

Class Plan

Recap

Reminders

GPU

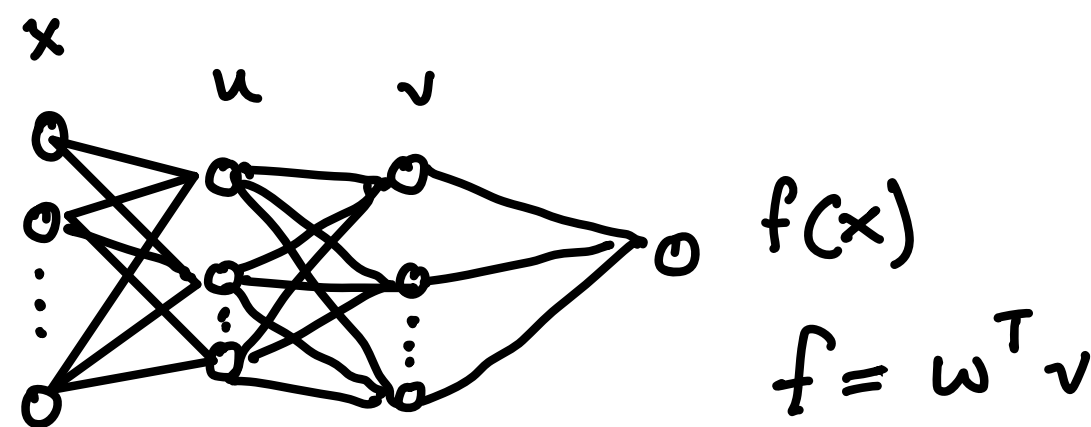
Variants of Gradient Descent

Initialization

Generalization

Neural Network

1)



$$u = \sigma(Wx) \quad v = \sigma(W'u)$$

2) $\mathcal{L}(w) =$ how bad our predictions are

3)

$$\nabla \mathcal{L}(w) = \left[\frac{\partial \mathcal{L}(w)}{\partial w_i} \right]$$

$w \in \mathbb{R}^d$

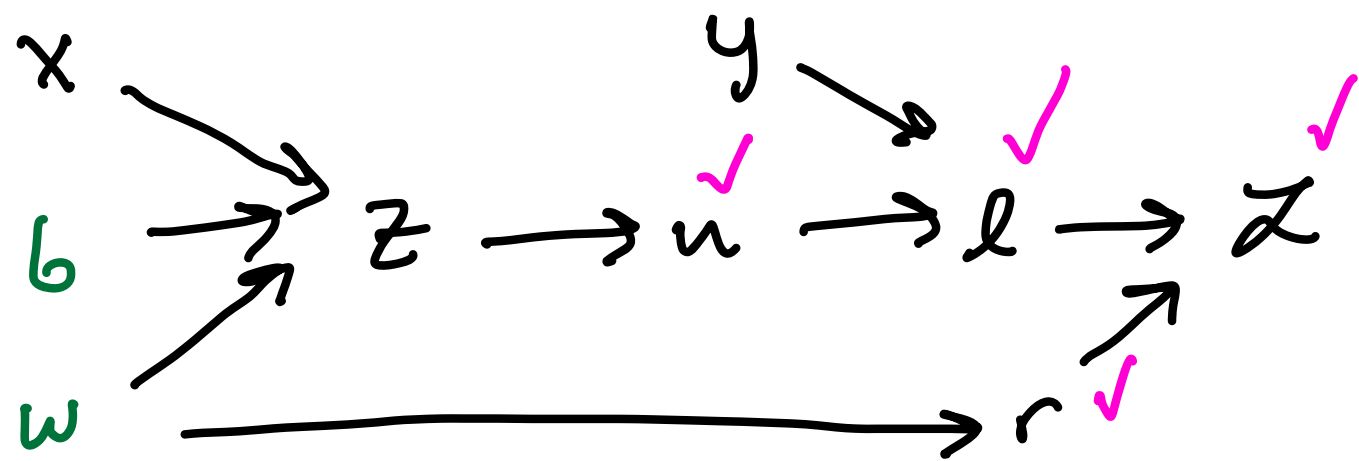
$$w \leftarrow w - \overset{\text{eta}}{\eta} \nabla \mathcal{L}(w)$$

$$\mathcal{L}(w, b) = \frac{1}{2}(y - \sigma(wx + b))^2 + \lambda w^2$$

$$z = wx + b \quad r = w^2$$

$$u = \sigma(z) \quad \mathcal{L} = l + \lambda r$$

$$l = \frac{1}{2}(y - u)^2$$



$$\frac{\partial \mathcal{L}}{\partial b}$$

$$\frac{\partial \mathcal{L}}{\partial w}$$

Forward pass
plug it in!

