

Clustering – An example





Content

- 1 Machine Learning basics
- 2 Classification and logistic regression
- 3 Text Classification



Content

2

Classification and logistic regression

- Classification - an example
- Decision boundary
- Model definition
- Cost function
- Stochastic gradient descend
- Multiclass classification



Classification - an example



A power generator

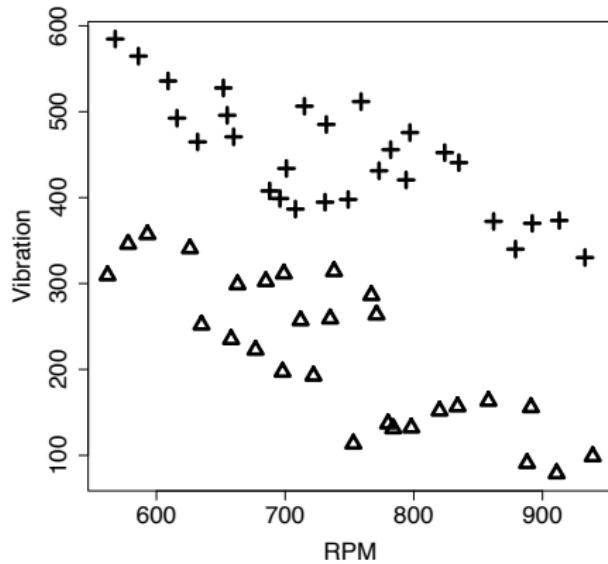
**Table:** A dataset listing features for a number of generators.

ID	RPM	VIBRATION	STATUS	ID	RPM	VIBRATION	STATUS
1	568	585	good	29	562	309	faulty
2	586	565	good	30	578	346	faulty
3	609	536	good	31	593	357	faulty
4	616	492	good	32	626	341	faulty
5	632	465	good	33	635	252	faulty
6	652	528	good	34	658	235	faulty
7	655	496	good	35	663	299	faulty
8	660	471	good	36	677	223	faulty
9	688	408	good	37	685	303	faulty
10	696	399	good	38	698	197	faulty
11	708	387	good	39	699	311	faulty
12	701	434	good	40	712	257	faulty
13	715	506	good	41	722	193	faulty
14	732	485	good	42	735	259	faulty
15	731	395	good	43	738	314	faulty
16	749	398	good	44	753	113	faulty
17	759	512	good	45	767	286	faulty
18	773	431	good	46	771	264	faulty
19	782	456	good	47	780	137	faulty
20	797	476	good	48	784	131	faulty
21	794	421	good	49	798	132	faulty
22	824	452	good	50	820	152	faulty
23	835	441	good	51	834	157	faulty
24	862	372	good	52	858	163	faulty
25	879	340	good	53	888	91	faulty
26	892	370	good	54	891	156	faulty
27	913	373	good	55	911	79	faulty
28	933	330	good	56	939	99	faulty

John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



Classification - an example



John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



Content

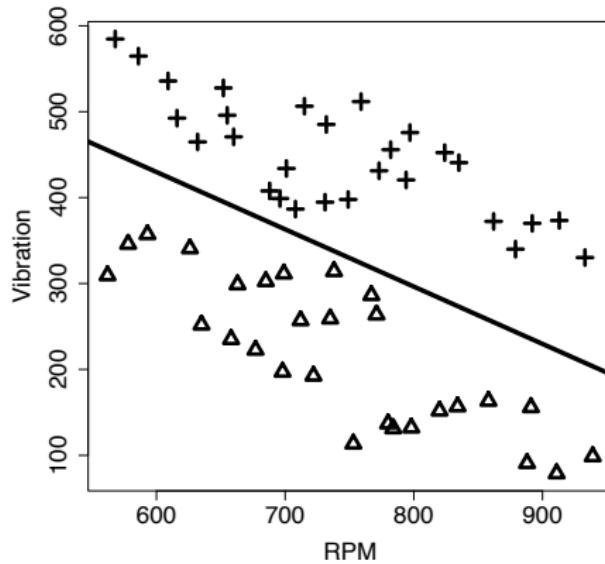
2

Classification and logistic regression

- Classification - an example
- **Decision boundary**
- Model definition
- Cost function
- Stochastic gradient descend
- Multiclass classification



Decision boundary



John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



Decision boundary

$$830 - 0.667 \times \text{RPM} - \text{Vibration} = 0$$

$$\theta_0 + \theta_1 x_1 + x_2 = 0$$

- For all “good” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 \geq 0$
- For all “faulty” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 < 0$



Decision boundary

$$830 - 0.667 \times \text{RPM} - \text{Vibration} = 0$$

$$\theta_0 + \theta_1 x_1 + x_2 = 0$$

- For all “good” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 \geq 0$
- For all “faulty” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 < 0$



Decision boundary

$$830 - 0.667 \times \text{RPM} - \text{Vibration} = 0$$

$$\theta_0 + \theta_1 x_1 + x_2 = 0$$

- For all “good” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 \geq 0$
- For all “faulty” generators, we have: $\theta_0 + \theta_1 x_1 + x_2 < 0$



