

pca_summary

pmcdevitt

29 juillet 2017

Statement of Problem

Estimating market value of a home for sale has significant implications for all parties involved in the transaction : seller, buyer, agents, mortgage providers and even local taxing authorities. Getting it right can improve local economies. Inefficiencies associated to historical methods of value assessments create hesitation on the part of buyers and lenders, and potential loss of revenue for sellers and agents. Developing a model that considers all available factors and provides a transparent valuation that can be shared among all parties in the transaction can enable the participants to proceed with increased confidence, thus increasing the velocity of the local real estate market.

That is the purpose of this evaluation : use all available contributing factors for the residential real estate market in Ames, Iowa and create a predictive model to better estimate market valuation for future properties to be proposed for sale.

Data Available & Utilized

For this evaluation, there are seventy-nine explanatory variables available for exploitation, based on residential sales in the years 2006 through 2010, comprising approximately 1500 sales. The explanatory variables include traditional expected characteristics, such as : neighborhood, square footage, number of bedrooms, number of bathrooms, etc. and also several factors that are perhaps considered secondary or tertiary, but are included in the modeling to increase predictive capability. Some of these additional factors include : heating type, number of fireplaces, qualitative assessment of the kitchen condition.

Model Construction

In order to build the model, the following steps are taken :

- * read in the raw training data set provided,
 - * basic cleaning of the data, including removing significant outliers (for this purpose, more than 5 std deviations from mean)
 - * imputing values for features where none was provided (for this purpose, setting to mean value for numeric features, and creating a new factor level “None” for categorical features),
 - * plot and visually examine each feature in relation to $\log(\text{SalePrice})$...
 - + this provides a basis for removing some features from consideration based on inspection
 - + some features may have 1400 / 1460 within same category, thus not providing variability worthwhile including in a model
 - + some numerical features are sparsely populated, and the few values visually exhibit zero slope in relation to $\log(\text{SalePrice})$
 - + a new feature was created “saledate” from the “year sold” and “month sold” features. Upon visual examination, there was no obvious trend in the time series view for $\log(\text{SalePrice})$ s, so this was eventually discarded
 - + this visual examination then results in eliminating approximately 25 of the features from consideration in the model.
 - + (All of the plots are available for review at the referenced GitHub site : “homes_train_plots.pdf”)
-

Prior Models Considered and Results

In all cases, the basic data set consists of 52 predictor variables and the dependent output variable $\log(\text{SalePrice})$

Four different models are built :

* Stepwise - modeled in SAS proc glmselect

* Forward - modeled in SAS proc glmselect

* Backward - modeled in SAS proc glmselect

* CUSTOM - model based on the above three models ... just average the results of these models and see if this improves the Kaggle score.

Model	Adj R ²	CV Press	Kaggle Score
Forward	0.91	22.10	0.189
Backward	0.92	23.17	0.226
Stepwise	0.91	22.13	0.133
Custom	0.90		0.225
PCA_1			0.127

Prior Conclusion:

- For this effort, the 4 models each provide good predictive capability for estimating market valuation of residential real estate to be proposed for offering in the Ames, Iowa market.
- The stepwise model outperformed the other selection methods for the features chosen in this case. In addition, the stepwise selection also resulted in the least number of features (14) in comparison, also providing a simpler model. The table above shows the characteristics relative to model fit, along with the Kaggle scores when the model is applied to the test case data set. Clearly, the stepwise model selection is the preferred model among those evaluated.

The retained features in the final model include :

Feature	Feature	Feature	Feature	Feature
bsmtfinsf1	centralair	fireplaces	garagecars	kitchenqual
log(grlivarea)	log(lotarea)	mszoning	neighborhood	overallcond
overallqual	totalbsmtsf	yearbuilt	yearremodadd	

2. Principal Components (40%)

- address the assumptions of principal components:

The assumptions are the same as those used in regular multiple regression ¹:

* linearity,

* constant variance (no outliers), and

* independence.

* Since PC regression does not provide confidence limits, normality need not be assumed.

From the modeling that was completed previously, the linearity of the model and (mostly) constant variance were demonstrated.

