

Lecture 15

Polynomial regression work sheet:

Q1 : Degree 4 polynomial model :

$$Y = C_0 + C_1 X_1 + C_2 X_1^2 + C_3 X_1^3 + C_4 X_1^4$$

Q2 $\vec{Y} = A \vec{C}$

$$Y = C_0 + C_1 X_1 + C_2 X_1^2$$

happiness

icecream

$$\text{icecream} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

$$\text{Chocolate} = \begin{bmatrix} 99 \\ 999 \\ 9999 \\ 1999 \end{bmatrix}$$

$$\vec{Y} = A \vec{C}$$

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} \begin{bmatrix} C_0 \\ C_1 \\ C_2 \end{bmatrix}$$

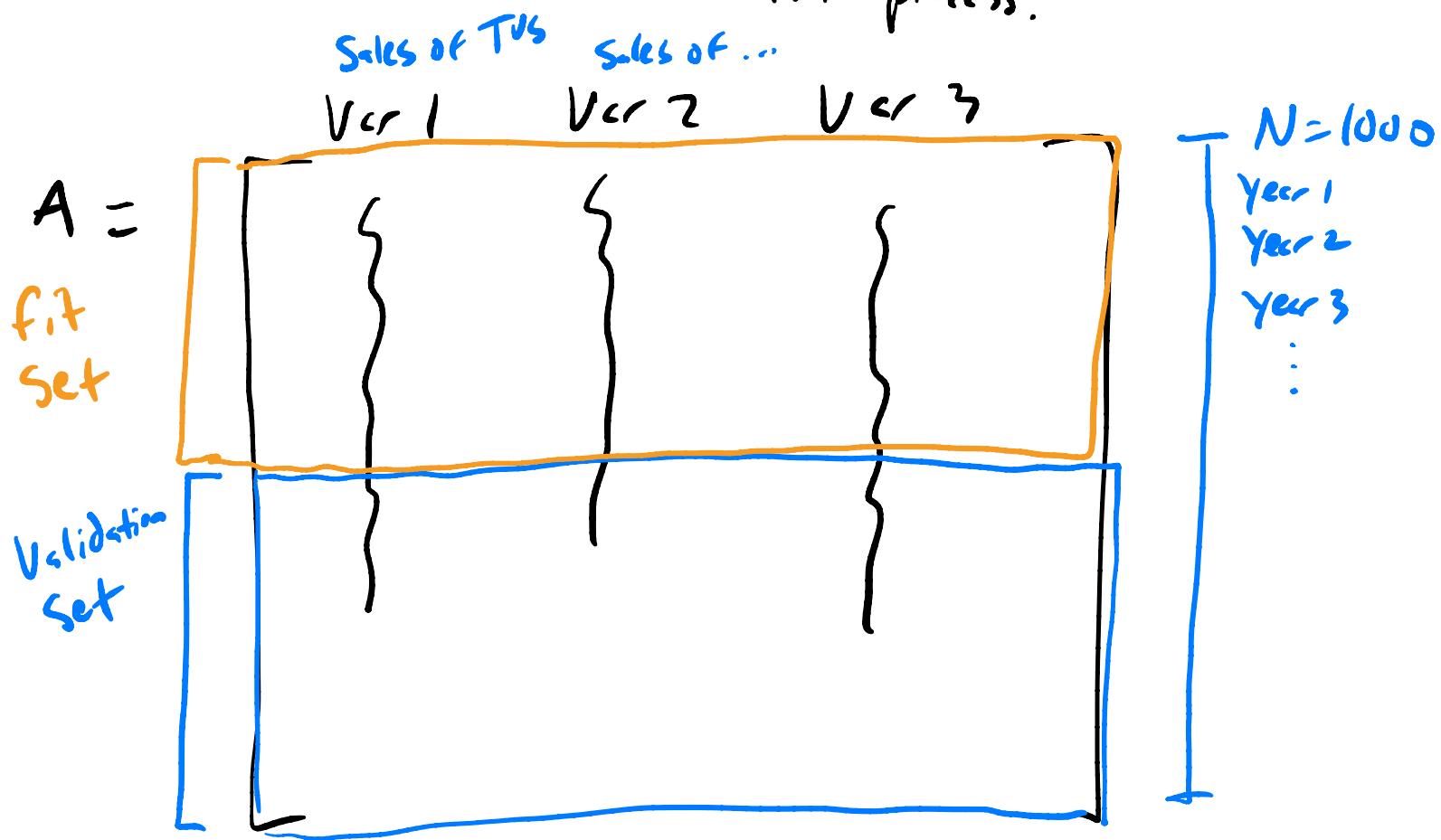
A

$$\text{happiness} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Fit and Validation Sets:

To help detect overfitting / prevent it, you can split your dataset into 2 parts:

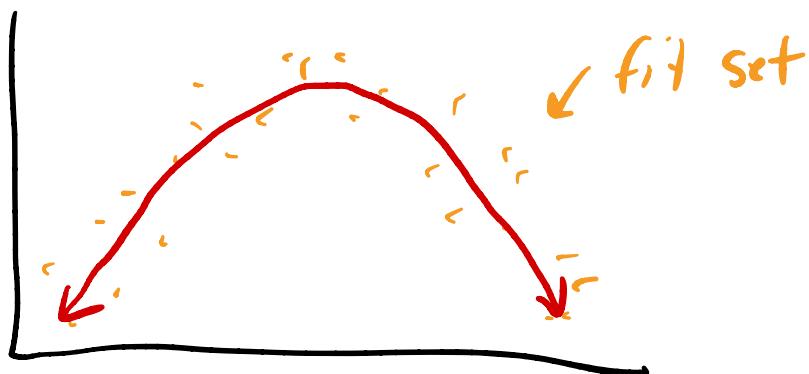
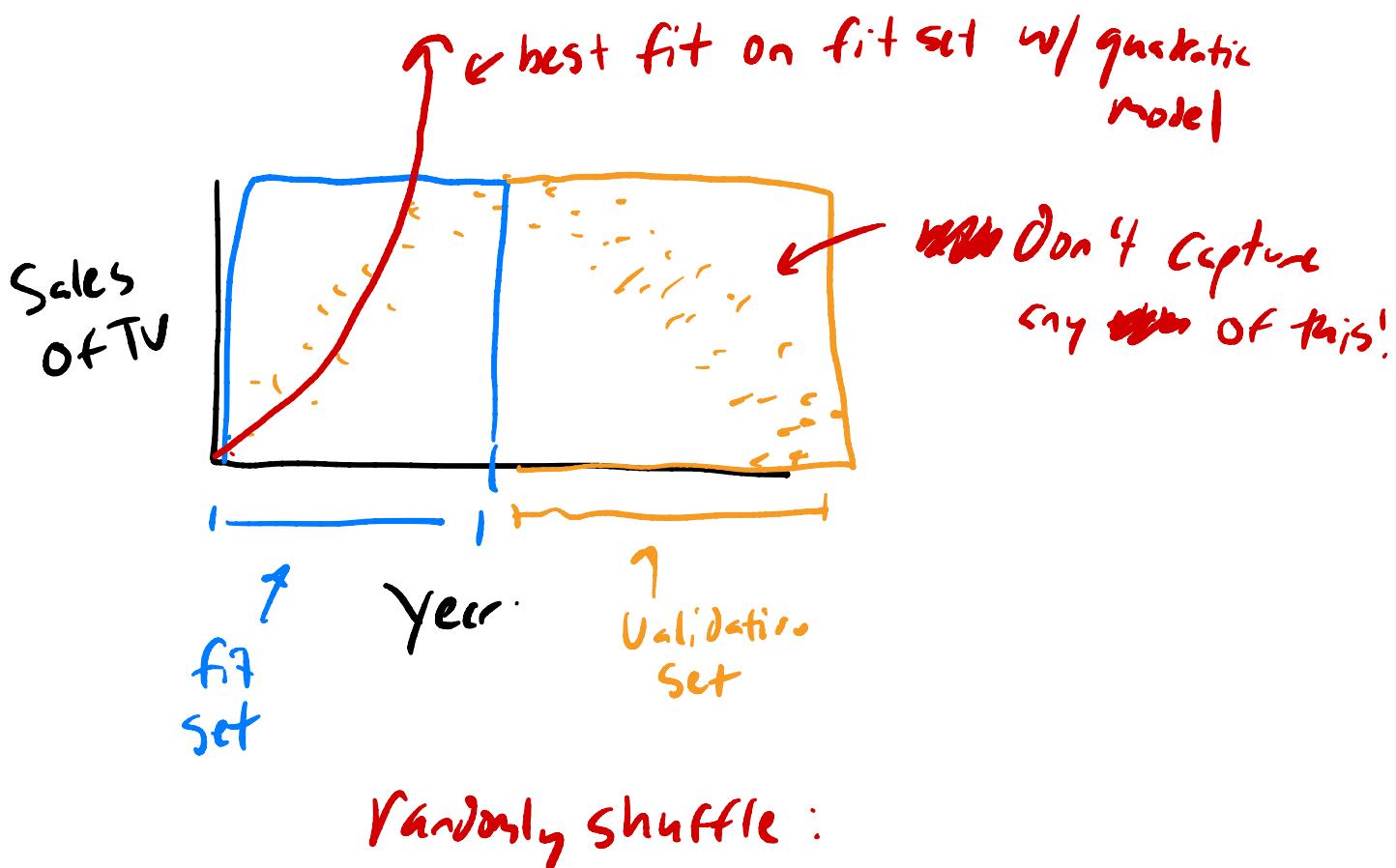
- 1) **Fit set**: data used to fit linear regression
- 2) **Validation set**: data used to check for overfitting:
how well does regression fit
generalize to data not used
in fit process.



