

Data Preprocessing : One Hot Encoding

[creating new features
- feature engineering]

Suppose a data set that consists of four columns.

- | | | |
|------------------|---|-----------|
| 1. locality | - | text data |
| 2. area-sq.-ft | - | number |
| 3. bedrooms | - | " |
| 4. price - lakhs | - | " |

property price depends upon locality also.

computer understands numbers.

text can be converted into numbers - feature engineering

Method 1 - Label Encoding

Assign numbers to text data

eg:

Locality (categorical data)

Kollur - 1

Banjarahills - 2

Mankhal - 3

categorical data

1. ~~categorical~~ nominal (no order).

2. ~~categorical~~ ordinal

(categories have some kind of order)

eg: green & red & blue.

Kollur <

Banjarahills <

Mankhal

bachelors < masters < phd

dissatisfied < neutral < satisfied

low, medium, high

One hot encoding is used.

Label encoding is used



Locality A	Locality B	Locality C	Area sq-ft	Bed rooms	Price lakhs
1	0	0	656	3	39
1	0	0	1260	2	83
1	0	0	1057	3	86
0	1	0	2400	3	300
0	1	0	1008	2	50
0	0	1	1125	2	64
0	0	1	1700	3	100

Locality A, Locality B, Locality C - Dummy variables
Original column - locality

On that column, we apply onehot encoding and we get 3 dummy columns.

In this one of the column is redundant.

If we know two set of columns, from that we can infer the third one.

If it is neither A or B then it will be C.

Multicollinearity:- occurs when two or more independent variables are highly correlated, making it difficult to distinguish their individual effects on the dependent variable.

To overcome this problem, we drop any one column.

pd.get_dummies(df, columns = ["locality"])

due to multicollinearity column one col is dropped.

pd.get_dummies(df, columns = ["locality"],
drop_first = True)
↳ drop 1st dummy variable

df_encoded = pd.get_dummies(df, columns =
["locality"], drop_first = True)

df_encoded.sample(5)

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

X = df_encoded.drop('price_lakhs', axis=1)
y = df_encoded['price_lakhs']

X_train, X_test, y_train, y_test = train_test_split(

X, y, test_size = 0.2, random_state = 42)

model = LinearRegression
model.fit(X_train, y_train)

model.score(X_test, y_test)

test = pd.DataFrame([
{'area_sq_ft': 1600, 'bedrooms': 2,
'locality_kollu': False, 'locality-
Mankhali': False}
)

Date ____ / ____ / ____



model.predict(test)



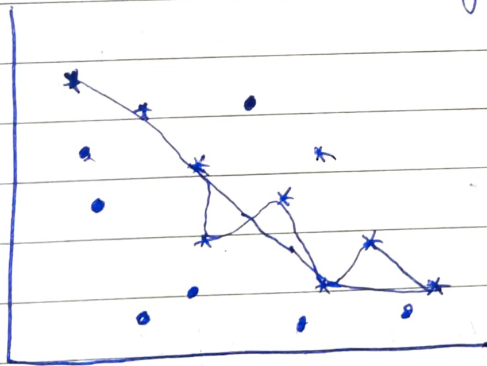
Overfitting :- Occurs when a model learns too much detail and noise from the training data, affecting its performance on new data.

Underfitting :- Happens when a model is too simple and can't learn the data pattern, leading to poor performance on all data.

Balanced fit - Achieved when a model accurately ~~known~~ learns the training data patterns and performs well on unseen data.

Overfitting

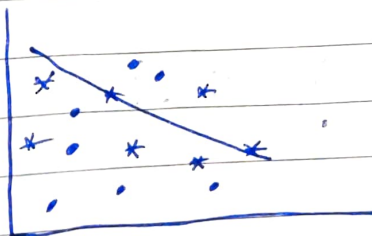
- Train error - low
- Test error - High



- - test datapoints
- * - train "

Learning from noise/
fluctuations
not trying to learn
general patterns

underfit



test & train error - high

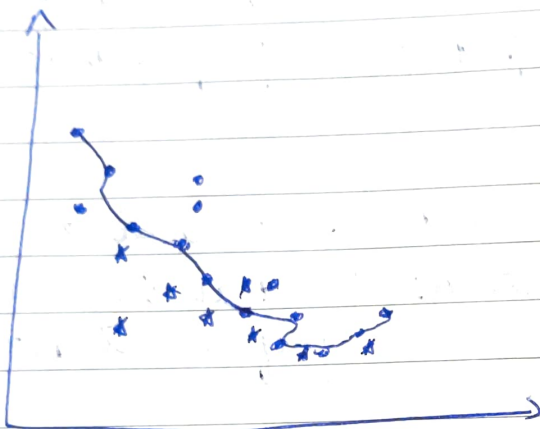
Date

Regularization

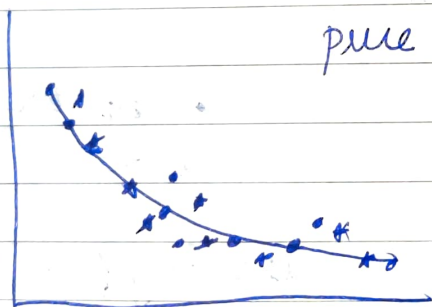
→ used to reduce overfitting

Let an overfit model using polynomial regression

$$\text{price} = \beta_0 + \beta_1 * \text{mileage} + \beta_2 * \text{mileage}^2 + \beta_3 * \text{mileage}^3 + \beta_4 * \text{mileage}^4 + \beta_5 * \text{mileage}^5$$



right model here is this



$$\text{price} = \beta_0 + \beta_1 * \text{mileage} + \beta_2 * \text{mileage}^2$$

If we try to reduce β_3, β_4 & β_5 as zero we will get second eqn.

