

The housing crisis

Data challenges and opportunities

Hadley Wickham, Garret Grolemund,
Dex Gannon, Gabi Quart, Barret Schloekre

Disclaimer

Some statistics, but more of a focus on the data. What you need **before** you can do statistics

You'll see some extreme, but representative, problems. Almost every analysis has a big data preparation component.

1. Motivation & data questions.
2. Three representative data sources.
3. What is a metropolitan area?
4. Collaboration & reproducibility

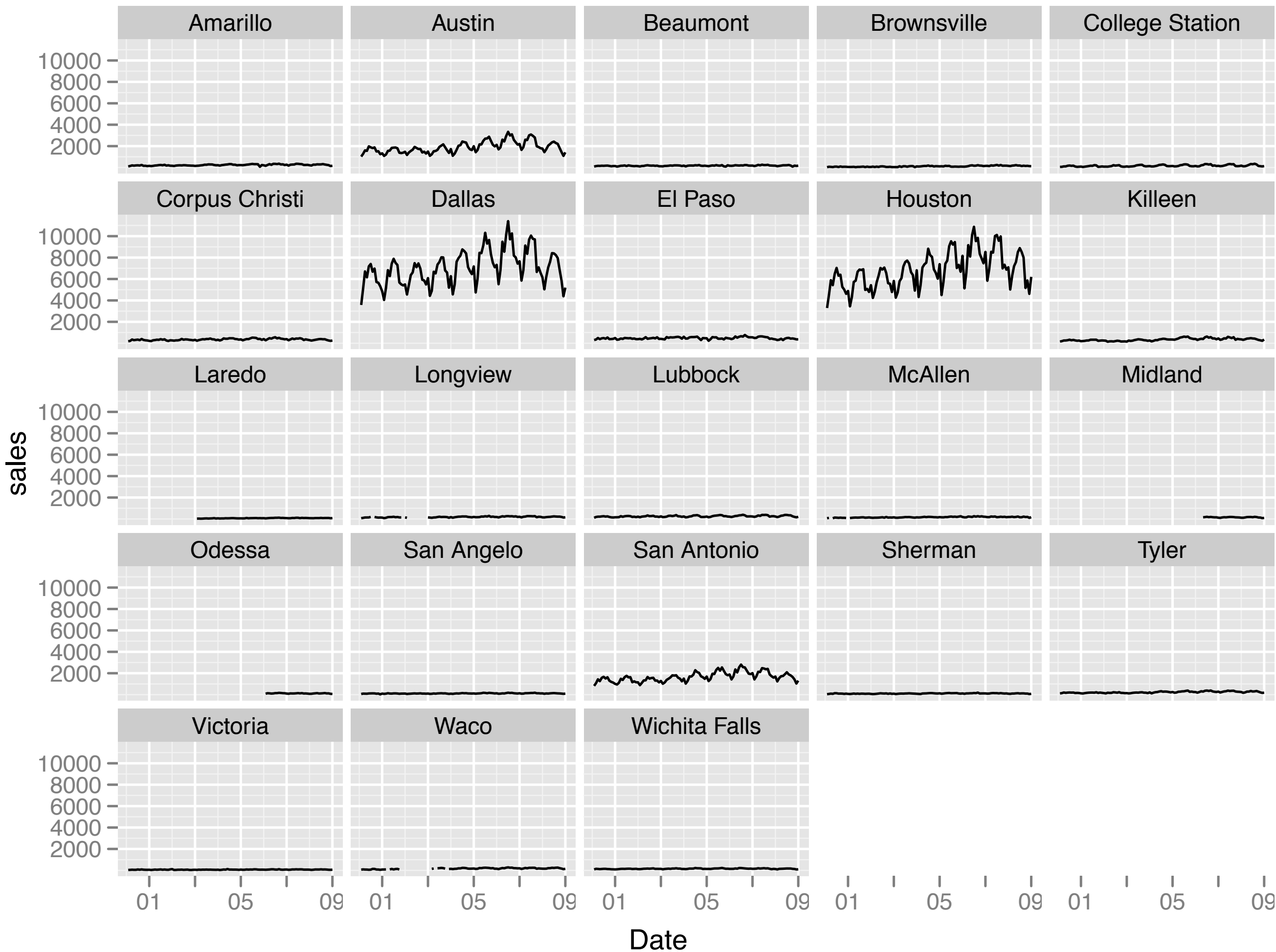


Motivation

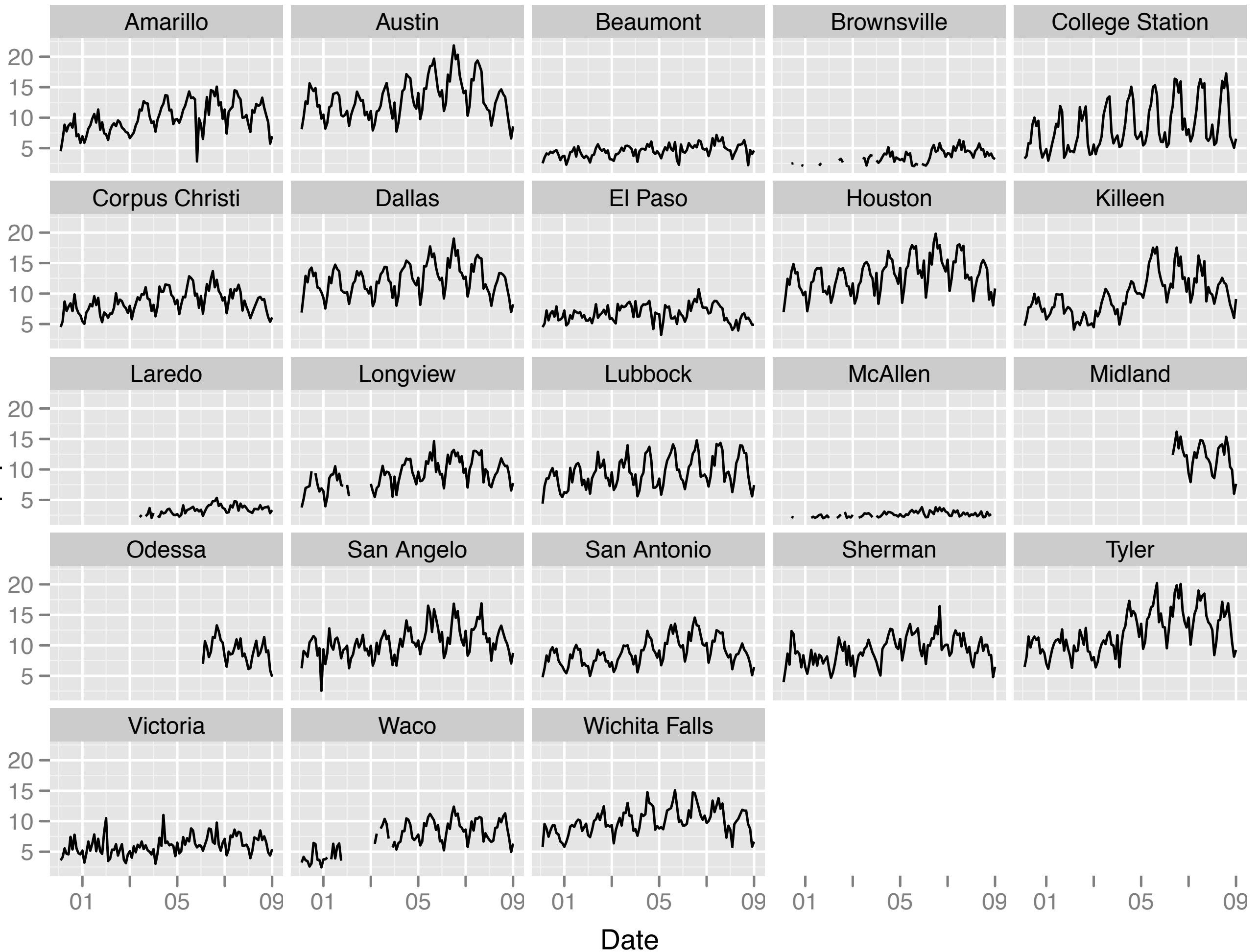
Pretty obvious! **But** the data is hard to find, hidden behind pay walls, hard to use, hard to combine.

