

# Neural nets

Nets: transform our data into a latent featurespace

$\bar{w}^T f(\bar{x})$  replace  $f(\bar{x})$  with a nonlinear func of the  
of  $f(\bar{x})$

Define  $\bar{z} = g(Vf(\bar{x}))$   
 nonlinearities  $\underbrace{\quad}_{\text{d} \times n \text{ matrix } X}$   
 $\underbrace{V}_{n \text{-dimensional feature vector}}$

How can  $V + g$  give us useful  
latent features?

Classification with  $\bar{w}^T \bar{z}(\bar{x})$

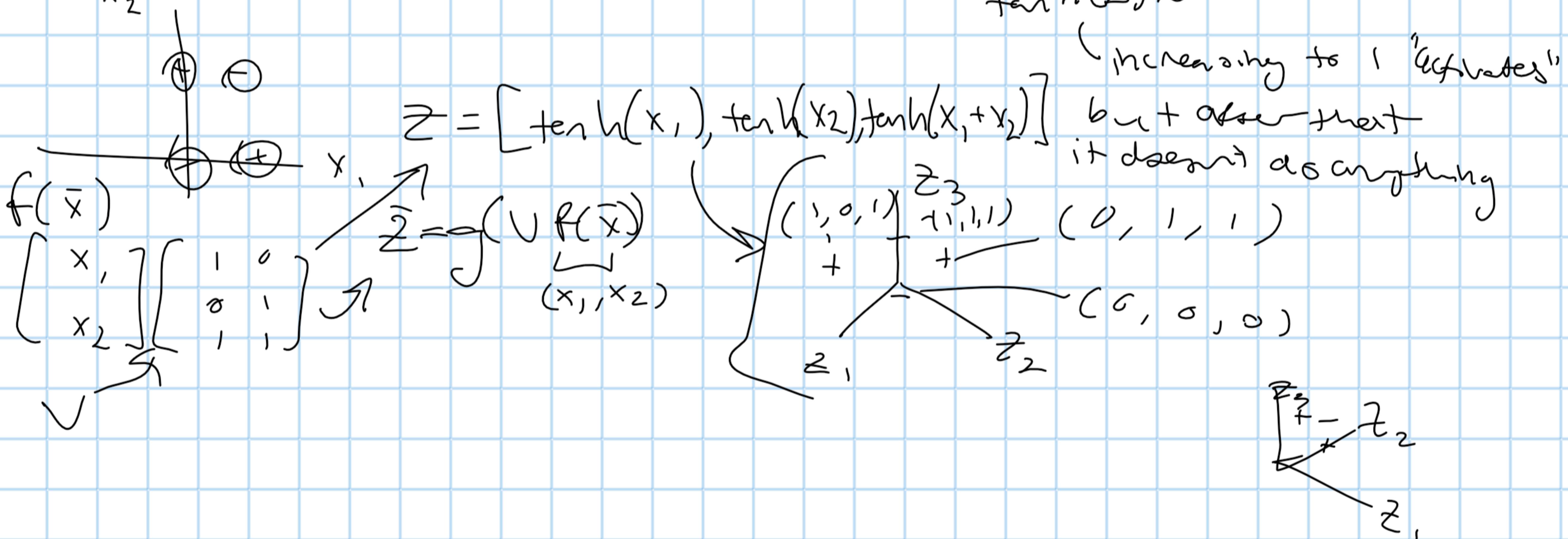
$$g = \tanh \cancel{\text{+}}$$

$$g = \text{ReLU} \cancel{\text{+}}$$

$$\cancel{\text{+}}$$

Ex:

Suppose  $V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$   $g = \tanh$   
 $x_1 \quad \tanh(0) = 0, \tanh(1) \approx 1$   
 $x_2 \quad \tanh(2) \approx -1$



