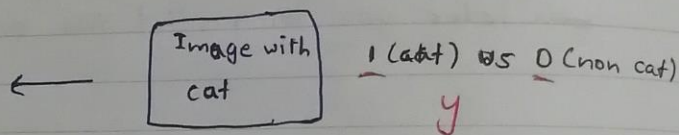
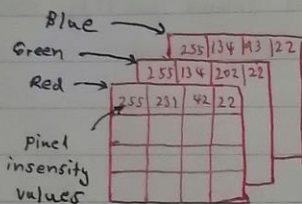


Week 2 Neural Networks Basics

Supervised Learning

Binary Classification

- * Logistic regression is an algorithm for binary classification.
- * Let's assume we have a image in our computer. That image always stored in three separated metrics. each metrics represent by one color (Red, Green, Blue)



- * if you store image 64x64 pixels image, then you will have three 64x64 pixels metrics.
- * Now we have to convert those pixel intensity values into a feature vector. There for we have to assign those values to our input feature vector x .

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$$

pixel x pixel x no. of metrics

$$64 \times 64 \times 3 = 12288$$

$$n_x = 12288$$

- * In binary classification, our goal is learn a classifier that can input an image represented by feature vector x and predict whether given label is 1 or 0. (In my case cat or not cat)

* A single training example is represented by (x, y)

$$(x, y) \quad x \in \mathbb{R}^{n_x}, \quad y \in \{0, 1\}$$

x is an n_x dimensional feature

y is the label either 0 or 1,

* Your training set comprise: m training examples

* first training example: $(x^{(1)}, y^{(1)})$

* second training example: $(x^{(2)}, y^{(2)})$

* m training example: $(x^{(m)}, y^{(m)})$

$$m \text{ training example: } \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

* To define those training set create matrix using X

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$\xleftarrow{\quad m(\text{column}) \quad} \xrightarrow{\quad}$

* In dimensional metrics

$$X \in \mathbb{R}^{n_x \times m}$$

* In Python we use command

$$X.\text{shape} = (n_x, m)$$

* In output (Y) $\rightarrow Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

Logistic Regression

Given x , want $\bar{y} = P(y=1|x)$

* Let's assume x is a picture and you want \bar{y} to tell you what is the chance that this is a cat picture

* Given that parameter of logistic regression will be w which also an n dimensional vector ~~$x \in \mathbb{R}^n$~~ and together with b which is just a real number

Parameters : $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

* Then, given x, w and b how can we generate output \bar{y}

* What if we use

$$\text{Output } \bar{y} = w^T x + b$$

* Above formula represent linear regression, but it is not good algorithm for binary classification because we want \bar{y} hat to be $(y=1)$.

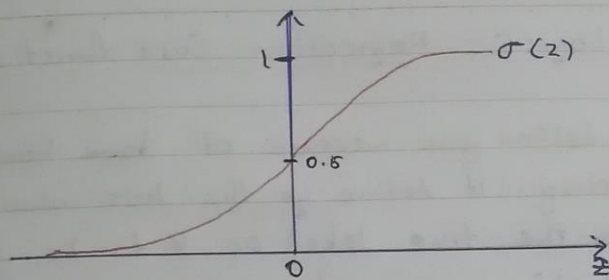
$\therefore \bar{y}$ should be in between 1 and 0 ($0 \leq \bar{y} \leq 1$)

* $w^T x$ can give more than 1 or negative value.

* Therefore, In Logistic regression we use σ (sigmoid function) to put \bar{y} in between 1 and 0.

$$\text{Output } \bar{y} = \sigma(\underbrace{w^T x + b}_z)$$

$$\bar{y} = \sigma(z)$$



sigmoid function

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

$$E_n: X_0 = 1, \quad x \in \mathbb{R}^{n_x + 1}$$

$$\bar{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{matrix} \left. \vphantom{\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix}} \right\} b \\ \left. \vphantom{\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix}} \right\} w \end{matrix}$$

Logistic Regression Cost function

* Loss function tells you measure of ~~how~~ how well our algorithm is doing. it define ~~you~~ the loss when your algorithm outputs \bar{y} and the true label as y to be the square error or one half a square error.

