

Deep learning: Resources

- Goodfellow, Courville and Bengio: *Deep Learning*, MIT Press 2016 http://www.deeplearningbook.org/ https://github.com/janishar/mit-deep-learning-book-pdf
- P. Mehta, M. Bukov, C.-H. Wang, A.G.R. Day, C. Richardson, C.K. Fisher, D.J. Schwab: A high-bias, low-variance introduction to Machine Learning for physicists https://arxiv.org/abs/1803.08823













WaveNet







WaveNet



Describes without errors



A person riding a motorcycle on a dirt road.





Two dogs play in the grass.

Somewhat related to the image



A skateboarder does a trick on a ramp.

A little girl in a pink hat is blowing bubbles.





A refrigerator filled with lots of food and drinks.



A yellow school bus parked in a parking lot.

Freie Universität



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.

A dog is jumping to catch a frisbee.



Supervised learning with feedforward neural networks















Activation functions

https://en.wikipedia.org/wiki/Activation_function





Activation functions

https://en.wikipedia.org/wiki/Activation_function





Activation functions

https://en.wikipedia.org/wiki/Activation_function



Currently most widely used. Empirically easier to train and results in sparse networks.

Nair and Hinton: Rectified linear units improve restricted Boltzmann machines ICLM'10, 807-814 (2010) Glorot, Bordes and Bengio: Deep sparse rectifier neural networks. PMLR 15:315-323 (2011) Freie Universität

Perceptron (Rosenblatt 1957)

Used as a binary classifier - equivalent to support vector machine (SVM)





Perceptron (Rosenblatt 1957)





Freie Universität

Multilayer perceptrons (~1985)



- Network architecture is typically the most complex model decision:
 - Number of hidden layers L-1 and number of neurons $(n_1, ..., n_{L-1})$.
 - Type of activation functions σ.
 - Layer types (dense, convolutional, etc. more details later)
 - Additional regularization terms, e.g. promoting sparsity.

• Choosing the network architecture:

- In principle: Hyperparameter selection problem.
- In practice: Engineering problem, problem specific architecture.
- Depends on amount and type of data and available computational resources.



Universal Representation Theorem

Hornik et al., 1989; Cybenko, 1989

Sloppy formulation: A neural network with **a single** hidden layer and a monotonically increasing nonlinearity σ can approximate any continuous function $F : \mathbb{R}^{n_0} \mapsto \mathbb{R}^{n_2}$ from inputs \mathbb{R}^{n_0} to outputs \mathbb{R}^{n_2} with arbitrary accuracy given that sufficiently many hidden neurons n_1 are provided.

- Proof for sigmoid activation functions: G. Cybenko: "Approximations by superpositions of sigmoidal functions", Mathematics of Control, Signals, and Systems 2, 303-314 (1989)
- Generalization: K. Hornik: "Approximation Capabilities of Multilayer Feedforward Networks", Neural Networks 4, 251-257 (1991)

Caveats:

- Existence of network parameters that approximate *F* does not mean these parameters can be efficiently found.
- Approximation of a function does not mean exact representation, error might be too large for practical purposes.
- For many complex functions, the number of hidden neurons required to achieve acceptable accuracy is unpractical. → deep neural networks.



Parameter optimization (weight training)



Gradient descent of loss function

Cost or **loss function** *C* quantifies performance of network with parameters θ to predict observations **X**.

Learning network weights

 $\hat{\boldsymbol{\theta}} = \operatorname{arg\,min}_{\boldsymbol{\theta}} C(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta})$

Minimizing cost $C \equiv \max inizing \ score - C$. Often we just write $C(\theta)$.

Neural networks are usually trained with some form of gradient descent:

Simple gradient descent algorithm

• Initialize θ_0

- For t = 0, ..., T 1 or until converged:
 - Compute gradient $\mathbf{g}(\theta_t) = \nabla_{\theta} C(\theta_t)$
 - **O** Update parameters: $\theta_{t+1} = \theta_t \eta_t \mathbf{g}(\theta_t)$
- Networks are functions of functions, so we will need to use the chain rule of differentiation.
- Key is to make differentiation fully automatic so that we can just define the network architecture without worrying about functions and gradients.
 - \rightarrow Backpropagation.



