

# 1 Data Prep, EDA, and Theory development

## 1.1 Variable Selection & Explanation

For the purpose of analyzing the determinants of prices of home sales in the US, the following variables were included in the analysis (Table 1):

- Sale Price (DV)
- Total Living Space (IV) & Years Since Remodeling (at point of Sale) (IV)
- Confounders & Controls: Quality, Lot Area, Condition

As can be seen in Table 1, a total of 1,460 house sales were recorded between 2006 and 2010 for the district of Ames, Iowa (USA). As can be observed in Table 1, the mean sale price of a house was (in 1000s) \$180.921 (SD = 79.443). Combined with the range [34,900, 755,000] a positive skewness was to be expected (skew = 1.881), considering that the outcome variable is of financial nature. The Total Living Area displays a mean of 2,572.89 square feet (SD = 823.598) in addition to a large range of values [334, 11,752]. Years Since remodeling (at time of sale) shows that the average property did not undergo renovations for 22.95 (23) years (SD = 20.950). Contrary to expectation, this variable distributes reasonably equally across the range, stopping out at a maximum of 60 years (See Figure 1B - quantiles). Furthermore, the variable Quality (and Condition) represents a rating from 1 to 10, similar to a Likert Scale. Quality has to be considered a categorical variable in this case i.e. because the distances between each rating level are not constant and the distribution is skewed (**SEE SUPPLEMENTARY APPENDIX REGRESSION AND PICTURE**). However, while strictly speaking there we cannot consider Quality a numerical variable, for the purpose of certain examples at a later point, this variable will be considered as both a categorical and numerical variable (no inference will be made is used as numerical). Finally, Lot Area, will be used as a confounder in the regressions to control for the association larger lot sizes creating larger houses (**AS AN INTERACTION**).

Table 1: DESCRIPTIVE STATISTICS OF NUMERIC VARIABLES

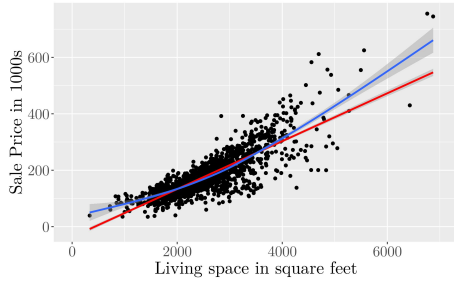
Statistic	Mean	St. Dev.	Min	Pctl(25)	Median	Pctl(75)	Max
SalePrice	180.921	79.443	34.900	129.975	163.000	214.000	755.000
Lot Area	10,516.830	9,981.265	1,300	7,553.5	9,478.5	11,601.5	215,245
Quality	6.099	1.383	1	5	6	7	10
Condition	5.575	1.113	1	5	5	6	9
Total Living Space	2,572.893	823.598	334	2,014	2,479	3,008.5	11,752
Years Since Remodeling	22.950	20.641	0	4	14	41	60

Notes: N = 1460

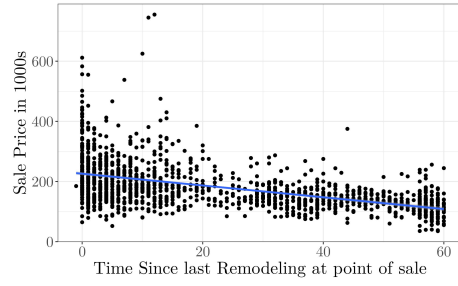
Beyond this table, there are multiple categorical variables (**see Appendix**), such as Zoning and Year of Sale. The original (MS)Zoning variable contains seven categories, of which five contain data; these zones correspond to the administrative classification of the ground on which the properties are constructed (Commercial, Floating Village, Low-Density, Moderate-Density, High-Density contain data; Residential Low Density Park, Agricultural, Industrial do not contain records). For the purpose of this analysis, this number was reduced to four categories based on the similar behaviour of Moderate and High Density properties (**SEE SUPPLEMENTARY APPENDIX PLOT**) in the data as Ames, Iowa, represents the stereotypical picture of a mid-western town in the US, thereby displaying fewer densely populated areas. Thus, the main question of this analysis section focuses on the difference between Low and higher density properties.<sup>1</sup> In addition, year of sale will be used to control for the effect of the 2008/2009 housing crisis.

