

The Ancient Secrets



Computer Vision

Logistics:

- Homework 3 is out!
 - Optical flow
 - Won't work on newest OpenCV 3.4
 - Should work on OpenCV 2, maybe early versions of 3
 - demo!

Previously
On



Ancient Secrets
of Computer Vision

Softmax: normalized exponential

Generalization of logistic

Input: vector of reals

Output: probability distribution

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

softmax([1,2,7,3,2]):

Calculate e^x : [2.72, 7.39, 1096.63, 20.09, 7.39]

Calculate $\text{sum}(e^x)$: $2.72+7.39+1096.63+20.09+7.39 = 1134.22$

Normalize: $e^x/\text{sum}(e^x) = [0.002, 0.007, 0.967, 0.017, 0.007]$

Multinomial logistic regression

Probability of that a data point belongs to a class is the normalized, weighted sum of the input variables with the learned weights.

$\text{softmax}(\mathbf{w}\mathbf{x} + \mathbf{b})$

$$P(y = j \mid \mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^\top \mathbf{w}_k}}$$

Multinomial logistic regression

Probability of that a data point belongs to a class is the normalized, weighted sum of the input variables with the learned weights.

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \text{softmax} \left(\begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \right)$$

MNIST: Handwriting recognition

50,000 images of handwriting

28 x 28 x 1 (grayscale)

Numbers 0-9

10 class softmax regression

Input is 784 pixel values

Train with SGD

> 95% accuracy

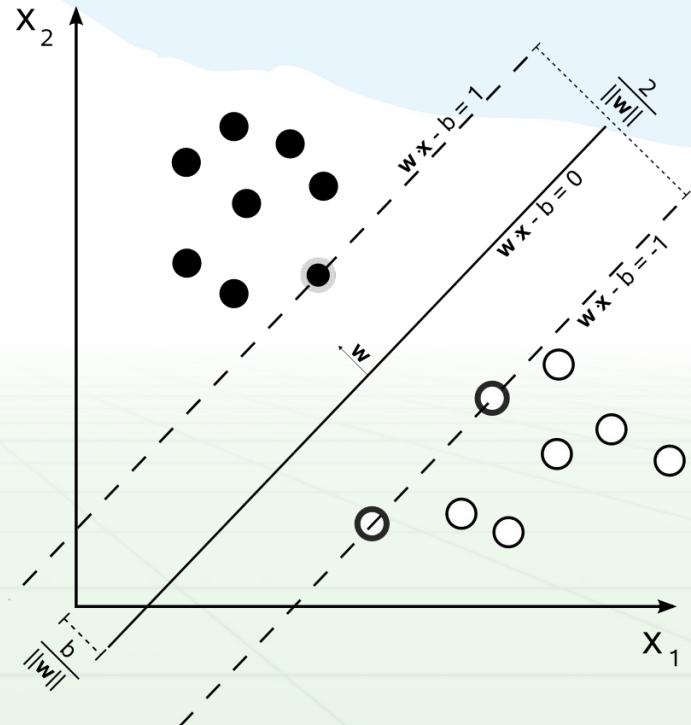


Support Vector Machine (SVM)

Find max-margin classifier. Examples on the margin are supporting data points, support vectors.

$$\begin{aligned} & \min ||\mathbf{w}||_2 \\ \text{s.t. } & y_n(\mathbf{w} \cdot \mathbf{x}_n - b) \geq 1, \quad n = 1, 2 \dots \end{aligned}$$

Or: minimize weights such that margin for each point is at least 1



Case study: Person detection

Dalal and Triggs '05:

Train SVM on HOG features of image

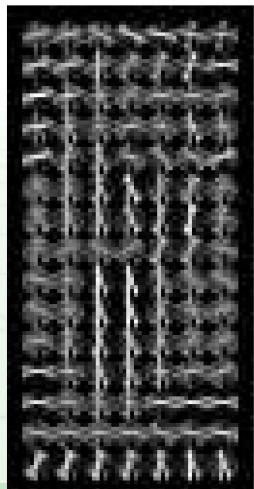
2 classes, person/not person

At test time:

Extract HOG features at many scales

Run SVM classifier at every location

High responses = person?



Case study: Person detection

Dalal and Triggs is a sliding window detector

Many scales

Every location

10k+ classifier
evaluations per
image.

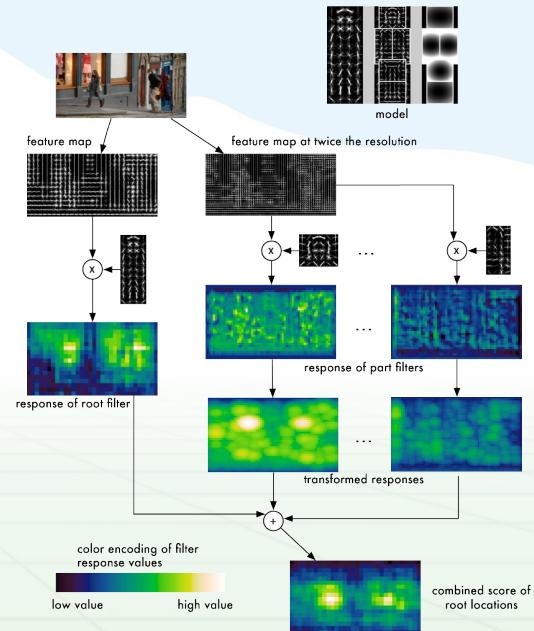


Case study: Deformable parts models

Objects have parts, learn to recognize parts and where they are

Latent SVM: Learn part appearances and locations without explicit data

Hard negative mining: rebalance classes for sliding window detectors



Case study: Image classification

Given an image, what's in it?

Old state-of-the-art:

Extract features from image

SIFT and Fisher Vectors

Train Linear SVM

On 1000 different classes, 54% accurate

What's wrong with this?

Machine learning needs features!!

What are the right features?

HOG?

SIFT?

FV?

Why not let the algorithm decide

Neural networks: Feature extraction + linear model

Success of Neural Networks

Image classification:

54% -> 80% accuracy on 1000 classes

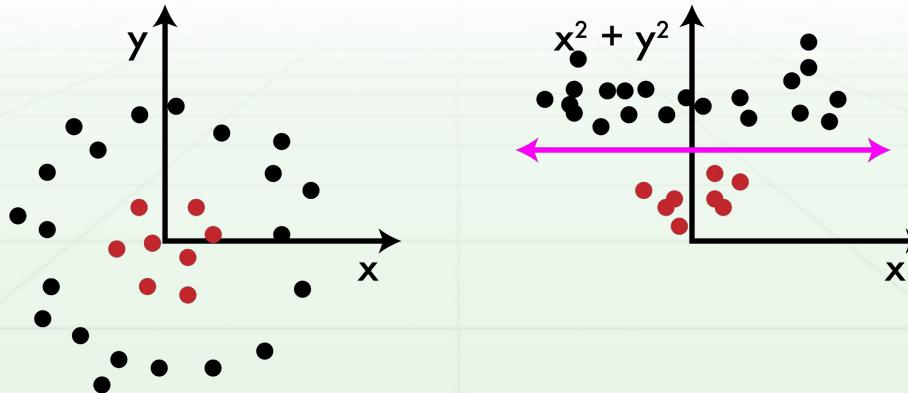
Object detection:

33% mAP (DPM) -> 88% mAP on 20 classes

Feature engineering

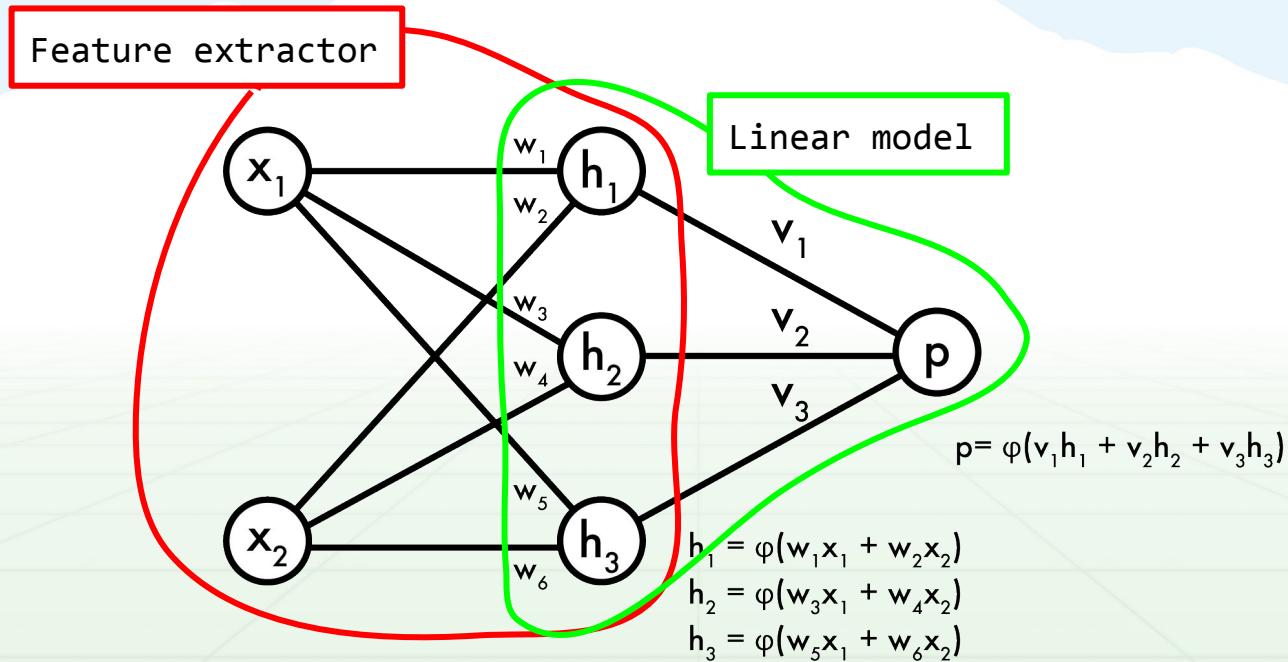
Arguably the **core** problem of machine learning
(especially in practice)

ML models work well if there is a clear relationship between the inputs and outputs of the function you are trying to model



What if we added more processing?

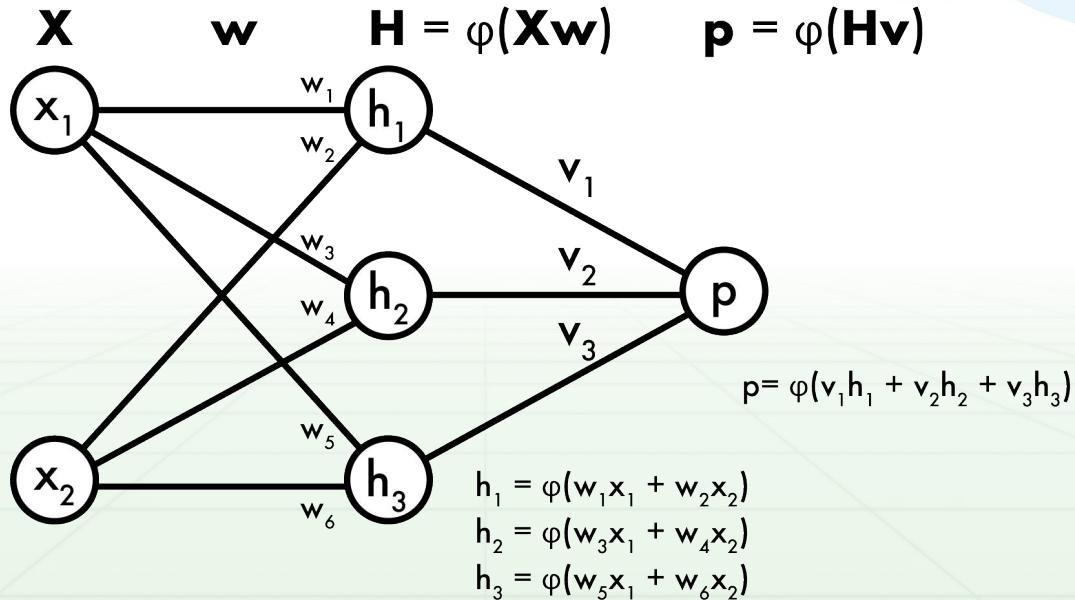
Now our prediction p is a function of our hidden layer



This is a neural network!

This one has 1 hidden layer, but can have **way** more

Each layer is just some function φ applied to linear combination of the previous layer



Universal approximation theorem

What if ϕ not linear?

Universal approximation theorem (Cybenko 89, Hornik 91)

ϕ : any nonconstant, bounded, monotonically increasing function

I_m : m -dimensional unit hypercube (interval [0-1] in m -d)

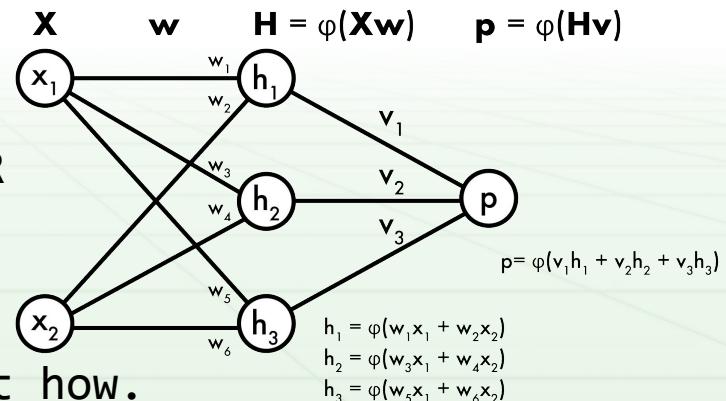
Then 1-layer neural network with ϕ as activation can model any continuous function $f: I_m \rightarrow R$
(no bound on size of hidden layer)

By extension, works on $f: \text{bounded } R^m \rightarrow R$

What can we learn? What can't we?

UAT just says it's possible to model, not how.

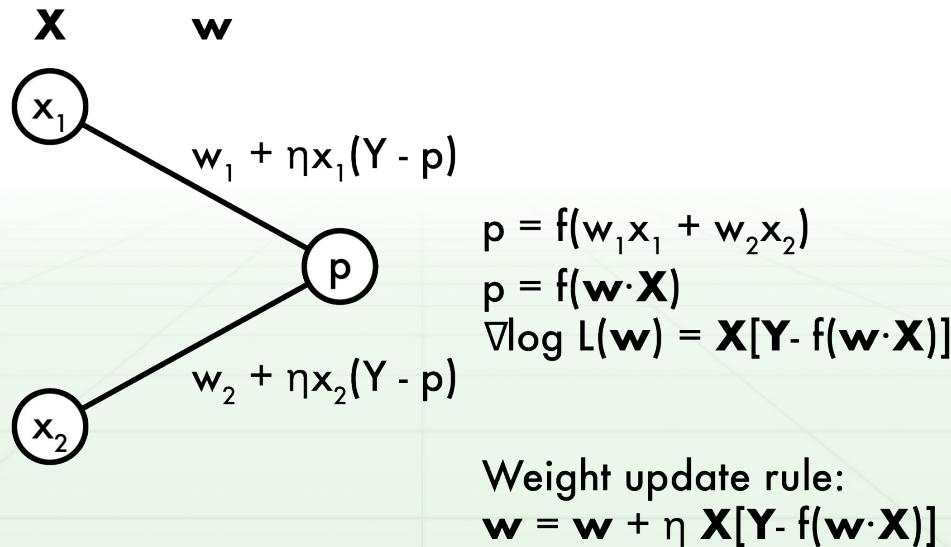
https://en.wikipedia.org/wiki/Universal_approximation_theorem



How do we learn it?

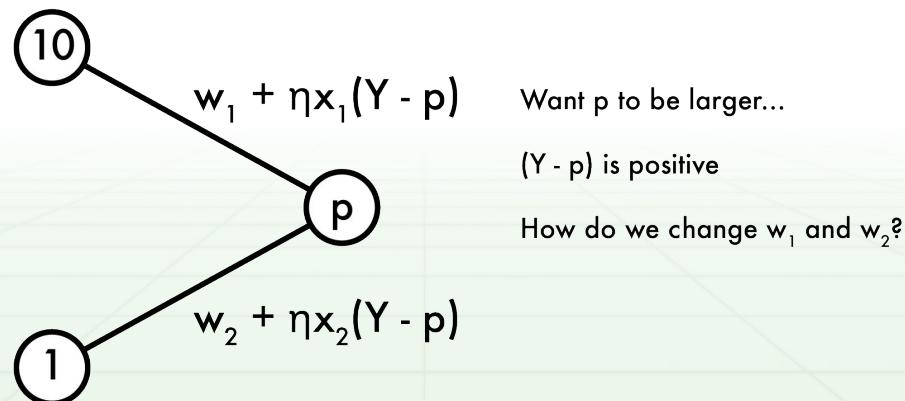
With gradient descent we calculate the partial derivatives of the loss (or likelihood) function for every weight: $\partial/\partial w_i \log L(w)$

Then do gradient descent (or ascent) by adding gradient to weight



How do we learn it?

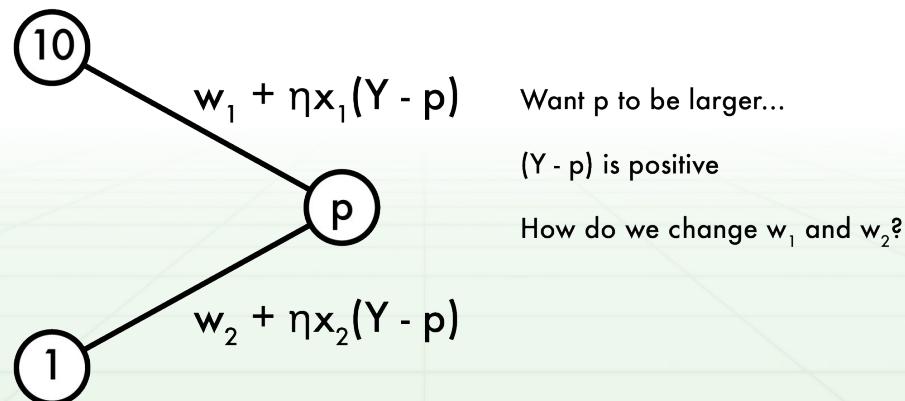
Simple example, say we have a data point $[10, 1]$ and we predict some p . We also know the “correct” label Y . Maybe our prediction p is too small and we want to make it larger. How do we adjust w ?