Backprop

• Activation of a single neuron:

$$x_i^l = \sigma\left(z_i^l
ight)$$

 $z_i^l = \sum_j w_{ij}^l x_j^{l-1} + b_i^l$

$$egin{aligned} \mathcal{C}(\mathbf{X}, m{ heta}) &= rac{1}{N} \sum_{k=1}^N c(\mathbf{x}_k, \mathbf{y}_k, m{ heta}) \ rac{\partial \mathcal{C}(\mathbf{X}, m{ heta})}{\partial \hat{y}_{k,i}} &= rac{1}{N} \sum_{k=1}^N rac{\partial c(\mathbf{x}_k, \mathbf{y}_k, m{ heta})}{\partial \hat{y}_{k,i}} \end{aligned}$$

• Error: Loss gradient with respect to the *i*th weighted input

$$\mathbf{e}_{i}^{L} = \frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_{i}^{L}} = \underbrace{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial \hat{y}_{i}}}_{\text{loss derivative}} \frac{\partial \hat{y}_{i}}{\partial z_{i}^{L}} = \frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial \sigma(z_{i}^{L})} \underbrace{\frac{\sigma'(z_{i}^{L})}{\sigma(z_{i}^{L})}}_{\text{activation derivative}}$$

• Backpropagate error to earlier layers

$$\mathbf{e}_i^l = \frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l} = \boldsymbol{\sigma}'(z_i^l) \sum_j \mathbf{e}_j^{l+1} w_{ji}^{l+1}$$

• Weight gradients at each layer

$$\frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial b_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial w_{ij}^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}} = \frac{\frac{\partial c(\mathbf{x}, \mathbf{y}, \theta)}{\partial z_i^l}}$$

• Weight update:

$$egin{aligned} b_i^{l} \leftarrow b_i^{l} - \eta \, \mathbf{e}_i^{l} \ w_{ij}^{l} \leftarrow w_{ij}^{l} - \eta \, \mathbf{e}_i^{l} x_j^{l-1} \end{aligned}$$

with learning rate η .

















Stochastic gradient descent

Many cost functions and their gradients are of the form

$$C(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} c_i(\mathbf{x}_i, \theta) \qquad \nabla_{\theta} C(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} c_i(\mathbf{x}_i, \theta)$$

with data samples \mathbf{x}_i . Computing such gradients is expensive, involving $\mathcal{O}(N)$ operations, where N may be millions.

Idea: rewrite gradient as an expectation $\mathbb{E}_{\mathbf{x}}$:

$$\nabla_{\theta} C(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} c_i(\mathbf{x}_i, \theta) \approx \mathbb{E}_{\mathbf{x}} [\nabla_{\theta} c(\mathbf{x}, \theta)]$$

Subdivide the N-sample average into M averages over n samples each. Call index sets B_m :

$$\nabla_{\theta} C(\mathbf{X}, \theta) = \frac{1}{M} \sum_{m=1}^{M} \frac{1}{n} \sum_{i \in B_m} \nabla_{\theta} c_i(\mathbf{x}_i, \theta) \approx \mathbb{E}_{\mathbf{x}} [\nabla_{\theta} c(\mathbf{x}, \theta)]$$

Stochastic gradient descent

- Initialize θ_0
- For epoch e = 1, ..., E or until converged:
 - For minibatch m = 1, ..., M:
 - Compute minibatch gradient $\mathbf{g}_m(\theta) = \frac{1}{n} \sum_{i \in B_m} \nabla_{\theta} c_i(\mathbf{x}_i, \theta)$
 - **2** Update parameters: $\theta_{t+1} = \theta_t \eta_t \mathbf{g}_m(\theta_t)$



Stochastic gradient descent

SGD update with Momentum

Using the momentum parameter $0 \le \gamma \le 1$, we can implement SGD as:

$$egin{aligned} \mathbf{v}_t &= \gamma \mathbf{v}_{t-1} + \eta_t
abla_ heta \, C(heta_t) \ heta_{t+1} &= heta_t - \mathbf{v}_t \end{aligned}$$

or, equivalently, written with parameter updates $\Delta \theta_t = \theta_t - \theta_{t-1}$:

$$\Delta \theta_{t+1} = \gamma \Delta \theta_t - \eta_t \mathbf{g}(\theta_t)$$

SGD with momentum is equivalent to numerically solving a dynamical equation of a particle with mass m and position θ moving in a viscous medium under the dimensionless potential $C(\theta)$:



Langevin form: When approximating the gradient $\nabla_{\theta} C(\theta)$ with minibatches, it can be written as a deterministic part (gradient in the limit $M \to N$), plus a noise term depending on minibatch size. Then, Eq. (1) is a Langevin equation.



