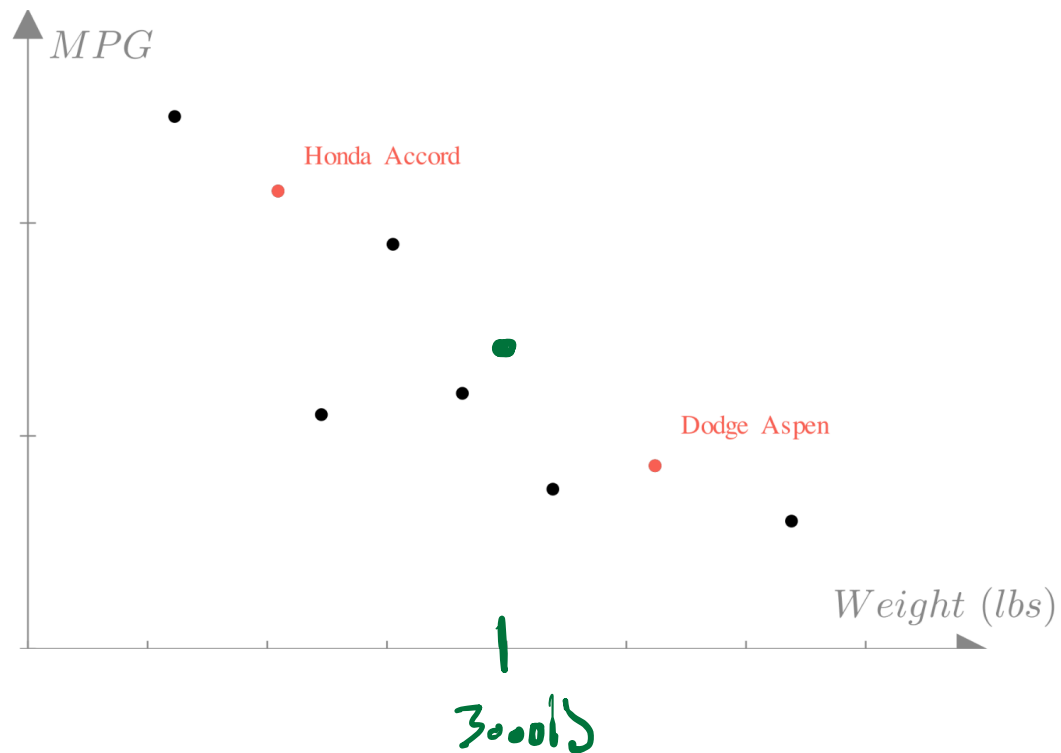Honda Accord: $\begin{bmatrix} \text{Weight:} & \text{2500 lbs} \\ \text{Horsepower:} & \text{123 HP} \\ \text{Displacement:} & \text{2.4 L} \\ \text{0-60mph:} & \text{7.8 Sec} \end{bmatrix} \longrightarrow$ MPG: 33mpg

Dodge Aspen: $\begin{bmatrix} \text{Weight:} & \text{3800 lbs} \\ \text{Horsepower:} & \text{155 HP} \\ \text{Displacement:} & \text{3.2 L} \\ \text{0-60mph:} & \text{6.8 Sec} \end{bmatrix} \longrightarrow$ MPG: 21mpg

$\vdots \quad \vdots$

$$\begin{bmatrix} \text{Weight} \\ hP \\ \text{Displct} \\ \text{Acc.} \end{bmatrix} \longrightarrow MPG \ ?$$

Dataset: $\mathcal{D} = \{(\mathbf{x}_i, y_i) \text{ for } i \in 1 \ldots N\}$

$$\text{Input: } \mathbf{x}_i = \begin{bmatrix} \text{Weight} \\ \text{Horsepower} \\ \text{Displacement} \\ \text{0-60mph} \end{bmatrix}, \quad \text{Output: } y_i = MPG$$

$$mpg = f(\text{weight, horsepower} \ldots)$$

function

$$y = f(x)$$

$x_i$: Vector input

$y_i$: Scalar output

$i$: Index of obs

$D$: Data set

$N$: size of $D$

# Vectors

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \Big\} \, d = 3$$

$\mathbf{x}$ — bolded

$x_3$ — scalar

```
x = np.array([2, 5, 1])
```
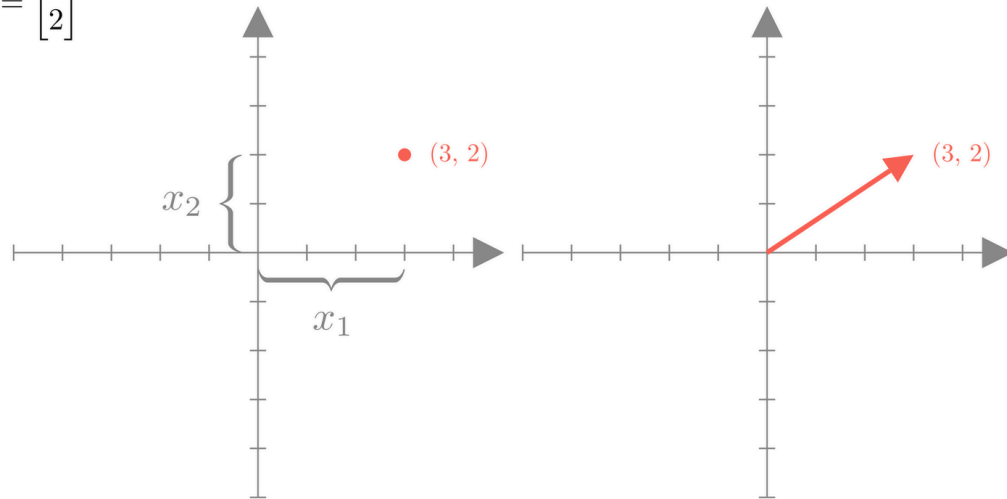
$$\mathbf{x} = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$