How do we learn it?

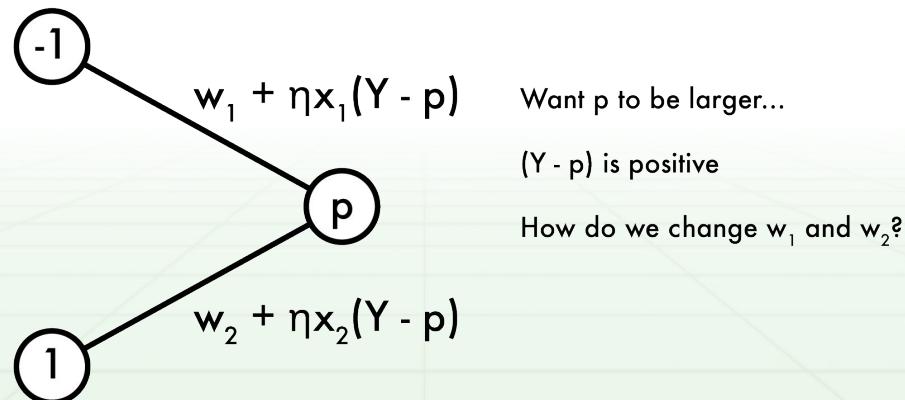
Simple example, say we have a data point $[10, 1]$ and we predict some p . We also know the “correct” label Y . Maybe our prediction p is too small and we want to make it larger. How do we adjust w ?

We adjust w_1 much more than w_2 , why?



How do we learn it?

Simple example, say we have a data point $[-1, 1]$ and we predict some p . We also know the “correct” label Y . Maybe our prediction p is too small and we want to make it larger. How do we adjust w ?



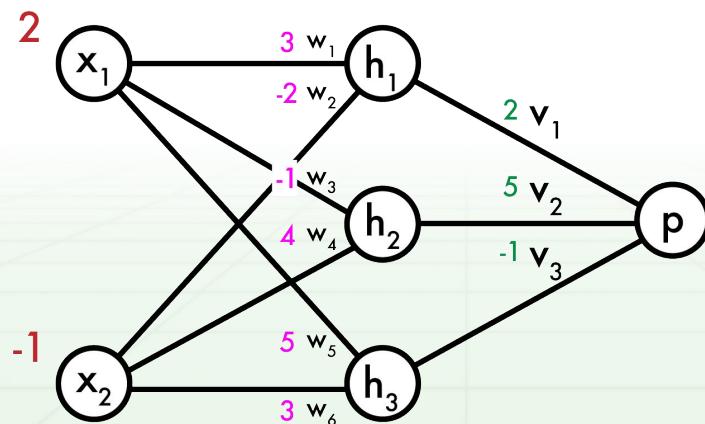
Chapter Twelve



Neural Networks

How do we learn it?

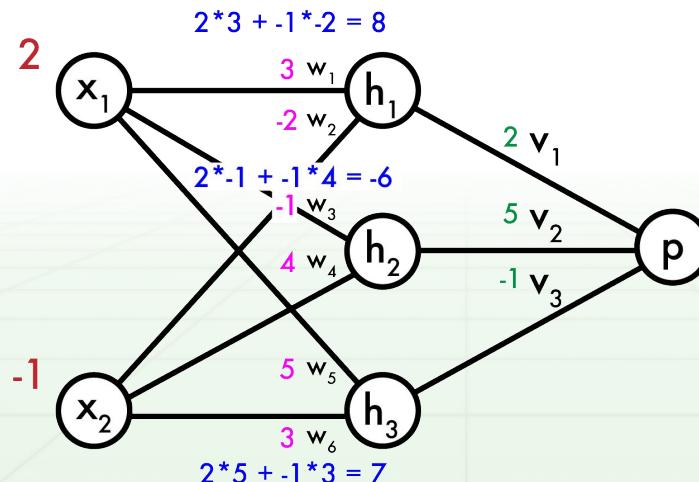
Now we have a “real” neural network (using linear activation for simplicity). How do we predict p ?



How do we learn it?

Now we have a “real” neural network (using linear activation for simplicity). How do we predict p?

Calculate hidden layer neurons

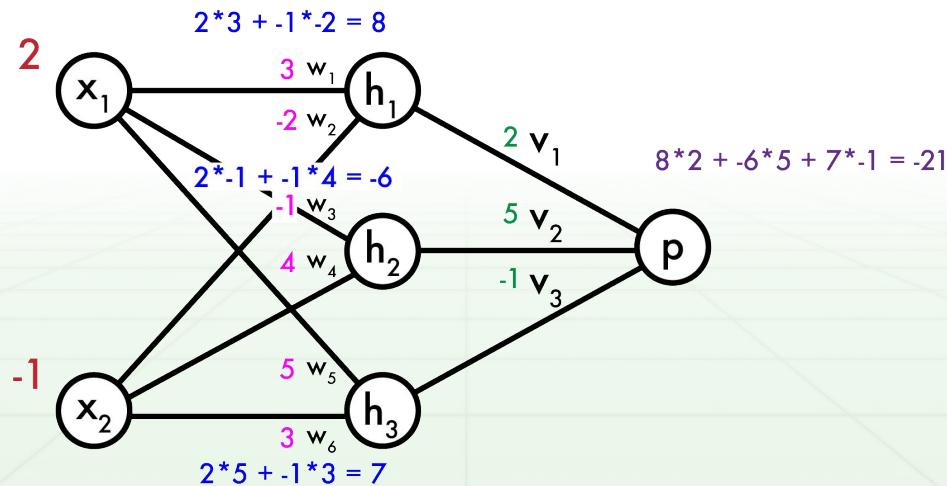


How do we learn it?

Now we have a “real” neural network (using linear activation for simplicity). How do we predict p?

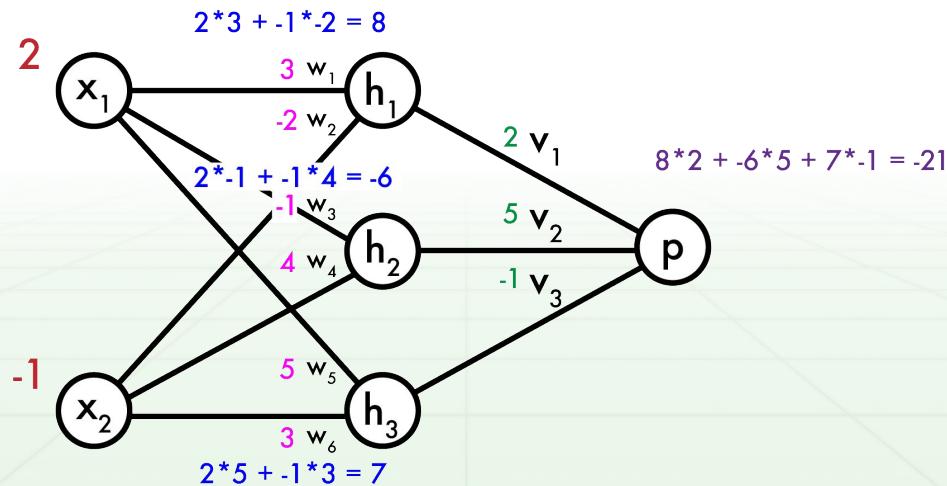
Calculate hidden layer neurons

Calculate output p



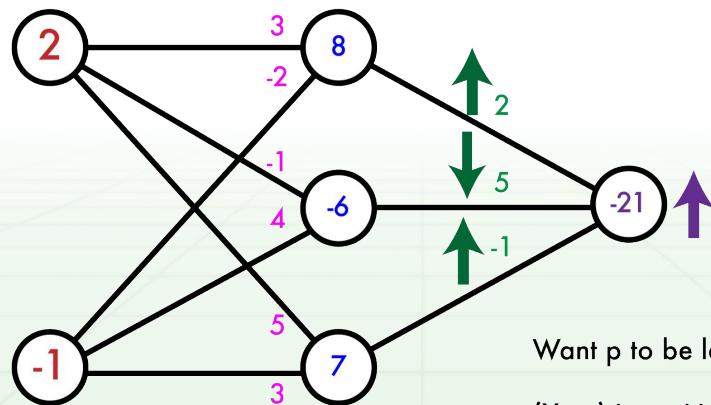
How do we learn it?

Say we want to make p larger. How do we modify the weights? The first layer is easy, same as normal linear model:



How do we learn it?

Say we want to make p larger. How do we modify the weights? The first layer is easy, same as normal linear model:



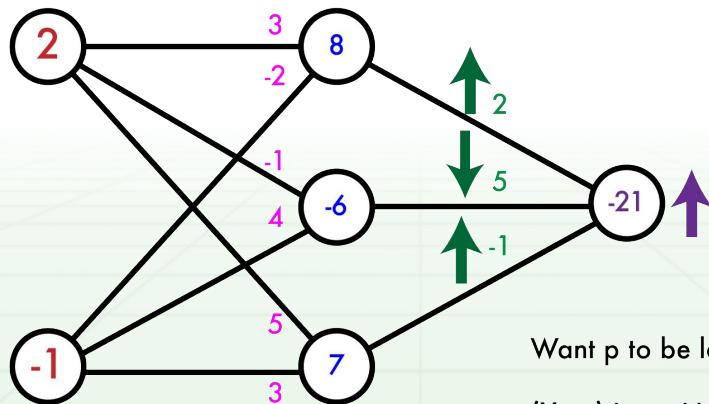
Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?

How do we learn it?

Now what? Let's calculate the “error” that the hidden layer makes. We want p to be larger, given current weights how should we adjust the hidden layer output to do that?



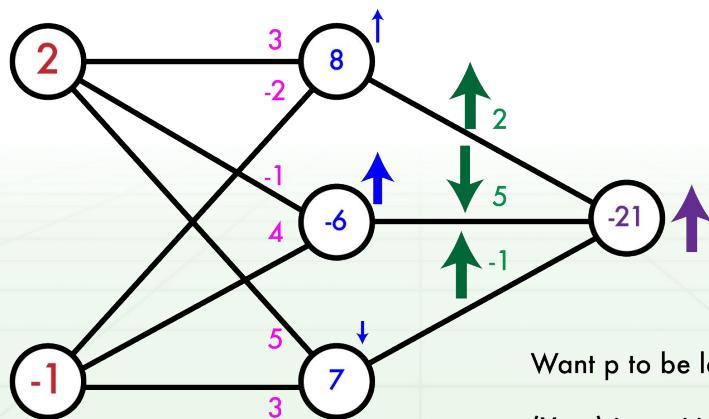
Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?

How do we learn it?

Now what? Let's calculate the “error” that the hidden layer makes. We want p to be larger, given current weights how should we adjust the hidden layer output to do that?



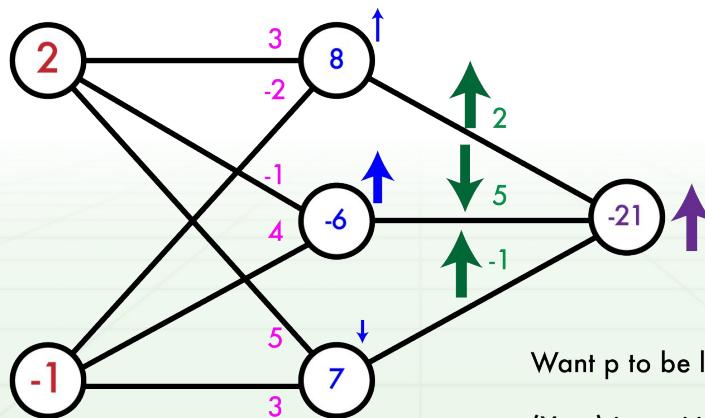
Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?

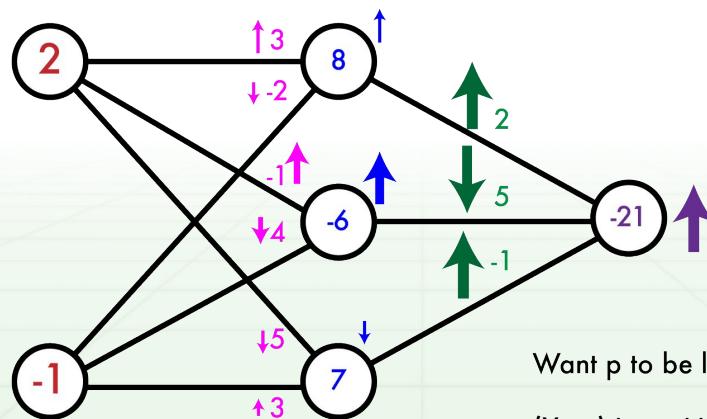
How do we learn it?

Now that we have an “error” in our hidden layer, want to modify the previous weights. Easy again, just like our linear model.



How do we learn it?

Now that we have an “error” in our hidden layer, want to modify the previous weights. Easy again, just like our linear model.



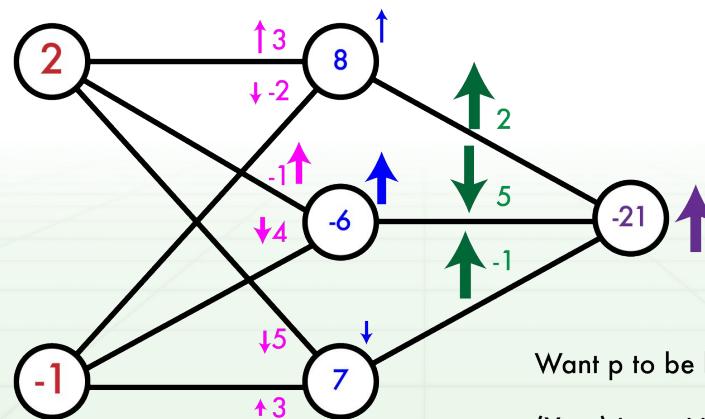
Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?

Backpropagation: just taking derivatives

This is the backpropagation algorithm. It's really just an easy way to calculate partial derivatives in a neural network. We forward-propagate information through the network, calculate our error, then backpropagate that error through network to calculate weight updates.



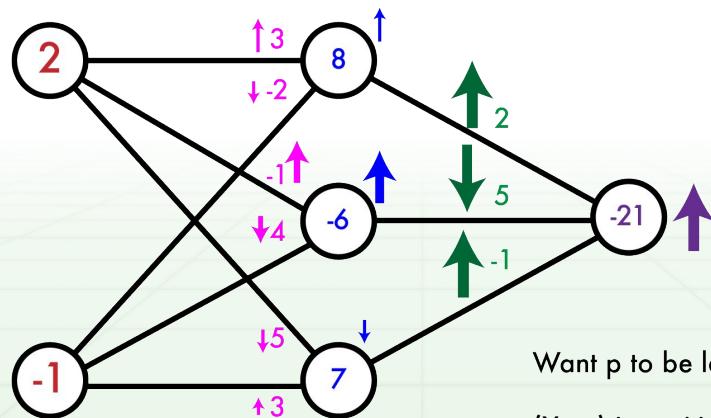
Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?

Backpropagation: just taking derivatives

This was with linear activations but the process is the same for any φ , just have to calculate $\varphi'(x)$ for that neuron as well.



Want p to be larger...

$(Y - p)$ is positive

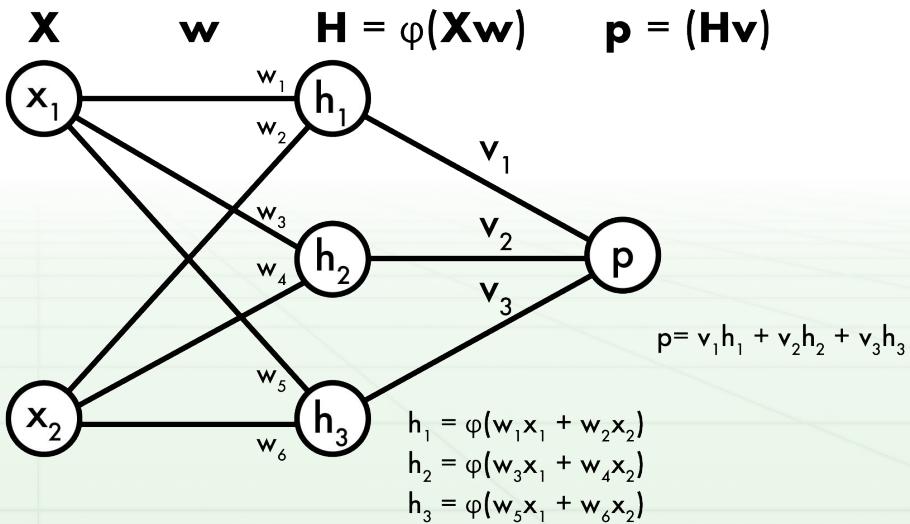
How do we change our weights?

