

**Data Cleaning Challenge: Handling missing values**

Python notebook using data from [multiple data sources](#) · 78,787 views · dailychallenge, multiple data sources

602

Fork

13429



Version 8

8 commits

Notebook

Data

Log

Comments






This kernel has been released under the [Apache 2.0](#) open source license.


**Did you find this Kernel useful?**  
Show your appreciation with an upvote


  
602





Data Sources


▼  Detailed NFL Play-by-...


 NFL... 362k x 102

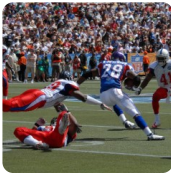
 NF... 408k x 102

 NF... 449k x 255

▼  San Francisco Buildin...

 Build... 199k x 43

 DataDict... 43 x 3



Detailed NFL Play-by-Play Data  
2009-2018

nflscrapR generated NFL dataset wiith expected  
points and win probability

Last Updated: a month ago (Version 6)

About this Dataset

Introduction

The lack of publicly available National Football League (NFL) data sources has been a major obstacle in the creation of modern, reproducible research in football analytics. While clean play-by-play data is available via open-source software packages in other sports (e.g. nhlscrapr for hockey; PitchF/x data in baseball; the Basketball Reference for basketball), the equivalent datasets are not freely available for researchers interested in the statistical analysis of the NFL. To solve this issue, a group of [Carnegie Mellon University statistical researchers](#) including Maksim Horowitz, Ron Yurko, and Sam Ventura, built and released [nflscrapR](#) an R package which uses an API maintained by the NFL to scrape, clean, parse, and output clean datasets at the individual play, player, game, and season levels. Using the data outputted by the package, the trio went on to develop reproducible methods for building expected point and win probability models for the NFL. The outputs of these models are included in this dataset and can be accessed using the nflscrapR package.

Content

The dataset made available on Kaggle contains all the regular season plays from the 2009-2016 NFL seasons. The dataset has 356,768 rows and 100 columns. Each play is broken down into great detail containing information on: game situation, players involved, results, and advanced metrics such as expected point and win probability values. Detailed information about the

Run Info

Succeeded	True	Run Time	24.7 seconds
Exit Code	0	Queue Time	0 seconds
Docker Image Name	<a href="#">kaggle/python(Dockerfile)</a>	Output Size	0
Timeout Exceeded	False	Used All Space	False
Failure Message			

```
Time   Line #   Log Message
3.7s   1         [NbConvertApp] Converting notebook script.ipynb to html
3.7s   2         [NbConvertApp] Executing notebook with kernel: python3
23.9s  3         [NbConvertApp] Writing 294163 bytes to __results__.html
23.9s  4
23.9s  6         Complete. Exited with code 0.
```

Comments (317)

Filter/sort >

Please [sign in](#) to leave a comment.

Wei Chen • 6 months ago

1

Very useful, thank you!

Options

perkleeee • 6 months ago

1

Thank you for great tutorial !

Options

Cyringa • 6 months ago

1

Very useful challenge for novice. Thank you!

Options

VSS • 6 months ago

1

Completed Day 1 Challenge.  
Thankyou Rachael . It is really helpful and informative :) .

Options

Joshua • 6 months ago

1

Rachel, thanks for putting this challenge together! It is perfect for beginners like myself.

Options

shreya • 7 months ago

3

Hi Rachael! I think that was a wonderful introduction to data cleaning.  
Could you please explain the below syntax and its important uses?  
np.random.seed(0)  
Thanks!

Options	Rachael Tatman · 7 months ago		<div><div>^</div><div>3</div><div>▼</div></div>	
Thank you Rachael!				
Options	JamesZhou · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
THX ! :) I solved it				
Options	Marina · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
thanks! I solved it				
Options	tecnoholic · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Good tutorial Rachel. Thank you.				
Options	Imran Iskandar · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Thank you! Waiting for tommorow's challange!				
Options	Jake Nocentini · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Thanks for this tutorial!				
Options	Artem Prosvirnin · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Thank you Rachael. The tutorial was very useful.				
Options	Lim Chia Hoon · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Thanks for this tutorial!				
Options	Mandy · 7 months ago		<div><div>^</div><div>1</div><div>▼</div></div>	
Thanks Rachael! Here is my work: <a href="https://www.kaggle.com/tcmandyw/data-cleaning-challenge-handling-missing-values">https://www.kaggle.com/tcmandyw/data-cleaning-challenge-handling-missing-values</a> Keep cracking down other challenges^^				

Options	r w	
	e	
<b>Parichita Nandi</b> • 7 months ago	d o s o	<div><div>^</div><div>1</div><div>v</div></div>
Options	m e	
	t	
<b>Mohd Akhtarh</b> • 7 months ago	n i	<div><div>^</div><div>1</div><div>v</div></div>
Options	g w	
	i	
<b>shreya</b> • 7 months ago	t h a n	<div><div>^</div><div>1</div><div>v</div></div>
Options	e l	
	e	
<b>zhouleyuan</b> • 7 months ago	m e	<div><div>^</div><div>1</div><div>v</div></div>
Options	d o	
	n	
<b>Rachael Tatman</b> • 7 months ago	e p s s ( n t - k e g s o n u r f n - i y n s g d t m h e c c o r	<div><div>^</div><div>0</div><div>v</div></div>
<b>Tarak</b> • 8 months ago	s i s ( n t - k e g s o n u r f n - i y n s g d t m h e c c o r	<div><div>^</div><div>1</div><div>v</div></div>
Options	n k e g s o n u r f n - i y n s g d t m h e c c o r	
<b>shonafeng41</b> • 8 months ago	n u r f n - i y n s g d t m h e c c o r	<div><div>^</div><div>1</div><div>v</div></div>
Options	n u r f n - i y n s g d t m h e c c o r	
<b>Mahbub</b> • 8 months ago	n u r f n - i y n s g d t m h e c c o r	<div><div>^</div><div>1</div><div>v</div></div>
Options	n u r f n - i y n s g d t m h e c c o r	

<div>Berk Atik</div> <div>Options</div>		<div> <div>9 months ago</div> <div>Your 5-day challenges are awesome for the beginners like me, it helps to wrap my head around the concept. Thank you a lot :)</div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>
<div>Akash Agarwal</div> <div>Options</div>		<div> <div>9 months ago</div> <div>Thanks a lot! this was very helpful.</div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>
<div>DavinderKumar</div> <div>Options</div>		<div> <div>10 months ago</div> <div>Great learning ....waiting eagerly for tomorrow's session :)</div> </div>	<div> <div>^</div> <div>3</div> <div>v</div> </div>
<div>Prableen Kaur</div> <div>Options</div>		<div> <div>10 months ago</div> <div>           Nice tutorial Rachel!            I like how you are trying to get us to build an intuition of how to do "imputation".            Looking forward to tomorrow's lesson!            Here's my solution: <a href="https://www.kaggle.com/prableen20/data-cleaning-challenge-handling-missing-values">https://www.kaggle.com/prableen20/data-cleaning-challenge-handling-missing-values</a> </div> </div>	<div> <div>^</div> <div>4</div> <div>v</div> </div>
<div>Hannah Catri</div> <div>Options</div>		<div> <div>10 months ago</div> <div>Also, I hope to participate in more challenges in future :)</div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>
<div>Hannah Catri</div> <div>Options</div>		<div> <div>10 months ago</div> <div>Thanks for creating the 5 day challenge, I learned a lot!</div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>
<div>Gopesh Dwivedi</div> <div>Options</div>		<div> <div>10 months ago</div> <div>This was a great tutorial, thanks :)</div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>
<div>pawan bansal</div> <div>Options</div>		<div> <div>10 months ago</div> <div> <a href="https://github.com/hackbansu/kaggle-5-day-data-cleaning-challenge">https://github.com/hackbansu/kaggle-5-day-data-cleaning-challenge</a>            All my notebooks link.            Thanks Rachael for this great challenge.         </div> </div>	<div> <div>^</div> <div>1</div> <div>v</div> </div>



Options

gcoolt • 10 months ago

This was very helpful. Thanks

Options

Zayed Shah • 10 months ago

Thanks Rachael

Options

Praxitelis-Nikolaos Kouroupetroglou • 10 months ago

thank you Rachael for the data-cleaning challenge day-1, capital work!

