

ML day 15

High Dimensional Data

- n - number of rows and p - number of parameters that is mean column or variables
- Most traditional $n \gg p$, it is low dimensional and it is ok
- but if there is $n < p$, n is too much greater than p there will be high dimensional data .
- Gap between the n and p must be less for high Dimensional Data .

High-Dimensional Data

Most traditional statistical techniques for regression and classification are intended for the low-dimensional setting in which n is much greater than p .

Let's consider predicting the blood pressure based on age, gender and body mass index (BMI). There are three or four (if an intercept is included) predictors in the model and thousands of patients with the predictor values. Hence $n \gg p$ and the problem is low-dimensional.

But now it is common to collect a large number of feature measurements (in fields such as Finance, marketing and medicine). While p can be extremely large n is often limited due to cost, sample availability or other considerations.

Datasets containing more features than observations are called as High-dimensional. Classical approaches such as least squares linear regression are not appropriate in this case.

What goes wrong in High Dimensions ?

- So consider the chart 1, in here there is only one p and near there are 100 of rows, So think if we increase the number of parameters the best fit line

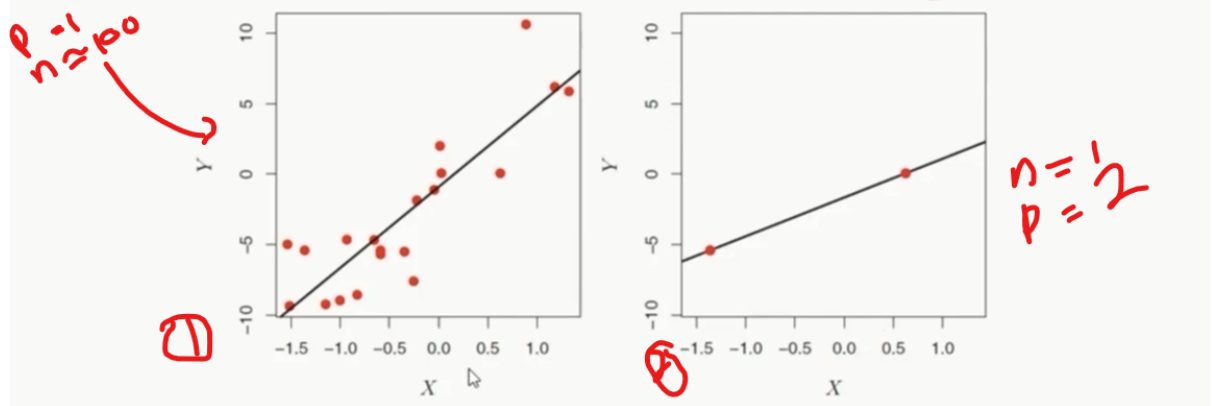
will get curved

- when we consider the image number 2 , there is $n = 1$ and $p = 2$, so $n > p$ it is true but this will cause the overfitting because this data set will not suite for the unseen data

What Goes Wrong in High Dimensions?

When p is as large as or larger than n , least squares cannot or should not be performed. Because regardless of whether or not there is a relationship between the features and the response, least squares will yield a set of coefficient estimates that result in a perfect fit to the data such that the residuals are zero.

Consider following example.



Curse of Dimensionality

- example :
 - If there is a child we will give cookies to test , first think if there is two flavors at least child must eat two biscuit to select best , if there is three shapes and two flavors child must eat $2 \times 3 = 8$ cookies
- When the number of features grows the number of data also getting high
- This one also not a good thing .
- What can be do in here ? So we have to reduce the dimensions .

Dimension Reduction

- There are two main ways to do dimensionality Reduction
 1. Feature Selection
 2. Dimensionality Reduction Algorithms

Feature Selection : Supervised Learning

- Best Subset Selection
- Forward Selection
- backward Selection

Best Subset Selection

- We have to create model from null model to Full Model
- After That we have to get the less RSS value (Residual Sum of Squares)
 -

1. Definition:

- RSS (Residual Sum of Squares):

- The sum of the squared differences (residuals) between the actual and predicted values.

- Formula:

$$RSS = \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

Where $y^{(i)}$ is the actual value, $\hat{y}^{(i)}$ is the predicted value, and m is the number of data points.

- MSE (Mean Squared Error):

- The average of the squared differences between the actual and predicted values. It is essentially the RSS divided by the number of data points.

- Formula:

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

- 2 to the power P will give number of Models
- If the P is increasing there will be more models so we have to find the error in more models .

Handwritten notes illustrating the progression of linear models:

- Features: x_1, x_2, x_3, x_4, x_5
- Null model: $y = \beta_0$
- Model 1: $y = \beta_0 + \beta_1 x_1$
- Model 2: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- Model 3: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$
- Model 4: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$
- Full model: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$

Forward Selection

- start with null model
- Fit P simple linear regression models and get the lowest RSS model .
- Increase that model with rest of the variables (P-1) with the lowest RSS .
- Construct until some stopping rule is satisfied .
-

Handwritten notes illustrating the Forward Selection process for three variables x_1 , x_2 , and x_3 :

- Start with the null model: $y = \beta_0 \leftarrow \text{Null}$
- Fit simple linear regression models for each variable:
 - $y = \beta_0 + \beta_1 x_1$ (RSS)
 - $y = \beta_0 + \beta_1 x_2$ (RSS) — This model is selected as the best among the single-variable models.
 - $y = \beta_0 + \beta_1 x_3$ (RSS)
- Fit a multiple linear regression model with the selected variable and the next variable:
 - $y = \beta_0 + \beta_1 x_2 + \beta_2 x_1$ (RSS) — This model is selected as the best among the two-variable models.
 - $y = \beta_0 + \beta_1 x_2 + \beta_2 x_3$ (RSS)
- Construct the full model: $y = \beta_0 + \beta_1 x_2 + \beta_2 x_1 + \beta_3 x_3 \leftarrow \text{Full}$

Backward Selection

Backward Selection

- Start with the model with all P variables.
- Remove one variable and select the most significant model ($P-1$).
- Continue until a stopping rule is reached