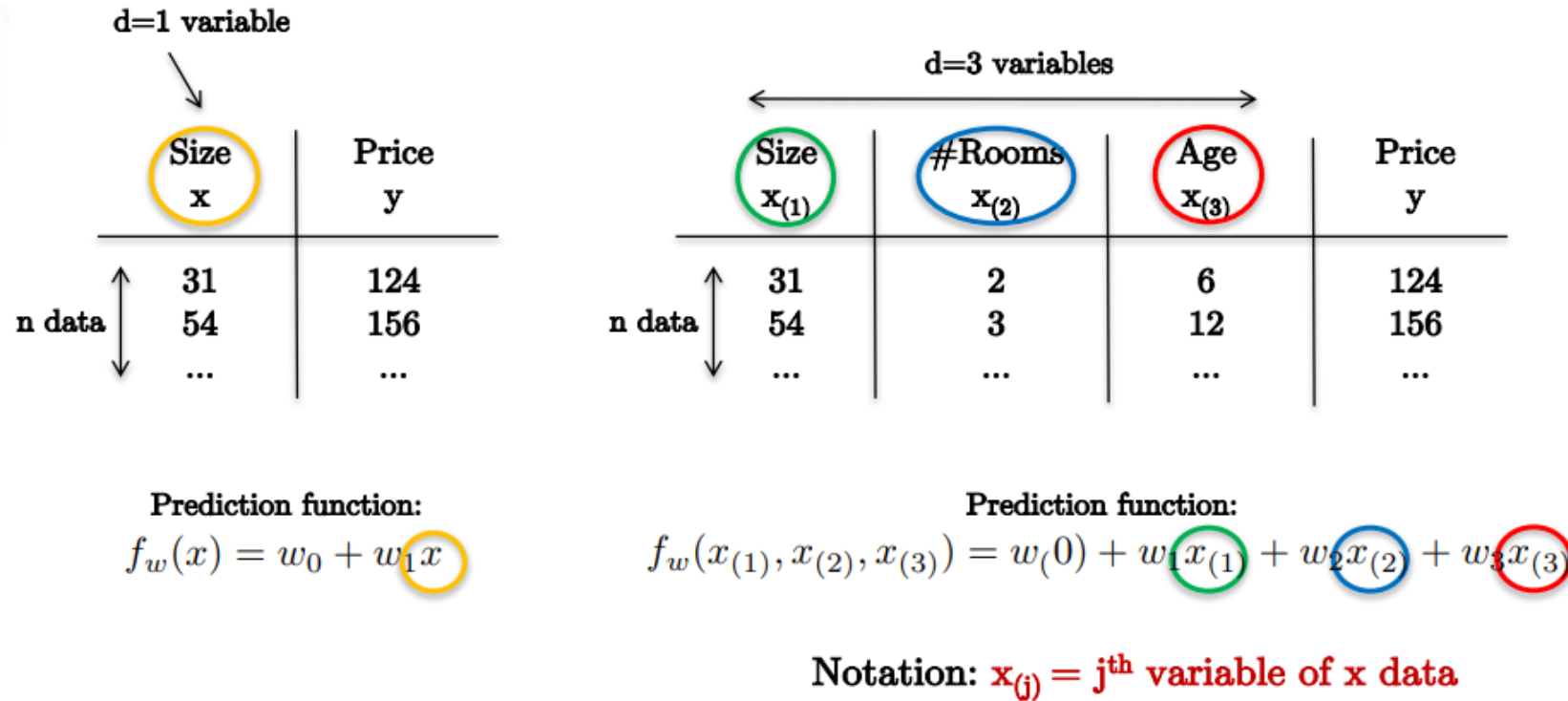# 인공지능과 수학적 배경

# Linear Regression
## Multiple variables

◆ **Data may have K even M features/attributes** (curse of dimensionality).

d=1 variable

| Size x | Price y |
|--------|---------|
| 31 | 124 |
| 54 | 156 |
| ... | ... |

n data

d=3 variables

| Size $x_{(1)}$ | #Rooms $x_{(2)}$ | Age $x_{(3)}$ | Price y |
|-----------|-----------|----------|---------|
| 31 | 2 | 6 | 124 |
| 54 | 3 | 12 | 156 |
| ... | ... | ... | ... |

n data

Prediction function:

$$f_w(x) = w_0 + w_1 x$$

Prediction function:

$$f_w(x_{(1)}, x_{(2)}, x_{(3)}) = w_{(0)} + w_1 x_{(1)} + w_2 x_{(2)} + w_3 x_{(3)}$$

