# Case 1

# Consumption Prediction of a detached

# Residential building
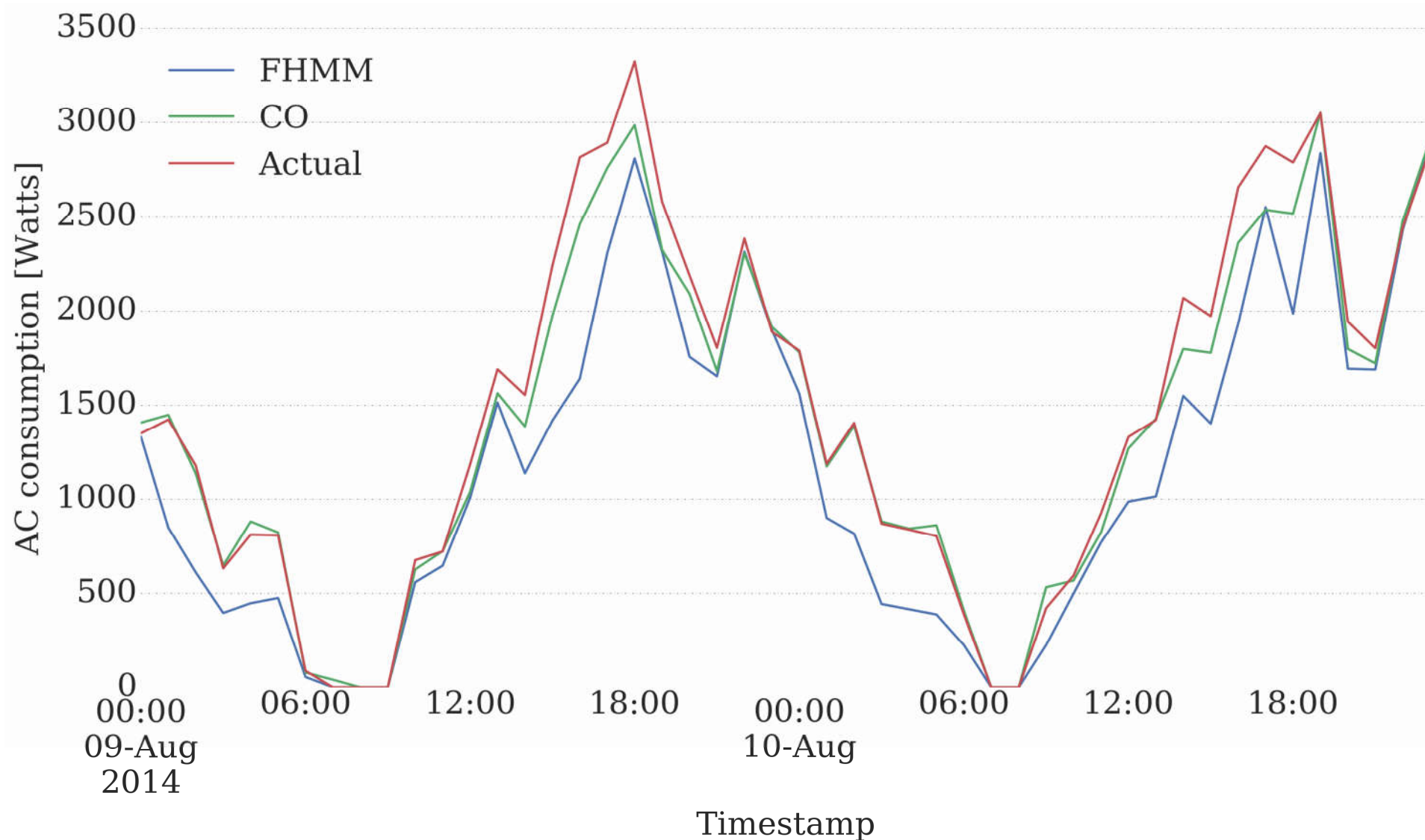
❖ Location: Austin, Texas

❖ Time-stamped data including:

✓ Aggregate consumption

✓ Appliance-by-appliance consumptions including Air conditioner

✓ Ambient temperature

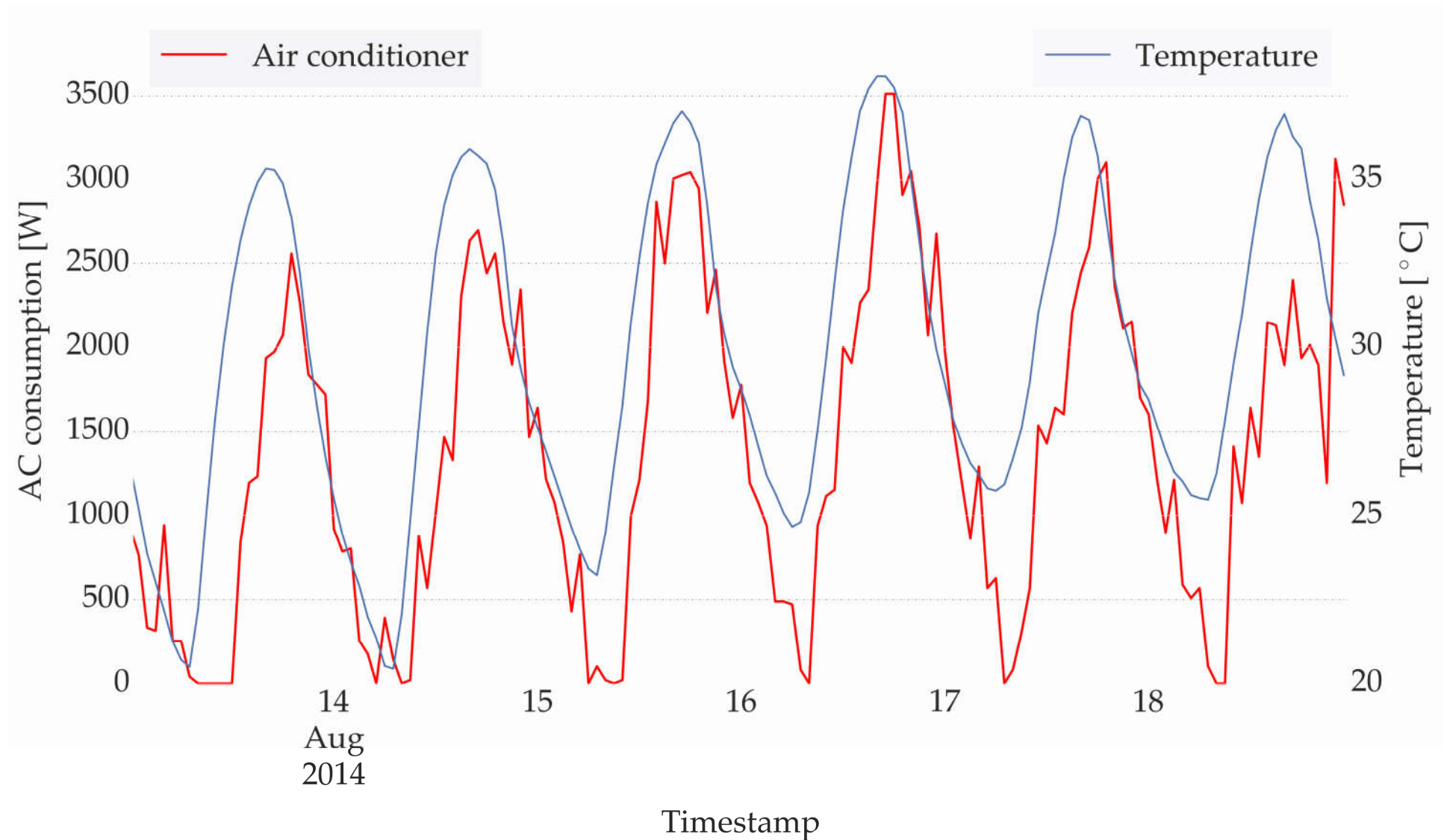➤ Time-stamped PV generation from buildings: represents solar irradiation

