# Predicting Property Value

Sakava Kiv | Alex Thibeaux

William Jones | Chris Mathew

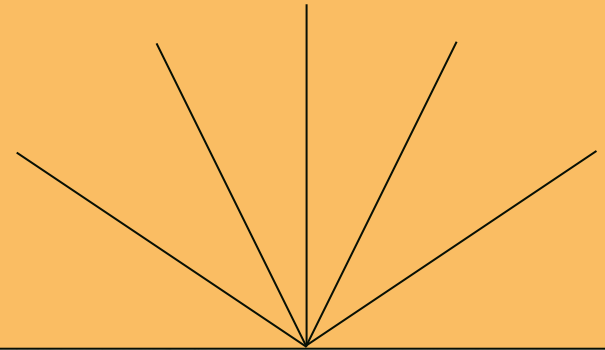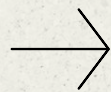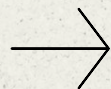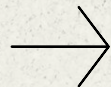# TABLE OF CONTENTS

# Purpose & Background

# 01.

# Purpose

To gain experience applying data science techniques on real-world datasets and implementing a database
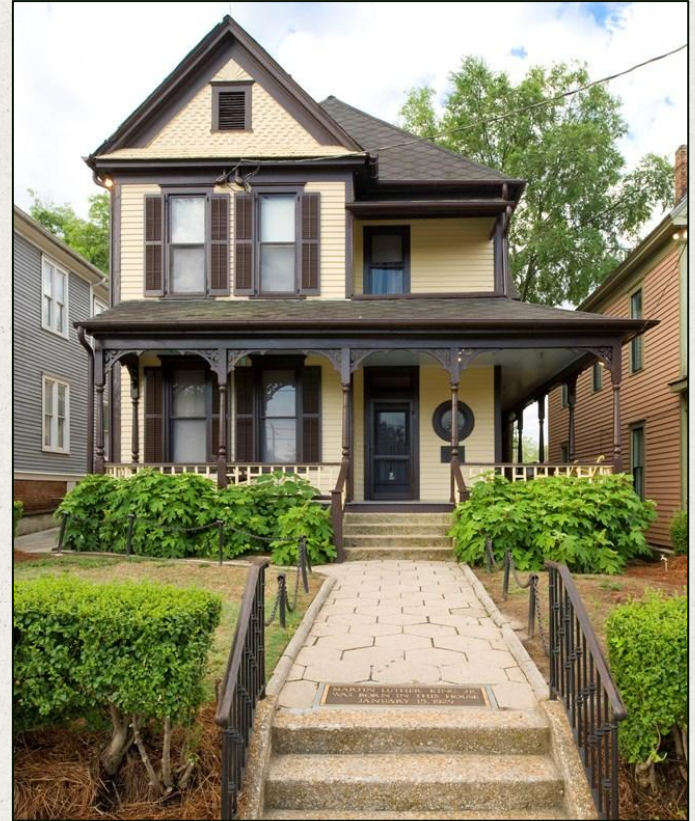
Gain insights in the real estate market

Predicting the property value

# What is property value?

- Probable price of a given property at a given time

- It is important for real estate deals and assessing the property tax

- There are many factors that affect the property value i.e. location, size, condition of building

# Available Data

## MONETARY

taxamount
landtaxvaluedollarcnt
taxvaluedollarcnt
structuretaxvalue-
dollarcount
saleprice

## PROPERTY

bathroomcnt
bathroomcnt
calculatedbathnbr
calculatedfinshed-
squarefeet
finishedsquarefeet12
fips

lotsizesquarefeet
roomcnt
yearbuilt
assessmentyear
latitude
longitude
fullbathroomcnt

