The Information Theory of Deep Learning: What do the layers represent?

Naftali Tishby InterSpeech, Hyderabad, September 2018

School of Engineering and Computer Science

The Edmond & Lily Safra Center for Brain Sciences











Noga Zaslavsky Ravid Schwartz-Ziv



The Deep Learning revolution

edureka!

ARTIFICIAL INTELLIGENCE

Engineering of making Intelligent Machines and Programs

1960's

1970's

1980's

+

1950's

Ability to learn without being explicitly programmed

edurel



1990's

reka! e

2000's

2006's

DEEP LEARNING

2012's 🎽

2017's

Learning based on Deep Neural Network

2010's

Mimic Biology- Brain like Neurons





Simple Neural Networks





Deep Learning: Neural-Nets strike back



IS18 Tutorial, September 2018 - Tishby



We begin to obtain new theoretical understanding...

We combine 3 different ingredients:

- Rethinking Statistical Learning Theory
 - Worse case PAC bounds → TYPICAL data dependent model free bounds...
 - From expressivity/Hypothesis class \rightarrow Input Compression bounds
- Information Theory (statistical mechanics...)
 - Large scale learning *Typical* input patterns
 - \rightarrow Concentration of the Mutual Information values
 - \rightarrow Huge parameter space exponentially many optimal solutions
- Stochastic dynamics of the training process
 - Convergence of SGD to locally-Gibbs (Max Entropy) weight distribution
 - \rightarrow The mechanism of representation compression in Deep Learning
 - > Convergence times explains the benefit of the hidden layers



The match between DL and the Information Bottleneck Main results:

- Optimality: The layers converge to the [finite-sample] IB bound
 - DL can achieve optimal (model free, rule dependent) sample complexity-accuracy tradeoff
 - Through the diffusion/noisy phase of the Stochastic Gradient Descent optimization
 - Which compresses the representation by "forgetting" irrelevant details
- Benefit of the Hidden Layers
 - The benefit is mostly computational boosting the compression!
 - The location of the optimal layers is determined by the problem
- Interpretability
 - Full layers can have clear problem specific meaning, NOT single neurons (in general)!
- Design principles
 - DL is good for stochastic, compressible rules.
 - Layers final position is related to critical points of the Information Bottleneck



Known issues & reservations

Objections to the theory:

- Information estimation [requires quantization or noise, not scalable? ...]
 - Not needed for training, only as a tool for understating!
 - Binning is done with the actual known resolution. It should not affect network performance.
 - Requires finite precision or quantization CORRECT!
 - Mutual Information values concentrate & become MORE stable the larger the problem!
- Input Compression/Information loss not necessary [ResNets, RevNets, i-RevNets,...]
 - Compression comes from unit saturation, not seen with RelU's (Saxe 2018) WRONG!
 - Indeed, good generalization can be achieved without apparent layer compression.
 - Similar to the classical physics paradox of reversible microscopic laws & entropy increase...
 - No "forgetting" of non-informative features (really?)
- Stochastic Gradients not needed [no convergence to local weight Gibbs distribution]
 - Good generalization achieved without stochastic gradients in INFINITE TIMES! How?
 - Convergence to Gibbs (MaxEnt) distribution is only local (in each layer).
 - The benefits of the stochasticity is dynamical (computational), but also in saving training data!
 - There is important INFORMATION in the mini-batch fluctuations!
- Is the IB bound relevant?
 - It actually gives concreate predictions and interpretation of the layers & weights.
 - May explain biological neural network organization... our ultimate motivation.





Deep Neural Nets and Information Theory ??





