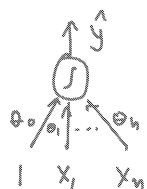


[illegible]

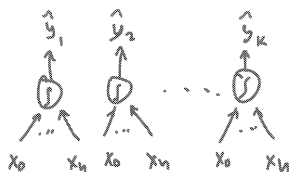
- [illegible]



$$\begin{aligned} \hat{y} &= f(\theta_1 x_1 + \dots + \theta_n x_n + \theta_0) \\ &= f(\theta^T x + \theta_0) \\ &= \begin{cases} 1 & \text{if } \theta^T x > \theta_0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

## This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 20 lines visible. The paper has a slight shadow on its right side, suggesting it's resting on a surface.

- [illegible]



$$\begin{matrix} \hat{y}_1 = f(\theta_1^T x + \theta_{10}) \\ \vdots \\ \hat{y}_k = f(\theta_k^T x + \theta_{k0}) \end{matrix}$$

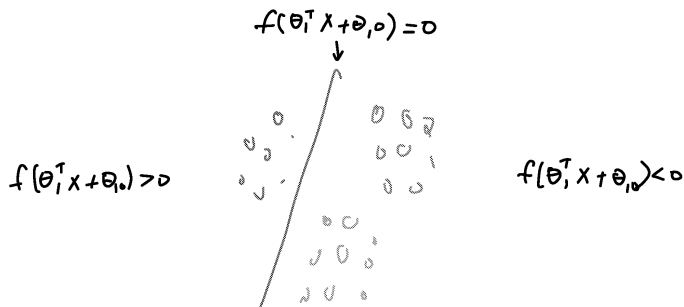
$$\hat{y} = f(x; \theta) = \theta^T x + \theta_0$$

[illegible]

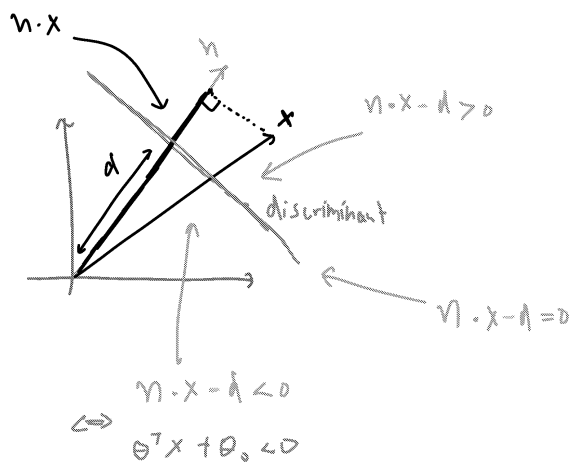
- \* Per form linear discrimination  
⇒ each unit defines a linear decision boundary
- \* Perform template matching  
⇒ each unit defines a template
- \* perform projection into a new basis  
⇒ each unit defines a basis vector

## Linear discrimination interpretation

- The rows of the weight matrix  $\theta^T$  represent parameters of K linear discriminant functions
- Each linear discriminant function separates one class from all others
- The discriminant function measures distance from a linear decision boundary. The distance is positive for examples belonging to the class and negative for all other examples

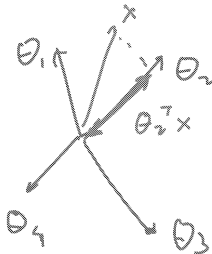


## Linear discriminants



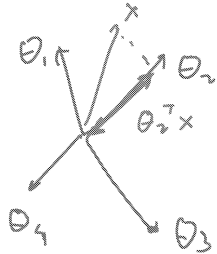
## Template matching interpretation

- The rows of the weight matrix  $\theta^T$  represent templates
- With K rows we have K templates (one template per class)
- The product  $\theta^T x$  measures how well The input vector matches the template (dot product between vectors)
- High similarity to the template of a particular class indicates membership in this class



## Projection into a new basis interpretation

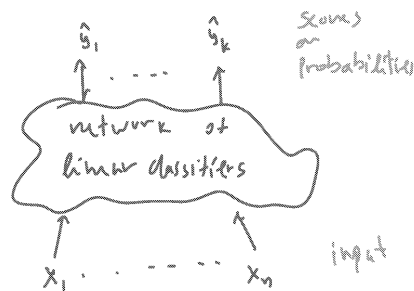
- The rows of the weight matrix  $\Theta^T$  represent basis vectors
- By taking the dot product of  $X$  with each of the rows we move to a representation in a new basis and so transform the data



## Linear classifiers

- what if decision boundaries are non-linear?
- what if more than one template is needed per class?

