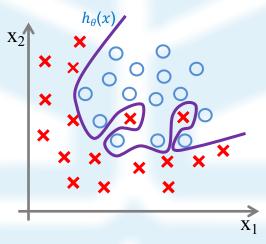
### **Machine Learning**

# Regularization

The problem of overfitting



Dr. Ali Valinejad

valinejad.ir valinejad@umz.ac.ir



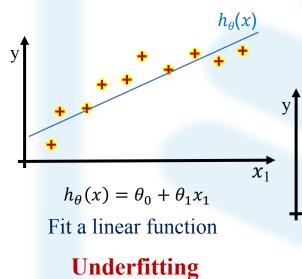
University of Mazandaran

#### Outlines:

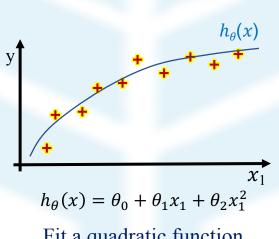
- Which model is a good model?
- Definition of a good model
- Underfitting?
- Overfitting
- Regularization
- Regularized linear regression
  - Normal equation
  - oGradient Decent
- Regularized logistic regression



### Which model is a good model?

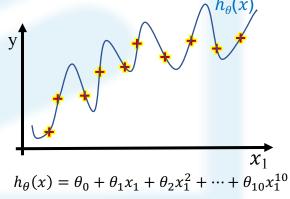


(High Bias)



Fit a quadratic function

**Work Well** 



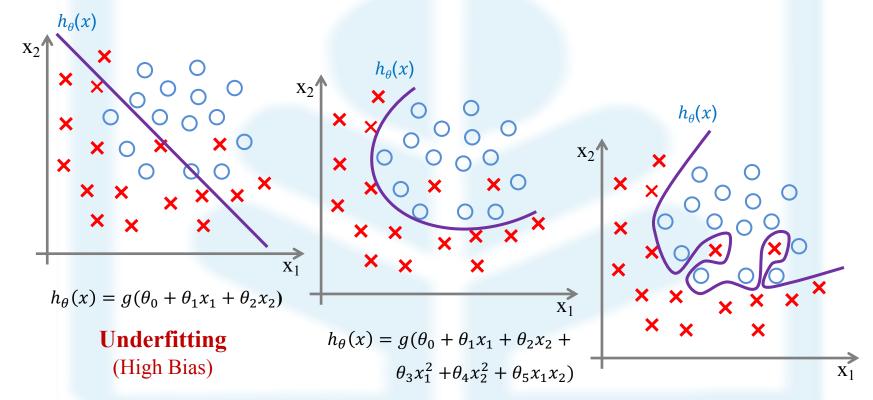
Fit a 10th order polynomial

#### **Overfitting**

(High variance)



### Which model is a good model?



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

**Work Well** 

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

#### **Overfitting**

(High variance)



### Definition of a good model

The definition of a 'good' model with *least training error*:

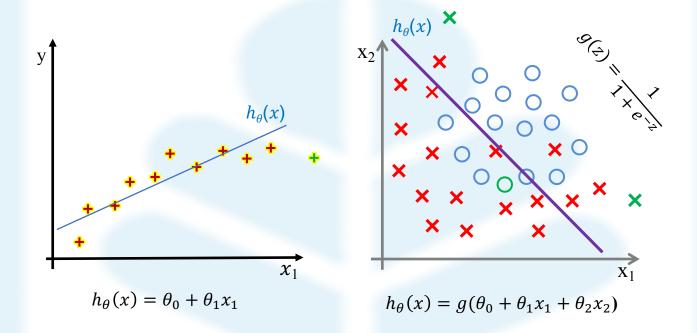
- **Step 1.** define a hypothesis function  $h_{\theta}(x)$  (i.e. model).
- **Step 2.** define a Cost Function  $J(\theta)$  for error measuring.
- **Step 3.** learn parameters of  $h_{\theta}(x)$  by minimizing  $J(\theta)$ .



# **Underfitting**



### **Underfitting**

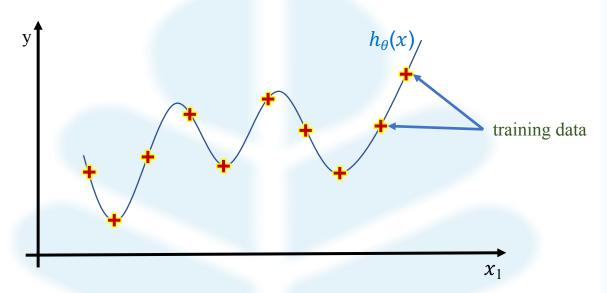


- This model has huge training error,
- ❖ When new data (green point) comes in, this model could has huge *prediction* error.





