

# Initialization

Gabriel Hope

# Adam

Can we combine adaptive scaling and momentum?

# Adam

Update *velocity*

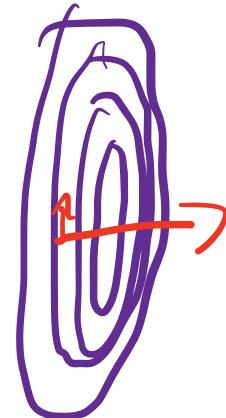
$$\mathbf{v}^{(k+1)} \leftarrow \beta_1 \mathbf{v}^{(k)} + (1 - \beta_1) \nabla_{\mathbf{w}} \text{Loss}(\mathbf{w}^{(k)}, \mathbf{X}, \mathbf{y})$$

Update *scaling*

$$\mathbf{s}^{(k+1)} \leftarrow \beta_2 \mathbf{s}^{(k)} + (1 - \beta_2) (\nabla_{\mathbf{w}} \text{Loss}(\mathbf{w}^{(k)}, \mathbf{X}, \mathbf{y}))^2$$

Update weights

$$\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha \frac{\mathbf{v}^{(k+1)}}{\sqrt{\mathbf{s}^{(k+1)} + \epsilon}}$$



# Adam

Update velocity

$$\mathbf{v}^{(k+1)} \leftarrow \beta_1 \mathbf{v}^{(k)} + (1 - \beta_1) \nabla_{\mathbf{w}} \text{Loss}(\mathbf{w}^{(k)}, \mathbf{X}, \mathbf{y})$$

Update scaling

$$\mathbf{s}^{(k+1)} \leftarrow \beta_2 \mathbf{s}^{(k)} + (1 - \beta_2) (\nabla_{\mathbf{w}} \text{Loss}(\mathbf{w}^{(k)}, \mathbf{X}, \mathbf{y}))^2$$

Modified weight update:

$$\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \frac{\alpha \mathbf{v}^{(k+1)}}{\sqrt{\mathbf{s}^{(k+1)}} + \epsilon}$$

*learning rate*  
*velocity*  
*scale*

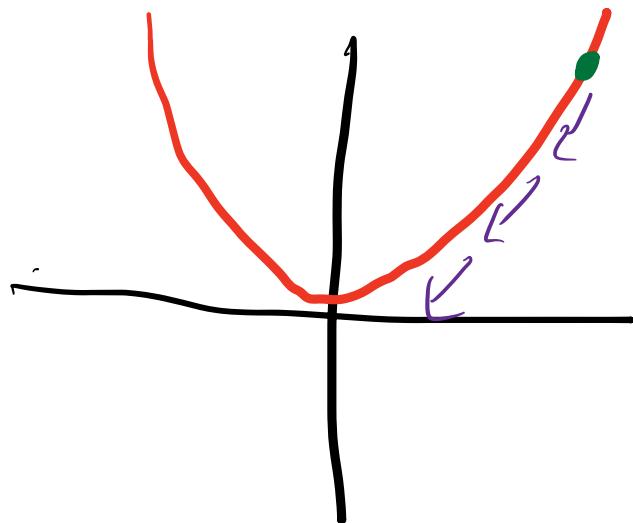
# Adam

At step 0:

$$v^{(k+1)} = \beta_1 \frac{v^{(k)}}{D} + (1 - \beta_1) \nabla_w \text{Loss}(w^{(k)}, X, y)$$

$$v^{(0)} = 0, \quad s^{(0)} = 0$$

$$\frac{v^{(k+1)}}{(1 - \beta_1^k)} = \frac{\beta_1 0 + (1 - \beta_1) \nabla_w \text{Loss}(w^{(k)}, X, y)}{(1 - \beta_1^k)} = \nabla_w \text{Loss}(w^{(k)}, X, y)$$



as  $k \rightarrow \infty$   
 $(1 - \beta^k) \rightarrow 1$   
e.g.  $\beta = 0.9$

$$(0.9)^{100} \approx 0$$
$$\rightarrow (1 - 0.9^{100}) \approx 1$$

# **Summary of gradient descent issues**

**Updates are too slow**

- Stochastic/minibatch gradient descent

**SGD gradients are very noisy (high variance)**

- Increase batch size, use momentum

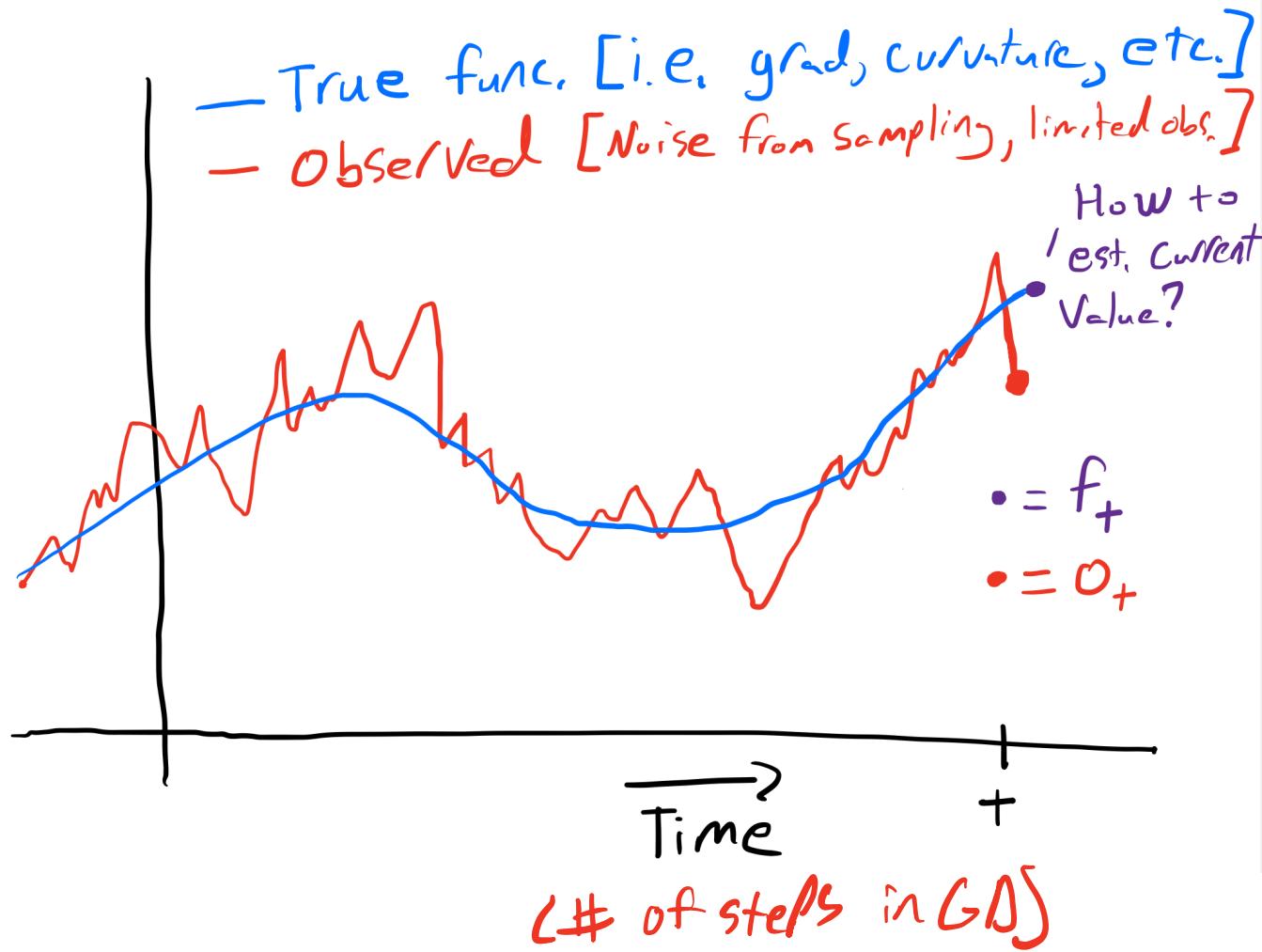
**Stuck at saddle points or shallow optima**

