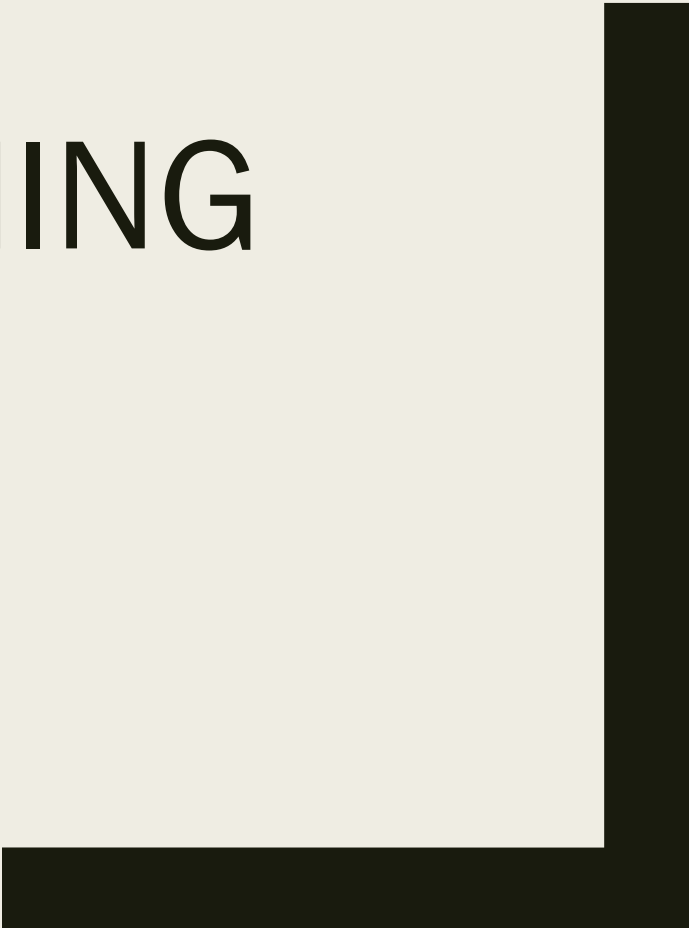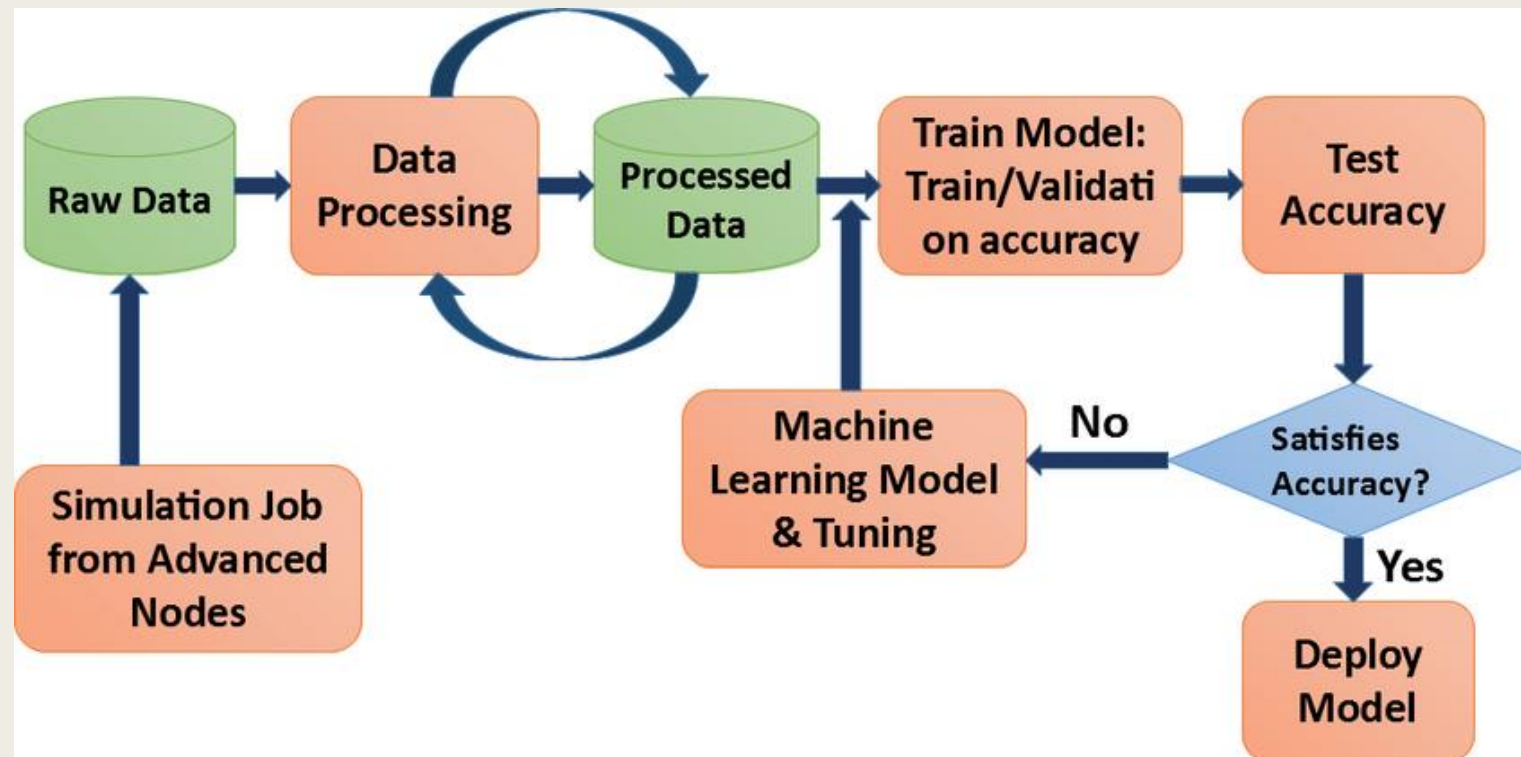# MACHINE