

Get started

This is your home for all data science and engineering work.

We'll show you how to set up clusters, data and users.

Set up your workspace

- Create a cluster
- Import data
- Build a Delta Live Tables pipeline
- Invite your team

Next steps

- Explore data with Python
- Read documentation

titanic Python

File Edit View Run Help Last edit was 7 minutes ago Give feedback

⚠️ Free trial ends in 14 days. Continue with a pay-as-you-go subscription by [providing your billing information](#).

```

1 # linear algebra
2 import numpy as np
3
4 # data processing
5 import pandas as pd
6
7 # data visualization
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 from matplotlib import style
11 from matplotlib import rc
12
13 # Algorithm
14 from sklearn import linear_model
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.svm import LinearSVC
17 from sklearn.linear_model import Perceptron
18 from sklearn.linear_model import SGDClassifier
19 from sklearn.tree import DecisionTreeClassifier
20 from sklearn.neighbors import KNeighborsClassifier
21 from sklearn.svm import SVC, LinearSVC
22 from sklearn.naive_bayes import GaussianNB

```

Command took 0.06 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 train_df.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   PassengerId 891 non-null   int64  
 1   Survived    891 non-null   int64  
 2   Pclass      891 non-null   int64  
 3   Name        891 non-null   object  
 4   Sex         891 non-null   object  
 5   Age         891 non-null   float64 
 6   SibSp       891 non-null   int64  
 7   Parch       891 non-null   int64  
 8   Ticket      891 non-null   object  
 9   Fare        891 non-null   float64 
 10  Cabin        284 non-null   object  
 11  Embarked    889 non-null   object  
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

```

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 train_df.describe()

```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000	891.000000		
mean	446.000000	0.353838	2.309642	29.699118	0.523008	0.351594	32.204208			
std	257.353842	0.465992	0.836071	14.526497	1.102743	0.806057	49.693429			
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000			
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400			
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200			
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000			
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200			

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 train_df.head(15)

```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	Nan	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	Nan	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	33.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	Nan	S
5	6	0	3	Moran, Mr. James	male	Nan	0	0	330877	8.4583	Nan	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17483	51.8625	E46	S
7	8	0	3	Paisson, Master Gosta Leonard	male	2.0	3	1	349909	21.0750	Nan	S
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	Nan	S
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	Nan	C
10	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 3549	16.7000	G6	S
11	12	1	1	Bonnell, Miss. Elizabeth	female	58.0	0	0	113783	26.5500	C103	S
12	13	0	3	Seendercock, Mr. William Henry	male	20.0	0	0	A/5 2151	8.0500	Nan	S
13	14	0	3	Andersson, Mr. Anders Johansen	male	39.0	1	5	347082	31.2750	Nan	S
14	15	0	3	Vestrom, Miss. Hilda Amanda Adolfsen	female	14.0	0	0	350406	7.8542	Nan	S

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 total = train_df.isnull().sum().sort_values(ascending=False)
2 percent_1 = train_df.isnull().sum()/train_df.isnull().count()*100
3 percent_2 = (round(percent_1, 1)).sort_values(ascending=False)
4 missing_data = pd.concat([total, percent_2], axis=1, keys=['Total', '%'])
5 missing_data.head(5)

```

	Total	%
Cabin	687	77.1
Age	177	19.9
Embarked	2	0.2
PassengerId	0	0.0
Survived	0	0.0

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 train_df.columns.values

```

