

# Neural Networks: Backpropagation & Regularization

Benjamin Roth, Nina Poerner

CIS LMU München

# Outline

1 Backpropagation

2 Regularization

# Backpropagation

- Forward propagation: Input information  $x$  propagates through network to produce output  $\hat{y}$ .
- Calculate cost  $J(\theta)$ , as you would with regression.
- Compute gradients w.r.t. all model parameters  $\theta$ ...
- ... how?
  - ▶ We know how to compute gradients w.r.t. parameters of the output layer (just like regression).
  - ▶ How to calculate them w.r.t. parameters of the hidden layers?

# Chain Rule of Calculus

- Let  $x, y, z \in \mathbb{R}$ .
- Let functions  $f, g : \mathbb{R} \rightarrow \mathbb{R}$ .
- $y = g(x)$
- $z = f(g(x))$
- Then

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

# Chain Rule of Calculus: Vector-valued Functions

- Let  $\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n, z \in \mathbb{R}$
- Let functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$
- $\mathbf{y} = g(\mathbf{x})$
- $z = f(g(\mathbf{x})) = f(\mathbf{y})$
- Then

$$\frac{\partial z}{\partial x_i} = \sum_{j=1}^n \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

- In order to write this in vector notation, we need to define the Jacobian matrix.

# Jacobian

- The Jacobian is the matrix of all first-order partial derivatives of a vector-valued function.

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_m} \\ \frac{\partial y_2}{\partial x_1} & & \frac{\partial y_2}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

- How to write in terms of gradients?
- We can write the chain rule as:

$$\nabla_{\mathbf{x}} z = \begin{pmatrix} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \end{pmatrix}^T \nabla_{\mathbf{y}} z$$

$m \times 1$                        $n \times m$                        $n \times 1$

## Viewing the Network as a Graph

- Nodes are function outputs (can be scalar or vector valued)
- Arrows are functions
- Example:
  - $\hat{y} = \mathbf{v}^T \text{relu}(\mathbf{W}^T \mathbf{x})$
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}; \mathbf{r} = \text{relu}(\mathbf{z})$

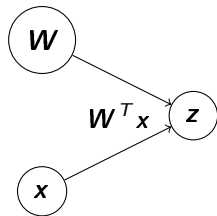
$\mathbf{W}$

$\mathbf{x}$

$\mathbf{v}$

## Viewing the Network as a Graph

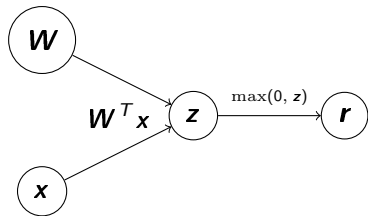
- Nodes are function outputs (can be scalar or vector valued)
- Arrows are functions
- Example:
  - $\hat{y} = \mathbf{v}^T \text{relu}(\mathbf{W}^T \mathbf{x})$
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}; \mathbf{r} = \text{relu}(\mathbf{z})$





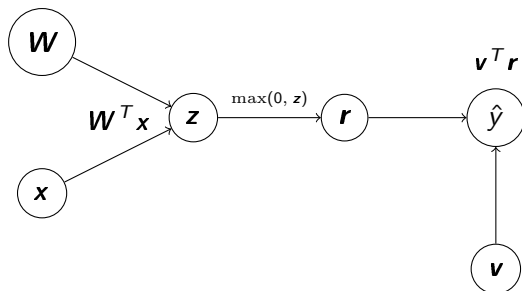
## Viewing the Network as a Graph

- Nodes are function outputs (can be scalar or vector valued)
- Arrows are functions
- Example:
  - $\hat{y} = \mathbf{v}^T \text{relu}(\mathbf{W}^T \mathbf{x})$
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}; \mathbf{r} = \text{relu}(\mathbf{z})$