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \varphi(Xw)v$$

$$F(X) = \varphi(x_1 w_1 + x_2 w_2) v_1 + \varphi(x_1 w_3 + x_2 w_4) v_2 + \varphi(x_1 w_5 + x_2 w_6) v_3$$



Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \varphi(Xw)v$$

$$F(X) = \varphi(x_1 w_1 + x_2 w_2)v_1 + \varphi(x_1 w_3 + x_2 w_4)v_2 + \varphi(x_1 w_5 + x_2 w_6)v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \varphi(Xw)v)^2$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2)v_1 + \phi(x_1 w_3 + x_2 w_4)v_2 + \phi(x_1 w_5 + x_2 w_6)v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\partial/\partial v_2 L_{x,y}(w, v) = \partial/\partial v_2 \frac{1}{2}(Y - \phi(Xw)v)^2$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2)v_1 + \phi(x_1 w_3 + x_2 w_4)v_2 + \phi(x_1 w_5 + x_2 w_6)v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\begin{aligned}\partial/\partial v_2 L_{x,y}(w, v) &= \partial/\partial v_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial v_2 \phi(Xw)v]\end{aligned}$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2)v_1 + \phi(x_1 w_3 + x_2 w_4)v_2 + \phi(x_1 w_5 + x_2 w_6)v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\begin{aligned}\partial/\partial v_2 L_{x,y}(w, v) &= \partial/\partial v_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial v_2 \phi(Xw)v] \\ &= (Y - \phi(Xw)v) * -[\partial/\partial v_2 \phi(x_1 w_3 + x_2 w_4)v_2]\end{aligned}$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2)v_1 + \phi(x_1 w_3 + x_2 w_4)v_2 + \phi(x_1 w_5 + x_2 w_6)v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

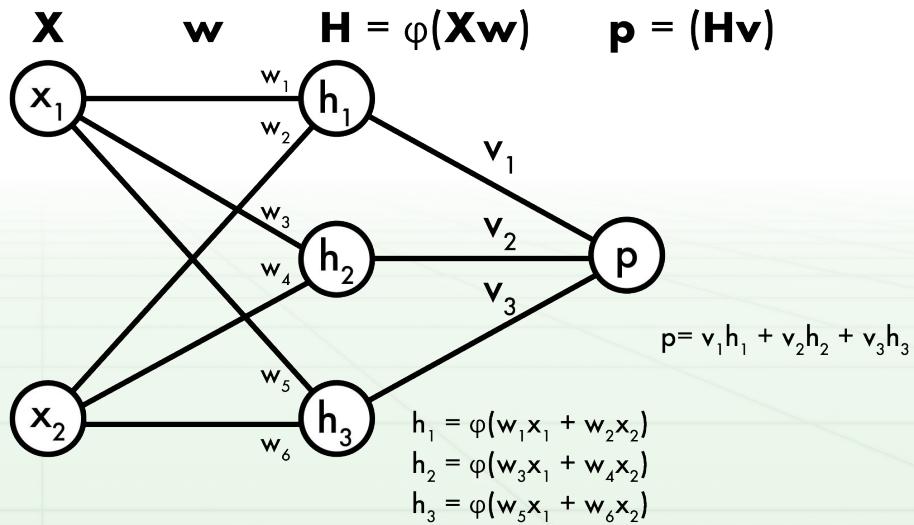
Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\begin{aligned}\partial/\partial v_2 L_{x,y}(w, v) &= \partial/\partial v_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial v_2 \phi(Xw)v] \\ &= (Y - \phi(Xw)v) * -[\partial/\partial v_2 \phi(x_1 w_3 + x_2 w_4)v_2] \\ &= (Y - \phi(Xw)v) * -\phi(x_1 w_3 + x_2 w_4)\end{aligned}$$

Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say v_2 :

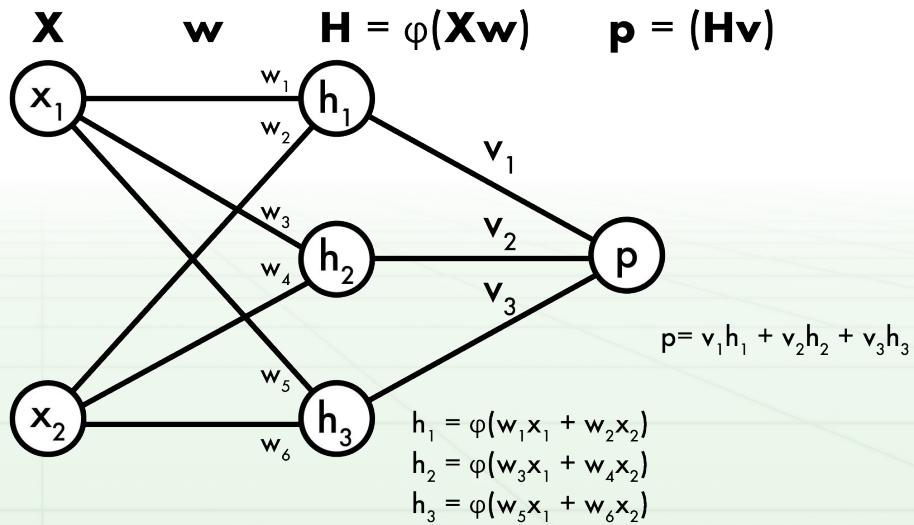
$$\frac{\partial}{\partial v_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -\varphi(x_1w_3 + x_2w_4)$$



Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\frac{\partial}{\partial v_2} L_{x,y}(w, v) = -(Y - \varphi(\mathbf{X}w)v) * h_2$$



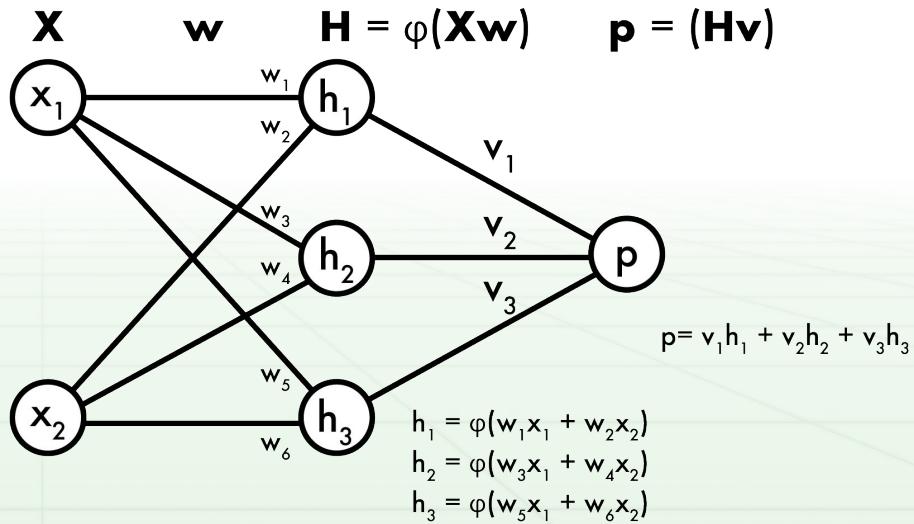
Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say v_2 :

$$\frac{\partial}{\partial v_2} L_{x,y}(w, v) = -(Y - \varphi(\mathbf{X}w)v) * h_2$$

Weight update rule (remember descend on loss):

$$v_2 = v_2 + \eta(Y - \varphi(\mathbf{X}w)v) * h_2$$



Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2) v_1 + \phi(x_1 w_3 + x_2 w_4) v_2 + \phi(x_1 w_5 + x_2 w_6) v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say w_2 :

$$\begin{aligned} \partial/\partial w_2 L_{x,y}(w, v) &= \partial/\partial w_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(Xw)v] \end{aligned}$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2) v_1 + \phi(x_1 w_3 + x_2 w_4) v_2 + \phi(x_1 w_5 + x_2 w_6) v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say w₂:

$$\begin{aligned}\partial/\partial w_2 L_{x,y}(w, v) &= \partial/\partial w_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(Xw)v] \\ &= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(x_1 w_1 + x_2 w_2)v_1]\end{aligned}$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2) v_1 + \phi(x_1 w_3 + x_2 w_4) v_2 + \phi(x_1 w_5 + x_2 w_6) v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say w₂:

$$\begin{aligned}\partial/\partial w_2 L_{x,y}(w, v) &= \partial/\partial w_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(Xw)v] \\ &= (Y - \phi(Xw)v) * -v_1 [\partial/\partial w_2 \phi(x_1 w_1 + x_2 w_2)]\end{aligned}$$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2) v_1 + \phi(x_1 w_3 + x_2 w_4) v_2 + \phi(x_1 w_5 + x_2 w_6) v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

Want partial derivative of Loss w.r.t. weights, say

$$\begin{aligned}\partial/\partial w_2 L_{x,y}(w, v) &= \partial/\partial w_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\ &= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(Xw)v] \\ &= (Y - \phi(Xw)v) * -v_1 [\partial/\partial w_2 \phi(x_1 w_1 + x_2 w_2)] \\ &= (Y - \phi(Xw)v) * -v_1 \phi'(x_1 w_1 + x_2 w_2) [\partial/\partial w_2 (x_1 w_1 + x_2 w_2)]\end{aligned}$$

Chain rule!
If $F(x) = f(g(x))$
 $F'(x) = f'(g(x))g'(x)$

Backpropagation: the math

1-layer NN, sigmoid activation at hidden layer, linear output:

$$F(X) = \phi(Xw)v$$

$$F(X) = \phi(x_1 w_1 + x_2 w_2) v_1 + \phi(x_1 w_3 + x_2 w_4) v_2 + \phi(x_1 w_5 + x_2 w_6) v_3$$

Say regression, Loss function is $\frac{1}{2}$ L₂ norm, expected output is Y:

$$L_{x,y}(w, v) = \frac{1}{2}(Y - F(X))^2 = \frac{1}{2}(Y - \phi(Xw)v)^2$$

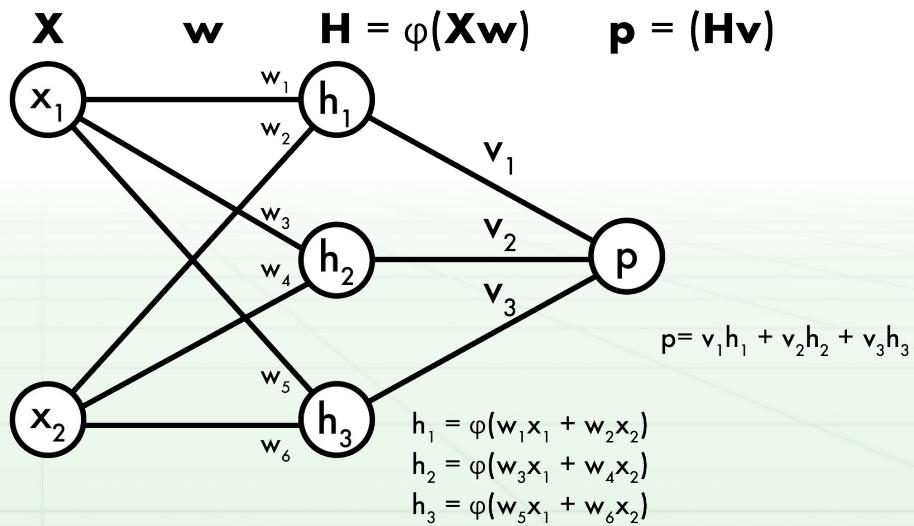
Want partial derivative of Loss w.r.t. weights, say w₂:

$$\begin{aligned}\partial/\partial w_2 L_{x,y}(w, v) &= \partial/\partial w_2 \frac{1}{2}(Y - \phi(Xw)v)^2 \\&= (Y - \phi(Xw)v) * -[\partial/\partial w_2 \phi(Xw)v] \\&= (Y - \phi(Xw)v) * -v_1 [\partial/\partial w_2 \phi(x_1 w_1 + x_2 w_2)] \\&= (Y - \phi(Xw)v) * -v_1 \phi'(x_1 w_1 + x_2 w_2) [\partial/\partial w_2 (x_1 w_1 + x_2 w_2)] \\&= (Y - \phi(Xw)v) * -v_1 \phi'(x_1 w_1 + x_2 w_2) * x_2\end{aligned}$$

Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say w_2 :

$$\frac{\partial}{\partial w_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -v_1\varphi'(x_1w_1 + x_2w_2) * x_2$$

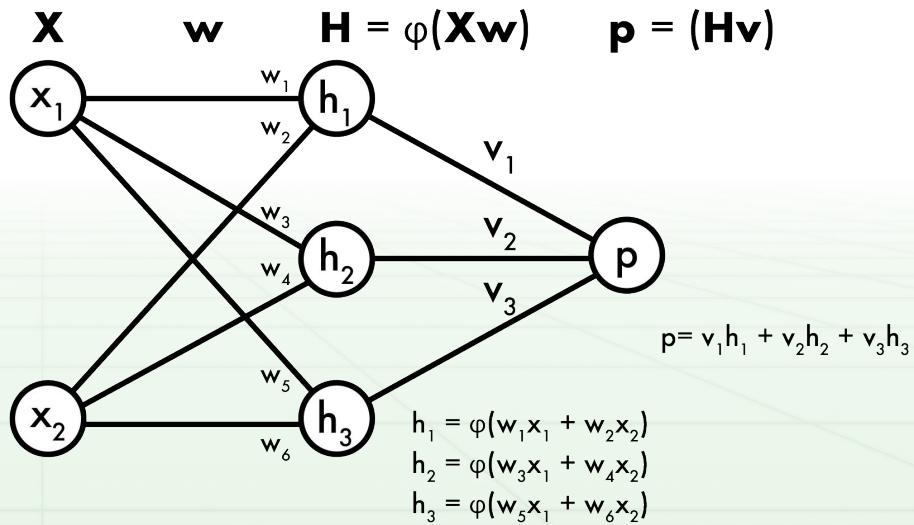


Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say w_2 :

$$\frac{\partial}{\partial w_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -v_1\varphi'(x_1w_1 + x_2w_2) * x_2$$

Model error at p



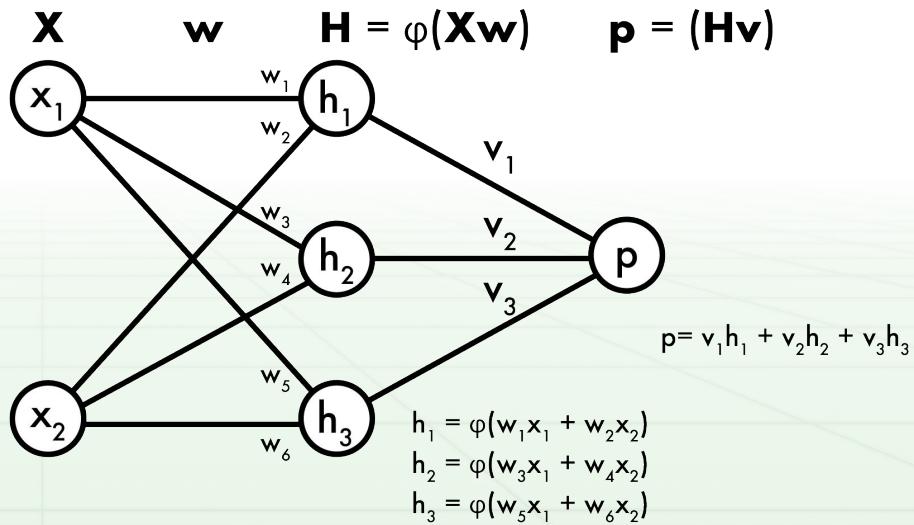
Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say w_2 :

$$\frac{\partial}{\partial w_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -v_1\varphi'(x_1w_1 + x_2w_2) * x_2$$

Model error at p

Backpropagate through v_1



Backpropagation: the math

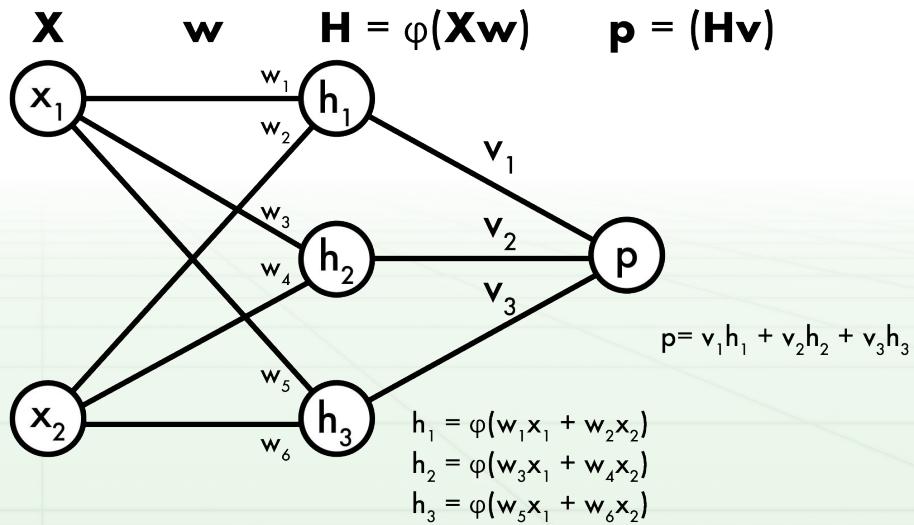
Want partial derivative of Loss w.r.t. weights, say w_2 :

$$\frac{\partial}{\partial w_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -v_1\varphi'(x_1w_1 + x_2w_2) * x_2$$

Model error at p

Model error at h_1

Backpropagate through v_1



Backpropagation: the math

Want partial derivative of Loss w.r.t. weights, say w_2 :

