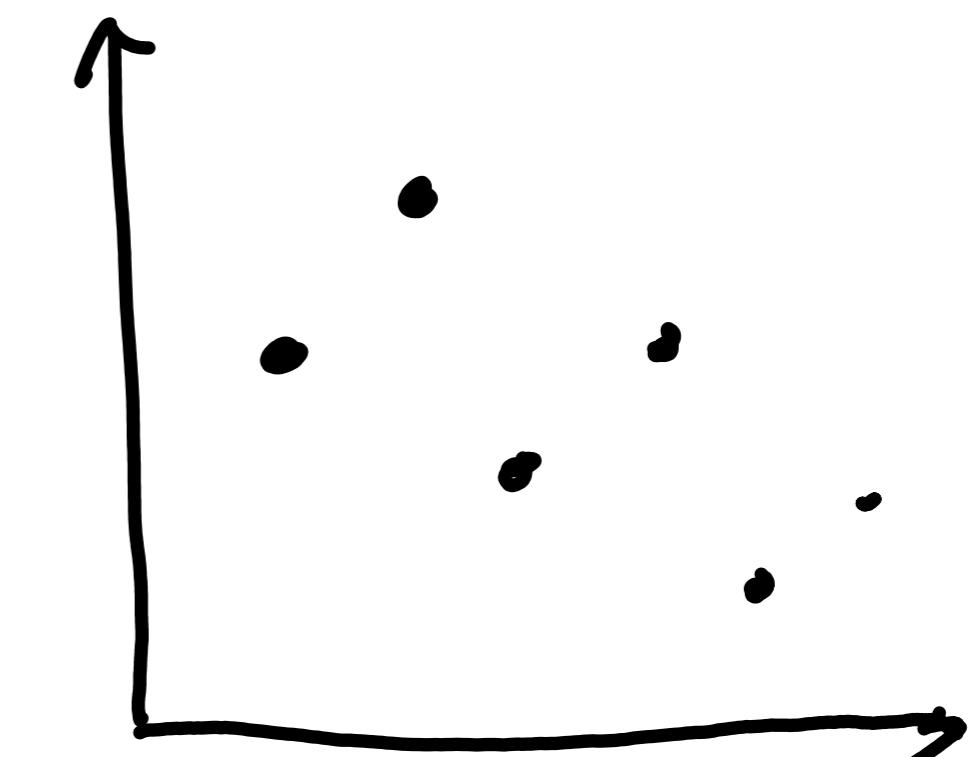
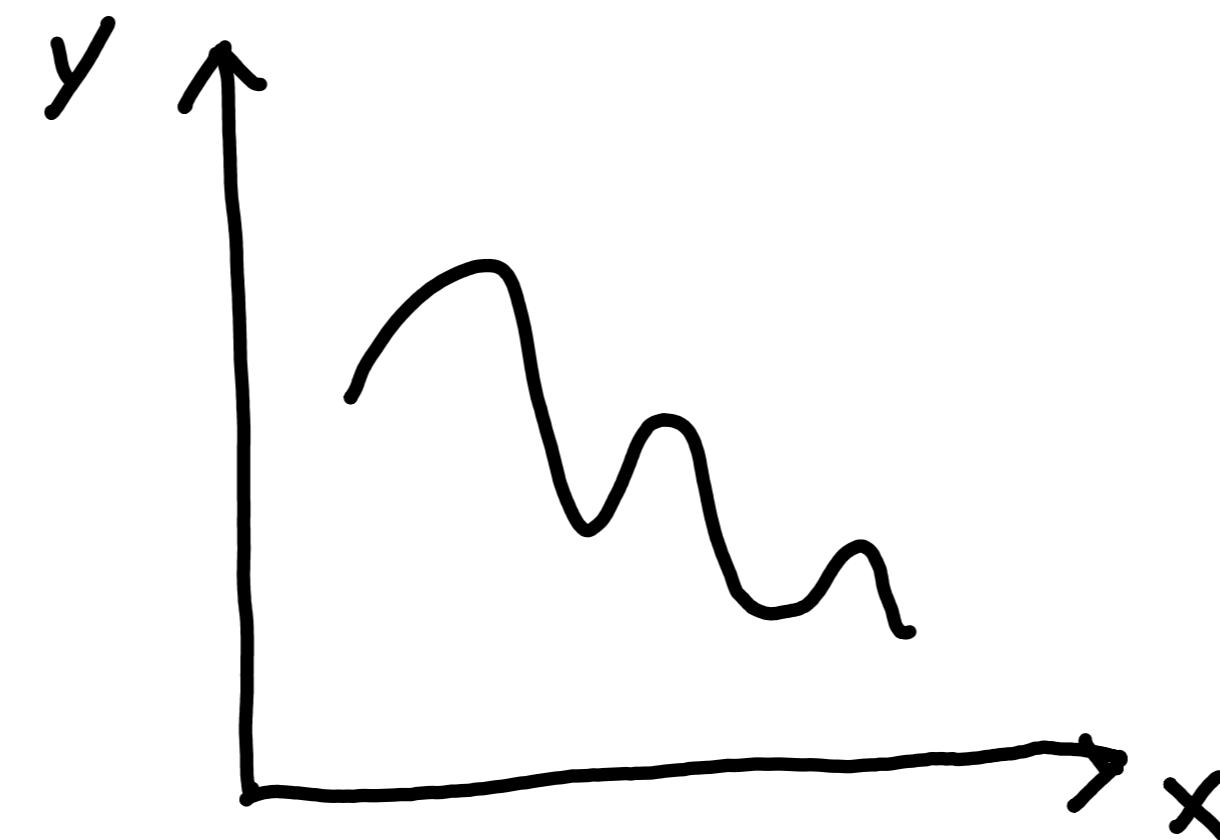
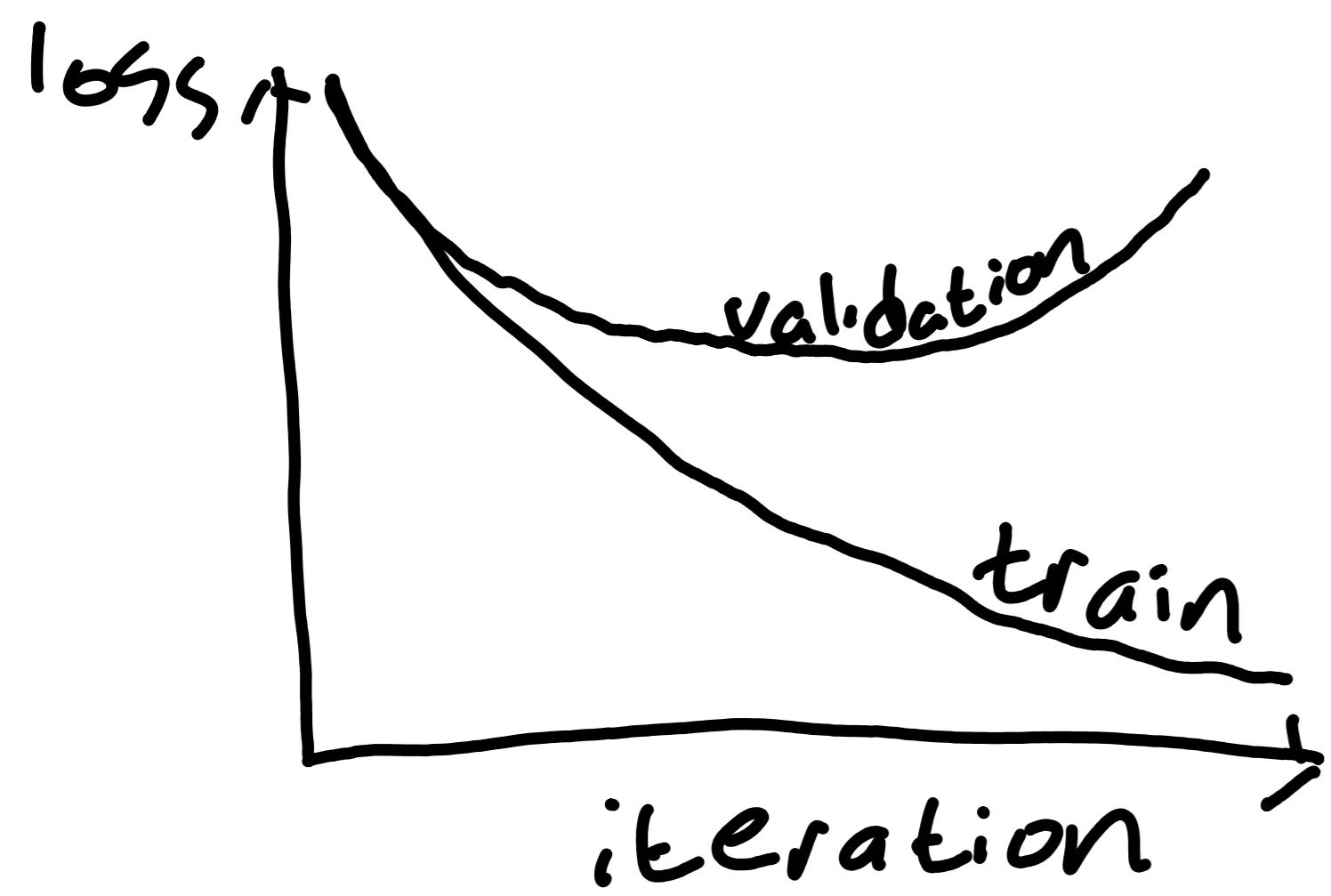


