

Training, Validation and Testing Set

Sept 2021

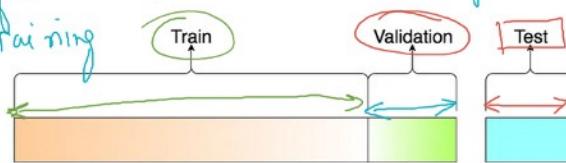
Training set: Sample of data used to fit/train the model and make the model learn the hidden features/pattern in the data

For 3x1

Validation set: It is a set of data separate from training set that is used to validate the model performance during training & update hyperparameters

Final Model

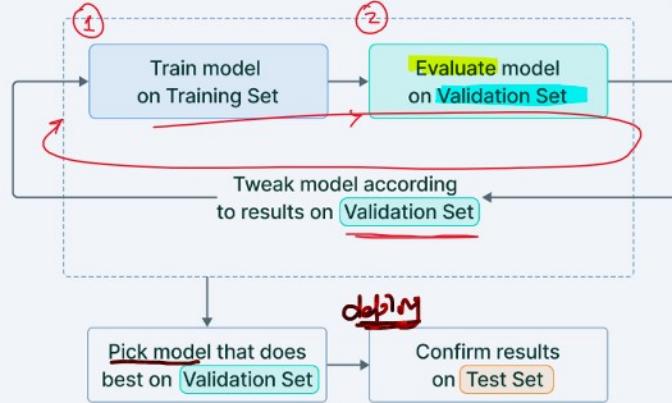
Test set: It is a separate set of data used to test the model after completing the final training



deploy the model

Training data/validation/test

9 Dec



Forecasting Sales

Prod A	Month	Actual	
		Model	Actual
	JAN	100	150
	FEB	80	90
	MAR	90	100

$$\begin{aligned}
 (\text{Error}) &= \frac{50}{150} \times 100 = 33.3\% \quad \text{Feb} \\
 &= \frac{10}{90} \times 100 = 11.1\% \quad \text{Mar} \\
 &= \frac{10}{100} \times 100 = 10\% \quad \text{Avg} \\
 &= \frac{90}{3} = 30
 \end{aligned}$$

100/9=11.1111

85%

Data Training Needs

$$\underline{500} \times 80\% = 400$$



$$\begin{cases}
 \text{Training} = 80\% \\
 \text{Validation} = 10\% \\
 \text{Testing} = 10\%
 \end{cases}$$

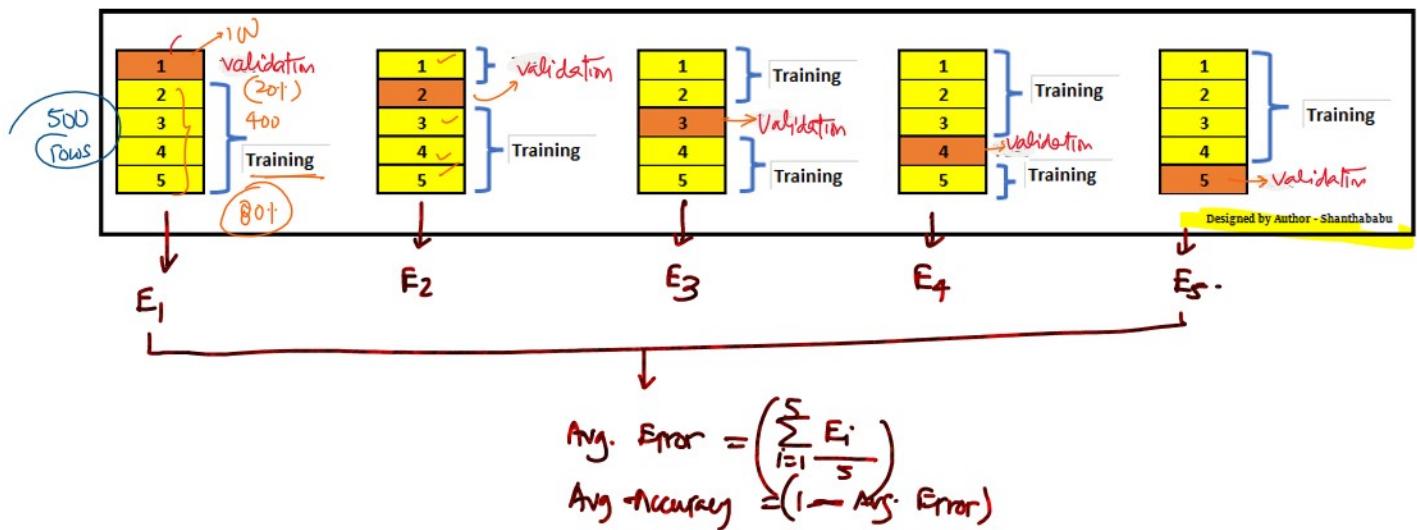
max
data quality is bad

min

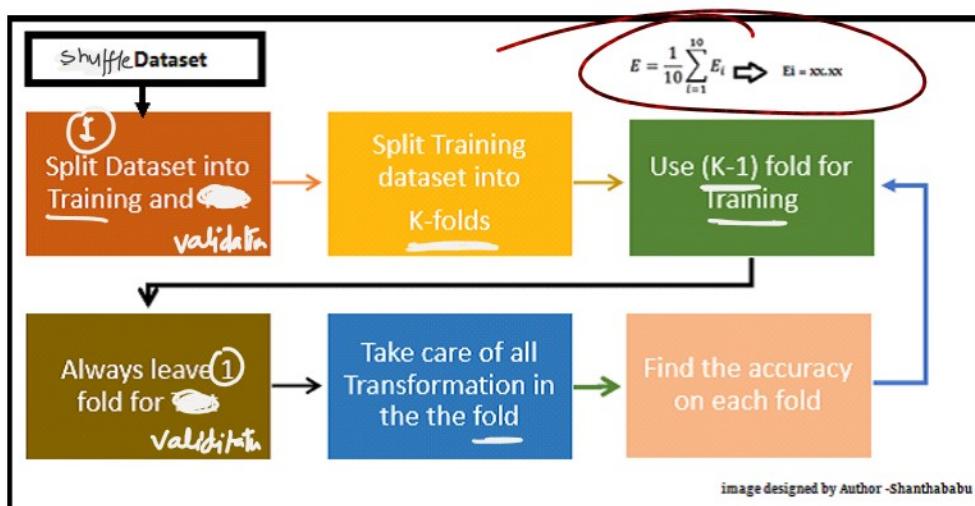
Validation Techniques

K-fold CV : K-fold cross validation

- it is a technique for evaluating the predictive models during training.
- dataset is divided into **k**-subsets or folds randomly.
- Model is trained and evaluated **K-times**, using a different fold as the validation set each time.



Life-cycle of K-fold CV



a Mo $\beta - \gamma$ M del

