CMSC 472/672

Lecture 9

Neural Network Optimization





CMSC 472/672



Objective

Data

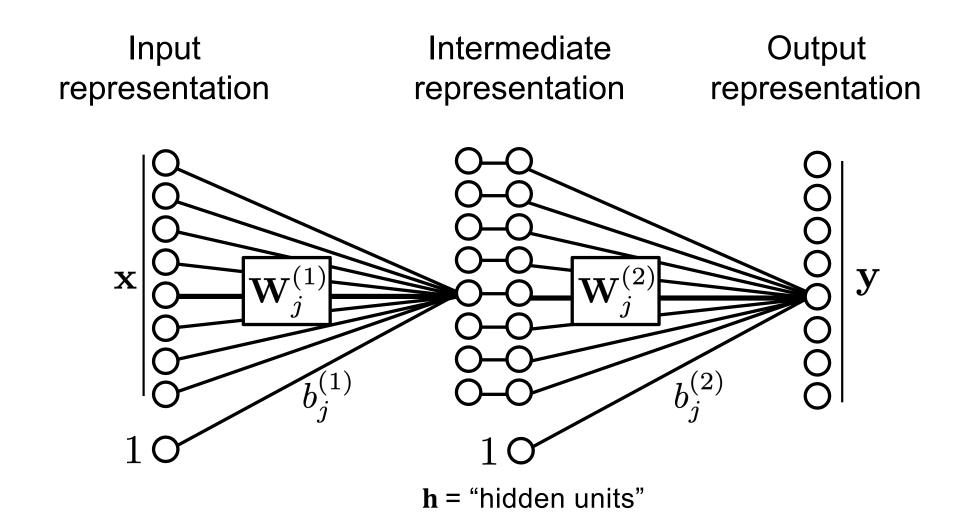
Hypothesis space

Optimizer

 \uparrow

Compute

Stacking layers

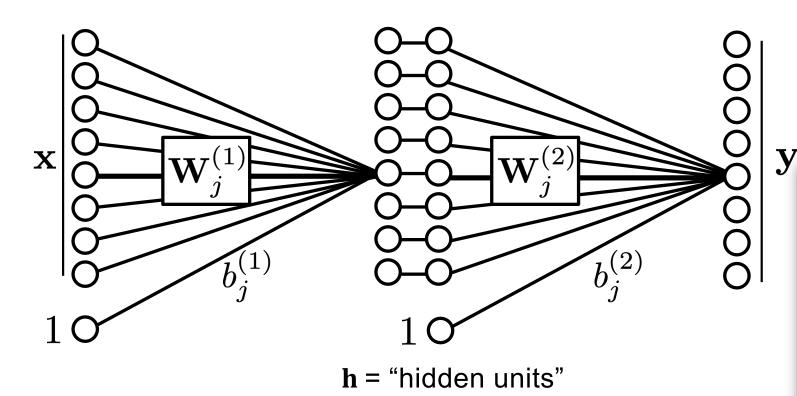


Stacking layers

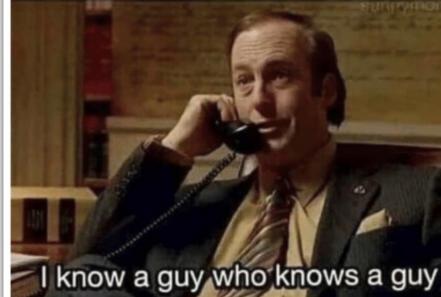
Input representation

Intermediate representation

Output representation



How Neural Networks work? Neurons:



An Incremental Learning Strategy

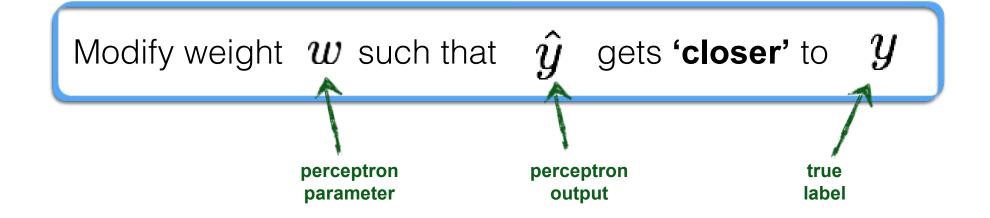
(gradient descent)

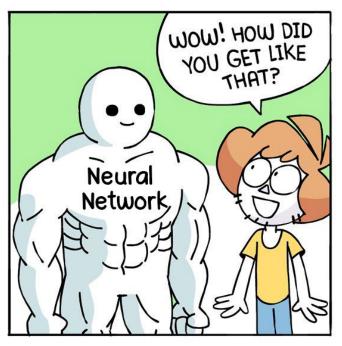
Given several examples

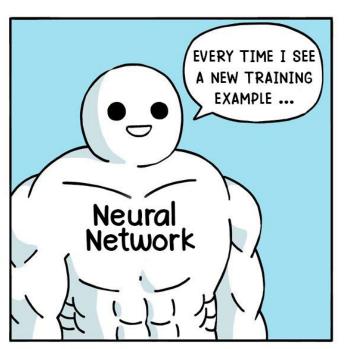
$$\{(x_1,y_1),(x_2,y_2),\ldots,(x_N,y_N)\}$$

