

Regression, Prediction and Classification

Jacob LaRiviere
(some content from Justin Rao)

Terminology

X : *features*.

y : outcomes

Goal is to model outcomes as a function of features.

Width of X is p

Length of X and y is n

Terminology cont.

X : *features*.

y : outcomes

$$\text{Obs 1} \longrightarrow \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

Feature 1

$$\text{Obs 1} \longrightarrow \mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

Estimating equations

$$Y = f(X) + \epsilon$$

$$y_i = f(x_i) + \epsilon_i$$

Where ϵ_i is call the error term.

Example: coin flipping

$$Y = f(X) + \epsilon$$

$$flip_i = p + \epsilon_i$$

if outcome = 1 then $\epsilon_i = 1 - p$
if outcome = 0 then $\epsilon_i = -p$

p is a constant, called the “success rate”

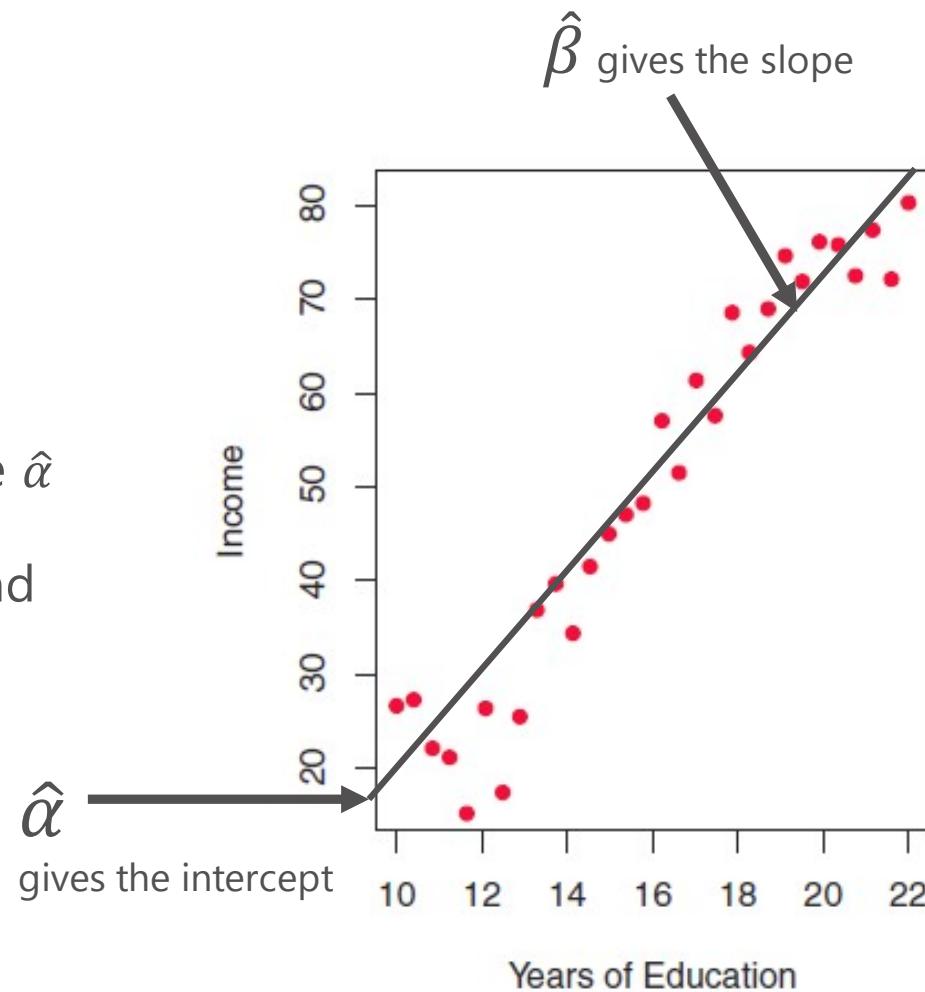
Linear regression

$$y_i = \alpha + \beta_1 x_i + \epsilon_i$$

Ordinary least squares will find the $\hat{\alpha}$ and $\hat{\beta}$ to minimize the squared distance between the fitted line and the observations

$$\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$$

Minimizes $(\hat{y} - y)^2$



How does this work intuitively?

Assume $\alpha = 0$

$$\begin{aligned} \min_{\{\beta\}} (y - \hat{y})^2 &= (y - \hat{\beta}x)^2 \\ &= y^2 - 2\hat{\beta}xy + (\hat{\beta}x)^2 \\ FOC \rightarrow -2xy + 2\hat{\beta}x^2 &= 0 \\ \rightarrow \hat{\beta}x^2 &= xy \\ \hat{\beta} &= \frac{xy}{x^2} \end{aligned}$$

$$\hat{\beta} = (X'X)^{-1}(X'Y)$$

How does this work for inference?

What about hypothesis testing? (Now assume $\alpha \neq 0$)

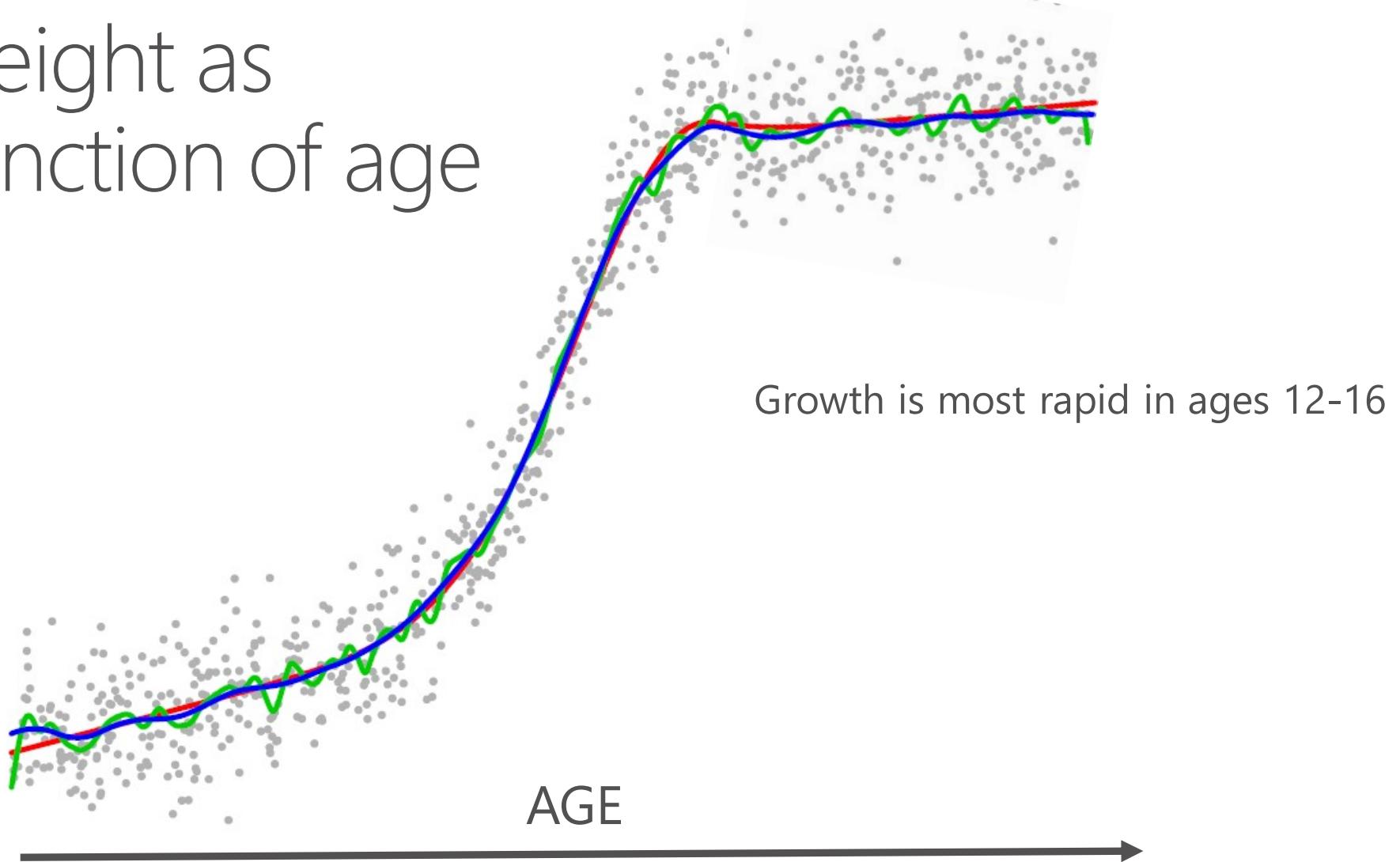
$$se(\hat{\beta}) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2}$$

$$= \frac{1}{n} \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{\epsilon}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

