

MS&E 226: Fundamentals of Data Science

Lecture 3: More on linear regression

Ramesh Johari

`ramesh.johari@stanford.edu`

Recap: Linear regression

The linear regression model

Given:

- ▶ n outcomes Y_i , $i = 1, \dots, n$
- ▶ n vectors of covariates (X_{i1}, \dots, X_{ip}) , $i = 1, \dots, n$, $p < n$

let \mathbf{X} be the design matrix where the i 'th row is $(1, X_{i1}, \dots, X_{ip})$.

OLS solution (with intercept) is:

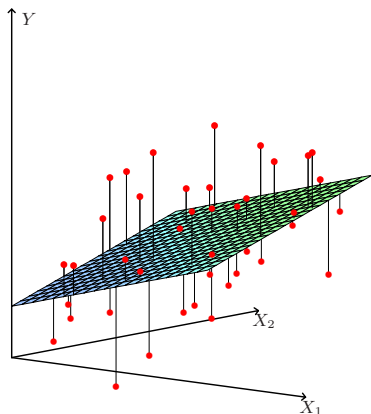
$$\hat{Y}_i = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_{ij},$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

Residuals

The OLS solution minimizes squared error; this is $\|\hat{\mathbf{r}}\|^2$, where $\hat{\mathbf{r}} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$.

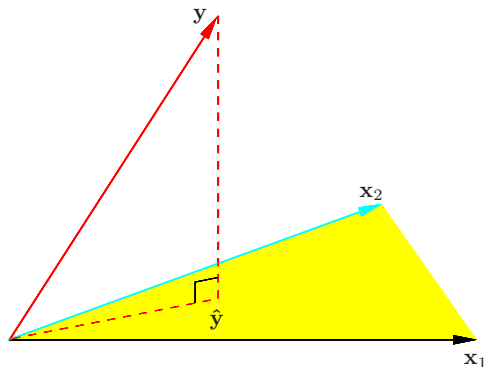
E.g., with two covariates:¹



¹Figure courtesy of *Elements of Statistical Learning*.

Geometry

This picture summarizes the geometry of OLS with two covariates:²



It explains why the residual vector $\hat{\mathbf{r}}$ is *orthogonal to every column of \mathbf{X}* .

²Figure courtesy of *Elements of Statistical Learning*.

Relating the previous two pictures

The first picture is in $p + 1$ -dimensional space: p covariates, and an additional dimension for the outcome.

- ▶ It displays a regression model fitted with an intercept term.
- ▶ The fitted value at each data point is the value on the regression plane associated to that data point.
- ▶ The black lines associated to each data point are the residuals.

The second picture is in n -dimensional space: each vector is a vector of the same length as the number of observations.

- ▶ The vector $\hat{\mathbf{Y}}$ is the vector of fitted values corresponding to the observations \mathbf{Y} .
- ▶ The residual vector is $\hat{\mathbf{r}} = \mathbf{Y} - \hat{\mathbf{Y}}$.

Key assumptions

We assumed that $p < n$ and \mathbf{X} has *full rank* $p + 1$.

What happens if these assumptions are violated?

Collinearity and identifiability

If \mathbf{X} does not have full rank, then $\mathbf{X}^\top \mathbf{X}$ is *not invertible*.

In this case, the optimal $\hat{\beta}$ that minimizes SSE is *not unique*.

The problem is that if a column of \mathbf{X} can be expressed as a linear combination of other columns, then the coefficients of these columns are not uniquely determined.³

We refer to this problem as *collinearity*. We also say the resulting model is *nonidentifiable*.

³In practice, \mathbf{X} may have full rank but be *ill conditioned*, in which case the coefficients $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ will be very sensitive to the design matrix.