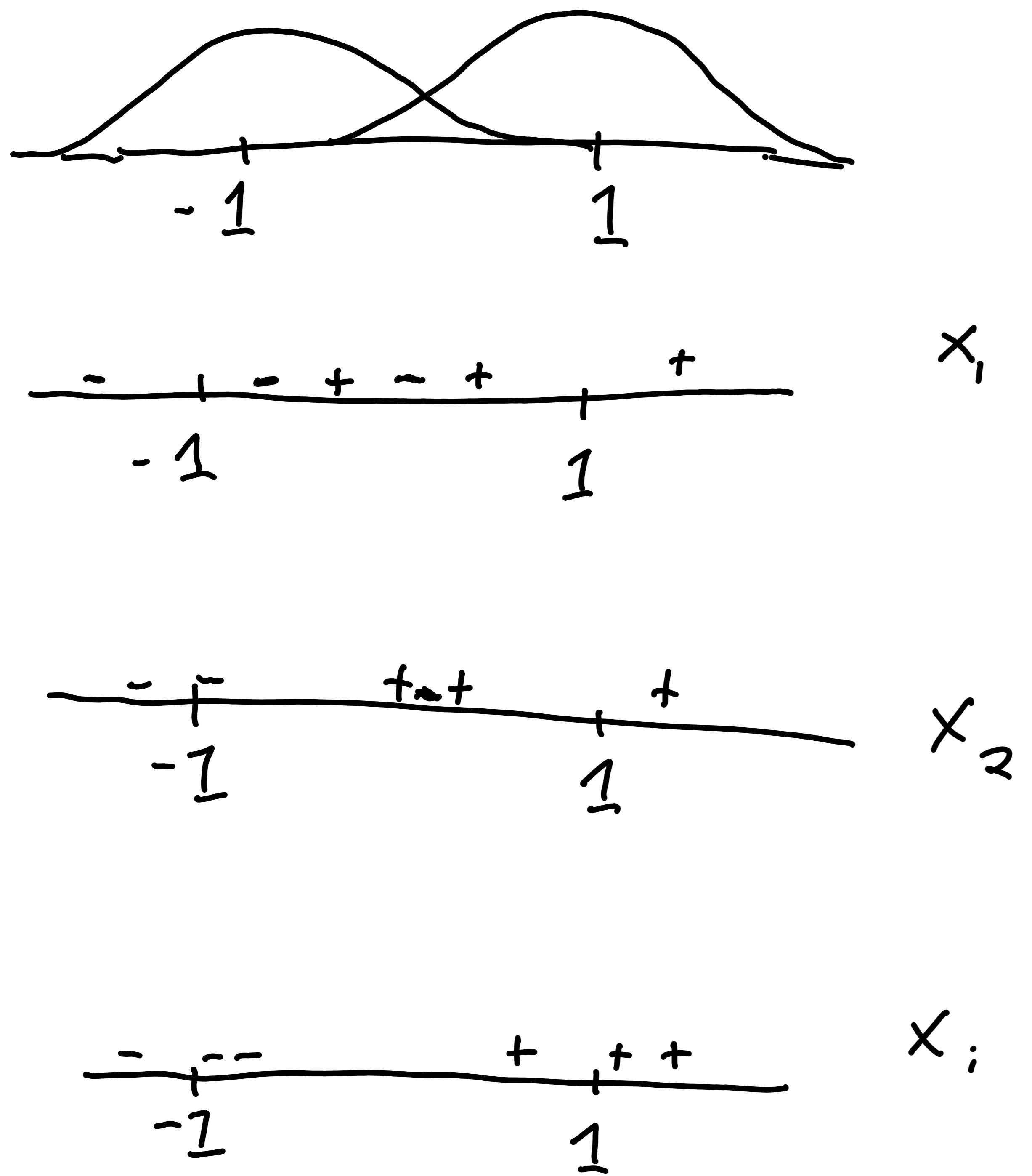
Overfitting and Regularization



How can we overfit in a linear model?

Say $x_i \sim N(y_i, \sigma)$ ($y_i \in \{-1, 1\}$)

We want to use $\hat{y} = w x$



As $\dim(x)$ grows
and the training data
shrinks, it becomes
increasingly likely that
some feature x_i
allows overfitting!

How can we avoid overreliance on a feature?