- Use momentum

**Inconsistent scaling of the gradient**

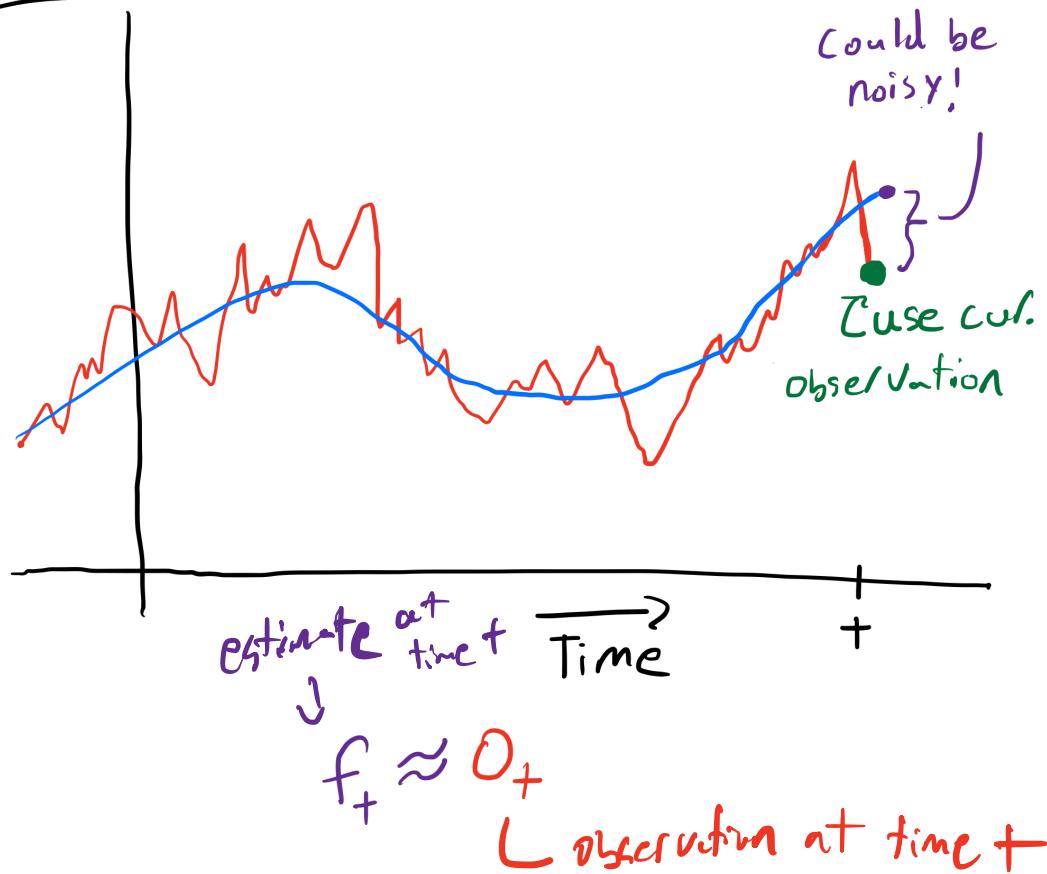
- Use RMSProp scaling

## Exponential moving average (EMA)



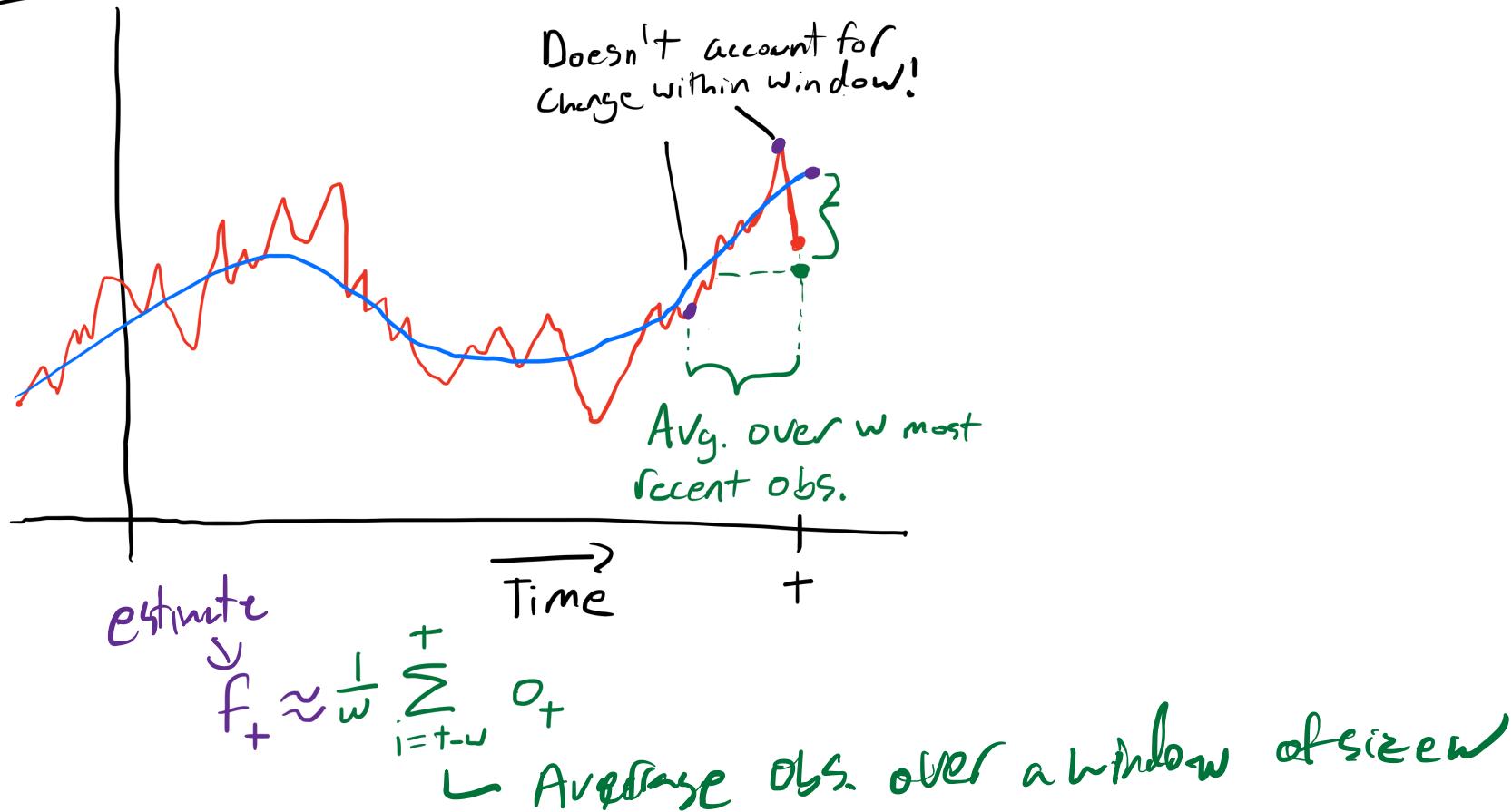
# Exponential moving average (EMA)

Idea !



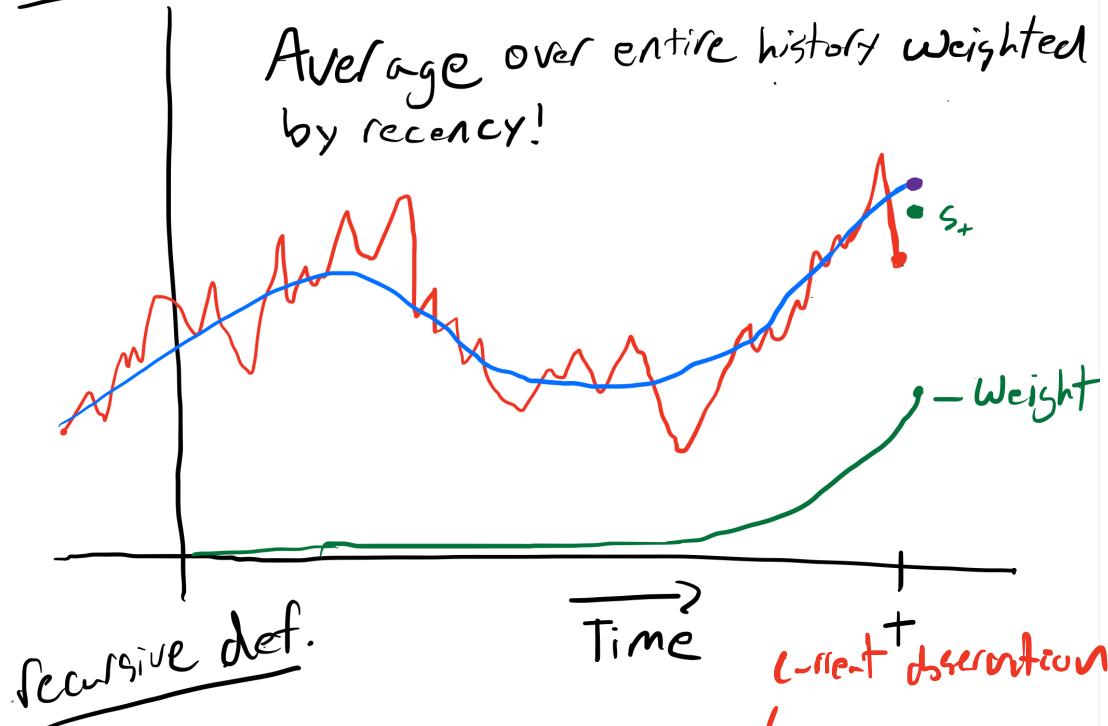
## Exponential moving average (EMA)

Idea 2 [Box avg.]