## Identification

parcelid
propertycountyland-
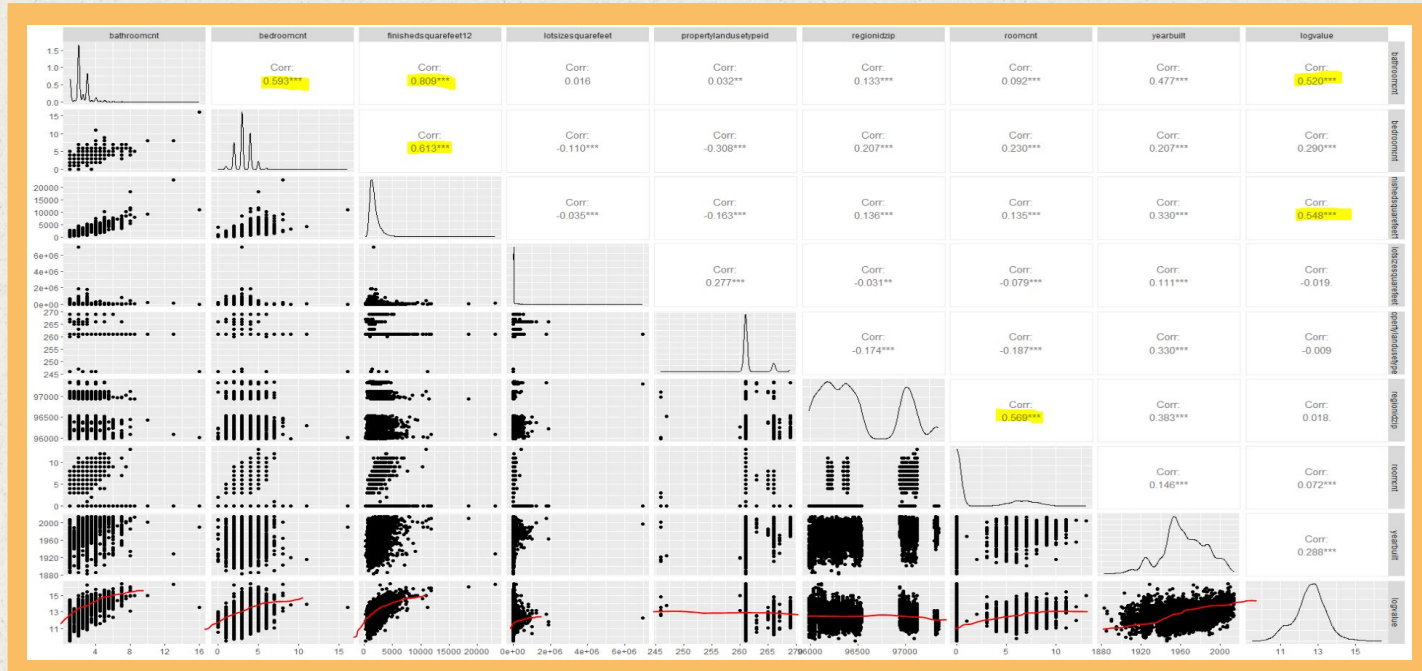usecode
propertycountyland-
usetypeid
rawcensustractandblock
regionidcity
regionidcounty
regionidzip

# Exploratory Data Analysis

# 02.

# Missing Values

From a random sample of 10,000 observations, it was determined that there were 30 parameters that could not be used in the analysis due to missing data. Imputation is not helpful for parameters that are missing more than 10% of data.
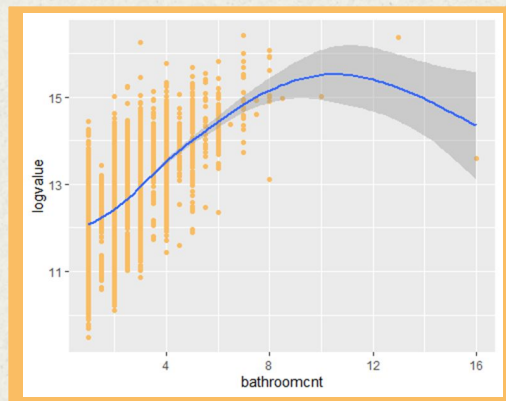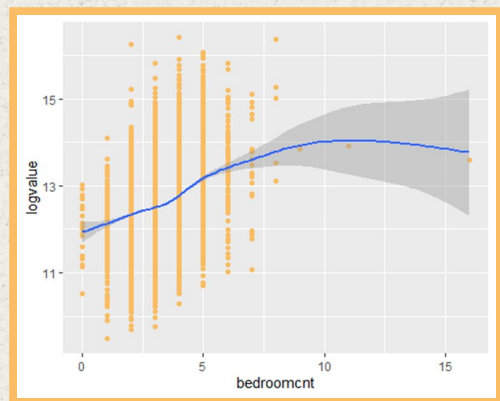
# Assumptions: Multi-Linear Regression

