# Chapter 2

*DJM*

*4 February 2020*

## What is this chapter about?

Problems with regression, and in particular, linear regression

A quick overview:

1. The truth is almost never linear.
2. Collinearity **can** cause difficulties for numerics and interpretation.
3. The estimator depends strongly on the marginal distribution of $X$.
4. Leaving out important variables is bad.
5. Noisy measurements of variables can be bad, but it may not matter.

## Asymptotic notation

- The Taylor series expansion of the mean function $\mu(x)$ at some point $u$

$$\mu(x) = \mu(u) + (x - u)^\top \frac{\partial \mu(x)}{\partial x}|_{x=u} + O(\|x - u\|^2)$$

- The notation $f(x) = O(g(x))$ means that for any $x$ there exists a constant $C$ such that $f(x)/g(x) < C$.

- More intuitively, this notation means that the remainder (all the higher order terms) are about the size of the distance between $x$ and $u$ or smaller.

- So as long as we are looking at points $u$ near by $x$, a linear approximation to $\mu(x) = \mathbb{E}[Y \mid X = x]$ is reasonably accurate.

## What is bias?

- We need to be more specific about what we mean when we say **bias**.

- Bias is neither good nor bad in and of itself.

- A very simple example: let $Z_1, \ldots, Z_n \sim N(\mu, 1)$.

    - We don't know $\mu$, so we try to use the data (the $Z_i$'s) to estimate it.
    - I propose 3 estimators:
        1. $\widehat{\mu}_1 = 12$,
        2. $\widehat{\mu}_2 = Z_6$,
        3. $\widehat{\mu}_3 = \overline{Z}$.
    - The **bias** (by definition) of my estimator is $\mathbb{E}[\widehat{\mu}] - \mu$.
    - Calculate the bias and variance of each estimator.

## Regression in general

- If I want to predict $Y$ from $X$, it is almost always the case that

$$\mu(x) = \mathbb{E}[Y \mid X = x] \neq x^\top \beta$$

- There are always those errors $O(\|x - u\|)^2$, so the **bias** is not zero.

- We can include as many predictors as we like, but this doesn't change the fact that the world is **non-linear**.

## Covariance between the prediction error and the predictors

- In theory, we have (if we know things about the state of nature)

$$\beta^* = \arg\min_{\beta} \mathbb{E}\left[\|Y - X\beta\|^2\right] = \text{Cov}\left[X,\ X\right]^{-1}\text{Cov}\left[X,\ Y\right]$$

- Define $v^{-1} = \text{Cov}\left[X,\ X\right]^{-1}$.

- Using this optimal value $\beta^*$, what is $\text{Cov}\left[Y - X\beta^*,\ X\right]$?

$$
\begin{aligned}
\text{Cov}\left[Y - X\beta^*,\ X\right] &= \text{Cov}\left[Y,\ X\right] - \text{Cov}\left[X\beta^*,\ X\right] &&\text{(Cov is linear)}\\
&= \text{Cov}\left[Y,\ X\right] - \text{Cov}\left[X(v^{-1}\text{Cov}\left[X,\ Y\right]),\ X\right] &&\text{(substitute the def. of } \beta^*)\\
&= \text{Cov}\left[Y,\ X\right] - \text{Cov}\left[X,\ X\right]v^{-1}\text{Cov}\left[X,\ Y\right] &&\text{(Cov is linear in the first arg)}\\
&= \text{Cov}\left[Y,\ X\right] - \text{Cov}\left[X,\ Y\right] = 0.
\end{aligned}
$$

## Bias and Collinearity

- Adding or dropping variables may impact the bias of a model
  - Suppose $\mu(x) = \beta_0 + \beta_1 x_1$. It **is** linear. What is our estimator of $\beta_0$?
  - If we instead estimate the model $y_i = \beta_0$, our estimator of $\beta_0$ will be biased. How biased?
  - But now suppose that $x_1 = 12$ always. Then we don't need to include $x_1$ in the model. Why not?
  - Form the matrix $[1\ x_1]$. Are the columns collinear? What does this actually mean?