1. minimize $\|w\|_2^2$ by adding $\lambda \|w\|_2^2$ to our loss function

$$w - \eta \nabla_w (L + \lambda \|w\|_2^2) = \eta \nabla_w L + (1 - \eta \lambda \lambda) w$$

so we are performing "weight decay"

2. "Dropout" features at random:

During training,

$$x'_i = \begin{cases} 0 & \text{with probability } p \\ \frac{x_i}{1-p} & \text{with probability } 1-p \end{cases}$$

hyperparameter

$$\mathbb{E}[x'_i] = x_i$$

Multilayer Perceptron (MLP)

"linear model" \rightarrow an increase or decrease in a feature must always increase or decrease the output.

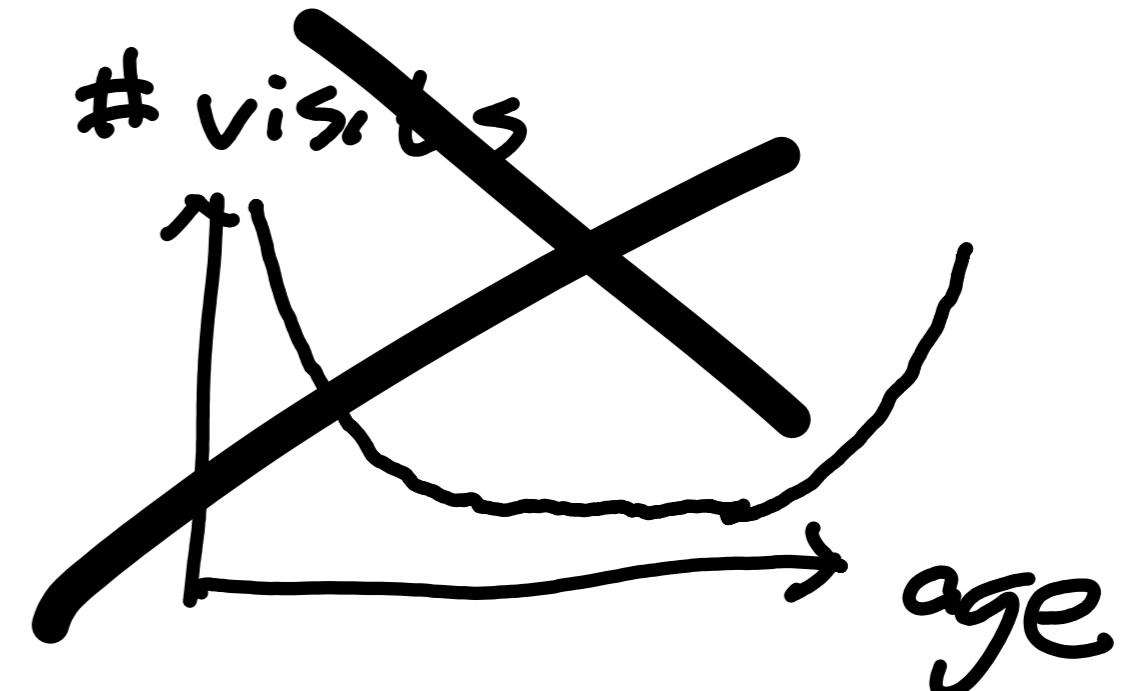
$$\hat{y} = w x + b$$

w positive: $x \uparrow, y \uparrow$

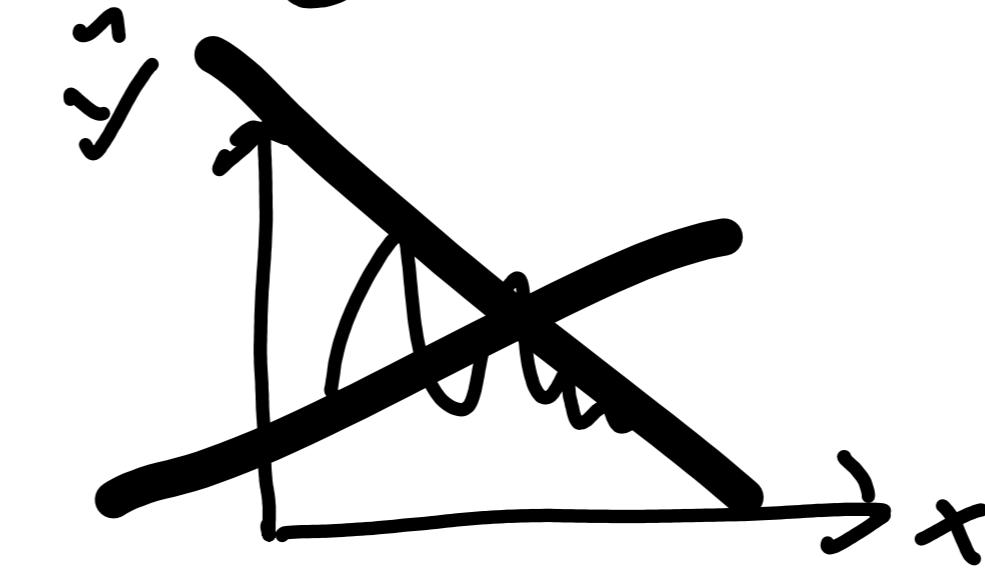
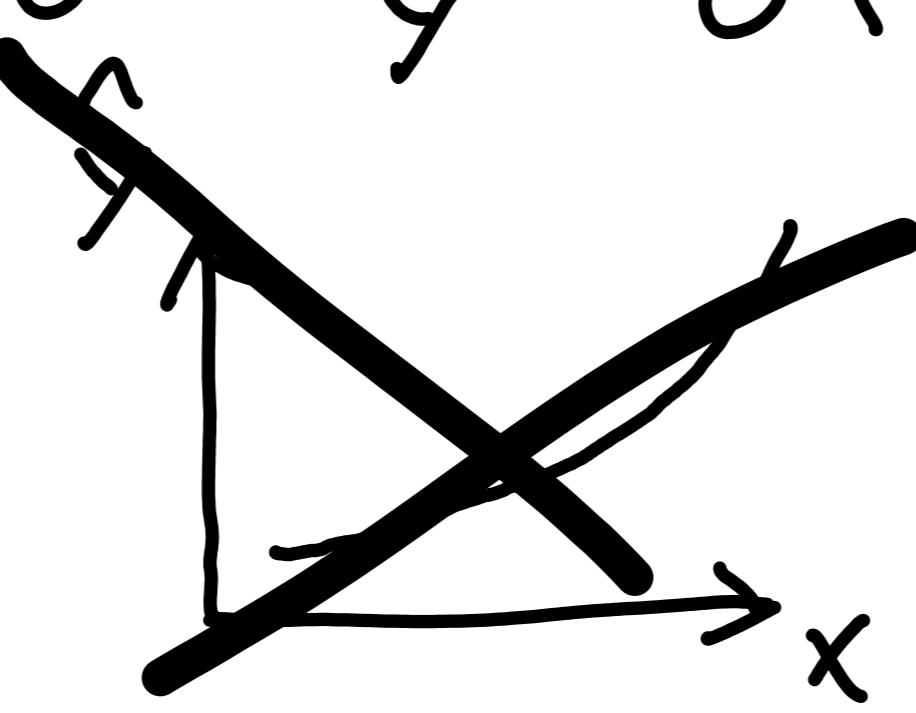
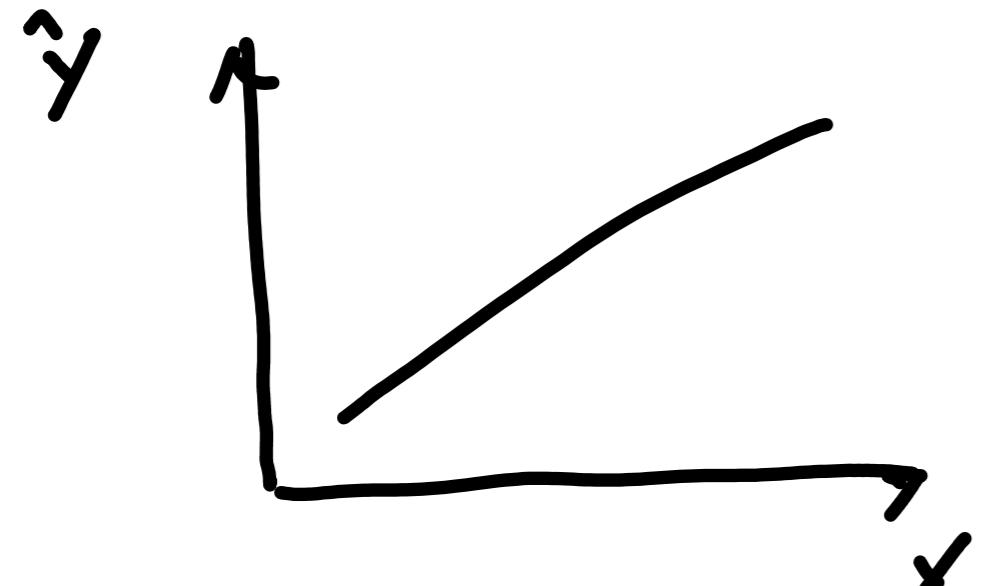
w negative: $x \downarrow, y \uparrow$