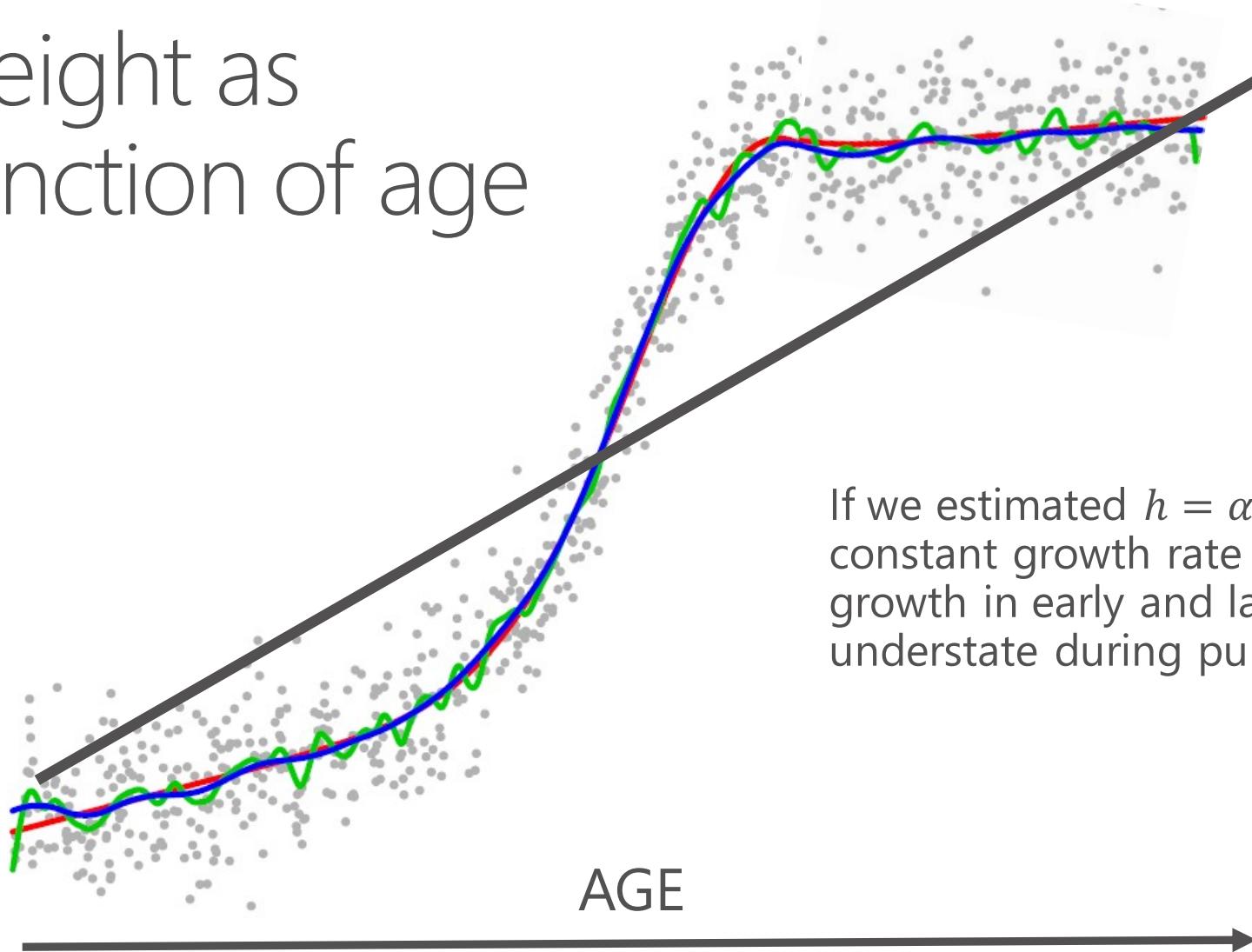
$$\rightarrow t = \frac{\hat{\beta}_1 - \beta_{1,0}}{se(\hat{\beta})}$$

Both of these are themselves sample statistics from data y and X. As a result, they both have their own distributions (e.g., normal and Chi Squared). The results ration has the Student t distribution which is used for hypothesis testing for a given sample.

Height as function of age



Height as function of age



If we estimated $h = \alpha + \beta * \text{age}$, the constant growth rate β would overstate growth in early and later years, and understate during puberty

How would we estimate this relationship?

General formula:

$$h_i = f(\text{age}_i) + \epsilon_i$$

ϵ_i allows people of the same age to have different heights. Even if we have the correct growth model using age alone, we expect some variation in height conditional on age.

To use a linear model, polynomial features allow for a non-linear relationship between the feature, age, and the outcome, height.

$$h_i = \alpha + \beta_1 \text{age}_i + \beta_2 \text{age}_i^2 + \cdots + \beta_p \text{age}_i^p + \epsilon_i$$

Linear regression

Linear regression allows you to use non-linear functions of the features, provided they enter the function as an additive term, with weight β

$$y_i = \alpha + \beta_1 g_1(x_{i,1}) + \cdots + \beta_p g_p(x_{i,p}) + \epsilon_i$$

The simplest example is where all features simply enter without any transformations

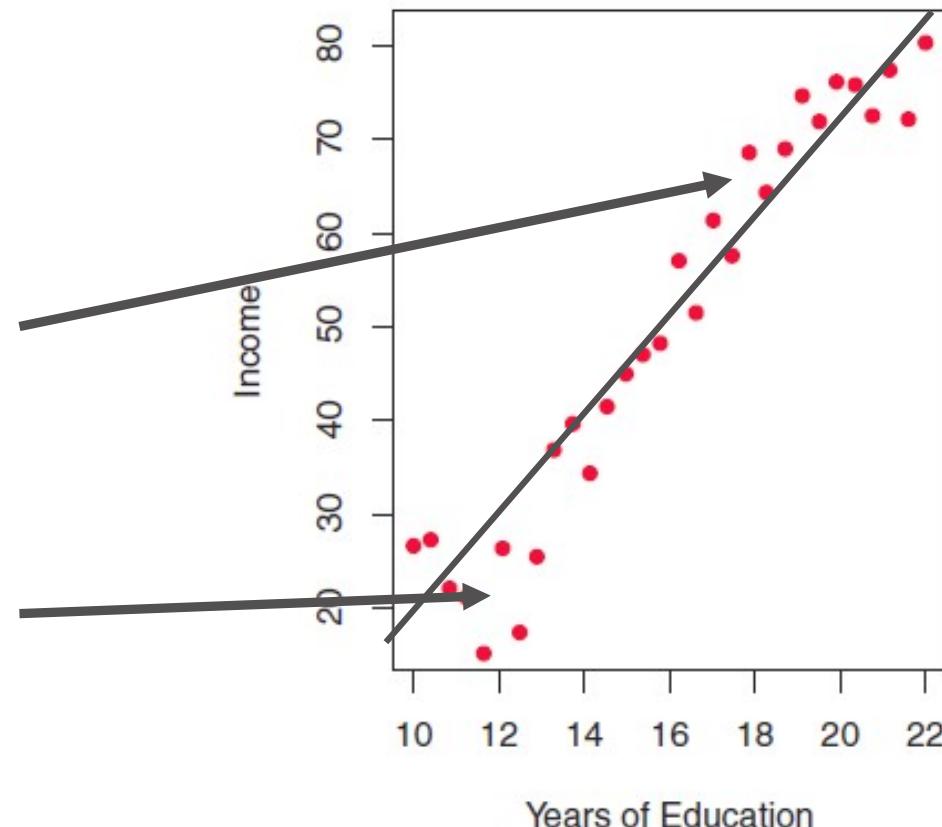
$$y_i = \alpha + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p} + \epsilon_i$$