Notation: $x_{(j)} = j^{th}$ variable of x data

# Linear Regression
## Data matrix $X$

◆ **Notation:**

- ▪ $n$ = number of training data

- ▪ $d$ = number of variables/features, dim(x)

- ▪ $x_i$ = i$^{th}$ training data, size($x_i$) = $d$ x 1

- ▪ $x_{(j)}$ = j$^{th}$ feature of training data x, scalar

- ▪ X = data matrix, size(X)= $n$ x $d$

- ▪ $X_{ij}$ = $x_{i(j)}$ = j$^{th}$ variable of i$^{th}$ training data

◆ **Example:**

$$x_1 = \begin{bmatrix} 2 \\ 1 \\ 7 \end{bmatrix} \updownarrow \text{d variables}$$

$d$ x 1

$$x_1^T = \begin{bmatrix} 2 & 1 & 7 \end{bmatrix}$$

1 x $d$

$$\xleftarrow{\hspace{1cm}} \text{d variables} \xrightarrow{\hspace{1cm}}$$

$$X = \begin{bmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{bmatrix} \updownarrow \text{n data}$$

Data matrix

$n$ x $d$

All training data in a single matrix
(very efficient for linear algebra computations)

# Linear Regression
## Prediction function with $d$ variables

◆ **$d=1$ variable:** $\qquad f_w(x) = w_0 + w_1 x$

◆ **$d$ variables:** $\qquad f_w(x_{(1)}, ..., x_{(d)}) = w_{(0)} + w_1 x_{(1)} + ... + w_d x_{(d)}$

◆ **Example: $d=2$** $\qquad f_w(x_{(1)}, x_{(2)}) = w_{(0)} + w_1 x_{(1)} + w_d x_{(2)}$



$$f_w(x_{(1)}, x_{(2)}) = -0.3 + 1.2 x_{(1)} + 9.3 x_{(2)}$$

# Linear Regression
## Vector representation

◆ **Define the vectors:**

$$x = \begin{bmatrix} 1 \\ x_{(1)} \\ x_{(2)} \\ \vdots \\ x_{(d)} \end{bmatrix} \qquad w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$(d+1)$ x 1 $\qquad\qquad$ $(d+1)$ x 1

◆ **Re-write the prediction function (as vector-vector multiplication):**

$$f_w(x) = w_0.1 + w_1 x_{(1)} + ... + w_d x_{(d)}$$

$$= w^T x \quad \text{One line of code}$$

1 x $(d+1)$ $\qquad$ $(d+1)$ x 1

$f$ is called multivariate linear regression function.