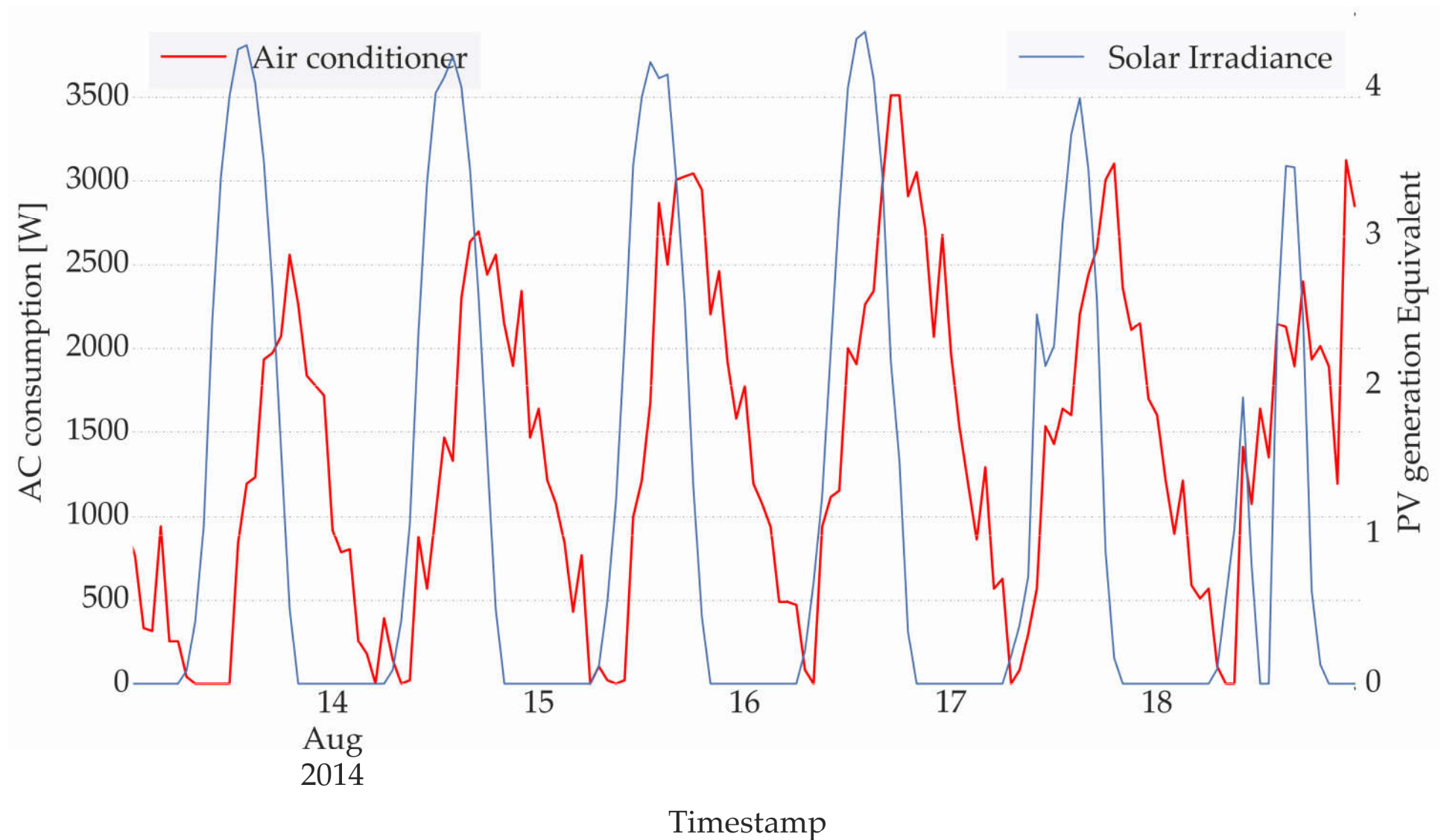# Air conditioner disaggregation results

Independent Variables

Dependent Variable

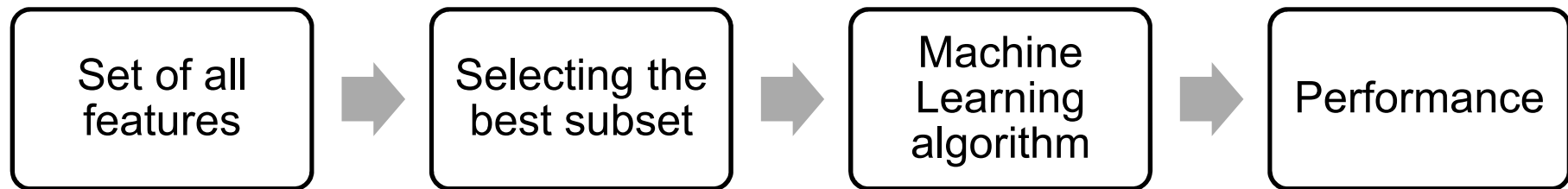| S.No | Date | Time | Hour Of day (0-23) | Day or Night (1/0) | Temperature (°C) | Irradiance (0-1) | Relative Humidity (0-100) | Energy consumed by AC (W) |
|------|------|------|--------------------|--------------------|------------------|------------------|---------------------------|---------------------------|
| 1 | 01-01-2017 | 14:00:00 | 14 | 1 | 44 | 0.83 | 60 | 2400 |
| 2 | 01-01-2017 | 15:00:00 | 15 | 1 | 44.5 | 0.85 | 60 | 2500 |
| 3 | … | … | … | … | … | … | … | … |

