

Adaptive Gradient Methods

Rather than using the same learning rate for all params, use a different learning rate for each computed "adaptively," based on past gradients.

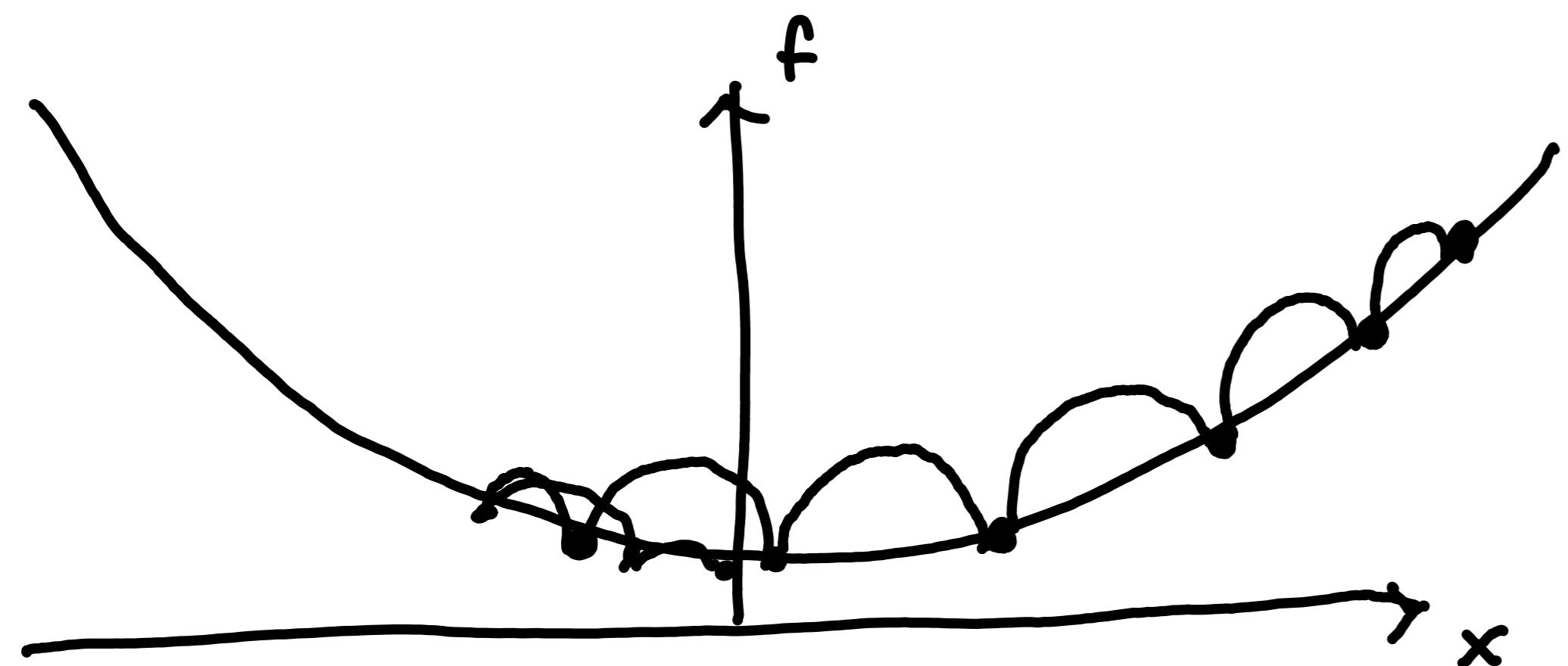
Simplest: Momentum

$$\begin{aligned} v_t &\leftarrow \beta v_{t-1} + g_t & \xrightarrow{\text{new param}} & \text{gradient vector} \\ && & \xrightarrow{\text{at current timestep}} \nabla_{x_t} f(x) \\ x_t &\leftarrow x_{t-1} - \eta v_t & \xrightarrow{\text{learning rate}} & \in \mathbb{R}^d \end{aligned}$$

