

## multiple linear regression

features : 'n' number of features

label : 1 label  $\rightarrow$  numeric

### Venture Capitalist Use-case :

R&DSpend  $\rightarrow$  numeric  $f_1$

Administration  $\rightarrow$  numeric  $f_2$

marketing  $\rightarrow$  numeric  $f_3$

State  $\rightarrow$  Categorical  $f_4$

Profit  $\rightarrow$  numeric label

No missing  
data issue

After handling categorical data,

California

Florida

New York

R&D Spend

Administration

marketing

Profit

Equation of line goes:

$$\begin{aligned} \text{Profit} = & b_0 + b_1 [\text{California}] \\ & + b_2 [\text{Florida}] + b_3 [\text{NY}] \\ & + b_4 [\text{R\&D Spend}] + b_5 [\text{Adm}] \\ & + b_6 [\text{market}] \end{aligned}$$

# FEATURE ENGINEERING

## Feature Selection

① RFE

② Select by model

③ ANOVA

(Analysis of Variance)

## Feature Extraction

① PCA (Principal Component Analysis)

② LDA (Linear Discriminant Analysis)

## Feature Selection analysis

### Regression :-

→ ① Correlation analysis (corr) ✓

② Elimination using OLS ←

③ Feature selection ✓

④ Feature extraction

## Backward elimination using OLS

Hypothesis  
Testing

R<sup>2</sup>D  $\longleftrightarrow$  Profit  
P value

Calc<sup>n</sup>  
using OLS

$P \leq 0.05 \rightarrow$  Rel<sup>n</sup>ship exists so keep the feature

else  $\rightarrow$  Eliminate feature

Significance level

p-value

(the amount of  
error you can  
tolerate)

confidence

(the trust you can  
have on the model)

$\rightarrow 0.01 \leftarrow$

$\rightarrow 0.05 \leftarrow$

$\rightarrow 0.1 \leftarrow$

$\rightarrow 0.99 \leftarrow$

$\rightarrow 0.95 \leftarrow$

$\rightarrow 0.9 \leftarrow$

my significance level is

$P = 0.05$

$\leftarrow$  Data Scientist

Role of ML engg:

o Eliminate the features whose p value  
is greater than 0.05

Ideal steps we perform in feature elimination  
using OLS / AIC

---

- ① Perform All-IN (ensure you provide all coeff and intercept coeff)

$$\text{profit} = \textcolor{red}{b_0} + \textcolor{yellow}{b_1} (\text{California}) + \textcolor{yellow}{b_2} (\text{Florida}) + \textcolor{yellow}{b_3} (\text{NY}) + \textcolor{yellow}{b_4} (\text{R\&D}) + \textcolor{yellow}{b_5} (\text{adm}) + \textcolor{yellow}{b_6} (\text{market})$$

- ② Decide the Significance level (0.05)
- ③ select the feature in summary which has highest p-value.
- ④ If p value > SL(0.05), eliminate that feature.
- ⑤ Repeat step 3 & 4 till condition ④ is false. If condition ④ is false, go to step 6.
- ⑥ consider all remaining features as selected feature for model building.