Finally, the plots are to be considered in the context of the hypothesis in the next section.

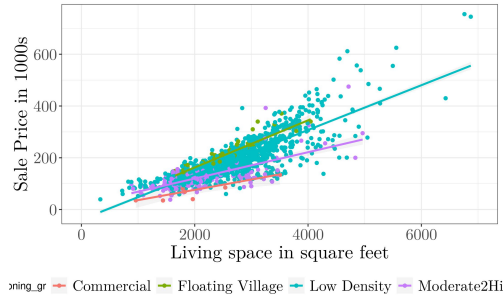
<sup>1</sup>Floating Village and Commercial behave too differently to be merged



(a) Positive association of total Living Space & Sale Price & optimal line; range [0, 7000]



(b) Positive association of total Living Space & Sale Price



(c) Positive association of total Living Space & Sale Price; range [0, 7000]

Figure 1: Three Hypothesis Graphs displaying their repective association with the outcome variable

## 2 Theoretical model and OLS assumptions

### 2.1 Hypotheses

Based on the plots generated during the EDA, a mini theory was created to explain the variation in the sales price of properties (Figure 2).

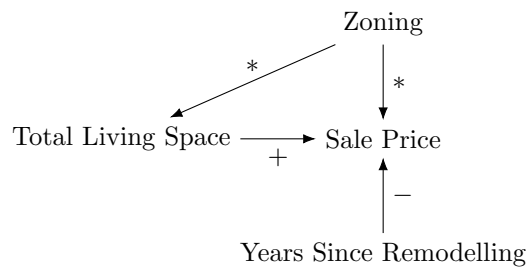


FIGURE 2: Causal relationship Scheme

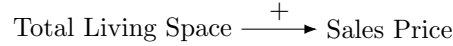
Based on this model, the following three hypothesis were created.

#### 2.1.1 Hypotheses 1

**Figure 1A** displays a potential direct positive association between Total Living Space (IV) and Sale Price (DV), including an optimal fit, showing a small. Thus, one expects that larger houses have a higher sale price. Consequently, we assume that:

**Hypothesis 1 (H1):** *Total living space (IV) has a direct postive association with Sales Price (DV)*

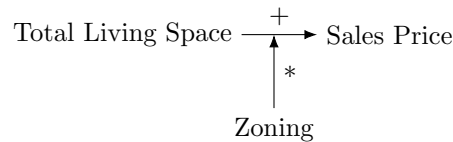
**Figure 3: Hypothesis 1**



Subsequently, when taking the Zoning (MSZoning or MSZoning\_grouped - IV) into account to reflect the administrative borders of the Ames districts, larger houses in more densely populated areas of the city appear to have a lower price when compared to houses of same size in less densely populated areas as can be seen in **Plot XXXX2**. This suggests a separation between "downtown less affluent areas" and "suburban affluent areas"<sup>2</sup> and concludes in the hypothesis that

**Hypothesis 2 (H2):** *Zoning moderates (MIV) the direct postive association of Total Living Space (IV) and Sale Price (DV). The association of Space and Sale Price is proposed to be weaker for more densly populated areas than for more rural areas.*

**Figure 4: Hypothesis 2**



As can be observed on Figure 1c, older the houses 8not renovated (IV), are associated with lower sales prices (DV). Thus, the final hypothesis corresponds to:<sup>3</sup>

**Hypothesis 3 (H3):** *Years Since Remodeling (M-IV) or Construction has an amplyfying effect on the direct postive effect of Quality (IV) and Sale Price (DV)*

**Figure 5: Hypothesis 3**



**IMPORTANT: DO YOU INCLUDE A QUADRATIC TERM IN THE POPULATION REGRESSION MODEL???? NO! BUt make more equations explaining each quadartic term and represent the interaction effects etc adn what not**

$$(2 - 1) \text{ SalePrice} = \alpha + \beta_1 \text{TotalLivingSpace} + \beta_2 \text{Zoning} - \beta_3 \text{YearsSinceRemodeling} + \epsilon$$

## 2.2 Assumptions

A1: The linearity of model parameters and error term assumption suggests that a) the functional form of the underlying population regression model is linear and additive. Thus, this assumption is generally assumed to hold, i.a. not having introduced quadratic or polynomial (by parameter)

<sup>2</sup>As an extention we will test whether the groups of Neighborhoods generally stay in the same zoning category; if Neighborhoods and Zoning are not related (so eg 50 % of one neighborhood is in rural zone while the other part is in moderately populated zone) then we have a problem that this woudl induce a bias. Otherwise, we can just proceed.

