

Architectures

Strategy

1. Select Network Structure appropriate for problem

2. Check for implementation bugs with gradient checks

3. Parameter initialization

4. Optimization

5. Check if the model is powerful enough to overfit

Structure: Single words, fixed window, sentence based, context level; bag of words, recursive vs. recurrent, CNN

Nonlinearity (Activation Functions)

1. Implement your gradient

2. Implement a finite difference computation by looping through the parameters of your network, adding and subtracting a small epsilon (~10^-4) and estimate derivatives

3. Compare the two and make sure they are almost the same

$$f'(\theta) \approx \frac{J(\theta^{(i+1)}) - J(\theta^{(i-1)})}{2\epsilon}$$

$$\theta^{(i+1)} = \theta + \epsilon \times \epsilon_i$$

If you gradient fails and you don't know why?

Simplify your model until you have no bug!

What now? Create a very tiny synthetic model and dataset

Only softmax on fixed input

Backprop into word vectors and softmax

single unit single hidden layer

Add multi unit single layer

Add second layer single unit, add multiple units, bias - Add one softmax on top, then two softmax layers

Add bias

Example: Start from simplest model then go to what you want:

Initialize hidden layer biases to 0 and output (or reconstruction) biases to optimal value if weights were 0 (e.g., mean target or inverse sigmoid of mean target).

Initialize weights ~ Uniform(-r, r), r inversely proportional to fan-in (previous layer size) and fan-out (next layer size):

$$\sqrt{6/(\text{fan-in} + \text{fan-out})}$$

Is a first-order iterative optimization algorithm for finding the minimum of a function. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point. If instead one takes steps proportional to the positive of the gradient, one approaches a local maximum of that function; the procedure is then known as gradient ascent.

Gradient descent uses total gradient over all examples per update, SGD updates after only 1 or few examples:

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_t(\theta)$$

Ordinary gradient descent as a batch method is very slow, should never be used. Use 2nd order batch method such as L-BFGS

On large datasets, SGD usually wins over all batch methods. On smaller datasets L-BFGS or Conjugate Gradients win. Large-batch L-BFGS extends the reach of L-BFGS [Le et al. ICML 2011].

Gradient descent uses total gradient over all examples per update, SGD updates after only 1 example

Most commonly used now, Size of each mini batch B: 20 to 1000

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J_{t:t+B}(\theta)$$

Helps parallelizing any model by computing gradients for multiple elements of the batch in parallel

Idea: Add a fraction v of previous update to current one. When the gradient keeps pointing in the same direction, this will increase the size of the steps taken towards the minimum.

Reduce global learning rate when using a lot of momentum

Update Rule

$$v = \mu v - \alpha \nabla_{\theta} J_t(\theta)$$

$$\theta^{new} = \theta^{old} + v$$

v is initialized at 0

Momentum often increased after some epochs (0.5 to 0.99)

Adaptive learning rates for each parameter!

Learning rate is adapting differently for each parameter and rare parameters get larger updates than frequently occurring parameters. Word vectors!

Let $g_{t,i} = \frac{\partial}{\partial \theta_i} J_t(\theta)$, then: $\theta_{t,i} = \theta_{t-1,i} - \frac{\alpha}{\sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}} g_{t,i}$

If not, change model structure or make model "larger"

Simple first step: Reduce model size by lowering number of units and layers and other parameters

Standard L1 or L2 regularization on weights

Early Stopping: Use parameters that gave best validation error

Sparsity constraints on hidden activations, e.g., add to cost:

Training time: at each instance of evaluation (in online SGD-training), randomly set 50% of the inputs to each neuron to 0.

Test time: halve the model weights (now twice as many!) This prevents feature co-adaptation: A feature cannot only be useful in the presence of particular other features

Dropout

In a single layer: A kind of middle-ground between Naive Bayes (where all feature weights are set independently) and logistic regression models (where weights are set in the context of all others)

Can be thought of as a form of model bagging

It also acts as a strong regularizer

Feed Forward

Kinds

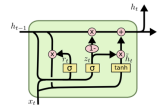
Single-Layer Perceptron

The inputs are fed directly to the outputs via a series of weights. By adding an Logistic activation function to the outputs, the model is identical to a classical Logistic Regression model.

Multi-Layer Perceptron

This class of networks consists of multiple layers of computational units, usually interconnected in a feed-forward way. Each neuron in one layer has directed connections to the neurons of the subsequent layer. In many applications the units of these networks apply a sigmoid function as an activation function.

LSTMs

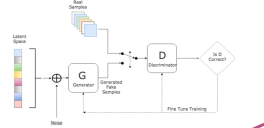


$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$
$$\tilde{h}_t = \tanh(W \cdot [r_t \cdot h_{t-1}, x_t])$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

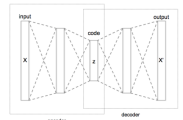
Long short-term memory - It is a type of recurrent (RNN), allowing data to flow both forwards and backwards within the network.

GANs

GANs or Generative Adversarial Networks are a class of artificial intelligence algorithms used in unsupervised machine learning, implemented by a system of two neural networks contesting with each other in a zero-sum game framework.

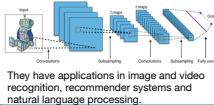


Auto-Encoders



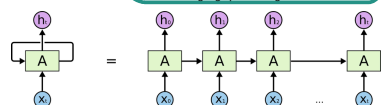
Is an artificial neural network used for unsupervised learning of efficient codings.

Convolutional Neural Networks (CNN)



They have applications in image and video recognition, recommender systems and natural language processing.

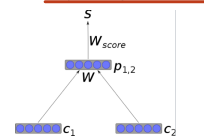
RNNs (Recurrent)



Is a class of artificial neural network where connections between units form a directed cycle. This allows it to exhibit dynamic temporal behavior. Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs.

This makes them applicable to tasks such as unsegmented, connected handwriting recognition or speech recognition.

RNNs (Recursive)



Is a kind of deep neural network created by applying the same set of weights recursively over a structure, to produce a structured prediction over variable-size input structures, or a scalar prediction on it, by traversing a given structure in topological order.

RNNs have been successful for instance in learning sequence and tree structures in natural language processing, mainly phrase and sentence continuous representations based on word embedding.

