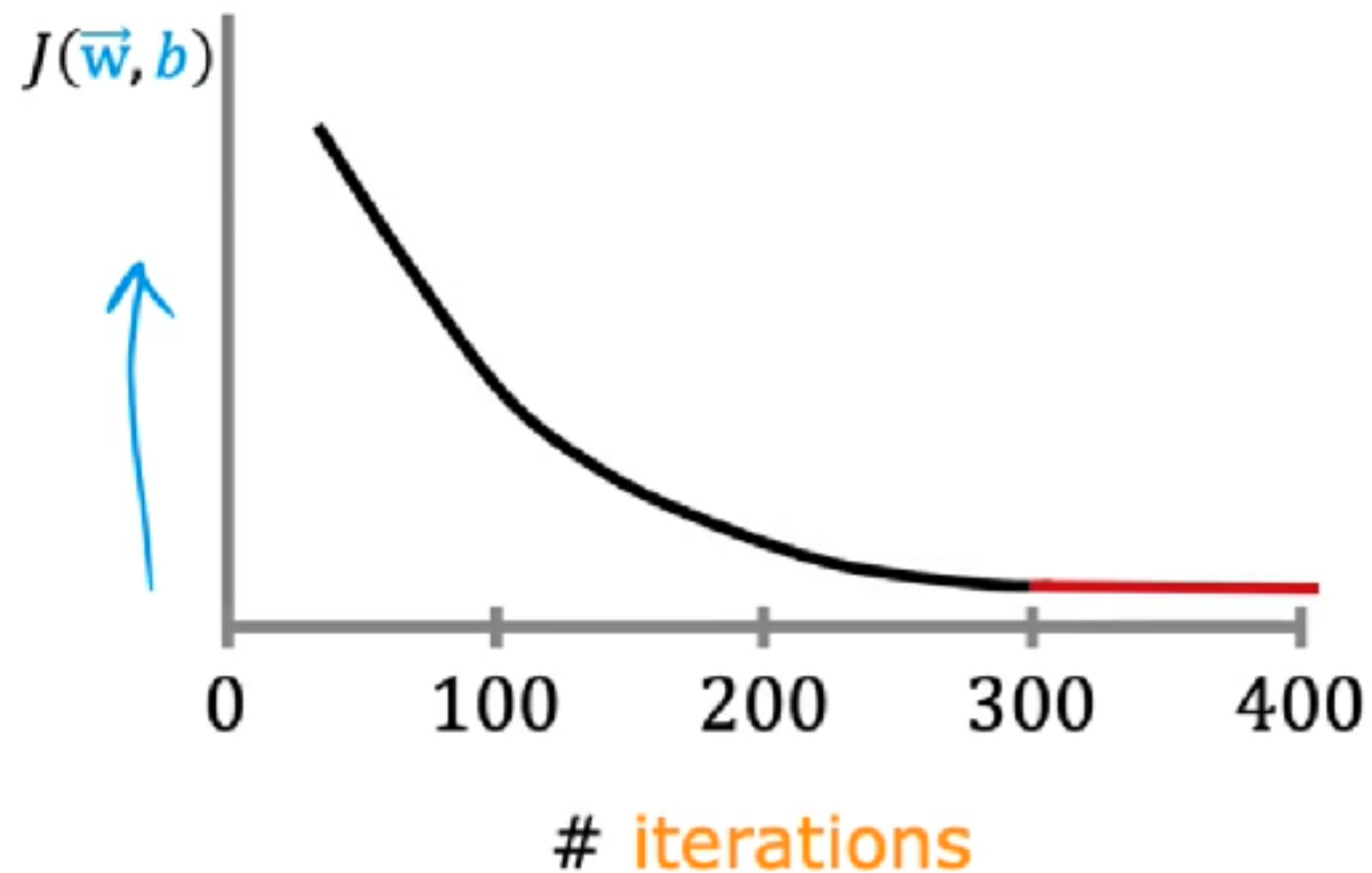


3_Checking Gradient Descent for Convergence

Make sure gradient descent is working correctly

