

Contents

1 Representation and unsupervised learning: Introduction (Sanjeev Arora)	1
2 Apocryphal Results, Recent Results, and Open Problems in Neural Network Representations, Matus Telgarsky (UIUC)	3
2.1 Representation	4
2.2 Apocryphal/subspaces	5
2.3 Separation: k to k^2 via Δ	6
2.4 Algorithmic open questions	7
2.5 k to $k + 1$	8
2.6 Recurrent networks	9

1 Representation and unsupervised learning: Introduction (Sanjeev Arora)

I'll describe ways in which the current frameworks are not up to the task.

In statistical learning theory, we observe examples x_1, \dots, x_n with corresponding labels y_1, \dots, y_n . ERM is finding

$$\operatorname{argmin}_C \sum_{i=1}^n L_C(x_i, y_i) + R(C)$$

where L_C is the loss function. If x_1, \dots, x_n are iid distributed according to D and C comes from a class of “low complexity” (Rademacher, etc.) then the test loss is approximately the training loss.

This framework doesn't explain many things about deep nets. The following are mysteries.

1. Generalization in deep nets. See Moritz, Recht [Zha+16].¹ Take image data with *random* labels. Training a neural net with sufficient capacity (M nodes) achieves 0 training error. The neural net can fit even random labels.

But neural nets trained with backprop still generalizes.

2. Unrelated classes seem to help. ImageNet has 1000 classes, each with 1000 examples. Knowing additional classes helps training for a specific class.
3. Transfer learning. Train a deep net on images. The top 2 layers of the deep net has great features for many other visual task.
4. Zero-shot learning: BM Lake, R Salakhutdinov, JB Tenenbaum [LST15].² Train on data set with letters from 50 languages. Learn a new character. The program can learn well with one or two examples.

¹<https://arxiv.org/abs/1611.03530>

²<https://staff.fnwi.uva.nl/t.e.j.mensink/zsl2016/zslpubs/lake15science.pdf>

When you see a character, what does your mind do to remember it? Remember it as distinct strokes.

In theory papers we ignore the distribution. The statistical learning theorem works for every distribution. But the classifier is very related to the distribution.

Regularization can incorporate a prior: this is one way to relate with Bayesian approaches.

5. Domain adaptation: If the distributions are very different there are no guarantees. But the distributions can be different: we can train on ImageNet and use for X-ray classification.

The classifier is very related to the distribution. You can't think of choosing them independently. Clustering, classification with margin are examples. We want a more sophisticated language.

In text, we work with bag-of-words models.

Representation learning is finding a map from a data space (raw pixels, etc.) to a representation space, $x \mapsto h$. This could be a many-to-one map.

We need a language to talk about such maps.

We can use simple models like cluster models, mixtures of Gaussians.

The Bayesian view has been to describe representation learning as distribution learning. The task of learning h is finding the MLE for what generated x . Their notion of representation learning is distribution learning: minimize KL divergence. Why should KL divergence lead to good learning?

We have a loglinear topic model with dynamics which comes up with a vector for each word. How to do this for fMRI data, or small corpora (100 documents, too small to do topic model)? Once you've understood how to understand meanings using Wikipedia, you can do a simple domain adaptation.

We had fMRI data (50k dimensional vector every time step) of people watching a movie (Sherlock). The semantic content of the movie was replaced by human annotation every 5 seconds, 2000 in total. Now we correlate.

Evaluation: Presented a 50k vector, given 5 annotations, choose the correct one.

If we want to classify according to any linear classifier; then we have to preserve all the bits. The human has a semantic model to generate the text; there is a mapping from the semantic representation to the text; we try to find the reverse map.

Some discussions:

- Example: Clustering should take into account the task: how they are going to be used. Representations can be profoundly affected by what you want to do. In clustering there is the concrete task of matching the clustering algorithm with the task. Coming up with tools about how to pick the clustering algorithm is an important problem.
- Early stopping in NN gives better generalization. This is more general than NN. Which algorithms have this property? You can also do early stopping with gradient descent with kernel methods. In principle you can overfit but you do not. (Note: The surrogate might never reach 0, but the classification error can be 0.)

- Can we phrase GANs as a problem rather than an algorithm?
- Neural networks act a little like our brain. This is fascinating. If you look at details in the brain, it's different. Are we functionally capturing what's happening? Once you have a good implementation, you can steal any brain. Ex. Train a neural net to repeat blood flow and use it as a feature.
- Human-in-the-loop: Practitioners have a problem where they want to achieve human accuracy. If they fail to get training error to 0, make network bigger. Look at the generalization error on the validation set. If that's too high, add regularization. Play with network structure, get more data, etc.