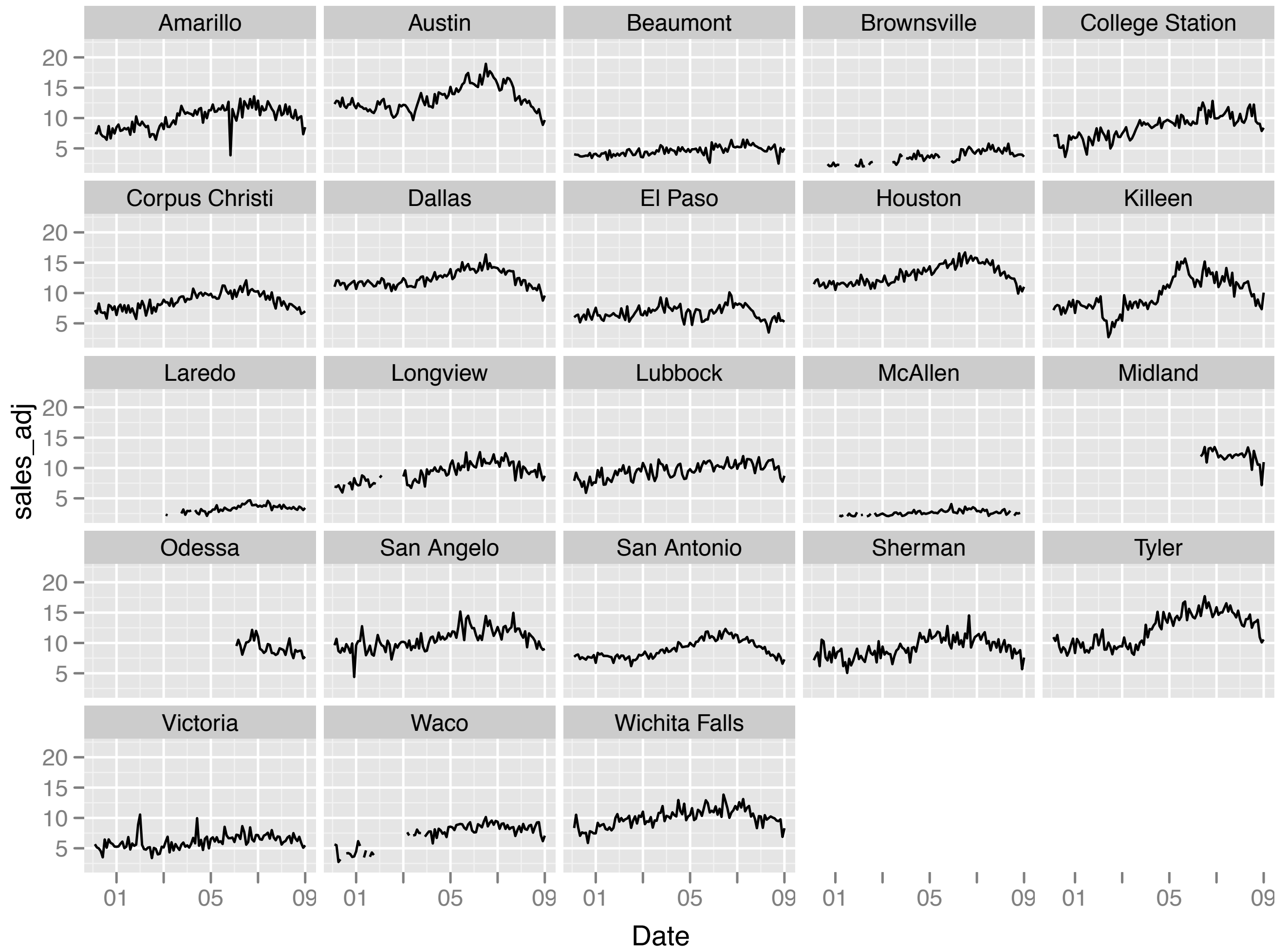
Makes it very difficult to make decisions based on fact, not anecdotes.

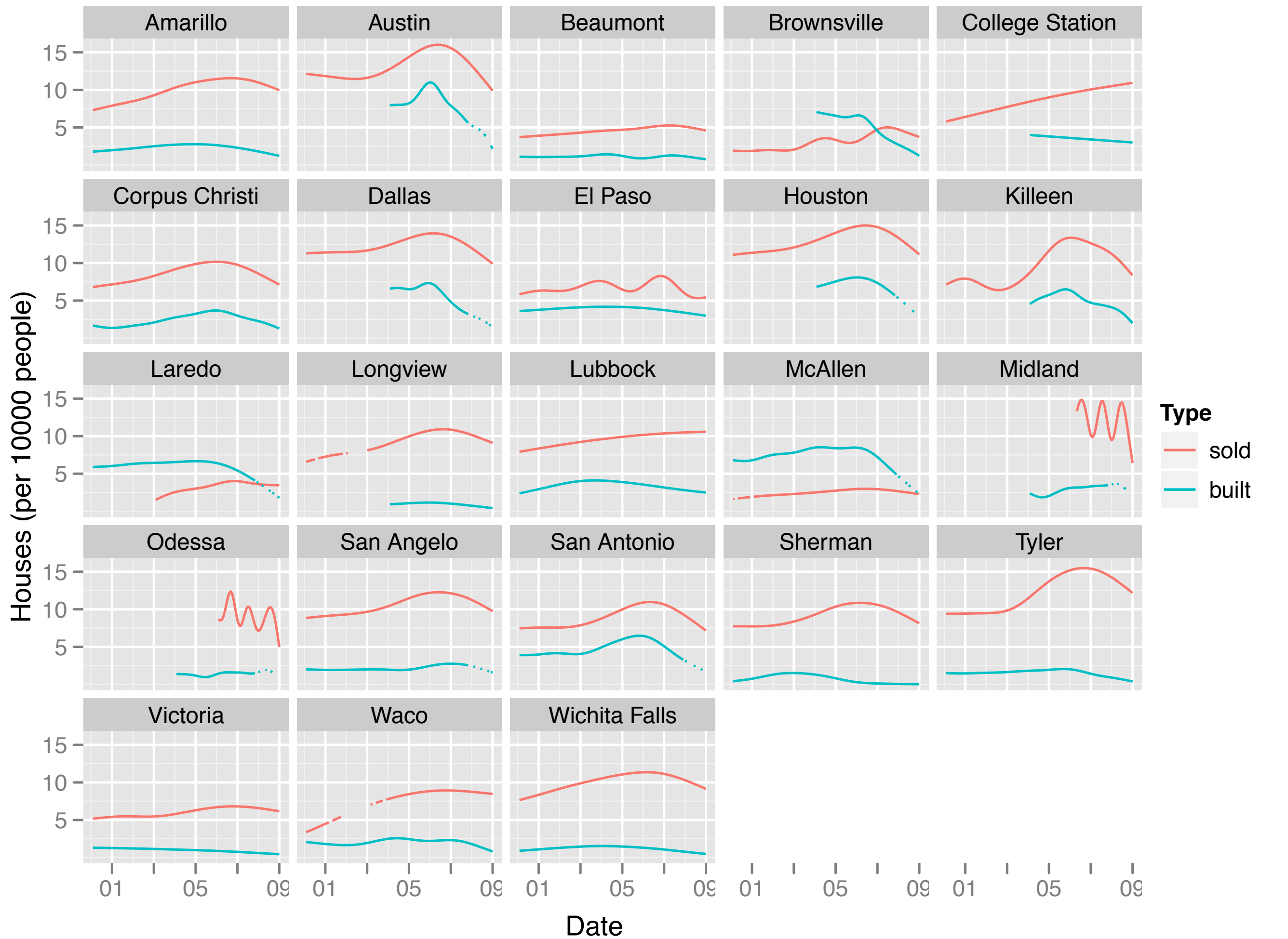
Next: a few examples of the types of things you want to be able to explore