Linear regression

$$y_i = \alpha + \beta_1 x_i + \epsilon_i$$

A linear function tends to underestimate for medium-high education

And overestimate income for low education

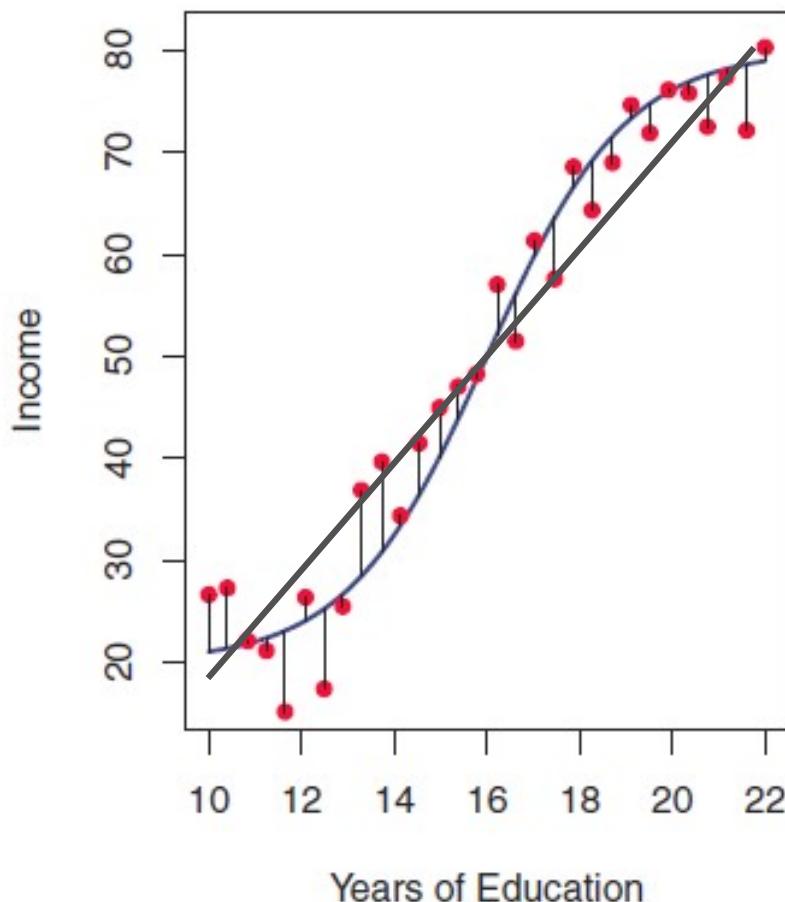


Linear regression

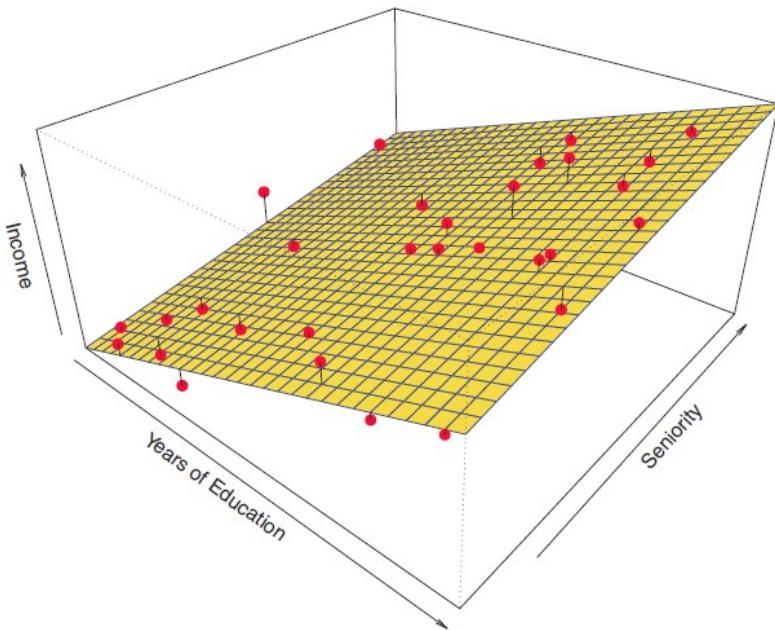
A cubic polynomial fits the data much better

$$\hat{y}_i = \hat{\alpha} + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \hat{\beta}_3 x_i^3$$

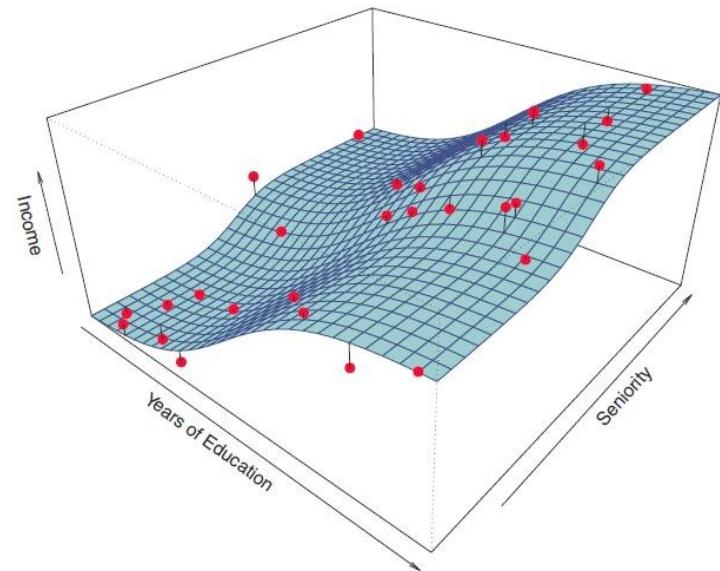
$(\hat{y} - y)^2$ is now much smaller



Visualizing in two dimensions



$$y_i = \alpha + \beta_1 \text{educ}_i + \beta_2 \text{seniority} + \epsilon_i$$



Higher order polynomial in educ and seniority could approximate this function

Steps for linear regression

1. Define the relationship we are trying to model. What is the outcome? What raw features do we have? E.g.

$$h_i = f(\text{age}_i) + \epsilon_i$$

2. Define a linear model to approximate this relationship, e.g.

$$h_i = \alpha + \beta_1 \text{age}_i + \beta_2 \text{age}_i^2 + \dots + \beta_p \text{age}_i^p + \epsilon_i$$

3. Create the necessary features eg. Age2=age² (or use poly)
4. Estimate the model:

```
mymodel=lm(y~age+age2...+ageP)  
summary(mymodel)
```

5. Evaluate the model with an evaluation metric.
6. Repeat steps 2-5 to improve model fit.

Model prediction

1. After we have fit a model, we can predict the outcome for any input.

$$\hat{h}_i = \alpha + \hat{\beta}_1 age_i + \hat{\beta}_2 age_i^2 + \cdots + \hat{\beta}_p age_i^p$$

2. For any given observation $\hat{h}_i - h_i = prediction\ error$

$(\hat{h}_i - h_i)^2$: squared loss
 $|\hat{h}_i - h_i|$: absolute loss

Model Fit

Fit refers to how well our model, an approximation of reality, matches what we actually observe in the data

Mean squared error is the average of squared loss for each data point (row) we evaluate

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

Training vs. Test

We typically want to “train” the model on the data “we have” and test the model on “new data”

This is a realistic test of how well the model predicts outcomes.

