

## New Era University College of Informatics and Computing Studies Department of Computer Science



### **FREE ELECTIVE 3**

CCSEL3-18



# Supervised Machine Learning: Regression and Classification

### coursera

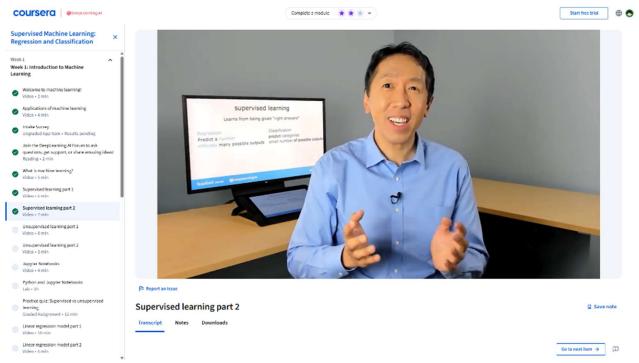
**ONLINE COURSE DOCUMENTATION** 



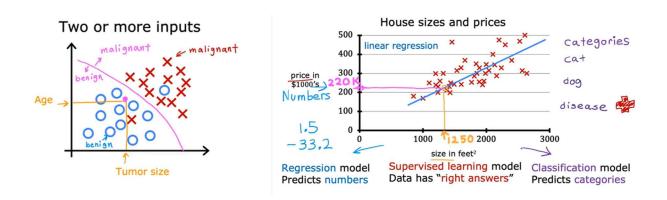
Marasigan, Vem Aiensi A.
BS Computer Science | 21-11295-310

#### Module 1: Introduction to Machine Learning

[Completed by September 8 – 9]

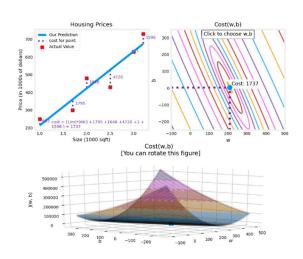


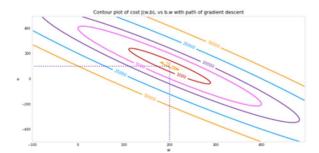
The first module/week for the course introduces machine learning concepts and foundational knowledge on how machine learning works as well as its applications. It showed different forms of machine learning which are; Regression, for predicting numeric values, and Classification, for predicting a category, and, how supervised learning differs from unsupervised learning.



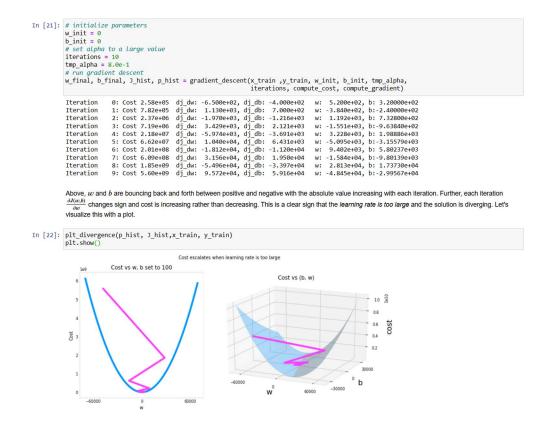
Proceeding, this module's topic is mostly around Regression, covering the most basic form of how predictions are achieved using Linear Regression Algorithm. This module discussed primary knowledge, from the linear function, which is also called the model, up to choosing the correct w

and b, and learning rate for the gradient descent function. This explained the functions needed in Linear Regression Algorithm which are: the model's f(x), the cost function that determines whether the w and b is good, and finally the gradient descent function which helps determine the best w and b that achieve the minimum value of cost function.



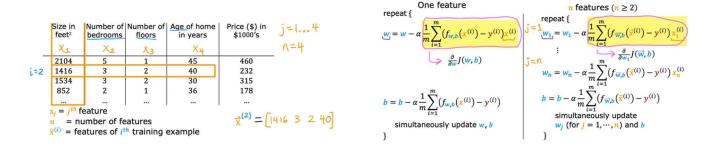


Also, learning all these functions was quite overwhelming at first. But with the aid of the laboratory Jupiter notebooks, with visualizations and interactivity, learning the concepts became much easier.