## Neural Network Visualization

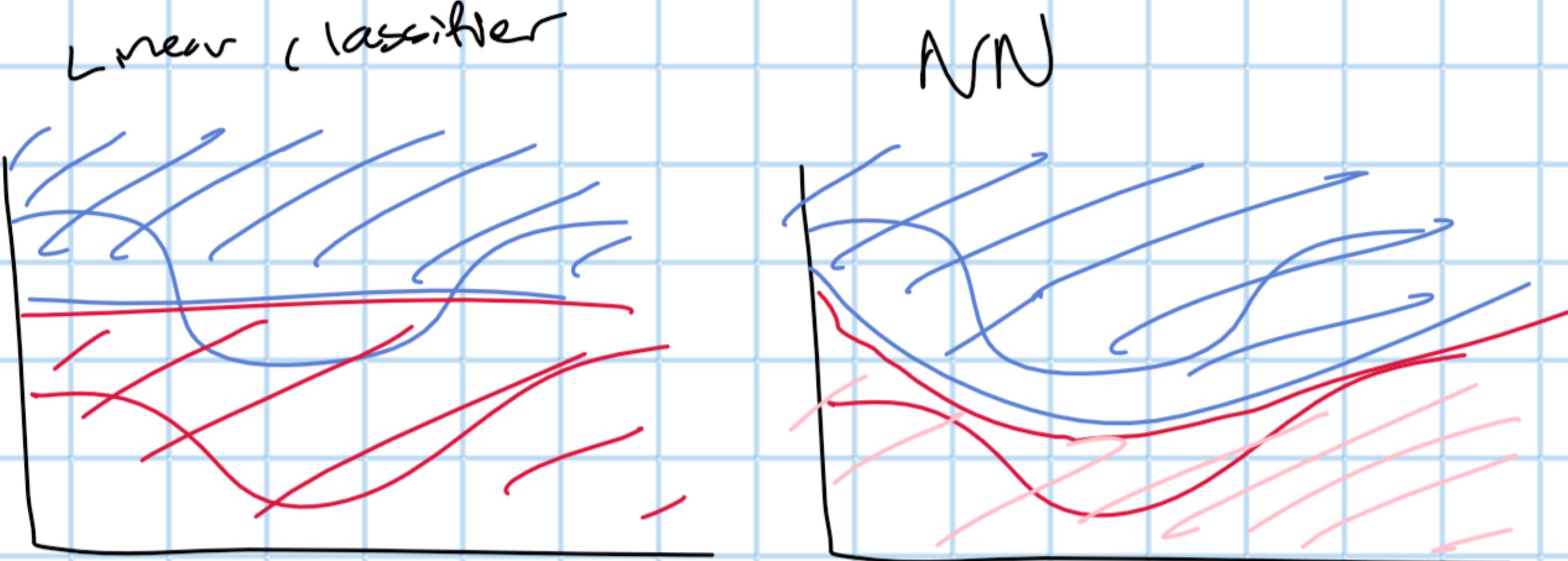
$$z = \tilde{g}(\nabla f(x) + b)$$

non-linear warp shift

$$y_{\text{pred}} = \arg \max_y w_y^T z$$

Linear classifier

we can ignore  $+b$



## Deep Neural Networks

$$z_1 = g(V_1 f(x))$$

$$z_2 = g(V_2 z_1)$$

...

$$y_{\text{pred}} = \arg \max_g w_g^T z_n$$

multiple layers

learn transformation with backpropagation

## Feedforward Neural Networks, Backpropagation

$$LR: P(y|x) = \frac{\exp(w_y^T f(x))}{\sum_{y \in \gamma} \exp(w_y^T f(x))}$$

single scalar probability

3 classes "diff weights"	$w_1^T f(x)$	-1.1	softmax	0.036	
	$w_2^T f(x)$	2.1	→	0.89	class
	$w_3^T f(x)$	-0.4		0.07	probs

