

Neural networks: an introduction

UE de Master 2, AOS2

Fall 2023

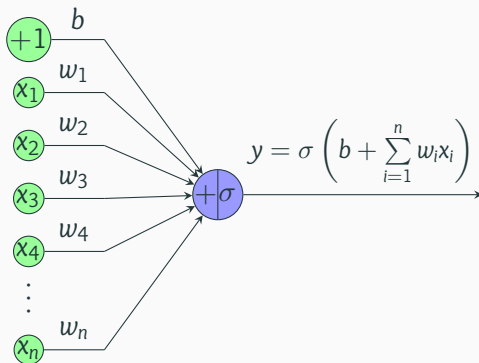
S. Rousseau

Feed forward neural network

What is a neuron?

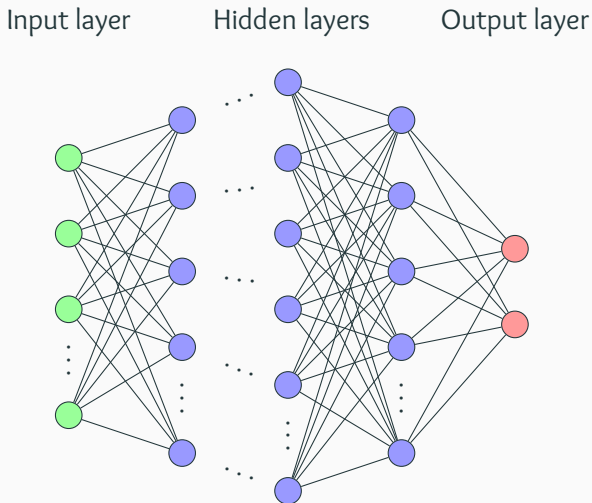
Neuron = linear transform of the input followed by a non-linearity

- Input: $\mathbf{x} = (x_i)_{i=1\dots n} \in \mathbb{R}^n$
- Scalar output: $y \in \mathbb{R}$
- +1 denotes the intercept (or bias) b
- The w_i 's are called **weights**
- σ is nonlinear function called an **activation function**



Feed forward neural network

- Stacked collection of neurons arranged in layers
- Input layer nodes are not neurons
- Output layer neurons do not have necessarily an activation function



Input layer and first hidden layer

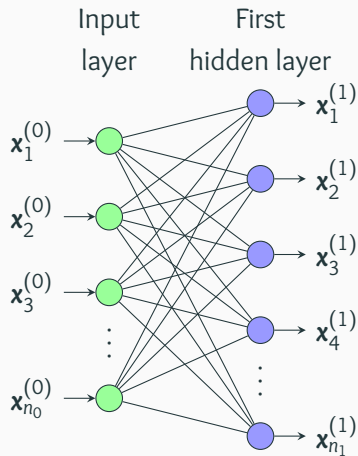
- Input: $\mathbf{x}^{(0)} = \left(\mathbf{x}_i^{(0)} \right)_{i=1 \dots n_0} \in \mathbb{R}^{n_0}$
- Output: $\mathbf{x}^{(1)} = \left(\mathbf{x}_i^{(1)} \right)_{i=1 \dots n_1} \in \mathbb{R}^{n_1}$

We have

$$\mathbf{x}_i^{(1)} = \sigma \left(\sum_{j=1}^{n_0} \mathbf{w}_{ij}^{(1)} \mathbf{x}_j^{(0)} + \mathbf{b}_i^{(1)} \right)$$

where

- $\mathbf{b}_i^{(1)}$ is the bias of neuron i (not displayed)
- $\mathbf{w}_{ij}^{(1)}$ is the j -th coefficient of i -th neuron
- Vector version $\mathbf{x}_i^{(1)} = \sigma \left(\left\langle \mathbf{w}_i^{(1)}, \mathbf{x}^{(0)} \right\rangle + \mathbf{b}_i^{(1)} \right)$

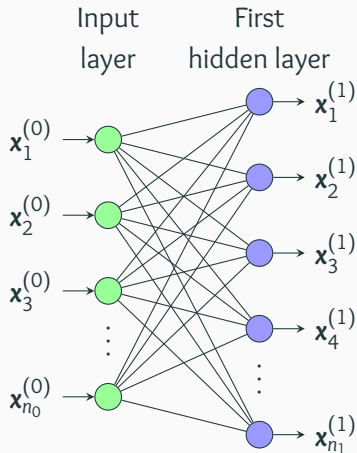


Input layer and first hidden layer

- Matrix version:

$$\mathbf{x}^{(1)} = \sigma \left(\left(\mathbf{w}^{(1)} \right)^T \mathbf{x}^{(0)} + \mathbf{b}^{(1)} \right)$$

- σ is applied element-wise
- $\mathbf{w}^{(1)}$ is of size $n_0 \times n_1$
- $\mathbf{w}_i^{(1)} \in \mathbb{R}^{n_0}$ columns of $\mathbf{w}^{(1)}$
- $\mathbf{b}^{(1)} \in \mathbb{R}^{n_1}$ groups all the biases
- The parameters are $\Theta^{(1)} = \{\mathbf{w}^{(1)}, \mathbf{b}^{(1)}\}$



Two hidden layers

- n_k is the number of neurons at layer k

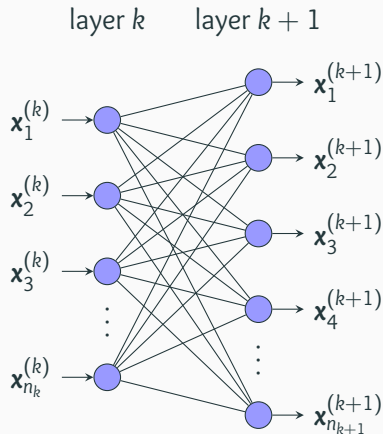
$$\mathbf{x}_i^{(k+1)} = \sigma \left(\sum_{j=1}^{n_k} \mathbf{w}_{ij}^{(k+1)} \mathbf{x}_j^{(k)} + \mathbf{b}_i^{(k+1)} \right)$$

- Condensed version

$$\mathbf{x}^{(k+1)} = \sigma \left(\left(\mathbf{w}^{(k+1)} \right)^T \mathbf{x}^{(k)} + \mathbf{b}^{(k+1)} \right)$$

- The parameters at layer $k + 1$ are

$$\Theta^{(k+1)} = \{ \mathbf{w}^{(k+1)}, \mathbf{b}^{(k+1)} \}$$



All layers

- Transformation at layer k :

$$F_{\Theta_k}(\mathbf{x}) = \sigma\left(\left(\mathbf{w}^{(k)}\right)^T \mathbf{x} + \mathbf{b}^{(k)}\right)$$