## Exponential moving average (EMA)

Idea 3 [ EMA ]



$$\alpha \in [0, 1]$$

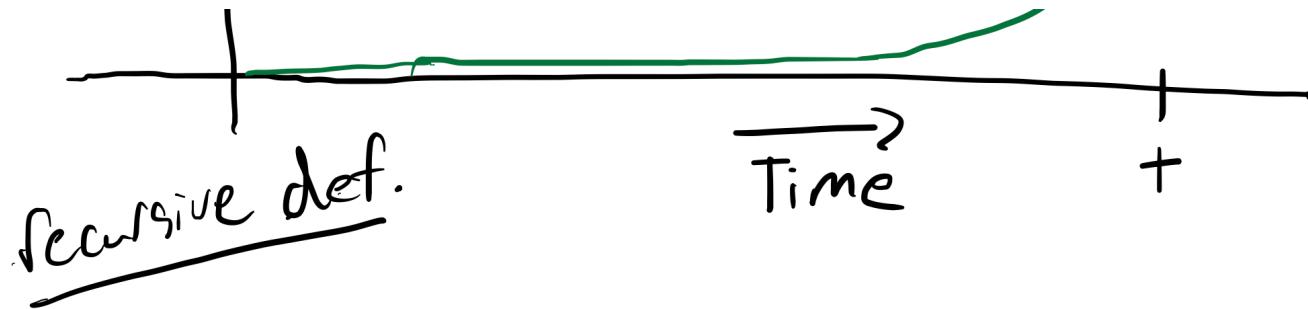
$$f_+ \approx s_+ \quad s_+ = (1-\alpha) o_t + \alpha s_{t-1}$$

$\downarrow$

estimate at time  $t$

$\hookrightarrow$  prev. at  $t-1$

## Exponential moving average (EMA)



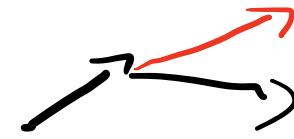
$$f_t \approx S_t \quad S_t = (1-\alpha) O_t + \alpha S_{t-1}$$

$$\rightarrow f_t \approx (1-\alpha) O_t + (1-\alpha)\alpha O_{t-1} + (1-\alpha)\alpha^2 O_{t-2} + \dots$$

Ex.  $\alpha = 0.9$

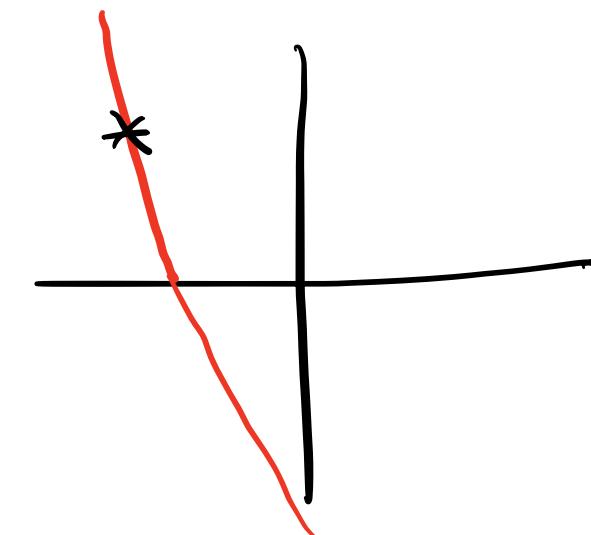
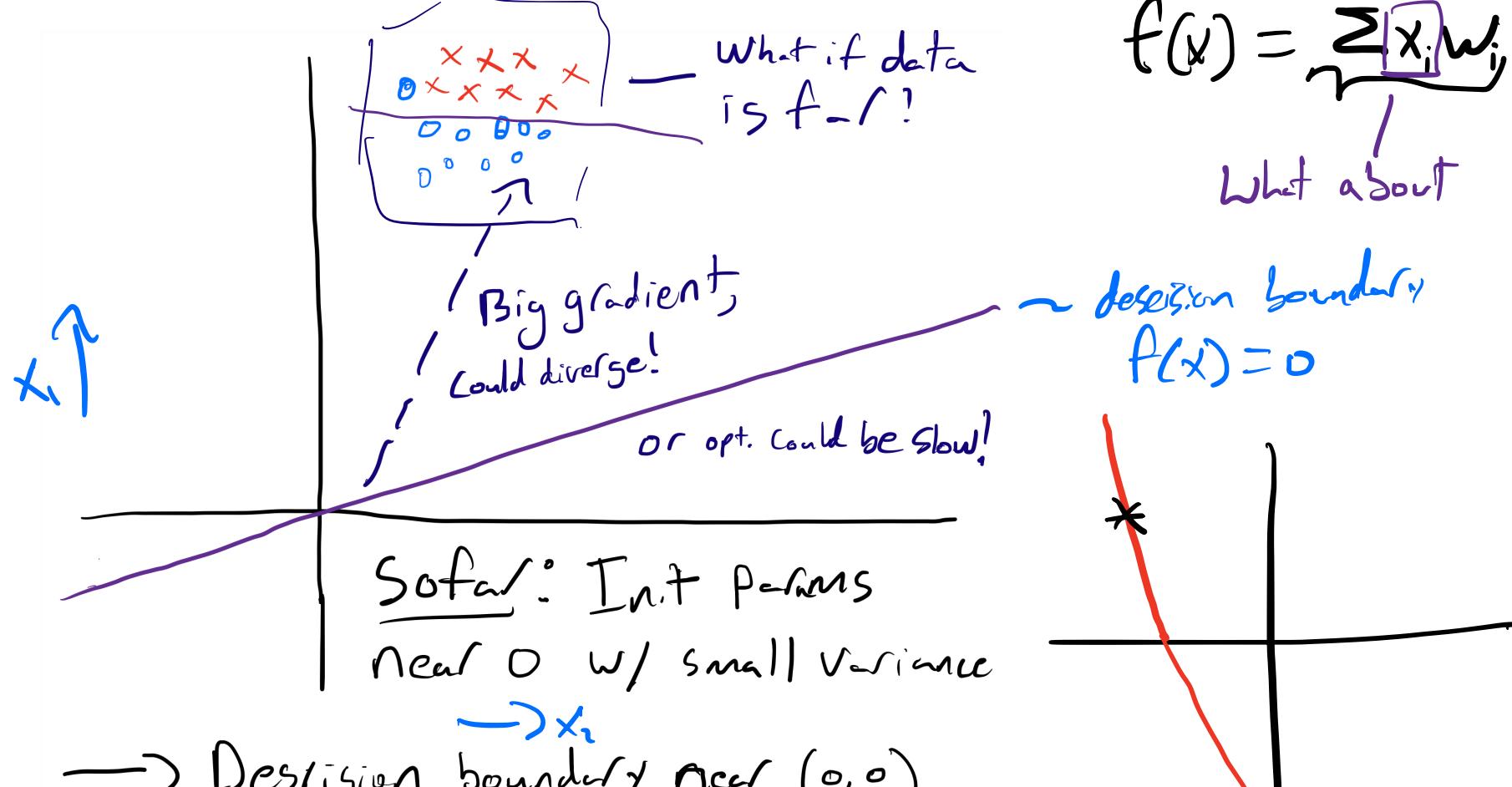
$$f_t \approx 0.1 O_t + 0.09 O_{t-1} + 0.081 O_{t-2} + 0.0729 O_{t-3} + \dots$$

