

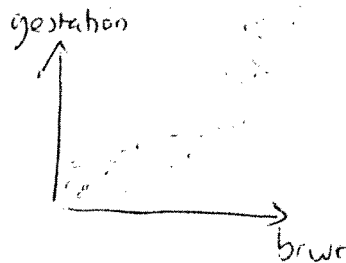
Polynomial Regression

①

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i1}^2 + \beta_3 x_{i1}^3 + \beta_4 x_{i2} + \dots + \varepsilon_i$$

polynomial in x_1

Example



Note - poly-reg is not the same as transform models

$$y_i = \beta_0 + \beta_1 (x_{i1})^2 + \beta_2 x_{i4} + \beta_3 (\log(x_{i5})) + z_i$$

transform x 's prior to model fitting

② a transform model suffices \rightarrow you save some parameters
Comp. poly-reg and avoid possible collinearities
between (x, x^2, x^3, \dots)

Note - poly-reg is not a nonlinear model, since $E(y)$ still linear in the parameters β 's. (and that makes life a lot easier)

Nonlinear models: $y_i = \beta_0 + \beta_1 \exp(\beta_2 x_{i1}) + \varepsilon_i$

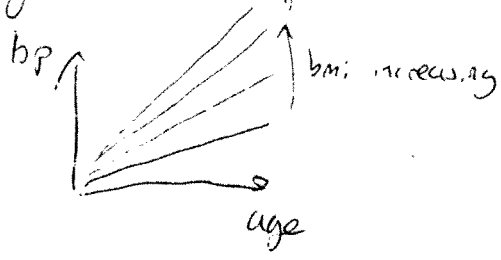
(even $y_i = \beta_0 \exp(\beta_2 x_{i1}) + \varepsilon_i$ since error additive
 $\phi(y) = \beta_0 \exp(\beta_2 x)$ $\log E(y_i) = \log(\beta_0) + \beta_2 x_{i1} \Rightarrow$ Generalized Linear Models)

Interactions

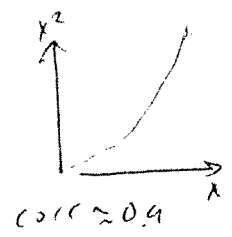
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_{12} x_{i1} x_{i2} + \beta_3 x_{i1}^2 + \beta_4 x_{i2}^2 + \epsilon_i$$

interaction term.

example: age x bmi → effect on blood pressure (bp)



Problems? (x, x^2, x^p) are often highly correlated

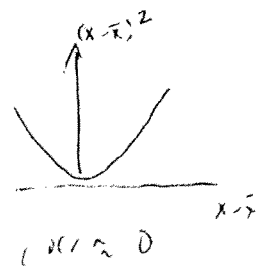


⇒ $(X'X)^{-1}$ unstable, near singular

Is there any way around this?

centering x before fitting the polyreg!

$$\tilde{x}_{i1} = (x_{i1} - \bar{x}_1) \quad \tilde{x}_{i1}^2 = (x_{i1} - \bar{x}_1)^2 \quad \dots$$



The correlation between $\tilde{x}, \tilde{x}^2, \dots, \tilde{x}^p$ much reduced compared to corr between x, x^2, \dots, x^p .

Fit $y_i = \beta_0 + \beta_1 \tilde{x}_{i1} + \beta_2 \tilde{x}_{i1}^2 + \dots + \epsilon_i$

Note → β 's have slightly different interpretation now.

A word of warning

- it's best to keep the order of the polynomial as low as possible
- Always try transform models first.
- Keep it reasonable → limit the order & number of interaction terms.

The idea behind centering, i.e. preprocessing x 's & higher order terms to reduce the correlation can be extended to other "bases" → e.g. splines, wavelets (more later)

Practical

- To fit poly-reg → incorporate $\tilde{x}_{i1}, \tilde{x}_{i2}, \dots$ in \tilde{X} (design matrix)
- $\hat{\beta} = (\tilde{X}'\tilde{X})^{-1}\tilde{X}'y$ as before
- You can use the t-test, F-test, but be aware of remaining collinearities
⋮
- Carefully examine residual plots & rely on data context to determine the order of the polynomial regression model.

Detecting interactions

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + (\beta_3 x_{i1} x_{i2}) + \epsilon_i$$