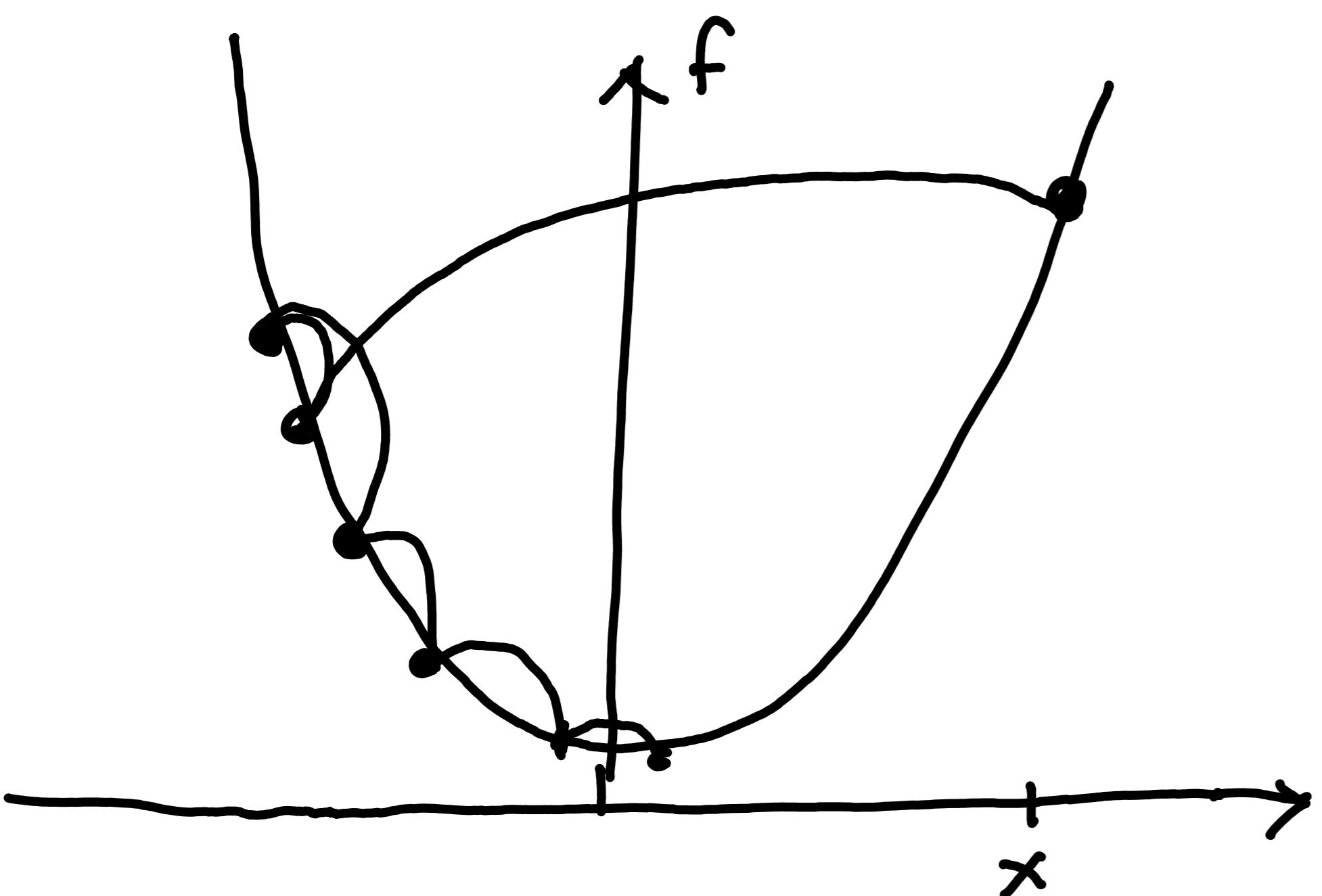
v_t will change more slowly than g_t . If g_t is consistent, v_t will grow in that direction. If g_t is inconsistent, it might "cancel out" changes.

$$v_t \leftarrow \beta v_{t-1} + g_t$$

(assume β is close to 1, 0.9 is typical)



gradient points in the
same direction \rightarrow
take bigger steps



gradient direction is
inconsistent \rightarrow
take smaller steps

Adam:

$$v_t \leftarrow \beta_1 v_{t-1} + (1 - \beta_1) g_t \quad (\text{"first moment"})$$
$$s_t \leftarrow \beta_2 s_{t-1} + (1 - \beta_2) g_t^2 \quad (\text{"second moment"})$$

$$g_t' = \frac{m v_t}{\sqrt{s_t} + \epsilon} \xrightarrow{\text{small constant}}$$

$$x_t \leftarrow x_{t-1} - g_t'$$

$$v_1 = \beta_1 v_0 + (1 - \beta_1) g_1 = (1 - \beta_1) g_1$$

$$v_2 = \beta_1 (1 - \beta_1) g_1 + (1 - \beta_1) g_2$$

$$\beta_1 (1 - \beta_1) + (1 - \beta_1) = 1 - \beta_1^2 \quad (\text{total gradient rescaling})$$

$$v_3 = \underbrace{\beta_1 (\beta_1 (1 - \beta_1) g_1 + (1 - \beta_1) g_2)}_{1 - \beta_1^3} + \underbrace{(1 - \beta_1) g_3}_{1 - \beta_1^3}$$

Total gradient scale is $1 - \beta_1^t$

Typically, β_1 and β_2 are close to 1.

$1 - \beta_1^t$ and $1 - \beta_2^t$ will start near 0 and grow slowly. So v_t and s_t will "biased" towards zero.

To correct this bias: Divide out $1 - \beta_1^t$ and $1 - \beta_2^t$

$$v_t \leftarrow \beta_1 v_{t-1} + (1 - \beta_1) g_t \quad \hat{v}_t = \frac{v_t}{1 - \beta_1^t}$$

$$s_t \leftarrow \beta_2 s_{t-1} + (1 - \beta_2) g_t^2 \quad \hat{s}_t = \frac{s_t}{1 - \beta_2^t}$$

$$g_t' = \frac{\eta \hat{v}_t}{\sqrt{\hat{s}_t} + \epsilon}$$

$$x_t \leftarrow x_{t-1} - g_t'$$

Adam rescales gradients so that you use
(approximately) the average gradient divided
by the average magnitude.

