Machine learning for particle physicists II. Introduction to neural networks

Anja Butter

26th Vietnam School of Physics

Recap I - Linear Regression

Task: Model calorimeter response

Model:

Loss function:

Optimized parameters:

 $h(\mathbf{x}) = \mathbf{x}\mathbf{w} = \tilde{\mathbf{y}} \leftarrow \text{ prediction}$ $\mathcal{L} = \sum_{i}^{n_{data}} (y_i - \tilde{\mathbf{y}})^2$ $\mathbf{w}_{opt} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Recap II - Logistic Regression

Task: Classify jets

Model:
$$h(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w}) = P(y = 1|\mathbf{x}) \leftarrow$$
 probabilityLoss function: $\mathcal{L} = \sum_{i} -y_i \log(h(\mathbf{x}_i)) - (1 - y_i) \log(1 - h(\mathbf{x}_i))$ Parameter optimization: $\mathbf{w}_n \rightarrow \mathbf{w}_{n+1} = \mathbf{w}_n + \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_n)$

Limitations

So far we were limited to linear separations:



 \rightarrow Need more complex models to learn complex structures

Many possibilities:

boost decision trees $\rightarrow\,$ popular, close to intuitive cut & count analysis

Decision Trees

Example from single top measurement



 \Rightarrow Further information: boosting, random forest

Limitations

So far we were limited to linear separations:



 \rightarrow Need more complex models to learn complex structures

Many possibilities:

boost decision trees $\rightarrow\,$ popular, close to intuitive cut & count analysis

Limitations

So far we were limited to linear separations:



 \rightarrow Need more complex models to learn complex structures

Many possibilities:

boost decision trees \rightarrow popular, close to intuitive cut & count analysis support vector machines \rightarrow often very good baseline neural networks \rightarrow more flexibility, excellent performance

Why neural?

Inspiration from brain



- 1 Dendrites get incoming signal
- 2 Soma processes signal (fire?!)
- 3 Axon transports signal to cells
- 4 Synapse connects to other dendrites

Intelligence is *somehow* the product of the connection of many neurons.

An artificial neuron



$$\rightarrow output = f(\sum_{i} x_i w_i)$$
$$= f(\mathbf{x} \mathbf{w})$$

Combination of linear mapping & non-linear activation function



Every line represents a free parameter called weight



$$\overset{\mathsf{input}}{\boldsymbol{X}} = [n_{data} \times d_{feat}]$$

















If we combine several linear layers...

$$\begin{aligned} \mathbf{Y} &= \mathbf{X} \cdot W_1 \cdot W_2 \cdot w_{out} \\ &= \mathbf{X}_{ij} W_{1,jk} W_{2,kl} W_{out,l} & \leftarrow \text{ Einstein summation convention} \\ &= \mathbf{X}_{ij} \tilde{\mathbf{w}}_j & \text{ with } \tilde{\mathbf{w}} = W_{1,jk} W_{2,kl} W_{out,l} \end{aligned}$$

... we obtain a linear layer!

Not more expressiveness than a simple scalar product!

 \rightarrow Include non-linearities

Activation functions

The right choice can facilitate the training:

- \rightarrow smooth/sharp, limited/unlimited, computing time efficient, \ldots
- Limited (typically for classification)
 - Sigmoid: $\sigma(x) = \frac{1}{1+e-x}$
 - Step: $\theta(x) = \operatorname{sign}(x)$
 - tanh: tanh(x)

1.0 0.5 0.0 -0.5 -1.0 -4 -2 0 2 4

- Unlimited
 - ReLU: max(0, x)
 Leaky ReLU: max(αx, x)

• ELU:
$$\begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x < 0 \end{cases}$$













Finally we have obtained an expressive network! But how can we train it? Gradient descent?!

How to minimize the loss

Remember gradient descent for logistic regression:

$$\rightarrow \boldsymbol{w}_{n+1} = \boldsymbol{w}_n + \alpha \nabla_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}_n)$$

We start from the last layer: Let $\tilde{\pmb{X}}$ be the latent representation after the 2. layer

$$\mathcal{L}(\boldsymbol{w}_{out}) = \sum_{j=1}^{n_{data}} (\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j)^2$$

$$\frac{\partial \mathcal{L}(\boldsymbol{w}_{out})}{\partial w_{out,i}} = \sum_{j=1}^{n_{data}} 2(\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j) \cdot \frac{\partial \sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j}{\partial w_{out,i}} \leftarrow \text{ chain rule}$$

$$= \sum_{j=1}^{n_{data}} 2(\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j) \cdot \frac{\partial \sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j}{\partial (\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j} \cdot \frac{\partial (\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j}{\partial w_{out,i}}$$

How to minimize the loss

Remember gradient descent for logistic regression:

$$\rightarrow \boldsymbol{w}_{n+1} = \boldsymbol{w}_n + \alpha \nabla_{\boldsymbol{w}} \mathcal{L}(\boldsymbol{w}_n)$$

We start from the last layer: Let $\tilde{\pmb{X}}$ be the latent representation after the 2. layer

$$\mathcal{L}(\boldsymbol{w}_{out}) = \sum_{j=1}^{n_{data}} (\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j)^2$$

$$\frac{\partial \mathcal{L}(\boldsymbol{w}_{out})}{\partial w_{out,i}} = \sum_{j=1}^{n_{data}} 2(\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j) \cdot \frac{\partial \sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j}{\partial w_{out,i}} \leftarrow \text{ chain rule}$$

$$= \sum_{j=1}^{n_{data}} 2(\sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j - Y_j) \cdot \frac{\partial \sigma(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j}{\partial(\tilde{\boldsymbol{X}} \boldsymbol{w}_{out})_j} \cdot \tilde{\boldsymbol{X}}_{ji}$$