IS18 Tutorial, September 2018 - Tishby

Some Information Theory basics

• The KL-distribution divergence:

for any two distributions p(x) & q(x) over X:

$$D[p(x) || q(x)] = \sum_{x} p(x) \log \frac{p(x)}{q(x)} \ge 0$$

• The Mutual Information:

for any two random variables, X, Y: I(X;Y) = D[p(x,y) || p(x)p(y)] = D[p(x|y) || p(x)] = D[p(y|x) || p(y)] = H(X) - H(X|Y)

• Data Processing Inequality (DPI) & Invariance:

for any Markov chain: $X \rightarrow Y \rightarrow Z$: $I(X;Y) \ge I(X;Z)$ Reparametrization Invariance, for invertible ϕ, ψ : $I(X; Y) = I(\phi(X); \psi(Y))$

What do the DNN Layers represent?



Data Processing Inequalities:

$$H(X) \ge I(X; h_{i}) \ge I(X; h_{i+1}) \ge I(X; h_{i+2}) \ge \dots$$

$$I(X; Y) \ge I(h_{i}; Y) \ge I(h_{i+1}; Y) \ge I(h_{i+2}; Y) \ge \dots$$

- A Markov chain of topologically distinct [soft] partitions of the input variable X.
- Successive Refinement of Relevant Information
- Individual neurons can be easily "scrambled" within each layer



Each layer is characterized by its Encoder & Decoder Information



<u>Theorem (Information Plane):</u> For large typical X, the sample complexity of a DNN is completely determined by the encoder mutual information, I(X;T), of the last hidden layer; the accuracy (generalization error) is determined by the decoder information, I(T;Y), of the last hidden layer.

The complexity of the problem shifts from the decoder to the encoder, across the layers... 12 8/12/18 IS18 Tutorial, September 2018 - Tishby

100 DNN Layers in Info-Plane without averaging



Inside Deep Learning

New experiments reveal how deep neural networks evolve as they learn.





Learning From Experience

Deep neural networks learn by adjusting the strengths of their connections to better convey input signals through multiple layers to neurons associated with the right general concepts.



When data is fed into a network, each artificial neuron that fires (labeled "1") transmits signals to certain neurons in the next layer, which are likely to fire if multiple signals are received. The process filters out noise and retains only the most relevant features.

Inside Deep Learning

New experiments reveal how deep neural networks evolve as they learn.



A INITIAL STATE: Neurons in Layer 1 encode everything about the input data, including all information about its label. Neurons in the highest layers are in a nearly random state bearing little to no relationship to the data or its label.

B FITTING PHASE: As deep learning begins, neurons in higher layers gain information about the input and get better at fitting labels to it.

© PHASE CHANGE: The layers suddenly shift gears and start to "forget" information about the input.

© COMPRESSION PHASE: Higher layers compress their representation of the input data, keeping what is most relevant to the output label. They get better at predicting the label.

FINAL STATE: The last layer achieves an optimal balance of accuracy and compression, retaining only what is needed to predict the label.



The role of stochasticity: How do we measure Mutual Information?

- The representation invariance of the mutual information raises an interesting question.
- Obviously, the computational complexity of learning is not representation invariant (think about learning from encrypted patterns). Thus, Information measures can't tell the whole story.
- Our experiments crucially depends on how we estimate information. We consider 3 types of estimations: (1) binning the variables. (2) adding noise / stochasticity (3) parametric approximations.
 In our experiments we quantized/bin the neuronal output values.
- All assume compressibity/refineability of the variables. They are not robust to arbitrary invertible transformations!
- The assertion that the layers are invertible transformations of the input is NOT robust to small noise and misleading. Binning or assuming stochastic mapping is essential for our information theoretic approach.
- Moreover, the IB is trivial (uninteresting) for completely deterministic rules! I argue that our theory predicts that <u>without additional structural information on the patterns</u>, DL can't work for completely deterministic rules, as they can't be distinguished from random (fully mixing) rules!



