

Estimating Costal Property Values in Florida

Langyi Tian

Executive summary

Objective

- Organizational mission: Quantify the financial risk of tidal flooding to address stakeholder concern
- Task for this research: Make the numbers more accurate with market value estimation

Data

- Administrative property records (3 million), transaction records from ATTOM Data Solutions
- Demographic data from census data
- Flooding projections from National Oceanic and Atmospheric Administration (NOAA)

Methodology

- Build separate models within each city and county
- Regularization models (Ridge, LASSO) as baseline
- Regression trees e.g. random forest and gradient boosting as comparison for trial models

Findings

- Trees work better than baseline
- Random forest algorithm consistently outperformed by gradient boosting
- City-level models have varying performance
- Ability to predict within 10% deviation in some cities

Project roadmap

1. Data preparation
2. Exploratory analysis
3. Separate modeling: one model for each city/county data
4. Feature selection with individual models
5. Build separate predictive models for 85 city-level subsamples
6. Functionalities to test and parse out performance metrics for regularized models and regression trees
7. Personalize data filtering parameters
8. Iteration through cities/counties and view cross-validated model performance, map county-level model performance

1. Data preparation

Select features from property records in a real estate broker's database.

```
#Select features to import, subset and save
home_dta <- select(
  home_dta_original,
  attomid, #Matching ID
  deedlastsaleprice, #Transaction price last sale
  situsstatecode, #State code
  situscounty, #County code
  propertyaddresscity, #City code
  ownertypedescription1, #First owner is individual/company?
  ownertypedescription2, #Second owner is individual/company?
  deedlastsaledate, #Date of market sale
  yearbuilt, #Year when built
  propertyusegroup, #Commercial/residential?
  areabuilding, #Living area in sq. feet
  censustract, #Census tract division
  propertylatitude, #Lat of property
  propertylongitude, #Lon of property
  roomsatticflag, #See below
  parkinggarage:communityrecroomflag #A series variable measuring physical attributes of the property, including room
)
```

1. Data preparation

Select features from environmental risk and demographic data set constructed by Porter

```
risk_dta <-  
  risk_dta_original %>% select(attomid = ATTOM_ID, #ID  
                                dist_coast, #Distance to coast  
                                mdkt32, #Flooding probability estimate in next years  
                                Totpopbg:near_reading_rates, )#a set of demographic information varying by census tra
```

Get a simple feature of transaction frequence from transaction records

```
trans_dta <-  
  trans_dta_original %>% group_by(attomid) %>% summarize(trans_times = n())#Number of transactions  
  
## Up to here, the data dimension is 3327923 * 203
```

1. Data preparation

- Drop all unary features

```
## Up to here, the data dimension is 3327923 * 91
```

- Drop variables with too many levels, besides those are numerical

```
## [1] "hvacheatingdetail"  
## [1] "exterior1code"  
## [1] "roofmaterial"  
## [1] "roofconstruction"
```

```
## Up to here, the data dimension is 3327923 * 87
```

- Recoded characters to factors

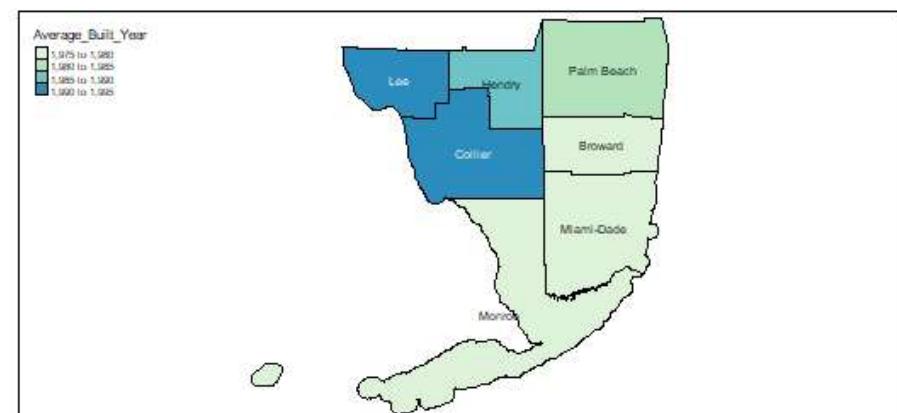
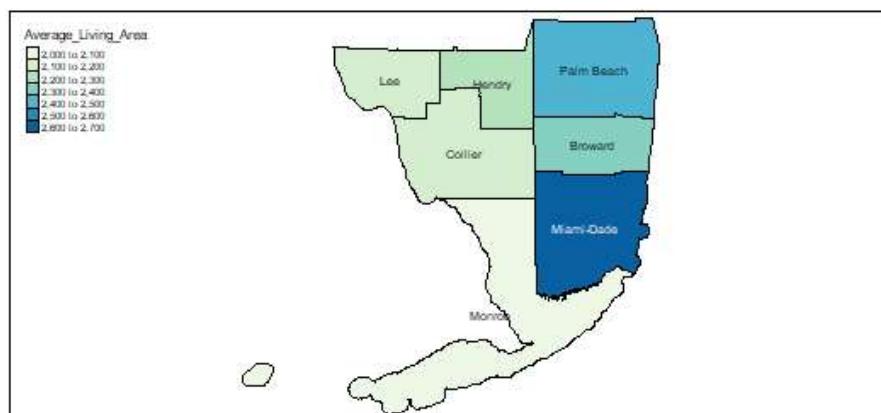
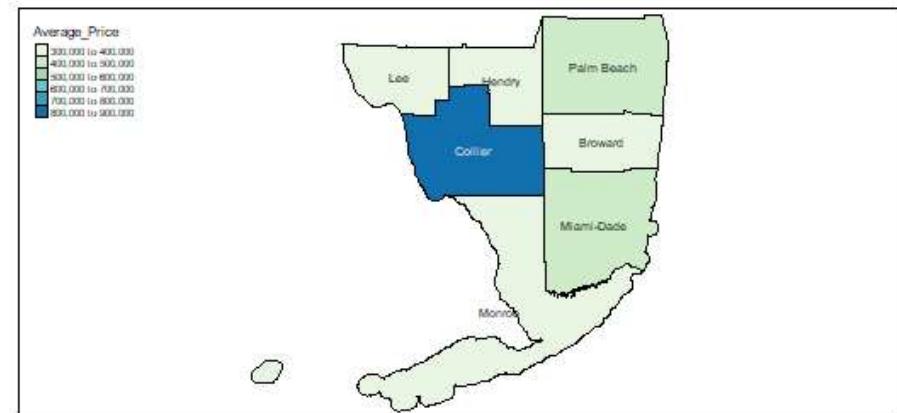
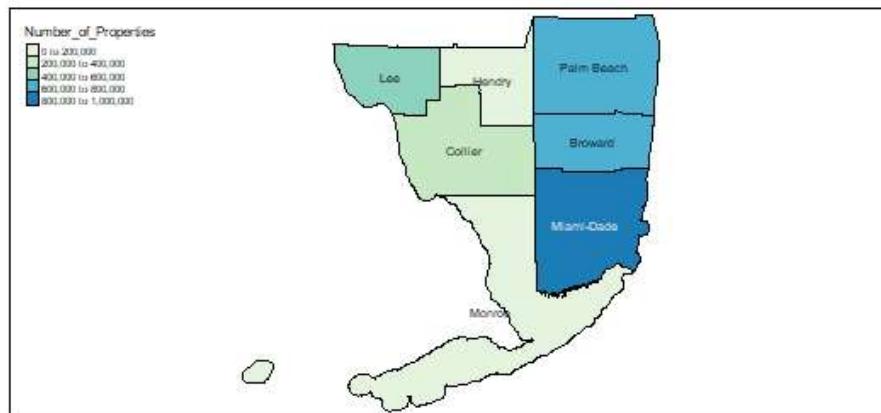
```
## These variables are recoded to factors
```

```
## [1] "situsstatecode"  
## [1] "situscounty"  
## [1] "propertyaddresscity"  
## [1] "ownertypedescription1"  
## [1] "ownertypedescription2"  
## [1] "propertyusegroup"  
## [1] "viewdescription"  
## [1] "porchcode"
```