Loss (Error) function: $L(\bar{y}, y) = \frac{1}{2}(\bar{y} - y)^2$
loss function

$$L(\bar{y}, y) = - (y \log \bar{y} + (1-y) \log(1-\bar{y})) \quad 0 \leq y \leq 1$$

~~$y=0$~~ $y=1$

$$L(\bar{y}, y) = -\log \bar{y} \quad y=1. \text{ because we want } -\log \bar{y} \text{ to be small as possible}$$

$y=0$

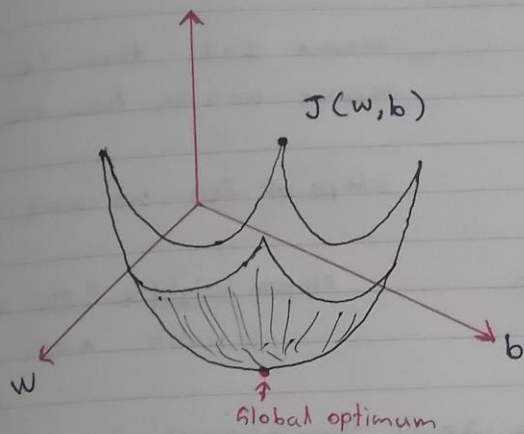
$$L(\bar{y}, y) = -\log(1-\bar{y}) \quad y=0. \text{ because we want } -\log(1-\bar{y}) \text{ to be large as possible}$$

* Cost function tells you measure of how ~~at~~ well our entire training set doing. Cost function denote by J .

Cost function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\bar{y}^{(i)}, y^{(i)})$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\bar{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \bar{y}^{(i)} + (1-y^{(i)}) \log(1-\bar{y}^{(i)})]$$

Gradient Descent



w : single real number
 b : single real number
 $J(w, b)$: height of the surface



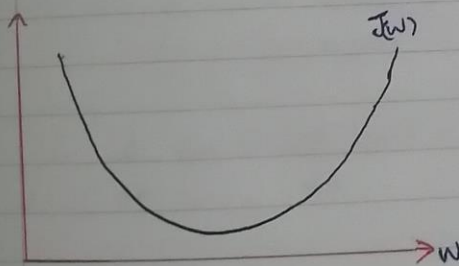
convex



non convex

convex : have one optima
 nonconvex : have multiple optima

Let's remove "b" for a moment



Repeat {

$$w : w - \alpha \frac{dJ(w)}{dw}$$

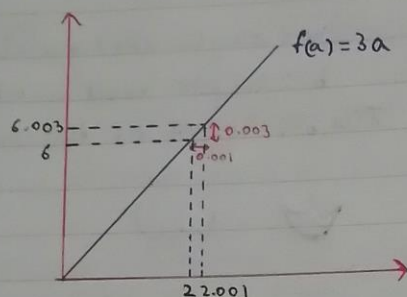
}

$J(w, b)$

$$w : w - \alpha \frac{dJ(w, b)}{dw}$$

$$b : b - \alpha \frac{dJ(w, b)}{db}$$

Derivatives



assume $a=2$ then $f(a)=6$
assume $a=2.001$ then $f(a)=6.003$

slope of $f(a)$ at $a=2$ is

$$\text{slope} = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.001} = \underline{\underline{3}}$$

for another example assume $a=5$ then $f(a)=15$

assume $a=5.001$ then $f(a)=15.003$

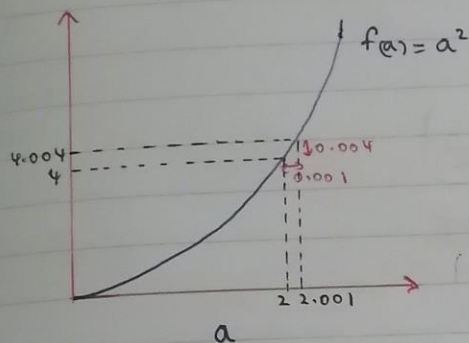
Then slope of $f(a)$ at $a=5$ is

$$\text{slope} = \frac{0.003}{0.001} = \underline{\underline{3}}$$

$$\therefore \frac{df(a)}{da} = 3 = \frac{d}{da} f(a)$$

* if you increase value of " a " ~~slope~~ by 0.001, the value of $f(a)$ goes up by three times as much. Therefore, we can tell this function has same slope everywhere.

