

Lecture 13 (ch.3)

Summary

Given data (x_i, y_i) , minimize $SSE = \sum_i (y_i - \hat{y}_i)^2 \Rightarrow \hat{y} = \hat{\alpha} + \hat{\beta} x + \dots$

We will add things here, below

The ANOVA decomposition in regression:

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 = SSE.$$

$SST = SS_{exp.} + SS_{unexp.}$

OLS estimates of α, β .

Interpretation: Regression "reduces" the variability in y from SST to SSE, by accounting for (filtering out) the variability due to x .

$R^2 = \frac{SS_{expl}}{SST} \sim$ percent variability in y due to (or explained by) x .

$s_e = \sqrt{\frac{SSE}{n-2}} \sim$ typical variability in y NOT due to x .
(ie. typical error).

It may help to see ANOVA "By hand": "By R": $\begin{cases} \text{lm.1} = \text{lm}(y \sim x) \\ \text{anova}(\text{lm.1}) \end{cases}$

$\text{lm.1} = \text{lm}(y \sim x)$

$\text{sum}((\text{predict}(\text{lm.1}) - \text{mean}(y))^2) = SS_{expl.} = SS_{\text{model}}$

$\text{sum}((y - \text{predict}(\text{lm.1}))^2) = SS_{unexp} = SS_{\text{Residual}}$

Q: What's the "best" number for predicting y ?

— if we have data on variable y , only.

A: Sample mean of y , ie. \bar{y} report $\bar{y} \pm s_y$ ← We will improve on these, later. in ch.11

— if we also have data on x , which is related to y .

A: The fitted value $\hat{y}(x) = \hat{\alpha} + \hat{\beta} x$ report $\hat{y}(x) \pm s_e \leftarrow s_y$

← Here, the prediction depends on x

A Couple of more comments about ANOVA:

⇒ When we write $SST = SS_{\text{expl}} + SS_{\text{unexpl}}$,

these SS quantities are generally formatted in an ANOVA Table. Look at p.124 and learn how to read the outputs to identify what you need. For example, some computer outputs may call R^2 , Coeff. of determ., or r^2 , R-sqd, ... Also, they may give RMSE, instead of Se :

$$Se = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2}$$

Root → funny mean error ← squared

see examples in Lab.

⇒ Why is R^2 written as R^2 (or even r^2 , as in our book) +

IMPORTANT! R^2 — not a square of anything; at least not generally.
= symbol.

as in our / many books

To see why it is written as R^2 (or even r^2), consider our example: from previous lecture.

$$r = \frac{\overline{xy} - \bar{x}\bar{y}}{\sqrt{(\bar{x}^2 - \bar{x}^2)(\bar{y}^2 - \bar{y}^2)}} \quad \text{or} \quad \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = 0.88916$$

Note $(.88916)^2 = \underline{0.79}$ (see R^2 in prev lecture)

I.e. coeff. of deter. (R^2) = $(r)^2$

But only in simple linear regression.

i.e. $y = \alpha + \beta x$.

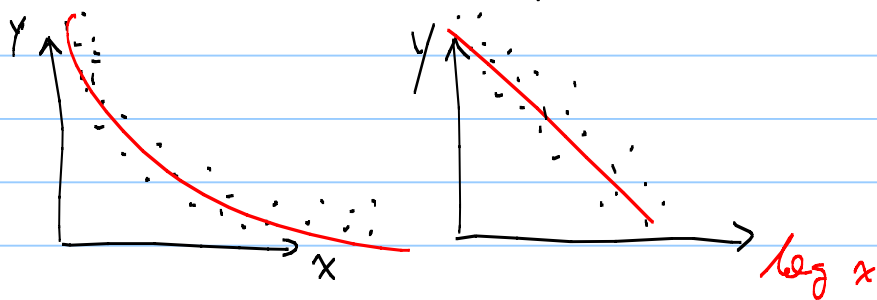
In everything else we will do next, $R^2 \neq (r)^2$.

Non linear relations

So far, we've considered situations where x & y are linearly related. If the relationship (in scatterplot) is nonlinear, then 2 options:

1) If monotonic, then transform data:

For example, $x \rightarrow \log(x)$ often straightens scatterplots that look like this



Then, we do regression on y vs. $\log(x)$.

I.e. $y = \alpha + \beta(\log x)$ not $y = \alpha + \beta x$

and decompose (i.e. Anova) as before.