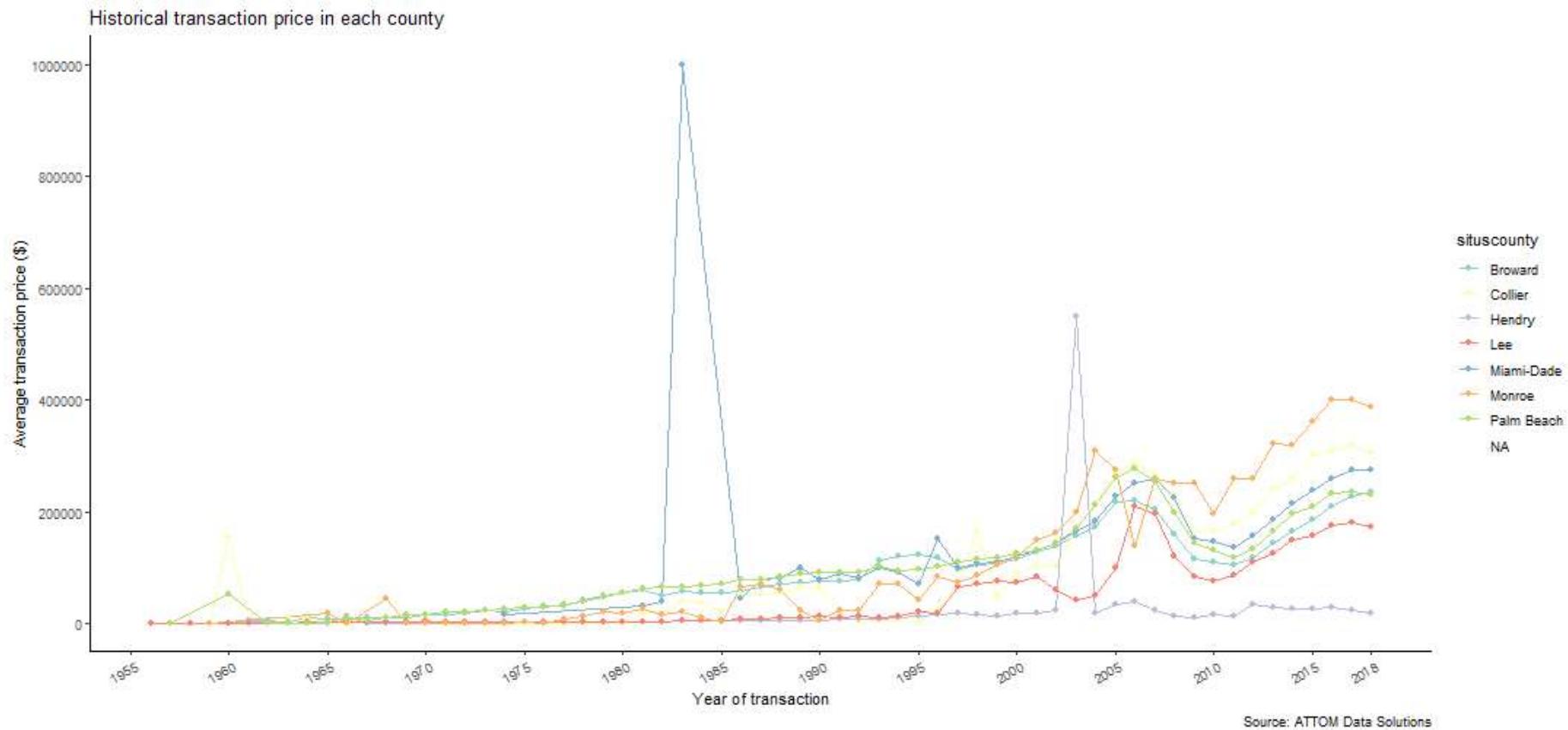
- Hold out cities with sample size too small to go into tree model.

2. Exploratory analysis

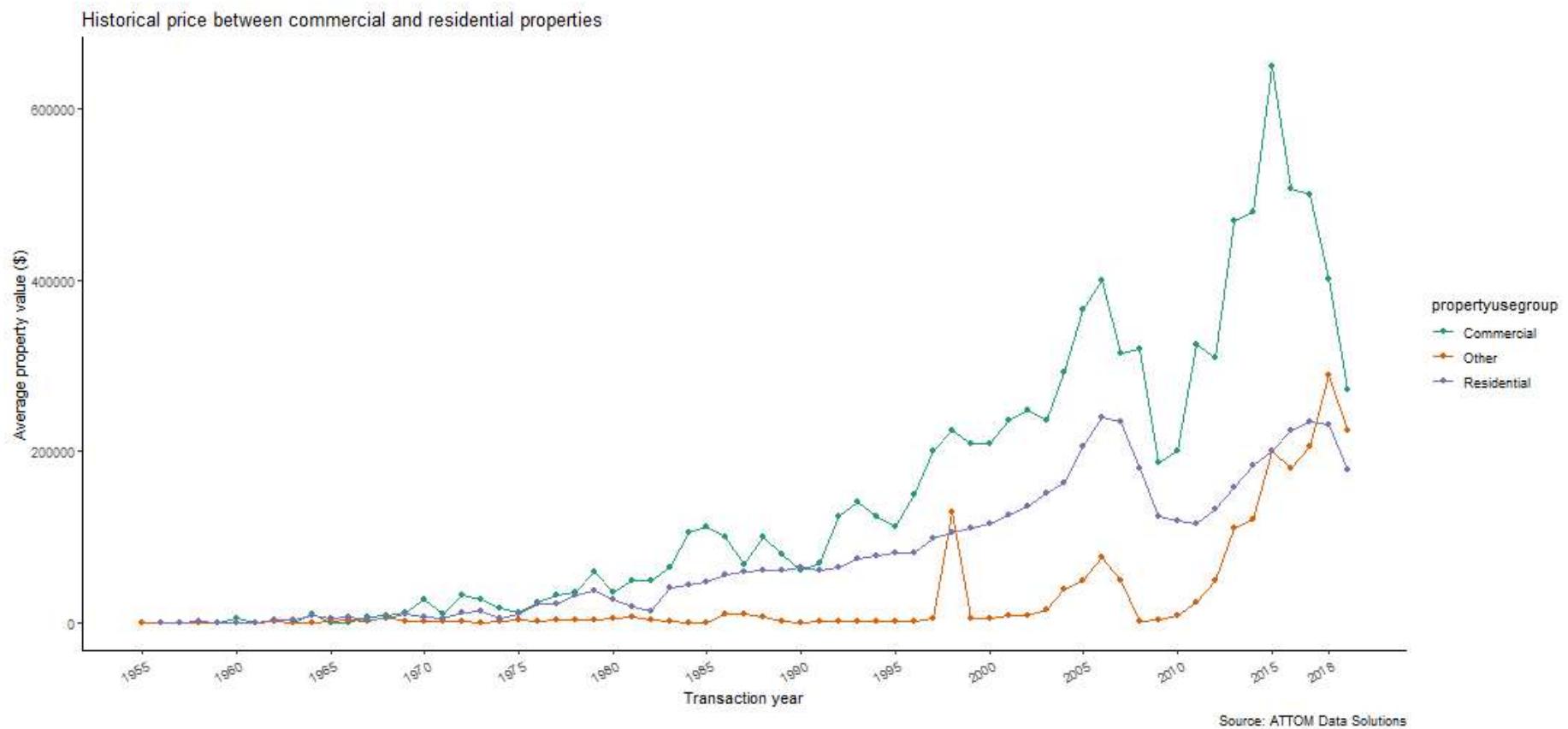
- Summarize a few key variables by county to see the geographical variation