- Desirable to adapt the learning rate by taking large steps in shallow, flat directions
- Second-order methods do this by computing the Hessian matrix, but this is computationally expensive.
- Alternative: methods that track not only the gradient but also its second moment.
- Examples:
 - AdaGrad (Duchi et al., 2011)
 - AdaDelta (Zeiler, 2012)
 - RMS-Prop (Tieleman and Hinton, 2012)
 - ADAM (Kingma and Ba, 2014).



What does a neural network learn?



Feature detectors



Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.





VERSIL

What is this unit doing?



Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.





VERSIL













Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.

What features might you expect a good NN to learn, when trained with data like this?





Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.





Figure 1.2: Examples of handwritten digits from U.S. postal envelopes.





But what about position invariance? Our example unit detectors were tied to specific parts of the image


Deep Networks



Successive layers can detect higher-level features



Freie Universität





2006: the deep learning breakthrough



- Hinton, Osindero & Teh « <u>A Fast Learning</u> <u>Algorithm for Deep</u> <u>Belief Nets</u> », Neural Computation, 2006
- Bengio, Lamblin, Popovici, Larochelle « <u>Greedy Layer-Wise</u> <u>Training of Deep</u> <u>Networks</u> », *NIPS'2006*
- Ranzato, Poultney, Chopra, LeCun
 « Efficient Learning of Sparse Representations with an Energy-Based Model », NIPS'2006

2019 Turing award





Until 2009, deep networks were rarely used as training them seemed too difficult. Key advances that enabled the training of deep networks include:

- Rectifiers: Changing from activation functions such as logistic, tanh or arctan to rectifiers (ReLu, ELU, SoftPlus) avoids the vanishing gradient problem, makes the loss surface less "frustrated" and thus improves convergence.
- Stochastic Gradient Descent: Changing from full gradient descent to stochastic gradient descent resulted in significant reduction of the computational cost to convergence (cheaper iterations and ability to escape flat local minima and saddle points)
- GPU implementations of neural networks (Theano, TensorFlow), the resulting computational speedup, and the flexibility of automatic differentiation.



Applications to molecular systems



Physical Systems — example: O₂ molecule

Invariance to reference frame: U(x) is invariant when translating or rotating the molecule. We can thus choose an arbitrary position and orientation of the molecule, e.g.:

$$\mathbf{x}_1 = (0,0,0)$$
 $\mathbf{x}_2 = (d,0,0),$

 \rightarrow energy becomes a function of the interatomic distance, E(d).

2 Energy conservation: $U(\mathbf{x})$ and $\mathbf{f}(\mathbf{x})$ are related by

$$\mathbf{f}(\mathbf{x}) = -\nabla U(\mathbf{x}).$$

With the choice, we can compute the components of f(x) as:

$$\mathbf{f}_1 = \left(-\frac{\partial U(d)}{\partial d}, 0, 0\right) \qquad \mathbf{f}_2 = \left(\frac{\partial U(d)}{\partial d}, 0, 0\right)$$

Our Content of Physical laws: Same physical laws apply everywhere in the universe. Energy is unchanged if exchange the labels "1" and "2" (permutation invariance).







Physical Systems — example: O₂ molecule

- *O*₂ **example**: Energy is invariant wrt rotation, translation and permutation.
- **Equivariance**: A function is equivariant if it transforms the same way as its argument. For example, the force is defined by

$$\mathbf{f}(\mathbf{x}) = -\nabla U(\mathbf{x}).$$

If we rotate the molecule with \mathbf{R} , the force will rotate in the same way as the gradient is equivariant wrt rotation





Building Physics into the ML model

- Incorporating physical symmetries into ML model avoid learning them "by heart" via data augmentation.
- Reduces the dimensionality of the problem. In the O₂ example above, we have reduced the learning of:

$$U(\mathbf{x}): \mathbb{R}^6 o \mathbb{R}, \, \mathbf{f}(\mathbf{x}): \mathbb{R}^6 o \mathbb{R}^6$$

to

$$U(d): \mathbb{R} \to \mathbb{R}, \mathbf{f}(d) = -\nabla_{\mathbf{x}} U(d).$$

• We only operate on the manifold of physically meaningful solutions.