Notation

$$\text{Let: } \theta = \begin{bmatrix} 1 \\ \theta_1 \\ \theta_0 \end{bmatrix}, x = \begin{bmatrix} x_2 \\ x_1 \\ 1 \end{bmatrix}$$

Then we have the decision boundary:

$$d_\theta(x) = \theta^T x = 0$$



Content

2

Classification and logistic regression

- Classification - an example
- Decision boundary
- **Model definition**
- Cost function
- Stochastic gradient descend
- Multiclass classification

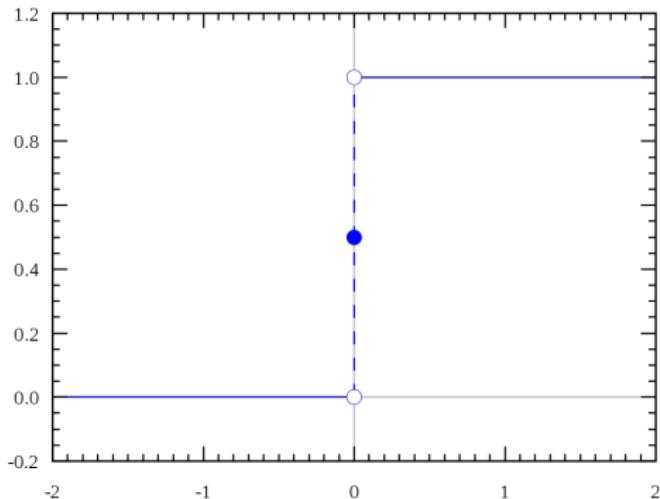


Heaviside step function

$$\text{Heaviside}(x) = \begin{cases} 1 & \text{if: } x \geq 0 \\ 0 & \text{if: } x < 0 \end{cases}$$



Heaviside step function



By Omegatron - Own work, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=801382>



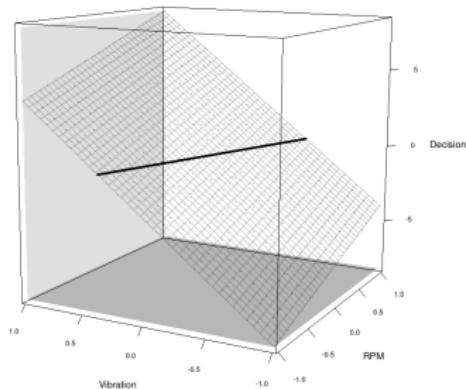
Model as a Heaviside step function

$$h_{\theta}(x) = \text{Heaviside}(d_{\theta}(x))$$

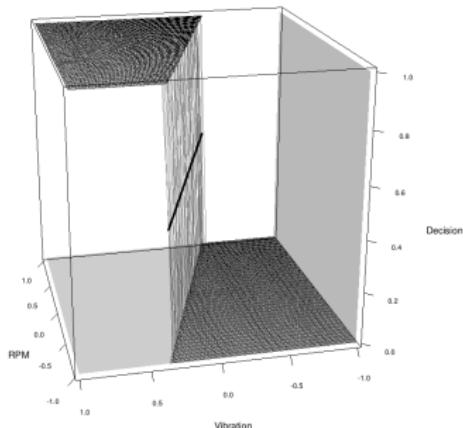
$$= \begin{cases} 1 & \text{if: } d_{\theta}(x) = \theta^T x \geq 0; \text{ for good generators} \\ 0 & \text{if: } d_{\theta}(x) = \theta^T x < 0; \text{ for faulty generators} \end{cases}$$



Model as a Heaviside step function



(a)



(b)

$$d_\theta(x)$$

$$h_\theta(x) = \text{Heaviside}(d_\theta(x))$$

John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



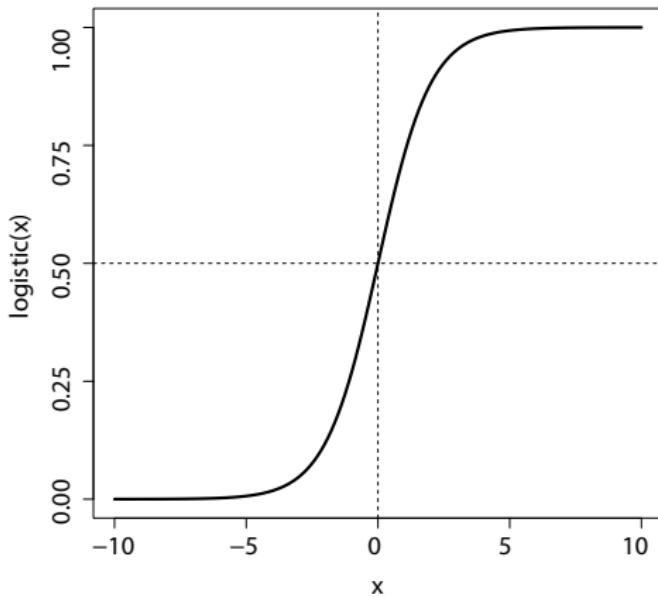
Problem of the Heaviside step function

- Heaviside step function is not derivable and hard to optimize.
- An alternative is using a Logistic function (sigmoid function):