More Derivative Exempler



$$a = 2 \quad \therefore f(a) = 4$$

$$a = 2.001 \quad \therefore f(a) \approx 4.004 \quad (4.004001)$$

slope of $f(a)$ at $a = 2$ is

$$\text{slope} = \frac{0.004}{0.001} = 4 = \frac{d f(a)}{d a}$$

for another example $a = 5 \quad \therefore f(a) = 25$

$$a = 5.001 \quad \therefore f(a) \approx 25.010$$

$$\therefore \frac{d f(a)}{d a} = \frac{0.01}{0.001} = 10, \text{ when } a = 5$$

$$\therefore \frac{d f(a)}{d a} = \frac{d}{d a} (a^2) = \underline{\underline{2a}}$$

* slope will change time to time in these functions.

Computation Graph

Example

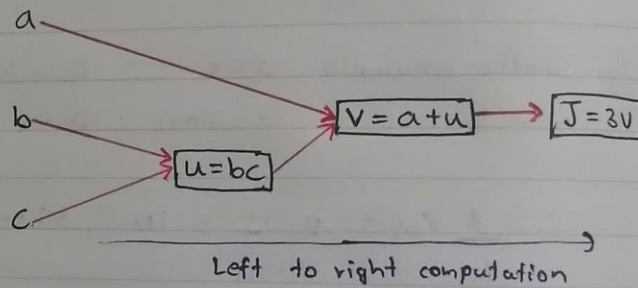
$$J(a,b,c) = 3(a + \underbrace{bc}_u)$$

$\underbrace{\quad}_v$
 $\underbrace{\quad}_J$

$$u = bc$$

$$v = a + u$$

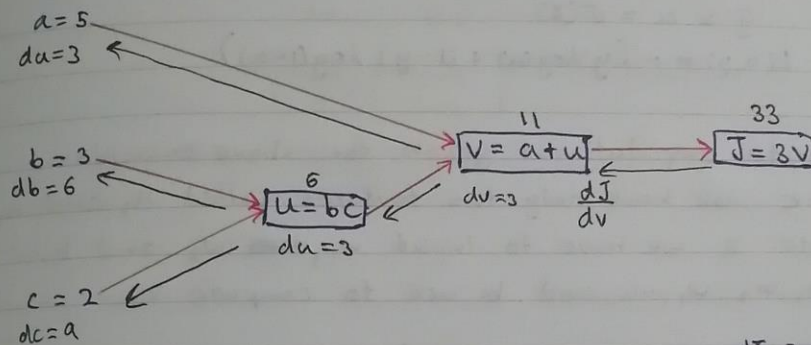
$$J = 3v$$



* In here "J" is our cost function. because we are trying to minimize it.

* we saw in here is left to right computation (forward propagation)

Derivatives with a Computation Graph



$$\begin{aligned}
 u=6 &\rightarrow 6.001 & \frac{dJ}{du} &= \frac{dJ}{dv} \cdot \frac{dv}{du} & \frac{dJ}{dv} &= \frac{0.003}{0.001} = 3 \\
 v=1 &\rightarrow 11.001 & \frac{dv}{du} &= \frac{1}{3} & \frac{dv}{du} &= \frac{0.001}{0.001} = 1 \\
 J=33 &\rightarrow 33.003 & \therefore \frac{dJ}{du} &= 3
 \end{aligned}$$

$$\begin{aligned}
 b=3 &\rightarrow 3.001 & \frac{dJ}{db} &= \frac{dJ}{du} \cdot \frac{du}{db} & \frac{dJ}{du} &= \frac{0.006}{0.002} = 3 \\
 u=bc=6 &\rightarrow 6.002 & \frac{du}{db} &= \frac{1}{2} & \frac{du}{db} &= \frac{0.002}{0.001} = 2 \\
 J=33 &\rightarrow 33.006 & \therefore \frac{dJ}{db} &= 6
 \end{aligned}$$

$$\begin{aligned}
 c=2 &\rightarrow 2.001 & \frac{dJ}{dc} &= \frac{dJ}{du} \cdot \frac{du}{dc} & \frac{dJ}{du} &= \frac{0.004}{0.003} = 3 \\
 u=bc=6 &\rightarrow 6.003 & \frac{du}{dc} &= \frac{1}{3} & \frac{du}{dc} &= \frac{0.003}{0.001} = 3 \\
 J=33 &\rightarrow 33.009 & \therefore \frac{dJ}{dc} &= 9
 \end{aligned}$$