- **R** · rotates each atom with random 3D rotation matrix,**T** translates each atom with a translation vector.
- Potential or free energy of a molecule without external field is invariant to roto-translation:

$$U(\mathbf{Rx}+\mathbf{T})=U(\mathbf{x}).$$

 \rightarrow easily achieved by transforming the x into rototranslationally invariant features y (e.g., distances, angles).

• Force is equivariant to rotation, but invariant to translation

$$-\nabla U(\mathbf{Rx} + \mathbf{T}) = -\mathbf{R}\nabla U(\mathbf{x})$$

achieved by computing force explicitly in network (SchNet, CGnet)





Permutation invariance

• Energy/force a molecule are invariant/equivariant with permutation of equivalent particles:

$$U(\mathbf{Px}) = U(\mathbf{x})$$

 $-\nabla U(\mathbf{Px}) = -\mathbf{P}\nabla U(\mathbf{x})$

- Number of permutations increases exponentially with the number of identical particles → learning permutation invariance by data augmentation is hopeless.
- **DeepSets**¹, Theorem 2: any permutation symmetric function f(r) can be represented in the form $\rho(\sum_i \phi(\mathbf{r}_i))$, with functions ρ, ϕ .
- **Common choice**: compute potential energy as a sum

$$E(\mathbf{x}) = \sum_{i} E_{i}(\mathbf{x}), \qquad (4)$$

 E_i is the contribution of the energy by the *i*th atom in its chemical environment.

Freie Universität

¹Zaheer et al, *NIPS* 2017

Learning to represent energy functions



- Rototranslational and permutation invariances; parameter sharing.
- a) Network for computing the atomic energy of a single atom *i*. The system coordinates are mapped to rototranslationally-invariant features describing the chemical environment of atom *i*.
- b) Molecular system, e.g., H₂O, is composed by employing a network copy for each atom. Parameters are shared between networks for same elements.
- Total energy equals sum of atomic energies \rightarrow permutation

Behler and Parrinello, PRL 98, 146401 (2007)



Learning coarse-grained free energy functions

CGnet⁹



- Rototranslationally invariant features $\mathbf{y} \rightarrow \text{invariant energy } U(\mathbf{x})$.
- Energy conserving: $\mathbf{f}(\mathbf{x}) = -\nabla_{\mathbf{x}} U(\mathbf{x}) \rightarrow \text{equivariant force } \mathbf{f}(\mathbf{x}).$
- Prior energy: defines correct asympotic behavior where $U(\mathbf{x}) \rightarrow \infty$.
- No permutation invariance.
- No parameter sharing \rightarrow not scalable, not transferable.

⁹Wang, Olsson, Wehmeyer, Pérez, Charron, De Fabritiis, Noé, Clementi: ACS Cent. Sci. 5, 755-767 (2019)



Convolutional Neural Networks



LeNet 5

Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998

Convolutional filters

Discrete convolution:

$$y_i = (\mathbf{x} * \mathbf{w})_i = \sum_j x_j w_{i-j} = \sum_j w_j x_{i-j}$$

Convolutional Kernel / Filter

Apply convolutions

30	3,	2_2	1	0
02	02	10	3	1
30	1,	2_2	2	3
2	0	0	2	2
2	0	0	0	1

12.0	1 2 .0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

0	1	2	
2	2	0	
0	1	2	

3	30	21	1_{2}	0
0	02	1_2	30	1
3	10	2_1	2_{2}	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
00	01	1_2	3	1
32	1_2	20	2	3
20	0,	02	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	00	1_1	32	1
3	1_{2}	2_2	20	3
2	00	01	22	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



Convolutions are translation-equivariant



- Convolutions combine parameter sharing and equivariance.
- Standard convolutions are translation-equivariant.
- Rotation-equivariant convolutions exist (spherical CNNs, SchNet).