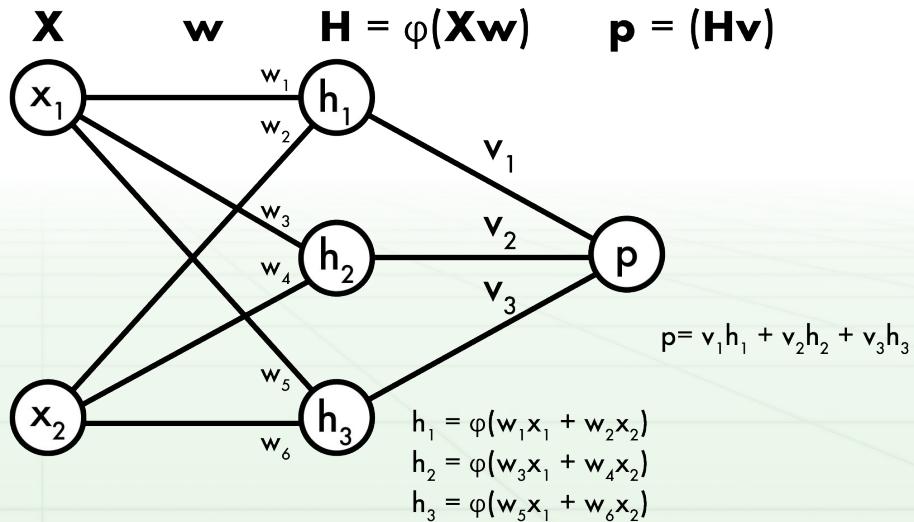
$$\frac{\partial}{\partial w_2} L_{x,y}(w, v) = (Y - \varphi(\mathbf{X}w)v) * -v_1 \varphi'(x_1 w_1 + x_2 w_2) * x_2$$

