# Logistic Regression

Md Arafat Hussain

Chennupalli Sathwik

Ramesh Chandra Soren

Aryan D Rozario

Avishikta Das

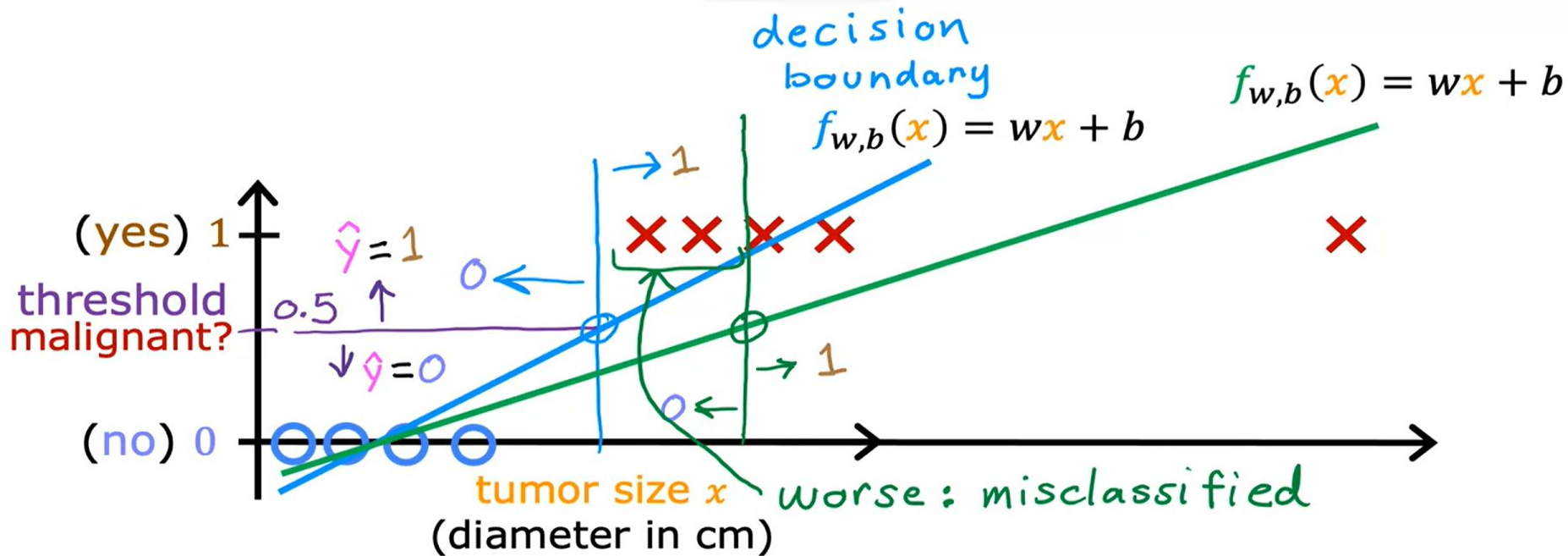Sanjoy Garai

# *Objective*

- Overview of Logistic Regression
- Sigmoid Function
- Squared error cost
- Gradient descent
- Overfitting Issue
- Regularization

# Overview

# Problem with linear regression



if $f_{w,b}(x) < 0.5 \rightarrow \hat{y} = 0$

if $f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$

- Linear regression is used when the outcome (dependent variable) is continuous. It predicts a value based on a linear relationship between the independent variables and the outcome.
- It cannot accurately binarily classify an entity according to categorical outcome variables.
- Logistic regression is specifically designed for binary classification problems, which are prevalent in various fields such as healthcare (disease diagnosis), finance (fraud detection), and marketing (customer segmentation). Linear regression cannot handle binary outcomes directly.

# Difference between linear and regression model

- **Linear Regression**:
  - Predicts continuous numerical values.
  - Models the relationship between dependent and independent variables using a linear equation.
  - Used for regression tasks.
  - Solved by minimizing the least squares error.

- **Logistic Regression**:
  - Predicts categorical outcomes.
  - Models the probability of a binary outcome using a logistic function.
  - Used for classification tasks.
  - Uses the logistic loss function to penalize large errors.
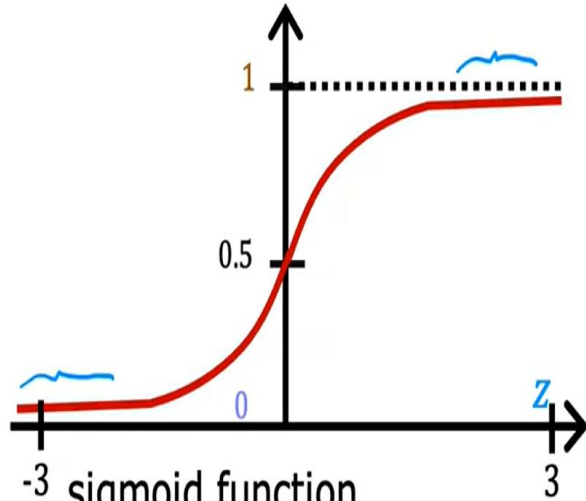
# Sigmoid function

# Example of Logistic regression

A Linear function cannot be used to classify accurately the given data.

