# Logistic Regression Part 2

**Dr Mehrdad Ghaziasgar**

Senior Lecturer in Computer Science
Department of Computer Science
University of the Western Cape
mghaziasgar@uwc.ac.za

# Content

- Part 1:
    - Classification With Logistic Regression
    - Model Evaluation For Classification

- Part 2:
    - Overfitting and Underfitting

- Part 3:
    - Practical Issues With Classification

# Content – Part 2

- Overfitting and Underfitting
    - Introduction
    - The Concept of Regularization
    - Regularized Linear Regression
    - Regularized Logistic Regression
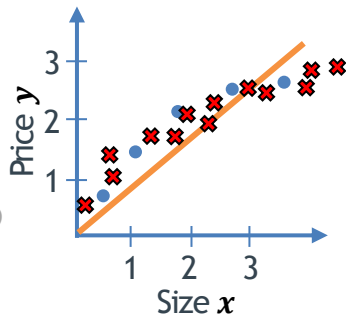
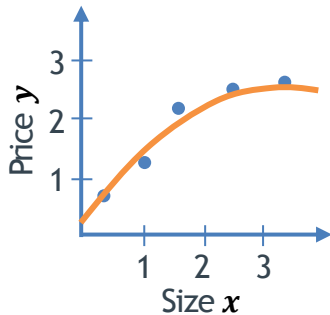# Overfitting and Underfitting
## Introduction

# Introduction

- Underfitting and overfitting are problems that can arise that cause a model (either regression or classification) to perform very poorly on test (new) examples
- We'll describe these problems and then describe how we can mitigate them using a technique called "Regularization"

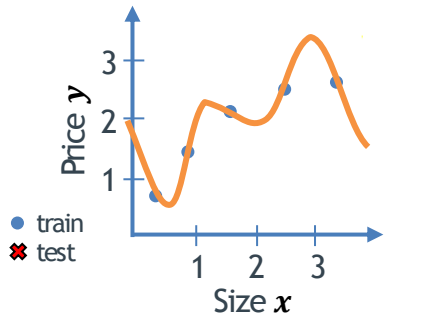# Introduction – Overfitting and Underfitting in Linear Regression



$$h_\theta(x) = \theta_0 + \theta_1 x$$

Underfitting
"High Bias"

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfitting
"High Variance"

- Underfitting: The model $h_\theta(x)$ is too simple (too few features) and doesn't fit the training examples very well or even at all (so $J(\theta)$ is very large) and so the model will also not reflect reality and doesn't fit new (test) examples well either
- The model has high "bias" i.e. it has a biased (rigid, inflexible) nature

# Introduction – Overfitting and Underfitting in Linear Regression



$h_\theta(x) = \theta_0 + \theta_1 x$

Underfitting
"High Bias"

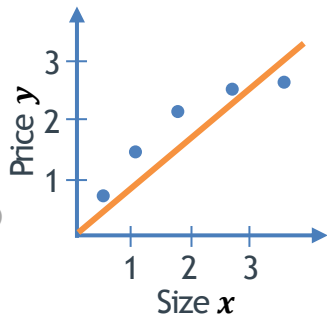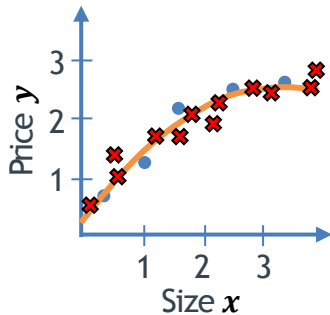$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

"Just right"

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

Overfitting
"High Variance"

- Overfitting: The model $h_\theta(x)$ is too complex (too many features) and fits the training examples very well or even perfectly (so $J(\theta) \approx 0$) but the model doesn't reflect reality and doesn't fit the test examples very well (or at all)
- The model has "high variance" i.e. it has a "too flexible" / "too variant" nature

