# Universidad Politécnica de Yucatán

## MACHINE LEARNING

## SOLUTION TO MOST COMMON PROBLEMS IN ML

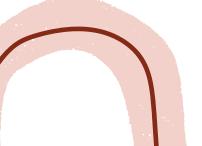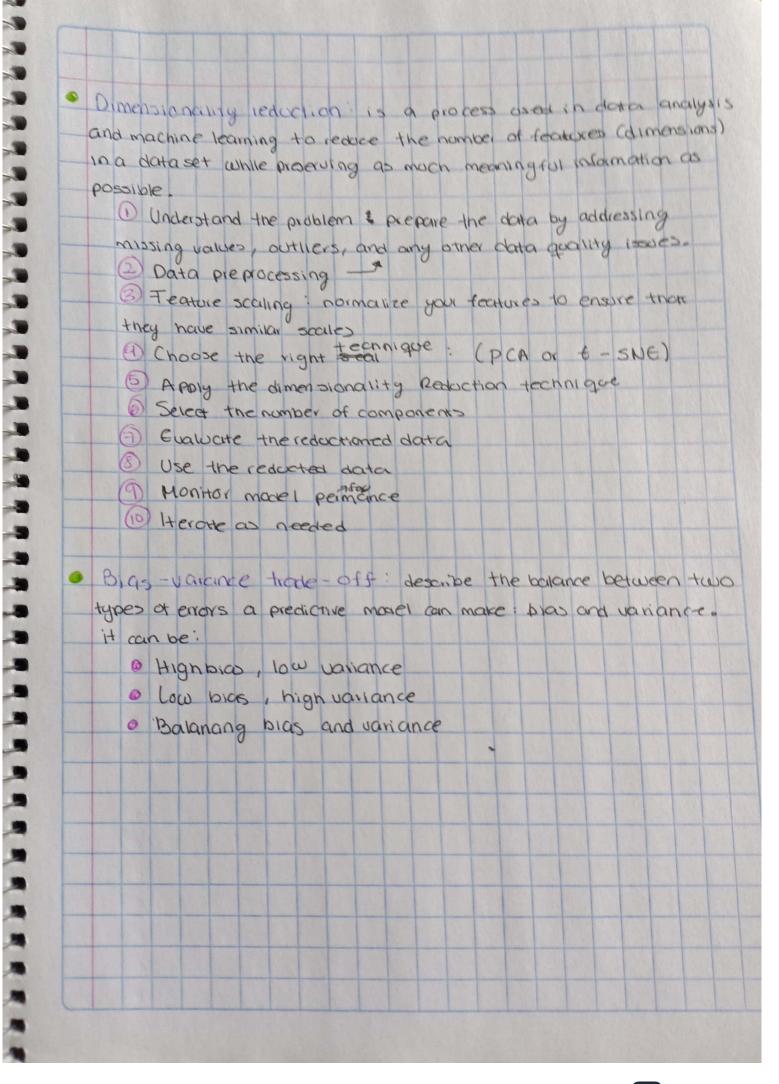GABRIELA ELIZABETH AVILA CHAN

2009003

IRC 9°B

- **Overfiting**: occurs when the machine learning model learns the training data too well, to the point that it captures noise, random fluctuations, or outliers in the data, rather than the underlying patterns. As a result, an overfiting model perfoms very well on training data but poorly on unseen data because it has essentially memorized the training examples instead of learning the true underlying relationships. Signs of it:
  - Low training error (the model fits the training data very closely)
  - High validation or test error (poor performance on new data)
  - The model's predictions are too sensitive to small changes in the input data.

- **Underfiting**: occurs when a machine learning model is too simple to capture the underlying patterns in the training data. It fails to learn the relationships between the features and the target variable, resulting on poor performance on both the training data and new data. It lacks of the capacity to represent the complexity of the data. Signs of it:
  - High training error (the model cannot even fit training data)
  - High validation or test error
  - The mode's predictions do not capture the true patterns in the data.

- **Outliers**: data points that significantly deviate from the majority of the data in a dataset. Some of its key characteristics
  - Extreme values
  - Unusual a rare
  - Impact on summary statistics (can significantly affect summary statistics)
  - Visual identification (are often visually identifiable on plots)
  - Data errors (can sometimes be the result of data entry errors, measurement errors)
  - Influence on models
  - Domain knowledge (to be sure when is an outlier or when should be treated as an error)

- Common solutions for them
  - Overfiting
    - Regularization = (apply techniques like L1 or L2 regularization to penalize larg models coef).
    - Cross-validation (use techniques like K-fold cross-validation to asses model peisormance)
    - Reduce model complexity.
    - Early stopping (monitor the models peisormance during training data and stop when peisormance degrade)
    - Ensemble methods

  - Underfiting
    - Increase model complexity
    - Feature engineering (add more relevant features)
    - Collect more data set
    - Hyperparameter tuning
    - Ensemble methods

  - Presence of outliers
    - transformations (log-transformations or winsorization)
    - Robust models (use models that are less sensitive to outliers)
    - Imputation
    - Data cleaning
    - Contextual understanding

- Dimensionality problem: it happens when dealing with high-dimentional data sets or feature spaces. It encompasses various issues, here's an overview of the dimensionarity problem:
  - Increased computational complexity
  - Data sparsity
  - Increased risk of overfitting
  - Difficulty in visualization
  - Loss of intuition

- **Dimensionality reduction** is a process used in data analysis and machine learning to reduce the number of features (dimensions) in a dataset while preserving as much meaningful information as possible.

  ① Understand the problem & prepare the data by addressing missing values, outliers, and any other data quality issues.

  ② Data preprocessing ⟶

  ③ Feature scaling: normalize your features to ensure that they have similar scales

  ④ Choose the right technique: (PCA or t-SNE)

  ⑤ Apply the dimensionality Reduction technique

  ⑥ Select the number of components

  ⑦ Evaluate the reductioned data

  ⑧ Use the reducted data

  ⑨ Monitor model performance

  ⑩ Iterate as needed

- **Bias-variance trade-off**: describe the balance between two types of errors a predictive model can make: bias and variance. it can be:

  - High bias, low variance
  - Low bias, high variance
  - Balanang bias and variance

Hertel, P. (n.d.). *Common problems with machine learning that companies face*. Hyperon.Io. Retrieved September 15, 2023, from https://www.hyperon.io/blog/common-problems-with-machine-learning-that-companies-face

*Issues in Machine Learning*. (n.d.). Www.javatpoint.com. Retrieved September 15, 2023, from https://www.javatpoint.com/issues-in-machine-learning