$$Energy\ consumed(t) = F(Date \quad time(t), Hour(t), Temperature(t), IRR(t), RH(t))$$

Question:
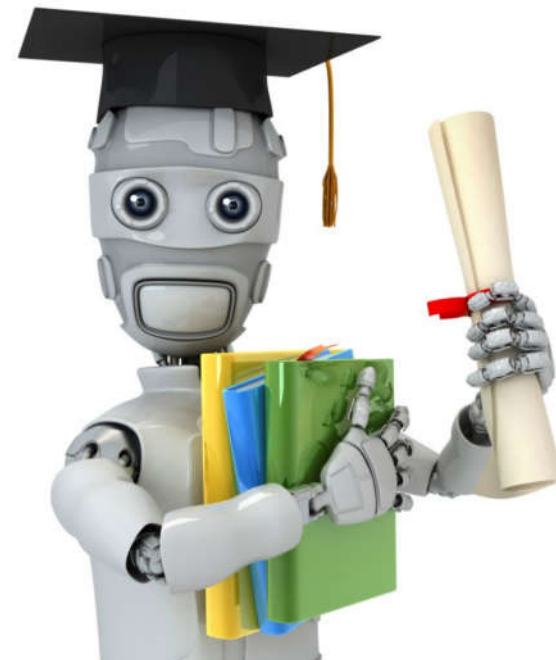What other independent variables affect Energy consumed by an Air conditioner ?

| Set of all features | → | Selecting the best subset | → | Machine Learning algorithm | → | Performance |
|---|---|---|---|---|---|---|

1. Identifying set of all features >> Domain knowledge
2. Selecting the best subset      >> Correlation (Person, Spearman)
3. Learning algorithm             >> Performance of the model

# A Rapid Introduction to Machine Learning*



*Important Note: This introduction is a simplified presentation directly extracted from Andrew Ng's Stanford ML course

❖ **Simple informal Definition (Arthur Samuel)**

✓ the field of study that gives computers the ability to learn without being explicitly programmed

❖ **A more accurate definition (Tom Mitchel)**

✓ A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.

❑ **Example: playing checkers.**

✓ E = the experience of playing many games of checkers, T = the task of playing checkers.

✓ P = the probability that the program will win the next game.

❖ **Main Machine Learning Categories**

✓ Supervised Learning

✓ Unsupervised Learning

❖ **Supervised Learning**

✓ in supervised learning, we are given a data set and already know what our correct output should look like, having the idea that there is a relationship between the input and the output. Supervised learning problems are categorized into :

➢ **Regression**

✓ In a regression problem, we are trying to predict results within a continuous output, meaning that we are trying to map input variables to some continuous function.

❑ Example: Given data about the size of houses on the real estate market, try to predict their price. Price as a function of size is a continuous output

➢ **Classification**

✓ In a classification problem, we are instead trying to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories.

❑ Example: given a picture of tumour we would like to decide whether it is malignant or benign
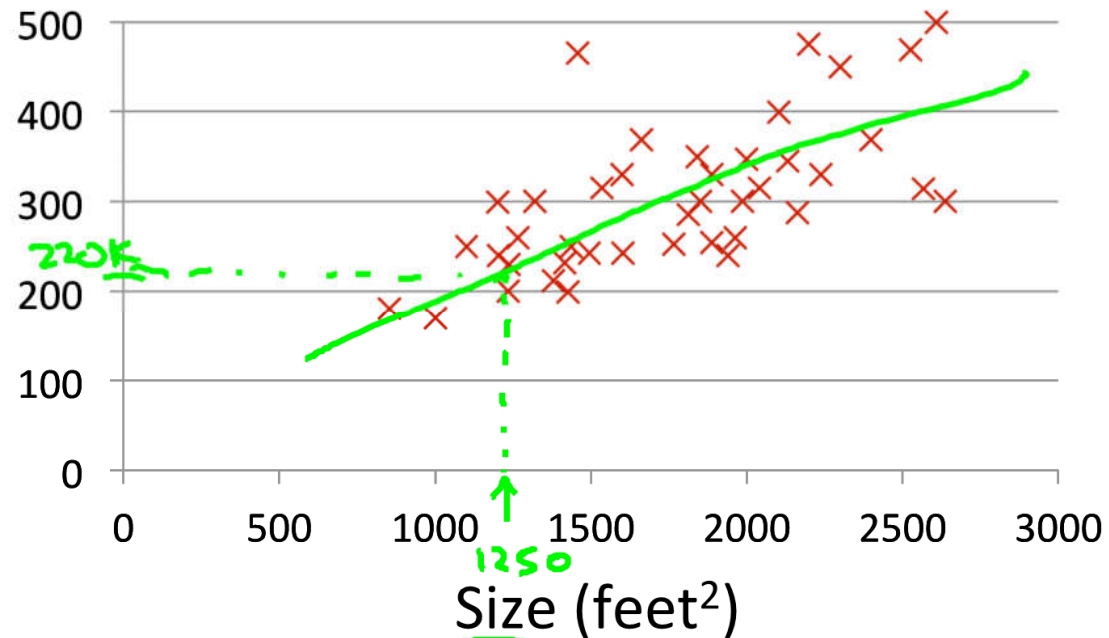
❖ **Unsupervised learning**

✓ Unsupervised learning, allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

✓ We can derive this structure by clustering the data based on relationships among the variables in the data.

✓ With unsupervised learning there is no feedback based on the prediction results, i.e., there is no teacher to correct you.