Rethinking Learning Theory

X

"Old" Generalization bounds: $\epsilon^{2} < \frac{\log |H_{\epsilon}| + \log \frac{1}{\delta}}{2m}$

- ϵ generalization error ρ_{μ}
- δ confidence
- *m*-number of training examples
- H_{ϵ} ϵ -cover of the Hypothesis class

typically we assume: $|H_{\varepsilon}| \sim \left(\frac{1}{\varepsilon}\right)^{a}$

d - the class (VC,...) dimension

... Don't work for Deep Learning! Higher expressivity - worse bound! New: Input Compression bound $|H_{\epsilon}| \sim 2^{|X|} \rightarrow 2^{|T_{\epsilon}|}$

 $T_{t} - \varepsilon$ -partition of the input variable X with respect to the distortion:

 $d_{B}(x,t) = D[p(y|x) \parallel p(y|t)]$

 $\geq \frac{1}{2\ln 2} \| p(y|x) - p(y|t) \|_{1}^{2}$

when
$$p(y|t) = \sum_{x} p(y|x) p(x|t)$$

 $\langle d_{B} \rangle = I(X;Y) - I(T;Y)$
small IB distortion, or high $I(T;Y)$,

⇒ small [typical] generalization error



Rethinking Learning Theory...

What are "large typical" patterns?

Typicality emerges when the underling pattern distribution can be asymptotically expressed as a long product of localized conditional probabilities.

E.g. Markov Random Fields, Hidden Markov Models, *pairwise* interaction Hamiltonians in physics, all common Graphical models, etc.

In our case it includes images, speech & text, long molecular sequences, signals generated by localized dynamic systems, etc. Then, the Shannon-McMillen limit for the entropy exists:

$$\lim_{n\to\infty} -\frac{1}{n}\log p(x_1,\dots,x_n) = H(X)$$

and almost all patterns are typical with probability:

 $p(x_1,...,x_n) \simeq 2^{-nH(X)}$

and also for large enough typical partitions, T:

$$p(x_1, ..., x_n | T) \simeq 2^{-n E(Z|T)}$$



Concentration of Mutual Information

$$I(X;T) = \left\langle \log \frac{p(x|t)}{p(x)} \right\rangle_{x,T} = \left\langle \log \prod_{i} \frac{p(x_i | Pa(x_i), t)}{p(x_i | Pa(x_i))} \right\rangle_{x,T} = \left\langle \sum_{i} \log \frac{p(x_i | Pa(x_i), t)}{p(x_i | Pa(x_i))} \right\rangle_{x,T}$$
$$I(T;Y) = \left\langle \log \sum_{x} p(y|x) p(x|t) - \log p(y) \right\rangle_{Y,T} = \left\langle \log \left[\sum_{x} p(y|x) \prod_{i} p(x_i | Pa(x_i), t) \right] - \log p(y) \right\rangle_{Y,T}$$

Proposition:

1. Both I(T;X) and I(T;Y), as defined, concentrate, uniformly, under the partition typicality assumption.

2. Both can be estimated uniformly well (over the partitions) from a sample of p(X,



Rethinking Learning Theory

 $\gamma^{H(X|T_{\epsilon})}$

X

 $\mathbf{2}^{H(X)}$

"Old" Generalization bounds: $\log |H_{\varepsilon}| + \log \frac{1}{\delta}{2m}$

- $\boldsymbol{\epsilon}$ generalization error
- δ confidence
- *m* number of training examples
- H_{ϵ} ϵ -cover of the Hypothesis class

typically we assume: $|H_{\varepsilon}| \sim \left(\frac{1}{\varepsilon}\right)^{a}$

d - the class (VC,...) dimension... Don't work for Deep Learning!Higher expressivity - worse bound!

New: Input Compression bound: $|H_{\varepsilon}| \sim 2^{|X|} \rightarrow 2^{|T_{\varepsilon}|}$

> $T_{ε}$ - ε-partition of the input variable *X* Information Theory: $|T_{ε}| \sim 2^{I(T_{ε};X)}$



... K bits of compression of X are like a factor of 2^{κ} training examples!



The Information Bottleneck (IB) Method

(Tishby, Pereira, Bialek, 1999)

- The Information Bottleneck method was born out of the Speech Recognition problem:
 - What are the simplest (efficient) representations

of the (high entropy) Acoustic Signal that yield good prediction of the (low entropy) phonemes ?