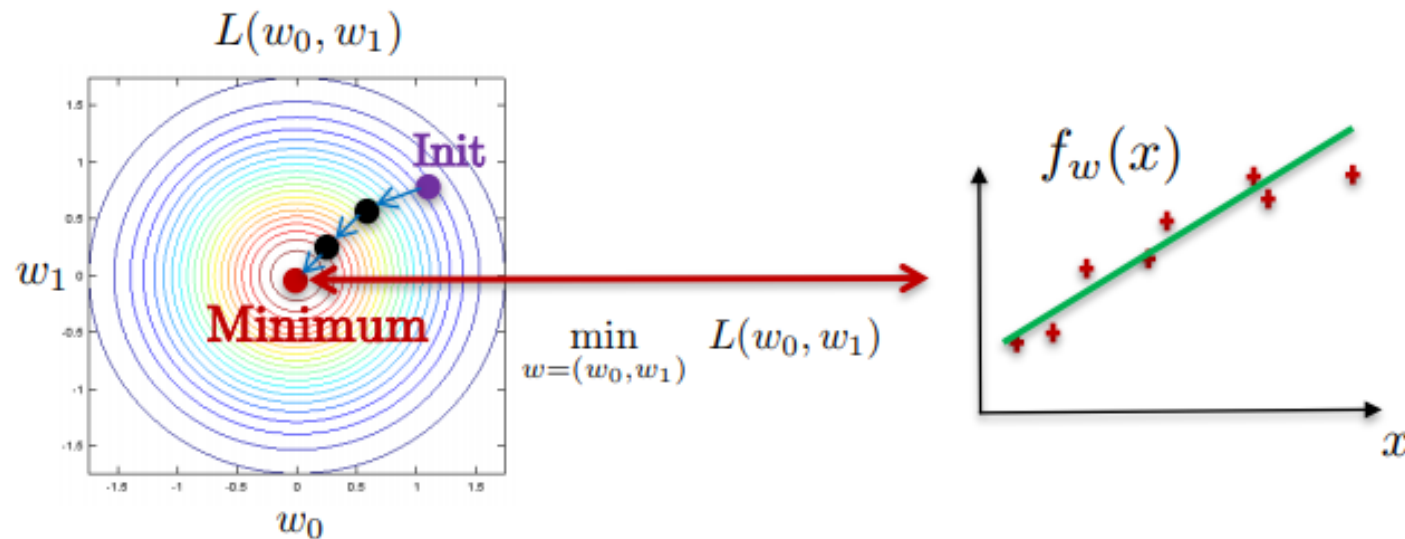# Linear Regression
## Gradient descent with $d=1$ variable

◆ **Prediction function:** $\quad f_w(x) = w_0 + w_1 x$

◆ **Parameters:** $\quad w_0, w_1$

◆ **Loss function:** $\quad L(w_0, w_1) = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( w_0 + w_1 x_i - y_i \right)^2$

◆ **Optimization:** $\quad \min\limits_{w=(w_0,w_1)} L(w_0, w_1)$

◆ **Gradient descent:** $\quad w_j \leftarrow w_j - \tau \dfrac{\partial}{\partial w_j} L(w)$



$L(w_0, w_1)$

Init

$w_1$

Minimum

$w_0$

$\min\limits_{w=(w_0,w_1)} L(w_0, w_1)$

$f_w(x)$

$x$

# Linear Regression
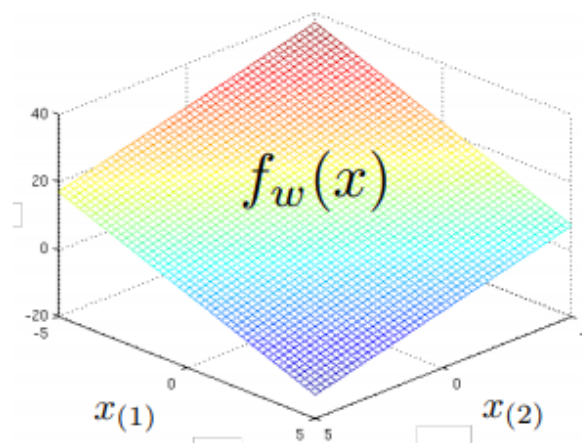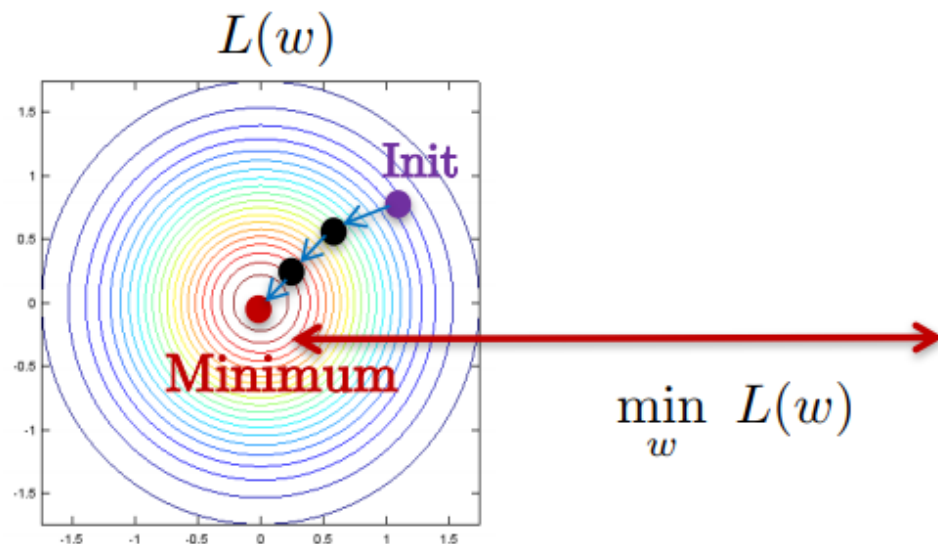# Gradient descent with *d* variables