objective: $\min_{\vec{w}, b} J(\vec{w}, b)$

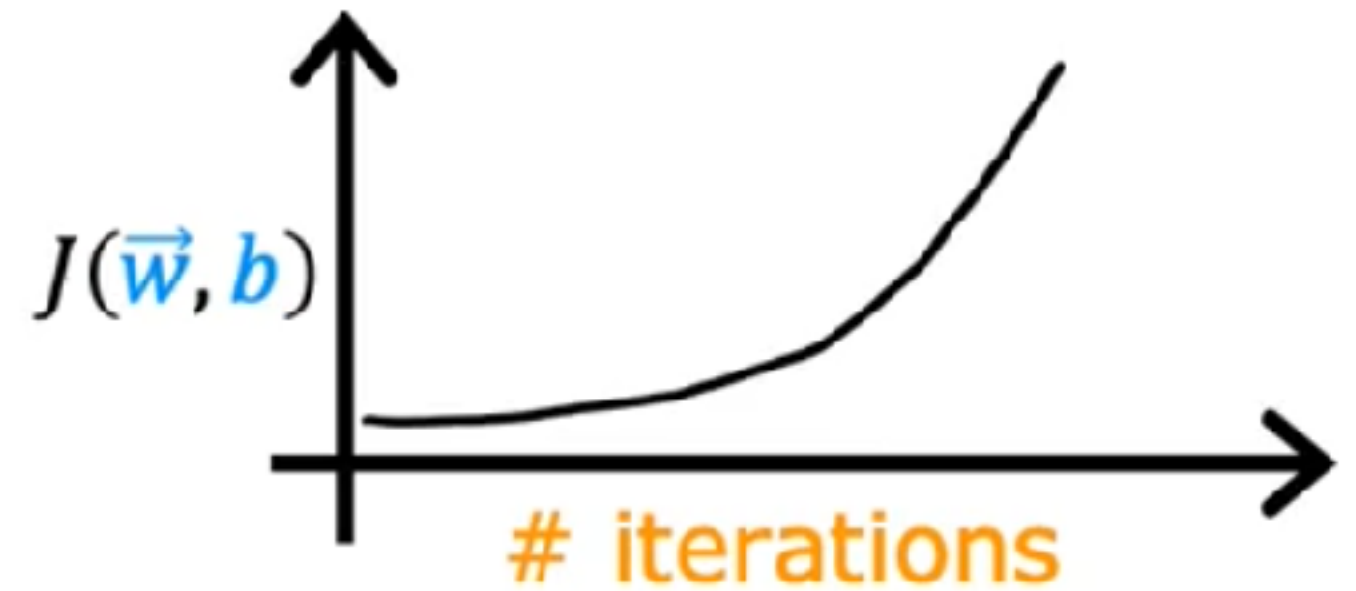
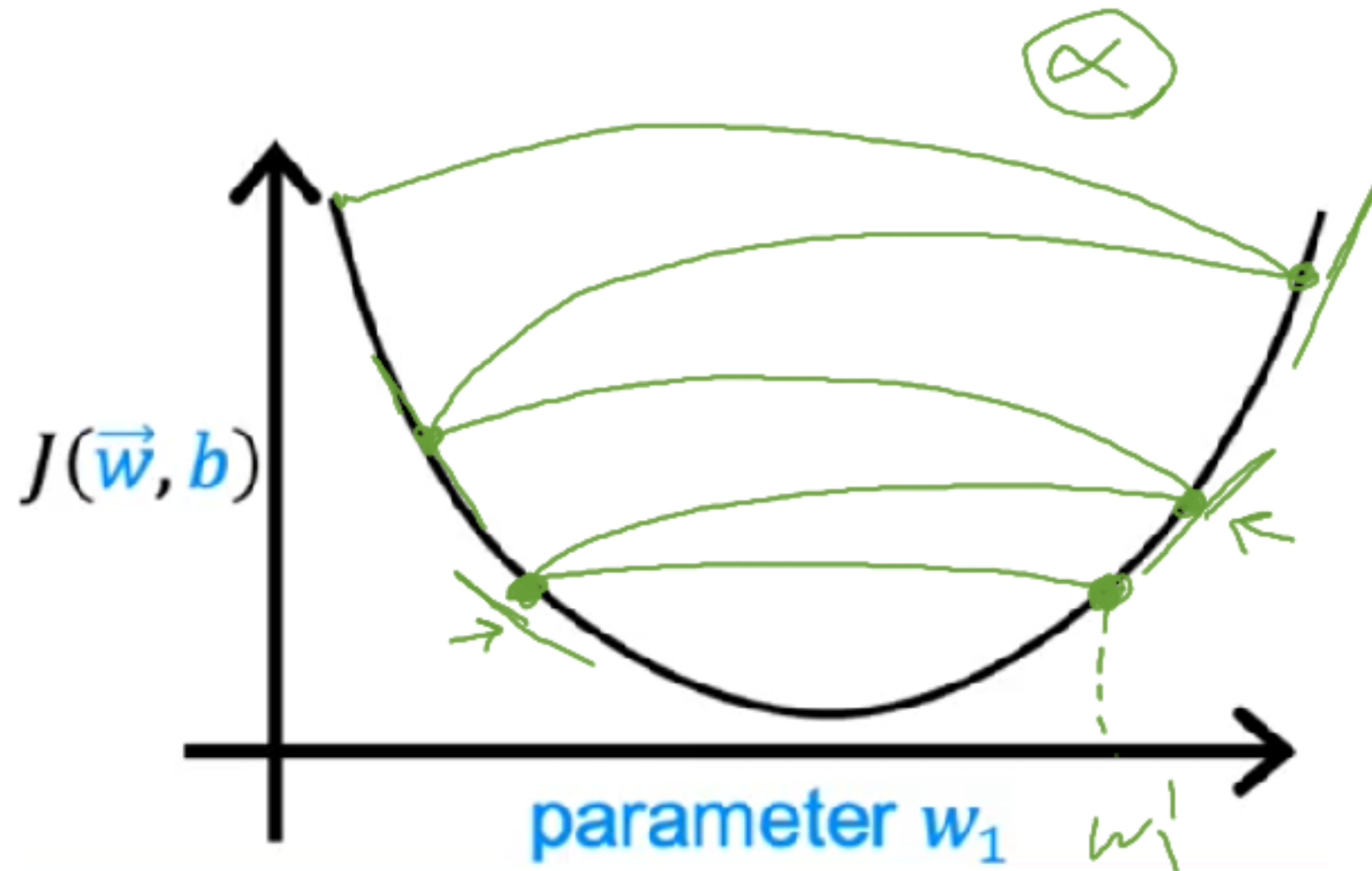


