



# Deep learning: a non-alchemical view

Yuesong Shen

Computer Vision & Artificial Intelligence Technical University of Munich July 22, 2020



Tur Uhrenturm

## Outline

- Deep learning: the original series
- Deep learning: the next generation?
- Neural networks as chain graphs
  - $\circ~$  Core ideas
  - Generality of the CG interpretation
  - Selected case studies of NN designs
- Conclusion and discussions

### Deep learning: the original series



<sup>&</sup>lt;sup>1</sup>From https://www.cbc.ca/news/technology/turing-award-ai-deep-learning-1.5070415





## The deep learning boom

Success stories of DL on ImageNet:

- AlexNet<sup>2</sup> (2012): significantly outperforms other ML methods
  → Deep CNN (8 layers), max-pooling, ReLU, LRN, dropout
- ResNet<sup>3</sup> (2015): reaches human-level
  → Very deep CNN (100+ layers), Residual block, BatchNorm

NNs are widely used, achieving SOTA results in many domains:

Computer visionNatural language processingReinforcement learningGenerative modelsGraph learningUncertainty estimation

Many cool works from our chair use DL!

<sup>3</sup>He et al. "Deep Residual Learning for Image Recognition". In: CVPR. 2016.

Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM

<sup>&</sup>lt;sup>2</sup>Krizhevsky et al. "ImageNet Classification with Deep Convolutional Neural Networks". In: *NIPS*. 2012, 1097–1105.





# Mystical arts of deep learning at a glance

We have an **overwhelming** amount of NN designs and practices:

- Architecture: ResNet, Inception, (bi-)RNN, Siamese, Auto-encoder, GAN ....
- Layer: Convolution, Pooling, LSTM, Attention, GCN, GAT, skip connection ...
- Activation: (leaky) ReLU, Sigmoid, tanh, softmax, maxout, ELU ...
- Loss: L1, L2, Hinge, Cross Entropy, NLL, KL Div., artisanal heuristics ...
- Optimizer: SGD, momemtum, Adagrad, Adam, RMSProp, AMSGrad ...
- Initialization: Xavier, Kaiming, Fixup, ReZero, Isometric ...
- **Regularization**: weight decay, dropout, Batch/Layer/Group/Weight Norm, gradient clipping, DropConnect ...
- And many, many, many more (input pre-processing, Ir scheduling ...)!







### The alchemy age of deep learning

### As Ali Rahimi pointed out in the NIPS 2017 "Test of Time" talk:

"Machine learning has become alchemy."



#### We have scalable model, but non-scalable progress!

#### We need the "science" of deep learning

<sup>&</sup>lt;sup>5</sup>From https://medium.com/@Synced/lecun-vs-rahimi-has-machine-learning-become-alchemy-21cb1557920d Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM

### Deep learning: the next generation?



<sup>&</sup>lt;sup>6</sup>Adapted from http://yvettecandraw.blogspot.com/2011/09/speaking-of-outer-spacethe-top-ten.html





# Going beyond the deep learning alchemy

- DL does have a theory: Universal approximation theorem
  ~> Establish repr. power of NNs, Little insight, still "black-box".
- Several prior works aim to provide theoretical insights to NNs:
  - Gaussian process, kernel theory: NN as GP<sup>7</sup>, NTK<sup>8</sup>
  - Continuous counterpart: Neural ODE<sup>9</sup>
  - Optimal transport for generative models<sup>10</sup>
  - Signal propagation<sup>11</sup>, Lottery ticket hypothesis<sup>12</sup>...
- Our approach tries to make the link between PGM and DL.

<sup>&</sup>lt;sup>7</sup>Lee et al. "Deep Neural Networks as Gaussian Processes". In: *ICLR*. 2018.

<sup>&</sup>lt;sup>8</sup>Jacot et al. "Neural Tangent Kernel: Convergence and Generalization in Neural Networks". In: *NeurIPS*. 2018. <sup>9</sup>Chen et al. "Neural Ordinary Differential Equations". In: *NeurIPS*. 2018, 6572–6583.

<sup>&</sup>lt;sup>10</sup>Genevay et al. "GAN and VAE from an optimal transport point of view". In: arXiv:1706.01807 (2017).

<sup>&</sup>lt;sup>11</sup>Poole et al. "Exponential expressivity in deep neural networks through transient chaos". In: *NIPS*. 2016.