momentum



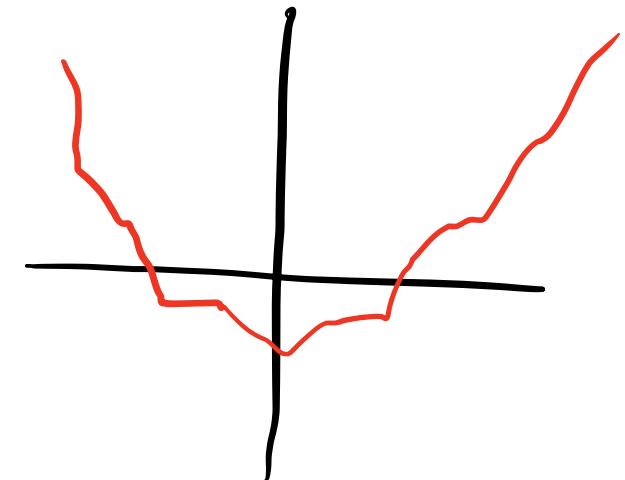
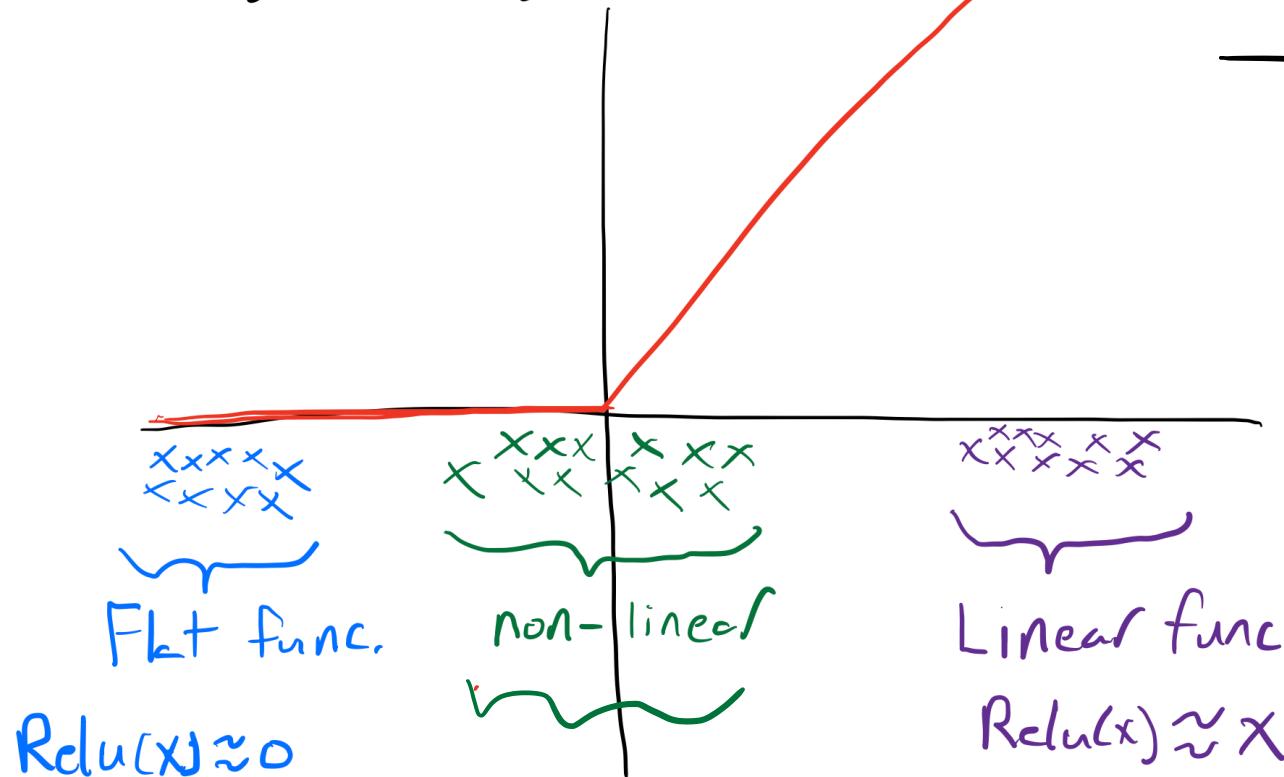
$$\approx 0 \\ + \boxed{O} O_{t-100}$$