What algorithms would work well with human-in-the-loop?

For nonconvex problem, often human creativity is required to get good initialization. To prove anything nontrivial, initialization plays a role.

See Curves dataset.

Some general topics:

- Nonconvex optimization
- Generalization vs. optimization for nonconvex models.
- Representation power of depth (measure other than number of nodes)
- Models for representations (GANs)
- Differentiable computing. (Real numbers, logic)

Layers are like function calls; recurrent nets are loops. There are primitive operations which you can compose: loop operator, multiply, divide. Kernel methods cannot represent these efficiently. You get power by reusing computation.

Many things you say about neural networks you can say about polynomials. Why don't people use polynomials in ML? There are certain polys you can write as a neural network.

How to interpret "Turing-complete" in the context of NN? Have a read-write memory. It can read things it wrote last time, write things for next time.

2 Apocryphal Results, Recent Results, and Open Problems in Neural Network Representations, Matus Telgarsky (UIUC)

3

³We use "Avi" notation:

10 is a small constant.

100 is big-Oh.

1000 is big-Oh with scary constant.

- Representation? If you use a certain class to represent your predictor, what are the benefits of that choice?
- Apocryphal/subspaces: Ancient results: subspaces can fit every function. Many people can state these results but are not comfortable with the proofs. Ex. Every continuous function can be approximated with polynomials. This can be used as a lemma to prove that a neural net with two layers can fit continuous functions.
- k to k^2 via Δ : There's always a benefit to adding more layers. We will see this by constructing a strange, pathological function. You can use this function to compute many polynomials, parity, etc.
- Applications of Δ
- Algorithmic interlude: I'll comment about algorithmic question: how to determine whether you should add a $k + 1$ layer. Even for polynomials this is hard.
- k to $k + 1$: Initially I thought this was the kind of technical problem that wastes the time of many people, but I don't even know the hard function looks like—my guesses have failed.
- Rational functions: are also approximable.
- Recurrent networks: there's a trivial construction of a recurrent network that cannot be computed by a deep net, because you can have loops.
Eduardo Sontag and Siegelmann: Turing-complete model. Empirical work has missed this.
- Generalization (?)
- Two-layer NN

4

2.1 Representation

How can we represent our function class?

For example there are 3 ways to represent polynomials.

1. By degree

$$\bigcup_{r \geq 1} \left\{ x \mapsto \sum_{|\alpha| \leq r} c_\alpha x^\alpha : c_\alpha \in \mathbb{R}^{\binom{r+d}{r}} \right\}.$$

⁴ Keep in mind:

1. View theory as leading the way, not just analyzing existing algorithms.
2. We shouldn't lock ourselves in rooms.

2. By kernel

$$\bigcup_{r \geq 1} \bigcup_{n \geq 1} \left\{ x \mapsto \sum_{i=1}^n c_i k_r(x, x_i) : (x_i)_{i=1}^n \subseteq \mathbb{R}^{nd}, c \in \mathbb{R}^d \right\}$$

where $k_r(x, z)$ is the inner product of $[1, x_1, \dots, x_1^2, x_1 x_2, \dots, x_n^r]$ with $[1, z_1, \dots]$.

3. Sum-product network.

There are different tradeoffs in ease of representability, generalization, etc.

What is the set of low-complexity polys? A sparse function in one representation can be dense in another.

Things we want to prove:

1. Equivalence: One representation can be approximated by another.
2. Inequivalence: There are functions in one class that cannot be approximated by another. Depending on prior, one class really is better.

Compare and contrast programming languages and ML languages.

- In PL, we want human readable/writeable. In ML we want machine learnable human interpretable? (But neural nets are not human interpretable...)
- In PL we have (debugging) tools to diagnose errors. In ML we don't have this, but generalization theory.
- In PL we have Turing completeness. In ML we haven't looked at this. (This isn't just about function approximation—Turing-complete machine can take variable length input and time.)

My conjecture on why NN are successful is because of representation: the things they represent are what we see in nature.

Challenge: Come up with a function class better than neural networks in some aspect that you can quantify.

2.2 Apocryphal/subspaces

Theorem 2.1. Let $RECT = \{x \mapsto \mathbb{1}[x \in \prod_{i=1}^d [a_i, b_i]]\}$. Then

$$\sup_{f: [0,1]^d \rightarrow \mathbb{R}, f \text{ continuous}} \inf_{g \in \text{span}(RECT)} \|f - g\|_1 = 0.$$

This implies “boosted decision trees” also suffice. (An indicator of a box is an intersection of halfspaces.) The size is $(\frac{L}{\epsilon})^d$. (Good luck doing anything with this.) (Information-theoretically this is the bound; pack bits into a grid.)