**Module 2:** Regression with multiple input variables

[Completed by September 16]



Taking a step further, learning multiple linear regression involved a lot of concepts. Considering multiple features, optimizations are introduced in this module. These are **vectorization**, **feature scaling**, **feature engineering**, and a little bit of polynomial regression example.

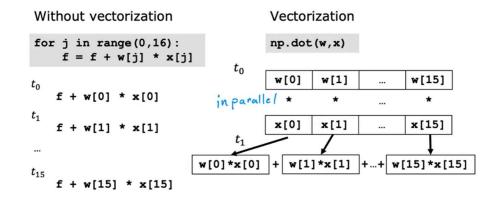
#### Vectorization

Due to the nature of samples having multiple features  $(x_1, x_2, x_3 ... x_n)$ , the resulting model would be:  $f_{\overrightarrow{w},b}(\overrightarrow{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b \quad \text{which can range up to thousands of features.}$ Say 1000 features, this approach would be inefficient considering that it must iterate in every cost function.

Hence, this is where Vectorization comes in by using functions for sets, on of which is the dot product function. This can be achieved in code as well using the NumPy library's dot function

$$f_{\overrightarrow{\mathbf{w}},b}(\overrightarrow{\mathbf{x}}) = \overrightarrow{\mathbf{w}} \cdot \overrightarrow{\mathbf{x}} + b$$
  $\mathbf{f} = \text{np.dot}(\mathbf{w},\mathbf{x}) + \mathbf{b}$ 

This makes the formula representation cleaner and especially execute the code much faster.



#### **CCSEL3-18 | FREE ELECTIVE 3**

Online Course Documentation

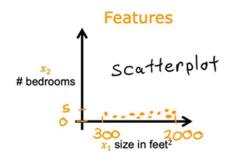
Supervised Machine Learning: Regression and Classification

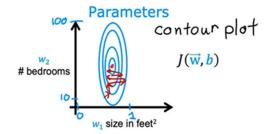
#### **Feature Scaling**

Another thing to consider in multiple linear regression would be the range of values for each feature. Having a large number in range in other features affects how gradient descent works to achieve minimum. Say:

$$300 \le x_1 \le 2000$$
  $0 \le x_2 \le 5$ 

If the difference in the ranges is far from each other, these might introduce inefficiency or overshooting when performing gradient descent to reach minimum.





So, the thing to do here is to transform feature values so that features can be comparable to other features' range of values. This is called feature scaling, can be normalization as well, and when features are rescaled first, this would be the gradient descent:



There are many ways to rescale features, these are:

**Basic Feature scaling** 

Mean normalization

**Z-score normalization** 

$$x_{1,scaled} = \frac{x_1}{2000}$$

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

Rescaling all features not only makes a better representation of a graph but most importantly help gradient descent find the minimum costs for each  $w_s$  considering multiple features.

#### **Feature Engineering**

Having multiple features to train a model, another we can do is to add custom features based on existing ones. These relies on our intuition whether if some features are combined can have significance for predictions.



As a basic example from the course, say existing features are:

x₁ is frontage and,

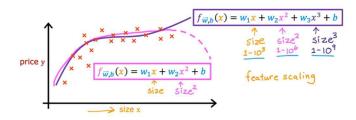
x₂is depth

$$area = frontage \times depth$$

Thinking that the area of the lot has a significance, say where predicting its price, we may create another feature "area" to be our  $x_3$ 

$$f_{\vec{\mathbf{w}},b}(\vec{\mathbf{x}}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

By the end of the module, we get a glimpse of **Polynomial Regression**. This is used for scenarios where there is a tragectory where a curve is much more suitable instead of a straight line, linear regression.