Freie Universität

Convolutional filters perform image processing

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	



Example: face recognition







First Layer Representation

Second Layer Representation Third Layer

Third Layer Representation



540,000 artificial distortions

+ 60,000 original

Test error: 0.8%

60,000 original datasets Test error: 0.95%







Convolution is a linear operation

Convolving $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{w} \in \mathbb{R}^m$, $m \le n$ can be written as linear operation

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

with $\mathbf{W} \in \mathbb{R}^{n-m+1 \times n}$ being a **Toeplitz matrix**:

$$\mathbf{W} = \left(\begin{array}{ccccc} w_1 & \cdots & w_m & & \\ & \ddots & & \ddots & \\ & & w_1 & \cdots & w_m \end{array}\right)$$

Convolutional Networks:







Motivation from neuroscience

- Pioneering neuroscientists: Hubel and Wiesel (1959, 1962, 1968):
 - Discrovered basic functionality of mammalian vision system by recording individual neuron activity in cats
 - Neurons in the early visual system respond mostly to specific patterns of light, such as precisely oriented bars.
- Primary visual cortex (V1): first brain area that performs advanced processing of visual input \rightarrow inspires ConvNets
 - V1 is arranged in 2*d* spatial map, mirroring the image in the retina.
 → inspires 2*d* structure of ConvNets.
 - V1 simple cells respond approximately linearly to a small, spatially localized receptive field.

 \rightarrow inspires ConvNet detector units

• Most simple cells seem to perform convolutions whose weights are described by **Gabor functions**.

	1 1 1 1 1 1 1 1	
///	1 1 1 1 1 1 1 1	
1111==000		
0 0 0 0 0 0 0 0 0		
0 0 0 8 0 0 0 0 0		
* * * = = // // //		
* * * = = = // //		
*****		Freie Universität

Berlin

Unsupervised learning



Principal component analysis

O Compute the covariance matrix

$$\mathbf{C}_0 = \frac{1}{T-1} \mathbf{X}^\top \mathbf{X},$$

which is just a scaled version of $\boldsymbol{X}^{\top}\boldsymbol{X}$

Olve the Eigenvalue problem:

$$\mathbf{C}_0 \mathbf{w}_i = \sigma_i^2 \mathbf{w}_i$$

with normalization $\mathbf{w}_i^{\top} \mathbf{w}_i = 1$.

- Select m eigenvectors with largest eigenvalues.
- **3** Reduce dimension from *n* to *m* with $\mathbf{W}_m = [\mathbf{w}_1, ..., \mathbf{w}_m]$:

$$\mathbf{Y}_m = \mathbf{X}\mathbf{W}_m$$
.





Autoencoder



encoder

$$\operatorname{Loss}(\mathbf{x};\theta) = \sum_{i} [\hat{y}_{i}(\mathbf{x};\theta) - x_{i}]^{2}$$



Recurrent neural networks



RNNs



- **One-to-one**: Classical (dense or convolutional) feedforward nets
- One-to-many: Given Image, generate sequence of words.
- Many-to-one: Given sequence of words, make classification (e.g., sentiment)
- Many-to-many: Text translation, on-line video segmentation



RNNs

• Specialized for:

- processing a sequence of values [x(1),...,x(t),...,x(T)].
 (t not necessarily physical time)
- Sequences of variable length.
- Example: "I went to Chile in 2013" and "In 2013, I went to Chile."
 - Extract the year the narrator went to Chile.
 - Traditional fully connected net would have separate parameters for each input feature, i.e. learn all of the rules of the language separately at each position in the sentence.
- Key idea: parameter sharing.
 - related idea: Convolution across a 1-D temporal sequence (time-delay neural networks – see Lang and Hinton, 1988; Waibel et al., 1989; Lang et al., 1990).



• **Example**: dynamical system driven by external signal $\mathbf{x}(t)$:

$$\mathbf{h}(t) = f(\mathbf{h}(t-1), \mathbf{x}(t); \theta)$$

- Train network to predict future from the past.
- h(t) is summary of present and past. It is lossy as it represents an arbitrary-length sequence (x₁,...,x_t) with fixed length.



• **Unfolding**: represent recurrence after t steps with a function g_t :

$$\mathbf{h}(t) = g_t(\mathbf{x}(t), ..., \mathbf{x}(1)) = f(\mathbf{h}(t-1), \mathbf{x}(t); \theta)$$

- Folding factorizes g_t into repeated application of a function f.
- f has same input size regardless of sequence length.
- Parameter sharing: use same transition function f every time step.



Basic universal RNN



- Can with finite size compute any function computable by a Turing machine (Siegelmann and Sontag, 1991; Siegelmann, 1995; Siegelmann and Sontag, 1995; Hyotyniemi, 1996).
- Loss L measures distance of predictions \mathbf{o}_t from training targets \mathbf{y}_t .



RNNs are deep — gradient annihilation

short "time" dependency: The clouds are in the sky.



long "time" dependency: I was born in France. Over the years, I have learned many foreign languages, but my native language is french.