❑ Example: (Clustering) Take a collection of 1000 essays written on the US Economy, and find a way to automatically group these essays into a small number that are somehow similar or related by different variables, such as word frequency, sentence length, page count, and so on.

**Housing Prices
(Portland, OR)**

Price
(in 1000s
of dollars)



Size (feet$^2$)

220K

1250

**Supervised Learning**

Given the "right answer" for
each example in the data.

**Regression Problem**

Predict real-valued output

Classification: Discrete-valued output

*Source: Andrew Ng's Machine Learning Course

**Training set of housing prices (Portland, OR)**

| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

*m = 47*

Notation:

**m** = Number of training examples

**x**'s = "input" variable / features

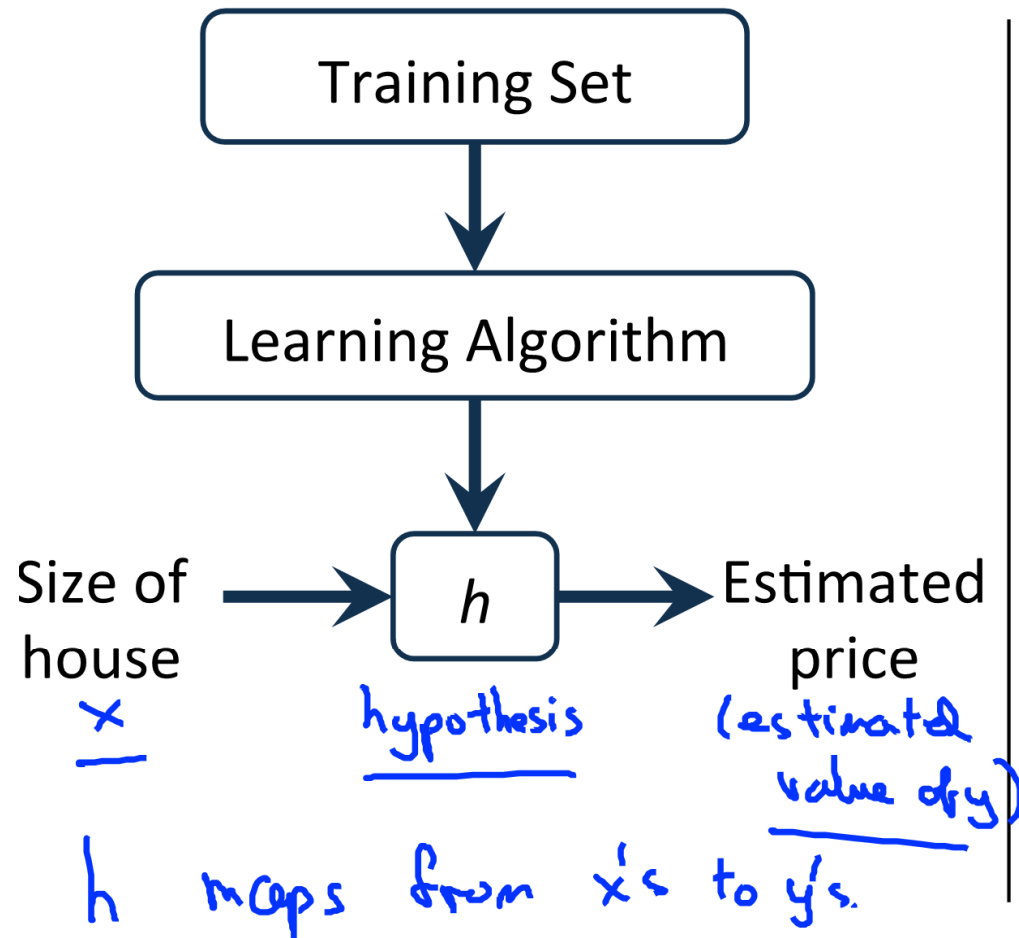**y**'s = "output" variable / "target" variable

$(x, y)$ – one training example

$(x^{(i)}, y^{(i)})$ – $i^{th}$ training example

$x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$

Training Set

↓

Learning Algorithm

↓

Size of house → $h$ → Estimated price

$x$

hypothesis

$h$ maps from $x$'s to $y$'s.

(estimated value of $y$)

**How do we represent $h$ ?**

$$h_\theta(x) = \theta_0 + \theta_1 x$$

Shorthand: $h(x)$

$h(x) = \theta_0 + \theta_1 x$

Linear regression with one variable. $(x)$
Univariate linear regression.

└ one variable

*Source: Andrew Ng's Machine Learning Course

Training Set

| Size in feet² (x) | Price ($) in 1000's (y) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$M = 47$

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s: <u>Parameters</u>

How to choose $\theta_i$'s ?
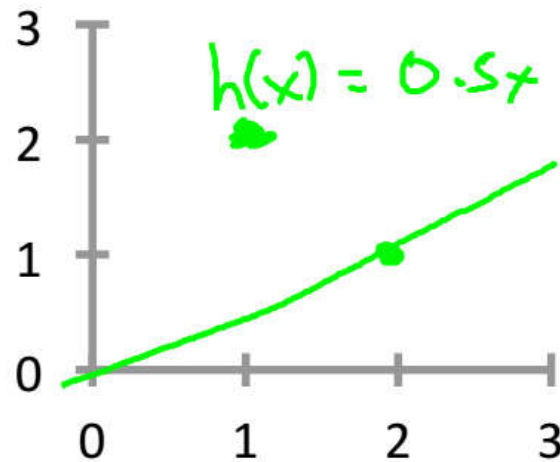
*Source: Andrew Ng's Machine Learning Course

$$h_\theta(x) = \theta_0 + \theta_1 x$$



$h(x) = 1.5 + 0 \cdot x$

$\theta_0 = 1.5$
$\theta_1 = 0$

$h(x) = 0.5x$

$\theta_0 = 0$
$\theta_1 = 0.5$

$h_\theta(x)$

$h(x)$

$\theta_0 = 1$
$\theta_1 = 0.5$

**\*Source: Andrew Ng's Machine Learning Course**

Training Set

| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$M = 47$

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s: Parameters

How to choose $\theta_i$'s ?

