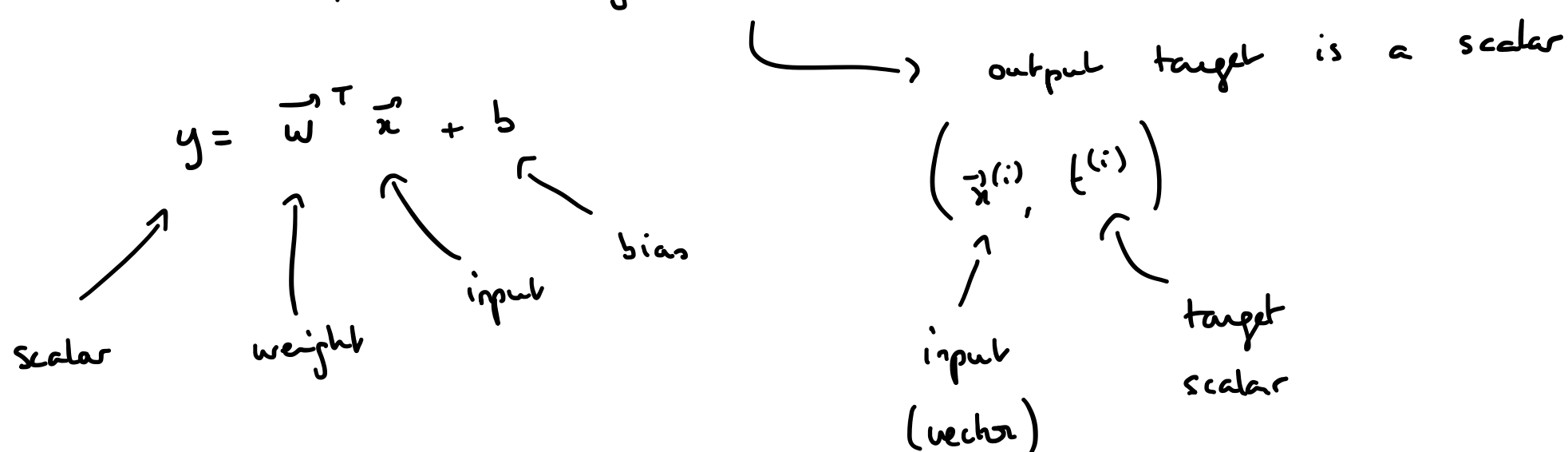


Gradient descent for linear regression



Loss function

(quantifies the error between the model's prediction & the target)

For one training example:

$$L(y, t) = \frac{1}{2} (y - t)^2$$

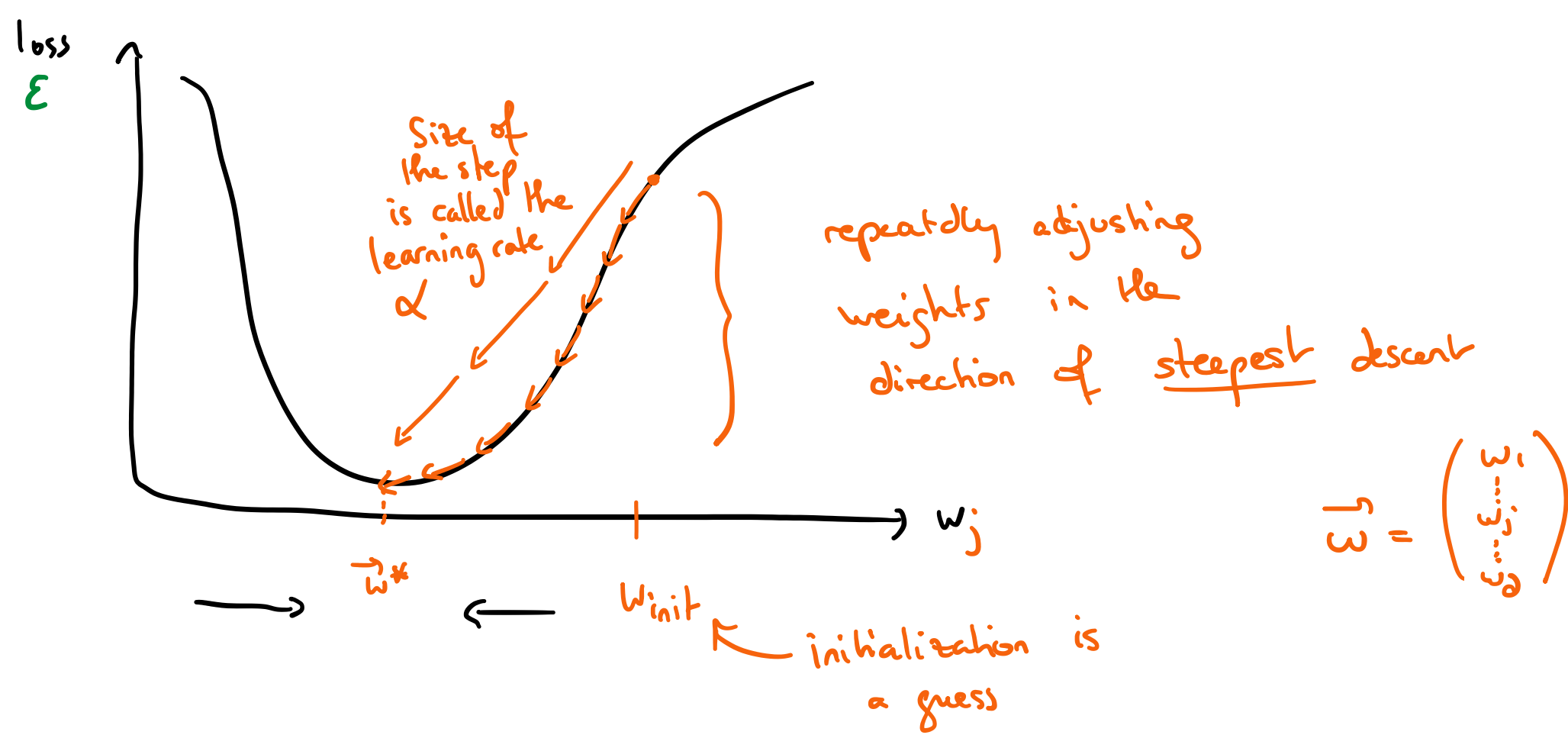
y : model prediction
 t : target
 for notation convenience

For entire dataset: N training examples

$$\mathcal{E}(w, b) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2$$

Gradient descent

- will allow us to minimize the loss function
- more broadly applicable (e.g., can be used w/ non-linear models)
- Grad descent is an iterative algorithm
 - \Rightarrow we are going to apply an update to the weights repeatedly
 - \Rightarrow until some criterion is met



$$\frac{\partial \mathcal{E}}{\partial w} < 0$$

If we increase w , we will decrease loss \mathcal{E}

$$\frac{\partial \mathcal{E}}{\partial w} > 0$$

If we increase w , we will increase the loss \mathcal{E}

direction of steepest descent

Update rule

$$w_j \leftarrow w_j - \alpha \frac{\partial \mathcal{E}}{\partial w_j}$$

If $\frac{\partial \mathcal{E}}{\partial w} > 0$, $-\frac{\partial \mathcal{E}}{\partial w} < 0$, we decrease w , we decrease the loss \mathcal{E}

If $\frac{\partial \mathcal{E}}{\partial w} < 0$, $-\frac{\partial \mathcal{E}}{\partial w} > 0$, we increase w , we decrease the loss \mathcal{E}

Learning rate

α : the larger it is, the faster \vec{w} changes.

α is a hyperparameter which we need to tune using our training/validation splits

Typically we use small values of $\alpha \approx 0.01$ or 0.0001

Update rule in vectorized form:

$$\vec{w} \leftarrow \vec{w} - \alpha \frac{\partial \mathcal{E}}{\partial \vec{w}}$$

$$\frac{\partial \mathcal{E}}{\partial \vec{w}} = \begin{pmatrix} \frac{\partial \mathcal{E}}{\partial w_1} \\ \frac{\partial \mathcal{E}}{\partial w_2} \\ \vdots \\ \frac{\partial \mathcal{E}}{\partial w_n} \end{pmatrix}$$

here model has n weights.

direction of fastest increase of \mathcal{E}