Collinearity: Example

If we run `lm` on a less than full rank design matrix, we obtain NA in the coefficient vector:

```
> sh$livingArea_copy = sh$livingArea
> fm = lm(data = sh, price ~ 1 + livingArea + livingArea_copy)
> coef(fm)
      (Intercept)      livingArea livingArea_copy
      13439.3940         113.1225              NA
```

High dimension

If $p \approx n$, then the number of covariates is of a similar order to the number of observations.

Assuming the number of observations is large, this is known as the *high-dimensional* regime.

When $p + 1 \geq n$, we have enough *degrees of freedom* (through the $p + 1$ coefficients) to perfectly fit the data. Is this a good model?

Note that if $p \geq n$, then in general the model is nonidentifiable.

Interpreting regression coefficients

Coefficients

How to interpret the coefficients?

- ▶ $\hat{\beta}_0$ is the fitted value when all the covariates are zero.
- ▶ $\hat{\beta}_j$ is the change in the fitted value for a one unit change in the j 'th covariate, *holding all other covariates constant*.
- ▶ What language do we use to talk about coefficients?

Language

Suppose we completely believe our model. Different things we might say:

- ▶ “A one unit change in X_{ij} is *associated* (or *correlated*) with a $\hat{\beta}_j$ unit change in Y_i .”
- ▶ “Given a particular covariate vector (X_{i1}, \dots, X_{ip}) , we *predict* Y_i will be $\hat{\beta}_0 + \sum_j \hat{\beta}_j X_{ij}$.”
- ▶ “If X_{ij} changes by one unit, *then* Y_i will be $\hat{\beta}_j$ units higher.”

This course focuses heavily on helping you understand conditions under which these statements are possible (and more importantly, when they are not ok!).

Example in R

Recall `new` is Yes or No depending on whether the house is new construction.

```
> fm = lm(data = sh, price ~ 1 + new)
> summary(fm)
...
Coefficients:
              Estimate ...
(Intercept)  282307 ...
newNo        -73800 ...
...
```

Example in R

Note that:

```
> mean(sh$price[sh$new == "Yes"])  
[1] 282306.8  
> mean(sh$price[sh$new == "No"])  
[1] 208507.4  
> mean(sh$price[sh$new == "Yes"]) -  
+   mean(sh$price[sh$new == "No"])  
[1] 73799.46
```

Another example with categorical variables

Recall heating:

- ▶ electric
- ▶ hot water/steam
- ▶ hot air

Does it make sense to build a model where:

price $\approx \hat{\beta}_0 + \hat{\beta}_1$ heating?

(Not really.)

Example in R: Categorical variables

Example with heating:

```
> fm = lm(data = sh, price ~ 1 + heating)
```

```
> summary(fm)
```

```
...
```

```
Coefficients:
```

	Estimate	...
(Intercept)	226355	...
heatinghot water/steam	-17223	...
heatingelectric	-64467	...

```
...
```

Example in R: Categorical variables

Effectively, R creates two new binary variables:

- ▶ The first is 1 if heating is hot water/steam, and zero otherwise.
- ▶ The second is 1 if heating is electric, and zero otherwise.

What if they are all zero? Why is there no variable heatinghot air?

What do the coefficients mean?

```
> mean(sh$price[sh$heating == "hot water/steam"])  
[1] 209132.5  
> mean(sh$price[sh$heating == "electric"])  
[1] 161888.6  
> mean(sh$price[sh$heating == "hot air"])  
[1] 226355.4
```

Conditional means

We are observing something that will pop up repeatedly:

The regression model tries to estimate the conditional average of the outcome, given the covariates.

In the simple regression setting $Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$ with a binary covariate $X_i \in \{0, 1\}$, we find that:

- ▶ $\hat{\beta}_0$ is the sample mean of the Y_i 's where $X_i = 0$.
- ▶ $\hat{\beta}_0 + \hat{\beta}_1$ is the sample mean of the Y_i 's where $X_i = 1$.

In other words, $\hat{\beta}_0 + \hat{\beta}_1 X_i$ is the *conditional sample mean* of the Y_i 's given X_i .

