CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 12, 2021







2 Review of last lecture



Outline



2 Review of last lecture





• Sign-up with your group members for the project!

Outline







Linear Discriminant Analysis

The main bottleneck is *not knowing* $\mathcal{P}(X = \boldsymbol{x}|y = c)$

$$\mathcal{P}(y = c | X = \boldsymbol{x}) = \frac{\mathcal{P}(X = \boldsymbol{x} | y = c) \mathcal{P}(y = c)}{\mathcal{P}(X = \boldsymbol{x})}$$

LDA makes two simplifying assumptions:

• Let
$$\mathcal{P}(X = oldsymbol{x} | y = c) \sim \mathcal{N}(oldsymbol{\mu}_c, \Sigma_c)$$
, and

• Let all class covariances be the same i.e. $\Sigma_c = \Sigma$ for all $c \in [C]$

If so, the decision boundary (for binary classification) is given by

$$\mathcal{P}(y=0|X=\boldsymbol{x}) = \mathcal{P}(y=1|X=\boldsymbol{x})$$
$$\mathcal{P}(X=\boldsymbol{x}|y=0)\mathcal{P}(y=0) = \mathcal{P}(X=\boldsymbol{x}|y=1)\mathcal{P}(y=1)$$

What do the decision boundaries look like?



The decision boundaries are a quadratic when Σ 's are not the same, this is known as *Quadratic Discriminant Analysis*!

Outline





3 Neural Nets

- Definition
- Backpropagation
- Preventing overfitting

Linear models are not always adequate



We can use a nonlinear mapping as discussed:

$$oldsymbol{\phi}(oldsymbol{x}):oldsymbol{x}\in\mathbb{R}^{\mathsf{D}}
ightarrowoldsymbol{z}\in\mathbb{R}^{\mathsf{M}}$$

Linear models are not always adequate



We can use a nonlinear mapping as discussed:

$$oldsymbol{\phi}(oldsymbol{x}):oldsymbol{x}\in\mathbb{R}^{\mathsf{D}} ooldsymbol{z}\in\mathbb{R}^{\mathsf{M}}$$

But what kind of nonlinear mapping ϕ should be used? Can we actually learn this nonlinear mapping?

Linear models are not always adequate



We can use a nonlinear mapping as discussed:

$$oldsymbol{\phi}(oldsymbol{x}):oldsymbol{x}\in\mathbb{R}^{\mathsf{D}} ooldsymbol{z}\in\mathbb{R}^{\mathsf{M}}$$

But what kind of nonlinear mapping ϕ should be used? Can we actually learn this nonlinear mapping?

THE most popular nonlinear models nowadays: neural nets

Linear model as a one-layer neural net



h(a) = a for linear model

Linear model as a one-layer neural net



h(a) = a for linear model

To create non-linearity, can use

- Rectified Linear Unit (ReLU): $h(a) = \max\{0, a\}$
- sigmoid function: $h(a) = \frac{1}{1+e^{-a}}$
- TanH: $h(a) = \frac{e^a e^{-a}}{e^a + e^{-a}}$
- many more

More output nodes



 $oldsymbol{W} \in \mathbb{R}^{4 imes 3}$, $oldsymbol{h}: \mathbb{R}^4 o \mathbb{R}^4$ so $oldsymbol{h}(oldsymbol{a}) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$

More output nodes



 $m{W} \in \mathbb{R}^{4 imes 3}$, $m{h}: \mathbb{R}^4 o \mathbb{R}^4$ so $m{h}(m{a}) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$

Can think of this as a nonlinear basis: ${f \Phi}({m x})={m h}({m W}{m x})$

More layers



Becomes a network:

More layers



Becomes a network:

More layers



Becomes a network:

- h is called the activation function
 - can use h(a) = 1 for one neuron in each layer to *incorporate bias term*
 - output neuron can use h(a) = a

More layers



Becomes a network:



- h is called the activation function
 - can use h(a) = 1 for one neuron in each layer to *incorporate bias term*
 - output neuron can use h(a) = a
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)

More layers



Becomes a network:



- h is called the activation function
 - can use h(a) = 1 for one neuron in each layer to *incorporate bias term*
 - output neuron can use h(a) = a
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters

More layers



Becomes a network:



- h is called the activation function
 - can use h(a) = 1 for one neuron in each layer to *incorporate bias term*
 - output neuron can use h(a) = a
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters
- this is a feedforward, fully connected neural net, there are many variants

How powerful are neural nets?