Backpropagation



output			
w _{out} =	[d ₁₂	×	1]

 $ightarrow rac{\partial \mathcal{L}}{\partial \sigma_3} rac{\partial \sigma_3}{\partial \textit{lin3}}$

$$[W_2 = [d_{l1} \times d_{l2}] \leftarrow \texttt{output} \\ w_{out} = [d_{l2} \times 1]$$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \sigma_{3}} \frac{\partial \sigma_{3}}{\partial lin3} \cdot \frac{\partial lin3}{\partial \sigma_{2}} \frac{\partial \sigma_{2}}{\partial lin2}$$

$$\begin{matrix} \text{layer 1} \\ W_1 = [d_{feat} \times d_{l1}] \end{matrix} \xleftarrow{} \begin{matrix} \text{layer 2} \\ W_2 = [d_{l1} \times d_{l2}] \end{matrix} \xleftarrow{} \begin{matrix} \text{output} \\ w_{out} = [d_{l2} \times 1] \end{matrix}$$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \sigma_{3}} \frac{\partial \sigma_{3}}{\partial lin3} \cdot \frac{\partial lin3}{\partial \sigma_{2}} \frac{\partial \sigma_{2}}{\partial lin2} \cdot \frac{\partial lin2}{\partial \sigma_{1}} \frac{\partial \sigma_{1}}{\partial lin1}$$

$$\begin{array}{c} \text{input} \\ \textbf{X} = [n_{data} \times d_{feat}] \end{array} \begin{array}{c} \text{layer 1} \\ \textbf{W}_1 = [d_{feat} \times d_{l1}] \end{array} \begin{array}{c} \text{layer 2} \\ \textbf{W}_2 = [d_{l1} \times d_{l2}] \end{array} \begin{array}{c} \text{output} \\ \textbf{w}_{out} = [d_{l2} \times 1] \end{array}$$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \sigma_{3}} \frac{\partial \sigma_{3}}{\partial lin3} \cdot \frac{\partial lin3}{\partial \sigma_{2}} \frac{\partial \sigma_{2}}{\partial lin2} \cdot \frac{\partial lin2}{\partial \sigma_{1}} \frac{\partial \sigma_{1}}{\partial lin1} \cdot \frac{\partial lin1}{\partial W_{1}}$$

$$\begin{tabular}{|c|c|c|c|} \hline layer 1 & layer 2 \\ \hline \textbf{X} = [n_{data} \times d_{feat}] & \textbf{W}_1 = [d_{feat} \times d_{l1}] & \textbf{W}_2 = [d_{l1} \times d_{l2}] & \textbf{w}_{out} = [d_{l2} \times 1] \\ \hline \end{tabular}$$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \sigma_{3}} \frac{\partial \sigma_{3}}{\partial lin3} \cdot \frac{\partial lin3}{\partial \sigma_{2}} \frac{\partial \sigma_{2}}{\partial lin2} \cdot \frac{\partial lin2}{\partial \sigma_{1}} \frac{\partial \sigma_{1}}{\partial lin1} \cdot \frac{\partial lin1}{\partial W_{1}}$$

Latest winner of ImageNet 'ViT-H/14' has 32 layers with a total of 632M parameters Let's start calculating...?

Automatic gradient computations

Backpropagation implemented via computation graphs in dedicated frameworks



Ayoosh Kathuria

High-level ML frameworks



- TensorFlow & PyTorch \rightarrow DNN
- Keras \rightarrow user friendly TF interface
- Scikit-learn \rightarrow SVM, BDT, clustering
- Spark ML \rightarrow part of Spark, basic ML
- Hugging face \rightarrow NLP, transformers

Which framework to choose for NN?

Biggest players: PyTorch (facebook) and TensorFlow (google)



Weltweit. Letzte 5 Jahre. Websuche.

Before you start training your own neural network

Let's talk about:

 \rightarrow efficient training

 \rightarrow overtraining (overfitting)

Stochastic gradient descent

- Loss defined on entire dataset
- Each weight update: gradient for full dataset
 - \rightarrow very computing expensive!



Stochastic gradient descent

- Loss defined on entire dataset
- Each weight update: gradient for full dataset → very computing expensive!
- SGD: 1 iteration: gradient for 1 random data point
- Compromise: batch gradient descent (batch size: 32/62/.../1024)



- + faster
- + less sensitive to local minima
- + profit from vectorization

Overtraining



... when networks start to learn "noise"

How to control overtraining?



- Training loss \nearrow validation loss $\searrow \Rightarrow$ overfitting
- Why do we need test data?
- How would you split your dataset? 1:1:1? 8:1:1?

How to control overtraining?



- Training loss \nearrow validation loss $\searrow \Rightarrow$ overfitting
- Why do we need test data?
 - parameter tests = training on validation dataset
- How would you split your dataset? 1:1:1? 8:1:1?

How to control overtraining?



- Training loss \nearrow validation loss $\searrow \Rightarrow$ overfitting
- Why do we need test data?
 - parameter tests = training on validation dataset
- How would you split your dataset? 1:1:1? 8:1:1?
 - dataset dependent
 - val/test data large enough to test performance

Today's summary

We saw:

- How to build a network with multiple layers
- Why we need activation functions in each layer
- How to train a deep neural network \rightarrow Backpropagation
- Stochastic/batch gradient descent
- Overtraining \rightarrow Train validate test

Now you can start training your own neural network!

Procrastination over the course of time



credits: XKCD & u/AmpyeriDracula