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{predicted}})^2$$



here $y_{\text{predicted}}$ is

$$\text{price} = \beta_0 + \beta_1 * \text{mileage} + \beta_2 * \text{mileage}^2 + \beta_3 * \text{mileage}^3 + \beta_4 * \text{mileage}^4 + \beta_5 * \text{mileage}$$

In general

$$y_{\text{predicted}} = \beta_0 + \beta_1 x + \beta_2 * x^2 + \dots + \beta_n * x^n$$

~~we~~ we need to reduce β_3, β_4 & β_5 close to zero so that we will get a balanced fit model.

for that we are adding a new term to the equation.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{predicted}})^2 + \lambda \sum_{i=1}^n \beta_i^2$$

$\lambda \rightarrow$ a hyperparameter that we can configure
 $\sum_{i=1}^n \beta_i^2 = \beta_0^2 + \beta_1^2 + \beta_2^2 - \dots$

If larger β_3 value is there let it be 100.
 $100 \times 100 * \lambda = 10000$
 a big number.

$\beta_3 * x^3$ gives again a big no.

If we substitute this value in equation for $y_{\text{predicted}}$ MSE will also increase

so Gradient Descent try to reduce the MSE
 So it will configure β_3 or β_5 value

Thus we get a balanced fit
(all large polynomial values $\rightarrow 0$)

This is known as L_2 regularization
This is also called ridge regression

L_1 regularization / Lasso regression

penalty term adding is $|\beta_i|$
absolute value of β

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{predicted}})^2 + \sum_{i=1}^n |\beta_i|$$

$$y_{\text{predicted}} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_n x^n$$

Let $\beta_3 = 100$

$$\frac{100 \times x^3}{\text{large value}}$$

MSE increases

GD tries to minimize error.

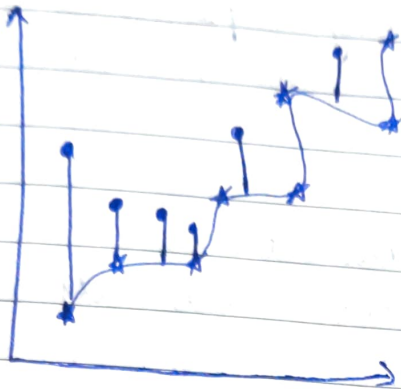
to minimize error

it minimize higher polynomial
 β terms.



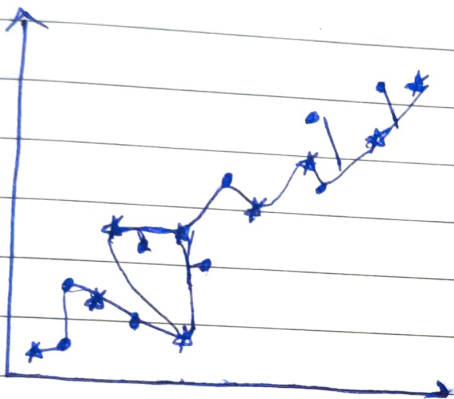
Bias Variance Tradeoff

Eg:



Train dataset error = 0
Test " " = 100
higher test error

high
variance
in terms of test data



Train dataset error = 0
Test " " = 27
lower test error

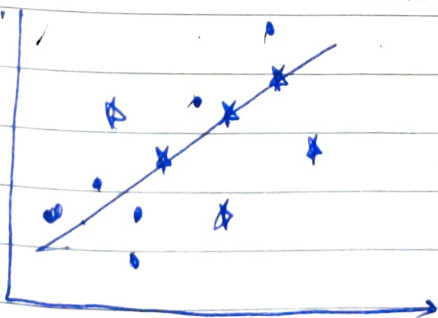
same data set.

overfitting

test error varies

usually overfit model have high variance
Based on the selection of the training dataset,
the prediction errors vary greatly (i.e. showing
high variance among errors)

underfitting



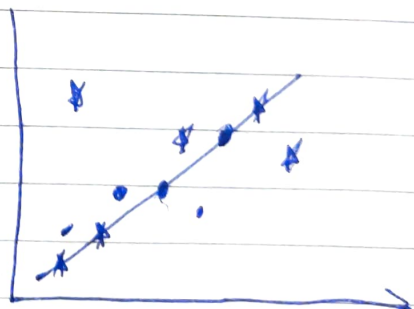
Test error = 47

Train " = 43

high bias

error high

high
bias



Test error = 37

Train error = 41

Bias is a measurement of how accurately a model
can capture a pattern in a training dataset.

In above case since error is big it is said to have high bias

When you try to build a very simple model, you will get high errors in testing & training data \Rightarrow high bias

In case of balanced data set

\Rightarrow low variance, low bias.