Backward pass

$$\frac{\partial \mathcal{L}}{\partial \mathcal{L}} = 1 \quad \frac{\partial \mathcal{L}}{\partial r} = \lambda \quad \frac{\partial \mathcal{L}}{\partial l} = 1$$

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial l} \cdot \frac{\partial l}{\partial u} = 1 \cdot (u - y)$$

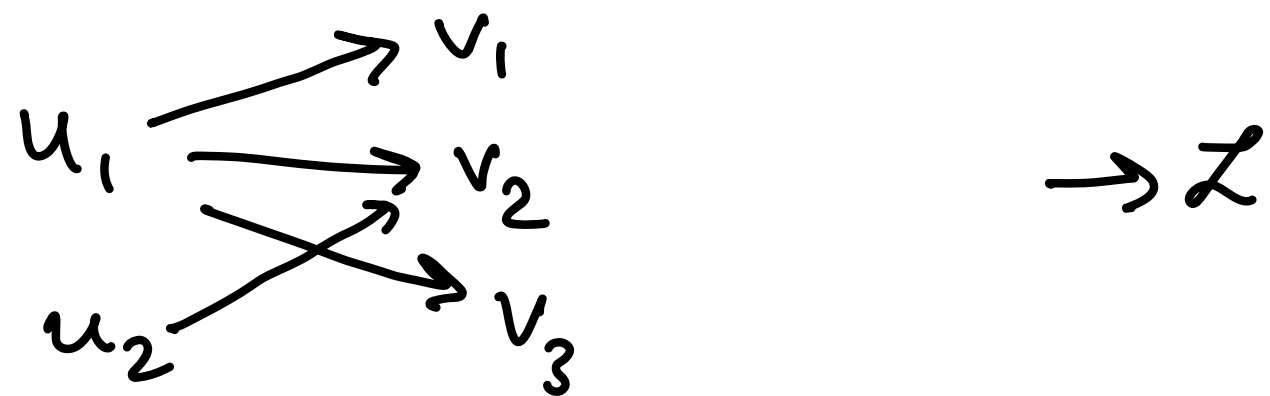
$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial u} \cdot \frac{\partial u}{\partial z} = \frac{\partial \mathcal{L}}{\partial u} \cdot \sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial w} + \frac{\partial \mathcal{L}}{\partial r} \cdot \frac{\partial r}{\partial w}$$

\nwarrow \nearrow
 x w

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial b}$$

\nwarrow
 1



Forward pass

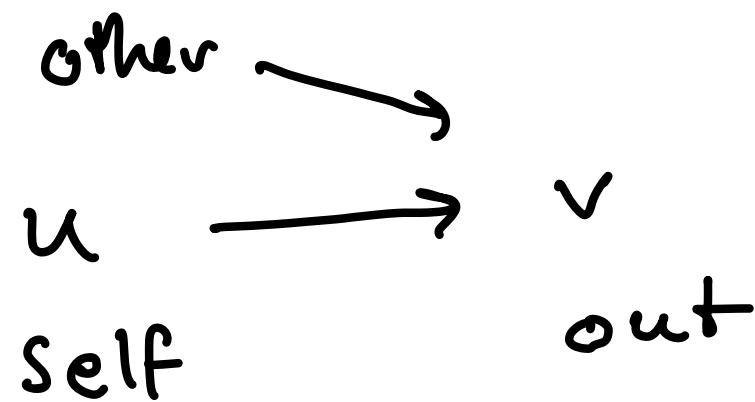
$$v_i = v_i(u_1, u_2)$$

Backward pass

$$\frac{dZ}{dout} = out.grad$$

$$\frac{\partial \mathcal{L}}{\partial u_i} = \sum_v \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