# Is this is a good model?

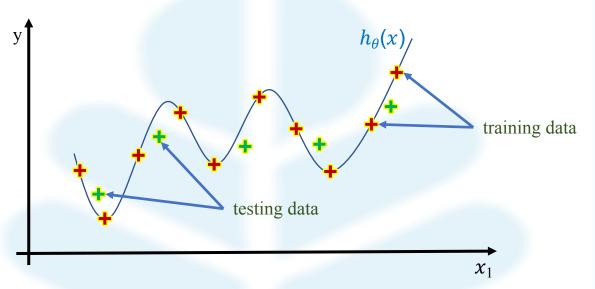


- This model has zero training error,
- This model fits training data perfectly,
- ❖ The predict label of all training data is equal to their corresponding truth label.

Is this is a good model?



#### How come the performance of a good model is terrible?

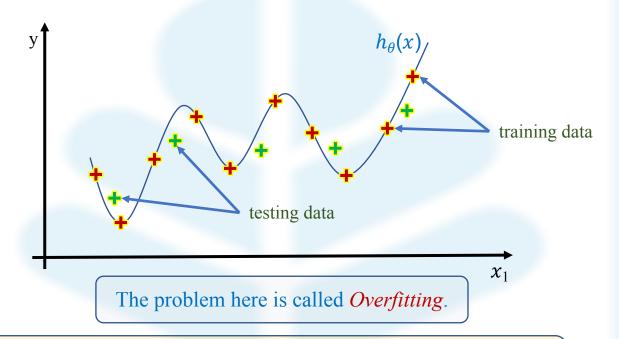


- The model has little *training error*.
- ❖ When new data (green point) comes in, this model has huge *prediction* error.

The learned model works well for training data but *terrible* for testing data (unknown data).

A truly good model must have both little training error and little prediction error.



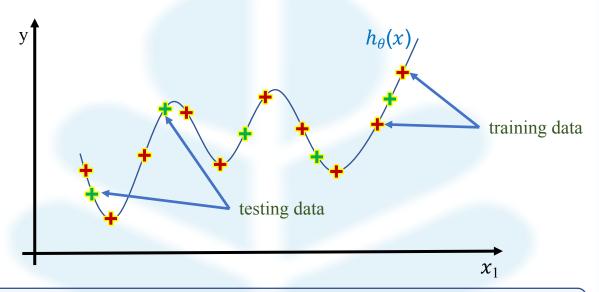


Overfitting: If we have *too many features*, the learned hypothesis may fit the training set very well, but fail to generalize to new examples.

#### Algorithm has high variance:

- Variance, in the context of Machine Learning, is a type of error that occurs due to a model's sensitivity to small fluctuations in the training set.
- Thise means that hypothesis can basically fit any training data.
- High variance would cause an algorithm to model the noise in the training set.





If we have a whole dataset that captured *all possibilities* of the problem, we don't need to worry about overfitting.

When new data comes in, it shall fall into one possibility, hence, the model can predict it perfectly.

A truly good model must have both little training error and little prediction error.



#### **Addressing overfitting:**

- \* Reduce number of features
  - o Manually select which features to keep.
  - o Model selection algorithm.
- \* Regularization
  - $\circ$  Keep all the features, but reduce magnitude/values of parameters  $\theta_i$ .
  - Works well when we have a lot of features, each of which contributes a bit to predicting *y*.



# Regularization



## Regularization

When *overfitting* occurs, we get an over complex model with too many features. One way to avoid it is to apply *Regularization* and then we can get a better model with proper features.

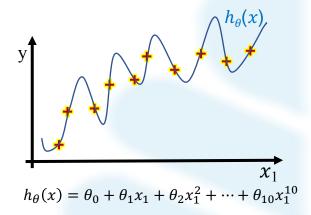
**Regularization** is a technique applied to Cost Function  $J(\theta)$  in order to avoid Overfitting.

The *core idea* in *Regularization* is to keep more important features and ignore unimportant ones. The importance of feature is measured by the value of its parameter  $\theta_j$ .



