## **Notation**

- Supervised Learning involves supplying a number (m) of examples.
- Each example is a pair consisting of
  - vector  $\mathbf{x}$  consisting of  $n \geq 1$  features (attributes)
  - scalar (sometimes a vector) **y** 
    - $\circ$  referred to as the *target* value or *label* associated with  ${f x}$
- we use **bold face** to indicate a vector (e.g,  $\mathbf{x}$ )

- We use superscript (i) to index examples, when we have more than one
  - $\mathbf{x^{(i)}}, \mathbf{x}^{(i')}, i \neq i'$  are two distinct examples
- $\bullet \ \, \text{denote an element } i \text{ of a collection of } m \text{ examples (e.g., } \mathbf{x^{(i)}}) \\ \bullet \ \, \text{We use subscript } j \text{ to index element } j \text{ of a vector, e.g., } \mathbf{x}^{(i)}_j \\$

• So  $\mathbf{x}^{(i)}$  is

$$\mathbf{x^{(i)}} = egin{pmatrix} \mathbf{x}_1^{(i)} \ \mathbf{x}_2^{(i)} \ dots \ \mathbf{x}_n^{(i)} \end{pmatrix}$$

Each element of  $\mathbf{x^{(i)}}$  is a "feature"

•  $\mathbf{x}_{j}^{(\mathbf{i})}$  is the  $j^{th}$  feature of example i

### Training set

• The collection of examples used for fitting (training) a model is called the *training* set:

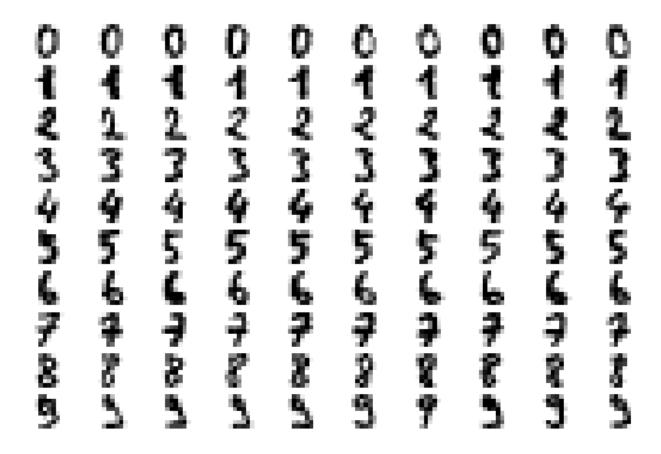
$$\langle \mathbf{X}, \mathbf{y} \rangle = [\mathbf{x^{(i)}}, \mathbf{y^{(i)}} | 1 \le i \le m]$$

where m is the size of training set and each  $\mathbf{x^{(i)}}$  is a feature vector of length n.

- By seeing many (m) pairs of feature vectors and associated labels we will try to infer the correct label  $\mathbf{y^{(i)}}$  from the features in  $\mathbf{x^{(i)}}$
- ${f X}$  is an (m imes n) matrix and  ${f y}$  is an (m imes 1) vector of targets.

$$\mathbf{X} = egin{pmatrix} (\mathbf{x}^{(1)})^T \ (\mathbf{x}^{(2)})^T \ dots \ (\mathbf{x}^{(m)})^T \end{pmatrix} = egin{pmatrix} \mathbf{x}_1^{(1)} \dots \mathbf{x}_n^{(1)} \ \mathbf{x}_1^{(2)} \dots \mathbf{x}_n^{(2)} \ dots \ \mathbf{x}_1^{(m)} \dots \mathbf{x}_n^{(m)} \end{pmatrix}$$

#### **Training set**



• We will sometimes add a "constant" feature by setting  $\mathbf{x}_0^{(i)} = 1, 0 \le i \le m$  so that the first column of  $\mathbf{x}$  is 1:

$$\mathbf{X} = egin{pmatrix} 1 & \mathbf{x}_1^{(1)} & \dots & \mathbf{x}_n^{(1)} \ 1 & \mathbf{x}_1^{(2)} & \dots & \mathbf{x}_n^{(2)} \ dots & dots & \dots & dots \ 1 & \mathbf{x}_1^{(m)} & \dots & \mathbf{x}_n^{(m)} \end{pmatrix}$$

ullet So each of the m rows is an example and each of the n columns is a feature.