4. Choosing the Learning Rate

$$w^{\text{new}} = w^{\text{old}} - \alpha \frac{\partial J}{\partial w}$$

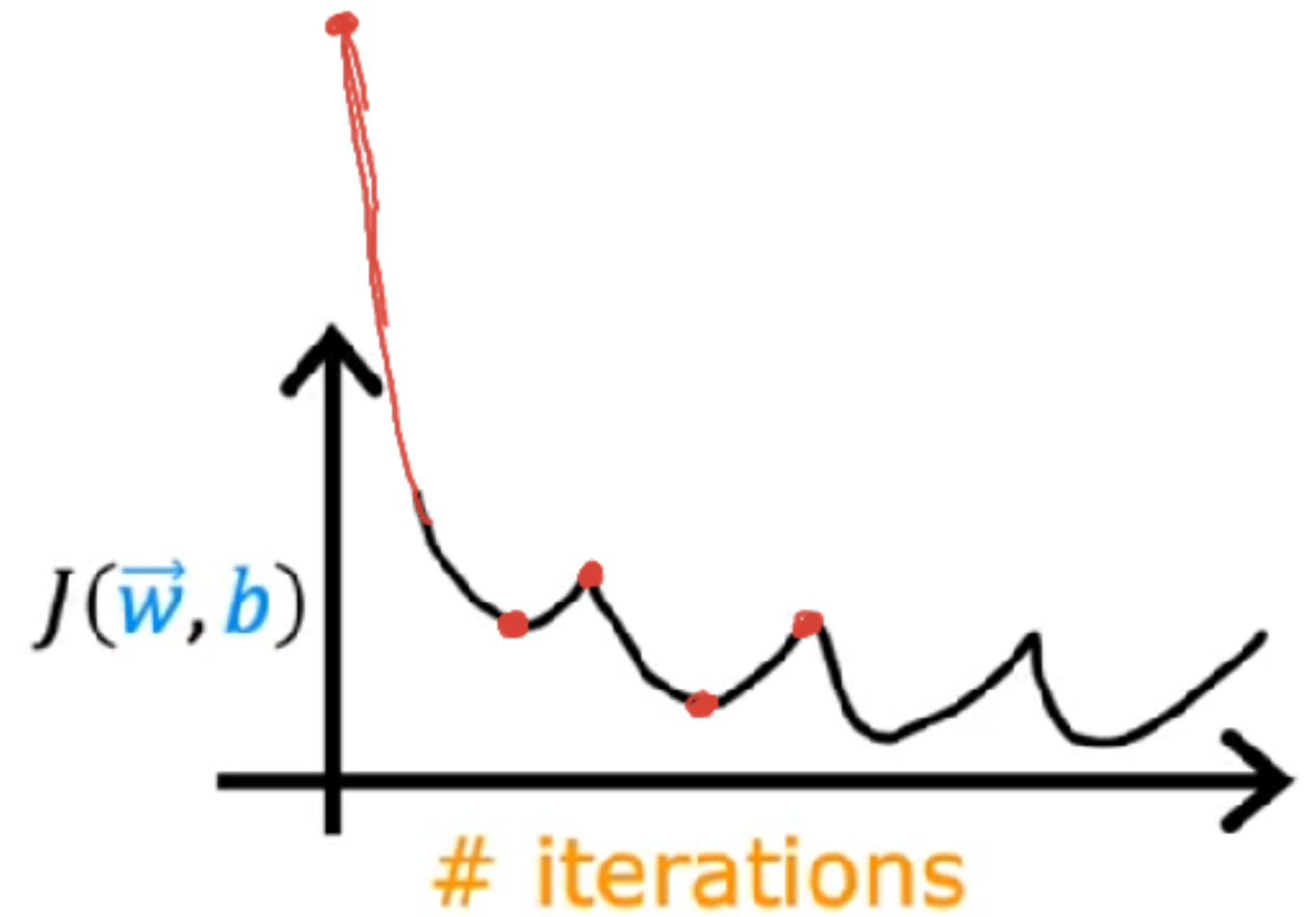
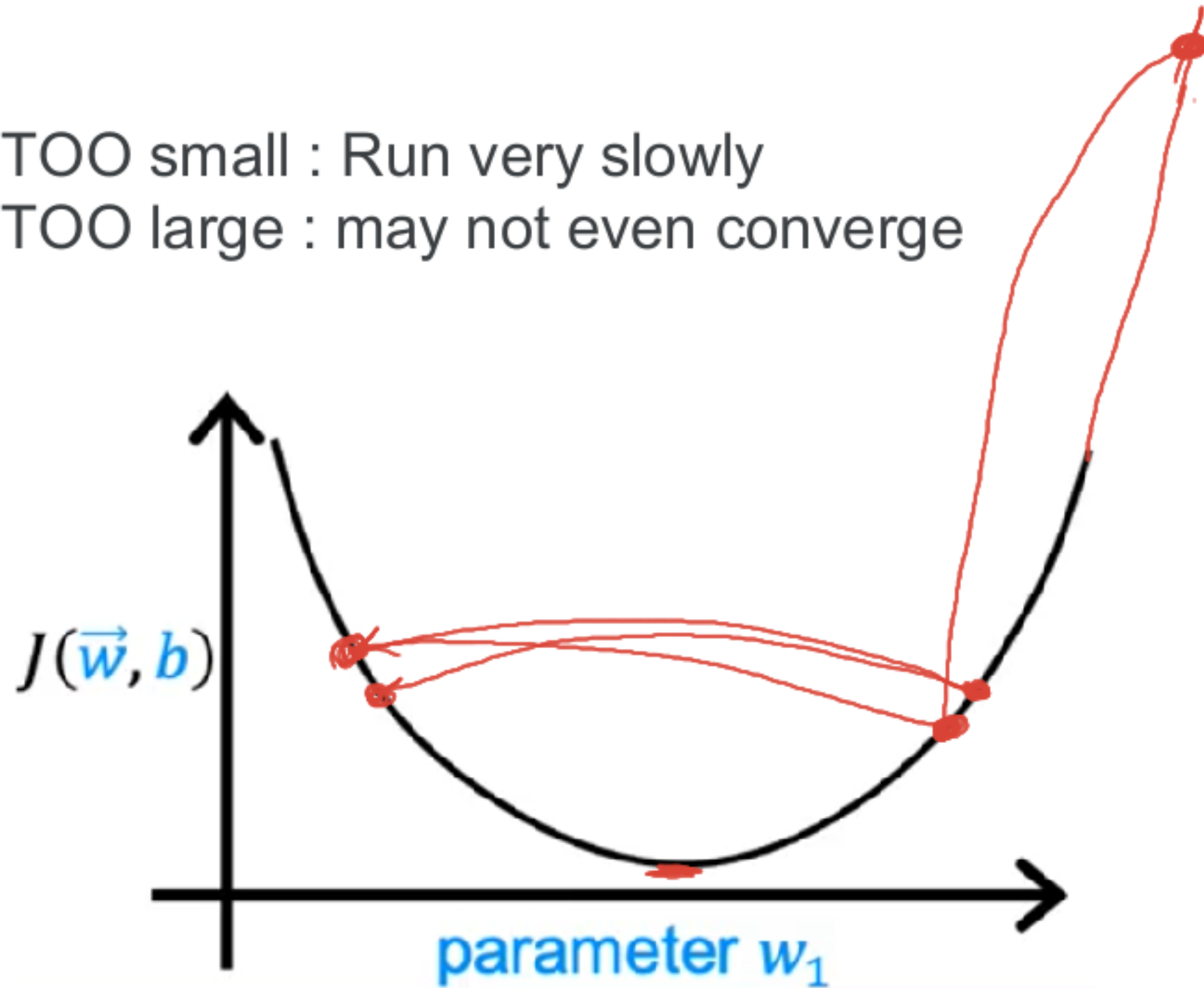
TOO small : Run very slowly