example: if $N = 1000$, could make 1st $N = 500$
fit set, remaining 500 be validation set.

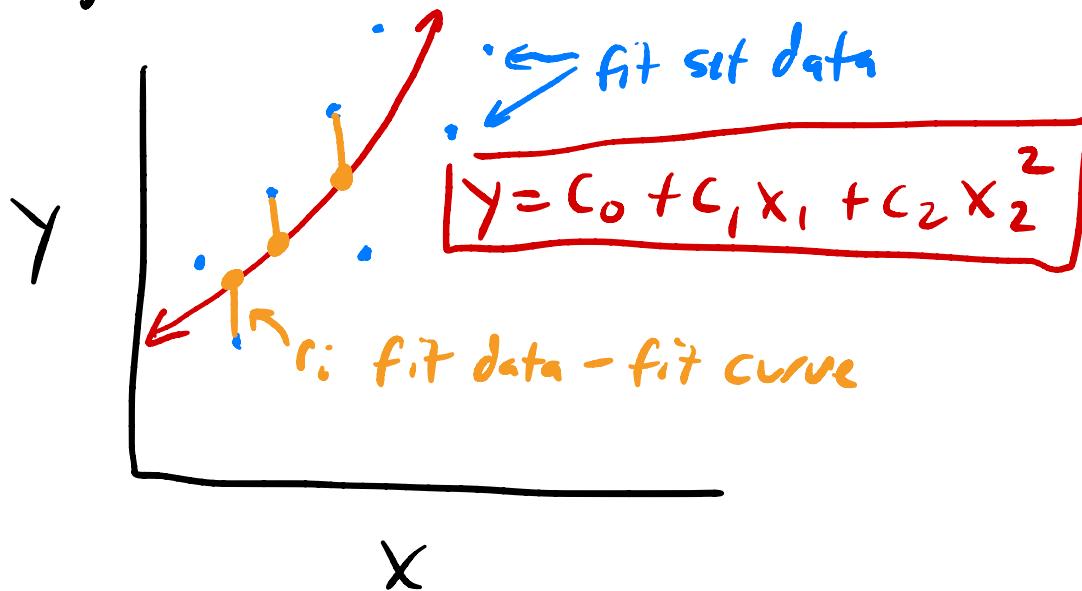
★ Each sample (A) must be in either the fit set or Validation, not both.

★ Usually best ^{randomly} shuffle order of N samples (rows) before splitting A into fit/validation sets.