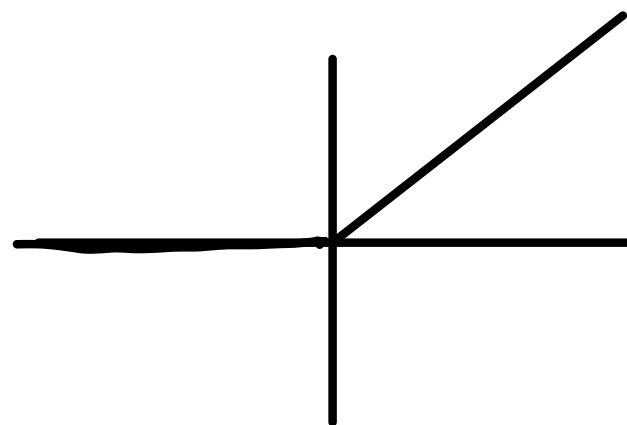
$\frac{dout}{dself}$



out $\frac{dout}{dself}$
 forward backward

| | | |
|------|--------------|-------------------------|
| add | self+other | 1 |
| mul | self*other | other |
| pow | self**other | other*self**(other-1) |
| relu | max(0, self) | 1 if self > 0 0 else |

Sigmoid



Reminders

- Form: 23 / 26

 - ↳ timestamp

 - ↳ I will not respond individ.

- Homework due Friday!

 - ↳ five questions

 - ↳ try LaTeX

 - ↳ solve before

 - ↳ update source

- Canvas discussion!

GPU

video games

↳ parallel cores

↳ large peak memory

Forward "embarrassingly parallel"

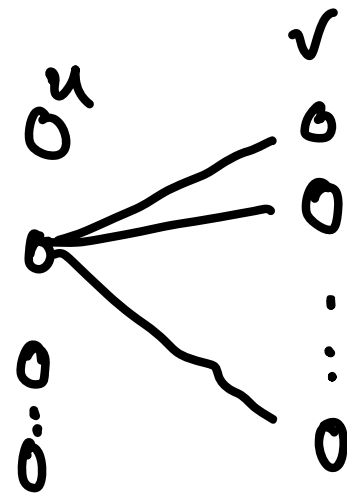
$$\sigma(Wu) = \left[\sigma(kw^{(1)}, w) \right]$$

Backward

"embarrassingly parallel"

$$\frac{\partial \mathcal{L}}{\partial u} = \left[\frac{\partial \mathcal{L}}{\partial u_j} \right] = \begin{bmatrix} \dots & \frac{\partial v_i}{\partial u_j} & \dots \end{bmatrix} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial v_i} \end{bmatrix}$$

↑
vector



1 hr → 20 seconds

$S \sim$ sample

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i) \quad \text{vs.} \quad \frac{1}{|S|} \sum_{i \in S} \ell(y_i, \hat{y}_i)$$

$\mathcal{L}(w)$ $\tilde{\mathcal{L}}(w)$

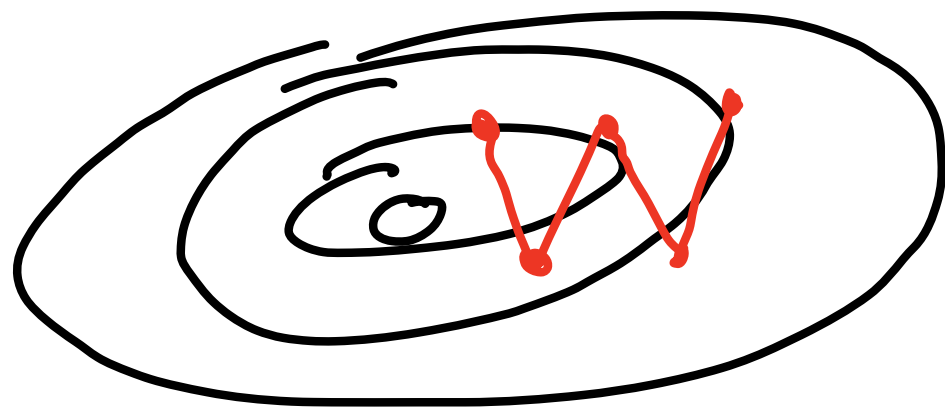
$$\nabla \tilde{\mathcal{L}}(w)$$

Variants of SGD

$$w_{t+1} \leftarrow w_t - \eta \nabla \mathcal{L}(w_t)$$

Problems:

- oscillate a lot
- get stuck

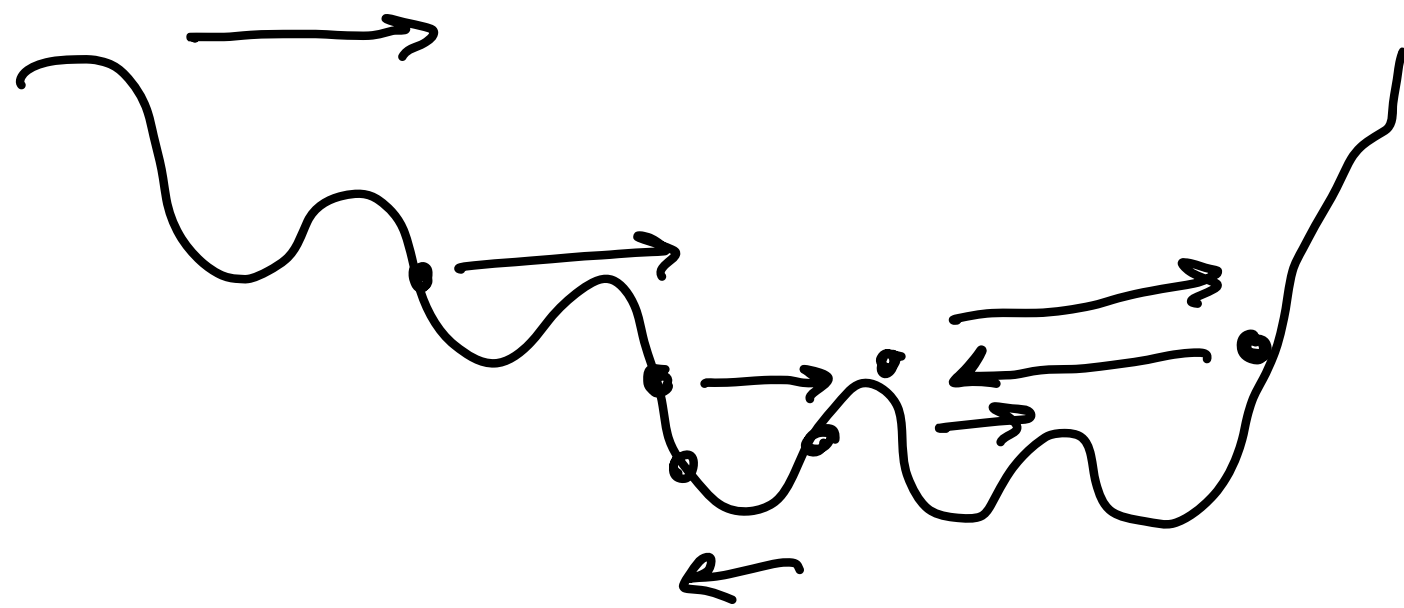


Momentum

keeps memory

$$v_t \leftarrow \mu \overset{\text{"mu"}}{v_{t-1}} + (1-\mu) \nabla \mathcal{L}(w_t)$$

$$w_{t+1} \leftarrow w_t - \eta v_t$$



Adaptive Learning Rate

$$g_t \leftarrow \nabla L(w)$$

$$s_t \leftarrow s_{t-1} + \|g_t - \text{mean}(g_t)\|_2^2$$

$$w_{t+1} \leftarrow w_t - \eta \frac{1}{\sqrt{s_t} + \epsilon} g_t$$

Adam (Adaptive Moment Estimate)

$$v_t \leftarrow \frac{\mu v_{t-1} + (1-\mu)g_t}{1-\mu^t}$$

$$1-\mu^t \leftarrow \text{Idk}$$

$$s_t \leftarrow \frac{\beta \overset{\text{"beta"}}{s_{t-1}} + (1-\beta) \|g_t - \text{mean}(g_t)\|_2^2}{1-\beta^t}$$

$$1-\beta^t$$

$$w_{t+1} \leftarrow w_t - \frac{\eta}{\sqrt{s_t} + \epsilon} v_t$$

Both!!
:)