```

Out[34]: array(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
   'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype=object)

```

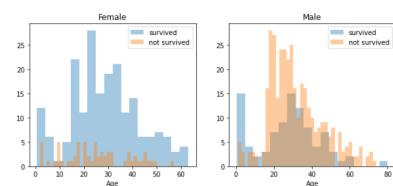
Command took 0.09 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

1 survived = 'survived'
2 not_survived = 'not survived'
3 fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(10, 4))
4 women = train_df[train_df['Sex']=='female']
5 men = train_df[train_df['Sex']=='male']
6 ax = sns.distplot(women[women['Survived']==1].Age.dropna(), bins=18, label = survived, ax = axes[0], kde = False)
7 ax = sns.distplot(women[women['Survived']==0].Age.dropna(), bins=40, label = not_survived, ax = axes[0], kde = False)
8 ax.legend()
9 ax.set_title('Female')
10 ax = sns.distplot(men[men['Survived']==1].Age.dropna(), bins=18, label = survived, ax = axes[1], kde = False)
11 ax = sns.distplot(men[men['Survived']==0].Age.dropna(), bins=40, label = not_survived, ax = axes[1], kde = False)
12 ax.legend()
13 _ = ax.set_title('Male')

```

```
/databricks/python/lib/python2.9/site-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
  warnings.warn(msg, FutureWarning)
```

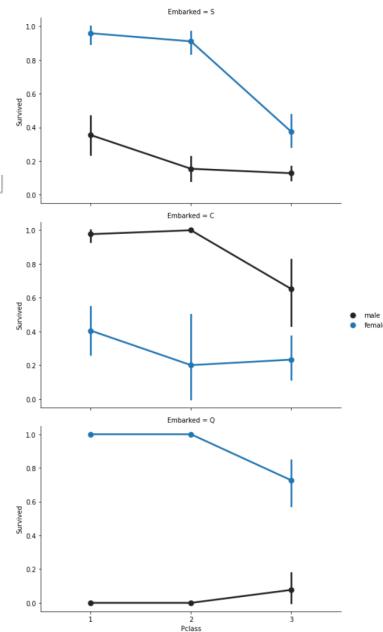


```
Command took 0.90 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

```
Cmd 9
```

```
1 | FacetGrid = sns.FacetGrid(train_df, row='Embarked', size=4.5, aspect=1.6)
2 | FacetGrid.map(sns.pointplot, 'Pclass', 'Survived', 'Sex', palette=None, order=None, hue_order=None)
3 | FacetGrid.add_legend()
```

```
/databricks/python/lib/python2.9/site-packages/seaborn/axisgrid.py:337: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
  warnings.warn(msg, UserWarning)
Out[36]: <seaborn.axisgrid.FacetGrid at 0x7f38e9ec620>
```

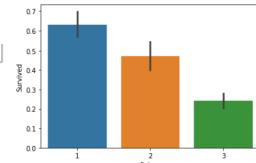


```
Command took 1.80 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

```
Cmd 10
```

```
1 | sns.barplot(x='Pclass', y='Survived', data=train_df)
```

```
Out[37]: <AxesSubplot:xlabel='Pclass', ylabel='Survived'>
```

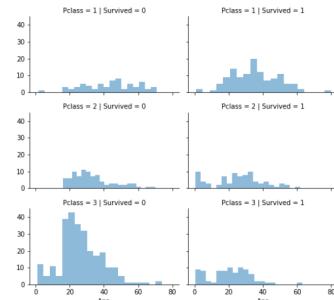


```
Command took 0.40 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

```
Cmd 11
```

```
1 | grid = sns.FacetGrid(train_df, col='Survived', row='Pclass', size=2.2, aspect=1.6)
2 | grid.map(plt.hist, 'Age', alpha=.5, bins=20)
3 | grid.add_legend()
```

```
/databricks/python/lib/python3.9/site-packages/seaborn/axisgrid.py:337: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
  warnings.warn(msg, UserWarning)
```



```
Command took 1.80 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

```
Cmd 12
```

```
1 | data = [train_df, test_df]
2 | for dataset in data:
3 |     dataset['relatives1'] = dataset['SibSp'] + dataset['Parch']
4 |     dataset.loc[(dataset['relatives1'] > 0, 'not_alone')] = 0
5 |     dataset.loc[(dataset['relatives1'] == 0, 'not_alone')] = 1
6 |     dataset['not_alone'] = dataset['not_alone'].astype(int)
```

```
Command took 0.09 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

