside the circle satisfy $w \cdot \phi(x) + b \leq 0$.

→ Points outside circle satisfy $w \cdot \phi(x) + b \geq 0$.

→ Points on the boundary: $w \cdot \phi(x) + b = 0$.

∴ We can conclude that in the feature space (x_1, x_2, x_1^2, x_2^2) , every circular region is separated from the rest of the plane by a hyperplane. Thus circular regions are linearly separable in this space.

#

#

Task 6

$$\phi(x_1, x_2) = (z_1, z_2) = (x_1, x_1 x_2)$$

Given the points and labels: $z_1 = x_1$, $z_2 = x_1 x_2$

(x_1, x_2)	label	(z_1, z_2)
$(-1, -1)$	-	$(-1, +1)$
$(-1, +1)$	+	$(-1, -1)$
$(+1, -1)$	+	$(+1, -1)$
$(+1, +1)$	-	$(+1, +1)$

From the table, we can see that:

All positives lie at $z_2 = +1$

All negatives lie at $z_2 = -1$

∴ We can use the best separating line as halfway between them.

Now we can use:

if $z_2 > 0 \rightarrow$ predicts negative.

if $z_2 < 0 \rightarrow$ predicts positive.

All 4 points lie at distance 1 from the boundary so all four are support vectors.

After mapping (x_1, x_2) to $(x_1, x_1 x_2)$ positives are at $z_2 = -1$ and negatives at $z_2 = +1$.

→ the maximal margin separator is the horizontal line $z_2 = 0$.

→ $w = (0, 1)$, $b = 0$ gives the boundary.

→ Margin = 1 and all the 4 points are support vectors.

Task 8

$$\text{Ellipse equation: } c(x_1 - a)^2 + d(x_2 - b)^2 - 1 = 0$$

$$cx_1^2 - 2ax_1 + ca^2 + dx_2^2 - 2bx_2 + bd^2 - 1 = 0$$

$$\text{the given feature vector: } \phi(x) = \begin{pmatrix} 1, x_1, x_2, x_1^2, \\ x_2^2, x_1 x_2 \end{pmatrix}.$$

If we use expanded ellipse equation, we see that it is in quadratic form in (x_1, x_2) . Thus we can express it as $\phi(x) \cdot w \cdot \phi(x) = 0$

here with suitable weights $w \cdot x_1, x_2$
 weight can be 0 for axis aligned ellipses and non-zero for not axis aligned. Thus we can say that any elliptic boundary is a hyperplane in this feature space.

The degree -2 kernel is $k(u, v) = (1 + u \cdot v)^{-2}$
 This equals a dot product in quadratic feature above: $(1 + u \cdot v)^{-2} = \phi(u) \cdot \phi(v)$

For a mapping equivalent to $(1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$
 scaling doesn't change linear separability.

Therefore a SVM using $k(u, v) = (1 + u \cdot v)^{-2}$ is same as a linear SVM in the quadratic features $(1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

Github code link:

<https://github.com/AtharvaBOT7/CSE572-DM-HW2-1233380696>