1

Denoising Autoencoders

Sargur Srihari srihari@buffalo.edu

Topics in Autoencoders

- What is an autoencoder?
- 1. Undercomplete Autoencoders
- 2. Regularized Autoencoders
- 3. Representational Power, Layout Size and Depth
- 4. Stochastic Encoders and Decoders
- 5. Denoising Autoencoders
- 6. Learning Manifolds and Autoencoders
- 7. Contractive Autoencoders
- 8. Predictive Sparse Decomposition
- 9. Applications of Autoencoders

5. Denoising Autoencoders

- An autoencoder that receives a corrupted data point as input and is trained to predict the original, uncorrupted data point as its output
 - Traditional autoencoders minimize L(x, g(f(x)))
 - where *L* is a loss function penalizing g(f(x)) for being dissimilar from *x*, such as L^2 norm of difference: mean squared error
 - A DAE minimizes $L(\boldsymbol{x}, g(f(\tilde{\boldsymbol{x}})))$
 - where \tilde{x} is a copy of x that is corrupted by some form of noise
 - The autoencoder must undo this corruption rather than simply copying their input

Example of Noise in a DAE

- An autoencoder with high capacity can end up learning an identity function (also called null function) where input=output
- A DAE can solve this problem by corrupting the data input
- How much noise to add?
 - Corrupt input nodes by setting 30-50% of random input nodes to zero





Original input, corrupted data, reconstructed data

DAE Training procedure

- Computational graph of cost function below
 - DAE trained to reconstruct clean data point *x* from the corrupted Accomplished by minimizing loss *L*=-log *p*_{encoder}(*x*|*h*=*f*(*x*))



Corruption process, $C(\tilde{x} | x)$ is a conditional distribution over corrupted samples \tilde{x} given the data sample x

The autoencoder learns a reconstruction distribution $p_{\text{reconstruct}}(\mathbf{x} | \ \tilde{\mathbf{x}})$ estimated from training pairs $(\mathbf{x}, \tilde{\mathbf{x}})$ as follows:

1 Sample a training sample *x* from the training data 2. Sample a corrupted version \tilde{x} from $C(\tilde{x} | x)$ 3. Use (x, \tilde{x}) as a training example for estimating the autoencoder distribution $p_{\text{recoconstruct}}(x | \tilde{x}) = p_{\text{decoder}}(x | h)$ with *h* the output of encoder $f(\tilde{x})$ and p_{decoder} typically defined by a decoder g(h)

• DAE performs SGD on the expectation $E_{\tilde{\boldsymbol{x}} \sim p^{\wedge} data(\boldsymbol{x})} \log p_{decoder}(\boldsymbol{x}|\boldsymbol{h}=\boldsymbol{f}(\tilde{\boldsymbol{x}}))$

Deep Learning

DAE for MNIST data

Python/Theano

•	import theano.tensor as T	7616
•	from opendeep.models.model import Model	1400
•	from opendeep.utils.nnet import get_weights_uniform, get_bias	
•	from opendeep.utils.noise import salt_and_pepper	D69(
•	from opendeep.utils.activation import tanh, sigmoid	5070
•	from opendeep.utils.cost import binary_crossentropy	24/1
•	# create our class initialization!	9665
•	class DenoisingAutoencoder(Model):	
•	nn	Unsuperv
•	A denoising autoencoder will corrupt an input (add noise) and try to reconstruct it.	
•	definit(self):	
•	# Define some model hyperparameters to work with MNIST images!	
•	input_size = 28*28 # dimensions of image	
•	hidden_size = 1000 # number of hidden units - generally bigger than input size for	DAE
•	# Now, define the symbolic input to the model (Theano)	
•	# We use a matrix rather than a vector so that minibatch processing can b	e done in parallel.
•	x = T.fmatrix("X")	
•	self.inputs = [x]	
•	# Build the model's parameters - a weight matrix and two bias vectors	
•	W = get_weights_uniform(shape=(input_size, hidden_size), name="W")	
•	b0 = get_bias(shape=input_size, name="b0")	
•	b1 = get_bias(shape=hidden_size, name="b1")	
•	self.params = [W, b0, b1]	
•	# Perform the computation for a denoising autoencoder!	
•	# first, add noise (corrupt) the input	
•	corrupted_input = salt_and_pepper(input=x, corruption_level=0.4)	
•	# next, compute the hidden layer given the inputs (the encoding function)	
•	hiddens = tanh(T.dot(corrupted_input, W) + b1)	
•	# finally, create the reconstruction from the hidden layer (we tie the weight	s with W.T)
•	reconstruction = sigmoid(T.dot(hiddens, W.T) + b0)	-
•	# the training cost is reconstruction error - with MNIST this is binary cross-	entropy
•	self.train_cost = binary_crossentropy(output=reconstruction, target=x)	



Dervised Denoising Autoencoder Left: original test images Center: corrupted noisy images Right: reconstructed images

Srihari

Estimating the Score

- An autoencoder can be based on encouraging the model to have the same score as the data distribution at every training point x
 - The *score* is a particular gradient field is: $\nabla_x \log p(x)$
 - Learning the gradient field of $\log p_{\rm data}$ is one way to learn the structure of $p_{\rm data}$ itself
- Score Matching works by fitting the slope (score) of the model density to the slope of the true underlying density at the data points
- DAE, with conditionally Gaussian p(x|h), estimates this score as (g(f(x)-x)
 - The DAE is trained to minimize $||g(f(\hat{x}-x))||^2$
- DAE estimates a vector fields as illustrated next

Deep Learning

DAE learns a vector field

- Training examples x lie on a low-dimensional manifold
 - Training examples *x* are red crosses
- Gray circle is equiprobable corruptions
- The vector field (g(f(x)-x), indicated by green arrows, estimates the score ∇_x log p(x) which is the slope of the density of data



Vector field learnt by a DAE

- 1-D curved manifold near which the data concentrate
- Each arrow proportional to reconstruction minus input vector of DAE and points towards higher probability
- Where probability is maximum arrows shrink