Options

Vibha • 10 months ago

Really good introductory tutorial Rachel! Having done these basics and watching them execute is a tiny first step in my Data Science journey. Thank you :)

Options

Options

vica • 10 months ago

Nice work thank you.

Options

Arvind Kumar • 10 months ago

This Topic really makes a sense at Kaggle...

Options

Rich King • 10 months ago

Thanks for this tutorial. It's super helpful!

Options

Ryan Juliano • 10 months ago

to echo many other comments, this was a nice tutorial. Thanks for investing the time to make it. I look forward to completing the subsequent parts (which I've sadly fallen a bit behind on!)

Options

Aman Saxena

• 10 months ago

1

I'm not able to find an R notebook for this, can anyone help me out? @Rachael, you've only uploaded the Python notebook, right ?

Options

Rachael Tatman

• 10 months ago

1

Shehjar Kaul

• 10 months ago

1

Great tutorial to start with! [Here](#) is my code!

Options

Vipin Kumar

• 10 months ago

1

Thank you for the great tutorial, and providing the suggestion for further learning from Dan's notebook.

Options

jaludac

• 10 months ago

1

It's great!  
Thank you Rachael

Options

konichuvak

• 10 months ago

1

Handling NaN's is crucial! Thanks for doing the series!

Options

Samarendra P

• 10 months ago

1

Thanks Rachel !! The content was great and the datasets were good !!

Options

Rahul Singh

• 10 months ago

1

Is anyone trying these out in R?

Options

Rachael Tatman

• 10 months ago

0

RohithB

• 10 months ago

1

Thanks for the great tutorial @Rachael. Please check my Day1 solution : [Handling missing values](#)

Options		
rory81 • 10 months ago		<div>^1v</div>
Thanks Rachel! Not only am I learning a lot about data cleaning, but as a Python newbie also a lot about using Python. So here is my <a href="#">Solution</a>		
Options		
Masum Rumia • 10 months ago		<div>^1v</div>
Thank you so much <a href="#">@Rachel</a> . Here is my <a href="#">Solution</a> . Doing day#2 now and looking forward to the rest of them.		
Options		
John Olayisade • 7 months ago		<div>^0v</div>
Dany Venero • 10 months ago		<div>^1v</div>
Learning a lot...Thanks Rachel!		
Options		
Yong • 10 months ago		<div>^1v</div>
thanks		
Options		
Dare • 10 months ago		<div>^1v</div>
Nice!		
Options		
Hernan Goldenberg • 10 months ago		<div>^1v</div>
Thanks a lot Rachel, I'm taking my first steps in data science..... <a href="https://www.kaggle.com/ruso72/data-cleaning-challenge-handling-missing-v-ea4b91">https://www.kaggle.com/ruso72/data-cleaning-challenge-handling-missing-v-ea4b91</a>		
Options		
Banjo • 10 months ago		<div>^1v</div>
Thanks Rachel! I really enjoyed this first part and hope other parts will be as simple as you did to this first part. Your turn <a href="#">solution link</a>		
Options		
ColaCole • 10 months ago		<div>^1v</div>
My Teacher thank you very much for your education Ms. Rachael Tatman. Great work.		

Options		Options	
Said Aspen		10 months ago	1
Good stuff, thanks!			
Options			
Sabu Joseph		10 months ago	1
Hi Racheal, Updated the notebook with analysis around Street Number Suffix and Zipcode. Analysis around zipcode can be productive since it may be connected with missing Location or inherent issues in the source system that creates sf_permits dataset while collecting these tow fields or if zipcodes are extracted based on Location using wrapper applications. Another assumption is the growing importance of Location coordinates in many applications Will post finding from datasets about San Fransisco on the Datasets listing			
Options			
hariprasaath		10 months ago	1
hope it helps. <a href="https://www.kaggle.com/hariprasaath/data-cleaning-challenge-handling-missing">https://www.kaggle.com/hariprasaath/data-cleaning-challenge-handling-missing</a>			
Options			
Bellamy		10 months ago	1
This is great, thanks for this tutorial. I'm looking forward to the next few days.			
Options			
Jessica Chopra		10 months ago	1
Great Tutorial! I have used tools like Angoss and Enterprise Miner to impute missing values. These tools can create decision trees based upon certain conditions and we can impute missing values based upon that data in the neighborhood records.			
Options			
WangYong		10 months ago	1
Thanks for sharing, it is a great kickoff. <a href="#">Here is my Kernel</a> Features:  1. Visualization for missing data 2. Ignorable & Non-ignorable concept			
Options			

Semloh

• 10 months ago

^

1

v

Thank you Rachael was really a great tutorial. Looking forward for the next lesson.

Options

naveensaliyan

• 10 months ago

^

1

v

Thank you Rachel, good refresher on data cleaning, waiting for advanced cleaning techniques in next tutorials.  
I never knew there is a **sample** method for DF, thanks for introducing :)

Options

Deep Thought

• 10 months ago

^

1

v

Great tutorial. Thanks Rachel!

Options

GSD

• 10 months ago

^

1

v

Great tutorial Rachel..Thanks

Options

Elvirasun28

• 10 months ago

^

1

v

Pretty interesting and useful to a new fresher. Please check my solution  
<https://www.kaggle.com/elvirasun28/data-cleaning-challenge-handling-missing-values>

Options

Punyajoysaha

• 10 months ago

^

1

v

nice way to learn data cleaning

Options

Anthony Wu

• 10 months ago

^

1

v

Very helpful for beginners like me. Thanks!

Options

LuizGermano

• 10 months ago

^

1

v

Thanks Rachel!!! It's an amazing tutorial. I'll try the lessons!

Options

Javed Sheikh

• 10 months ago

^

1

v

Options

Completed!

Options

M.G.Hunter • 10 months ago

1

I'm still new to data science and have little knowledge of coding. Yet, this was easy to follow and understand. Thanks Rachael

Options

Manojna Kalapala • 10 months ago

1

Needed this challenge!!!

Options

Peiqin • 10 months ago

1

Thanks Rachael for the nice tutorial! My solution: <https://www.kaggle.com/peiqin/data-cleaning-challenge-handling-missing-values>

Options

Arash Koupaei • 10 months ago

1

Thank you Rachael - interesting competition - [my solution to the challenge](#)

Options

KotsirasDimitris • 10 months ago

1

Nice tutorial!thank you! <https://www.kaggle.com/jimintel/data-cleaning-challenge-handling-missing-values>

Options

Imron1 • 10 months ago

1

This is very interesting. Thanks Rachel.  
this is my solution:  
<https://www.kaggle.com/imron1/data-cleaning-challenge-handling-missing-v-4ba114>  
feel free to have any feedback.

Options

Ace Palmero • 10 months ago

1

Great Tutorial! I learned a lot. A quick question about imputation. Doesn't such a broad fix add a lot of incorrect information to your dataset. For example, with the street number suffix, a lot of my values in that column are filled as 'A'. I assume the street number suffixes should be different so isn't it preferable to have Nan instead of incorrect values.

