Cognitive-Assisted Interactive Labeling

Francois Luus Research Scientist *IBM Research - Africa*





Overview

- **1. SETI RFI signal classification**
- 2. RFI signal clustering
- 3. Toolchain for Software 2.0
- 4. Cognitive-assisted interactive labeling





SETI RFI Signal Classification

IBM Research | Africa







SETI Institute – RFI Classification

Problem

Identify common RFI and prioritize follow-up measurements for anomalous RFI – while having no labels

Solution

Unsupervised feature extraction + clustering
 Cognitive-assisted interactive labeling

Methods Convolutional AutoEncoding, t-SNE

~1 Million

Signal events

IBM Confidential | Do not distribute | © 2018 IBM Corporation

Advanced and Applied AI (IBM Research - Africa)

Unsupervised RFI Classification

- 14 Million subband spectrogram/images
- 90-95% no signal
- Reoccurring RFI is prevalent
- ~20 common types of RFI spectra
- No labels



RFI Spectrograms



RFI Spectrograms



RFI Spectrograms





RFI Signal Clustering

IBM Research | Africa





RFI Clustering



Spectrogram samples - Original density



Spectrogram samples – Equalized density (Stat+PCA)



Spectrogram Fourier features



Spectrogram samples – Equalized density (Feature engineering)



RFI Clustering



Convolutional Autoencoding: Architectures



Autoencoding Comparison: Encoding Capacity

Dense	Convolutional	Dense Variational	Convolutional Variational		
			And the second s		

RFI Clustering





Laurens van der Maaten

Research scientist in machine learning and computer vision.

- 🖸 Email
- **F**acebook
- & Google+
- **O** Github

t-SNE

t-Distributed Stochastic Neighbor Embedding (t-SNE) is a (<u>prize-winning</u>) technique for dimensionality reduction that is particularly well suited for the visualization of high-dimensional datasets. The technique can be implemented via Barnes-Hut approximations, allowing it to be applied on large real-world datasets. We applied it on data sets with up to 30 million examples. The technique and its variants are introduced in the following papers:

- L.J.P. van der Maaten. Accelerating t-SNE using Tree-Based Algorithms. Journal of Machine Learning Research 15(Oct):3221-3245, 2014. PDF [Supplemental material]
- L.J.P. van der Maaten and G.E. Hinton. Visualizing Non-Metric Similarities in Multiple Maps. Machine Learning 87(1):33-55, 2012.
 <u>PDF</u>
- L.J.P. van der Maaten. Learning a Parametric Embedding by Preserving Local Structure. In Proceedings of the Twelfth International Conference on Artificial Intelligence & Statistics (AI-STATS), JMLR W&CP 5:384-391, 2009. Apple
- L.J.P. van der Maaten and G.E. Hinton. Visualizing High-Dimensional Data Using t-SNE. Journal of Machine Learning Research 9(Nov):2579-2605, 2008. PDF [Supplemental material] [Talk]

SNE (2002) of MNIST







Low-dimensional embedding



Gradient comparison

Low-dimensional distance >	1: 1: 1: <th></th><th>1 0.5 0 -0.8</th>		1 0.5 0 -0.8
High-dimensional distance >		High-dimensional distance >	

(a) Gradient of SNE.

(c) Gradient of t-SNE.

From: van der Maaten, 2008

High-dimension & Similarity reproduction

<u>High-dimensional input</u>

$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)},$$

Reproduce similarities

$$C = \sum_{i} KL(P_i||Q_i) = \sum_{i} \sum_{j} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}},$$

Curse of dimensionality



Autoencoding: tSNE plots



RFI Clustering



Mean-shift DBSCAN Affinity Average Mini-batch Birch Spectral Ward Propagation clustering K-means linkage 0.1 0.3 0.6 10 10 10 10 10 0.08 0.25 0.7 20 20 20 20 20 are of 100 17 -0.06 0.2 0.8 30 30 30 30 30 0.15 0.04 0.9 40 40 40 40 40 -53 0.02 0.1 0.95 50 50 50 50 50 15

tSNE of comparison Clustering

Clustering: Qualitative inspection



RFI Clustering



For n_clusters = 80 The average silhouette_score is : 0.45



Cluster sample 1



Clusters sample 2





Toolchain for Software 2.0

IBM Research | Africa





Building the Software 2.0 Stack by Andrej Karpathy from Tesla



Software 2.0



(Karpathy, 2017)