Hopefully this helps us use the same learning
rate for all parameters and all problems!

$$\beta_1 = 0.9$$

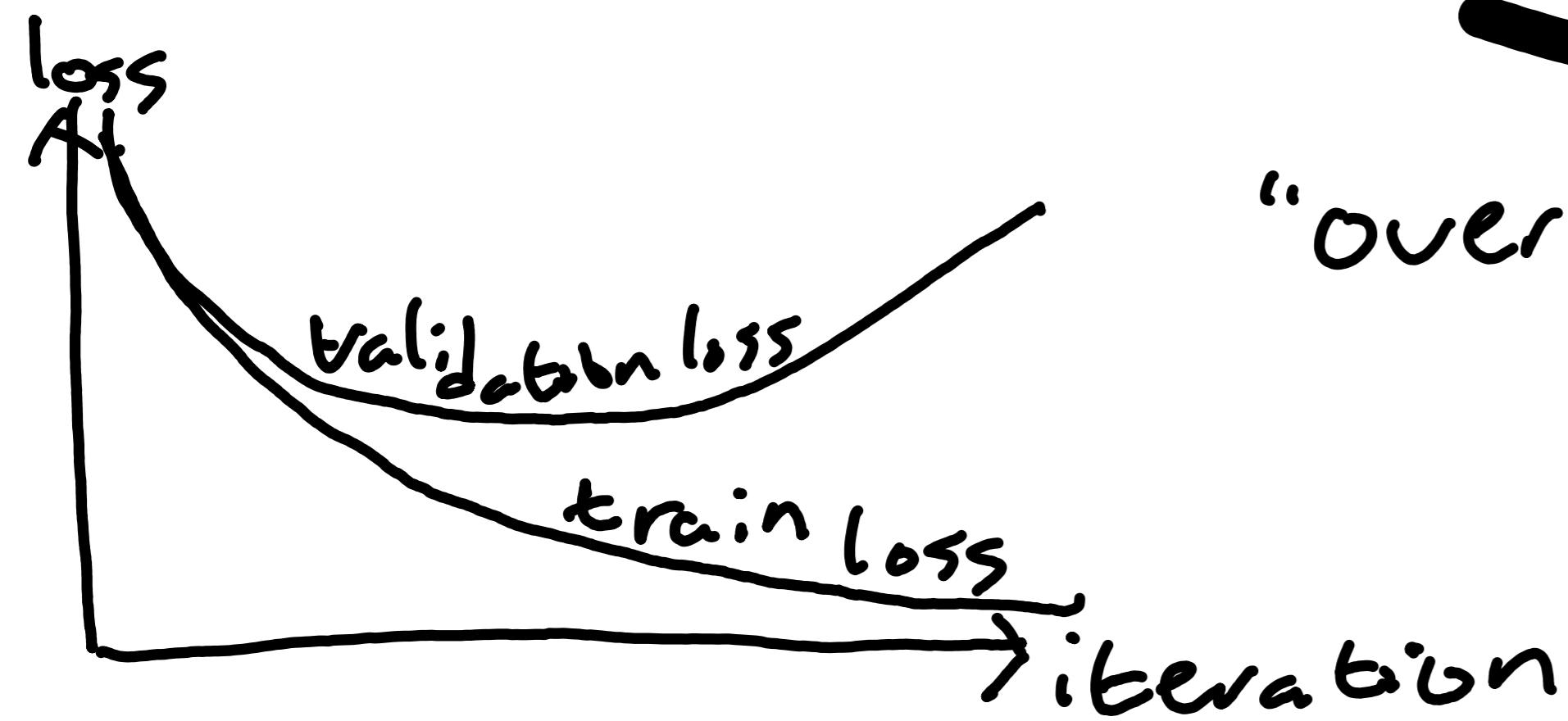
$$\beta_2 = 0.999$$

$$\eta = 0.001$$

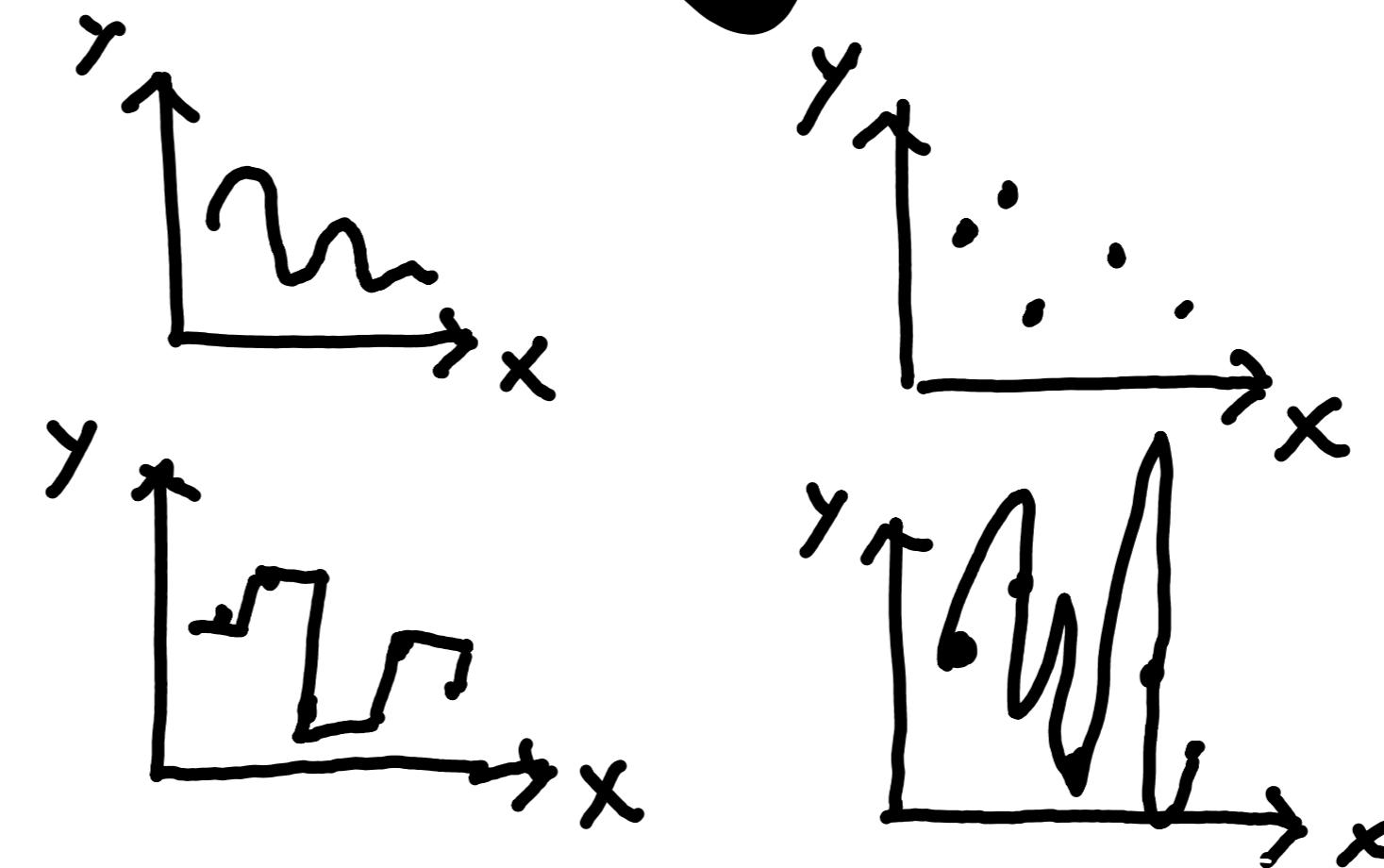
$$\epsilon = 10^{-8}$$

Note: Have to store v_t and s_t somewhere...

Overfitting and regularization

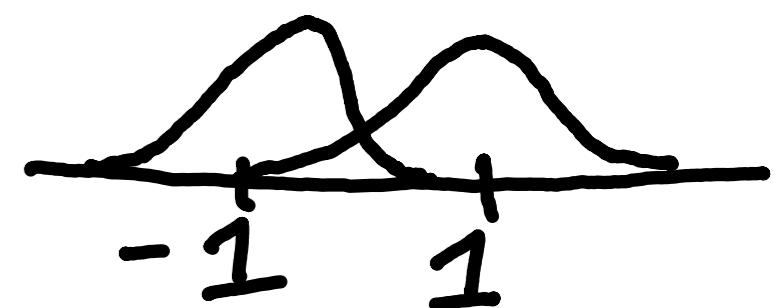


"overfitting"



How can we overfit in a linear model?

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



We use $\hat{y} = wx$

As $\dim(x)$ grows and size of training dataset shrinks, it becomes increasingly likely that for some feature x_i , we can achieve zero loss \rightarrow overfitting!

How can we prevent overreliance on a feature?

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

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

So we are performing "weight decay"

2. "Dropout" features at random :