- Suppose the neural network has K layers (input excluded, output included)

$$\mathbf{x}^{(K)} = F_{\Theta_K}(\mathbf{x}^{(K-1)})$$

\vdots

$$\mathbf{x}^{(K)} = F_{\Theta_K}\left(F_{\Theta_{K-1}}\left(\dots F_{\Theta_1}\left(\mathbf{x}^{(0)}\right)\right)\right) = F_{\Theta}\left(\mathbf{x}^{(0)}\right)$$

- Parameters are

$$\Theta = \left\{\mathbf{b}^{(1)}, \mathbf{w}^{(1)}, \mathbf{b}^{(2)}, \mathbf{w}^{(2)}, \dots, \mathbf{b}^{(K)}, \mathbf{w}^{(K)}\right\}$$

Loss function

- Measure the discrepancy between
 - the prediction $F_{\Theta}(\mathbf{x})$
 - the expected output $\mathbf{y} = (y_1, \dots, y_{n_K})$
- Denoted $\ell(F_{\Theta}(\mathbf{x}), \mathbf{y})$
- Output $\mathbf{x} = (x_1, \dots, x_{n_K})$
- Expected output $\mathbf{y} = (y_1, \dots, y_{n_K})$
- Mean squared error (MSE) loss:

$$\frac{1}{n_K} \|\mathbf{x} - \mathbf{y}\|^2 = \frac{1}{n_K} \sum_{i=1}^{n_K} (x_i - y_i)^2$$

Cross entropy

- Adapted to expected output in range (0, 1)
- In classification task, all the y_i 's are zero except one (one-hot encoding)
- The output **has to** take values in that range
- Use normalizing transform if necessary (softmax)
- Cross entropy error

$$\text{CE}(\mathbf{x}, \mathbf{y}) = - \sum_{i=1}^{n_k} y_i \log x_i$$

- If \mathbf{y} is one-hot encoded, $\text{CE}(\mathbf{x}, \mathbf{y}) = -\log x_k$ with $y_k = 1$

- Parameter-free normalizing transform

$$\text{softmax} : \mathbb{R}^{n_K} \rightarrow \{(p_1, \dots, p_{n_K}) \in \mathbb{R}^{n_K}, p_i \geq 0, p_1 + \dots + p_{n_K} = 1\}$$

- If $\mathbf{z} = \text{softmax}(\mathbf{x})$ we have

$$z_i = \frac{\exp x_i}{\sum_{j=1}^{n_K} \exp x_j}$$

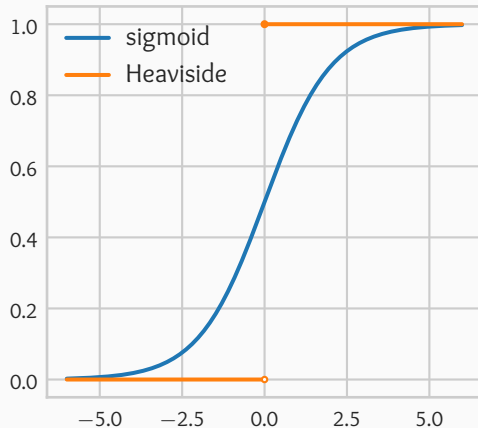
- $\text{softmax}(\mathbf{x} + c\mathbf{1}) = \text{softmax}(\mathbf{x})$, with $c \in \mathbb{R}$
- If $x_i \geq x_j$ then $\text{softmax}(\mathbf{x})_i \geq \text{softmax}(\mathbf{x})_j$
- Really an “arg softmax” rather than a “softmax”

Activation functions

Logistic function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Saturates at 0 and 1 when $x \rightarrow \pm\infty$
- Smooth version of Heaviside function
- $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- Killing gradients



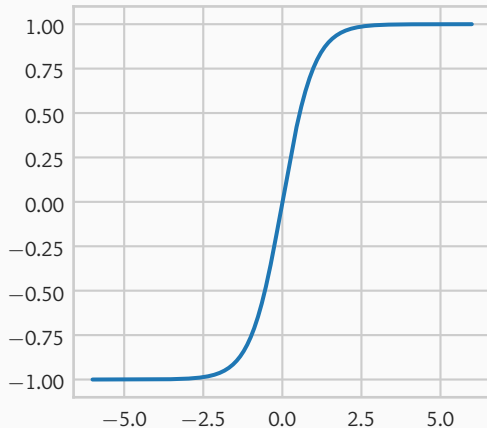
Hyperbolic tangent

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

- Symmetric sigmoid
- Saturates at ± 1 when $x \rightarrow \pm\infty$
- $\tanh'(x) = 1 - \tanh^2(x)$
- Linearly related to the sigmoid function by

$$\tanh(x) = 2\sigma(2x) - 1$$

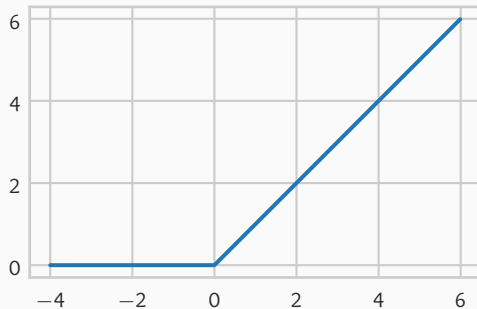
- Killing gradients



Rectified linear unit (ReLU)

$$\text{ReLU}(x) = \max(x, 0)$$

- Learns faster than sigmoid-like activation function
- Simpler to compute
- Provide sparsity of activations
- Dead ReLU: never activated across whole training set.
- Non negative activation function: zig-zag learning (equation (4))