Options



Options		
Saranjit saikia • 10 months ago	1	
nice tutorial learned a lot		
Options		
Hugh Zabriskie • 10 months ago	1	
Thanks for the tutorial, Rachel! It seems that the NFL data doesn't have the column descriptions in the documentation as mentioned. The link to the data in the tutorial is: <a href="https://www.kaggle.com/maxhorowitz/nflplaybyplay2009to2016">https://www.kaggle.com/maxhorowitz/nflplaybyplay2009to2016</a> . I went to Data > Column Metadata, but didn't see the column descriptions.		
Options	Options	
	Rachael Tatman • 10 months ago	0
uday • 10 months ago	1	
An excellent tutorial looking forward to next lesson. Thank you, Rachael, for this series. My solution: <a href="https://www.kaggle.com/udaykiran3207/data-cleaning-challenge-handling-missing-values-8af847">https://www.kaggle.com/udaykiran3207/data-cleaning-challenge-handling-missing-values-8af847</a>		
Options		
Prashant Brahmabhatt • 10 months ago	1	
As a rookie this was really really helpful. Thanks!		
Options		
ameer.khosrofi • 10 months ago	1	
Excellent job.. a very helpful tutorial for beginner.		
Options		
snehanshuseengupta • 10 months ago	1	
Nice article		
Options		
Prashant Kishore • 10 months ago	1	
Nlce tutorial. Thank you! <a href="https://www.kaggle.com/prashant001/data-cleaning-challenge-handling-missing-values">https://www.kaggle.com/prashant001/data-cleaning-challenge-handling-missing-values</a>		
Options		



Anand Panchal • 10 months ago

Thank you for the wonderful tutorial Rachael. I am looking forward to the next challenge. Here's my notebook for day 1: <https://www.kaggle.com/apanchal2/data-cleaning-challenge-handling-missing-values>

Options

Jordan • 10 months ago

edit -solved

Options

Atpaul • 10 months ago

Thank you Rachael. Looking forward for Day 2 of the Challenge

Options

Monalisa Mishra • 10 months ago

Really useful in refreshing the concepts. Thanks, Rachael!!  
<https://www.kaggle.com/mmishra90/data-cleaning-challenge-handling-missing-values>

Options

huangkx • 10 months ago

It is easy to have a handle on data cleaning by following the notebook, thanks!!

Options

TheRunningPhysicist • 10 months ago

As a new guy on Kaggle I find this kind of competitions/lessons to be really useful. I really appreciate the work and commitment in this 5-days challenge!  
  
If somebody wants to take a look at my notebook you are more than welcome:  
<https://www.kaggle.com/therunningphysicist/data-cleaning-challenge-handling-missing-values/notebook>  
  
Can't wait for next lesson!

Options

Sidharth • 10 months ago

Great work !!  
I joined Kaggle recently and was skeptical about signing up for this tutorial, but really enjoyed and learnt a lot.  
Looking forward to upcoming tutorials.  
  
Glad to share my work : <https://www.kaggle.com/sidharthsuman/data-cleaning-challenge-handling-missing-values>

Options

lbrxk • 10 months ago

Great Tutorial here is my solution :-)  
<https://www.kaggle.com/ibrahimkochbati/data-cleaning-challenge-handling-missing-values>

Options

Steven Hu • 10 months ago

Thank you, Rachael!! This hands-on practice is super helpful for beginners, especially people who have a dream of data scientist. Thank you! Cannot wait for day two challenge! [here is my notebook for day1](#)

Options

Anthony Jatobá • 10 months ago

Nice tutorial, Rachael! The examples were fine and my poor python wasn't a obstacle.  
Here's my notebook: <https://www.kaggle.com/anthonyjatoba/data-cleaning-challenge-handling-missing-values>

Options

Jonas Conde • 10 months ago

Nice tutorial!  
Is there any way to import the csv files from the OpenAddresses kernel?

Options

[Rachael Tatman](#) • 10 months ago

Options

Saroj • 10 months ago

Thanks Rachael  
<https://www.kaggle.com/saroj13278/data-cleaning-challenge-handling-missing-values>

Options

Patrick Kakou • 10 months ago

Nice challenge, I am just for Day 1 and ready for tomorrow. Thanks Rachael

Options

AhmedSultan • 10 months ago

Thanks for the tutorial! ... waiting for tomorrow's tutorial

Options

Wei Chun Chang • 10 months ago

^

1

▼

Hi Rachael, thanks for the tutorial. I love your challenge :)

Here are my fork notebook and some review. Welcome for any suggestion and comment.

<https://www.kaggle.com/justjun0321/data-cleaning-challenge-day-1-done-review>

Options

SarahXu • 10 months ago

^

1

▼

Thank you ! :)

Options

Sabu Joseph • 10 months ago

^

1

▼

Great experience..thanks for explaining the concepts in deep.

How can I access iowa housing data to try with methods DanB suggested ? i tried with path `pd.readcsv("../input/(iowa-housing-snapshot/iowadata.csv')` but getting file not found error

Options

Rachael Tatman • 10 months ago

^

1

▼

Demongolem • 10 months ago

^

1

▼

I think maybe I am not so familiar with notebooks. When I click in the first cell that I see that loads all the data, I get a `NameError: name 'np' is not defined` error. It says [2] to the left, so is there a [1] somewhere that I am missing?

**EDIT** It works now, but for me it seems I have to run several cells multiple times for the correct action to be taken. For example, I just printed out the pct of missing values for the `sf_permits` data set, but it would not spit out the correct number until I ran it 3 times.

Options

FSBDS\_AndreMarquesLeite • 10 months ago

^

1

▼

Thanks Rachael, Very detailed and clear. I'm also looking fwd for tomorrow's lesson ;)

Options

GranVisir • 10 months ago

^

1

▼

Thank you Rachael! Very Socratic lesson to learn: "Know your data"!

Happy to share my work: [Cleaning Data Challenge: Day 1](#)

Options

John Xu • 10 months ago

^

1

▼

Very enjoyable study. Thanks. Here is my notebook [John's Notebook](#)

Options

Options

Sergey Proshuta • 10 months ago

^

1

▼

Great tutorial and cool start to the first day!

I tried to add some extra things to my solution:

<https://www.kaggle.com/proshuta/data-cleaning-challenge-handling-missing-extra/>

Cannot wait for the next day :)

Options

Tanveer Baba • 10 months ago

^

1

▼

It was great. Honestly, After first lesson i'm excited for tomorrow's lesson

Thank you Rachael

Thank you very much

Options

Shridhar • 10 months ago

^

1

▼

Great learning, thank you.

Regards,

Shridhar

Options

TFu • 10 months ago

^

1

▼

I really like the way reasoning behind the missing value, I usually just remove the rows or columns before this tutorial

Options

SK • 10 months ago

^

1

▼

Thanks for the detailed explanations. Looking fwd to tomorrow's session

Options

AimeShangula • 10 months ago

^

1

▼

thanks for t(o)uts Rachel

Options

manasa • 10 months ago

^

1

▼

Great Tutorial!! Thanks Rachael. Looking forward for more.

Options

Abdelrahman Hamza • 10 months ago

Thank you! amazing day :) Looking forward to tomorrow's lesson.

Options

Xenophontas Psichis • 10 months ago

Really nice and helpful tutorial!

Options

dilex • 10 months ago

Thanks! Very good first steps tutorial.