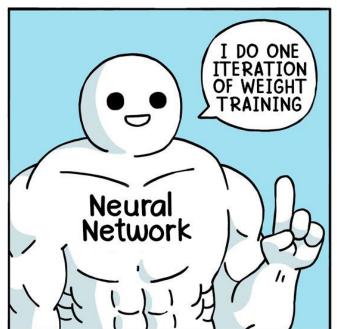
and a perceptron

$$\hat{y} = wx$$



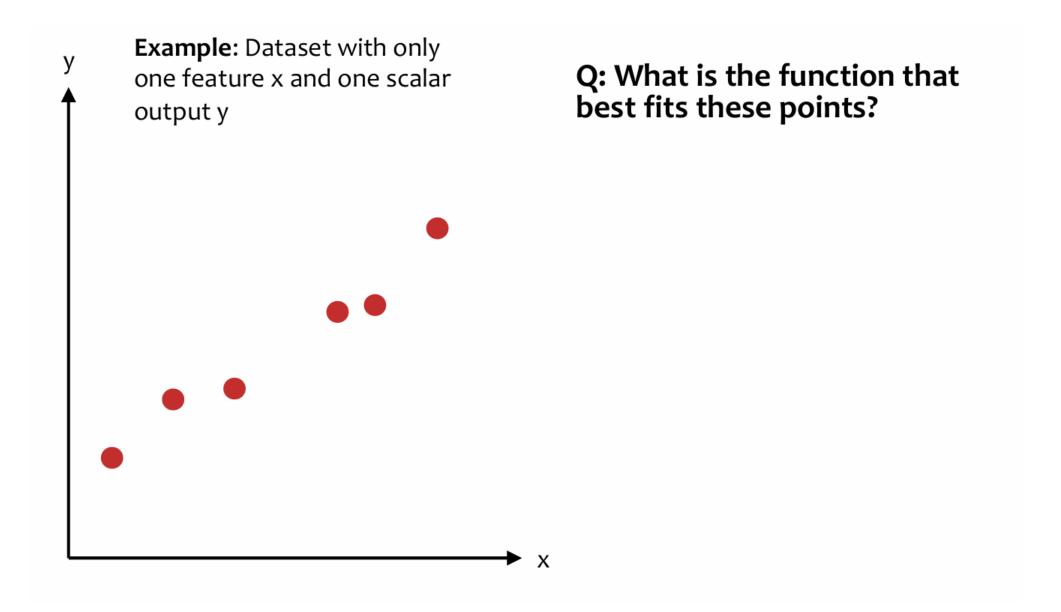




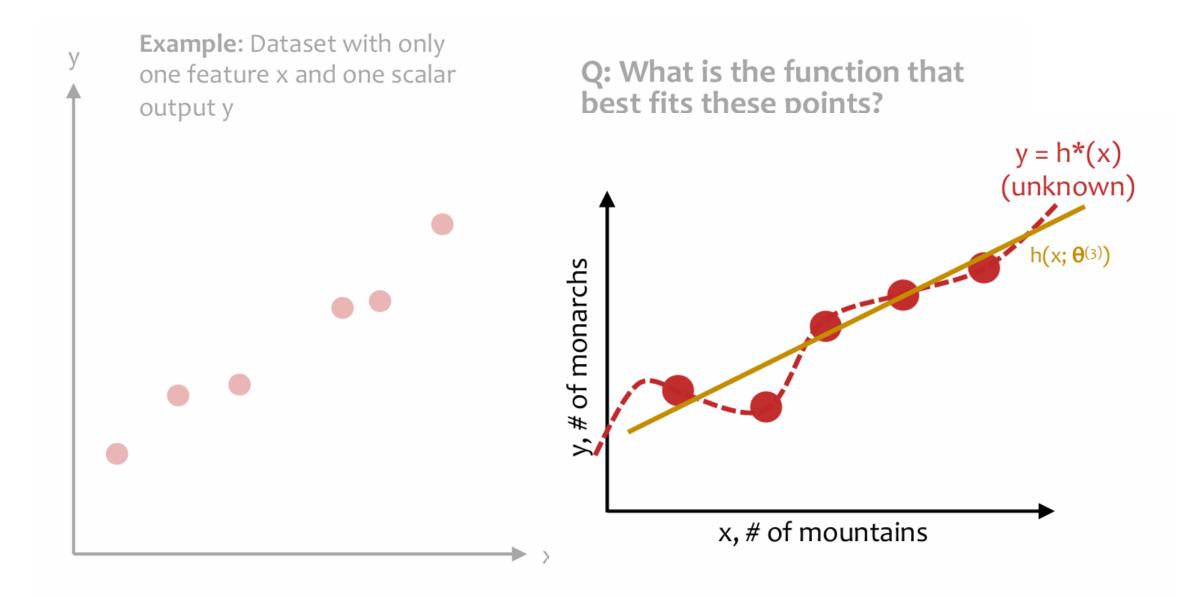




Recap: Linear Regression



Recap: Linear Regression



Naïve Idea:

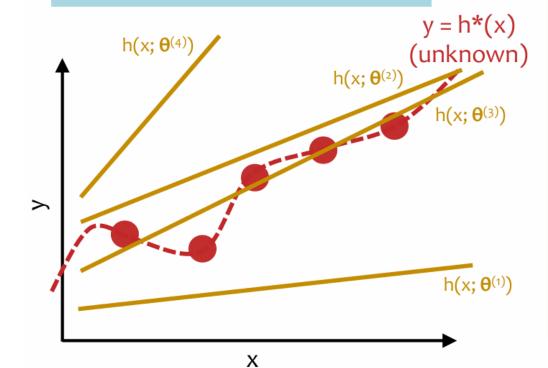
Linear Regression by Random Guessing

!!!

Linear Regression by Rand. Guessing

Optimization Method #0: Random Guessing

- 1. Pick a random θ
- 2. Evaluate $J(\theta)$
- Repeat steps 1 and 2 many times
- 4. Return θ that gives smallest $J(\theta)$



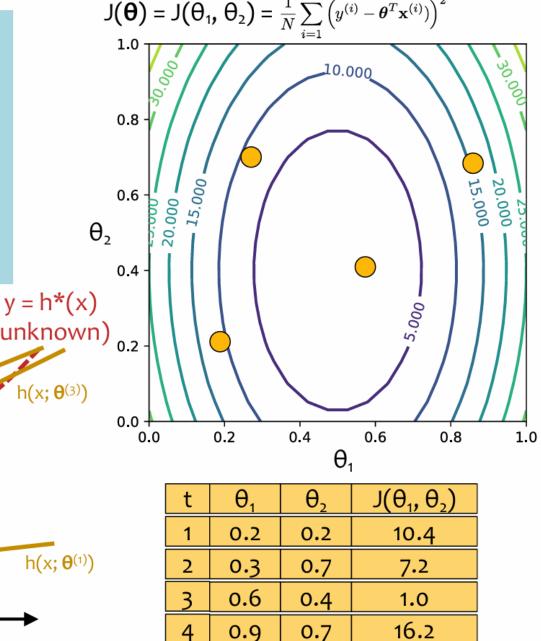
For Linear Regression:

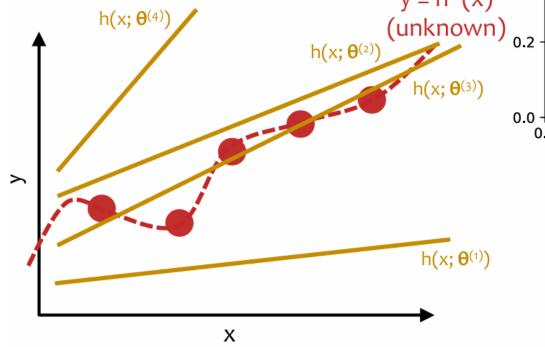
- target function h*(x) is unknown
- only have access to h*(x) through training examples (x⁽ⁱ⁾,y⁽ⁱ⁾)
- want h(x; θ^(t)) that best approximates h*(x)
- enable generalization w/inductive bias that restricts hypothesis class to linear functions

Linear Regression by Rand. Guessing $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - \theta^T \mathbf{x}^{(i)} \right)^2$

Optimization Method #0: Random Guessing

- 1. Pick a random θ
- 2. Evaluate $J(\theta)$
- 3. Repeat steps 1 and 2 many times
- 4. Return θ that gives smallest $J(\theta)$





Better Idea:

An optimization algorithm called

"Gradient Descent"

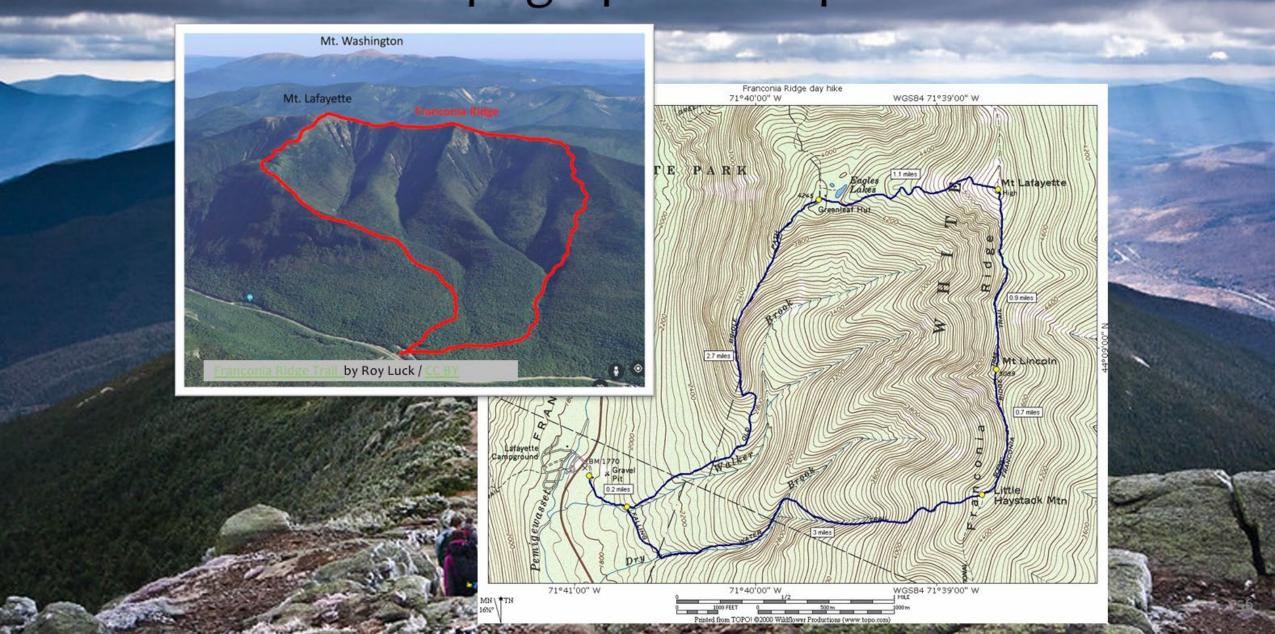
Recall: Gradient of a vector

Def: The gradient of $J:\mathbb{R}^M \to \mathbb{R}$ is

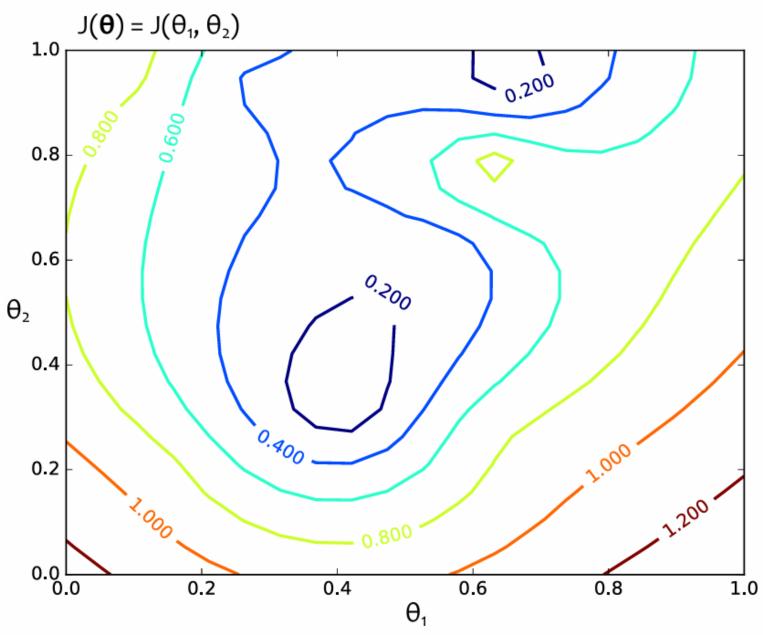
$$abla J(oldsymbol{ heta}) = egin{bmatrix} rac{\partial J(oldsymbol{ heta})}{\partial heta_1} rac{\partial J(oldsymbol{ heta})}{\partial heta_2} \ dots \ rac{\partial J(oldsymbol{ heta})}{\partial heta_M} \end{bmatrix}$$