ex sqft \uparrow , value \uparrow

ex age \uparrow , # doctor visits



also, effect on \hat{y} of changes in x is always the same.



How can we get more complex functions?

Stack multiple transformations.
(compose multiple functions)

Let's try:

$$\text{hidden} = \text{input} \times w_i + b_i$$

$$\text{output} = H w_2 + b_2$$

$$O = (X_w + b) w_2 + b_2$$

A diagram illustrating a mathematical relationship between two parallel horizontal lines. The top line features two wavy lines labeled w_1 and w_2 positioned above it, and one wavy line labeled w positioned below it. The bottom line features two wavy lines labeled b_1 , w_2 and $+ b_2$ positioned above it, and one wavy line labeled b positioned below it.

$$= x^w + b$$

Need nonlinearity! Apply nonlinear functions
(almost always elementwise)

$$H = \sigma(XW_1 + b_1) \quad O = Hw_2 + b_2$$

↳ nonlinear function!

"element wise"

$$h_i = \sigma(x_i)$$

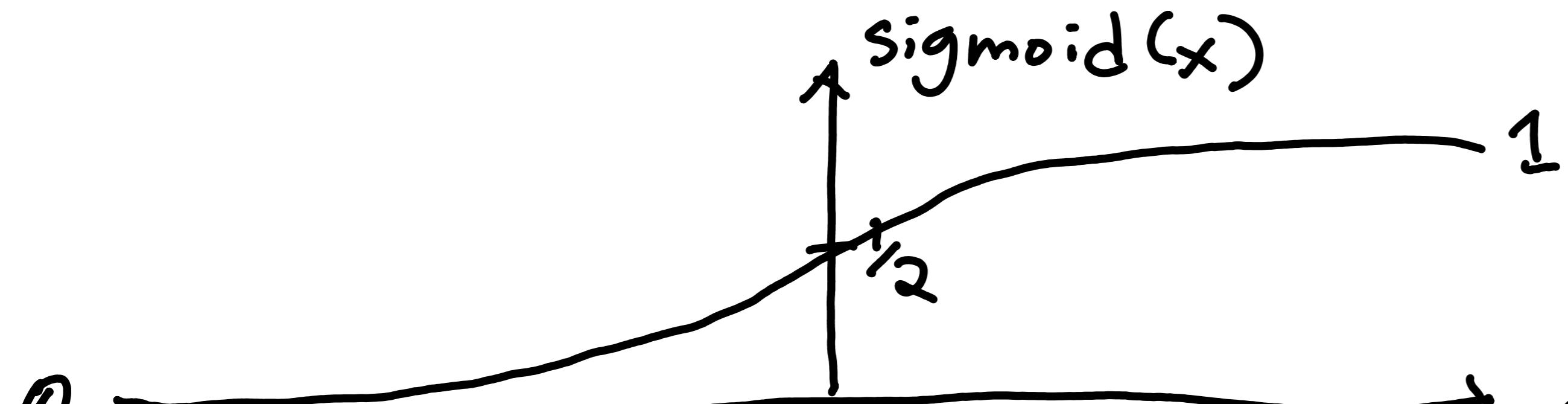
"Neural network"

"Multi-layer perceptron"

"Dense / feed-forward / fully connected"
(layer or network)