→ if this is the true model then

(1) if x_2 is held constant & x_1 increases by 1

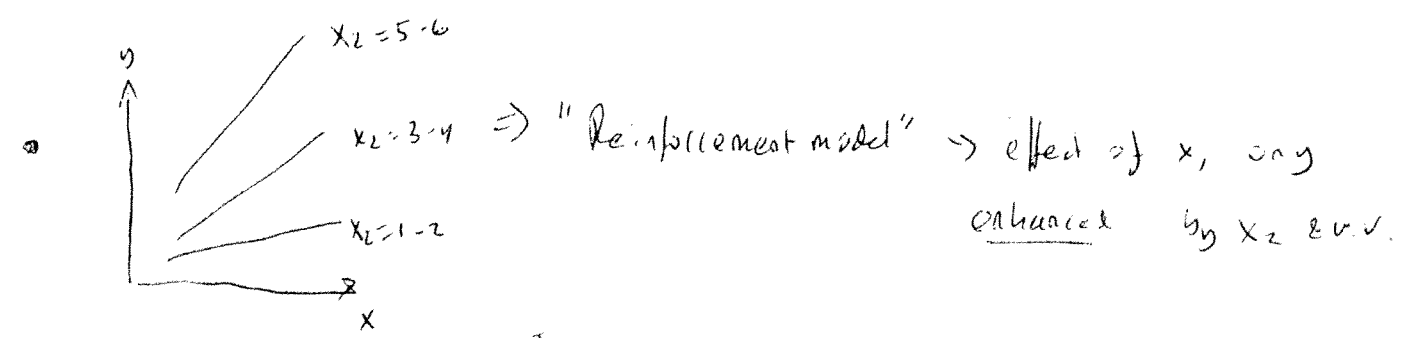
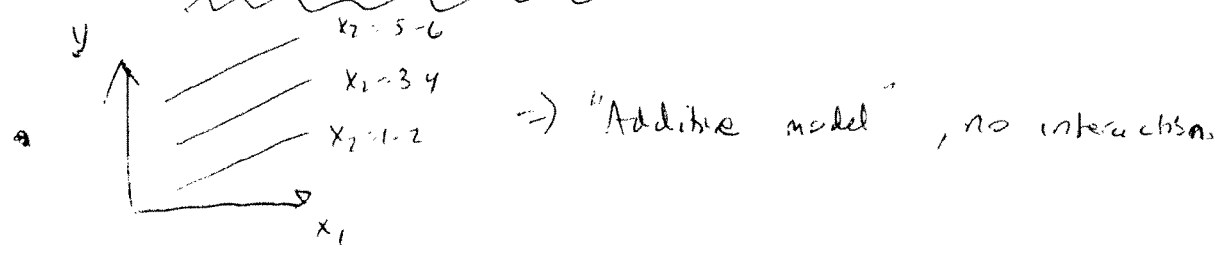
→ expect y to increase by $\beta_1 + \beta_3 x_2$

(2) if x_1 is held constant & x_2 increases by 1

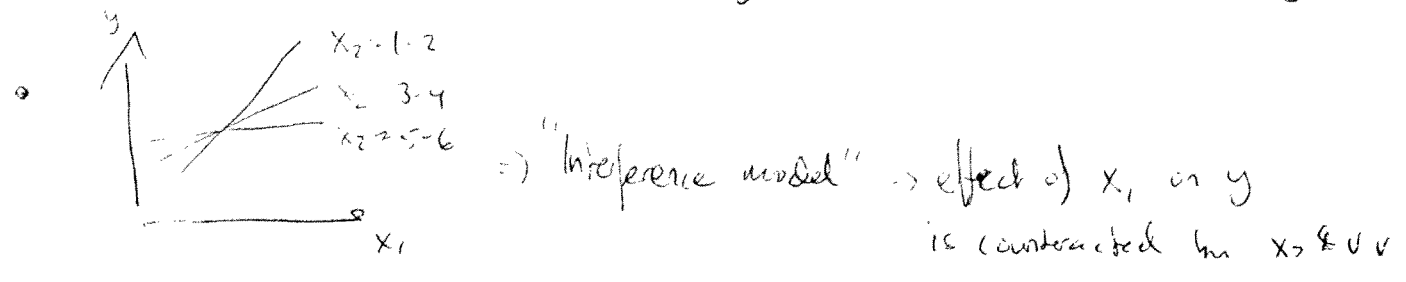
→ expect y to increase by $\beta_2 + \beta_3 x_1$

→ so effect of x_1 on y depends on x_2 & v.v.

Use conditional effects plots to detect



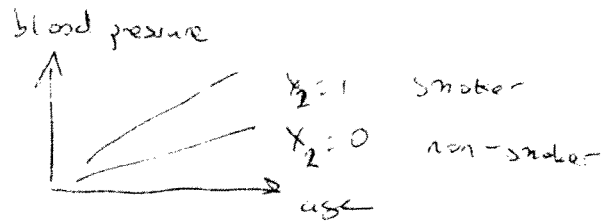
Ex: Salary ~ yr employment * yr of schooling



Qualitative variables

E.g. smoker, non-smoker

$$X_2 = \{0, 1\}$$



It's easy to incorporate qualitative variables in regression

using Indicator variables

$$X_2 = \begin{Bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{Bmatrix}$$

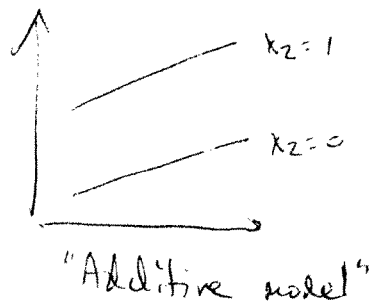
\Rightarrow add to \bar{X}
(the design matrix)



an additive model

$$E(y) = \underbrace{\beta_0}_{\beta_0} + \beta_1 x_1 \quad \text{if } X_2 = 0$$

$$E(y) = \underbrace{\beta_0 + \beta_2}_{\beta_0'} + \beta_1 x_1 \quad \text{if } X_2 = 1$$



So including an indicator in \bar{X} allows for a different intercept - but what if being a smoker/non-smoker changes the relationship between age and blood pressure?