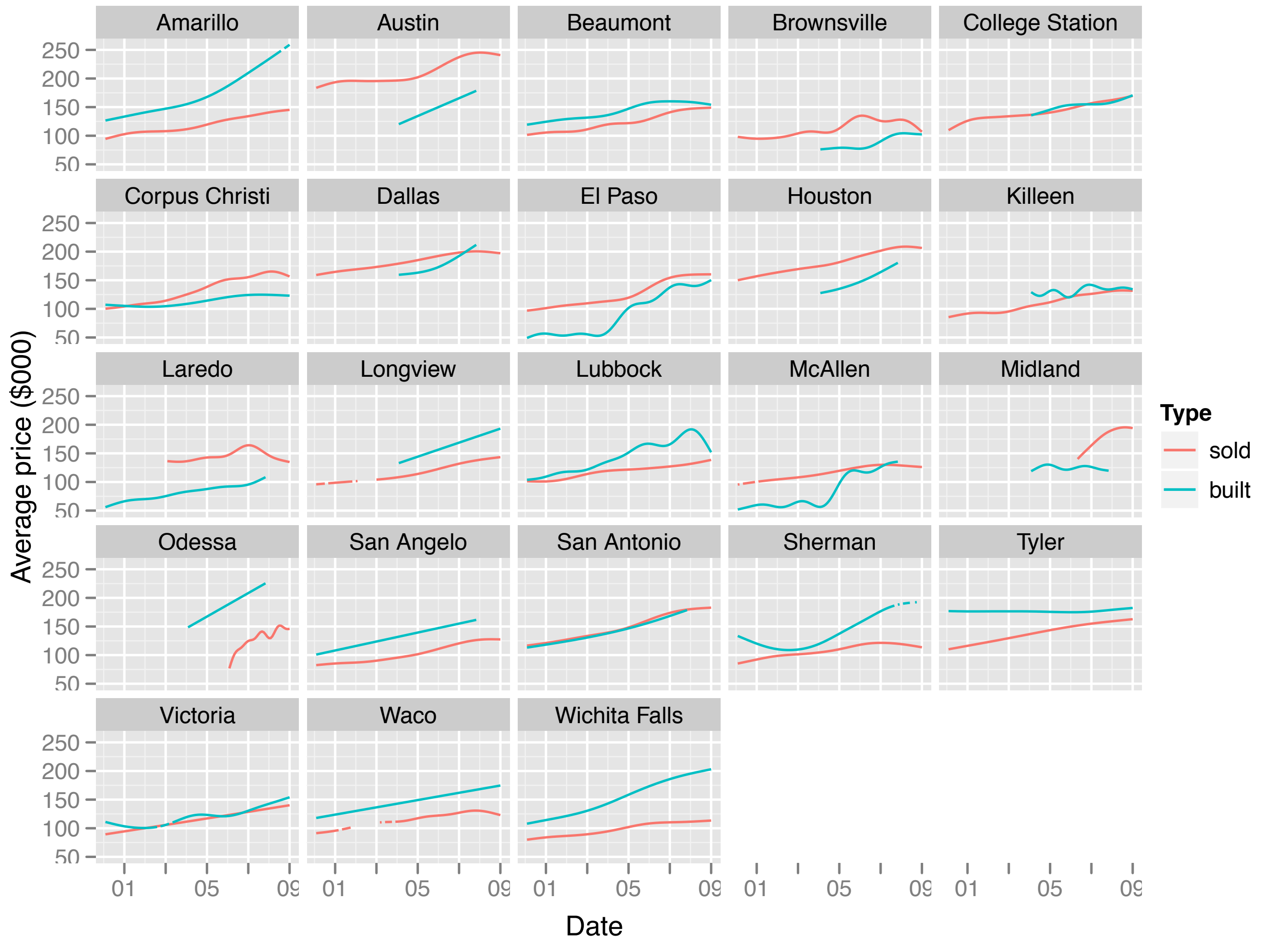


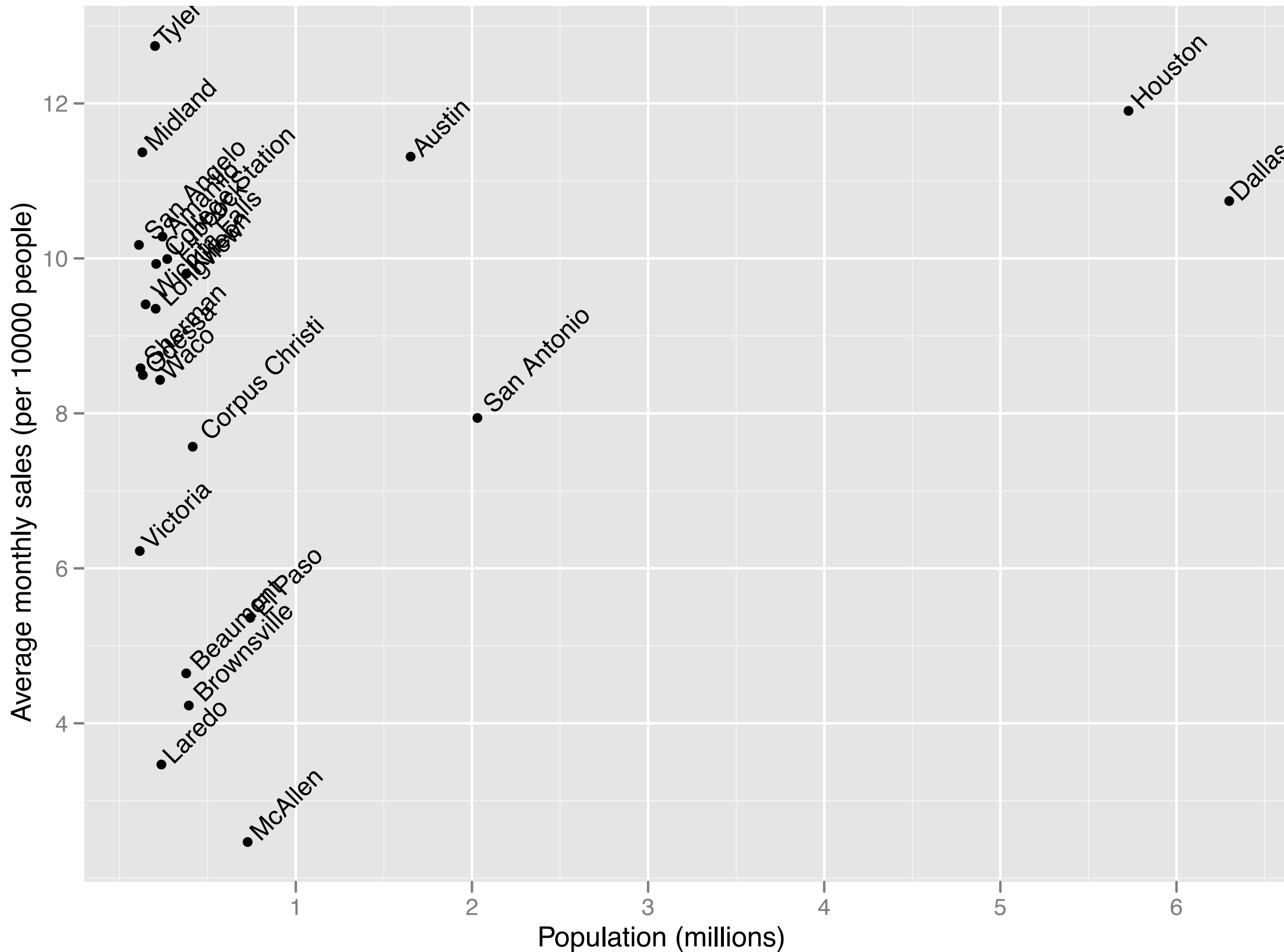
sales/pop * 10000











Data sources

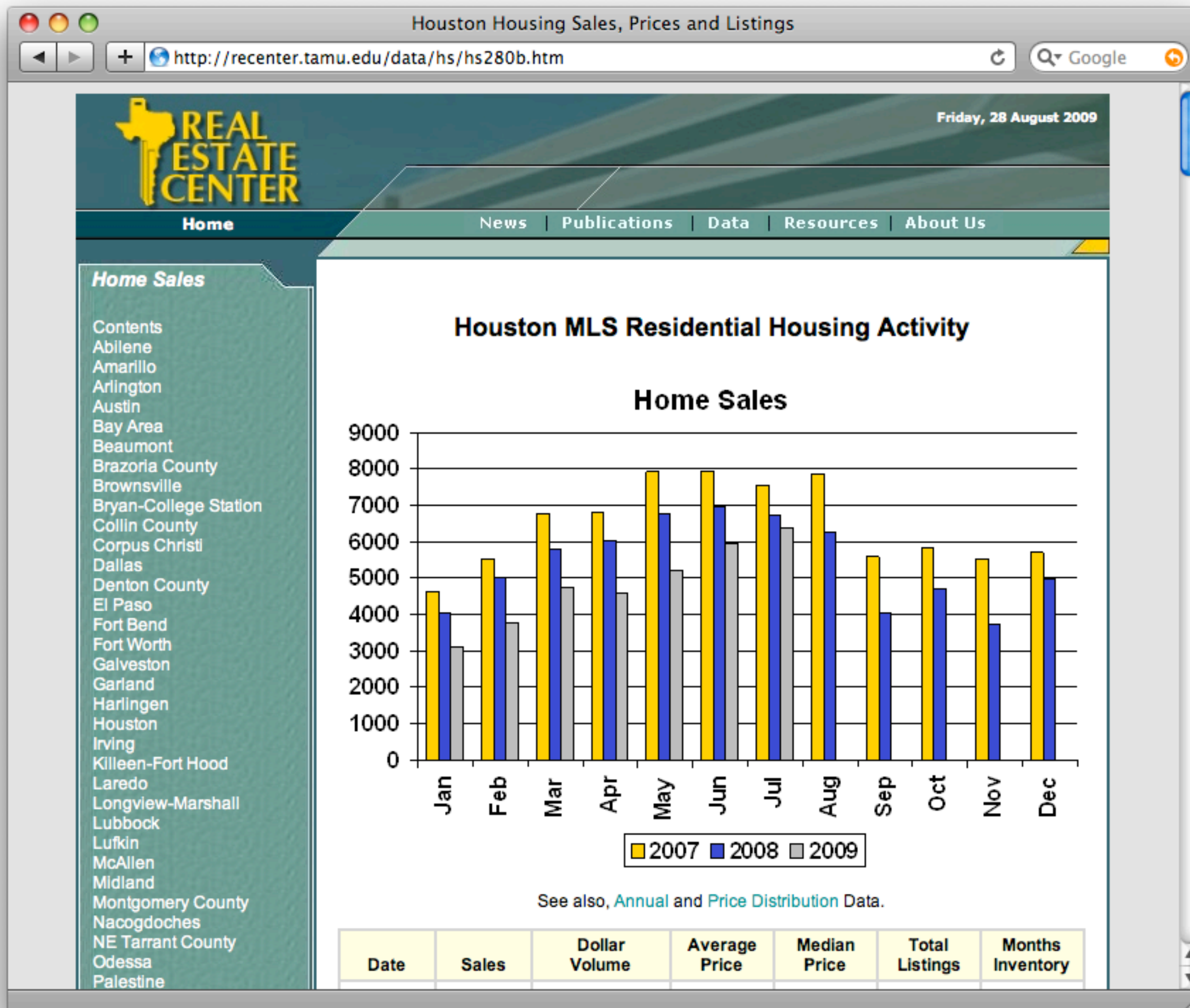
Texas multiple listing data from the Real Estate Center at A&M (sales and average sale price).

New construction data from the census (number of single unit dwellings and average price).

Population data, also from the census.

Combined by metropolitan statistical area.

Multiple listing service data



Houston Housing Sales, Prices and Listings						
http://recenter.tamu.edu/data/hs/hs280b.htm						
Nacogdoches NE Tarrant County Odessa Palestine Paris Port Arthur San Angelo San Antonio San Marcos Sherman-Denison Temple-Belton Texarkana Tyler Victoria Waco Wichita Falls Texas Totals	Date	Sales	Dollar Volume	Average Price	Median Price	Total Listings
	Months Inventory					
	1990-Jan	2,587	228,313,659	88,300	47,300	23,817
	Feb	2,003	182,314,615	91,000	67,500	24,689
	Mar	2,643	234,158,254	88,600	68,300	25,744
	Apr	2,519	234,787,170	93,200	68,200	26,206