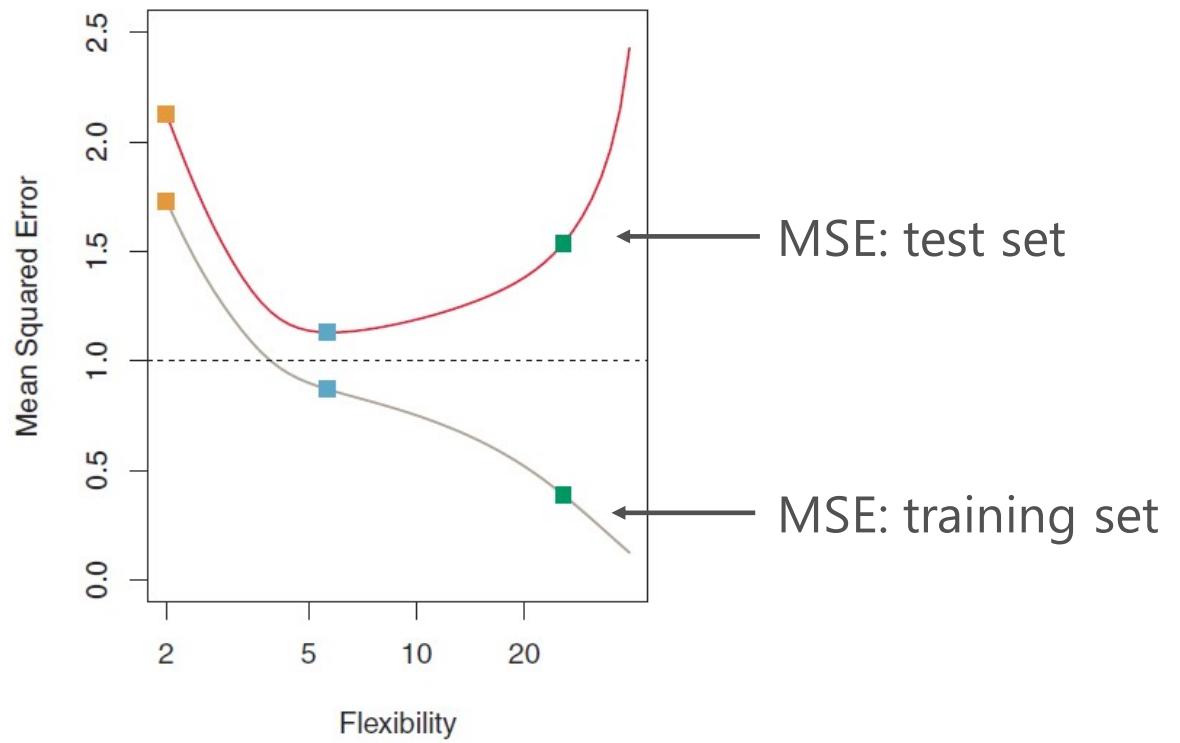
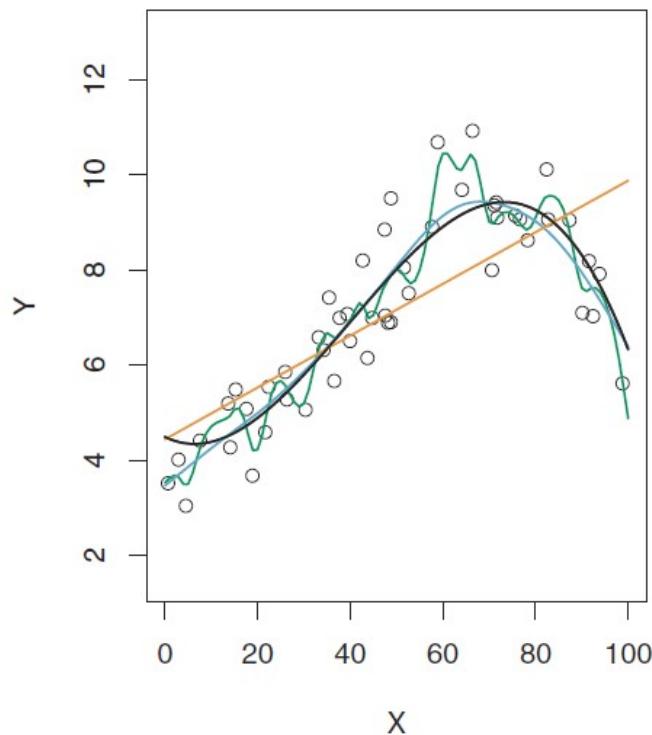
In practice, we will subset our data:

X%: Training

Y%: Test

80% is commonly used for training.

Why do we need training vs. test?

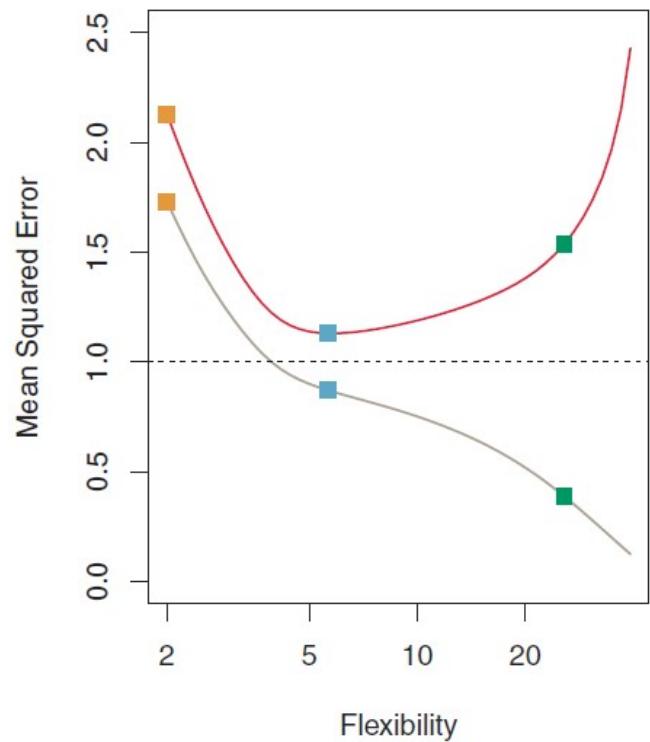
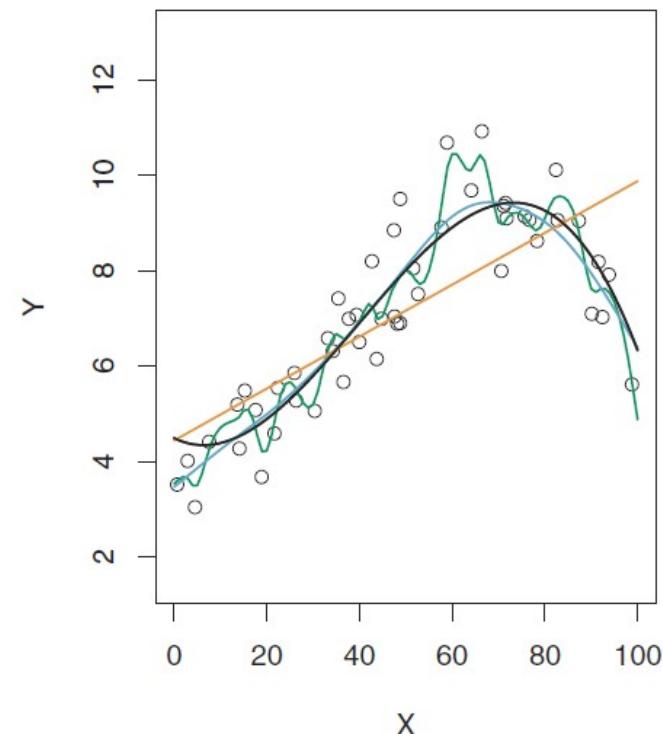


Why do we need training vs. test?

The green curve “overfits” the data. It fits the noise.

Using a “out of sample” test set guards against “overfitting”

The blue curve gets closest to the black (truth) curve



Training vs. test in practice

- Randomly select observations (rows) to be in the test set. The remainder are the training set.
- Why random: want complete coverage. E.g. don't want to train in February and test on July.
- Subset on the training set for all model estimation. E.g. any `data.frame` you pass to `lm`
- Use the `predict` command to make predictions (\hat{y}) for the test set
- Compute the ***evaluation metric*** using test set (you have the observe value and the predicted or modeled value)

R-squared

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Residual sum of squares

$$\text{TSS} = \sum (y_i - \bar{y})^2$$

Total variation in outcome variable

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

Fraction of the variation in the outcome variable captured by the model

Adjusted R-squared

- R-squared is reported for the training data. That is, it is “in sample”
- The “in sample” model fit weakly improves with more explanatory variables
- $Adjusted\ R^2 = 1 - \frac{n-1}{n-k-1} \frac{SSR}{TSS}$
- n is number of observations and k is number of explanatory variables.
 - penalizes adjusted r^2 when add variables without explanatory power.

A “Fair” R-squared Measure

- R-squared is reported for the training data. That is, it is “in sample”
- We learned that we typically want to use out-of-sample evaluation metrics, what do you we do
- It turns out that $R^2 = \text{corr}(\hat{y}_i, y_i)^2$
- This means, $\text{corr}(\hat{y}_i, y_i)^2$, for the test set can be interpreted just like R-squared (% of variation explain) and is a “fair test”

Related to an important concept called over-fitting.