An example with continuous covariates

Recall our linear model with the SaratogaHouses data:

```
> fm = lm(data = sh, price ~ 1 + livingArea)
```

```
> summary(fm)
```

```
...
```

```
Coefficients:
```

```
                Estimate ...
```

```
(Intercept) 13439.394 ...
```

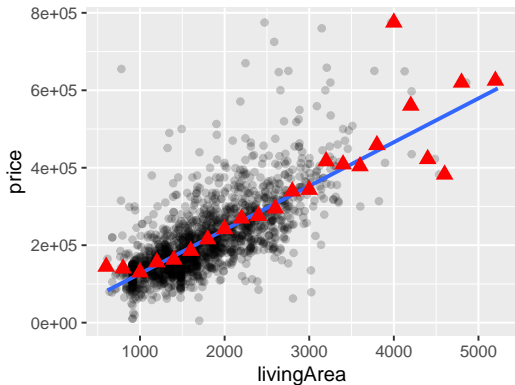
```
livingArea   113.123 ...
```

```
...
```

What does 13439.394 represent? (Is it sensible?) What does 113.123 represent?

Example in R

Let's plot the model, together with the conditional sample mean of price given different values of livingArea:



In this figure, livingArea is rounded to the nearest multiple of 200, and the resulting price values are averaged to produce the triangles.

Simple linear regression

We can get more intuition for the regression coefficients by looking at the case where there is only one covariate:

$$Y_i \approx \hat{\beta}_0 + \hat{\beta}_1 X_i.$$

(Note we dropped the second index on the covariate.)

It can be shown that:

$$\hat{\beta}_1 = \hat{\rho} \frac{\hat{\sigma}_Y}{\hat{\sigma}_X}; \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X},$$

where $\hat{\rho}$ is the *sample correlation*:

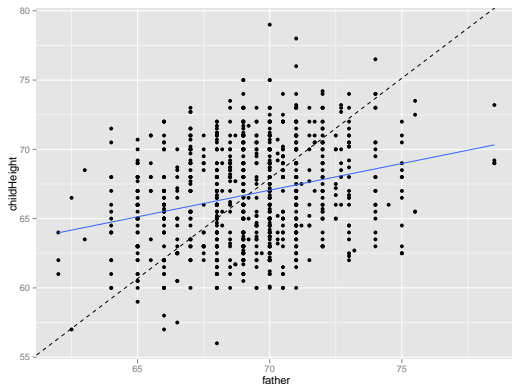
$$\hat{\rho} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\hat{\sigma}_X \hat{\sigma}_Y}.$$

Simple linear regression

To get intuition for $\hat{\beta}_1 = \hat{\rho}\hat{\sigma}_Y/\hat{\sigma}_X$, note that $-1 \leq \hat{\rho} \leq 1$.

We can compare to the *SD line*, that goes through (\bar{X}, \bar{Y}) and has slope $\text{sign}(\hat{\rho}) \times \hat{\sigma}_Y/\hat{\sigma}_X$.

We do this in the following graph, using data from Galton on the heights of children and their fathers [SM]:



“Reversion” to the mean

Note that correlation $\hat{\rho}$ is typically such that $|\hat{\rho}| < 1$.

E.g., if $\hat{\rho} > 0$:

- ▶ Suppose a covariate X_i is A s.d.'s larger than \bar{X}
- ▶ The fitted value \hat{Y}_i will only be $\hat{\rho}A$ s.d.'s larger than \bar{Y} .

On average, fitted values are closer to their mean than the covariates are to their mean. This is sometimes referred to as “mean reversion” or “regression to the mean” (terrible terminology).

Beyond linearity

Linearity

The linear regression model projects the outcomes into a hyperplane, determined by the covariates.

This fit might have systematic problems because the relationship between \mathbf{Y} and \mathbf{X} is inherently *nonlinear*.

Sometimes looking at the residuals will suggest such a problem, but not always!