## Not just numbers!

The features aren't restricted to be numeric!

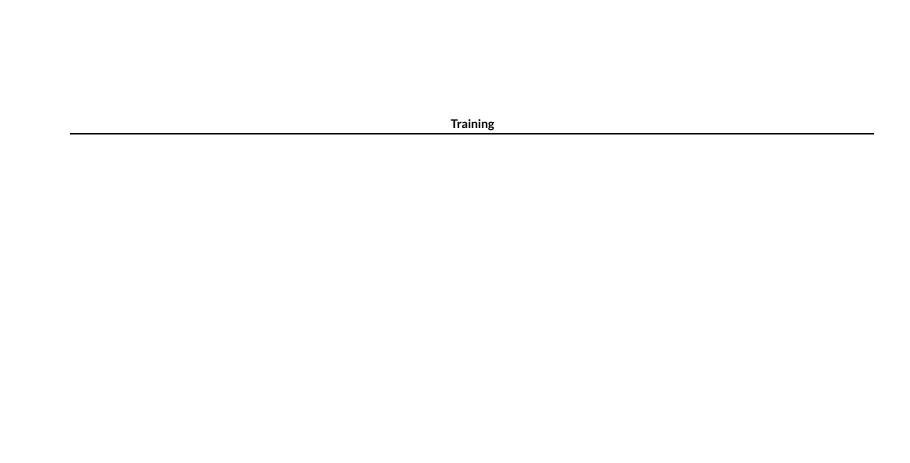
In this course, we will deal with data that is

- numeric
- categorical
- text
- image
- sound (not this course)

Of course, you'll have to encode this data as numbers in order for numerical algorithms to handle them.

### **Prediction**

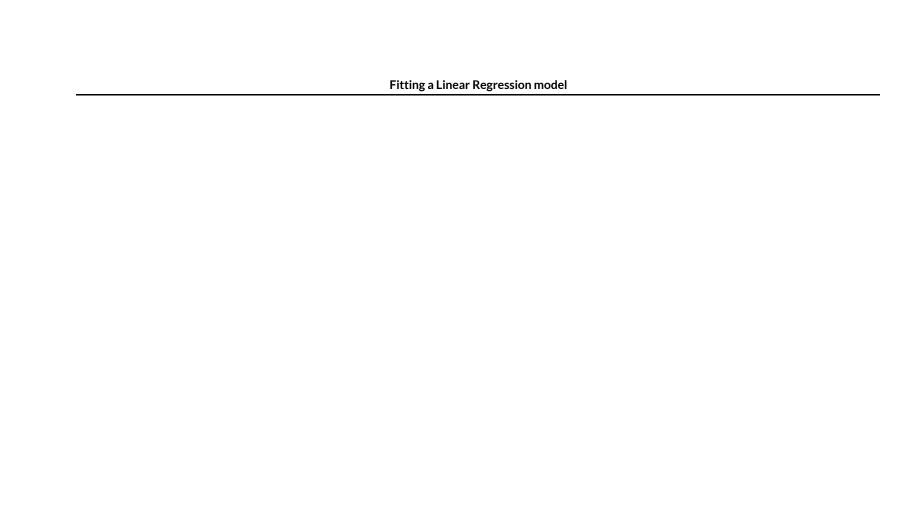
- Given training example  $\mathbf{x^{(i)}}$ , we construct a function h to predict its label  $\hat{\mathbf{y}^{(i)}} = h(\mathbf{x^{(i)}})$
- We use a "hat" to denote predictions:  $\hat{\mathbf{y}}^{(i)}$
- The function h will often be parameterized (by  $\Theta$ ) so, for clarity, we should write  $\hat{\mathbf{y}}^{(\mathbf{i})} = h(\mathbf{x}^{(\mathbf{i})}; \Theta)$
- We will often drop  $\Theta$  for ease of reading.
- Since h is a function, it should also be possible to make a prediction for a vector  $\mathbf{x}$  that is **not** part of the training set.
- That is, we are able *generalize* to non-training examples: to make out of sample predictions



The key task of Machine Learning is finding the "best" values for parameters  $\Theta$ .