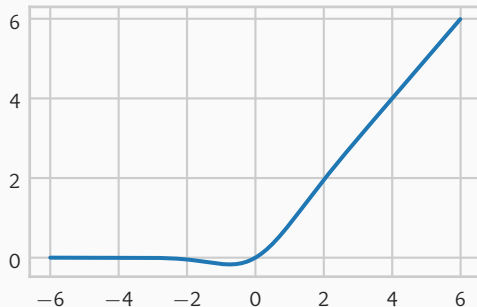


Gaussian Error Linear Unit (GELU), see Hendrycks and Gimpel 2023

$$\text{GELU}(x) = x\Phi(x)$$

where Φ is the standard Gaussian cumulative distribution function

- Smoother version of ReLU
- Mostly used in transformers



Summary

- Use ReLU
- Sigmoid not used anymore
- Prefer hyperbolic tangent to sigmoid

Gradient descent algorithms

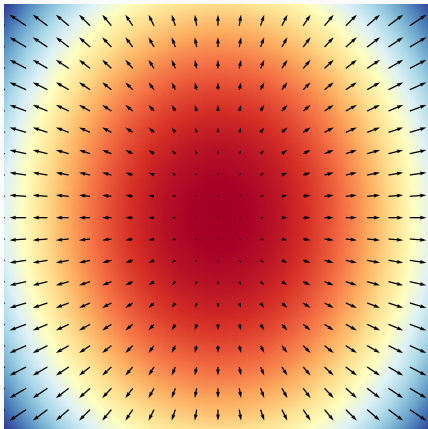
Gradient of a (scalar) function

- Differentiable function

$$f(x,y) = 1.5x^2 + y^2 + 2$$

- The gradient

$$\nabla_{(x,y)} f = \begin{pmatrix} 3x \\ 2y \end{pmatrix}$$



Gradient descent: motivation I

Gradient descent step improves current solution:

- Suppose \mathcal{L} is a differentiable function we want to minimize

$$\arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta})$$

- Starting from the first order Taylor expansion. For a small $\|\mathbf{h}\|$ we have

$$\mathcal{L}(\boldsymbol{\theta} + \mathbf{h}) \approx \mathcal{L}(\boldsymbol{\theta}) + \langle \nabla_{\boldsymbol{\theta}} \mathcal{L}, \mathbf{h} \rangle$$

- Choose $\mathbf{h} = -\eta \nabla_{\boldsymbol{\theta}} \mathcal{L}$ (the gradient descent step), we have

$$\mathcal{L}(\boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}) \approx \mathcal{L}(\boldsymbol{\theta}) - \eta \|\nabla_{\boldsymbol{\theta}} \mathcal{L}\|^2 < \mathcal{L}(\boldsymbol{\theta})$$

- $\boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}$ is better than $\boldsymbol{\theta}$

Gradient descent: motivation II

Gradient descent step yields best update of linearized and regularized objective function

- Looking for the best $\theta' = \theta + \mathbf{h}$ around fixed θ
- Instead of minimizing $\mathcal{L}(\theta')$, we minimize

$$\mathcal{L}(\theta') \approx \mathcal{L}(\theta) + \langle \mathbf{h}, \nabla_{\theta} \mathcal{L} \rangle \quad (1)$$

- θ' should stay close to θ for (1) to hold, we penalize by $\|\theta - \theta'\| = \|\mathbf{h}\|$

$$\mathcal{L}(\theta) + \langle \mathbf{h}, \nabla_{\theta} \mathcal{L} \rangle + \frac{1}{2\eta} \|\mathbf{h}\|^2$$

- Minimizing w.r.t \mathbf{h} gives

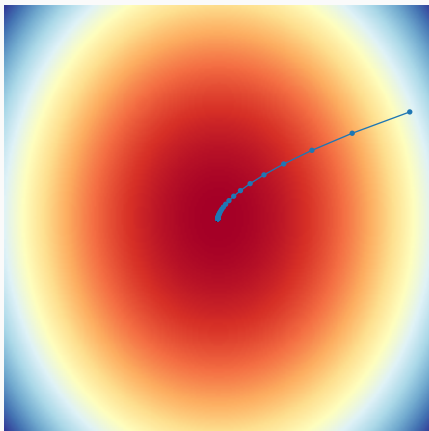
$$\theta' = \theta - \eta \nabla_{\theta} \mathcal{L}$$

Gradient descent algorithm

- Starting point θ_0
- Gradient descent step

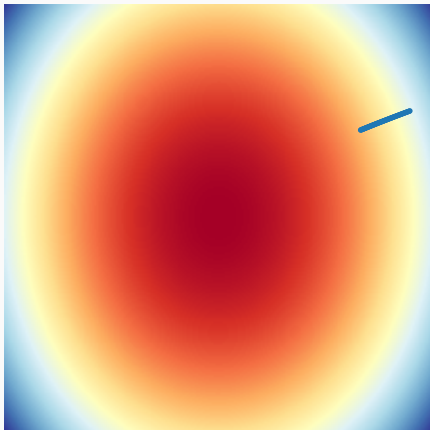
$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \mathcal{L}$$

- η is the *learning rate*
- Learning rate too small: slow convergence
- Learning rate too high: fluctuate around minimum or even diverge

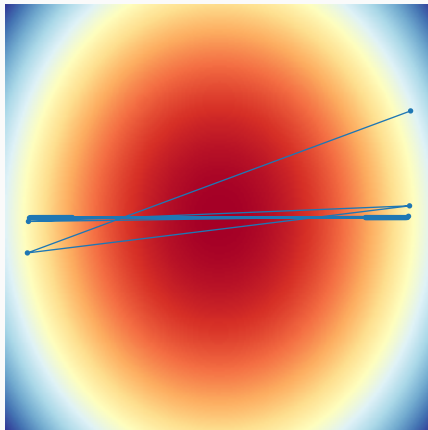


(a) $\eta = 0.1$

Gradient descent: illustration



(a) $\eta = 0.001$



(b) $\eta = 0.665757$

- Empirical risk minimization (ERM)

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(F_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i)$$

- Computing $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$ requires computing $\nabla_{\boldsymbol{\theta}} \ell(F_{\boldsymbol{\theta}}(\mathbf{x}_i), y_i)$ for all $i = 1, \dots, n$!
- Efficient when $n < 10000$
- Impractical when $n > 10000$

Stochastic gradient descent (SGD)

- Directly minimizing the expected risk

$$\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} \ell(F_{\boldsymbol{\theta}}(\mathbf{x}), y)$$

- Applying the gradient operator $\nabla_{\boldsymbol{\theta}}$ yields

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} \ell(F_{\boldsymbol{\theta}}(\mathbf{x}), y) \\ &= \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} \nabla_{\boldsymbol{\theta}} \ell(F_{\boldsymbol{\theta}}(\mathbf{x}), y)\end{aligned}$$

- $\nabla_{\boldsymbol{\theta}} \ell(F_{\boldsymbol{\theta}}(\mathbf{x}), y)$ is an unbiased estimator of $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta})$

