

Week 3

Motivations

Classification

Question	Answer "y"
Is this email <u>spam</u> ?	no yes
Is the transaction <u>fraudulent</u> ?	no yes
Is the tumor <u>malignant</u> ?	no yes

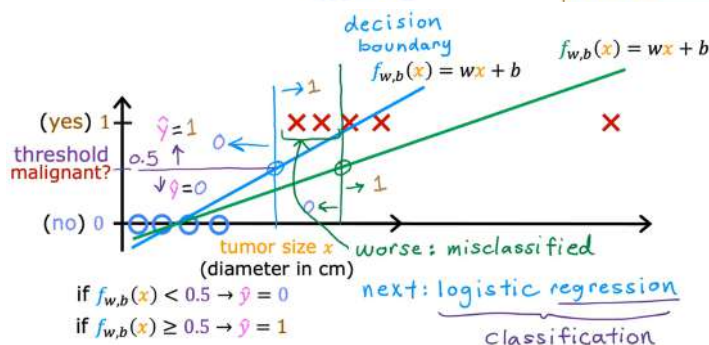
y can only be one of two values

"binary classification"

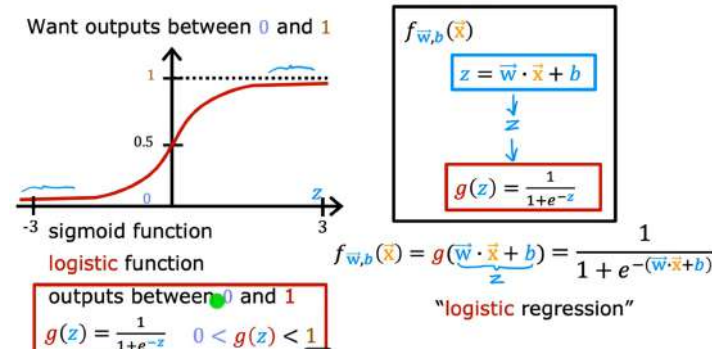
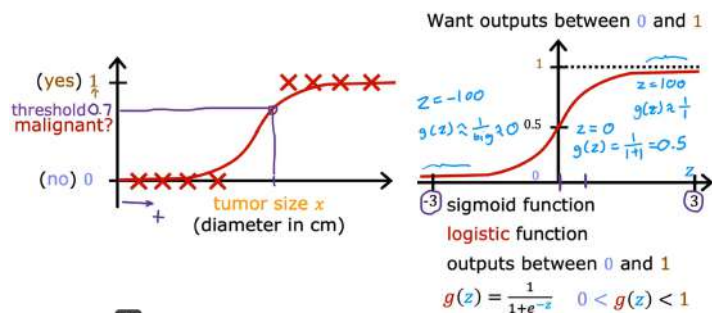
class = category

false true
0 1
useful for classification

"negative class" "positive class"
≠ "bad" ≠ "good"
absence presence



Logistic regression



Interpretation of logistic regression output

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

"probability" that class is 1

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

Probability that y is 1, given input \bar{x} , parameters \bar{w}, b

Example:

x is "tumor size"
y is 0 (not malignant)
or 1 (malignant)