## Data normalization



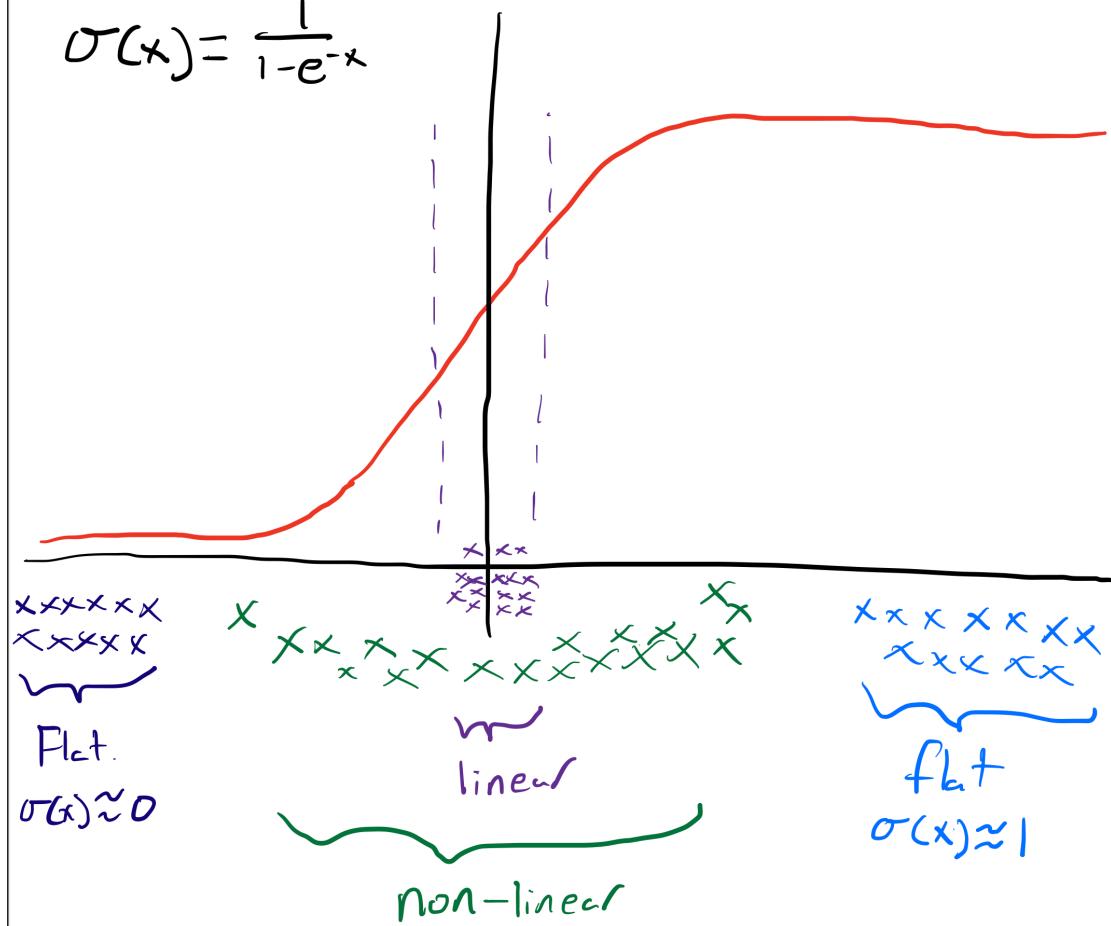
# Data normalization

$$\text{ReLU}(x) = \max(0, x)$$

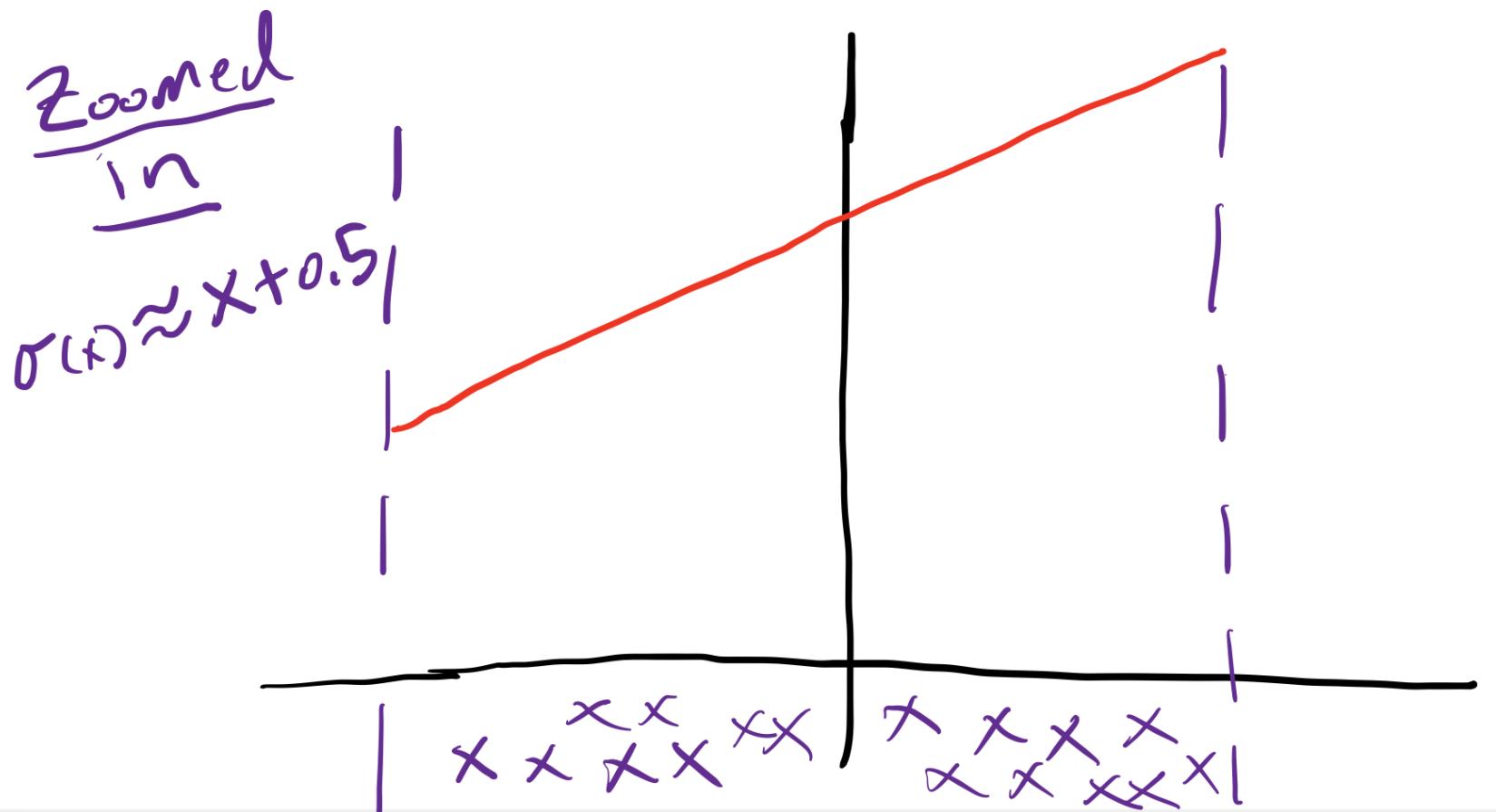


# Data normalization

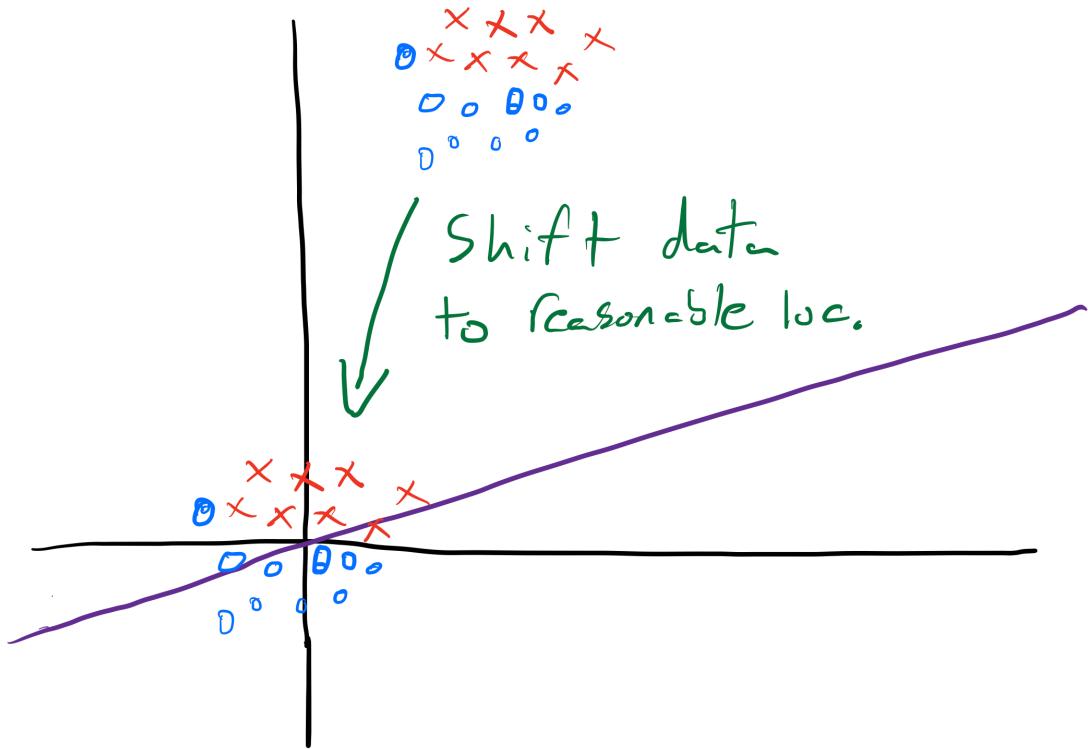
$$\sigma(x) = \frac{1}{1-e^{-x}}$$



## Data normalization

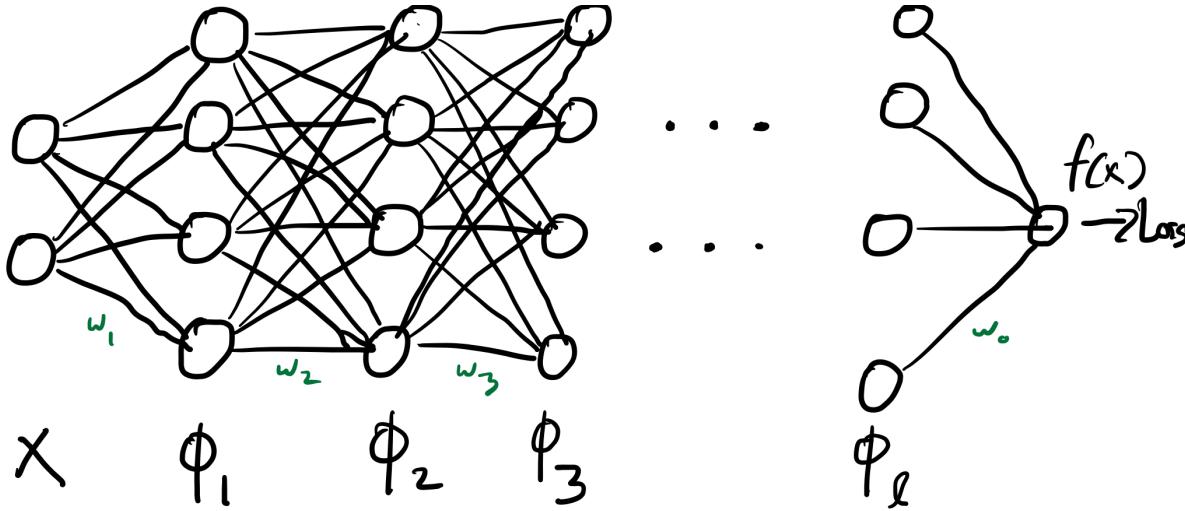


# Data normalization



Decision boundary near  $(0,0)$

# Vanishing and exploding gradients



$$\phi_1 = \sigma(x^T w_1 + b)$$

$$\phi_2 = \sigma(\phi_1^T w_2 + b)$$

$$\phi_3 = \sigma(\phi_2^T w_3 + b)$$

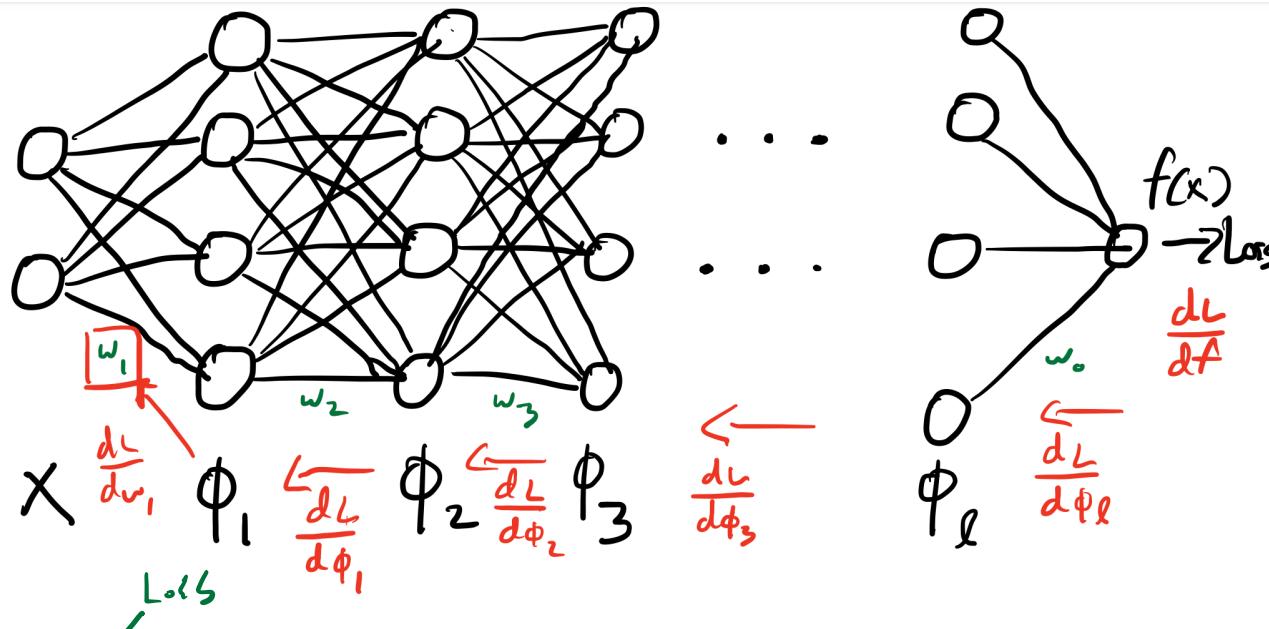
⋮

$$f = \phi_l^T w_l + b \quad L = \text{Loss}(f)$$

Many layers

$L \rightarrow \infty$

# Vanishing and exploding gradients



$$\frac{dL}{dw_i} = \underbrace{\frac{d\phi_1}{d w_i} \cdot \frac{d\phi_2}{d\phi_1} \cdot \frac{d\phi_3}{d\phi_2} \cdots}_{\text{Loss}} \underbrace{\frac{d\phi_\ell}{d\phi_{\ell-1}} \frac{df}{d\phi_\ell} \cdot \frac{dL}{df}}$$