- Linearity: The predicted vs. actual model from the prior linear regression modeling shows a very good fit, and strong linear response characteristics. In addition, the residuals plot do not show any obvious tendency towards non-linearity or increasing / decreasing variance across the range of values evaluated. The assumption of linearity appears respected for this model. When the final model is constructed again here, these will be re-assessed, but we can state that the available features used to construct a multi-variate linear regression model previously demonstrated satisfactory characteristics with regards to linearity.
- Equal variance: Similarly, from the prior modeling, the residuals plot and the studentized residuals plots both show that, within reasonable bounds, equal variances here are acceptable.
- Cook's D and the Leverage Plots showed a few influential points. The leverage plot, in particular, shows 2 points with relatively high leverage; however in this model, with 1400+ additional data points, the influence of these 2 points is not substantial.
- Relative to independence of the dependent responses (sale prices), we are assuming that with 1400+ single family residential home sales that the actions are reasonably independent. One can imagine a scenario where a developer acquires a large number of properties in a short period of time for a specific development project or a government project forces local project, but with 1400+ sales registered for this training set, and with no evidence to the contrary, we will make the assumption that independence is sufficiently respected with this data set.

produce the interpretations of the eigenvalues

Pearson Correlation Coefficients, N = 1460 Prob > r under H0: Rho=0														
	Prin1	Prin2	Prin3	Prin4	bsmtfinsf1	fireplaces	garagecars	log_griivarea	log_lotarea	overallcond	overallqual	totalbsmtsf	yearbuilt	yearremodd
Prin1	1.00000 1.0000	0.00000 1.0000	0.00000 1.0000	0.00000 1.0000	0.44226 1.0000	0.53105 1.0000	0.77555 1.0000	0.69525 1.0000	0.40059 1.0000	-0.23372 1.0000	0.84060 1.0000	0.73180 1.0000	0.69567 1.0000	0.60679 1.0000
Prin2	0.00000 1.0000	1.00000 1.0000	0.00000 1.0000	0.00000 1.0000	0.20267 1.0000	0.48189 1.0000	-0.09498 1.0000	0.29298 1.0000	0.63426 1.0000	0.30741 1.0000	-0.11972 1.0000	0.14575 1.0000	-0.56538 1.0000	-0.44583 1.0000
Prin3	0.00000 1.0000	0.00000 1.0000	1.00000 1.0000	-0.00000 1.0000	-0.39083 1.0000	0.05315 1.0000	0.02872 1.0000	0.21167 1.0000	-0.08450 1.0000	0.76442 1.0000	0.22398 1.0000	-0.24602 1.0000	-0.13709 1.0000	0.45292 1.0000
Prin4	0.00000 1.0000	0.00000 1.0000	-0.00000 1.0000	1.00000 1.0000	0.69819 1.0000	-0.12663 1.0000	-0.13617 1.0000	-0.37352 1.0000	-0.11067 1.0000	0.46163 1.0000	-0.05574 1.0000	0.25862 1.0000	0.03541 1.0000	0.17955 1.0000
bsmtfinsf1	0.44226 1.0000	0.20267 1.0000	-0.39083 1.0000	0.69819 1.0000	1.00000 1.0000	0.24210 1.0000	0.23221 1.0000	0.11691 1.0000	0.19829 1.0000	-0.04420 1.0000	0.22801 1.0000	0.47513 1.0000	0.25147 1.0000	0.12541 1.0000
fireplaces	0.53105 1.0000	0.48189 1.0000	0.05315 1.0000	-0.12663 1.0000	0.24210 1.0000	1.00000 1.0000	0.30079 1.0000	0.45622 1.0000	0.29770 1.0000	-0.02382 1.0000	0.39677 1.0000	0.32543 1.0000	0.14772 1.0000	0.11258 1.0000
garagecars	0.77555 1.0000	-0.09498 1.0000	0.02872 1.0000	-0.13617 1.0000	0.23221 1.0000	0.30079 1.0000	1.00000 1.0000	0.49306 1.0000	0.28113 1.0000	-0.18576 1.0000	0.60067 1.0000	0.45323 1.0000	0.53785 1.0000	0.42062 1.0000
log_griivarea	0.69525 1.0000	0.29298 1.0000	0.21167 1.0000	-0.37352 1.0000	0.11691 1.0000	0.45622 1.0000	0.49306 1.0000	1.00000 1.0000	0.37469 1.0000	-0.10084 1.0000	0.59835 1.0000	0.40442 1.0000	0.23015 1.0000	0.30911 1.0000
log_lotarea	0.40059 1.0000	0.63426 1.0000	-0.08450 1.0000	-0.11067 1.0000	0.19829 1.0000	0.29770 1.0000	0.28113 1.0000	0.37469 1.0000	1.00000 1.0000	-0.00637 1.0000	0.17772 1.0000	0.33028 1.0000	0.02314 1.0000	0.02876 1.0000
overallcond	-0.23372 1.0000	0.30741 1.0000	0.76442 1.0000	0.46163 1.0000	-0.04420 1.0000	-0.02382 1.0000	-0.18576 1.0000	-0.10084 1.0000	-0.00637 1.0000	1.00000 1.0000	-0.09193 1.0000	-0.17517 1.0000	-0.37598 1.0000	0.07374 1.0000
overallqual	0.84060 1.0000	-0.11972 1.0000	0.22398 1.0000	-0.05574 1.0000	0.22801 1.0000	0.39677 1.0000	0.60067 1.0000	0.59835 1.0000	0.17772 1.0000	-0.09193 1.0000	1.00000 1.0000	0.54071 1.0000	0.57232 1.0000	0.55068 1.0000
totalbsmtsf	0.73180 1.0000	0.14575 1.0000	-0.24602 1.0000	0.25862 1.0000	0.47513 1.0000	0.32543 1.0000	0.45323 1.0000	0.40442 1.0000	0.33028 1.0000	-0.17517 1.0000	0.54071 1.0000	1.00000 1.0000	0.40050 1.0000	0.29600 1.0000
yearbuilt	0.69567 1.0000	-0.56538 1.0000	-0.13709 1.0000	0.03541 1.0000	0.25147 1.0000	0.14772 1.0000	0.53785 1.0000	0.23015 1.0000	0.02314 1.0000	-0.37598 1.0000	0.57232 1.0000	0.40050 1.0000	1.00000 1.0000	0.59285 1.0000
yearremodd	0.60679 1.0000	-0.44583 1.0000	0.45292 1.0000	0.17955 1.0000	0.12541 1.0000	0.11258 1.0000	0.42062 1.0000	0.30911 1.0000	0.02876 1.0000	0.07374 1.0000	0.55068 1.0000	0.29600 1.0000	0.59285 1.0000	1.00000 1.0000