$\Rightarrow$  Combine multiple linear classifiers



## Activation functions

- Combining linear classifiers in sequence can be reduced to a single layer (linear combination of linear combinations is just a different linear combination)
- Therefore we need a nonlinear activation function between linear units

non-linear activation

$$g = h(\underbrace{\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n}_{\text{linear combination}})$$

A diagram of a single neuron. Inputs  $x_1, \dots, x_n$  are shown entering a circle containing a summation symbol  $\Sigma$ . Each input is multiplied by a weight  $\theta_1, \dots, \theta_n$ . A bias  $\theta_0$  is also added. The output of the summation is passed through a non-linear activation function  $h$  to produce the final output  $g$ .

### Activation functions

step:  $h(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$



sigmoid:  $h(a) = \frac{1}{1 + \exp(-a)}$



step<sup>x</sup>:  $h(a) = \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{otherwise} \end{cases}$



tanh:  $h(a) = \tanh(a) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

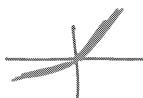


### Activation functions

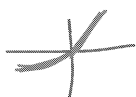
ReLU:  $h(a) = \max(0, a)$   
(Rectified Linear Units)



leaky ReLU:  $h(a) = \max(2a, a)$



ELU:  $h(a) = \begin{cases} a & a \geq 0 \\ a(e^a - 1) & a < 0 \end{cases}$



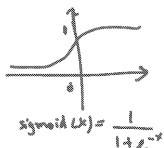
### Converting similarity to probability at the output layer

- Given similarity score (or decision boundary distance  $s_j$ ) from a binary classifier (higher better), compute:

$$\hat{y}_j^{(i)} \equiv P(y=j | X^{(i)})$$

- In a 2-class classification case:

$$\hat{y}_j^{(i)} \equiv P(y=j | X^{(i)}) = \text{Sigmoid}(s_j^{(i)})$$



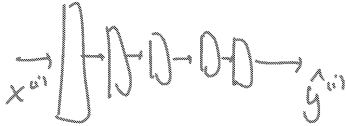
- In a k-class classification case

$$\hat{y}_j^{(i)} \equiv P(y=j | X^{(i)}) = \frac{\exp(s_j^{(i)})}{\sum_{l=1}^k \exp(s_l^{(i)})} \leftarrow \text{Softmax}$$

## Training the network

- Training: find weights for all network layers
- Loss: difference between obtained expected output
- Optimization: find the network parameters that will minimize the loss
- Select:
  - Loss function
  - Optimization framework

Given Training examples  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$



$$L(\theta) = \sum_{i=1}^n d(y^{(i)}, \hat{y}^{(i)}) \quad \theta^* = \arg\min_{\theta} L(\theta)$$

## Loss function

The weights  $\theta$  are set to minimize loss:

$$\begin{aligned} \theta^* &\equiv \arg\min_{\theta} L(\theta) \\ &\equiv \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^n L_i(x^{(i)}, y^{(i)}; \theta) \quad \text{sample loss} \\ &\equiv \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^n L_i(\underbrace{f(x^{(i)}; \theta)}_{\hat{y}^{(i)}}, y^{(i)}) \end{aligned}$$

## L1 and L2 loss

In regression problems minimize distance between known and predicted values:

$$L_1: L_i(\theta) = \sum_{j=1}^k |\hat{y}_j^{(i)} - y_j^{(i)}|$$

$$L_2: L_i(\theta) = \sum_{j=1}^k (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

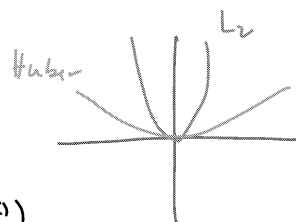
$\hat{y}_j^{(i)}$  is predicted value:  $\hat{y}_j^{(i)} \in \mathbb{R}$

### Huber loss

- Cap loss of outliers (robust regression)

$$\rho_{\delta}(d) = \begin{cases} \frac{1}{2} d^2 & \text{if } |d| \leq \delta \\ \delta \left( d - \frac{1}{2} \delta \right) & \text{otherwise} \end{cases}$$

{ quadratic for small  $d$   
linear for large  $d$



$$L_i(\theta) = \sum_{j=1}^k \rho_{\delta}(\hat{y}_j^{(i)} - y^{(i)})$$

### Cross entropy loss

- \* In classification problems assume that the true class label  $y^{(i)}$  is one-hot-encoded

$$x^{(i)} \in C_k \Rightarrow y^{(i)} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