During training,

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

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

Multi-layer Perceptrons

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

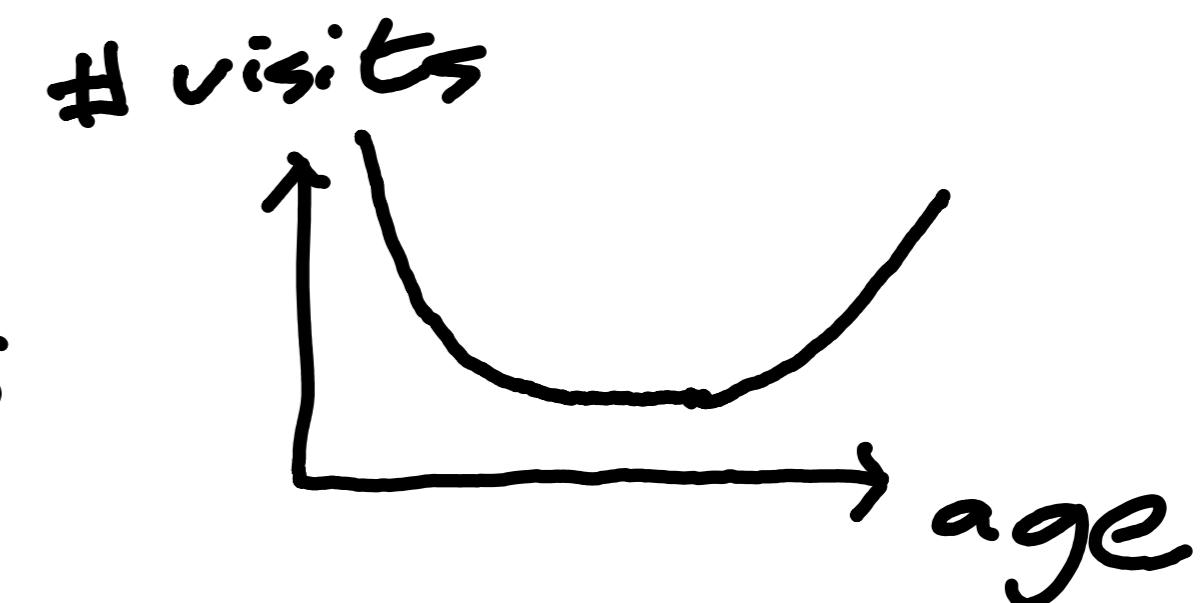
$$\hat{y} = w x + b$$

w pos: $x \uparrow, \uparrow \hat{y}$

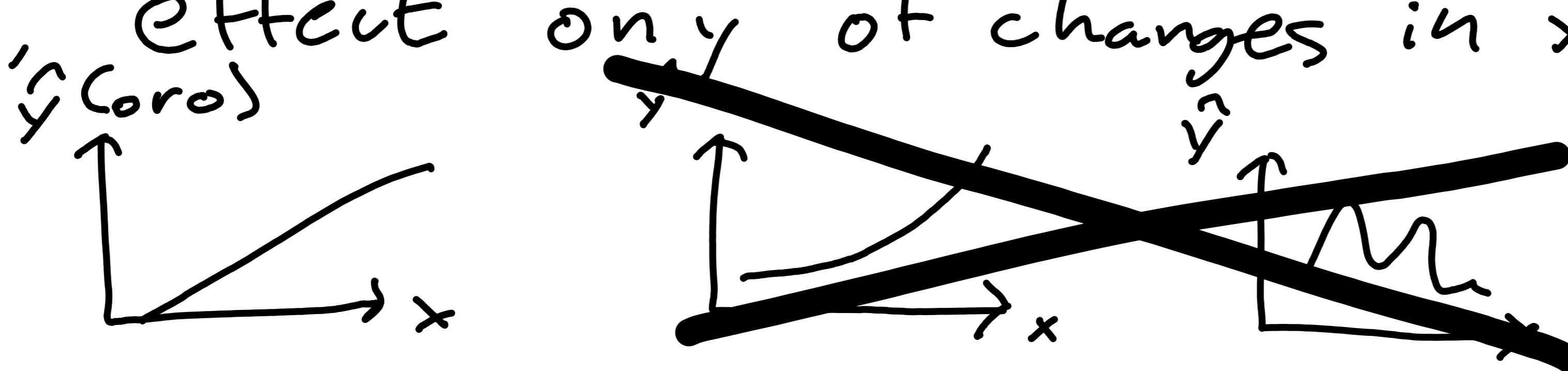
w neg: $x \uparrow, \downarrow \hat{y}$

ex sqft \uparrow , value \uparrow

ex age \uparrow , # of 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

$$\begin{aligned} H &= Xw_1 + b_1 \\ \text{hidden features} &\quad \text{input} \\ O &= Hw_2 + b_2 \\ \text{output} & \\ O &= (Xw_1 + b_1)w_2 + b_2 \\ &= X(w_1 w_2 + b_1 w_2 + b_2) \\ &= Xw + b \end{aligned}$$

The diagram shows a large black 'X' drawn over the equations, indicating that the proposed approach is incorrect or invalid.

Need non linearity! Apply nonlinear functions.

$$H = \sigma(XW_1 + b_1) \quad O = H W_2 + b_2$$

↳ nonlinear function!

One "layer" is a linear transformation and a non linearity.

σ is usually "element wise"