Figure 1: Explanation of eigenvectors

The table above provides a basis to interpret the principal components of one of the sets of principal components evaluated in this model. As an example, this table includes a PCA set for a model with just 10

of the features included and the Pearson's correlation coefficient for the 1st four principal components. With more features, the explanation becomes a bit more cumbersome, so we will provide details for this case, and the method of interpretation for larger set models follows the same method. The following observations can be made about this set of principal components :

- **PC_1** : Strongest contributions from Overall_Quality, Number_of_Cars_Garage, and Total_Basement_Sq_Ft. In a sense, these are the high value characteristics of the home price market.
- **PC_2** : Strong (negative) contributors from Year_Built and Year_Remodelled. This principal component is essentially the age feature of the home value - the newer the home and the more recently remodelled adds incremental value accounted for by PC_2.
- **PC_3** : Strongest contribution from Overall_Condition. We observed that Overall_Quality is not well correlated with Overall_Condition. This PC_3 is adding the value associated to whatever is the difference between _Condition and _Quality, in contrast to the _Quality value that is associated to PC_1.
- **PC_4** : Strong contribution from Basement_Finish_Sq-Ft. We observed that basement and home size increase home sale prices in this market. This PC is accounting for a contribution from that type home feature.
- For **PC_5** through **PC_10**, the relative weights of the individual independent variables are depicted in the below bar graph, with the relative correlation of each feature to that principal component indicated by the height of the associated color in the bar graph. In this case, we can see that among the ten features included they form natural sub-groups with two features each :
 - age related (year built, year remodeled)
 - amenities (garage size, fireplace)
 - basement features
 - size, area
 - quality and condition indicators.

The different hues of the same color identify features that can be logically associated (e.g., year_remodelled and year_built can both be associated to the idea of *age*).

