

I 3344-Digital Modeling & Simulation

problems given a set of x 's and their respectful y 's. Estimate y_n for a given x_n .

Solved by interpolation & extrapolation.

$$0 \leq n \leq i \quad n \text{ outside range}$$

Old way: polynomial regression (fitting a curve that goes through all points)

not suitable for cases of high noise/error

Better way: Least-Square Approximation by Linear Regression

$$n \begin{bmatrix} x_i & y_i & x_i^2 & x_i y_i & \hat{y}_i & e_i & e_i^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum x & \sum y & \sum x^2 & \sum xy & & & \end{bmatrix}$$

$(\sum x)^2$

$$a_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$
$$= r \frac{\sigma_y}{\sigma_x}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

$$e_i = y_i - \hat{y}_i$$

(Sum of the Squares of the Residual) $SSE = \sum (e_i)^2$

Normalize the errors: $MSE = \frac{SSE}{n}$

$$RMSE = \sqrt{MSE}$$

$$NRMSE = \frac{RMSE}{\bar{y}}$$

close to 0 \Rightarrow Perfect model
close to 1 \Rightarrow Null model

k-Nearest Neighbors

For a given x_n , we can approximate \hat{y}_n by:

- take k (usually $k = \sqrt{n}$)

- for every x_i , calculate $d_i = \sqrt{(x_n - x_i)^2}$

- take the y 's of the k least d 's and $\hat{y}_n = \frac{y_1 + y_2 + \dots + y_k}{k}$

k-fold cross validation: taking the set as k subsets. Each iteration we set one of them as the validation set and the rest as a training set. We fit the function on the latter (training) and test it on the former (validation).

* KNN: per iteration, calculate \hat{y} for each point in the validation set using the k nearest neighbors from the training set only and get NRMSE.

* Linear Regression: get a_1 & a_0 using the training set and calculate NRMSE in the validation set using these values per iteration

Cross validation: $CV = \frac{\sum NRMSE}{k}$

the method with the lowest CV is better.

Multiple Linear Regression: One y and multiple xs.

$$\hat{Y} = X * \beta \quad \text{with} \quad \beta = (X^T X)^{-1} X^T Y$$

$$R^2 = 1 - \frac{SSE}{TSS} \rightarrow \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad \sim 1 \text{ Perfect Model}$$

$$\sim 0 \text{ Null Model.}$$

Problem: R^2 always \uparrow as we add variables (xs) but that doesn't mean a better model.

fix: adjusted R^2 : $R_a^2 = 1 - (1 - R^2) \frac{n-1}{n-p-1}$
 where p = number of predictors.

ex)

A	B	C	D
1	2	14	3
3	5	23	4
7	9	17	8
4	5	19	8

$$Y = A = [1, 3, 7, 4]$$

$$\text{for } X = B = \begin{pmatrix} 2 \\ 5 \\ 9 \\ 5 \end{pmatrix}, R_a^2 = 0.973$$

$$X = C, R_a^2 = -0.4$$

$$X = D, R_a^2 = 0.95$$

pick B for second iteration:

$$X = B C = \begin{pmatrix} 2 & 14 \\ 5 & 23 \\ 9 & 17 \\ 5 & 19 \end{pmatrix}, R_a^2 = 0.94$$

$$X = B D, R_a^2 = 1 \checkmark$$

Notes for solving by hand:

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{pmatrix}, A^T = \begin{pmatrix} 1 & 5 \\ 2 & 6 \\ 3 & 7 \\ 4 & 8 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & 10 & 10 \\ 10 & 30 & 26 \\ 10 & 26 & 26 \end{pmatrix}, \begin{pmatrix} 4 & 10 & 10 \\ 10 & 30 & 26 \\ 10 & 26 & 26 \end{pmatrix}$$

3000 + 2704 + 2600 = top

3180 + 2600 + 2600 = bottom

$$|A| = \text{top} - \text{bottom} = 120 - 104 = 16$$

$$A = \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix}, \text{cof}(A) = \begin{pmatrix} |e f| & -|d f| & |d e| \\ -|b c| & |a c| & -|a b| \\ |b e| & -|a e| & |a b| \end{pmatrix}$$

$$A^{-1} = \frac{1}{|A|} \text{cof}(A)$$