```
Cmd 13
```

```
1 | train_df['not_alone'].value_counts()
```

```
Out[40]: 1    537
```

```
0    354
```

```
Name: not_alone, dtype: int64
```

```
Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 14

1 axes = sns.factorplot('relatives','Survived',
2                         data=train_df, aspect = 2.5, )
3
4 #databricks/python/lib/python3.9/site-packages/seaborn/categorical.py:3717: UserWarning: The 'factorplot' function has been renamed to 'catplot'. The original name will be removed in a future release. Please update your code. Note that the default 'kind' in 'factorplot' ('point') has changed 'strip' in 'catplot'.
5 #warnings.warn(msg)
6 #databricks/python/lib/python3.9/site-packages/seaborn/_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
7 #warnings.warn()
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
279
280
281
282
283
284
285
286
287
288
289
289
290
291
292
293
294
295
296
297
298
299
299
300
301
302
303
304
305
306
307
308
309
309
310
311
312
313
314
315
316
317
318
319
319
320
321
322
323
324
325
326
327
328
329
329
330
331
332
333
334
335
336
337
338
339
339
340
341
342
343
344
345
346
347
348
349
349
350
351
352
353
354
355
356
357
358
359
359
360
361
362
363
364
365
366
367
368
369
369
370
371
372
373
374
375
376
377
378
379
379
380
381
382
383
384
385
386
387
388
389
389
390
391
392
393
394
395
396
397
398
399
399
400
401
402
403
404
405
406
407
408
409
409
410
411
412
413
414
415
416
417
418
419
419
420
421
422
423
424
425
426
427
428
429
429
430
431
432
433
434
435
436
437
438
439
439
440
441
442
443
444
445
446
447
447
448
449
449
450
451
452
453
454
455
456
457
458
459
459
460
461
462
463
464
465
466
467
468
469
469
470
471
472
473
474
475
476
477
478
479
479
480
481
482
483
484
485
486
487
487
488
489
489
490
491
492
493
494
495
496
497
497
498
499
499
500
501
502
503
504
505
506
507
508
509
509
510
511
512
513
514
515
516
517
517
518
519
519
520
521
522
523
524
525
526
527
528
529
529
530
531
532
533
534
535
536
537
538
539
539
540
541
542
543
544
545
546
547
547
548
549
549
550
551
552
553
554
555
556
557
558
559
559
560
561
562
563
564
565
566
567
568
569
569
570
571
572
573
574
575
576
577
578
579
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
597
598
599
599
600
601
602
603
604
605
606
607
608
609
609
610
611
612
613
614
615
616
617
617
618
619
619
620
621
622
623
624
625
626
627
628
629
629
630
631
632
633
634
635
636
637
638
639
639
640
641
642
643
644
645
646
647
647
648
649
649
650
651
652
653
654
655
656
657
658
659
659
660
661
662
663
664
665
666
667
668
669
669
670
671
672
673
674
675
676
677
678
679
679
680
681
682
683
684
685
686
687
687
688
689
689
690
691
692
693
694
695
696
697
697
698
699
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
787
788
789
789
790
791
792
793
794
795
796
797
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
888
889
889
890
891
892
893
894
895
896
897
897
898
899
899
900
```

```

15     # filling NaN with 0, to get safe
16     dataset['Title'] = dataset['Title'].fillna(0)

Command took 0.09 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 25

1 train_df = train_df.drop(['Name'], axis=1)
2 test_df = test_df.drop(['Name'], axis=1)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 26

1 genders = {"male": 0, "female": 1}
2 data = [train_df, test_df]
3
4 for dataset in data:
5     dataset['Sex'] = dataset['Sex'].map(genders)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 27

1 train_df['Ticket'].describe()

Out[54]: count    891
unique      681
top     347082
freq       7
Name: Ticket, dtype: object

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 28

1 train_df = train_df.drop(['Ticket'], axis=1)
2 test_df = test_df.drop(['Ticket'], axis=1)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 29