Initialization

- vanishing/exploding gradient

$$f_w(x) = W_L W_{L-1} \dots W_1 x$$

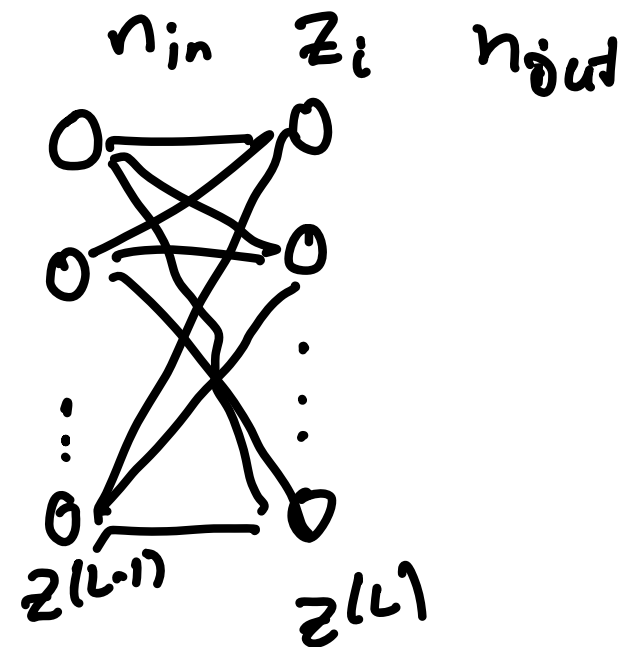
ignoring activation

$$\mathcal{L}(w) = \frac{1}{2} \|f_w(x) - y\|_2^2$$

$$\frac{\partial \mathcal{L}(w)}{\partial w}$$

$$= (f_w(x) - y) W_L W_{L-1} \dots W_{r+1} W_{r-1} \dots W_1 x$$

✓ • symmetry



✓ randomly initialize

from distribution with "good" norm

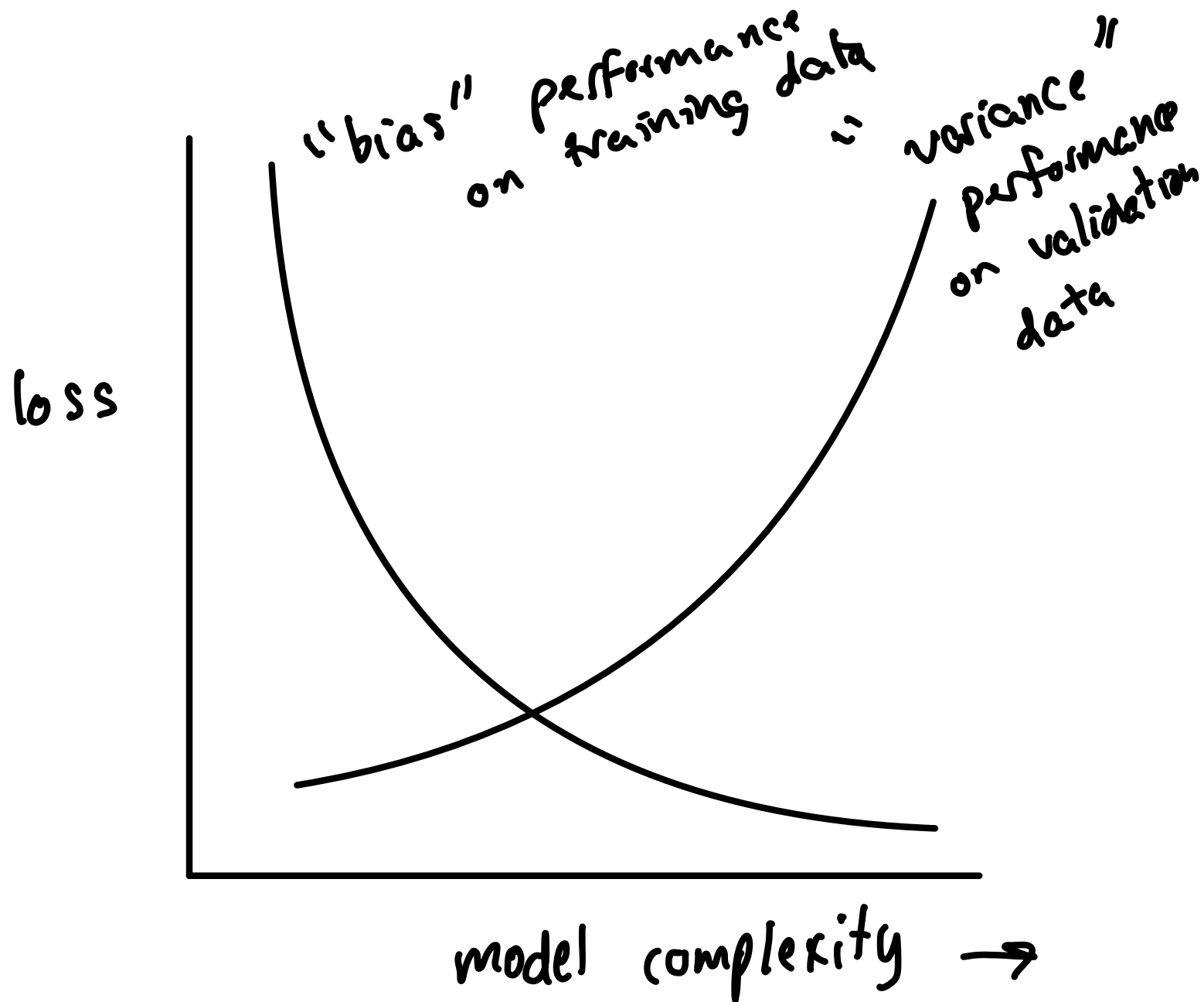
$$z_i^{(L)} = \sum_{j=1}^{n_{in}} W_{ij}^{(L)} \cdot z_j^{(L-1)}$$