$$P(y = 0) + P(y = 1) = 1$$

$f_{\bar{w},b}(\bar{x}) = 0.7$
70% chance that y is 1

C1_W3_Lab02_Sigmoid_function_Soln

- we would like the predictions of our classification model to be between 0 and 1 since our output variable y is either 0 or 1.
- This can be accomplished by using a "sigmoid function" which maps all input values to values between 0 and 1.

```
# Input is an array.
input_array = np.array([1,2,3])
exp_array = np.exp(input_array)
print("Input to exp:", input_array)
print("Output of exp:", exp_array)
# Input is a single number
input_val = 1
exp_val = np.exp(input_val)
print("Input to exp:", input_val)
print("Output of exp:", exp_val)
#Input to exp: [1 2 3]
#Output of exp: [ 2.72  7.39 20.09]
#Input to exp: 1
#Output of exp: 2.718281828459045

def sigmoid(z):
    g = 1/(1+np.exp(-z))

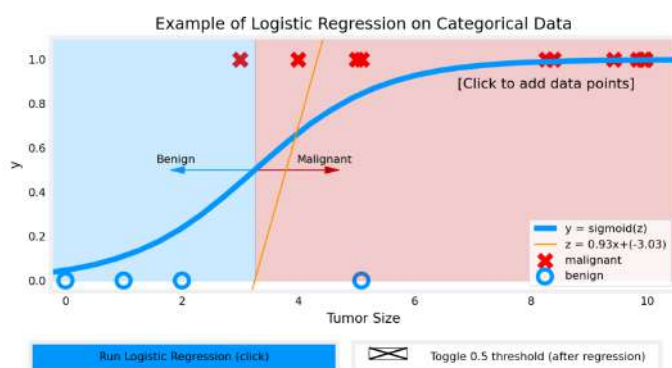
    return g

# Generate an array of evenly spaced values between -10 and 10
z_tmp = np.arange(-10,11)
# Use the function implemented above to get the sigmoid values
y = sigmoid(z_tmp)
#Input (z), Output (sigmoid(z))
#[[-1.000e+01  4.540e-05]
#[[-9.000e+00  1.234e-04]
#[[-8.000e+00  3.354e-04]
#[[-7.000e+00  9.111e-04].....

x_train = np.array([0., 1, 2, 3, 4, 5])
y_train = np.array([0, 0, 0, 1, 1, 1])

w_in = np.zeros((1))
b_in = 0
```

here we are giving a bunch of x having values 0 or 1 and we make a linear regression model and fit it in a sigmoid function to get a logistical regression model which is then plotted to give this curve where y is btw 0 and 1 and x is your normal x (input data x)



Decision boundary

- sigmoid fn is also called logistic fn

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

$$= P(y = 1 | \vec{x}; \vec{w}, b) \quad 0.7 \quad 0.3$$

0 or 1? threshold

$$\text{Is } f_{\vec{w},b}(\vec{x}) \geq 0.5?$$

$$\text{Yes: } \hat{y} = 1$$

$$\text{No: } \hat{y} = 0$$

$$\text{When is } f_{\vec{w},b}(\vec{x}) \geq 0.5?$$

$$g(z) \geq 0.5$$

$$z \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

$$\hat{y} = 1$$

$$\vec{w} \cdot \vec{x} + b < 0$$

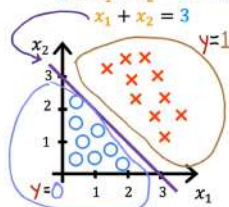
$$\hat{y} = 0$$

the decision boundary may not be a linear boundary always but also non linear/different shapes for other complex functions.

Decision boundary

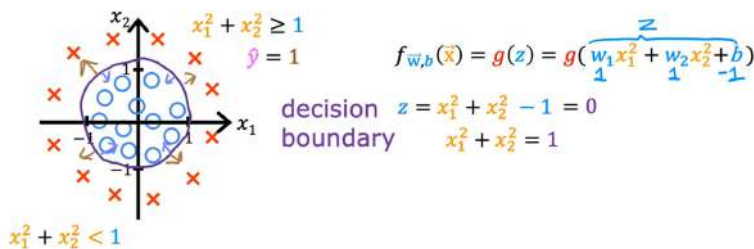
$$f_{\vec{w},b}(\vec{x}) = g(z) = g\left(\underbrace{w_1 x_1 + w_2 x_2 + b}_{z}\right)$$

$$\text{Decision boundary } z = \vec{w} \cdot \vec{x} + b = 0$$

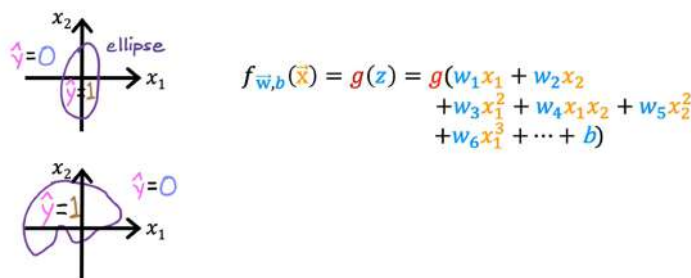


For a decision boundary line, $z=0$!!! the line for boundary is found at $z=0$

Non-linear decision boundaries



Non-linear decision boundaries



Cost function for logistic regression

- the loss function measures how well you're doing on one training example and is by summing up the losses on all of the training examples that you then get, the cost function, which measures how well you're doing on the entire training set

