簡報閱讀

重要知識點

什麼是 kaggle?

鐵達尼號生存預測

學習心得(完成)

AI共學社群 > Python資料科學程式馬拉松 > 探索性資料分析(EDA)_數據理解與重覆和遺失值處理



囯

問題討論

探索性資料分析(EDA)_數據理解與重覆和遺失值處理

範例與作業







- Kaggle 早期因與各領域的公司單位合作舉辦數據分析競賽而出名;
- 競賽的舉辦由與其合作的公司或單位來定義想解決的問題並提供相關數據資料,然後開放給各路好手 建立解決問題的預測模型。

Compete Datasets Notebooks Discuss Courses ***

Notebooks environment. Access free GPUs and a huge repository of community published data &



from sklearn.metrics import confusion_matrix from sklearn.utile.multiclass import unique_labels

Join Competition

More



訓練資料集,訓練模型

Kaggle 上傳的資料格式。

Detail

About this file

意義

測試資料集,以此資料透過模型,產生預測值

Discussion Leaderboard Rules Team

Compact

< gender_submission.csv (3.18 KE

Column

離散

離散

Ahoy, welcome to Kaggle! You're in the right place.

competitions and familiarize yourself with how the Kaggle platform works.

survived the Titanic shipwreck.

This is the legendary Titanic ML competition – the best, first challenge for you to dive into ML

The competition is simple: use machine learning to create a model that predicts which passengers

Overview Data Notebooks **Data Explorer**

90.9 KB

test.csv

train.csv

資料架構。

train.cvs

test.csv



大此和土的流体法 不应先大海州传山田 遺

為什麼要處理遺失值?

載入鐵達尼資料,發現有遺失值

Cabin

Embarked

(Missing not an Random)

處理遺失值的方法 - 刪除

處理遺失值的方法,可以分為刪除和補值兩類。

非隨機缺失

刪除

 刪除有缺失資料的樣本 刪除有過多缺失資料的變數 (通常超過 60% 就會刪除) 缺點: # 會導致資訊丟失。 # 不是完全隨機缺失,若採用簡單刪除法就會使得估計係數出現偏誤。 				
處理遺失值的方法 - 補值				
 補值 給定一個固定值去填補遺失值 由後往前補值 (時間性相關適用) 由前往後補值 (時間性相關適用) 用現有的資料取平均值、中位數、眾數等進行補值 用預測方法補值,迴歸或機器學習 				
回到今天的課程範例				

運用 python , 偵測與判斷是否有重複或遺失值 ,以實際資料 ,如何進行資料遺失值處理 。

#顯示有重覆的資訊:

df_train['Age']=df_train['Age'].fillna(df_train['Age'].mean())

知識點回顧

df_train['Age']=df_train['Age'].fillna(method='pad')

條件隨機缺失與非隨機缺失 • 運用 python ,處理資料中重覆和遺失值處理,可透過刪除和補值的技巧處理遺失值。

延伸閱讀

1. 下載資料(資料與 gender_submission.csv) 2. 撰寫程式建立模型,產生預測結果 3. 上傳 submission.csv and "Make Submission"

- Getting Started Prediction Competition Titanic: Machine Learning from Disaster Start here! Predict survival on the Titanic and get familiar with ML basics

Overview Data Notebooks Discussion Leaderboard Rules Team

Kaggle · 21,860 teams · Ongoing

My Submissions

Submit Predictions

網站: <u>kaggle</u>

Predict Malicious Websites: XGBoost
Pytina roteback using data have Melicious Start with more than import numpy as no import pandas as pd a blinking cursor Kaggle offers a no-setup, customizable, Jupyter

Home ♀ Compete □ Data Start here! Predict survival on the Titanic and get familiar with ML basics <> Notebooks Discuss Kaggle · 18,310 teams · Ongoing ○ Courses Overview Data Notebooks Discussion Leaderboard Rules 🖆 Jobs Overview

Description

Evaluation

Frequently Asked Questions

資料描述 • 鐵達尼號沈船事件,發生在1912年4月15日 • 船上共 2224 名乘客, 共 1502 名死亡。 • 在過程中,發現某些族群容易存活下來。

資料

Gender_submission.csv

gender_submission.csv

Fare 字串 乘客票價

船艙號碼

登船港口

有些程:	式與演算法,不容許在	有遺失值出現			
		的估計,讓程式執行下			
遺失值特性的]分類				
遺失值根據遺失的特性,大致上可以分成以下三種情形:					
遺失值根據遺失的	的特性,大致上可以分	·成以下三種情形:			
遺失值根據遺失的 遺失類		成以下三種情形: 定義		例子	
	型 缺失資料			,不小一掉了幾張考 的遺失和成績無關,	

缺失資料依賴於該變數本身。

別相關。

有分收入的問卷調查,通常薪資高

的人,不喜歡填有關收入資訊,所

以收入資料缺失和收入高低有關。

Python 語法

• 觀察-是否有重覆

• 由前往後補植:

df_train.duplicated() • 用平均值補值

df_train['Age']=df_train['Age'].fillna(method='bfill') • 由後往前補:

• 以 kaggle 鐵達尼資料集為例,理解資料科學數據議題

• 遺失值有三大種類,包含完全隨機缺失

Kaggle 參與比賽的流程

資料來源: kaggle

Overview Description Evaluation

下一步:閱讀範例與完成作業