

- where will we get data from
- offline or batch mode
- what are we trying to solve
- who are our target audience
- cost

1. Frame the problem

- team size
- How will end product look
- is the model supervised or unsupervised
- what type of algorithms will help us

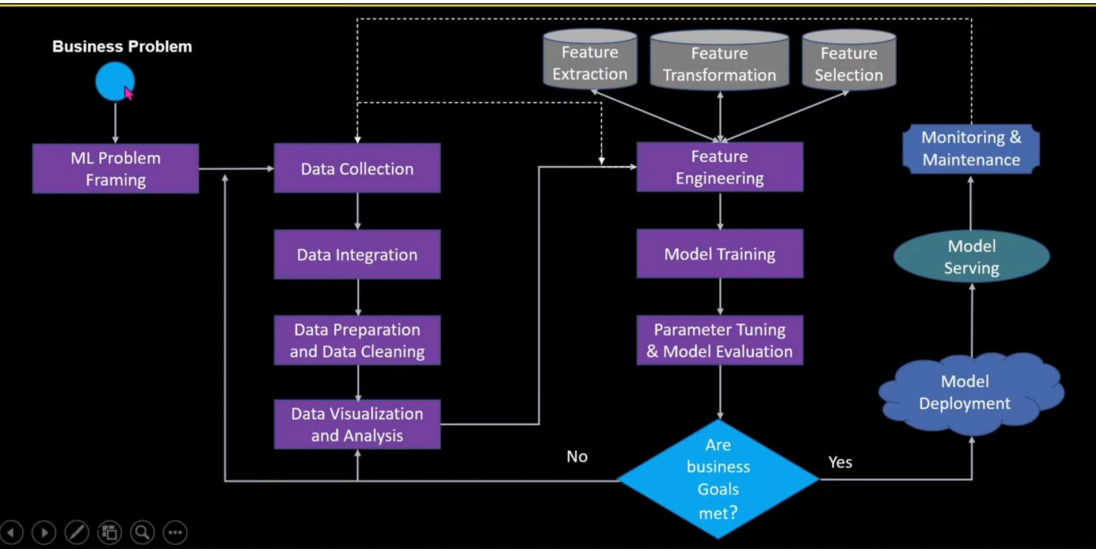
2. Gathering Data

- CSV
- API
- web scraping
- Database → Datawarehouse [ETL]
- Spark clusters

3. Data preprocessing

- Remove duplicates
- Remove missing values
- outliers
- scaling

4. EDA - study input & output relation, experiment with data vizualization, univariate / Bivariate / Multivariate Imbalanced data.



5. Feature Engineering / Selection

6. Model Training, Evaluation & selection

- Performance metrics
- Hyperparameter tuning

7. Model Deployment

8. Testing (Beta testing) A/B testing

9. Optimize.