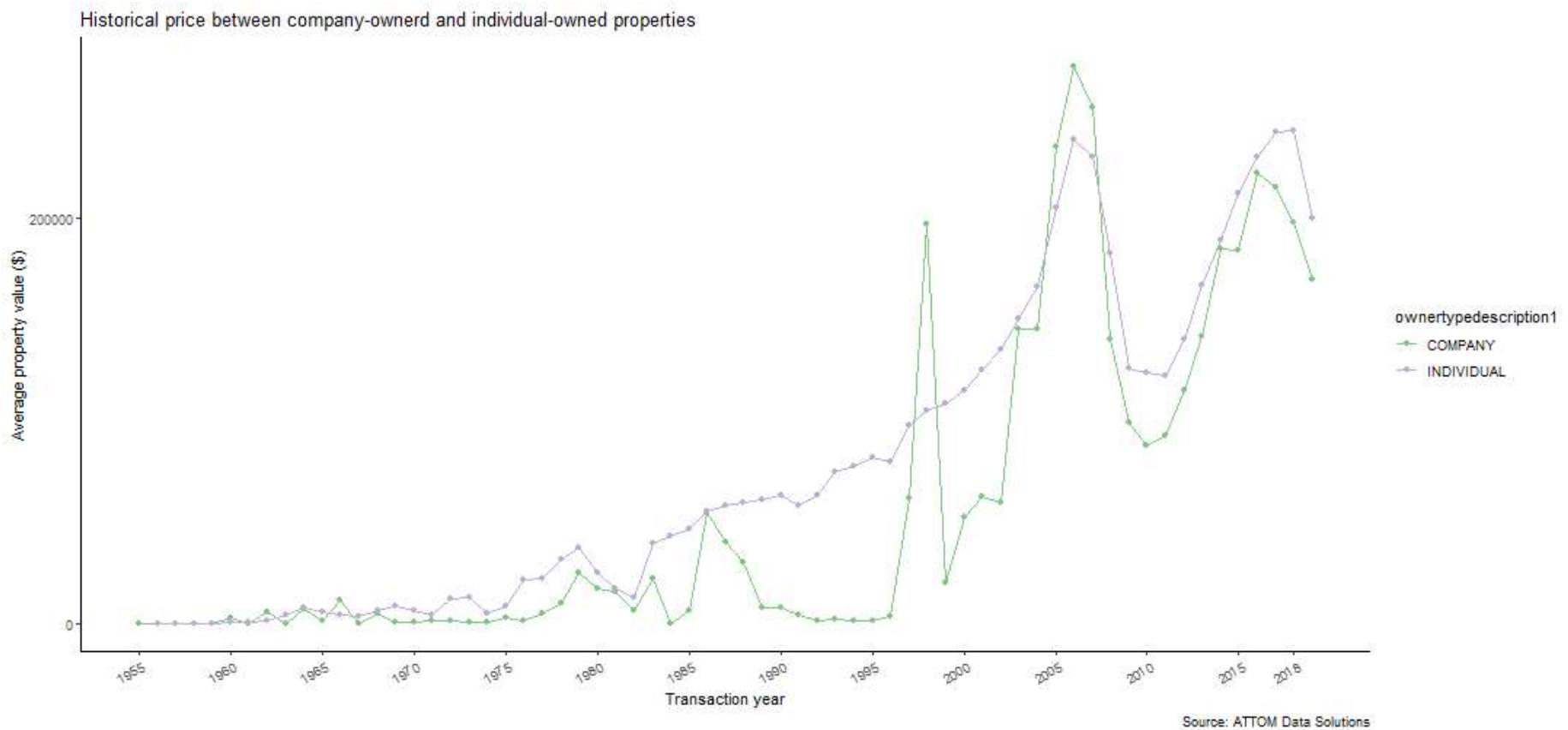
2. Exploratory analysis



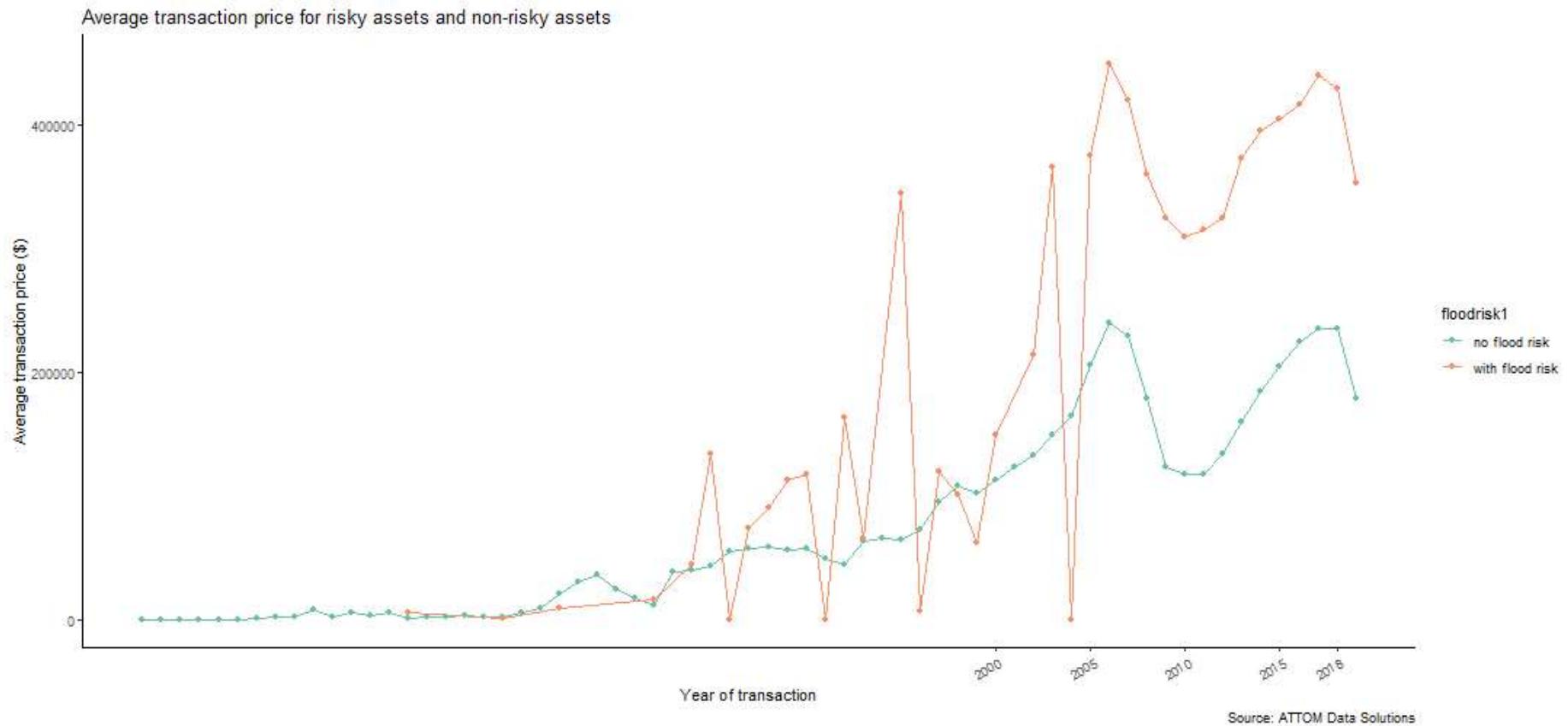
2. Exploratory analysis



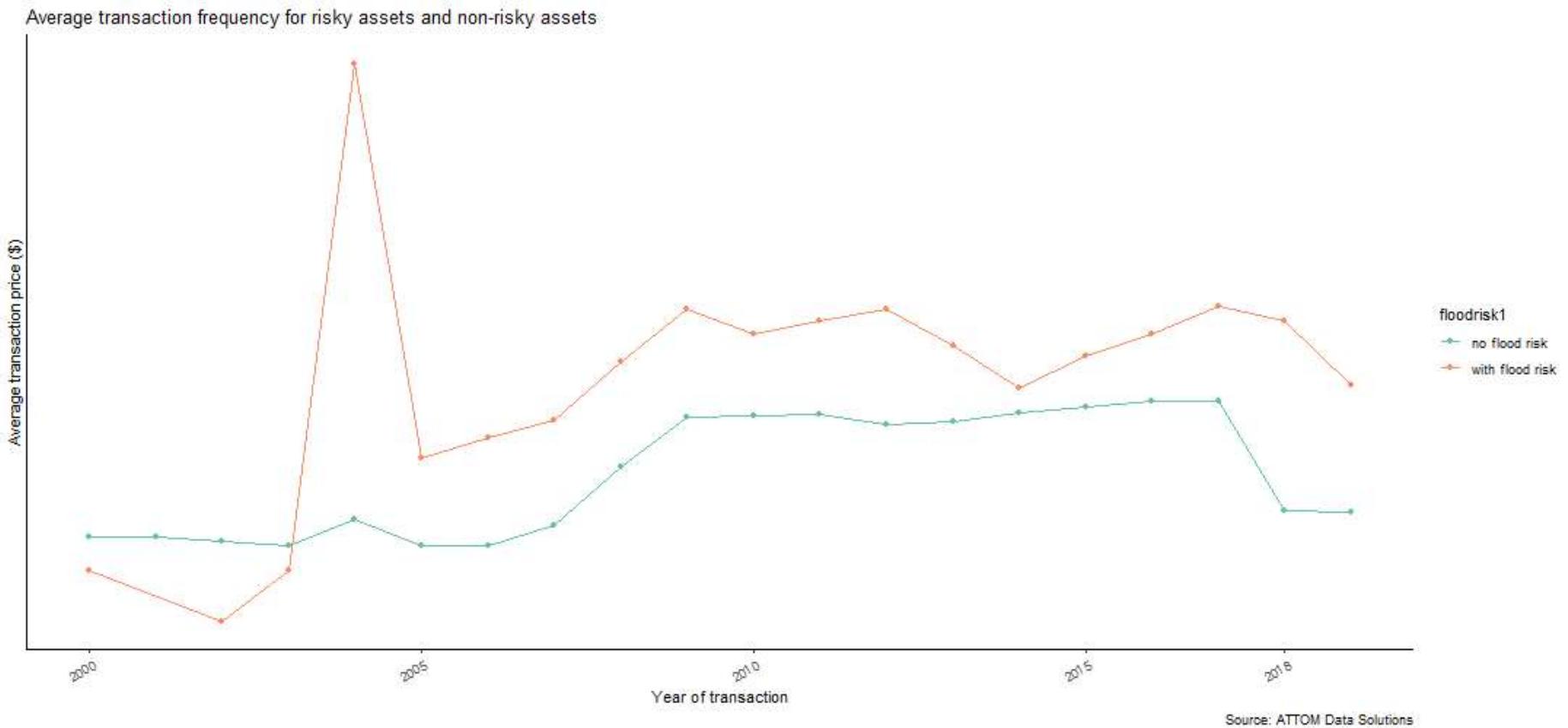
2. Exploratory analysis



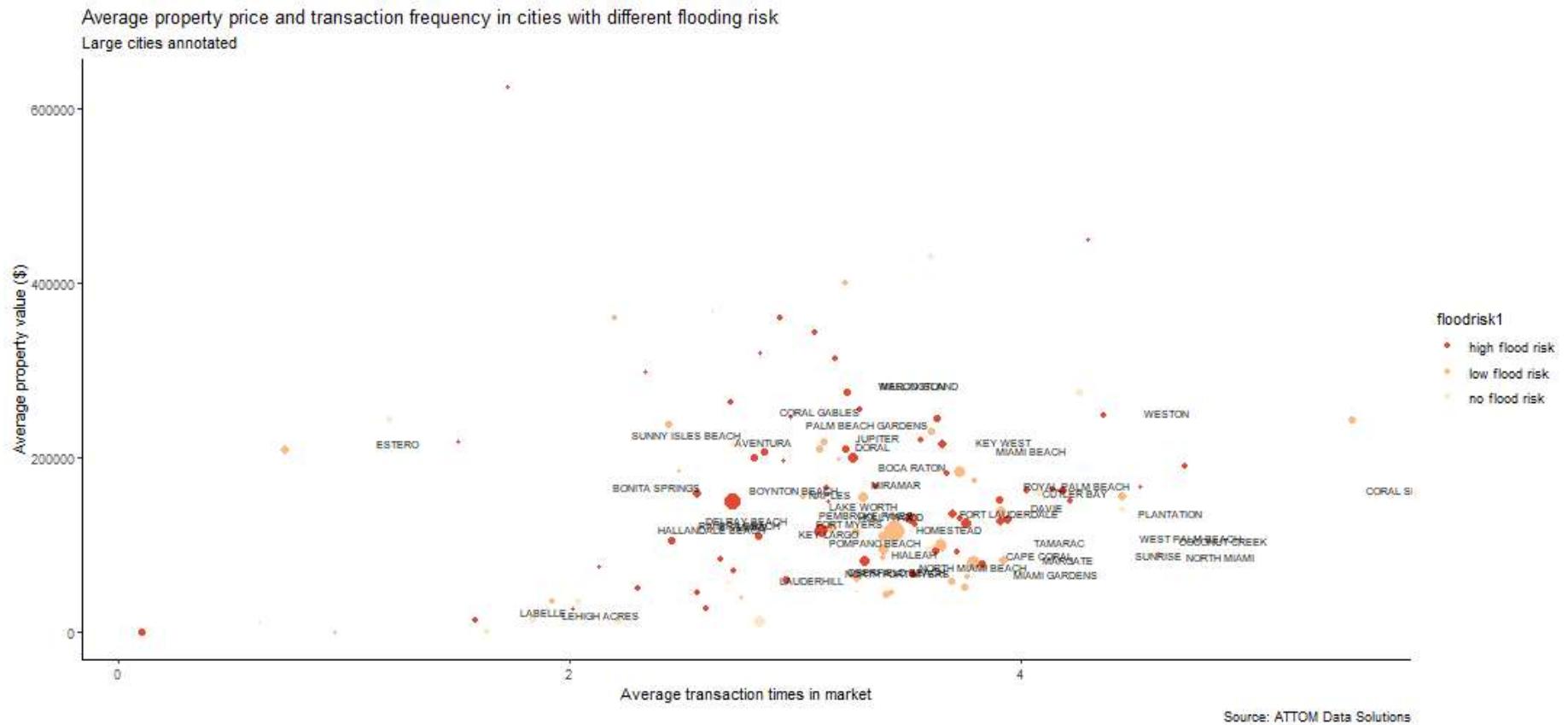
2. Exploratory analysis



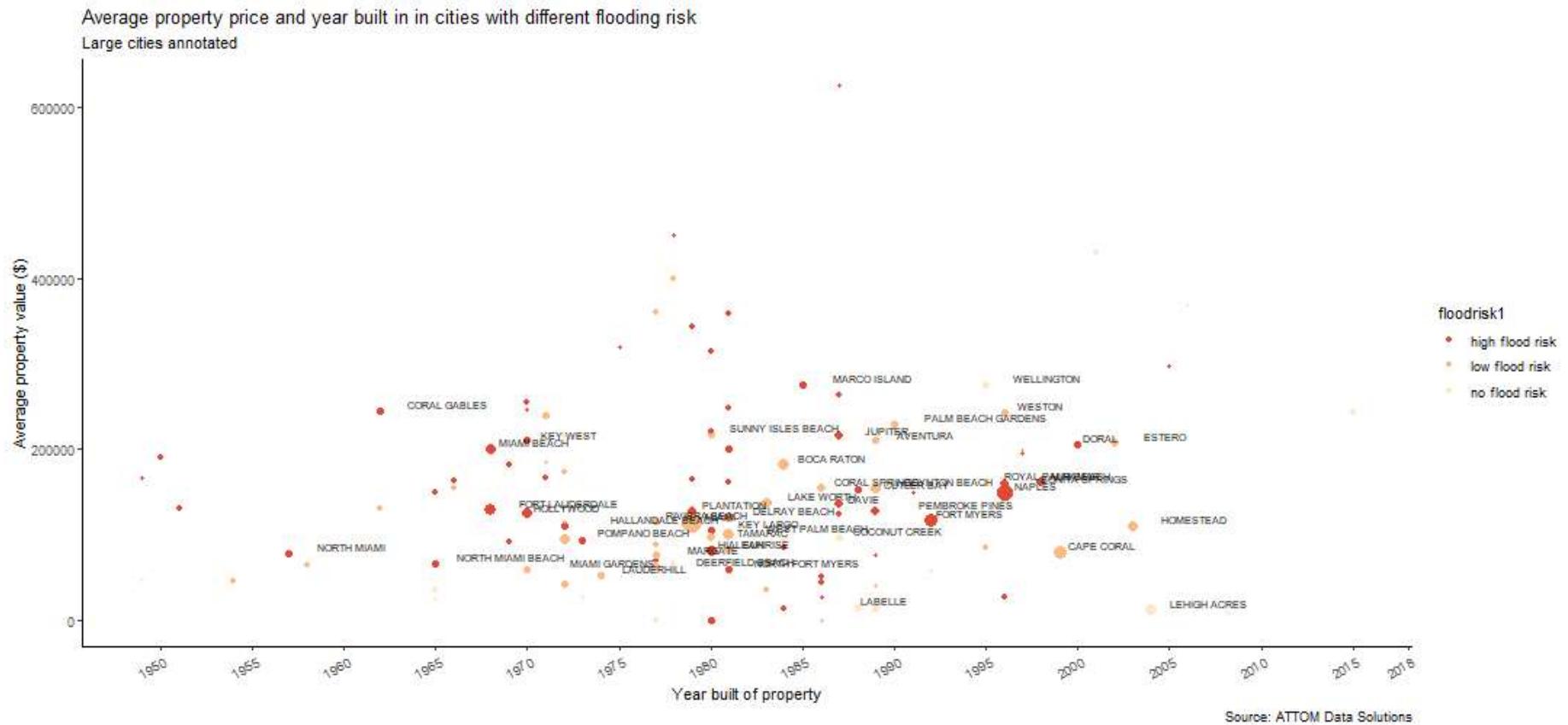
2. Exploratory analysis



2. Exploratory analysis



2. Exploratory analysis



3. Seperate modeling: one model for each city/county data

- For every city/county in South Florida, we fit an individual model to it seperately.
- Inspired by the "submarket" notion (the housing prices between neibourhoods, cities and counties vary a lot)
- Feature structure vary by county
- There are some variables that represent housing attributes that miss in data with different cases between counties.
- The process for county-level modeling is rather straightforward.

4. Feature selection with individual models

- Build function to automatically drop variables that are not important and are missing over 10% values, as we don't wish too many obs are omitted in tree model due to missing value.
- Apply the selection functions built just now to all subsamples
- Divided subsamples by county

```
summary(home_county_cleaned_dta)
```

```
##          Length Class      Mode
## Broward     79   data.frame  list
## Collier     72   data.frame  list
## Hendry      71   data.frame  list
## Lee         79   data.frame  list
## Miami-Dade 82   data.frame  list
## Monroe      78   data.frame  list
## Palm Beach  83   data.frame  list
```

4. Feature selection with individual models

- First 5 divided subsamples by city

```
summary(home_city_cleaned_dta)%>%head()
```

```
##          Length Class      Mode
## ALVA           70  data.frame list
## AVE MARIA      62  data.frame list
## AVENTURA       73  data.frame list
## BAL HARBOUR    68  data.frame list
## BANYAN VILLAGE 15  data.frame list
## BAY HARBOR ISLANDS 69  data.frame list
```