- **Prediction function:** $f_w(x) = w^T x = w_0 + w_1 x_{(1)} + ... + w_d x_{(d)}$

- **Parameters:** $w = [w_0, w_1, ..., w_d]$

- **Loss function:** $L(w) = \dfrac{1}{n} \sum_{i=1}^{n} \left( w^T x_i - y_i \right)^2$

- **Optimization:** $\min_{w} L(w)$

- **Gradient descent:** $w_j \leftarrow w_j - \tau \dfrac{\partial}{\partial w_j} L(w)$

$L(w)$



$$\min_{w} L(w)$$

$f_w(x)$

$x_{(1)}$     $x_{(2)}$

# Linear Regression
## Gradient descent equations

- *d*=1 (one variable):

$$\frac{\partial}{\partial w_j} L(w_0, w_1) = \frac{\partial}{\partial w_j} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( w_0.1 + w_1 x_i - y_i \right)^2 \right]$$

$$w_0 \leftarrow w_0 - \tau \frac{2}{n} \sum_{i=1}^{n} (w_0 + w_1 x_i - y_i).1$$

$$w_1 \leftarrow w_1 - \tau \frac{2}{n} \sum_{i=1}^{n} (w_0 + w_1 x_i - y_i).x_i$$

- *d* variable

$$\frac{\partial}{\partial w_j} L(w) = \frac{\partial}{\partial w_j} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + ... + w_d x_{i(d)} - y_i \right)^2 \right]$$

# Linear Regression
## Gradient descent equations

◆ **Gradient:**

$$\frac{\partial}{\partial w_j} L(w) = \frac{\partial}{\partial w_j} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right)^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w_j} \left[ \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right)^2 \right]$$

$$= \frac{2}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right) \cdot$$

$$\frac{\partial}{\partial w_j} \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right)$$

$$= \frac{2}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right) x_{i(j)}$$

◆ **Gradient descent:**

$$w_j \leftarrow w_j - \tau \frac{2}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i \right) \cdot x_{i(j)}$$

Data matrix $X_{ij}$

# Linear Regression
## Matrix-vector representation

◆ **Vectorize gradient descent scheme:**

$$w_j \leftarrow w_j - \tau \frac{2}{n} \sum_{i=1}^{n} (w_0 + w_1 x_{i(1)} + \ldots + w_d x_{i(d)} - y_i).x_{i(j)}$$

$$w_j \leftarrow w_j - \tau \frac{2}{n} \sum_{i=1}^{n} X_{ij}.(x_i^T w - y_i)$$

$$x_i = \begin{bmatrix} 1 \\ x_{i(1)} \\ \vdots \\ x_{i(d)} \end{bmatrix} \quad (d+1) \text{ x } 1$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad (d+1) \text{ x } 1$$

**for** $w_0$

$$w_j \leftarrow w_j - \tau \frac{2}{n} X_j^T (Xw - y)$$

$$X = \begin{bmatrix} 1 & -x_1^T- \\ & \vdots \\ 1 & -x_n^T- \end{bmatrix} \quad n \text{ x } (d+1)$$

