



Python for Data Science

Dr. Charlie H. Apigian
Interim Director - Data Science Institute
Middle Tennessee State University
April 6, 2019


<https://github.com/capigian/charleston>

I AM *true*
BLUE

`mt_dsi['using data for good'].max()`



Overview of Presentation

- Focus for today
 - Learn a little about , the Data Science Institute and what we are doing at MTSU.
 - Help debunk misperceptions of how and who can learn machine learning.
 - Understand the generalities of ML and how it can be applied
 - Establish a process for building machine learning solution that actually solves a problem.
 - Work through a problem showing one process

mtsu.edu/dsi

<https://github.com/capigian/charleston>

What a Data Scientist?

MODERN DATA SCIENTIST

Data Scientist, the second job of 21st century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand what a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- Machine learning
- Statistical modeling
- Experiment design
- Bayesian inference
- Supervised learning: decision trees, random forests, logistic regression
- Unsupervised learning: clustering, dimensionality reduction
- Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- Computer science fundamentals
- Scripting languages e.g. Python
- Statistical computing package e.g. R
- Databases SQL and NoSQL
- Relational algebra
- Parallel databases and parallel query processing
- MapReduce concepts
- Hadoop and Hadoop Pig
- Custom reducers
- Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- Passionate about the business
- Curious about data
- Influence without authority
- Hacker mindset
- Problem solver
- Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- Able to engage with senior management
- Story telling skills
- Translate data driven insights into decisions and actions
- Visual art design
- R packages like ggplot2 or lattice
- Knowledge of any of visualization

DATA SCIENCE AS A TEAM SPORT



ANACONDA

© 2017 Continuum Analytics, Inc. All rights reserved.

4

Start to think of Data Science as a team sport.

Anaconda - Continuum Analytics

mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

What skills do you need?

Data Scientist
also known as Data Managers, statisticians.

A data scientist will be able to take data science projects from end to end. They can help store large amounts of data, create predictive modelling processes and present the findings.

Skills: Mathematics, Programming, Communication

Will use programmes such as: SQL, Python, R

Data Engineers
also known as database administrators and data architects.

They are versatile generalists who use computer science to help process large datasets. They typically focus on coding, cleaning up data sets, and implementing requests that come from data scientists.

Skills: Programming, Mathematics, Big data

Will use programmes such as: Hadoop, NoSQL, and Python

Data Analysts
also known as business Analysts.

They typically help people from across the company understand specific queries with charts.

Skills: Statistics, Communication, Business knowledge

Will use programmes such as: Excel, Tableau, SQL

Data Scientists

- Knowledge of database like MySQL.
- Knowledge of Java, Python, R Programming.
- Knowledge of various analytical functions.
- Expertise in Math, Statistics, correlation, data mining, and predictive analysis.
- Knowledge of Machine Learning, Clustering

Data Analysts

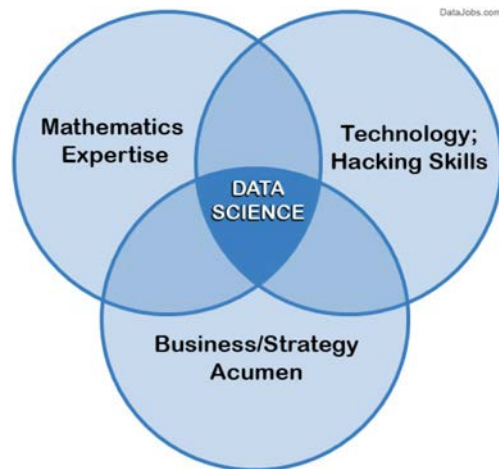
- Hands-on experience with data warehousing and business intelligence concepts.
- In-depth exposure to SQL and analytics.
- Data Storing and retrieving skills and tools.
- Proficient in data architecture.
- Proficiency in decision making
- Familiar with various ETL tools

Which skills should you learn first?

Depends on your objective and your interest.

Regardless of interest, you must have an education and skills.

Skills alone is not good enough.