Interactions

\Rightarrow Create a new X_i variable X_i indicator variable and add to \bar{X} .

A few alternatives

① 2 separate regression models

- use 2 indicators

$$\begin{matrix}
 \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix} & \begin{pmatrix} 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\
 X_S & X_{NS}
 \end{matrix}$$

, for smokers & nonsmokers

- create 2 new age variables

$$\text{age} \times X_S = \begin{pmatrix} 0 \\ 0 \\ ? \\ 0 \\ \text{age} \\ \text{age} \\ \vdots \end{pmatrix}, \quad \text{age} \times X_{NS} = \begin{pmatrix} \text{age} \\ \text{age} \\ \text{age} \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

- add the new variables to the design matrix X

- run regression \rightarrow get $E(y) = \beta_0^S + \beta_1^S \cdot \text{age}$; smoker
 $E(y) = \beta_0^{NS} + \beta_1^{NS} \cdot \text{age}$; nonsmoker

But we usually want to make inference about the contrasts between smokers & non-smokers

② One indicator

- baseline model

$$X_2 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{pmatrix}$$

$$, \text{ age} \cdot X_2 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \text{age} \\ \text{age} \\ \vdots \end{pmatrix}$$

→ add to \bar{X}

- run regression & get

$$\begin{cases} E(y) = \beta_0 + \beta_1 \cdot \text{age} & ; X_2 = 0 \\ E(y) = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \cdot \text{age} & ; X_2 = 1 \end{cases}$$

- Now we can test

a) $\beta_3 = 0$ → additive model

b) $\beta_2 = \beta_3 = 0$ → no effect due to smoking.

- Note, this formulation implicitly treats non-smokers as the baseline population for comparisons.

③ Contrast

$$X_2 = \begin{pmatrix} -1 \\ -1 \\ -1 \\ +1 \\ +1 \\ +1 \end{pmatrix}$$

instead of 0 & 1

2 new variables

- $X_2 \cdot \text{age}$

- add the new variables to \bar{X}

2-Var regression

$$\Rightarrow \begin{cases} E(y) = (\beta_0 - \beta_2) + (\beta_1 - \beta_3) \cdot \text{age} & \text{if nonsmoker} \\ E(y) = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \cdot \text{age} & \text{if smoker} \end{cases}$$

Now, can test $\beta_3 = 0$ (additive model) or

$\beta_2 = \beta_3 = 0$ (no effect due to smoking)

and the baseline is now the average of both populations (smokers & nonsmokers)

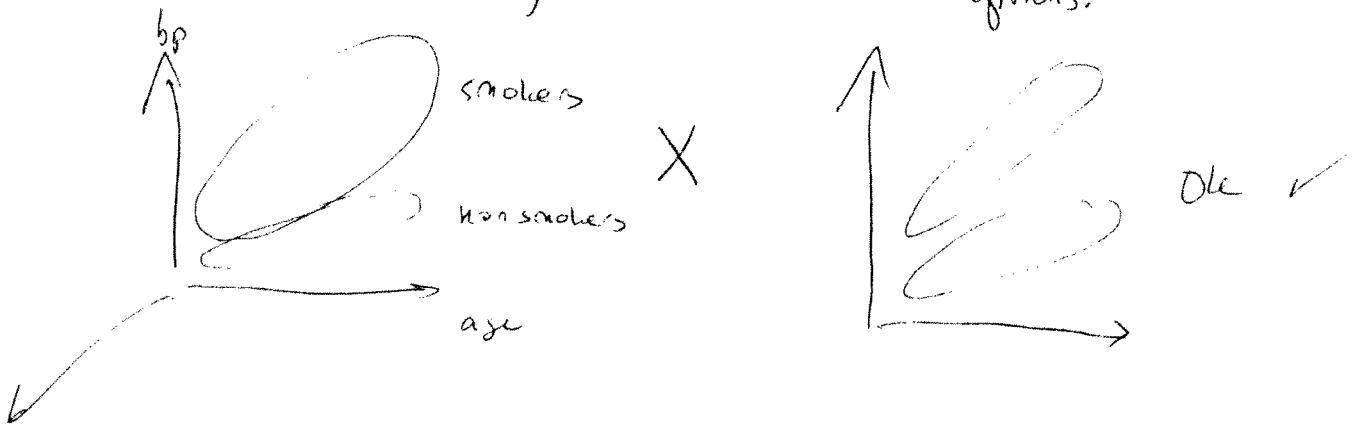
A few words of warning:

9

① It is common practice to always include the main terms if interactions are included - makes the interpretation of the interactions easier to deal with

② The basic assumptions still need to hold - e.g. constant variance.

④ The variance is different for smokers & non-smokers, it is a violation of the basic assumptions.



What to do here?

- a) weighted Least Squares
- b) 2 regression models

