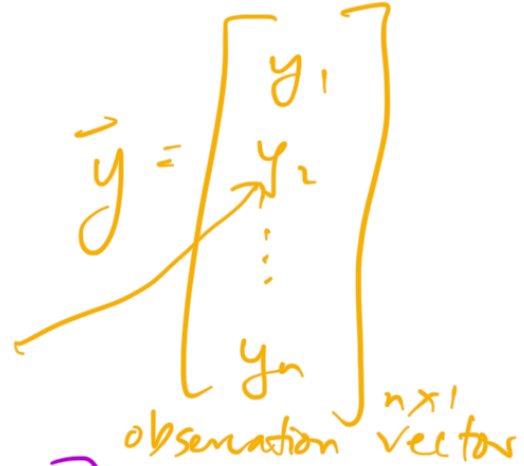
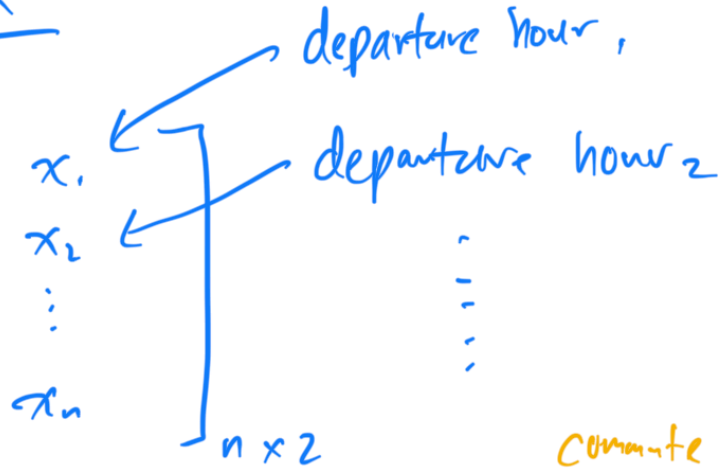




$X$  design matrix

$$X = \begin{bmatrix} | \\ \vdots \\ | \end{bmatrix}$$



$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

parameter vector

$$\vec{h} = X\vec{w} = \begin{bmatrix} H(x_1) \\ H(x_2) \\ \vdots \\ H(x_n) \end{bmatrix}$$

hypothesis vector

compute time<sub>2</sub>

mean squared error

Goal: choose the  $\vec{w}^*$  that minimizes  $\frac{1}{n} \|\vec{y} - X\vec{w}\|^2$

$R_{sq}(\vec{w})$  < 6.1 >

?



Equivalent to the other formulas:

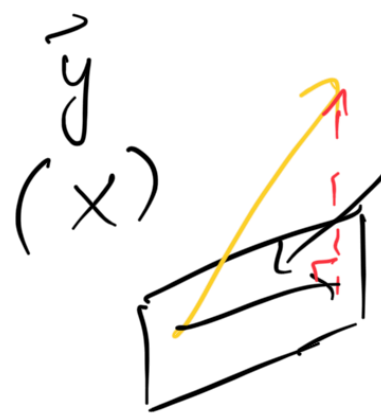
$$X^T X \vec{w}^* = X^T \vec{y}$$

normal equations

if  $X^T X$  is invertible, unique solution

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

got this by projecting onto span



best predictions

the same!!!

$$w_1^* = r \frac{\sigma_y}{\sigma_x}$$

$$w_0^* = \bar{y} - w_1^* \bar{x}$$

?



$$X = \begin{bmatrix} 1 & dh_1 & dm_1 \\ 1 & dh_2 & dm_2 \\ \vdots & \vdots & \vdots \\ 1 & dh_n & dm_n \end{bmatrix} \quad n \times 3$$

departure hour<sub>n</sub> (pointing to  $dh_n$ )  
day of month<sub>n</sub> (pointing to  $dm_n$ )

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

unchanged!

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad 3 \times 1$$



## Feature vectors

- Suppose we have the following dataset.

|     | departure_hour | day_of_month | minutes |
|-----|----------------|--------------|---------|
| row |                |              |         |
| 1   | 8.45           | 22           | 63.0    |
| 2   | 8.90           | 28           | 89.0    |
| 3   | 8.72           | 18           | 89.0    |

not a feature!

$$\vec{x}_1 = \begin{bmatrix} 8.45 \\ 22 \end{bmatrix}_{2 \times 1}$$

$$\vec{x}_2 = \begin{bmatrix} 8.90 \\ 28 \end{bmatrix}$$

$$\vec{x}_3 = \begin{bmatrix} 8.72 \\ 18 \end{bmatrix}$$

# ⊗ Augmented feature vector

$\text{Aug}(\vec{x}) =$  append a 1 to the beginning of  $\vec{x}$

$$\text{Aug}(\vec{x}_1) = \begin{bmatrix} 1 \\ 8.45 \\ 22 \end{bmatrix}$$

In general, if  $\vec{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix}_{d \times 1}$

then

$$\text{Aug}(\vec{x}) = \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix}_{(d+1) \times 1}$$

⊗ Why make an augmented feature vector?  
 ⇒  $Aug(\vec{x})$  and  $\vec{w}$  have same dimensions!

e.g. two features

$$\vec{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

intercept term

$$Aug(\vec{x}) = \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \end{bmatrix}$$

now, our predictions for a single data point are a dot product!

$$\begin{aligned}
 H(\vec{x}) &= w_0 + w_1 x^{(1)} + w_2 x^{(2)} \\
 &= \vec{w} \cdot Aug(\vec{x})
 \end{aligned}$$

↑ upgraded  $y = w_0 + w_1 x$





In general, with  $d$  features,  $X$  is of shape  $n \times d+1$ !

$$X = \begin{bmatrix} 1 & x_1 & x_1 & \dots & x_1 \\ 1 & x_2 & x_2 & \dots & x_2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n & \dots & x_n \end{bmatrix}_{n \times (d+1)}$$





# Linear in the parameters

- Using linear regression, we can fit rules like:

$$w_0 + w_1x + w_2x^2$$

$$w_1 e^{-x^{(1)^2}} + w_2 \cos(x^{(2)} + \pi) + w_3 \frac{\log 2x^{(3)}}{x^{(2)}}$$

no parameters!

- This includes arbitrary polynomials.
- These are all linear combinations of (just) features.

$w_0 + w_1 \boxed{\phantom{x}} + w_2 \boxed{\phantom{x}} + \dots + w_d \boxed{\phantom{x}}$   
 should not have  $w!$

- For any of the above examples, we **could** express our model as a product of a design matrix and parameter vector, and that's all that LinearRegression in sklearn needs.

What we put in the X argument to model.fit is up to us!

$w_1$  is part of  $e!$

- Using linear regression, we **can't** fit rules like:

$$w_0 + e^{w_1x}$$

$$w_0 + \sin(w_1x^{(1)} + w_2x^{(2)})$$

- These are **not** linear combinations of just features.

$w_1, w_2$  are in the sin!

$$X \vec{w}^*$$