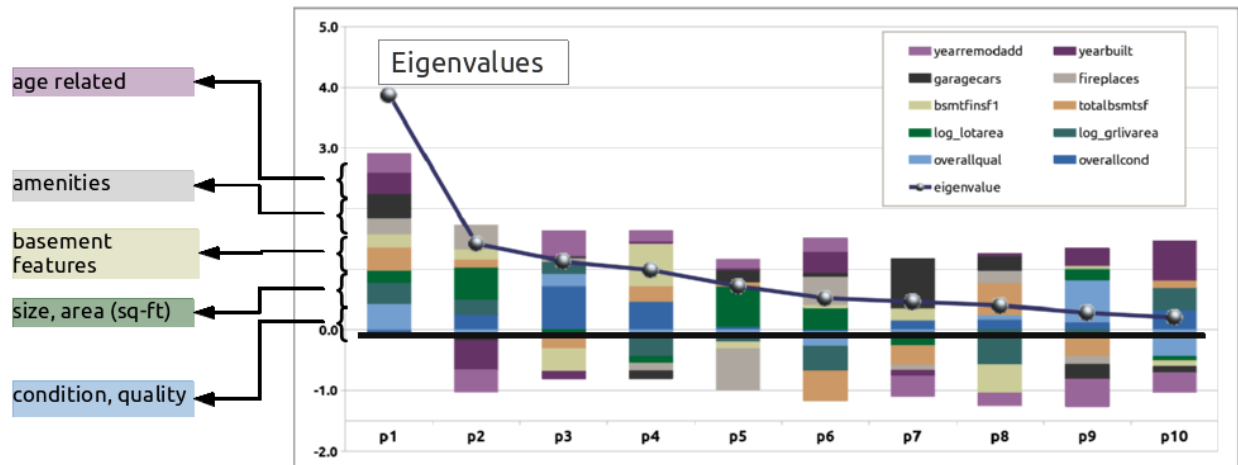


Figure 2: Visual view of principal components correlation with each feature

- screeplots,

The scree plot associated to these principal components is shown below. This relationship indicates relative weight of the eigenvalues for the the principal components. It can be observed that 1st principal component has value of 3.8; the second principal component has value of 1.4 and each succeeding value progressively less. The sum of these eigenvalues is 10, corresponding to the number of principal components. The plot on the right shows the relative increase in variance explained (among the independent variables) for the successive addition of each principal component. The first five principal components can account for 80% of the variation in the data and 7 components can account for 90% of the variation.

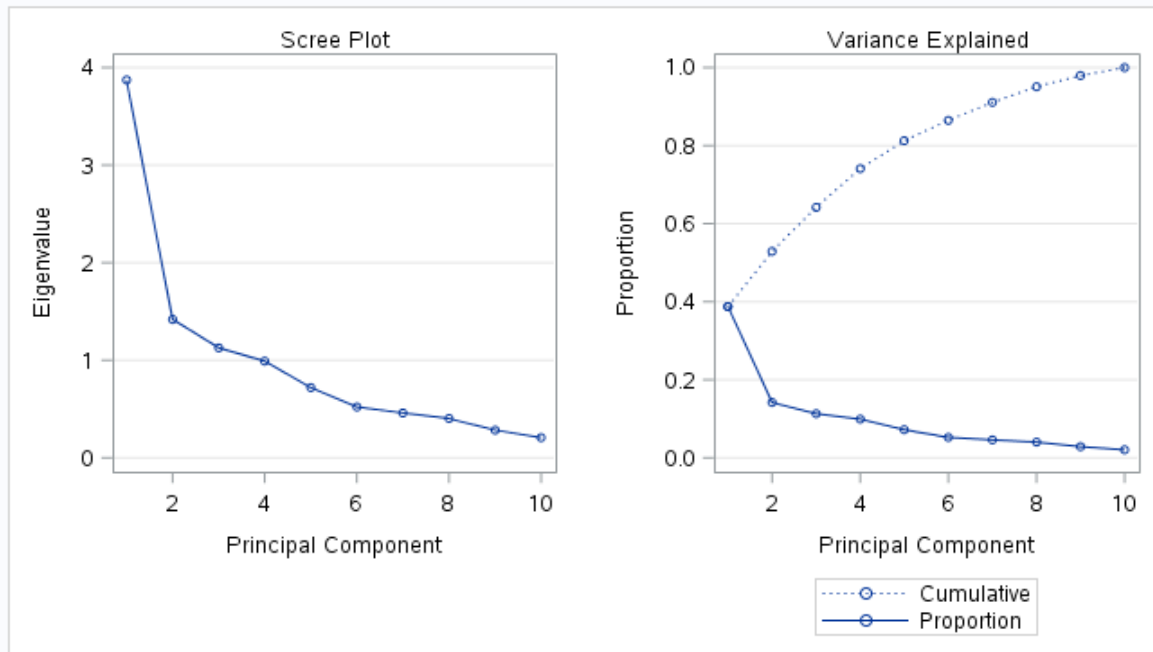


Figure 3: Scree Plot

This needs to be done →

Furthermore, the investigation of the use of PCA in this project should be used in conjunction with the regression techniques you have already been using (adding categorical variables, investigating OLS, LASSO and coefficients, cross validation, variable selection, etc.)