<sup>&</sup>lt;sup>12</sup>Frankle and Carbin. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks". In: *ICLR*. 2019. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





# Probabilistic graphical models: a new hope?

- **Representation**: graphical description of RVs and their relations Graph  $\leftrightarrow$  Overall distrib. Node  $\leftrightarrow$  Random var. Edge  $\leftrightarrow$  Direct influence
- Many PGM variants:



$$P(A, B, C, D) = P(A)P(B|A)P(C|A)P(D|B, C)$$
(BN)

$$=\frac{1}{Z}\phi(A,B)\phi(A,C)\phi(B,D)\phi(C,D)$$
(MN)

- = P(A, B)P(C, D|A, B)(CG)
- Flexible modeling: generative, unsupervised, domain-knowledge ...

Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





## The intractability strikes back

How efficient are **inference** and **learning** on general PGM?

~> Bad news: exact inference is proven to be intractable!<sup>13</sup>

How about approximations?

 $\rightarrow$  Bad news: Any with better-than-random guarantee is intractable!<sup>13</sup>

Don't panic. All is not lost!

- Tractable for special cases: tree, binary submodular ...
- Many (somewhat) practical approx. inference / learning methods:
  Mean field, Loopy belief propagation, Gibbs sampling ...
  Pseudo-likelihood, Contrastive divergence ...

### PGM offers flexible modeling, but is hard to do scalable learning!

<sup>&</sup>lt;sup>13</sup>Koller and Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





### PGM and NN: a tale of two models

- PGM: elegant theory, sub-optimal scalability
- NN: strong practical performance, weak theoretical support
- → Would be nice to have the best of both worlds!

### Shocking revelation: NN can actually be interpreted as PGM!



<sup>&</sup>lt;sup>14</sup>Adapted from https://knowyourmeme.com/memes/lets-see-who-this-really-is Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM

### Neural networks as chain graphs

A.K.A.: Deriving Neural Network Design and Learning from the Probabilistic Framework of Chain Graphs



<sup>&</sup>lt;sup>15</sup>Adapted from http://www.dearcastandcrew.com/content/2015/12/10/how-to-avoid-the-star-wars-syndrome.html

# Outline

- Deep learning: the original series
- Deep learning: the next generation?
- Neural networks as chain graphs
  - $\circ~$  Core ideas
  - Generality of the CG interpretation
  - Selected case studies of NN designs
- Conclusion and discussions





### From a binary sequential layered CG ...

• L layers of binary RVs  $X = (X^1, \ldots, X^L), X'_i \in \{\alpha', \beta'\}$ .



- Chain graph:  $P(X^2, ..., X^L | X^1) = \prod_{l=2}^{L} P(X^l | X^{l-1});$
- Bipartite chain component:  $P(X'|X'^{-1}) = \prod_{i=1}^{N'} P(X'_i|X'^{-1});$
- Pairwise CRF:  $P(X'_i|X'^{-1}) = \frac{1}{Z'_i(X'^{-1})} \exp\left(b'_i X'_i + \sum_{j=1}^{N'^{-1}} w'_{j,i} X'_j X'_j \right);$
- Binary RV:  $P(X'_i|X'^{-1}) = \sigma((2X'_i (\alpha' + \beta'))(b'_i + (X'^{-1})^\top w'_{\cdot,i})).$





### ... To a neural network!

Given an input  $\tilde{x}^1$ , we denote for each RV node  $X'_i$ :

- Marginal distribution:  $Q_i^{\prime}(x_i^{\prime}|\tilde{x}^1) = P(X_i^{\prime} = x_i^{\prime}|X^1 = \tilde{x}^1);$
- Expected value:  $q'_i = \mathbb{E}_{Q'_i}[X'_i] \ (q^1 = \tilde{x}^1).$

We have the following recursive expression for Q $Q'_i(x'_i|\tilde{x}^1) = \sum_{x'^{-1} \in \{\alpha'^{-1}, \beta'^{-1}\}^{N'^{-1}}} P(x'_i|x'^{-1})Q'^{-1}(x'^{-1}|\tilde{x}^1) = \mathbb{E}_{Q'^{-1}}[P(x'_i|X'^{-1})].$  $\rightsquigarrow$  Looks intractable ... How to proceed?

