In this cause we will need to learn the weights of the neural network. As we have discussed in lecture that we will be using gradient descent to learn the network weights.  $M(K+t) = M(K) - E \Delta^{M} \chi$ In the gradient descent step we will need to compute Vwol. In a neural network the weights in the earlier layers are connected to the 1055 function through a composition of functions, and in such cases computing the gradients require repeated application of the chain rule.

In this discussion we will go over  
a computational mechanism that will  
enable us to compute the gradients  
in an efficient and modular manner.  
Chain rule revisited:  
Let's consider the following regularized  
linear classification model  
$$Z = WX + b$$
  
 $y = \sigma(Z)$   
 $d = \frac{1}{2}(y - t)^2$   
 $R = \frac{1}{2}W^2$   
 $\chi(reg = \alpha' + \beta R)$ 

-



$$\frac{\partial \sqrt{reg}}{\partial b} = \frac{\partial}{\partial b} \left[ \frac{1}{2} (t(wx+b)-t) + \frac{\partial}{2} w^{2} \right]$$

$$= \frac{1}{2} \frac{\partial}{\partial b} \left[ t(wx+b)-t \right] + \frac{\partial}{2} \frac{\partial}{\partial b} w^{2}$$

$$= \left[ t(wx+b)-t \right] \frac{\partial}{\partial b} \left[ t(wx+b)-t \right] + 0$$

$$= \left[ t(wx+b)-t \right] t(wx+b) \frac{\partial}{\partial b} (wx+b) \right]$$

$$= \left[ t(wx+b)-t \right] t(wx+b) \frac{\partial}{\partial b} (wx+b)$$
The above computation involves
(a) Lots of copies of terms from one line to next
(b) Lots of redundant work. For
(c) Instance, the first 3 steps in the two derivations above are nearly identical.

The idea behind backprop is to Share repeated computations wherever possible. Computational graphs and backprop: Computational graphs help to visualize and contextualize how we aim to backpropagate Derivatives. It is the mechanism by which deeplearning libraries compute gradients. A computational graph is a directed acyclic graph where each node in the graph denote a mathematical operation.

Let's draw the computational graph for our running example,





The backprop algorithm consists of two steps:

(i) In forward pass we take
the input X and the initial values
of W and b and propagate
in the forward Direction through
the network to compute the quantities
in red.
(ii) In backward pass we compute
the derivatives starting at the
output and propagating it backward
to the Input.





$$\begin{array}{l} z = b + m \\ \hline \frac{\partial z}{\partial b} = 1, \quad \frac{\partial z}{\partial m} = 1 \\ \hline \frac{\partial x reg}{\partial b} = \frac{\partial z}{\partial b} \quad \frac{\partial x reg}{\partial z} \\ = \sigma'(z)g = \sigma'(\omega x t b) [\sigma(\omega x t b) - b] \\ = \sigma'(z)g = \sigma'(\omega x t b) [\sigma(\omega x t b) - b] \end{array}$$



$$m = q_{1} \times$$

$$\frac{\partial m}{\partial q_{1}} = \chi \quad \frac{\partial m}{\partial \chi} = q_{1}$$

$$\frac{\partial dreg}{\partial \chi} = \frac{\partial m}{\partial \chi} \quad \frac{\partial dreg}{\partial m}$$

$$= q_{1} \sigma^{1}(z) g$$

$$\frac{\partial dreg}{\partial q_{1}} = \frac{\partial m}{\partial q_{1}} \quad \frac{\partial dreg}{\partial m}$$

$$= \chi \sigma^{1}(z) g$$





So, 
$$\mathcal{L} = ||m||_2^2$$
  
 $\nabla_m \mathcal{L} = 2m \in \mathbb{R}^n (\epsilon a^n |31)$   
 $m = g - x$   
 $\nabla_g m = I \in \mathbb{R}^{n \times n}$ 





..