$x^T$

Column vector: $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$, Row vector: $\mathbf{x}^T = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$
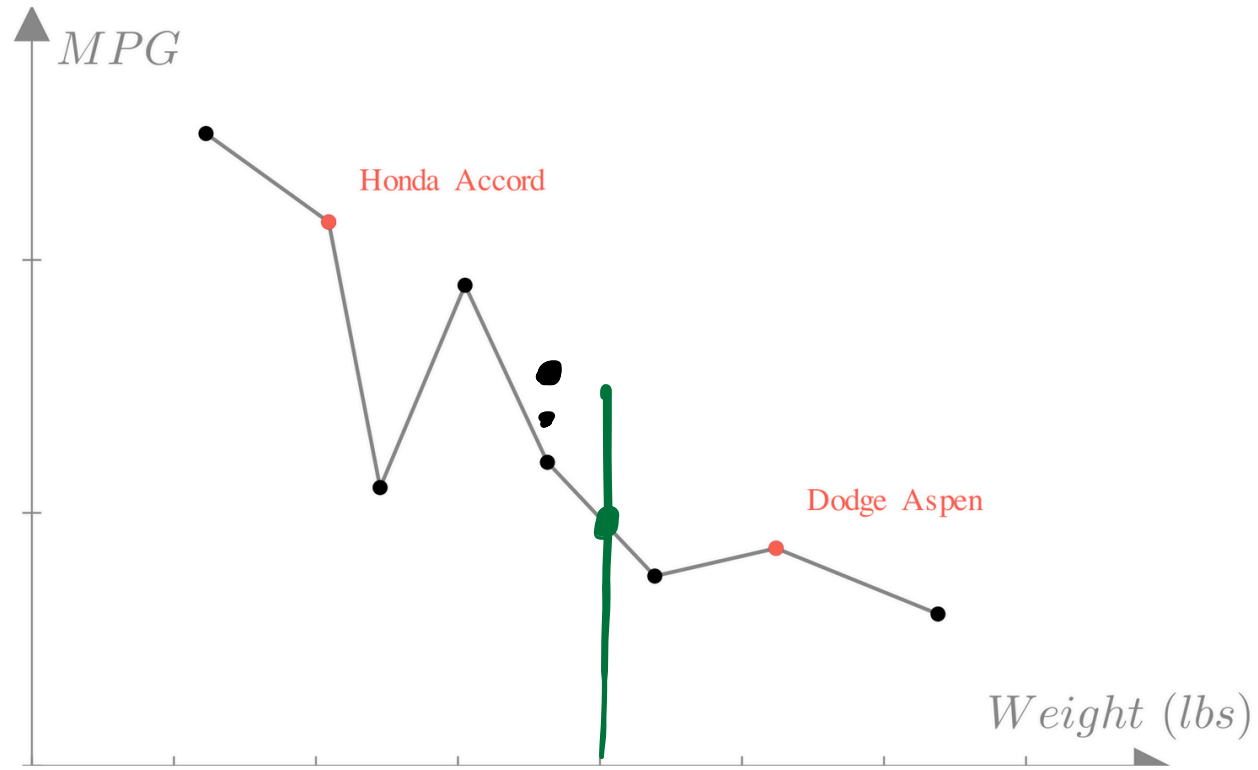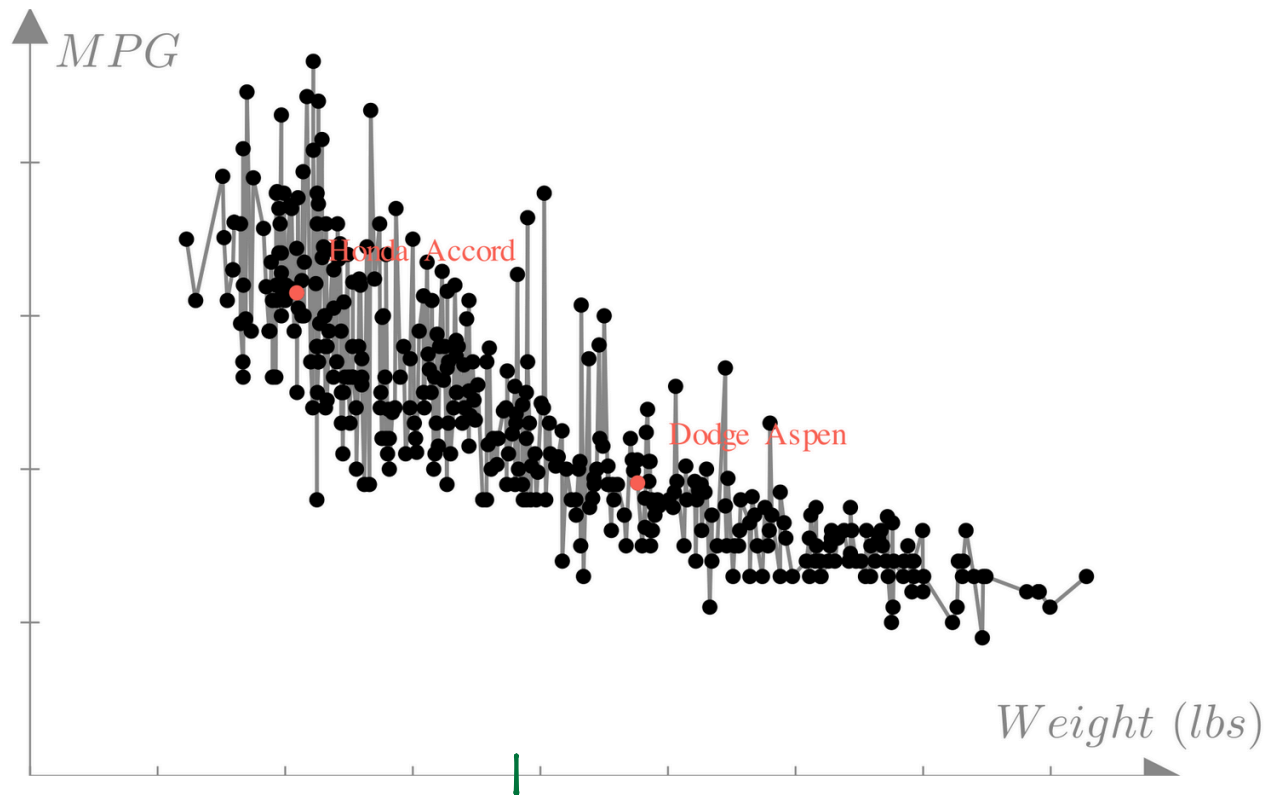
|||