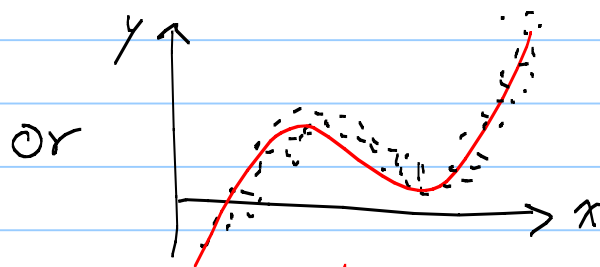
$$SST = SS_{exp} + SS_{unexp} \quad \text{explained by } \log(x).$$

Usually, one of the following transformations straightens the scatterplot:

$\log x$, e^x , \sqrt{x} , $(x)^{1/3}$, same for y .

The best rule is to try different ones, and check the scatterplot.

2) If the relationship is not monotonic?



quadratic

$$\hat{y} = \alpha + \beta_1 x + \beta_2 x^2$$

cubic

$$\hat{y} = \alpha + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

These are examples of polynomial regression.

in R $\text{lm}(y \sim x + I(x^2) + I(x^3) + \dots)$
 for R reasons.

$\hat{\alpha}, \hat{\beta}_1$, etc. are obtained as before, by $\frac{\partial}{\partial \alpha}, \frac{\partial}{\partial \beta_1}, \frac{\partial}{\partial \beta_2}$

See
hw

As in simple linear regression, we can still decompose The total variability in y into explained and un-explained, \dots, R^2, S_e, \dots

The only difference is that

$$SST = SS_{\text{expl}} + SS_{\text{unexpl}}^{\text{SSE}}$$

$R^2 \neq r^2$

$$S_e^2 = \frac{1}{n - (k+1)} \sum (y_i - \hat{y}_i)^2$$

of β 's \nearrow The α .

Note that with The same basic ideas we have learned so far, we can now fit (almost) any data.

But, These tools are dangerous (in the wrong hands).
 In particular, look below, for "overfitting."

Summary

When you see data on (x, y)

→ Look at their scatter plot (and histograms, and ...)

→ If linear, do $y = \alpha + \beta x$

Assess performance with ANOVA (R^2 , se, ^{skipped} residual plots, ...)

→ If non linear, but monotonic,

Then transform x and/or y . E.g. $y = \alpha + \beta \log x$

Assess performance with ANOVA.

→ If non monotonic, Then polynomial regression:

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots$$

Assess performance with ANOVA.

→ Extrapolate Cautiously! Remember The -755 pound person!

→ And Do not overfit! It's bad for predictions (See below)

Learn how to manipulate eqns like These:

$$\text{E.g. } y = \alpha + \beta \ln x$$

$$\hookrightarrow (y - \alpha) / \beta = \ln x \Rightarrow x = e^{(y - \alpha) / \beta}$$

$$\hookrightarrow y = \ln e^\alpha + \ln x^\beta = \ln(e^\alpha x^\beta) \Rightarrow e^y = e^\alpha x^\beta$$

Additive/multiplicative errors:

Additive $y = \alpha + \beta x + \epsilon$

Mult. $y = \alpha e^{\beta x} + \epsilon \rightarrow \ln y = \alpha + \beta \ln x + \epsilon$

So a problem with multiplicative errors can be handled by doing linear regression on the log of all data.

Some details of The above summary

Given data $(x_i, y_i) \quad i=1, 2, \dots, n$

assume $y = \alpha + \beta_1 x + \beta_2 x^2 + \dots$ (or transform, \sqrt{x} , $\log y$, ...)

which means $y_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \dots + \epsilon_i$

minimize $SSE = \sum_{i=1}^n \epsilon_i^2$

to get $\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \dots$ OLS estimates of $\alpha, \beta_1, \beta_2, \dots$

predict: $\hat{y}(x) = \hat{\alpha} + \hat{\beta}_1 x + \hat{\beta}_2 x^2$ OLS fit to data.

ANOVA: $S_{yy} = SST = SS_{exp} + SS_{unexp}$

$$R^2 = \frac{SS_{exp}}{SST}$$

\sim goodness of fit.

$$S_e = \sqrt{\frac{SS_{unexp}}{n-(k+1)}} = \text{std. dev. of errors.}$$

\sim typical error.