<sup>3</sup>It is notable, that more complex analysis will be considered in part 4 & 5, such as the moderating effect of Years Since remodeling on the association of Quality and Sale price

terms into the population regression equation (2-1). However as Figure 1b shows, there might be a quadratic relationship present between Total Living Area and Sale Price. The reason why this might be the case may be that with increase values for  $x$  ( $[X = x]$ ), the effect of  $X$  on  $Y$  increases, which would be represented in later regression models by an additional quadratic term in the regression equation. However, as will be shown in later parts, this can be (partially) remedied.

A2: Full rank assumes that no independent variable can be a linear function of other independent variables; e.g. a 2-dimensional sphere might be compressed to a line, figuratively speaking, meaning that there is no optimal (or none at all) solution for the parameters in question. This first part of the assumption might easily be violated when not dropping a "comparative" category of a categorical variable. However, violations of full rank might also come in the form of multicollinearity. Multicollinearity is the "almost" violation of the full rank violation as a given variable may, for instance, be highly correlated with another given variable. This way the sphere is almost compressed to a single line, which leads to a multitude of problems regarding inference of the model: primarily, the standard errors of the model become unstable to the inclusion of other explanatory or control variables. Multicollinearity is almost always present with observational data (actually making regression interesting in the first place); however, the extent is more relevant. For example might be if we included the Total Living Space and the TotalNumberOfRooms, both of which will be strongly correlated. Further tests will be conducted to probe for this assumption violation.

#### **in what direction is the estimator biased when this variable is left out**

A3: This assumption is referred to as mean independence (or exogeneity); assuming no correlation/systematic association between independent variables and residuals (or conditional mean error is dependent on independent variables). Thus this assumption might induce a bias/ or reduce consistency of the estimator and is, thus, critical to the model itself. A common way this variable is violated is via unobserved confounders or omitted variables. If a given variable is left out, part of the error term can be explained by a given variable which suffers under the influence of the confounder. An example in this case will be the missing information regarding crime by the area a property is located. It is obvious that higher crime levels will reduce the value in a given area. Thus, given a certain Neighborhood in question, part of the variance that would have been explained by the Crime variable is now falsely attributed to Neighborhood, biasing the estimate. Omitted variables bias the estimate of the population coefficient either up or down, thereby reducing the consistency of the estimate.

A4: The error term has a constant variance for each observation expected! — heteroscedasticity

The assumption of homoscedasticity assumes constant variance of each observation. However, if the variance of the estimate is varying we face heteroscedastic variances. While this is not necessarily problematic regarding the estimated parameter (coefficient), heteroscedasticity impacts the standard errors of the estimate, leading to problems regarding inference; heteroscedasticity can lead to both type I and II error. In this research, this violation might occur if the eg. sales prices vary stronger for large houses than for small houses. The resulting (averaged) standard error would neither be representative for small and large house standard errors and, thus, inference itself.

#### **NOTE AUTOCORRELATION AND YEARS SINCE BUILT**

Generally, as the sample size increases, this assumption becomes less important for the estimate itself, considering that OLS is a consistent estimator; which is the case here ( $N$  is quite large). However, heteroscedasticity may lead to problems regarding inference, as the standard error of the estimate may still be biased. As such, solutions will be later implemented to control for such violations.

A5: Data generation:  $x_i$  can be random or fixed; the data was collected for predictive purposes so we might not be able to verify this. As the data was not generated via a random experiment, none of the variables at hand are **random variables**. Due to the data originating from sales in a town in Iowa (Ames), we have to assume that the data at hand is "fixed". Thus, inference regarding a wider population cannot be made from this data, as its fixed nature only applies to small towns in the mid-western USA. However, we can assume that the fixed variables collected are measured without error; particularly as the sample collected can be somewhat representative of the wider population of small towns in the mid-western USA.

generally, due to the nature of the data, being a

#### **in order to**

A6: This assumption regards inference; if the residuals do not follow a standard normal distribution

bution, this might result in incorrect decisions as the distribution might for instance be "fatter" in the tails, resulting in potentially Type I or Type II errors. This might be the case if for instance we do not have a lot of observations for e.g. subcategories. Subsequently, the resulting inference might be biased. However, as the sample size increases, most distributions approach the normal form. However, in order to answer the question, this assumption is violated via i.a. large amounts of outliers. As can be seen in Figure 1a the range goes from 0 to 7000 square feet, thus, ommits two extreme outliers.