We need a certain sigmoid function to vary between 0 and 1 for any features of input x to better classify the input according to a threshold value.

Want outputs between 0 and 1



sigmoid function

logistic function

outputs between 0 and 1

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

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

$$z = \vec{w} \cdot \vec{x} + b$$

$$z$$

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

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

"logistic regression"

w and b are the parameters.

w1,w2,w3.......wn; for x1,x2,x3......xn attributes.

b is the bias term.

# Squared Error Cost

# Squared Error Cost

- **Squared Error**:
  - Measures the difference between predicted and actual values in regression.
  - Calculated by squaring the difference between predicted and actual values, then summing across all observations.
  - Possesses nice mathematical properties, such as being everywhere differentiable.
  - Induced by an inner product on the underlying space, facilitating geometric interpretation.
- **Choice of Loss Function**:
  - Depends on the nature of the data and the problem being solved. This calculates the loss in each attribute of the given data.
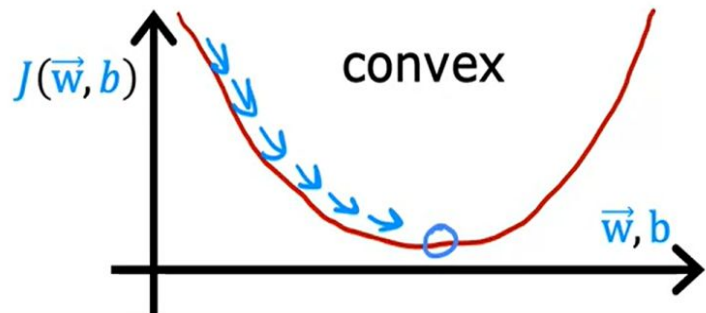
# Squared error cost

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

loss $\boxed{L\left( f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)} \right)}$

## linear regression

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

convex

$J(\vec{w}, b)$

$\vec{w}, b$

## logistic regression

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

non-convex

$J(\vec{w}, b)$

$\vec{w}, b$

# Cost

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

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

Convex
↳ can reach a
global minimum

# Simplified cost function

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

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

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

# Gradient Descent

# Necessity of Gradient Descent

- **Parameter Reduction**:
  - Gradient descent iteratively adjusts model parameters (coefficients and bias) to minimize the logistic loss function.
  - At each step, gradients of the loss function are computed to determine the direction of parameter updates.
  - Parameters are updated in the opposite direction of the gradients, scaled by the learning rate.
  - This process repeats until convergence, refining parameters to best fit the data and optimize classification accuracy.

  **Learning Rate(alpha)**: As in other variants of gradient descent, the learning rate determines the step size of each update. It's a hyperparameter that needs to be carefully chosen to balance convergence speed and stability.

# Gradient descent

(Alpha) is the learning rate.

cost

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

repeat {

$j = 1 \ldots n$

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

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

} simultaneous updates

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

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

# 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

Same concepts:
- Monitor gradient descent (learning curve)
- Vectorized implementation

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

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

# Overfitting Issue

- **Overfitting**:
  - **Definition**: Model memorizes noise in training data rather than learning underlying patterns, leading to poor performance on unseen data.
  - **Causes**: Model complexity, insufficient data, noise in data.
  - **Effects**: Reduced generalization, poor performance on unseen data, unreliable insights.
  - **Prevention and Remedies**: Simplify model, cross-validation, regularization, feature selection, ensemble methods, early stopping.
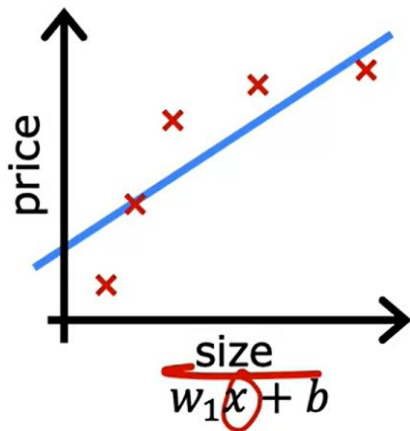
Overfitting can be addressed by ensuring models are appropriately complex for the data available, validating performance on unseen data, and employing techniques to reduce noise and complexity.
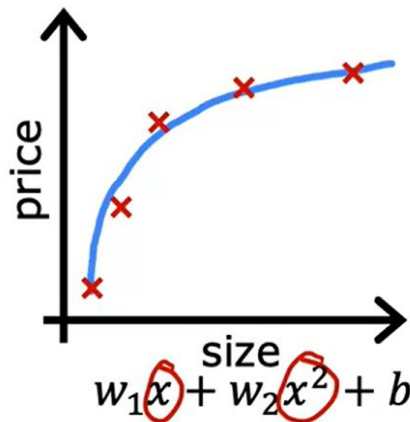
For linear regression