$$\text{var}(z_i^{(L)}) = \underbrace{n_{in} \cdot \text{var}(w_{ij})}_{\approx 1} \cdot \text{var}(z_i^{(L-1)})$$

$$\text{var}(w_{ij}) = \frac{2}{n_{in} + n_{out}} \quad w_{ij} \sim \mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right)$$

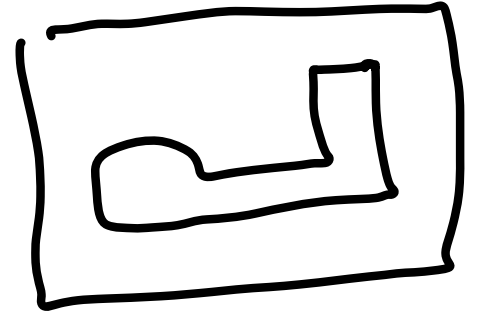
Generalization

↳ want good performance on unseen data

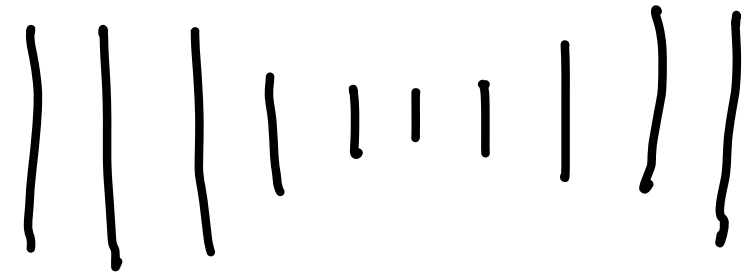


"Regularize"