6. Build predictive models for individual subgroups

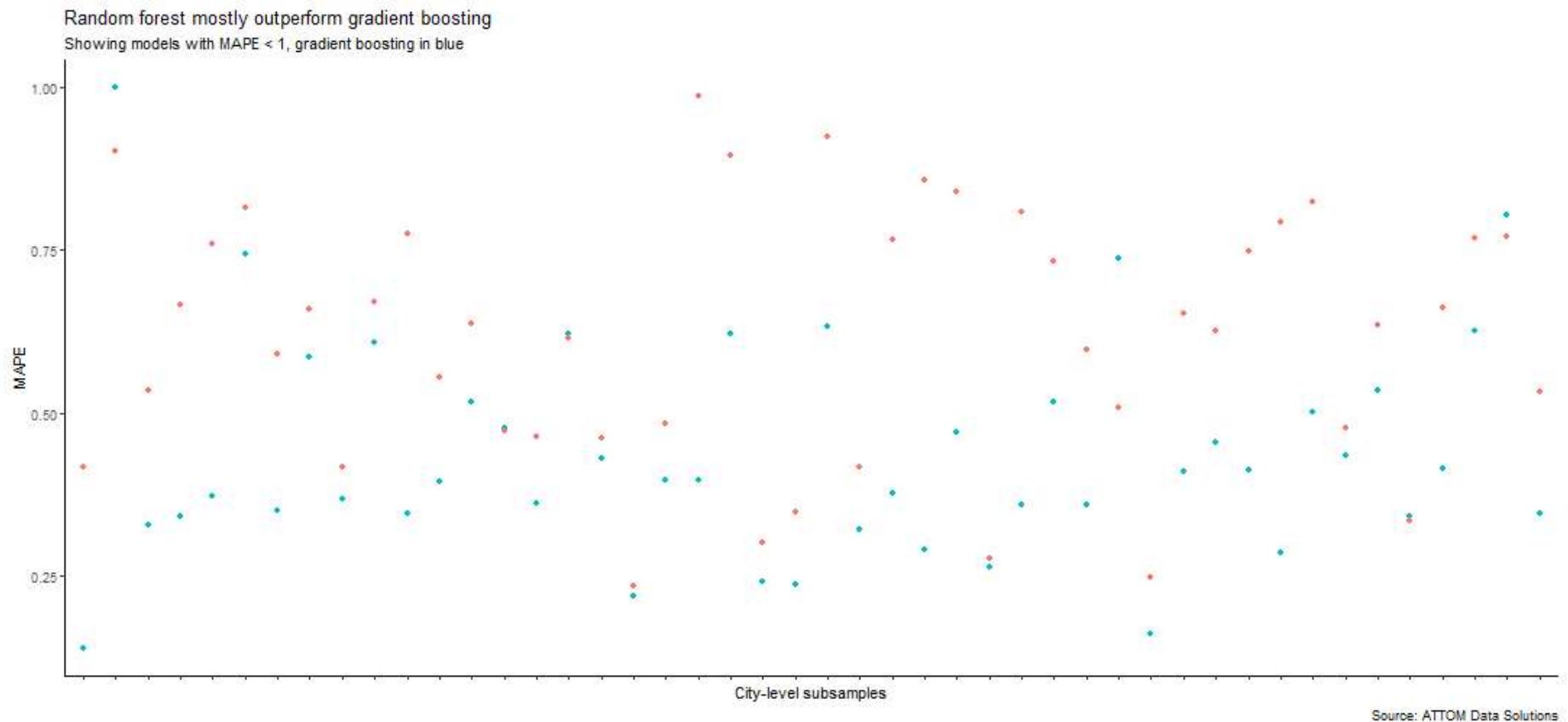
Final Preprocessing of individual model before modeling

- Develop functionality for final preprocessing before feeding data into model
- Here we have within each model the control of missing values and division of training/test set
- Therefore, the wrangling process was automated and run in each subsample

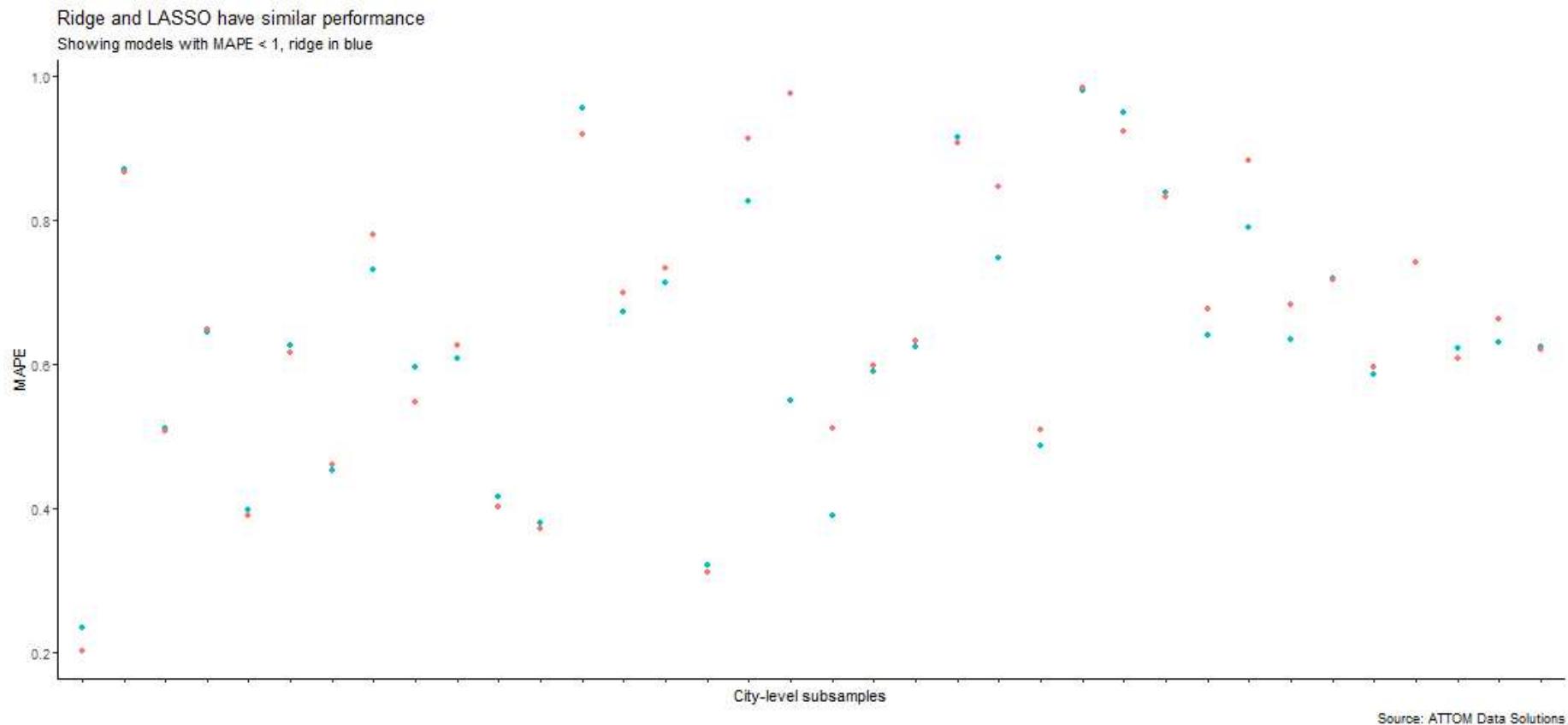
Modeling functionalities

- Regularization (Ridge and LASSO) as baseline, try also trees
- Random forest deal with data with large variations well
- Gradient boosting is good at learn specific models
- A model comparison with output from random forest and gradient boosting (also Ridge and LASSO)
- The function will parse validated key performance metrics and variables of importance of the models