$d \times 1$ matrix

|||

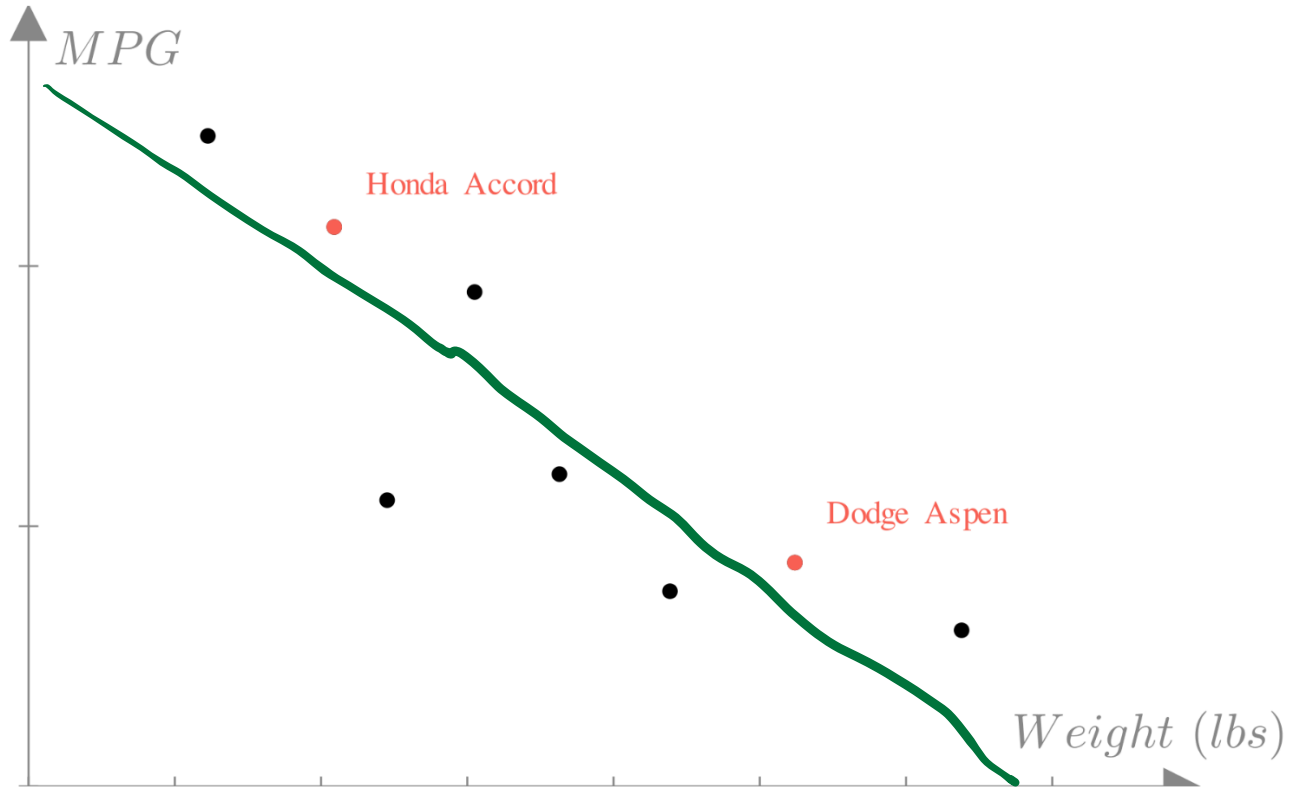$1 \times d$ matrix

linear interpolation

MPG

Honda Accord

Dodge Aspen

Weight (lbs)

Linear Regression

A *linear function* is any function $f$ where the following conditions always hold:

$$f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$$

and

$$f(a\mathbf{x}) = af(\mathbf{x})$$

For a linear function, the output can be defined as a weighted sum of the inputs. In other words a linear function of a vector can always be written as:
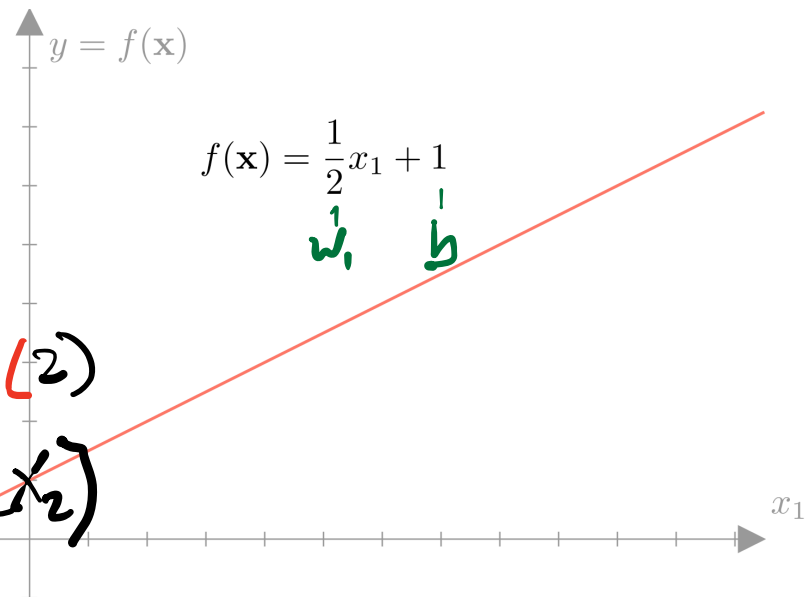
$$f(\mathbf{x}) = \sum_{i=1}^{n} x_i w_i$$

We can add an offset $b$ to create an *affine* function:

$$f(\mathbf{x}) = \sum_{i=1}^{n} x_i w_i + b$$

$$f\left(\begin{bmatrix} 1 \\ 2 \end{bmatrix}\right) = 3(1) + 5(2)$$
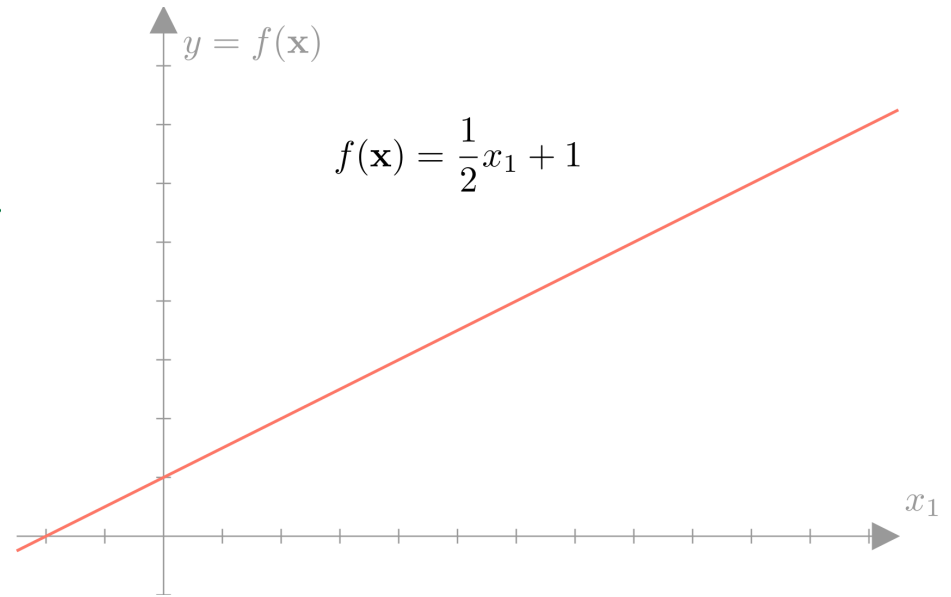
$$f(x) = 3(x_1) + 5(x_2)$$

$$f(\underline{x}) = \sum_{i=1}^{n} \overset{x}{x_i} \overset{w}{w_i} = x \cdot w = x^T w = [\longrightarrow] \begin{bmatrix} \downarrow \end{bmatrix}$$