Instead of $g \in \text{span}(RECT)$,

- Can take g to be 3-layer neural nets. Such g can compute functions in RECT.

- For g polynomials, two proofs:
 1. use Bernstein polynomials. Control the L^∞ norm: break into 2 cases, close to point and far from point.
Weierstrass's original proof: convolve with a Gaussian, look at Taylor approximation. $x \mapsto \mathbb{E}_{\tau \sim N(0, \sigma^2)} f(x + \tau)$ is analytic; it has Taylor expansion.

“Discrete-time vs. continuous-time diffusion.”

There is a proof for 2-layer NN that reduces to polynomial approximation. See

- Hornik-Stinchcombe-White [HSW89]
- Cybenko [Cyb89]
- Funahashi [Fun89]

Reverse engineering the size estimate, I get $(\frac{L}{\epsilon})^{2d}$.

I don't think this is a neural net theorem; this is a theorem about a linear subspace $x \mapsto \sum_i c_i \sigma(a_i^T x + b_i)$.

2.3 Separation: k to k^2 via Δ

Theorem 2.2 ([Tel16]). • For all $k \geq 1$,

- there exists $f : [0, 1] \rightarrow [0, 1]$, f has $10k^2$ layers, 2 nodes per layer, and
- for all $g : \mathbb{R} \rightarrow \mathbb{R}$, with $\leq k$ layers, $\leq 2^k$ nodes,

$$\int_{[0,1]} |f(x) - g(x)|^2 dx \geq \frac{1}{100}.$$

Why care about number of nodes? All the complexity measures in terms of generalization scale with number of nodes. VC dimension scales as number of layers times parameters.

This works for many classes of functions; we prove it with ReLU. The fact that I'm using the uniform measure is amazing. The ReLU allows me to use it. If I make this statement for the sigmoid, I only know how to prove it with a special hand-crafted measure, and I approximate sigmoids with ReLUs.

If we fix the complexity in the correct matter, each of {shallow networks, deep networks} represents functions that can't be represented by the other.

Open: Find a function easy to write down as a shallow network and is hard to represent by a deep net. You can express it with d blowup in size; the claim is that factor is necessary. I have a candidate function I don't know how to show; it seems to be a coding theory problem.

5

⁵“It's what not what you teach, it's what they learn.”

Proof. Step 1: Find the hard function.

Observation: In a linear representation, the number of bumps is about the number of basis functions.

Define the triangle function

$$\Delta(x) = \begin{cases} 2x, & x \in [0, \frac{1}{2}] \\ 2(1-x), & x \in [\frac{1}{2}, 1] \\ 0, & \text{otherwise.} \end{cases} = \sigma(2\sigma(x) - 4\sigma(x - \frac{1}{2}))$$

where $\sigma(x) = \max\{0, x\}$ is ReLU.

This function has incredible properties in composition: It copies functions, the second copy in reverse. If a function has k bumps, then composing with Δ gives $2k$ bumps.

Thus $\Delta^{(k)}$ has 2^{k-1} bumps.

Take $f = \Delta^{(10k^2)}$.

Step 2: Upper bound the number of bumps for shallow networks.

Define a s -affine function as a piecewise linear function with s pieces.

The sum of s -affine and t -affine function is $s + t - 1$ affine. The composition of s and t -affine functions is at most st affine, and this is tight.

Adding functions slowly increases complexity; composing quickly increases complexity.

Thus a neural network computes $\leq (\text{nodes})^{(\text{layers})}$ affine pieces.

The shallow network can have at most $(2^k)^k$ affine pieces, compared to 2^{10k^2} .

Step 3: L_1 bound.

This is a counting argument. Break the shallow function into pieces based on when they cross $y = \frac{1}{2}$. On average $\Delta^{(10k^2)}$ has 10 times as many oscillations as the shallow function. In any such piece, about half of the triangles are on the wrong side. \square

This proof is robust: I can prove separations about polynomials, etc.

2.4 Algorithmic open questions

- Representations reachable by SGD or some other efficient method.

Rewrite the $k \rightarrow k^2$ theorem in terms of a function that is learnable efficiently by a net with k^2 layers such that the same algorithm for fewer layers means you need exponentially many nodes.

- Polytime algorithm to see if you want to add layers.

This is a pain. Consider a poly $x_1x_2x_3$, correlation with any other monomial is 0.

You need structural assumptions on the network to make such a theorem go through.

Can you make a boosting algorithm on layers? Residual networks are like boosting.