Each entry is a first-order partial derivative

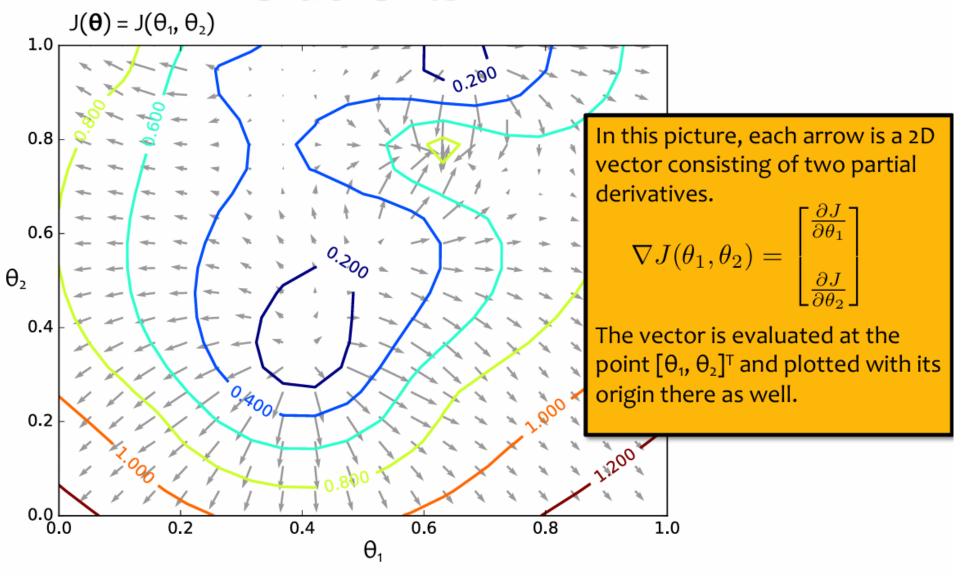
Topographical Maps



Gradients

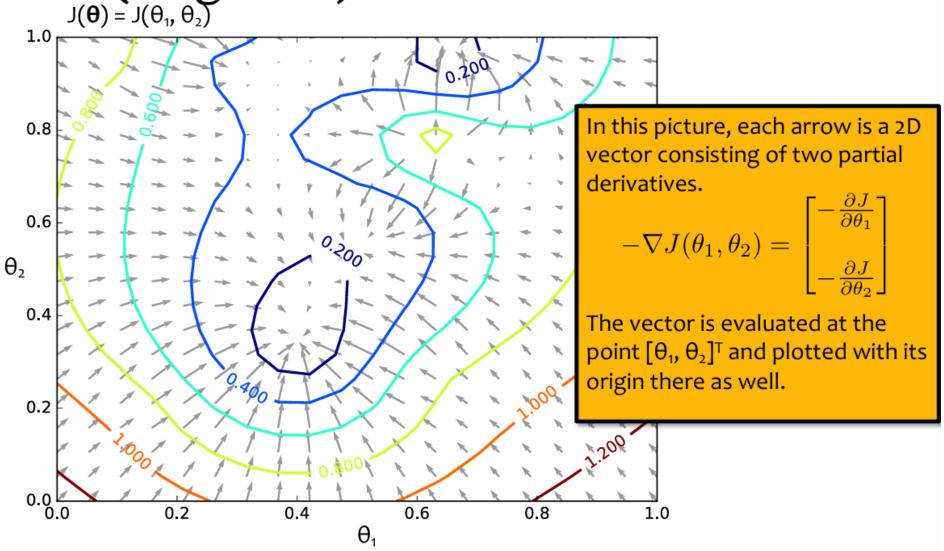


Gradients



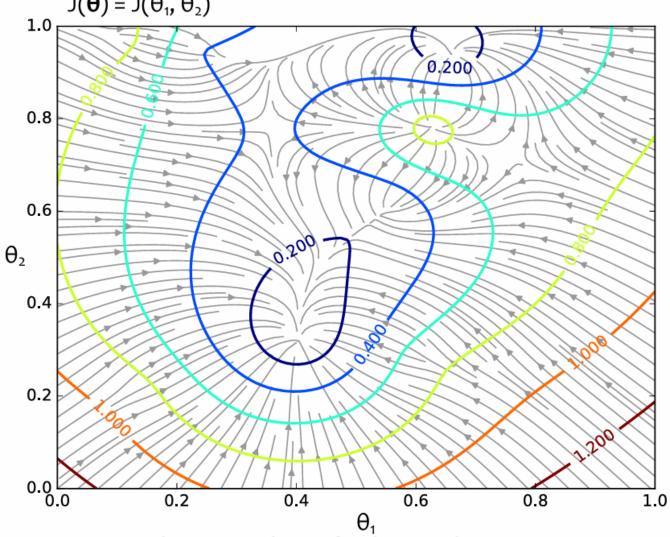
These are the **gradients** that Gradient **Ascent** would follow.

(Negative) Gradients $J(\theta) = J(\theta_1, \theta_2)$



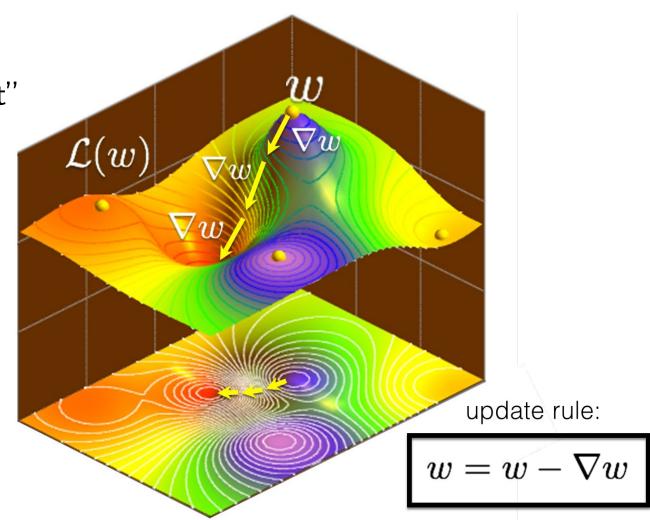
These are the **negative** gradients that Gradient **Descent** would follow.

(Negative) Gradient Paths $\int_{J(\theta)=J(\theta_1,\theta_2)}^{J(\theta)=J(\theta_1,\theta_2)}$



Shown are the **paths** that Gradient Descent would follow if it were making **infinitesimally small steps**.

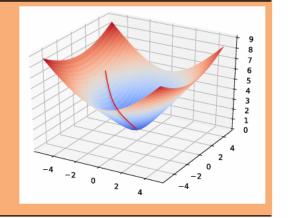
Go down the path of "steepest descent"



• Go down the path of steepest descent

Algorithm 1 Gradient Descent

- 1: **procedure** $\mathrm{GD}(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
- 2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
- 3: while not converged do
- 4: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- 5: return θ



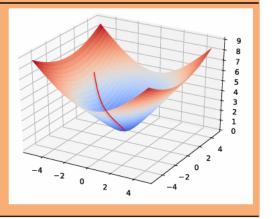
In order to apply GD to Linear Regression all we need is the gradient of the objective function (i.e. vector of partial derivatives).

$$abla_{m{ heta}} J(m{ heta}) = egin{bmatrix} rac{d heta_1}{d} J(m{ heta}) \ rac{d}{d heta_2} J(m{ heta}) \ dots \ rac{d}{d heta_M} J(m{ heta}) \end{bmatrix}$$

Go down the path of steepest descent

Algorithm 1 Gradient Descent

- 1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
- 2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
- 3: **while** not converged **do**
- 4: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- 5: return θ



There are many possible ways to detect **convergence**. For example, we could check whether the L2 norm of the gradient is below some small tolerance.