↑ reduce complexity

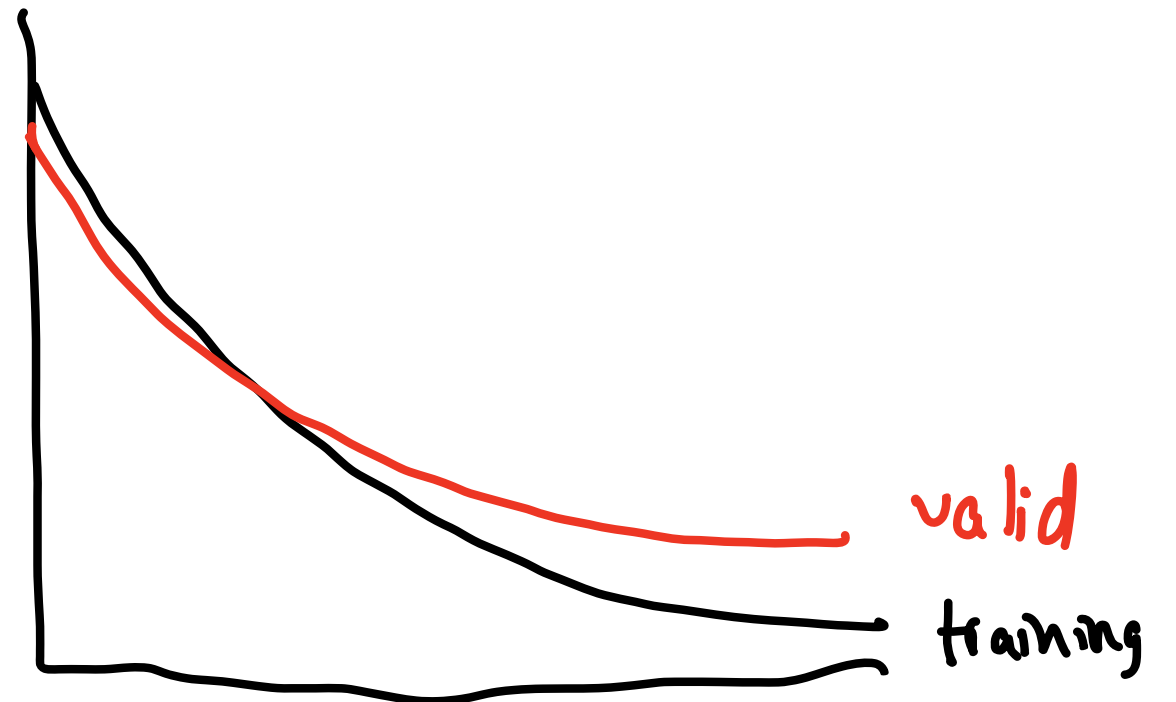


• bottleneck



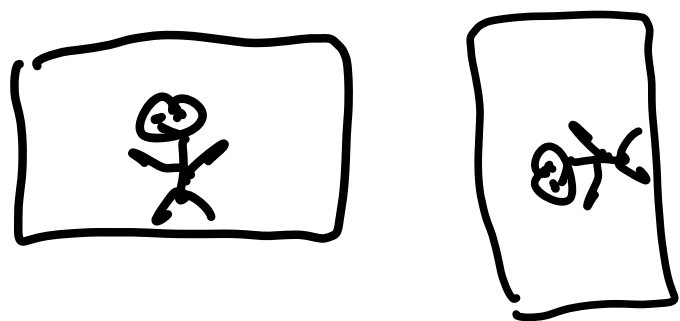
Stop early

•

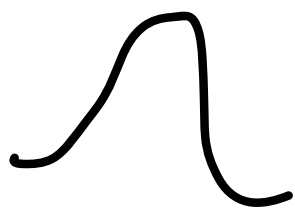
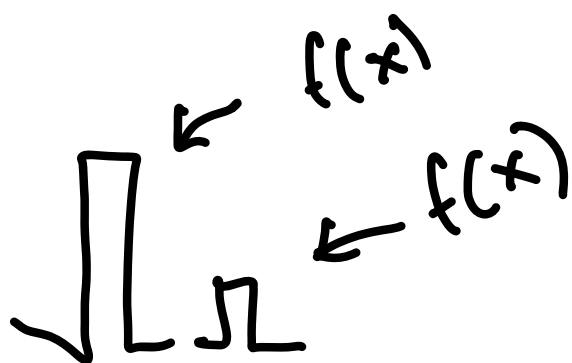


key: move from "memorization" to "learning"

- augment data



- weight decay $\mathcal{L}(w) + \lambda \|w\|_2^2$



- dropout $\in \{0, 1\}$
 $v = m_i \sigma(u)$

- transfer learning

Hint

```
import torch.nn as nn
```

```
nn.Sequential
```

```
nn.Linear
```

```
nn.ReLU
```

```
nn.Softmax
```