# Midterm Review

Autograd: Decomposing model into atomic operations, making a computational graph

ResNets: Know what a residual connection is and design considerations

# Linear Regression

$$\hat{y} = w^T \underset{\text{input}}{x} + b$$

$$\underset{\text{parameters}}{}$$

$$L = \frac{1}{2}(\hat{y} - y)^2$$

$$\underset{\substack{\text{true} \\ \text{value}}}{}$$

$$w \leftarrow w - \underset{\substack{\text{learning} \\ \text{rate}}}{\eta} \nabla_w L$$

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dw} = (\hat{y} - y)x$$

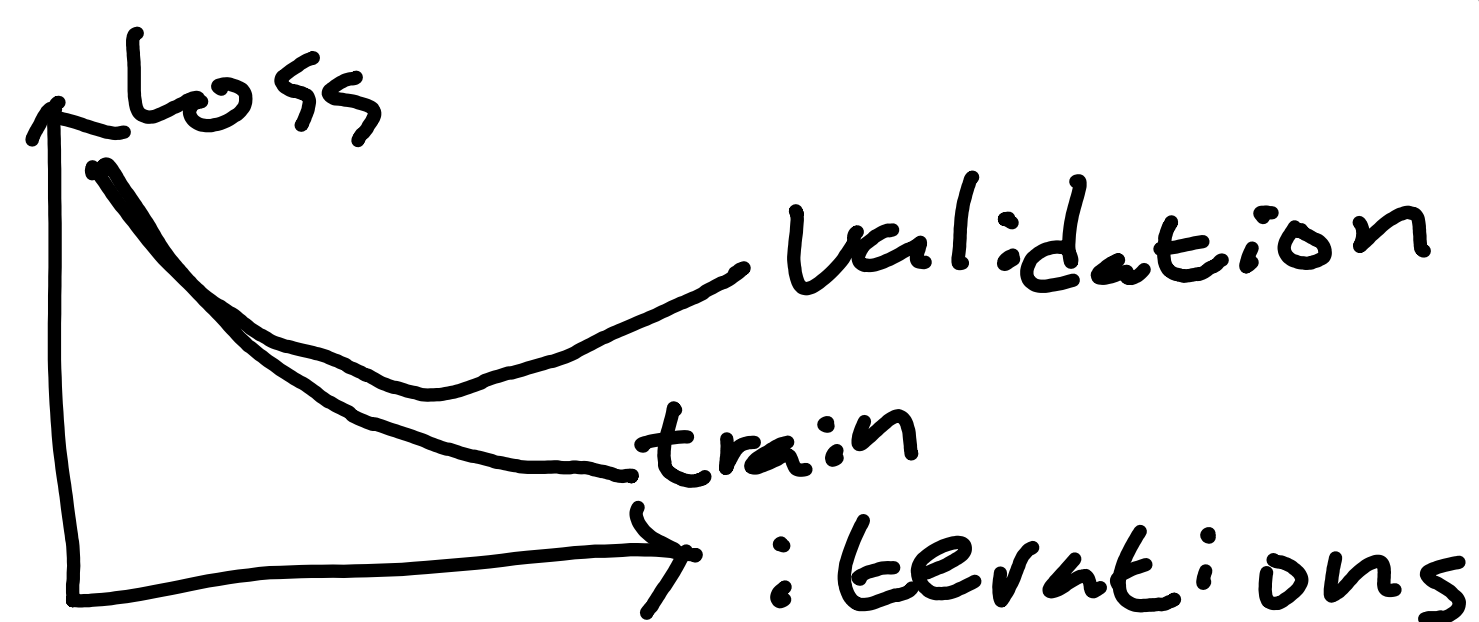# Softmax Regression

(logistic)

$$o = Wx + b$$

"scores"

$$\hat{y} = \text{softmax}(o) \qquad \hat{y}_i = \frac{\exp(o_i)}{\sum_k \exp(o_k)}$$

$$L = -\sum_j y_j \log \hat{y}_j = -\log \hat{y}_i \longrightarrow \text{correct class}$$

label
(one-hot)

$$\frac{dL}{do_j} = \text{softmax}(o)_j - y_j$$

# Overfitting and regularization

loss

validation

train

iterations

Weight decay:
add loss term $\lambda \|w\|^2$

$\hookrightarrow$ parameter vector

Dropout:

$$h' = \begin{cases} 0 & \text{with probability } p \\ n/1-p & \text{with probability } 1-p \end{cases}$$

# Adaptive gradient descent methods

Momentum:

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

where $g_t$ is the gradient at iteration $t$: $\nabla_\theta L$

$$\theta_t \leftarrow \theta_{t-1} - \eta v_t$$

Adam:

(complex update)

$$g_t' = \frac{\eta v_t}{\sqrt{\hat{s}_t} + \varepsilon}$$

# MLP
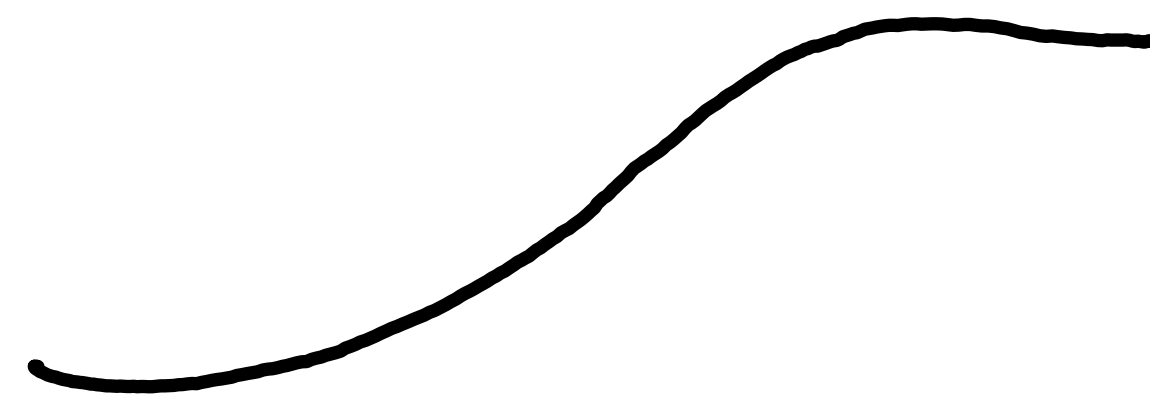
feed forward, dense, fully-connection, "DNN"

$$h = \phi(W_1 x + b_1)$$

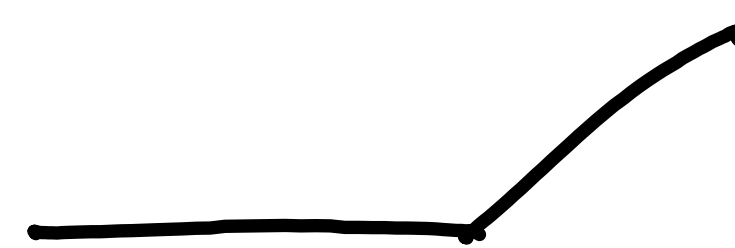hidden rep.    non-linearity

$$o = W_2 h + b_2$$

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

$$\frac{d\, \text{sigmoid}}{dx} = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$$

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

$$\frac{d\, \text{ReLU}}{dx} = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

# Backpropagation

$$z_m = W_m a_{m-1} + b_m$$

$$a_m = \phi_m(z_m)$$

$$C = \text{cost function of } a_L \text{ and } y$$

$$\frac{dC}{dW_m} = \frac{dC}{dz_m} a_{m-1}^T$$

$$\frac{dC}{dz_m} = \left( \underbrace{W_{m+1}^T \underbrace{\frac{dC}{dz_{m+1}}}} \right) \circ \frac{da_m}{dz_m}$$

recursion!

# Initialization

- Assume $x \sim \mathcal{N}(0, I)$
- Assume initial $w \sim \mathcal{N}(0, \sigma^2 I)$
- Ignore nonlinearities