### Do a linear approximation! ~> Move the expectation inside.

$$q'_{i} = \sum_{x'_{i} \in \{\alpha',\beta'\}} x'_{i} \mathbb{E}_{Q'^{-1}}[P(x'_{i}|X'^{-1})] \approx \frac{\beta'-\alpha'}{2} \tanh\left(\frac{\beta'-\alpha'}{2}(b'_{i} + (q'^{-1})^{\top}w'_{\cdot,i})\right) + \frac{\alpha'+\beta'}{2}$$
  
Hey, this is just a feed-forward with  $q!$ 

$$\alpha' = \mathbf{0}, \beta' = \mathbf{1} \Longrightarrow$$
 Sigmoid,  $\alpha' = -\mathbf{1}, \beta' = \mathbf{1} \Longrightarrow$  Tanh.

Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





# Summarizing the core ideas

Overview of some correspondences:

NNLayerNeuronFeed-forwardNon-linearityCGChain componentRV nodeApprox. inferenceRV distrib.

Learning?

→ Same! Loss via MLE / SRM + SGD + Backprop.

When is the approximation accurate?

- Small pairwise weight
  weight decay is good!
- Linear neighborhood

→ might explain why ReLU is good.

### Clear modeling, assumptions and approximations!

Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM

# Outline

- Deep learning: the original series
- Deep learning: the next generation?
- Neural networks as chain graphs
  - $\circ$  Core ideas
  - Generality of the CG interpretation
  - Selected case studies of NN designs
- $\circ~$  Conclusion and discussions



# Let's generalize ...

- Network architecture: general (Partially) Directed Acyclic Graph
  - Multi-branch: Inception, ResNet
  - Time series: RNNs
- Layer connection: sparse connection, shared / fixed weight
  - NN layers: dense, convolution, average pooling, skip connection ...
  - Intra-connection: non-bipartite CRF as the last layer<sup>161718</sup>
- Modular components: transfer learning via layer reuse.
- $\rightsquigarrow$  How about activation functions?

<sup>17</sup>Chen et al. "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs". In: *TPAMI* 40.4 (2018), pp. 834–848.

<sup>&</sup>lt;sup>16</sup>Zheng et al. "Conditional Random Fields as Recurrent Neural Networks". In: *ICCV*. 2015.

<sup>&</sup>lt;sup>18</sup>Lample et al. "Neural Architectures for Named Entity Recognition". In: *NACCL*. 2016. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM



# General form of activation function

- We define statistics T' that extract features from each node  $X'_i$ .
- General nodewise CPD: (positive unit-sum) generalized linear model  $P(X_i^l|X^{l-1}) = f_i^l(X_i^l, e_i^l(T^{l-1}(X^{l-1})))$

with

$$e'_i(T^{l-1}(X^{l-1})) = b'_i + \sum_{j=1}^{N^{l-1}} w'_{j,i} T^{l-1}(X^{l-1}_j).$$

• General feed-forward expression:

$$q'_i = \mathbb{E}_{Q'_i}[T'(X'_i)] pprox \int_{x'_i} T'(x'_i) f'_i(x'_i, e'_i(q'^{-1})) dx'_i.$$

- Example: multi-label node with label indicator statistics  $\implies$  softmax.
- Batch/Layer/Weight/Group Norm: reparametrization of *e*, as part of *f*. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





# Rectified Gaussian distribution family The question to address: Can we get ReLU from this?

Yes we can. With rectified Gaussian nodes, defined as:

$$T'_i(X'_i) = X'_i = \max(0, Y'_i); \quad Y'_i \sim \mathcal{N}(e'_i(T^{l-1}(X^{l-1})), (s'_i)^2)$$

In fact, several non-linearities can be approximated from this family:

- softplus:  $s_i^{\prime} \approx 1.776$  (constant std);
- ReLU:  $s'_i = \tanh(e'_i(T^{l-1}(X^{l-1})));$
- $\epsilon$ -Leaky ReLU:  $X'_i = \max(\epsilon Y'_i, Y'_i)$ .

 $\rightsquigarrow$  Without rectification (or  $\epsilon = 1$ ), we get the identity activation!



## Outline

- Deep learning: the original series
- Deep learning: the next generation?
- Neural networks as chain graphs
  - Core ideas
  - Generality of the CG interpretation
  - Selected case studies of NN designs
- Conclusion and discussions





# Going deeper: residual block as refinement module

Residual block is effective for going deeper  $\Rightarrow$  better results!  $\rightsquigarrow$  Can we interpret it as CG? Yes we can!



Introducing the refinement module: side-chain augmentation!

- base submodule:  $X^{l-1} \rightarrow X^{l}$ ;
- refining submodule:  $ilde{X}' o Z' o X'$ .