Training set

	tumor size (cm)	...	patient's age	malignant?	
	x_0		x_n	y	
$i=1$	10		52	1	$i = 1, \dots, m \leftarrow \text{training examples}$
\vdots	2		73	0	$j = 1, \dots, n \leftarrow \text{features}$
\vdots	5		55	0	
$i=m$	12		49	1	
	

target y is 0 or 1

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

here we are given n features with m training examples, the objective is to find the right set of $[w_1, w_2, w_3 \dots]$ and b using the cost function for logistic regression. (not squared error cost fn!!!)

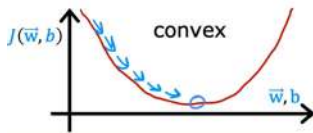
Squared error cost

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

loss $L(f_{\bar{w}, b}(\bar{x}^{(i)}), y^{(i)})$

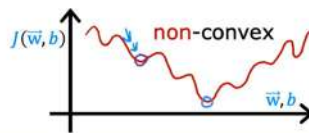
linear regression

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



logistic regression

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



here we see that we have modified the cost fn b taking 1/2 inside the summation and defined the loss L to be equal to the term inside \sum

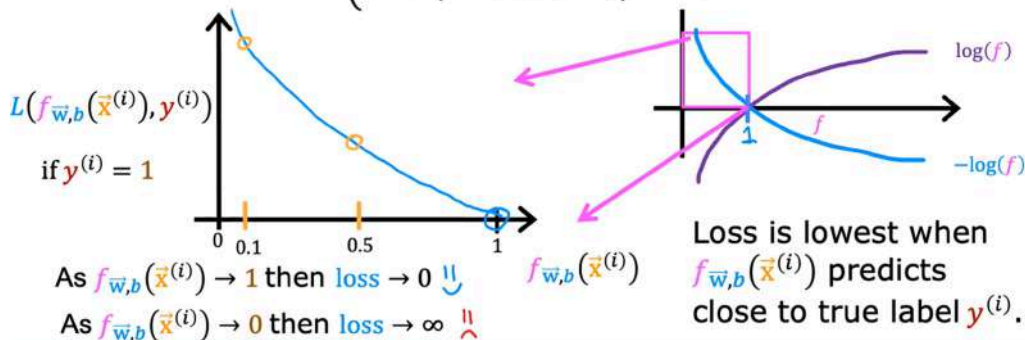
$$\text{loss } L(f_{\bar{w}, b}(\bar{x}^{(i)}), y^{(i)}) = \frac{1}{2} (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2$$

the loss fn takes to parameters, f(x) and output y (0 or 1)

also in the image, in linear regression the curve is convex when plotting cost fn against weights but in logistic regression, it is a non-convex with multiple local minimums hence its hard to do GD, thus we need to make it a good curve.