*Source: Andrew Ng's Machine Learning Course

Idea: Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

$x, y$

$(x^{(i)}, y^{(i)})$

$\theta_0, \theta_1$

#training examples

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

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

Cost function

Squared error function

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)



Price ($) in 1000's

Size in feet$^2$ (x)

$\theta_0 = 50$

$\theta_1 = 0.06$

$$h_\theta(x) = 50 + 0.06x$$

$\theta_1$

$\theta_0, \theta_1$

*Source: Andrew Ng's Machine Learning Course

Contour plots
Contour figures -

$J(\theta_0, \theta_1)$

$\theta_1$

$\theta_0$

$J(\theta_0, \theta_1)$

*Source: Andrew Ng's Machine Learning Course

Have some function $J(\theta_0, \theta_1)$   $J(\theta_0, \theta_1, \theta_2, \ldots, \theta_n)$

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

**Outline:**

- Start with some $\theta_0, \theta_1$   $(\text{say } \theta_0 = 0, \theta_1 = 0)$

- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

  until we hopefully end up at a minimum

*Source: Andrew Ng's Machine Learning Course

## Gradient descent algorithm

repeat until convergence {

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

(simultaneously update $j = 0$ and $j = 1$)

}

learning rate

derivative

$$\min_{\theta_1} J(\theta_1)$$

$$\theta_1 \in \mathbb{R}.$$

**\*Source: Andrew Ng's Machine Learning Course**

## Gradient descent algorithm

repeat until convergence {

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

(for $j = 1$ and $j = 0$)

}

## Linear Regression Model

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

**\*Source: Andrew Ng's Machine Learning Course**

**Gradient descent algorithm**

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

}

update $\theta_0$ and $\theta_1$ simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

**\*Source: Andrew Ng's Machine Learning Course**

*Source: Andrew Ng's Machine Learning Course

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

*Source: Andrew Ng's Machine Learning Course

## Multiple features (variables).

| Size (feet$^2$) $x_1$ | Number of bedrooms $x_2$ | Number of floors $x_3$ | Age of home (years) $x_4$ | Price ($1000) $y$ |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$m = 47$

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$x_3^{(2)} = 2$

Notation:

$n$ = number of features    $n = 4$

$x^{(i)}$ = input (features) of $i^{th}$ training example.

$x_j^{(i)}$ = value of feature $j$ in $i^{th}$ training example.

*Source: Andrew Ng's Machine Learning Course

Hypothesis:

Previously: $h_\theta(x) = \theta_0 + \theta_1 x$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

e.g. $h_\theta(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$

age

*Source: Andrew Ng's Machine Learning Course

Hypothesis:
$$h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

$\rightarrow x_0 = 1$

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$ $\theta$

$n+1$ - dimensional vector

Cost function:
$$J(\theta_0, \theta_1, \ldots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$J(\theta)$

Gradient descent:

Repeat {
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n) \quad J(\theta)$$
}

(simultaneously update for every $j = 0, \ldots, n$)

Hypothesis: $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

$\rightarrow x_0 = 1$

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$

$n+1 - $ dimensional vector

Cost function:

$$J(\theta_0, \theta_1, \ldots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$
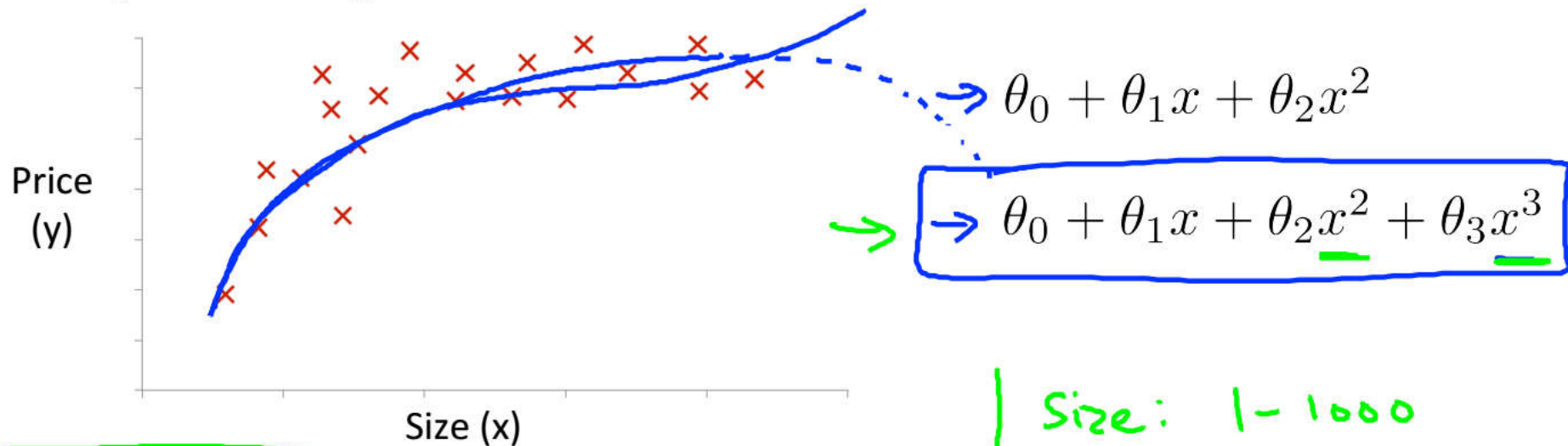
$J(\theta)$

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n) \quad J(\theta)$$

}

(simultaneously update for every $j = 0, \ldots, n$)

## Polynomial regression



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$
$$= \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$$