Feature types

Binary features: equal to 1 or 0. For example:

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male,} \end{cases}$$

Continuous features: can take an real value

Categorical features: can take one of N values. E.g. Saturday, Sunday, Monday... declaring as `.factor` then R will treat it as N binary variables for each category (ex. =1 if Saturday, =0 otherwise)

Running Example: Orange Juice

Data: 83 Chicago-area stores



At weekly level:

- Sales ("log move")
- Average sales price
- Whether advertised ("feat")

At store level:

- Various demographics of the shoppers

Data taken from:

"Determinants of Store-Level Price Elasticity" Stephen J. Hoch, Byung-Do Kim, Alan L. Montgomery and Peter F. Rossi
Journal of Marketing Research Vol. 32, No. 1 (Feb., 1995), pp. 17-29



Data

Log of quantity sold
each week

Price is not logged

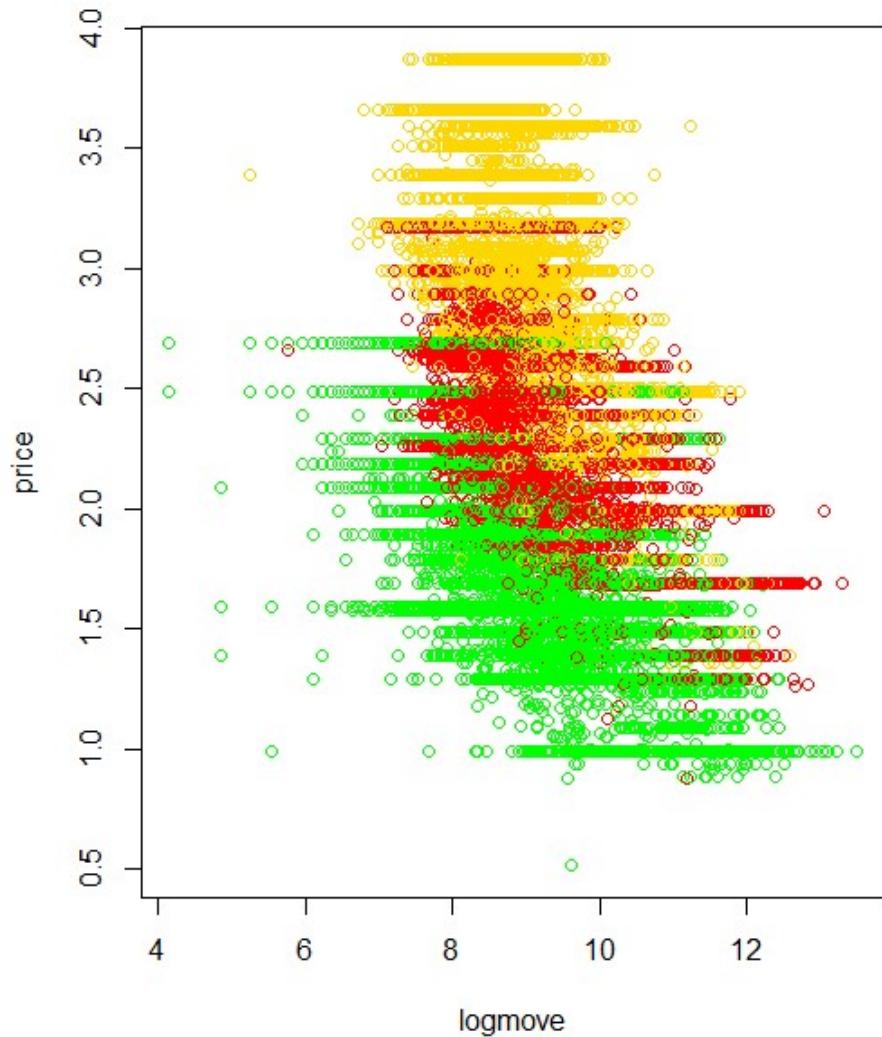
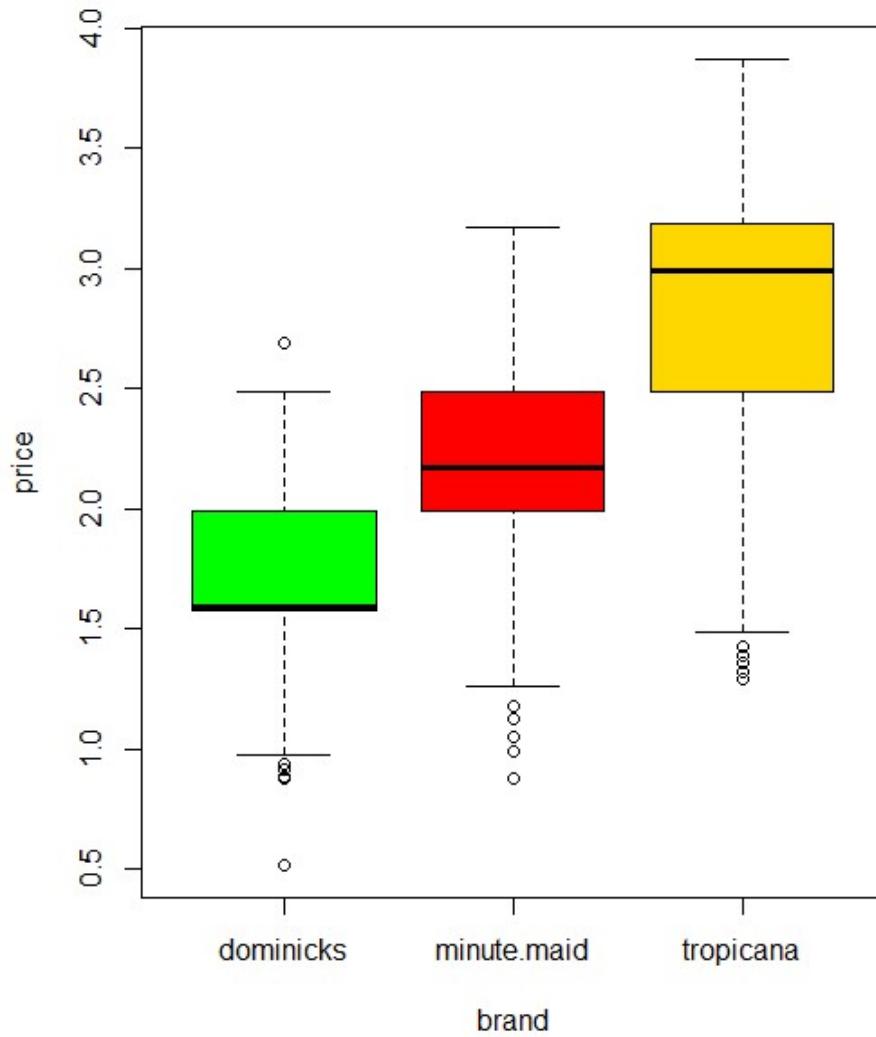
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	store	brand	week	logmove	feat	price	AGE60	EDUC	ETHNIC	INCOME	HHLARGE	WORKWOM	HVAL150	SSTRDIST	SSTRVOL	CPDIST5	CPWVOL5
2	2 tropicana		40	9.018695488	0	3.87	0.232864734	0.248934934	0.11428	10.55321	0.103953406	0.303585347	0.463887065	2.110122129	1.142857143	1.92728	0.376927
3	2 minute.maid		40	8.407378325	0	3.17	0.232864734	0.248934934	0.11428	10.55321	0.103953406	0.303585347	0.463887065	2.110122129	1.142857143	1.92728	0.376927
4	2 dominicks		40	9.264828557	1	1.59	0.232864734	0.248934934	0.11428	10.55321	0.103953406	0.303585347	0.463887065	2.110122129	1.142857143	1.92728	0.376927
5	5 tropicana		40	8.68067166	0	3.66	0.117368032	0.32122573	0.053875	10.92237	0.103091585	0.410568032	0.535883355	3.801997814	0.681818182	1.600573	0.736307
6	5 minute.maid		40	8.348537825	0	2.99	0.117368032	0.32122573	0.053875	10.92237	0.103091585	0.410568032	0.535883355	3.801997814	0.681818182	1.600573	0.736307
7	5 dominicks		40	7.491087594	1	1.59	0.117368032	0.32122573	0.053875	10.92237	0.103091585	0.410568032	0.535883355	3.801997814	0.681818182	1.600573	0.736307
8	8 tropicana		40	9.093357017	0	3.19	0.252394035	0.095173274	0.035243	10.59701	0.131749698	0.283074736	0.054227156	2.636332801	1.5	2.905384	0.641016
9	8 minute.maid		40	8.712759975	0	2.62	0.252394035	0.095173274	0.035243	10.59701	0.131749698	0.283074736	0.054227156	2.636332801	1.5	2.905384	0.641016
10	8 dominicks		40	9.570529135	1	1.59	0.252394035	0.095173274	0.035243	10.59701	0.131749698	0.283074736	0.054227156	2.636332801	1.5	2.905384	0.641016
11	9 tropicana		40	8.987196821	0	3.19	0.269119018	0.222172318	0.032619	10.78715	0.096830474	0.358994501	0.505747126	1.103278736	0.666666667	1.820474	0.441268
12	9 minute.maid		40	8.449342525	0	2.62	0.269119018	0.222172318	0.032619	10.78715	0.096830474	0.358994501	0.505747126	1.103278736	0.666666667	1.820474	0.441268
13	9 dominicks		40	9.056722883	1	1.59	0.269119018	0.222172318	0.032619	10.78715	0.096830474	0.358994501	0.505747126	1.103278736	0.666666667	1.820474	0.441268
14	14 tropicana		40	8.764053269	0	3.87	0.213949275	0.348293024	0.034179	11.04393	0.107894294	0.362305746	0.750769908	3.320947571	0.75	2.359891	0.73181
15	14 minute.maid		40	7.82244473	0	3.17	0.213949275	0.348293024	0.034179	11.04393	0.107894294	0.362305746	0.750769908	3.320947571	0.75	2.359891	0.73181
16	14 dominicks		40	8.886270902	1	1.59	0.213949275	0.348293024	0.034179	11.04393	0.107894294	0.362305746	0.750769908	3.320947571	0.75	2.359891	0.73181
17	18 tropicana		40	8.886270902	0	3.19	0.272313368	0.072246456	0.074417	10.39198	0.097937631	0.313907719	0.093351946	1.980307092	0.666666667	1.684315	0.154225
18	18 minute.maid		40	8.502688505	0	2.62	0.272313368	0.072246456	0.074417	10.39198	0.097937631	0.313907719	0.093351946	1.980307092	0.666666667	1.684315	0.154225
19	18 dominicks		40	8.528330936	1	1.59	0.272313368	0.072246456	0.074417	10.39198	0.097937631	0.313907719	0.093351946	1.980307092	0.666666667	1.684315	0.154225