Logistic loss function

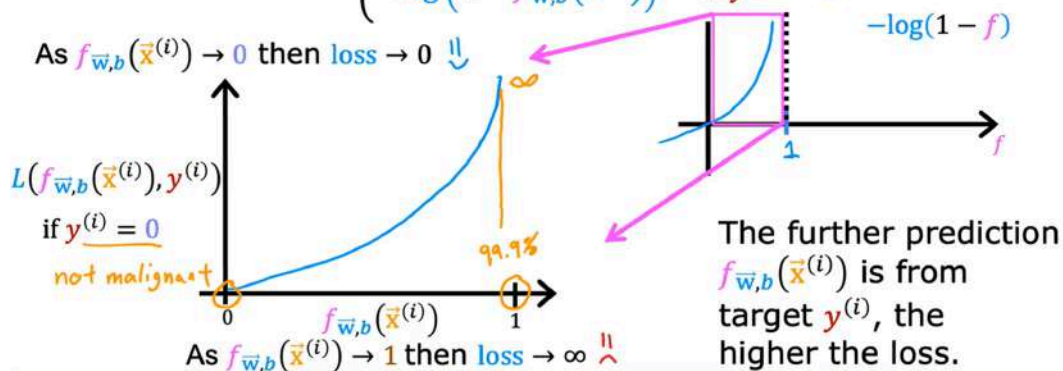
$$L(f_{\bar{w}, b}(\bar{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\bar{w}, b}(\bar{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\bar{w}, b}(\bar{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$



- here we see we have defined Loss in such a way that the curve is smooth. also the fn is different for $y=0$ and $y=1$.
- the x axis of the graph on the left is $f_{wb}(xi)$ (the prediction), this graph tells us **that given real output is 1 and the prediction should be 1** we are plotting loss against predicted output (between 0 and 1 obv)
- here for $y=1$, we see that as the model predicts output (y) closer to 1, the cost tends to 0. this means that lesser the cost, more the model pushes to adapt to the answer having loss tending to 0. when output tends to 0, loss tends to ∞ .
- also in the graph on the right, since we are dealing with $0 < y < 1$, we only consider the part/curve of the graph between 0 and 1. we penalize the model if the loss is very high! (discourage)

Logistic loss function

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



- the x axis of the graph on the left is $f_{wb}(xi)$ (the prediction), this graph tells us **that given real output is 0 and the prediction should be 0** we are plotting loss against predicted output (between 0 and 1 obv)
- here for $y=0$, we see that as the model predicts output (y) closer to 0, the cost tends to 0. this means that lesser the cost, more the model pushes to adapt to the answer having loss tending to 0. when output tends to 1, loss tends to ∞ .
- here above the loss is defined like this

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

the loss is defined for every training example (i) and the cost is defined as such

Cost

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

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

convex
↳ can reach a global minimum

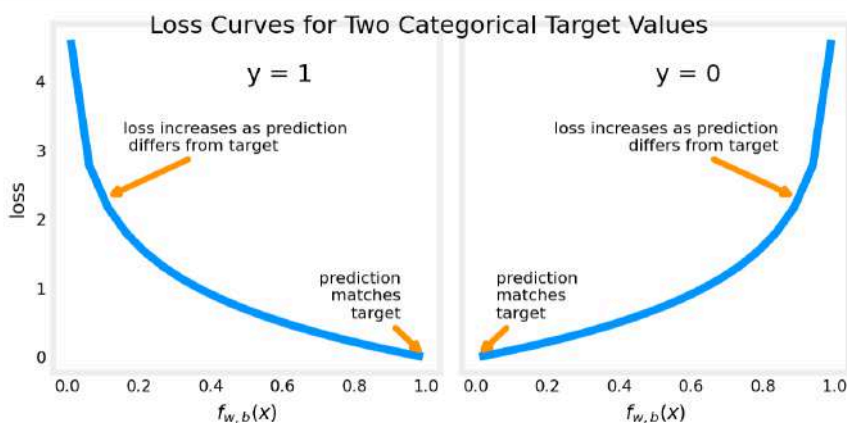
find w, b that minimize cost J

so we know the real output (0 OR 1), we know the $f_{w,b}(x(i))$ (the prediction), for every iteration we will choose the correct formula of loss depending on whether the real answer is 0 or 1 and then find the loss by putting $f_{w,b}(x(i))$ in the equation and finally sum all the losses and divide by m (no. of training examples) to get total cost. objective is now to reduce cost to finally get some set of optimized w, b that minimizes cost! remember the model is $z = w \cdot x + b$ and logistic regression is $1/(1+e^{-(w \cdot x + b)})$ and this is the actual predicted output!

$$z = w \cdot x + b$$

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

Loss is a measure of the difference of a single example to its target value while the **Cost** is a measure of the losses over the training set