1) linearity
2) homoskedasticity
3) independence of errors
4) normality
5) independence of independent variables
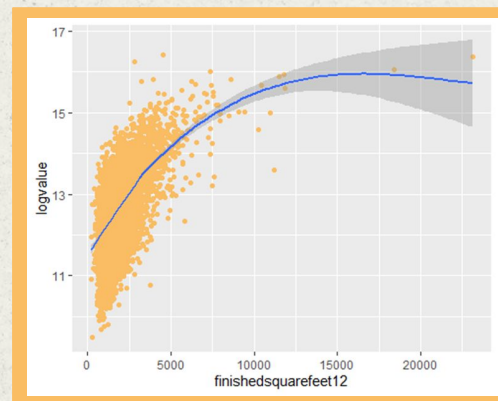
# Simple Linear Regression



# Bathrooms          # Bedrooms          Finished Sq. Ft.
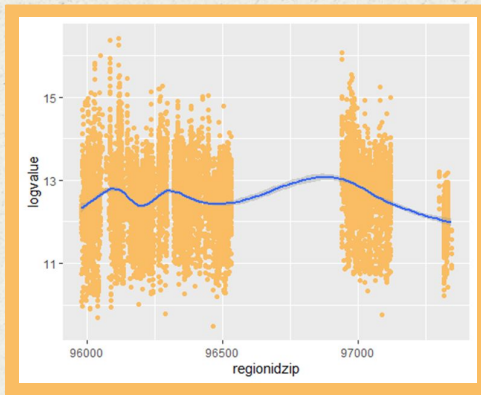
# Simple Linear Regression



Longitude



Zip Code



Finished Sq. Ft.

# Location & History Data



Zip Code

County

Year Built

# Categorical Parameters



Zip Code

County

Property Land use Type

# Model RMSE



## Full Model

RMSE: 0.674



## Logged Model, 2016

RMSE: 0.678



## Logged Model, 2017

RMSE: 0.669



## K-Nearest Neighbors

RMSE: 0.725

# Relational Database

## 03.

# SQL Schema Design: Raw vs Cleaned Data

## Four Tables

## 58 Columns

Number of columns in *raw* data tables

## 24 Columns

Number of columns in *cleaned* data tables

# Importing the Data into MySQL Database



## Table Data Import Wizard

- Easier user-interface
- Slower loading speed

## LOAD DATA Statement

- More configurable
- Faster loading speed

# Errors and Obstacles during Database Implementation

1. Managing Load Speeds

2. Determining Column Data Types

3. Error while loading Data using LOAD DATA statement

LOAD DATA LOCAL INFILE

Convert CSV to SQL

Error Code: 3948. Loading local data is disabled; this must be enabled on both the client and server sides

# Summary Statistics with SQL

```
-- Summary Statistics:
-- Retrieve summary statistics about the home values to get an overview of the data:
-- => Total properties in three counties (Los Angeles, Orange and Ventura, California) data in 2016 = 2985217
-- => min home value = 0, max home value = 282786000, Average home value = 414485.6639, Standard deviation home value = 722871.73
SELECT
  COUNT(*) AS total_properties,
  MIN(taxvaluedollarcnt) AS min_home_value,
  MAX(taxvaluedollarcnt) AS max_home_value,
  AVG(taxvaluedollarcnt) AS avg_home_value,
  STDDEV(taxvaluedollarcnt) AS stdev_home_value
  FROM properties_2016;
```

| total_properties | min_home_value | max_home_value | avg_home_value | stdev_home_value |
|---|---|---|---|---|
| 2985217 | 0 | 282786000 | 414485.6639 | 722871.7307041432 |

