itanic import numby as up import pandes as pol import Jeobon as S/S OB Set G'dorkeril') Style = from skleum, ensemble import Random Forest Clasistier. from skloom preprocessing import Onellot Encoder, Label Encoder, Standard Scoler from skleam netrics import roccure, auc from Sklean Model-Selection import Stratified Ktold inport String input womings. Wornings filter warnings (ignore) JEED : 42. def concat_df (train_dota, test_dota); Petum pd. corcat ([train_data, test_data], Sort=True). reset-index (dop = True) def divide-df (all-dota) return all dota 2. /oc[:890], all-dota. /oc[891:]. drop ([:Survival] OXIS=1)

```
of train . pd. read ou ( 'puth')
df_test = pd. road_csu ( 'poth')
df_all = Concat_df (df_tmin, df_text)
of train name = 'Traing det'
of test name = 7st Set'
df-all. name = 'All Set'
Plint (Number of training Sets : { 3', format (df.train, Stope [0]))
print ( ' ' Test " ' '
print ( Training Sets X Slape : { 3', farmet( * *))
1. Exploratory Data Analysis.
 1.1 Overview
 print (df-train. Info())
 of-train, sample (3)
  1.2 missing values.
  def display_missing (df):
  for cal in df. colums, tolist():
          print ( fool, afford, is not). Jano) 3 column misses why : {
         Print ('In')
```

for df in dfs:

Print (df. name)

display - prissing (df)

1.2.1 Age. df_all_corr - df_all.corre). also). Unstacke), Sort-values Ctokind = 'quickbort', ascending = False), reset_index() df_all_corr. rename (calums = { '/axel_0' = 'Forme /', -- 3, inplace = T) df_all-corr [df_all-corr [feature 1] = Age] index 4 age_by-pclass_sex)=df_all, graphy (['Sex', 'pclass']), median()['Age'] nediant & 熟题. df-all ['Age'] = df-all. grouply (['Sex', 'poloss']) ['Age'], apply (landa x: X. Alha (X. Median ())) 1.22 Embarked defall [df_all ['Embarked'], isnull()] df-all ['Emborked'] = df-all ['Emborked']. filma('S') 1.2.3 Fore ['fare'] med-Fore · df-all grouply (['pals', 'parch', 'sibsp'])['fore']. median()[3][0][1] [Fare'] 1.24 Colin det of all ['secki] - of all ['Cabin'], apply (hambda s: S[0] if pol motroll(s) else 'm') df-all-decks = df-all grouply (['Deck', 'Palus']), count(). drop((olumo ·[' -- ']

DESIGN K

def get-pchss-dist (df);

Greating a dictionary for every passenger class court in every deck.

decks = {'A': {}, 'B': {}, --- 'T': {}}}

decks = df, columns, levels [0]

for deck in decks:

for polass in rape (1,4):

try:

Court = df[deck][pclass][0]

deck_courts [deck][pclass] = court

except Kex[mor;

deck_courts [deck][pclass] = 0

Parantyc.

df-decks = pd. Potytrane (deck_conts.)

deck_percentures = { }

for col in df_decks, columns;

deck_parantypes[aol] = [(count/df_decks[col].Sun()) x/00

for count in df_decks[col]]

return deck-courts, deck-percentures.

def display-polass-dist (perconleges): Anthone df-percentages - pd. (percentages), transposec) deck-nowes - ('A', ···· 'T') bor-count = 119. orange (/en(deck rames)) bor- width = 0.85 pclus = df_parentyes [0] pc/as 2 - of percentyes [1] polass 3 - df- percentycs [2] P/t. figure (figsize = (20,10)) pt. bor (bor-count, pclas), color: #65469', edecolor='white', width - bar-width, label . posseger clas 1') plt. for (" botton · pobs 1 " Plt. har (" botton - poliss 1 + poliss 2 - 2 3 Jake 1019 219 pt. x/atel (Peck', Size 15, (latel pad)=20) At. Ylabel ("Passoyer Clas Penange", Size-15, hareland > 20) plt. Xticks (box-court, dock-runes) -> x= label 423 PH. tick-params (axis = 'X', labelsize = 15) " (" ") plt./good (/ac · 'upper /eft', bbx_to_ anchor: (1-1), prop · f'size' · 15 }) 1, Size-18, /= 1.05) PH. title (' pH. Show ()

DESIGN K

all_deck_court, all_deck_per = get_pchss_dist (df_all_decks)
display_pchss_dist (all_deck_per)

Passesor in the I deck is changed to A.

idx = df_all[df_all[Deck] == 'T']. index.

df_all.loc[idx, 'Deck'] = 'A'

df_all_decks_Durvived = df_all_grouply (['Deck', 'Survived']). Count(),

drop(colorns = [---]).

renone(colorns = [1/bre' : Count']).

transpose() = 1/2 colorne. Watel colorne.

def get_survived_dist(cdf);

Greating a dictionary for every survival court in every deck.