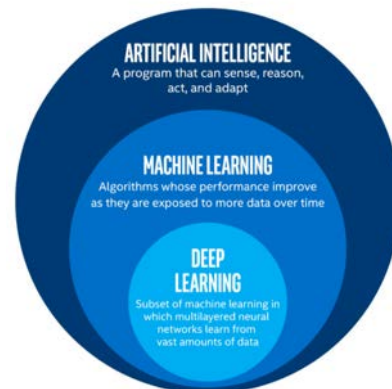


mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

What is machine Learning?

- **Machine learning** is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed.

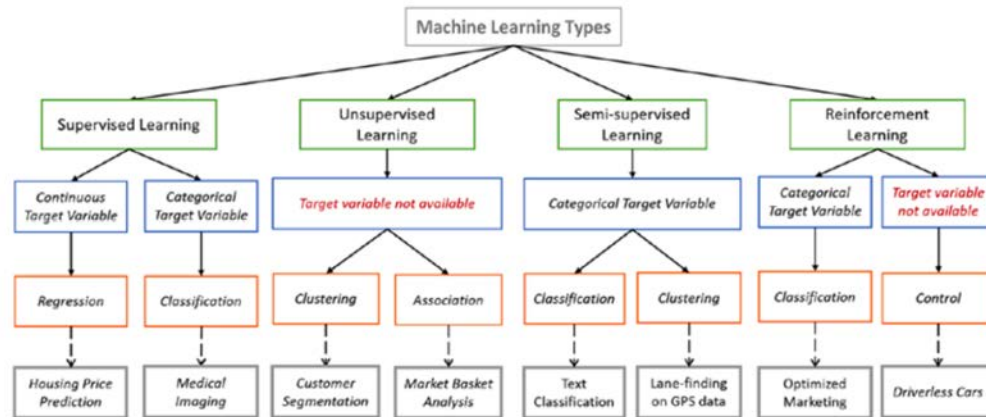


Machine learning focuses on the development of computer programs that can access data and use it learn for themselves.

mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

Different Types of Machine Learning

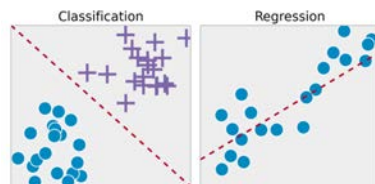


towardsdatascience.com

mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

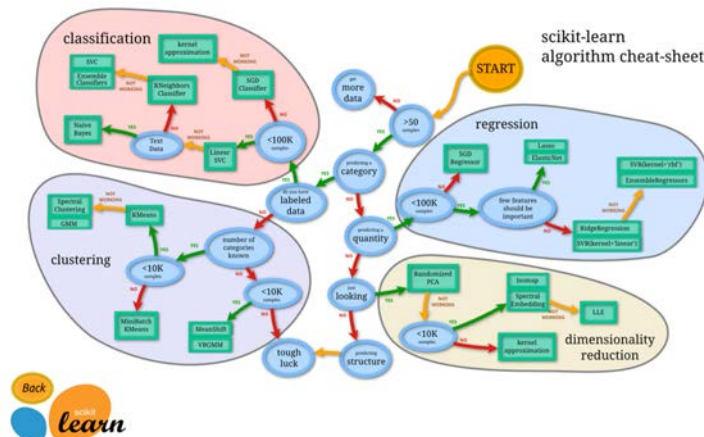
Machine Learning Algorithms Mind-Map



mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

Scikit-Learn for Machine Learning Algorithm Selection



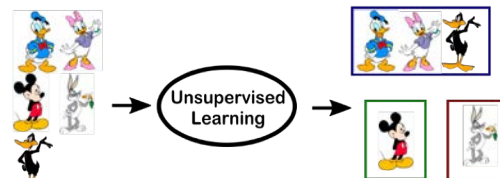
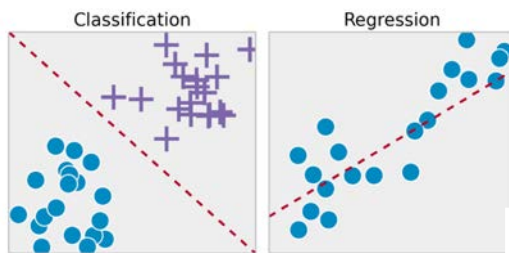
mtsu.edu/dsi

```
mt_dsi['using data for good'].max()
```