1 ports = {"S": 0, "C": 1, "Q": 2}
2 data = [train_df, test_df]
3
4 for dataset in data:
5     dataset['Embarked'] = dataset['Embarked'].map(ports)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 30

1 data = [train_df, test_df]
2 for dataset in data:
3     dataset['Age'] = dataset['Age'].astype(int)
4     dataset.loc[dataset['Age'] <= 11, 'Age'] = 0
5     dataset.loc[(dataset['Age'] > 11) & (dataset['Age'] <= 18), 'Age'] = 1
6     dataset.loc[(dataset['Age'] > 18) & (dataset['Age'] <= 22), 'Age'] = 2
7     dataset.loc[(dataset['Age'] > 22) & (dataset['Age'] <= 27), 'Age'] = 3
8     dataset.loc[(dataset['Age'] > 27) & (dataset['Age'] <= 33), 'Age'] = 4
9     dataset.loc[(dataset['Age'] > 33) & (dataset['Age'] <= 40), 'Age'] = 5
10    dataset.loc[(dataset['Age'] > 40) & (dataset['Age'] <= 66), 'Age'] = 6
11    dataset.loc[ dataset['Age'] > 66, 'Age'] = 6

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 31

1 # let's see how it's distributed
2 train_df['Age'].value_counts()

Out[58]: 6    165
4   153
5   147
3   143
2   117
1   98
0    68
Name: Age, dtype: int64

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 32

1 train_df.head(18)

   Survived Pclass Sex Age SibSp Parch Fare Embarked relatives not_alone Deck Title
0         0     1.0   3  0.0  2.0  1.0  0.0    7.0        0  1.0  0.0  8.0  1.0
1         1     1.0   1  1.0  5.0  1.0  0.0   71.0       1  1.0  0.0  3.0  3.0
2         1     1.0   3  1.0  3.0  0.0  0.0    7.0       0  0.0  1.0  8.0  2.0
3         1     1.0   1  1.0  5.0  1.0  0.0   53.0       0  1.0  0.0  3.0  3.0
4         0     1.0   3  0.0  5.0  0.0  0.0    8.0       0  0.0  1.0  8.0  1.0
5         0     3.0   0  0.0  3.0  0.0  0.0    8.0       2  0.0  1.0  8.0  1.0
6         0     1.0   1  0.0  6.0  0.0  0.0   51.0       0  0.0  1.0  5.0  1.0
7         0     3.0   0  0.0  3.0  0.0  1.0   21.0       0  4.0  0.0  8.0  4.0
8         1     3.0   1  1.0  3.0  0.0  2.0   11.0       0  2.0  0.0  8.0  3.0
9         1     2.0   1  1.0  1.0  0.0  0.0   30.0       1  1.0  0.0  8.0  3.0

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 33

1 data = [train_df, test_df]
2
3 for dataset in data:
4     dataset.loc[dataset['Fare'] <= 7.91, 'Fare'] = 0
5     dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1
6     dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare'] = 2
7     dataset.loc[(dataset['Fare'] > 31) & (dataset['Fare'] <= 99), 'Fare'] = 3
8     dataset.loc[(dataset['Fare'] > 99) & (dataset['Fare'] <= 250), 'Fare'] = 4
9     dataset.loc[ dataset['Fare'] > 250, 'Fare'] = 5
10    dataset['Fare'] = dataset['Fare'].astype(int)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 34

1 data = [train_df, test_df]
2
3 for dataset in data:
4     dataset['Age_Class']= dataset['Age']* dataset['Pclass']

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 35

1 for dataset in data:
2     dataset['Fare_Per_Person'] = dataset['Fare']/(dataset['relatives']+1)
3     dataset['Fare_Per_Person'] = dataset['Fare_Per_Person'].astype(int)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 36

1 # let's take a last look at the training set, before we start training the models.
2 train_df.head(28)