Notation:

I say $\hat{\alpha}, \hat{\beta}, \hat{y}, SSE$, book says $a, b, \hat{y}, SS_{Resid}$ (and SSE)

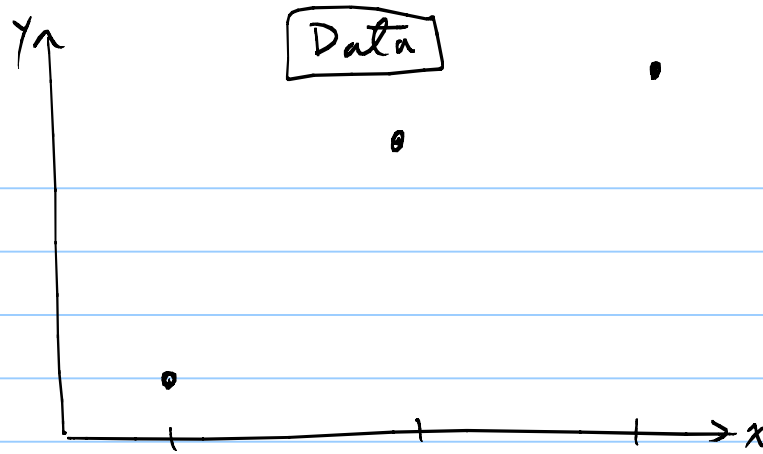
FYI

It may seem like The idea of minimizing SSE should be equivalent to maximizing SS_{exp} , because $SS_{exp} + SSE = \text{Constant (SST)}$.

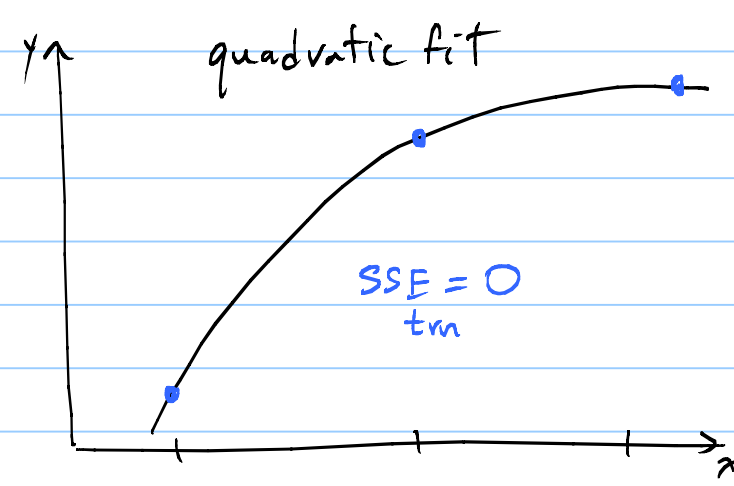
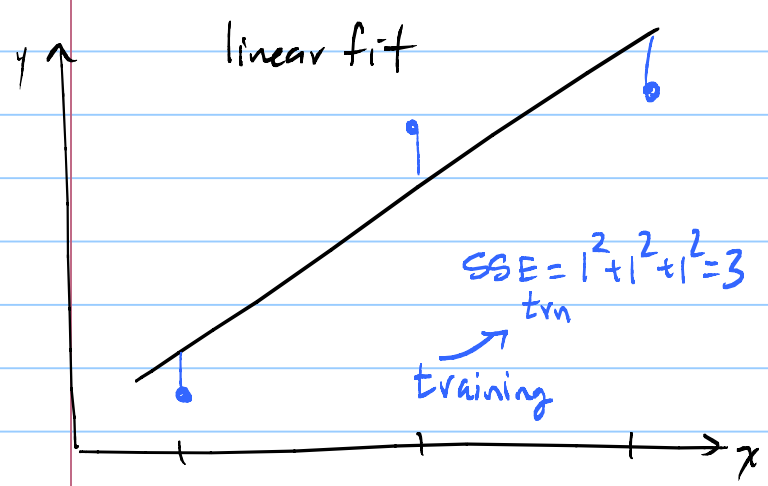
However, That argument is flawed because The decomposition $SS_{exp} + SSE = SST$ is itself true only when SSE is minimized. If SSE is not minimized, Then There is a non-zero cross-term.

\Rightarrow With The tools you have learned now, you can fit any data (x, y) regardless of how complicated (nonlinear) The relationship between x & y .