Interesting remarks:

- Pre-activation variant ~> found as the preferred choice<sup>19</sup>!
- **Recursive refinement** ~> sequence of consecutive res. blocks!

<sup>&</sup>lt;sup>19</sup>He et al. "Identity Mappings in Deep Residual Networks". In: *ECCV*. 2016. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





## Learning sequential data

Simple (Elman) RNN  $\Rightarrow$  Markov assumption

 $\rightsquigarrow$  component  $P(X^{l,t}|X^{l-1,t},X^{l,t-1})$  + dense connect.  $X^{l,t-1} \rightarrow X^{l,t}$ 

→ Vanishing / exploding gradient! alleviated with LSTM, GRU ...



Rethinking dependency: conditional independence through time!

$$P(X_i^{l,t}|X^{l-1,t},X^{l,t-1}) = P(X_i^{l,t}|X^{l-1,t},X_i^{l,t-1})$$

→ We have rediscovered the independently RNN<sup>20</sup>! Yay!

- Simpler than Elman RNN, better than LSTM!
- Can use ReLU, residual block, etc and go deeper.

<sup>&</sup>lt;sup>20</sup>Li et al. "Independently Recurrent Neural Network (IndRNN): Building a Longer and Deeper RNN". . In: *CVPR*. 2018. Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





# Stochastic regularization with dropout

We can also replicate dropout within the CG framework:

- For each  $X_i^{l-1}$ , introduce an auxiliary RV  $D_i^{l-1} \sim \text{Bernoulli}(p^{l-1})$ ;
- Pairwise connection  $X_j^{l-1} \to X_i^l \Longrightarrow$  ternary  $X_j^{l-1} \to X_i^l \leftarrow D_j^{l-1}$ .

Node-wise CPD with dropout:

$$P(X_i^{l}|X^{l-1}, D^{l-1}) = f_i^{l}(X_i^{l}, e_i^{l}(T^{l-1}(X^{l-1}), D^{l-1}))$$
  
with  $e_i^{l}(T^{l-1}(X^{l-1}), D^{l-1}) = b_i^{l} + \sum_{j=1}^{N^{l-1}} D_j^{l-1} w_{j,i}^{l} T^{l-1}(X_j^{l-1}).$ 

Reproducing the behavior of dropout:

- During training: sample  $D_i^{l-1}$  during each feed-forward;
- During test: marginalize  $D_i^{l-1} \Rightarrow$  constant scale  $p^{l-1}$ .





## Stochasticity without dropout?

We have a PGM, stochastic behavior can simply be achieved via sampling!

For directed graphs, we can use forward (ancestral) sampling!

- need a single pass like feed-forward;
- produce unbiased samples.

Issue: sampling doesn't go well with backprop!

- Continuous RV: reparametrization trick
- Discrete RV: Gumbel-softmax trick? (continuous relaxation ...)
- High variance

What can we do?

∽→ what made dropout work?





# Something new: partially collapsed feed-forward

Main idea: sampling only a part of the nodes. (with proba.  $p \in [0, 1]$ )

→ we simply "mix up" feed-forward and forward sampling!

 $q_i^l = \begin{cases} E_{Q_i^l}[T^l(X_i^l)] & ext{if collapsed (feed-forward);} \\ T^l(x_i^l), \ x_i^l \sim Q_i^l & ext{if uncollapsed (forward sampling).} \end{cases}$ 

Following "collapsed sampling", I dub thee partially collapsed feed-forward !

Cool things about PCFF:

- gradient flows freely (can still use reparam. trick)
- offers a bias-variance trade-off
- generalizes over feed-forward and forward sampling





### Experimental evaluations of PCFF

### Does PCFF actually work in practice? Yup!

### Experiments with Conv. ResNet20 on CIFAR-10:



- · Comparable to dropout, but more consistent improvement!
- More exp. in the paper: PCFF also work for IndRNN!

Yuesong Shen, Deep learning: a non-alchemical view, CVAI - TUM





### Conclusion and discussions

• The take-home message:

### Neural networks are chain graphs!

• I repeat (because it is the take-home message!):

### Neural networks are chain graphs!

- Some possible tracks for future work:
  - 1. Analyze existing designs using the NNasCG framework;
  - 2. Propose improvements / new designs guided by the framework;
  - 3. Go beyond the FeedForward + BackProp learning strategy;
  - 4. And more!

### Many work can be done: collaborations are welcome!

### Thank you! Questions?



<sup>&</sup>lt;sup>21</sup>From http://barkingalien.blogspot.com/2017/02/i-choose-you.html