

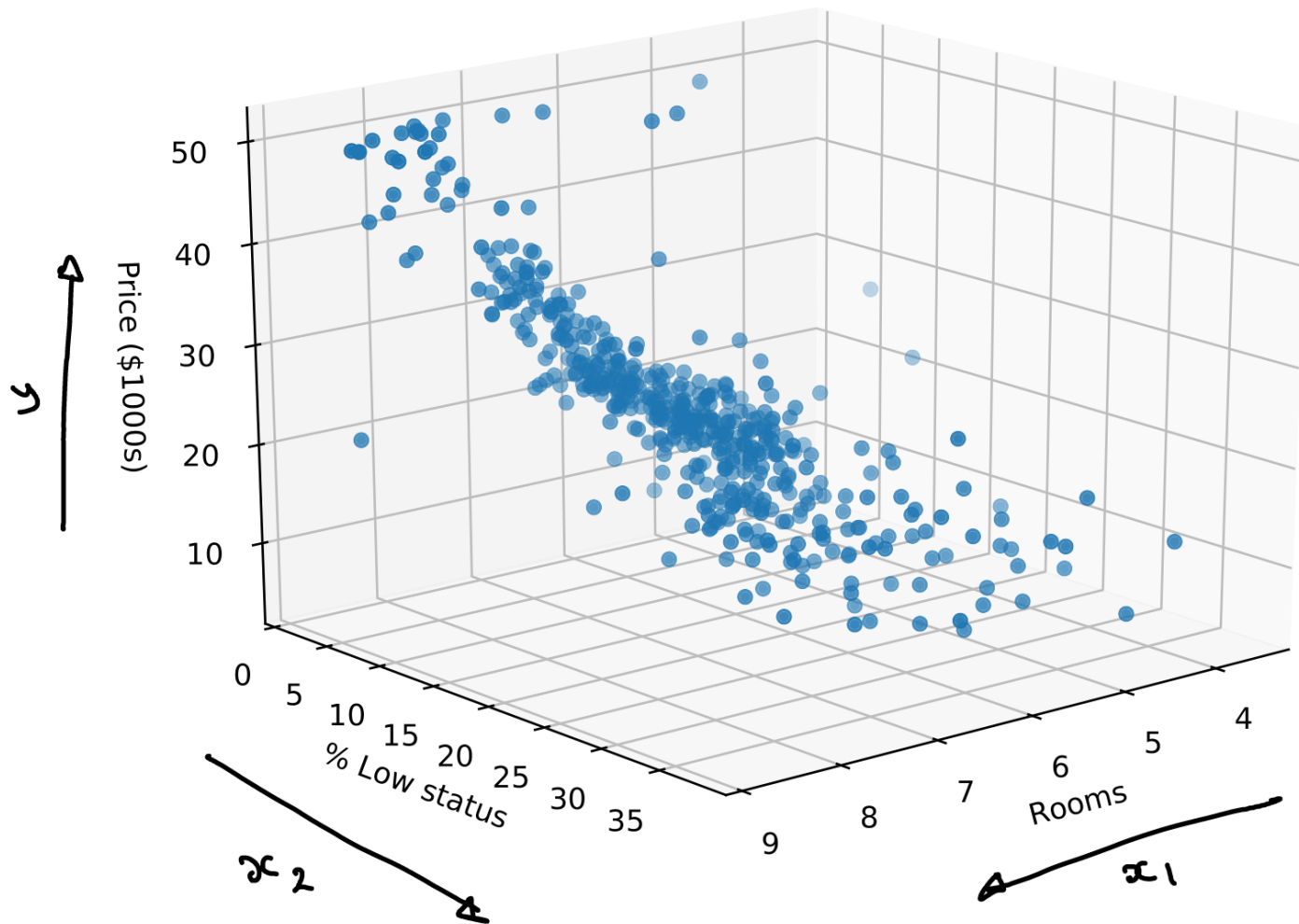
Multiple linear regression

Model and loss

Herman Kamper

<http://www.kamperh.com/>

Boston house prices



Multiple linear regression

The model:

$$f(x_1, x_2, \dots, x_D; w_0, w_1, w_2, \dots, w_D)$$

$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

"bias"

$$f(\underline{x}; \underline{w}) = \underline{w}^T \underline{x} \quad [\text{Pretending } x_0 = 1]$$