- \* The outputs of the network are assumed to be the probabilities of belonging to each class

$$\hat{y}^{(i)} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\hat{y}_j^{(i)} = p(y=j | x^{(i)})$$

### Cross entropy loss

- Probability of the true class:  $\prod_{j=1}^k (\hat{y}_j^{(i)})^{y_j^{(i)}}$

$$\text{Likelihood: } L(\theta) = \prod_{i=1}^m \prod_{j=1}^k (\hat{y}_j^{(i)})^{y_j^{(i)}}$$

↑  
maximize

$$\text{negative Log-likelihood: } \ell(\theta) = -\log L(\theta) = -\sum_{i=1}^m \underbrace{\sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})}_{\text{sample loss}}$$

↑  
minimize (negative LL)

### Cross entropy loss

- possible cross-entropy loss values:

$$L_i(\theta) \in [0, \infty] \quad \begin{array}{ll} 0 & \text{when } p \rightarrow 1 \\ \infty & \text{when } p \rightarrow 0 \end{array}$$

- For random class assignment:

$$\hat{y}_j^{(n)} = P(y=j | x^{(n)}) = \frac{1}{K}$$

$$\ell(\theta) = - \sum_{i=1}^m \sum_{j=1}^K y_j^{(i)} \log\left(\frac{1}{K}\right) = - \sum_{i=1}^m -K = mK$$

↑  
bad loss value

### Regularization

- Occam's razor: simpler explanations are better

- A simpler solution is when the weights  $\theta$  are lower (e.g. when  $\theta_{ij} = 0$  we remove our coefficient).

- Smaller coefficients  $\rightarrow$  more stable solution that will generalize better

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m L_i(f(x_i; \theta), y_i) + \lambda R(\theta)$$

$$R(\theta) = \sum_i \sum_j \theta_{ij}^2$$

↑ regularization term  
weight of regularization (hyper parameter)

### Regularization

- Alternative Regularization ( $L_1$ ):

$$R(\theta) = \sum_i \sum_j |\theta_{ij}|$$

- $L_2$  Regularization minimizes weights while spreading them (e.g.  $0.5^2 + 0.5^2 < 1^2 + 0^2$ )

- $L_1$  Regularization does not have this property and may concentrate weights (e.g.  $|10.5| + |0.5| = |11| + |0|$ )

## Optimization

- To Minimize the loss (objective) solve:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta)$$

$$\Rightarrow \nabla L(\theta) = 0 \Rightarrow \theta^*$$

- If  $\nabla L(\theta)$  is not linear it is hard to find an explicit solution

$\Rightarrow$  numerical iterative solution  
(e.g. Gradient descent)

## Gradient descent

\* Iterative optimization algorithm

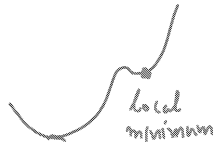
\* start with guess  $\theta_0$

\* Repeat:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \underset{\text{learning rate}}{\eta} \nabla L(\theta^{(i)})$$

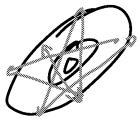
\* stop condition:

$$|L(\theta^{(i+1)}) - L(\theta^{(i)})| < \epsilon$$



## Gradient descent

\* The learning rate has to be selected correctly



$\eta$  too large



$\eta$  too small



good  $\eta$

\* Does not work on piecewise constant functions  
(e.g. empirical loss)





## Stochastic gradient descent (SGD)

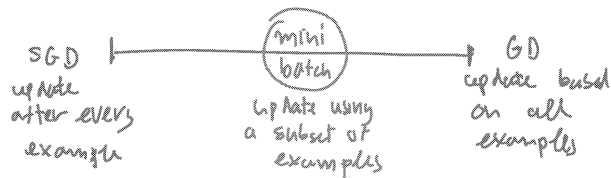
\* Randomly order examples

\* For  $j = 1..m$

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L_j(\theta^{(i)})$$

loss of j-th example

\* More frequent but less accurate updates



## Data normalization

- Features with high values have more influence and different features require a different learning rate.