	May	3,103	274,674,699	88,500	67,200	26,886
	Jun	3,315	313,691,112	94,600	71,600	26,519
	Jul	3,230	302,314,553	93,600	70,800	26,649
	Aug	3,752	357,123,243	95,200	70,500	25,777
	Sep	2,678	231,788,408	86,600	66,100	24,883
	Oct	2,902	252,502,440	87,000	66,600	24,573
	Nov	2,587	228,313,659	88,300	47,300	23,817
	Dec	2,298	206,631,235	89,900	66,800	23,331
	1991-Jan	1,656	148,683,429	89,800	67,900	22,565
	Feb	2,009	160,113,604	79,700	65,500	23,543
	Mar	2,268	216,367,964	95,400	73,400	24,260
	Apr	2,732	249,169,462	91,200	68,800	25,205
	May	3,345	321,692,198	96,200	78,100	25,854
	Jun	3,294	320,638,203	97,300	75,400	26,811
	Jul	3,229	321,592,804	99,600	77,600	26,068
	Aug	3,401	326,937,004	96,100	68,100	25,750
	Sep	2,747	248,639,729	90,500	80,100	25,596
	Oct	2,568	243,506,648	94,800	72,100	24,983
	Nov	2,789	200,412,654	71,900	66,300	24,032
	Dec	2,453	238,511,083	97,200	73,400	21,470
	1992-Jan	1,599	144,276,401	90,200	75,700	24,110
	Feb	1,737	192,976,831	111,100	74,800	25,027
	Mar	2,505	260,784,136	104,100	78,300	26,162