On a side note : there are missing "s" in sf\_permits in the end of paragraphs 2 &5 ;)

Options

Rachael Tatman • 10 months ago

LennartGrosser • 10 months ago

Cool notebook on an important topic which can really make a difference in terms of accuracy when done well.

I'd further suggest looking at the distribution of a column when replacing missing values.

Basic example: if the values of a column are equally distributed you want to impute differently than if they were normally distributed.

For categorical values one can look at the histogram.

Options

Rachael Tatman • 10 months ago

Andreas Paulin • 10 months ago

Great lesson, thanks a lot! Looking forward to more!

Options

Marco Rigo • 10 months ago

Nice tutorial!  
very clear

Options

Hui Miao • 10 months ago

Thank you Rachael! Great lesson for someone who is new to data cleaning!

Options

Eladio Rego

10 months ago

0

1

Great!!! I do like it!

Options

Rohit Singh Adhikari

10 months ago

1

Thanks for the explanation....will follow the coming sessions!!!

Options

Sunny Chidambaram

10 months ago

1

Great.. :)

Options

Chayan Shrang Raj

10 months ago

1

Awesome Challenge! Simple and effective data cleaning techniques for beginners. Happy to share my notebook <https://www.kaggle.com/chayanraj/data-cleaning-challenge-handling-missing-values>.

Never Stop Learning...

Options

LakshyaRoy

10 months ago

1

Great tutorial

Options

Arunava

10 months ago

1

It was a nice walkthrough. Thanks for this :)

Here's the link to my fork for [Day 1 of Data Cleaning Challenge](#)

Options

ChrisBow

10 months ago

1

Thanks Rachael, I'd been looking a bit of a more structured way to work on my Python (not that I'll forget about you, R!), and this was spot on! Looking forward to the rest of the week.

[Here's my notebook for day 1](#) (I hope I imputed the missing value at the end of one of the questions correctly!)

Options

Ibio

10 months ago

1

useful skills

Options		
CRK • 10 months ago		
Great tutorial Rachel. Thank you. Excited to move on to Day 2! This course reminded me of an interesting project that i was working in 2011 where we had customers in Dublin,Ireland. I will not go in detail about the project but part of it was to deliver points cards to the customers address. I did not knew until then that Dublin,Ireland did not have any postcode... the records in the database showed 'NaN' for few postal codes and few others had 'text and numbers' mixed together. I was trying to figure out for a long time if it was "does not exist" or "wasnt recorded" But later in 2014,ireland started to introduce the postal code. <a href="https://qz.com/272332/ireland-is-just-now-getting-around-to-introducing-postal-codes/">https://qz.com/272332/ireland-is-just-now-getting-around-to-introducing-postal-codes/</a>		
Options	Options	Options
Sachin	Rachael Tatman • 10 months ago	0
Sabu Joseph	• 10 months ago	0
Feiqi Zhang • 10 months ago		
nice teaching style that is well organized and super easy to do and follow...		
Options		
Somrik Banerjee • 10 months ago		
Thank you Rachel for pointing out the subtle difference between data not being recorded and nonexistence. I am familiar with most of this content otherwise. Looking forward to the next session. <a href="https://www.kaggle.com/somrikbanerjee/data-cleaning-challenge-handling-missing-values">https://www.kaggle.com/somrikbanerjee/data-cleaning-challenge-handling-missing-values</a>		
Options		
NehaS • 10 months ago		
Was amazing lesson... thanks a lot. Looking forward to more lessons.		
Options		
Moni Gill • 10 months ago		
Nicely commented and explained. Looking fwd to other sessions.		
Options		
时间刺客 • 10 months ago		
I'm a novice.我觉得这是很好的学习教程。没有什么练习的机会，刚好可以试试手。很期待明天的挑战！！！！		
Options		

nengkuan2

• 10 months ago

Options

1

very good reasoning on why/when/what to drop/impute. But I am still confused about NAN and None

1. It looks to me that dropna() only remove NAN. It does not remove None. Am I right?

2. Is there a way to remove None?

3. My understanding about NAN is that it is a result from floating operation like 0/0. We are also using NAN to represent missing or empty. Does None also mean missing or empty?

Options

junnu

• 10 months ago

Options

3

Niyamat Ullah

• 10 months ago

Options

1

Actual learning is doing in practice. And Kaggle is doing the same thing. Thanks Kaggle!

Options

Jonas Conde

• 10 months ago

Options

0

Weimin Yu

• 10 months ago

Options

1

Tkanks! ,it's very useful for learner! It's my solution :

<https://www.kaggle.com/yuweiming70/data-cleaning-challenge-handling-missing-values>

Options

Djaballah

• 10 months ago

Options

1

Well-explained tutorial, Looking forward to the next challenge.

Here is my notebook

[Data Cleaning Challenge: Handling missing values Day 1](#)

Options

Kelechi

• 10 months ago

Options

2

Great job, 🙌 Rachael. I just published an article about my experience here -

[https://medium.com/@Kelechukwu\\_/a-practical-guide-to-data-cleaning-for-beginners-bd4f45967825](https://medium.com/@Kelechukwu_/a-practical-guide-to-data-cleaning-for-beginners-bd4f45967825)

Options

Richard Allen

• 10 months ago

Options

2

Thanks for the kernel Rachel, just a quick note - for the 'More Practicle!' section, I've found the data for San Francisco addresses to be missing in what seems like a deliberate exclusion. I have a brief analysis here:

<https://www.kaggle.com/bigironsphere/missing-san-francisco-data>

And the missing data itself can be directly downloaded here:

[http://results.openaddresses.io/sources/us/ca/san\\_francisco](http://results.openaddresses.io/sources/us/ca/san_francisco)

Options



Options

Rachael Tatman • 10 months ago

0

Kristen O Anderson • 10 months ago

2

Everything is so well explained and the steps are so logical. Thanks!

Options

Lewis Tunstall • 10 months ago

2

In addition to calculating the total fraction of missing values in a dataset, I often find it's useful to visualise the percentage of missing values per feature. In case this is useful to others, the following code snippet produces a bar plot for the `sf_permits` dataframe:

```
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

# get percentage of missing values
df_na = (sf_permits.isnull().sum() / len(sf_permits)) * 100

# drop columns without 1 values and sort
df_na = df_na.drop(df_na[df_na == 0].index).sort_values(ascending=False)

# create plot
f, ax = plt.subplots(figsize=(12, 8))
plt.xticks(rotation='90')
sns.barplot(x=df_na.index, y=df_na)
ax.set(title='Percentage of missing data by feature',
        ylabel='Percentage missing')
plt.show()
```

Options

aydi • 10 months ago

0

Chirag Sehra • 10 months ago

2

Great Work Rachel! The doubts and comment section shows what a great teacher you are. Thanks for the efforts to make everyone understand such important topics that people mostly skip.