Supervised vs. Unsupervised

Supervised

Unsupervised



<https://www.kdnuggets.com/2018/04/supervised-vs-unsupervised-learning.html>

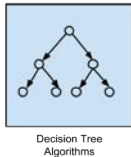
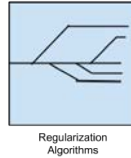
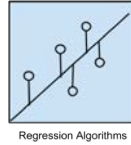
	Supervised Learning	Unsupervised Learning
Discrete	classification or categorization	clustering
Continuous	regression	dimensionality reduction

mtsu.edu/dsi

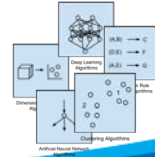
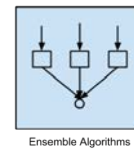
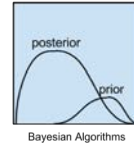
```
mt_dsi['using data for good'].max()
```


Types of Algorithms

- Regression Algorithms
 - Ordinary Least Squares Regression (OLSR)
 - Linear Regression
 - Logistic Regression
 - Stepwise Regression
- Regularization Algorithms
 - Ridge Regression
 - Least Absolute Shrinkage and Selection Operator (LASSO)
 - Elastic Net
- Decision Tree Algorithms
 - Classification and Regression Tree (CART)
 - Iterative Dichotomiser 3 (ID3)
 - Chi-squared Automatic Interaction Detection (CHAID)



- Bayesian Algorithms
 - Naive Bayes
 - Gaussian Naive Bayes
 - Multinomial Naive Bayes
- Ensemble Algorithms
 - Boosting
 - Bootstrapped Aggregation (Bagging)
 - AdaBoost
 - Gradient Boosting Machines (GBM)
 - Random Forest
- Clustering Algorithms
- Association Rule Learning Algorithms
- Dimensionality Reduction Algorithms
- Deep Learning Algorithms
- Artificial Neural Network Algorithms
- ...and many others



mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

When do you build a model?

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.

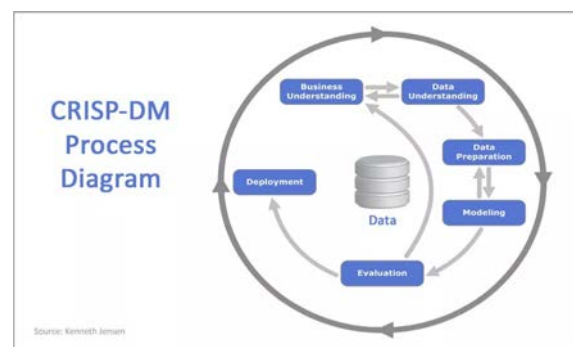
Modeling

Model Building Process

- | | |
|--|---------------------------|
| 1. Select a model | 1. Business Understanding |
| 2. Identify and select the data that fits that model | 2. Data Understanding |
| 3. Transform the data | 3. Data Preparation |
| 4. Identify a business problem | 4. Modeling |
| 5. Train/Test Split | 5. Evaluation |
| 6. Run Model | 6. Deployment |