LEARNING WORKFLOW
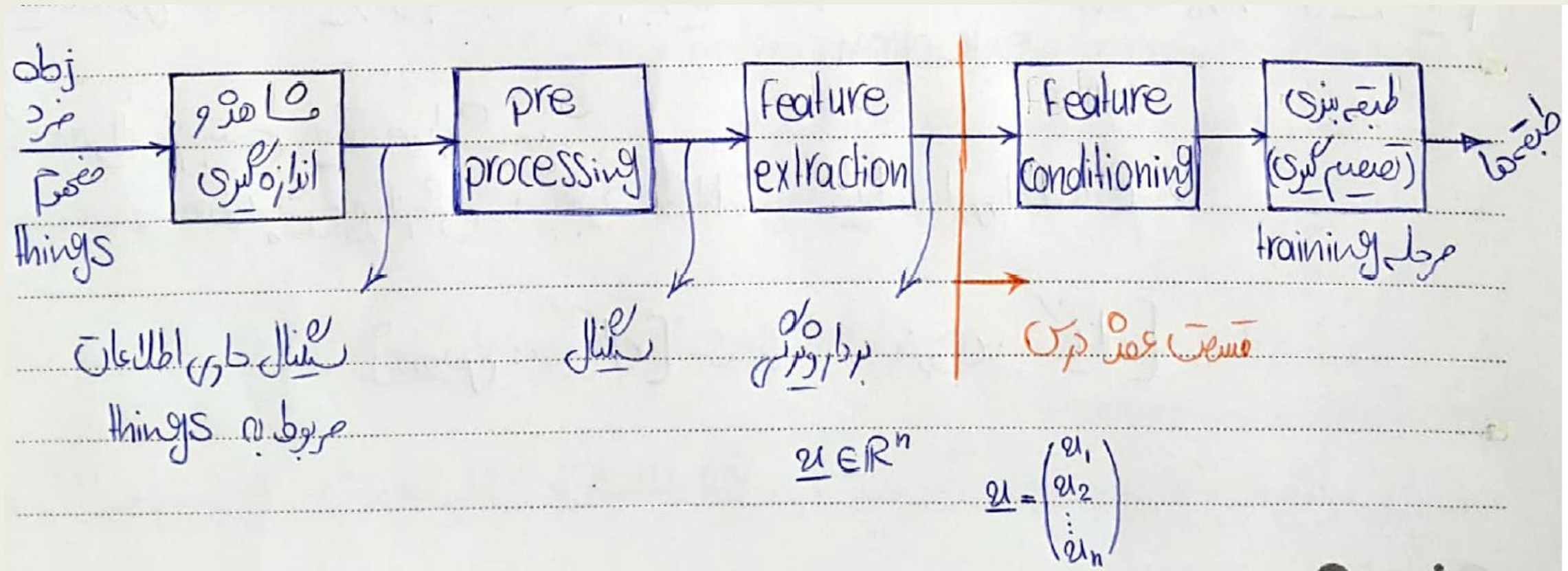
AI Summer School