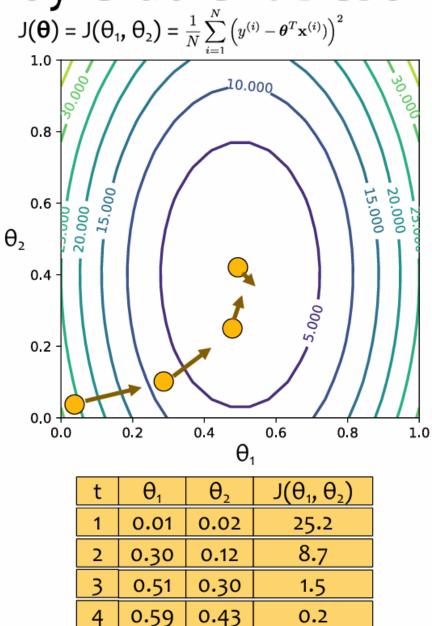
$$||\nabla_{\boldsymbol{\theta}}J(\boldsymbol{\theta})||_2 \leq \epsilon$$

Alternatively we could check that the reduction in the objective function from one iteration to the next is small.

Linear Regression by Gradient Desc. $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

Optimization Method #1: Gradient Descent

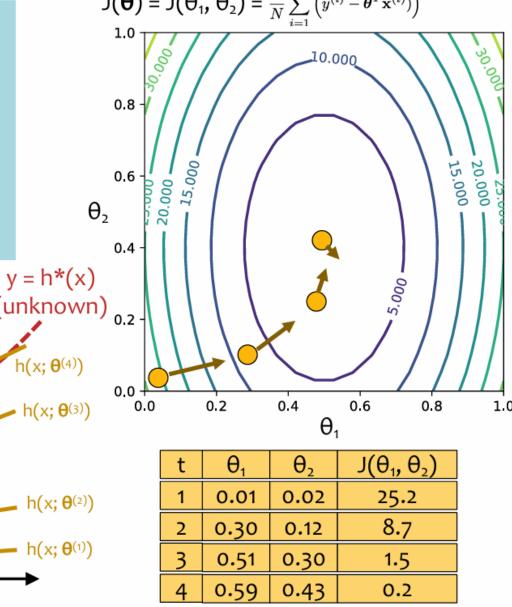
- 1. Pick a random $\boldsymbol{\theta}$
- 2. Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- 3. Return θ that gives smallest $J(\theta)$

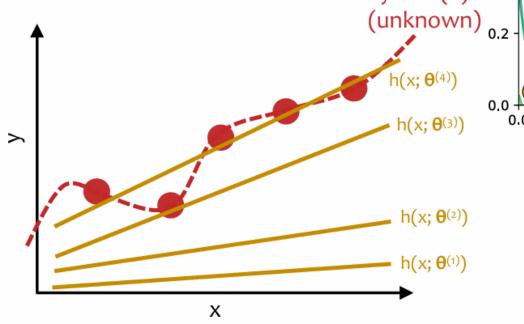


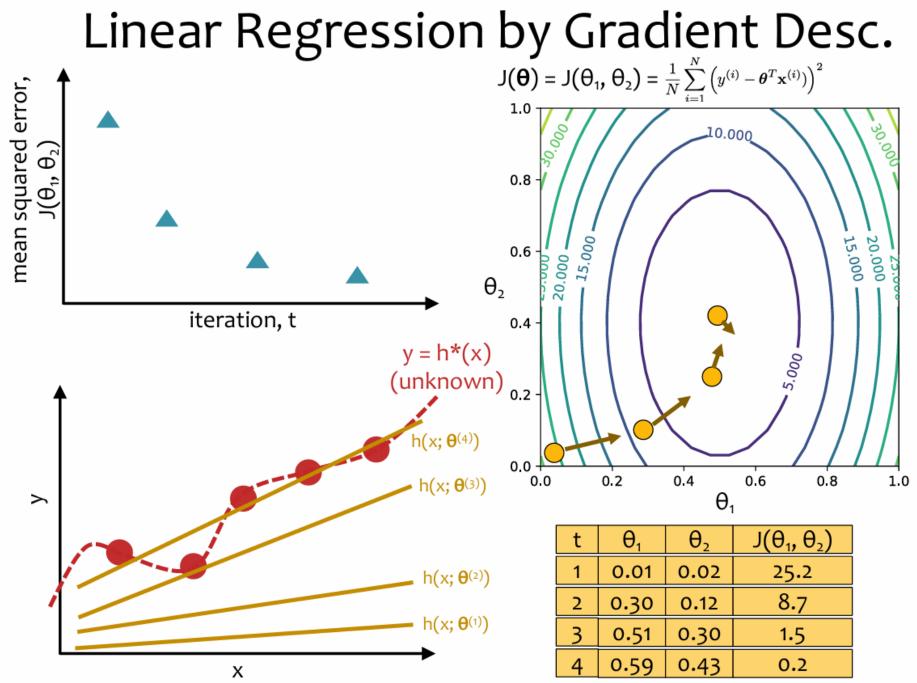
Linear Regression by Gradient Desc. $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

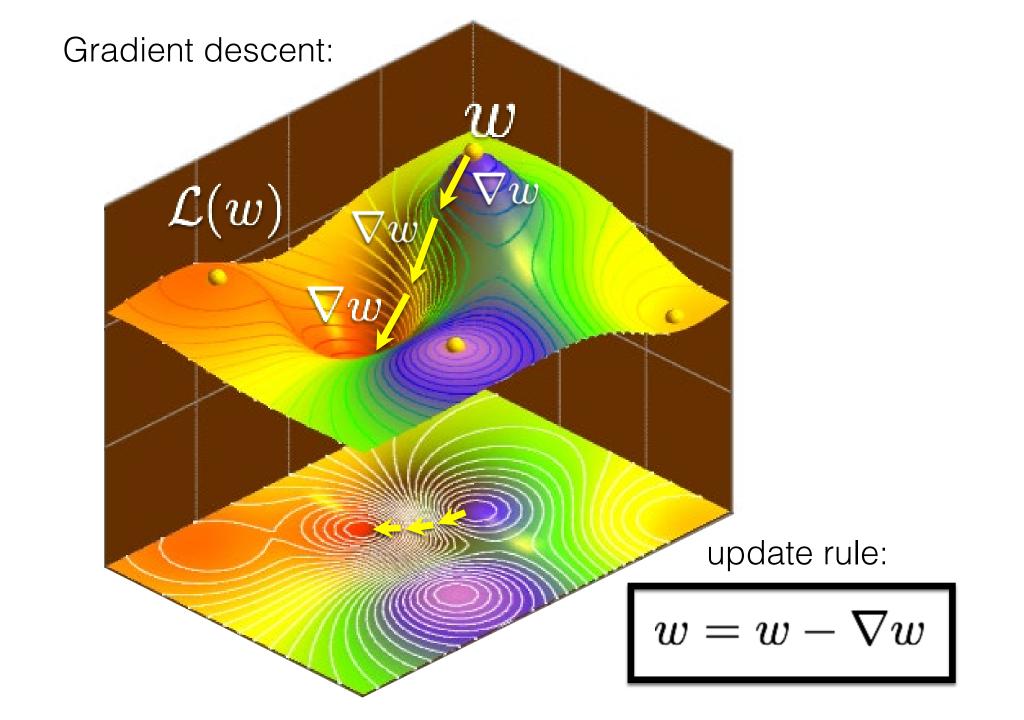
Optimization Method #1: Gradient Descent

- 1. Pick a random θ
- 2. Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- 3. Return θ that gives smallest $J(\theta)$