Fit and validation data workflow

- 1) Fit regression model to fit set data \Rightarrow get coefficients \vec{c} :

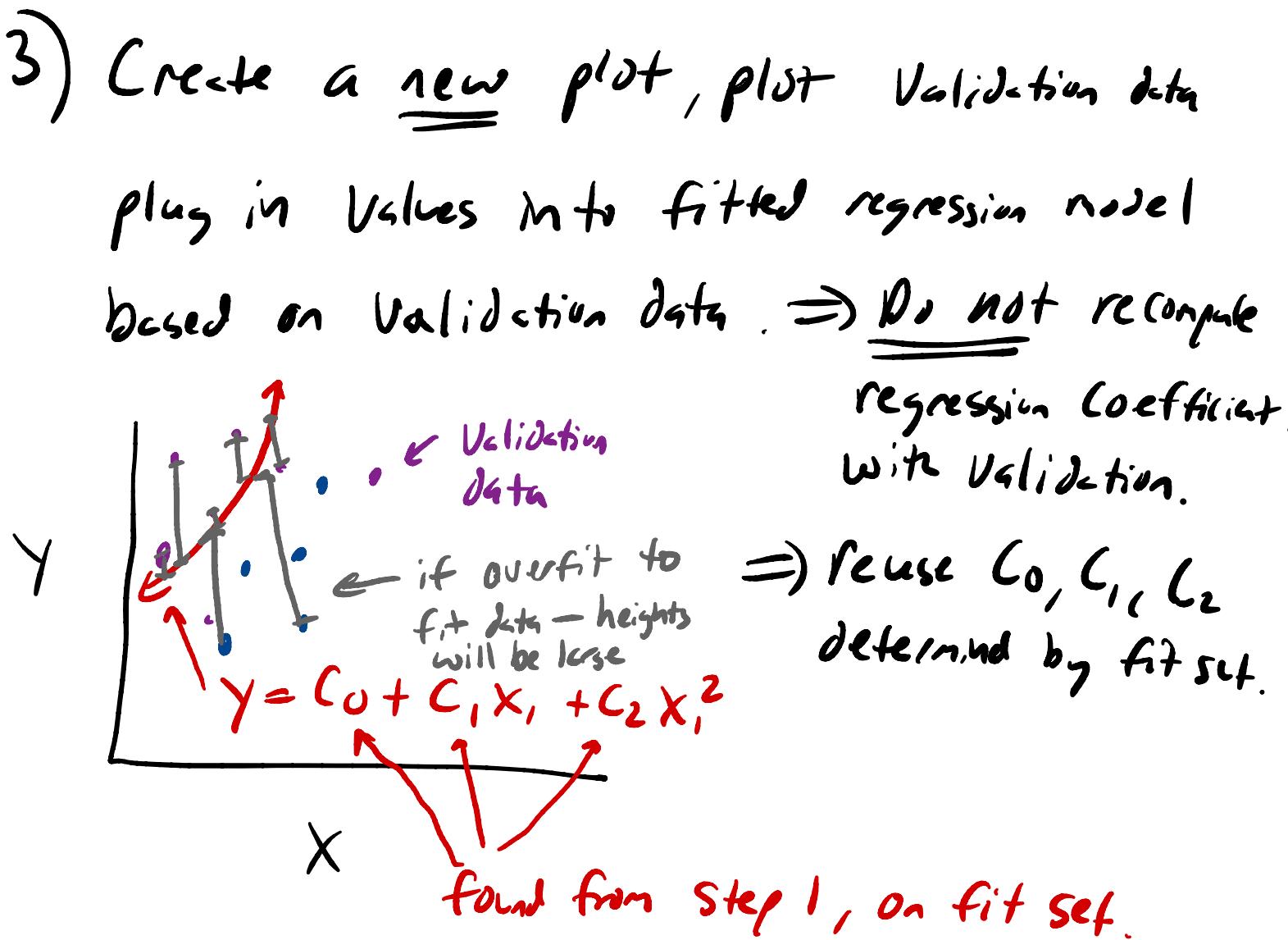


- 2) See how well I fit the fit data

Compute residuals between

- quadratic model
- fit data

$\Rightarrow R^2, \text{MSSE}, \dots$ get quality of fit.