The team with the best Kaggle score this time around will again earn an extra 3 points. Simply provide the same table you filled out before with your new scores as well as the code you used to generate it. As before, only techniques we have learned in this class can be used.

Conclusion/Discussion Required

The conclusion should reprise the questions and conclusions of the introduction,

Submission and Description	Public Score	Use for Final Score
kaggle_submit_pca_step1.csv 3 minutes ago by bici.sancta add submission details	0.12730	<input type="checkbox"/>
kaggle_submit_nca2_w_imp_discr_2017.07.27.csv	0.15735	<input type="checkbox"/>

Figure 4:

perhaps augmented by some additional observations or details gleaned from the analysis section. New questions, future work, etc., can also be raised here.

Appendix

SAS Code for principal components regression analysis

```
/* ----- */
/* ... Principal Components Regression ... */

/* kaggle home prices data set
/* patrick mcdevitt
/* 29-jul-2017
/* ----- */

/* ----- */
/* ... start with clean memory ... */
/* ----- */

proc datasets lib=work kill nolist memtype=data;
quit;

/* ----- */
/* ... read in training data set ... */
/* ----- */

FILENAME REFFILE '/folders/myfolders/stats_ii/training_set_cleaned.csv';

PROC IMPORT DATAFILE = REFFILE
    DBMS = CSV
    OUT = home_prices;
    GETNAMES = yes;
RUN;

/* ----- */
/* ... read in test data set ... */
/* ----- */

filename reffile '/folders/myfolders/stats_ii/test_set_cleaned.csv';

proc import datafile = REFFILE
    DBMS = csv
    OUT = test_set;
    GETNAMES = yes;
RUN;

/* ----- */
/* ... combine train and test data sets ... */
/* ----- */

data train_test;
    set home_prices test_set;
run;

title 'PCA for all (selected) independent numeric variables in training set';
proc princomp data = train_test out = pc_home_prices;
var bsmtfinsf1
    bsmtfullbath
```



```

fireplaces
fullbath
garagearea
garagecars
halfbath
log_grlivarea
log_lotarea
log_lotfrontage
overallcond
overallqual
totalbsmtsf
totrmsabvgrd
x1stflrsf
x2ndflrsf
yearbuilt
yearremodadd;
run;

proc print data = pc_home_prices;
run;

/*****
      model with principal components + categoricals
*****/

title 'Regression (stepwise) with full PC set + categorical variables';
proc glmselect data = pc_home_prices plots = (criteria) seed = 3;
class bsmtfintype1
    bsmtqual
    centralair
    electrical
    exterior1st
    exterior2nd
    exterqual
    fireplacequ
    foundation
    garagefinish
    garagetype
    heatingqc
    housestyle
    kitchenqual
    lotshape
    masvnrtype
    mszoning
    neighborhood
    saletype;
model log_saleprice =
    /*      continuous variables      */
    prin1-prin18
    /*      categorical variables      */
    bsmtfintype1
    bsmtqual
    centralair

```

```

electrical
exterior1st
exterior2nd
exterqual
fireplacequ
foundation
garagefinish
garagetype
heatingqc
housestyle
/* kitchenqual */
lotshape
masvnrtype
/* mszoning */
neighborhood
saletype / selection = stepwise(choose = CV select = cv stop = aic);
output out = result p = Predict;
run;

/* ----- */
/* create kaggle submission file */
/* two columns with appropriate labels. */
/* ----- */

proc means data = result Min Max;
run;

proc means data = result noprint;
var Predict;
output out = means mean(Predict) = mean_predict;
run;

/* ----- */
/* in case any missing values in predicted set, */
/* impute with mean of predicted sale prices */
/* ----- */

data kaggle_submit;
set result;
SalePrice = exp(Predict);
if Predict = . then SalePrice = exp(12.024);
keep id SalePrice;
where id > 1460;
run;

proc export data = kaggle_submit replace
outfile = '/folders/myfolders/stats_ii/kaggle_submit_pca_step.1.csv'
dbms = csv;
run;

/* ----- */
/* ... end_of_file */
/* ----- */

```