

COMS 474 HW 1

Haadi Majeed

February 2, 2022

Contents

1	Problem 1	3
2	Problem 2	5
3	Problem 3	6

1 Problem 1

[16 points; 4 each part] For each of parts (a) through (d), indicate whether we would generally expect the performance of fitted models on future data to be better or worse if the model class has many degrees of freedom (e.g. class of high degree polynomials) compared to a model class with few degrees of freedom (e.g. constant or linear functions). Briefly explain why (1-3 sentences).

(a) the number of samples n is large, and the number of features p is small.

- When working with large amounts of data, and fitting it with a small number of features, the model that comes out of it will not accurately represent the dataset as it may be treating trends and their variations as just noise that it ignores. Having p be too small would result, generally, in under-fitting the data. This would result in a poorer model in general.

(b) the number of features p is large, and the number of samples n is small.

- When working with small amounts of data and many features, or higher flexibility, it would over-fit the data where it would effectively treat the data as if there was no noise. Such can be seen in the image below, where the model to the 16th power successfully hits every point of data, however is not very useable for future data, which is what we care about. This would result in a poorer model in general.

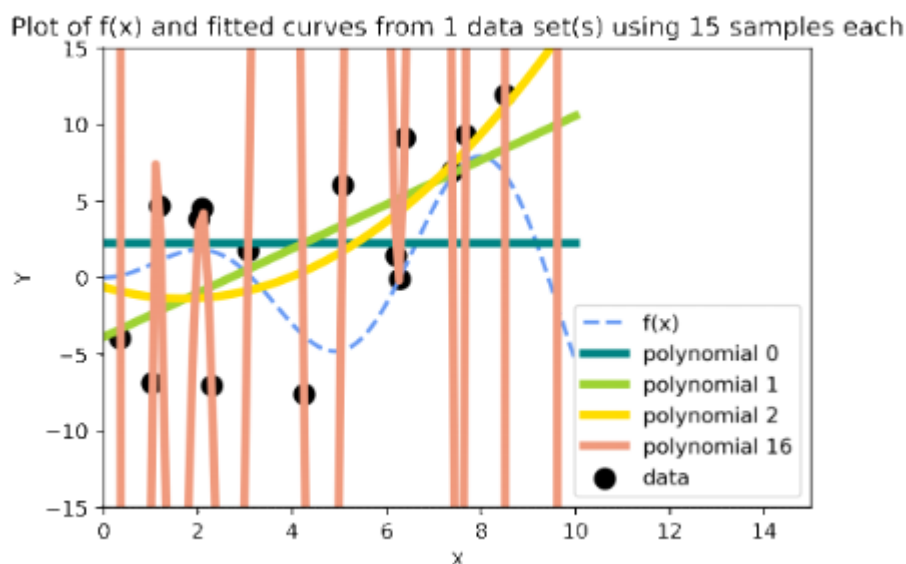


Image from page 8 of January 27th's Notes

(c) The relationship between the features $\{X_1 \dots X_p\}$ and the response (Y) is highly non-linear.

- As the predicted model has more features, the response will look less and less like a linear model, this is because the model will start to overfit the data. For example, if there is 16 data points, with 20 features vs 16 data points but only 5 features, the 20 feature one will overfit the data and most likely take in the noise value as is instead of as variance like the 5 feature model may.

(d) The variance of the noise terms, $\{X_1 \dots X_p\}$, is extremely high

- When the noise of the data is extremely high, it becomes more difficult to fit a model to it as the future data may not share the variance to such a magnitude. This is because all noise is absolutely independent and irreducible, and predicting future noise is impossible. This model could be either Over or Under fitted as the approach does not set which type of behaviour is taken, regardless fitting a model to this would risk being askewed and would result in a poorer model.

2 Problem 2

[18 points total; 6 each part] You will now think of some real-life applications for machine learning.

- (a) Describe three real-life applications in which classification might be useful. Describe the response, as well as the predictors
- Qualitative, these would be things that are attributes or properties that something could have
 - Determining the integrity of a bridge in a simulator
Predictors: Load on the bridge and materials used
Response: Stress test results if bridge breaks or not
 - Is this picture a bird or not
Predictors: Input Photo and previous data
Response: Yes or No based on algorithm
 - Does this email look like spam
Predictors: Email and previous data
Response: Yes or no based on algorithm
- (b) Describe three real-life applications in which regression might be useful. Describe the response, as well as the predictors.
- Quantitative, these are often numerical values that something has
 - Resistance of a electrical component
Predictors: Voltage, Current
Response: Voltage, Current
 - What will the temperature be tomorrow?
Predictors: Previous Data, Date
Response: A temperature value
 - Credit Score
Predictors: Previous Data, Purchase History,
Response: A Score that is carefully calculated
- (c) Describe three real-life applications in which cluster analysis might be useful
- Behavioural Clustering - Seeing how different people with various mindsets react to the same situation. Predictors may include past events, gender, ethnicity, and backgrounds
 - Determine Fruit/Veggie ripeness - Based off past data, determine how under/over/correctly ripe the input is. Predictors may include size and colour.
 - Search Engine Results - Determine how close some results are than others to the user's query and keywords

3 Problem 3

We want to learn a model to predict Y , let N denote the number of samples of data. Using calculus, derive the optimal constant function under mean square error

$$\beta_0^* = \arg \min_{\beta_0 \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (Y(i) - \beta_0)^2.$$

make sure to check you found a minimizing β_0 , not a maximizing β_0 .

$$\beta_0^* = \arg \min_{\beta_0 \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (Y(i) - \beta_0)^2.$$

$$\hat{y} = \frac{d}{db} \left(\frac{1}{n} ((Y(1) - \beta_1)^2 + (Y(2) - \beta_2)^2 + \dots + (Y(n) - \beta_n)^2) \right)$$

$$\hat{y}' = \frac{1}{n} ((2(Y(1) - \beta_1)) + (2(Y(2) - \beta_2)) + \dots + (2(Y(n) - \beta_n)))$$

$$\hat{y}' = \frac{1}{n} (2Y(1) - 2\beta_1 + 2Y(2) - 2\beta_2 + \dots + 2Y(n) - 2\beta_n)$$

$$\hat{y}' = \frac{1}{n} (2 \sum_{i=1}^n (Y(i)) - 2(\beta_1 + \beta_2 + \dots + \beta_n))$$

$$\hat{y}' = \frac{2}{n} \sum_{i=1}^n (Y(i)) - \frac{2}{n} (\beta_1 + \beta_2 + \dots + \beta_n)$$

$$\hat{y}'' = \frac{d}{db} \hat{y}' \Rightarrow 2 > 0 \text{ for all } \beta$$

Therefore

$$\beta_0^* = \frac{1}{n} (\beta_1 + \beta_2 + \dots + \beta_n)$$

Is when \hat{y} is minimized