Data

Building permits
 Employment
 Home Sales
 Housing Affordability
 Population
 Rural Land

```

    <!-- Insert Main content below -->
<p align="center" class="maintitle">Houston MLS Residential Housing Activity</p>
<p align="center"></p>
<p align="center">See also, <a href="hs280a.htm">Annual</a> and <a href="hs280c.htm">Price Distribution</a> Data.</p>
<table border="1" cellspacing="0" cellpadding="3" align="CENTER" bordercolor="#D0D0D0">
<tr align="center" valign="bottom">
<td width="60" bgcolor="#FFFFDD"><b>Date</b></TD>
<td width="60" bgcolor="#FFFFDD"><b>Sales</b></TD>
<td width="90" bgcolor="#FFFFDD"><b>Dollar<br>Volume</b></TD>
<td width="65" bgcolor="#FFFFDD"><b>Average<br>Price</b></TD>
<td width="65" bgcolor="#FFFFDD"><b>Median<br>Price</b></TD>
<td width="65" bgcolor="#FFFFDD"><b>Total<br>Listings</b></TD>
<td width="65" bgcolor="#FFFFDD"><b>Months<br>Inventory</b></TD>
</tr>
<tr align="right">
<TD>1990-Jan</TD><TD>2,587</TD><TD>228,313,659</TD><TD>88,300</TD><TD>47,300</TD><TD>23,817</TD><TD>9.3</TD>
</TR>
<tr align="right">
<TD>Feb</TD><TD>2,003</TD><TD>182,314,615</TD><TD>91,000</TD><TD>67,500</TD><TD>24,689</TD><TD>9.6</TD>
</TR>
<tr align="right">
<TD>Mar</TD><TD>2,643</TD><TD>234,158,254</TD><TD>88,600</TD><TD>68,300</TD><TD>25,744</TD><TD>10.0</TD>
</TR>
<tr align="right">
<TD>Apr</TD><TD>2,519</TD><TD>234,787,170</TD><TD>93,200</TD><TD>68,200</TD><TD>26,206</TD><TD>10.1</TD>
</TR>
<tr align="right">
<TD>May</TD><TD>3,103</TD><TD>274,674,699</TD><TD>88,500</TD><TD>67,200</TD><TD>26,886</TD><TD>10.2</TD>
</TR>
<tr align="right">
<TD>Jun</TD><TD>3,315</TD><TD>313,691,112</TD><TD>94,600</TD><TD>71,600</TD><TD>26,519</TD><TD>9.9</TD>
</TR>
<tr align="right">
<TD>Jul</TD><TD>3,230</TD><TD>302,314,553</TD><TD>93,600</TD><TD>70,800</TD><TD>26,649</TD><TD>9.8</TD>
</TR>
<tr align="right">
<TD>Aug</TD><TD>3,752</TD><TD>357,123,243</TD><TD>95,200</TD><TD>70,500</TD><TD>25,777</TD><TD>9.4</TD>
</TR>

```

Strategy

Locate elements on page.

Use script to extract each occurrence.

Then: firebug + xpath + ruby.

Now: selectorgadget + css selectors + csvget.

Sunday, 30 August 2009



Home

News

Publications

Data

Resources

About Us

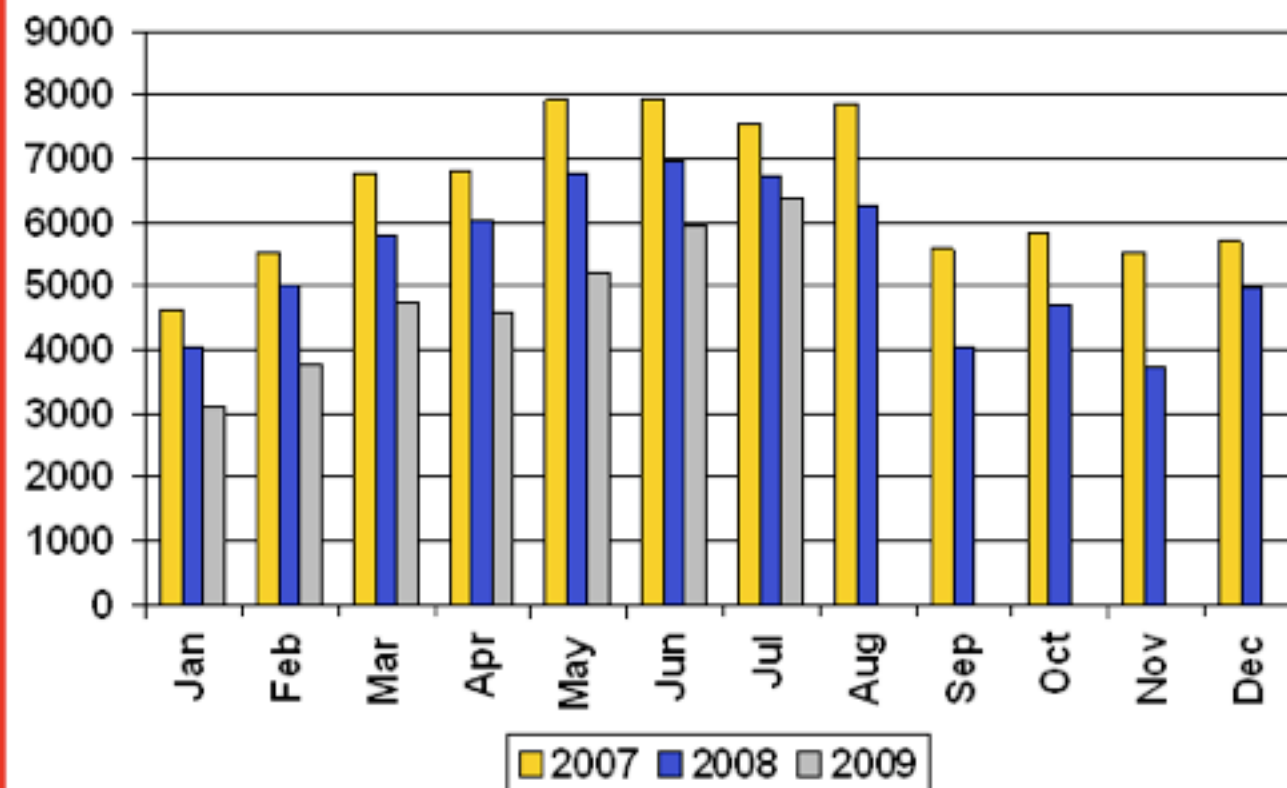
Home Sales

Contents

- Abilene
- Amarillo
- Arlington
- Austin
- Bay Area
- Beaumont
- Brazoria County
- Brownsville
- Bryan-College Station
- Collin County
- Corpus Christi
- Dallas
- Denton County
- El Paso
- Fort Bend
- Fort Worth
- Galveston
- Garland
- Harlingen
- Houston
- Irving
- Killeen-Fort Hood
- Laredo
- Longview-Marshall
- Lubbock
- Lufkin
- McAllen
- Midland
- Montgomery County
- Nacogdoches
- NE Tarrant County
- Odessa
- Palestine
- Paris
- Port Arthur
- San Angelo

Houston MLS Residential Housing Activity

Home Sales

See also, [Annual](#) and [Price Distribution](#) Data.

Date	Sales	Dollar Volume	Average Price	Median Price	Total Listings	Months Inventory
1990-Jan	2,587	228,313,659	88,300	47,300	23,817	9.3
Feb	2,003	182,314,615	91,000	67,500	24,689	9.6
Mar	2,642	234,158,254	88,600	68,300	25,744	10.0
Apr	2,642	234,158,254	88,600	68,300	25,744	10.0
May	3,103	274,674,699	88,500	67,200	26,886	10.2

table:nth-child(5) td

Clear (1652)