For this reason context and domain expertise is critical in building a good model.

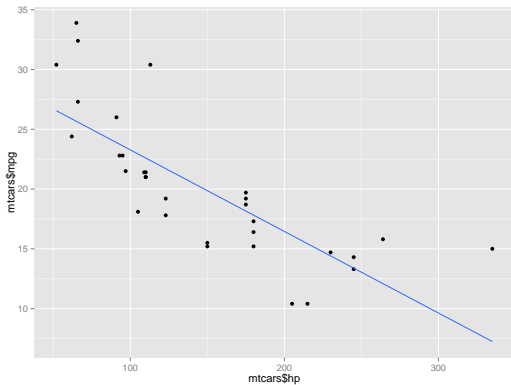
Higher order terms

For this example we use the Motor Trend dataset on cars (built into R):

```
> data(mtcars)
> fm = lm(data = mtcars, formula = mpg ~ 1 + hp)
> summary(fm)
...
Coefficients:
              Estimate ...
(Intercept) 30.09886 ...
hp          -0.06823 ...
...
Multiple R-squared:  0.6024 ...
...
```

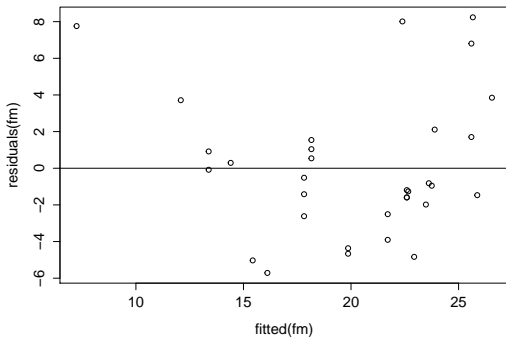
Higher order terms

Visualization reveals that the line is not a great fit...



Higher order terms

...and the residuals look suspicious:



Higher order terms

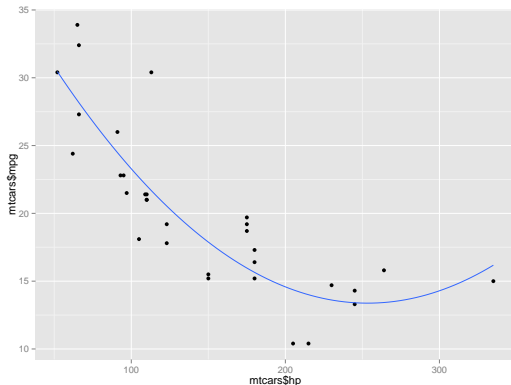
The fit suggests that we might benefit by modeling mpg as being a *quadratic* function of hp.

```
> fm = lm(data = mtcars, formula = mpg ~ 1 + hp + I(hp^2))
> summary(fm)
...
Coefficients:
              Estimate ...
(Intercept)  4.041e+01 ...
hp           -2.133e-01 ...
I(hp^2)      4.208e-04 ...
...
Multiple R-squared:  0.7561 ...
...
```

Note increase in R^2 with new model.

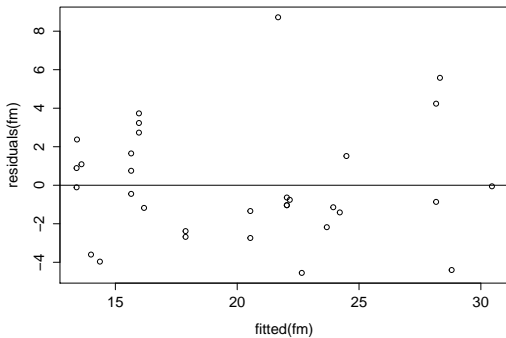
Higher order terms

Visualization suggests better fit with quadratic model...



Higher order terms

...and the residuals look (a little) better:



Higher order terms

Consider a model with one covariate, with values (X_1, \dots, X_n) .