   Survived Pclass Sex Age SibSp Parch Fare Embarked relatives not_alone Deck Title Age_Class Fare_Per_Person
0         0     1.0   3  0.0  2.0  1.0  0.0    0.0       0  1.0  0.0  8.0  1.0          6            0
1         1     1.0   1  1.0  5.0  1.0  0.0   3.0       1  1.0  0.0  3.0  3.0          5            1
2         1     1.0   3  1.0  3.0  0.0  0.0    0.0       0  0.0  1.0  8.0  2.0          9            0
3         1     1.0   1  1.0  5.0  1.0  0.0   3.0       0  1.0  0.0  3.0  3.0          5            1
4         0     1.0   3  0.0  5.0  0.0  0.0    0.0       1  0.0  0.0  1.0  8.0  1.0          15           1
5         0     3.0   0  0.0  3.0  0.0  0.0    0.0       1  2.0  0.0  1.0  8.0  1.0          9            1
6         0     1.0   1  0.0  6.0  0.0  0.0    3.0       0  0.0  1.0  5.0  1.0          6            3
7         0     3.0   0  0.0  3.0  1.0  2.0    0.0       4  0.0  0.0  8.0  4.0          0            0
8         1     3.0   1  1.0  3.0  0.0  2.0    1.0       0  2.0  0.0  8.0  3.0          9            0
9         1     2.0   1  1.0  1.0  1.0  0.0    2.0       1  1.0  0.0  8.0  3.0          2            1
10        1    1.0   3  1.0  0.0  1.0  1.0    2.0       0  2.0  0.0  7.0  2.0          0            0
11        1    1.0   1  1.0  6.0  0.0  0.0    2.0       0  0.0  1.0  3.0  2.0          6            2
12        0    1.0   0  0.0  2.0  0.0  0.0    1.0       0  0.0  1.0  8.0  1.0          6            1
13        0    3.0   0  0.0  5.0  1.0  5.0    2.0       0  6.0  0.0  8.0  1.0          15           0
14        0    3.0   1  1.0  0.0  0.0  0.0    0.0       0  0.0  1.0  8.0  2.0          3            0
```

```

15   1   2   1   6   0   0   2   0   0   1   8   3   12   2
16   0   3   0   0   4   1   2   2   5   0   8   4   0   0
17   1   2   0   5   0   0   1   0   0   1   8   1   10   1
18   0   3   1   4   1   0   2   0   1   0   8   3   12   1
19   1   3   1   5   0   0   0   1   0   1   8   3   15   0

```

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

Cmd 37

```

1 X_train = train_df.drop("Survived", axis=1)
2 Y_train = train_df["Survived"]
3 X_test = test_df.drop("PassengerId", axis=1).copy()

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 38

```

1 # stochastic gradient descent (SGD) learning
2 sgd = linear_model.SGDClassifier(max_iter=5, tol=None)
3 sgd.fit(X_train, Y_train)
4 Y_pred = sgd.predict(X_test)
5
6 sgd.score(X_train, Y_train)
7
8 acc_sgd = round(sgd.score(X_train, Y_train) * 100, 2)
9
10 print(round(acc_sgd,2), "%")

77.89 %

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 39

```

1 # Random Forest
2 random_forest = RandomForestClassifier(n_estimators=100)
3 random_forest.fit(X_train, Y_train)
4
5 Y_prediction = random_forest.predict(X_test)
6
7 random_forest.score(X_train, Y_train)
8 acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)
9 print(round(acc_random_forest,2), "%")

92.48 %

Command took 0.39 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 40

```

1 # Logistic Regression
2 logreg = LogisticRegression()
3 logreg.fit(X_train, Y_train)
4
5 Y_pred = logreg.predict(X_test)
6
7 acc_log = round(logreg.score(X_train, Y_train) * 100, 2)
8 print(round(acc_log,2), "%")

81.59 %
/databricks/python/lib/python3.9/site-packages/sklearn/linear_model/_logistic.py:763: ConvergenceWarning: lbfsgs failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result()

Command took 0.20 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 41