We'll study 3 brands,
"dominicks" is the
store's brand

1 if advertised that
week by that store, 0 if
not

Demographics of
shoppers by store

Office

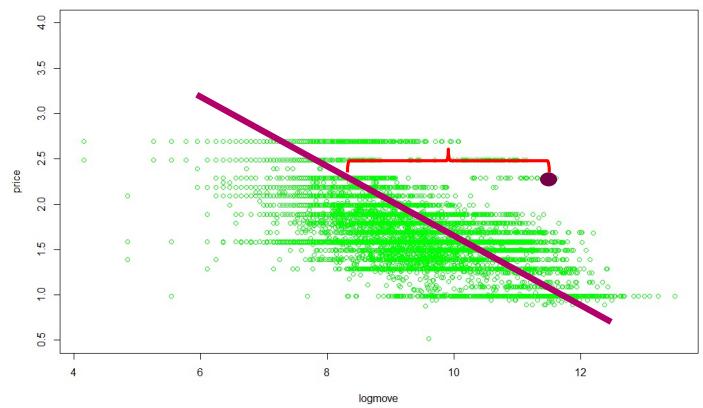


Estimating equations

$$Y = f(X) + \epsilon \rightarrow Y = \ln(Q) = \alpha + \ln(X) + \epsilon$$

$$\rightarrow \log move_{it} = \alpha + \log(price_{it}) + \epsilon_{it}$$

ϵ_i is called the error term.

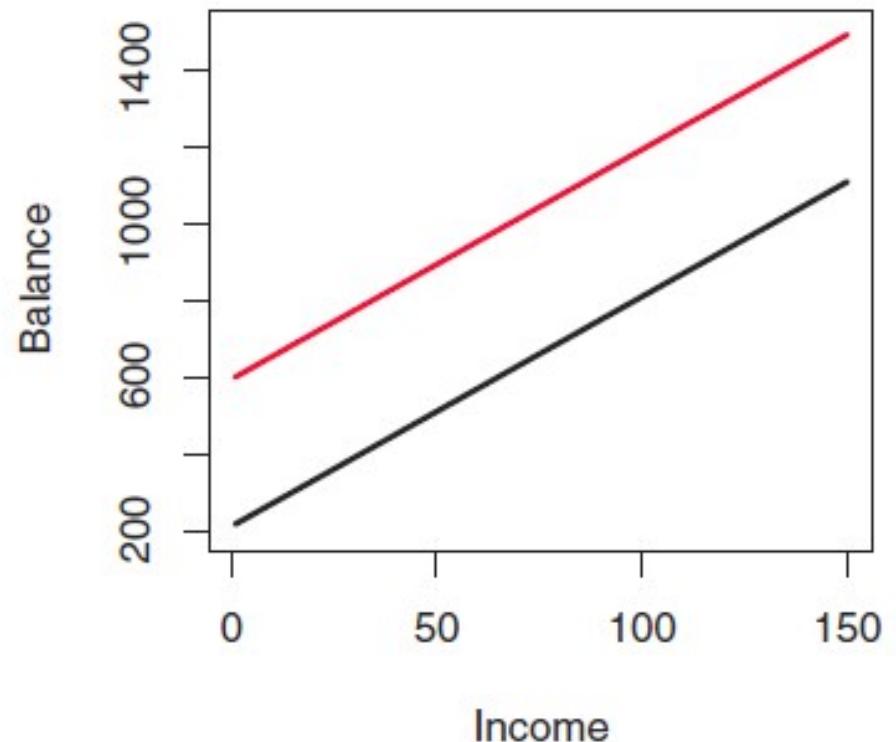


Interpretation of Coefficients

Understand binary features

β_1 allows for a different y-intercept for students

$$y_i = \alpha + \beta_1 \{1 \text{ if student}\} + \beta_2 \text{income} + \epsilon_i$$

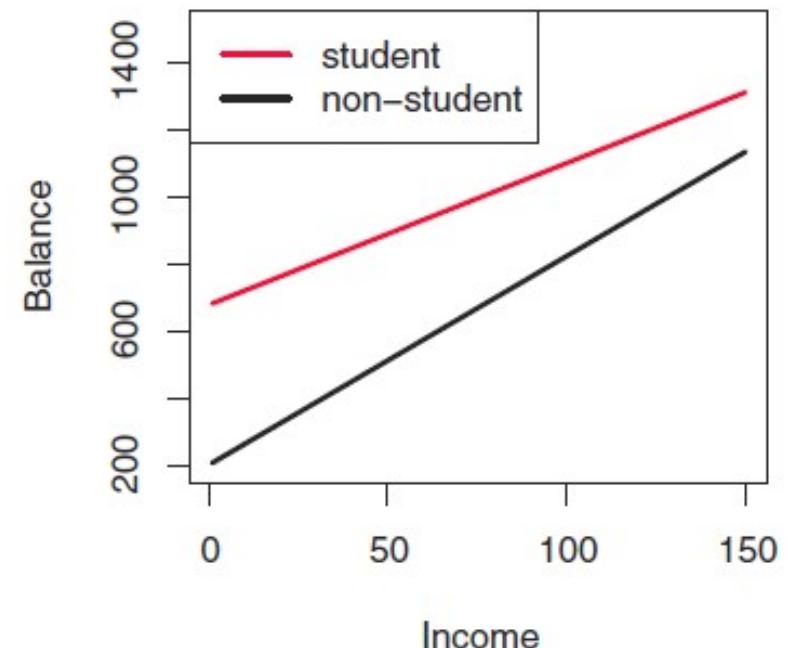


Understand binary features

β_1 allows for a different y-intercept for students

β_3 allows for a different Slope for students

$$y_i = \alpha + \beta_1 \{1 \text{ if student}\} + \beta_2 \text{income} + \beta_3 \{1 \text{ if student}\} * \text{income} + \epsilon_i$$



Interaction terms

Interaction term: multiplying two features by each other. When done as such:
 $\{\text{continuous feature}\} * \{\text{binary feature}\}$

it allows for a different slope for the group represented by the binary feature. Note, be sure to include the continuous feature without the interaction as well.