$x_1 = (size)$
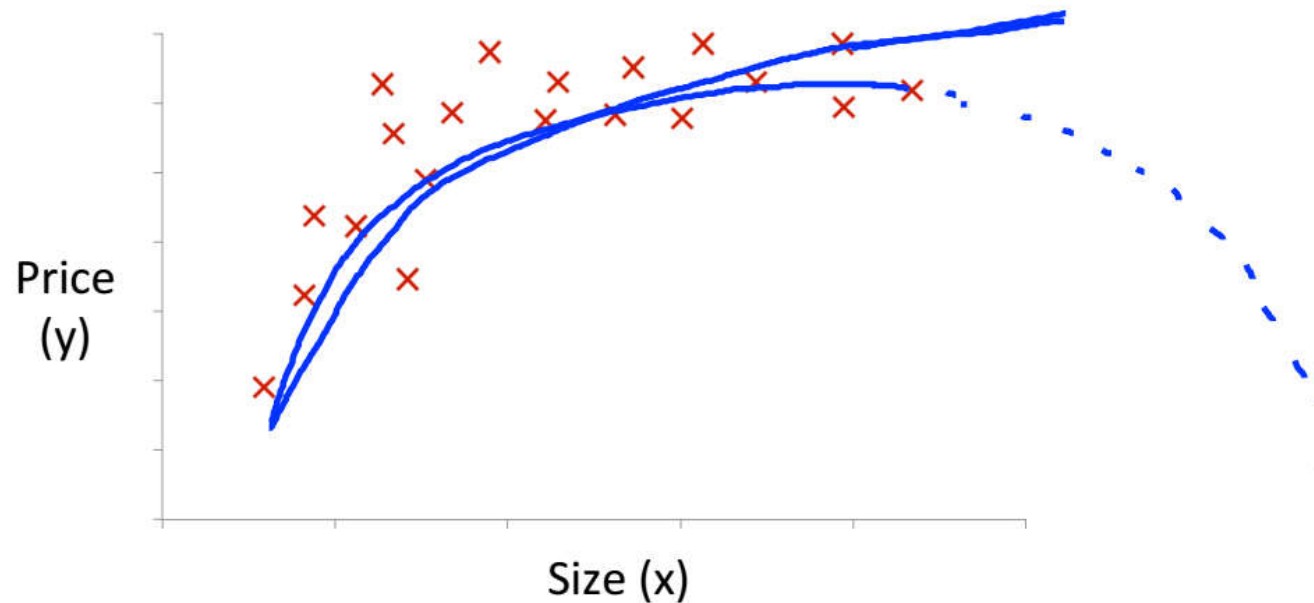$x_2 = (size)^2$
$x_3 = (size)^3$

Size: 1 - 1000

Size$^2$: 1 - 1000,000

Size$^3$: 1 - 10$^9$

*Source: Andrew Ng's Machine Learning Course

## Choice of features



$$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2(size)^2$$

$$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2\sqrt{(size)}$$

**\*Source: Andrew Ng's Machine Learning Course**

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2$$
$$+\theta_3 x_1 x_2 + \theta_4 x_1^2 x_2$$
$$+\theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \ldots)$$

$$x_1 = \text{size}$$
$$x_2 = \text{\# bedrooms}$$
$$x_3 = \text{\# floors}$$
$$x_4 = \text{age}$$
$$\ldots$$
$$x_{100} -$$

$$\to x_1^2, \; x_1 x_2, \; x_1 x_3, \; x_1 x_4 \ldots x_1 x_{100}$$
$$x_2^2, \; x_2 x_3 \ldots$$

$$\approx 5000 \; \text{feature} \qquad O(n^2)$$

$$n = 100$$

$$\to x_1^2, x_2^2, x_3^2, \ldots, x_{100}^2 \qquad \approx \frac{n^2}{2}$$

$$\to x_1 x_2 x_3, \; x_1^2 x_2, \; x_{10} x_{11} x_{17}, \ldots$$
$$O(n^3)$$
$$170,000$$

*Source: Andrew Ng's Machine Learning Course
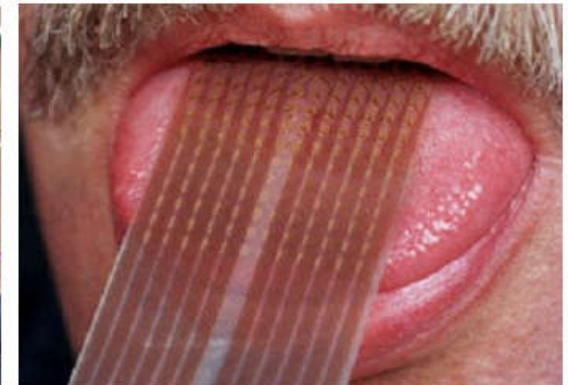
❖ **First Solution**

✓ Let's have a look at what happens in the nature! → Human brain !

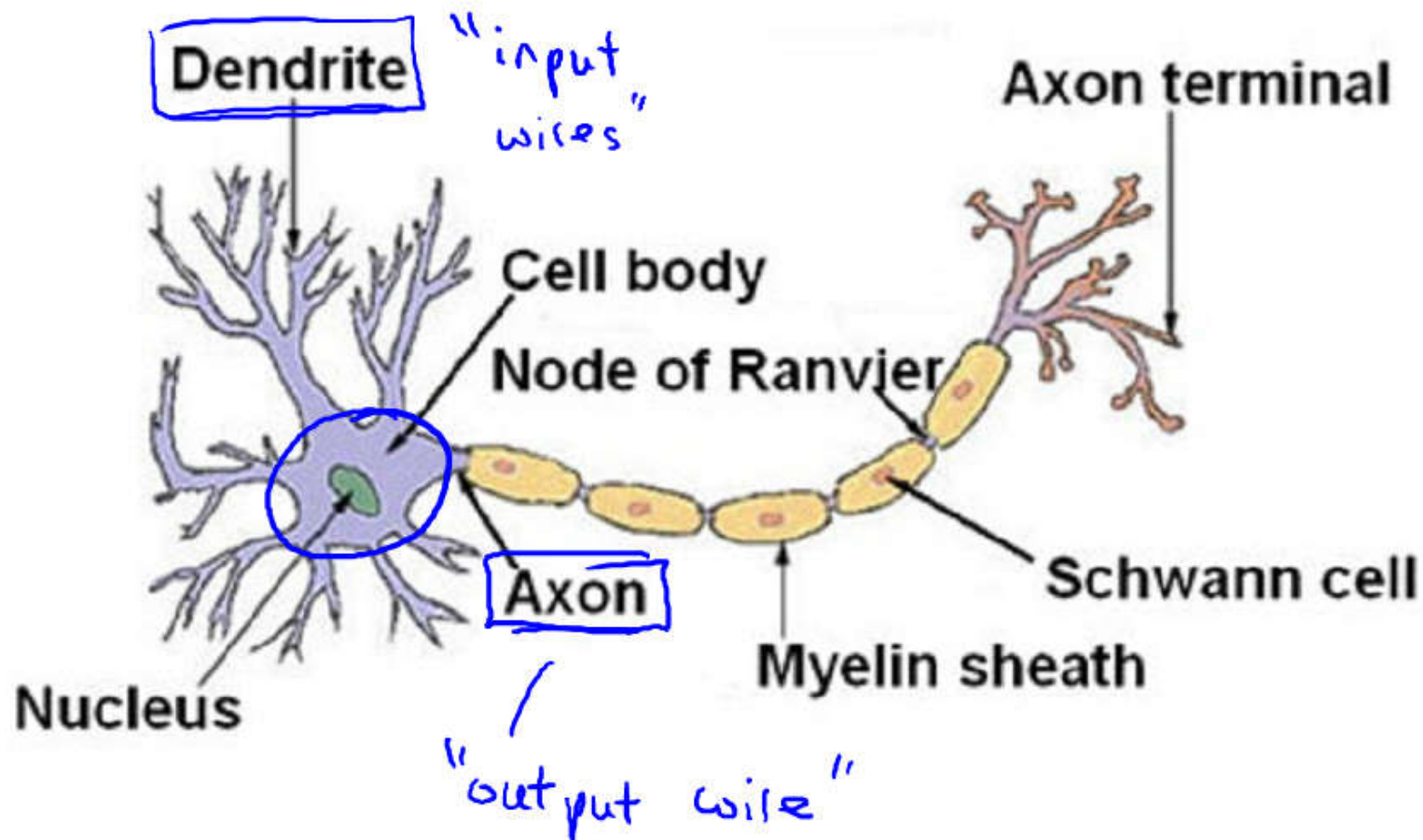**The "one learning algorithm" hypothesis**