# Average Value for Property types SQL

```sql
-- Average Property Value by Property Type:
-- This query calculates the average property value
-- for each property type (e.g., residential, commercial, etc.).
SELECT
  propertylandusetypeid,
  CASE propertylandusetypeid
    WHEN 0 THEN 'Unknown/Not Specified'
    WHEN 31 THEN 'Commercial/Office/Residential Mixed Used'
    WHEN 47 THEN 'Industrial'
    WHEN 246 THEN 'Triplex (3 units)'
    WHEN 247 THEN 'Quadruplex (4 units)'
    WHEN 248 THEN 'Double Wide'
    WHEN 260 THEN 'Residential General'
    WHEN 261 THEN 'Single Family Residential'
    WHEN 263 THEN 'Mobile Home'
    WHEN 264 THEN 'Townhouse'
    WHEN 265 THEN 'Cluster Home'
    WHEN 266 THEN 'Condominium'
    WHEN 267 THEN 'Multi-Family (2-4 units)'
    WHEN 269 THEN 'Cooperative'
    WHEN 270 THEN 'Condominium, Duplex (2 units)'
    WHEN 275 THEN 'Planned Unit Development (PUD)'
    ELSE 'Unknown' -- Handles any other value not listed above
  END AS property_type,
  AVG(taxvaluedollarcnt) AS avg_property_value
FROM properties_2016
GROUP BY propertylandusetypeid;
```

| propertylandusetypeid | property_type | avg_property_value |
|---|---|---|
| 47 | Industrial | 2297121.0582 |
| 260 | Residential General | 1960105.8671 |
| 31 | Commercial/Office/Residential Mixed Used | 691333.4876 |
| 264 | Townhouse | 475094.8578 |
| 248 | Double Wide | 471168.1080 |
| 269 | Cooperative | 470533.0984 |
| 261 | Single Family Residential | 437828.2185 |
| 247 | Quadruplex (4 units) | 379340.5007 |
| 266 | Condominium | 348295.8650 |
| 246 | Triplex (3 units) | 344641.4541 |
| 267 | Multi-Family (2-4 units) | 323832.4518 |
| 265 | Cluster Home | 273258.3506 |
| 275 | Planned Unit Development (PUD) | 173421.2614 |
| 263 | Mobile Home | 19593.9638 |
| 0 | Unknown/Not Specified | 0.0000 |
| 270 | Condominium, Duplex (2 units) | 0.0000 |

# Los Angeles, California Vs Dallas, Texas



**$495,000**  3 bd | 2 ba | 1,305 sqft
15109 Leadwell St, Van Nuys, CA 91405
● **Auction** | View Zestimate ®

| Contact agent |

Overview | Facts and features | Home value | Price and tax h ›

🏠 Single family residence
📅 Built in 1947
🌡 No data
❄ See remarks
🅿 2 Garage spaces
📐 8,108 sqft
📊 $379 price/sqft
💲 2% buyers agency fee



**$492,450**  5 bd | 3 ba | 2,586 sqft
3813 Weisenberger Dr, Dallas, TX 75212
● For sale

**Est.:** $3,407/mo  Ⓢ **Personalize this payment**

| Request a tour<br>as early as tomorrow at 9:00 am | Contact agent |

Overview | Facts and features | Price and tax history | Month ›

🏠 Single family residence
📅 Built in 2022
🌡 No data
❄ No data
🅿 2 Attached garage spaces
📐 0.34 Acres
📊 $190 price/sqft
💲 3.00% buyers agency fee

Average Single Family home of $437,828 for California counties combined

# Average Value for Property Trends (Year Built) SQL

```sql
-- Average Property Value Trend Over Years:
-- This query calculates the average property value based on the year
-- the properties were built, providing insights into how property values
-- have changed over the years.
SELECT
  yearbuilt,
  AVG(taxvaluedollarcnt) AS avg_property_value
FROM properties_2016
GROUP BY yearbuilt
ORDER BY yearbuilt;
```

| yearbuilt | avg_property_value |
|-----------|--------------------|
| 0         | 275432.5542        |
| 1801      | 537285.0000        |
| 1805      | 268630.0000        |
| 1806      | 287696.5000        |
| 1807      | 66851.0000         |
| 1808      | 109836.5000        |
| 1810      | 167898.0000        |
| 1812      | 468763.8000        |
| 1815      | 288879.0000        |
| 1819      | 172396.0000        |
| 1821      | 154795.0000        |
| 1823      | 35449.0000         |
| 1824      | 226361.0000        |
| 1825      | 42240.0000         |
| 1827      | 94978.0000         |
| 1828      | 503556.0000        |
| 1829      | 96150.0000         |
| 1831      | 301264.0000        |
| 1833      | 285864.0000        |
| 1834      | 199836.0000        |

List continues to 2015...