Example, suppose temperature has a different impact on # of bikeshare trips taken on weekdays vs. weekends. We can create a binary variable `is_weekend`.

If we add `is_weekend` to the model, we allow for a different baseline ridership on weekends. If we add `is_weekend * temp`, we allow temperature to have a different effect for weekends.

We can “interact” `is_weekend` with a polynomial in `temp`. This effectively gives us a new model for temperature for the weekends vs. weekdays

Feature explosion

If we have two raw features, A and B, how many models can we make?

Without interactions (4): none, just A, just B, A + B

With interactions (8): A*B, A+B + A*B, A + A*B, B+ A*B

If we have p features then we have 2^{p+1} possible models!

$$p = 29 \rightarrow 1,073,741,824$$

We cannot possibly run all these models... so we'll learn methods to guide us.

Human intelligence is often a great guide as well.

Interpreting regression output

`Summary(my_model)` gives the coefficient estimates (beta's), t-statistics, standard errors and p-values.

t-statistics evaluate the hypothesis that the coefficient is equal to zero. Thus t greater than 1.96 in absolute value allows us to reject this hypothesis at the 95% level. P-values simply convert t-stats into the probability we would get something this far from zero due to sampling chance alone.

Note that t-statistics evaluate features "one by one". Overall model fit is a better guide for explanatory power (e.g. we might be better for leaving insignificant features in sometimes)

t-stats can be a good guide to throw out irrelevant features

Parametric vs. non-parametric regression

$$Y = f(X) + \epsilon$$

Parametric: f defined by a model we write down with *parameters that we have to estimate* (e.g. $\alpha, \beta_1, \text{etc.}$)

Non-parametric: we directly fit the data

Local averaging

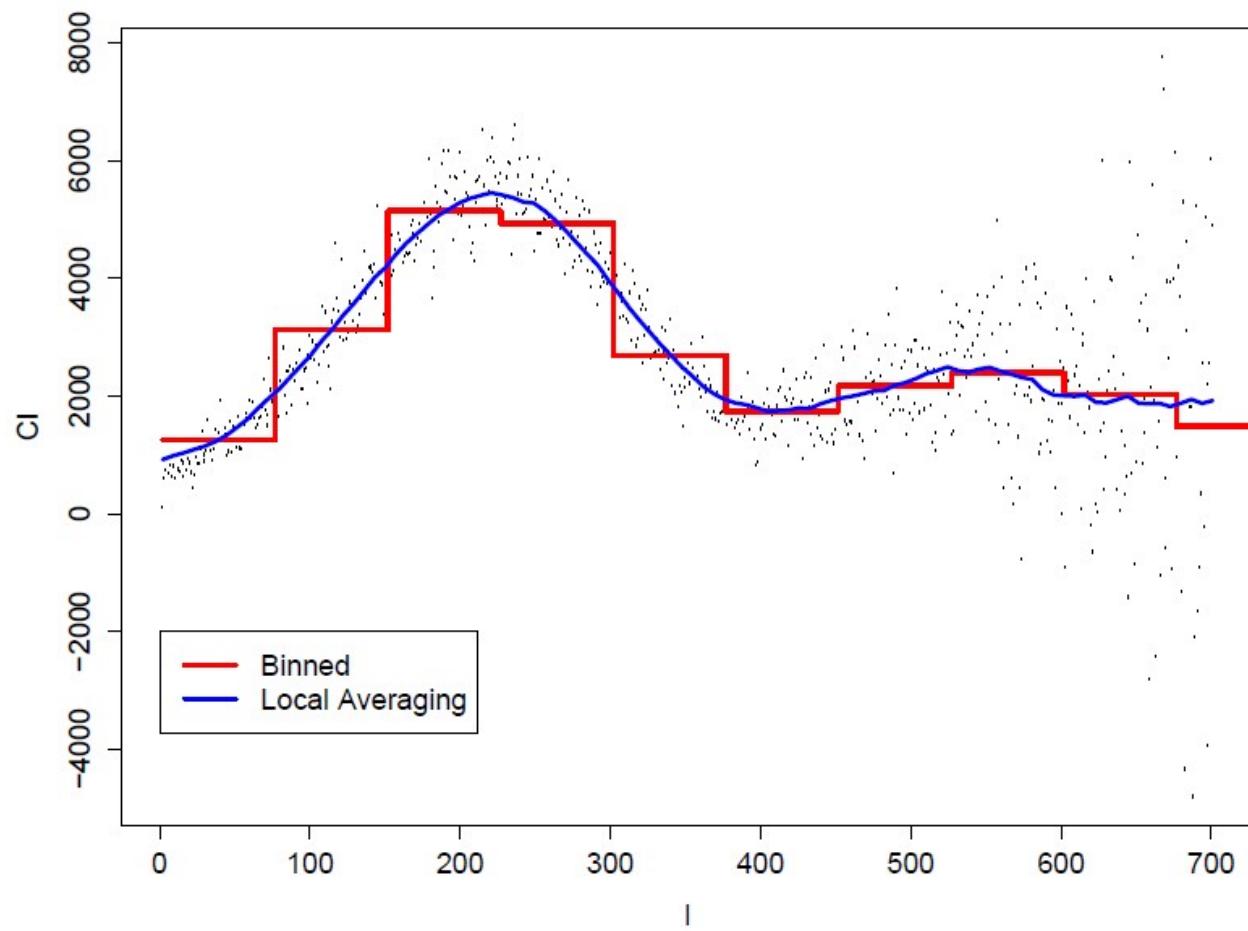
For each x let

$$N_x = \left[x - \frac{h}{2}, x + \frac{h}{2} \right].$$

This is a moving window of length h , centered at x . Define

$$\hat{f}(x) = \text{mean} \left\{ Y_i : X_i \in N_x \right\}.$$

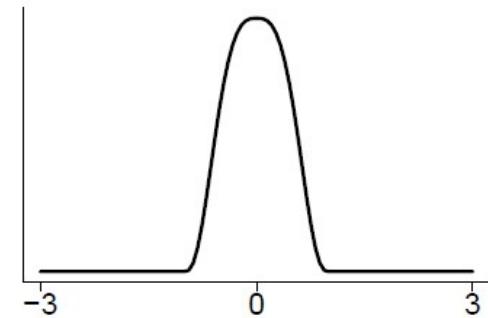
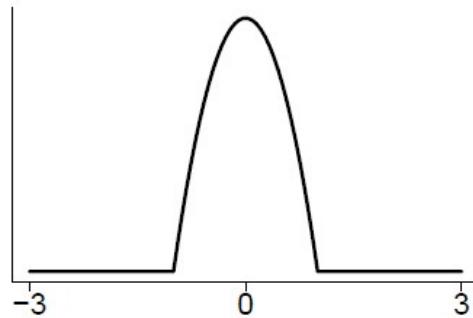
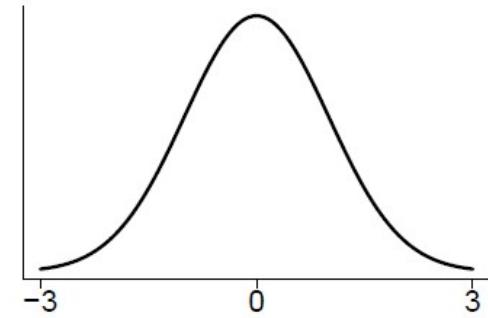
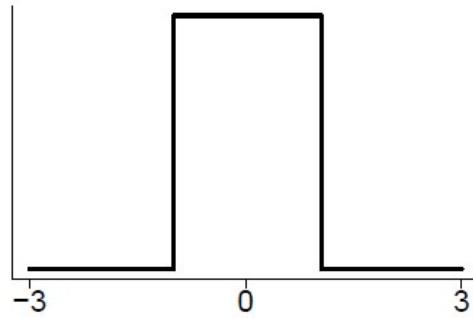
Local averaging vs. bins



Kernels

A method of assigning higher weight to the points closer to the target point I am trying to fit.

Choice of kernel usually does not matter much.



Bandwidths

Bandwidths tell us how wide to make our kernel (window)

Larger bandwidths → smoother functions because there is more data used to make the averages

Bandwidth is sometimes called the smoothing parameter.

Bandwidths matter!

