# Cost Function

* **Training Set :** $m$ - examples $\{(x^{(1)}, y^{(1)}), \text{--------}, (x^{(m)}, y^{(m)})\}$

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_0 = 1 \quad y \in = \{0, 1\}$$
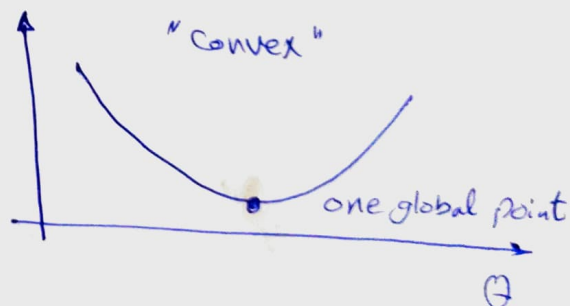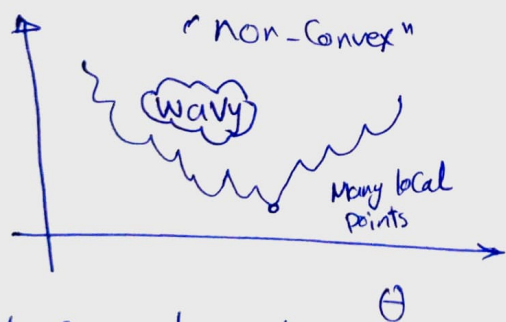
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose Parameters $\theta$ ?

→ **For Linear Regression :**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2}_{\text{Cost } (h(x^{(i)}) - y^{(i)}) \quad \text{(loss)}}$$

$$\text{Cost Function} = \frac{1}{2} \left[ h_\theta(x) - y \right]^2$$

If we apply this equ to logistic Regression where $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$ it will produce a non-Convex fn.



"non-Convex"

wavy

Many local points

$\theta$

"Convex"

one global point

$\theta$

→ No guarantee to Converge at min. point if we applied gradient decsent.
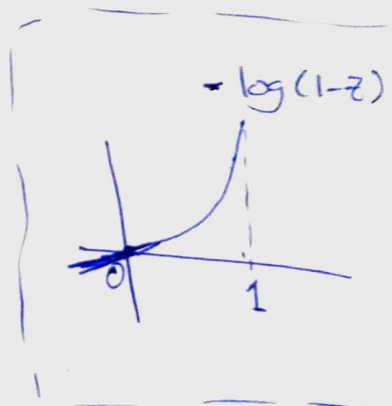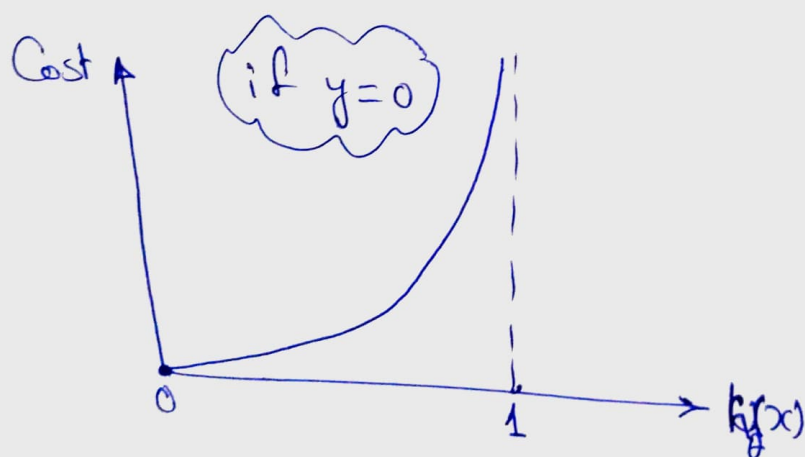
# * Logistic Regression Cost Function:

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

if $y=1$

$y=1$, $h_\theta(x) = 1 \rightarrow \boxed{Cost = 0}$

as $h_\theta(x) < 1$

$\log(z)$

$-\log(z)$

$h_\theta(x)$

<u>But</u>, as $h_\theta(x) \rightarrow 0$ $(P(y=1|x;\theta) = 0)$
    $Cost \rightarrow \infty$

Captures intuition that $h_\theta(x) = 0$ $(P(y=1|x;\theta) = 0)$
but actually $\underline{y=1}$, so we penalize learning algorithm
by a very large Cost.

if $y=0$

$-\log(1-z)$

## • Notes:

→ If $h_\theta(x) = y \rightarrow Cost = 0$ (for $y=0$ and $y=1$)

→ If $y=0$, then $(Cost \rightarrow \infty$ as $h_\theta(x) \rightarrow 1)$

→ Regardless of whether $y=0$ or $y=1$, if $h_\theta(x) = 0.5$, then
    $Cost > 0$

# Logistic Regression Cost function

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

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

for single example

↳ Compressing it in one eqn :

$$Cost(h_\theta(x), y) = -y\log(h_\theta(x)) - (1-y)\log(1-h_\theta(x))$$

So,

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

And why do we choose this particular function, while it may there could other cost functions. This cost function can be derived from statistics using the principle of Maximum Likelihood Estimation which is an idea in statistics for how to efficiently find parameters data for different models. And it also has a nice property that it is Convex

So, we want to fit parameter $\theta$ to make $J(\theta)$ is minimum to make a prediction given new $x$ :

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

→ We will use Gradient Descent

# Gradient Descent

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

Simulataneously update all $\theta_j$

➤ **For Logistic Regression**

Repeat {
$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^{m} \left[ (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right]$$
}

Simulataneously update all $\theta_j$

➡ It's Actually the same algorithm of gradient Descent for Linear Regression, but :

    ✳ <u>Linear Regression</u> : $h_\theta(x) = \theta^T x$

    ✳ <u>Logistic Regression</u> : $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$

➡ <u>A vectorized Implementation</u> :    (Run faster)

$$h = g(X\theta)$$
$$J(\theta) = \frac{1}{m} \cdot [-y^T \log(h) - (1-y^T) \log(1-h)]$$

<u>Gradient Descent</u> : $\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$

➔ <u>Also, Feature Scaling</u>

    ↳ helps speed up gradient descent for the log-regression algorithm.