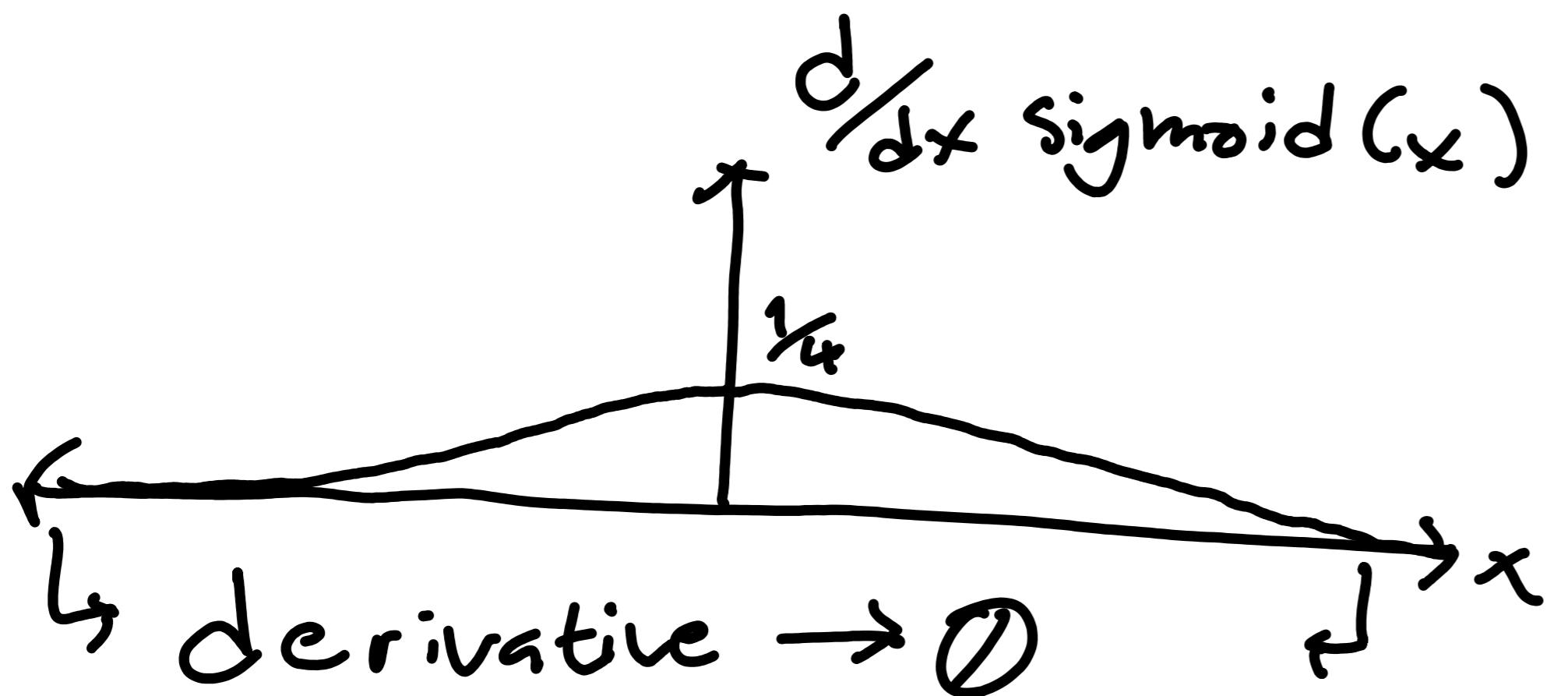
What to use for $\sigma(x)$?

$$\text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$



$$\frac{d}{dx} \text{sigmoid}(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2}$$

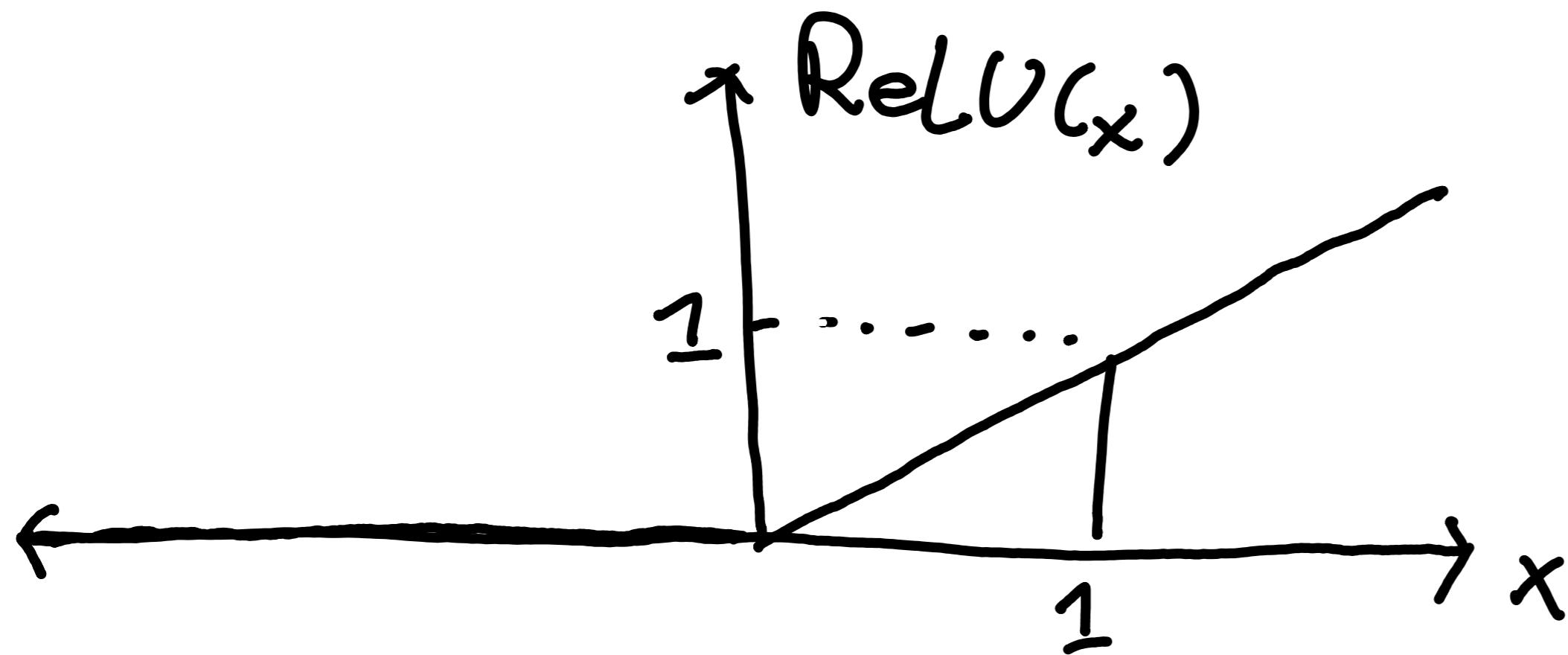
$$= \text{sigmoid}(x)(1 - \text{sigmoid}(x))$$