$$X_j = \begin{bmatrix} x_{1(j)} \\ \vdots \\ x_{n(j)} \end{bmatrix} \quad n \text{ x } 1$$
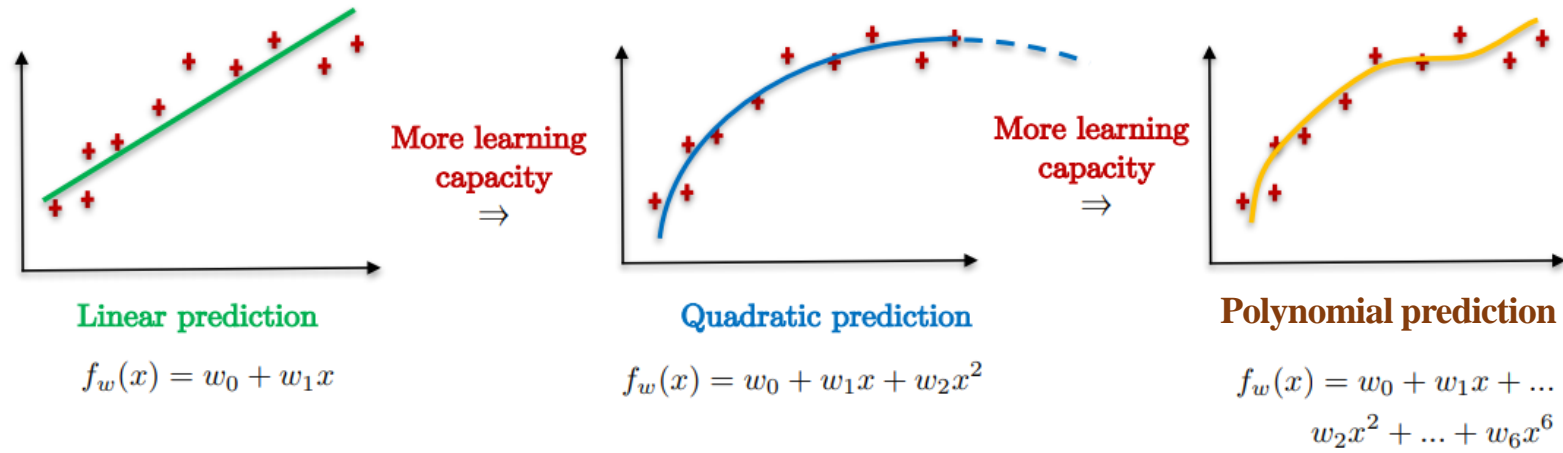
**Data matrix**

$$w \leftarrow w - \tau \frac{2}{n} X^T (Xw - y) \quad \text{1 line of code}$$

$(d+1)$ x 1

$(d+1)$ x $n$        $n$ x 1

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad n \text{ x } 1$$

## Linear Regression
## Beyond linear regression

◆ **Example:** Housing prices prediction function



Linear prediction

$$f_w(x) = w_0 + w_1 x$$

Quadratic prediction

$$f_w(x) = w_0 + w_1 x + w_2 x^2$$

Polynomial prediction

$$f_w(x) = w_0 + w_1 x + ...$$
$$w_2 x^2 + ... + w_6 x^6$$

◆ **Handcrafted prediction function:** Domain expertise allows to define better families of predictive regression functions. Example:

$$f_w(x) = w_0 + w_1 \sqrt{x} + w_2 e^{-x}$$

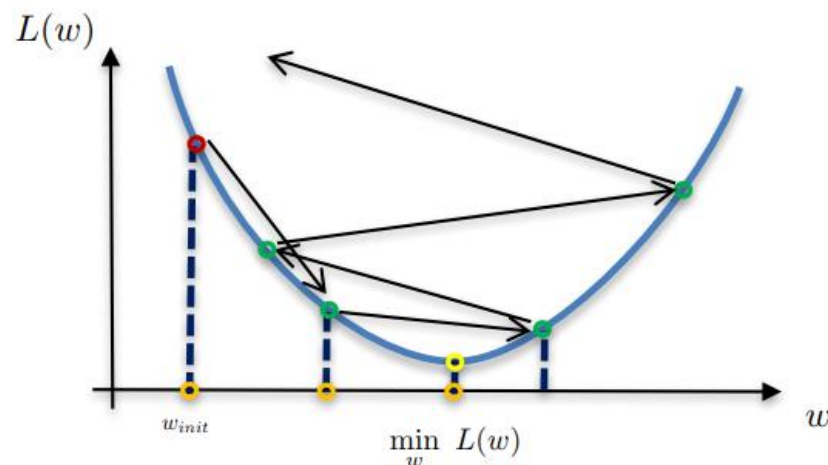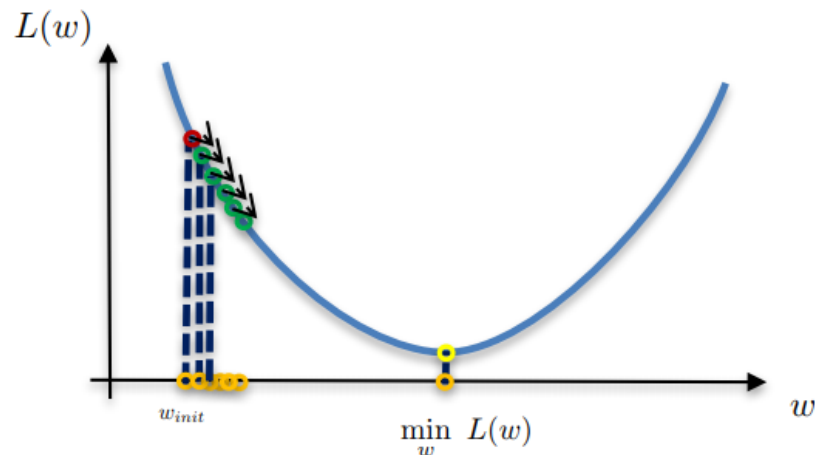◆ **Best non-linear regression technique:** Neural networks.

