Deep Learning and Its Application to Signal and Image Processing and Analysis

CLASS III - SPRING 2022

TAMMY RIKLIN RAVIV, ELECTRICAL AND COMPUTER ENGINEERING

Today's topics

- Stochastic Gradient Descent
- Backpropagation
- •Vanishing/Exploiding gradients
- Optimization & Optimizers
- •Training, Evaluation & Test
- Batch Normalization
- •Dropout









Optimization

Cost function values: the discrepancies between the outputs (NN estimations) and the training set data points.

Goal : find the set of weights for which a global minimum Is obtained

the cost function is parameterized by the network's weights — we control our loss function by changing the weights.

In reality the cost function is not convex.

















$$\label{eq:Gradient descent} \begin{split} & W\colon = W - \alpha \frac{\partial L(W)}{\partial W} \quad \text{Update rule} \\ & W\colon = W - \alpha \nabla_W L(W) \\ & & \text{Learning rate} \end{split} \\ \\ & \text{Learning rate: An important hyperparameter} \\ & \text{too small - very slow convergence or gets stuck in local minima} \\ & \text{Too Big - may "skip" the target minimum; may go in the wrong direction} \end{split}$$

Backpropagation to train multilayer architectures

The backpropagation procedure to computes the gradient of an objective function with respect to the weights of a multilayer stack of modules.

Practical application of the chain rule.

The key insight is that the gradient of the objective with respect to the input of a module can be computed by working backwards from the gradient with respect to the output of that module (or the input of the subsequent module).

The backpropagation equation can be **applied repeatedly** to **propagate gradients** through all modules, starting from the output at the top (where the network produces its prediction) all the way to the bottom (where the external input is fed).

Once these gradients have been computed, it is straightforward to compute the gradients with respect to the weights of each module.

16

Chain rule of calculus

x - a real number

 $f() \ {\rm and} \ g()$ - functions mapping from a real number to a real number.

$$y = g(x)$$
 and $z = f(g(x)) = f(y)$

chain rule:

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

















Backpropagation - Limitations

Gradient descent with backpropagation is not guaranteed to find the global minimum of the error function, but only a local minimum

It has trouble crossing plateaux in the error function landscape.

This issue, caused by the non-convexity of error functions in neural networks, was long thought to be a major drawback, but in a 2015 review article, Yann LeCun *et al.* (Deep Learning, Nature) argue that in many practical problems, it is not.

 Backpropagation learning does not require normalization of input vectors; however, normalization could improve performance

Modes of learning

There are two modes of learning to choose from: stochastic and batch.

Stochastic learning: each propagation is followed immediately by a weight update.

Batch learning: many propagations occur before updating the weights, accumulating errors over the samples within a batch.

Stochastic learning introduces "noise" into the gradient descent process, using the local gradient calculated from one data point; this reduces the chance of the network getting stuck in a local minima. Yet batch learning typically yields a faster, more stable descent to a local minima, since each update is performed in the direction of the average error of the batch samples. In modern applications a common compromise choice is to use "mini-batches", meaning batch learning but with a batch of small size and with stochastically selected samples.

Training data collections

Online learning is used for dynamic environments that provide a continuous stream of new training data patterns.

Offline learning makes use of a training set of static patterns.

Internal Covariate Shift

Internal Covariate Shift is defined as the change in the distribution of network

activations due to the change in network parameters during training.

Batch normalization is a method intended to mitigate internal covariate shift for neural networks.



https://medium.com/analytics-vidhya/internal-covariate-shift-anoverview-of-how-to-speed-up-neural-network-training-3e2a3dcdd5cc

Internal Covariate Shift

Internal Covariate Shift is defined as the change in the distribution of network activations due to the change in network parameters during training.

Deeper networks are more affected by the internal covariate shift.

Key idea: stabilize the input values for each layer (defined as **z** = **Wx** + **b**, where

z is the linear transformation of the W weights/parameters and the biases).













Batch Normalization

we demand from our features to follow a Gaussian distribution

、

with zero mean and unit variance.

$$BN(x) = \gamma(rac{x-\mu(x)}{\sigma(x)}) + eta \ N = rac{W}{W}$$

$$\mu_{\mathcal{C}}(x) = rac{1}{NHW}\sum_{n=1}^{N}\sum_{h=1}^{H}\sum_{w=1}^{W}x_{nchw}$$

$$\sigma_{c}(x) = \sqrt{rac{1}{NHW}\sum\limits_{n=1}^{N}\sum\limits_{h=1}^{H}\sum\limits_{w=1}^{W}(x_{nchw}-\mu_{c}(x))^{2}}$$



Layer Normalization

In BN, the statistics are computed across the batch and the spatial dims.

In contrast, in **Layer Normalization** (**LN**), the statistics (mean and variance) are **computed across all channels and spatial dims**. Thus, the statistics are independent of the batch. This layer was initially introduced to handle vectors (mostly the RNN outputs).



Instance Normalization (IN) is computed only across the features' spatial dimensions. It is independent for each channel and sample.



Instance Normalization

The affine parameters in **IN can completely change the style** of the output image. As opposed to BN, **IN can normalize the style of each individual sample to a target style** (modeled by γ and β). For this reason, training a model to transfer to a specific style is easier.