The process of using training examples  ${f X}$  to find  ${f \Theta}$ 

- is called *fitting* the model
- is solved as an optimization problem (to be described)



#### Summary

- a training example is a pair  $(\mathbf{x^{(i)}}, \mathbf{y^{(i)}})$  drawn from training set  $\langle \mathbf{X}, \mathbf{y} \rangle$  consisting of
  - a feature vector  $\mathbf{x}^{(i)}$  of length n
  - lacktriangle the associated label (target)  $\mathbf{y^{(i)}}$
  - **X** is of dimension  $m \times n$
  - ullet  ${f y}$  is dimension m imes 1, i.e., target is a single, continuous value per example
- predictions are indicated with a "hat:
  - $\quad \ \ \, \hat{y}^{(i)}$  is the prediction made given  $x^{(i)}$  as input

# Loss/Cost, Utility

- The prediction  $\hat{\mathbf{y}}^{(i)}$  for example  $\mathbf{x^{(i)}}$  is perfect if it matches the true label  $\mathbf{y^{(i)}}$   $\hat{\mathbf{y}}^{(i)} = \mathbf{y^{(i)}}$
- Perfection is hard (at least at first) so we need a measure for "how far off" the prediction is.
- We will call the distance between  $\hat{\mathbf{y}}^{(i)}, \mathbf{y}^{(i)}$  the Loss (or Cost) for example i:

$$\mathcal{L}_{\Theta}^{(\mathbf{i})} = L(\ h(\mathbf{x^{(i)}}; \Theta), \mathbf{y^{(i)}}\ ) = L(\hat{\mathbf{y}^{(i)}}, \mathbf{y^{(i)}})$$

where L(a,b) is a function that is 0 when a=b and increasing as a increasingly differs from b.

Two common forms of  ${\cal L}$  are Mean Squared Error (for Regression) and Cross Entropy Loss (for classification).

The Loss for the entire training set is simply the average (across examples) of the Loss for the example

$$\mathcal{L}_{\Theta} = rac{1}{m} \sum_{i=1}^{m} \mathcal{L}_{\Theta}^{(\mathbf{i})}$$



Whereas Loss describes how "bad" our prediction is, we sometimes refer to the converse -- how "good" the prediction is.

We call the "goodness" of the prediction the *Utility*  $U_{\Theta}$ .

So we could state the optimization objective either as "minimize Cost" or "maximize Utility".

By convention, the DL optimization problem is usually framed as one of minimization (of cost or loss) rather than maximization of utility.

Since Cost is inversely related to Utility, you will sometimes see the minimization objective written as "minimize -1 times Utility".
So be forewarned that you will often see Loss function with leading "negation" signs.

## Creating Loss functions is a key part of Deep Learning

As you will come to see, particularly for Deep Learning, the essence of many problems is in creating a Loss Function that captures the objective of your problem.

This is far from a trivial part of the process.

# Fitting/Training a Model

The best (optimal)  $\Theta$  is the one that minimizes the Average (across training examples) Loss

$$\Theta^* = \operatorname*{argmin}_{\Theta} \mathcal{L}_{\Theta}$$

•	The goal of fitting/training is to solve for the $\Theta$ that minimizes the training set le	oss
	$L_{\Theta}$	

 $\bullet \;$  The method for finding  $\Theta$  is called optimization.

# The dot product: Template matching

- The "dot product" (special case of inner product) is one function that often appears in template matching
- It measures the similarity of two vectors

$$\mathbf{v}\cdot\mathbf{v}'=\sum_{i=1}^n\mathbf{v}_i\mathbf{v}_i'$$

• As a similarity measure (rather than as a distance) high dot product means "more similar".

- There are several intuitions for the dot product
- The dot product is maximized when large (resp., small) values appear in similar positions in both vectors
  - this becomes even more obvious if we 0-center both vectors such that "small" values become negative
  - this looks like the statistical formula for covariance
    - if we normalize both vectors to unit length, then this looks like correlation

We can generalize dot product to higher dimensions

- Compute pair-wise product of corresponding entries
- Reduce to a scalar by summing

```
In [2]: print("Done")
```

Done