Options

Made in Russia • 10 months ago

2

nengkuan2 • 10 months ago

1

Thanks for tutorial. I'm kind of late bird, but anyways that was a good start. Thats my notebook <https://www.kaggle.com/madeinrussia/data-cleaning-challenge-handling-missing-values/notebook>

Options

Aahan Singh • 10 months ago

^

2

▼

Thank you for the tutorial Rachael. Here's my solution  
<https://www.kaggle.com/aahansingh/data-cleaning-challenge-handling-missing-values>

Options

jaishiva • 10 months ago

^

2

▼

Guys we can get the StreetNumberSuffix >> based on description ; and Street. Below are the examples

In the description we can find the Address information based on which we can get the StreeNumber Suffix

based on below concept

A number suffix is a letter that might come after an address if there aren't enough numbers for all the buildings on a street. (For example, if your address is 9A Main Street, the suffix would be "A".)

below is the link you can have a look at  
[http://www.mississauga.ca/portal/helpfeedback/faq?paf\\_gear\\_id=2000021&itemId=104500575n&action=faqAnswer](http://www.mississauga.ca/portal/helpfeedback/faq?paf_gear_id=2000021&itemId=104500575n&action=faqAnswer)

we can update the data  
StreetnumberSuffix / Street

A >> Minna/Dolores/Pearl/Rausch >> Street space >> for blanks

V >> Hampshire / Hallam / 06th >> Street space >> for blanks

C >> Hayes / Montgomery / Stockton

Options

SHAF • 10 months ago

^

2

▼

In case someone want's to fill the NA's in both direction (forward and backward in chain) can use this code:  
subsetnfldata.fillna(method='ffill',axis=0).fillna(method='bfill',axis=0)

Also, thank you for the tutorial Rachael.

Options

Docty • 10 months ago

^

2

▼

Hello Rachel, I thank you for your tutorial, However, how we can impute missing categorical values in the data???

Options

Rachael Tatman • 10 months ago

^

0

▼

Joseph Romani • 10 months ago

^

2

▼

This is a great tutorial, but I do have to echo Ace Palmero's comments. How you handle missing values in a dataset is entirely dependant on what you are intending to use the data

for and what techniques you are going to apply. Solution is here  
<https://www.kaggle.com/attackgnome/data-cleaning-challenge-handling-missing-values>

Options

Ahmed Alqam • 10 months ago

^ 2 v

Awesome competition, looking forward to continuing the rest of it.

Options

Gabriel Montañola • 10 months ago

^ 2 v

## For those trying to use OpenAddresses dataset

I'm new to Python (2 months) and Pandas (3 days...I guess..) - so mind any newbie coding

San Francisco is not listed as a City. Don't ask me why:

I googled the range of zipcodes and used this (first part is converting the POSTCODE to numeric values, then selecting the range with between the values I want to.)

```
ca_addresses['POSTCODE'] = pd.to_numeric(ca_addresses['POSTCODE'],
errors='coerce')
```

```
ca_addresses['NUMBER'] = pd.to_numeric(ca_addresses['NUMBER'],
errors='coerce')
```

```
new = ca_addresses[ca_addresses['POSTCODE'].between(94102,
94177)].sort_values(by=['POSTCODE'])
```

I found a way to get the Zipcode using the closest number and the street name but I can't get it to work with pandas apply() - mostly because I suck at it.

But it goes like this:

### Defining a function to concat the street name

```
def street_concat(row):
    return row['Street Name'].upper() + " " + row['Street
Suffix'].upper()
```

I know that there is some *AV* vs *AVE* problems, but I can work it later.

### Applying this to the dataframe to create a new column with the concat street name

```
sf_permits['COMBO_ST'] = sf_permits.apply(street_concat, axis=1)
```

### Def a function to get the nearest number

LennartGrosser • 10 months ago

^ 0 v

```
def find_nearest(array, value):
    idx = (np.abs(array-value)).idxmin()
    return array[idx]
```

But then I'm stuck. I can get a ZIPCODE value using this:

```
new['POSTCODE'][(new['STREET'] == 'JERROLD AVE') &
(new['NUMBER'] == find_nearest(new[new['STREET'] == 'JERROLD AVE']
['NUMBER'], 2241))].values[0]
```

This returns me **94124**. But I'm hardcoding the *street name* and *number*. I don't know how to build a function to automate this process. Tried this:

```
def zip_finder(row):
    if row['Zipcode'].isna():
        return new['POSTCODE'][(new['STREET'] == row['COMBO_ST'])
&
        (new['NUMBER'] == find_nearest(new[new['STREET'] ==
row['COMBO_ST']]['NUMBER'],
        row['NUMBER']))].values[0]
    else:
        return row['Zipcode']
```

But it fails in so many ways that I'm ashamed to post it here. Someone willing to improve this?

Options

Killdary Aguiar

• 10 months ago

2

Very cool, thank you Rachael, I learned a lot. So I have a question, in filling in the data, would not it be better to sort the data through the columns related to the addresses before filling with the dataset's own data?

Here is my solution: <https://www.kaggle.com/killdary/data-cleaning-challenge-handling-missing-values>

Options

Rachael Tauman

• 10 months ago

0

Surekha

• 10 months ago

2

Thank you Rachael for this useful tutorial.

I have a question about the fillna() method. If we are replacing all NA's with the value that comes directly after it in the same column, why do we still fill the values with "0"? subsetn/data.fillna(method = 'bfill', axis=0).fillna("0"). Is this for the last row, where there are no values after?

Here is my solution: <https://www.kaggle.com/surekha09/data-cleaning-challenge-handling-missing-values>

Options

Rachael Tauman

• 10 months ago

1

Thomas Zoeller

• 10 months ago

2

Hi @Rachael, funny and interesting to walk trough. Thanks!

Btw.: On the Mac CTRL+ENTER works for executing the code or SHIFT+ENTER to execute and jump to next section.

Options

Rachael Tatman

• 10 months ago

0

Nitin

• 10 months ago

0

mememem • 10 months ago

2

"More practice!" section is a must read, because strategies presented in main notebook are not that advanced.

Options

DEBJIT GHOSH • 10 months ago

2

This was my **first kernel** on Kaggle. Not only did I learn about handling missing data in Python, but I also learnt how to write, edit and publish kernels. Great learning I must say. Thank you Rachael for the tutorial! Looking forward to Day 2 of the challenge. Here's my notebook: <https://www.kaggle.com/debjitghosh/data-cleaning-challenge-handling-missing-values>

Options

rush • 10 months ago

2

Thanks a lot Rachael, nice tutorial!

My fork is [here](#).

But [OpenAddresses dataset](#) doesn't contain data for San Francisco.

Options

Rachael Tatman

• 10 months ago

0

Mihai • 10 months ago

2

I was getting this error when reading the data

```
/opt/conda/lib/python3.6/site-packages/IPython/core/interactiveshell.py:2698: DtypeWarning: Columns (25,51) have mixed types. Specify dtype option on import or set low_memory=False.
```

`interactivity=interactivity, compiler=compiler, result=result)`

I've fixed it by using

```
nfl_data = pd.read_csv("../input/nflplaybyplay2009to2016/NFL Play by Play 2009-2017 (v4).csv", low_memory=False)
sf_permits = pd.read_csv("../input/building-permit-applications-data/Building_Permits.csv", low_memory=False)
```

Options

Tamanna Patil • 10 months ago

Great Learning, Rachael. waiting for the coming sessions!!

Options

Lakshmipath Kakarla • 10 months ago

Really nice and helpful tutorial! Thanks you Rachael!

Options

Faiz Ahmed Khan • 10 months ago

An amazing way to learn glad to share my notebook, looking forward to the challenge and hopefully such competitions won't stop coming in  
<https://www.kaggle.com/faizkhan/data-cleaning-challenge-handling-missing-values>

Options

Gwin • 10 months ago

Excellent job, Rachael ! ... looking forward to the tomorrow's lesson :)

Options

Arunkumar V Ramanan • 10 months ago

Wonderful job, Rachael! Well, Done on Day 1 of the 5-Day Data Challenge! Looking forward to getting started Day 2 :)

Options

ANURAG GUPTA • 10 months ago

Do we have this exercise in R ?