Momentum

- Accelerate learning, especially in the face of high curvature, small but consistent gradients, or noisy gradients.
- Accumulates an exponentially decaying moving average of past gradients and continues to move in their direction.
- Addresses **ill-conditioning, i.e.,** when SGD gets "stuck" in the sense that even very small steps increase the cost function





SGD with Momentum

Require: Learning rate ϵ , momentum parameter α **Require:** Initial parameter $\boldsymbol{\theta}$, initial velocity \boldsymbol{v} while stopping criterion not met do Sample a minibatch of m examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding targets $\boldsymbol{y}^{(i)}$. Compute gradient estimate: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$. Compute velocity update: $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \boldsymbol{g}$. Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$. end while

The size of the step depends on how large and how aligned a sequence of gradients are





$$oldsymbol{v} \leftarrow lpha oldsymbol{v} - \epsilon
abla_{oldsymbol{ heta}} \left[rac{1}{m} \sum_{i=1}^m L\left(oldsymbol{f}(oldsymbol{x}^{(i)};oldsymbol{ heta} + oldsymbol{lpha} oldsymbol{v},oldsymbol{y}^{(i)}
ight)
ight]$$

 $oldsymbol{ heta} \leftarrow oldsymbol{ heta} + oldsymbol{v},$

- The gradient is evaluated after the current velocity is applied
- Add a correction factor to the standard method of momentum.
- Improves the rate of convergence in Convex batch gradient case

SGD with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α Require: Initial parameter θ , initial velocity vwhile stopping criterion not met do Sample a minibatch of m examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding labels $\boldsymbol{y}^{(i)}$. Apply interim update: $\tilde{\boldsymbol{\theta}} \leftarrow \boldsymbol{\theta} + \alpha \boldsymbol{v}$. Compute gradient (at interim point): $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\boldsymbol{\theta}}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \tilde{\boldsymbol{\theta}}), \boldsymbol{y}^{(i)})$. Compute velocity update: $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \boldsymbol{g}$. Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$. end while

Algorithms with Adaptive Learning Rates

- Assumption: directions of sensitivity are somewhat axis aligned
- Use a separate learning rate for each parameter
- · Automatically adapt these learning rates throughout the course of learning





The AdaGrad Algorithm

Require: Global learning rate ϵ **Require:** Initial parameter $\boldsymbol{\theta}$ **Require:** Small constant δ , perhaps 10^{-7} , for numerical stability Initialize gradient accumulation variable $\boldsymbol{r} = \boldsymbol{0}$ **while** stopping criterion not met **do** Sample a minibatch of m examples from the training set $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ with corresponding targets $\boldsymbol{y}^{(i)}$. Compute gradient: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$. Accumulate squared gradient: $\boldsymbol{r} \leftarrow \boldsymbol{r} + \boldsymbol{g} \odot \boldsymbol{g}$. Compute update: $\Delta \boldsymbol{\theta} \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot \boldsymbol{g}$. (Division and square root applied element-wise) Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$. **end while**





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 r = 0while stopping criterion not met do Sample a minibatch of *m* examples from the training set { $x^{(1)}, \ldots, x^{(m)}$ } with corresponding targets $y^{(i)}$. Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i} L(f(x^{(i)}; \theta), y^{(i)})$. Accumulate squared gradient: $r \leftarrow \rho r + (1 - \rho)g \odot g$. Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$. $(\frac{1}{\sqrt{\delta + r}}$ applied element-wise) Apply update: $\theta \leftarrow \theta + \Delta \theta$. end while

RMSProp algorithm with Nesterov momentum

Require: Global learning rate ϵ , decay rate ρ , momentum coefficient α **Require:** Initial parameter θ , initial velocity vInitialize accumulation variable r = 0while stopping criterion not met do Sample a minibatch of m examples from the training set $\{x^{(1)}, \ldots, x^{(m)}\}$ with corresponding targets $y^{(i)}$. Compute interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$. Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_{i} L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$. Accumulate gradient: $r \leftarrow \rho r + (1 - \rho)g \odot g$. Compute velocity update: $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$. $(\frac{1}{\sqrt{r}}$ applied element-wise) Apply update: $\theta \leftarrow \theta + v$. end while

70

Adam Optimizer (Adaptive Moment)

- A variant on the combination of RMSProp and momentum with a few important distinctions.
- · Momentum is incorporated directly as an estimate of the first-order moment

(with exponential weighting) of the gradient.

- Momentum is applied to the rescaled gradients
- Includes bias corrections to the estimates of both the first-order moments (the momentum term) and the (uncentered) second-order moments to account for their initialization at the origin

| | 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 (Suggested default: |
| ADAM | Require: Initial parameters $\boldsymbol{\theta}$ |
| | Initialize 1st and 2nd moment variables $s = 0, r = 0$ |
| Algorithm | Initialize time step $t = 0$ |
| | while stopping criterion not met do |
| | Sample a minibatch of m examples from the training set $\{x^{(1)}, \ldots, x^{(m)}\}$ with corresponding targets $y^{(i)}$. |
| | Compute gradient: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$ |
| | $t \leftarrow t + 1$ |
| | Update biased first moment estimate: $\boldsymbol{s} \leftarrow \rho_1 \boldsymbol{s} + (1 - \rho_1) \boldsymbol{g}$ |
| | Update biased second moment estimate: $\boldsymbol{r} \leftarrow \rho_2 \boldsymbol{r} + (1 - \rho_2) \boldsymbol{g} \odot \boldsymbol{g}$ |
| | Correct bias in first moment: $\hat{s} \leftarrow \frac{s}{1-\rho_1^t}$ |
| | Correct bias in second moment: $\hat{r} \leftarrow \frac{1}{1-\rho_2^t}$ |
| | Compute update: $\Delta \theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r}+\delta}}$ (operations applied element-wise) |
| | Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$ |
| | end while |
| | |

NADAM

Nadam is an extension of the Adam version of gradient descent that incorporates Nesterov momentum