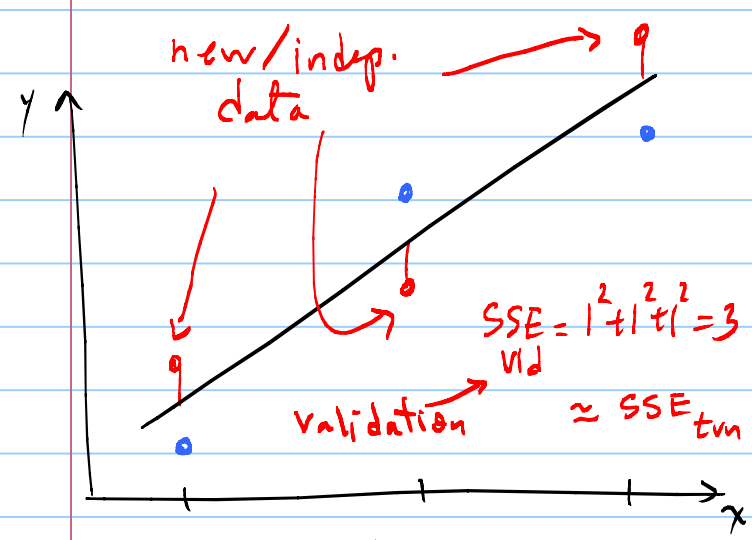
However, powerful tools can be dangerous \searrow



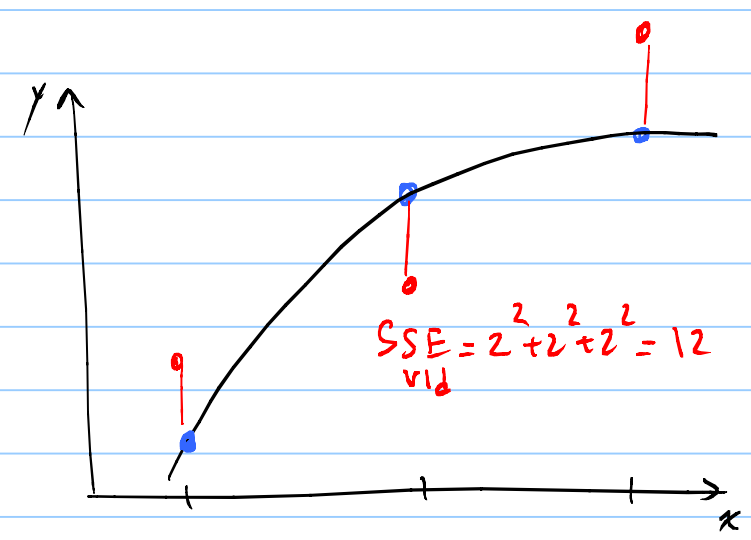
overfitting



If The true relationship between x, y is linear, Then new/independent data will be approximately random about The linear fit, e.g. red dots.