- The idea was to extract (approximate) Minimal-Sufficient Statistics - simplest features - of the complex signal (sound), that are informative for the simpler one (text).
- This was cast into a simple looking Information Theoretic tradeoff between compression and accuracy 8/14/18



The Information Bottleneck (IB) Method

(Tishby, Pereira, Bialek, 1999)

(1) Approximate Minimal Sufficient Statistics: Markov chain: $Y \rightarrow X \rightarrow S(X) \rightarrow \hat{X}$ $\hat{X} = \arg \min_{S(X) \rightarrow I(X;Y) - I(X;Y)} I(S(X); X)$ Relaxation - given p(X, Y): $\hat{X} = \arg \min_{p(X|X)} I(\hat{X}; X) - \beta I(\hat{X}; Y), \beta > 0$ (Shamir, Sabato,T., TCS 2010)

IS18 Tutorial, September 2018 - Tishby



The Information Bottleneck optimality bound (Tishby, Pereira, Bialek, 1999)

The IB bound optimality equations:

$$\min_{\rho(\hat{\boldsymbol{x}}|\boldsymbol{x}) \to \boldsymbol{X} \to \hat{\boldsymbol{y}}} I(\hat{\boldsymbol{X}}; \boldsymbol{X}) - \beta I(\hat{\boldsymbol{X}}; \boldsymbol{Y}) , \beta > 0$$

$$p(x \mid \hat{x}) = \frac{p(x)}{Z(x,\beta)} \exp(-\beta D[p(y \mid x) \parallel p(y \mid \hat{x}])]$$

$$Z(x,\beta) = \sum_{\hat{x}} p(\hat{x}) \exp(-\beta D[p(y \mid x) \parallel p(y \mid \hat{x}])]$$

$$p(\hat{x}) = \sum_{x} p(\hat{x} \mid x) p(x)$$

$$p(y \mid \hat{x}) - \sum_{x} p(y \mid x) p(x \mid \hat{x})$$

Solved by Arimoto-Blahut like iterations,

but with possibly sub-optimal solutions, bifurcations (!),





IS18 Tutorial, September 2018 - Tishby

Rethinking Learning Theory...

... but we need to guarantee the label homogeneity of the -partition with finite samples. Without additional structural information on the inputs (stability, robustness, topology), we must use the stochasticity of the rule and the IB distortion measure:

measure: The ε - partition, T_{ε} , is with the empirical distortion

 $d_{IB}(x,t) = D[\rho_{emp}(y|x) \parallel p(y|t)]$

as
$$\langle d_{IB} \rangle_{emp} = I(X;Y) - \hat{I}_{emp}(T;Y)$$

with a finite sample there is another information loss:

$$I(T;Y) \leq \hat{I}_{emp}(T;Y) + O_{\sqrt{\left(\frac{2^{I(T;X)}|Y|}{m}\right)}},$$

both should remain small for good generlaization! 24 8/12/18





- Layers of optimal DNN converge to [a successively refineable approximation of] the optimal finite-sample IB limit information-curve
- Layers must be in "different topological phases" of the IB solutions
- The DNN encoder & decoder for each layer satisfy the IB self-consistent equations







ICRI-CI

Layers paths with training/generalization error









In the noisy phase the weights diffuse and grow (ket)

IS18 Tutorial, September 2018 - Tishby



Gradients SNR, Diffusion & Compression - all layers





Layer - 4

Layer - 5

Layer - 6







Iterations

Compariantian of Inconting

The role of the batch size





CRI-CI

Break



Relevant and Irrelevant local dimensions

• The covariance matrix of the gradients is very narrow in the relevant local dimensions and very wide in the many other dimensions.



Is it the general picture? Yes!



6 layer committee machine

IS18 Tutorial, September 2018 - Tishby



... and for "Real-world" problems? Yes!



