

# Midterm Review

## 1. Machine Learning Understanding

- Supervised Learning vs. unsupervised Learning

Supervised learning has a clear objective such detecting spam email, predicting prices

Unsupervised learning aims to find a pattern, anomalies, or clusters.

- Regression vs. classification

Regression predicts a value, while classification predicts a label.

- Training set vs. test set

The training set is used to train the model, while the test set is used to evaluate the model's performance. We should split the data set into a training set and a test set before the training process. The test set is not used to train the model, so it represents new data.

- How to evaluate a regression model?

MSE (mean-squared-error), MAE (mean-absolute-error), Visualize the model and see how it fits the data.

- How to evaluate a classification model?

Accuracy, Precision, Recall, Confusion Matrix, F-1 score, Visualization of the decision regions..

- Cross validation

Cross validation eliminates the possibility that the good performance only occurs for a specific training set.

- Which model is a proper choice for a regression task: linear regression, polynomial regression
- Which model is a proper choice for a classification task: logistic regression

## 2. Machine Learning Models

- Linear regression
- Polynomial regression
- Logistic regression

For each model, you are expected to explain:

- The mathematical expression for that model
- What cost function can be used to measure the model's performance during the training phase?
- How to find parameter values that correspond to an optimal cost?
- Work on a simple example

Index	(1)	(2)	(3)	(4)
X (input variable)	1	2	3	4

Y (output variable)	4	4.9	6.2	8
---------------------	---	-----	-----	---

## Linear Regression:

Intuition: We use a line (a plane, or a hyper-plane) to describe the relationship between the input variables  $X_1, X_2, \dots, X_n$  and the output variable  $y$

Model expression

$$Y = mX + b$$

If there is a single input variable

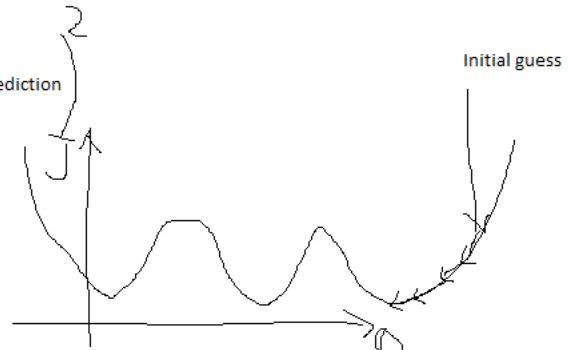
$$Y = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n$$

MSE cost function:

$$J(\theta_0, \dots, \theta_n) = \frac{1}{N} \sum (\text{actual output} - \text{prediction})^2$$

How to find the minimal cost?

1. Use the normal equation
2. Use the gradient descent method



## Polynomial Regression:

Intuition: Use a polynomial curve to describe the relationship between the input and the output.

Model expression: single variable, degree 4

$$y = \theta_0 + \theta_1 X + \theta_2 X^2 + \theta_3 X^3 + \theta_4 X^4$$

Model expression: suppose there are two variables  $x_1$  and  $x_2$ , and the largest degree is 3.

$$y = \theta_{00} + \theta_{11} X_1 + \theta_{12} X_2 + \theta_{21} X_1^2 + \theta_{22} X_1 X_2 + \theta_{23} X_2^2 + \theta_{31} X_1^3 + \theta_{32} X_1^2 X_2 + \theta_{33} X_1 X_2^2 + \theta_{34} X_2^3$$

MSE cost function  $J(\theta) = \frac{1}{N} \sum (\text{actual output} - \text{model prediction})^2$


Training algorithm:

1. Using normal equation (with polynomial features)
2. Gradient descent method

## Logistic Regression

Intuition: Predict probability distribution over all possible classes. The probability distribution is represented as a linear combination of all inputs, following by a transformation that guarantees the output to be probabilities (logistic function or the softmax function).

Model expression: binary classifier

$$P = \sigma(\theta_0 + \theta_1 X_1 + \dots + \theta_n X_n) \quad \sigma(t) = \frac{1}{1 + e^{-t}}$$


Model expression: multiple classes

$t_1, \dots, t_n =$  linear combinations of input variables

$$p_1, \dots, p_n = \sigma(t_1, \dots, t_n) \quad \sigma: \text{softmax function}$$

Cross-Entropy Cost function

Training algorithm: gradient descent

## 3. Python Libraries

- NumPy: NumPy array, mathematical functions, matrix functions
- Pandas: Pandas data frame, data handling, data transformation
- Matplotlib: line plot, scatter plot, histogram
- Sklearn: models, model evaluations

TODO: go over the class notes and summarize the functions we have used.