### 3 OLS regression and model fit

#### 3.1 Normal Regression

**NOTE: This regression table shifts from including interaction effects, scaling, etc from one model to the other; I chose to The choice to not include incremental models was done to keep the analysis concise. In the appendix one can see that the  $Adj - R$  increases also only due to the inclusion of further control variables without interaction terms. Nonetheless, also the interaction terms improve on the model fit (with and without control variables). Finally, Model (3) & (4) include all interaction effects and scaled variables (such as LN SalePrice) as no model shows a lower issues in this regard. Thus, only the simple multi-variate model and an extensive model was displayed.**

Table 2 reports a simplified (1) and complex model (3), containing all corresponding interaction effects and control variables <sup>4</sup>

Table 2 reportes the (unstandardized) regression results for the Model (1) excluding any interaction effects & and including those interaction effects in Model (3).<sup>5</sup> The corresponding standardized coefficients can be found in Std.Model (2) & (3).<sup>6</sup>

he study of Sales Price of property in Ames, Iowa, has shown that the Total Living Area has a significant positive effect on the Sales Price

Table 3 reports the regression results of for the four successively extensive models starting from only the independent variable from Figure 2, then adding interaction terms in model (2), control variables in model (3) without interactions, and a complete model (4) including interaction terms, which will be discussed below. All four models have been reported to demonstrate that for the main independent variables (see Figure 2) the coefficients and Standard Errors are stable to the inclusion of control variables and interaction terms.<sup>7</sup> Considering model (4) as the Complete model Specifically, this research shows that the association Total Living Space and Sale Price is positive and statistically significant ( $\hat{\beta} = 0.040, p = 0.01$ )

he study of Sales Price of property in Ames, Iowa, has shown that the Total Living Area has a significant positive effect on the Sales Price

Overall, the inclusion of the interaction of Zone (M-IV) and Total Living Area (IV), in addition to introducing a quadratic term of Total Living Area (IV) improves the model fit slightly from

#### 3.2 Standardized Regression Coefficients and Effect Size

### 4 Diagnostic checking

#### 4.1 A1

Linearity

<sup>4</sup>It is notable, that more complex analysis will be considered in part 4 & 5, such as the moderating effect of Years Since remodeling on the association of Quality and Sale price

<sup>5</sup>It is notable, that more complex analysis will be considered in part 4 & 5, such as the moderating effect of Years Since remodeling on the association of Quality and Sale price

<sup>6</sup>It may be noted that the Sale Price will be log scaled in Part 4; the interaction term of Total Living Space with itself was already included in this section as I assume that part 3 is mainly about the independent variables demonstrating once the normal interpretation; as such a complete model will be provided in part 4, containing the log scaled version of the independent

<sup>7</sup>Please note that log scaling of the dependent variable was not done on purpose at this point due to this task falling into part 4 of this assignment

Table 2: MERGED STANDARDIZED REGRESSON RESULTS CONTAINING BOTH STANDARDIZED AND UNSTANDARDIZED COEFFICIENTS FINAL VERSION 3 MODELS COMPLETE YAY

