
Loss functions, optimizers and tricks to train a NN

Komal Teru
Mila, McGill University



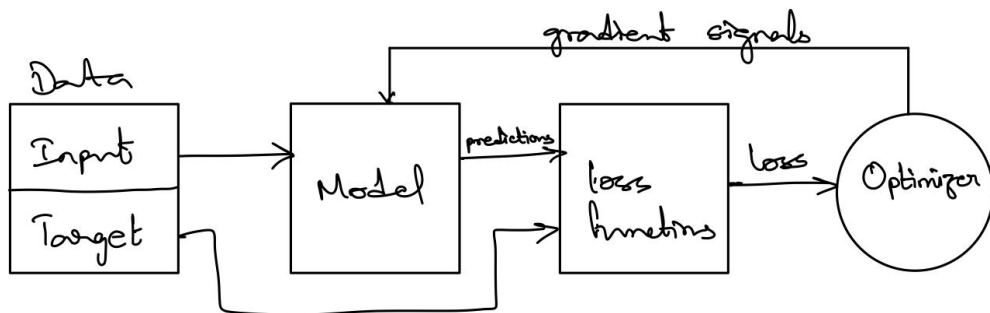
Topics to cover today

- **Review of MLP**
- **More types of optimizers and loss functions**
- **Tricks to train an MLP**

Review of MLP

Key ingredients of ML pipeline

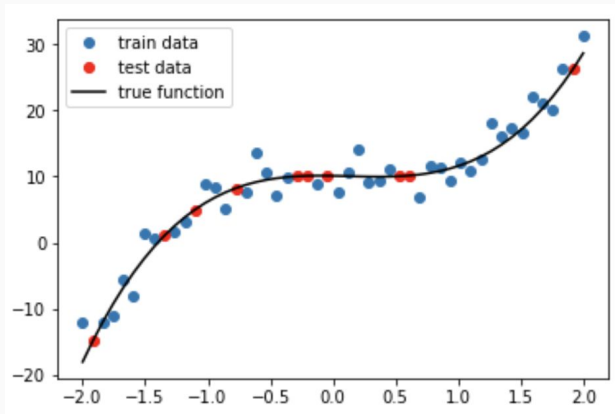
- Data : image, text, speech, graph etc.
- Model : MLP, CNN, RNN, etc.
- Loss function : MSE, cross-entropy, etc.
- Optimization algorithm (optimizer) : SGD, Adam, Adagrad, etc.



Review of MLP

Multi layer perceptron for regression

Data



$$x \in \mathbb{R}$$

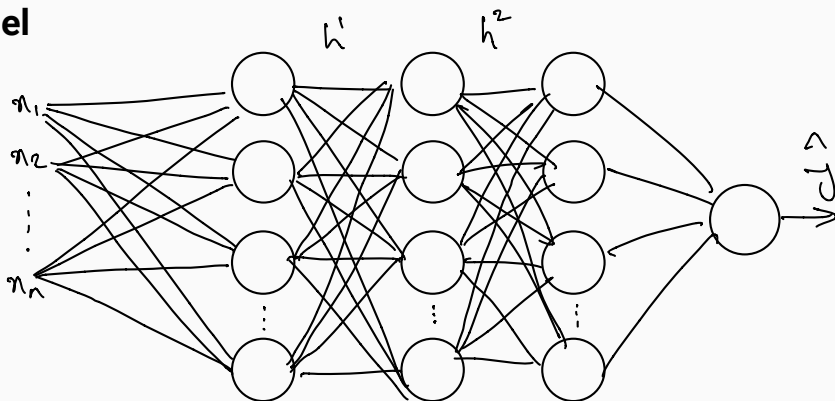
$$y \in \mathbb{R}$$

$$\mathcal{D} = (x, y)$$

Review of MLP

Multi layer perceptron for regression

Model



$$h^1 = f(x_i w^1 + b^1)$$

$$h^2 = f(h^1 w^2 + b^2)$$

$$\hat{y}_i = \underline{\underline{h^2 w^3 + b^3}}$$

$$mse = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Loss function : mean squared error **Optimizer:** SGD

Review of MLP

Multi layer perceptron for regression

Key points:

- The final layer doesn't have activation function. (why?)
- Bias is necessary to adjust any offset in the data.
- Activation function is necessary to fit to non-linear data.

