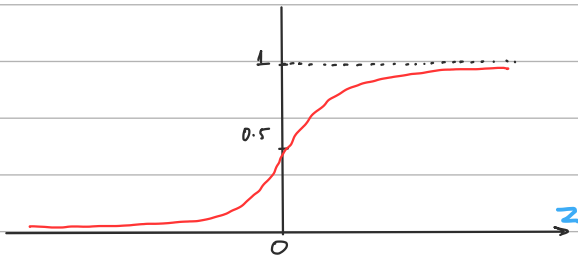


Week 3

0 Classification with Logistic Regression.

We use Sigmoid Function to classify values as 0 or 1



We need an equation which results value in the range 0 to 1, hence we use sigmoid function

$$g(z) = \frac{1}{1 + e^{-z}}$$

□ Logistic Regression Model

$$f_{\vec{w}, b}(\vec{x}) = z = \vec{w} \cdot \vec{x} + b$$

As,

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x} + b)$$

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

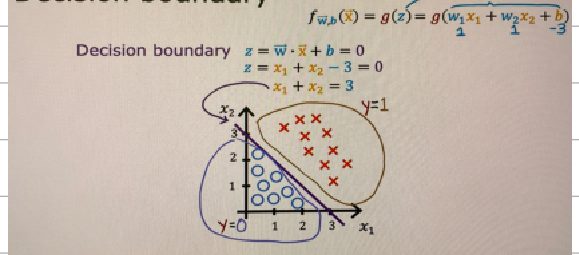
- It's also represented as,

$$f_{\vec{w}, b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w}, b)$$

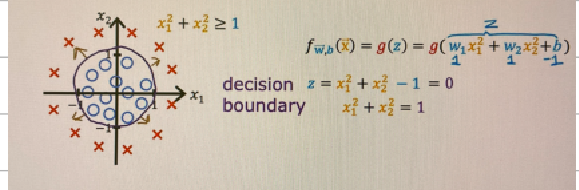
Which means, what is probability that y is 1 given input \vec{x} and parameters \vec{w}, b

0 Decision Boundary

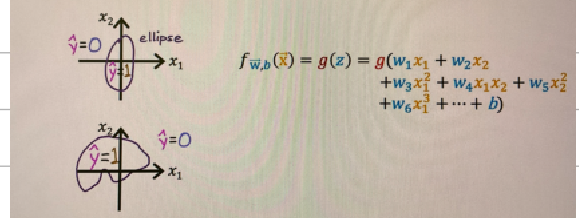
Decision boundary



Non-linear decision boundaries



Non-linear decision boundaries



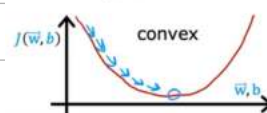
0 Why can't we use Squared error cost function for Classification?

Squared error cost

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

linear regression

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$



logistic regression

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



- As seen in the image above, Squared error cost function would create a non-convex curve with multiple local minima.

o Logistic Loss Function

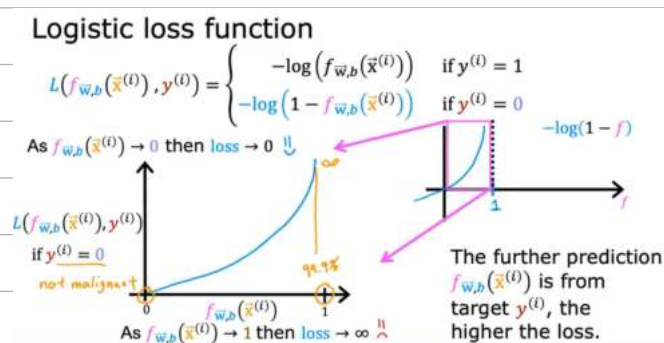
- we find loss for each example & then get the average of all the losses

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})}_{\text{Loss}}$$

Here Loss Equals,

$$L = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

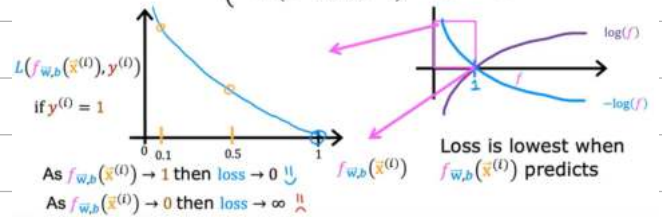
o Logistic Loss for $y^{(i)} = 0$



o Logistic Loss for $y^{(i)} = 1$

Logistic loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$



o Simplified Loss Function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

- This is a simple extension of the given conditional value which further helps in simplifying the cost function.