Consider the linear regression model:

$$Y_i \approx \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 X_i^2 + \dots + \hat{\beta}_k X_i^k.$$

What happens to R^2 as you increase k ?

Interactions

Consider the following example:

```
> fm = lm(data = sh, formula = price ~ 1 + new + livingArea)
> summary(fm)
...
Coefficients:
              Estimate ...
(Intercept) -3680.83 ...
newNo       15394.22 ...
livingArea   114.52 ...
...
```

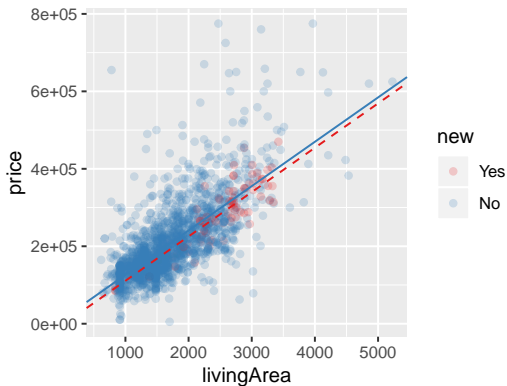
Interpretation:

- ▶ $\text{new} = \text{Yes} \implies \text{price} \approx -3681 + 115 \times \text{livingArea}.$
- ▶ $\text{new} = \text{No} \implies \text{price} \approx 11713 + 115 \times \text{livingArea}.$

Note that both have the *same slope*.

Interactions

Visualization:



The plot suggests *higher* slope when $\text{New} = \text{Yes}$.

Interactions

When changing the *value* of one covariate affects the *slope* of another, then we need an *interaction* term in the model.

E.g., consider a regression model with two covariates:

$$Y_i \approx \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2}.$$

The model with an interaction between the two covariates is:

$$Y_i \approx \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \hat{\beta}_2 X_{i2} + \hat{\beta}_{1:2} X_{i1} X_{i2}.$$

Interactions

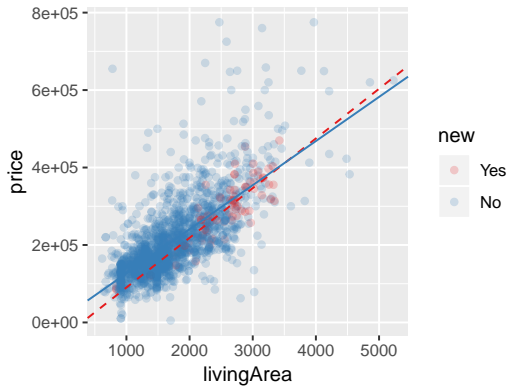
```
> fm = lm(data = sh,
           formula = price ~ 1 + new +
           livingArea + new:livingArea)
> summary(fm)
...
Coefficients:
                Estimate ...
(Intercept)    -38045.90 ...
newNo           50804.51 ...
livingArea      128.28 ...
newNo:livingArea -14.37 ...
...
```

Interpretation:

- ▶ When new = Yes,
then $\text{price} \approx -38046 + 128 \times \text{livingArea}$.
- ▶ When new = No,
then $\text{price} \approx 12759 + 114 \times \text{livingArea}$.

Interactions

Visualization:



Rules of thumb

When should you try including interaction terms?

- ▶ If a particular covariate has a large effect on the fitted value (high coefficient), it is also worth considering including interactions with that covariate.
- ▶ Often worth including interactions with covariates that describe *groups* of data (e.g., `new` or `heating` in this example), since coefficients of other covariates may differ across groups.

Summary

Higher order terms and interactions are powerful techniques for making models much more flexible.

A warning: now the effect of a single covariate is captured by multiple coefficients!

Beyond minimizing SSE

Minimizing SSE

Ordinary least squares minimizes the sum of squared errors:

$$\text{SSE} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2.$$

What about other objective functions?

Motivation: Adding covariates

Adding covariates to a model can only make R^2 increase.