$$\text{Logistic}(x) = \frac{1}{1 + e^{-x}}$$



Logistic function

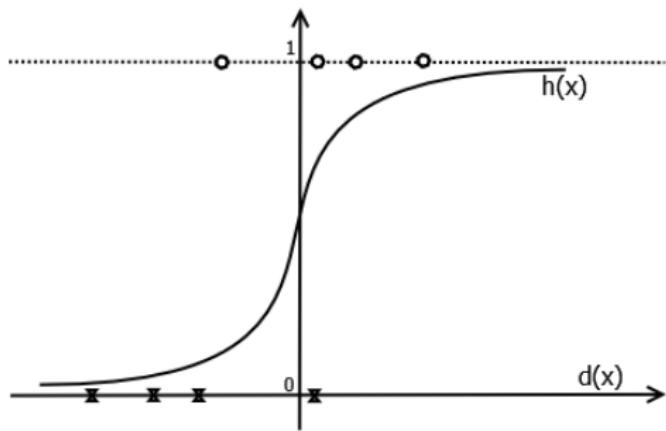


John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



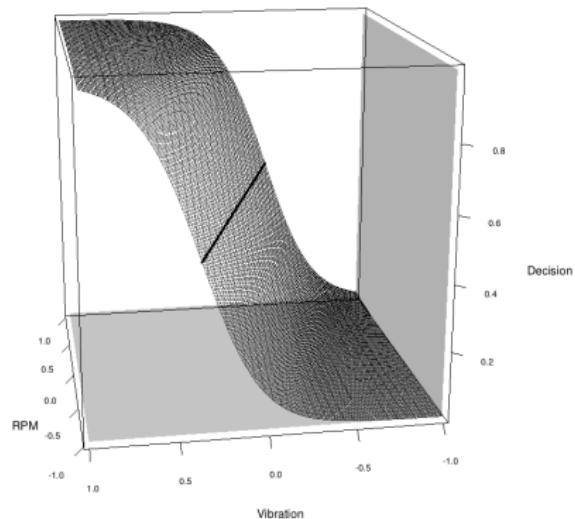
Model as a logistic function

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





Model as a logistic function



John Kelleher and Brian Mac Namee and Aoife D'Arcy, Fundamentals of Machine Learning for Predictive Data Analytics



Content

2

Classification and logistic regression

- Classification - an example
- Decision boundary
- Model definition
- **Cost function**
- Stochastic gradient descend
- Multiclass classification



Cost function

- The objective of a learning algorithm is to search a model to best predict the target feature on the training data given the inductive bias.
- One way to achieve this goal is to minimize the errors of the prediction over the training data.
- A cost function (loss function) is defined as the sum of the errors over the training samples.



Cost function

- A possible cost function:

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

- It is not a good cost function for Logistic regression because it is non-convex for a Logistic function.



Convex vs. non convex functions

Convex Vs Non-Convex

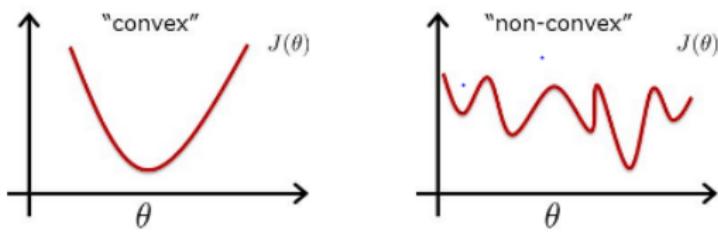


figure source: <https://www.fromthegenesis.com/artificial-neural-network-part-7/>



Cost function for Logistic regression

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n Cost(h_\theta(x^{(i)}), y^{(i)})$$

where:

$$\begin{aligned} Cost(h_\theta(x), y) &= \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases} \\ &= -[y \log(h_\theta(x)) + (1 - y) \log(1 - h_\theta(x))] \end{aligned}$$

$$y^{(*)} \in \{0, 1\}$$



Cost function for Logistic regression

$$\text{Cost}(h_\theta(x), y) = -[y \log(h_\theta(x)) + (1 - y) \log(1 - h_\theta(x))]$$