Universal approximation theorem (Cybenko, 89; Hornik, 91):

A feedforward neural net with a single hidden layer can approximate any continuous functions.

How powerful are neural nets?

Universal approximation theorem (Cybenko, 89; Hornik, 91):

A feedforward neural net with a single hidden layer can approximate any continuous functions.

It might need a huge number of neurons though, and *depth helps!*

How powerful are neural nets?

Universal approximation theorem (Cybenko, 89; Hornik, 91):

A feedforward neural net with a single hidden layer can approximate any continuous functions.

It might need a huge number of neurons though, and *depth helps!*

Designing network architecture is important and very complicated

• for feedforward network, need to decide number of hidden layers, number of neurons at each layer, activation functions, etc.

Math formulation

An L-layer neural net can be written as

$$\boldsymbol{f}(\boldsymbol{x}) = \boldsymbol{h}_{\mathsf{L}}\left(\boldsymbol{W}_{L}\boldsymbol{h}_{\mathsf{L}-1}\left(\boldsymbol{W}_{L-1}\cdots\boldsymbol{h}_{1}\left(\boldsymbol{W}_{1}\boldsymbol{x}
ight)
ight)$$



input layer hidden layer 1 hidden layer 2 output layer

Math formulation

An L-layer neural net can be written as

$$oldsymbol{f}(oldsymbol{x}) = oldsymbol{h}_{\mathsf{L}}\left(oldsymbol{W}_{L}oldsymbol{h}_{\mathsf{L}-1}\left(oldsymbol{W}_{L-1}\cdotsoldsymbol{h}_{1}\left(oldsymbol{W}_{1}oldsymbol{x}
ight)
ight)$$



nput layer hidden layer 1 hidden layer 2 output layer

To ease notation, for a given input x, define recursively

$$oldsymbol{o}_0 = oldsymbol{x}, \qquad oldsymbol{a}_\ell = oldsymbol{W}_\ell oldsymbol{o}_{\ell-1}, \qquad oldsymbol{o}_\ell = oldsymbol{h}_\ell(oldsymbol{a}_\ell) \qquad \quad (\ell = 1, \dots, \mathsf{L})$$

where

- $m{W}_\ell \in \mathbb{R}^{\mathsf{D}_\ell imes \mathsf{D}_{\ell-1}}$ is the weights between layer $\ell-1$ and ℓ
- $\bullet \ D_0 = D, D_1, \ldots, D_L$ are numbers of neurons at each layer
- $a_\ell \in \mathbb{R}^{\mathsf{D}_\ell}$ is input to layer ℓ
- $o_\ell \in \mathbb{R}^{\mathsf{D}_\ell}$ is output to layer ℓ
- $h_{\ell} : \mathbb{R}^{\mathsf{D}_{\ell}} \to \mathbb{R}^{\mathsf{D}_{\ell}}$ is activation functions at layer ℓ

Learning the model

No matter how complicated the model is, our goal is the same: minimize

$$\mathcal{E}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \frac{1}{N} \sum_{n=1}^{\mathsf{N}} \mathcal{E}_n(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}})$$

Learning the model

No matter how complicated the model is, our goal is the same: minimize

$$\mathcal{E}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \frac{1}{N} \sum_{n=1}^{\mathsf{N}} \mathcal{E}_n(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}})$$

where

$$\mathcal{E}_n(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \begin{cases} \|\boldsymbol{f}(\boldsymbol{x}_n) - \boldsymbol{y}_n\|_2^2 & \text{for regression} \\ \ln\left(1 + \sum_{k \neq y_n} e^{f(\boldsymbol{x}_n)_k - f(\boldsymbol{x}_n)_{y_n}}\right) & \text{for classification} \end{cases}$$

Same thing: apply SGD! even if the model is nonconvex.