$$h_i = \sigma(x_i)$$

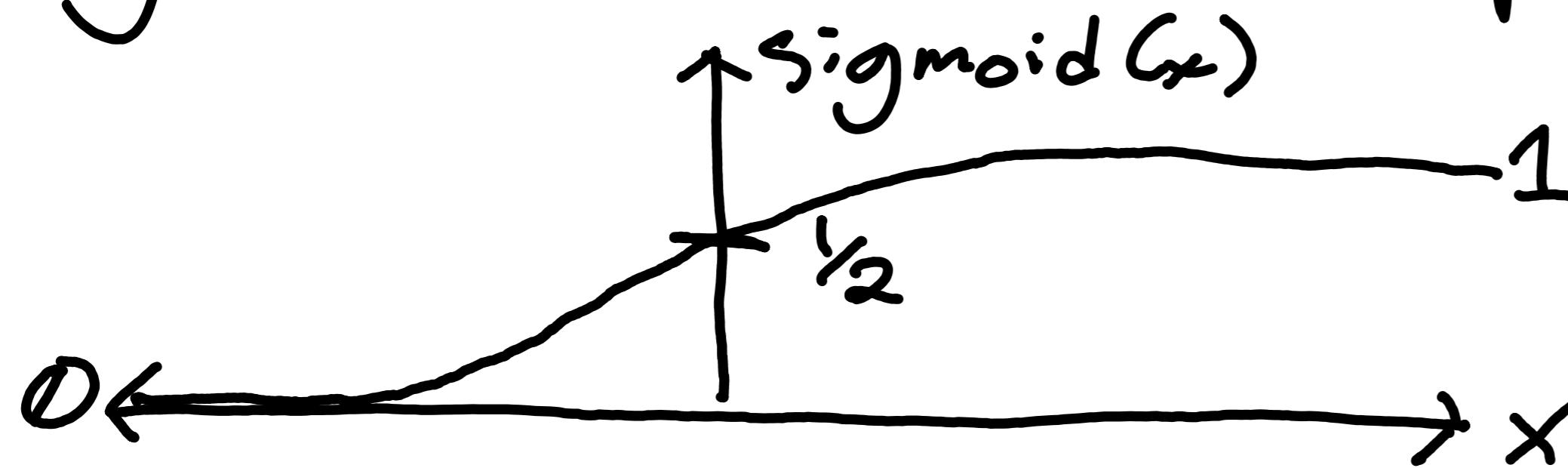
"Neural network"

"Multi-layer perceptron"

"Dense/feed-forward/fully connected"

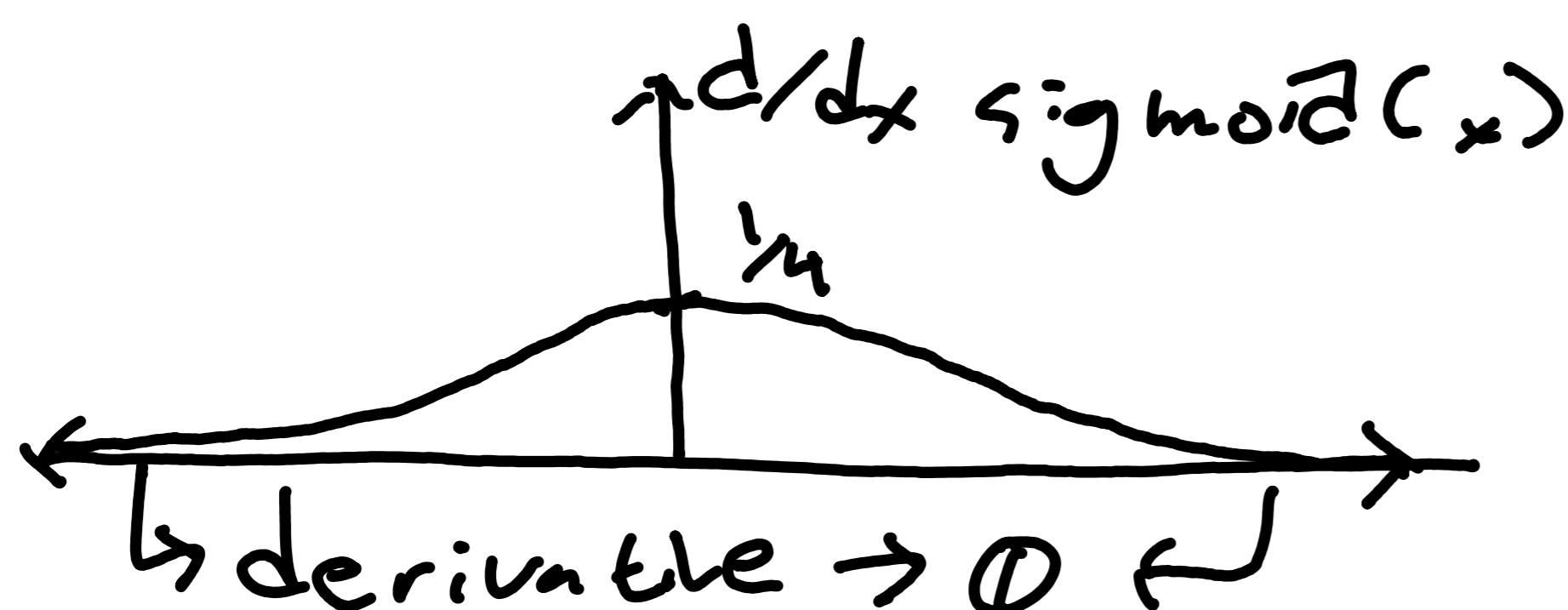
What to use for $\sigma(x)$?