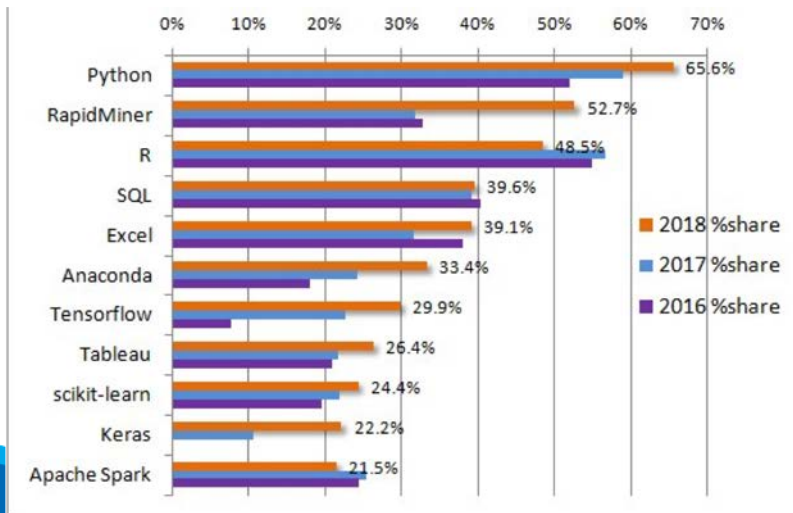
Data Analysis Process

1. Business Understanding
 - a. Frame the problem and the REAL pain point
 - b. Available resources, problems, goals
2. Data Understanding
 - a. What data do you have available to you?
 - b. Setup your workspace with tools or applications
 - Programming – Jupyter notebooks for Python or R Studio for R
 - BI/spreadsheets – Excel – PowerPivot - Tableau
 - c. Import or download the data
 - d. View, explore, and summarize the data
3. Data Preparation
 - a. Clean up null values, outliers, mistakes
 - b. Construct new data, transform or feature engineering
 - c. Integrate and merge data
 - d. Format data (strings, integers, floats, etc.)
 - e. Create you X and y
4. Modeling
 - a. Split your data (Train/Test Split)
 - b. Setup models for machine learning/AI processes
 - c. Can include visuals, dashboards or reports



5. Evaluation
 - a. Fine tune your model
 - b. Create a report of the findings
6. Deployment of models

Why Python?



KDnuggets Analytics/Data Science 2018 Software Poll:
Top tools in 2018, and their share in the 2016-7 polls

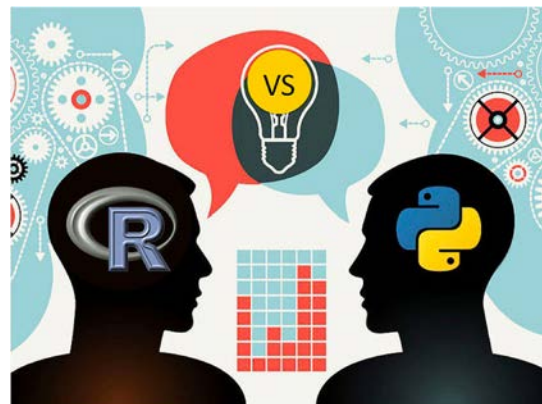
<https://www.kdnuggets.com/2016/06/r-python-top-analytics-data-mining-data-science-software.html>

mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

What should you use?

- Whatever you and your community want to learn and use.



What will we use?