Simplified Cost Function for Logistic Regression

Simplified loss function

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

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

if $y^{(i)} = 1$:

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

if $y^{(i)} = 0$:

$$L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)}) = - (1-0) \log(1-f(\vec{x}))$$

the loss function can be written in 1 single line like shown above and plugging in 0 or 1 will give you the appropriate eqn

Simplified cost function

$$\begin{aligned}
 \text{loss} \quad L(f_{\bar{w},b}(\vec{x}^{(i)}), y^{(i)}) &= -y^{(i)} \log(f_{\bar{w},b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\bar{w},b}(\vec{x}^{(i)})) \\
 \text{cost} \quad J(\bar{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L(f_{\bar{w},b}(\vec{x}^{(i)}), y^{(i)})] \quad \text{convex (single global minimum)} \\
 &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\bar{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\bar{w},b}(\vec{x}^{(i)}))] \\
 &\quad \text{maximum likelihood (don't worry about it!)}
 \end{aligned}$$

the reason we chose this particular cost fn is because of the statistical method of "maximum likelihood" and also that this gave us a convex curve

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

$$\text{loss} \quad L(f_{\bar{w},b}(\vec{x}^{(i)}), y^{(i)}) \approx \frac{1}{2} (f_{\bar{w},b}(\vec{x}^{(i)}) - y^{(i)})^2$$

$$\text{cost} \quad J(\bar{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\bar{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\bar{w},b}(\vec{x}^{(i)}))]$$

C1_W3_Lab05_Cost_Function_Soln

Cost function

In a previous lab, you developed the *logistic loss* function. Recall, loss is defined to apply to one example. Here you combine the losses to form the **cost**, which includes all the examples.

Recall that for logistic regression, the cost function is of the form

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=0}^{m-1} [\text{loss}(f_{\mathbf{w},b}(\mathbf{x}^{(i)}), y^{(i)})] \quad (1)$$

where

- $\text{loss}(f_{\mathbf{w},b}(\mathbf{x}^{(i)}), y^{(i)})$ is the cost for a single data point, which is:

$$\text{loss}(f_{\mathbf{w},b}(\mathbf{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\mathbf{w},b}(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\mathbf{w},b}(\mathbf{x}^{(i)})) \quad (2)$$

- where m is the number of training examples in the data set and:

$$f_{\mathbf{w},b}(\mathbf{x}^{(i)}) = g(z^{(i)}) \quad (3)$$

$$z^{(i)} = \mathbf{w} \cdot \mathbf{x}^{(i)} + b \quad (4)$$

$$g(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}} \quad (5)$$

```

import numpy as np
from lab_utils_common import plot_data, sigmoid, dlc

X_train = np.array([[0.5, 1.5], [1, 1], [1.5, 0.5], [3, 0.5], [2, 2], [1, 2.5]]) # (m,n)
y_train = np.array([0, 0, 0, 1, 1, 1]) # (m,)

def compute_cost_logistic(X, y, w, b):
    # here the linear model z = x.w + b
    m = X.shape[0]
    cost = 0.0
    for i in range(m):
        z_i = np.dot(X[i], w) + b
        f_wb_i = sigmoid(z_i) # returns the sigmoid of z, this is not a built in but imported, see top
        cost += -y[i]*np.log(f_wb_i) - (1-y[i])*np.log(1-f_wb_i)

    cost = cost / m
    return cost

w_tmp = np.array([1, 1])
b_tmp = -3
print(compute_cost_logistic(X_train, y_train, w_tmp, b_tmp))

```

Gradient Descent Implementation

in GD the algo remains same, simultaneous update of w and b and that same formula we used, shown below

aim: to find w (vector) and b to be able to predict output for a fresh new dataset.

Gradient descent

cost

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

repeat {
 $j = 1 \dots n$
 $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$
 $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$
 } simultaneous updates

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

remember we have j no. of w 's, thus while finding out derivative of cost wrt j , we will use $x(i, j)$

Gradient descent for logistic regression

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

looks like linear regression!