# Regression example



Press Esc to exit full screen

**Underfit** (price vs size)
$$w_1 x + b$$
underfit

- Does not fit the training set well

high bias

**Just right** (price vs size)
$$w_1 x + w_2 x^2 + b$$
just right

- Fits training set pretty well

generalization

**Overfit** (price vs size)
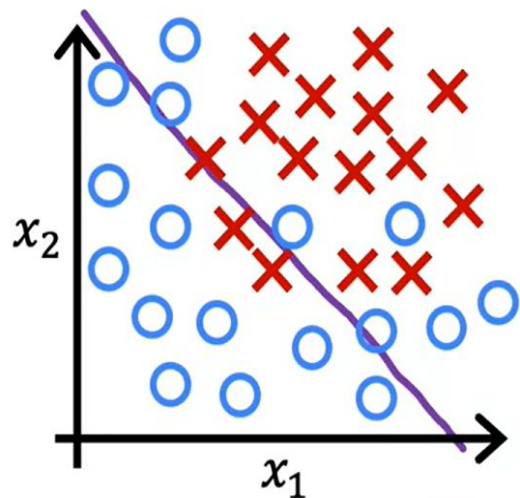$$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$
overfit

- Fits the training set extremely well
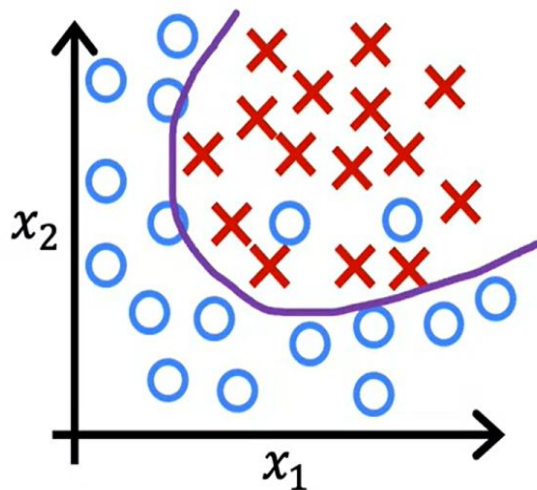
high variance

# Classification



$z = w_1 x_1 + w_2 x_2 + b$

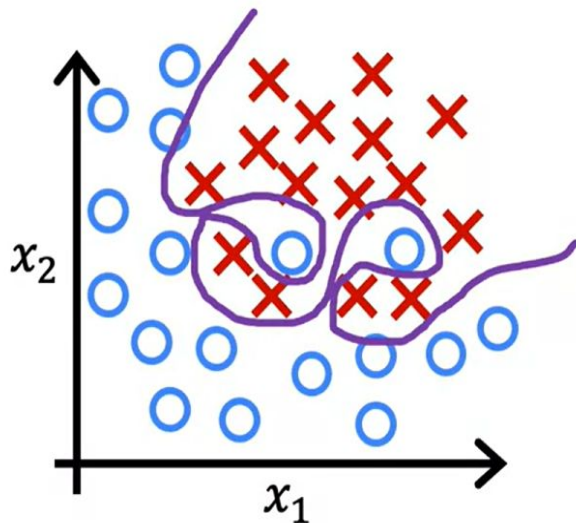$f_{\vec{w},b}(\vec{x}) = g(z)$

$g$ is the sigmoid function

underfit    high bias

$z = w_1 x_1 + w_2 x_2$
$+ w_3 x_1^2 + w_4 x_2^2$
$+ w_5 x_1 x_2 + b$

just right

$z = w_1 x_1 + w_2 x_2$
$+ w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2$
$+ w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2$
$+ \cdots + b$

overfit

# Addressing overfitting

Options

1. Collect more data

2. Select features
   - Feature selection

3. Reduce size of parameters
   - "Regularization"

Regularization

# Need for Regularization

- Regularization is a technique to prevent overfitting by adding a penalty term to the model's loss function, discouraging overly complex models. It comes in two main forms:
- Regularization helps in improving model generalization by penalizing complexity, although it may introduce some bias. The choice of regularization strength is crucial and typically requires tuning. Regularization is widely used in various models like linear regression, logistic regression, and neural networks.

# Regularization

small values $w_1, w_2, \cdots, w_n, b$

simpler model

less likely to overfit

$w_3 \hat{\approx} 0$

$w_4 \approx 0$

| size $x_1$ | bedrooms $x_2$ | floors $x_3$ | age $x_4$ | avg income $x_5$ | ... | distance to coffee shop $x_{100}$ | price $y$ |
|---|---|---|---|---|---|---|---|
| | | | | | | | |

$n = 100$

$w_1, w_1, w_2, \cdots, w_{100}, b$

$n$ features

regularization term

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

can include or exclude $b$

"lambda"
regularization parameter   $\lambda > 0$

# Regularization

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

fit data

← Keep wj small

$\lambda$ balances both goals

choose $\lambda = 10^{10}$

$$f_{\vec{w},b}(\vec{x}) = \underset{\approx 0}{w_1 x} + \underset{\approx 0}{w_2 x^2} + \underset{\approx 0}{w_3 x^3} + \underset{\approx 0}{w_4 x^4} + b$$

$f(x) = b$

choose $\lambda$

price

$\lambda = 0$

b

# Thank You