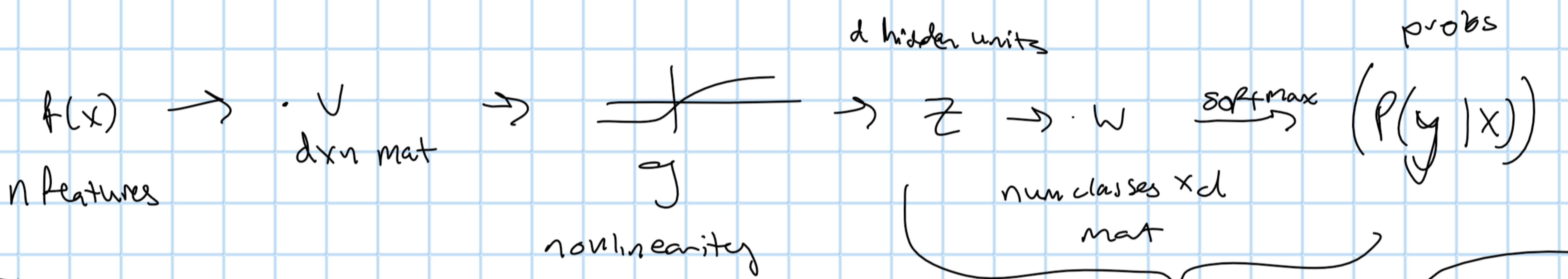
Softmax operation = "exp and normalize"

$$\text{softmax}(w^T f(x)) = P(y|x)$$

vector weight vector per class  
w is [num classes x num feats]

$$P(y|x) = \text{softmax}(w_g(Vf(x)))$$

1 hidden layer



Training

$$P(y|x) = \text{softmax}(wz) \quad z = g(Vf(x))$$

$$\mathcal{L}(x, i^*) = \log P(y=i^*|x) = \log (\text{softmax}(wz)e_{i^*})$$

maximize log likelihood of training data  
in a 3 class problem, if the  $i^*$  is in the 2nd pos,  
 $i^*$ : index of the gold label      use  $e_2$

$e_i$ : 1 in the  $i$ th row, zero elsewhere      dotting by this  
    \ num classes x 1 vector

$$\mathcal{L}(x, i^*) = wz \cdot e_{i^*} - \log \sum_j \exp(wz) \cdot e_j$$

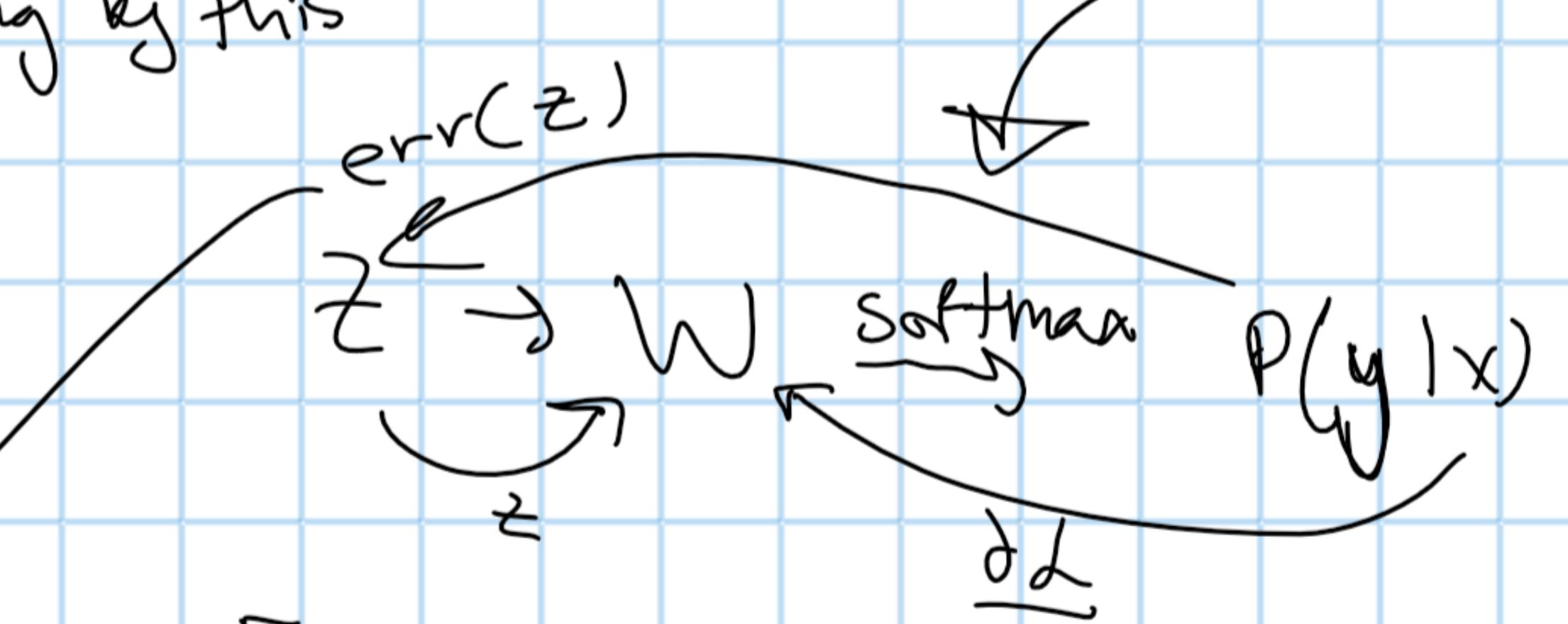
can forget everything  
after  $z$ , treat it as the  
output and do bp