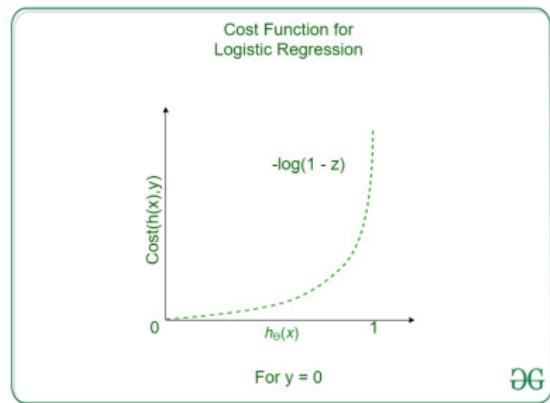
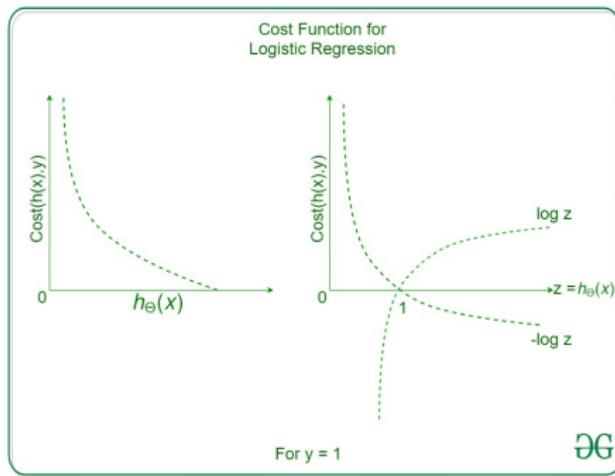


figure source: <https://www.geeksforgeeks.org/ml-cost-function-in-logistic-regression/>



Content

2

Classification and logistic regression

- Classification - an example
- Decision boundary
- Model definition
- Cost function
- Stochastic gradient descend**
- Multiclass classification



Gradient descend

- Now we have a cost function which defines the errors of a model over the training data.
- Next we need to search all the possible models to find a model with the minimal cost.
- We use a gradient descend algorithm for this purpose.



Gradient descend

Have some function $J(\theta_0, \theta_1)$

Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

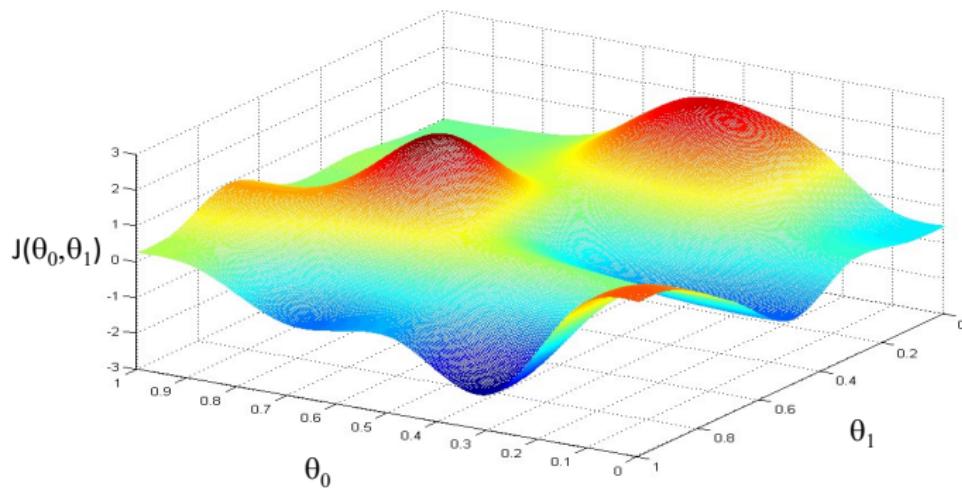
Outline:

- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum

Andrew Ng, Machine Learning, Coursera course



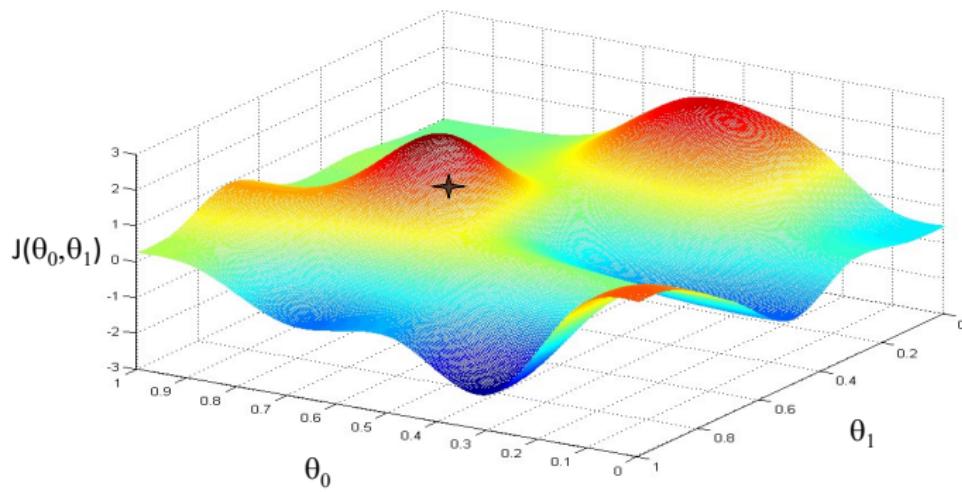
Gradient descend



Andrew Ng, Machine Learning, Coursera course



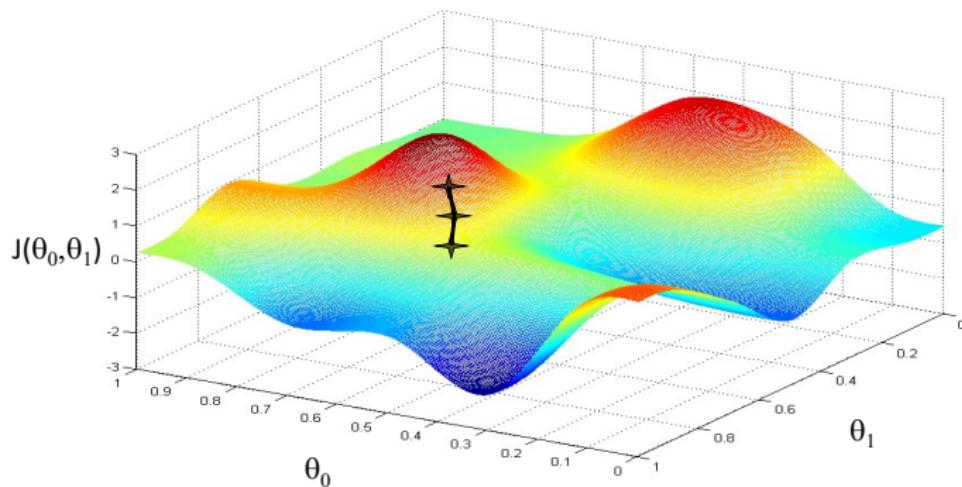
Gradient descend



Andrew Ng, Machine Learning, Coursera course



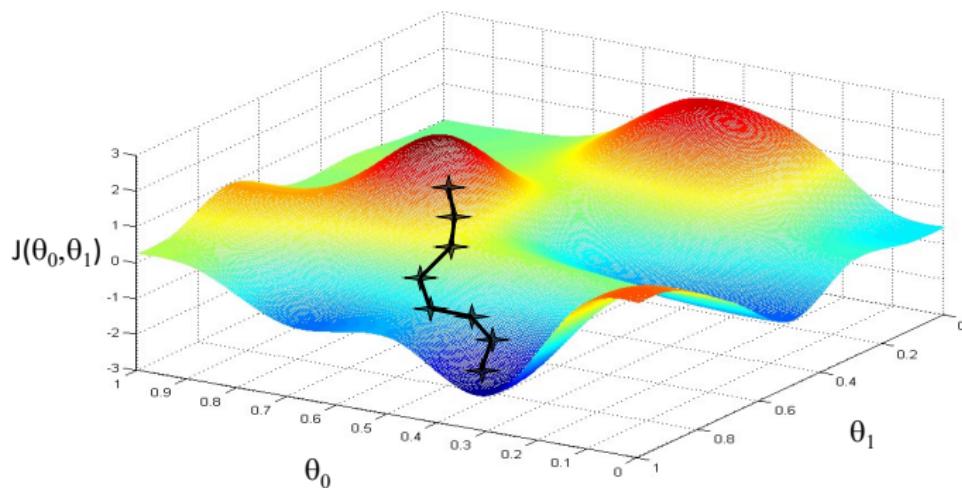
Gradient descend



Andrew Ng, Machine Learning, Coursera course



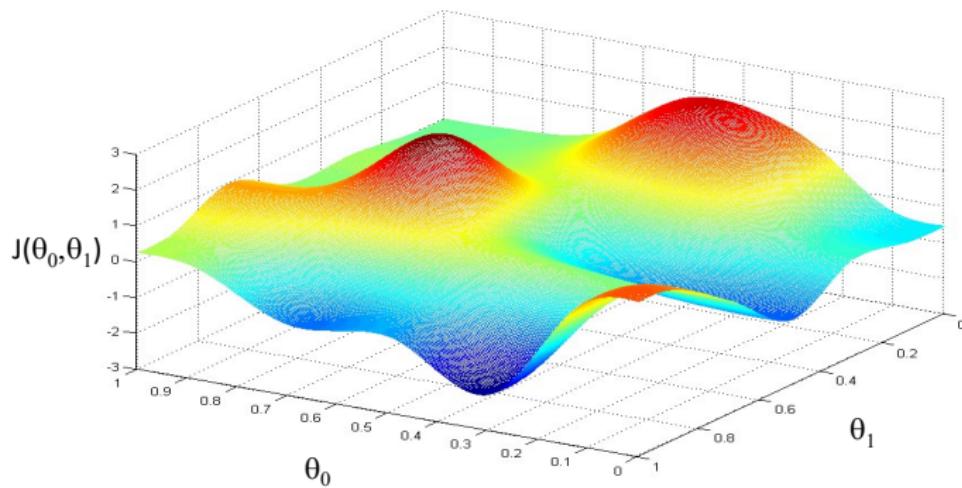
Gradient descend



Andrew Ng, Machine Learning, Coursera course



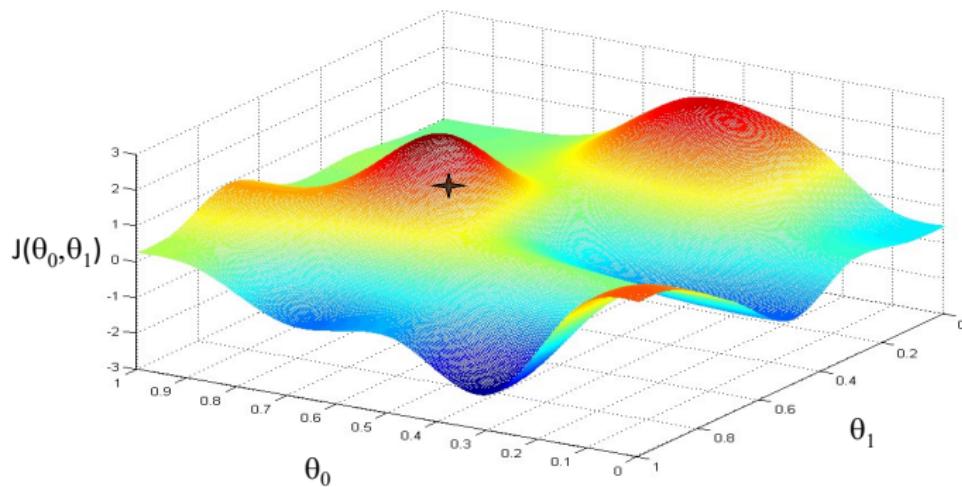
Gradient descend