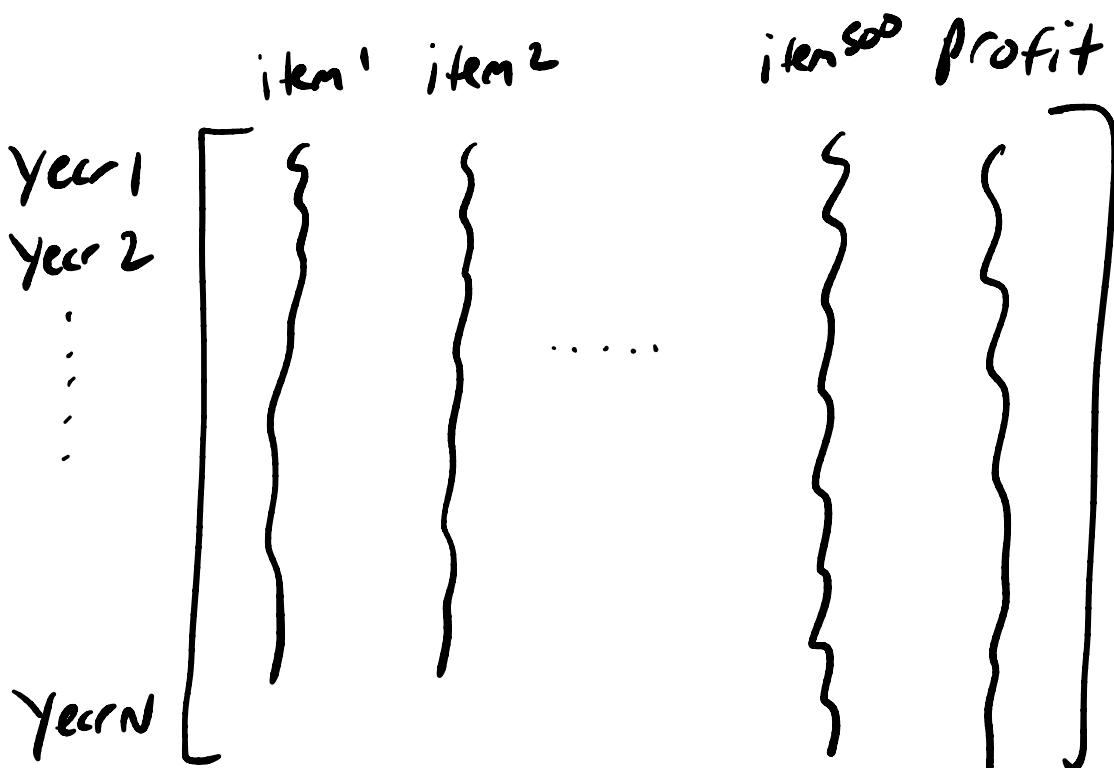
- 4) Compute residuals between
- regression curve fit using fit data
 - Validation data
- \Rightarrow Calculate R^2 , MSSE using

Overfitting

fit data : $R^2 \approx 1$, MSSE ≈ 0

Validation data : $R^2 \approx 0$, MSSE big

Stepwise linear regression



Stepwise linear regression : greedy algorithm.

- 1) Start regression model w/ intercept only: $A = [\vec{1}]$
- 2) Try to add each one of the 500 vars to your regression model:

Try: $A = [\vec{1}, \vec{x}_1]$

\uparrow
1st var

\Rightarrow Calc R^2 to measure how predictive model is of \vec{y}

Next: $A = [\vec{1}, \vec{x}_2]$

\Rightarrow keep track of which added var produces the highest improvement in R^2 value.

\Rightarrow After visiting all 500 vars, officially add var that improved R^2 the most.

e.g. \vec{x}_{30}

$$A = [\vec{1}, \vec{x}_{30}]$$

3) With updated model, try again to add all remaining vars (499) to regression

e.g.

$$A = [\vec{1}, \vec{x}_{30}, \vec{x}_1] \Rightarrow R^2$$

$$A = [\vec{1}, \vec{x}_{30}, \vec{x}_2] \Rightarrow R^2$$

:

:

4) Keep going until all variables added to large enough regression