$y = f(\mathbf{x})$

$$f(\mathbf{x}) = \frac{1}{2} x_1 + 1$$
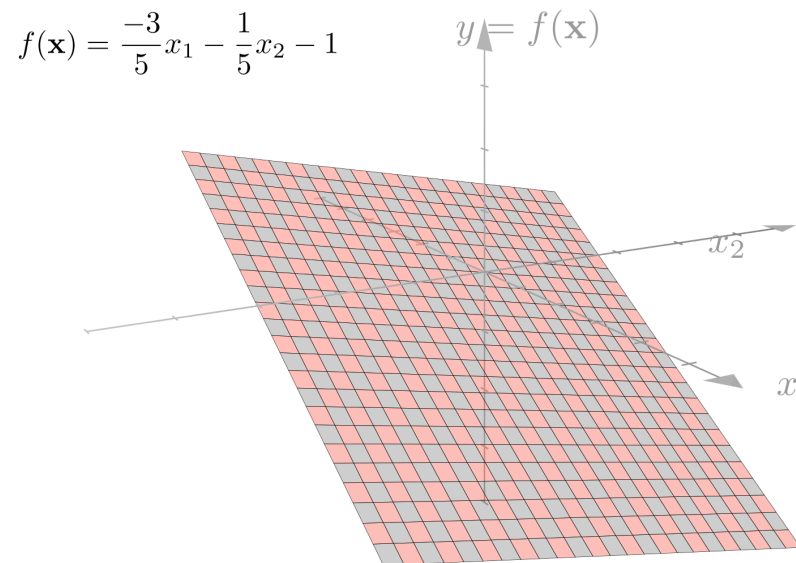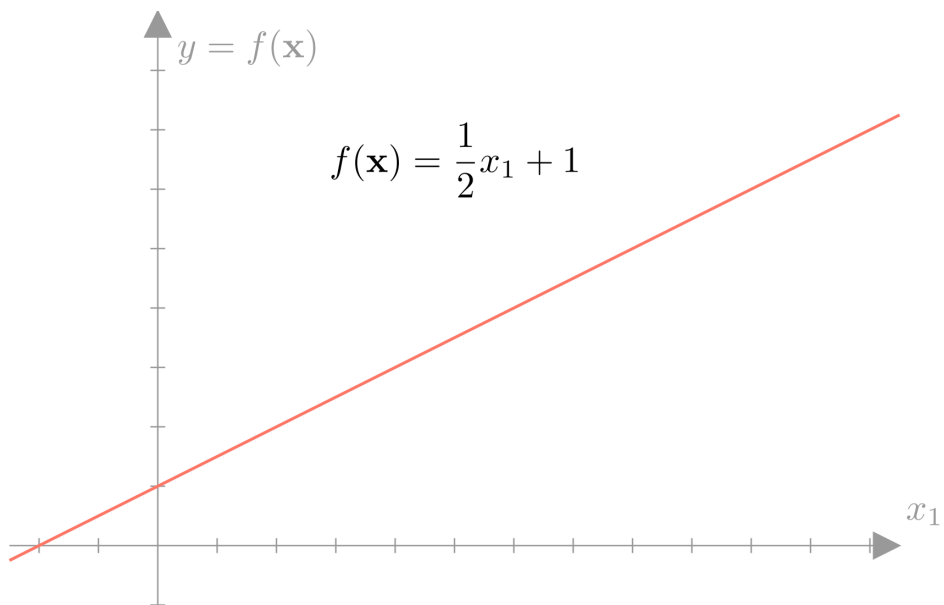
$w_1 \quad b$

$x_1$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \sum_{i=1}^{n} x_i w_i$$

```
def f(x):
    w = np.array([-0.6, -0.2])
    b = -1
    return np.dot(x, w) + b
```

Numpy

$y = f(\mathbf{x})$

$$f(\mathbf{x}) = \frac{1}{2} x_1 + 1$$

$x_1$

$$f(\mathbf{x}) = \frac{1}{2}x_1 + 1$$

$$f(\mathbf{x}) = \frac{-3}{5}x_1 - \frac{1}{5}x_2 - 1$$

**Linear regression** is the approach of modeling an unknown function with a linear function. From our discussion of linear functions, we know that this means that we will make predictions using a function of the form:

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \sum_{i=1}^{n} x_i w_i$$

Meaning that the output will be a weighted sum of the *features* of the input. In the case of our car example, we will make predictions as:
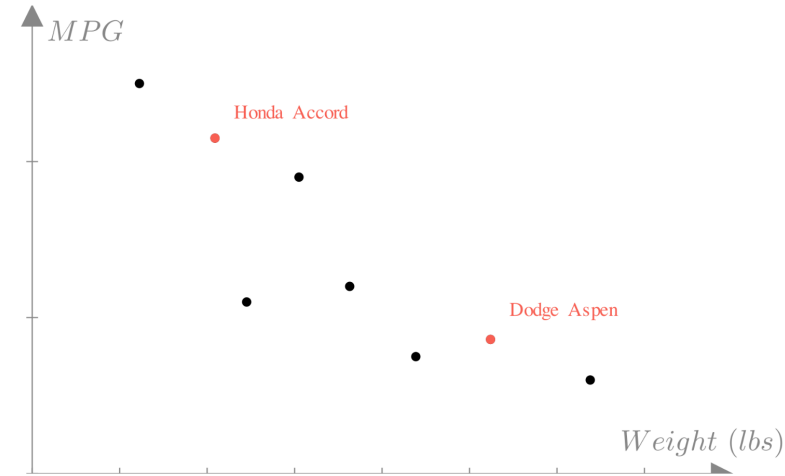
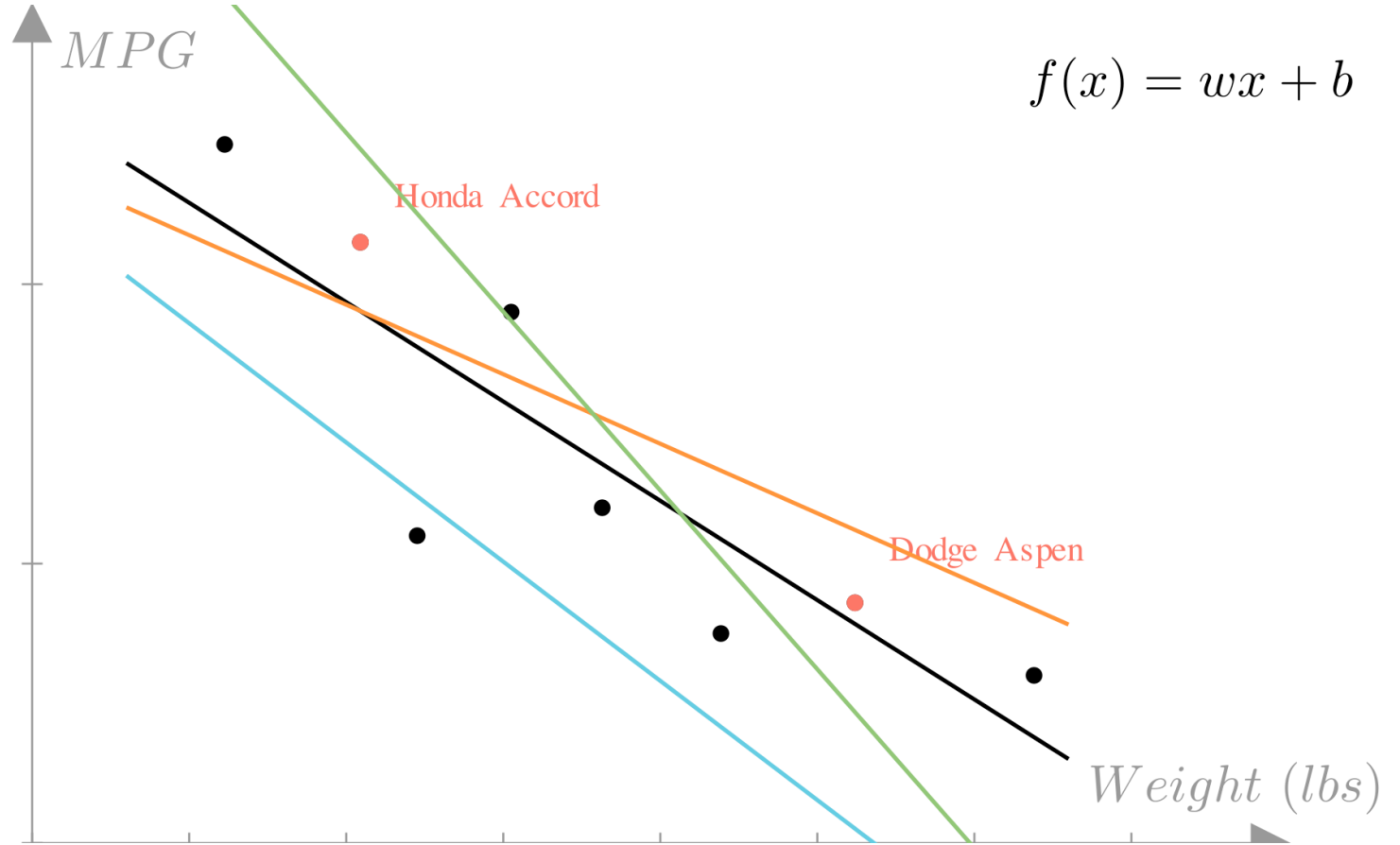$$\text{Predicted MPG} = f(\mathbf{x}) =$$

$$(\text{weight})w_1 + (\text{horsepower})w_2 + (\text{displacement})w_3 + (\text{0-60mph})w_4 + b$$