K-nearest neighbors

An alternative to bandwidths is instead of specifying a fixed window, we can say “use the nearest K nearest neighbors”

This ensures each bin has the same amount of data and can be a useful tool

Functionally it will often be quite similar to using a bandwidth, except when there are “sparse data” issues (then K-NNN is preferred)

Often expressed as a fraction of data. E.g. NN=.1 says “each bin is 10% of data”

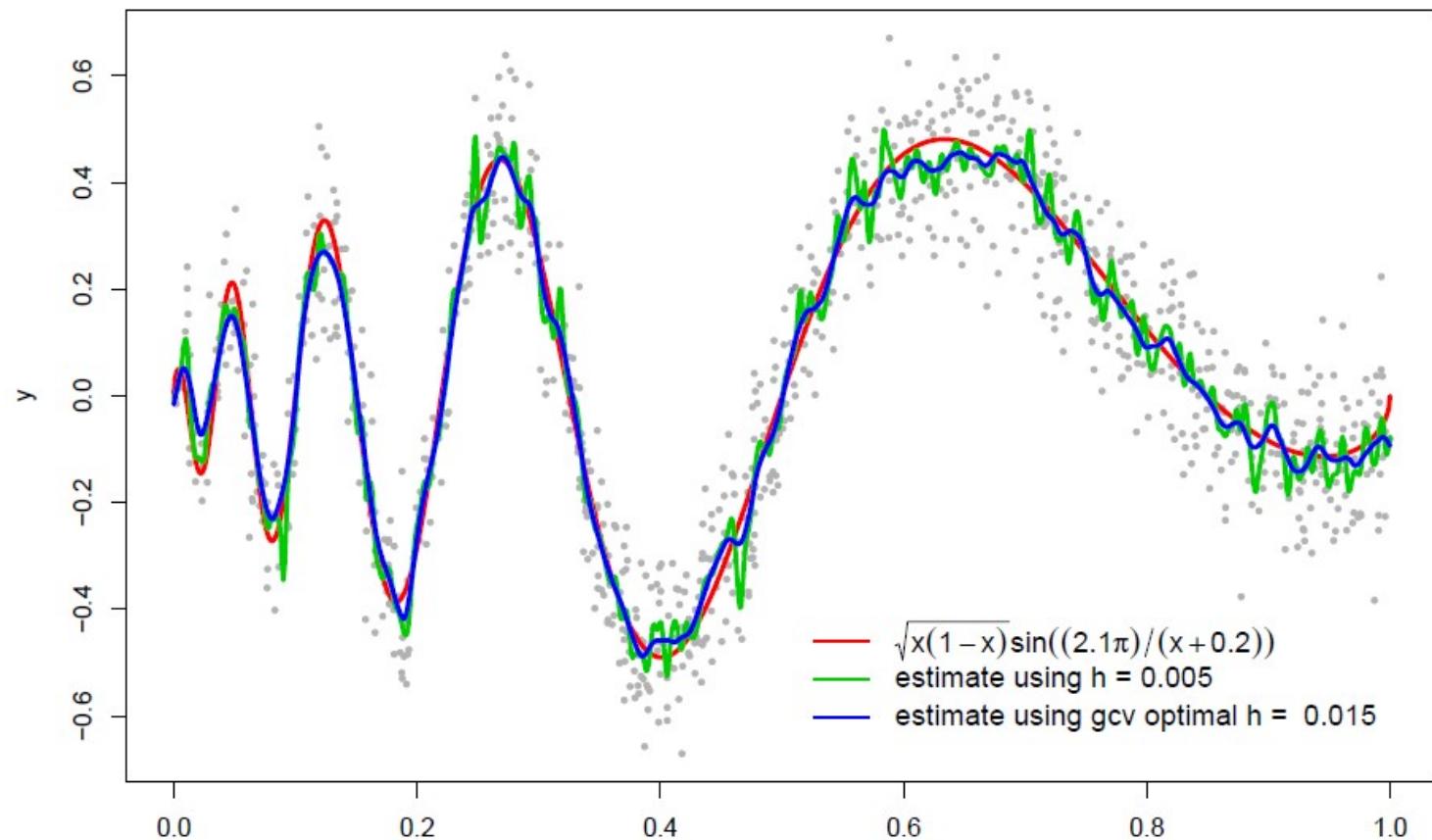
Local polynomial regression

Thing of local averaging as fitting a constant for a given window of data

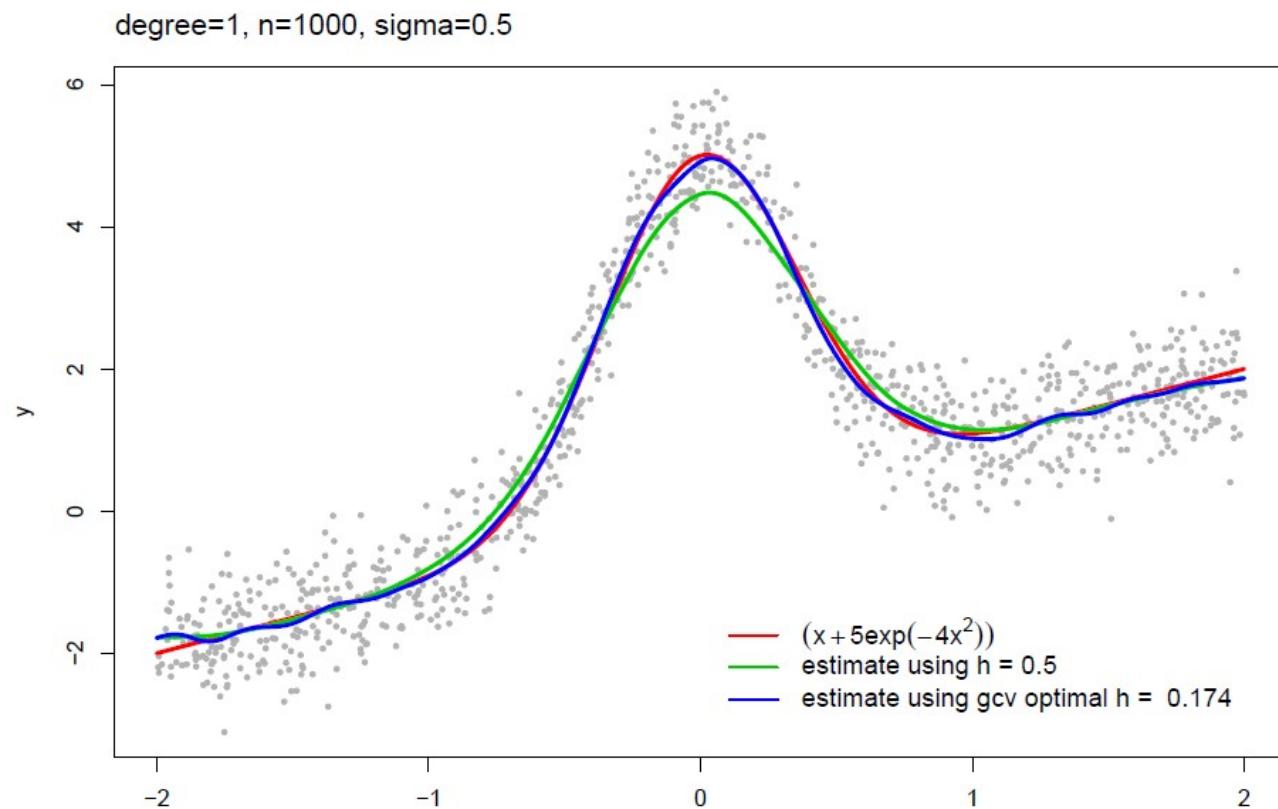
Instead of fitting a constant, we can fit a line ($y=a*x + b$) or a polynomial ($y=a*x + b*x^2 + c$), etc.

It is thus more general than local averaging, but very similar

Some examples



Some examples



Local polynomial in R locfit

```
install.packages("locfit")
library(locfit)
```

Modeler needs to specify: bandwidth (or fraction of nearest neighbors), degree of the local polynomial

Contrast: in parametric regression, we had to write down a model

Aside: we'll learn that non-parametric methods struggle in higher dimensions

When to use each model?

If there are a relatively small number of features (e.g. <4) then non-parametric makes a lot of sense.

With many features, the “curse of dimensionality” sets in and non-parametric methods fall apart (the “neighborhood” is generally empty)

Parametric models using interaction terms and polynomials are the preferred method with many features

“Semi-parametric” methods combine elements of both