Why? Each additional covariate “explains” at least some of the variance in the outcome variable.

So is a model with more covariates “better”?

Regularization

Sometimes, with many covariates, we only want our regression to pick out coefficients that are “meaningful” .

One way to achieve this is to *regularize* the objective function.

Regularization: Ridge regression

Instead of minimizing SSE, minimize:

$$\text{SSE} + \lambda \sum_{j=1}^p |\hat{\beta}_j|^2.$$

where $\lambda > 0$. This is called *ridge regression*.

In practice, the consequence is that it penalizes $\hat{\beta}$ vectors with “large” norms.

Regularization: Lasso

Instead of minimizing SSE, minimize:

$$\text{SSE} + \lambda \sum_{j=1}^p |\hat{\beta}_j|$$

where $\lambda > 0$.

This is called the *Lasso*.

In practice, the resulting coefficient vector will be “sparser” than the unregularized coefficient vector.

Regularization

Both lasso and ridge regression are “shrinkage” methods for covariate selection:

- ▶ Relative to OLS, both lasso and ridge regression will yield coefficients $\hat{\beta}$ that have “shrunk” towards zero.
- ▶ The most explanatory covariates are the ones that will be retained.
- ▶ Lasso typically yields a much smaller subset of nonzero coefficients than ridge regression or OLS (i.e., fewer nonzero entries in $\hat{\beta}$).

Intuition for lasso

Why does lasso tend to “truncate” more coefficients at zero than ridge?

Intuition for lasso

Regularization

Often regularized regression is used for “kitchen sink” data analysis:

- ▶ Include every covariate you can find.
- ▶ Run a regularized regression to “pick out” which covariates are important.

What are some pros and cons of this approach?

Robustness to “outliers”

Linear regression can be very sensitive to data points far from the fitted model (“outliers”).

One source of this sensitivity is that $(Y_i - \hat{Y}_i)^2$ is *quadratic*.

- ▶ Derivative w.r.t. Y_i is $2(Y_i - \hat{Y}_i)$.
- ▶ So if Y_i is far away from the fitted value, small changes in Y_i cause large changes in SSE.
- ▶ This in turn makes the optimal model very sensitive to large Y_i .

Robustness to “outliers”

An alternative approach is to minimize the sum of *absolute deviations*:

$$\sum_{i=1}^n |Y_i - \hat{Y}_i|.$$

Though computationally more challenging to optimize, this approach tends to be more stable in the face of outliers:

The resulting linear model approximates the *conditional median* (instead of the conditional mean).

(Note that regularized regression is also less vulnerable to outliers.)

Data transformations

Transformations of the data

Often it is useful to work with transformed versions of the data, to make regressions more meaningful.

We already saw one example of this: creating indicators for each value of a categorical variable.

We'll discuss two more transformations in particular:

- ▶ Logarithms of positive variables
- ▶ Centering and standardizing

Logarithmic transformations

In many contexts, outcomes are *positive*: e.g., physical characteristics (height, weight, etc.), counts, revenues/sales, etc.

For such outcomes linear regression can be problematic, because it can lead to a model where \hat{Y}_i is negative for some \mathbf{X}_i .

One approach to deal with this issue is to take a *logarithmic* transformation of the data before applying OLS:

$$\log Y_i \approx \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_{ij}.$$

Exponentiating, note that this becomes a model that is *multiplicative*:

$$Y_i \approx e^{\hat{\beta}_0} e^{\hat{\beta}_1 X_{i1}} \dots e^{\hat{\beta}_p X_{ip}}.$$

Note that holding all other covariates constant, a *one unit change in X_{ij} is associated with a proportional change in the fitted value by $e^{\hat{\beta}_j}$.*

Logarithmic transformations

Note that $e^{\hat{\beta}_j} \approx 1 + \hat{\beta}_j$ for small $\hat{\beta}_j$.