Options

Rachael Tatman • 10 months ago

Options

Kannan Veeramani • 10 months ago

Very good learning. Well explained. Thanks!!!  
<https://www.kaggle.com/kannanv/data-cleaning-challenge-handling-missing-values?scriptVersionId=2929955>

Options

Jojo Xu • 10 months ago

Options

shafiquekhan • 10 months ago

how we can go to R notebook

Options

Completed



how many total missing values do we have?

```
totalcells=np.product(nfldata.shape)
```

dear kagglers what does that line mean ??

Apoorv Patne • 10 months ago

^

0

▼

I've got a doubt in the last section. Here's my solution :  
<https://www.kaggle.com/apoorvwatsky/data-cleaning-challenge-handling-missing-v-f07f3f>  
Can someone please tell me why the missing values in the sf\_permits are not replaced by values directly after them?

Options

Betty Holyk • 10 months ago

^

0

▼

Would you guys give me more comment on this, pl?:

**replace all NA's the value that comes directly after it in the same column,**

**then replace all the reamining na's with 0**

```
subsetnfl_data.fillna(method = 'bfill', axis=0).fillna("0")
```

Options

Rachael Tayman • 10 months ago

^

0

▼

Vaughn Shideler • 10 months ago

^

0

▼

Can someone please explain how the slicing works in the example under "Filling in missing values automatically"? I'm assuming that ["EPA" : "Season"] prints all the columns between "EPA" and "Season," but what does the first colon do? And what if you wanted to print specific rows instead?

Options

Jonas Conde • 10 months ago

^

3

▼

aseemmyada • 10 months ago

^

0

▼

can any one help me where i can write and run the code written above

Options

Chirag Senra • 10 months ago

^

1

▼

Sabu Joseph • 10 months ago

^

0

▼

How can I retrieve the code I saved ? I executed it, versioned, made the code public too

Options

Rachael Tayman • 10 months ago

^

0

▼

GremlinB • 10 months ago

^

0

▼

Is there a shutdown procedure, similar to working locally, when we are finished with a day's



work?

Options

Rachael Tatman

• 10 months ago

^

0

▼

Aman Saxena • 10 months ago

^

0

▼

This was good and easy! Why is it not in R too @Rachael ?

Here's my notebook .. <https://www.kaggle.com/amansaxena/data-cleaning-challenge-handling-missing-v-c12d60>

Options

Rachael Tatman

• 10 months ago

^

1

▼

SeanGalloway • 10 months ago

^

0

▼

How do I download the data listed here?

Options

Rachael Tatman

• 10 months ago

^

0

▼

Rishabhjain • 10 months ago

^

0

▼

Hi Rachel

I have also signed up for this course but i am not getting any mails regarding this

My mail id is [rishabh15virgo@gmail.com](mailto:rishabh15virgo@gmail.com)

Please have a look

Options

Rachael Tatman

• 10 months ago

^

0

▼

russell luttrel • 10 months ago

^

0

▼

Thanks for the tutorial. The last cell seems to have its comment cut off mid sentence though.

Options

Rachael Tatman

• 10 months ago

^

0

▼

mgiy • 10 months ago

^

0

▼

Thank you Rachael!

One question Is there any special reason for filling the remaining NAs' with string "0", not number 0 in the last line: `subsetnfl/data.fillna(method = 'bfill', axis=0).fillna("0")` ?

Options

Rachael Tatman

• 10 months ago

^

0

▼

Kaustubh • 10 months ago

^

0

▼

Thank you for simple and intelligible challenge!  
Can you create a similar notebook where the missing values can be filled using various Regression Analysis techniques? In some specific cases, that would be more useful than filling the values with constant or based on a value in the same column of next row.

Options

[Rachael Tatman](#)

10 months ago

^0v

hklee · 10 months ago

^0v

Thank you for your clear and really helpful instruction.  
How about

`subsetnfldata.fillna(method = 'bfill', axis=0).fillna(method = 'ffill', axis=0)`  
instead of

`subsetnfldata.fillna(method = 'bfill', axis=0).fillna("0")` ?

Options

[Rachael Tatman](#)

10 months ago

^0v

Siddharth · 10 months ago

^0v

Hi rachael...Amazing tutorial...really got to learn a lot... I knew the commands but never got to practice on the data sets. Thanks for taking up this topic.  
Btw I haven't got the link for Data cleaning challenge, Day 2 yet.  
When can I expect to get it? PLEASE LET ME KNOW!

Options

[Rachael Tatman](#)

10 months ago

^0v

Sami · 10 months ago

^0v

Hi Rachael, thank you for this tutorial.  
Have you ever had issues with columns types ? I remember having issues with columns types because when you convert object type to float, it doesn't handle correctly `NaN` and `None`.

Options

[Rachael Tatman](#)

10 months ago

^0v

Preety Missir · 10 months ago

^0v

Thanks Rachel!  
Can you please explain the below :  
`print("Columns in original dataset: %d \n" % sfermits.shape[1]) print("Columns with na's dropped: %d" % columnswithnadropped.shape[1])`

Options

[Rachael Tatman](#)

10 months ago

^0v

Chris · 10 months ago

^

0

▼

Thanks @Rachael ! I have a question by the way, what's the difference between:

```
subset_nfl_data.fillna(1) = 'bfill', axis=0).fillna("0")
subset_nfl_data.fillna(method = 'bfill').fillna("0")
```

Because both give me the same result. Thanks!

Options

LeonnartGrosser · 10 months ago

^

1

▼

Chathura · 10 months ago

^

0

▼

Hi Rachael, I didn't get an email for the second day Challenge?

Options

Rachael  
Tatman · 10 months ago

^

0

▼

dejavuk2 · 10 months ago

Diana Petuhova · 10 months ago

^

0

▼

Hi, Rachael! D  
Why did you use subsetnfl/data.fillna(0), filling missing values with zeros? Is it because the data is already normalised?

Options

Rachael  
Tatman · 10 months ago

^

0

▼

SNEHITHA TIGER · 10 months ago

^

0

▼

Hi@ RachaelTatman ,  
Nice tutorial,, Thanks for helping how to code for beginner's like me.  
I'm new to python. When I'm trying to run this data an error is shown as initializing from file failed. How to sort it.

Options

Rachael  
Tatman · 10 months ago

^

0

▼

Vivi · 10 months ago

^

0

▼

Hey Rachael, thanks for the tutorial! :D I have one question regarding missing data during training. Suppose there is one real-valued feature that might not exist for every sample (=value missing because it doesn't exist). How do I then train a model if I really want to include this feature(a classical regression model cannot deal with "NaN"s, while trees could perhaps still deal with it)? What methods do you usually choose?

Options

Rachael  
Tatman · 10 months ago

^

1

▼

boosting\_75

10 months ago

^

0

v

Easy and fast way to learn Python;-) Here my code for review:  
<https://www.kaggle.com/zeus75/data-cleaning-challenge-handling-missing-values>

Options

Tyler James Martin

10 months ago

^

0

v

Another common strategy is to replace with mean.  
You can use:  
from sklearn.preprocessing import Imputer  
imputer = Imputer(missing\_values = 'NaN', strategy = 'mean', axis = 0)  
  
And then fit & transform your data.

Options

Abdoulaye Diou

10 months ago

^

0

v

SNEHITHA  
TIGER

10 months ago

^

2

v

Hello, I'm so late for this chalenge. thnk you Rachel. But could someone explain me what  
.loc[] refer to? and Axis=0?  
I'm beginner in python and im still learning basics.  
  
Thanks

Options

jaishiva

10 months ago

^

0

v

Hey Rachel is there any way to find out the Street Number Suffix ?