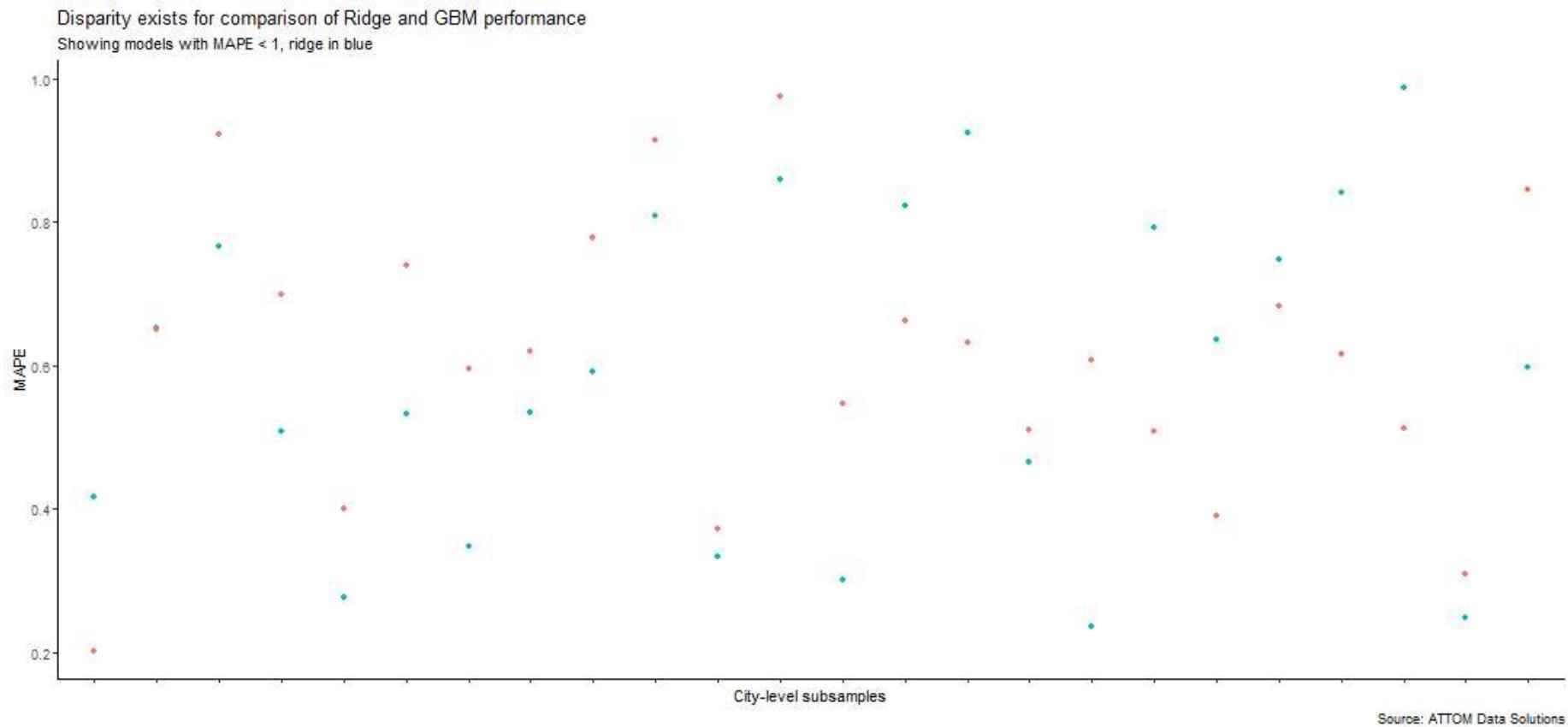
7 Visualize model performance



7 Visualize model performance



7 Visualize model performance



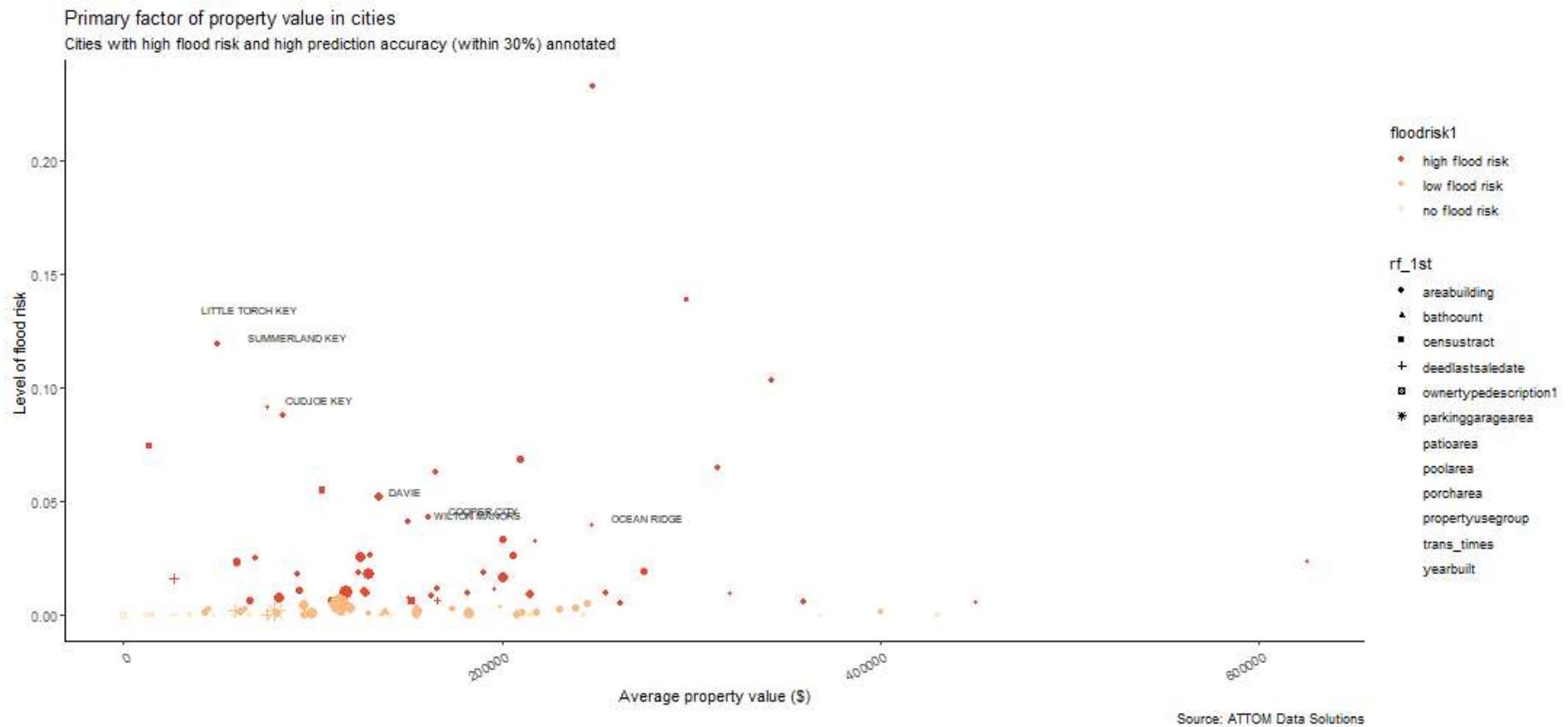
7 Visualize model performance

Show 10 ▾ entries

Search:

	subsample	averageprice	floodrisk	rf_MAPE	gb_MAPE
1	FL_Lee_ALVA	74.597892401553	0.016044061302682	7.63040335953268	10.4672964918621
2	FL_Collier_AVE MARIA	152.508728179551	0	0.139168804531982	0.417172647871918
3	FL_Miami-Dade_AVVENTURA	173.075149838635	0.00141786510026332	7.1957843595276	11.6980898259232
4	FL_Miami-Dade_BAL HARBOUR	295.389048991354	0.000250878073256397	1.03109217165162	1.04122521856096
5	FL_Hendry_BANYAN VILLAGE		0	1.53315834349759	1.43580133177156
6	FL_Miami-Dade_BAY HARBOR ISLANDS	135.947496884088	0.00353982300884956	0.412165184066806	0.748069804780997
7	FL_Palm Beach_BELLE GLADE	26.8028188286809	0	0.853550432457862	1.74058852655098

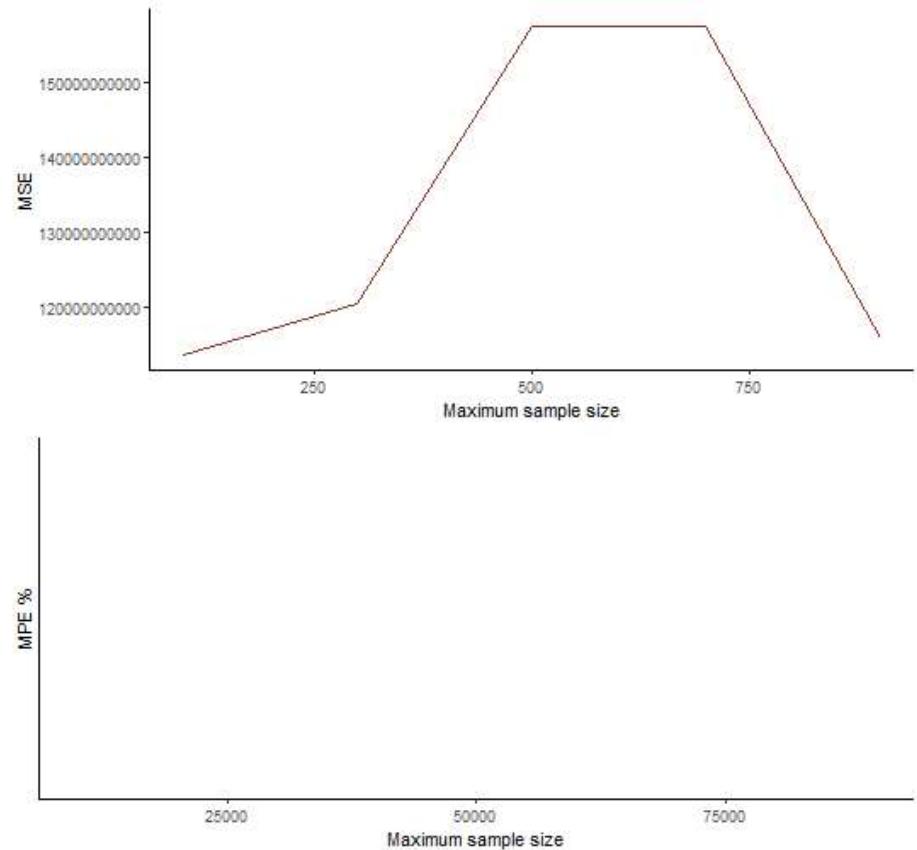
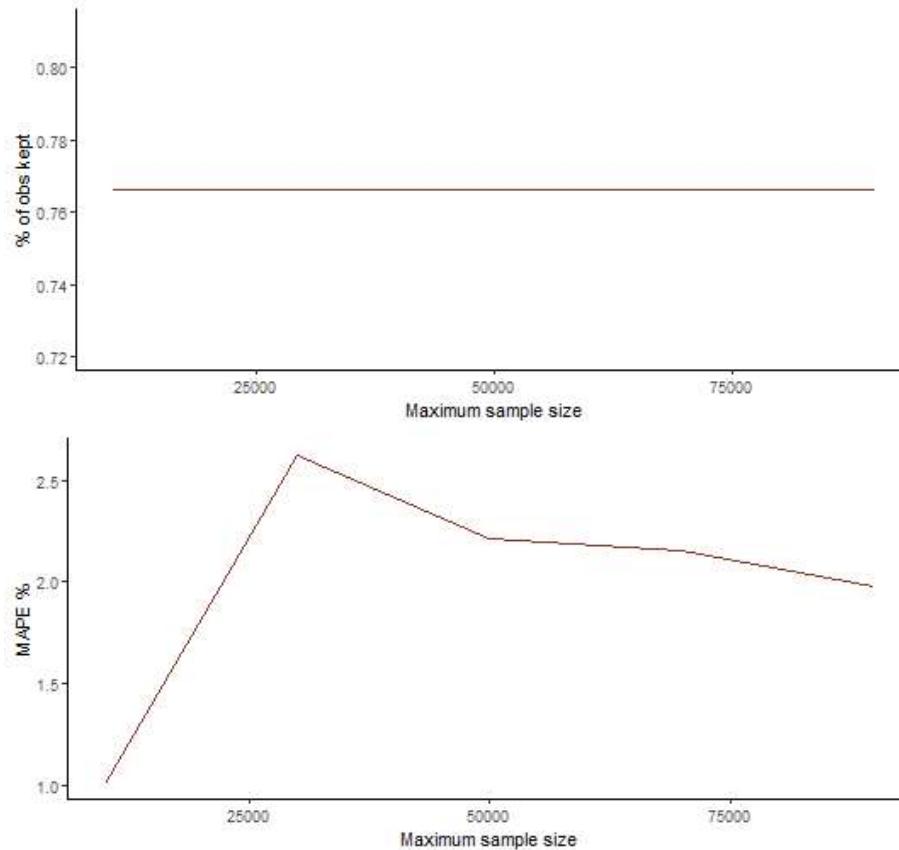
7.5 Visualize model performance