# Introduction – Overfitting and Underfitting in Logistic Regression

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$
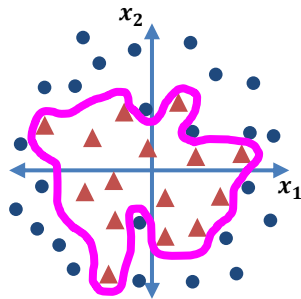
$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \cdots)$$

Underfitting
"High Bias"

"Just right"

Overfitting
"High Variance"

# Introduction – Addressing High Bias and Variance

- To address high bias:
    - Increase the number of features i.e. new types of features, high-order features, combination of features

- To address high variance, there are two main methods:
    1. Reduce the Features
        - Manually (pain-stakingly) select and remove features
        - Use train-cv-test sets to compare various models and select the best one
        - Works well when you have relatively fewer features and/or you have specific ideas about models to use/compare

    2. Apply Regularization
        - Use as complex of a model as you like (number of features, high-order terms, combination of features)
        - Apply regularization to reduce the impact of these features towards predicting $y$ by adjusting/reducing all of the $\theta$ parameters
        - Works well when you have a LOT of features and manual selection will be hard

# Overfitting and Underfitting
## The Concept of Regularization

# Regularization - Concept

- If we minimize $J(\theta)$ with this hypothesis, we'll get a $h_\theta(x)$ that overfits
- Suppose that we penalize $\theta_3$ and $\theta_4$ to make them really small by modifying the cost function as follows:

$$\underset{\theta}{\text{minimize}}\ \frac{1}{2m}\sum_{i=1}^{m}\big(h_\theta(x^{(i)}) - y^{(i)}\big)^2 + 100\ \theta_3^2 + 100\ \theta_4^2$$



$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$
$$+ \theta_3 x^3 + \theta_4 x^4$$

- With everything else the same, the only way to reduce (and minimize) the cost in the same way as before is to ensure that $\theta_3$ and $\theta_4$ become really really small
  - Has the effect of setting $\theta_3$ and $\theta_4$ such that $+ 100\ \theta_3^2 + 100\ \theta_4^2$ is set to 0 to get back the original cost i.e. setting $\theta_3$ and $\theta_4$ close to 0
  - Also the effect of almost excluding the third and fourth order terms in $h_\theta(x)$

# Regularization - Intuition

- Intuition 2: Hamburger shop:
    - As the manager of a hamburger joint, you're selling 10 200g hamburgers per day @ R10 per hamburger
    - The "cost" of a burger is:

    $$\text{Cost} = P_{Meat} + P_{Veggies} + P_{Roll}$$
    $$= 3.50 + 1.20 + 1.50 = R6.20$$

    - At this cost, the boss is making a handsome R2000 per day in profits
    - Life is good

# Regularization - Intuition

- Intuition 2: Hamburger shop:
  - Suddenly, the centre announces that they are supporting "animal rights" or some weird thing like that
  - They are imposing an extra "cost" on every kg of meat you sell
  - They will charge you 1cent for every 1g of meat you sell (basically the weight of the meat)
  - The "cost" of a burger now changes:

$$Cost = P_{Meat} + P_{Veggies} + P_{Roll}$$
$$= \quad 3.50 + \quad 1.20 + 1.50$$

  - Unfortunately, the boss says he's not willing to accept that much less profit and orders you to reduce the new cost down to R6.50 **somehow.**
  - You've got to reduce your cost somehow to keep profits constant

  - The prices at which you're buying meat, veggies and rolls are also still the same
  - What can you change to keep the cost the same as it was before?
    - Reducing the "weight" of the burger