Options

David Nichols

10 months ago

^

0

v

For imputing missing values, NA's can be imputed by first determining the likely distribution  
of the variable to be imputed, then sampling from its distribution and randomly imputing  
each NA.  
<https://epic2020datascience.blogspot.com/2018/01/>

Options

Thi An

10 months ago

^

0

v

Hi all, I'm new to this and I'd like to know if I can Fork this notebook and make it an R  
network. Thank you very much in advanced.

Options

Rachael  
Tatman

• 10 months ago

^

0

▼

blackspider007

• 7 months ago

^

0

▼

Thanks Rachael for the tutorial! But when I am running the codes in the console here, first the `nf1data` *did not get defined*. I started working with `sfpermit` but this is even not running. I'm getting an error message "`sfpermit is not defined`". *This is the same error I was getting for `nf1data`*. Will you please help me?

Options

Rachael  
Tatman

• 7 months ago

^

0

▼

John Olayisade

• 7 months ago

^

0

▼

Thank you Rachael. Will this work for a ,R' environment?

Options

Rachael  
Tatman

• 7 months ago

^

0

▼

Srishti Jaiswal

• 7 months ago

^

0

▼

Thanks a lot for this simple and helpful tutorial. I have a question that when I ran my last code:-

"sf\_permits.fillna(method='bfill',axis=0).fillna(0)"

In the output I got NaN values which was not supposed to come. Can you please help clearing my doubt?

Options

Imran  
Islahdar

• 7 months ago

^

0

▼

Prabesh123

• 6 months ago

^

0

▼

Hi Rachael I am trying my hands on R and was wondering if you run a similar daily challenge with R as well ? I could not find it anywhere. Many thanks :)

Options

Rachael  
Tatman

• 6 months ago

^

0

▼

Bhumika Lamba

• 6 months ago

^

0

▼

Hi,

for the `fillna(method='bfill',axis=0).fillna(0)` . Why do we need to set `axis=0`? I am not able to understand how this method works. If it befill method then there won't be any NaN value left to be filled as a zero.

Options

Rachael Tatman

• 6 months ago

^

1

▼

aljqa2

• 6 months ago

^

0

▼

Hi Rachel,

Thanks - this lesson is awesome! I do have a question - I don't understand why the following head() command doesn't return any rows:

```
rows_with_na_dropped = sf_permits.dropna(axis=0,how = 'any')  
  
rows_with_na_dropped.head()
```

My notebook is located @ <https://www.kaggle.com/aljqa2/data-cleaning-challenge-handling-missing-values/edit>

Thank you!

Options

Joshua

• 6 months ago

^

1

▼

Abdoulaye

• 9 months ago

^

0

▼

Yo Ram

• 6 months ago

^

0

▼

There are a lot of typo's in the text, but really interesting challenge!

Options

Aniket Gulhane

• 5 months ago

^

0

▼

Hi,

I am new to data science, undergoing a PG diploma.

I have a doubt while treating the NA values. The data I have consist of 15 variable, of which for 2 variables, there are around 988 NA observations out of the 1232 total observations. I am not sure how I should deal with this situation.

Any help is much appreciated!!

Thanks!

Options

Alex Poulain

• 5 months ago

^

0

▼

Thank you! Very nice cleaning tutorial.

Options

Aditya Agrawal

• 5 months ago

^

0

▼

Hi Rachael and everyone here! Thanks a lot for this lesson. I took away some vital points from this notebook- one of the major ones being that every column would need to be dealt with individually! I have just published [my first kernel](#) with some pre-processing and EDA on a



	<div>clothing-fitting feedback dataset. I would really appreciate it, if you guys could check it out and give me some feedback on it!</div> <div>Thanks!</div> <div>Options</div>	
	<div><b>Emad Tolba</b> • 4 months ago</div> <div>thanks</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>Oumaima Hourrane</b> • 4 months ago</div> <div>Very interesting and helpful, thank you!</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>Ewelina Szymakowicz-Grabania</b> • 4 months ago</div> <div>Thank you Rachel, it is very useful challenge. I am learning a lot.</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>Farshad</b> • 4 months ago</div> <div>The course is highly appreciated. Because data cleaning is always one of my headaches!</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>Rachael Tauman</b> • 10 months ago</div> <div><b>Micheal_Mike</b> • 3 months ago</div> <div>Very helpful</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>selin</b> • 3 months ago</div> <div>Very helpful Rachael thanks for sharing!</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>ganesh</b> • 16 days ago</div> <div>how do you handle missing value in categorical data set ?</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>
	<div><b>jors</b> • 8 days ago</div> <div>Thanks for this Rachael!</div> <div>Options</div>	<div><div>^</div><div>0</div><div>v</div></div>

10 months ago <div>This Comment was deleted.</div> <div>Options</div>		
10 months ago <div>This Comment was deleted.</div>		
10 months ago <div>This Comment was deleted.</div>		
10 months ago <div>This Comment was deleted.</div> <div>Options</div>		
10 months ago <div>This Comment was deleted.</div> <div>Options</div>		
10 months ago <div>This Comment was deleted.</div> <div>Options</div>		
10 months ago <div>This Comment was deleted.</div>		
10 months ago <div>This Comment was deleted.</div>		
7 months ago <div>This Comment was deleted.</div>		
7 months ago <div>This Comment was deleted.</div> <div>Options</div>		
7 months ago <div>This Comment was deleted.</div>		




7 months ago

This Comment was deleted.


7 months ago

This Comment was deleted.


Similar Kernels




SQL Scavenger Hunt Handbook




Data Cleaning Challenge: Scale And Normalize Data



SQL Scavenger Hunt: Day 1



Data Cleaning Challenge: Parsing Dates






Data Cleaning Challenge: Inconsistent Data Entry

© 2019 Kaggle Inc

Loading [Contrib]/a11y/accessibility-menu.js

Team Terms Privacy Contact/Support



ju

10 months ago

^

1

▼

```
elmitsd = False, StreetNumber suffix, name_upper + " " + surname
```

row [street suffix] . . .  
super ( )

b b  
 K E  
 e n  
 n p  
 p n  
 n i  
 g h  
 h t  
 c a  
 a w  
 o m  
 t a  
 h h  
 a p  
 e a  
 p A  
 p p  
 o n  
 s s  
 e n  
 c i  
 o n  
 t r  
 i b  
 u t  
 e s  
 o o  
 t t  
 h h  
 e m  
 a a  
 n n  
 a o  
 n s  
 e e  
 n n  
 e e  
 f f  
 a a  
 r r  
 e e  
 s s  
 e e  
 m m  
 o o  
 n n  
 e e  
 y y

get a new version on the market. Although it is a bid to be the best of the best, it is a challenge to the established players in the market. The new version is a challenge to the established players in the market.

Rachael  
Tatman

• 6 months ago

^

0

▼

Rachael  
Tatman

• 10 months ago

^

0

▼

Options

alpha2

• 6 months ago

^

0

▼

g  
a  
t  
h  
h  
a  
n  
h  
a  
e  
a  
e  
e  
a  
n  
t  
n  
o  
o  
o  
b  
b  
e  
a  
g  
m  
e  
a  
h  
h  
a  
n  
t  
e  
a  
a  
n  
e  
e  
a  
b  
a  
n  
r  
e  
e  
s  
s  
t  
i  
n  
g  
r  
a  
y  
(  
l  
h  
a  
s  
a  
d  
o  
c  
y  
e  
m  
s  
e  
t  
h  
t  
e  
s  
h

jaishiva • 10 months ago

^ 2 v

an  
ap  
y.  
eab  
tbs  
hs  
e(  
ear  
e  
o  
o  
n  
g  
w  
h  
i  
s  
h  
a  
h  
e  
i  
n  
g  
x  
h  
m  
i  
n  
e  
t  
o  
s  
a  
e  
p  
e  
r  
n  
b  
a  
w  
e  
d  
s  
d  
a  
g  
i  
x  
o  
t  
b  
B  
B  
t  
t  
p  
p  
e  
n  
e  
a



stakeholders and the board of directors.

new website [P.O.S.T.C.O.D.E.]

nengkuan2 • 10 months ago

the following estimates of the number of people who have been exposed to the virus since the start of the outbreak are based on the number of people who have been exposed to the virus since the start of the outbreak.