Image of a single family home built in the 1900's in Los Angeles California
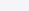
# Feature Analysis SQL

```sql
-- To gain insights about the price of properties with pools vs. properties without pools,
-- you can use a SQL query to calculate various statistics such as the average, minimum,
-- and maximum property prices for each group. Here's how you can do it:
SELECT
  poolcnt,
  COUNT(*) AS property_count,
  AVG(taxvaluedollarcnt) AS avg_property_value,
  MIN(taxvaluedollarcnt) AS min_property_value,
  MAX(taxvaluedollarcnt) AS max_property_value
FROM properties_2016
GROUP BY poolcnt;
```
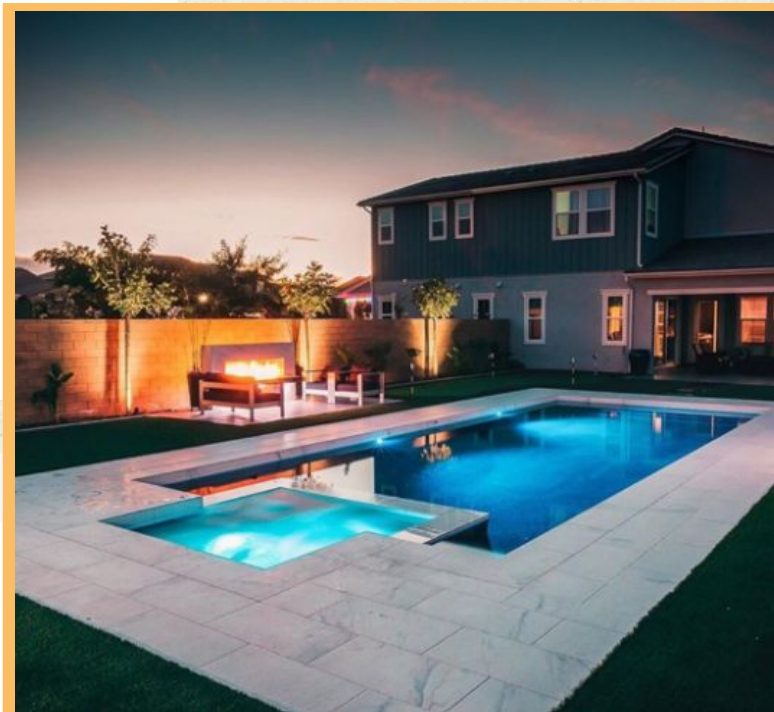
Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| poolcnt | property_count | avg_property_value | min_property_value | max_property_value |
|---------|----------------|--------------------|--------------------|--------------------|
| 1 | 517534 | 633229.1619 | 0 | 149613482 |
| 0 | 2467683 | 368609.7562 | 0 | 282786000 |

- Our model was not able to fit the data set well due to missing data in some crucial columns.
- Having a home built closer to current date increases property value
- When building the RDBMS loading the data using the LOAD DATA statement was faster than using the data import wizard.
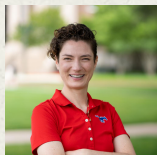
## Future Works

- Inquire to data collectors about missing data.
- Develop a stronger model by integrating different data science techniques together
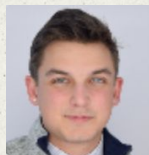
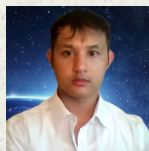# Conclusion

# THANKS!

## DO YOU HAVE ANY QUESTIONS?

Alex Thibeaux
athibeaux@smu.edu

William Jones
wfjones@smu.edu

Chris Mathew
mathewc@smu.edu

Sakava Kiv
skiv@smu.edu