Computing Gradients with backprop:

$$\mathcal{L}(x, i^*) = wz \cdot e_{i^*} - \log \sum_j \exp(wz) \cdot e_j$$

gradient wrt  $w$ : apply chain rule

$$\frac{\partial \mathcal{L}(x, i^*)}{\partial w_{ij}} = \underbrace{\frac{\partial \mathcal{L}(x, i^*)}{\partial z}}_{\text{error}} \cdot \underbrace{\frac{\partial z}{\partial w_{ij}}}_{\text{error}}$$



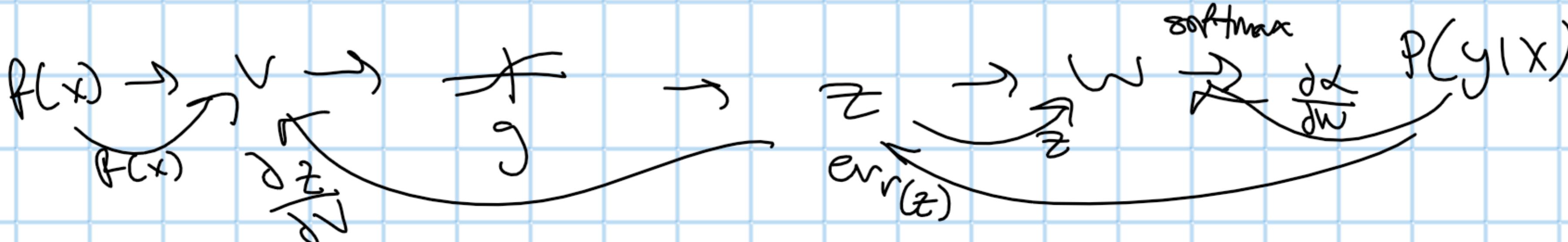
$\nabla$  wrt  $w$ :  
looks like LR func, can be  
computed treating  $z$  as the features

$$z = g(Vf(x))$$

activations at hidden layer

$$\frac{\partial z}{\partial w_{ij}} = \frac{\partial g(a)}{\partial a} \cdot \frac{\partial a}{\partial w_{ij}}$$

$a = Vf(x)$



# Neural Network Implementation

automatic differentiation

$$y = x^* x \rightarrow (y, dy) = (x^* x, 2 \cdot x \cdot dx)$$

`torch.nn.Module`

takes an example  $x$  and computes result  
`forward(x)`.

Computes gradient after `forward` is called

`backward()` is produced automatically

Training and Optimization

```
class FFNN(nn.Module):
    def __init__(self, input_size, hidden_size, out_size):
        super(FFNN, self).__init__()
        self.V = nn.Linear(input_size, hidden_size)
        self.g = nn.Tanh() # or nn.ReLU(), sigmoid()...
        self.W = nn.Linear(hidden_size, out_size)
        self.softmax = nn.Softmax(dim=0)

    def forward(self, x):
        return self.softmax(self.W(self.g(self.V(x))))
            # syntactic sugar for forward
```

```
P(y|x) = softmax(Wg(Vf(x))) one-hot vector
        of the label
        (e.g., [0, 1, 0])
ffnn = FFNN(inp, hid, out)
optimizer = optim.Adam(ffnn.parameters(), lr=lr)
for epoch in range(0, num_epochs):
    for (input, gold_label) in training_data:
        ffnn.zero_grad() # clear gradient variables
        probs = ffnn.forward(input)
        loss = torch.neg(torch.log(probs)).dot(gold_label)
        loss.backward() negative log-likelihood of correct answer
        optimizer.step()
```

Training:

Define modules, etc.