with

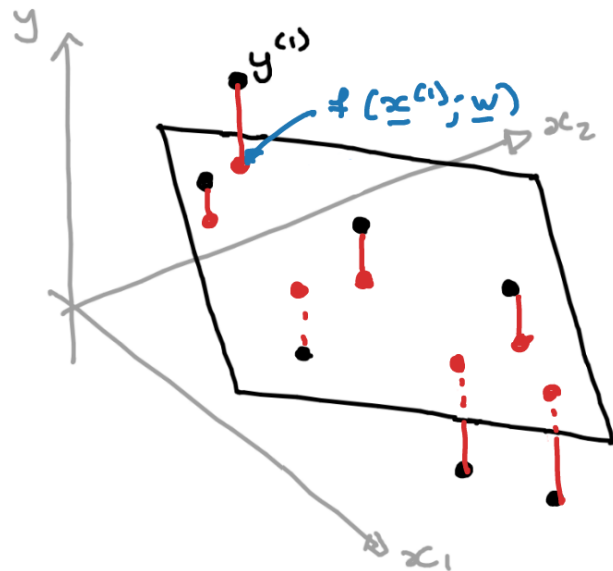
$$\underline{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \text{and} \quad \underline{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

Loss function:

Squared loss:

$$J(\underline{w}) = \sum_{n=1}^N (y^{(n)} - f(\underline{x}^{(n)}; \underline{w}))^2$$

$$= \sum_{n=1}^N (y^{(n)} - (w_0 + w_1 x_1^{(n)} + w_2 x_2^{(n)} + \dots + w_D x_D^{(n)}))^2$$



Optimization:

Derive $\frac{\partial J}{\partial w_0}$, $\frac{\partial J}{\partial w_1}$, \dots , $\frac{\partial J}{\partial w_D}$

and set equal to zero.

Painful to do term-by-term.

Idea: Rather write in vector form and find $\frac{\partial J}{\partial \underline{w}}$.

Interlude: Watch vector and matrix derivatives.

Writing the loss in matrix form

We want to minimize:

$$J(\underline{w}) = \sum_{i=1}^N (y^{(i)} - f(\underline{x}^{(i)}; \underline{w}))^2$$

$$= \sum_{i=1}^N (y^{(i)} - \underline{w}^T \underline{x}^{(i)})^2 \dots \textcircled{1}$$

"design matrix"

Define:

$$\underline{X} = \begin{bmatrix} - (\underline{x}^{(1)})^T - \\ - (\underline{x}^{(2)})^T - \\ \vdots \\ - (\underline{x}^{(N)})^T - \end{bmatrix}; \quad \underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

Now we can write $\textcircled{1}$ as:

$$J(\underline{w}) = (\underline{y} - \underline{X}\underline{w})^T (\underline{y} - \underline{X}\underline{w}) \dots \textcircled{2}$$

To see this, you can define an error vector:

$$\underline{e} = \begin{bmatrix} y^{(1)} - \underline{w}^T \underline{x}^{(1)} \\ y^{(2)} - \underline{w}^T \underline{x}^{(2)} \\ \vdots \\ y^{(N)} - \underline{w}^T \underline{x}^{(N)} \end{bmatrix}$$

Then note that $\textcircled{1}$ can be written as $J(\underline{w}) = \underline{e}^T \underline{e}$, and $\underline{e} = \underline{y} - \underline{X}\underline{w}$, which leads to $\textcircled{2}$. We now use the form in $\textcircled{2}$ to determine $\frac{\partial J}{\partial \underline{w}}$.

$$J(\underline{w}) = \underline{y}^T \underline{y} - \underbrace{\underline{y}^T \underline{X}}_{\underline{1}^T \underline{X}} \underline{w} - \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w}$$

$$= \underline{y}^T \underline{y} - \underline{w}^T \underline{X}^T \underline{y} - \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w}$$

$$= \underline{y}^T \underline{y} - 2 \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w}$$

Multiple linear regression

Optimization

Herman Kamper

<http://www.kamperh.com/>

The normal equations

Now we set $\frac{\partial J}{\partial \underline{w}} = \underline{0}$, i.e. $\frac{\partial J}{\partial w_0} = 0, \dots, \frac{\partial J}{\partial w_D} = 0$

$$\begin{aligned} \frac{\partial J}{\partial \underline{w}} &= \frac{\partial}{\partial \underline{w}} \left[\underline{y}^T \underline{y} - 2 \underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w} \right] \\ &= -2 \underline{X}^T \underline{y} + \left(\underline{X}^T \underline{X} + (\underline{X}^T \underline{X})^T \right) \underline{w} \dots \textcircled{3} \\ &= -2 \underline{X}^T \underline{y} + 2 \underline{X}^T \underline{X} \underline{w} \end{aligned}$$