Model error at p

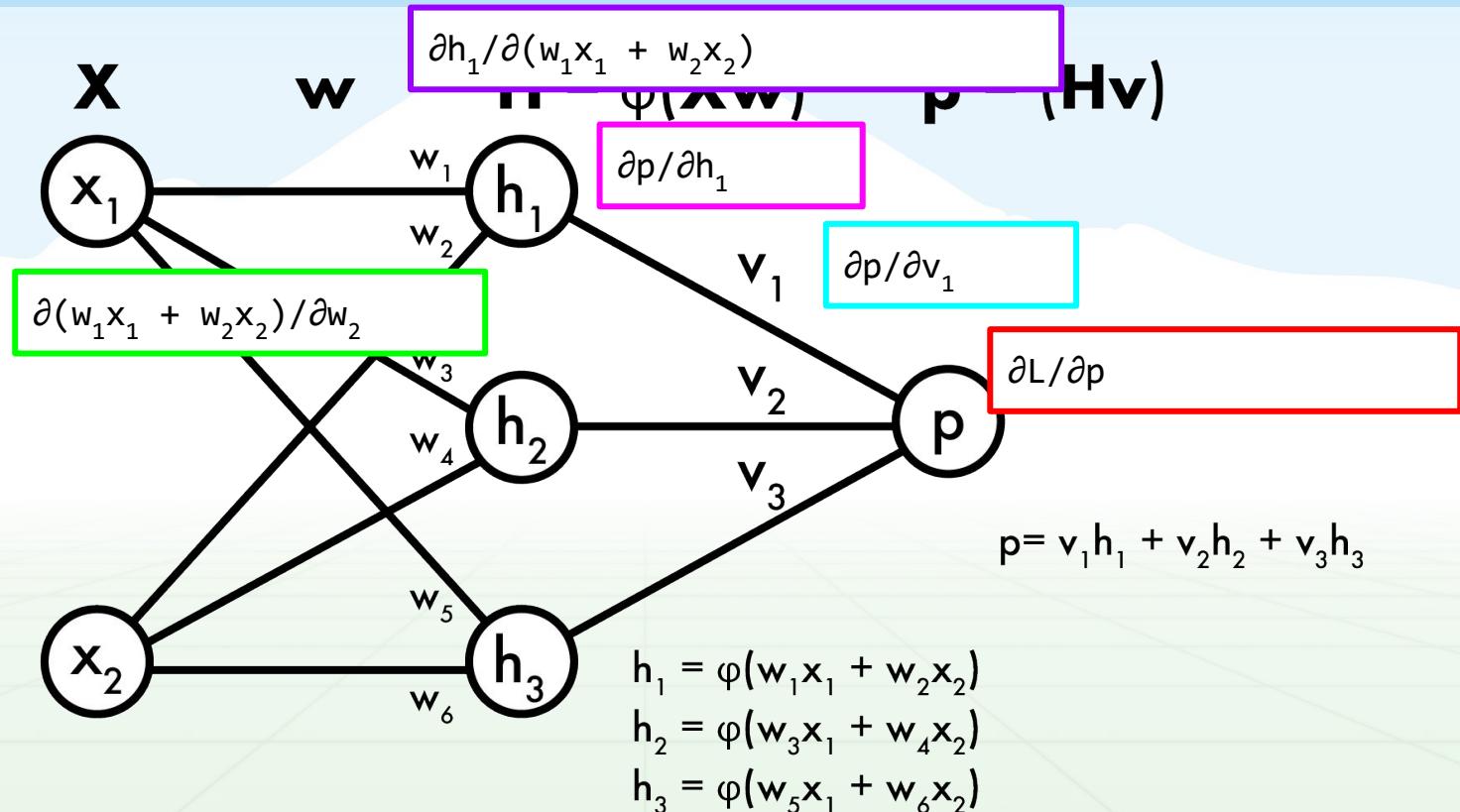
Model error at h_1

Multiply by x_2 :
gradient w.r.t. w_2

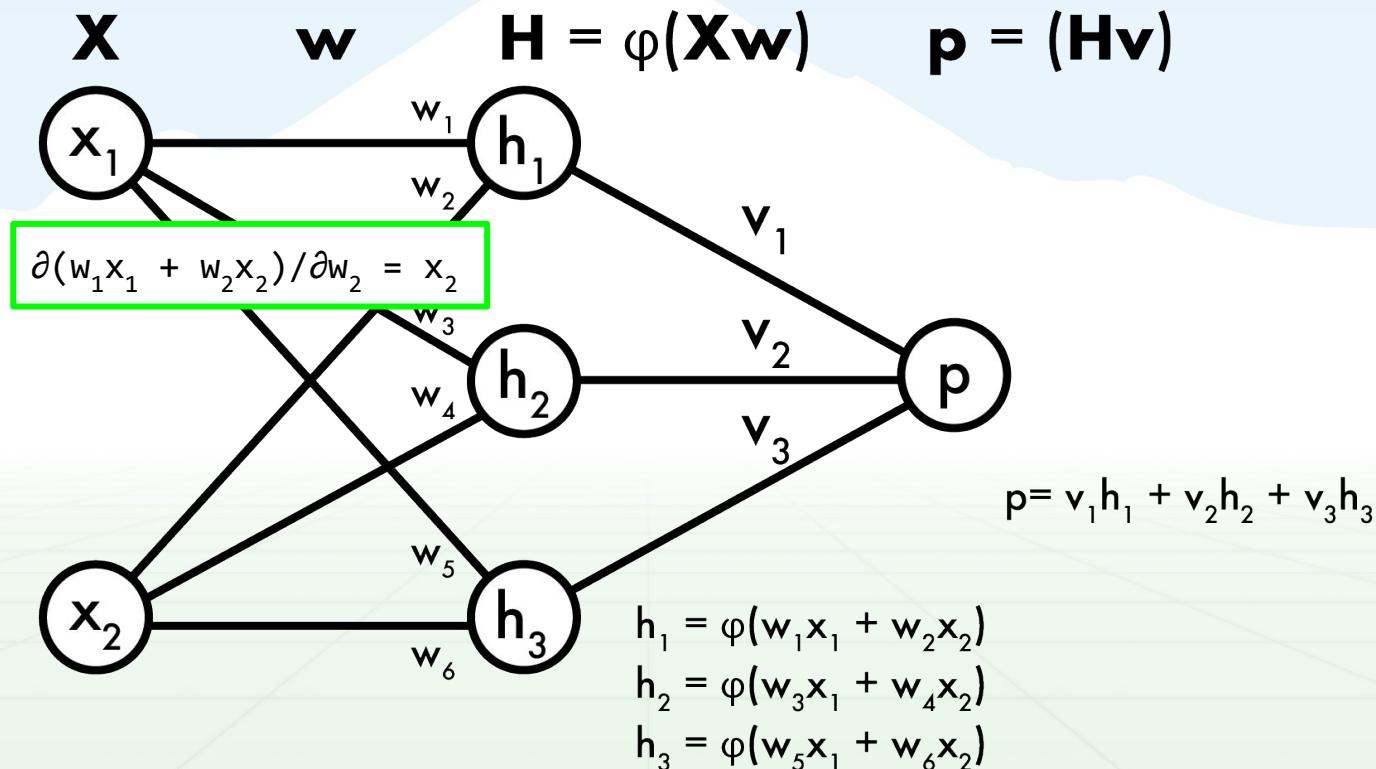
Backpropagate through v_1



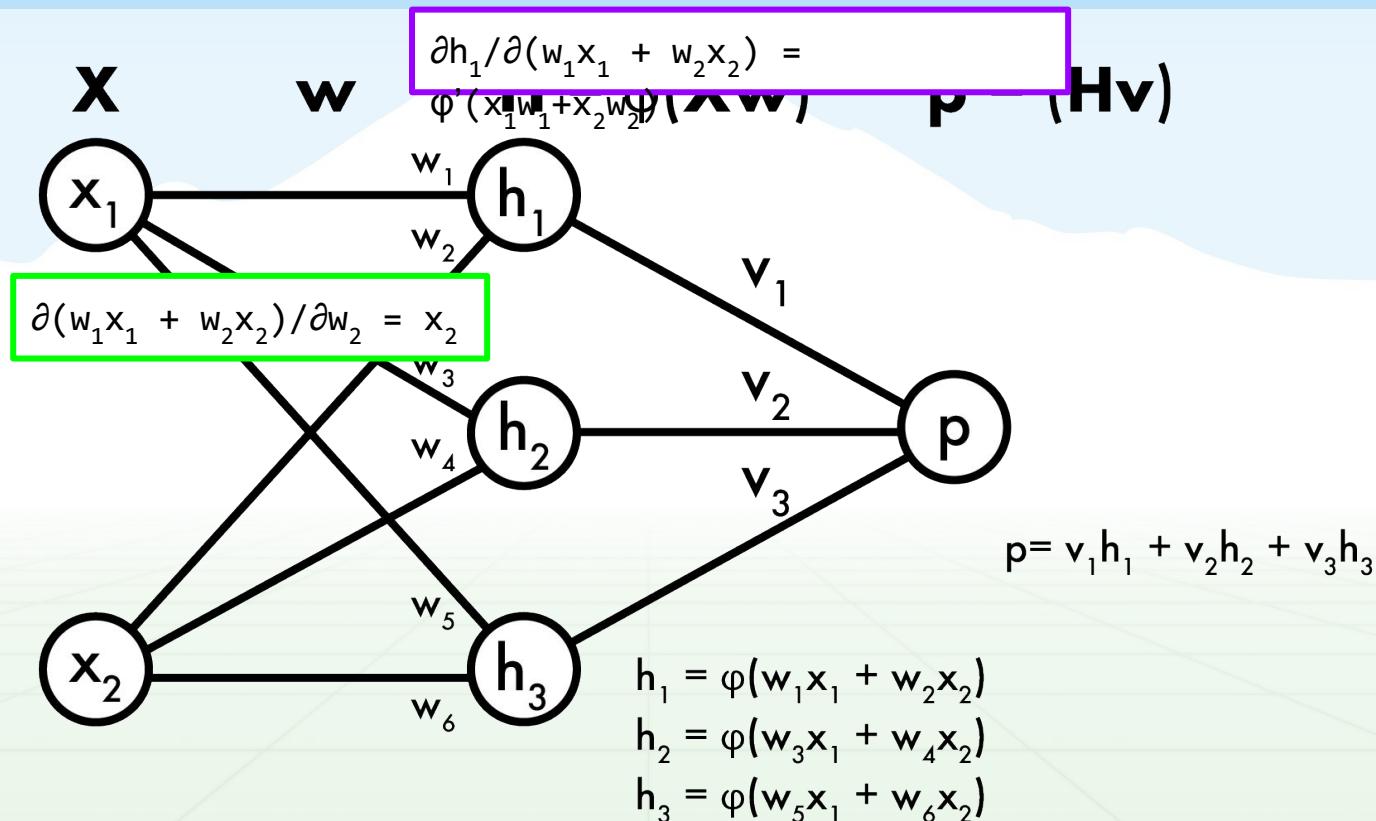
Backpropagation: the math



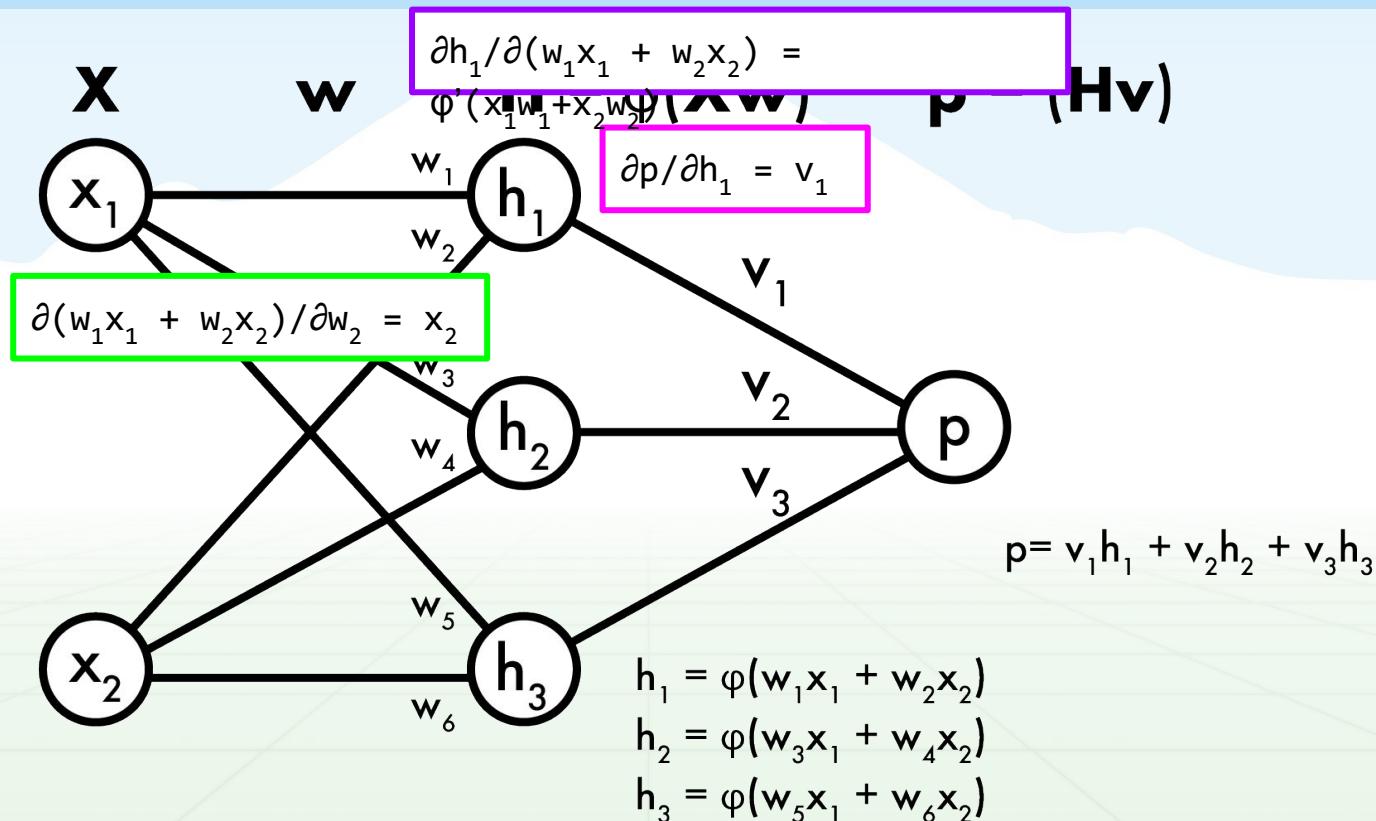
Backpropagation: the math



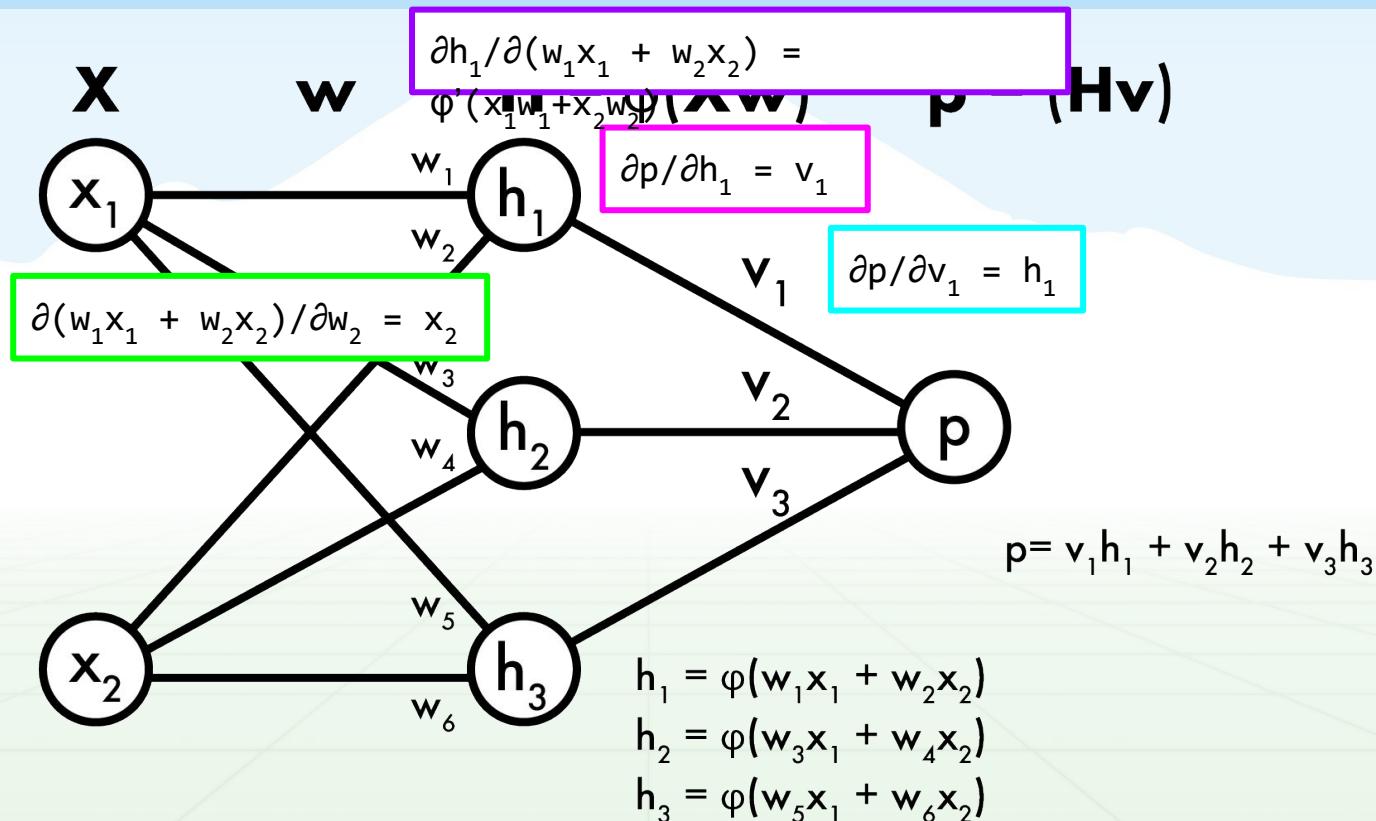
Backpropagation: the math



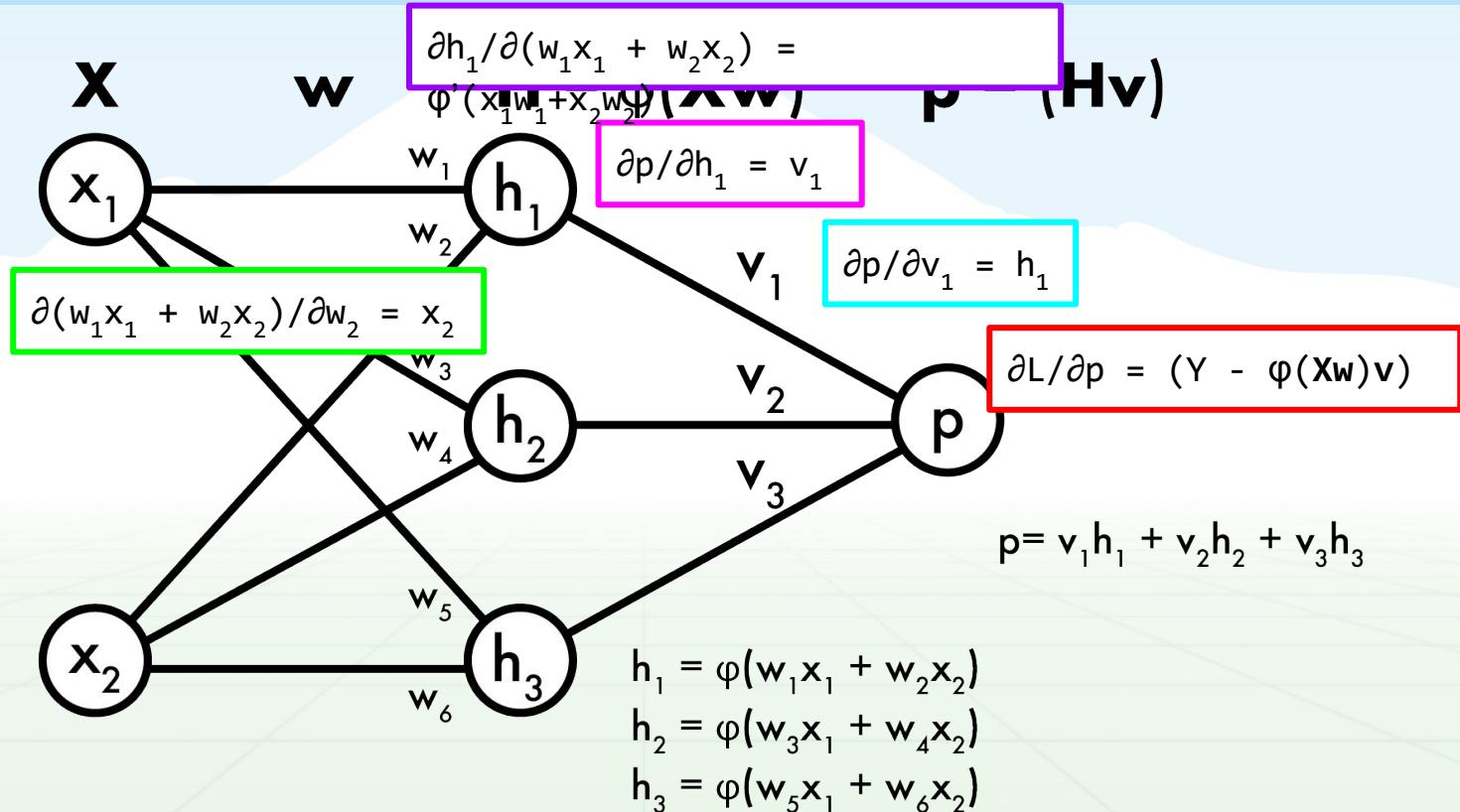
Backpropagation: the math



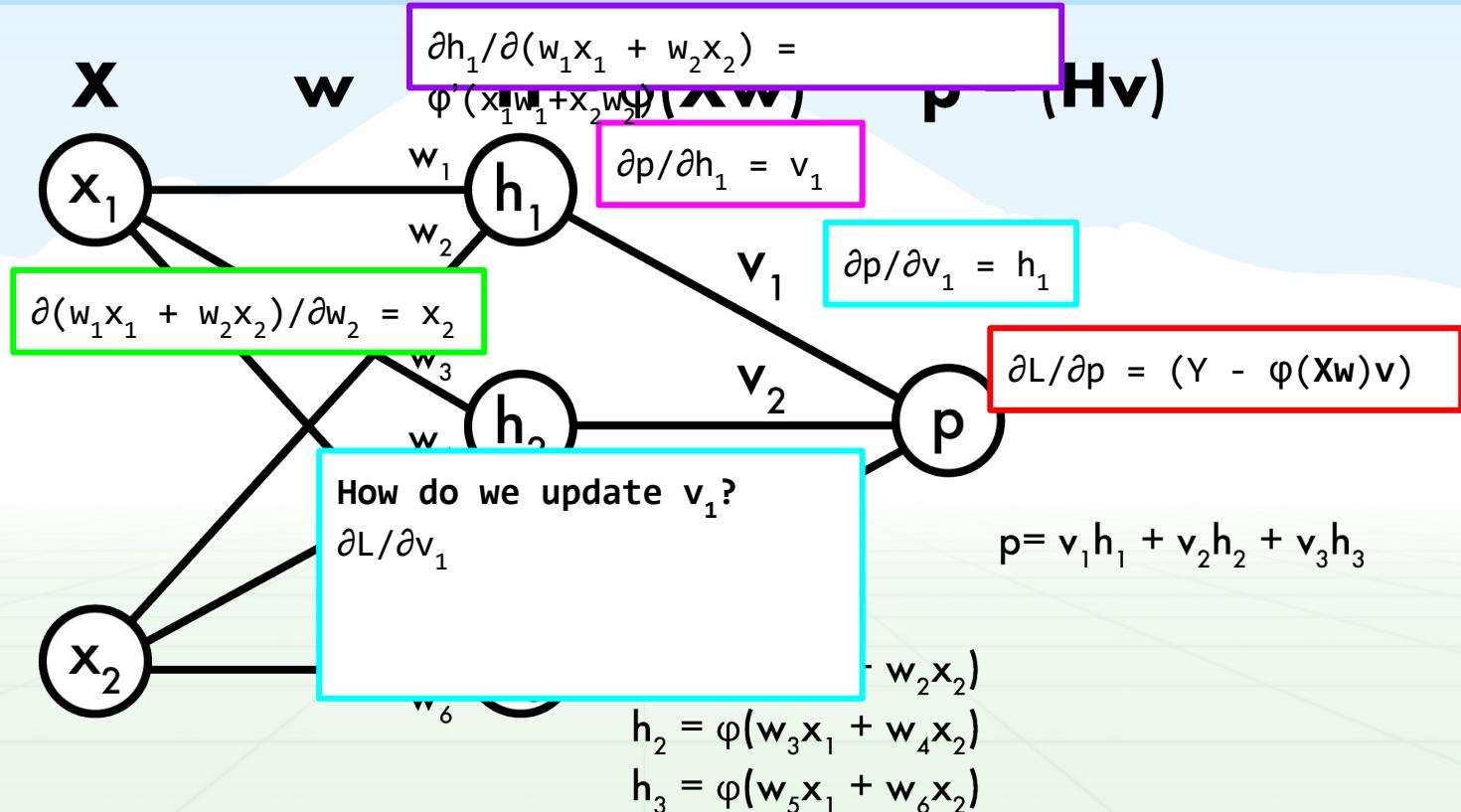
Backpropagation: the math



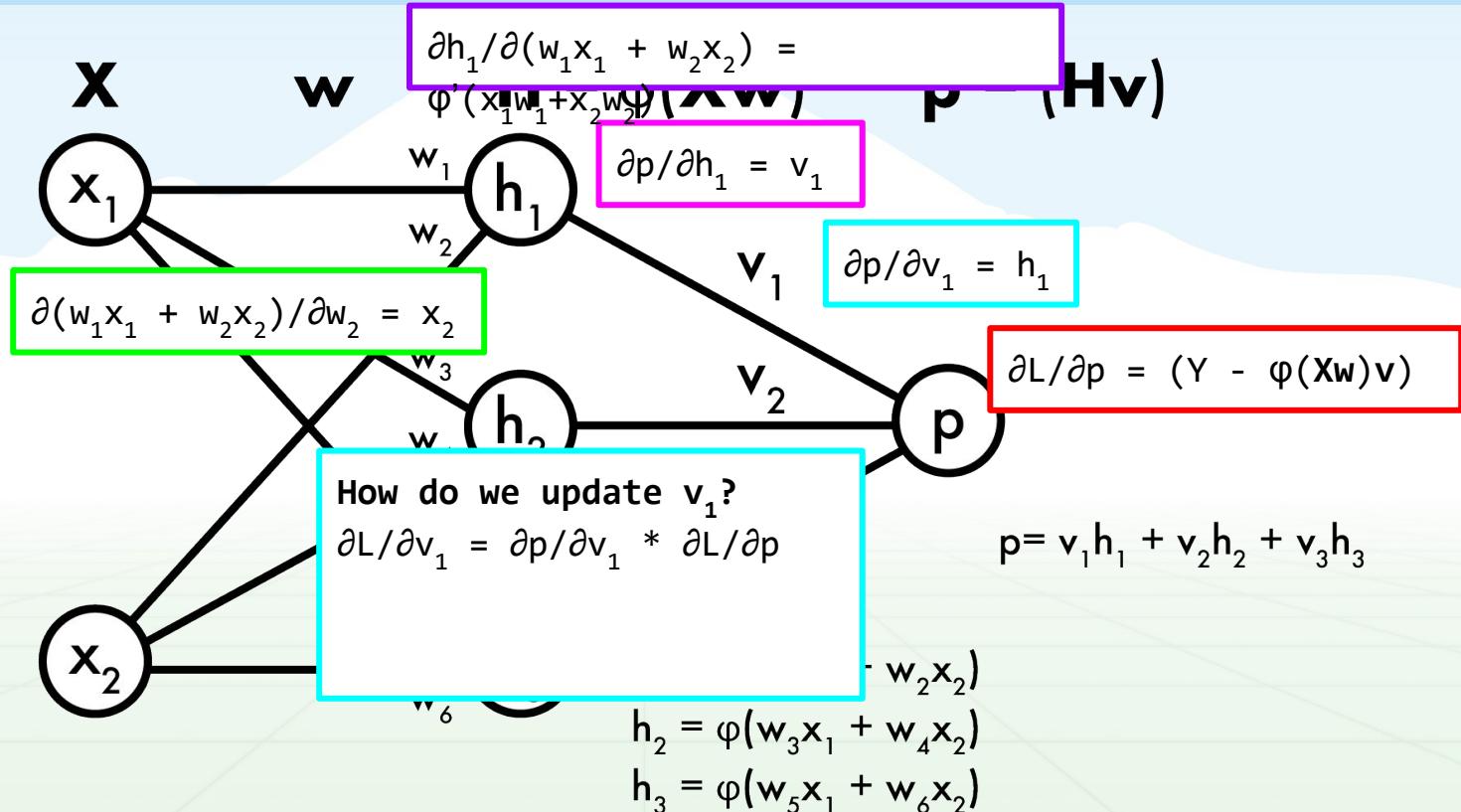