```

1 # KNN
2 knn = KNeighborsClassifier(n_neighbors = 3)
3 knn.fit(X_train, Y_train)
4
5 Y_pred = knn.predict(X_test)
6
7 acc_knn = round(knn.score(X_train, Y_train) * 100, 2)
8 print(round(acc_knn,2), "%")

86.31 %

Command took 0.19 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 42

```

1 # Gaussian Naive Bayes
2 gaussian = GaussianNB()
3 gaussian.fit(X_train, Y_train)
4
5 Y_pred = gaussian.predict(X_test)
6
7 acc_gaussian = round(gaussian.score(X_train, Y_train) * 100, 2)
8 print(round(acc_gaussian,2), "%")

77.44 %

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 43

```

1 # Perceptron
2 perceptron = Perceptron(max_iter=5)
3 perceptron.fit(X_train, Y_train)
4
5 Y_pred = perceptron.predict(X_test)
6
7 acc_perceptron = round(perceptron.score(X_train, Y_train) * 100, 2)
8 print(round(acc_perceptron,2), "%")

63.75 %
/databricks/python/lib/python3.9/site-packages/sklearn/linear_model/_stochastic_gradient.py:574: ConvergenceWarning: Maximum number of iteration reached before convergence. Consider increasing max_iter to improve the fit.
warnings.warn("Maximum number of iteration reached before convergence. Consider increasing max_iter to improve the fit.")

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 44

```

1 # Linear SVC
2 linear_svc = LinearSVC()
3 linear_svc.fit(X_train, Y_train)
4
5 Y_pred = linear_svc.predict(X_test)
6
7 acc_linear_svc = round(linear_svc.score(X_train, Y_train) * 100, 2)
8 print(round(acc_linear_svc,2), "%")

81.48 %
/databricks/python/lib/python3.9/site-packages/sklearn/svm/_base.py:985: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
warnings.warn("Liblinear failed to converge, increase ")

Command took 0.19 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 45

```

1 # Decision Tree
2 decision_tree = DecisionTreeClassifier()
3 decision_tree.fit(X_train, Y_train)
4
5 Y_pred = decision_tree.predict(X_test)
6
7 acc_decision_tree = round(decision_tree.score(X_train, Y_train) * 100, 2)
8 print(round(acc_decision_tree,2), "%")

92.48 %

Command took 0.09 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster

```

Cmd 46

```

1 results = pd.DataFrame({
2     'Model': ['Support Vector Machines', 'KNN', 'Logistic Regression',
3               'Random Forest', 'Naive Bayes', 'Perceptron',
4               'Stochastic Gradient Descent',
5               'Decision Tree'],
6     'Score': [acc_linear_svc, acc_knn, acc_log,
7               acc_random_forest, acc_gaussian, acc_perceptron,
8               acc_sgd, acc_decision_tree])
9 result_df = results.sort_values(by='Score', ascending=False)
10 result_df = result_df.set_index('Score')
11 result_df.head(8)


```

```

Model
Score
92.48 Random Forest
92.48 Decision Tree
86.31 KNN
81.59 Logistic Regression
81.48 Support Vector Machines
77.89 Stochastic Gradient Descent
77.44 Naive Bayes
63.75 Perceptron

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 47

1 from sklearn.model_selection import cross_val_score
2 rf = RandomForestClassifier(n_estimators=100)
3 scores = cross_val_score(rf, X_train, Y_train, cv=10, scoring = "accuracy")

Command took 2.20 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 48

1 print("Scores:", scores)
2 print("Mean:", scores.mean())
3 print("Standard Deviation:", scores.std())

Scores: [0.77777778 0.82622472 0.74157303 0.84269663 0.85393258 0.85393258
0.79775281 0.7752809 0.83146967 0.83146967]
Mean: 0.8126902384519351
Standard Deviation: 0.03600879831824795

Command took 0.09 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 49

1 importances = pd.DataFrame({'feature':X_train.columns,'importance':np.round(random_forest.feature_importances_,3)})
2 importances = importances.sort_values('importance',ascending=False).set_index('feature')