For each example sample $\{x_i,y_i\}$

1. Predict

$$rac{\partial \mathcal{L}}{\partial heta}$$

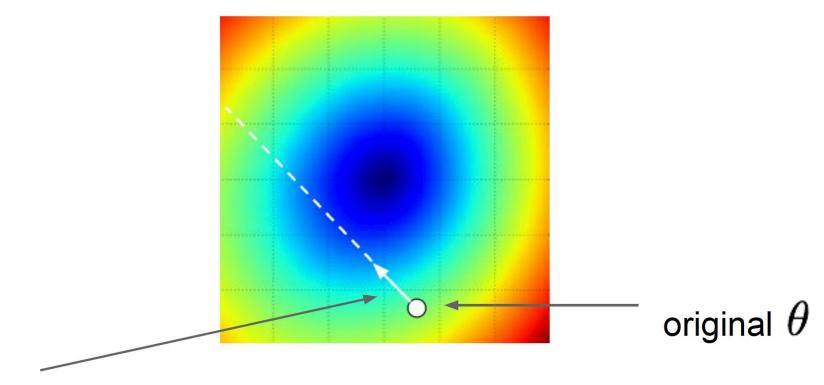
 \mathcal{L}_i

vector of parameter partial derivatives

 $\hat{y} = f_{\text{MLP}}(x_i; \theta)$

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

Learning rates



negative gradient direction

$$heta \leftarrow heta$$
 - $\eta rac{\partial \mathcal{L}}{\partial heta}$

Step size: learning rate

Too big: will miss the minimum

Too small: slow convergence

ac - Mr. as

Late - Ave - TC

et - - Profit



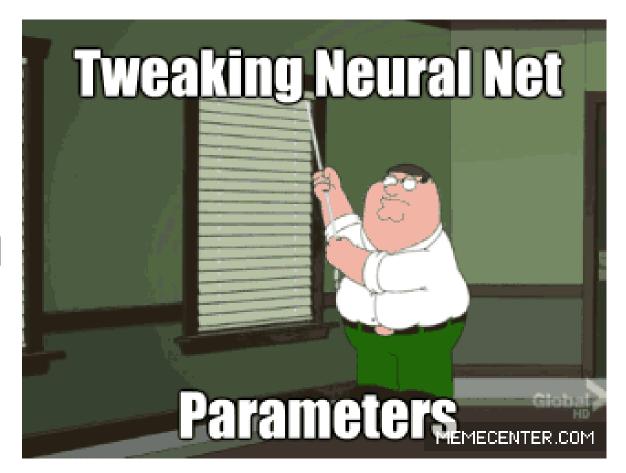


GD for LR: Python Step-by-Step Example

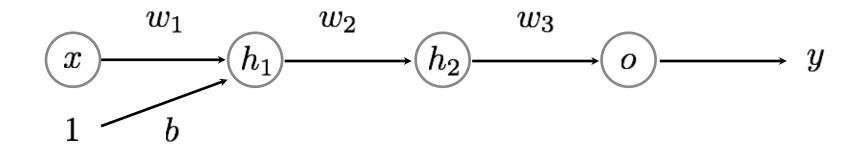
 https://colab.research.google.com/drive /17dK6cynECzk2ObyCqDk5gKcUyN1k MjSR?usp=sharing



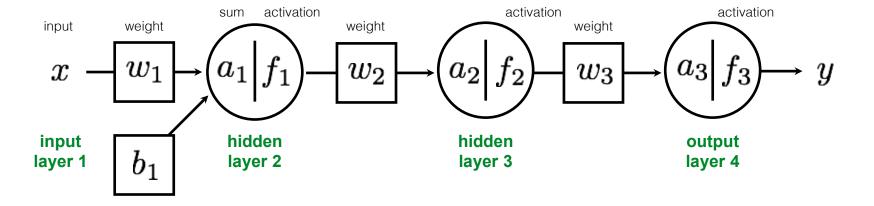
Backpropagation

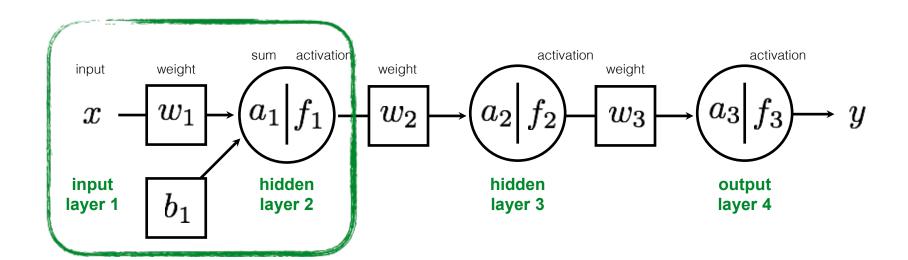


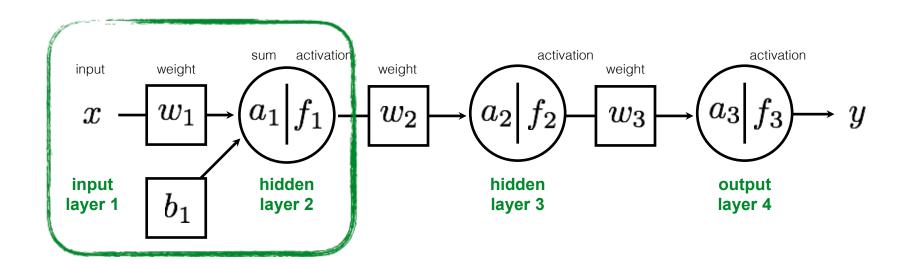
Example of a multi-layer perceptron



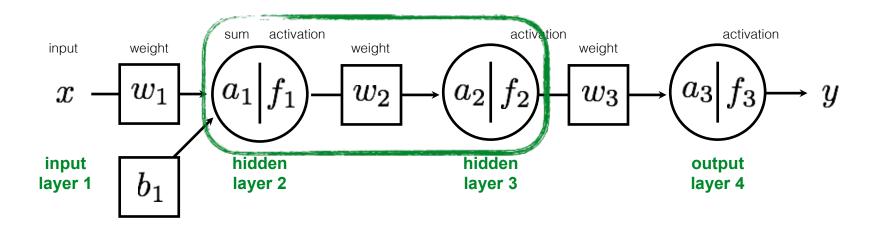
function of FOUR parameters and FOUR layers!



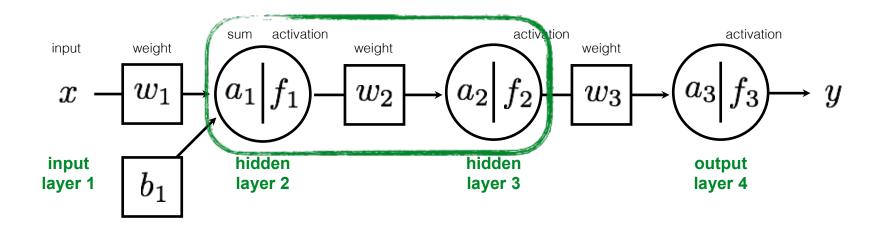




$$a_1 = w_1 \cdot x + b_1$$

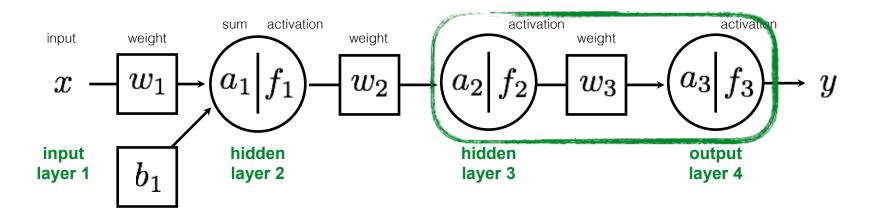


$$a_1 = w_1 \cdot x + b_1$$



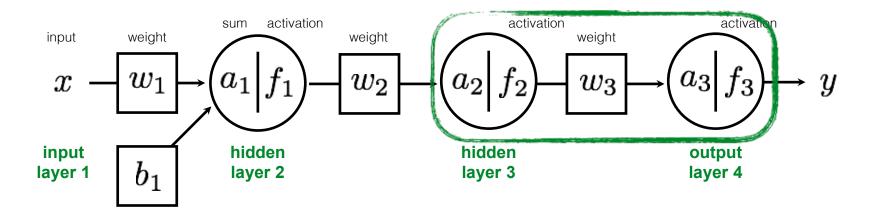
$$a_1 = w_1 \cdot x + b_1$$

 $a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$



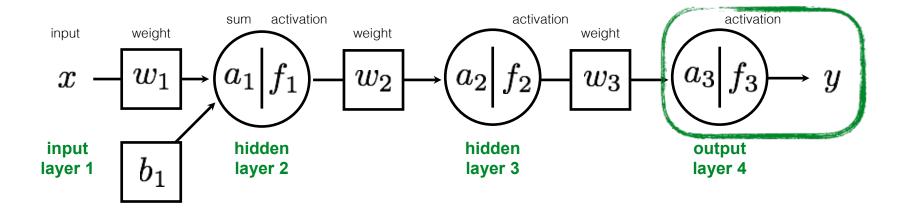
$$a_1 = w_1 \cdot x + b_1$$

 $a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$



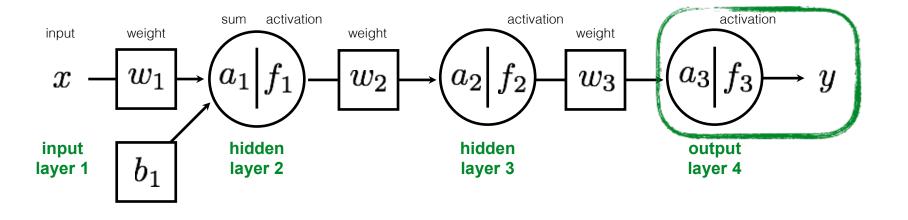
$$a_1 = w_1 \cdot x + b_1$$

 $a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$
 $a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$



$$a_1 = w_1 \cdot x + b_1$$

 $a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$
 $a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$



$$a_1 = w_1 \cdot x + b_1$$

 $a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$
 $a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$
 $y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

We need to train the network:

What is known? What is unknown?