predictions are good

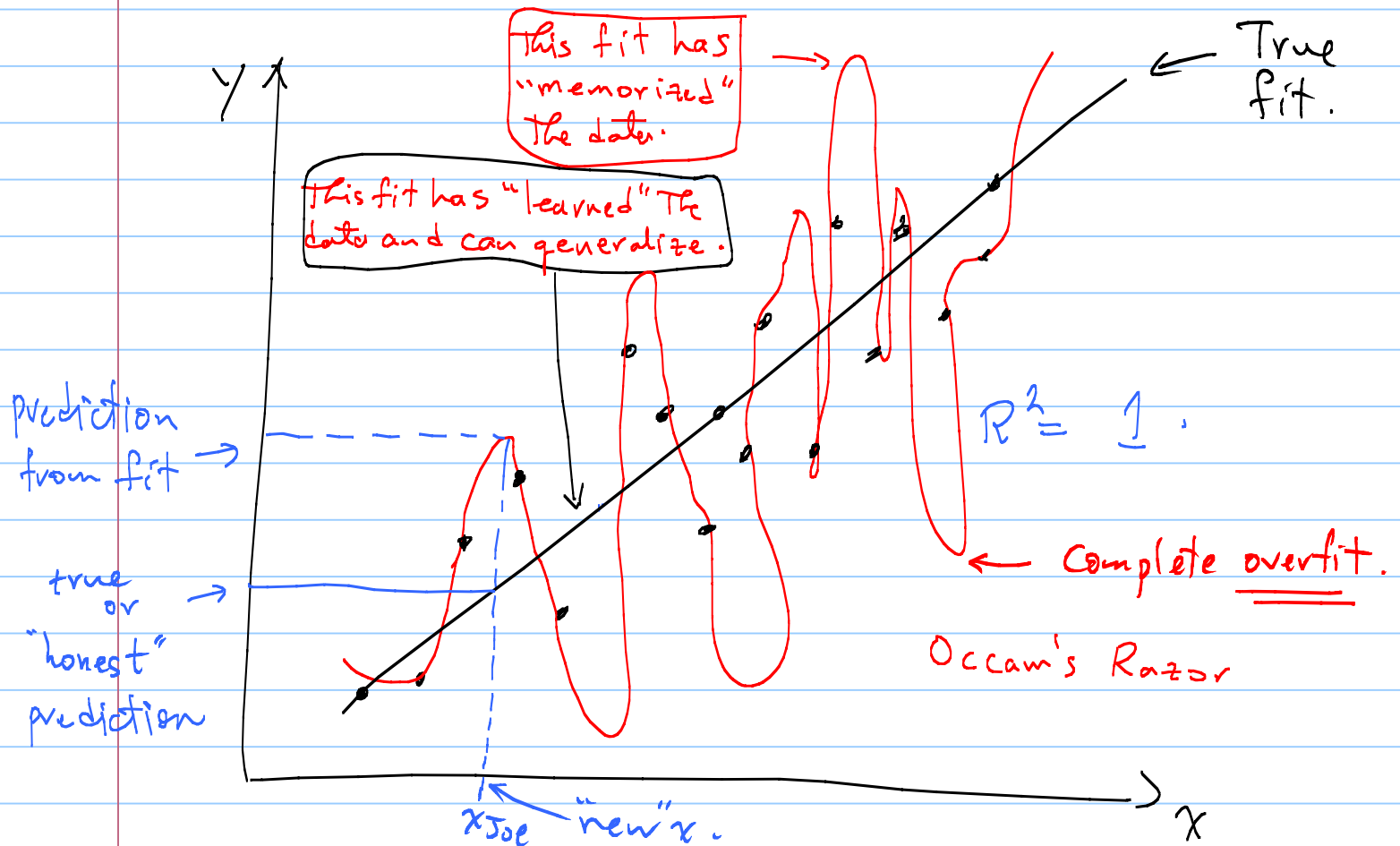


predictions are poor.

ie. it overfit the trn set.

⌈ The quadratic model "memorized" the training data ($SSE_{trn} = 0$), but did not learn/generalize.

The picture we usually see for overfitting is something like this: But this is a bit wrong because it assumes a "new" case has a new x value. As you can see above, overfitting can occur even when the set of x values is fixed (across old and new data).



A Overfitting will yield a model that does really well on the data used for developing the model (that data is called the training set) but it will do poorly on "new" data (called the test/validation set) slightly different but never mind

Overfitting is a "gray" thing. The above red curve shows complete overfitting. But lower levels of overfitting are still bad (in predictions). Moral: Don't overfit!

Q: How will you know if/when you have overfit? Hard Question

A: Try testing your model on indep./new data.

Google "cross-validation" or "bootstrap"! Check Lab.

hw_lect13_1

- Read the data file transform_dat.txt from the course website into R, and
- Make a scatterplot of y vs. x.
- Transform x and/or y to linearize the relationship.
- Perform regression on the transformed data, i.e., do (lm),
- Overlay the corresponding line on the scatterplot
- What percentage of the variability in the transformed y is explained by the transformed x, and what's the typical error in the prediction of the transformed y.

hw_lect13_2

The procedure for estimating the regression coefficients in polynomial regression is the same as before, i.e. by minimizing MSE with respect to alpha, beta_1, beta_2, Each derivative leads to a linear equation, and the system of equations can be uniquely solved to get alpha_hat, beta_1_hat, etc. For this hw, consider a quadratic regression, and derive the linear equations that must be satisfied by alpha_hat, beta_1_hat, etc. Write these equations in terms of the following means:

Do not solve the system of equations.

$$\bar{x}, \overline{x^2}, \overline{x^3}, \overline{x^4}, \overline{xy}, \overline{x^2y}, \bar{y}$$

hw_lect13_3 (By R)

Return to the data you collected at the start of the quarter. Call the 2 continuous variables x and y, depending on which variable you want to predict from the other, and

- Perform simple linear regression to estimate the regression coefficients, and interpret them.
- Draw the regression line on the scatterplot of y vs. x
- Compute R^2 and interpret it
- Compute s_e and interpret it
- Do you need to consider polynomial regression? Or transforming variables? If so, do it!