Set $\frac{\partial J}{\partial \underline{w}} = 0$, then:

$$\begin{aligned} \underline{X}^T \underline{X} \underline{w} &= \underline{X}^T \underline{y} \\ \underline{w} &= (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \end{aligned}$$

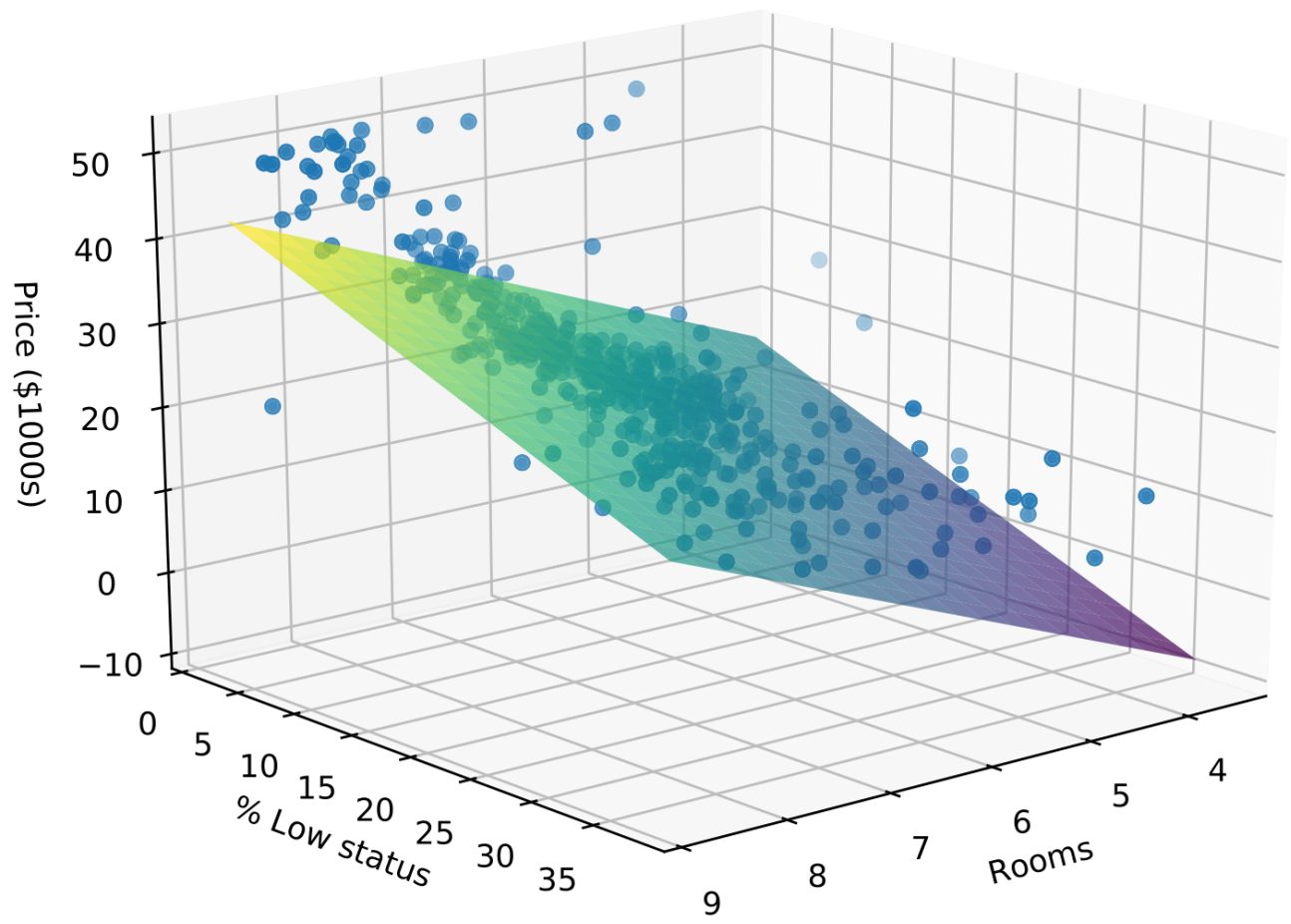
This is called the normal equations. With one line of Python, can get the estimates of $D+1$ parameters.