Optimizers

SGD algorithm

Algorithm 1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$
 with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{h}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{h}}$

end while

Optimizers

Momentum method

Algorithm 1 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\mathbf{h} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{h}$

 Apply update: $\theta \leftarrow \theta + \mathbf{v}$

end while

Optimizers

Adagrad algorithm

Algorithm 1 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient: $\mathbf{h} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{h} \odot \mathbf{h}$

Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{h}$. (Division and square root applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

Optimizers

RMSProp method

Algorithm 1 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ .

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers.

Initialize accumulation variables $\mathbf{r} = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{h} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{h} \odot \mathbf{h}$

 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{h}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied elem-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Optimizers

Adam method

Algorithm 1 The Adam algorithm

17

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggestion: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$
 with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{h} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

 Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{h}$

 Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{h} \odot \mathbf{h}$

 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

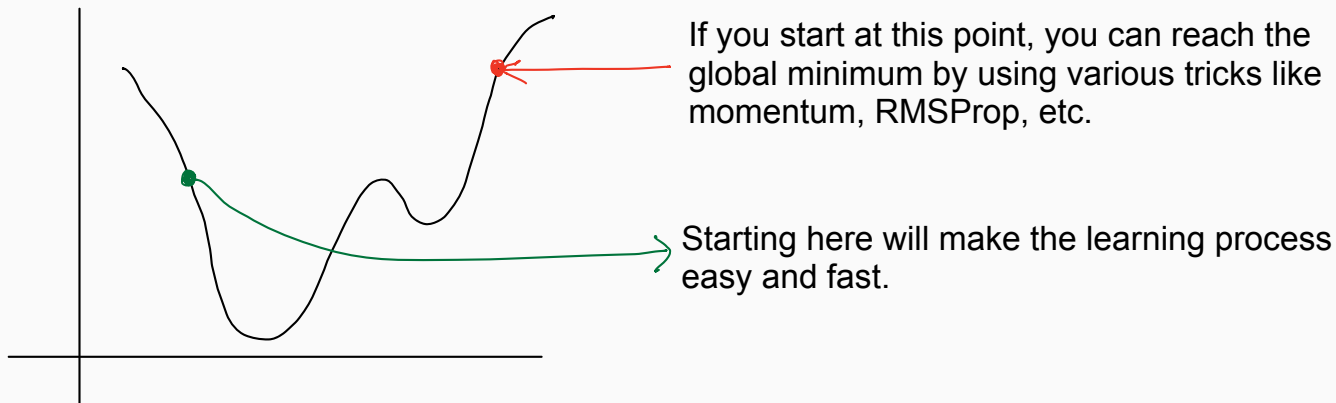
 Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Xavier/Glorot Initialization

Initialization of weights is key to the training dynamics.



Xavier/Glorot Initialization

These initializations put you in a better position to start your optimization

Xavier uniform initialization

$$\theta \sim \mathcal{U}(-a, a)$$
$$a = \sqrt{\frac{6}{f_{in} + f_{out}}} \approx \sqrt{\frac{3}{f_{in}}}$$

Xavier normal initialization

$$\theta = \mathcal{N}(0, \sigma^2)$$
$$\sigma = \sqrt{\frac{2}{f_{in} + f_{out}}} \approx \sqrt{\frac{1}{f_{in}}}$$

In class exercise

We feed a 4 dimensional data to a 1 hidden layer MLP with 30 nodes/perceptrons. The output has a dimension of 4. How many trainable parameters does the model have?

L2 Regularization

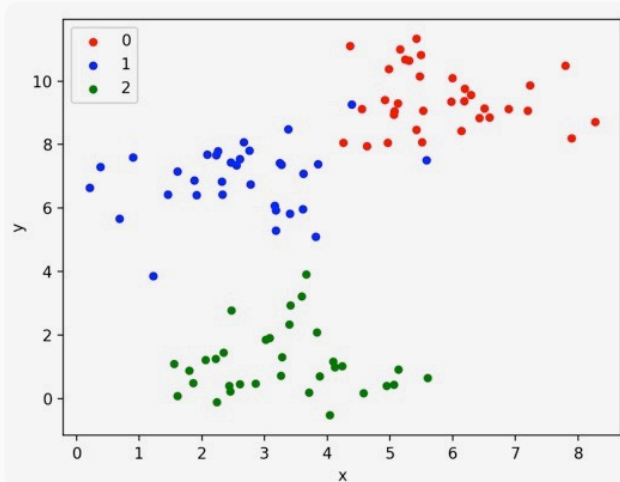
$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$$

↓
weight penalty.