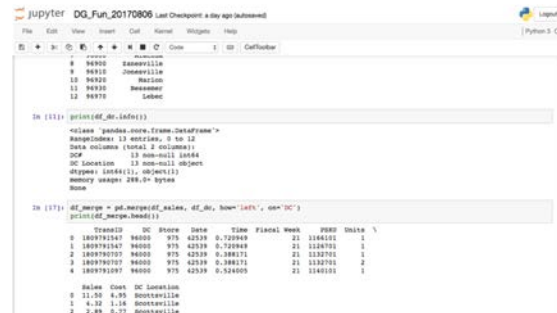
Anaconda

<https://anaconda.org/>



Jupyter Notebooks

<http://jupyter.org/>



mtsu.edu/dsi

`mt_dsi['using data for good'].max()`











Use of Python Libraries




http://chris35wills.github.io/courses/pydata_stack/

mtsu.edu/dsi

`mt_dsi['using data for good'].max()`

Library	Category	2017	2018	2019	2020	2021	2022	2023	2024	2025
 Orange3	Machine learning	22 753	1 084	86	2 114	28 098	14 005	21	265	26
XGBoost	Machine learning	3277	280	9	868	11 991	5 425	12	364	43
LightGBM	Machine learning	1083	79	14	363	5 488	1 467	14	77	69
CatBoost	Machine learning	1509	61	20	157	2 780	369	25	75	46
 eli5	Machine learning	922	6	22	39	672	89	154	42	112
 SciPy	Data wrangling	19 150	608	99	301	4 447	2 318	31	193	7
 NumPy	Data wrangling	17 911	641	136	390	7 215	2 766	28	132	11
 pandas	Data wrangling	17 144	1 165	93	858	14 294	5 788	15	184	12
 StatsModels	Statistics	10 067	153	21	234	2 868	1 240	66	479	19
 TensorFlow	Deep learning	33 339	1 469	58	7 968	99 664	62 952	23	575	68
 PyTorch	Deep learning	11 306	635	16	816	15 512	3 483	18	707	24
 Keras	Deep learning	4 539	671	41	1 673	29 444	10 964	7	1111	44
dist-keras	Distributed deep learning	1125	5	7	41	431	106	225	161	86
elephas	Distributed deep learning	170	13	5	97	913	189	13	34	70
spark-deep-learning	Distributed deep learning	67	11	3	116	920	206	6	22	84
 Natural Language Toolkit	NLP	13 041	736	74	467					

<https://activewizards.com/blog/top-20-python-libraries-for-data-science-in-2018/>



What can you use?

Technique	
Data Import	Pandas, Regular Expressions
Visualization	Matplotlib, Seaborn
Cleansing Data	Pandas, Numpy
Transforming Data	Pandas or ScikitLearn
Basic Statistics	Pandas, ScikitLearn, StatsModels
Unsupervised Learning	ScikitLearn
Supervised Learning - Regression	StatsModels or ScikitLearn
Supervised Learning - Classification	ScikitLearn
Deep Learning	Keras or TensorFlow

mtsu.edu/dsi
`mt_dsi['using data for good'].max()`

Python example

Let's get started with Python and
Jupyter notebooks



mtsu.edu/dsi

```
mt_dsi['using data for good'].max()
```

Appleton Lending Co

- **Operations**
 - Over 80% of the loans provided by Appleton are personal. These loans are mostly made by borrowers in order to consolidate debt or pay off credit cards, but they may be provided for numerous reasons such as weddings, vacations, and for small businesses.
- **Strategy**
 - Over the past two years, Appleton has provided over 3 billion dollars in loans. The company provides personal loans for amounts between \$1,000 and \$40,000 that can be repaid over time periods of 3 or 5 years. Appleton approves loans based on credit history, credit score, debt to income ratio (dti), and the amount of the loan applied for. Appleton is highly selective with the loans it accepts, with over an 80% denial rate over the past four years. This ensures that Appleton provides high quality opportunities for itself and for lenders.

Appleton Lending Co

- After some negative publicity at the board level, Appleton is looking to refocus its efforts on providing high quality loans. They are wanting to better understand their customers and most importantly, the difference between good loans and bad loans.
- After understanding the type of customers that they serve, they would like to improve the company's ability to predict borrowers who will default on loans. Additionally, Appleton is interested in predicting how much a borrower would be able to pay back, regardless of how large of a loan they have applied for.

Business Understanding

1. **Business Understanding**
 - Available resources, problems, goals
2. Data Understanding
3. Data Preparation
4. Modeling
5. Evaluation
6. Deployment of models

- What are the available resources?
 - What are the key performance indicators (variables)?
- What are Appleton's expressed problems?
- What are Appleton's expressed and underlying goals?

• What data do you have available?

Why is this data too much?