Toggle Position

XPath

Help

X

Temple-Belton

Sunday, 30 August 2009



Home

News | Publications

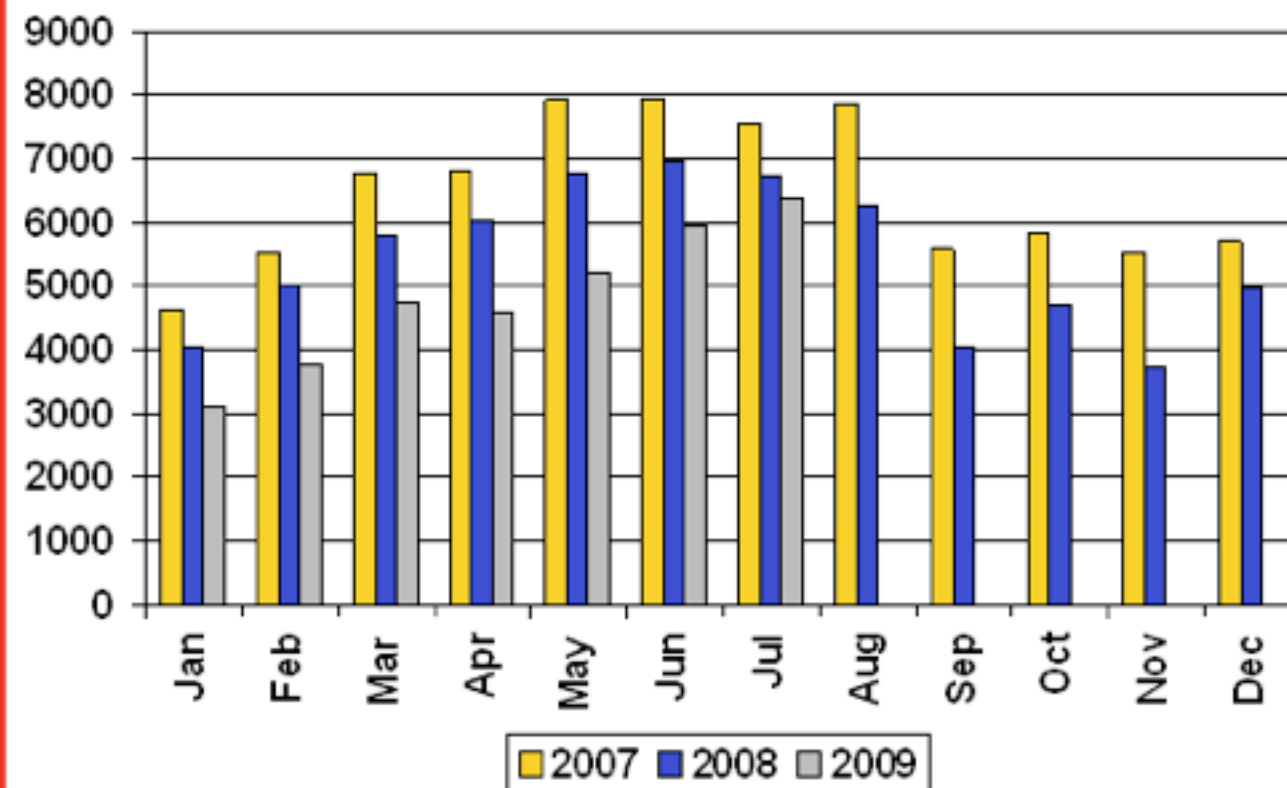
tr td

Home Sales

Contents
 Abilene
 Amarillo
 Arlington
 Austin
 Bay Area
 Beaumont
 Brazoria County
 Brownsville
 Bryan-College Station
 Collin County
 Corpus Christi
 Dallas
 Denton County
 El Paso
 Fort Bend
 Fort Worth
 Galveston
 Garland
 Harlingen
 Houston
 Irving
 Killeen-Fort Hood
 Laredo
 Longview-Marshall
 Lubbock
 Lufkin
 McAllen
 Midland
 Montgomery County
 Nacogdoches
 NE Tarrant County
 Odessa
 Palestine
 Paris
 Port Arthur
 San Angelo

Houston MLS Residential Housing Activity

Home Sales

See also, [Annual](#) and [Price Distribution](#) Data.

Date	Sales	Dollar Volume	Average Price	Median Price	Total Listings	Months Inventory
1990-Jan	2,587	228,313,659	88,300	47,300	23,817	9.3
Feb	2,003	182,314,615	91,000	67,500	24,689	9.6
Mar	2,642	234,158,254	88,600	68,300	25,744	10.0

table:nth-child(5) td

```
{  
  "listings(table:nth-child(5) tr)": [{  
    "cell": "td"  
  }]  
}
```


Final steps

Repeat for all 46 “MSAs” and combine into single csv file.

Fix dates: add missing years and convert months to numeric.

Tidy up column names

Construction data

Table 3u. New Privately Owned Housing Units Authorized
Unadjusted Units by Metropolitan Area

January 2000

					Num of Struc- tures With		
	Total	1 Unit	2 Units	3 and 4 Units	5 Units or More	5 Units or More	
Abilene* TX MSA	16	16	0	0	0	0	
Albany* GA MSA	138	42	0	0	96	12	
Albany-Schenectady-Troy* NY MSA	85	75	0	0	10	1	
Albuquerque* NM MSA	371	337	0	4	30	2	
Alexandria* LA MSA	29	29	0	0	0	0	
Allentown-Bethlehem-Easton* PA MSA	98	70	0	4	24	2	
Altoona* PA MSA	4	4	0	0	0	0	

Table 3u. New Privately Owned Housing Units Authorized
Unadjusted Units by Metropolitan Area

January 2008

Monthly Coverage Percent	Total	1 Unit	2 Units	3 & 4 Units or more		5 Units or more		Num of Struc- tures With 5 Units
				Units	or more	or more	or more	
Abilene, TX	14	10	4	0	0	0	91	
Akron, OH	93	46	0	3	44	7	69	
Albany, GA	24	22	2	0	0	0	84	
Albany-Schenectady-Troy, NY	39	39	0	0	0	0	59	
Albuquerque, NM	204	163	0	0	41	2	100	
Alexandria, LA	41	41	0	0	0	0	97	
Allentown-Bethlehem-Easton, PA-NJ	118	113	0	0	5	1	100	
Altoona, PA	3	3	0	0	0	0	7	

Table 3u. New Privately Owned Housing Units Authorized
Unadjusted Units by Metropolitan Area

January 2008

Different headers

Different column widths

Monthly Coverage

Percent

	Total	1 Unit	2 Units	3 & 4 Units	5 Units or more	Num of Structures With 5 Units or more	
Abilene, TX	14	10	4	0	0	0	91
Akron, OH	93	46	0	3	44	7	69
Albany, GA	24	22	2	0	0	0	84
Albany-Schenectady-Troy, NY	39	39	0	0	0	0	59
Albuquerque, NM	204	163	0	0	41	2	100
Alexandria, LA	41	41	0	0	0	0	97
Allentown-Bethlehem-Easton, PA-NJ	118	113	0	0	5	1	100
Altoona, PA	3	3	0	0	0	0	7

Different wrapping conventions

Different variables

Strategy

Identify consistent patterns across all data sets and turn into code. Heavy use of regular expressions.

Patch up other errors as found.

Apply to all 224 files.


```
# The first line of data is the second line with  
# one or more characters, and a non-blank in the  
# second column
```

```
first <- which(  
  nchar(raw) > 1 & substr(raw, 2, 2) != " ")[2]
```

```
# The last line of data is the first line with  
# less than two characters
```

```
last <- which(  
  nchar(raw[-seq_len(first)]) < 2)[1] + first - 1
```

```
name <- trim(name)
name <- gsub("[*,]", " ", name)
name <- gsub(" (CMSA|MSA|PMSA|P MSA|PM SA|PMS|PMS A)",
  "", name)
name <- gsub("- | -", "-", name)
name <- gsub(" {2,}", " ", name)
```

```
# Random fixes
```

```
name <- gsub("Bea ch", "Beach", name)
name <- gsub("Bernar dino", "Bernardino", name)
name <- gsub("dAlene|d\"Alene", "d'Alene", name)
name <- gsub("Murfreesboro-Franklin",
  "Murfreesboro--Franklin", name)
```

Population

Strategy

Data in single csv file. Phew!