Long-short time memory

Simple RNN



Long Short Time Memory (LSTM, Hochreiter & Schmidhuber 1997)





- 4.4 MB text input,
- 3-layer RNN, 512 hidden nodes per layer.

VIOLA:

Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

KING LEAR:

O, if you were a feeble sight, the courtesy of your law, Your sight and several breath, will wear the gods With his heads, and my hands are wonder'd at the deeds, So drop upon your lordship's head, and your opinion Shall be against your honour."



Generate Math

- Source: 16 MB LaTeX code *http://stacks.math.columbia.edu/*
- Sample: almost compiles, minor corrections needed.

For $\bigoplus_{n=1,\dots,m}$ where $\mathcal{L}_{m_{\bullet}} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X, U is a closed immersion of S, then $U \to T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \operatorname{Spec}(R) = U \times_X U \times_X U$$

and the comparison in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps M along the set of points Sch_{fppf} and $U \to U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ??. Hence we obtain a scheme S and any open subset $W \subset U$ in Sh(G) such that $Spec(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\operatorname{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}|_U$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for i > 0 and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^{ullet} = \mathcal{I}^{ullet} \otimes_{\operatorname{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$Arrows = (Sch/S)_{fppf}^{opp}, (Sch/S)_{fppf}$$



Pretending Math¹

- Source: 16 MB LaTeX code *http://stacks.math.columbia.edu/*
- Sample: almost compiles, minor corrections needed.

Proof. This is form all sheaves of sheaves on X. But given a scheme U and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,...,n} U_i$ be the scheme X over S at the schemes $X_i \to X$ and $U = \lim_i X_i$.

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{\mathcal{X},...,0}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.

Lemma 0.3. In Situation ??. Hence we may assume q' = 0.

Proof. We will use the property we see that \mathfrak{p} is the mext functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where K is an F-algebra where δ_{n+1} is a scheme over S.


Commonly-used Tricks



Multi-task learning

Alpha Go Zero

The network learns `tabula rasa' (from a blank slate)

At no point is the network trained using human knowledge or expert moves

The network

residual layer

residual layer

residual layer

residual layer

residual layer

residual layer

residual layer residual layer

residual layer residual layer residual layer

residual layer

residual layer residual layer

residual layer residual layer

residual layer residual layer

residual layer residual layer

residual laver

residual layer residual layer

residual layer

residual layer

residual layer

residual layer residual layer residual layer

residual layer

residual layer residual layer

residual layer

residual layer

residual laver

residual layer

residual layer

residual layer residual layer

convolutional layer

Input: The game state (see below)

3

policy head

value head

The value head game value for current player [-1, 1] tanh non-linearity 🔲 scalar Fully connected layer Rectifier non-linearity Hidden layer; size 256 Fully connected layer Rectifier non-linearity Batch normalisation 1 convolutional filter (lxl) Input



The policy head 19 x 19 + 1 (for pass) move logit probabilities Fully connected layer Rectifier non-linearity Batch normalisation 2 convolutional filters (1x1) Input A residual layer Rectifier non-linearity Skip connection Batch normalisation 256 convolutional filters (3x3) Rectifier non-linearity Batch normalisation 256 convolutional filters (3x3) Input

Adversarial attacks and training

- **Adversarial attacks** (trying to fake the network):
 - Intended perturbation: Accuracy reduced to random with test examples with small $\|\mathbf{x}_i^{\text{test}} \mathbf{x}_j^{\text{train}}\|$ but large $\|y_i^{\text{test}} y_j^{\text{train}}\|$ for given i, j.
 - Random perturbation: Add Gaussian noise with a very small amplitude. Picture indistinguishable for humans, but breaks neural network performance (Szegedy et al. 2014).



 Adversarial training: include adversarial attacks in the training data to force predictions to be locally constant near training data.



ResNets

- With most network architectures, when adding layers (increasing depth), the training loss first reduces but then increases.
- Indicates training problem adding layers make the network more expressive, so training loss should be non-increasing.
 → also affects the test loss.
- Residual Networks (He, Zhang, Ren, Sun, 2015) enable training of very deep networks.



Freie Universität

Berlin

- As in LSTMs, ResNets introduces a short-cut path that can carry gradients deep.
- ResNets are state-of-the-art in many problems (CIFAR, ImageNet, AlphaGo Zero).





Generative networks ==> Monday