- Polytime test to remove layers.

- In practice, obtain $((x_i, y_i))_{i=1}^n$.
- Fit f_0 of high complexity.

- Now fit f_1 of lower complexity $((x_i, f_0(x_i)))_{i=1}^n$. Claim: this gets lower (generalization?) error.
- Now fit $((x_i, f_1(x_i)))_{i=1}^n$, etc.

6

People have prescribed sizes of networks they try.

2.5 k to $k + 1$

Find a function in $k + 1$ layers which if you try to approximate with k layers you need exponentially many nodes.

Why should this be true? Evidence:

- Eldan-Shamir 2016, $k = 2$. [ES15]
- Boolean circuits: known forever.
- An average-case depth hierarchy theorem for boolean circuits, Servedio et al. 2016 [RST15]

I need a candidate hard function.

Ideas:

1. Another open problem: characterize many hard functions.
2. Kolmogorov 1957. (Superposition theorem⁷) Consider (almost) bijections $[0, 1]^d \rightarrow [0, 1]$ —space filling curve. Neural nets can build bijections.

We need multivariate functions here.

This is much more approachable: Characterize many hard functions.

Given a function (e.g. the sawtooth function), we can write $+1, -1$ when it is above/below the midline.

Define a mapping from sequences to lines. Given $\{-1, 1\}^d$, consider the function that is a sequence of triangles going above ($+1$) or below (-1) the $\frac{1}{2}$ line. A NN **certifies** this sequence if it's equal to this graph.

What is the set of all sequences achievable with $10k^2$ nodes?

- $\Lambda^{(10k^2)}$ certifies the sequence $\underbrace{-1, +1, \dots}_{2^{10k^2-1}}$
- Any sequence of length $\leq k$.
- Can only do exponentially few sequences of length $\geq 100k^4$.

“Find a graphical grammar with what you can do with neural networks.” (Rong Ge)

When you go to high dimensions, this problem becomes nasty: how do you even carve up the space and define the language? If I could answer this, I think it would give enough intuition to solve k vs. $k + 1$.

⁶S Arora: This can be used in a semisupervised setting. Use humans to label a small set; use the trained network f_0 to label more images.

⁷https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Arnold_representation_theorem

2.6 Recurrent networks

- Computational model by Siegelmann-Sontag.

RNN's are Turing-complete, so there exists $O(1)$ RNN's that is not approximated by any deep net.

The RNN gets an input, and gives output and whether it's thinking (whether it's done), and passes state. It can churn as long as it wants—it determines when to stop thinking and get the next input.

To encode a Turing machine, f is a state transition function for the TM; the state contains the tape, head, etc. (you can pack it into a real number)

This has loops so it has crazy power. How do they work in practice?

Here is a function that cannot be approximated by a deep net of any size: $f(x, k) = \Delta^{\lg(k)}(x)$. The RNN computes $k/2$ (the “output” bit) and asks if $k \approx 1$ (the “done” bit). RNN's can do loops of unbounded length.

This model of computation is absurdly powerful.

This theorem is trivial; the intelligence part is the model.

Another open question: generalization for RNN's. Can you remove the $\sqrt{\cdot}$ dependence on length?

I'm most interested in the $k \rightarrow k + 1$ problem and recurrent nets (learning algorithms, etc.).

References

- [Cyb89] George Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of Control, Signals, and Systems (MCSS)* 2.4 (1989), pp. 303–314.
- [ES15] Ronen Eldan and Ohad Shamir. “The Power of Depth for Feedforward Neural Networks”. In: *arXiv preprint arXiv:1512.03965* (2015).
- [Fun89] Ken-Ichi Funahashi. “On the approximate realization of continuous mappings by neural networks”. In: *Neural networks* 2.3 (1989), pp. 183–192.
- [HSW89] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. “Multilayer feedforward networks are universal approximators”. In: *Neural networks* 2.5 (1989), pp. 359–366.
- [LST15] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. “Human-level concept learning through probabilistic program induction”. In: *Science* 350.6266 (2015), pp. 1332–1338.
- [RST15] Benjamin Rossman, Rocco A Servedio, and Li-Yang Tan. “An average-case depth hierarchy theorem for boolean circuits”. In: *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*. IEEE. 2015, pp. 1030–1048.
- [Tel16] Matus Telgarsky. “Benefits of depth in neural networks”. In: *arXiv preprint arXiv:1602.04485* (2016).

- [Zha+16] Chiyuan Zhang et al. “Understanding deep learning requires rethinking generalization”. In: *arXiv preprint arXiv:1611.03530* (2016).