L2 regularization limits the models capacity by limiting the range of values it's parameters can take.

Classification model

Data



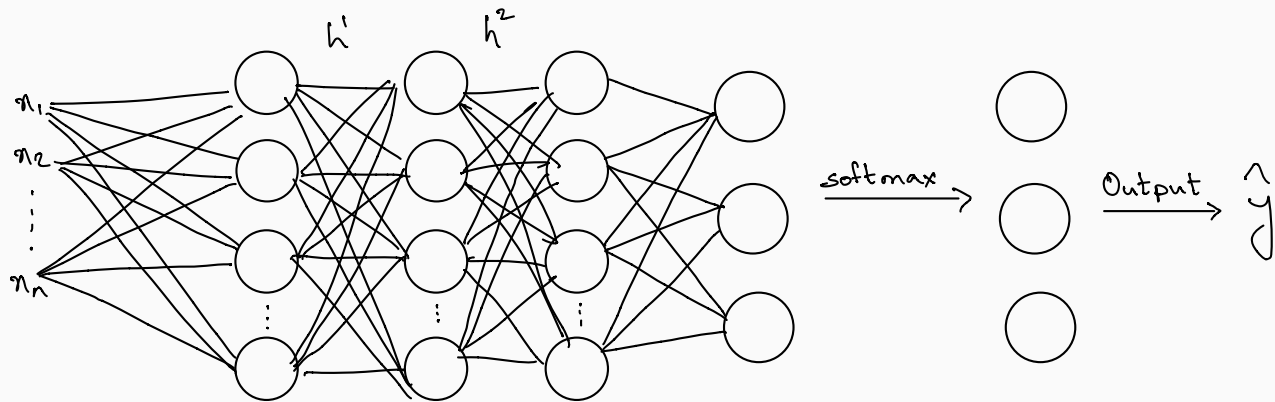
$$X_i \in \mathbb{R}, O_i \in \{0, 1, 2\}$$

$$\text{and } \gamma_{i-hot} \in \mathbb{B}^3 \text{ where } \mathbb{B} = \{0, 1\}$$

$$\begin{aligned} \text{Eg: IF } O_i = 0 &\Rightarrow y_i^{1-hot} = [1, 0, 0] \\ O_i = 1 &\Rightarrow y_i^{1-hot} = [0, 1, 0] \\ O_i = 2 &\Rightarrow y_i^{1-hot} = [0, 0, 1] \end{aligned}$$

Classification model

Model



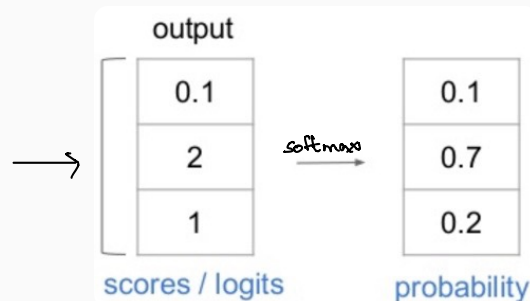
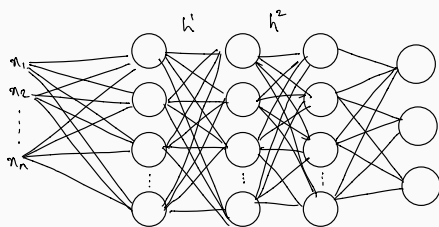
Loss function : cross-entropy **Optimizer:** SGD

Classification model

Softmax function

scores (logits)

$$S(l_i) = \frac{e^{l_i}}{\sum_k e^{l_k}}$$



Classification model

Cross-entropy

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^m -\log \hat{y}_i^{(o_i)}$$