Init weights and optimizer

For each epoch

    For each batch of data

        Zero gradient

        Compute loss on batch

        Autograd to compute gradients and take step on optimizer

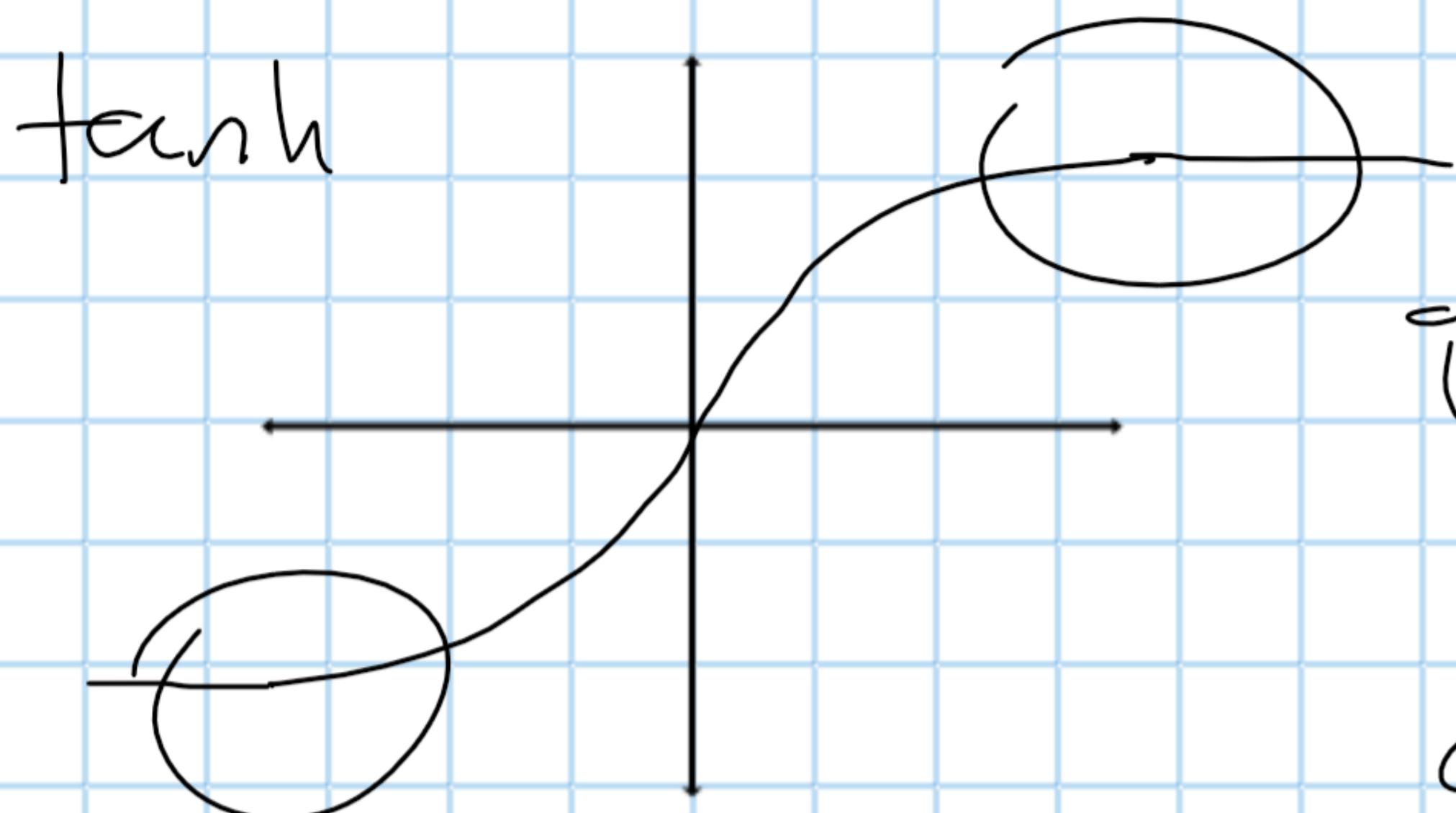
(Optional) check performance on dev set so no overfitting