Andrew Ng, Machine Learning, Coursera course



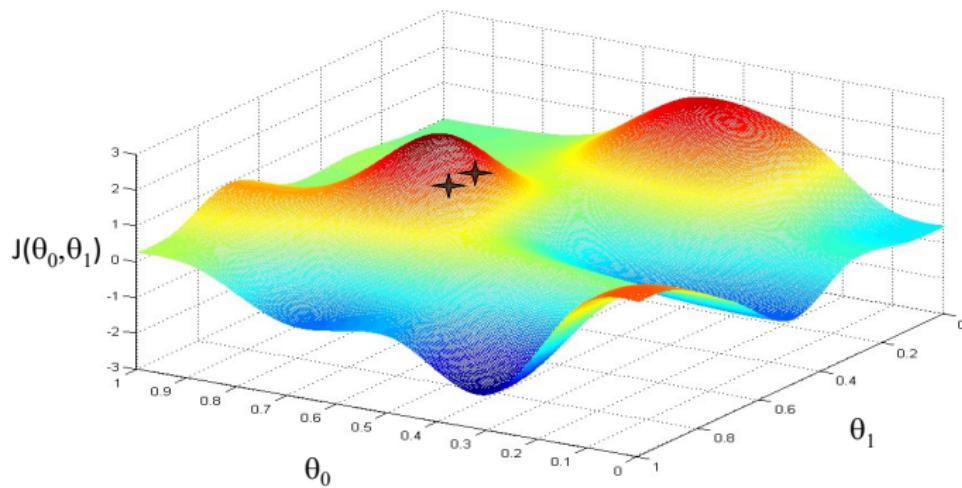
Gradient descend



Andrew Ng, Machine Learning, Coursera course



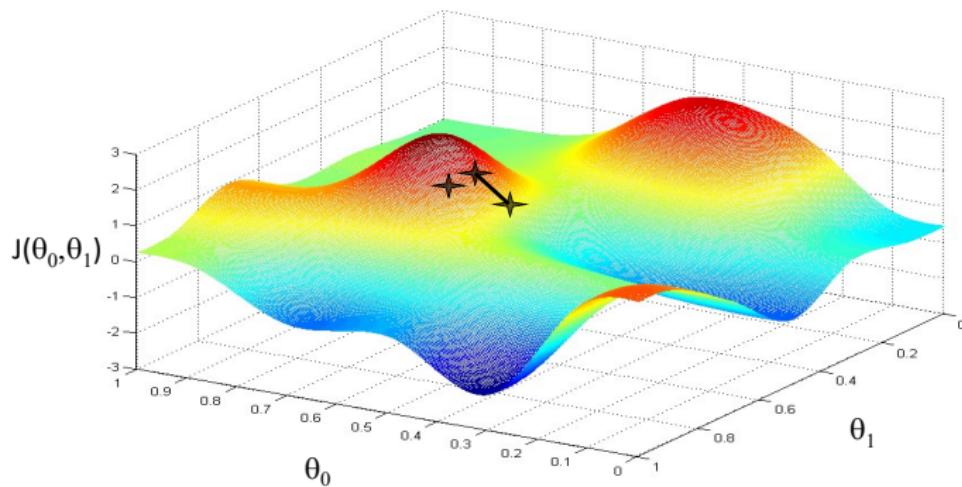
Gradient descend



Andrew Ng, Machine Learning, Coursera course



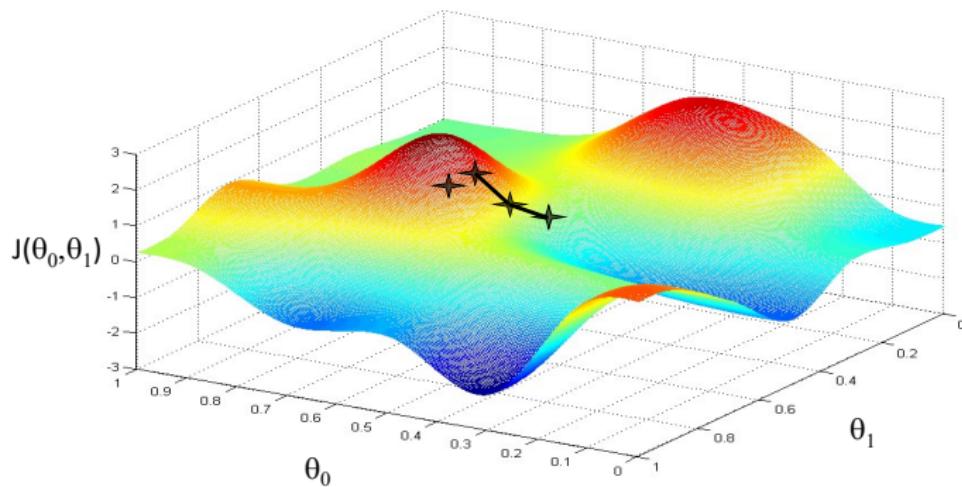
Gradient descend



Andrew Ng, Machine Learning, Coursera course



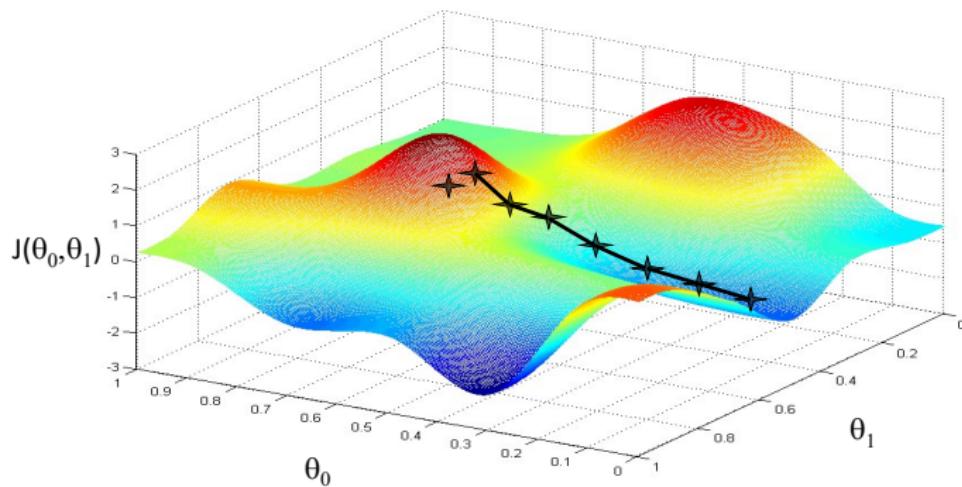
Gradient descend



Andrew Ng, Machine Learning, Coursera course



Gradient descend



Andrew Ng, Machine Learning, Coursera course



Gradient descend algorithm

Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \begin{matrix} \text{(simultaneously update} \\ j = 0 \text{ and } j = 1) \end{matrix}$$

Andrew Ng, Machine Learning, Coursera course



Gradient descend algorithm

Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

}

Correct: Simultaneous update

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Andrew Ng, Machine Learning, Coursera course



Gradient descend algorithm

Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

}

Correct: Simultaneous update

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect:

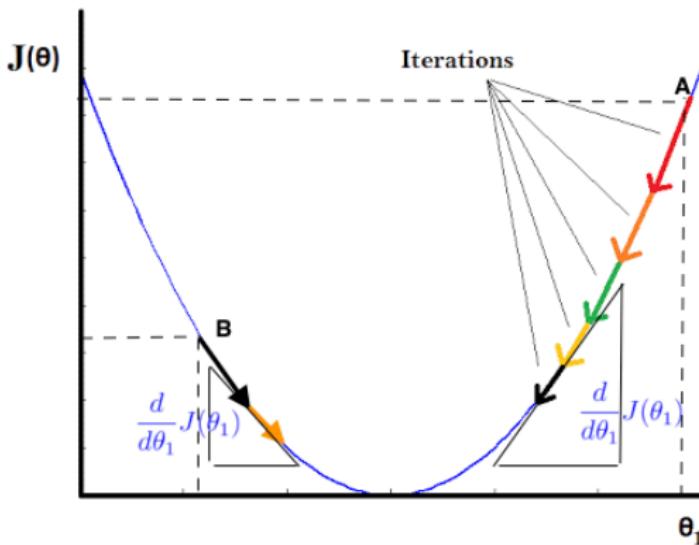
$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