* we saw in here is right to left computation.
 (backward propagation)

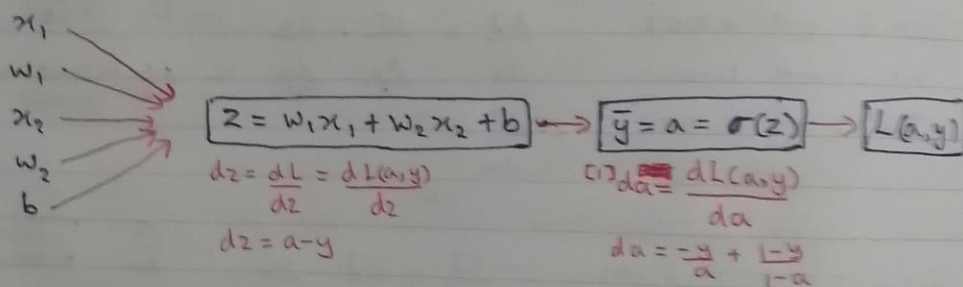
Logistic Regression Gradient Descent

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

- * Let's write computational graph for above example.
- * We assume we have only two features called x_1 and x_2 .
- * To compute z we have to input w_1 , w_2 and b .
- * Those x_1, x_2, w_1, w_2 and b use to compute z .



- * In logistic regression, what we want to do is to modify the parameters " w 's and " b ", in order to reduce loss $L(a,y)$.
- * To reduce loss, first thing to do is going backwards to compute the derivative of this loss with respect to da .

$$\begin{aligned}
 \frac{da}{dL(a,y)} &= \frac{dL(a,y)}{da} \\
 &= [-y \log(a)] \frac{dL(a,y)}{da} + [(1-y) \log(1-a)] \frac{dL(a,y)}{da} \\
 &= -\frac{y}{a} + \frac{1-y}{1-a}
 \end{aligned}$$

$$dz = \frac{dL}{dz}$$

$$= \frac{dL}{da} \cdot \frac{da}{dz}$$

$$= \left[\frac{-y}{a} + \frac{1-y}{1-a} \right] \times [a(1-a)]$$

$$dz = \underline{\underline{a-y}}$$

$$\begin{aligned} \frac{da}{dz} &= \sigma(z) \times (1 - \sigma(z)) \\ &= a \times (1-a) \\ &= \underline{\underline{a(1-a)}} \end{aligned}$$

next we assume how much we have to change w s and b

$$\therefore \frac{dL}{dw_1} = dw_1 = x_1 dz \quad \left| \quad dw_2 = x_2 dz \quad \right| \quad db = dz$$

$$\therefore w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

Gradient Descent on m Examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{d}{dw_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{d}{dw_1} L(a^{(i)}, y^{(i)})}_{dw_1^{(i)} = (x^{(i)}, y^{(i)})}$$

Example

$$J=0, dw_1=0, dw_2=0, db=0$$

For $i=1$ to m

[1]

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log (1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

[2]

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J /= m \leftarrow J = J/m$$

$$dw_1 /= m$$

$$dw_2 /= m$$

$$db /= m$$

weaknesses above code

* to ~~imple~~ implement logistic regression this way, we have to write two for loops. first ~~loop~~ for loop is given for loop^[1] above which is use to iterate until ~~we~~ m training examples. The second for loop use for over all the features over here. [2]

* those for loops make your algorithm slow. As a solution deep learning proposed "vectorization". It can use to get rid of from for loops.