Minibatch Stochastic gradient descent

- Easy to get an observation of $\nabla_{\theta} \ell(F_{\theta}(\mathbf{x}), y)$
- Unbiased but possibly of high variance
- Combining unbiased estimator to reduce variance

$$\mathcal{L}_{\mathcal{B}}(\theta) = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \ell(F_{\theta}(\mathbf{x}), y)$$

- $\nabla_{\theta} \mathcal{L}_{\mathcal{B}}$ is an unbiased estimate of $\nabla_{\theta} \mathcal{L}_{\mathcal{D}}$
- $\mathbb{E}_{\mathcal{B}} \nabla_{\theta} \mathcal{L}_{\mathcal{B}} = \nabla_{\theta} \mathcal{L}_{\mathcal{D}}$
- Size of minibatch $|\mathcal{B}|$ is a tradeoff term between good estimate of the gradient and cheap computations

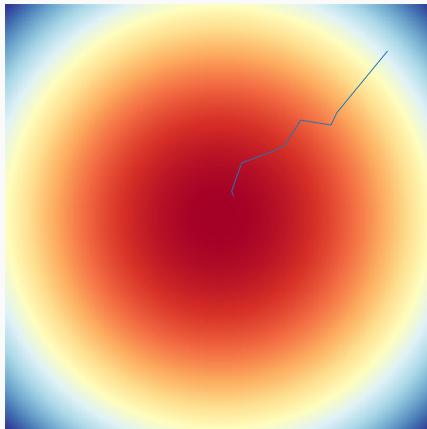
Illustrations

- Model $Y = \langle \mathbf{w}_0, \mathbf{X} \rangle + \varepsilon$ with:
 - $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \Sigma)$
 - $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- Minimization problem

$$\arg \min_{\mathbf{w} \in \mathbb{R}^p} \mathbb{E} (Y - \langle \mathbf{w}, \mathbf{X} \rangle)^2$$

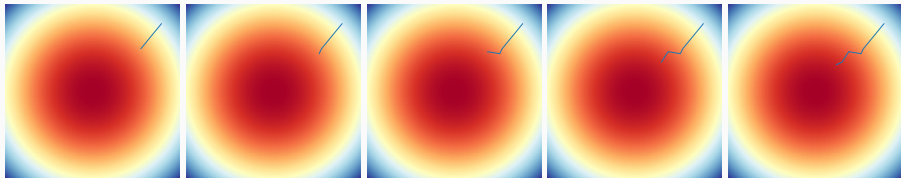
- One can show that

$$\mathbb{E} (Y - \langle \mathbf{w}, \mathbf{X} \rangle)^2 = (\mathbf{w} - \mathbf{w}_0)^T \Sigma (\mathbf{w} - \mathbf{w}_0)$$

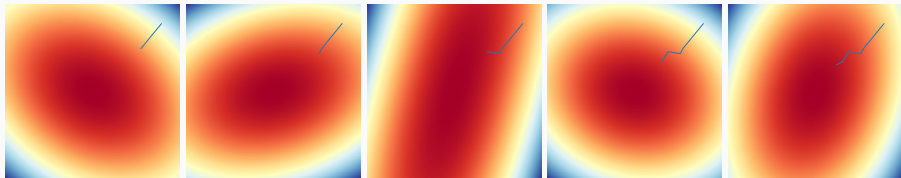


Stochastic gradient descent path

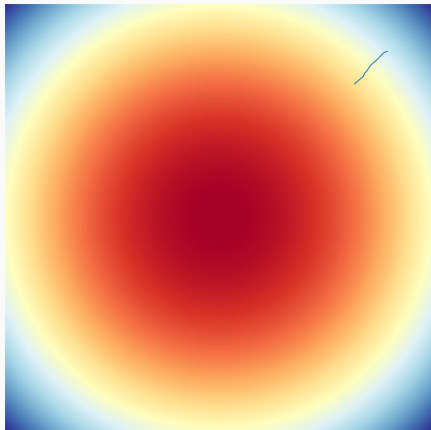
- Gradient descent for first minibatches, theoretical loss is $\mathbb{E} (Y - \langle \mathbf{w}, \mathbf{x} \rangle)^2$



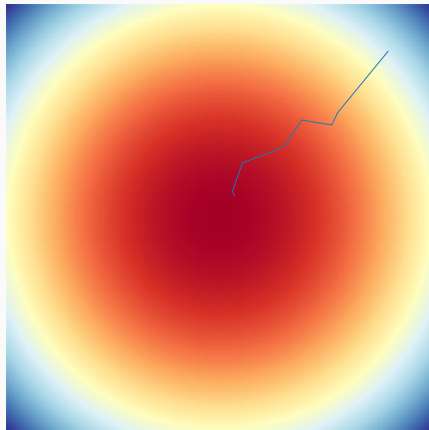
- What SGD is really seeing through minibatches, loss is $\frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} (y - \langle \mathbf{w}, \mathbf{x} \rangle)^2$



Influence of learning rate

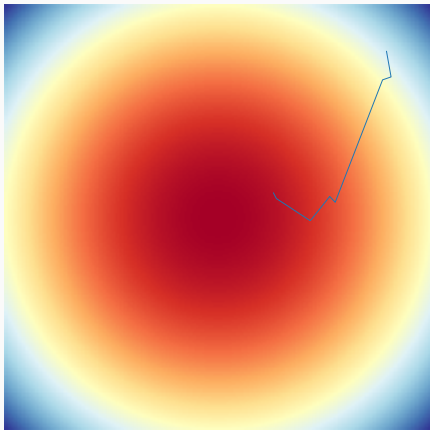


(a) $\eta = 0.01, |\mathcal{B}| = 10$

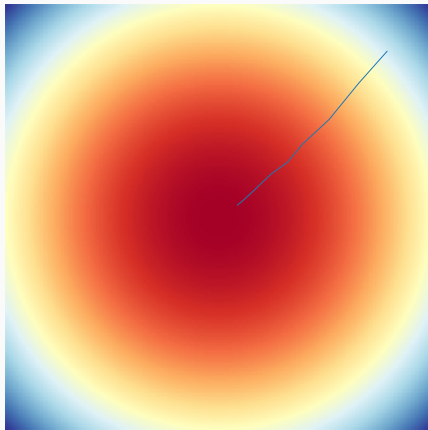


(b) $\eta = 0.1, |\mathcal{B}| = 10$

Influence of minibatch size



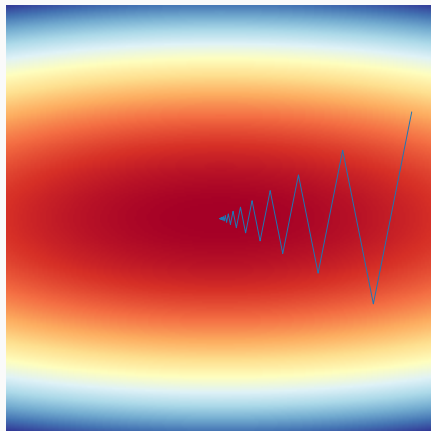
(a) $\eta = 0.1, |\mathcal{B}| = 1$



(b) $\eta = 0.1, |\mathcal{B}| = 100$

Gradient descent with momentum

- If learning rate is high or loss landscape is bumpy SGD (and GD) trajectory might be erratic