Same thing: apply **SGD**! even if the model is *nonconvex*.

What is the gradient of this complicated function?

Same thing: apply **SGD**! even if the model is *nonconvex*. What is the gradient of this complicated function?

Chain rule is the only secret:

 \bullet for a composite function f(g(w))

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

Same thing: apply **SGD**! even if the model is *nonconvex*. What is the gradient of this complicated function?

Chain rule is the only secret:

 \bullet for a composite function f(g(w))

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

• for a composite function $f(g_1(w), \ldots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^{d} \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

Same thing: apply **SGD**! even if the model is *nonconvex*. What is the gradient of this complicated function?

Chain rule is the only secret:

 \bullet for a composite function f(g(w))

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

• for a composite function $f(g_1(w), \ldots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^{d} \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

the simplest example $f(g_1(w), g_2(w)) = g_1(w)g_2(w)$

Drop the subscript ℓ for layer for simplicity.



Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}}$$

Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}}$$

Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} o_j$$

Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} o_j$$
$$\frac{\partial \mathcal{E}_n}{\partial a_i} = \frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i}$$

Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} o_j$$
$$\frac{\partial \mathcal{E}_n}{\partial a_i} = \frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} \frac{\partial a_k}{\partial o_i}\right) h'_i(a_i)$$

Drop the subscript ℓ for layer for simplicity.



$$\frac{\partial \mathcal{E}_n}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}} = \frac{\partial \mathcal{E}_n}{\partial a_i} o_j$$
$$\frac{\partial \mathcal{E}_n}{\partial a_i} = \frac{\partial \mathcal{E}_n}{\partial o_i} \frac{\partial o_i}{\partial a_i} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} \frac{\partial a_k}{\partial o_i}\right) h'_i(a_i) = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k} w_{ki}\right) h'_i(a_i)$$

Adding the subscript for layer:

$$\frac{\partial \mathcal{E}_n}{\partial w_{\ell,ij}} = \frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$
$$\frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki}\right) h'_{\ell,i}(a_{\ell,i})$$



Adding the subscript for layer:

$$\frac{\partial \mathcal{E}_n}{\partial w_{\ell,ij}} = \frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$
$$\frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki}\right) h'_{\ell,i}(a_{\ell,i})$$

For the last layer, for square loss

$$\frac{\partial \mathcal{E}_n}{\partial a_{\mathsf{L},i}} = \frac{\partial (h_{\mathsf{L},i}(a_{\mathsf{L},i}) - y_{n,i})^2}{\partial a_{\mathsf{L},i}}$$



Adding the subscript for layer:

$$\frac{\partial \mathcal{E}_n}{\partial w_{\ell,ij}} = \frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$
$$\frac{\partial \mathcal{E}_n}{\partial a_{\ell,i}} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki}\right) h'_{\ell,i}(a_{\ell,i})$$



For the last layer, for square loss

$$\frac{\partial \mathcal{E}_n}{\partial a_{\mathsf{L},i}} = \frac{\partial (h_{\mathsf{L},i}(a_{\mathsf{L},i}) - y_{n,i})^2}{\partial a_{\mathsf{L},i}} = 2(h_{\mathsf{L},i}(a_{\mathsf{L},i}) - y_{n,i})h'_{\mathsf{L},i}(a_{\mathsf{L},i})$$

Using matrix notation greatly simplifies presentation and implementation:

$$\frac{\partial \mathcal{E}_n}{\partial \boldsymbol{W}_{\ell}} = \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} \boldsymbol{o}_{\ell-1}^{\mathrm{T}}$$

$$\frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} = \begin{cases} \left(\boldsymbol{W}_{\ell+1}^{\mathrm{T}} \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell+1}} \right) \circ \boldsymbol{h}_{\ell}'(\boldsymbol{a}_{\ell}) & \text{ if } \ell < \mathsf{L} \\ 2(\boldsymbol{h}_{\mathsf{L}}(\boldsymbol{a}_{\mathsf{L}}) - \boldsymbol{y}_n) \circ \boldsymbol{h}_{\mathsf{L}}'(\boldsymbol{a}_{\mathsf{L}}) & \text{ else} \end{cases}$$

where $v_1 \circ v_2 = (v_{11}v_{21}, \cdots, v_{1D}v_{2D})$ is the element-wise product (a.k.a. Hadamard product).

