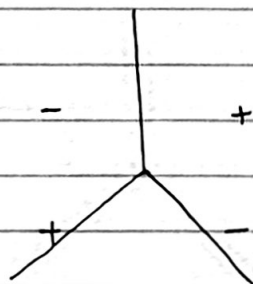


Neural Networks

	$f(\bar{x})$		
	G	B	N
good	1	0	0
bad	0	1	0
not good	1	0	1
not bad	0	1	1



No linear classifier
can classify
this

- Neural Nets : Transform data into latent feature space

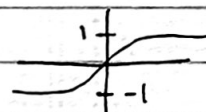
$\bar{w}^T f(\bar{x})$ \rightarrow replace $f(\bar{x})$ with a non-linear function of the original $f(\bar{x})$

Define $\bar{z} = g(\underbrace{V}_{d \times n} \underbrace{f(\bar{x})}_{n\text{-dimensional feature vector}})$

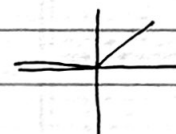
non-linearity \rightarrow g

Types of non-linearity

$g = \tanh$



$g = \text{ReLU}$

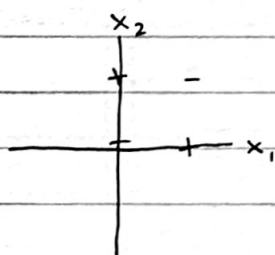


How can $V + g$ give us useful latent features?

Ex. Suppose $V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$ $g = \tanh$

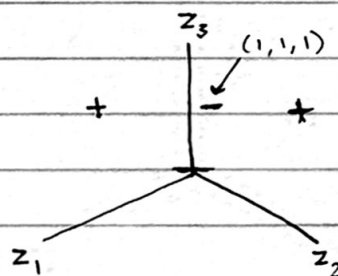
$\tanh(0) = 0$, $\tanh(1) \approx 1$, $\tanh(2) \approx 1$

$$\bar{z} = g(\underbrace{V f(\bar{x})}_{(x_1, x_2)}) = [\tanh(x_1), \tanh(x_2), \tanh(x_1 + x_2)]$$



OG Feature Space

Transform \rightarrow



New Feature Space

Features now
separable!

Feed Forward Neural Networks

• Vectorization and Softmax

Multiclass
Logistic Regression
"Diff. weights"

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

* Single scalar probability

3 classes	w_1^T	= 1.1	softmax	0.036	
"Different weights"	w_2^T	= 2.1	→	0.89	class probabilities
	w_3^T	= -0.4		0.07	

• Softmax operation: exponentiate and normalize

• Written as: $\text{softmax}(W f(x))$

↓
Divide by
their sum

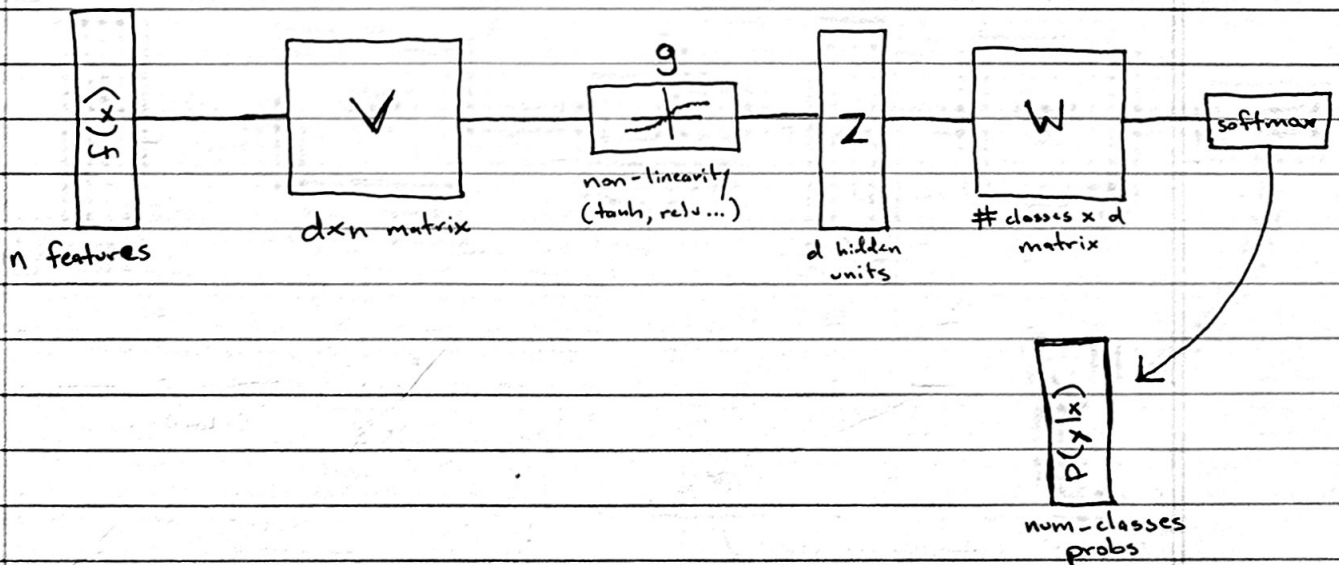
$$P(y|x) = \text{softmax}(W f(x))$$

* Weight vector per class

W is $[\# \text{ classes} \times \# \text{ feats}]$

$$P(y|x) = \text{softmax}(W g(v f(x)))$$

* Now one hidden layer



• Training Neural Nets

stage 2

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

stage 1

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

- Maximize log likelihood of training data :