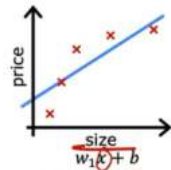
o Simplified Cost Function for Logistic Regression.

$$\begin{aligned} \text{loss } L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) &= -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \\ \text{cost } J(\vec{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})] \\ &= \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))] \end{aligned}$$

◦ Fitting issue in Linear Regression.

A. Underfitting :-

This happens when the model has few features.



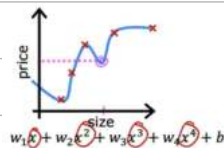
underfit

- Does not fit the training set well

high bias

B. Overfitting :-

This happens when the model has a lot of features



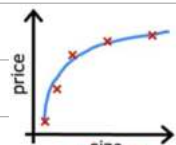
overfit

- Fits the training set extremely well

high variance

C. Generalization :-

This happens when the model has the perfect amount of features.

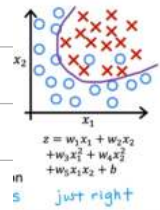


just right

- Fits training set pretty well

generalization

— Generalized Fitting →



◦ Ways to handle Overfitting

- Collect more data. If that's not possible then
- Improve feature selection. But it means you need to sacrifice data of few features. Hence you can use
- * Regularization :- Reduce size of parameters.

◦ Regularization in Regression.

- We penalise all the features by multiplying it with 'λ'.

$$\frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

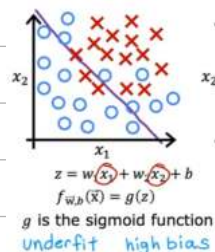
Here 2m is used so that regularization & cost function is scaled equally

Cost function after Regularization

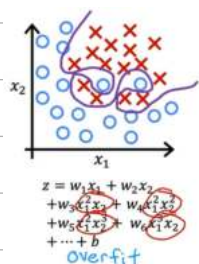
$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

◦ Fitting issue in Classification

— Underfitting →



— Overfitting →



0 Regularized Linear Regression & Gradient Descent

Regularized cost function is given as:-

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2}_{\text{Squared Error Cost function}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{Regularization}} \right]$$

Gradient Descent on Regularized cost function:-

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$
$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

} simultaneous update

don't have to regularize b

Let's expand w_j :-

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j \right]$$
$$= w_j - \alpha \frac{\lambda}{m} w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$
$$= w_j \left(1 - \alpha \frac{\lambda}{m} \right) - \text{usual update}$$

Let's assume $\alpha = 0.01$, $\lambda = 2$, $m = 50$

$$\therefore w_j \left(1 - 0.01 \times \frac{2}{50} \right) = w_j (1 - 0.0002) = w_j (0.9998)$$

Hence regularization basically makes sure that we slowly reduce w_j and reduce overfitting.

★ Derivation for $\frac{\partial}{\partial w_j} J(\vec{w}, b)$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} J(\vec{w}, b) &= \frac{\partial}{\partial w_j} \left[\frac{1}{2m} \sum_{i=1}^m \underbrace{(f(\vec{x}^{(i)}) - y^{(i)})^2}_{\vec{w} \cdot \vec{x}^{(i)} + b} + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \\
 &= \frac{1}{2m} \sum_{i=1}^m \left[(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \cancel{2x_j^{(i)}} \right] + \frac{\lambda}{2m} \cancel{2w_j} \quad \text{No } \sum_{j=1}^n \\
 &= \frac{1}{m} \sum_{i=1}^m \left[(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j \\
 &= \frac{1}{m} \sum_{i=1}^m \left[\underbrace{(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)})}_{f(\vec{x})} x_j^{(i)} \right] + \frac{\lambda}{m} w_j \\
 &= \frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j
 \end{aligned}$$

◦ Regularized Logistic Regression & Gradient Descent.

- we do the same thing we did for linear regression

★ Cost function looks like this:-

$$\min_{\vec{w}, b} J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

★ In the same way, updated gradient descent looks like this:-

Gradient descent

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

}

Looks same as for linear regression!

$$= \frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j$$

$$= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

don't have to re