But: data in strange format. Variable name and year combined in column:

POPESTIMATE2000, POPESTIMATE2001,
POPESTIMATE2002, POPESTIMATE2003,
POPESTIMATE2004, POPESTIMATE2005,
POPESTIMATE2006, POPESTIMATE2007,
POPESTIMATE2008

Use reshape package

"city","year","births","deaths","domesticmig","internationalmig","natural
inc","netmig","npopchg_","poestimate","residual","msa_code"

"Akron OH",2000,2252,1442,-96,223,810,127,999,695961,62,10420

"Akron OH",2001,8826,6540,-1112,709,2286,-403,2244,698205,361,10420

"Akron OH",2002,8527,6457,-1696,653,2070,-1043,1258,699463,231,10420

"Akron OH",2003,8352,6580,-1199,539,1772,-660,647,700110,-465,10420

"Akron OH",2004,8320,6669,-1920,521,1651,-1399,272,700382,20,10420

"Akron OH",2005,8272,6821,-1762,533,1451,-1229,21,700403,-201,10420

"Akron OH",2006,8124,6447,-3300,548,1677,-2752,-1148,699255,-73,10420

"Akron OH",2007,8518,6416,-2701,483,2102,-2218,-173,699082,-57,10420

"Akron OH",2008,8548,6527,-3079,485,2021,-2594,-529,698553,44,10420

"Albany GA",2000,568,312,-410,52,256,-358,-109,157759,-7,10500

"Albany GA",2001,2493,1335,-500,120,1158,-380,2055,159814,1277,10500

"Albany GA",2002,2276,1305,-1189,47,971,-1142,199,160013,370,10500

"Albany GA",2003,2223,1401,-96,-127,822,-223,1115,161128,516,10500

"Albany GA",2004,2228,1479,-537,196,749,-341,276,161404,-132,10500

"Albany GA",2005,2315,1335,-462,91,980,-371,508,161912,-101,10500

"Albany GA",2006,2447,1393,-85,150,1054,65,1062,162974,-57,10500

"Albany GA",2007,2454,1356,-184,70,1098,-114,967,163941,-17,10500

"Albany GA",2008,2392,1394,-117,101,998,-16,978,164919,-4,10500

Summary

Common problems

Spread over many files

In unhelpful formats (e.g. html)

Observations and variables confused

Format varies over time

Meaning varies across datasets

What is a
metropolitan area?

MLS data:

Houston

Construction data:

Houston-Galveston-Brazoria

Houston-Baytown-Sugar Land, TX

Houston-Sugar Land-Baytown, TX

Population data:

Houston-Sugar Land-Baytown TX

Metropolitan statistical area (MSA)

Contiguous urban area of at least 50,000 people.

In 2008, there were 362, containing 254 million people, 83% of the population

Defined by Office of Management and Budget. Updated every year.

Updated every year?!

Names change every year!
Fortunately, id codes don't.

But in 2003 they started from scratch.
No practical way to connect pre- and
post-2003 data.

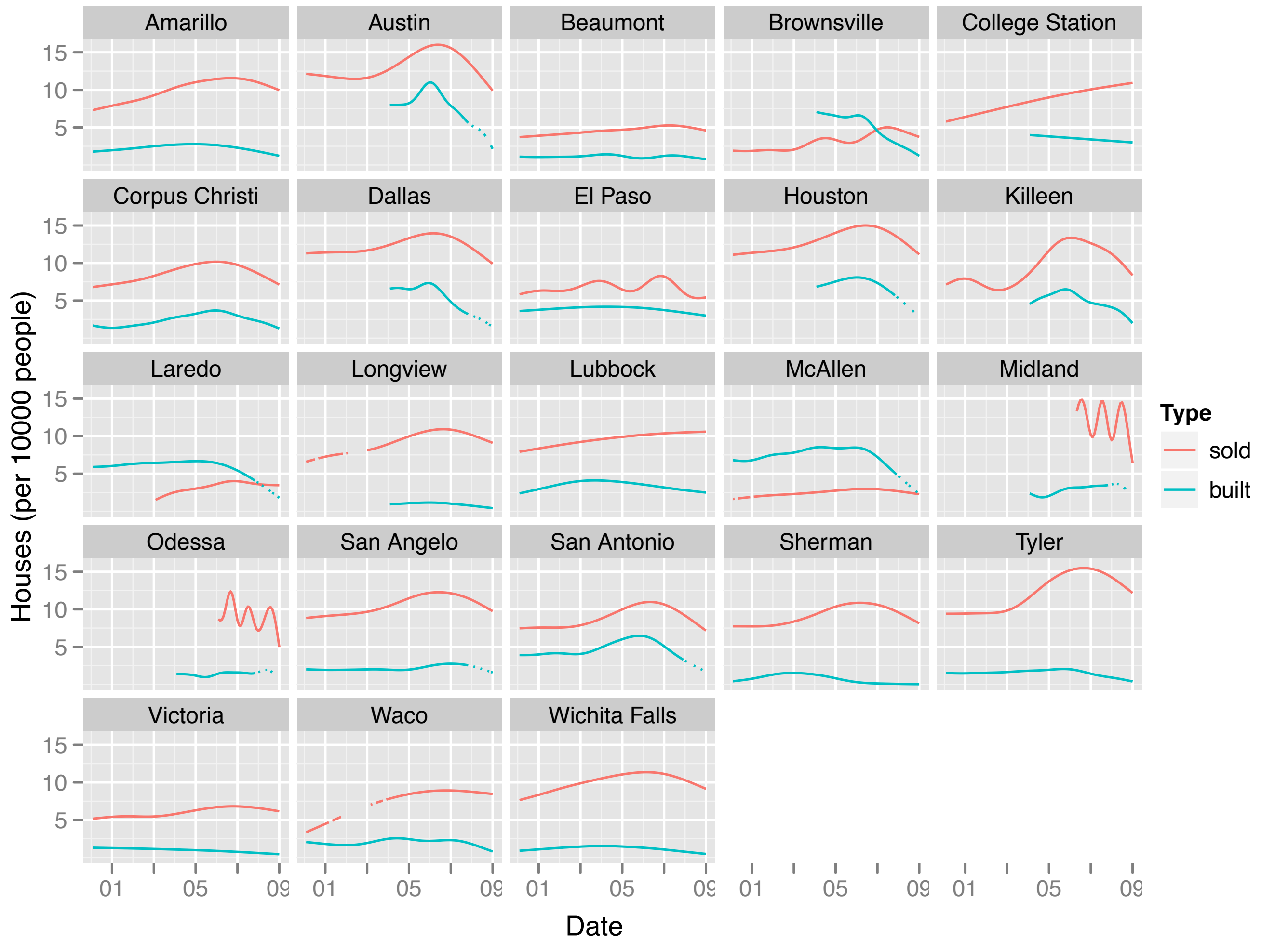
Plus, the Real Estate Centre doesn't use
real MSAs!

Strategy

Build list of all names ever used for an msa id. Apply common set of standardisations: remove commas etc.

Also prepare data set for labelling plots etc, that only contain major city.

Manually match Real Estate Centre areas to MSAs.



Collaboration & reproducibility

Reproducibility

Final results available: you can pick up and use the data

All working shown: you can see what we did and learn from it to help solve new problems.

Every part open licensed. You don't need to ask for permission to use it.

Even this talk!