University of Tehran

# Machine Learning Workflow
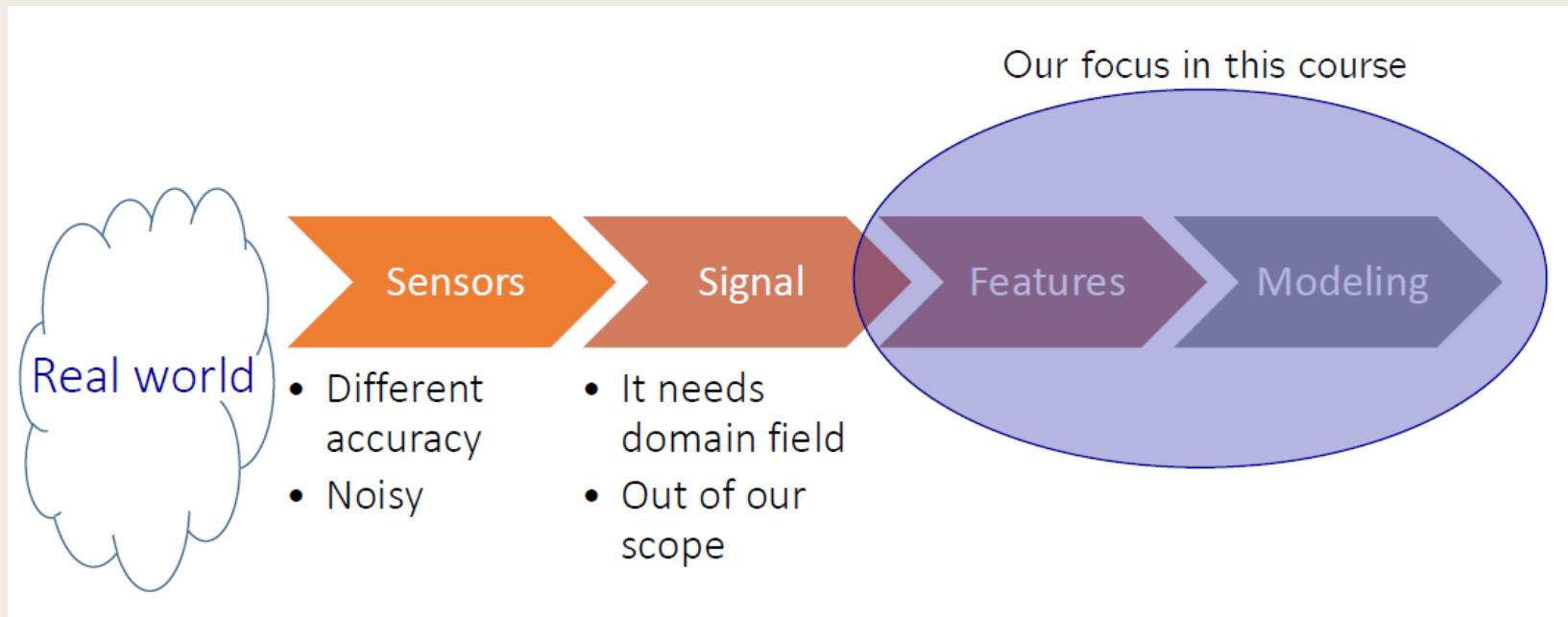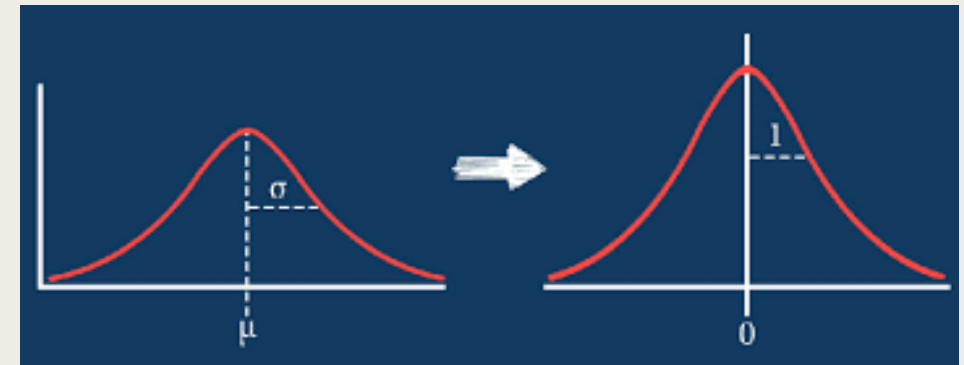
# Machine Learning Workflow
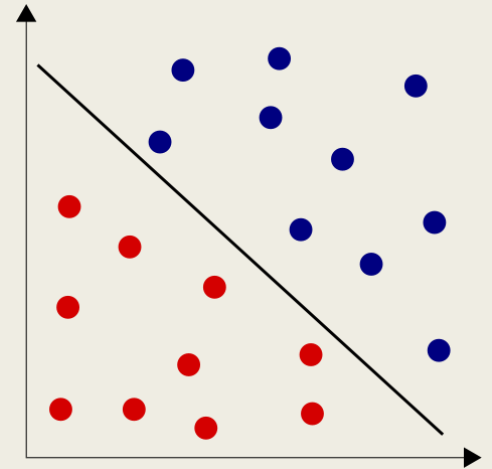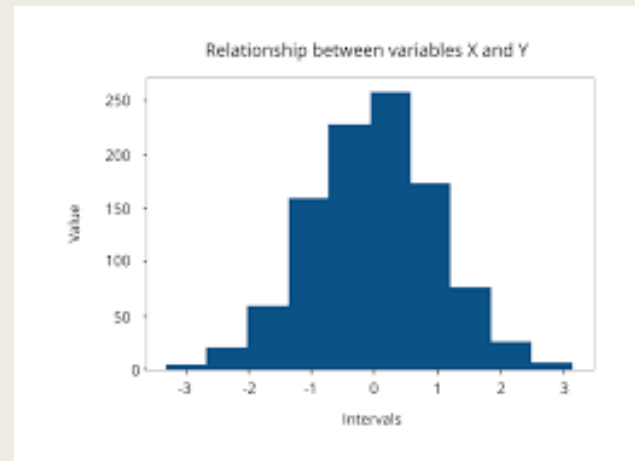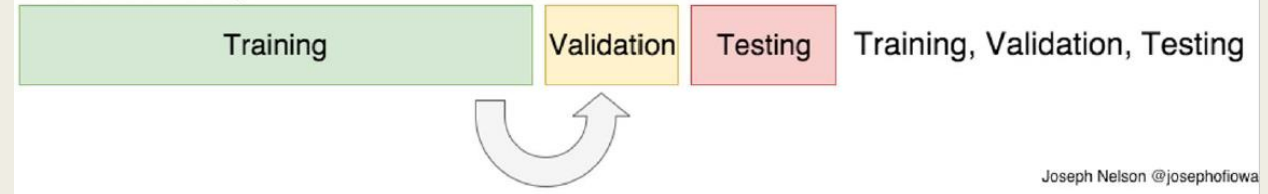
# Data Gathering