Or in matrix notation:

$$f(\mathbf{x}) = \begin{bmatrix} \text{Weight} \\ \text{Horsepower} \\ \text{Displacement} \\ \text{0-60mph} \\ 1 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ b \end{bmatrix}$$

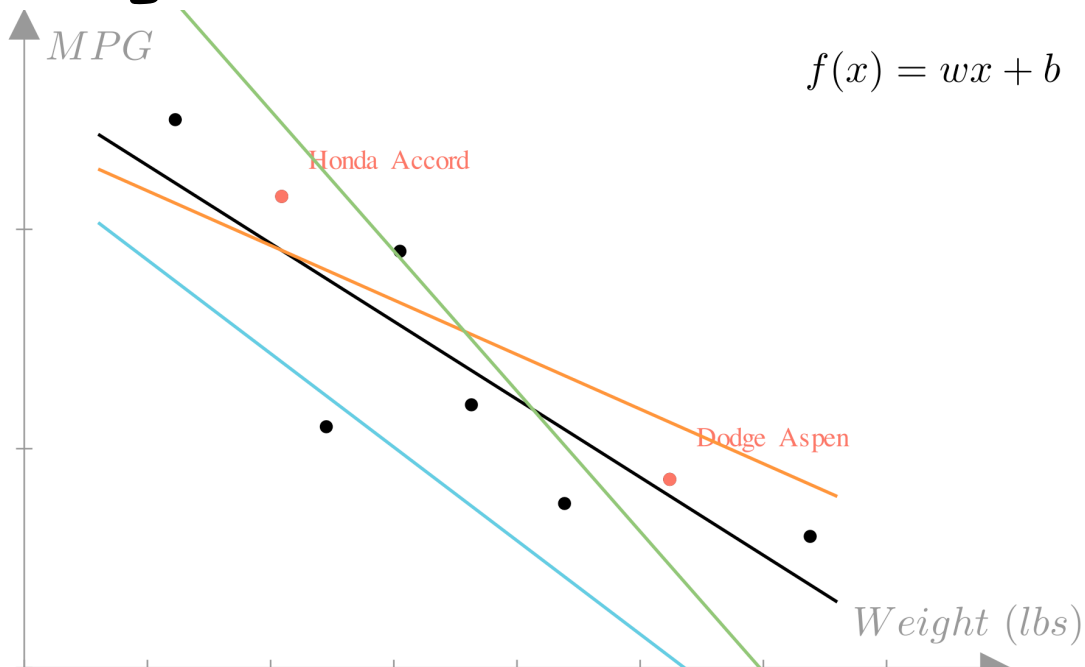Obs. cal

$$f(\mathbf{x}) = \sum_{i=1}^{n} x_i w_i + b$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_n \\ \vdots \end{bmatrix}$$

We typically refer to $\mathbf{w}$ specifically as the **weight vector** (or weights) and $b$ as the **bias**. To summarize:

**Affine function:** $f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$,   **Parameters:**   (Weights: $\mathbf{w}$, Bias: $b$)
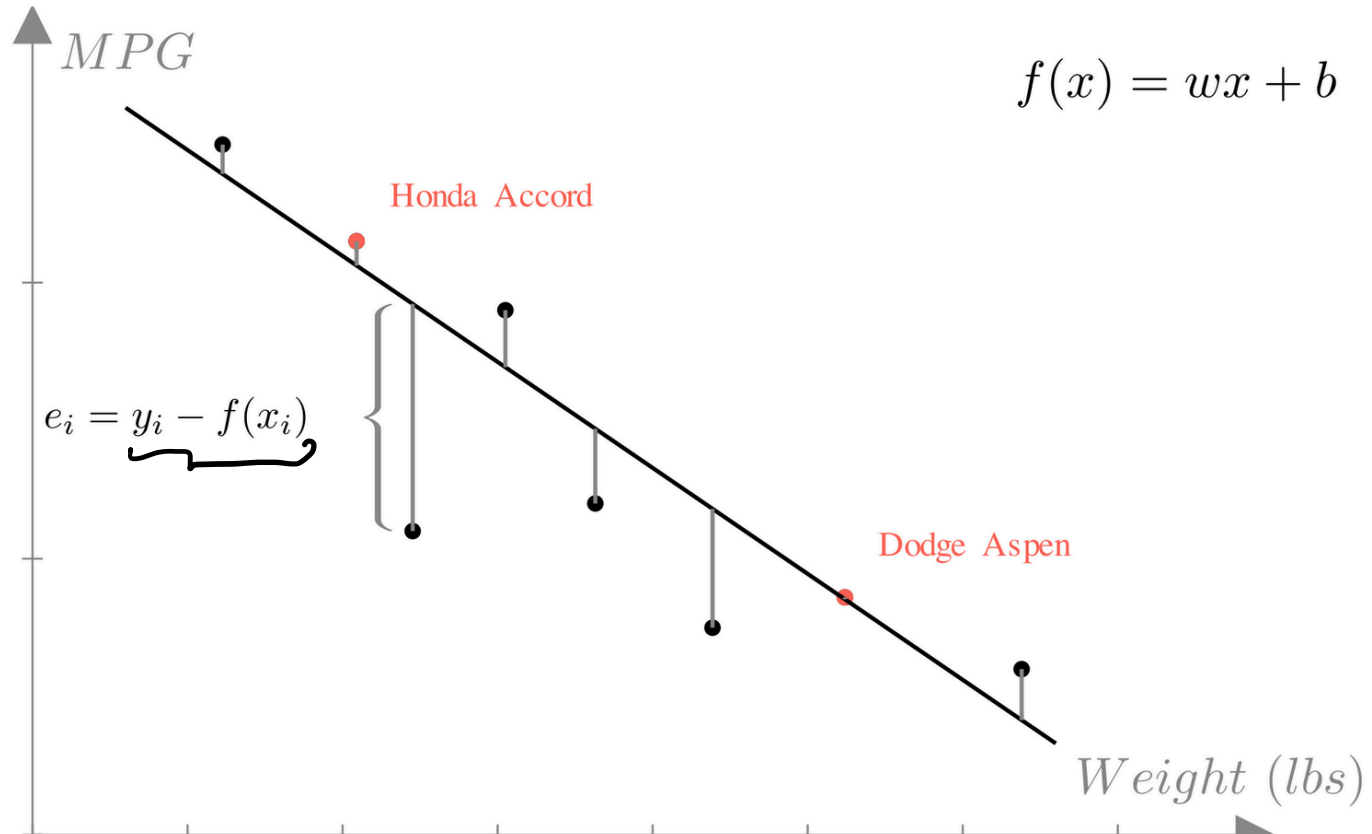
$$f(x) = wx + b$$

```
class Regression:
    def __init__(self, weights):
        self.weights = weights

    def predict(self, x):
        return np.dot(x, self.weights)
```



*MPG* — *Weight (lbs)* scatter plot with Honda Accord and Dodge Aspen points and multiple regression lines.