	<i>Dependent variable:</i>					
	SalePrice		SalePrice		ln_SalePrice	
	Model (1)	Std. Coef.	Model (2)	Std. Coef.	Model (3)	Std. Coef.
Constant	15.229 (25.697)		24.550 (24.541)		3.584*** (0.114)	
Living Space	0.035*** (0.002)	0.362***	0.038*** (0.005)	0.393***	0.0003*** (0.00002)	0.684***
Years Since Remodeling	-0.492*** (0.056)	-0.128***	-0.493*** (0.051)	-0.128***	-0.003*** (0.0002)	-0.170***
Low Dens. Zone	16.640*** (2.753)	0.086***	-31.118*** (8.368)	-0.160***	0.077** (0.039)	0.079**
Commercial Zone	-25.251** (11.881)	-0.026**	-35.782 (35.624)	-0.037	-0.673*** (0.165)	-0.139***
Floating Zone	18.710*** (5.258)	0.049***	-8.247 (22.841)	-0.021	0.127 (0.106)	0.065
Lot Area	0.001*** (0.0001)	0.074***	0.005*** (0.0005)	0.616***	0.00002*** (0.00000)	0.404***
I(Living Space^2)			-0.00000 (0.00000)	-0.002	-0.00000*** (0.000)	-0.234***
Living Space:Low Dens. Zone			0.019*** (0.004)	0.308***	0.00002 (0.00002)	0.058
Living Space:Commercial Zone			0.0001 (0.017)	0.0002	0.0002** (0.0001)	0.067**
Living Space:Floating Zone			0.013 (0.009)	0.088	0.00002 (0.00004)	0.022
Living Space:Lot Area			-0.00000*** (0.00000)	-0.647***	-0.000*** (0.000)	-0.396***
Year Sold	YES	YES	YES	YES	YES	YES
Overall Quality Rating	YES	YES	YES	YES	YES	YES
Building Type	YES	YES	YES	YES	YES	YES
R <sup>2</sup>		0.803		0.836		0.861
Adjusted R <sup>2</sup>		0.800		0.833		0.858
Residual Std. Error		35.548		32.488		0.150
df Residual Std. Error		(df = 1439)		(df = 1434)		(df = 1434)
F Statistic		292.393***		291.609***		354.941***
df F Statistic		(df = 20; 1439)		(df = 25; 1434)		(df = 25; 1434)

*Note:* N = 1460. All observations included in any model. OLS estimates, robust standard errors in parentheses.\*\*\* p<0.01, \*\* p<0.05, \* p<0.1. This table reportes the (unstandardized) regression results for the Model (1) excluding any interaction effects & and inclduing those interaction effects in Model (3). The corresponding standardized coefficients can be foudn in Std.Model (2) & (3). For complementary regression models **see appendix**. The comparison category in "Zone" (IV) is "Moderate-to-High Density". Years in date range is 2006 to 2010.

## 4.2 A2

Table 3

	GVIF	Df	GVIF <sup>1/(2*Df)</sup>
tot_liv_area	22.682	1	4.763
y_since_rem	1.535	1	1.239
low_density_zone	16.161	1	4.020
commercial_zone	11.942	1	3.456
floating_zone	30.700	1	5.541
I(tot_liv_area^2)	31.756	1	5.635
LotArea	32.949	1	5.740
YrSold	1.051	4	1.006
OverallQual_cat	4.301	9	1.084
Building_type	1.189	1	1.090
tot_liv_area:low_density_zone	31.262	1	5.591
tot_liv_area:commercial_zone	11.544	1	3.398
tot_liv_area:floating_zone	31.772	1	5.637
tot_liv_area:LotArea	46.292	1	6.804

In order to assess multicollinearity, the VIF was calculated using the "AER" package in R. To make the resulting VIF comparable, the  $GVIF^{1/(2*Df)}$  is used.<sup>8</sup>

While some variables display a VIF of more than 5, this does not automatically indicate multicollinearity. However, this necessitates further analysis. Considering i.a. Table 2 again, one can see that for the exception of some categories in Zoning, the standard errors stay stable. Additionally, the inclusion or exclusion of certain factors does not seem to induce a large change in the coefficients beyond what can be expected. Finally, the stability of the model and the significance of most of the terms in context of the significance of the overall model indicates that multicollinearity might not pose a considerable problem in this model.

However, what is notable is that (not reported) there seems to be some multicollinearity problems with the control variable categories of Quality. However, these do not seem to pose too big of a problem considering the main variables being not impacted.

## 4.3 A3

## 4.4 A4

## 4.5 A6

It is notable, that more complex analysis will be considered in part 4 & 5, such as the moderating effect of Years Since remodeling on the association of Quality and Sale price.