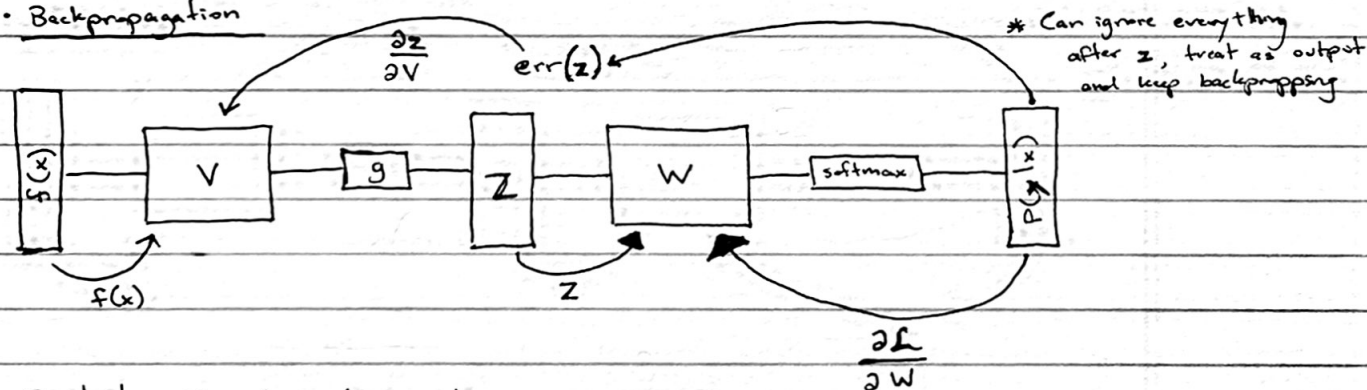
$$\mathcal{L}(x, i^*) = \log P(y = i^* | x) = \log(\text{softmax}(Wz) \cdot e_{i^*})$$

i^* = index of gold label

$e_i = 1$ in the i^{th} row, 0 elsewhere. Dot ^{by this} ~~matrix~~ = select i^{th} index

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

• Backpropagation



Gradient w.r.t. V : chain rule

$$\frac{\partial \mathcal{L}(x, i^*)}{\partial V_{ij}} = \underbrace{\frac{\partial \mathcal{L}(x, i^*)}{\partial z}}_1 \cdot \underbrace{\frac{\partial z}{\partial V_{ij}}}_2$$

* Gradient w.r.t. W :
looks like logistic regression,
can be computed using z
as the features

First Term (1) : $\text{err}(z)$ represents gradient w.r.t. z

Second Term (2) :

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

gradient of non-linear
activation function at a
(depends on current value)

gradient of
linear function

Implementing Neural Nets

- Computing gradients is hard!
- Automatic Differentiation: instrument code to keep track of derivatives

$$y = x \cdot x \xrightarrow[\text{code gen}]{\text{code gen}} (y, dy) = (x \cdot x, 2 \cdot x \cdot dx)$$

- PyTorch: Framework for defining computations that provides easy access to derivatives

Module: defines a neural network (can use wrap other modules which implement predefined layers)

* If `forward()` uses crazy math, you have to write `backward()` yourself

`torch.nn.Module`:

Takes example `x` and computes result
`forward(x)`:
...

Computes gradient after forward called
`backward()`:
... produced automatically!

Code and Pseudocode on "seg-15.pdf"

- Batching: Batching data gives speedups due to more efficient matrix operations
 - Batch sizes from 1-100 often work well

- How does initialization affect learning?

- How do we initialize `V` and `W`?
- Non-convex problem. Initialization matters!
- Must initialize to non-zero to get model to learn
- If you initialize too large, cells become saturated, gradient close to 0, no learning
- Can do random uniform/normal initialization w/ appropriate scale

`fan-in`: # inputs
`fan-out`: # outputs

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

* Want variance and gradient of inputs for each layer to be the same

• Dropout : Probabilistically zero out parts of network during training to prevent overfitting. Use whole network at test time

- Form of stochastic regularization
- Similar to benefits of ensembling (sub-networks)
- One line of code

• Gradient Clipping : Set max value for gradients