(e.g. price ~ 10k vs age ~ 10)

- Data normalization:

$$\bar{x}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \bar{x}_j)^2$$

$$\tilde{x}_j^{(i)} = \frac{x_j^{(i)} - \bar{x}_j}{\hat{\sigma}_j}$$

normalized features

## Learning rate decay

- Learning rate need not be fixed

- Make learning smaller as iterations progress

\* start begins:

1) step decay:

$$\text{every } k \text{ iterations } \eta \leftarrow \eta/2$$

2) exponential decay

$$\eta = \eta_0 e^{-k/t}$$

decay rate

iteration index


3) fractional decay

$$\eta = \eta_0 / (1 + k/t)$$

## Additional optimization methods

- Newton's method:  
compute learning rate instead of specifying it  
⇒ second order derivatives (Hessian matrix)
- AdaGrad:  
replace Hessian computation (expensive) with approximation
- RMSProp:  
add decay to AdaGrad
- ADAM:  
add gradient momentum to RMSProp
- Quasi-Newton methods (e.g. BFGS)  
approximate Hessian inverse (faster)

## Weight initialization

- $\Theta = \text{constant}$  → all outputs identical  
→ no learning  
→ random initialization
- $\theta = \text{large}$  → saturated activations,  
→ zero gradients  
→ no learning  
→ initialize close to zero
- 

⇒ Initialize at random close to zero

## Weight initialization

\* Glorot (Xavier) normal initialization

- Draw weights from normal distribution with:

$$\text{Var}[w_i] = \left( \frac{2}{f_{\text{in}} + f_{\text{out}}} \right)^{1/2} \quad \mu = 0$$

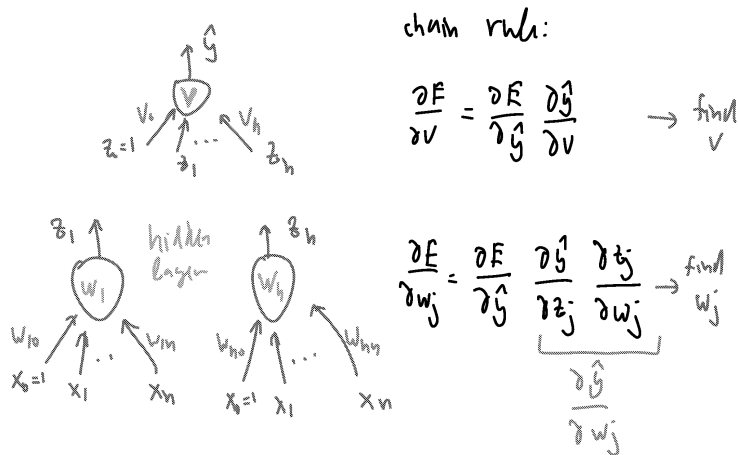
\* Glorot uniform initialization

- Draw weights from uniform distribution

within  $[-\text{lim}, \text{lim}]$

$$\text{lim} = \left( \frac{6}{f_{\text{in}} + f_{\text{out}}} \right)^{1/2}$$

## Computing gradients using Backpropagation

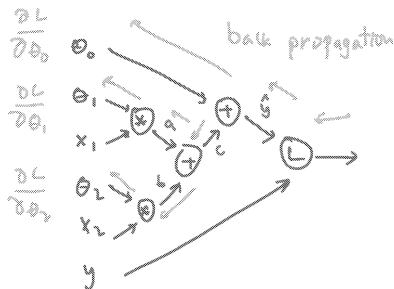


## Backpropagation

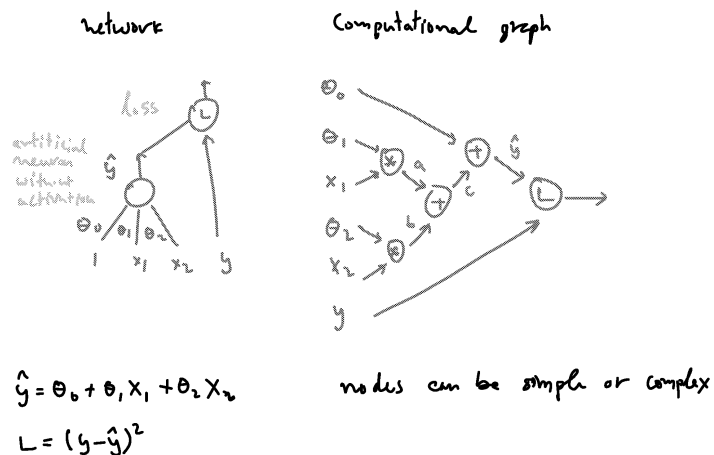
- Start with guess for parameters:  $v, \{w_j\}$  (at random close to zero)
- Use  $x^{(i)}$  and current parameters to compute outputs:  $z^{(i)}, \hat{y}^{(i)}$  forward pass
- Use outputs  $z^{(i)}, \hat{y}^{(i)}$  to update parameters:  $v, \{w_j\}$  backward pass (push back gradients)
- continue while loss changes