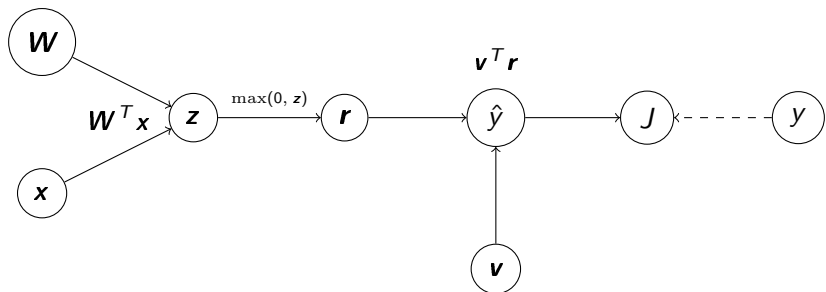
## Viewing the Network as a Graph

- Nodes are function outputs (can be scalar or vector valued)
- Arrows are functions
- Example:
  - $\hat{y} = \mathbf{v}^T \text{relu}(\mathbf{W}^T \mathbf{x})$
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}; \mathbf{r} = \text{relu}(\mathbf{z})$



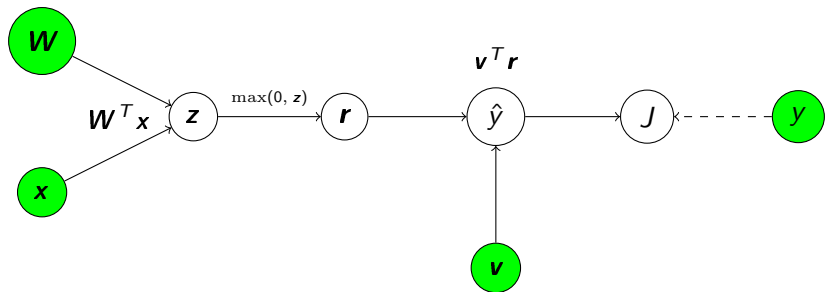
## Viewing the Network as a Graph

- Nodes are function outputs (can be scalar or vector valued)
- Arrows are functions
- Example:
  - $\hat{y} = \mathbf{v}^T \text{relu}(\mathbf{W}^T \mathbf{x})$
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}; \mathbf{r} = \text{relu}(\mathbf{z})$