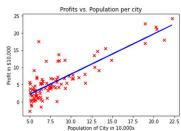


```
In [27]: # Plot the linear fit
plt.plot(x_train, predicted, c = "b")

# Create a scatter plot of the data.
plt.scatter(x_train, y_train, marker='x', c='r')

# Set the title
plt.title("Profits vs. Population per city")
# Set the y-axis label
plt.ylabel('Profit in $10,000')
# Set the x-axis label
plt.xlabel('Population of City in 10,000s')

Out[27]: Text(0.5, 0, 'Population of City in 10,000s')
```



And finally, a laboratory activity practicing what was learned so far in linear regression where we implement functions in code.

```
# You need to return the following variables correctly
dj_dw = 0
dj_db = 0

### START CODE HERE ###
for i in range(m):
   dj_dw += (w * x[i] + b - y[i]) * x[i]
   dj_db += w * x[i] + b - y[i]
dj_dw /= m
dj_db /= m
```

Module 3: Classification

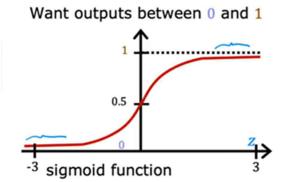
[Completed by September 17]

From the previous lessons about Linear Regression for predicting numerical values, this module now moves on to predicting classes or categories. This is done through Logistic Regression. Although it is slightly different from linear regression, most concepts and knowledge from previous lessons do carry-over in Logistic Regression, from computing gradient descent, rescaling, and even feature engineering. Still, these are the key differences:

**Logistic Function** model, or the **Sigmoid** function

$$f_{\overrightarrow{\mathbf{w}},b}(\overrightarrow{\mathbf{x}}) = g(\overrightarrow{\mathbf{w}} \cdot \overrightarrow{\mathbf{x}} + b) = \frac{1}{1 + e^{-(\overrightarrow{\mathbf{w}} \cdot \overrightarrow{\mathbf{x}} + b)}}$$

Considering classification, having only 2 values, either 0 or 1, true or false, or in short binary outputs, a linear regression model won't be reliable. Hence, we can use the sigmoid function that produces outputs ranging from 0 to 1, then setting up a threshold or decision boundary to determine if the prediction is positive, or if it is negative.



Also, having a new model also **Loss function** 

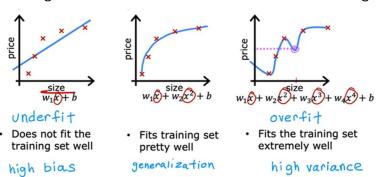
changes the cost function used in the gradient descent. This incorporates another function which is called the 
$$\begin{cases} -\log\left(f_{\overrightarrow{w},b}(\overrightarrow{x}^{(i)})\right) & \text{if } y^{(i)} = 1 \\ -\log\left(1 - f_{\overrightarrow{w},b}(\overrightarrow{x}^{(i)})\right) & \text{if } y^{(i)} = 0 \end{cases}$$

And for simplicity, for programming as well, this can also be written as:

$$\mathbf{y}^{(i)} \log \left( f_{\overrightarrow{\mathbf{w}},b}(\overrightarrow{\mathbf{x}}^{(i)}) \right) + \left( 1 - \mathbf{y}^{(i)} \right) \log \left( 1 - f_{\overrightarrow{\mathbf{w}},b}(\overrightarrow{\mathbf{x}}^{(i)}) \right)$$

**Decision Boundary** – This is the main objective of using logistic regression, this decision boundary is the resulting graph of the logistic regression model that will determine the category of the sample being predicted.

Another concept introduced in this module was **Generalization**, which applies to both Linear and Logistic Regression. This involves the decision of on what degree should the model be to create a



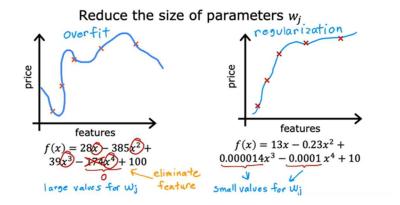
"just right" curve of predictions for a particular dataset. This generalization is done to avoid underfitting and overfitting and is achieved by regularization. This procedure involves many ways which are:

Eliminating other features,

features that seems to be "not that significant". This way, the model can reduce being overfitted. However, what if all features do mean significant factor for prediction, this is where regularization comes in.