Verify yourself!

The **backpropagation** algorithm (**Backprop**)

- Initialize W_1, \ldots, W_L . Repeat:
 - () randomly pick one data point $n \in [\mathsf{N}]$

The **backpropagation** algorithm (**Backprop**)

- Initialize W_1, \ldots, W_L . Repeat:
 - () randomly pick one data point $n \in [N]$
 - **(2)** forward propagation: for each layer $\ell = 1, \dots, L$
 - compute $oldsymbol{a}_\ell = W_\ell oldsymbol{o}_{\ell-1}$ and $oldsymbol{o}_\ell = oldsymbol{h}_\ell(oldsymbol{a}_\ell)$ $(oldsymbol{o}_0 = oldsymbol{x}_n)$

The backpropagation algorithm (Backprop)

Initialize W_1, \ldots, W_L . Repeat:

() randomly pick one data point $n \in [N]$

(a) forward propagation: for each layer $\ell = 1, \dots, L$

• compute $oldsymbol{a}_\ell = oldsymbol{W}_\ell oldsymbol{o}_{\ell-1}$ and $oldsymbol{o}_\ell = oldsymbol{h}_\ell(oldsymbol{a}_\ell)$ $(oldsymbol{o}_0 = oldsymbol{x}_n)$

(a) backward propagation: for each $\ell = L, \dots, 1$

• compute

$$\frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} = \begin{cases} \left(\boldsymbol{W}_{\ell+1}^{\mathrm{T}} \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell+1}} \right) \circ \boldsymbol{h}_{\ell}'(\boldsymbol{a}_{\ell}) & \text{ if } \ell < \mathsf{L} \\ 2(\boldsymbol{h}_{\mathsf{L}}(\boldsymbol{a}_{\mathsf{L}}) - \boldsymbol{y}_n) \circ \boldsymbol{h}_{\mathsf{L}}'(\boldsymbol{a}_{\mathsf{L}}) & \text{ else} \end{cases}$$

• update weights

$$\boldsymbol{W}_{\ell} \leftarrow \boldsymbol{W}_{\ell} - \eta \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{W}_{\ell}} = \boldsymbol{W}_{\ell} - \eta \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} \boldsymbol{o}_{\ell-1}^{\mathrm{T}}$$

The **backpropagation** algorithm (**Backprop**)

Initialize W_1, \ldots, W_L . Repeat:

() randomly pick one data point $n \in [N]$

(a) forward propagation: for each layer $\ell = 1, \dots, L$

- compute $oldsymbol{a}_\ell = oldsymbol{W}_\ell oldsymbol{o}_{\ell-1}$ and $oldsymbol{o}_\ell = oldsymbol{h}_\ell(oldsymbol{a}_\ell)$ $(oldsymbol{o}_0 = oldsymbol{x}_n)$
- **(a)** backward propagation: for each $\ell = L, \dots, 1$
 - compute

$$\frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} = \begin{cases} \left(\boldsymbol{W}_{\ell+1}^{\mathrm{T}} \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell+1}} \right) \circ \boldsymbol{h}_{\ell}'(\boldsymbol{a}_{\ell}) & \text{ if } \ell < \mathsf{L} \\ 2(\boldsymbol{h}_{\mathsf{L}}(\boldsymbol{a}_{\mathsf{L}}) - \boldsymbol{y}_n) \circ \boldsymbol{h}_{\mathsf{L}}'(\boldsymbol{a}_{\mathsf{L}}) & \text{ else} \end{cases}$$

update weights

$$\boldsymbol{W}_{\ell} \leftarrow \boldsymbol{W}_{\ell} - \eta \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{W}_{\ell}} = \boldsymbol{W}_{\ell} - \eta \frac{\partial \mathcal{E}_n}{\partial \boldsymbol{a}_{\ell}} \boldsymbol{o}_{\ell-1}^{\mathrm{T}}$$

Think about how to do the last two steps properly!