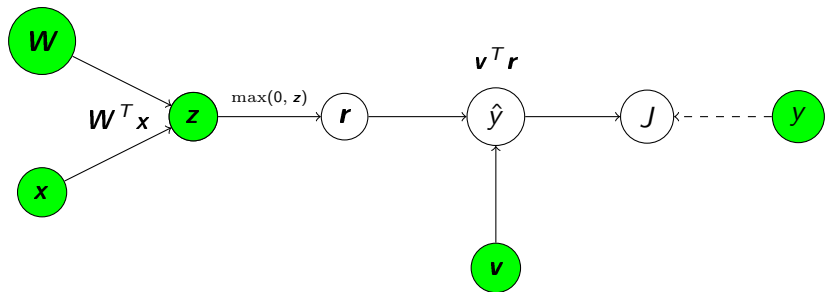
# Forward pass

Green: Known or computed node



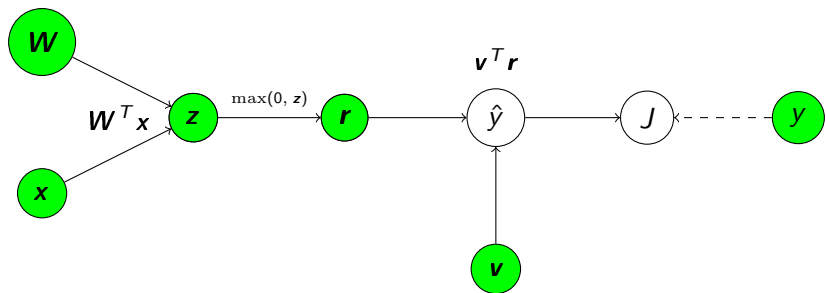
# Forward pass

Green: Known or computed node



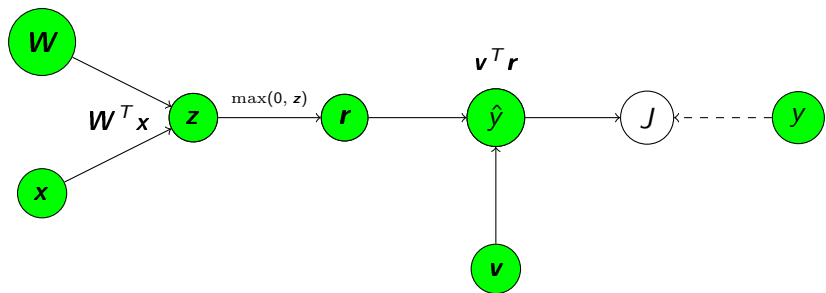
# Forward pass

Green: Known or computed node



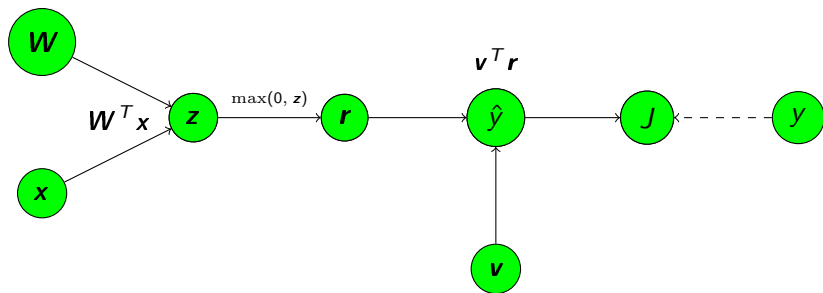
# Forward pass

Green: Known or computed node



# Forward pass

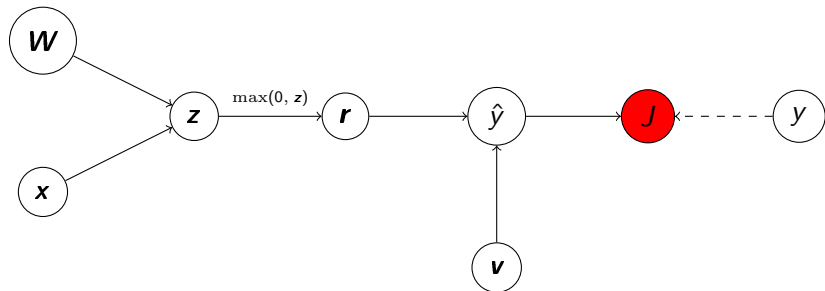
Green: Known or computed node





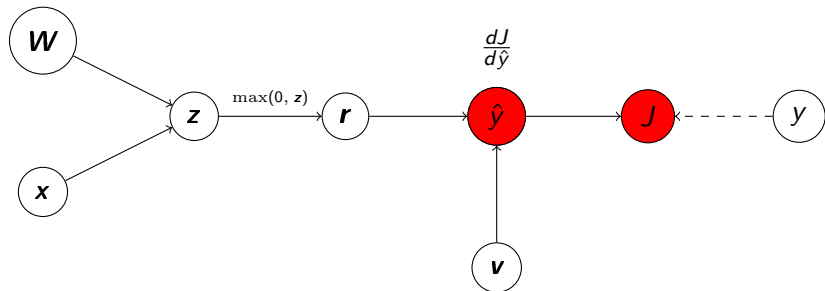
## Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed



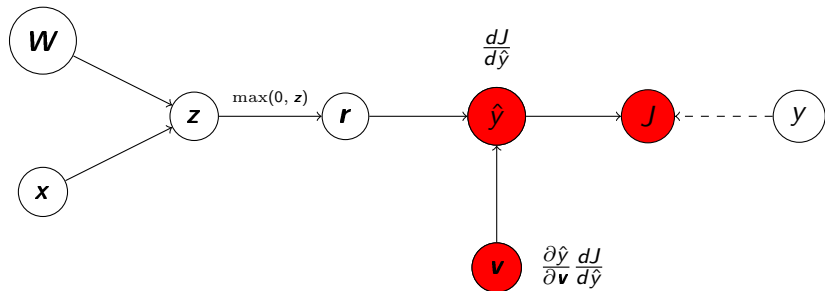
## Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed



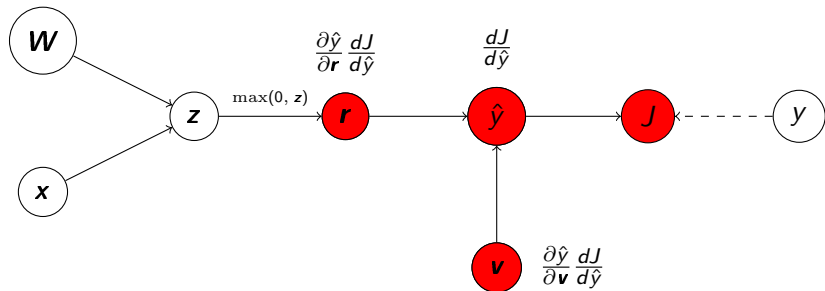
# Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed



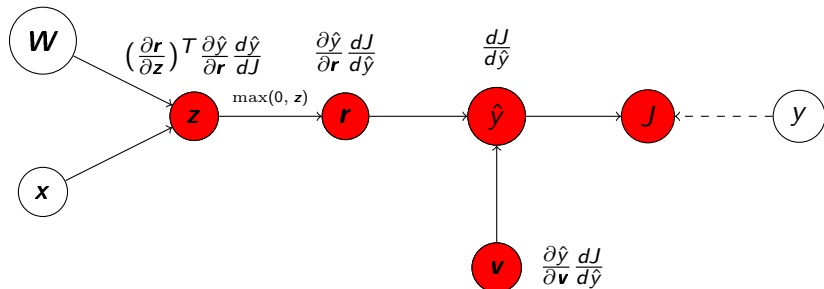
# Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed



# Backward pass

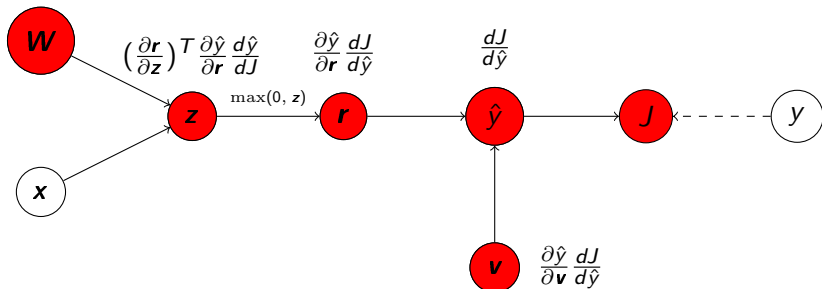
Red: Gradient of  $J$  w.r.t. node known or computed



# Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed

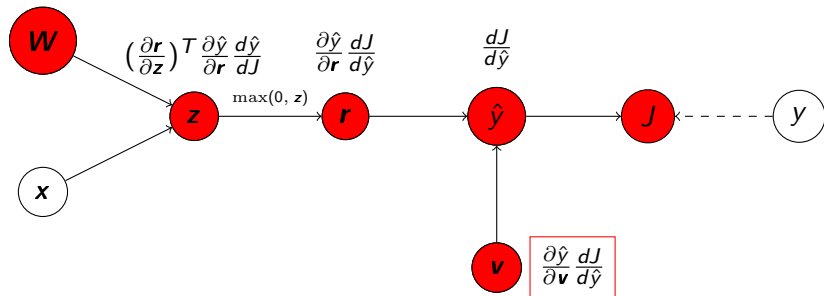
$$\left(\frac{\partial \mathbf{z}}{\partial \mathbf{W}}\right)^T \left(\frac{\partial \mathbf{r}}{\partial \mathbf{z}}\right)^T \frac{\partial \hat{y}}{\partial \mathbf{r}} \frac{dJ}{d\hat{y}}$$



# Backward pass

Red: Gradient of  $J$  w.r.t. node known or computed

$$\left(\frac{\partial z}{\partial W}\right)^T \left(\frac{\partial r}{\partial z}\right)^T \frac{\partial \hat{y}}{\partial r} \frac{dJ}{d\hat{y}}$$



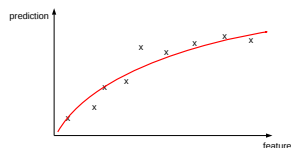
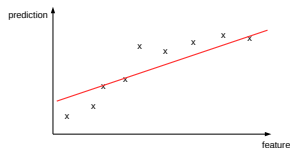
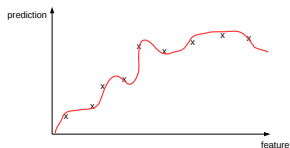
# Outline