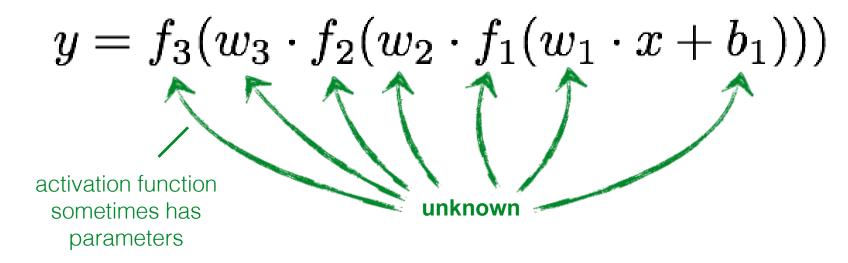
Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$
known

We need to train the network:

What is known? What is unknown?

Entire network can be written out as one long equation



We need to train the network:

What is known? What is unknown?

Learning an MLP

Given a set of samples and a MLP

$$\{x_i, y_i\}$$
$$y = f_{\text{MLP}}(x; \theta)$$

Estimate the parameters of the MLP

$$\theta = \{f, w, b\}$$

Gradient Descent

For each ${f random}$ sample $\{x_i,y_i\}$

- 1. Predict
 - a. Forward pass

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

- b. Compute Loss
- 2. Update
 - a. Back Propagation
 - b. Gradient update

$$rac{\partial \mathcal{L}}{\partial heta}$$

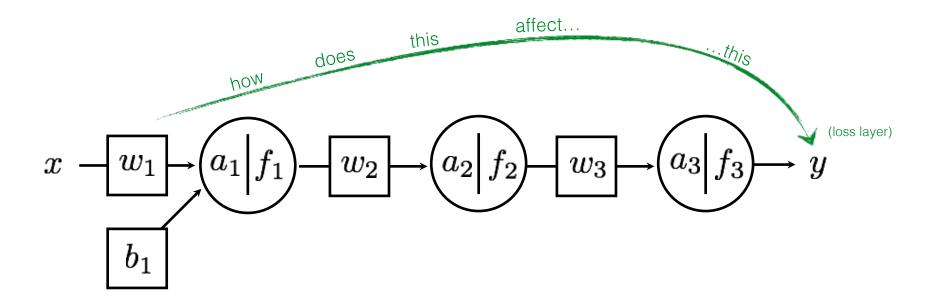
vector of parameter partial derivatives

 $\theta \leftarrow \theta - \eta \nabla \theta$

So we need to compute the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \left[\frac{\partial \mathcal{L}}{\partial w_3} \frac{\partial \mathcal{L}}{\partial w_2} \frac{\partial \mathcal{L}}{\partial w_1} \frac{\partial \mathcal{L}}{\partial b} \right]$$

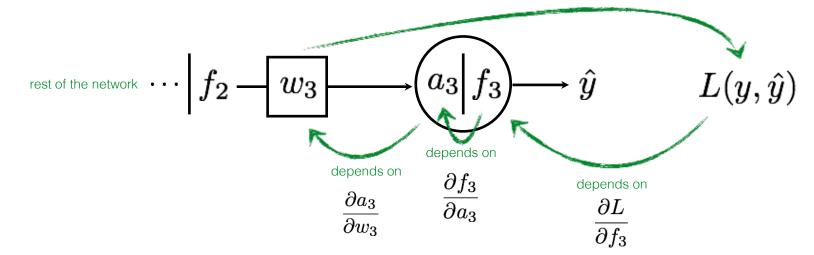
Remember, $\frac{\partial L}{\partial w_1} \ \ \text{describes...}$

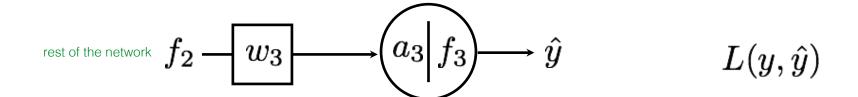


According to the chain rule...

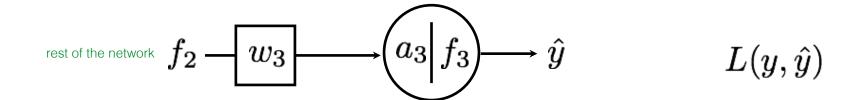
$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Intuitively, the effect of weight on loss function : $\frac{\partial L}{\partial w_3}$



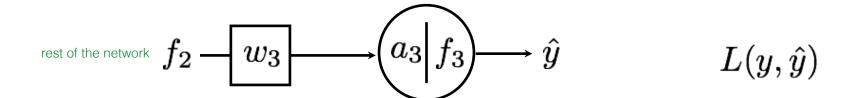


$$rac{\partial L}{\partial w_3} = rac{\partial L}{\partial f_3} rac{\partial f_3}{\partial a_3} rac{\partial a_3}{\partial w_3}$$
 Chain Rule!



$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$
$$= -\eta (y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

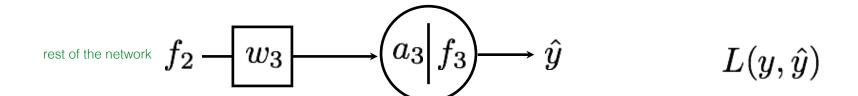
Just the partial derivative of L2 loss



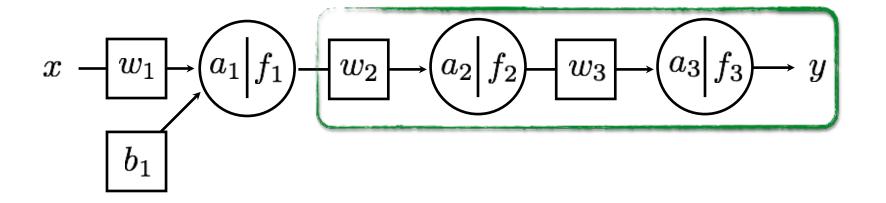
$$\begin{split} \frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ &= -\eta (y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \end{split}$$

Let's use a Sigmoid function

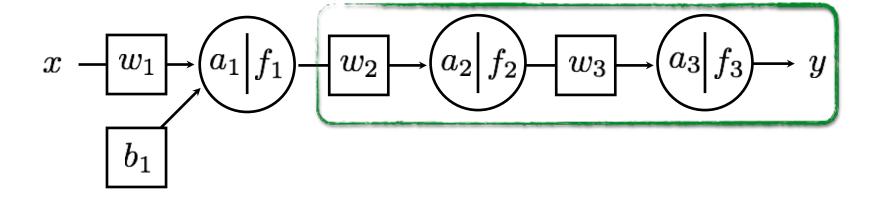
$$\frac{ds(x)}{dx} = s(x)(1 - s(x))$$



$$\begin{split} \frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ &= -\eta (y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ &= -\eta (y - \hat{y}) f_3 (1 - f_3) \frac{\partial a_3}{\partial w_3} \\ &= -\eta (y - \hat{y}) f_3 (1 - f_3) f_2 \end{split}$$



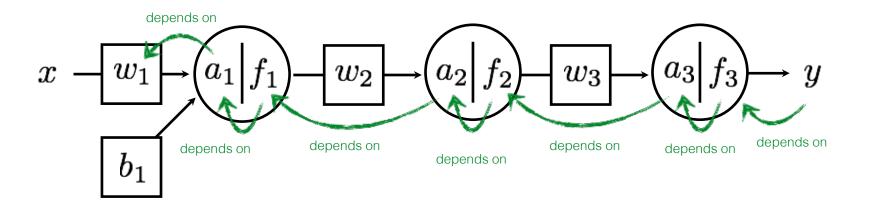
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$



$$\frac{\partial L}{\partial w_2} = \left[\frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \right]$$

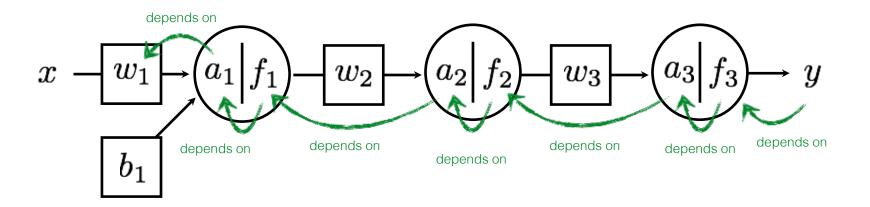
already computed. re-use (propagate)!