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 50

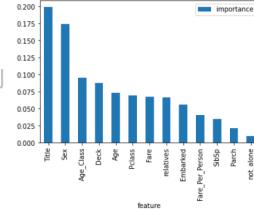
1 importances.head(15)

importance
feature
Title      0.199
Sex        0.174
Age_Class   0.096
Deck       0.088
Age        0.073
Pclass      0.070
Fare        0.068
relatives   0.067
Embarked    0.066
Fare_Per_Person  0.041
SibSp       0.035
Parch       0.022
not_alone   0.010

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 51

1 importances.plot.bar()

Out[78]: <AxesSubplot:xlabel='feature'>


Feature importance distribution. The x-axis lists 13 features: Title, Sex, Age_Class, Deck, Age, Pclass, Fare, relatives, Embarked, Fare_Per_Person, SibSp, Parch, and not_alone. The y-axis represents the importance score, ranging from 0.000 to 0.200. The bars are blue and show a decreasing trend from left to right. The most important feature is 'Title' at approximately 0.199, followed by 'Sex' at 0.174, and so on down to 'not_alone' at 0.010.

Command took 0.40 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 52

1 train_df = train_df.drop("not_alone", axis=1)
2 test_df = test_df.drop("not_alone", axis=1)
3
4 train_df = train_df.drop("Parch", axis=1)
5 test_df = test_df.drop("Parch", axis=1)

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 53

1 # Random Forest
2
3 random_forest = RandomForestClassifier(n_estimators=100, oob_score = True)
4 random_forest.fit(X_train, Y_train)
5 Y_prediction = random_forest.predict(X_test)
6
7 random_forest.score(X_train, Y_train)
8
9 acc_random_forest = round(random_forest.score(X_train, Y_train) * 100, 2)
10 print(round(acc_random_forest,2), "%")

92.48 %

Command took 0.68 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 54

1 print("oob score:", round(random_forest.oob_score_, 4)*100, "%")

oob score: 81.82000000000001 %

Command took 0.10 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 55

1 # Random Forest
2 random_forest = RandomForestClassifier(criterion = "gini",
3                                         min_samples_leaf = 1,
4                                         min_samples_split = 10,
5                                         n_estimators=100,
6                                         max_features="auto",
7                                         oob_score=True,
8                                         random_state=1,
9                                         n_jobs=-1)
10
11 random_forest.fit(X_train, Y_train)
12 Y_prediction = random_forest.predict(X_test)
13
14 random_forest.score(X_train, Y_train)
15
16 print("oob score:", round(random_forest.oob_score_, 4)*100, "%")

oob score: 82.49 %

Command took 0.49 seconds -- by kisung.park@jjsu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 56

1 from sklearn.model_selection import cross_val_predict
2 from sklearn.metrics import confusion_matrix
3 predictions = cross_val_predict(random_forest, X_train, Y_train, cv=3)
4 confusion_matrix(Y_train, predictions)

```

```

Out[83]: array([[487,  62,
   [101, 241]])

Command took 2.30 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 57

1 from sklearn.metrics import precision_score, recall_score
2
3 print("Precision:", precision_score(Y_train, predictions))
4 print("Recall:", recall_score(Y_train, predictions))

Precision: 0.7953795379537953
Recall: 0.7046783625738995

```

```

Command took 0.09 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 58

1 from sklearn.metrics import f1_score
2
3 f1_score(Y_train, predictions)

Out[85]: 0.7472868217854265

```

```

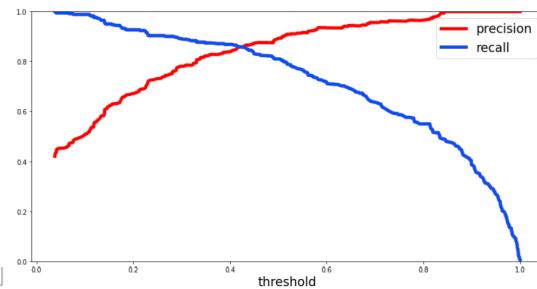
Command took 0.09 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 59