$\frac{d}{dw_i} \sigma(x^T w + b) = x \sigma'(x w + b)$

for all  $\frac{d\phi_i}{d\phi_{i-1}}$

$$\frac{dL}{dw_i} = x \sigma'(x w_i + b) \left( \prod_{j=L}^l \frac{d\phi_j}{d\phi_{j-1}} \right) \frac{df}{d\phi_\ell} \cdot \frac{dL}{df}$$

$\left| \frac{d\phi_i}{d\phi_{i-1}} \right| \approx M$

$$\left| \frac{dL}{dw_i} \right| \approx |x| |\sigma'(x w_i + b)| \left( \prod_{j=L}^l \left| \frac{d\phi_j}{d\phi_{j-1}} \right| \right)$$

$\left| \frac{d\phi_i}{d\phi_{i-1}} \right| \approx M$

$\approx M^2$

## Vanishing and exploding gradients

If  $M > 1 \rightarrow \left| \frac{dL}{dw_1} \right| \gg 1$

Gradient explodes!

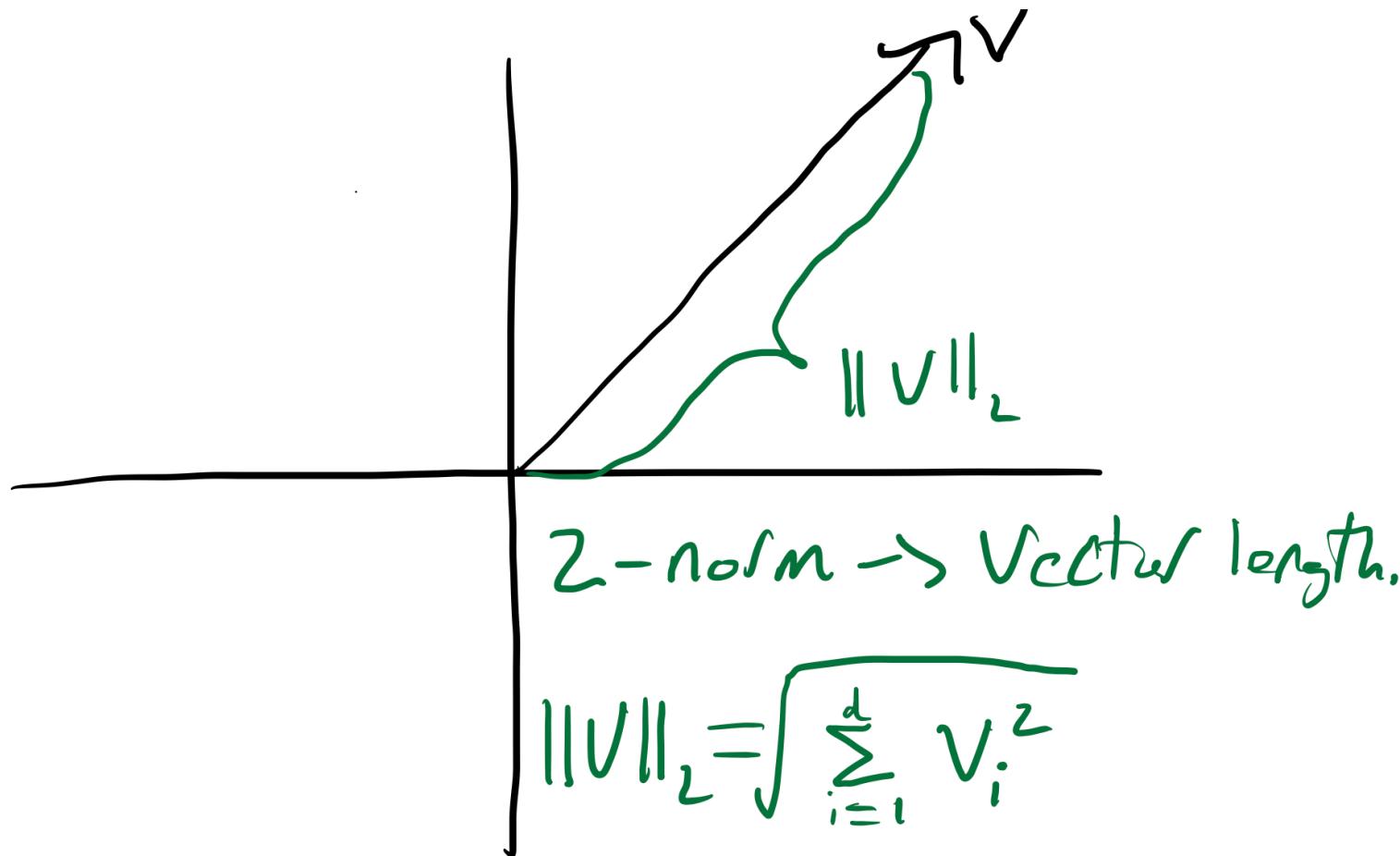
If  $M < 1 \rightarrow \left| \frac{dL}{dw_1} \right| \approx 0$

Gradient Vanishes

Both bad for gradient descent!