The chain rule says...



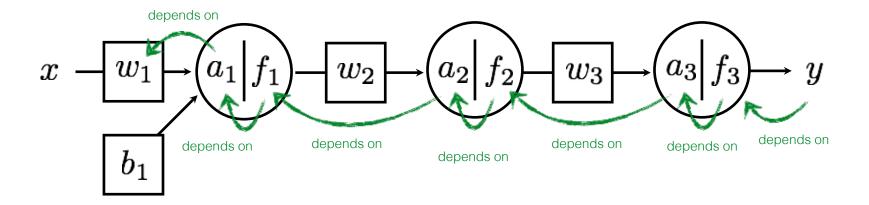
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

The chain rule says...

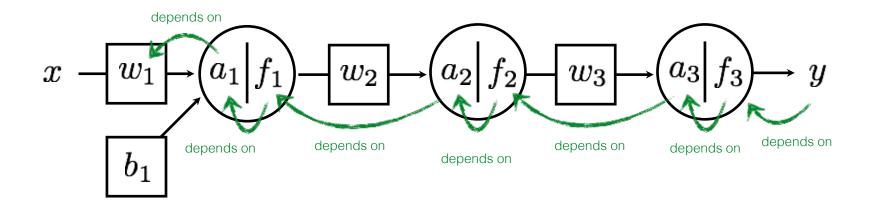


$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \frac{\partial a_1}{\partial w_1} \right)$$

already computed. re-use (propagate)!

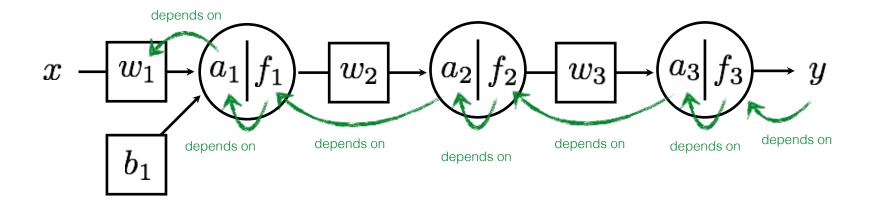


$$\frac{\partial \mathcal{L}}{\partial w_3} = \begin{bmatrix}
\frac{\partial \mathcal{L}}{\partial f_3} & \frac{\partial f_3}{\partial a_3} & \frac{\partial a_3}{\partial w_3} \\
\frac{\partial \mathcal{L}}{\partial w_2} & = \begin{bmatrix}
\frac{\partial \mathcal{L}}{\partial f_3} & \frac{\partial f_3}{\partial a_3} & \frac{\partial a_3}{\partial f_2} & \frac{\partial f_2}{\partial a_2} & \frac{\partial a_2}{\partial w_2} \\
\frac{\partial \mathcal{L}}{\partial w_1} & = \frac{\partial \mathcal{L}}{\partial f_3} & \frac{\partial f_3}{\partial a_3} & \frac{\partial f_2}{\partial f_2} & \frac{\partial f_2}{\partial a_2} & \frac{\partial f_1}{\partial a_1} & \frac{\partial f_1}{\partial w_1} \\
\frac{\partial \mathcal{L}}{\partial b} & = \frac{\partial \mathcal{L}}{\partial f_3} & \frac{\partial f_3}{\partial a_3} & \frac{\partial f_2}{\partial f_2} & \frac{\partial f_2}{\partial a_2} & \frac{\partial f_2}{\partial f_1} & \frac{\partial f_1}{\partial a_1} & \frac{\partial f_1}{\partial b}
\end{bmatrix}$$



$$\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \begin{bmatrix}
\frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\
\frac{\partial \mathcal{L}}{\partial w_1} = \begin{bmatrix}
\frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial f_2}{\partial a_2} \frac{\partial f_1}{\partial a_1} \frac{\partial f_1}{\partial a_1} \frac{\partial f_2}{\partial b_1} \\
\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial f_2}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial f_1}{\partial b_1}$$



$$\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial f_2}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial f_2}{\partial f_3} \frac{\partial f_1}{\partial a_1} \frac{\partial f_1}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial f_2}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial f_1}{\partial b}$$

Gradient Descent

For each example sample $\{x_i,y_i\}$

1. Predict

a. Back Propagation

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

$$\mathcal{L}_i$$

$$\begin{split} \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b} \end{split}$$

$$egin{aligned} w_3 &= w_3 - \eta
abla w_3 \ w_2 &= w_2 - \eta
abla w_2 \ w_1 &= w_1 - \eta
abla w_1 \ b &= b - \eta
abla b \end{aligned}$$

Gradient Descent

For each example sample
$$\{x_i,y_i\}$$

1. Predict

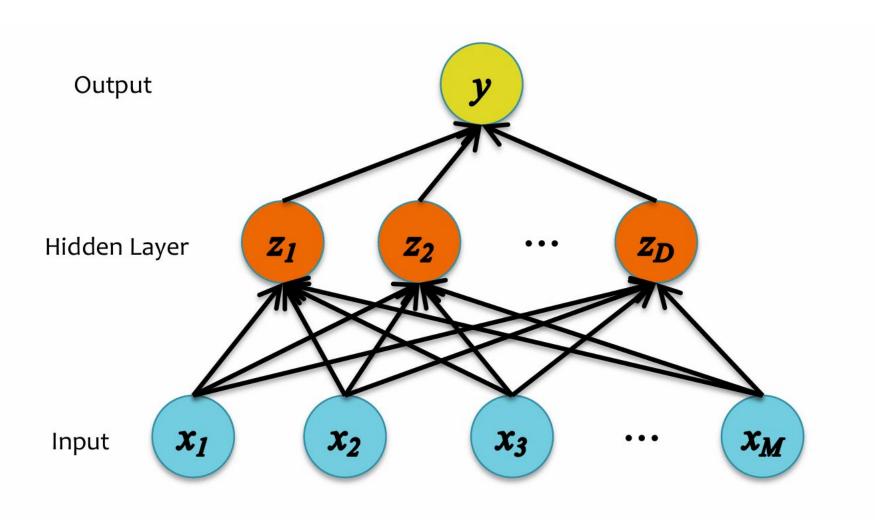
$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

$$\mathcal{L}_i$$

 $\frac{\partial \mathcal{L}}{\partial \theta}$

$$heta \leftarrow heta - \eta rac{\partial \mathcal{L}}{\partial heta}$$
 vector of update equations

Single Output Neural Networl Let's write the equation



Objective Functions for NNs

Quadratic Loss:

- the same objective as Linear Regression
- i.e. mean squared error

$$J = \ell_Q(y, y^{(i)}) = \frac{1}{2}(y - y^{(i)})^2$$
$$\frac{dJ}{dy} = y - y^{(i)}$$

2. Binary Cross-Entropy:

- the same objective as Binary Logistic Regression
- i.e. negative log likelihood
- This requires our output y to be a probability in [0,1]

$$J = \ell_{CE}(y, y^{(i)}) = -(y^{(i)} \log(y) + (1 - y^{(i)}) \log(1 - y))$$
$$\frac{dJ}{dy} = -\left(y^{(i)} \frac{1}{y} + (1 - y^{(i)}) \frac{1}{y - 1}\right)$$