(a) $\eta = 0.1$

Gradient descent with momentum

- Original update rule

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}_t} \mathcal{L}$$

- Gradient descent rule with momentum ($\mathbf{v}_0 = \nabla_{\boldsymbol{\theta}_0} \mathcal{L}$)

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{v}_t$$

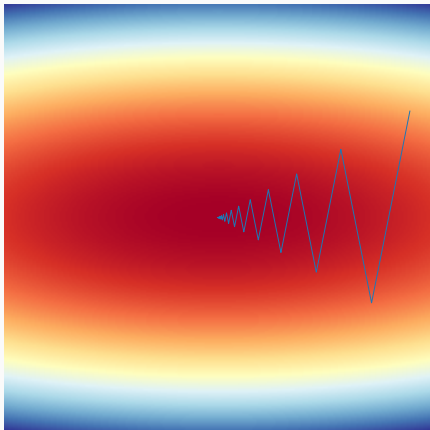
$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla_{\boldsymbol{\theta}_t} \mathcal{L}$$

- Average current gradient with previous one
- Exponentially weighted moving average on past gradients

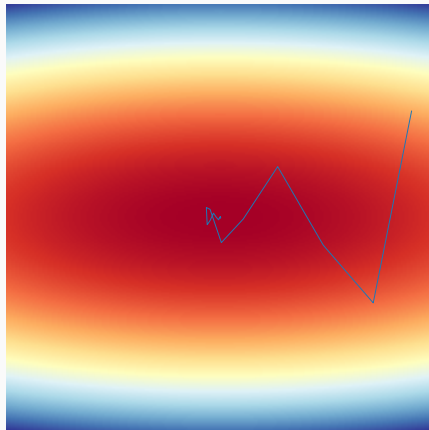
$$\mathbf{v}_t = (1 - \beta) \sum_{i=0}^{t-1} \beta^i \nabla_{\boldsymbol{\theta}_{t-i}} \mathcal{L} + \beta^t \nabla_{\boldsymbol{\theta}_0} \mathcal{L}$$

- Averaging so better estimate

Momentum in action (GD)



(a) $\eta = 0.1$, no momentum. Gradients are abruptly changing



(b) $\eta = 0.1$, momentum = 0.5. Gradients change is smoothed out

Other stochastic optimizers

- Adagrad, RMSProp, **Adam**: decreases the learning rate dynamically coordinate-wise
- Adadelta: adapt learning rate from rate of change in the parameters

Backpropagation algorithm

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta)$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta) \xrightarrow{f_2} f_2(f_1(\theta))$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\theta \xrightarrow{f_1} f_1(\theta) \xrightarrow{f_2} f_2(f_1(\theta)) \xrightarrow{f_3} f_3(f_2(f_1(\theta))) = F(\theta)$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ & & & & \downarrow & & \\ & & & & f_3'(f_2(f_1(\theta))) & & \end{array}$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f_3' \circ f_2 \circ f_1)(\theta) \cdot (f_2' \circ f_1)(\theta) \cdot f_1'(\theta) \\ &= f_3'((f_2 \circ f_1)(\theta)) \cdot f_2'(f_1(\theta)) \cdot f_1'(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ & & \downarrow & & \downarrow & & \\ & & f_1'(\theta) & \xleftarrow{\times} & f_2'(f_1(\theta)) & & \end{array}$$

Backpropagation

- Chain rule

$$(f \circ g)'(\theta) = (f' \circ g)(\theta) \cdot g'(\theta)$$

- Generalization to any number of functions ($F = f_3 \circ f_2 \circ f_1$)

$$\begin{aligned} F'(\theta) &= (f_3 \circ f_2 \circ f_1)'(\theta) = (f'_3 \circ f_2 \circ f_1)(\theta) \cdot (f'_2 \circ f_1)(\theta) \cdot f'_1(\theta) \\ &= f'_3((f_2 \circ f_1)(\theta)) \cdot f'_2(f_1(\theta)) \cdot f'_1(\theta) \end{aligned}$$

- Forward and backward pass

$$\begin{array}{ccccccc} \theta & \xrightarrow{f_1} & f_1(\theta) & \xrightarrow{f_2} & f_2(f_1(\theta)) & \xrightarrow{f_3} & f_3(f_2(f_1(\theta))) = F(\theta) \\ \downarrow & & \downarrow & & \downarrow & & \\ F'(\theta) = f'_1(\theta) & \xleftarrow{\times} & f'_2(f_1(\theta)) & \xleftarrow{\times} & f'_3(f_2(f_1(\theta))) & & \end{array}$$

- Generalizable to \mathbb{R}^n to \mathbb{R}^p functions using jacobians
- Generalizable to a computational graph (DAG)

Backpropagating the loss

Total loss on the minibatch \mathcal{B}

$$\mathcal{L}_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(F_{\Theta}(\mathbf{x}), \mathbf{y})$$

Differentiating w.r.t any scalar parameter θ (any $\mathbf{w}_{ij}^{(k)}$ or $\mathbf{b}_j^{(k)}$)

$$\frac{\partial \mathcal{L}_{\mathcal{B}}}{\partial \theta} = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \frac{\partial \ell(\mathbf{x}^{(k)}, \mathbf{y})}{\partial \theta} \quad \text{avec} \quad \mathbf{x}^{(k)} = F_{\Theta}(\mathbf{x})$$

$$\ell(\mathbf{x}^{(k)}, \mathbf{y}) = \frac{1}{n_K} \sum_{i=1}^{n_K} \ell(\mathbf{x}_i^{(k)}, \mathbf{y}_i)$$

It suffices to compute

$$\frac{\partial \ell(\mathbf{x}_i^{(k)}, \mathbf{y}_i)}{\partial \theta}$$

MSE loss

- MSE loss is defined by

$$\ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i) = (\mathbf{x}_i^{(K)} - \mathbf{y}_i)^2$$

- Differentiating with respect to some scalar parameter $\theta = \mathbf{b}_q^{(k)}$ or $\theta = \mathbf{w}_{pq}^{(k)}$
- $(\mathbf{x}_i^{(K)})$ and \mathbf{y}_i are scalars)

$$\begin{aligned}\frac{\partial \ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i)}{\partial \theta} &= \frac{\partial \ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i)}{\partial \mathbf{x}_i^{(K)}} \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \\ &= 2(\mathbf{x}_i^{(K)} - \mathbf{y}_i) \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}\end{aligned}$$

- We need to know $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}$ now!