## Backpropagation

- Represent the network using a computational graph
- Because each node is simple, it has a simple explicit expression for its derivative
- Forward pass: push input to compute all intermediate node values
- Backward pass: starting with end nodes, push gradients towards the beginning
- Multiply backpropagated gradients (front back) by current gradient and propagate this to the next node



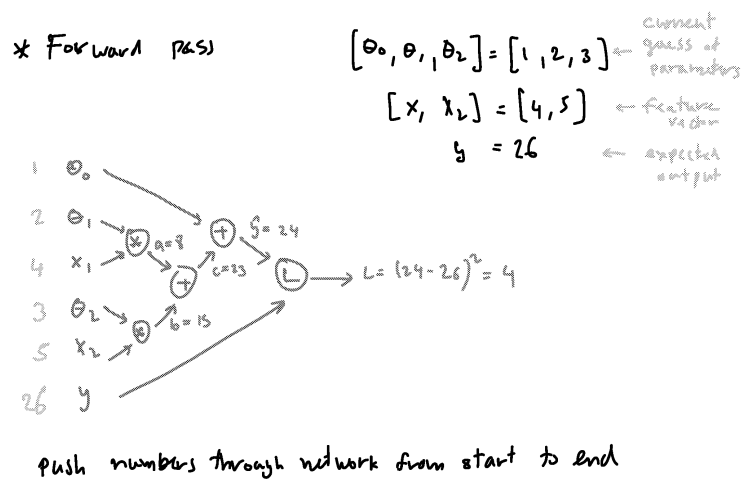
# Backpropagation example



$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

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

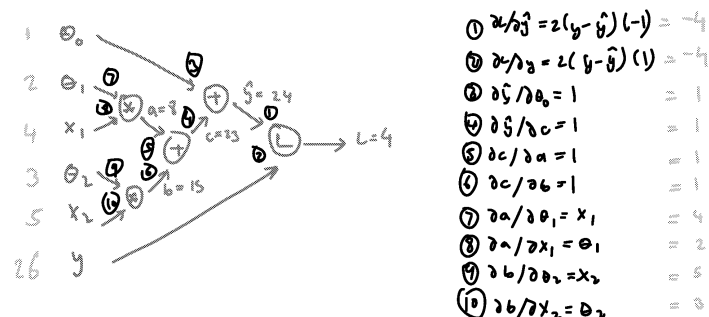
# Backpropagation example



# Backpropagation example

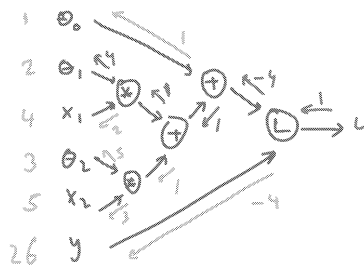
## \* Backward pass (compute gradients)

1. For each node: write derivative of output w.r.t. each input
2. Compute gradient at input of each node (using forward values)
3. Multiply gradients along paths



## Backpropagation example

\* Backward pass (contd.)



$$\frac{\partial L}{\partial \theta_0} = 1 \cdot (-4) \cdot 1 = -4$$

$$\frac{\partial L}{\partial \theta_1} = 1 \cdot (-4) \cdot 1 \cdot 1 \cdot 4 = -16$$

$$\frac{\partial L}{\partial \theta_2} = 1 \cdot (-4) \cdot 1 \cdot 1 \cdot 8 = -32$$

$$\textcircled{1} x_{\text{adj}} = 2(y - \hat{y})(-1) = -4$$

$$\textcircled{v} \partial_x \lambda_2 = 2(y - \hat{y})(1) = -4$$

$$\textcircled{2} \partial \hat{\gamma} / \partial \theta_0 = 1$$

$$\textcircled{4} \partial g / \partial c = 1$$

⑤  $\partial c / \partial a = 1$

$$(6) \partial c / \partial b = 1$$

$$\textcircled{7} \partial a / \partial \theta_1 = x_1 \quad = 4$$

$$\textcircled{8} \partial \alpha / \partial x_1 = 0, \quad \text{---} \quad 2$$

(10)  $\lambda_1 \lambda_2 = 2$

$$\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2}, \frac{\partial L}{\partial y} \quad \text{can}$$

be augmented but  
not needed

[illegible]