## When two variables are collinear, a few things happen.

1. We cannot **numerically** calculate $(\mathbf{X}^\top\mathbf{X})^{-1}$. It is rank deficient.
2. We cannot **intellectually** separate the contributions of the two variables.
3. We can (and should) drop one of them. This will not change the bias of our estimator, but it will alter our interpretations.
4. Collinearity appears most frequently with many categorical variables.
5. In these cases, software **automatically** drops one of the levels resulting in the baseline case being in the intercept. Alternately, we could drop the intercept!
6. High-dimensional problems (where we have more predictors than observations) also lead to rank deficiencies.
7. There are methods (regularizing) which attempt to handle this issue (both the numerics and the interpretability). We may have time to cover them slightly.

## White noise

**White noise** is a stronger assumption than **Gaussian**.

Consider a random vector $\epsilon$.

1. $\epsilon \sim \text{N}(0, \Sigma)$.
2. $\epsilon_i \sim \text{N}(0, \sigma^2(x_i))$.
3. $\epsilon \sim \text{N}(0, \sigma^2 I)$.

The third is white noise. The $\epsilon$ are normal, their variance is constant for all $i$ and independent of $x_i$, and they are independent.

## Asymptotic efficiency

This and MLE are covered in 420.

There are many properties one can ask of estimators $\widehat{\theta}$ of parameters $\theta$

1. Unbiased: $\mathbb{E}\left[\widehat{\theta}\right] - \theta = 0$
2. Consistent: $\widehat{\theta} \xrightarrow{n \to \infty} \theta$
3. Efficient: $\mathbb{V}\left[\widehat{\theta}\right]$ is the smallest of all unbiased estimators
4. Asymptotically efficient: Maybe not efficient for every $n$, but in the limit, the variance is the smallest of all unbiased estimators.
5. Minimax: over all possible estimators in some class, this one has the smallest MSE for the worst problem.
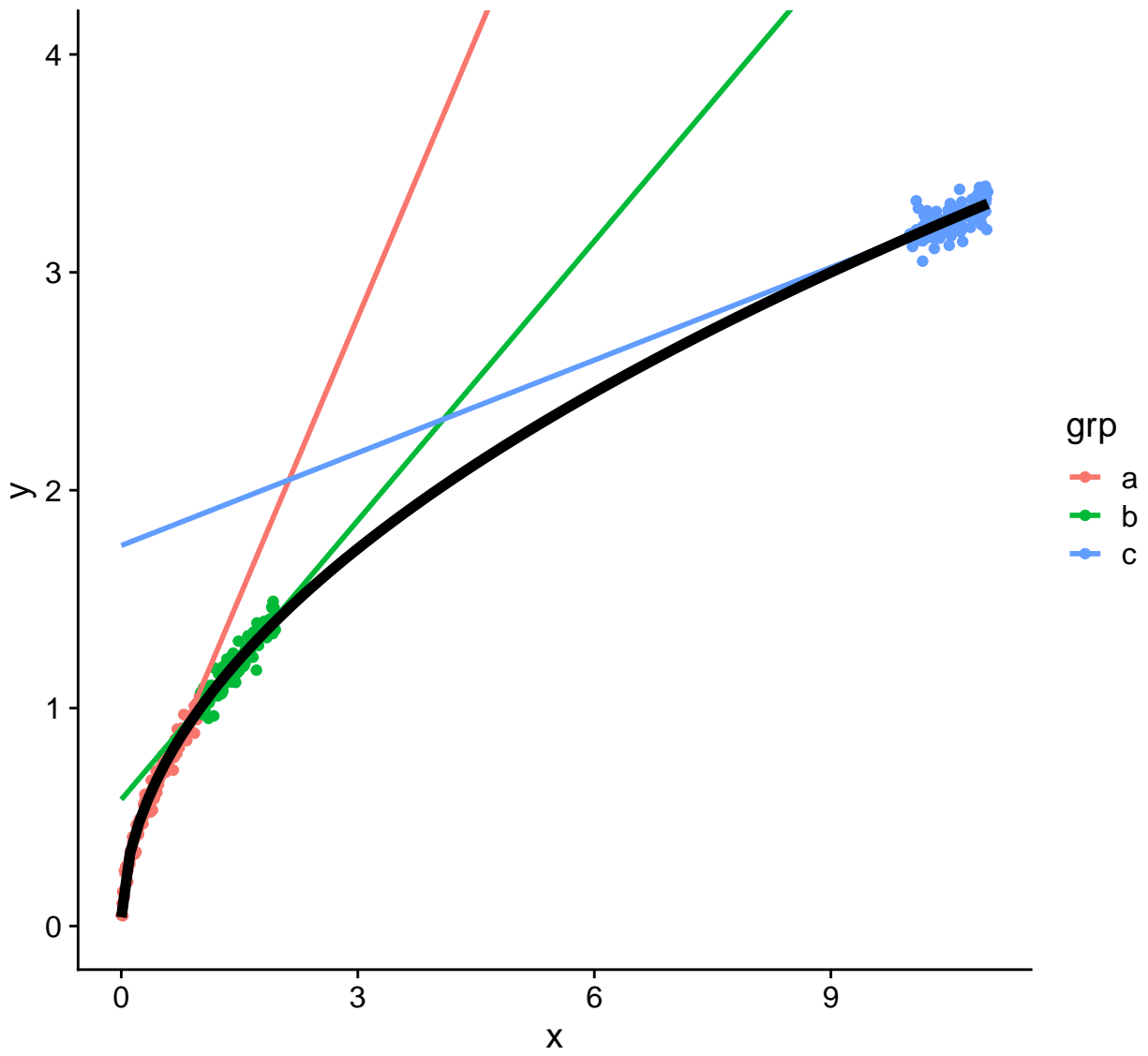6. ...