Cross entropy loss

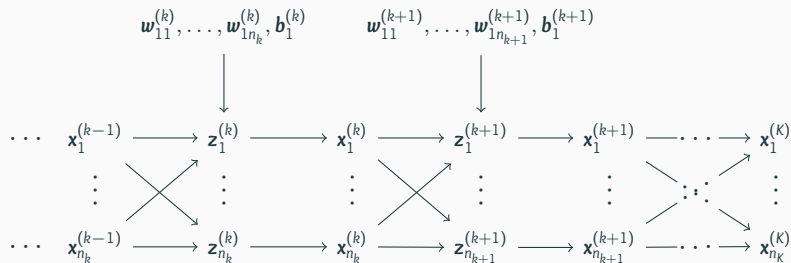
- Cross entropy loss: $\ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i) = -\mathbf{y}_i \log \mathbf{x}_i^{(K)}$
- Differentiating with respect to some scalar parameter $\theta = \mathbf{b}_q^{(k)}$ or $\theta = \mathbf{w}_{pq}^{(k)}$

$$\begin{aligned}\frac{\partial \ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i)}{\partial \theta} &= \frac{\partial \ell(\mathbf{x}_i^{(K)}, \mathbf{y}_i)}{\partial \mathbf{x}_i^{(K)}} \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta} \\ &= -\frac{\mathbf{y}_i}{\mathbf{x}_i^{(K)}} \cdot \frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}\end{aligned}$$

- We need to know $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \theta}$ now!

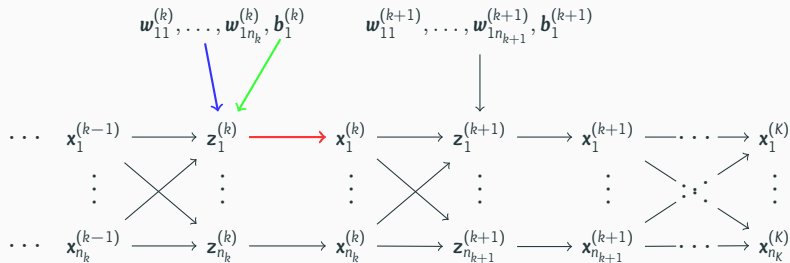
Computational graph

- Define $\mathbf{z}_i^{(k)} = \langle \mathbf{w}_i^{(k)}, \mathbf{x}^{(k-1)} \rangle + \mathbf{b}_i^{(k)}$ so that we have $\mathbf{x}_i^{(k)} = \sigma(\mathbf{z}_i^{(k)})$
- Computational graph for $\mathbf{x}_1^{(k-1)}, \mathbf{x}_1^{(k)}, \mathbf{x}_1^{(k+1)}$ only!



Gradient of last layer w.r.t parameters I

Computational graph for neuron 1 at layer k ($p = 1$):



$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial w_{pq}^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \frac{\partial \mathbf{x}_p^{(k)}}{\partial \mathbf{z}_p^{(k)}} \cdot \frac{\partial \mathbf{z}_p^{(k)}}{\partial w_{pq}^{(k)}}$$

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial b_p^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \frac{\partial \mathbf{x}_p^{(k)}}{\partial \mathbf{z}_p^{(k)}} \cdot \frac{\partial \mathbf{z}_p^{(k)}}{\partial b_p^{(k)}}$$

Gradient of last layer w.r.t parameters II

Given that $\mathbf{x}_p^{(k)} = \sigma(\mathbf{z}_p^{(k)})$, we have $\frac{\partial \mathbf{x}_p^{(k)}}{\partial \mathbf{z}_p^{(k)}} = \sigma'(\mathbf{z}_p^{(k)})$

And $\mathbf{z}_p^{(k)} = \langle \mathbf{w}_p^{(k)}, \mathbf{x}^{(k-1)} \rangle + \mathbf{b}_p^{(k)}$, so $\frac{\partial \mathbf{z}_p^{(k)}}{\partial \mathbf{w}_{pq}^{(k)}} = \mathbf{x}_q^{(k-1)}$ and $\frac{\partial \mathbf{z}_p^{(k)}}{\partial \mathbf{b}_p^{(k)}} = 1$

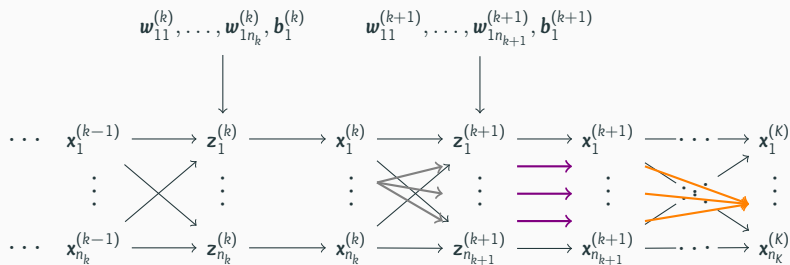
Finally

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{w}_{pq}^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \cdot \mathbf{x}_q^{(k-1)}$$

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{b}_p^{(k)}} = \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)})$$

Need to compute $\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}}$ now!