MNIST handwriting digit recognition with RelU's a CNN architecture 35 8/12/18 NASIT, May 2018 - Tishby

Is it the general picture? Yes!



CIFAR 10 object recognition task

Non decreasing layer widths - notice last hidden



Local weights Gibbs and optimal IB representations

Noisy relaxation (SGD) on training error (with *m* examples):

$$\frac{\partial W_k}{\partial t} = -\nabla E(W_k \mid X^{(m)}) + \beta_k^{-1} \xi(t) , \text{ layer } k, \quad \xi \sim N(0,1), \ \beta_k - decoder \ k - layer \ noise$$

$$\Rightarrow \text{Maximum Entropy:} \quad P_{Gibbe}(W_k \mid X^{(m)}) \propto \exp\left(-\beta_k E(W_k \mid X^{(m)})\right) \qquad \textbf{keely}$$

with additive [cross-entropy] training error and i.i.d. samples, using Bayes rule with the quenched W:

$$P_{Glbbs}(X | W_k) = P_{Glbbs}(X | T_k) \propto \exp\left(-\tilde{\beta}_k D_{KL}[p(Y | X) || p(Y | T)]\right)$$

This is precisely the IB optimal encoder with $\tilde{\beta}_k$ – the encoder k - layer noise

Since $I(X;T_k) = H(X) - H(X|T_k)$, Max Entropy of the weights \Rightarrow Min $I(X;T_k)$ thus SGD converges, layer by layer, to a maximally compressed representation, which is a SUCCESSIVELY REFINEABLE APPROXIMATION of the optimal IB bound! 37 8/12/18



The benefit of the hidden layers



IS18 Tutorial, September 2018 - Tishby

More layers take much FEWER training epochs for good generalization.

The optimization time depend super-linearly (exponentially?) on the compressed information, delta Ix, for each layer.



Relaxation times and the benefit of the hidden layers Noisy relaxation (SGD): $\frac{\partial W_k}{\partial t} = -\nabla E(W_k) + \beta_k^{-1}\xi(t)$, layer k, $\xi \sim N(0,1)$

 \Rightarrow Maximum Entropy (via Focker-Planck): $P_{Gibbs}(W_k) \propto \exp(-\beta_k E(W_k))$,

Relaxation time for non-strongly convex error: $\Delta t_k \sim \exp(\Delta S_k)$

Denote the layer compression be: $\Delta S_k = I(X; T_k) - I(X; T_{k-1})$

Since
$$\exp\left(\sum_{k} \Delta S_{k}\right) \gg \sum_{k} \exp\left(\Delta S_{k}\right) > \max_{k} \exp\left(\Delta S_{k}\right) \Rightarrow$$

Exponential boost in the relaxation time with K layers!

IS18 Tutorial, September 2018 - Tishby



Equilibration of Information Flow through the layers

The Information Capacity between two layers is bounded by the Gaussian capacity:

$$C_G(W_k) = \frac{1}{2} \log \left(1 + \frac{P_k}{N_k} \right) = \frac{1}{2} \log \left(1 + SNR_k \right)$$

The stochastic relaxation decreases the SNR of the irrelevant channels

- \Rightarrow In optimized DNN only the relevant information I(X;Y) flows through the network
- \Rightarrow SNR_k ~ const. we see this in the simulations.

This can determine the final layers locations in the Information Plane...

Unless the stochastic relaxation stops through critical slowing down near phase transitions !





- Fitting larger training data require more information in the hidden layers.
- It is the **mutual-information of the last hidden layer**, which determines generalization (unlike standard hypothesis class bounds)