Feature	Description	Feature	Description
member_id	A unique Appleton assigned id for the borrower member.	dti	A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested Appleton loan, divided by the borrower's self-reported monthly income.
loan_amnt	The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.	delinq_2yrs	The number of 30+ days past due incidences of delinquency in the borrower's credit file for the past 2 years.
funded_amnt	The total amount committed to that loan at that point in time.	earliest_cr_line	The month the borrower's earliest reported credit line was opened.
funded_amnt_inv	The total amount committed by investors for that loan at that point in time.	inq_last_6mths	The number of inquiries in past 6 months (excluding auto and mortgage inquiries).
term	The number of payments on the loan. Values are in months and can be either 36 or 60.	mths_since_last_delinq	The number of months since the borrower's last delinquency.
int_rate	Interest Rate on the loan.	mths_since_last_record	The number of months since the last public record.
installment	The monthly payment owed by the borrower if the loan originates.	open_acc	The number of open credit lines in the borrower's credit file.
grade	Appleton assigned loan grade: A, B, C, D, etc. with A being the best.	pub_rec	Number of derogatory public records.
sub_grade	Appleton assigned loan subgrade: A1, A2, A3, etc. with A1 being the best.	revol_bal	Total credit revolving balance.
emp_title	The job title supplied by the Borrower when applying for the loan.	revol_util	Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
emp_length	Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.	total_acc	The total number of credit lines currently in the borrower's credit file.
home_ownership	The home ownership status provided by the borrower during registration. Our values are RENT, OWN, MORTGAGE, OTHER.	initial_list_status	The initial listing status of the loan. Possible values are -W, F.
annual_inc	The self-reported annual income provided by the borrower during registration.	out_prncp	Remaining outstanding principal for total amount funded.
verification_status	Indicates if income was verified by Appleton, not verified, or if the income source was verified.	out_prncp_inv	Remaining outstanding principal for portion of total amount funded by investors.
issue_d	The month which the loan was funded.	total_pymnt	Payments received to date for total amount funded.
loan_status	Current status of the loan.	total_pymnt_inv	Payments received to date for portion of total amount funded by investors.
pymnt_plan	Indicates if a payment plan has been put in place for the loan.	total_rec_prncp	Principal received to date.
desc	Loan description provided by the borrower.	total_rec_int	Interest received to date.
purpose	A category provided by the borrower for the loan request.	total_rec_late_fee	Late fees received to date.
title	The loan title provided by the borrower.	recoveries	post charge off gross recovery.
zip_code	The first 3 numbers of the zip code provided by the borrower in the loan application.	collection_recovery_fee	post charge off collection fee.
addr_state	The state provided by the borrower in the loan application.	last_pymnt_d	Last month payment was received.
		last_pymnt_amnt	Last total payment amount received.
		next_pymnt_d	Next scheduled payment date.
		last_credit_pull_d	The most recent month Appleton pulled credit for this loan.
		collections_12_mths_ex_med	Number of collections in 12 months excluding medical collections.
		mths_since_last_major_derog	Months since most recent 90-day or worse rating.
		policy_code	publicly available policy_code=1 new products not publicly available policy_code=2
		application_type	Indicates whether the loan is an individual application or a joint application with two co-borrowers.
		acc_now_delinq	The number of accounts on which the borrower is now delinquent.
		tot_coll_amt	Total collection amounts ever owed.
		tot_cur_bal	Total current balance of all accounts.
		total_credit_rv	Total revolving high credit/credit limit.
		loan_status	True if borrower defaulted on loan. False if loan was good.

• What is the RIGHT available data?

```
df_loandata = pd.read_csv('data/Loan_Data.csv', index_col = 0, header = 0)
df_loandata.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 49870 entries, 149512 to 4076727
Data columns (total 20 columns):
loan_amnt          49870 non-null int64
term              49870 non-null int64
sub_grade         49870 non-null object
emp_length        49870 non-null int64
home_ownership    49870 non-null object
annual_inc        49870 non-null float64
purpose           49870 non-null object
delinq_2yrs       49870 non-null int64
mths_since_last_delinq  21790 non-null float64
open_acc          49870 non-null int64
pub_rec           49870 non-null int64
revol_bal         49870 non-null int64
total_acc         49865 non-null float64
collections_12_mths_ex_med  49870 non-null int64
mths_since_last_major_derog  49870 non-null int64
acc_now_delinq    49870 non-null int64
tot_coll_amt      49870 non-null int64
tot_cur_debt      49870 non-null int64
total_credit_rv    49870 non-null int64
loan_status       49870 non-null object
dtypes: float64(3), int64(13), object(4)
memory usage: 8.0+ MB
```

Contact Information

- Charlie Apigian, PhD.
- Interim Director - Data Science Institute
- Middle Tennessee State University
- charles.apigian@mtsu.edu
- @capigian
- www.mtsu.edu/dsi
- www.mtsu.edu/isa