Gradient of last layer w.r.t other layer



$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} = \sum_{j=1}^{n_{k+1}} \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_j^{(k+1)}} \cdot \frac{\partial \mathbf{x}_j^{(k+1)}}{\partial \mathbf{z}_j^{(k+1)}} \cdot \frac{\partial \mathbf{z}_j^{(k+1)}}{\partial \mathbf{x}_p^{(k)}}$$

Backpropagating from next layer

Given that $\mathbf{x}_j^{(k+1)} = \sigma(\mathbf{z}_j^{(k+1)})$, we have

$$\frac{\partial \mathbf{x}_j^{(k+1)}}{\partial \mathbf{z}_j^{(k+1)}} = \sigma'(\mathbf{z}_j^{(k+1)})$$

And $\mathbf{z}_j^{(k+1)} = \langle \mathbf{w}_j^{(k+1)}, \mathbf{x}^{(k)} \rangle + \mathbf{b}_j^{(k+1)}$, so

$$\frac{\partial \mathbf{z}_j^{(k+1)}}{\partial \mathbf{x}_p^{(k)}} = \mathbf{w}_{jp}^{(k+1)}$$

Replacing in last equation, we have

$$\frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_p^{(k)}} = \sum_{j=1}^{n_{k+1}} \frac{\partial \mathbf{x}_i^{(K)}}{\partial \mathbf{x}_j^{(k+1)}} \cdot \sigma'(\mathbf{z}_j^{(k+1)}) \cdot \mathbf{w}_{jp}^{(k+1)}$$

Backpropagation equations

- Backpropagation equations for parameters $\mathbf{w}_{pq}^{(k)}$ and $\mathbf{b}_p^{(k)}$

$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{w}_{pq}^{(k)}} = \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \cdot \mathbf{x}_q^{(k-1)} \quad (2)$$

$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{b}_p^{(k)}} = \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} \cdot \sigma'(\mathbf{z}_p^{(k)}) \quad (3)$$

- Backpropagation equation

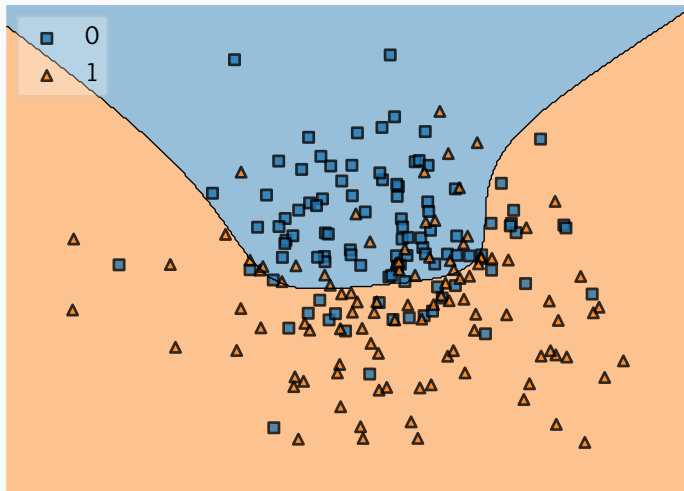
$$\frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_p^{(k)}} = \sum_{j=1}^{n_{k+1}} \frac{\partial \mathbf{x}_i^{(k)}}{\partial \mathbf{x}_j^{(k+1)}} \cdot \sigma'(\mathbf{z}_j^{(k+1)}) \cdot \mathbf{w}_{jp}^{(k+1)} \quad (4)$$

- Only one sweep backward is necessary to compute all gradients

Regularization

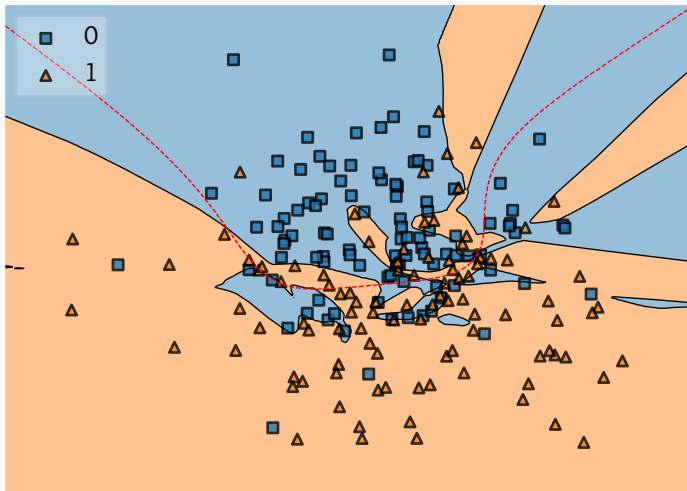
Toy dataset

- 200 samples, 2 classes
- Gaussian mixture model
- Bayes decision boundary



Standard neural network classification

- Neural network
 - 1 hidden layer with 50 units
 - ~ 200 parameters
- SGD algorithm
 - learning rate: 0.1
 - momentum: 0.9



Regularizing the ERM

- Empirical risk minimization (ERM)

$$\arg \min_{\theta \in \Theta} \mathcal{L}_{\mathcal{B}} = \arg \min_{\theta \in \Theta} \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(F_{\Theta}(\mathbf{x}), \mathbf{y})$$

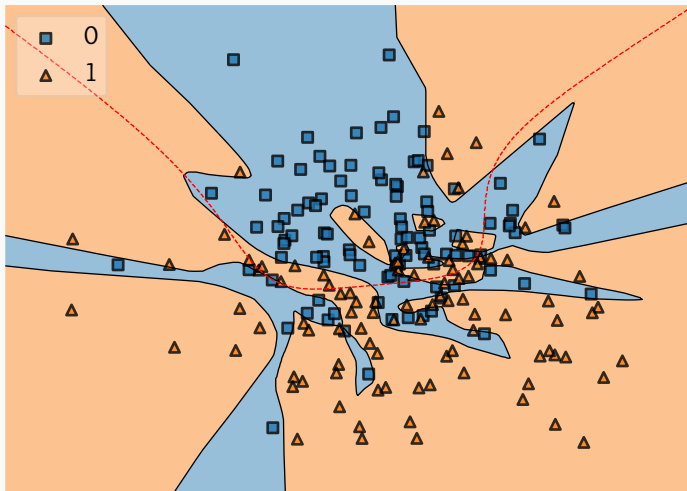
- L_2 penalizing term

$$\arg \min_{\theta \in \Theta} \mathcal{L}_{\mathcal{R}} = \arg \min_{\theta \in \Theta} \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \ell(F_{\Theta}(\mathbf{x}), \mathbf{y}) + \lambda \sum_{k=1}^K \left\| \omega^{(k)} \right\|_F$$

- Biases terms are not regularized
- λ is the tradeoff parameter called *weight decay*