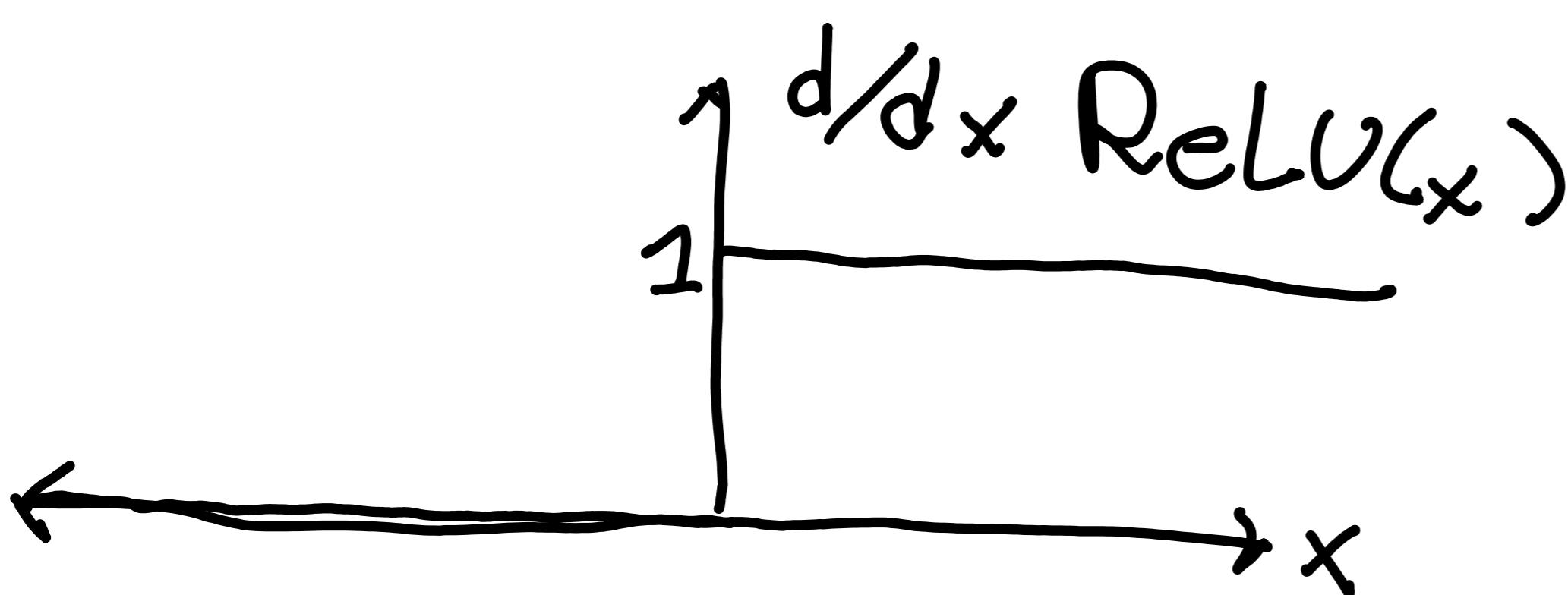
$$\begin{aligned} & 1 - \text{sigmoid}(x) \\ &= 1 - \frac{1}{1 + \exp(-x)} \\ &= \frac{1 + \exp(-x)}{1 + \exp(-x)} - \frac{1}{1 + \exp(-x)} \\ &= \frac{\exp(-x)}{1 + \exp(-x)} \end{aligned}$$

ReLU - "Rectified Linear Unit"

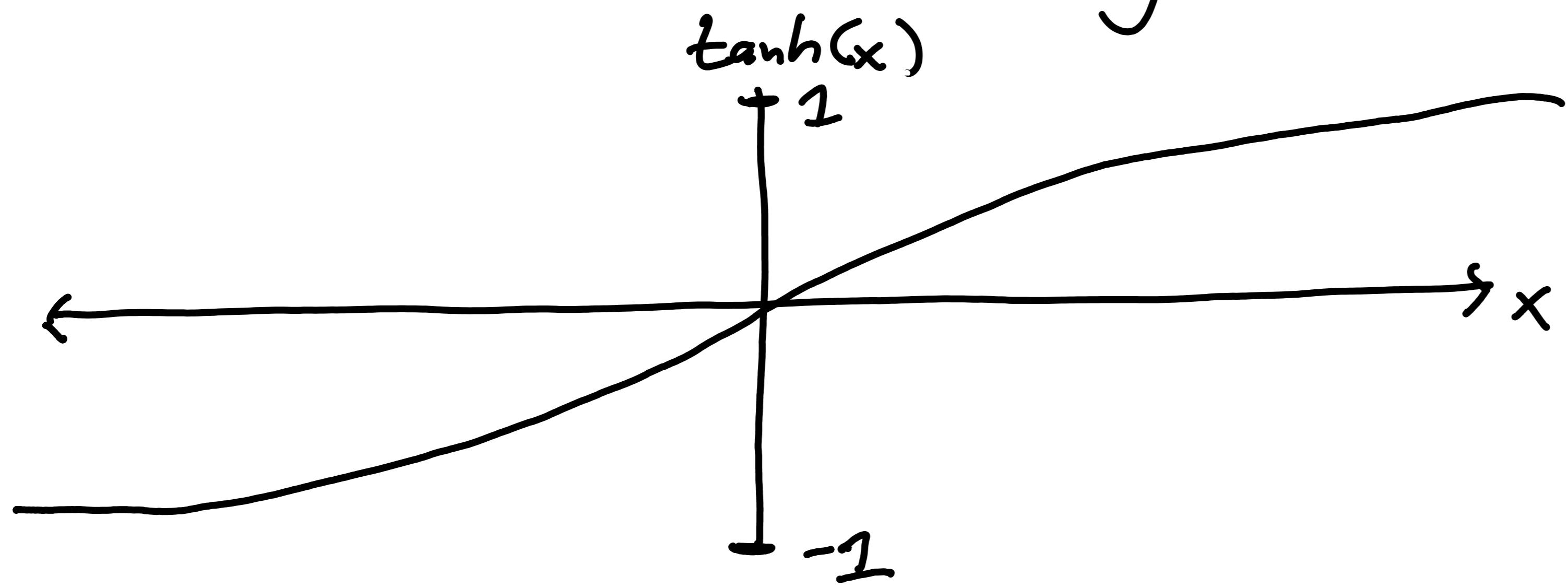
$$\text{ReLU}(x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases} = \max(0, x)$$



$$\frac{d \text{ReLU}}{dx} = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} = \text{step}(x)$$



$$\tanh(x) = 2 \text{sigmoid}(2x) - 1$$



Universal Approximation Theorem:

A sufficiently "wide" MLP with one hidden layer can approximate any function arbitrarily well.

("width" is the dimensionality of the hidden features)

Is universal approximation useful?

1. It means we can overfit. → regularization
pre-training
2. The model might be huge. → different
architectures
3. It doesn't tell us how. → gradient descent
(via backprop)

Gradient Descent

We want to minimize $f(x)$

$$f(x + \epsilon) = \sum_{n=0}^{\infty} \frac{\epsilon^n f^{(n)}(x)}{n!}$$

$$f(x + \underbrace{\epsilon}_{\text{small}}) = f(x) + \epsilon f'(x) + O(\epsilon^2)$$