More tricks to optimize neural nets

Many variants based on backprop

- SGD with **minibatch**: randomly sample a batch of examples to form a stochastic gradient
- SGD with momentum

• • • •

SGD with momentum

Initialize $oldsymbol{w}_0$ and velocity $oldsymbol{v}=oldsymbol{0}$

For t = 1, 2, ...

- form a stochastic gradient $oldsymbol{g}_t$
- update velocity $m{v} \leftarrow lpha m{v} \eta m{g}_t$ for some discount factor $lpha \in (0,1)$
- update weight $oldsymbol{w}_t \leftarrow oldsymbol{w}_{t-1} + oldsymbol{v}$

SGD with momentum

Initialize $oldsymbol{w}_0$ and velocity $oldsymbol{v}=oldsymbol{0}$

For t = 1, 2, ...

- form a stochastic gradient $oldsymbol{g}_t$
- update velocity $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} \eta \boldsymbol{g}_t$ for some discount factor $\alpha \in (0,1)$
- update weight $oldsymbol{w}_t \leftarrow oldsymbol{w}_{t-1} + oldsymbol{v}$

Updates for first few rounds:

•
$$w_1 = w_0 - \eta g_1$$

• $w_2 = w_1 - \alpha \eta g_1 - \eta g_2$
• $w_3 = w_2 - \alpha^2 \eta g_1 - \alpha \eta g_2 - \eta g_3$
• ...

Overfitting

Overfitting is very likely since the models are too powerful.

Methods to overcome overfitting:

- data augmentation
- regularization
- dropout
- early stopping
- • •

Data augmentation

Data: the more the better. How do we get more data?

Data augmentation

Data: the more the better. How do we get more data?

Exploit prior knowledge to add more training data

Affine Elastic Noise Distortion Deformation Random Horizontal Hue Shift flip Translation

Regularization

L2 regularization: minimize

$$\mathcal{E}'(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \mathcal{E}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) + \lambda \sum_{\ell=1}^{\mathsf{L}} \|\boldsymbol{W}_{\ell}\|_2^2$$

Regularization

L2 regularization: minimize

$$\mathcal{E}'(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \mathcal{E}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) + \lambda \sum_{\ell=1}^{\mathsf{L}} \|\boldsymbol{W}_{\ell}\|_2^2$$

Simple change to the gradient:

$$\frac{\partial \mathcal{E}'}{\partial w_{ij}} = \frac{\partial \mathcal{E}}{\partial w_{ij}} + 2\lambda w_{ij}$$

Regularization

L2 regularization: minimize

$$\mathcal{E}'(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) = \mathcal{E}(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}) + \lambda \sum_{\ell=1}^{\mathsf{L}} \|\boldsymbol{W}_{\ell}\|_2^2$$

Simple change to the gradient:

$$\frac{\partial \mathcal{E}'}{\partial w_{ij}} = \frac{\partial \mathcal{E}}{\partial w_{ij}} + 2\lambda w_{ij}$$

Introduce weight decaying effect

Dropout

Randomly delete neurons during training



Very effective, makes training faster as well

Early stopping

Stop training when the performance on validation set stops improving



Deep neural networks

• are hugely popular, achieving best performance on many problems

- are hugely popular, achieving *best performance* on many problems
- do need a lot of data to work well

- are hugely popular, achieving *best performance* on many problems
- do need a lot of data to work well
- take a lot of time to train (need GPUs for massive parallel computing)

- are hugely popular, achieving *best performance* on many problems
- do need a lot of data to work well
- take *a lot of time* to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters

- are hugely popular, achieving *best performance* on many problems
- do need a lot of data to work well
- take a lot of time to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters
- are still not well understood in theory