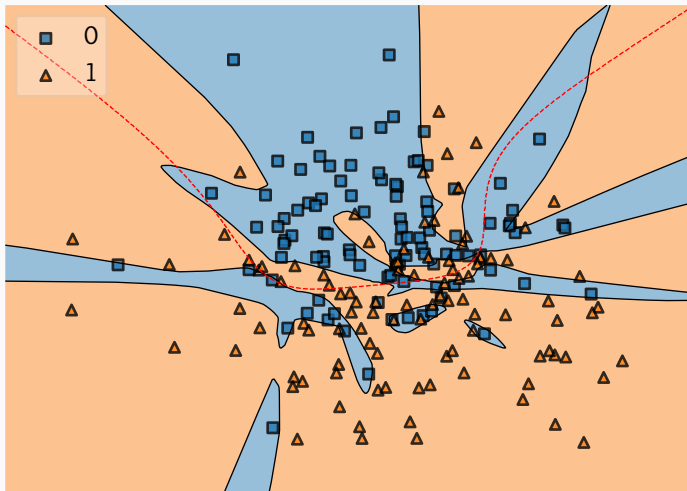
Weight decay: an example

- SGD algorithm
 - learning rate: 0.1
 - momentum: 0.9
- **weight decay:** 10^{-4}



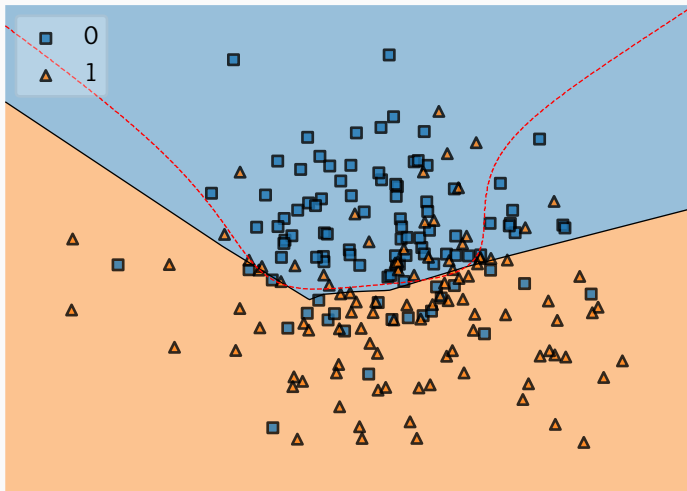
Weight decay: an example

- SGD algorithm
 - learning rate: 0.1
 - momentum: 0.9
- **weight decay:** 10^{-3}



Weight decay: an example

- SGD algorithm
 - learning rate: 0.1
 - momentum: 0.9
- **weight decay:** 10^{-2}



Weight decay: gradient

- One extra term in the loss

$$\mathcal{L}_{\mathcal{R}} = \mathcal{L}_{\mathcal{B}} + \lambda \sum_{k=1}^K \left\| \mathbf{w}^{(k)} \right\|_F$$

- Gradient is easy to get

$$\frac{\partial \mathcal{L}_{\mathcal{R}}}{\partial \mathbf{w}_{ij}^{(k)}} = \frac{\partial \mathcal{L}_{\mathcal{B}}}{\partial \mathbf{w}_{ij}^{(k)}} + 2\lambda \mathbf{w}_{ij}^{(k)}$$

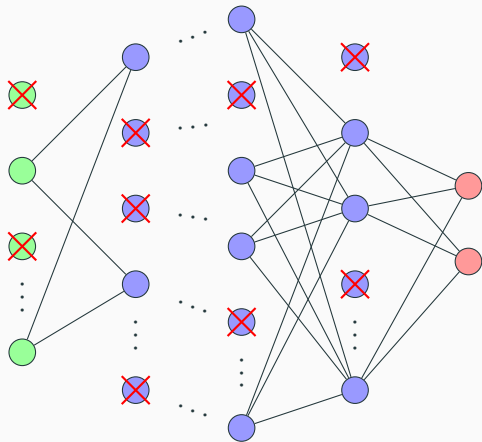
- Gradients w.r.t. $\mathbf{b}_i^{(k)}$ are unchanged

- Randomly kill nodes in layers during training time

$$\mathbf{x}^{(k)} = \sigma\left(\mathbf{W}^{(k)}\left(\mathbf{x}^{(k-1)} \odot \mathbf{h}^{(k)}\right) + \mathbf{b}^{(k)}\right)$$

with $\mathbf{h}_k \sim \mathcal{B}(h_k)^{\otimes n_k}$

- h_k is the rate of dropout at layer k
- Rescaling needed at test-time



Batch normalization Ioffe and Szegedy 2015

- Normalize each activation independently from the minibatch statistics

$$\mu^{(k)} = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}^{(k)} \in \mathcal{B}^{(k)}} \mathbf{z}^{(k)}$$

$$\sigma_i^{(k)} = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}^{(k)} \in \mathcal{B}^{(k)}} \left(\mathbf{z}_i^{(k)} - \mu_i^{(k)} \right)^2$$

- For each element in the minibatch, replace $\mathbf{z}_i^{(k)}$ and $\mathbf{x}_i^{(k+1)}$ by

$$\tilde{\mathbf{z}}_i^{(k)} = \frac{\mathbf{z}_i^{(k)} - \mu_i^{(k)}}{\sigma_i^{(k)}} \quad \tilde{\mathbf{x}}_i^{(k+1)} = \sigma \left(\gamma_i^{(k)} \tilde{\mathbf{z}}_i^{(k)} + \beta_i^{(k)} \right)$$

- $\gamma_i^{(k)}$ and $\beta_i^{(k)}$ are $2n_k$ extra parameters
- The \mathbf{b}_i 's from $\mathbf{z}_i^{(k)} = \langle \mathbf{w}_i^{(k)}, \mathbf{x}^{(k)} \rangle + \mathbf{b}_i^{(k)}$ are useless ($\tilde{\mathbf{z}}_i^{(k)}$ does not depend on \mathbf{b}_i)

References i

- [1] Nitish Srivastava et al. “**Dropout: a simple way to prevent neural networks from overfitting.**”. In: *Journal of machine learning research* 15.1 (1 2014), pp. 1929–1958.
- [2] Sergey Ioffe and Christian Szegedy. “**Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift**”. Feb. 10, 2015. arXiv: 1502.03167 [cs]. URL: <http://arxiv.org/abs/1502.03167> (visited on 11/06/2017).
- [3] Ian Goodfellow et al. ***Deep learning.*** Vol. 1. MIT press Cambridge, 2016.
- [4] Aston Zhang et al. ***Dive into deep learning.*** 2020.
- [5] Dan Hendrycks and Kevin Gimpel. ***Gaussian Error Linear Units (GELUs).*** June 5, 2023. DOI: 10.48550/arXiv.1606.08415. arXiv: 1606.08415 [cs]. URL: <http://arxiv.org/abs/1606.08415> (visited on 11/07/2023). preprint.