## Problems with R-squared

$$R^2 = 1 - \frac{SSE}{\sum_{i=1}^{n}(Y_i - \overline{Y})^2} = 1 - \frac{MSE}{\frac{1}{n}\sum_{i=1}^{n}(Y_i - \overline{Y})^2} = 1 - \frac{SSE}{SST}$$

- This gets spit out by software
- $X$ and $Y$ are both normal with (empirical) correlation $r$, then $R^2 = r^2$
- In this nice case, it measures how tightly grouped the data are about the regression line
- Data that are tightly grouped about the regression line can be predicted accurately by the regression line.
- Unfortunately, the implication does not go both ways.
- High $R^2$ can be achieved in many ways, same with low $R^2$
- You should just ignore it completely (and the adjusted version), and encourage your friends to do the same

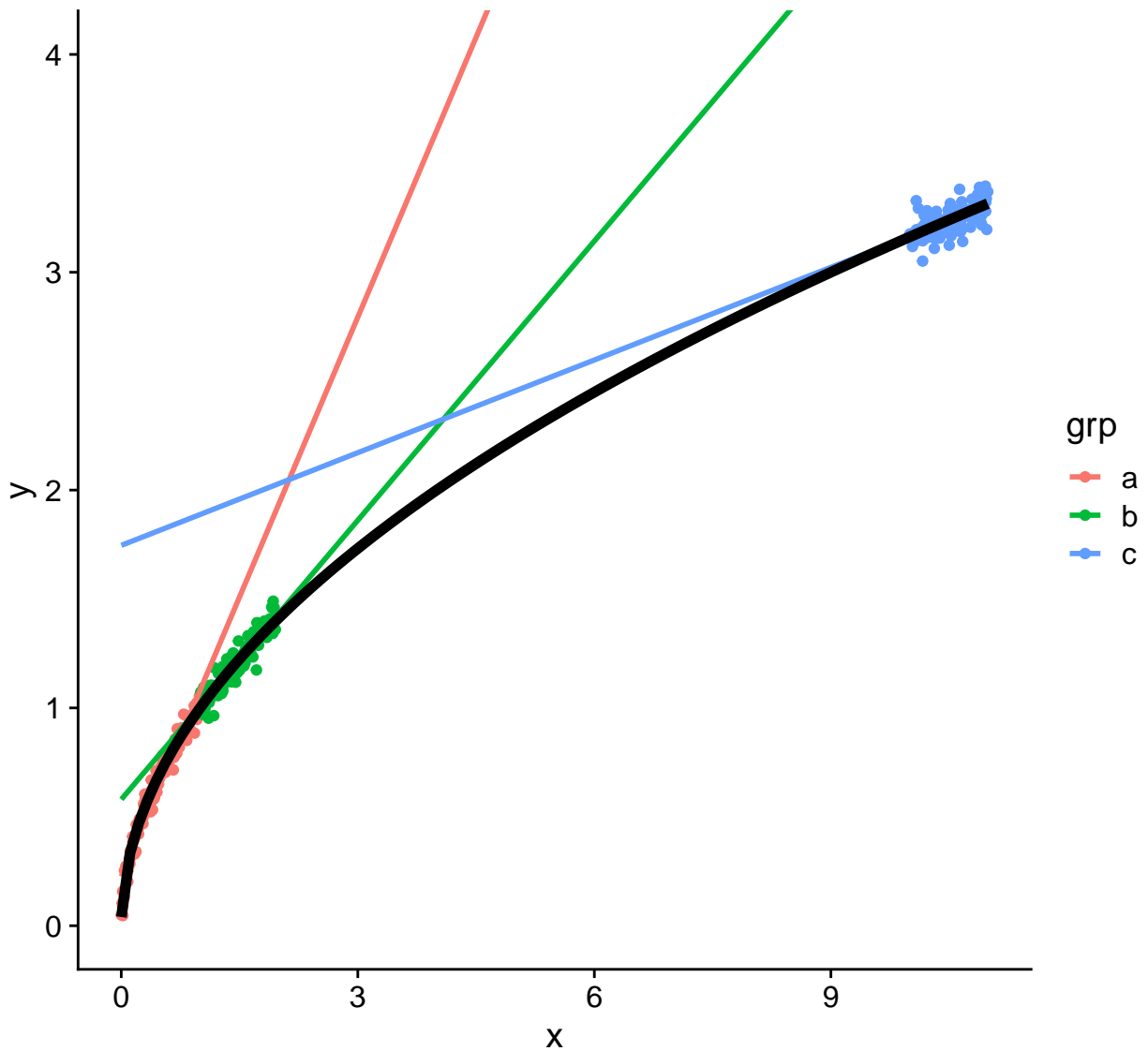## High R-squared with non-linear relationship

```
genY <- function(X, sig) Y = sqrt(X)+sig*rnorm(length(X))
sig=0.05; n=100
df = tibble(
  x = c(runif(n,0,1), runif(n,1,2), runif(n,10,11)),
  y = genY(x, sig),
  grp = rep(letters[1:3], each=n))

g1 = ggplot(df, aes(x, y, color=grp)) + geom_point() +
  geom_smooth(method = 'lm', fullrange=TRUE, se = FALSE) +
  coord_cartesian(ylim=c(0,4)) + stat_function(fun=sqrt,color='black',size=2) +
  theme_cowplot()
g1
```
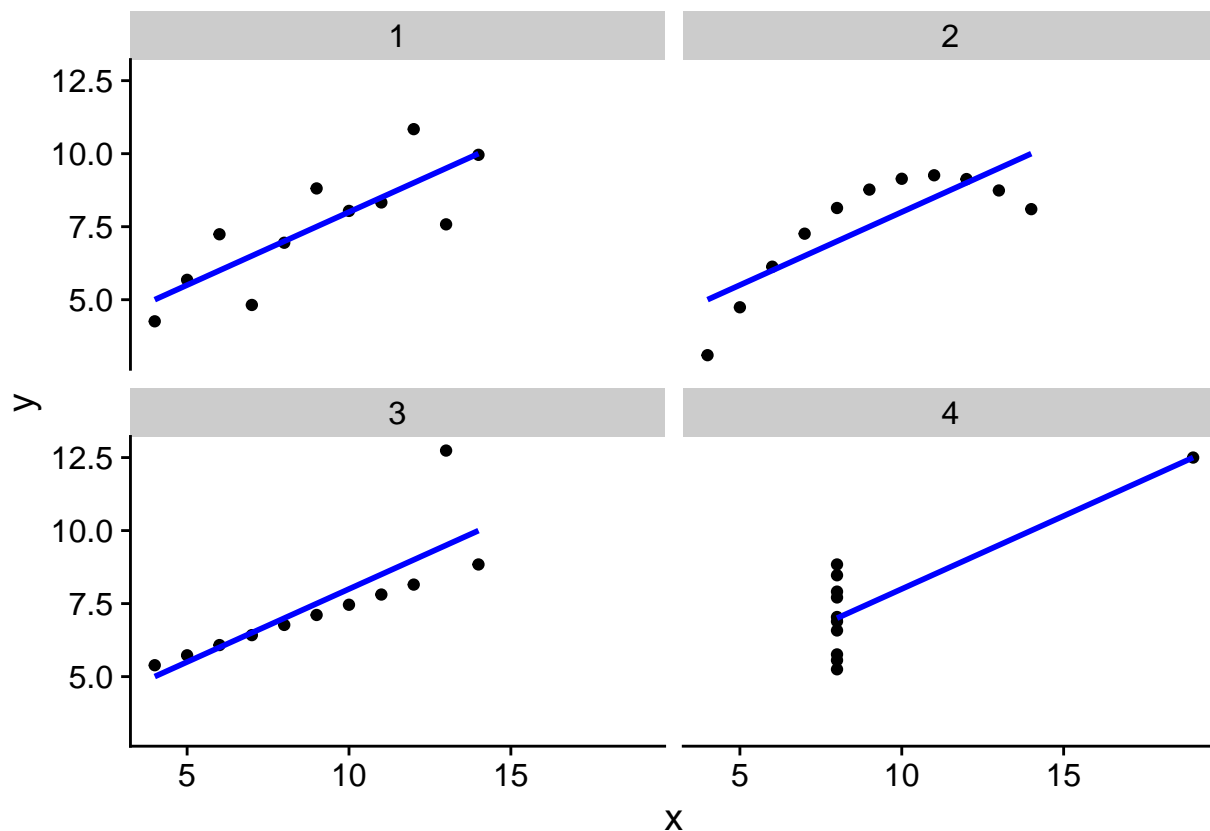
What are the numbers?