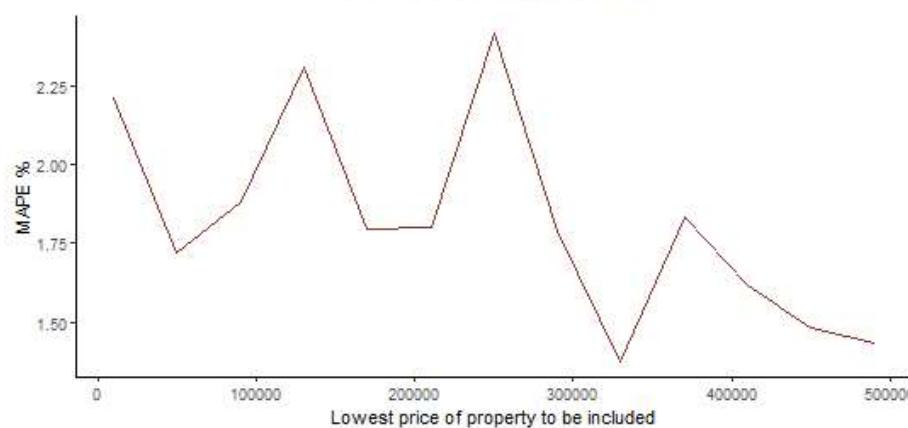
8. Personalize filtering parameters for random forest

- Filtering rule decides whether to include properties with extremely low/high price
- Crucial to model performance
- Build functionalities to show model performance depending on filtering parameters

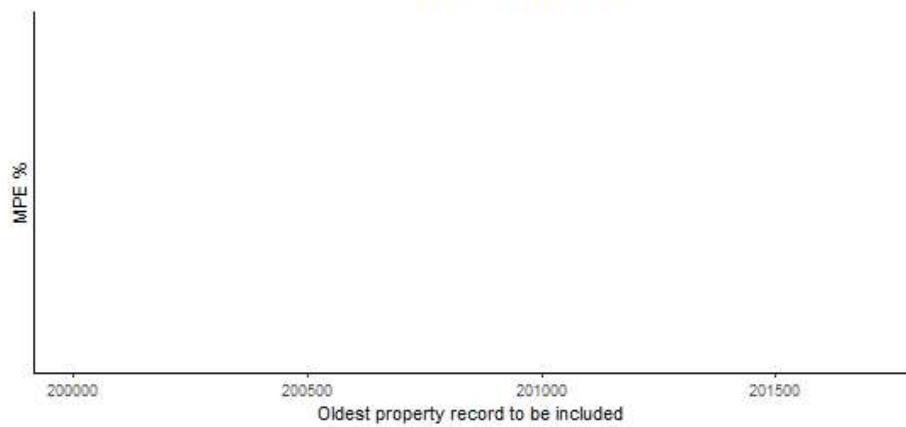
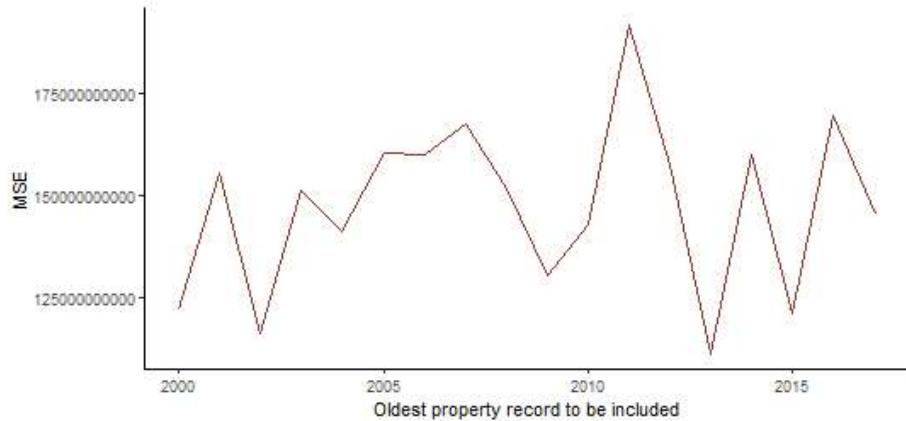
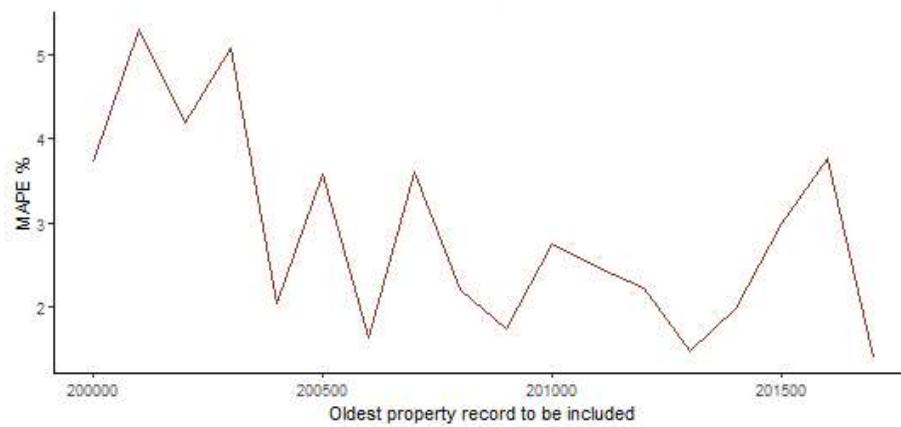
8. Personalize filtering parameters for random forest