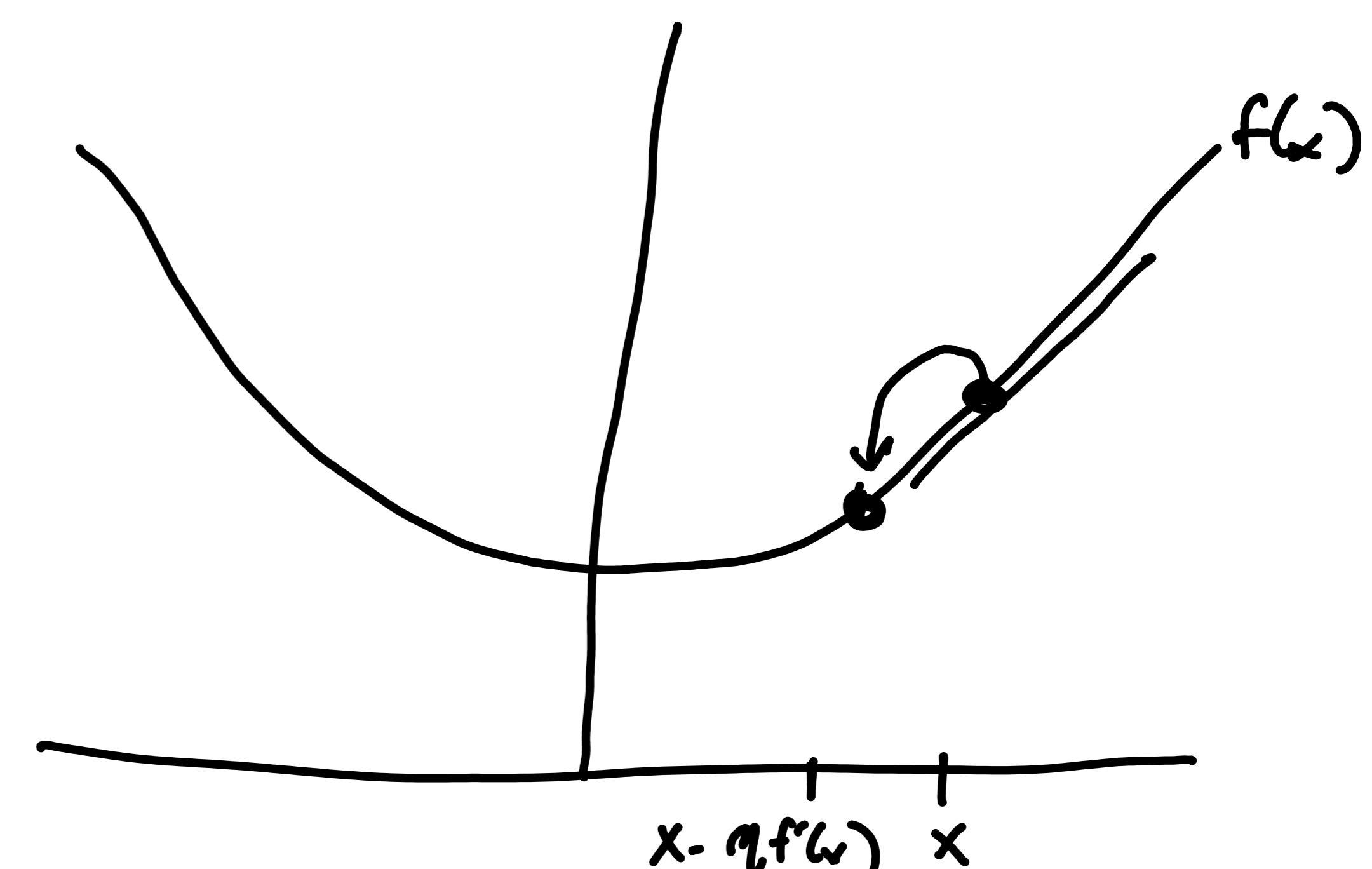
set $\epsilon = -\eta f'(x)$

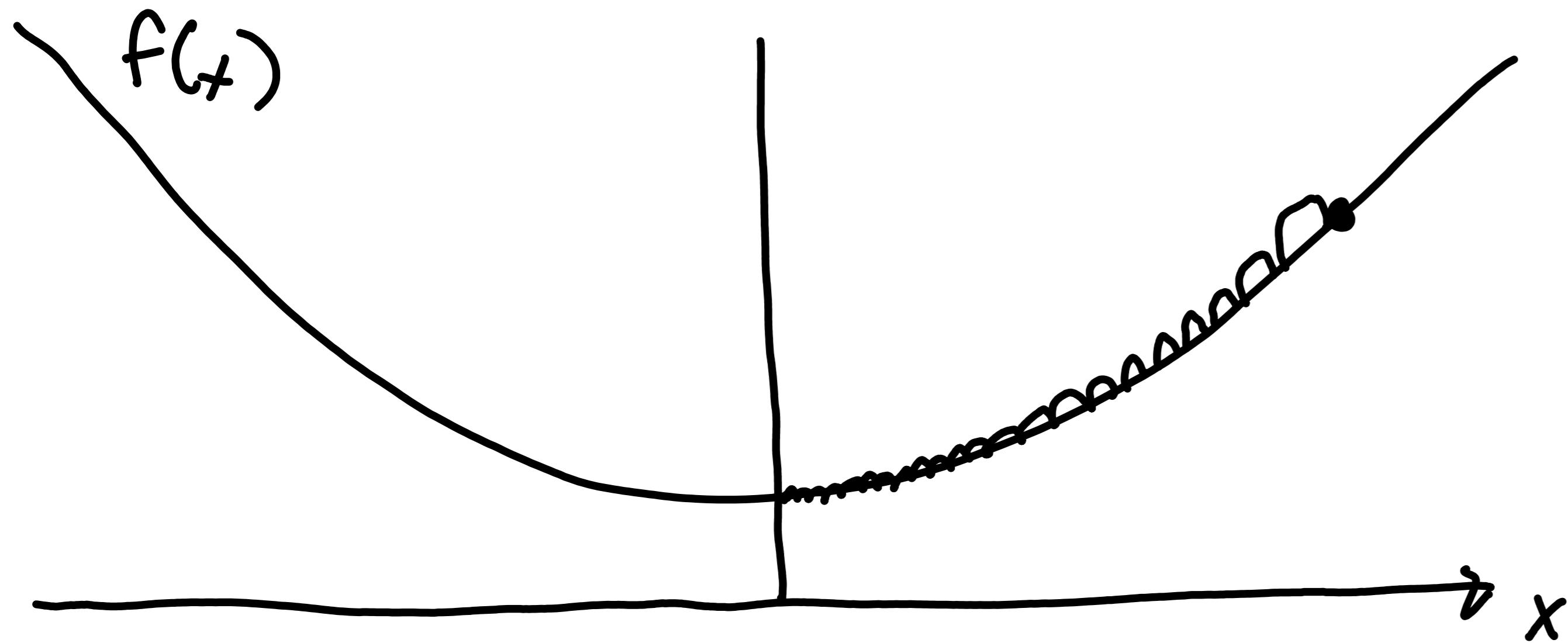
$$f(x - \eta f'(x)) = f(x) - \eta (f'(x))^2 + O(\eta^2 (f'(x))^2)$$