$$\sigma = \sqrt{\frac{2}{n_{in} + n_{out}}}$$

# Autograd

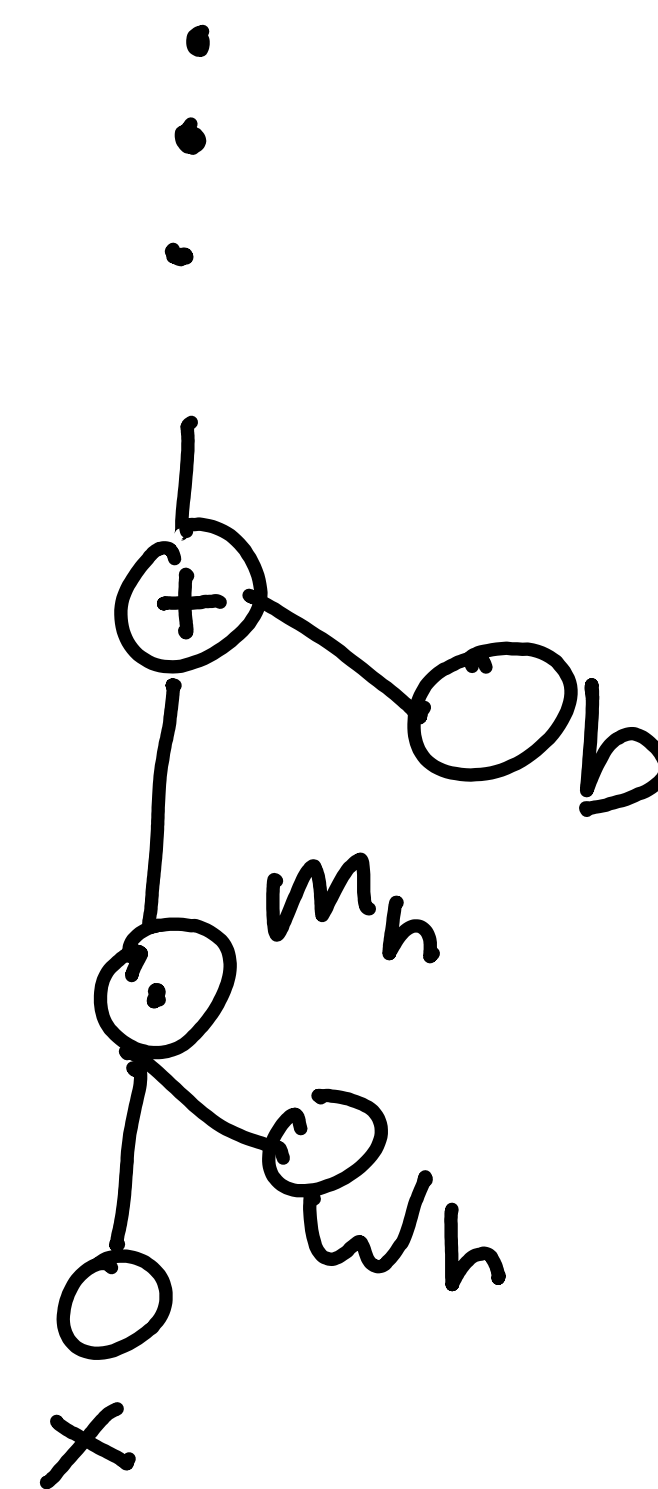$$h = ReLU(W_n x + b_n)$$

$$o = W_o h + b_o$$

$$L = (y - o)^2$$

---

$$m_h = W_n x$$
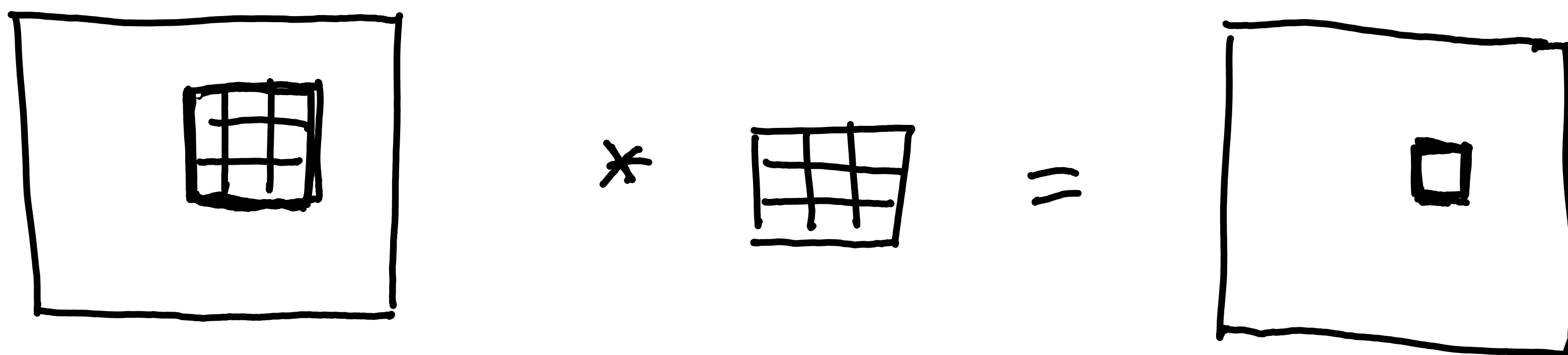
$$z_n = m_n + b_n$$

$$h = ReLU(z_n)$$

$$\vdots$$

$$\frac{dL}{dW_n} = \frac{dL}{de} \frac{de}{do} \frac{do}{m_o} \dots$$

# Convolution

$$H_{ijd} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} \sum_{c} V_{abcd} X_{i+a, i+b, c}$$



Other factors:
- padding
- striding
- multiple input and output channels
- receptive field

# ConvNet ingredients

- Convolutions
- Pooling - apply a reduction over a region of the input

- Dense layers

# Batch Norm

$$\text{Batch Norm}(x) = \gamma \odot \frac{x - \hat{\mu}_B}{\hat{\sigma}_B + \varepsilon} + \beta$$

$$\hat{\mu}_B = \frac{1}{|B|} \sum_{x \in B} x$$

$$\hat{\sigma}_B^2 = \frac{1}{|B|} \sum_{x \in B} (x - \hat{\mu}_B)^2$$

$\gamma, \beta$ learnable params, same shape as $x$

# Res Nets

Residual connection: $f(x) + x$
(optionally process $x$ to match shape)
1x1 conv