## Gradient clipping



# Gradient clipping

Explicitly clip the gradient to prevent it from becoming too large.

$$\text{clip}_{\text{value}}(\mathbf{x}, \epsilon) = \begin{bmatrix} \min(\max(x_1, -\epsilon), \epsilon) \\ \min(\max(x_2, -\epsilon), \epsilon) \\ \vdots \\ \min(\max(x_n, -\epsilon), \epsilon) \end{bmatrix}$$

*Clip by Value*

*$\epsilon = \text{limit}$*

*E.g. gradient*

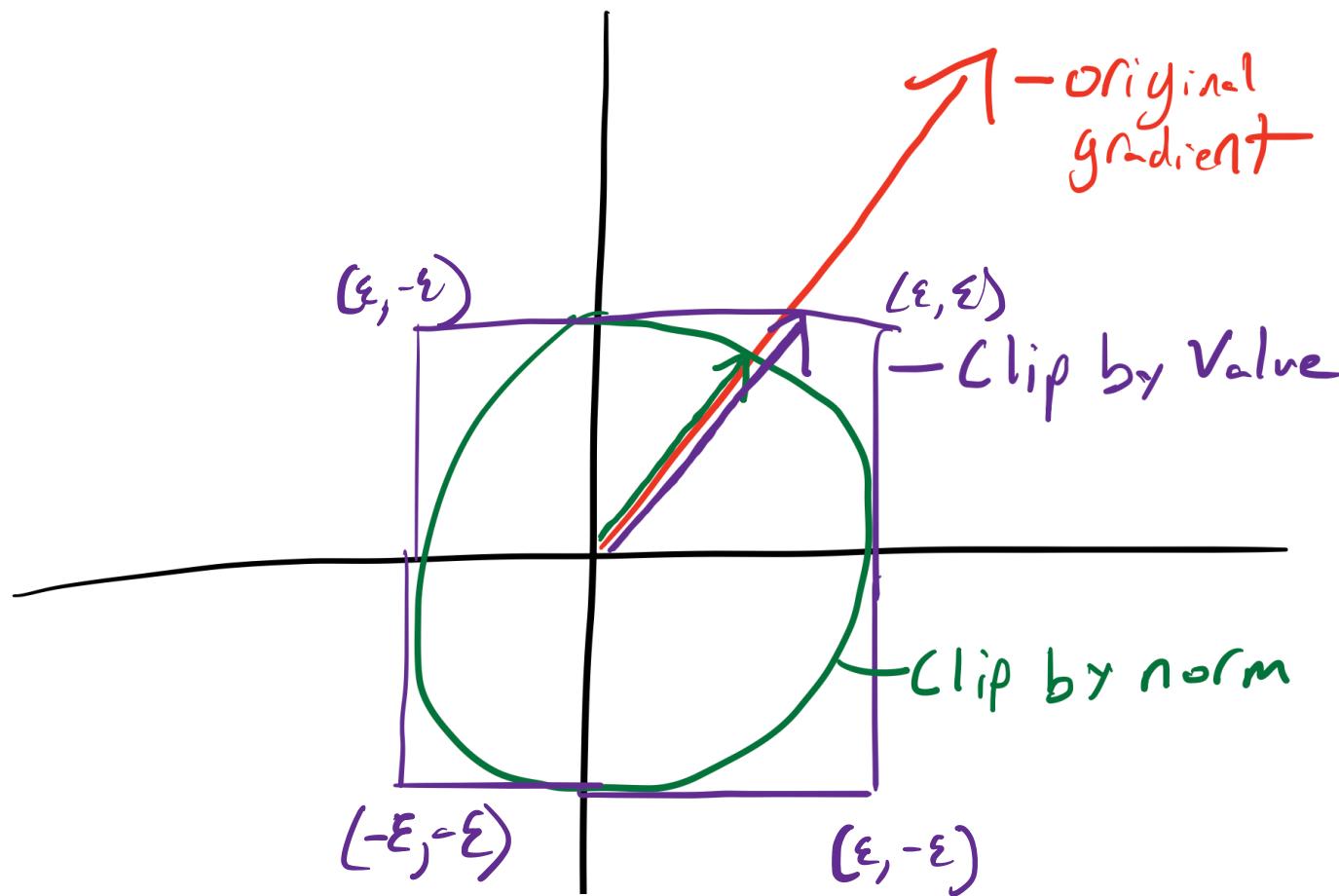
$$\text{clip}_{\text{norm}}(\mathbf{x}, \epsilon) = \begin{cases} \frac{\epsilon \mathbf{x}}{\|\mathbf{x}\|_2} & \text{if: } \|\mathbf{x}\|_2 > \epsilon \\ \mathbf{x} & \text{if: } \|\mathbf{x}\|_2 \leq \epsilon \end{cases}$$

*- enforce that  $\mathbf{x}$  has length  $\epsilon$*

Clipped  
gradient  
descent

$$\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \alpha \text{clip}(\nabla_{\mathbf{w}} \text{Loss}(\mathbf{w}^{(k)}, \mathbf{X}, \mathbf{y}))$$

# Gradient clipping



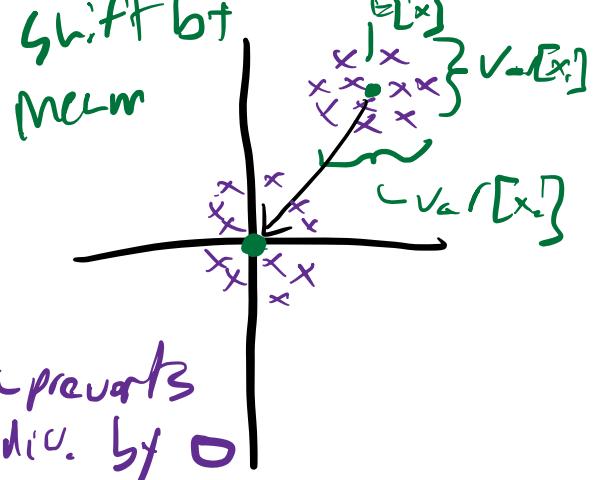
# Normalization

# Batch normalization

Normalize over the batch:

$$\text{BatchNorm}(x) = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}}$$

observation



Training time:

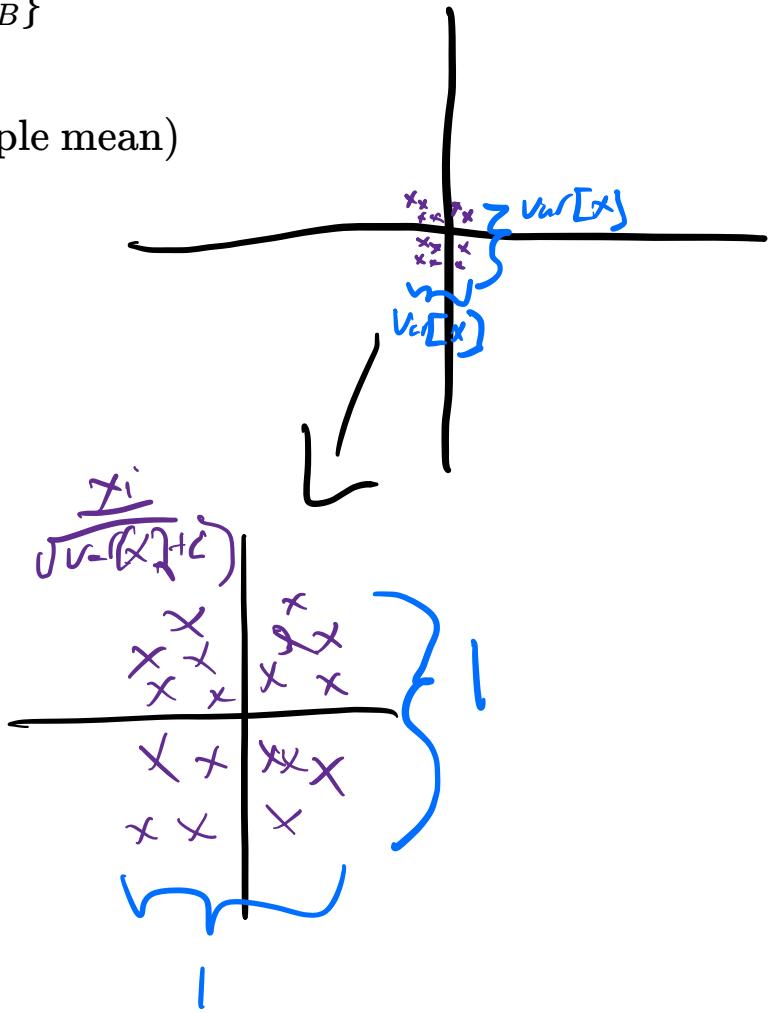
Batch:  $\{x_1, x_2, \dots, x_B\}$

$$\mathbb{E}[x] \approx \bar{x} = \frac{1}{B} \sum_{i=1}^B x_i \quad (\text{sample mean})$$

