

Feature Engineering

- is the process of creating new more informative features from existing raw data to improve a model's performance. (involves Domain knowledge + creativity)

Ex: instead of using raw Date values you might create features like Day of week / holiday. This can better explain something like sales (less or more on weekends?) helps model focus on relevant parts of data

Raw Data	Engineered Features
Date	Day, Is holiday
Jan 1, 2024	Sunday, Yes
Feb 3, 2024	Monday, No
Mar 5, 2024	Thursday, No

Feature Scaling (Normalization, standardization)

- is the process of Transforming numeric features to a similar scale, to prevent features with larger ranges from dominating the learning process

Ex: here the num's for salary are much larger than those for age and thus dominate model fitting

id	age	salary	id	norm age	norm salary
1	44	73000	1	0.8095	0.8387
2	27	47000	2	0	0.0
3	30	53000	3	0.1428	0.1935

Normalization
Ex: min = 27, max = 48 for age col
range (0-1)

Common methods: Standardization

$$x_{std} = \frac{x - \text{mean}(x)}{\text{standard deviation}(x)}$$
 \Rightarrow transforms mean = 0, std = 1

Normalization

$$x_{norm} = \frac{x - \min(x)}{\max(x) - \min(x)}$$
 \Rightarrow transforms range (0-1) or (-1-1)

- Why?

- in short many ML algo's are sensitive to scale of features by scaling we avoid bias in learning (larger feature = more important) and help some algo's like GD be faster (smaller vals)

- When to use what?

Scenario	norm	stand
Features have different scales	✓	✓
Data not gaussian	✓	✗
Data gaussian (normal dist)	✗	✓
Distance based model (KNN, SVM, GD)	✓	✓
NN	✓	✓ (sometimes)