TOO large : may not even converge

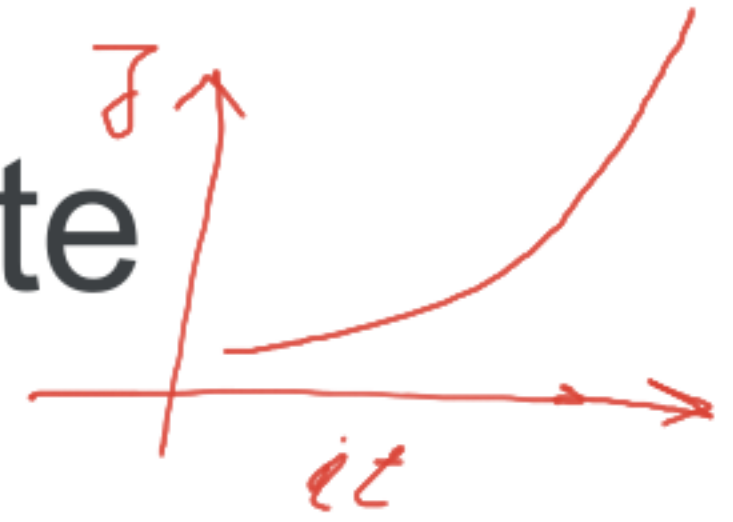


4. Choosing the Learning Rate

TOO small : Run very slowly
TOO large : may not even converge

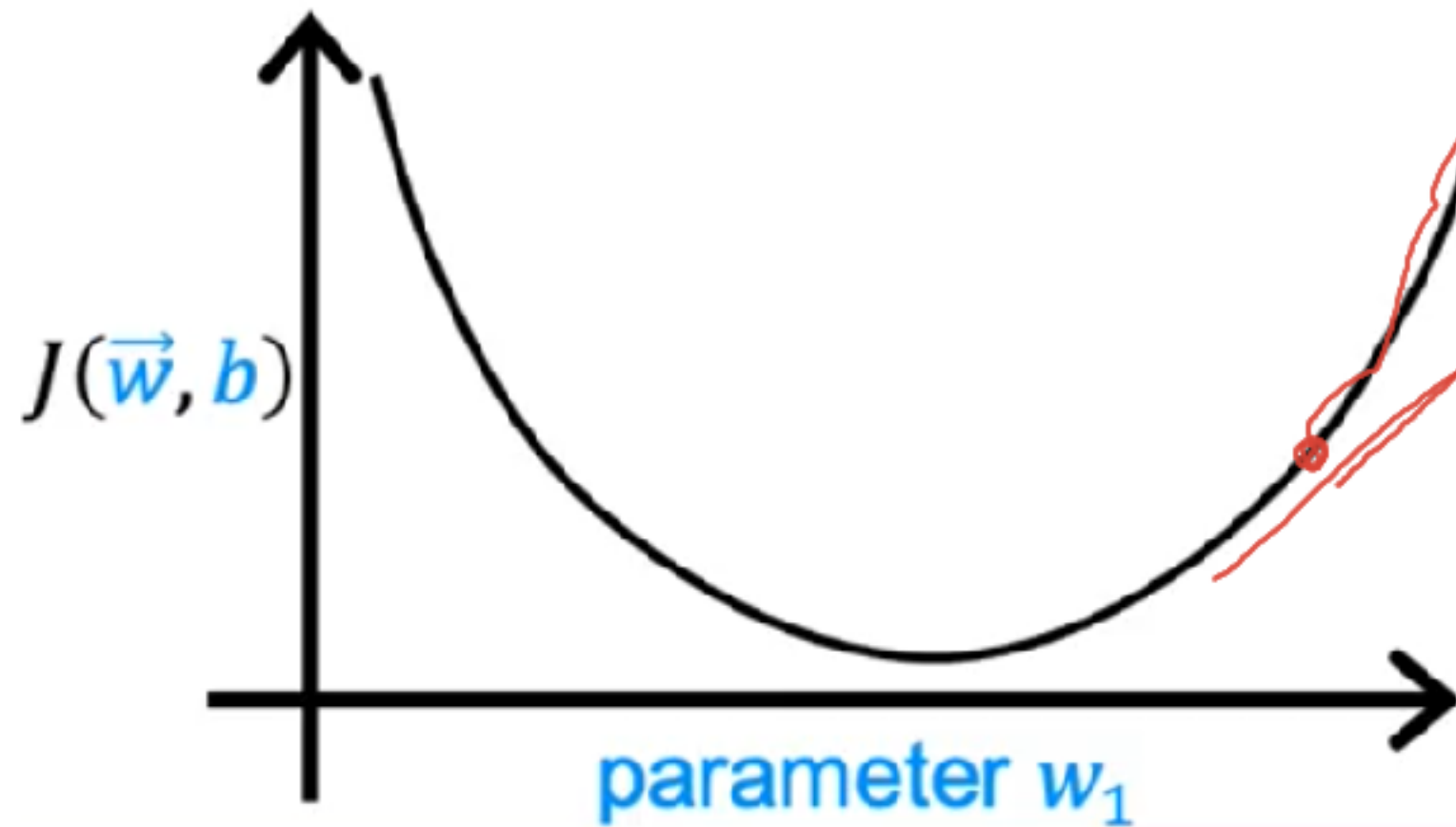


4. Choosing the Learning Rate



TOO small : Run very slowly