1 Backpropagation

2 Regularization

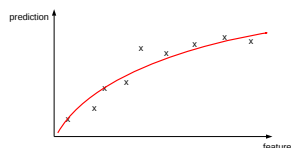
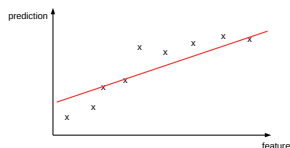
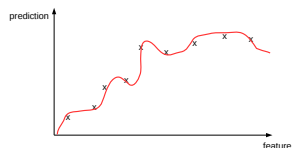


# Regularization



- Overfitting vs. underfitting

# Regularization



- Overfitting vs. underfitting
- Regularization: Any modification to a learning algorithm for reducing its generalization error but not its training error
- Build a “preference” into model for some solutions in hypothesis space
- Unpreferred solutions are penalized: only chosen if they fit training data much better than preferred solutions

# Regularization

- Large parameters  $\rightarrow$  overfitting
- Prefer models with smaller weights
- Popular regularizers:
  - ▶ Penalize large L2 norm (= Euclidian norm) of weight vectors
  - ▶ Penalize large L1 norm (= Manhattan norm) of weight vectors

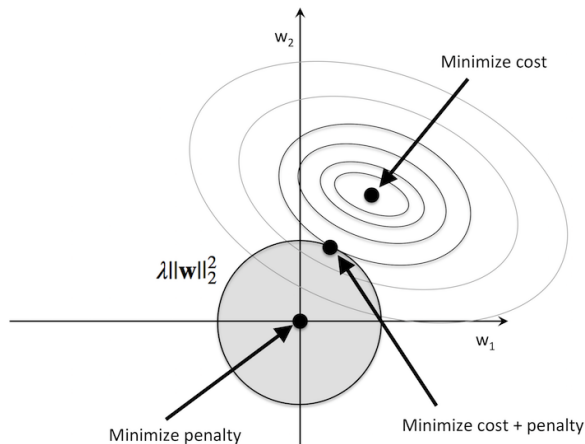
# L2-Regularization

- Add term that penalizes large L2 norm of weight vector  $\theta$
- The amount of penalty is controlled by a parameter  $\lambda$

$$J'(\theta) = J(\theta, \mathbf{x}, \mathbf{y}) + \frac{\lambda}{2} \theta^T \theta$$

## L2-Regularization

- The surface of the objective function is now a combination of the original loss and the regularization penalty.



# Summary

- Feedforward networks: layers of (non-linear) function compositions

# Summary

- Feedforward networks: layers of (non-linear) function compositions
- Non-Linearities for hidden layers: `relu`, `tanh`, ...

# Summary

- Feedforward networks: layers of (non-linear) function compositions
- Non-Linearities for hidden layers: `relu`, `tanh`, ...
- Non-Linearities for output units (classification):  $\sigma$ , `softmax`



# Summary

- Feedforward networks: layers of (non-linear) function compositions
- Non-Linearities for hidden layers:  $\text{relu}$ ,  $\text{tanh}$ , ...
- Non-Linearities for output units (classification):  $\sigma$ ,  $\text{softmax}$
- Training via backpropagation: compute gradient of cost w.r.t. parameters using chain rule

# Summary

- Feedforward networks: layers of (non-linear) function compositions
- Non-Linearities for hidden layers:  $\text{relu}$ ,  $\text{tanh}$ , ...
- Non-Linearities for output units (classification):  $\sigma$ ,  $\text{softmax}$
- Training via backpropagation: compute gradient of cost w.r.t. parameters using chain rule
- Regularization: penalize large parameter values, e.g. by adding L2-norm of parameter vector to loss

# Outlook

- “Manually” defining forward- and backward passes in numpy is time-consuming
- Deep Learning frameworks let you define forward pass as a “computation graph” made up of simple, differentiable operations (e.g., dot products).
- They do the backward pass for you
- tensorflow + keras, pytorch, theano, MXNet, CNTK, caffe, ...