Same concepts:

- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

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

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

The formula of dj_dw and dj_db is the same as LR

C1_W3_Lab06_Gradient_Descent_Soln

```
def compute_gradient_logistic(X, y, w, b):
    m, n = X.shape
    dj_dw = np.zeros((n,))
    dj_db = 0.

    for i in range(m):
        f_wb_i = sigmoid(np.dot(X[i], w) + b)
        err_i = f_wb_i - y[i]
        for j in range(n):
            dj_dw[j] = dj_dw[j] + err_i * X[i, j]
            dj_db = dj_db + err_i
    dj_dw = dj_dw / m
    dj_db = dj_db / m

    return dj_db, dj_dw

def gradient_descent(X, y, w_in, b_in, alpha, num_iters, compute_gradient_logistic):
    # An array to store cost J and w's at each iteration primarily for graphing later
    w = copy.deepcopy(w_in) #avoid modifying global w within function
    b = b_in

    for i in range(num_iters):
        # Calculate the gradient and update the parameters
        dj_db, dj_dw = compute_gradient_logistic(X, y, w, b)

        # Update Parameters using w, b, alpha and gradient
        w = w - alpha * dj_dw
        b = b - alpha * dj_db

    return w, b #return final w, b
```

C1_W3_Lab07_Scikit_Learn_Soln

implement GD for logistic regression using sci-kit learn

```
import numpy as np
from sklearn.linear_model import LogisticRegression

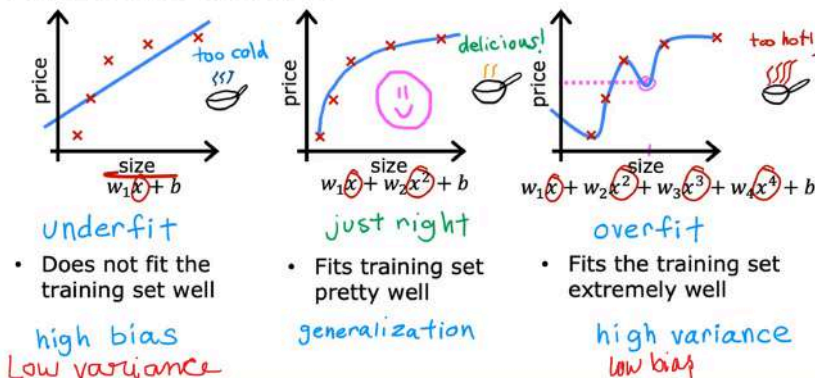
X = np.array([[0.5, 1.5], [1, 1], [1.5, 0.5], [3, 0.5], [2, 2], [1, 2.5]])
y = np.array([0, 0, 0, 1, 1, 1])

lr_model = LogisticRegression()
lr_model.fit(X, y)
y_pred = lr_model.predict(X)

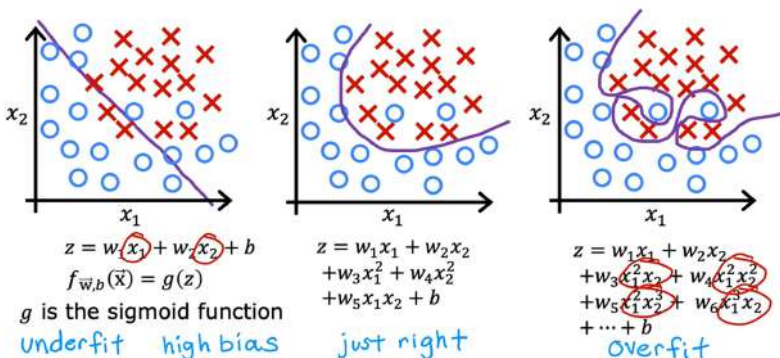
print("Prediction on training set:", y_pred)
print("Accuracy on training set:", lr_model.score(X, y))
```