# Regularization - Intuition

- Intuition 2: Hamburger shop:
  - Here's what needs to happen to get the new cost down from R8.20 to R6.50:

$$\text{Cost} = P_{Meat} + P_{Veggies} + P_{Roll} + \boldsymbol{\theta}_{Meat} \times 0.01$$

$$6.50 = 3.50 + 1.20 + 1.50 + \boldsymbol{\theta}_{Meat} \times 0.01$$

$$\boldsymbol{\theta}_{Meat} = \frac{6.50 - (3.50 + 1.20 + 1.50)}{0.01}$$

$$\boldsymbol{\theta}_{Meat} = \frac{6.50 - (3.50 + 1.20 + 1.50)}{0.01}$$

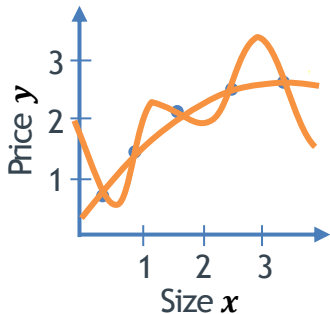$$\boldsymbol{\theta}_{Meat} = 30\text{g of meat}$$

# Regularization - Concept

- If we minimize $J(\theta)$ with this hypothesis, we'll get a $h_\theta(x)$ that overfits
- Suppose that we penalize $\theta_3$ and $\theta_4$ to make them really small by modifying the cost function as follows:

$$\underset{\theta}{\text{minimize}} \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2 + 100\, \theta_3^2 + 100\, \theta_4^2$$



$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$
$$+ \theta_3 x^3 + \theta_4 x^4$$
$$\quad\; \underset{0}{\approx} \qquad\; \underset{0}{\approx}$$

- With everything else the same, the only way to reduce (and minimize) the cost in the same way as before is to ensure that $\theta_3$ and $\theta_4$ become really really small
    - Has the effect of setting $\theta_3$ and $\theta_4$ such that $+ 100\, \theta_3^2 + 100\, \theta_4^2$ is set to 0 to get back the original cost i.e. setting $\theta_3$ and $\theta_4$ close to 0
    - Also the effect of almost excluding the third and fourth order terms in $h_\theta(x)$

# Regularization - Concept

- In the example we saw: we knew which terms were higher-order i.e. $\theta_3$ and $\theta_4$
- Adding those specific terms to the cost function reduced their weights
- In practice: we aren't always sure which features we need to regularize
  - We regularize ALL of the weights (except $\theta_0$) equally
  - We add all of them to the cost function

$$J_{\text{reg}}(\theta) = J(\theta) + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

- $J(\theta)$ is the normal cost function that we had before (either for linear or logistic regression)
- We don't regularize $\theta_0$ so $j$ starts from 1, not 0.

- The parameter $\lambda$ is the regularization parameter which helps control how much to regularize:
  - Setting $\lambda = 0$ means that the whole regularization term falls away i.e. no regularization at all – potential **overfitting** – **high variance**
  - Setting $\lambda > 0$ means increasing amount of regularization by penalizing all $\theta_j$ (except $\theta_0$) at the same time – moving towards **underfitting** – **high bias**

# Overfitting and Underfitting
## Regularized Linear Regression

# Regularized Linear Regression

- What we arrived at previously:

$$J_{\mathbf{reg}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

- So for linear regression:

$$J_{\mathbf{reg}}(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

$$J_{\mathbf{reg}}(\boldsymbol{\theta}) = \frac{1}{2m} \left[ \sum_{i=1}^{m} \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

# Regularized Linear Regression

- With the new cost function, we can simply apply gradient descent (or any other optimization/minimization technique) as before
- This involves making continuous updates to all $\boldsymbol{\theta_j}$ to minimize the cost

$$J(\boldsymbol{\theta}) = \frac{1}{2m}\left[\sum_{i=1}^{m}\left(h_{\boldsymbol{\theta}}\left(x^{(i)}\right) - y^{(i)}\right)^2 + \lambda\sum_{j=1}^{n}\theta_j^2\right]$$

**?**

- We'll be updating $\boldsymbol{\theta_j}$ using the same update rules as before

- One thing has changed here: $\frac{\partial}{\partial\theta_j}J(\boldsymbol{\theta})$ because the $J(\boldsymbol{\theta})$ has changed

Repeat until convergence:

