

Lab: Linear Regression and Bias-Variance Tradeoff

1. linear regression - Bias-Variance Tradeoff

Assume that X_1, X_2 are two independent variables but with a same distribution $N(1, 1)$. The true relationship between Y_i and X_{1i}, X_{2i} is $Y_i = 1 + X_{1i} + 0.001X_{2i} + \epsilon_i$, where $\epsilon_i \stackrel{i.i.d}{\sim} N(0, \sigma^2)$.

(a) Is the intercept/coefficient of X_1 biased if you only regress on X_1 ? What is your intuition?

(b) Follow-up: Show your conclusion in (a) mathmatically. Here are some hints:

Step 1: Based on the model we fit, we assume $E(Y|X_1) = \dots$

Step 2: We know the ‘true’ relationship between Y and X_1, X_2 , use it to replace Y in the above equation.

Step 3: Work on the expectation and reach your conclusion.

(c) Would you include X_2 to improve your model based on your intuition ?

(d) [Teamwork] Verify your conclusion using a simulation. Please follow the comments in the code chunk.

```
## Step 1: generate a training set
set.seed(263)
n = 50
x1 = rnorm(n,1,1)
x2 = rnorm(n,1,1)
eps = rnorm(n,0,0.2)
y = 1+x1+0.001*x2+eps
trainset=data.frame(cbind(y,x1,x2))
```

```
## Step 2: fit models on trainset: y~x1 and y~x1+x2
## fit1 = ...
## fit2 = ...
## Step 3: generate a test set
m=10000
x1 = rnorm(m,1,1)
x2 = rnorm(m,1,1)
eps = rnorm(m,0,0.2)
y = 1+x1+0.001*x2+eps
testset=data.frame(cbind(y,x1,x2))
## Step 4: get the predictions in test set:
## pred1 = ...
## pred2 = ...
## Step 5: compare the MSEs in test set
## MSE1 =
## MSE2 =
```

Try to answer the questions below and get the idea of bias-variance tradeoff:

- (1) In Model 1, the estimate of the intercept is biased/unbiased (choose one), the MSE on the test set is _____.
- (2) In Model 2, the estimate of the intercept is biased/unbiased (choose one), the MSE on the test set is _____.
- (3) Based on MSE, Model 1/2 (choose one) is better, so you can infer that the predictions using Model 1/2 (choose one) have a larger variance.

- (e) (Optional advanced problem) Let's go back to (a) and think, is the intercept/coefficient of X_1 biased if you only regress on X_1 , given that X_1 is correlated with X_2 ?

2. Predict House Price Using Regression

This dataset('kc_house_data.csv') contains house sale prices for King County, which includes Seattle. It includes homes sold between May 2014 and May 2015. There are 19 house features plus the price and the id columns, along with 21613 observations. The dictionary of the variables is listed in the next page.

```
## Load in the data, split the data into training set and test set
house = read.csv('kc_house_data.csv',header = T)
trainset = house[1:floor(nrow(house)*0.8),]
testset = house[-(1:floor(nrow(house)*0.8)),]
```

- (a) Fit a linear model on the training set: price = bedrooms + bathrooms + condition. Interpret the estimated coefficient of bathrooms and provide the corresponding 95% confidence interval.

- (b) Fit a linear model on the training set: $\text{price} = \text{bedrooms} + \text{bathrooms} + \text{condition} + \text{sqft_above}$. Compare the coefficients of bedrooms here with the one in (a), what do you find?
- (c) Using the model in (a), predict the price in the test set and calculate the mean square loss $L = \frac{1}{M} \sum_{i=1}^M (Y_i - \hat{Y}_i)^2$.
- (d) [Teamwork] Competition Time!
Use linear regression to get the best prediction! (Minimal square loss in test set). Think about:
1. Transformation: what is the best $\phi(x)$, e.g. $\log(\text{sqft_above})$? $\sqrt{\text{sqft_above}}$? or original?
2. Should we include all variables? How to combine different pieces of information, e.g. yr_built and yr_renovated ?
3. continuous or categorical?

Columns

- # id a notation for a house
- A date Date house was sold
- # price Price is prediction target
- # bedrooms Number of Bedrooms/House
- # bathrooms Number of bathrooms/House
- # sqft_living square footage of the home
- # sqft_lot square footage of the lot
- # floors Total floors (levels) in house
- A waterfront House which has a view to a waterfront
- A view Has been viewed
- A condition How good the condition is (Overall)
- A grade overall grade given to the housing unit, based on King County grading system
- # sqft_above square footage of house apart from basement
- # sqft_basement square footage of the basement
- # yr_built Built Year
- # yr_renovated Year when house was renovated
- # zipcode zip
- # lat Latitude coordinate
- # long Longitude coordinate
- # sqft_living15 Living room area in 2015(implies-- some renovations) This might or might not have affected the lotsize area
- # sqft_lot15 lotSize area in 2015(implies-- some renovations)

Figure 1: variable_dict

3. (Optional advanced problem) Distribution Theory - Matrix Representation

This question is beyond the scope of this class. It is here only for those who want more practice on matrix representations.

Let y be a $k \times 1$ multivariate normal random vector with mean μ and nonsingular variance-covariance matrix V , $y \sim N(\mu, V)$. Additionally, let A be a $k \times k$ matrix of constants and B be a $q \times k$ matrix. Then, the linear form $W = By$ and quadratic form $U = y^T Ay$ are independent if $BVA = 0$.

Prove that $\hat{\beta} = (X^T X)^{-1} X^T Y$ and $\hat{\sigma}^2$ are independent using the theorem above.