W r  
e e  
w s  
a t  
n (  
t n  
h e  
a w  
b [  
m n  
t e  
g w  
h [  
h .  
h S  
c T  
e R  
m E  
p E  
a T  
w .  
h ]  
e =  
n =  
a =  
h .  
a J  
d E  
d R  
i R  
g O  
b L  
s D  
y  
t A  
b V  
h E  
e .  
e ]  
e [  
a .  
t N  
e U  
d M  
w B  
t E  
e R  
h '  
b ]  
e ,  
g 2  
g 2  
h 4  
t 1  
m )  
b )  
g ]  
u .  
e v

sales  
[0]

a  
 a  
 h  
 b  
 a  
 m  
 e  
 e  
 t  
 s  
 v  
 d  
 b  
 a  
 'a  
 b  
 k  
 f  
 n  
 o  
 w  
 b  
 a  
 v  
 d  
 b  
 b  
 h  
 a  
 i  
 d  
 a  
 b  
 o  
 t  
 n  
 o  
 n  
 i  
 n  
 s  
 t  
 i  
 t  
 u  
 t  
 i  
 o  
 n  
 a  
 m  
 i  
 t  
 n  
 o  
 k  
 m  
 y  
 a  
 c  
 t  
 u  
 a  
 l  
 m  
 o  
 n  
 i  
 t  
 o  
 r  
 i  
 n  
 g  
 s  
 a  
 v  
 e  
 o  
 v  
 e  
 r  
 e  
 l  
 e

o  
b  
k  
a  
e  
h  
b  
d  
p  
h  
h  
e  
w  
w  
a  
r  
en  
o  
s  
s  
y  
a  
s  
a  
w  
g  
a  
t  
b  
a  
d  
p  
a  
v  
a  
a  
l  
a  
b  
e  
e  
p  
b  
e  
d  
b  
a  
c  
k  
y  
f  
a  
q  
p  
p  
a  
a

d  
e  
f  
z  
i  
p  
-  
f  
i  
n  
d  
e  
r  
(  
r  
o  
w  
)  
:  
i  
f  
r  
o  
w  
[  
'  
Z  
i  
p  
c  
o  
d  
e  
'  
]  
.  
i  
s  
n  
a  
(  
)

n :  
g  
A  
a  
n  
e  
i r  
d e  
1 t  
1 W  
2 or  
0 un  
0 l  
0 d  
0 n  
0 N  
2 ow  
1 l  
& n  
i e  
t p  
e c  
t o  
e a  
m T  
l S  
d e  
= D  
1 E  
0 A  
4 N  
5 "[  
0 v(  
0 an  
0 le  
5 w  
7 el  
5 i  
n n  
& t  
a h  
c E  
t f  
i u  
o t  
n u  
= r  
f e  
a o  
q p  
A e  
n rw  
s a  
w t  
e i  
r o  
t n  
h B  
a l  
t f  
i s  
i T  
o,

s,]  
t N  
ho  
en  
rea  
eim  
asp  
sp;  
or  
no  
wb(  
ha<sup>n</sup>  
y<sup>e</sup>  
wl<sup>w</sup>  
[  
ey,  
ha<sup>N</sup>  
al<sup>U</sup>  
v<sup>M</sup>  
e<sup>B</sup>  
S<sup>Æ</sup>  
u<sup>R</sup>  
ft'  
fe]  
im  
xp<sup>̄</sup>  
np<sup>̄</sup>  
ul  
ma<sup>f</sup>  
bc<sup>i</sup>  
ee<sup>n</sup>  
rh<sup>d</sup>  
so<sup>-</sup>  
fl<sup>n</sup>  
od<sup>e</sup>  
re<sup>a</sup>  
sr<sup>r</sup>  
se  
o.s  
mT<sub>t</sub>  
eh(  
rin  
ese  
ciw  
os[  
rm  
dœ  
stw  
af[  
ni'  
dn<sup>S</sup>  
fa<sup>T</sup>  
ol<sup>R</sup>  
r.<sup>E</sup>  
s<sup>E</sup>  
W<sup>T</sup>  
oe,  
ms]  
et'  
ri=  
el=



c l  
o rr  
r eo  
d ew  
s dl  
t t'  
h oC  
e cO  
d lM  
a eB  
t aO  
a nS  
i N  
s o,  
m nJ  
i eJ  
s bJ  
s e'  
i fN  
n uJ  
g rM  
. eB  
t E  
A rR  
c a'  
t iJ  
u n'  
a i  
l n'  
l gO  
y w  
a J  
m i'  
t n  
h gU  
i hM  
n tB  
k ?E  
2 R  
n l'  
g tJ  
i nJ  
f i)  
t gJ  
h h.  
e tV  
r bA  
e e1  
m aU  
i gE  
g oS  
h oJ  
t dJ  
b i  
e de  
a e1  
p as  
a te  
o:

t  
t  
e  
r  
n  
b  
a  
s  
e  
v  
d  
o  
n  
a  
w  
h  
i  
c  
h  
t  
h  
i  
s  
(  
d  
a  
t  
n'  
a  
h  
a  
e  
s  
b  
i  
e  
e  
p  
l  
p  
)  
a  
p  
r  
y  
m  
i  
s  
p  
e  
t  
g  
p  
o  
o  
m  
p  
a  
u  
n  
l  
y  
a  
w  
r  
a  
v  
y  
a  
s  
l  
t  
u  
h  
e  
(  
a  
t  
i  
f  
t  
m  
e  
a  
x  
t

Options

s )  
hf  
a o  
mr  
em  
di  
ts  
os  
pi  
on  
sg  
tv  
ia  
tl  
hu  
ee  
rc  
ee  
. l  
l

Options

. l  
**Jonas Conde**

• 10 months ago

^	1	v
---	---	---

t  
Hh  
i e  
Gr  
ae  
ba  
rh  
ia  
en  
ld  
, y  
n w  
na  
iy  
ct  
eo  
jd  
oo  
bt  
th  
ra  
yt  
i ?  
n

Options

o  
**Rachael Tatman**

• 10 months ago

^	0	v
---	---	---

i  
q  
u  
l  
e  
th  
b

u  
s  
d  
a  
h  
e  
c  
k  
o  
u  
e  
d  
h  
e  
l  
a  
g  
g  
o  
o  
g  
l  
e  
m  
a  
p  
e  
A  
P  
t  
h  
a  
b  
w  
e  
v  
e  
m  
p  
f  
t  
g  
a  
n  
e  
m  
e  
a  
b  
o  
k  
o  
f  
e  
a  
d  
d  
r

es  
sw  
b  
e  
a  
a  
t  
o  
t  
e  
d  
g  
e  
t  
s  
q  
m  
e  
a  
d  
d  
o  
e  
s  
s  
e  
s  
(  
A  
+  
l  
h  
s  
t  
y  
a  
g  
a  
s  
h  
t  
o  
m

Options

o  
r  
r  
o  
w  
w  
i  
t  
h  
m  
o  
r  
e  
p  
a

t  
i  
e  
n  
c  
e  
.

Options