**Regularization** is done by reducing the values of the parameters, this is the  $\boldsymbol{w}$  for every feature  $\boldsymbol{x}$ . This is done by incorporating the equation:  $\frac{\lambda}{2m}\sum_{i=0}^{n-1}w_j^2$ 

With this the resulting formula for cost function and gradient descent would be respectively:



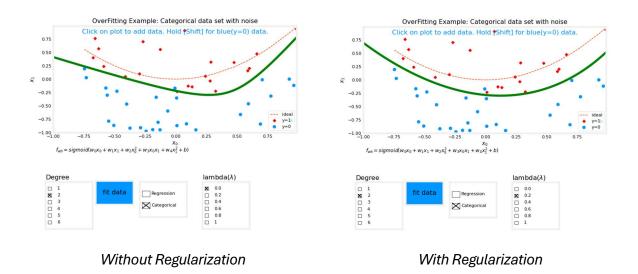
Cost Function: 
$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=0}^{m-1} \left[ -y^{(i)} \log \left( f_{\mathbf{w}, b} \left( \mathbf{x}^{(i)} \right) \right) - \left( 1 - y^{(i)} \right) \log \left( 1 - f_{\mathbf{w}, b} \left( \mathbf{x}^{(i)} \right) \right) \right] + \frac{\lambda}{2m} \sum_{j=0}^{m-1} w_j^2$$

Gradient Descent: 
$$\frac{\partial J(\mathbf{w},b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w},b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$
$$\frac{\partial J(\mathbf{w},b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w},b}(\mathbf{x}^{(i)}) - y^{(i)})$$

#### CCSEL3-18 | FREE ELECTIVE 3

Online Course Documentation

The key variable in this regularization is the **lambda**. Setting small values for lambda will reduce the features' influence in the model, and vise-versa, effectively reducing overfitting when a "just right value is used" achieving a generalized model. Here's a comparison below from the module:



This module also involved a practice lab to implement Logistic Regression. Here's a small snapshot of my take:

```
In [41]: # UNQ_C6
def compute_gradient_reg(X, y, w, b, lambda_ = 1):
                                                                                                                                                                  In [44]: plot_decision_boundary(w, b, X_mapped, y_train)
                                                                                                                                                                                   # Set the y-axis Label
plt.ylabel('Microchip Test 2')
                     Computes the gradient for logistic regression with regularization
                                                                                                                                                                                   # Set the x-axis label
plt.xlabel('Microchip Test 1')
                     Args:

X: (ndarray Shape (m,n)) data, m examples by n features
y: (ndarray Shape (m,)) target value
w: (ndarray Shape (n,)) values of parameters of the model
b: (scalar)
lambda : (scalar, float)
Returns
                                                                                                                                                                                   plt.legend(loc="upper right")
plt.show()
                        dj_{\cdot}db: (scalar) The gradient of the cost w.r.t. the parameter b. dj_{\cdot}dw: (ndarray Shape (n,)) The gradient of the cost w.r.t. the parameters w
                                                                                                                                                                                           1.0
                                                                                                                                                                                     Test 2
                                                                                                                                                                                          0.5
                     m, n = X.shape
                     dj_db, dj_dw = compute_gradient(X, y, w, b)
                                                                                                                                                                                           0.0
                      ### START CODE HERE ###
                     ### ### FNAN CODE MERE ###

for i in range(n):
    reg = w[i] * (lambda_ / m)
    dj_dw[i] += reg

### END CODE HERE ###
                                                                                                                                                                                         -0.5
                     return dj_db, dj_dw
```

Note: All these laboratory activities, as well as the optional labs, are also uploaded. You may check that out as well by scanning the QR code at the front page.

#### **CCSEL3-18 | FREE ELECTIVE 3**

Online Course Documentation

#### **COMPLETION RECORD**

#### **Machine Learning**