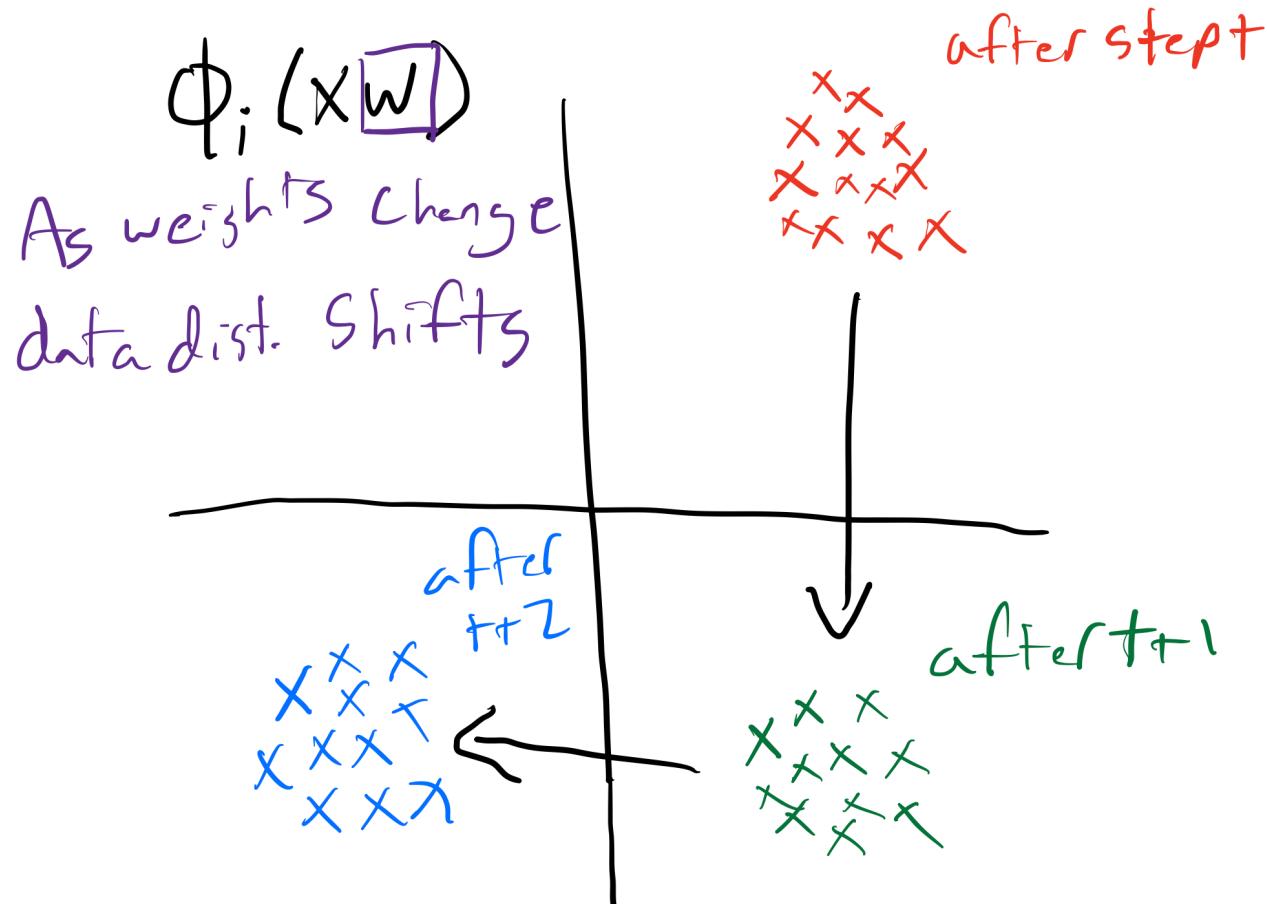


Our current estimate of  
The mean  $[\mathbb{E}[x]]$

Note  
Could also use full dataset

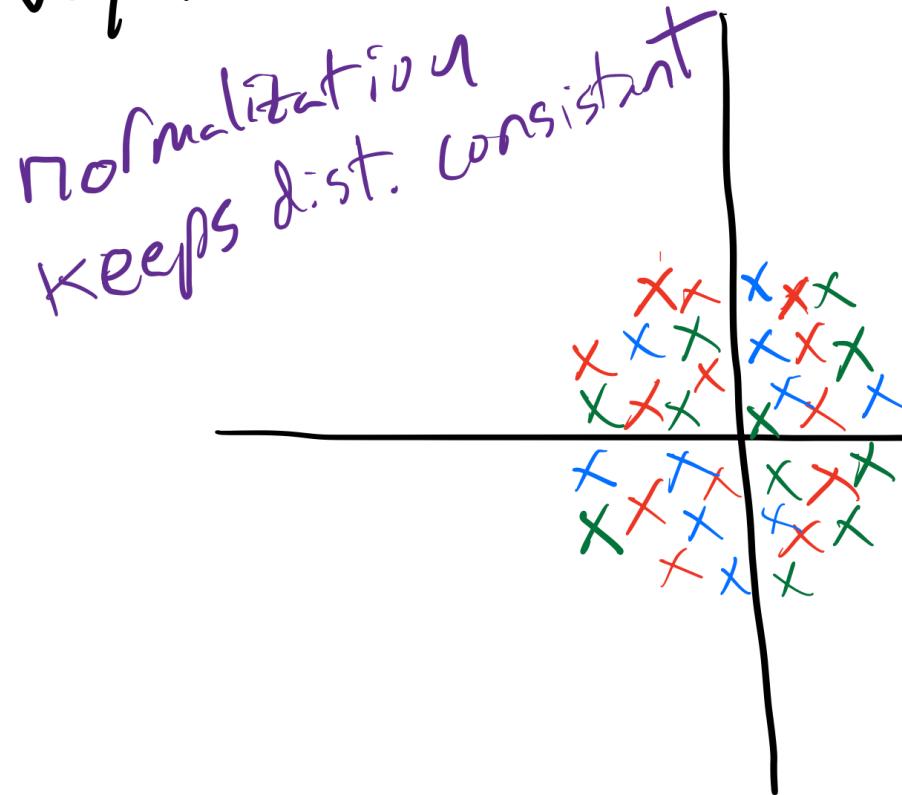


## Batch normalization



## Batch normalization

w/ normalization



# Batch normalization

Biased estimator:

$$\text{Var}[x] \approx s^2 = \frac{1}{B} \sum_{i=1}^B \left( x_i - \left( \frac{1}{B} \sum_{i=1}^B x_i \right) \right)^2 \quad (\text{sample var.})$$

Sample mean

Unbiased estimator:

$$\text{Var}[x] \approx s^2 = \frac{1}{B-1} \sum_{i=1}^B \left( x_i - \left( \frac{1}{B} \sum_{i=1}^B x_i \right) \right)^2 \quad (\text{sample var.})$$

$$\text{BatchNorm}_{\text{train}}(x) = \frac{x - \bar{x}}{\sqrt{s^2 + \epsilon}}$$

# Batch normalization

Running estimate:

$$\bar{\mu}^{(k+1)} \leftarrow \beta \bar{\mu}^{(k)} + (1 - \beta) \bar{x}^{(k)}$$

$$\bar{\sigma}^{2(k+1)} \leftarrow \beta \bar{\sigma}^{2(k)} + (1 - \beta) s^{2(k)}$$

$$\text{BatchNorm}_{\text{test}}(x) = \frac{x - \bar{\mu}}{\sqrt{\bar{\sigma}^2 + \epsilon}}$$

# Layer normalization

Normalize over the layer:

$$\text{LayerNorm}(\mathbf{x}) = \frac{\mathbf{x} - \bar{x}}{\sqrt{s^2 + \epsilon}}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Training & test time:

$$\bar{x} = \frac{1}{d} \sum_{i=1}^d x_i \quad (\text{output mean})$$

# Layer normalization

Biased estimator:

$$s^2 = \frac{1}{d} \sum_{i=1}^d \left( x_i - \left( \frac{1}{d} \sum_{i=1}^d x_i \right) \right)^2 \quad (\text{output var.})$$

Unbiased estimator:

$$s^2 = \frac{1}{d-1} \sum_{i=1}^d \left( x_i - \left( \frac{1}{d} \sum_{i=1}^d x_i \right) \right)^2 \quad (\text{output var.})$$

## Scaled normalization

$$\text{BatchNorm}(x) = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} \gamma + \kappa$$

$$\text{LayerNorm}(\mathbf{x}) = \frac{\mathbf{x} - \bar{x}}{\sqrt{s^2 + \epsilon}} \gamma + \kappa$$