(Karpathy, 2017)

Software 2.0

Andrej Karpathy (2017)

Software 2.0:

- "Software written by optimization in the language of weights"
- Engineered systems (1.0) -> Learned systems (2.0)
- End-to-end optimized systems
- Model skeletons + optimizers + supervised data
- Data labelers are the programmers



Toolchain for Software 2.0

Surrounding dataset/supervision infrastructure

IDE for Software 2.0 (Karpathy, 2017)

- Show full inventory/stats of current dataset
- Create/edit annotation layers for any datapoint
- Flag, escalate resolve discrepancies in multiple labels
- Flag, escalate datapoints thate are likely to be mislabeled
- Display predictions on an arbitrary set of test datapoints
- Autosuggest datapoints that should be labeled

Cognitive-assisted interactive labeling

Objectives

- 1. Visualize the quality of feature space
- 2. Define constraints/characteristics, could be through direct labeling
- 3. Update the feature space

<u>Requirements</u>

- 1. Iterative optimization
- 2. Dimensionality reduction for visualization
- 3. Fine-grained definition of sample characteristics

Cognitive-assisted interactive labeling

IBM Research | Africa











TensorFlov	V ™ Install	Develop	API r1.4	Deploy	Extend	Community	Versions	TFRC	Q	Search
Develop										
GET STARTED	PROGRAMMER'S GUIDE	TUTORIAL	.S PER	FORMANCE	MOBILE					

Getting Started

Getting Started With TensorFlow

MNIST For ML Beginners

Deep MNIST for Experts

TensorFlow Mechanics 101

tf.estimator Quickstart

Building Input Functions with tf. estimator

TensorBoard: Visualizing Learning

TensorBoard: Graph Visualization TensorBoard Histogram Dashboard

TensorFlow Versions

TensorBoard: Visualizing Learning

The computations you'll use TensorFlow for - like training a massive deep neural network - can be complex and confusing. To make it easier to understand, debug, and optimize TensorFlow programs, we've included a suite of visualization tools called TensorBoard. You can use TensorBoard to visualize your TensorFlow graph, plot quantitative metrics about the execution of your graph, and show additional data like images that pass through it. When TensorBoard is fully configured, it looks like this:



medium.com/tensorflow/interactive-supervision-with-tensorboard-9a101c91d3f7



Interactive supervision with TensorBoard

IBM Research AI



TensorBoard Projector showing labeled t-SNE

Originally published at www.ibm.com

TensorBoard

0



























Interactive Supervision with t-SNE



(k) 20% [CIFAR10] (l) 40% [CIFAR10] (m) 60% [CIFAR10]



Cognitive-assisted interactive labeling

- 1. Higher quality feature space means more homogeneous clusters
- 2. More homogeneous clusters means greater labeling efficiency
- 3. More labels means greater homogeneity
- 4. Edge cases and anomalies are obtained
- 5. Common classes defined



Questions





[D] Suggestion by Salesforce chief data scientist Discussion (i.redd.it)

submitted 21 hours ago by Prooffread3r to r/MachineLearning 72 comments share save hide report





Rather than spending a month figuring out an unsupervised machine learning problem, just label some data for a week and train a classifier.





Figure 1: Overview over the families of divergences and their relationship to each other. The shortcut *Prob.* denotes the special case of probability densities. For sake of clarity we show the most important relations only and do not claim completeness.

Digits



Fig. 6. Embeddings of all compared methods for Faces and Digits.