The **residual** or **error** of a prediction is the difference between the prediction and the true output:
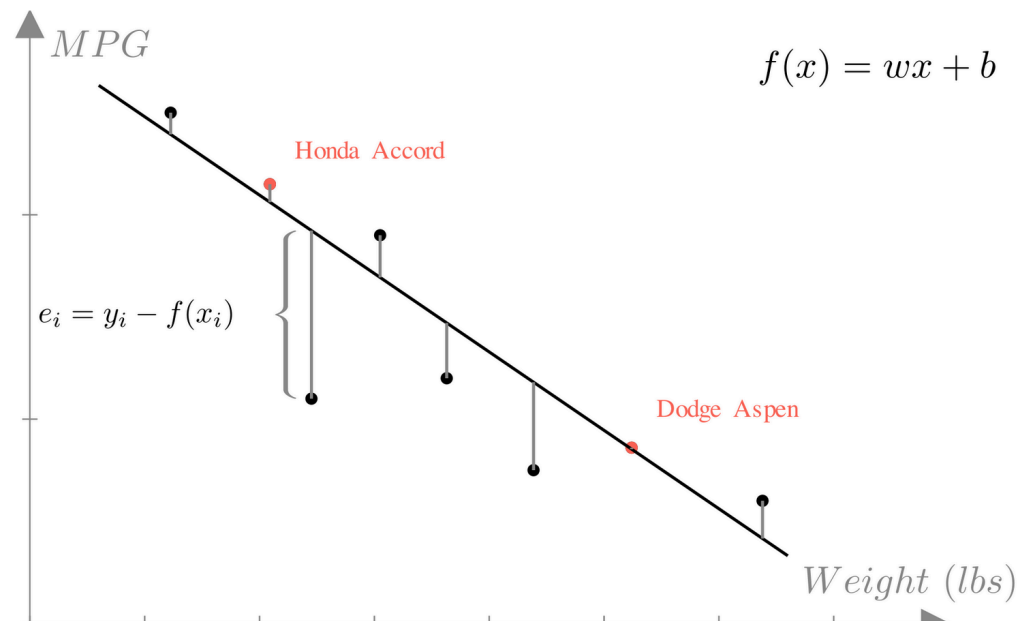
$$e_i = y_i - f(\mathbf{x}_i)$$

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots (\mathbf{x}_N, y_N)\}$$

$$f(x) = wx + b$$

Mean Squared error

$$MSE = \frac{1}{N} \sum_{i=1}^{N} (f(\mathbf{x}_i) - y_i)^2 = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i^T \mathbf{w} - y_i)^2$$

$N$: # of dos

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |f(x) - y_i|$$



MPG

Honda Accord

$e_i = y_i - f(x_i)$

Dodge Aspen

Weight (lbs)

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots (\mathbf{x}_N, y_N)\}$$

Vector
$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_m \end{bmatrix}$$

Weight    HP        Displacement

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{bmatrix} \begin{matrix} \text{Obs. 1} \\ \text{obs 2} \\ \\ \text{obs } N \end{matrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \begin{matrix} \text{Obs. 1} \\ \text{ob.2} \\ \\ \text{obs. } N \end{matrix}$$

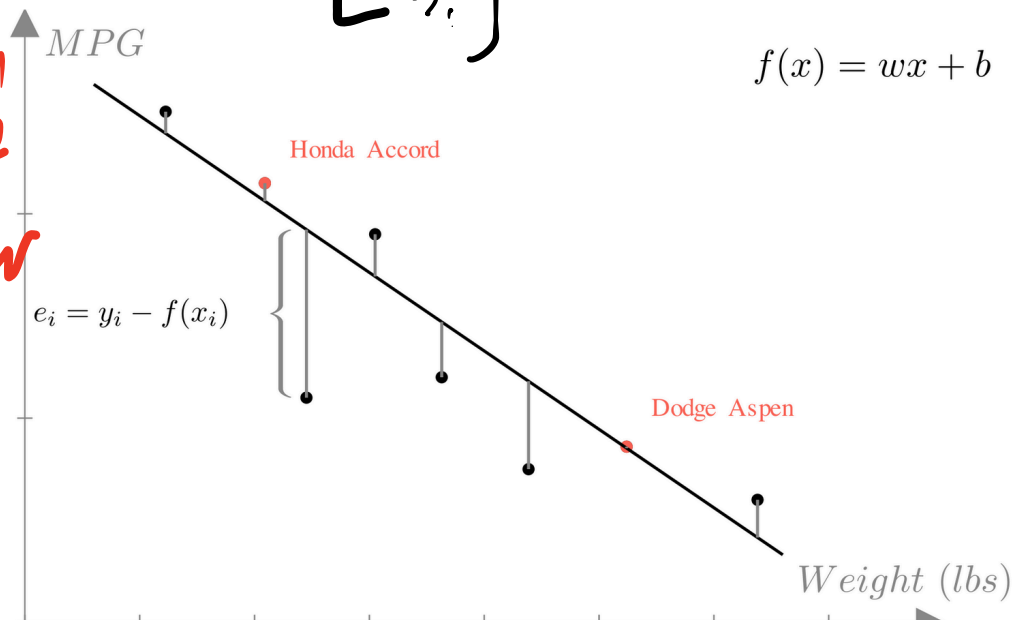Capital bold $X \implies$ matrix of all obs

$$MSE(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i^T \mathbf{w} - y_i)^2$$