To get to $\textcircled{3}$, we used the following identities:

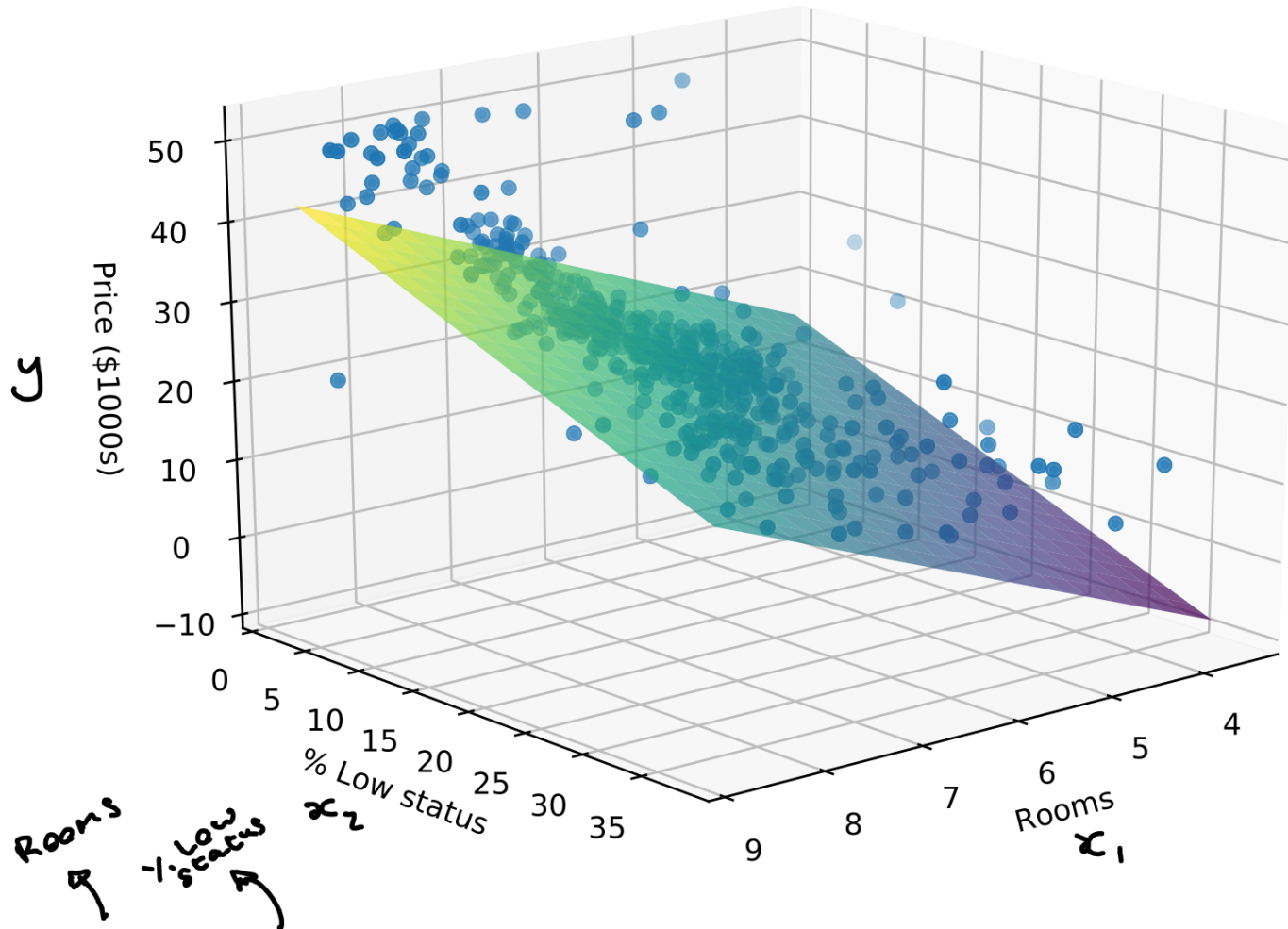
$$\frac{\partial \underline{x}^T \underline{A} \underline{x}}{\partial \underline{x}} = (\underline{A} + \underline{A}^T) \underline{x}$$
$$\frac{\partial \underline{x}^T \underline{a}}{\partial \underline{x}} = \underline{a}$$

You can find these in the Matrix calculus Wikipedia article.

Boston house prices fit



Boston house prices fit



$$f(\mathbf{x}; \hat{\mathbf{w}}) = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2 = -1.358 + 5.095x_1 - 0.642x_2$$