Backpropagation: the math



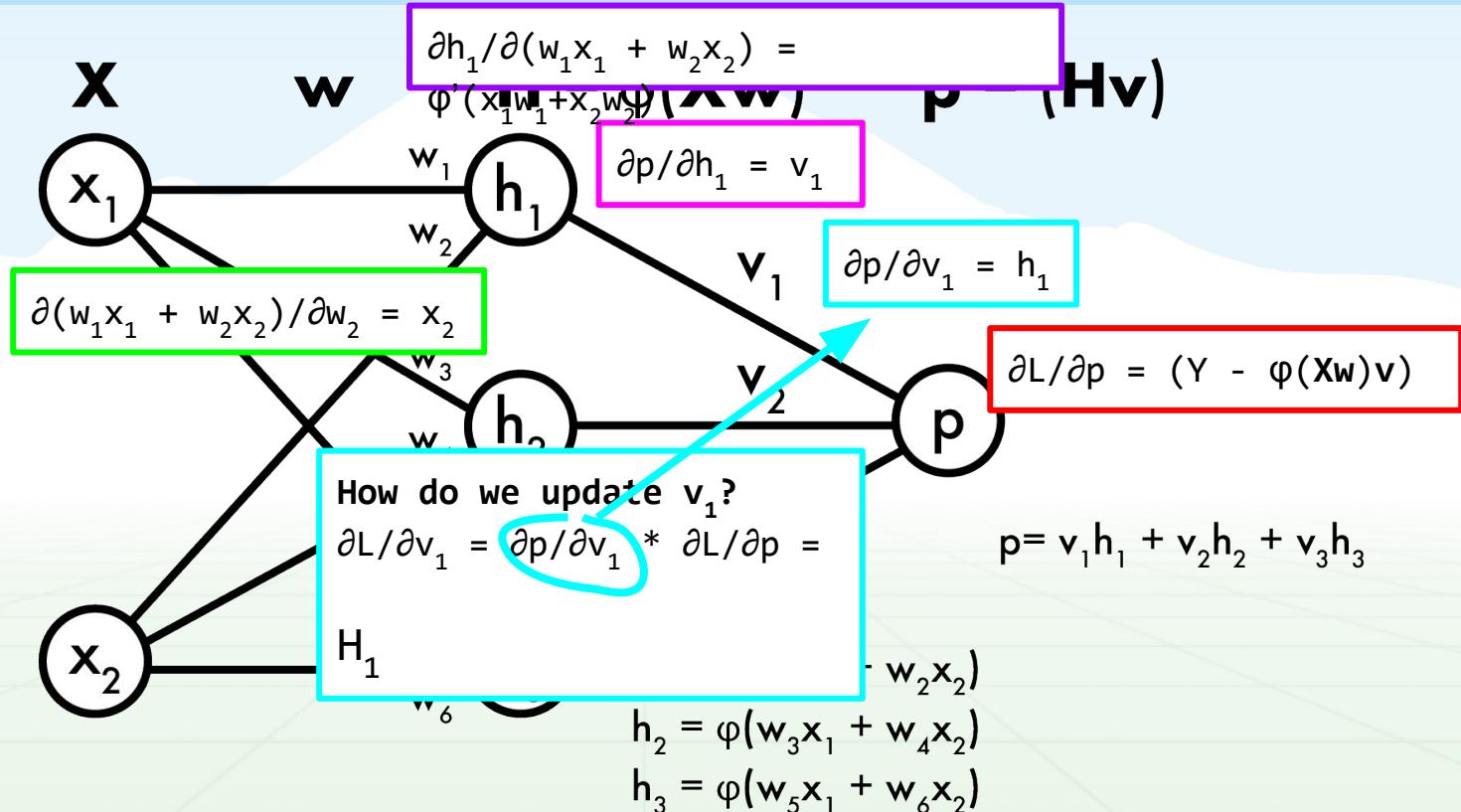
Backpropagation: the math



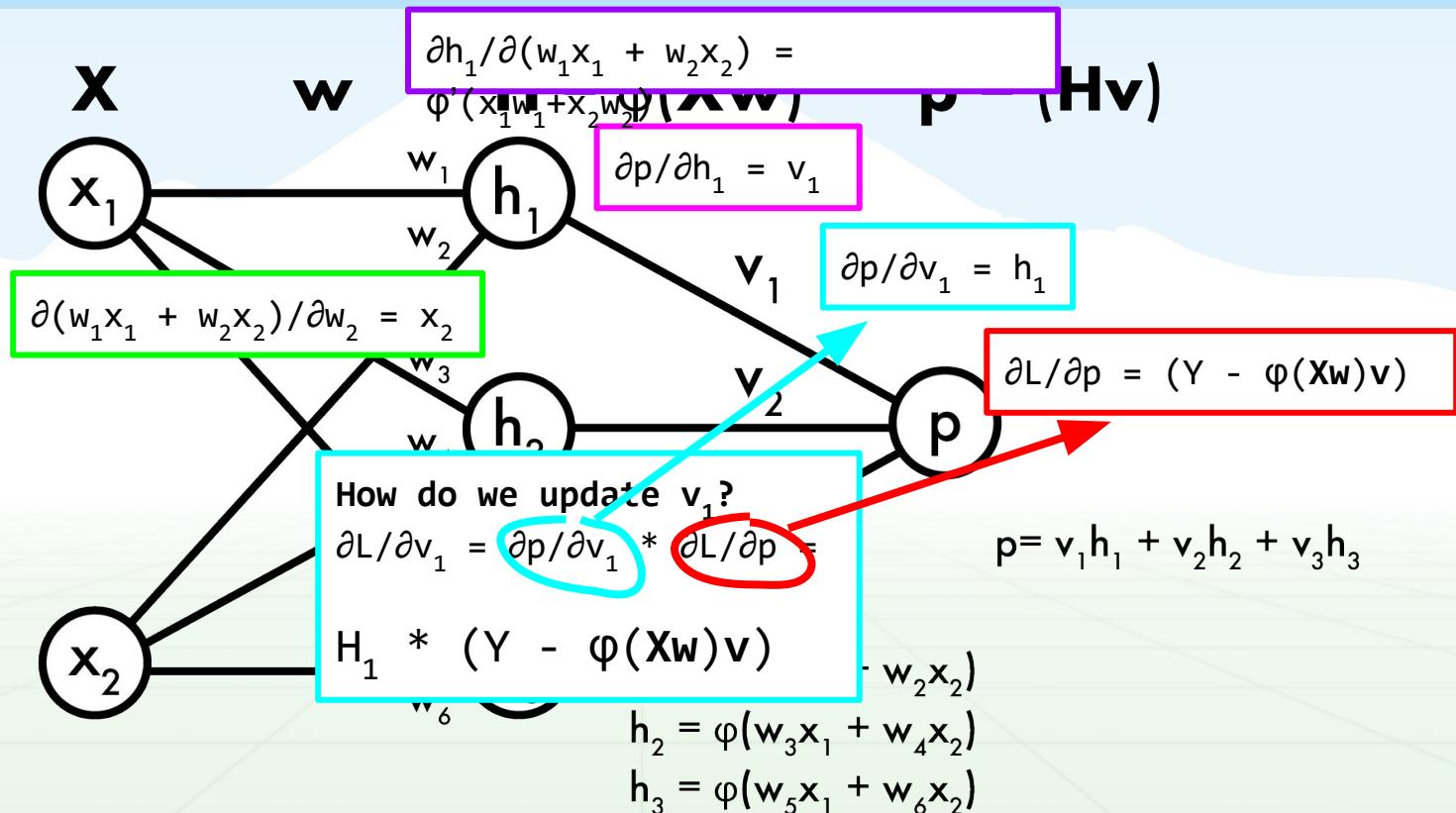
Backpropagation: the math



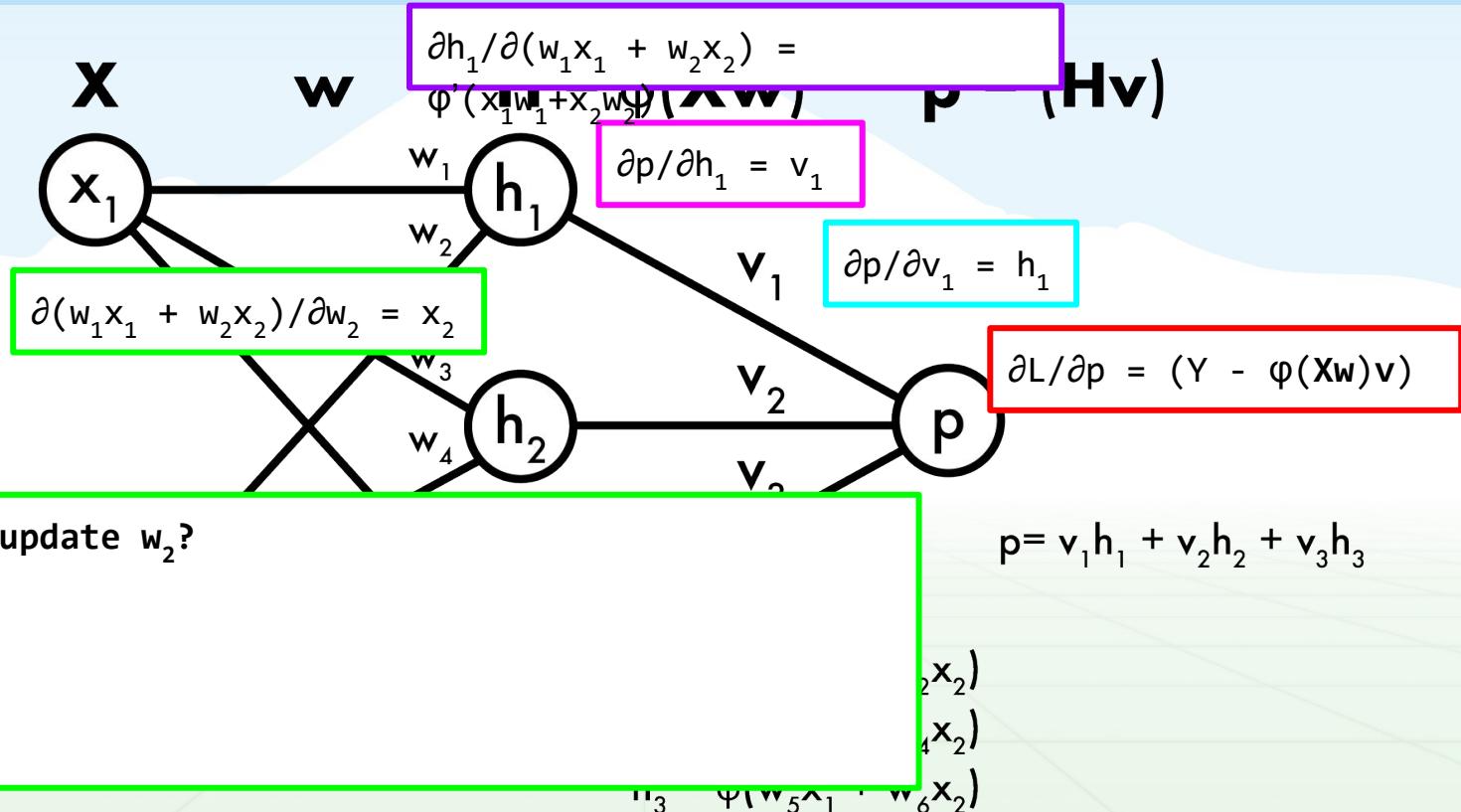
Backpropagation: the math



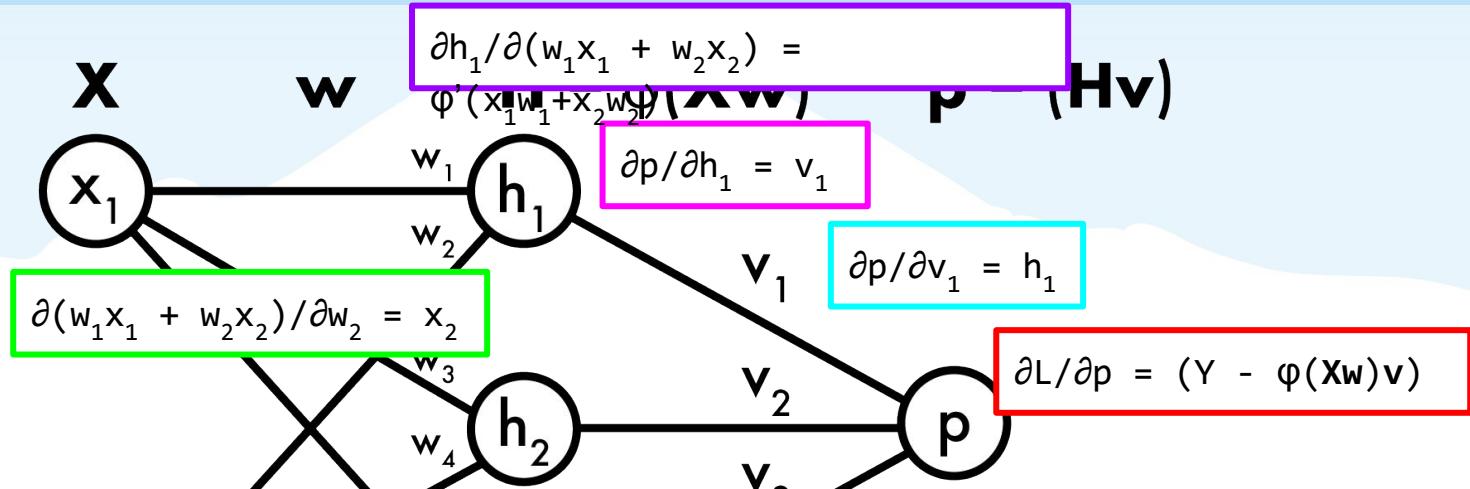
Backpropagation: the math



Backpropagation: the math



Backpropagation: the math



How do we update w_2 ?

$$\partial L / \partial w_2 =$$

$$\partial(w_1 x_1 + w_2 x_2) / \partial w_2 * \partial h_1 / \partial (w_1 x_1 + w_2 x_2) * \partial p / \partial h_1 * \partial L / \partial p$$