g1

```
df %>% group_by(grp) %>% summarise(rsq = summary(lm(y~x))$r.sq) %>% knitr::kable(digits = 2)
```

| grp | rsq |
|-----|------|
| a | 0.93 |
| b | 0.88 |
| c | 0.37 |

## Anscombe's quartet

```
ans = anscombe %>%
 pivot_longer(everything(),
   names_to = c(".value", "set"),
   names_pattern = "(.)(.)"
 )
ggplot(ans, aes(x,y)) + geom_point() +
  geom_smooth(method="lm",se=FALSE,color="blue") +
```

```
facet_wrap(~set, ncol=2) + theme_cowplot()
```



**Anscombe's quartet**

```
ans %>% group_by(set) %>%
  summarise(
    mx = mean(x), my = mean(y), sx = sd(x), sy = sd(y),
    int = coef(lm(y~x))[1], slope = coef(lm(y~x))[2],
    cor = cor(x,y), rsq = summary(lm(y~x))$r.sq) %>%
  knitr::kable(digits = 2)
```

| set | mx | my | sx | sy | int | slope | cor | rsq |
|-----|----|----|----|----|----|----|----|----|
| 1 | 9 | 7.5 | 3.32 | 2.03 | 3 | 0.5 | 0.82 | 0.67 |
| 2 | 9 | 7.5 | 3.32 | 2.03 | 3 | 0.5 | 0.82 | 0.67 |
| 3 | 9 | 7.5 | 3.32 | 2.03 | 3 | 0.5 | 0.82 | 0.67 |
| 4 | 9 | 7.5 | 3.32 | 2.03 | 3 | 0.5 | 0.82 | 0.67 |