# Linear Regression
## Beyond gradient descent

◆ **Gradient descent is the most generic optimization technique.**

$$\min_{w} \; L(w)$$

$$w_j \leftarrow w_j - \tau \frac{\partial}{\partial w_j} L(w)$$

◆ **But it has limitations:**

- **Choice of learning rate $\tau$**

- **Convergence speed (even with optimal $\tau$)**

# Linear Regression
## Beyond gradient descent

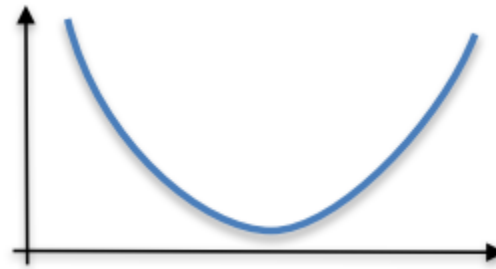◆ **We can leverage some mathematical properties to speed up the optimization.**

- **Prediction** function is **linear**:

$$f_w(x) = w^T x = w_0 + w_1 x_{(1)} + \ldots + w_d x_{(d)}$$

- **Loss** function is **convex** (quadratic):

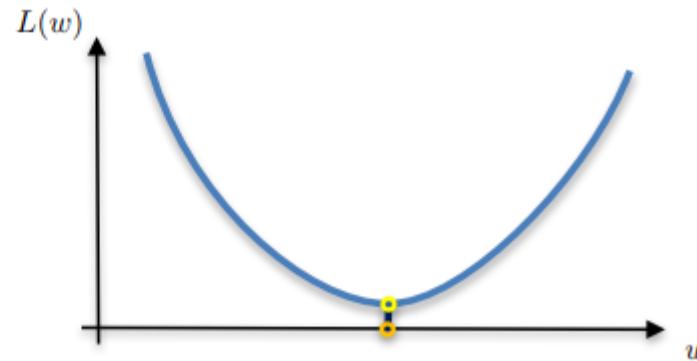$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \left( f_w(x_i) - y_i \right)^2$$

# Linear Regression
## Normal equation for $d$=1 and $n$=1

◆ **Normal equation:** Solution of mean square error loss

$$\min_{w} \left\{ L(w) = (wx - y)^2 \right\}$$



◆ **The minimum is obtained when the gradient/slope is zero:**

$$\frac{\partial}{\partial w} L(w) = 0$$

$$\frac{\partial}{\partial w}(wx - y)^2 = 2x(wx - y) = 0 \Rightarrow \boxed{w = x^{-1}y}$$

Solution                One line of code

# Linear Regression
## Normal equation for $d=1$ and $n$ data

- **Loss $L$:**

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \left( wx_i - y_i \right)^2$$