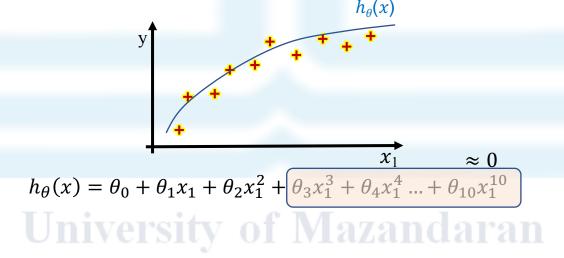
### Regularization

#### **Intuition**



Suppose we penalize and make  $\theta_3$ ,  $\theta_4$ , ...,  $\theta_{10}$  really small.

$$\min_{\theta} \ \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^{2} + 1000\theta_{3}^{2} + 1000\theta_{4}^{2} + \dots + 1000\theta_{10}^{2} \right]$$







#### **Standard linear regression:**

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^{2}$$

#### Regularized linear regression:

In linear regression, we modify its cost function by adding regularization term. The value of  $\theta_i$  is controlled by regularization parameter  $\lambda$ .

Note that m is the number of data and n is the number of features (parameters).



Normal equations to minimize  $J(\theta)$ 



**\diamond** Normal equations to minimize  $J(\theta)$ 

Given a training set, define the design matrix X

$$X := \begin{bmatrix} - & (x^{(1)})^T - \\ - & (x^{(2)})^T - \\ \vdots \\ - & (x^{(m)})^T - \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times (n+1)}$$



**\diamond** Normal equations to minimize  $J(\theta)$ 

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

if  $\lambda > 0$ :

$$\theta^* = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \right)^{-1} X^T y$$



Batch Gradient descent to minimize  $J(\theta)$ 



**A** Batch Gradient descent to minimize  $J(\theta)$ 

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

```
Repeat { oldsymbol{	heta} artheta \coloneqq oldsymbol{	heta} - oldsymbol{lpha} 
abla_{	heta} J(oldsymbol{	heta}) }
```

**α**: Learning rate



**Argonian Serice** Batch Gradient descent to minimize  $J(\theta)$ 

Repeat { 
$$\theta_0 \coloneqq \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$
 
$$\theta_j \coloneqq \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}\right]$$
 (simultaneously update for every  $j = 0, 1, ..., n$ ) }



# Regularized logistic regression



#### **Standard logistic regression:**

$$\mathbf{min}_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}) 
J(\boldsymbol{\theta}) = -\frac{1}{m} l(\boldsymbol{\theta}) = \frac{1}{m} \left[ \sum_{i=1}^{m} -y^{(i)} \log h_{\boldsymbol{\theta}}(x^{(i)}) - (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(x^{(i)})) \right]$$

#### **Regularized logistic regression:**

$$min_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta})$$

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

In logistic regression, we modify its cost function by adding regularization term. The value of  $\theta_j$  is controlled by regularization parameter  $\lambda$ .

Note that m is the number of data and n is the number of features (parameters).



## Regularized logistic regression

**A** Batch Gradient descent to minimize  $J(\theta)$ 

$$\begin{aligned} & \underset{\boldsymbol{\theta}}{\boldsymbol{min}} \boldsymbol{J}(\boldsymbol{\theta}) \\ & \boldsymbol{J}(\boldsymbol{\theta}) = -\frac{1}{m} l(\boldsymbol{\theta}) = \frac{1}{m} \left[ \sum_{i=1}^{m} -y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) - \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \boldsymbol{\theta}_{j}^{2} \end{aligned}$$

```
Repeat { oldsymbol{	heta} \coloneqq oldsymbol{	heta} - oldsymbol{lpha} 
abla_{	heta} J(oldsymbol{	heta}) }
```

**α**: Learning rate



## Regularized logistic regression

#### **A Patch Gradient descent to minimize** $J(\theta)$

Repeat { 
$$\theta_0 \coloneqq \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$
 
$$\theta_j \coloneqq \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \left[ \sum_{i=1}^m \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \right]$$
 (simultaneously update for every  $j = 0, 1, ..., n$ ) }



### Regularized

### Two advantages of using regularization:

- ❖ The *prediction error* of the regularized model is *lesser*, that is, it works well in testing data.
- ❖ The regularization model is *simpler* since it has less features (parameters).



#### References

- 1- https://www.geeksforgeeks.org
- 2- Andrew Ng, https://www.coursera.org/learn/machine-learning
- 3- https://medium.com/@qempsil0914/courseras-machine-learning-notes-week3-overfitting-and-regularization-partii-3e3f3f36a287