- Useful features

- Dataset size

- Data collection process

# Data Processing

- Train / Validation / Test split

- Data Visualization
  - *Class label distribution*
  - *Feature values distribution*
  - *Correlation between features*
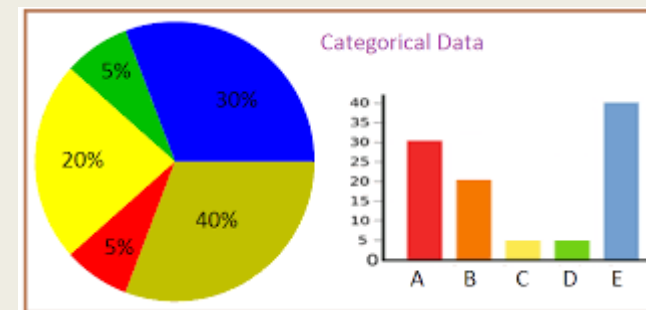  - *Separability of classes*

- Rescaling

# Data Processing

- Error correction
  - *Handling missing values*
  - *Handling inconsistent values*
  - *Handing outliers*

- Handling categorical features

- Feature engineering

# Data Processing

- Dealing with imbalanced classes:
  - *Oversampling*
  - *Undersampling*

- Data augmentation

# Data Processing

- Preprocessing for text:
  - *Normalization*
  - *Tokenization*
  - *Stop word deletion*
  - *Stemming*
  - *...*

- Preprocessing for image:
  - *Resizing*
  - *Converting the format: RGB, HSB, gray scale*
  - *Normalization*
  - *....*

# Training and Evaluation

- ■ Model selection

- ■ Hyper parameter tuning

- ■ Selecting proper evaluation metric

- ■ Fighting with overfitting

# Error Analysis