Objective Functions for NNs

Cross-entropy vs. Quadratic loss

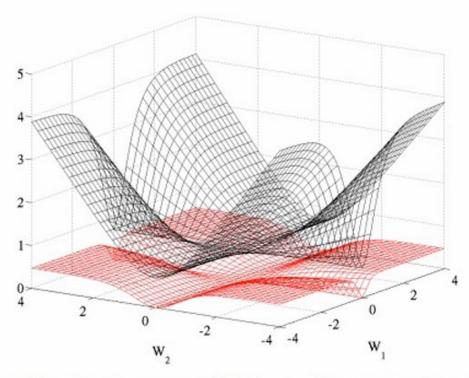
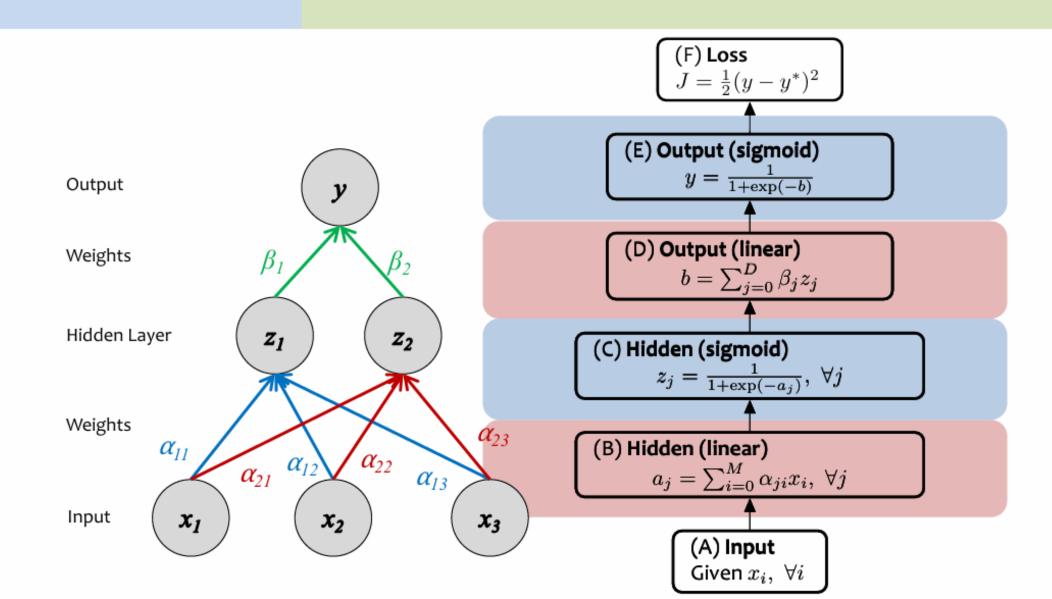


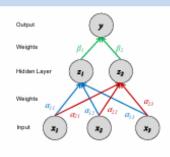
Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, W_1 respectively on the first layer and W_2 on the second, output layer.

Forward-Computation



Case 2: Neural Network

Backpropagation



Forward
$J = y^* \log y + (1 - y^*) \log(1 - y)$
$y = \frac{1}{y}$

$$b = \sum_{j=0}^{D} \beta_j z_j$$

Loss

Sigmoid

Linear

$$z_j = \frac{1}{1 + \exp(-a_j)}$$

$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i$$

$$g_y = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

$$g_b = g_y \frac{\partial y}{\partial b}, \ \frac{\partial y}{\partial b} = \frac{\exp(-b)}{(\exp(-b) + 1)^2}$$

$$g_{\beta_j} = g_b \frac{\partial b}{\partial \beta_j}, \ \frac{\partial b}{\partial \beta_j} = z_j$$

$$g_{z_j} = g_b \frac{\partial b}{\partial z_j}, \ \frac{\partial b}{\partial z_j} = \beta_j$$

$$g_{a_j} = g_{z_j} \frac{\partial z_j}{\partial a_j}, \ \frac{\partial z_j}{\partial a_j} = \frac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$$

$$g_{\alpha_{ji}} = g_{a_j} \frac{\partial a_j}{\partial \alpha_{ji}}, \ \frac{\partial a_j}{\partial \alpha_{ji}} = x_i$$

$$g_{x_i} = \sum_{j=0}^{D} g_{a_j} \frac{\partial a_j}{\partial x_i}, \ \frac{\partial a_j}{\partial x_i} = \alpha_{ji}$$

some history / controversy ...

Nature **volume 323**, pages533–536 (1986)

Nobel Prize in Physics 2024

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton† & Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California, San Diego, La Jolla, California 92093, USA † Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Philadelphia 15213, USA



- Precursors to backpropagation appeared in optimal control theory since 1950s.
 - LeCun et al 1985 credits 1950s work by Pontryagin.
- Modern backpropagation was first published by Seppo Linnainmaa as "reverse mode of automatic differentiation" (1970) for discrete connected networks of nested differentiable functions.
- Rumelhart "independently developed" backpropagation. He did not cite previous work as he was "unaware" of them.

https://people.idsia.ch/~juergen/who-invented-backpropagation.html



Why are NNs such good function approximators?

Universal Approximation Theorem:

"A neural network $\hat{f}(x)$ with a single hidden-layer can approximate any continuous real function f(x) to within any arbitrary degree of accuracy ϵ given a sufficient number of neurons in the hidden layer"

Nice Explanation:

https://www.deep-mind.org/2023/03/26/the-universal-approximation-theorem/

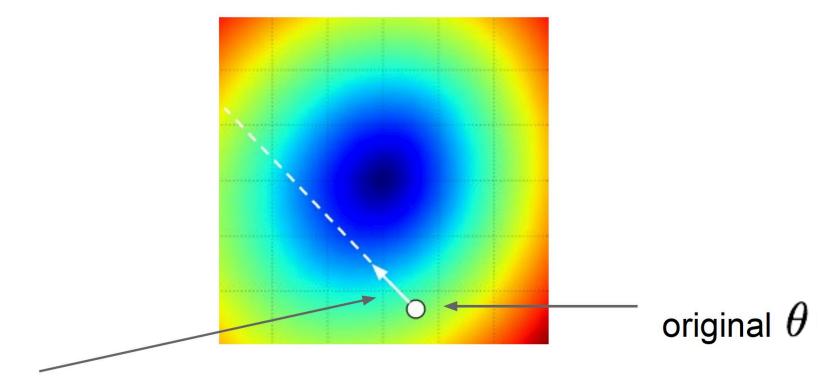
Why are NNs such good function approximators?

Universal Approximation Theorem (UAT) limitations:

- 1. Does not specify what the network size should be
- 2. Only guarantees *function approximation* on the given training data; does not guarantee generalization on new test data
- Only assures that approximation exists; does not provide insights on how to train the network to achieve that optimal function

"A neural network $\hat{f}(x)$ with a single hidden-layer can approximate any continuous real function f(x) to within any arbitrary degree of accuracy ϵ given a sufficient number of neurons in the hidden layer"

Learning rates



negative gradient direction

$$heta \leftarrow heta - \eta rac{\partial \mathcal{L}}{\partial heta}$$

Step size: learning rate

Too big: will miss the minimum

Too small: slow convergence

Learning rate scheduling

- Use different learning rate at each iteration.
- Most common choice:

$$\eta_t = \frac{\eta_0}{\sqrt{t}}$$

Need to select initial learning rate η_0

More modern choice: Adaptive learning rates.

$$\eta_t = G\left(\left\{\frac{\partial L}{\partial \theta}\right\}_{i=0}^t\right)$$

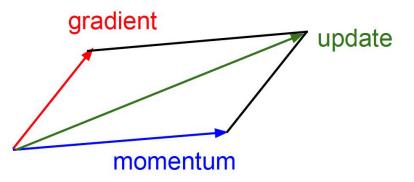
Many choices for G (Adam, Adagrad, Adadelta).

Momentum Update

$$heta \leftarrow heta$$
 - $\eta rac{\partial \mathcal{L}}{\partial heta}$

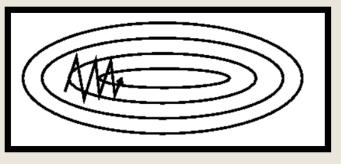
$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

$$\Delta \theta \leftarrow w \frac{\partial L}{\partial \theta} + (1 - w) \Delta \theta$$

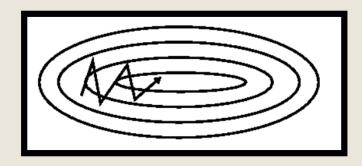


Take direction history into account!

```
weights grad = evaluate gradient(loss fun, data, weights)
vel = vel * 0.9 - step size * weights grad
weights += vel
```



(Fig. 2a)



Many other ways to perform optimization...

- Second order methods that use the Hessian (or its approximation): BFGS, **LBFGS**, etc.
- Currently, the lesson from the trenches is that well-tuned SGD+Momentum is very hard to beat for CNNs.
- No consensus on Adam etc.: Seem to give faster performance to worse local minima.