TOO large : may not even converge



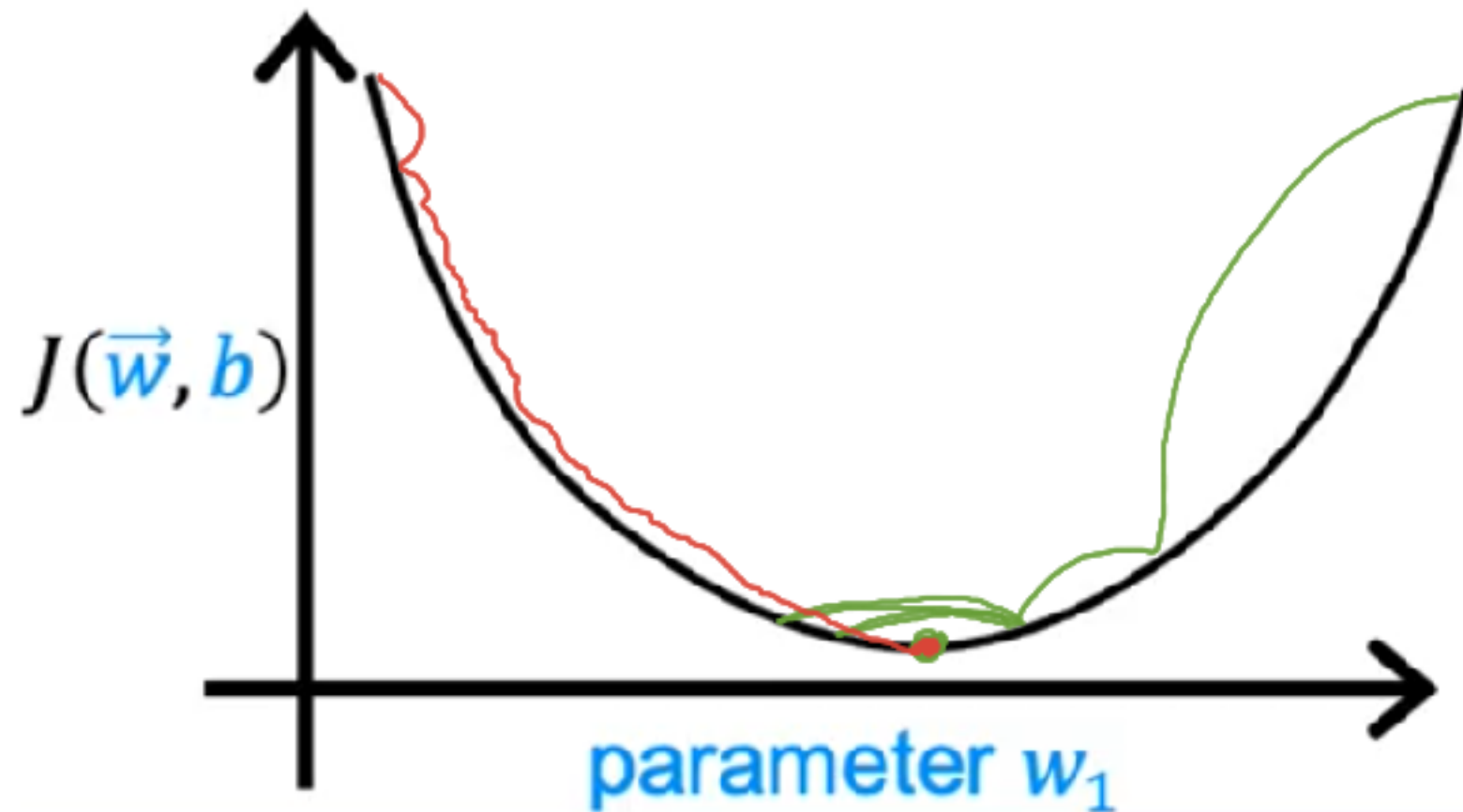
$$w^{new} = w^{old} - \alpha \frac{\partial J}{\partial w}$$

→ $w_1 = w_1 + \alpha d_1$ $\uparrow \downarrow$
use a minus sign
 $w_1 = w_1 - \alpha d_1$ $\uparrow \downarrow$