therefore

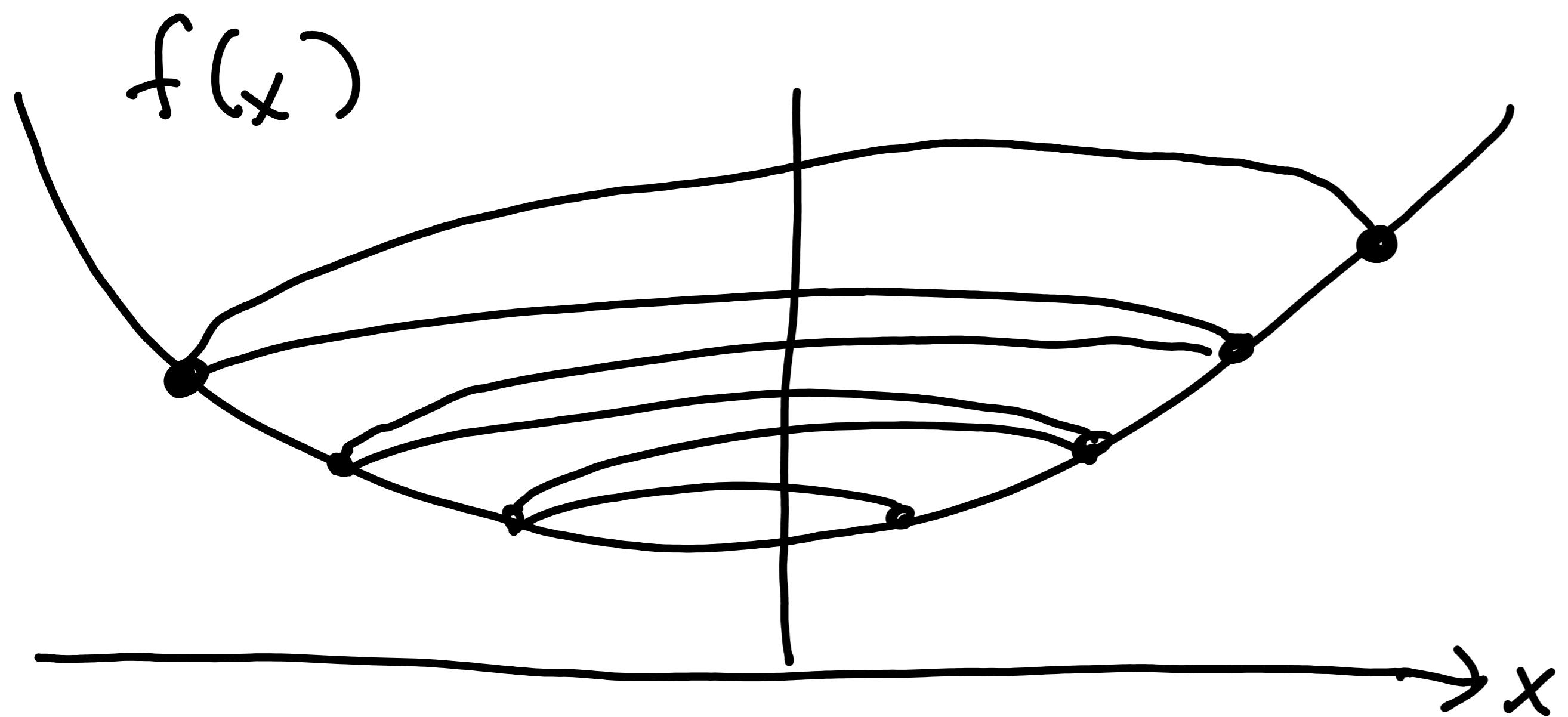
$$f(x - \eta f'(x)) \lesssim f(x)$$

$$\text{so: } x \leftarrow x - \eta f'(x)$$

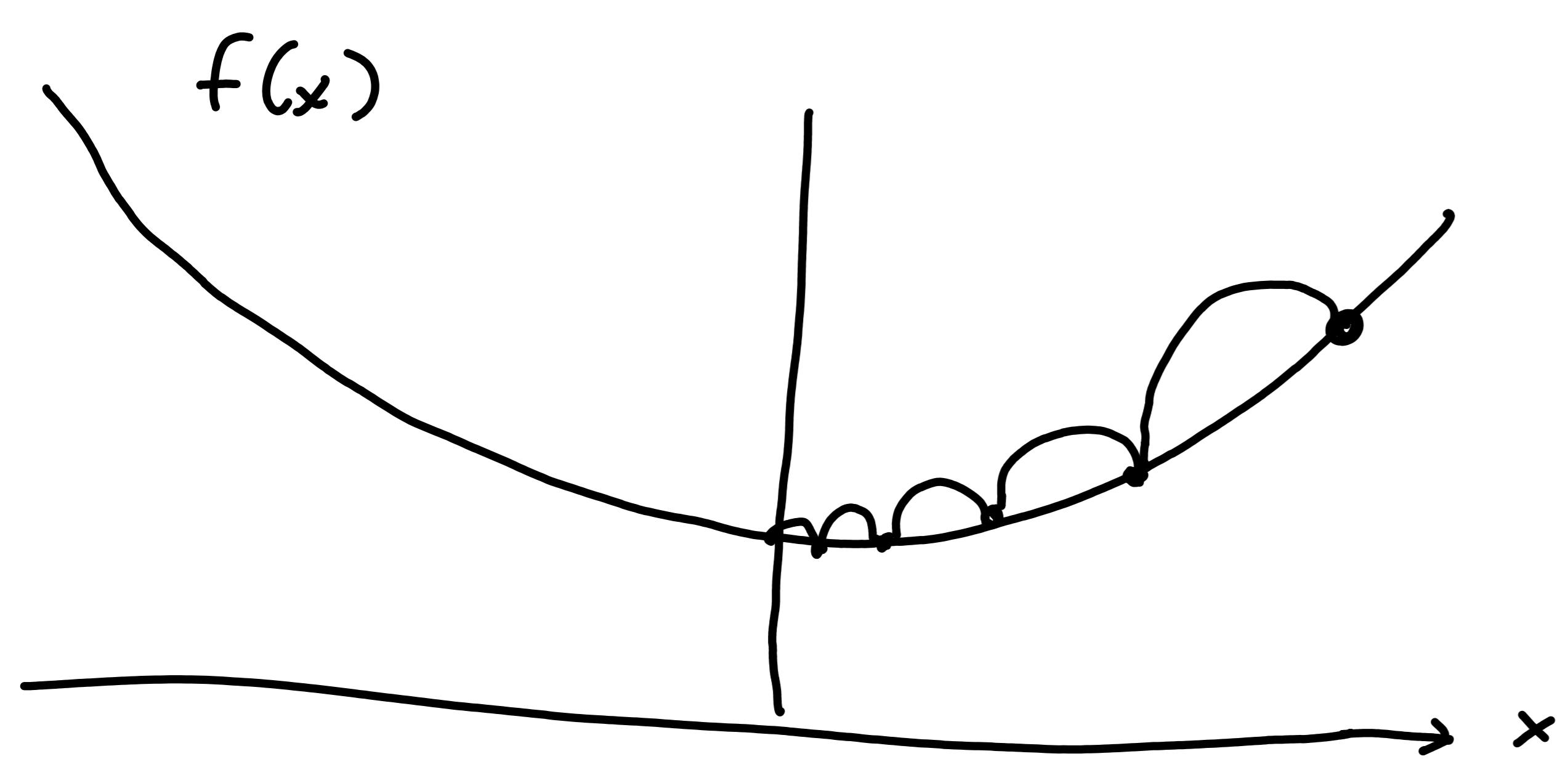




Small η



big η



"just right" η