Run on dev/test set

# Neural Net Training, Optimization

## Initialization

Nonconvex problem, so init matters (ex:  $0^\top$  for  $\nabla$  will not work)

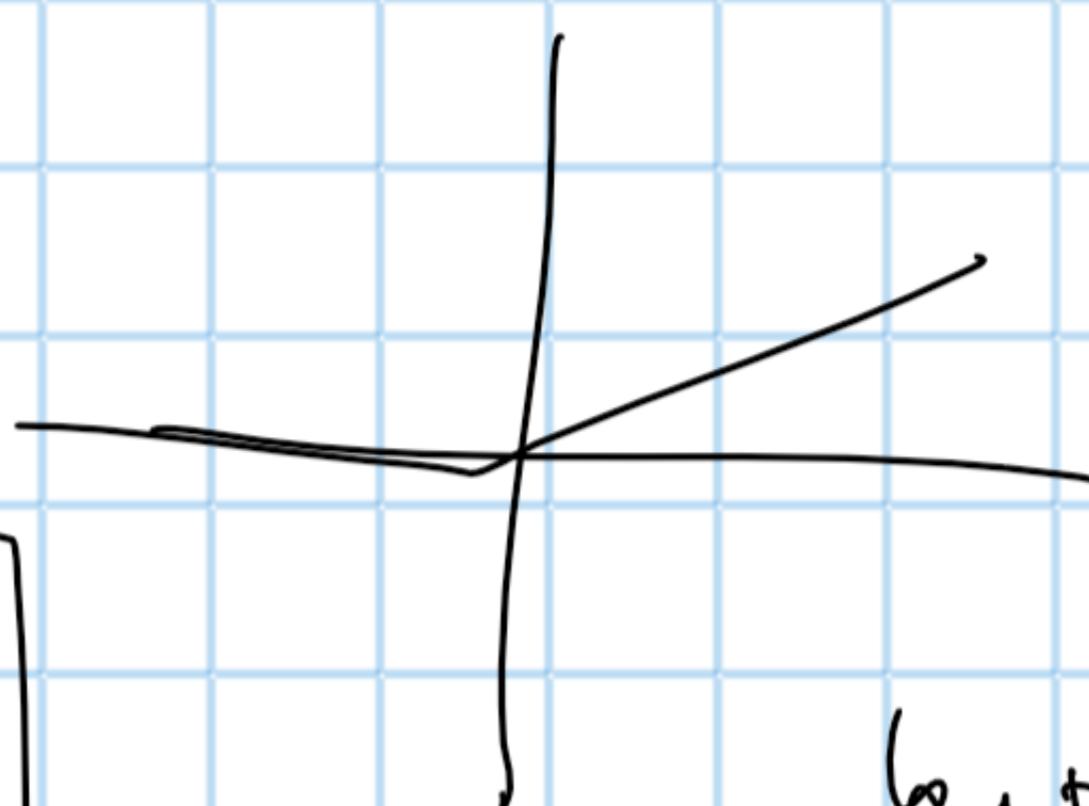


if gradients are small here, so model thinks that there is little change. So if inputs are too big (+ or -), the model won't change. So to avoid, you can use ReLU

can do random initialization with appropriate scale

$$\text{Glorot init: } U \left[ -\sqrt{\frac{6}{f_{\text{in}} + f_{\text{out}}}}, +\sqrt{\frac{6}{f_{\text{in}} + f_{\text{out}}}} \right]$$

want variance of inputs & gradients to be the same



but if many neg nums, some problem can happen

batch normalization is rescaling each layer to have mean 0 and var 1 over a batch  
this is useful if the net is deep (more applicable in CV than NLP) (Ioffe & Szegedy, 2015)

Dropout: zero out some neurons during training to prevent overfitting

too much dropout leads to bad performance, good dropout has some benefit

## Optimizer

Adam (Kingma and Ba, ICLR 2015) : adaptive step + momentum

gradient clipping sets max value for gradients