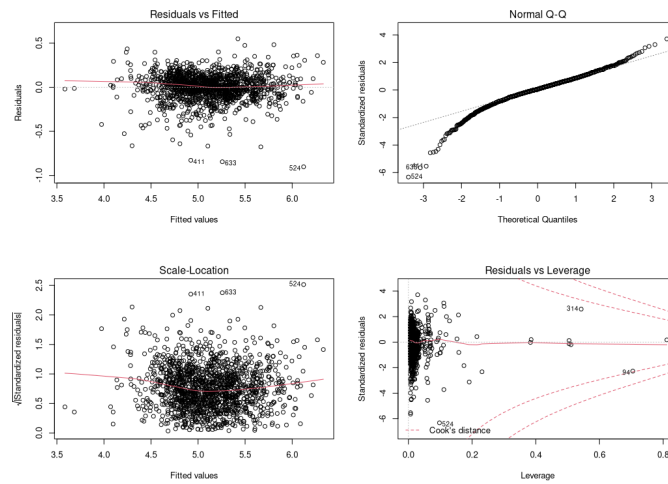
**IMPORTANT: FOR PART 4 CREATE A TABLE THAT CONTAINS OLS ROBUST TO STANDARD ERRORS!!!! (THIS EXACT TABLE!!)**

**IMPORTANT HERE WE WILL NOT ONLY COMPARE THE BASE MODEL WITHOUT SCALE SALE PRICE, But also THE FULL MODEL WITH SCALED EVERYTHING and corrections for standard errors**

to do: 1) log scale outcome variable because it is an important factor; 2) secondly, the quadratic term of Total living space is already in the part before, but here the rationale is different as we were only talking about the independent variables in that case. as such I did not consider the total living area to be an issue here

---

<sup>8</sup>IMPORTANT



(a) Positive association of total Living Space & Sale Price & optimal line; range [0, 7000]

Figure 2: QQ plots etc

## 5 Subset analyses

In this analysis we will examine further interaction and subsection analysis; such as quality as a numeric variable

MERGE THE STANDARDIZED AND UNSTANDARDIZED TABLE



Table 4: MERGED STANDARDIZED REGRESSON RESULTS CONTAINING BOTH STANDARDIZED AND UNSTANDARDIZED COEFFICIENTS FINAL VERSION BACKUP

	<i>Dependent variable:</i>			
	Sale Price		ln(Sale Price)	
	Model (1)	Std.Model (2)	Model (3)	Std.Model (4)
Constant	15.229 (25.697)	/	3.584*** (0.114)	/
Total Living Space	0.035*** (0.002)	0.362***	0.0003*** (0.00002)	0.684***
Years Since Remodelling	-0.492*** (0.056)	-0.128***	-0.003*** (0.0002)	-0.170***
Low Density Zone	16.640*** (2.753)	0.086***	0.077** (0.039)	0.079**
Commercial Zone	-25.251** (11.881)	-0.026**	-0.673*** (0.165)	-0.139***
Floating Zone	18.710*** (5.258)	0.049***	0.127 (0.106)	0.065
I(Total Living Area^2)			-0.00000*** (0.000)	-0.234***
Lot Area	0.001*** (0.0001)	0.074***	0.00002*** (0.00000)	0.404***
"Total Living Area": "Low Density Zone"			0.00002 (0.00002)	0.058
"Total Living Area": "Commercial Zone"			0.0002** (0.0001)	0.067**
"Total Living Area": "Floating Zone"			0.00002 (0.00004)	0.022
"Total Living Area": "Lot Area"			-0.000*** (0.000)	-0.396***
Year Sold	YES	YES	YES	YES
Overall Quality Rating	YES	YES	YES	YES
Building Type	YES	YES	YES	YES
R <sup>2</sup>	0.803	0.803	0.861	0.861
Adjusted R <sup>2</sup>	0.800	0.800	0.858	0.858
Residual Std. Error	35.548	35.548	0.150	0.150
df Residual Std. Error	(df = 1439)	(df = 1439)	(df = 1434)	(df = 1434)
F Statistic	292.393***	292.393***	354.941***	354.941***
df F Statistic	(df = 20; 1439)	(df = 20; 1439)	(df = 25; 1434)	(df = 25; 1434)

*Note:* N = 1460. All observations included in any model. OLS estimates, robust standard errors in parentheses.\*\*\* p<0.01, \*\* p<0.05, \* p<0.1. This table reportes the (unstandardized) regression results for the Model (1) excluding any interaction effects & and inclduing those interaction effects in Model (3). The corresponding standardized coefficients can be foudn in Std.Model (2) & (3). For complementary regression models **see appendix**. The comparison category in "Zone" (IV) is "Moderate-to-High Density". Years in date range is 2006 to 2010.