Parameters / obs. Inputs / obs. outputs/Labels

pred.   label

$MPG$

$f(x) = wx + b$

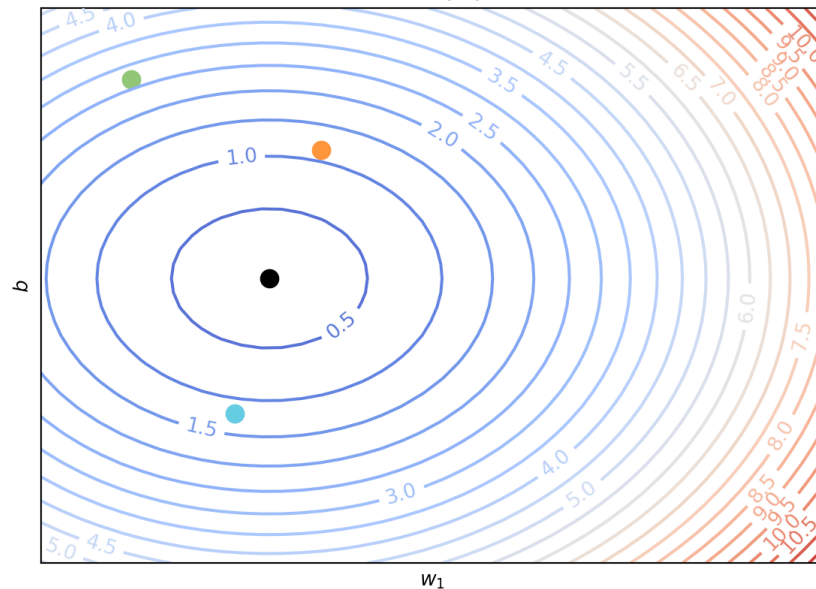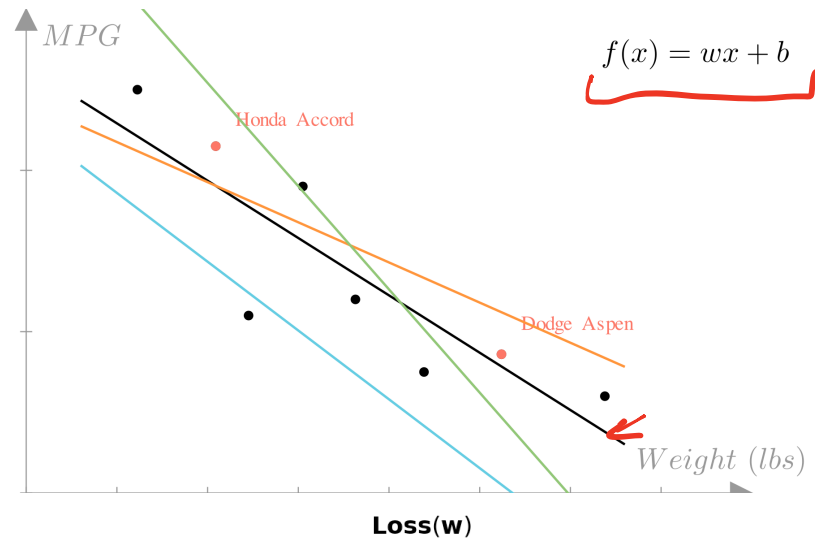Honda Accord

$e_i = y_i - f(x_i)$

Dodge Aspen

$Weight\ (lbs)$

$$\mathbf{Loss}(\mathbf{w}) = MSE(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \frac{1}{N}\sum_{i=1}^{N}(\mathbf{x}_i^T\mathbf{w} - y_i)^2$$

We can change
Fixed
$w$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \mathbf{Loss}(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{w}} \frac{1}{N}\sum_{i=1}^{N}(\mathbf{x}_i^T\mathbf{w} - y_i)^2$$
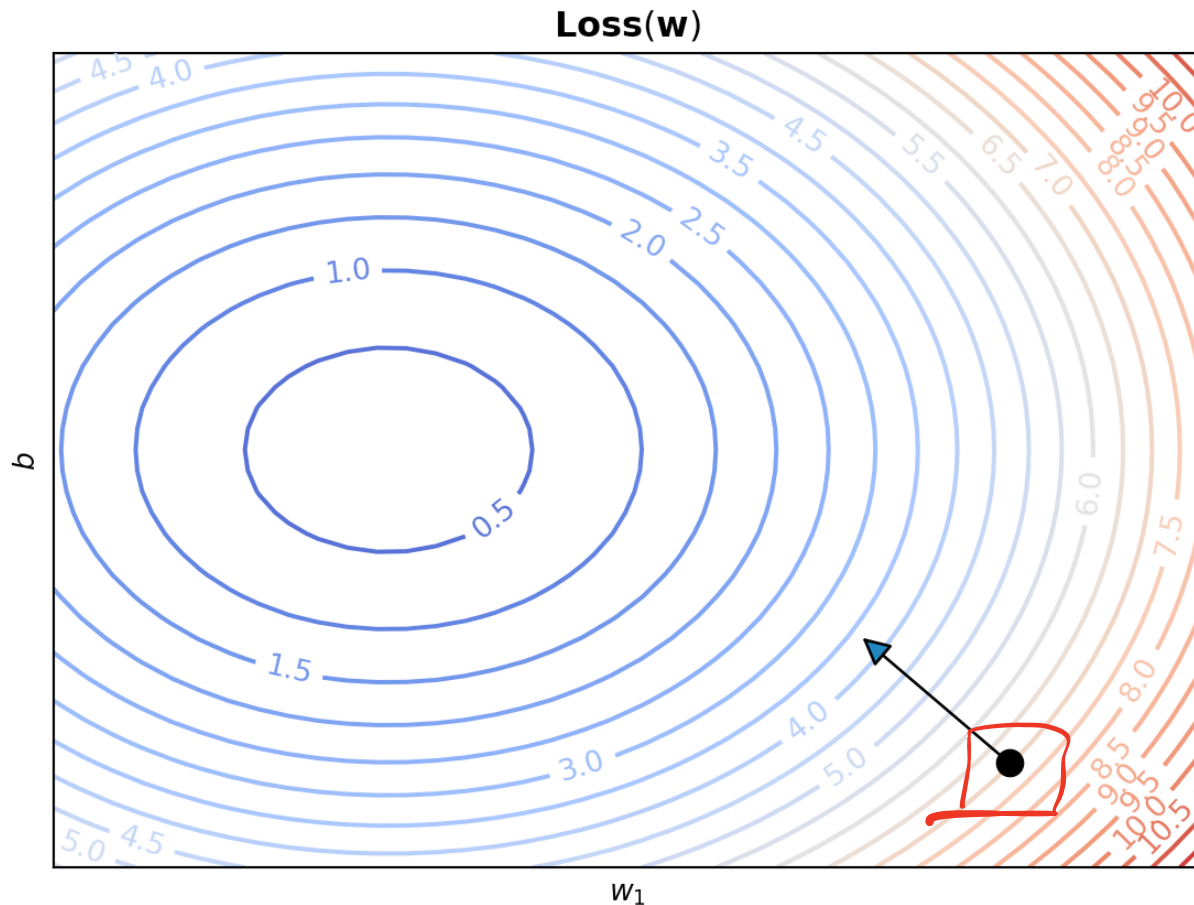