IS18 Tutorial, September 2018 - Tishby





-



Layer - 4 Epoch 0





Second order phase transitions on the IB curve



The IB bifurcation (phase-transitions) points

The IB bifurcation points can be found as follows: $p_{\beta}(\boldsymbol{x}|\boldsymbol{\hat{x}}) = \frac{p(\boldsymbol{x})}{Z(\boldsymbol{x},\boldsymbol{\beta})} \exp(-\beta D[p(\boldsymbol{y}|\boldsymbol{x}) \| \boldsymbol{p}_{\beta}(\boldsymbol{y}|\boldsymbol{\hat{x}}])$ $\ln p_{\beta}(\mathbf{X} | \mathbf{\hat{X}}) = \ln \frac{p(\mathbf{X})}{Z(\mathbf{X}, \beta)} - \beta D[p(\mathbf{y} | \mathbf{X}) || p_{\beta}(\mathbf{y} | \mathbf{\hat{X}})]$ or then: $\frac{\partial \ln \rho_{\beta}(\boldsymbol{x} \mid \hat{\boldsymbol{x}})}{\partial \hat{\boldsymbol{x}}} = \beta \sum_{\boldsymbol{y}} \rho(\boldsymbol{y} \mid \boldsymbol{x}) \frac{\partial \ln \rho_{\beta}(\boldsymbol{y} \mid \hat{\boldsymbol{x}})}{\partial \hat{\boldsymbol{x}}}$ similarly: $\rho_{\rm B}(\mathbf{y} | \hat{\mathbf{x}}) = \sum_{\mathbf{x}} \rho(\mathbf{y} | \mathbf{x}) \rho_{\rm B}(\mathbf{x} | \hat{\mathbf{x}})$ $\frac{\partial \ln \rho_{\beta}(\mathbf{y} \mid \hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}} = \frac{1}{\rho_{\alpha}(\mathbf{y} \mid \hat{\mathbf{x}})} \sum_{\mathbf{x}} \rho(\mathbf{y} \mid \mathbf{x}) \rho_{\beta}(\mathbf{x} \mid \hat{\mathbf{x}}) \frac{\partial \ln \rho_{\beta}(\mathbf{x} \mid \hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}}$



IS18 Tutorial, September 2018 - Tishby

The IB bifurcation (phase-transitions) points

Defining the matrices:

$$C_{xx}(\hat{x},\beta) = \sum_{y} \frac{p(y|x)}{p_{\beta}(y|\hat{x})} p_{\beta}(x'|\hat{x}) p(y|x')$$
$$C_{yy}(\hat{x},\beta) = \sum_{x} \frac{p(y|x)}{p_{\beta}(y|\hat{x})} p_{\beta}(x|\hat{x}) p(y'|x)$$

these equations can be combined into two (non-linear) eigenvalue problems:

$$\left[I - \beta C_{xx'}(\hat{x}, \beta)\right] \frac{\partial \ln P_{\beta}(x' \mid \hat{x})}{\partial \hat{x}} = 0$$
$$\left[I - \beta C_{yy'}(\hat{x}, \beta)\right] \frac{\partial \ln P_{\beta}(y' \mid \hat{x})}{\partial \hat{x}} = 0$$

These eigenvalue problems have non-trivial solutions (eigenvectors) only at the critical bifurcation points (second order phase transitions).

IS18 Tutorial, September 2018 - Tishby



Bifurcation diagrams in symmetric rule: layers diffusion slows down at phase transitions



Each layer encodes the information in the IB bifurcation from the previous layer.

$$W^{i}\!\!\approx\!\sum_{splits}rac{\partial\!\log p\!\left(\!x\!\mid\!\!t_{s}^{k ext{-}1}\!
ight)}{\partial t_{s}^{k ext{-}1}}$$



47 8/12/18

Summary

The Information Plane provides a unique visualization of DL

- Most of the learning time goes to compression
- Layers are learnt bottom up and "help" each other
- The layers converge to special (critical?) points on the IB bound
- The advantage of the layers is mostly computational
 - Relaxation times are super-linear (exponential?) in the Entropy gap
 - Hidden layers provide intermediate steps and boost convergence time
 - Hidden layers help in avoiding critical slowing down

Further directions

- Exactly solvable DNN models (through symmetry & group theory)
- New/better learning algorithms & design principles

- Predictions on the organization of biological layered networks ... 8/12/18



Thank you!