The problem of overfitting

Regression example



Classification



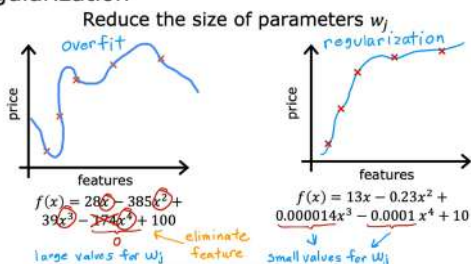
Addressing overfitting

method 1: collect more data, collecting more data would help the line fit data better and not be wiggly around data points ultimately not fitting new data properly

method 2: it may be possible that using all the features in the data set + having less data may be the reason causing overfitting. one solution to this maybe is to use a select few no. of those features and not all of them. the disadvantage to this is that maybe some imp. features be lost. this process is called feature selection.

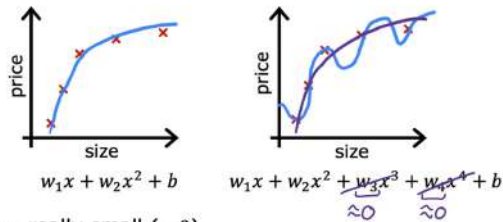
method 3: Regularization, sometimes we may not want to completely eliminate some features, thus we just reduce it (reduce the impact of that features) instead of completely removing it. this is called Regularization.

Regularization



in regularization we tend to regularize w 's only and not b . mostly not reqd. to regularize or make b small. in practice it would make very less diff. to make b small.

Intuition



make w_3, w_4 really small (≈ 0)

$$\min_{\bar{w}, b} \frac{1}{2m} \sum_{i=1}^m (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2 + \underbrace{1000 \omega_3}_{0.001} + \underbrace{1000 \omega_4}_{0.002}$$

here if we have to penalize the parameters, we can just add them to the cost fn and assign a large value of coeff to the, just their values would have to be very small to make cost fn small

Regularization

simpler model $w_3 \approx 0$
less likely to overfit $w_4 \approx 0$

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
x_1	x_2	x_3	x_4	x_5		x_{100}	y
$w_1, w_1, w_2, \dots, w_{100}, b$	n features					$n = 100$	

$$J(\bar{w}, b) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2 + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularization term}} + \underbrace{\frac{\lambda}{2m} b^2}_{\text{can include or exclude } b} \right]$$

"lambda" regularization parameter $\lambda > 0$

here λ is the **regularization parameter**, $2m$ is scaling the **regularization term**, w_j is every parameter, λ has to > 0 and $(\lambda/2m)b^2$ can be added or excluded.

Regularization

$$\min_{\bar{w}, b} J(\bar{w}, b) = \min_{\bar{w}, b} \left[\underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2}_{\text{mean squared error}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularization term}} \right]$$

fit data \rightarrow Keep w_j small

λ balances both goals

choose $\lambda = 10^{10}$

$f_{\bar{w}, b}(\bar{x}) = \underbrace{w_1x}_{\approx 0} + \underbrace{w_2x^2}_{\approx 0} + \underbrace{w_3x^3}_{\approx 0} + \underbrace{w_4x^4}_{\approx 0} + b$

$f(x) = b$ choose λ

we have to choose regularization parameter in such a way that its not enormous making the weight $\rightarrow 0$ or make it extremely small giving all the weights a chance to influence the model (outcome). thus an appropriate λ would make the model fit properly

Regularized linear regression

Regularized linear regression

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

Gradient descent

repeat {

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

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

} simultaneous update

don't have to regularize b

Implementing gradient descent

repeat {

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

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

} simultaneous update

some extra things

Implementing gradient descent

repeat {

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

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

} simultaneous update $j = 1 \dots n$

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

$$\alpha \frac{\lambda}{m} = 0.01 \frac{1}{50} = 0.0002$$

$$w_j (1 - 0.0002) = 0.9998$$

for every time w_j is updated, the w_j is reduced by a very small number which is $1 - \alpha \lambda/m$, thus the w_j is decreasing by a very small amount every iteration. hence we see the w_j is being updated and also the normal $-\alpha$ mse is there like in the LR GD. so this is basically the working/ how GD works in regularization