$f(x) = wx + b$

**Loss(w)**

Find: $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} f(\mathbf{w})$

Improved guess

$\mathbf{w}^{(1)} \leftarrow \boxed{\mathbf{w}^{(0)}} + \boxed{\mathbf{g}}$
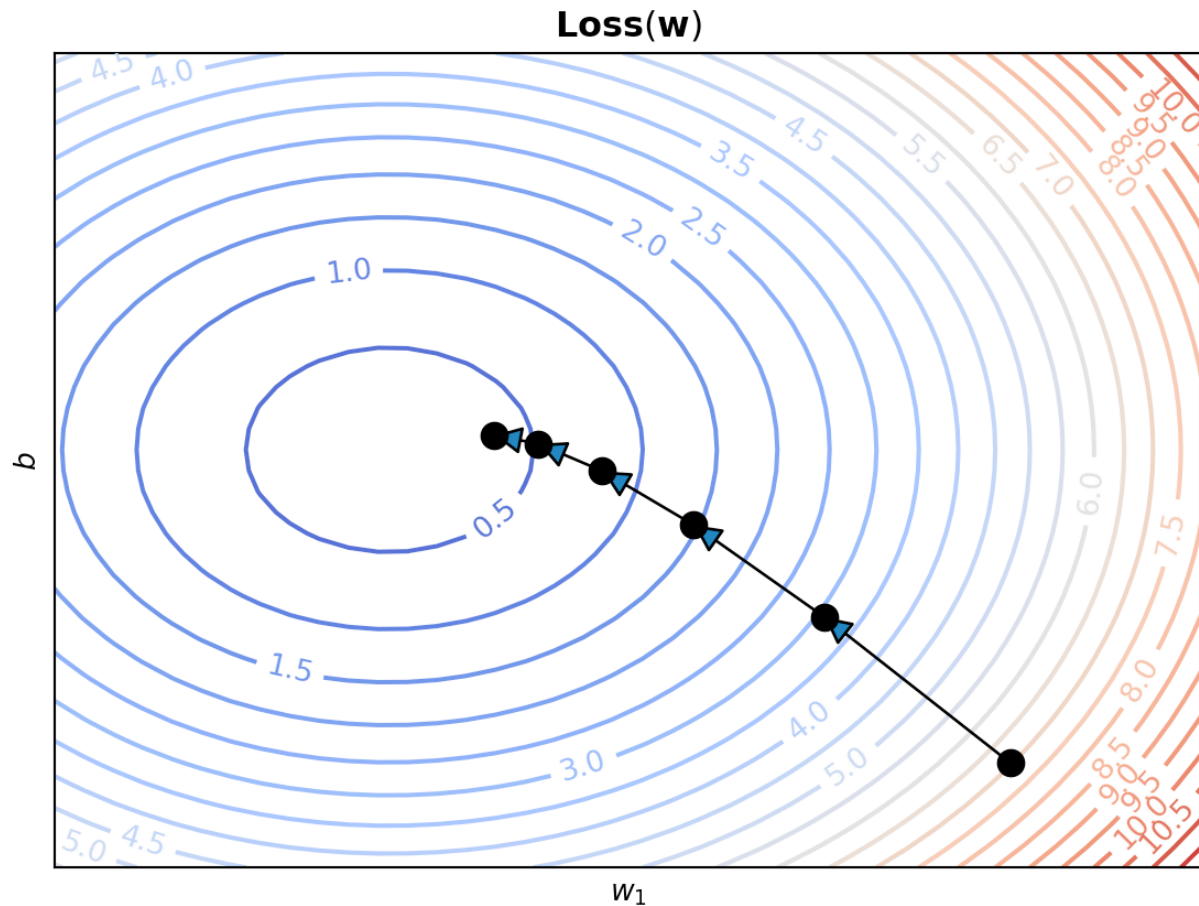
Initial guess

Direction of steepest descent



Loss(w)

Find: $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} f(\mathbf{w})$
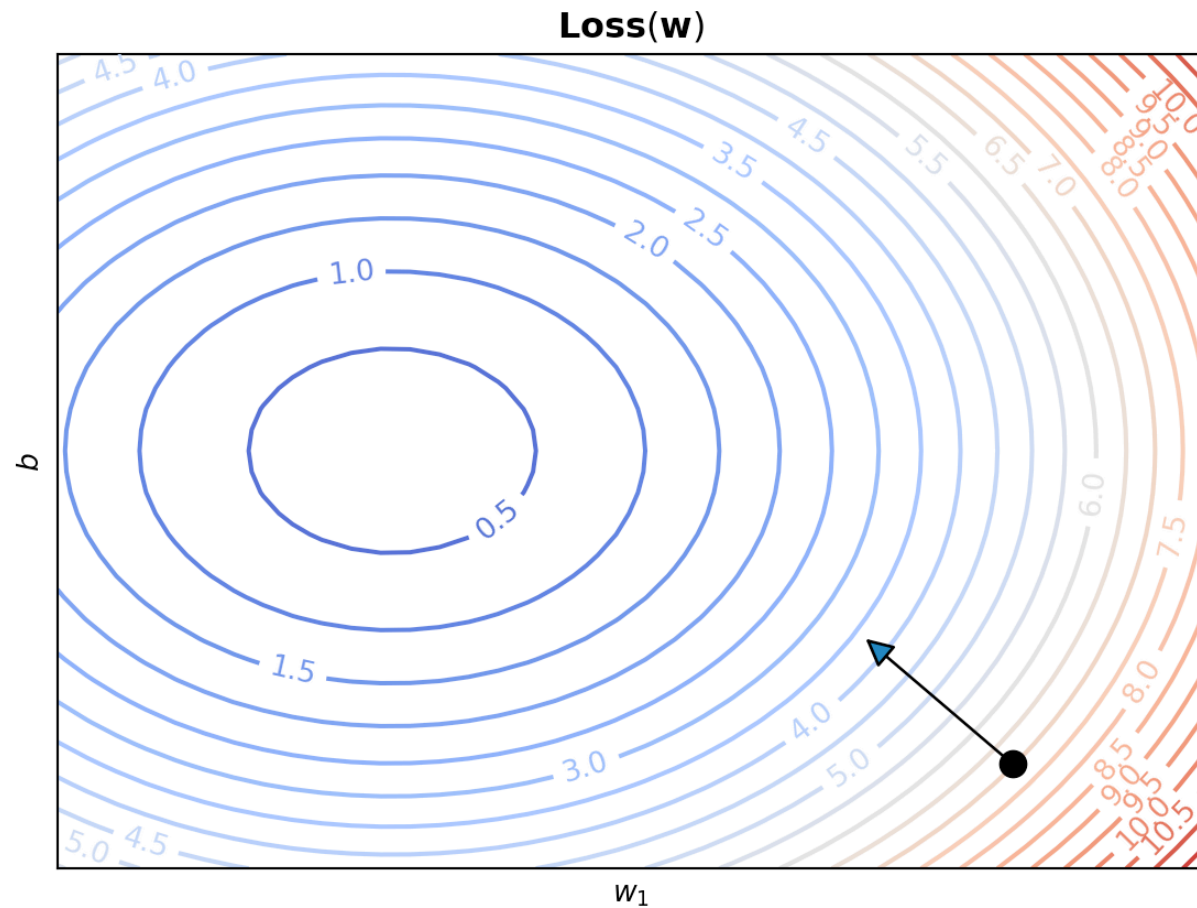
$$\mathbf{w}^{(1)} \leftarrow \mathbf{w}^{(0)} + \mathbf{g}$$

gradient descent



**Loss(w)**

Find: $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} f(\mathbf{w})$

$\mathbf{w}^{(1)} \leftarrow \mathbf{w}^{(0)} - \nabla f(\mathbf{w}^{(0)})$

**Loss(w)**

The **gradient** of a vector-input function is a vector such that each element is the partial derivative of the function with respect to the corresponding element of the input vector. We'll use the same notation as derivatives for gradients.

$$f(x)$$

$$\frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \\ \vdots \end{bmatrix}$$

$$f(x) = x^T x$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\frac{\partial}{\partial x_1} \mathbf{x}^T \mathbf{x} = \sum_{i=1}^{3} x_i \cdot x_i = \sum_{i=1}^{3} x_i^2 = x_1^2 + x_2^2 + x_3^2$$

$$= 2x_1$$