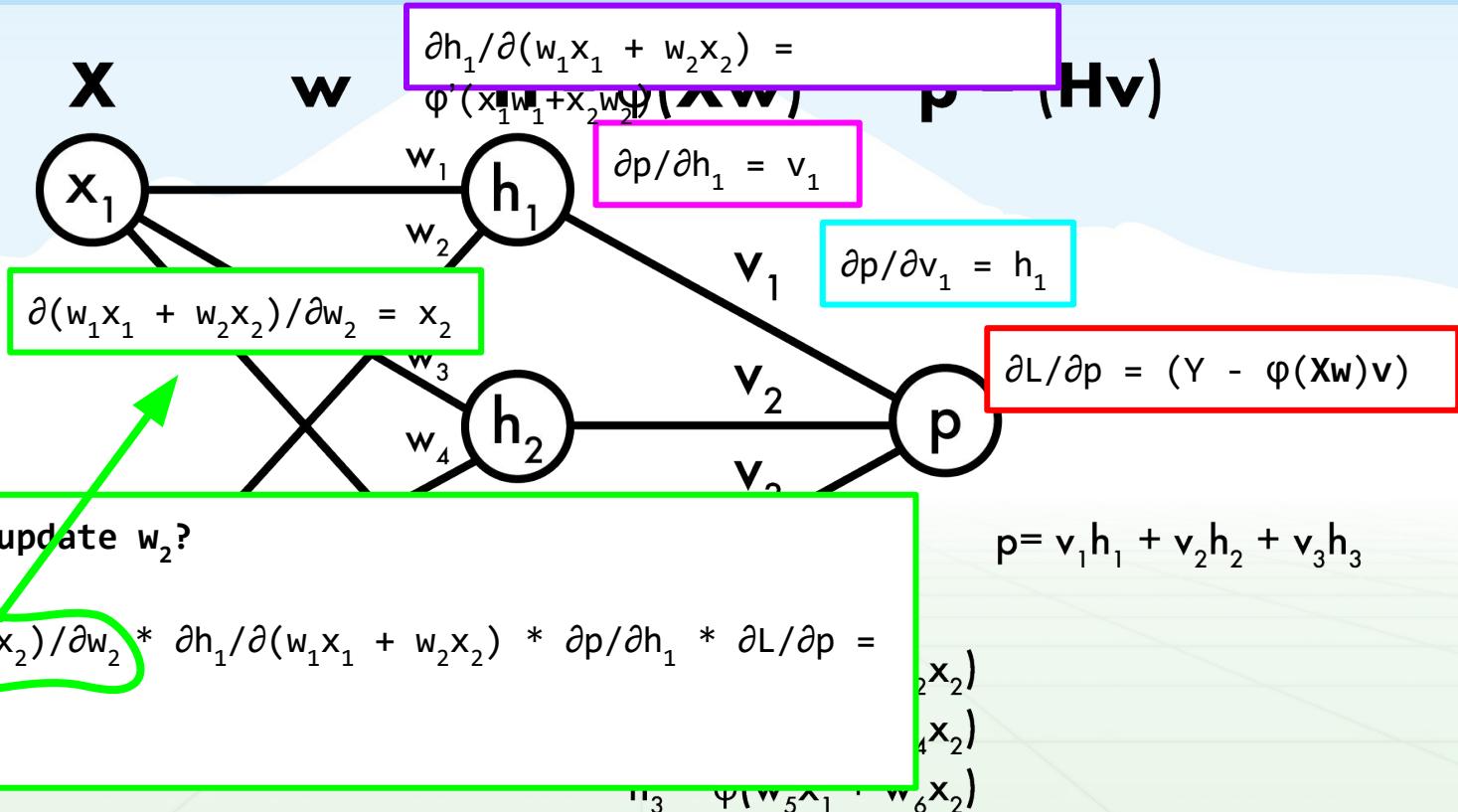
$$p = v_1 h_1 + v_2 h_2 + v_3 h_3$$

v_1
 v_2
 v_3
 h_1
 h_2
 h_3

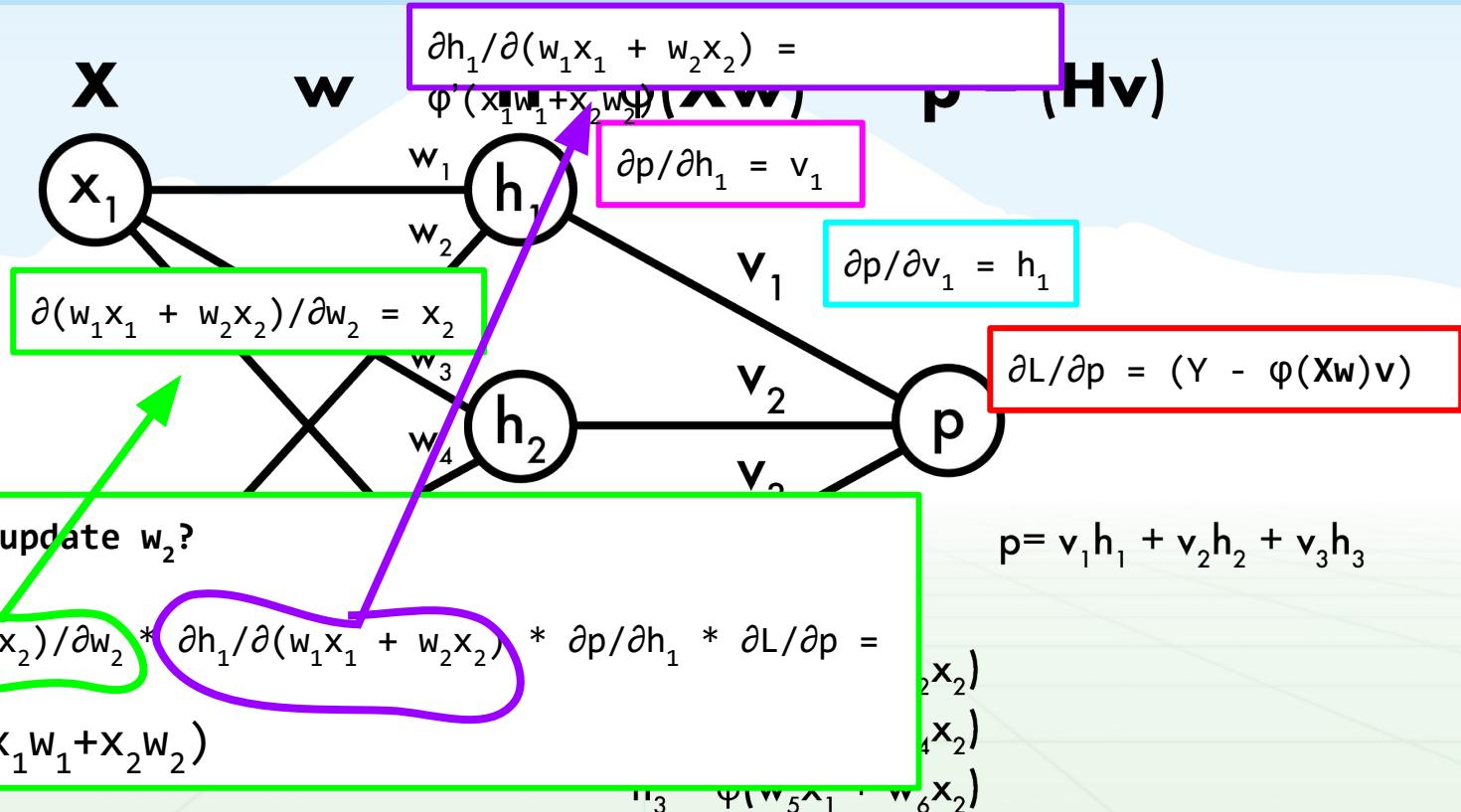
w_1
 w_2
 w_3
 w_4
 w_5
 w_6

x_1
 x_2
 x_3
 x_4
 x_5
 x_6

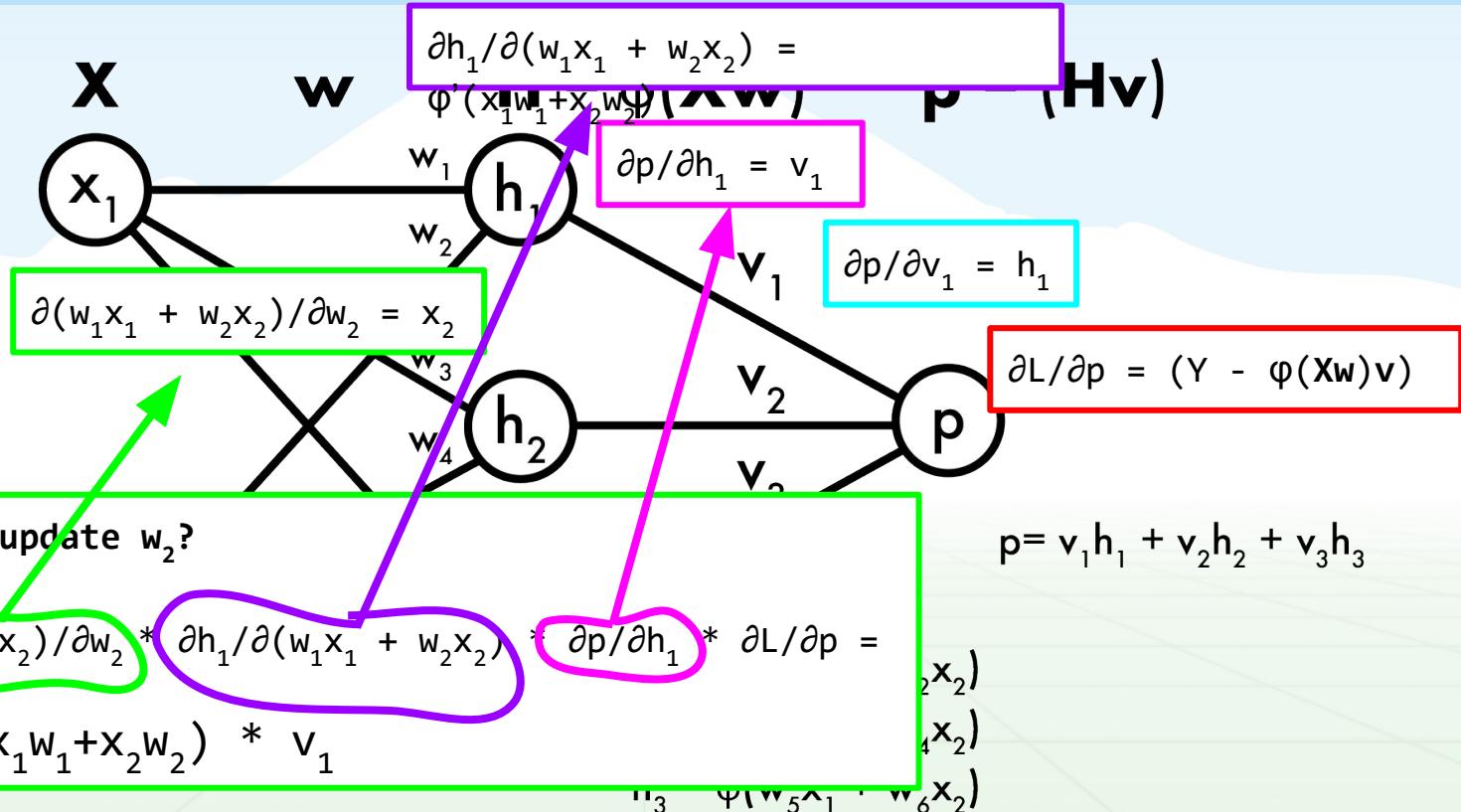
Backpropagation: the math



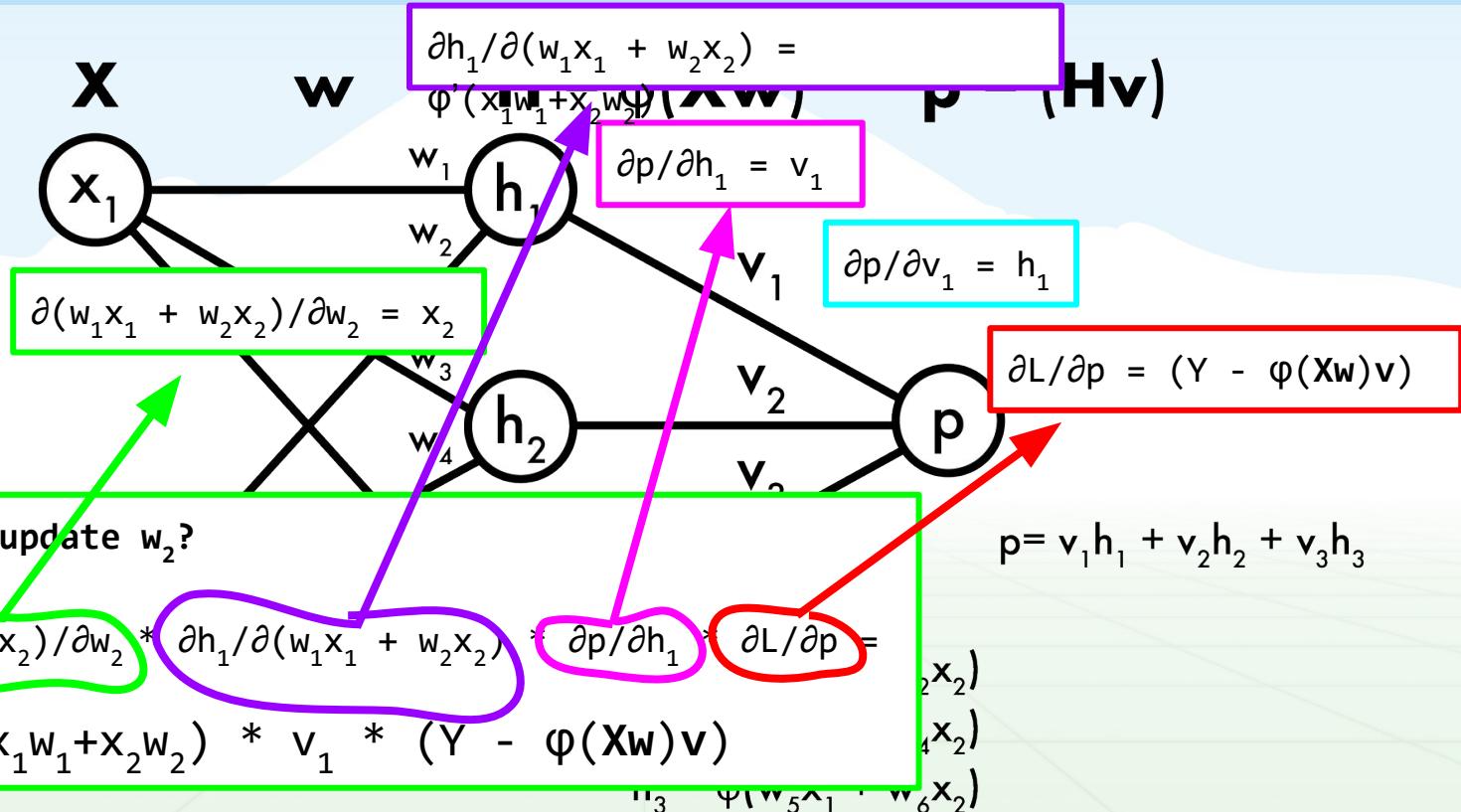
Backpropagation: the math



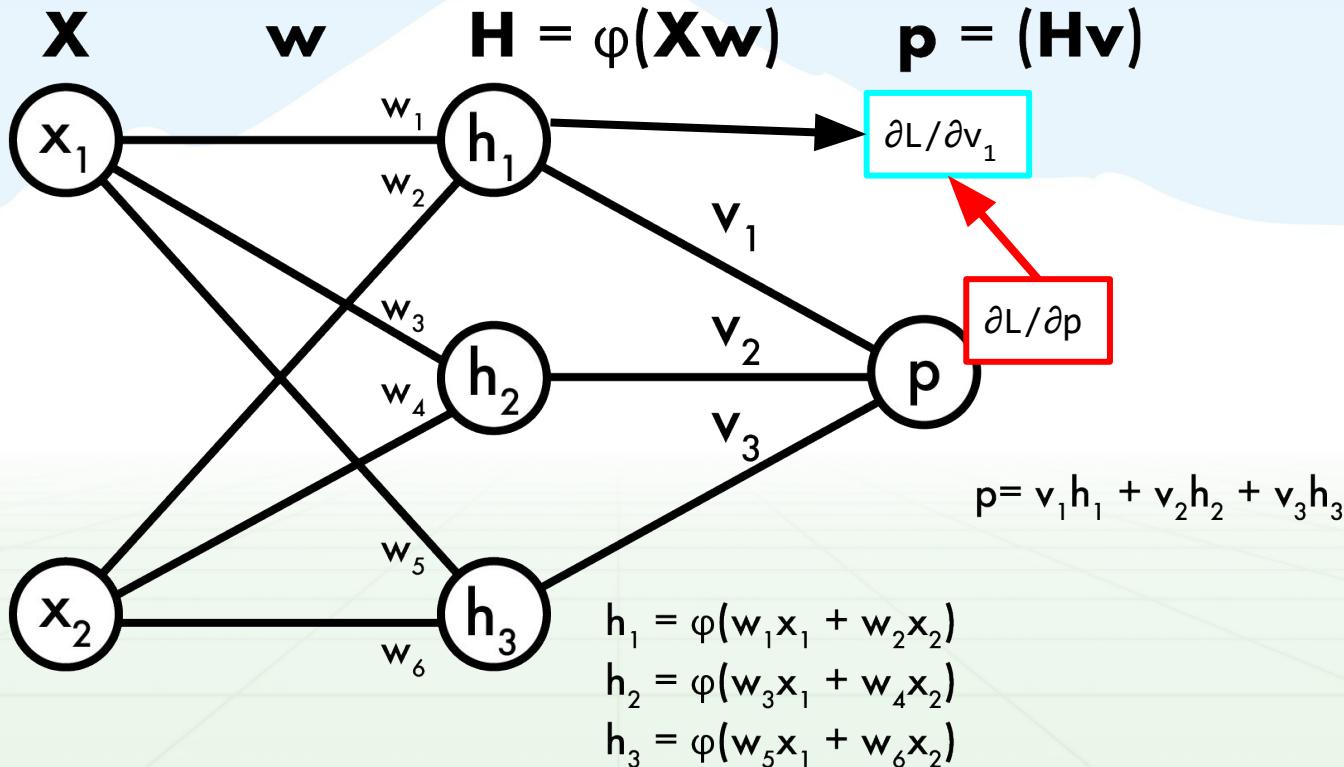
Backpropagation: the math



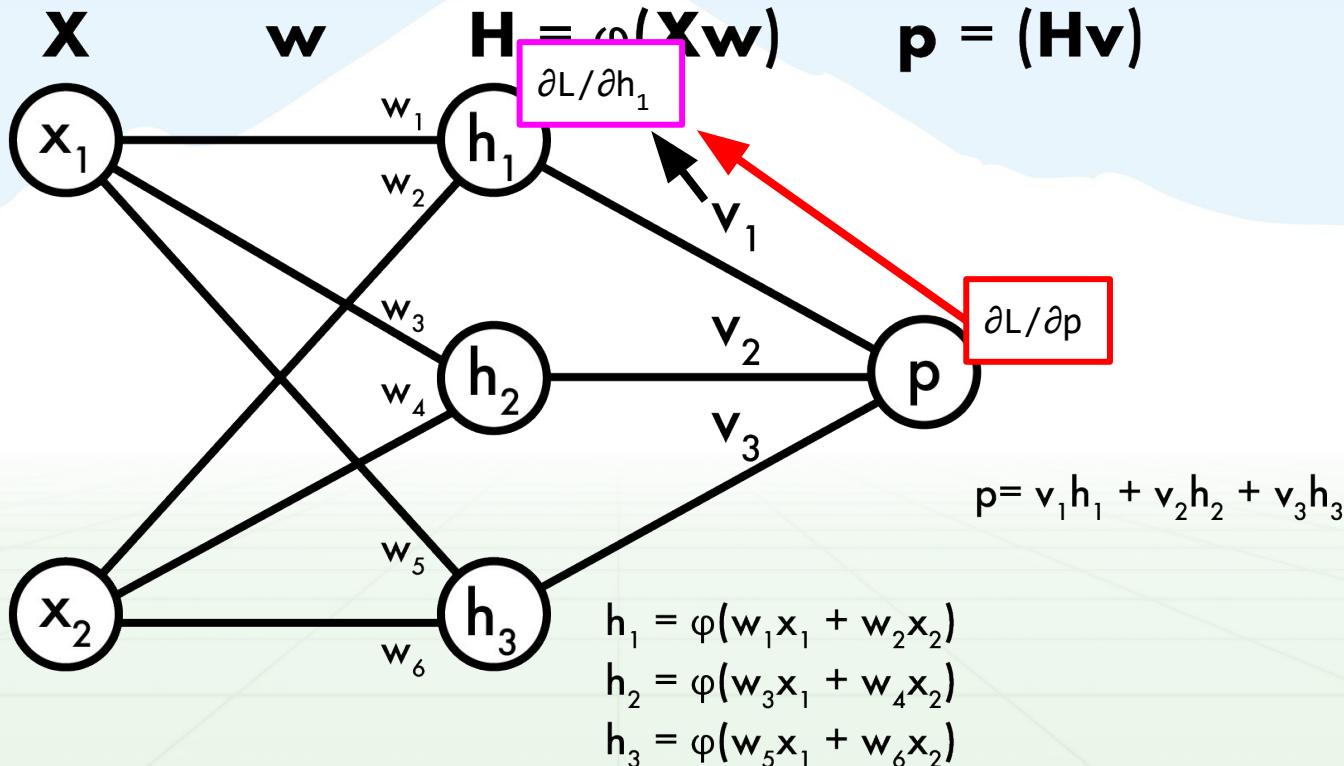
Backpropagation: the math



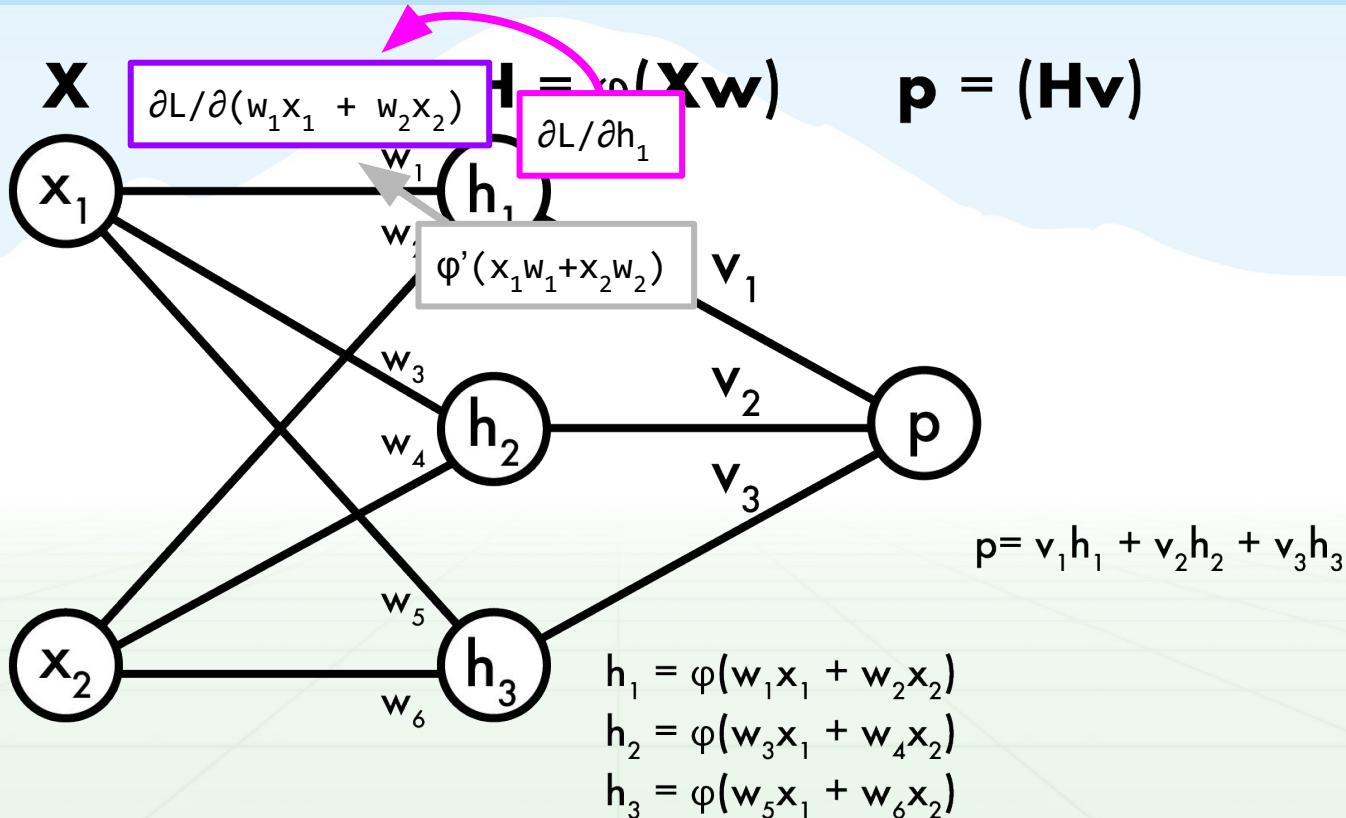
Backpropagation: the math



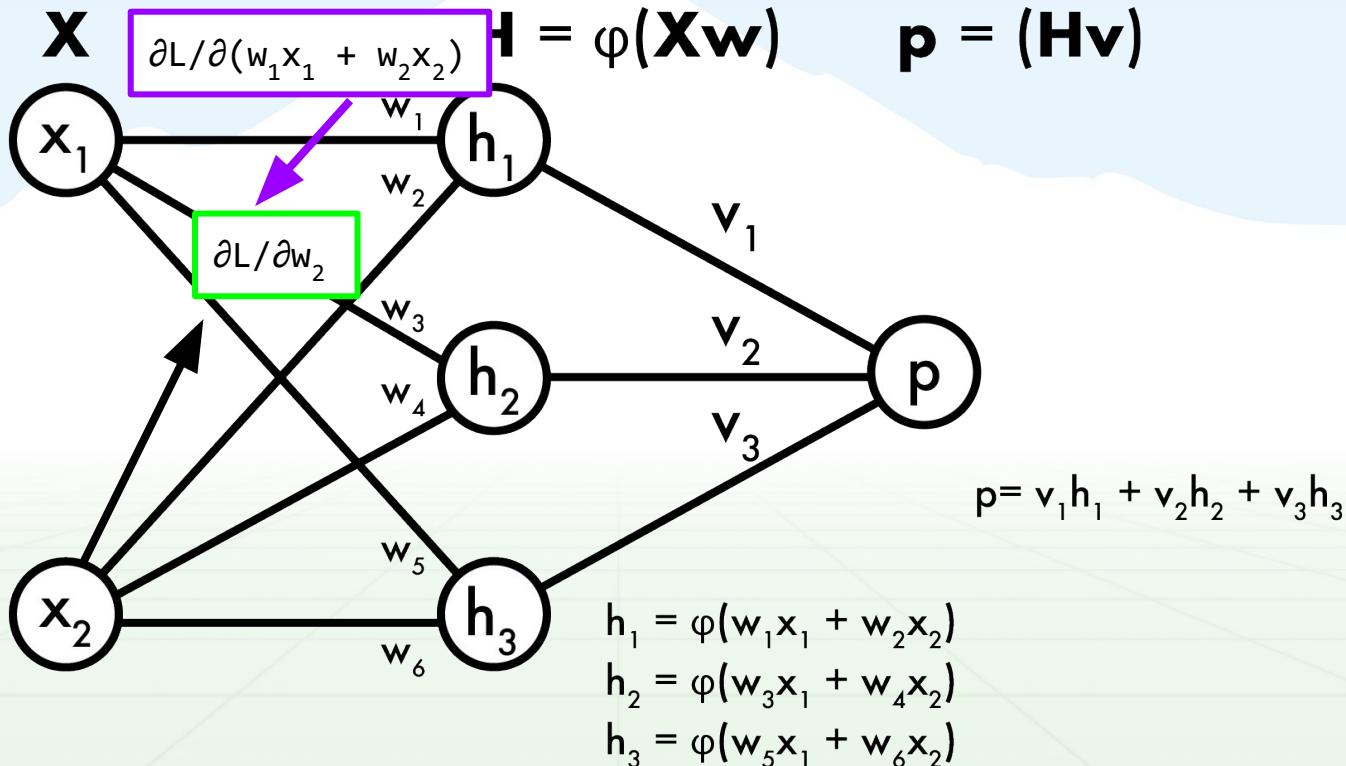
Backpropagation: the math



Backpropagation: the math

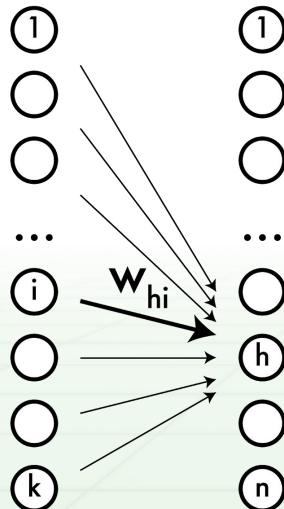


Backpropagation: the math

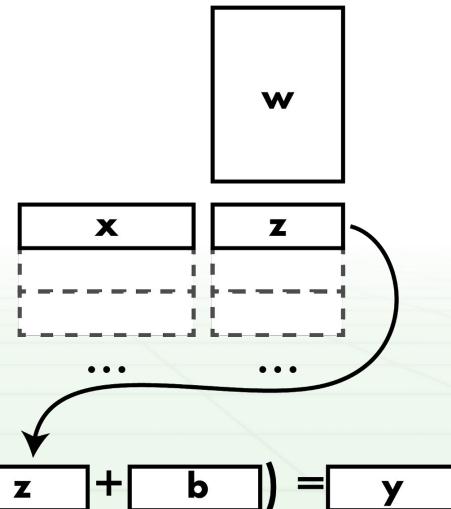


Forward propagation

$$y_h = \varphi(\sum x_i w_{hi} + b)$$



$$\mathbf{y} = \varphi(\mathbf{xw} + \mathbf{b})$$

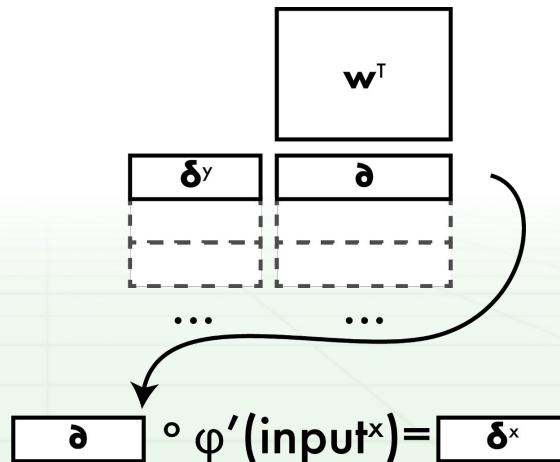
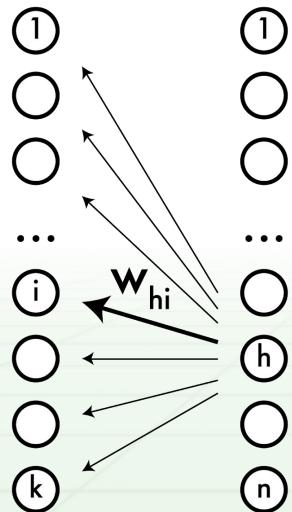


Layer x Layer y

$\mathbf{x}: m \times k$ $\mathbf{w}: k \times n$ $\mathbf{z}: m \times n$

Backward propagation

$$\delta_i^x = \left(\sum_j \delta_j^y w_{hi} \right) \varphi'(\text{input}_i^x) \quad \boldsymbol{\delta}^x = \boldsymbol{\delta}^y \mathbf{w}^T \circ \varphi'(\text{input}^x)$$



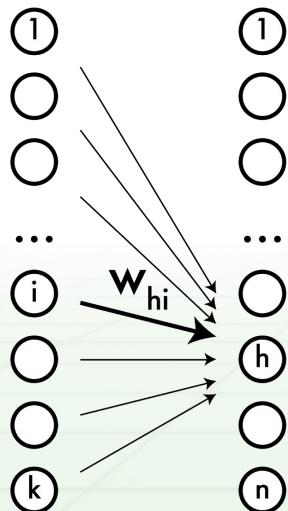
Layer x

Layer y

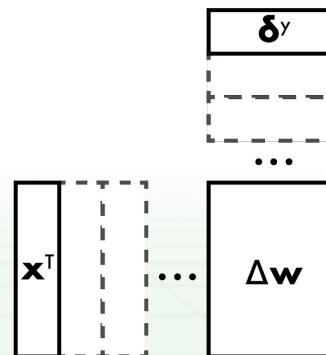
$\boldsymbol{\delta}^y: m \times n \quad \mathbf{w}^T: n \times k \quad \boldsymbol{\delta}^x: m \times k$

Weight updates

$$\Delta \mathbf{w}_{hi} = \delta_h^y \mathbf{x}_i$$



$$\Delta \mathbf{w} = \mathbf{x}^T \boldsymbol{\delta}^y$$



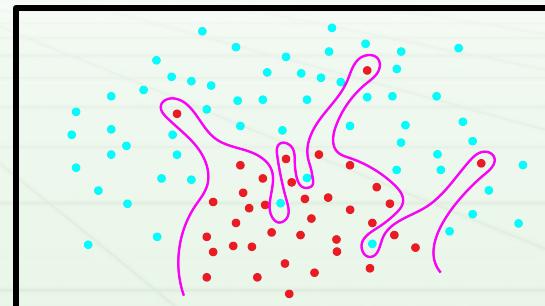
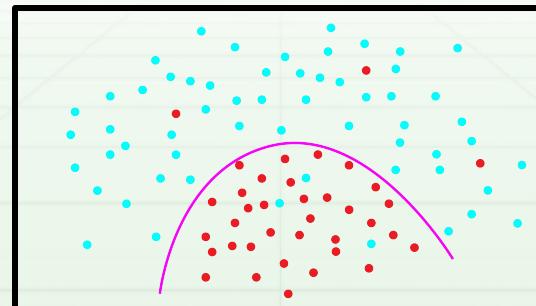
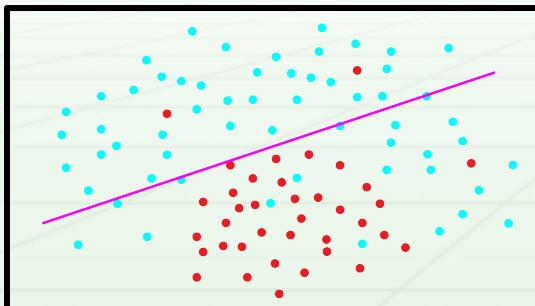
Layer x Layer y $\mathbf{x}^T: k \times m$ $\boldsymbol{\delta}^y: m \times n$ $\Delta \mathbf{w}: k \times n$

Under and Overfitting

Underfitting: model not powerful enough, too much bias

Overfitting: model too powerful, fits to noise, doesn't generalize well

Want the happy medium, how?



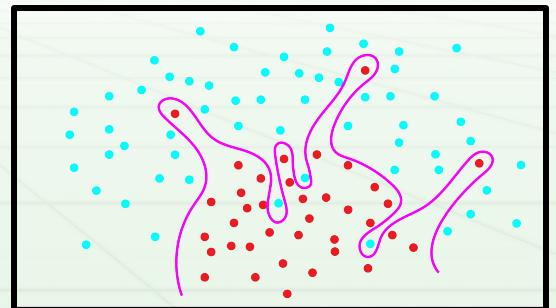
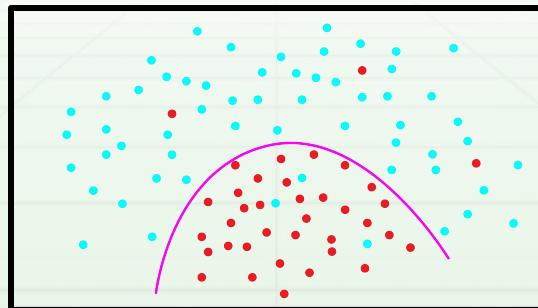
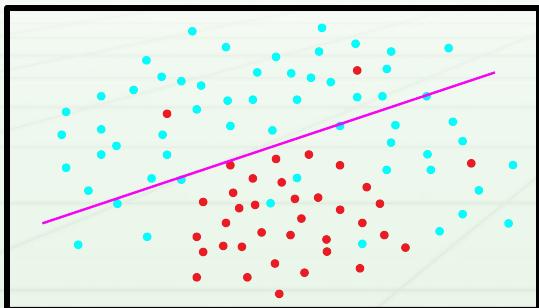
Under and Overfitting

Want the happy medium, how?

Pick the right model, but very hard to know a priori

Make weak model more powerful: boosting! (or other ways)

Make strong model less likely to overfit: *regularization*



With great power comes great *overfitting*

Neural networks are (sort of) all powerful! Which is not necessarily a good thing.

With great power comes great overfitting

Like SVMs, put limits on model that make it generalize better!

SVM:

$$\begin{aligned} & \min ||\mathbf{w}||_2 \\ \text{s.t. } & y_n(\mathbf{w} \cdot \mathbf{x}_n - b) \geq 1, \quad n = 1, 2 \dots \end{aligned}$$

Neural net:

Minimize loss function and weight magnitude

Before: $\operatorname{argmin}_{\mathbf{w}} L_x(\mathbf{w})$