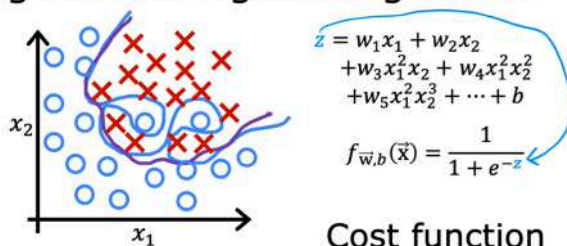
how the derivative term came for the regularization term

How we get the derivative term (optional)

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

Regularized logistic regression

Regularized logistic regression



Cost function

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

$\min_{\vec{w}, b} J(\vec{w}, b) \rightarrow w_j \downarrow$

Regularized logistic regression

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

$\min_{\vec{w}, b}$

Gradient descent

repeat {

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

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

}

Looks same as for linear regression!

logistic regression

you will see that the GD terms are exactly the same for Logistic regression as Linear regression but the only difference here is the meaning of $f_{\vec{w},b}(\vec{x})$ which here is a sigmoid fn unlike a linear fn in LinearR.

```

import numpy as np
import matplotlib.pyplot as plt

def compute_cost_linear_reg(X, y, w, b, lambda_ = 1):
# cost fn for REGULARISED LINEAR REGRESSION
    m = X.shape[0]
    n = len(w)
    cost = 0.
    for i in range(m):
        f_wb_i = np.dot(X[i], w) + b                #(n,)(n,)=scalar, see np.dot
        cost = cost + (f_wb_i - y[i])**2             #scalar
    cost = cost / (2 * m)                             #scalar

    reg_cost = 0
    for j in range(n):
        reg_cost += (w[j]**2)                        #scalar
    reg_cost = (lambda_/(2*m)) * reg_cost             #scalar

    total_cost = cost + reg_cost                      #scalar
    return total_cost

def compute_cost_logistic_reg(X, y, w, b, lambda_ = 1):
# cost fn for REGULARISED LOGISTIC REGRESSION
    m,n = X.shape
    cost = 0.
    for i in range(m):
        z_i = np.dot(X[i], w) + b                #(n,)(n,)=scalar, see np.dot
        f_wb_i = sigmoid(z_i)                     #scalar
        cost += -y[i]*np.log(f_wb_i) - (1-y[i])*np.log(1-f_wb_i) #scalar

    cost = cost/m                                    #scalar

    reg_cost = 0
    for j in range(n):
        reg_cost += (w[j]**2)                        #scalar
    reg_cost = (lambda_/(2*m)) * reg_cost             #scalar

    total_cost = cost + reg_cost                      #scalar
    return total_cost                                #scalar

def compute_gradient_linear_reg(X, y, w, b, lambda_):
# Gradient function for REGULARISED LINEAR REGRESSION
    m,n = X.shape                #(number of examples, number of features)
    dj_dw = np.zeros((n,))
    dj_db = 0.

    for i in range(m):
        err = (np.dot(X[i], w) + b) - y[i]
        for j in range(n):
            dj_dw[j] = dj_dw[j] + err * X[i, j]
        dj_db = dj_db + err
    dj_dw = dj_dw / m
    dj_db = dj_db / m

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

    return dj_db, dj_dw

def compute_gradient_logistic_reg(X, y, w, b, lambda_):
# Gradient function for REGULARISED LOGISTIC REGRESSION
    m,n = X.shape
    dj_dw = np.zeros((n,))        #(n,)
    dj_db = 0.0                   #scalar

    for i in range(m):
        f_wb_i = sigmoid(np.dot(X[i],w) + b)    #(n,)(n,)=scalar
        err_i = f_wb_i - y[i]                  #scalar
        for j in range(n):
            dj_dw[j] = dj_dw[j] + err_i * X[i,j] #scalar
        dj_db = dj_db + err_i

    dj_dw = dj_dw/m                #(n,)
    dj_db = dj_db/m                #scalar

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

```

```
return dj_db, dj_dw
```

do read the final lab assignment of week 3, very imp and useful