Update all $\boldsymbol{\theta_j}$ simultaneously:

$$\theta_j \leftarrow \theta_j - \alpha\frac{\partial}{\partial\theta_j}J(\boldsymbol{\theta})$$

# Regularized Linear Regression

- For any parameter $\theta_j$:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \left[ \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

Repeat until convergence:

Update all $\theta$ simultaneously:

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( \theta^T x^{(i)} - y^{(i)} \right) \cdot x_0^{(i)}$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} \left( \theta^T x^{(i)} - y^{(i)} \right) \cdot x_1^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

# Regularized Linear Regression

- The update rule for $\boldsymbol{\theta_0}$ is exactly the same as before: we're not regularizing it

- The update rules for all other $\boldsymbol{\theta_j}$ are different; they now have an added term

$$\theta_j \leftarrow \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (\boldsymbol{\theta}^T x^{(i)} - y^{(i)}) \cdot x_1^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

- If we slightly restructure this update rule, something interesting emerges:

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (\boldsymbol{\theta}^T x^{(i)} - y^{(i)}) \cdot x_1^{(i)} - \alpha \frac{\lambda}{m} \theta_j$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{\lambda}{m} \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (\boldsymbol{\theta}^T x^{(i)} - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_j \leftarrow \left(1 - \alpha \frac{\lambda}{m}\right) \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (\boldsymbol{\theta}^T x^{(i)} - y^{(i)}) \cdot x_1^{(i)}$$

- Given that $m$ will usually be larger than $\alpha\lambda$, the term $0 < \alpha \frac{\lambda}{m} < 1$ e.g 0.01

- Therefore, also $0 < \left(1 - \alpha \frac{\lambda}{m}\right) < 1$ e.g. 0.99

- Multiplying this by $\boldsymbol{\theta_j}$ has the effect of shrinking by e.g. 1% it before the update

# Regularized Linear Regression – Using sklearn

- Now that we know the details, we can use sklearn to fit a regularized model as follows:

  from sklearn.linear_model import Ridge

  model = Ridge(alpha= 10 )          $\lambda = 10$

  model.fit(X,y)

- Note that X should **not** contain the extra column of 1s for feature $x_0$.

- The Ridge class has other parameters: find out what they are/do

# Normal Equation With Regularization

- Given a feature matrix $X$, and the corresponding output matrix $y$
- The following equation solves for $\theta$ that best fits the data with regularization:

$$\theta = \left( X^T X \qquad\qquad\qquad \right)^{-1} X^T y$$

- Where
  - the matrix added is an $(n + 1) \times (n + 1)$ and
  - has a diagonal of 1s (apart from the top-left entry) and 0s everywhere else
- E.g.

$$X = \begin{bmatrix} 1 & 460 & 4 & 12 & 2 \\ 1 & 70 & 1 & 5 & 0 \\ 1 & 155 & 3 & 8 & 2 \\ 1 & 429 & 6 & 10 & 3 \end{bmatrix} \qquad y = \begin{bmatrix} 6639 \\ 1681 \\ 3969 \\ 5095 \end{bmatrix}$$

Size   #Rms   Age   #Grgs     Price

# Overfitting and Underfitting
## Regularized Logistic Regression

# Regularized Logistic Regression

- What we arrived at previously:

$$J_{\text{reg}}(\theta) = J(\theta) + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

- So for linear regression:

$$J_{\text{reg}}J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \ln\left(h_\theta(x^{(i)})\right) + (1 - y^{(i)}) \ln\left(1 - h_\theta(x^{(i)})\right) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

- Where:

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Regularized Logistic Regression

- With the new cost function, we can simply apply gradient descent (or any other optimization/minimization technique) as before
- This involves making continuous updates to all $\theta_j$ to minimize the cost

$$J(\theta) = -\frac{1}{m}\sum_{i=1}^{m}\left[y^{(i)}\ln(h_\theta(x^{(i)})) + (1-y^{(i)})\ln\left(1-h_\theta(x^{(i)})\right)\right] + \frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