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



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

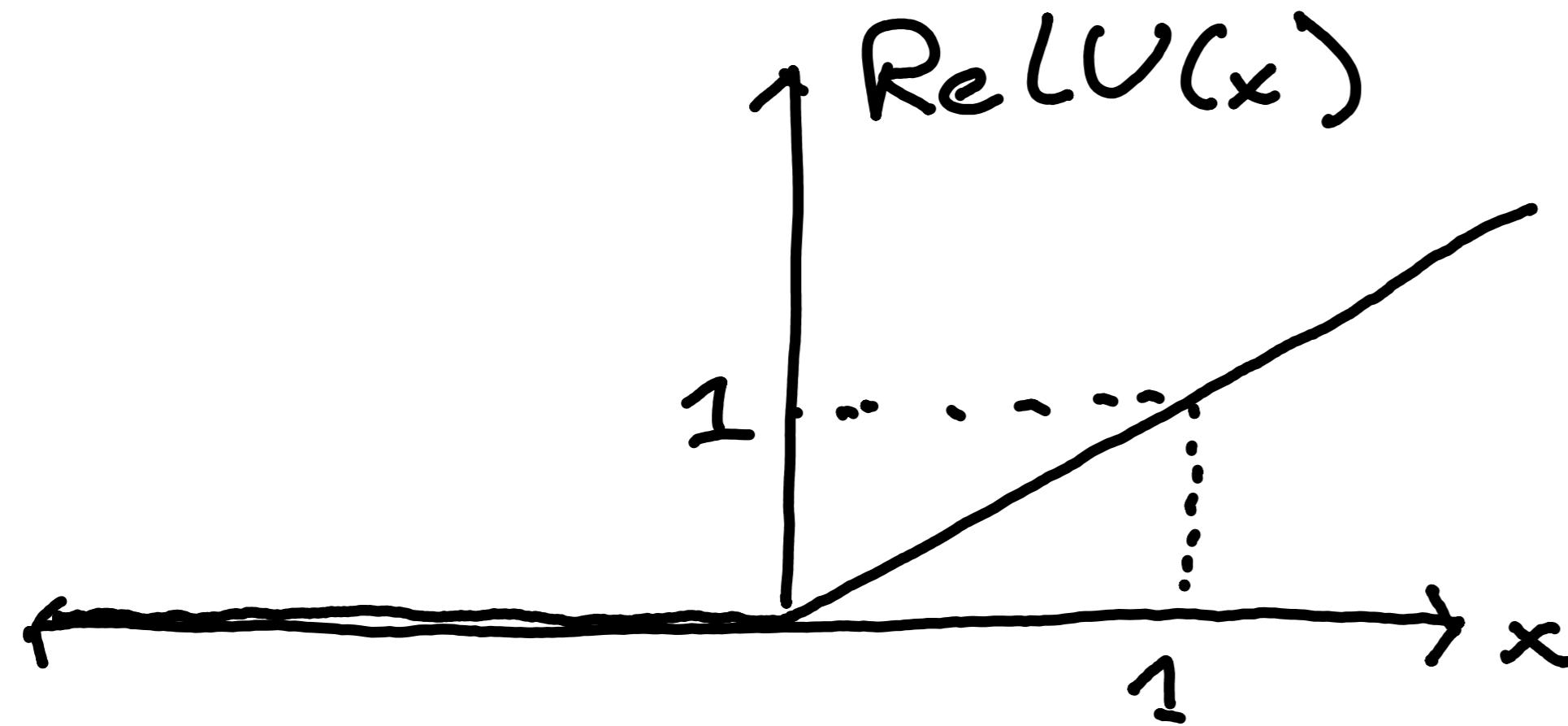
$$= \text{Sigmoid}(x)(1 - \text{Sigmoid}(x))$$



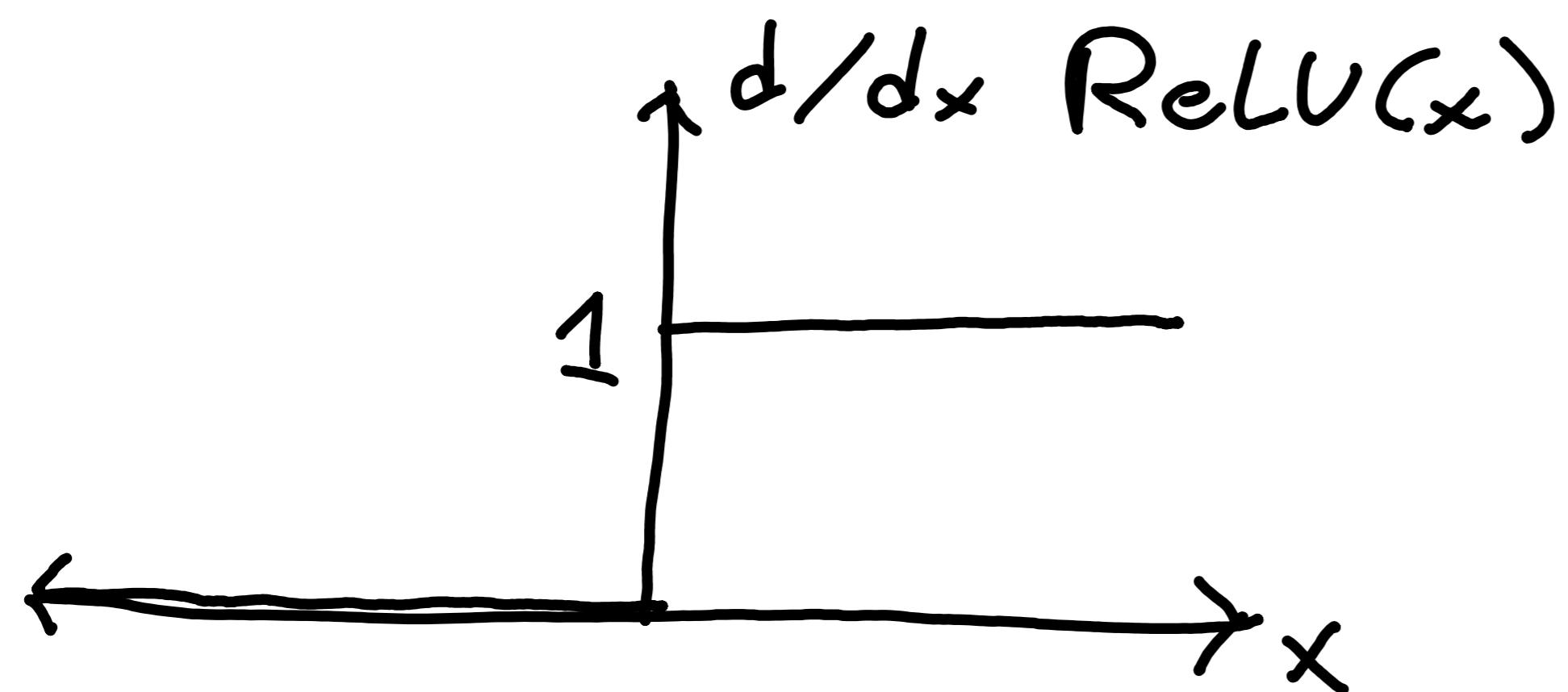
$$\begin{aligned} & 1 - \text{Sigmoid}(-x) \\ &= 1 - \frac{1}{1 + \exp(-x)} \\ &= \frac{1 + \exp(-x)}{1 + \exp(-x)} - \frac{1}{1 + \exp(-x)} \\ &= \frac{\exp(-x)}{1 + \exp(-x)} \end{aligned}$$