- **Gradient of the loss $L$ w.r.t. $w$:**

$$\frac{\partial}{\partial w} L(w) = 0 \Rightarrow \min_w \ L(w)$$

$$\frac{\partial}{\partial w} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( wx_i - y_i \right)^2 \right] = \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w} \left[ \left( wx_i - y_i \right)^2 \right]$$

$$= \frac{2}{n} \sum_{i=1}^{n} (wx_i - y_i) \frac{\partial}{\partial w} (wx_i - y_i)$$

$$= \frac{2}{n} \sum_{i=1}^{n} (wx_i - y_i) x_i$$

$$= \frac{2}{n} w \sum_{i=1}^{n} x_i^2 - \frac{2}{n} \sum_{i=1}^{n} y_i x_i$$

$$= 0 \Rightarrow w = \frac{\sum_{i=1}^{n} y_i x_i}{\sum_{i=1}^{n} x_i^2}$$

Solution

# Linear Regression
## Vectorization

◆ **Loss $L$:**

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \left( wx_i - y_i \right)^2$$

$$= \frac{1}{n} \underbrace{\left( wx - y \right)^T}_{1 \times n} \underbrace{\left( wx - y \right)}_{n \times 1}$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$
$n \times 1$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$
$n \times 1$

# Linear Regression
## Vectorization

◆ **Gradient of the loss $L$ w.r.t. $w$:**

$$\frac{\partial}{\partial w} L(w) = 0 \Rightarrow \min_{w} \; L(w)$$

$$\frac{\partial}{\partial w}\left[\frac{1}{n}(wx-y)^{T}(wx-y)\right] = \frac{1}{n}\frac{\partial}{\partial w}\left[(wx^{T}-y^{T})(wx-y)\right]$$

$$= \frac{2}{n}x^{T}(wx-y)$$

$$= \frac{2}{n}(wx^{T}x - x^{T}y)$$

$$= 0 \Rightarrow w = (x^{T}x)^{-1}x^{T}y \quad \text{One line of code}$$

# Linear Regression
## Normal equation for *d* features and *n* data

◆ **Loss *L*:**

$$L(w_0, ..., w_d) = \frac{1}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + ... + w_d x_{i(d)} - y_i \right)^2$$

◆ **Gradient of the loss *L* w.r.t. *w*:**

$$\frac{\partial}{\partial w_j} L(w_0, ..., w_d) = 0 \quad \forall j$$

$$\frac{\partial}{\partial w_j} L(w) = \frac{\partial}{\partial w_j} \left[ \frac{1}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + ... + w_d x_{i(d)} - y_i \right)^2 \right]$$

$$= \frac{2}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + ... + w_d x_{i(d)} - y_i \right) x_{i(j)}$$

$$= 0 \Rightarrow w_j = \frac{\sum_{k \neq j} \sum_i w_k x_{i(k)} x_{i(j)}}{\sum_i x_{i(j)}^2} \qquad \textbf{with} \ \ x_{i(0)} = 1$$

Solution

# Linear Regression
## Vectorization

◆ **Loss $L$:**

$$L(w_0, ..., w_d) = \frac{1}{n} \sum_{i=1}^{n} \left( w_0 + w_1 x_{i(1)} + ... + w_d x_{i(d)} - y_i \right)^2$$

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} \left( x_i^T w - y_i \right)^2$$

$$L(w) = \frac{1}{n} \left( Xw - y \right)^T \left( Xw - y \right)$$

$$w = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix} \quad x_i = \begin{bmatrix} x_{i(0)} \\ \vdots \\ x_{i(d)} \end{bmatrix} \quad X = \begin{bmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$(d+1) \times 1$      $(d+1) \times 1$      n x (d+1)      n x 1

**with** $x_{i(0)} = 1$      **Data matrix**

# Linear Regression
## Vectorization

◆ **Gradient of the loss $L$ w.r.t. $w$:**

$$\frac{\partial}{\partial w} L(w) = \frac{1}{n} \frac{\partial}{\partial w} \left[ \left( Xw - y \right)^T \left( Xw - y \right) \right]$$

$$= \frac{1}{n} \frac{\partial}{\partial w} \left[ \left( w^T X^T - y^T \right) \left( Xw - y \right) \right]$$

$$= \frac{2}{n} X^T \left( Xw - y \right)$$

$$= \frac{2}{n} \left( X^T Xw - X^T y \right)$$

$$= 0 \Rightarrow w = (X^T X)^{-1} X^T y$$

One line of code