From the computational graph,  

$$L_1 = -\frac{\gamma}{2} \log |K| = -\frac{\gamma}{2} \log (\det K)$$

Using matrix 
$$Gokbook$$
  
 $\frac{\partial L_{1}}{\partial k} = -\frac{1}{2} (k^{-1})^{T} = -\frac{1}{2} (k^{T})^{-1}$   
 $\frac{\partial L_{1}}{\partial k} = -\frac{1}{2} (k^{-1})^{T} = -\frac{1}{2} (k^{T})^{-1}$   
 $\frac{\partial L_{1}}{\partial k} = -\frac{1}{2} (k^{-1})^{T} = -\frac{1}{2} (k^{T})^{-1}$ 



Since AT has the same elements as A, so taking the gradient of Qi w. G. J AT we have

$$\frac{\Im L_{1}}{\Im A^{T}} = A^{T} \frac{\Im L_{1}}{\Im N}$$

$$= A^{T} \left[ -\frac{\sqrt{n}}{2} \left( k^{T} \right)^{-1} \right]$$

$$\frac{\Im L_{1}}{\Im A^{T}} = -\frac{\sqrt{n}}{2} \left( k^{T} A^{T} \right)^{T}$$
Now adding the gradient paths through  
A and AT we get
$$\frac{\Im L_{1}}{\Im A} = -\frac{\sqrt{n}}{2} \left( k^{T} \right)^{-1} A - \left[ \frac{\sqrt{n}}{2} \left( k^{-1} A \right)^{T} \right]^{T}$$
Since K is symmetrics
$$\frac{\Im L_{1}}{\Im A} = -\frac{\sqrt{n}}{2} \left( k^{-1} A \right)^{T}$$

5) The computational graph for  $\sqrt{2}$  is

grawn below:













Using the same reasoning as (a) to Compute the derivative of Ly with respect to AT DAT = - 2 AT P-1 DAT P-1  $= \frac{\alpha}{2} A^{T} P^{-1} B B^{T} P^{-1}$ Now adding the gradients parts through A and AT we get: OLZ + OLZ OA OA  $= \stackrel{\propto}{\gamma} p^{-1} B B^{T} p^{-1} A + \left[ \stackrel{\propto}{\gamma} A^{T} p^{-1} B B^{T} p^{-1} \right]^{T}$  $= \chi P' B B^T P' A + \chi (PT)' B B^T (PT)' A$ 



•



Using chain rule, backpropagating to P, we have







since 
$$\oplus$$
 gate passes the  
gradient, so  
 $dL = 2^{22} \ln 2 \oplus 29 \text{ ERK}$   
 $du = 2^{22} \ln 2 \oplus 29 \text{ ERK}$   
 $dL = 2^{22} \ln 2 \oplus 29 \text{ ERK}$   
 $dL = 2^{22} \ln 2 \oplus 29 \text{ ERK}$ 



$$\frac{dL}{dh_{1}} = \frac{du}{dh_{1}} \frac{dL}{du}$$

$$\frac{dL}{dh_{1}} = \frac{W_{2}T}{du} \frac{dL}{du} \in i\mathbb{R}^{m}$$
By the 3-D knsor derivative trick  
learned in class,  

$$\frac{dL}{dW_{2}} = \frac{dL}{du} h_{1} = \frac{dL}{dW_{2}} h_{1} = \frac{dL}{du} h_{1} = \frac{dL}{dW_{2}} h_{1} = \frac{dL}{du} h_{1} = \frac{dL}{dW_{2}} h_{1}$$



Now, t- Wix Computing the local gradient @ X  $dt = w_i T \in \mathbb{R}^{n \times m}$ SX Using chain rule, backpropagating to X we have  $\frac{dL}{dX} = \frac{dL}{dX} \frac{dL}{dt}$  $\frac{dL}{\partial X} = w_{1} \frac{dL}{dz_{1}} \in IR^{n}$ 

By the 3-D tensor derivative trick learned in class  $\frac{dL}{dw_{l}} = \frac{dL}{dt} \times^{T}$  $\frac{dL}{QW_{1}} = \frac{dL}{dZ_{1}} \times T \in \mathbb{R}^{m \times n}$