ReLU - "Rectified Linear Unit"

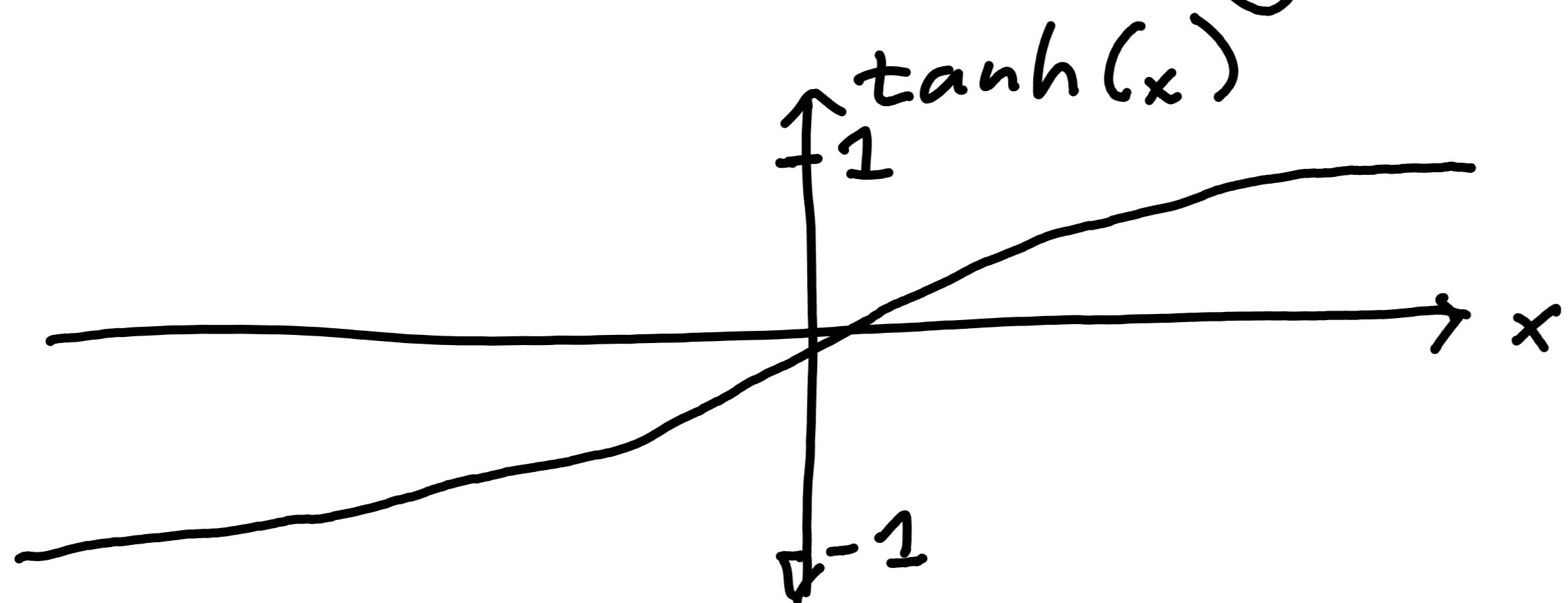
$$\text{ReLU}(x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases} = \max(0, x)$$



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



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



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. \rightarrow regularization
2. The model might be huge, \rightarrow different architectures
3. It doesn't tell us how. \rightarrow gradient descent
(via backprop)