Somatosensory Cortex

Seeing with your tongue

**\*Source: Andrew Ng's Machine Learning Course**

*Source: Andrew Ng's Machine Learning Course

## Neuron model: Logistic unit

"bias unit"

$x_0 = 1$

$x_0$

$x_1$

$x_2$

$x_3$

"input wires"

"output"

$h_\theta(x)$

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

"weights"

( parameters )

## Sigmoid (logistic) activation function.

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

*Source: Andrew Ng's Machine Learning Course

*Source: Andrew Ng's Machine Learning Course

## Neural Network



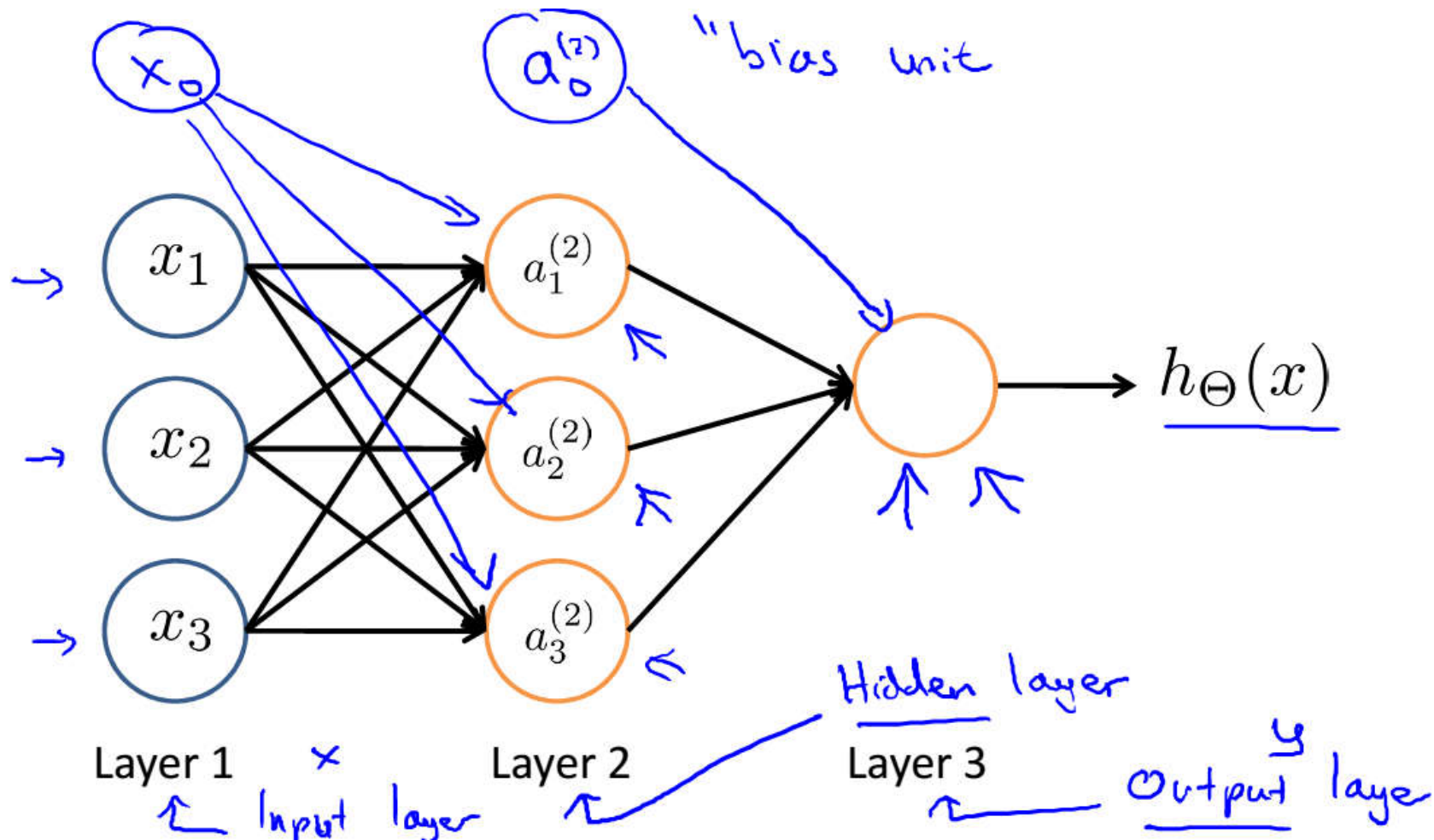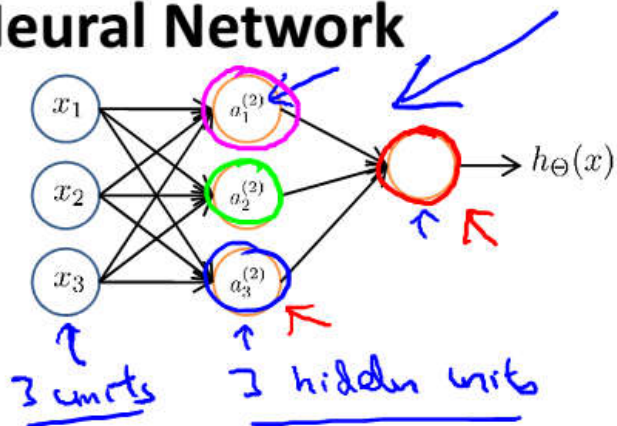$a_i^{(j)}$ = "activation" of unit $i$ in layer $j$