Different way to think about it:

$$\log(a) - \log(b) \approx \frac{a}{b} - 1$$

when a/b is close to 1; i.e., difference in logs gives (approximately) *percentage* changes.

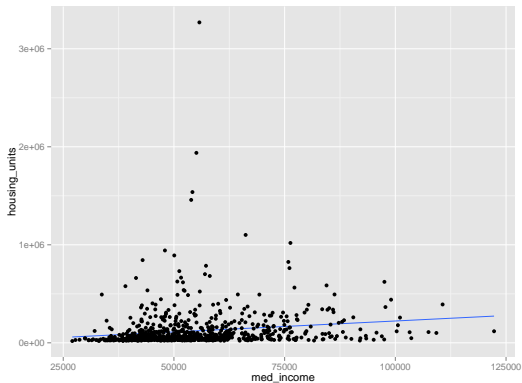
So we can interpret $\hat{\beta}_j$ as suggesting the proportional change in the outcome associated with a one unit change in the covariate.

If both data and outcome are logged, then $\hat{\beta}_j$ gives the proportional change in the outcome associated with a proportional change in the covariate.

Logarithmic transformations: Example

Data: 2014 housing and income by county, from U.S. Census

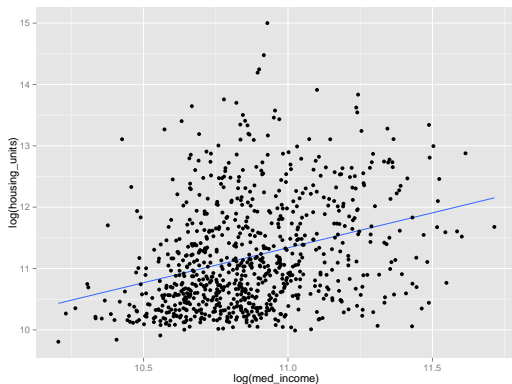
First plot number of housing units against median household income:



Logarithmic transformations: Example

Data: 2014 housing and income by county, from U.S. Census

Now do the same with logarithmically transformed data:



Logarithmic transformations: Example

Data: 2014 housing and income by county, from U.S. Census

The resulting model:

```
> fm = lm(data = income,
  formula = log(housing_units) ~ 1 + log(med_income))
> summary(fm)
...
Coefficients:
                Estimate ...
(Intercept)      -1.194 ...
log(med_income)   1.139 ...
...
```

The coefficient can be interpreted as saying that a 1% higher median household income is associated with a 1.14% higher number of housing units, on average.

Centering

Sometimes to make coefficients interpretable, it is useful to *center* covariates by removing the mean:

$$\tilde{X}_{ij} = X_{ij} - \bar{X}_j.$$

(Here $\bar{X}_j = \frac{1}{n} \sum_i X_{ij}$ denotes the sample mean of the j 'th covariate.)

In this case the regression model is:

$$Y_i \approx \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j \tilde{X}_{ij}.$$

Centering: Example

Consider our earlier regression, but now center `livingArea` first:

```
> sh$livingArea.c = sh$livingArea - mean(sh$livingArea)
> fm = lm(data = sh, formula = price ~ 1 + sh$livingArea.c)
> summary(fm)
```

...

Coefficients:

	Estimate	...
(Intercept)	211966.705	...
sh\$livingArea.c	113.123	...

...

Now the intercept is *directly interpretable* as (approximately) the average value of `price` around the average value of `livingArea`.

Centering

Two additional notes:

- ▶ Often useful to *standardize*: center and divide by sample standard deviation, i.e.

$$\tilde{X}_{ij} = \frac{X_{ij} - \bar{X}_j}{\hat{\sigma}_j},$$

where $\hat{\sigma}_j$ is sample standard deviation of j 'th covariate. This gives all covariates a normalized dispersion.

- ▶ Can also center the outcome $\tilde{Y}_i = Y_i - \bar{Y}$. Note that if all covariates and the outcome are centered, there will be no intercept term (why?).