1 from sklearn.metrics import precision_recall_curve
2
3 # getting the probabilities of our predictions
4 y_scores = random_forest.predict_proba(X_train)
5 y_scores = y_scores[:,1]
6
7 precision, recall, threshold = precision_recall_curve(Y_train, y_scores)

Command took 0.10 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 60

1 def plot_precision_and_recall(precision, recall, threshold):
2     plt.plot(threshold, precision[-1], "r-", label="precision", linewidth=5)
3     plt.plot(threshold, recall[-1], "b-", label="recall", linewidth=5)
4     plt.xlabel("threshold", fontsize=19)
5     plt.legend(loc="upper right", fontsize=19)
6     plt.ylim([0, 1])
7
8 plt.figure(figsize=(14, 7))
9 plot_precision_and_recall(precision, recall, threshold)
10 plt.show()

```

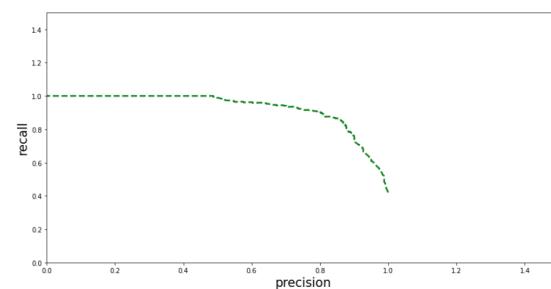


```

Command took 0.30 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 61

1 def plot_precision_vs_recall(precision, recall):
2     plt.plot(recall, precision, "g-", linewidth=2.5)
3     plt.xlabel("recall", fontsize=19)
4     plt.ylabel("precision", fontsize=19)
5     plt.axis([0, 1.5, 0, 1.5])
6
7 plt.figure(figsize=(14, 7))
8 plot_precision_vs_recall(precision, recall)
9 plt.show()

```



```

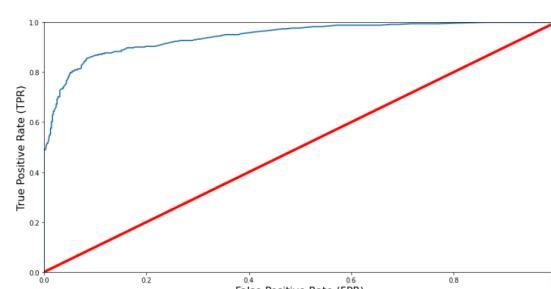
Command took 0.30 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 62

1 from sklearn.metrics import roc_curve
2
3 # compute true positive rate and false positive rate
4 false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_train, y_scores)

Command took 0.09 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 63

1 # plotting them against each other
2 def plot_roc_curve(false_positive_rate, true_positive_rate, label=None):
3     plt.plot(false_positive_rate, true_positive_rate, linewidth=2, label=label)
4     plt.plot([0, 1], [0, 1], 'r', linewidth=4)
5     plt.xlabel("False Positive Rate (FPR)", fontsize=16)
6     plt.ylabel("True Positive Rate (TPR)", fontsize=16)
7     plt.ylim([0, 1])
8
9 plt.figure(figsize=(14, 7))
10 plot_roc_curve(false_positive_rate, true_positive_rate)
11 plt.show()

```



```

Command took 0.30 seconds -- by kisung.park@jisu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
Cmd 64

```

```
1 from sklearn.metrics import roc_auc_score
2 r_a_score = roc_auc_score(y_train, y_scores)
3 print("ROC-AUC-Score:", r_a_score)

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

Cmd 65

```
1 submission = pd.DataFrame(
2     {"PassengerId": test_df["PassengerId"],
3      "Survived": Y_prediction
4      })
5 submission.to_csv('submission.csv', index=False)

Command took 0.10 seconds -- by kisung.park@sjtu.edu at 10/30/2022, 10:32:14 AM on [default]basic-starter-cluster
```

Shift+Enter to run