Vectorization

- * This is use to get rid of from explicit for loop from code
- * In deep learning we have to work with large datasets. Therefore, it is must to run our codes efficiently. Vectorization help us to reduce run time those codes.

$$Z = w^T x + b$$

Non-vectorized

$$Z = 0$$

for i in range(n-x)

$$Z += w[i] * x[i]$$

$$Z += b$$

vectorized

$$Z = \text{np.dot}(w, x) + b$$

$w^T x$

More Vectorization Examples

- * If you want to compute a vector "u" as the product of matrix "A" and another vector "v"

$$u = Av$$

- * Then definition of our matrix is

$$u_i = \sum_j A_{ij} v_j$$

* If we use non vectorized method

```
u = np.zeros(n, 1)
```

```
for i in range(n):
```

```
    for j in range(n):
```

```
        u[i] += A[i][j] * v[j]
```

* If we use vectorized method

```
u = np.dot(A, v)
```


Vectorizing Logistic Regression

* Let's assume you have M training examples

\therefore

first training example

$$z^{(1)} = w^T x^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

second training example

$$z^{(2)} = w^T x^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

third training example

$$z^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

* Likewise you have to do for m training examples.

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ \dots \ b] = [w^T x^{(1)} + b \quad w^T x^{(2)} + b \quad \dots \quad w^T x^{(m)} + b]$$

$$Z = \text{np.dot}(w.T, X) + b$$

* Python automatically take real number b and expands it out to $1 \times m$ row vector. This operation is called as "broadcasting"

* Likewise

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

Vectorizing Logistic Regression's Gradient Output

- * You have remembered in previous lessons we use $dz = a - y$. Let's try apply it to vectorization. Let's assume again you have m training examples.

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots \quad dz^{(m)} = a^{(m)} - y^{(m)}$$

- * We can vectorized the into like this.

$$dZ = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}]$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$\therefore dZ = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$

- * Now we have get rid of from one for loop which calculate dz . Now we have second for loop which contains below

$$dw = 0$$

$$dw += x^{(1)} dz^{(1)}$$

$$dw += x^{(2)} dz^{(2)}$$

\vdots

$$dw += x^{(m)} dz^{(m)}$$

$$dw /= m$$

$$db = 0$$

$$db += dz^{(1)}$$

$$db += dz^{(2)}$$

\vdots

$$db += dz^{(m)}$$

$$db /= m$$

- * Now we have to vectorize above for loop contains.

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \text{np.sum}(dZ)$$

$$dw = \frac{1}{m} X dZ^T$$

$$= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]_{n \times 1}$$

* Now Let's compare for 'loop version and vectorized version

$$J=0, dw_1=0, dw_2=0, db=0$$

for i=1 to m:

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 = x_1^{(i)} dz^{(i)}$$

$$dw_2 = x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/m$$

$$dw_1 = dw_1/m$$

$$dw_2 = dw_2/m$$

$$db = db/m$$

for item in range(1000) ←

$$z = w^T X + b$$

$$= \text{np.dot}(w.T, X) + b$$

$$A = \sigma(z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} \text{np.sum}(dZ)$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

but
* you want one for loop run
the vectorized model

Broadcasting in Python

* Broadcasting technique uses to make our Python code run faster.

	Apples	Beef	Eggs	Potatoes	
Carb	56	0	4.4	68	$= A_{(3,4)}$
Protein	1.2	104	52	8	
Fat	1.8	135	99	0.9	

*
$$\begin{bmatrix} \\ \\ \end{bmatrix}_{(m,n)} \begin{matrix} \rightarrow + \\ \rightarrow - \\ \rightarrow / \\ \rightarrow * \end{matrix} \begin{bmatrix} \\ \\ \end{bmatrix}_{(1,n)} = \begin{bmatrix} \\ \\ \end{bmatrix}_{(m,n)}$$

*
$$\begin{bmatrix} \\ \\ \end{bmatrix}_{(m,n)} \begin{matrix} \rightarrow + \\ \rightarrow - \\ \rightarrow / \\ \rightarrow * \end{matrix} \begin{bmatrix} \\ \\ \end{bmatrix}_{(m,1)} = \begin{bmatrix} \\ \\ \end{bmatrix}_{(m,n)}$$