Surv_courts = {'A': { }, ... }

decks = df. columns. /evels [o]

for deck in decks:

for survive in rappe(0,2):

SURV_ cants [deck][Survive] = df [deck][Survive][0]

df_surv = pd. DataFrame (surv_counts)

Surv-percentages = { 3

for cal in df_surv.colums:

Surv-percentages = [(count / df_surv[cal].sum()) x 100

for count in df_surv[col]]

DECIDAL IN

r.L.m	and courte	SUN_ percentages
return	JUIV- Carris	0010-100-10

det display - Surv_ dist (percentages):

别 神 智.

all_surv_count, all_surv_per = get_survived_dist (df_all_decks_surv)
display_surv_dist (all_surv_per)

volve 4771. df_all['Pack'] = df_all['Dack']. replace(['A', B', C'], 'ABC']

(,

df_all ['peck'], value_counts()

Dropping the Cabin feature

df_all.drop(['Cabin'], axis=1, inplace = T)

df_troin, df_test = divide_df (df_all)

dfs = [df_troin, df_test]

for df in dfs:

display_missing (df)

1.3	Sorvival Distribution	
Survived	= df_train['Survived'], value_courts()[1]	

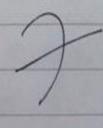
not_survived = " [0] Survived per - Survived / df_train. Shape [0] */00 not_ " not_ " "

Plt. figure (figsize = (10,8)) Shs. countplot (df_train ['Survived'])

1.4 Correlations

fig, axs = plt. subplots (nrow = 2, figsize=(20,20))

Sis. hertmap (df_train.drop(['fassengerID'], axis=1), corr(), ax = axs[o], anot = True, anot kus = {'size': 143, Square = True, chap='coolubm') df_test "



1.5 Survival Distribution in Features. Cont_features = ['Age', 'Fare'] SUN = df_train ['Survived'] == 1 fig. axs = plt. subplots (ncels=2, mas=2, figsize=(20, 20)) plt. sulplets - adjust (right = 1.5) for i, feature in enumerate (cont_features): # Distribution of Survival in feature Sns. distplot (df_train [~ surv][feature], /abel = 'Not Junited', hist = True, color = ', ax = axs[6][i]) SIS distplot (df-train[SURV][feature], /abel = "Survived", " axs[o][i]. Set_Xlobel(") } xi label Name. 15. [J " | MA A

for j in roge(2);

axs[i][j]. tick_parons (axs='x', lablesize=20)

""
y'

1.5.2 Cotoporical Features.

Cat_features = [~]

fig. axs = ptt. Sulphts (ncok = 2, now = 3, figsize = (20,20)) ptt. Subplots_addist (right = 1.5, tsp=1.25)

for i, feature in enumerate (cot-features, start = 1):

plt. Subptit (2,3, i) a Categorial lakel

SAS. count plot (X = feature, five) = 'Survived', data = df_train)

A CONTRACT VISION TO THE COMME

ME ACHAR CAR MAN STONE TOWN THAT HAVE SHOW THE

1.6 Conclusion

(DOTE OF A

df_all = Concot_df (df_train, df_test)

df_all head ()

2. Feature Engineering

2.1 Binning the continuous features

2.1.1 Fase

जिल्हा यहन व्यय

137年 724年 42.

df all [Fare'] = Pd. gat (df all [Fare'], 13)

fig. exes = plt. Jubplets (figsize = (22,9))

Sts. complet (X='Fae', bue = 'Sirvived', data=df_all)

) Option

2.1.2 Age

dfal['Aso'] = pl. qcot (dfal['Age'], 10)

fig. axs = pt. subplots (fisize (21, a))

Ins. complet (X='Age', hoe 'Survived', dota = df_all)

axs[0][0], Jet_title ('Family Size feature Ublue counts', Size=20, y=1.05)
axs[0][1] " ('Survival counts in family size', Size=20, y=1.05)

family_map = {1: 'Alone', 2: 'Alone', 3 df-all ['Fourily_Size_grap'] = df-all ['Family_Size'], map (family_map)

2.3 Title 6 Is mortied df_all ['Title'] = df_all ['Name'], Str. split (';', expand = True) [1]. \ Str. Split ('.', expand = True) [0] df-all ['Is-Married'] = 0 deal [Is Movied] = . /oc [deal [Title] == 'Ms] = 1 Bulptot) a Ens. burplot (x = df_all [Title], value_courts () index, ". Values, data day ax = axs[0]) A The Boy war to sent way (d. a [' A 24 Juroion rate. def extract_surname (data): families = [] for i in lepsth (data): name = data, iloc [i] if 'C' in name: name_10_bruket = name. Split('(')[0] else: Mane_no_bracket = hane.

family = MAR-10- brocket. Split (',')[0]
title = " [1]. Strip(). Splin('')[0]
755 Martinger 151 Har and Walland Andrew Thomas Hall
for c in String punctuation;
family = family, replace (C, "). Strip()
ALCHERON TO BE A THE PERSON OF
families, append (family)
the last select to the season of the selection of the sel
Veturn families
df_all ['Family'] = extract_surname (df_all ['Name'])
df_train = df_all./oc[:890]
df- test = " [891;]
dfs - [df_train, df_test]
Gentling a list of families and tickets
that are occurring in both training & test sets.
non_unique_families = [x for x in df-train ['Family'], uniquec)
if x in df. transfert ['family'], unique()]
101-Unique_tickets - [" ['Ticket'] " [,]]
grave
df-family_survival_rate = df_train ('Family') ['Survived', 'Family', 'Family size']
+ 1Kel (Median()
Heddan() Heddan() (Ticket) "Ticket"

family- moles = { 3 tidet_rates = { 3 for in name (len (df-family_Survival_rate)): it df_family_survival_rates, index[i] in non_uninve_families and df_family_survival_romes, ilec[i] 71: family_romes [df_family_ ~ index[i]] = df ~ ile[i,o] Tided > for