4. Choosing the Learning Rate

TOO small : Run very slowly

TOO large : may not even converge

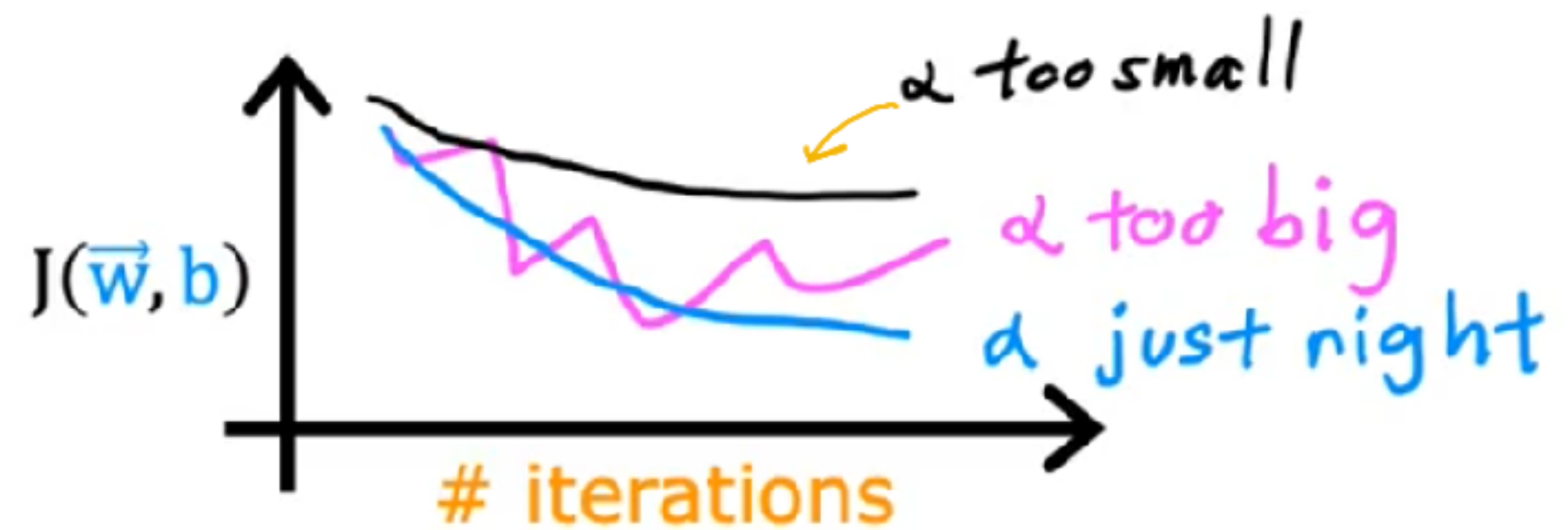
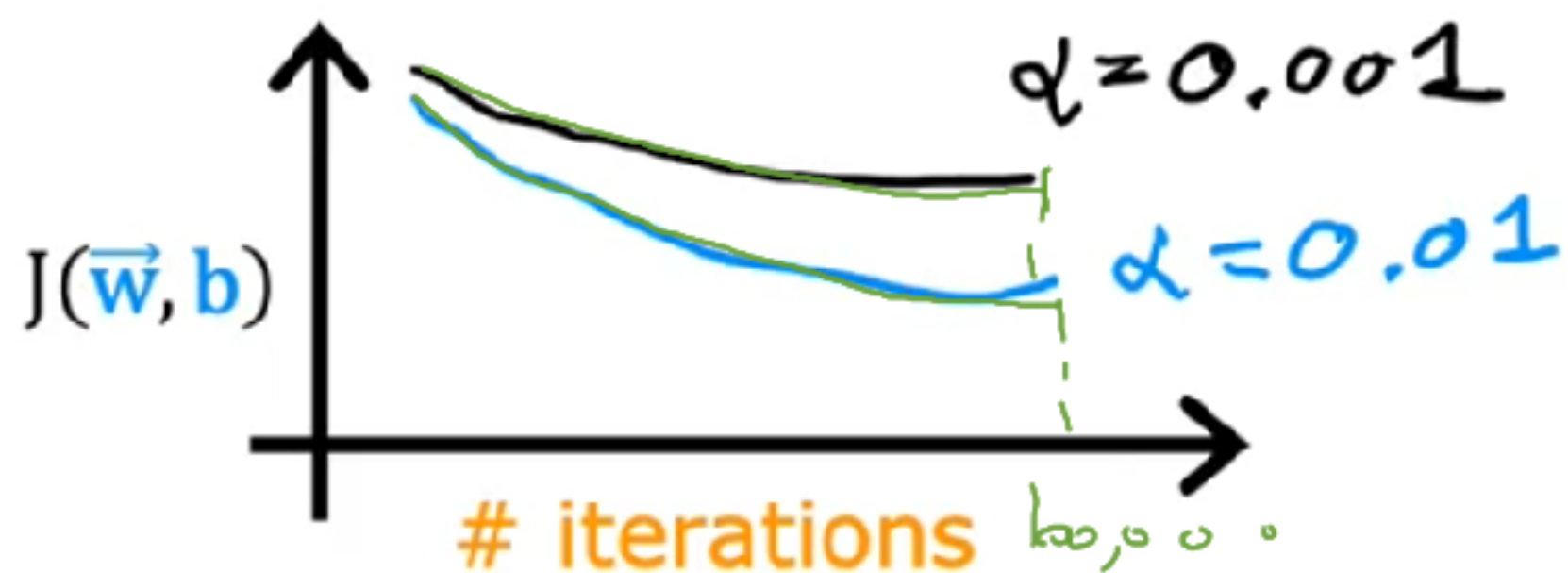