Barnes-Hut-SNE

Laurens van der Maaten Pattern Recognition and Bioinformatics Group, Delft University of Technology Mekelweg 4, 2628 CD Delft, The Netherlands lvdmaaten@gmail.com

Abstract

The paper presents an $\mathcal{O}(N \log N)$ -implementation of t-SNE — an embedding technique that is commonly used for the visualization of high-dimensional data in scatter plots and that normally runs in $\mathcal{O}(N^2)$. The new implementation uses vantage-point trees to compute sparse pairwise similarities between the input data objects, and it uses a variant of the Barnes-Hut algorithm to approximate the forces between the corresponding points in the embedding. Our experiments show that the new algorithm, called Barnes-Hut-SNE, leads to substantial computational advantages over standard t-SNE, and that it makes it possible to learn embeddings of data sets with millions of objects.

4.2 Approximating t-SNE Gradients

To approximate the t-SNE gradient, we start by splitting the gradient into two parts as follows:

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4(F_{attr} - F_{rep}) = 4\left(\sum_{j \neq i} p_{ij} q_{ij} Z(\mathbf{y}_i - \mathbf{y}_j) - \sum_{j \neq i} q_{ij}^2 Z(\mathbf{y}_i - \mathbf{y}_j)\right), \quad (8)$$

where F_{attr} denotes the sum of all attractive forces (the left sum), whereas F_{rep} denotes the sum of all repulsive forces (the right sum). Computing the sum of all attractive forces, F_{attr} , is computationally efficient; it can be done by summing over all non-zero elements of the sparse distribution P in $\mathcal{O}(uN)$. (Note that the term $q_{ij}Z = (1 + ||\mathbf{y}_i - \mathbf{y}_j||^2)^{-1}$ can be computed in $\mathcal{O}(1)$.) However, a naive computation of the sum of all repulsive forces, F_{rep} , is $\mathcal{O}(N^2)$. We now develop a Barnes-Hut algorithm to approximate F_{rep} efficiently in $\mathcal{O}(N \log N)$.

Consider three points \mathbf{y}_i , \mathbf{y}_j , and \mathbf{y}_k with $\|\mathbf{y}_i - \mathbf{y}_j\| \approx \|\mathbf{y}_i - \mathbf{y}_k\| \gg \|\mathbf{y}_j - \mathbf{y}_k\|$. In this situation, the contributions of \mathbf{y}_j and \mathbf{y}_k to F_{rep} will be roughly equal. The Barnes-Hut algorithm [1] exploits this by i) constructing a quadtree on the current embedding, ii) traversing the quadtree using a depth-first search, and iii) at every node in the quadtree, deciding whether the corresponding cell can be used as a "summary" for the gradient contributions of all points in that cell.





Quadtree. A quadtree is a tree in which each node represents a rectangular *cell* with a particular center, width, and height. Non-leaf nodes have four children that split up the cell into four smaller cells (quadrants) that lie "northwest", "northeast", "southwest", and "southeast" of the center of the parent node (see Figure 1 for an illustration). Leaf nodes represent cells that contain at most one point of the embedding; the root node represents the cell that contains the complete embedding. In each node, we store the center-of-mass of the embedding points that are located inside the corresponding cell, y_{cell} , and the total number of points that lie inside the cell, N_{cell} . A quadtree has $\mathcal{O}(N)$ nodes and can be constructed in $\mathcal{O}(N)$ time by inserting the points one-by-one, splitting a leaf node whenever a second point is inserted in its cell, and updating \mathbf{y}_{cell} and N_{cell} of all visited nodes.

Approximating the gradient. To approximate the repulsive part of the gradient, F_{rep} , we note that if a call is sufficiently small and suf-



Figure 1: Quadtree constructed on a twodimensional t-SNE embedding of 500 MNIST digits (the colors of the points correspond to the digit classes). Note how the quadtree adapts to the local point density in the embedding.