Git & github

Tools for open, collaborative development

data/msa-changes/3-states.r at 039707ca8f9913049ab493db558fc1030421ff53 from hadley's data-housing-crisis - GitHub

git http://github.com/hadley/data-housing-crisis/blob/039707ca8f9913049ab493db558fc1030421ff53/data/msa-chang RSS Google

github SOCIAL CODING Search Browse | Guides | Advanced


hadley account | profile | log out
132 dashboard | gists

Source Commits Network (0) Fork Queue Issues (4) Downloads (0) Wiki (2) Graphs Admin

hadley / data-housing-crisis edit pull request unwatch download 10 0

Description: Clean data related to the housing crisis [edit](#)
Homepage: Click to edit [edit](#)
Public Clone URL: [git://github.com/hadley/data-housing-crisis.git](https://github.com/hadley/data-housing-crisis.git)
Your Clone URL: git@github.com:hadley/data-housing-crisis.git

Create table of most recent MSA data suitable for labelling plots

 **hadley** (author)
about 6 hours ago

commit 039707ca8f9913049ab493db558fc1030421ff53
tree fa7e5a343d6e3beb981a3fcd613b68efea15b8a6
parent b08724f12e76f521e27f731ac1a859b09603860d

data-housing-crisis / data / msa-changes / 3-states.r

100644 17 lines (13 sloc) 0.5 kb edit raw blame history

```
1 # Separate states from city names and store with one record per state.
2 msa <- read.csv("msa-codes.csv")
3
4 divider <- as.numeric(regexpr(" [A-Z-]+$", msa$city))
5 city <- substr(msa$city, 1, divider - 1)
6 states <- strsplit(substr(msa$city, divider + 1, 100), "-")
7
8 lengths <- sapply(states, length)
9 rep <- rep(1:nrow(msa), lengths)
10 citystate <- data.frame(
11   city = city[rep],
12   state = unlist(states),
```

View code,

data/msa-changes at 039707ca8f9913049ab493db558fc1030421ff53 from hadley's data-housing-crisis - GitHub

git http://github.com/hadley/data-housing-crisis/tree/039707ca8f9913049ab493db558fc1030421ff53/data/msa-change

RSS

Google

name	age	message	history
..			
1-download.r	about 8 hours ago	Add 2008 data. Tidy up parsing code [hadley]	
2-clean.r	about 6 hours ago	Create table of most recent MSA data suitable f... [hadley]	
3-states.r	about 6 hours ago	Separate states processing into own file [hadley]	
msa-codes.csv	about 8 hours ago	Add 2008 data. Tidy up parsing code [hadley]	
msa-major.csv	about 6 hours ago	Create table of most recent MSA data suitable f... [hadley]	
msa-states.csv	about 6 hours ago	Fix type and correct state listing [hadley]	
original/	about 8 hours ago	Add 2008 data. Tidy up parsing code [hadley]	
readme.md	about 7 hours ago	Add a little info about msas and why we need th... [hadley]	

data/msa-changes/readme.md

Metropolitan statistical areas

A metropolitan statistical area is a collection of counties that encompass an urban centre of at least 50,000 people.

This dataset provides a mapping from MSA name to msa code for 2003-2008. MSA names can change on a yearly basis.

Unfortunately there is no easy way to match historical MSAs with current MSAs because prior to 2003 a completely different numbering standard was used. It may be possible to connect the two sets based on the counties that compose each MSA.

Blog | Support | Training | Contact | API | Status | GitHub on Twitter | Help

GitHub™ is Logical Awesome ©2009 | Terms of Service | Privacy Policy

Ruby Hosting by
Engine Yard

and documentation.

Commit History for hadley's data-housing-crisis - GitHub

<http://github.com/hadley/data-housing-crisis/commit/6c5d90a2c3a7f4008da3b41ca59035098644c833>

2009-07-20

<p>Palestine exploration for presentation</p> <p> hadley (author) July 20, 2009</p>	<p>commit 7a3976240a3e7ce04cbcffba36eca74a9ebfd3e3 tree d939d55a86653c77a47129fa4e0a421d7b9f8b84 parent 6c5d90a2c3a7f4008da3b41ca59035098644c833</p>
<p> garrettgman merged branch 'master' of git@github.com:hadley/data-housing-crisis</p>	
<p>modified hpi/explore.r</p> <p> garrettgman (author) July 20, 2009</p>	<p>commit 6f223c4752a88f3e80dd88f776898604c98fbfc3 tree 089cd2485f1130c7d3d4e9bdc0520754fba55e51 parent 15e2c5350bc2e18bc53f28b88fe2eda06f2d332b</p>
<p>moved all of the export (or most of them). FILES WILL NOT WORK. UPDATE THEM.</p> <p> bigbear (author) July 20, 2009</p>	<p>commit 6c5d90a2c3a7f4008da3b41ca59035098644c833 tree a97bd30b3aa44bf324989b97f598702158f35dee parent f6a775a9755a0adf0b790abc95f3ce2b56ba5607</p>
<p>massive file name and directory changing. plus minor things by me</p> <p> bigbear (author) July 20, 2009</p>	<p>commit f6a775a9755a0adf0b790abc95f3ce2b56ba5607 tree 348824d2fcb462d3c082d70f1f60cdc462a4dbdb parent 15e2c5350bc2e18bc53f28b88fe2eda06f2d332b</p>

2009-07-16

<p> garrettgman merged branch 'master' of git@github.com:hadley/data-housing-crisis</p>	
<p>hpi vs. out come by state and california</p> <p> garrettgman (author) July 16, 2009</p>	<p>commit ee52f1ca4d8df596f17b50230c87f172ebe5025e tree ca91447795542c06d8f775d639d9345f9d1be186 parent 15317a061b8bffc4473a28179d036fcb1fe36e</p>
<p>i hate readme</p> <p> gquart6 (author) July 16, 2009</p>	<p>commit fb4c92d899e0d5394cbb2784e21b47ce9208e532 tree a9b094da4a0efe7863a01f313d668d2d0fd01267 parent b69480c505cdb020944fac924bb79be887f20c3e</p>

Watch changes.

The data-housing-crisis Network - GitHub

git http://github.com/hadley/data-housing-crisis/network

github SOCIAL CODING

Browse | Guides | Advanced

hadley account | profile | log out
132 dashboard | gists

Source Commits Network (0) Fork Queue Issues (4) Downloads (0) Wiki (2) Graphs Admin

Graph Members Feed

hadley / data-housing-crisis edit pull request unwatch download 10 0

Description: Clean data related to the housing crisis edit
Homepage: Click to edit edit
Public Clone URL: git://github.com/hadley/data-housing-crisis.git
Your Clone URL: git@github.com:hadley/data-housing-crisis.git

The data-housing-crisis network graph

All branches in the network using **hadley/data-housing-crisis** as the reference point. [Read our blog post about how it works.](#)

Show Help

This graph has new data available. Reload the page to see it!

hadley

Combine work.

Issues - hadley/data-housing-crisis - GitHub

git http://github.com/hadley/data-housing-crisis/issues

github SOCIAL CODING

Search

hadley 132

account | profile | log out
dashboard | gists

Source Commits Network (0) Fork Queue Issues (4) Downloads (0) Wiki (2) Graphs Admin

hadley / data-housing-crisis edit pull request unwatch download 10 0

Description: Clean data related to the housing crisis edit

Create Issue

Select: All, None Actions

Search Issues

Sort by: Priority | Votes | Last Updated

Open (4)
Unread (6)
Closed

2. Gabi Get laborforce by msa
▲ 0 votes 0 comments Created about 1 month ago by hadley

10. Garrett compare second home information to housing crisis performance
▲ 0 votes 1 comment Created about 1 month ago by garrettgman

12. Write a 2-3 paragraph description of the project
▲ 0 votes 0 comments Created about 1 month ago by hadley

14. Reorganise explorations
▲ 0 votes 0 comments Created about 1 month ago by hadley

Labels New

- Barret
- Data
- Dex
- Gabi
- Garrett

Blog | Support | Training | Contact | API | Status | GitHub on Twitter | Help

GitHub™ is Logical Awesome ©2009 | Terms of Service | Privacy Policy

Ruby Hosting by Engine Yard

Track tasks.

Future work

More publicity.

Start actually using the data!

Add more data sources: foreclosures?
mortgage data?