# MNIST Classification: Multiple Linear Regression and Logistic Regression

CMPUT 328 Nilanjan Ray

#### **MNIST Dataset**



#### Classify images into digits

Each image is 28x28

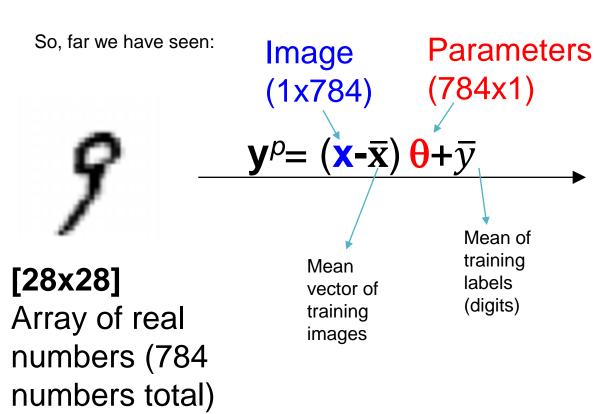
10 labels

**55,000** training images

**5,000** validation images

**10,000** test images.

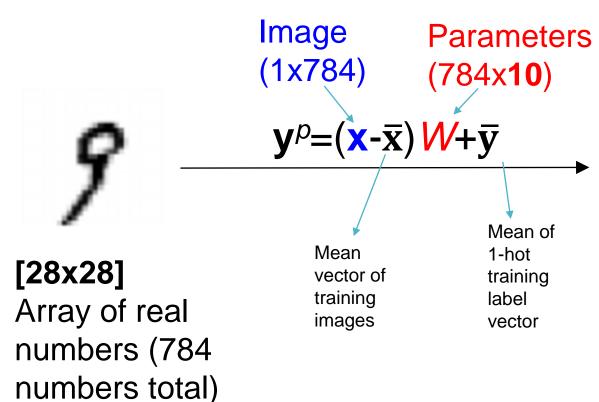
## Linear regression



1 number, indicating digit

	Digit		
<i>x</i> <sub>1</sub>	<i>X</i> <sub>2</sub>	 X <sub>784</sub>	v
0.1	0.3	 0.0	0
0.2	0.1	 0.5	1
0.0	0.98	 0.8	9
0.5	0.25	 0.36	?
0.1	0.95	 0.1	?

## Multiple or Vector Linear Regression



10 numbers, indicating class scores

Pixel values (feature) Digit: 1-not vector								
<b>x</b> <sub>1</sub>	<i>x</i> <sub>2</sub>		X <sub>784</sub>	<b>y</b> <sub>1</sub>		<b>y</b> <sub>10</sub>		
0.1	0.3		0.0	0		1		
0.2	0.1		0.5	1		0		
0.0	0.98		0.8	0		1		
0.5	0.25		0.36	?		?		
0.1	0.95		0.1	?		?		

Divolvalues (facture) Digit: 1 hot vector

## **Multiple** Linear Regression: PyTorch Implementation

 $\mathbf{y}^p = (\mathbf{x} - \bar{\mathbf{x}})W + \bar{\mathbf{y}}$ Prediction model:

Regularized loss function: 
$$L = \frac{1}{2} \sum_{i=1}^{n} ||\mathbf{y}_{i}^{p} - \mathbf{y}_{i}||^{2} + \frac{\gamma}{2} ||W||^{2}$$

https://en.wikipedia.org/wiki/Matrix calculus

This derivation requires matrix-vector differentiation

Gradient of loss function:

$$\nabla L = (X^T X + \gamma I)W - X^T Y$$

Equating gradient of *L* to zero matrix and solving for W gives us:

and matrix Y is defined as: 
$$Y = \begin{bmatrix} \mathbf{y}_1 - \overline{\mathbf{y}} \\ \vdots \\ \mathbf{y}_n - \overline{\mathbf{y}} \end{bmatrix}$$

where matrix X is defined as:  $X = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix}$ 

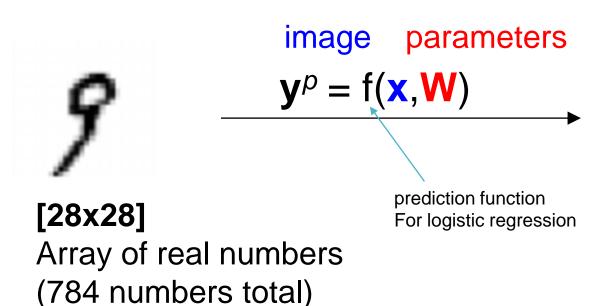
$$W = (X^T X + \gamma I)^{-1} X^T Y$$

and I is an identity matrix of size 784-by-784

We will "minimally" modify our linear regression scripts into multiple linear regression implementations!

## **Logistic** Regression

Would it not be nice if we can predict class probabilities instead of scores?



**10** numbers, indicating class probabilities

Pixel values (feature) Digit: 1-not vector							
<b>x</b> <sub>1</sub>	<i>x</i> <sub>2</sub>		X <sub>784</sub>	<b>y</b> <sub>1</sub>		<b>y</b> <sub>10</sub>	
0.1	0.3		0.0	0		1	
0.2	0.1		0.5	1		0	
0.0	0.98		0.8	0		1	
0.5	0.25		0.36	?		?	
0.1	0.95		0.1	?		?	

Divolvalues (feature) Digit, 1 hot voctor

## **Logistic** Regression

Can we modify scores from multiple regression function to output probabilities?

What is a suitable loss function for classification?

## Logistic regression: from multiple linear regression

Scores from multiple linear regression: 
$$\mathbf{s}_i = (\mathbf{x}_i - \bar{\mathbf{x}})W + \bar{\mathbf{y}}$$
 or  $\mathbf{s}_{i,k} = (\mathbf{x}_i - \bar{\mathbf{x}})W_{:,k} + \bar{\mathbf{y}}_k$   
Score for  $k^{\text{th}}$  class,  $k = 0, \dots, 9$ 

Predicted probability for 
$$k^{th}$$
 class:  $y_{i,k}^p = \frac{\exp(s_{i,k})}{\sum_{c=0}^9 \exp(s_{i,c})}$ 

"Softmax" function

## Logistic regression: loss function

Cross entropy loss: 
$$loss(\mathbf{y}^p, \mathbf{y}) = -\sum_{k=0}^{9} \mathbf{y}_k \log(\mathbf{y}_k^p)$$

Why this loss function? What does it mean? Why not use Euclidean loss as in MLR?