Now: $\operatorname{argmin}_{\mathbf{w}} L_x(\mathbf{w}) + \lambda ||\mathbf{w}||_2$

Weight decay: neural network regularization

$$\operatorname{argmin}_w L_x(w) + \lambda \|w\|_2$$

λ : regularization parameter

Higher: more penalty for large weights, less powerful model

Lower: less penalty, more overfitting

Commonly use L_2 norm to regularize, *weight decay*

Gradient descent update rule:

$$w_{t+1} = w_t - \eta [\partial / \partial w_t L(w_t) + \lambda w_t]$$

$$= w_t - \eta \partial / \partial w_t L(w_t) - \eta \lambda w_t$$

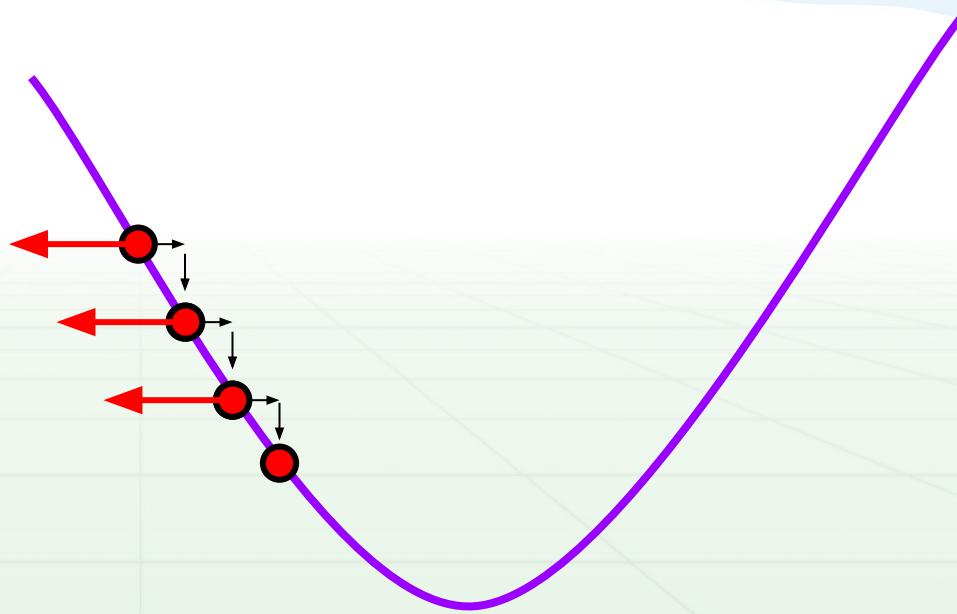
Subtract a little
bit of weight
every iteration

Sometimes training is SLOW

With SGD we make LOTS of little steps along the gradient

Sometimes we move in the same direction for a long time...

Maybe we should speed up in that direction!



Momentum: speeding up SGD

If we keep moving in same direction we should move further every round

Before:

$$\Delta w_t = -\partial/\partial w_t L(w_t)$$

Now:

$$\Delta w_t = -\partial/\partial w_t L(w_t) + m \Delta w_{t-1}$$

$$w_{t+1} = w_t + \eta \Delta w_t$$

Side effect: smooths out updates if gradient is in different directions

NN updates with weight decay and momentum

$$\Delta w_t = -\partial/\partial w_t L(w_t) - \lambda w_t + m \Delta w_{t-1}$$

$$w_{t+1} = w_t + \eta \Delta w_t$$

NN updates with weight decay and momentum

$$\Delta w_t = -\partial/\partial w_t L(w_t) - \lambda w_t + m \Delta w_{t-1}$$

Gradient of loss

$$w_{t+1} = w_t + \eta \Delta w_t$$

NN updates with weight decay and momentum

$$\Delta w_t = -\partial/\partial w_t L(w_t) - \lambda w_t + m \Delta w_{t-1}$$

Gradient of loss

Weight decay

$$w_{t+1} = w_t + \eta \Delta w_t$$

NN updates with weight decay and momentum

$$\Delta w_t = -\partial/\partial w_t L(w_t) - \lambda w_t + m \Delta w_{t-1}$$

Gradient of loss Weight decay Momentum

$$w_{t+1} = w_t + \eta \Delta w_t$$

NN updates with weight decay and momentum

$$\Delta w_t = -\partial/\partial w_t L(w_t) - \lambda w_t + m \Delta w_{t-1}$$

Gradient of loss Weight decay Momentum

$$w_{t+1} = w_t + n \Delta w_t$$

Learning rate

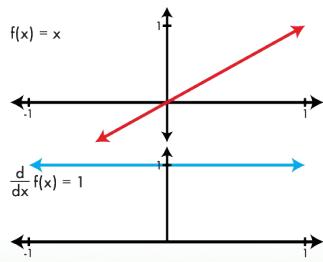
What about our activation functions φ

Many options, want them to be easy to take derivative

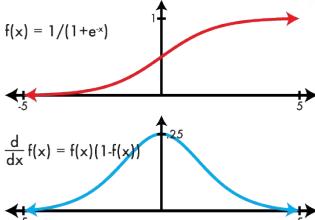
UAT holds when bounded, in practice bounds can be problematic

Common activation functions φ

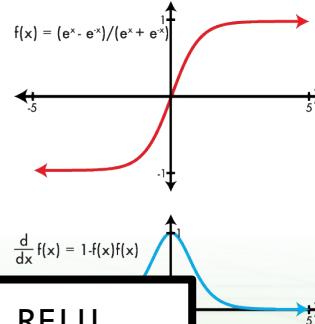
linear



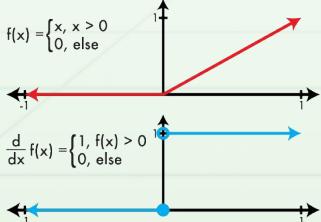
logistic



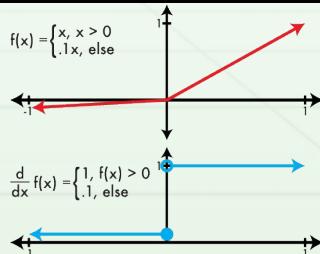
tanh



REctified Linear Unit (RELU)



Leaky RELU



So many hyper parameters!!

How do we know what to use??

Hyper Parameter Dark Magic

What follows are the one, true, correct, and only set of hyperparameters.

Praise be the NetLord!

$n = [.0001 - .01]$

$\lambda = .0005$

$m = .9$

$\phi = \text{leaky relu}$