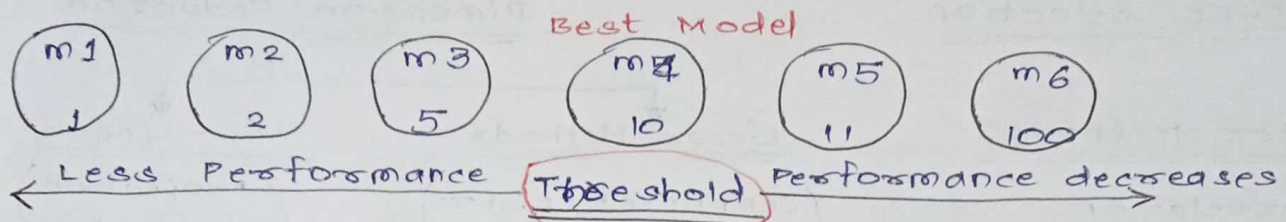


* Dimensionality Reduction *

① Dimensionality:- Features / columns in dataset.

② Curse of Dimensionality:-

↳ Generate models with increasing no. of features.



↳ Underfitting

↳ Less accurate

Threshold

Minimum num. of
required features for
the best model.

↳ May overfitting

↳ Increases complexity

↳ Decrease performance

↳ Feature not useful

↳ ↓ accuracy

↳ ↓ generalization performance

③ Dimensionality Reduction:-

↳ As dimensionality increases, accuracy increases but till a certain threshold, then after that start decreasing performance.

↳ Reduce features in dataset, while retain most of important info. as possible.

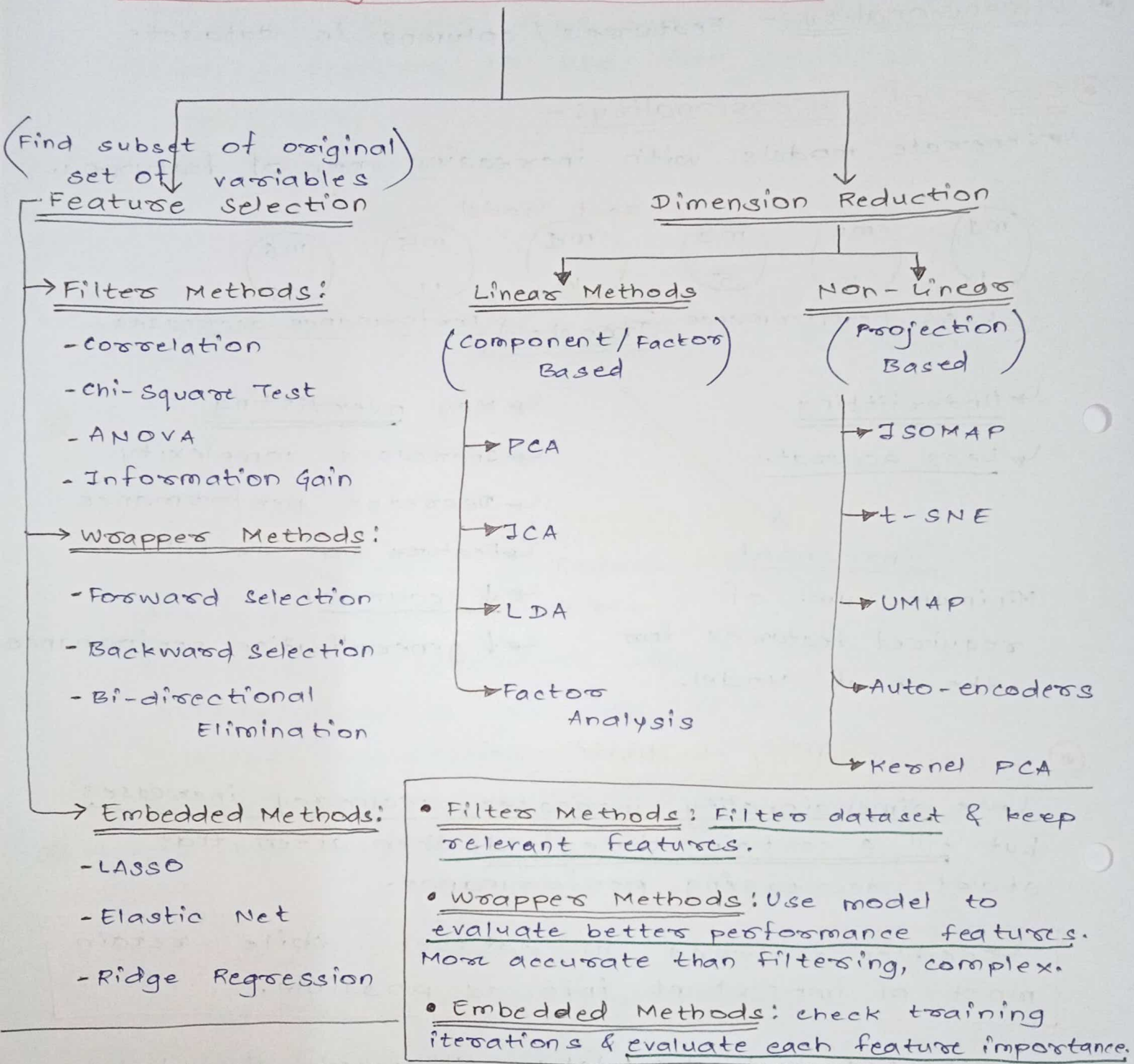
↳ convert high-dim. data \Rightarrow lower dim. data, still preserve essence of original data.

↳ Reduce complexity & improve generalization perf.

↳ Done during pre-processing, before building model to improve performance.

↳ Sometimes can discard useful info as well, so have to be careful.

Dimensionality Reduction Techniques



PCA (Principal Component Analysis)

LDA (Linear Discriminant Analysis)

ICA (Independent Component Analysis)

t-SNE (t-distributed stochastic Neighbor Embedding)

UMAP (Uniform Manifold Approximation & Mapping)

Kernel PCA (Kernel Principal Component Analysis)

② 2 components of Dimensionality Reduction:-

① Feature Selection: Find subset of original set of variables/features to use for training.

↳ Filter

↳ Wrappers

↳ Embedded

② Feature Extraction: Reduce data in high-dim space to low-dim space, with lesser no. of dimensions.

• Advantages of Dimensionality Reduction:

- ① Data compression \Rightarrow Reduce storage space.
- ② Reduce computation time.
- ③ Remove any redundant feature.
- ④ Easy visualize low-dim. data.
- ⑤ Reduce overfitting, Reduce complexity.
- ⑥ Improve model performance.
- ⑦ Reduce noise & irrelevant info.

• Limitations of Dimensionality Reduction:

- ① May loss some important data.
- ② Find linear correlations, sometimes undesirable.
- ③ May difficult understand relationship in original & reduced dimensions.
- ④ Sometimes may lead to overfitting.
- ⑤ Can be computational expensive, especially large datasets.