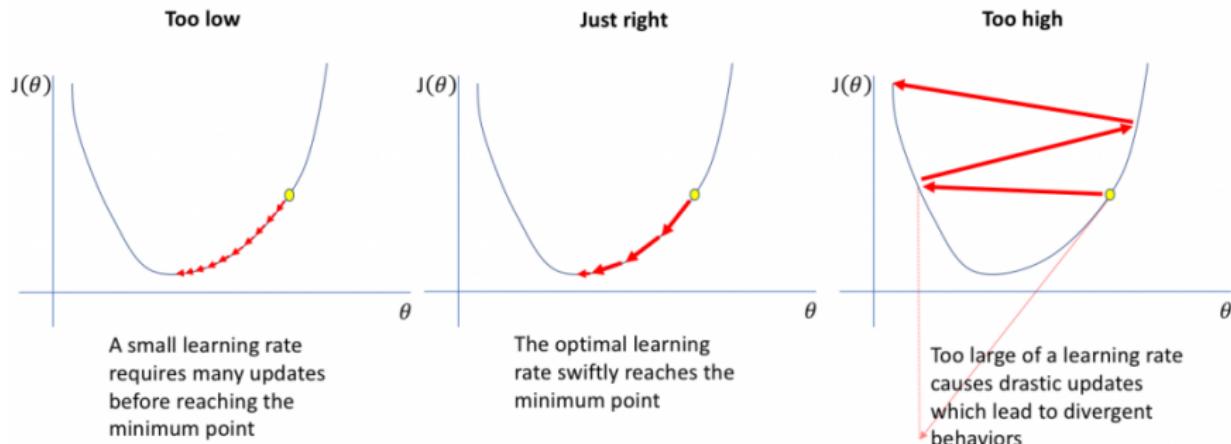
$$\theta_1 := \text{temp1}$$

Andrew Ng, Machine Learning, Coursera course



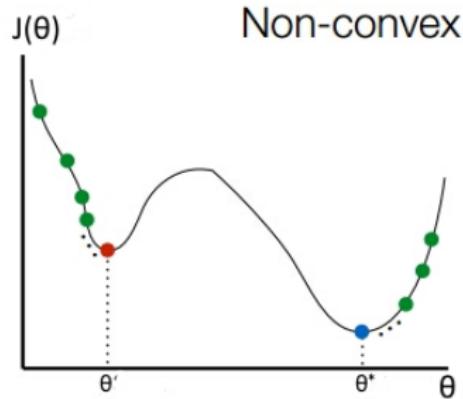
- The derivation $\frac{\partial J(\theta)}{\partial \theta}$ gives the direction of the movement.
- The learning rate α is used to adjust the size for each step.

figure source: <https://machinelearningmedium.com/2017/08/15/gradient-descent/>



The influence of the learning rate.

figure source: <https://ithelp.ithome.com.tw/m/articles/10204032>



A non-convex error surface may lead to local optima.



Gradient descend for Logistic regression

Model

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

Parameters

$$\theta_0, \theta_1$$

Cost Function

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

Goal

$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta)$$



Gradient descend for Logistic regression

Gradient descent:

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

Want $\min_{\theta} J(\theta)$:

Repeat {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta);$$

(simultaneously update all θ_j)

}



Gradient descend for Logistic regression

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

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= -\frac{1}{n} \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \left[y^{\{i\}} \log(h_\theta(x^{\{i\}})) + (1 - y^{\{i\}}) \log(1 - h_\theta(x^{\{i\}})) \right] \\ &= \dots \dots \\ &= \frac{1}{n} \sum_{i=1}^n (h_\theta(x^{\{i\}}) - y^{\{i\}}) x_j^{\{i\}}\end{aligned}$$



Gradient descend for Logistic regression

Gradient descent:

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

Want $\min_{\theta} J(\theta)$:

Repeat {

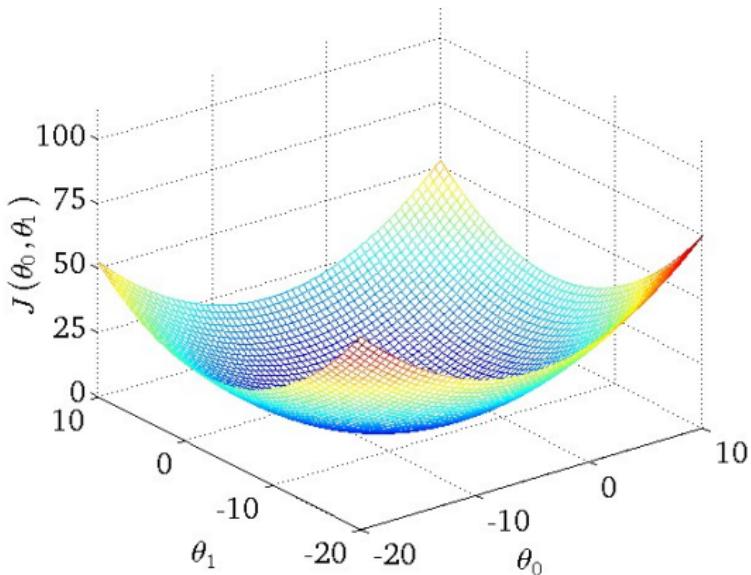
$$\theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(simultaneously update all θ_j)

}



Gradient descend for logistic regression



Fortunately the error surface for logistic regression is convex.

Andrew Ng, Machine Learning, Coursera course