**?**

- We'll be updating $\theta_j$ using the same update rules as before

- One thing has changed here: $\frac{\partial}{\partial\theta_j}J(\theta)$ because the $J(\theta)$ has changed

Repeat until convergence:

Update all $\theta_j$ simultaneously:

$$\theta_j \leftarrow \theta_j - \alpha\boxed{\frac{\partial}{\partial\theta_j}J(\theta)}$$

# Regularized Logistic Regression

- For any parameter $\theta_j$:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} + \frac{\lambda}{m} \theta_j$$

Repeat until convergence:

Update all $\theta$ simultaneously:

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x_0^{(i)}$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x_1^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

- Algorithm appears to be identical to gradient descent for regularized linear regression
  - It is NOT! The hypothesis here is that of logistic regression

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Regularized Logistic Regression Using Adv. Optim. Algorithms

- Implemented in SciPy's optimize library:

```
from scipy.optimize  import minimize

theta = np.zeros(X.shape[1])  #Initialize all thetas to zeros
result = minimize(costFunc,  theta,  args=(X,y),  method='BFGS',  jac=gradientFunc,
options={'maxiter'  : 400, 'disp': True})

print(result.x)
```

- costFunc is a function that returns the cost as per:

$$J(\boldsymbol{\theta}) = -\frac{1}{m}\sum_{i=1}^{m}\left[\boldsymbol{y}^{(i)}\ln(\boldsymbol{h_\theta}(\boldsymbol{x}^{(i)})) + (1 - \boldsymbol{y}^{(i)})\ln\left(1 - \boldsymbol{h_\theta}(\boldsymbol{x}^{(i)})\right)\right] + \frac{\lambda}{2m}\sum_{j=1}^{n}\boldsymbol{\theta}_j^2$$

- gradientFunc is a function that returns a list/vector of the gradients of all the $\boldsymbol{\theta}_j$ as per:

$$\frac{\partial}{\partial\boldsymbol{\theta}_j}J(\boldsymbol{\theta}) = \frac{1}{m}\sum_{i=1}^{m}(\boldsymbol{h_\theta}(\boldsymbol{x}^{(i)}) - \boldsymbol{y}^{(i)})\,\boldsymbol{x}_j^{(i)} + \frac{\lambda}{m}\boldsymbol{\theta}_j$$

# Regularized Logistic Regression Using Adv. Optim. Algorithms

- Implemented in SciPy's optimize library:

```
from scipy.optimize import minimize

theta = np.zeros(X.shape[1]) #Initialize all thetas to zeros
result = minimize(costFunc, theta, args=(X,y), method='BFGS', jac=gradientFunc,
options={'maxiter' : 400, 'disp': True})

print(result.x)
```

- result.x will contain the optimal $\boldsymbol{\theta}$ values
- Note that X must contain an extra column of zeros representing feature $x_0$

- Predictions can then be made by applying $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g(\boldsymbol{\theta}^T \boldsymbol{x})$

# Regularized Logistic Regression Using sklearn

- Implemented in sklearn specifically for Logistic Regression:

  from sklearn.linear_model import LogisticRegression

  clf = LogisticRegression(random_state=0, solver='lbfgs', multi_class='auto', C = 10 )

  clf.fit(X, y)

$C = \frac{1}{\lambda} = 10$

So here:

$\lambda = \frac{1}{C} = 0.1$

- Regularization in the LogisticRegression class is specified differently: it is specified via a parameter $C$ which is the inverse of $\lambda$.

$$C = \frac{1}{\lambda} \quad \text{so} \quad \lambda = \frac{1}{C}$$

- $C$ has the same goal as, but opposite effect to, $\lambda$ i.e.:
  - Setting $C$ very small means setting $\lambda$ very large i.e. **more** regularization
  - Setting $C$ very large means setting $\lambda$ very small i.e. **less** regularization

**THE END**
Of Logistic Regression Part 2