Optimization Equivalence of Divergences Improves Neighbor Embedding

Zhirong Yang² Jaakko Peltonen^{1,4} Samuel Kaski^{1,3}

ZHIRONG.YANG@AALTO.FI JAAKKO.PELTONEN@AALTO.FI SAMUEL.KASKI@AALTO.FI

¹Helsinki Institute for Information Technology HIIT, ²Department of Information and Computer Science, Aalto University, Finland, ³Department of Computer Science, University of Helsinki, and ⁴University of Tampere

Abstract

Visualization methods that arrange data objects in 2D or 3D layouts have followed two main schools, methods oriented for graph layout and methods oriented for vectorial embedding. We show the two previously separate approaches are tied by an optimization equivalence, making it possible to relate methods from the two approaches and to build new methods that take the best of both worlds. In detail, we prove a

1. Introduction

We address two research problems: nonlinear dimensionality reduction (NLDR) of vectorial data and graph layout. In NLDR, given a set of data points represented with highdimensional feature vectors or a distance matrix between such vectors, low-dimensional coordinates are sought for each data point. In graph layout, given a set of nodes (vertices) and a set of edges between node pairs, the task is to place the nodes on a 2D or 3D display. Solutions to both research problems are widely used in data visualization.

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4(F_{attr} - F_{rep}) = 4\left(\sum_{j \neq i} p_{ij} q_{ij} Z(\mathbf{y}_i - \mathbf{y}_j) - \sum_{j \neq i} q_{ij}^2 Z(\mathbf{y}_i - \mathbf{y}_j)\right)$$

amples). In contrast, ws-SNE uses edge repulsion to handle such cases of imbalanced degrees. When the d_i are not uniform, ws-SNE behaves differently from conventional s-SNE, which can be explained by its gradient $rac{\partial \mathcal{J}_{ ext{ws-SNE}}(Y)}{\partial y_i} \,=\, \sum_j p_{ij} q^{ heta}_{ij}(y_i-y_j) - c d_i d_j q^{1+ heta}_{ij}(y_i-y_j),$ where $c = \sum_{ij} p_{ij} / (\sum_{ij} d_i d_j q_{ij})$ is the connection scalar, and $\theta = 0$ for Gaussian q and $\theta = 1$ for Cauchy q. The first term in the summation is for attraction of nodes and the second for repulsion. Compared with the s-SNE gradient, the repulsion part is weighted by $d_i d_j$ in ws-SNE. That is, important nodes have extra repulsive force with the others and thus tend to be placed farther. This edge-repulsion strategy has been shown to be effective in graph drawing to overcome the "crowding problem", namely, many mapped points becoming crowded in the center of the display.

egy from Noack (2007), we insert weights M in the repulsive term: $\mathcal{J}_{\text{weighted-EE}}(Y) = \sum_{ij} p_{ij} ||y_i - y_j||^2 + \lambda \sum_{ij} M_{ij} \exp(-||y_i - y_j||^2)$, where $M_{ij} = d_i d_j$, and the vector d measures importance of the nodes. We use degree centrality as the measurement, i.e., d_i =degree of the *i*-th node. $\mathcal{J}_{\text{weighted-EE}}(Y)$ has downsides: it needs a userset edge repulsion weight λ , and is not invariant to scaling of p. By Theorem 1 we create a corresponding improved method, ws-SNE, minimizing a nonseparable divergence.

Proposition 5. Weighted EE is a separable divergence minimizing method and its non-separable variant is ws-SNE. **Proof:** Writing $q_{ij} = \exp(-||y_i - y_j||^2)$ with $q_{ii} = 0$, we have $\mathcal{J}_{\text{weighted-EE}}(Y) = D_{\mathrm{I}}(p||\lambda M \circ$ $q) + F(\lambda)$, where \circ denotes element-wise product and $F(\lambda) = C(\lambda) + \sum_{ij} p_{ij} \ln d_i d_j$ is a constant to Y. In the final and most important step, by a special case of Theorem 1: $\arg \min_Y D_{\mathrm{KL}}(p||M \circ q) =$ $\arg \min_Y [\min_{\lambda \ge 0} D_{\mathrm{I}}(p||\lambda M \circ q)]$, we obtain a new variant of SNE which minimizes $D_{\mathrm{KL}}(p||M \circ q)$ over Y:

$$\mathcal{T}_{ ext{ws-SNE}}(Y) = -\sum_{ij} p_{ij} \ln q_{ij} + \ln \sum_{ij} M_{ij} q_{ij} + ext{constant}$$

http://proceedings.mlr.press/v32/yange14.pdf