$\Theta^{(j)}$ = matrix of weights controlling function mapping from layer $j$ to layer $j+1$

$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$

3 units    3 hidden units

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

If network has $s_j$ units in layer $j$, $s_{j+1}$ units in layer $j+1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.
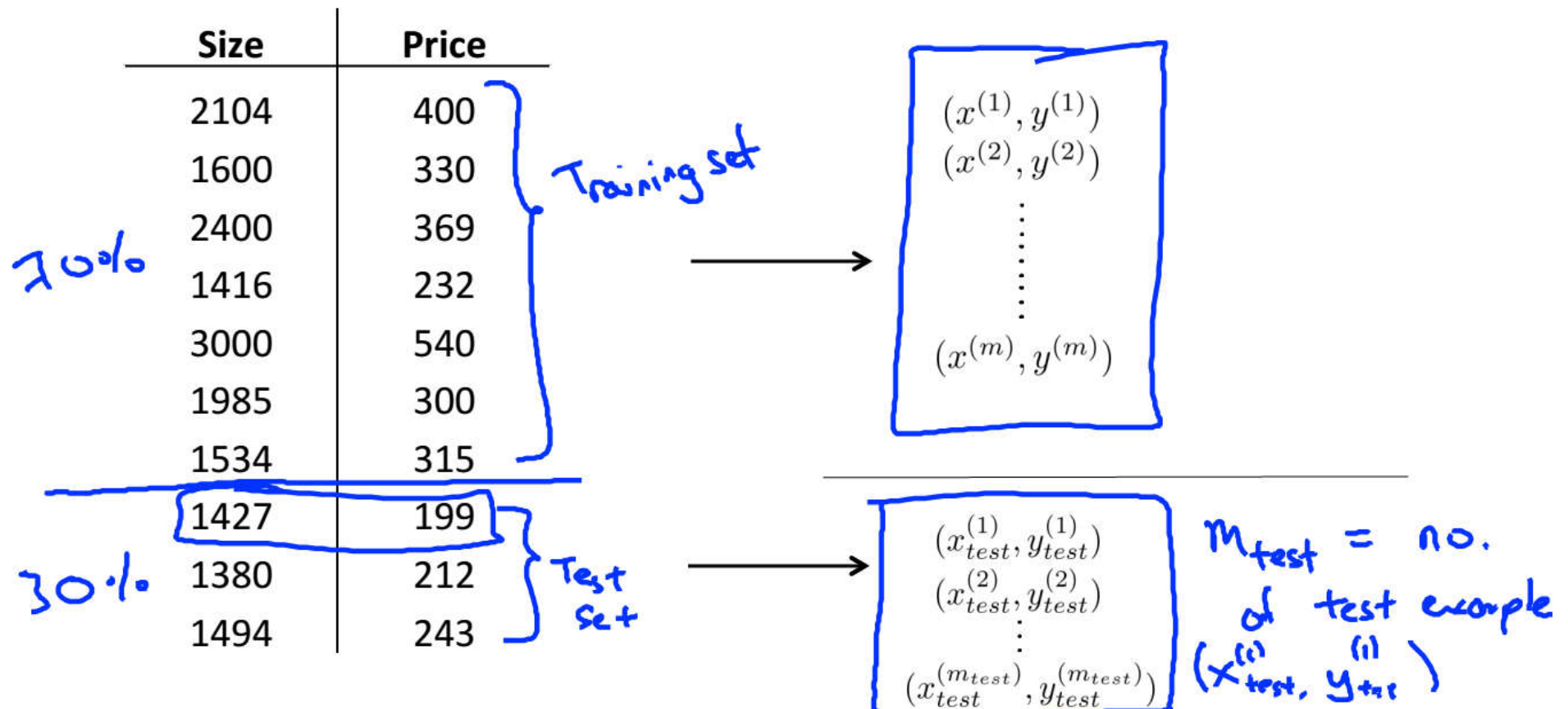
$s_{j+1} \times (s_j + 1)$

**\*Source: Andrew Ng's Machine Learning Course**

## Evaluating your hypothesis

Dataset:

| Size | Price |
|------|-------|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| 1985 | 300 |
| 1534 | 315 |
| 1427 | 199 |
| 1380 | 212 |
| 1494 | 243 |

*70%*  *Training set*

*30%*  *Test Set*

$$(x^{(1)}, y^{(1)})$$
$$(x^{(2)}, y^{(2)})$$
$$\vdots$$
$$(x^{(m)}, y^{(m)})$$

$$(x_{test}^{(1)}, y_{test}^{(1)})$$
$$(x_{test}^{(2)}, y_{test}^{(2)})$$
$$\vdots$$
$$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$$

*$m_{test}$ = no. of test example*
*$(x_{test}^{(i)}, y_{test}^{(i)})$*

*Source: Andrew Ng's Machine Learning Course

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2.$$

$$\text{MAE} = \frac{\sum_{i=1}^{n} |y_i - x_i|}{n} = \frac{\sum_{i=1}^{n} |e_i|}{n}.[1]$$

A data set has $n$ values marked $y_1,...,y_n$ (collectively known as $y_i$ or as a vector $y = [y_1,...,y_n]^T$), each associated with a predicted (or modeled) value $f_1,...,f_n$ (known as $f_i$, or sometimes $\hat{y}_i$, as a vector $f$).

If $\bar{y}$ is the mean of the observed data:

$$\bar{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$$

then the variability of the data set can be measured using three sums of squares formulas:

- The total sum of squares (proportional to the variance of the data):

$$SS_{tot} = \sum_i (y_i - \bar{y})^2,$$

- The regression sum of squares, also called the explained sum of squares:

$$SS_{reg} = \sum_i (f_i - \bar{y})^2,$$

- The sum of squares of residuals, also called the residual sum of squares:

$$SS_{res} = \sum_i (y_i - f_i)^2 = \sum_i e_i^2$$

The most general definition of the coefficient of determination is

$$R^2 \equiv 1 - \frac{SS_{res}}{SS_{tot}}.$$