If α is too small,
gradient descent takes
a lot more iterations to
converge

Values of α to try:



... 0.001 0.003 0.01 0.03 0.1 0.3 1 ...
 \nearrow \nearrow \nearrow \nearrow \nearrow
 $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$

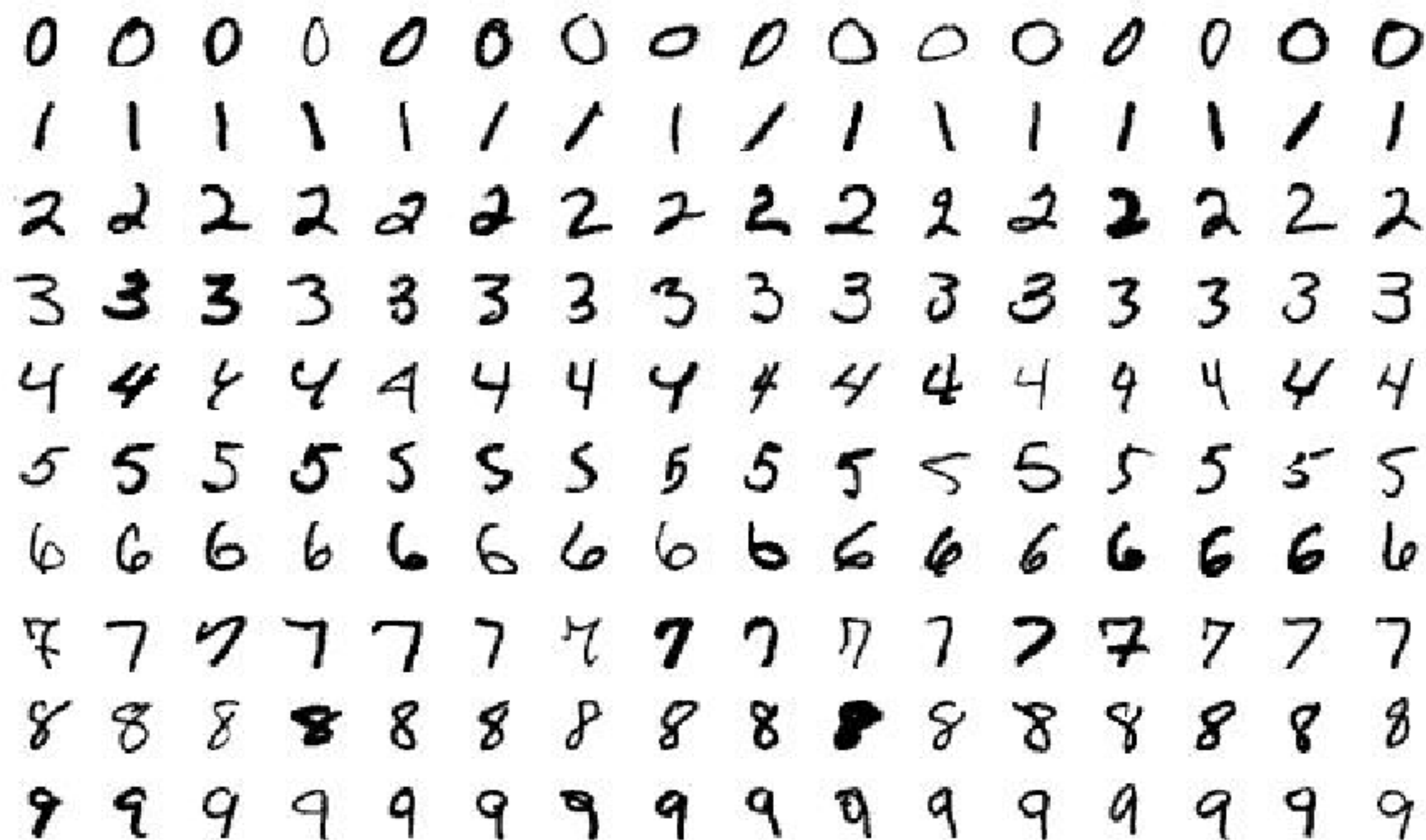


Feature Engineering

Choosing the right features is a critical step to making the algorithm work well.

MNIST Dataset

60,000 training images
10,000 testing images.



60,000 training images
10,000 testing images.

60,000 training images
10,000 testing images.

?

⑤ $28 \times 28 = 784$

③ $n \rightarrow 2000, 5$

[illegible]