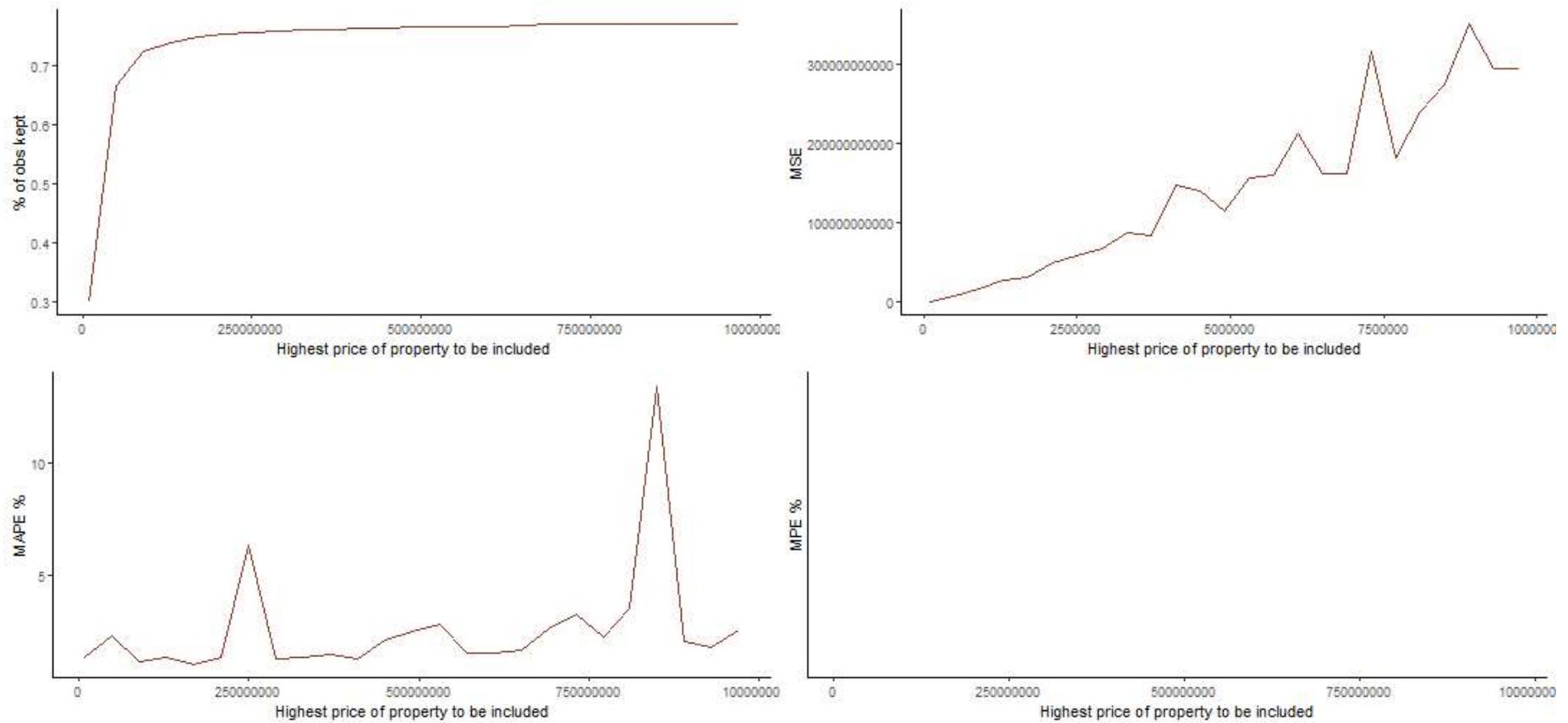
8. Personalize filtering parameters for random forest



8. Personalize filtering parameters for random forest



8. Personalize filtering parameters for random forest



9. Iteration through cities/counties and view cross-validated model performance

- Cities

Show 10 ▾ entries

Search:

	subsample	ridge_MPE	ridge_MAPE	LASSO_MPE	LASSO_MAPE
.1	FL_Collier_AVE MARIA	0.0203982135201874	0.125724668576108	0.0200407281535674	0.1250528702468
.116	FL_Broward_WESTON	0.0779097395957405	0.208957135411077	0.102455538102122	0.233765815936
.29	FL_Lee_ESTERO	0.00859072147337172	0.285669873444425	0.1243865835056	0.3146799064980
.96	FL_Lee_PUNTA GORDA	0.18024607522673	0.30964530747808	0.187859309819586	0.3216509541667
.2	FL_Miami-Dade_AVVENTURA	0.0357523896652135	0.338652338424877	0.0226081928721927	0.3345379769929
.41	FL_Palm Beach_HIGHLAND BEACH	0.112843021869129	0.302315305872082	0.181675875905562	0.3394697685121
.75	FL_Broward_MIRAMAR	0.215140738594086	0.39560114481853	0.166575928827484	0.3798334935117

9. Iteration through cities/counties and view cross-validated model performance

- Counties

Show 10 ▾ entries Search:

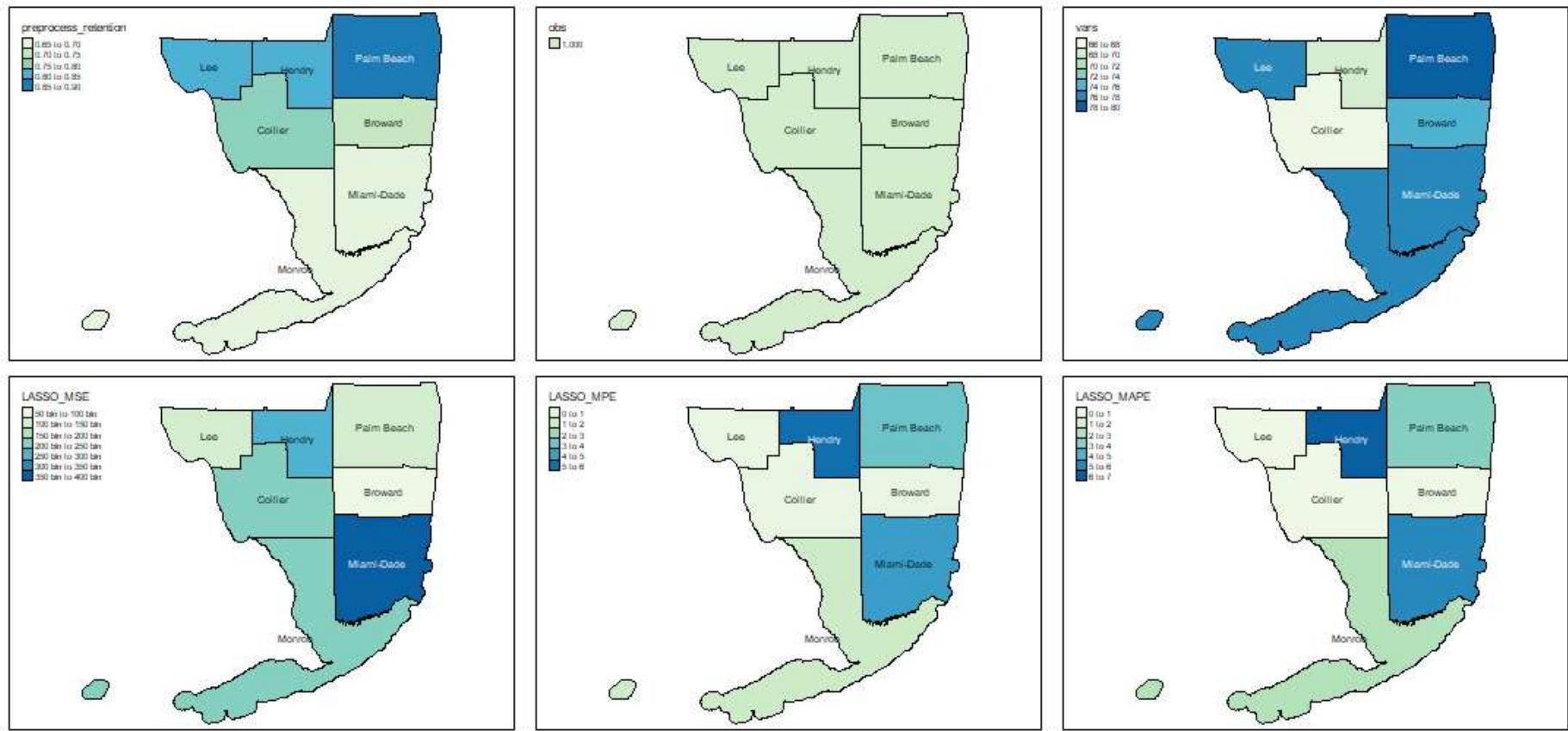
	subsample	ridge_MPE	ridge_MAPE	LASSO_MPE	LASSO_MAPE
..1	FL Collier	0.505128671588741	0.799125392180328	0.45493533077685	0.7900587809033
..3	FL Lee	0.658740045381184	0.833696739321913	0.59375836857304	0.80599816818285
.	FL Broward	0.709788388480453	0.899616779825107	0.689732506569334	0.902460244992195
.5	FL Monroe	2.06095987019894	2.26445715482043	1.90690092818672	2.12316217932807
.6	FL Palm Beach	2.33663845125845	2.48427375850428	3.57512033409527	3.76874391506075
..4	FL Miami-Dade	5.6224456642077	6.00813090214987	4.92200213878811	5.26792274384075
..2	FL Hendry	6.82166525108225	7.2131084650709	5.71983329622686	6.19049287168794

Showing 1 to 7 of 7 entries

Previous 1 Next

10. Map county-level model performance

We can see how each county performs on the map.



11. Brief Summary

- Random forest algorithm consistently outperformed by gradient boosting
- City-level models have varying performance and have potential for better model selection
- Ability to predict within 10% deviation in some cities
- Key filtering parameters optimized to improve prediction with reasonable constraint on sample size
- Transaction time and building area is 2 primary factor to determine property value

12. Next step

- Add census tract block medium price as control
- Can try tuning the filtering parameter for each model
- log scale price
- recode zip code to dummy (one-hot recoding)