



- (51) International Patent Classification:
G06N 3/08 (2006.01) G06N 3/02 (2006.01)
- (21) International Application Number:
PCT/US2019/013870
- (22) International Filing Date:
16 January 2019 (16.01.2019)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:
62/618,440 17 January 2018 (17.01.2018) US
62/792,648 15 January 2019 (15.01.2019) US
- (71) Applicant: UNLEARN AI, INC. [US/US]; 650 California Street, Seventh Floor, San Francisco, CA 94108 (US).
- (72) Inventors: FISHER, Charles, Kenneth; 1025 Lombard Street, #2, San Francisco, CA 94109 (US). SMITH, Aaron, Michael; 455 Euclid Ave., No. 306, San Francisco, CA 94118 (US). WALSH, Jonathan, Avenu; 513 Lexington Avenue, Ei Cerrito, CA 94530 (US).
- (74) Agent: LEE, Paul, J.; KPPB LLP, 2190 S. Towne Centre Place, Suite 300, Anaheim, CA 92806 (US).
- (81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA,

(54) Title: SYSTEMS AND METHODS FOR MODELING PROBABILITY DISTRIBUTIONS

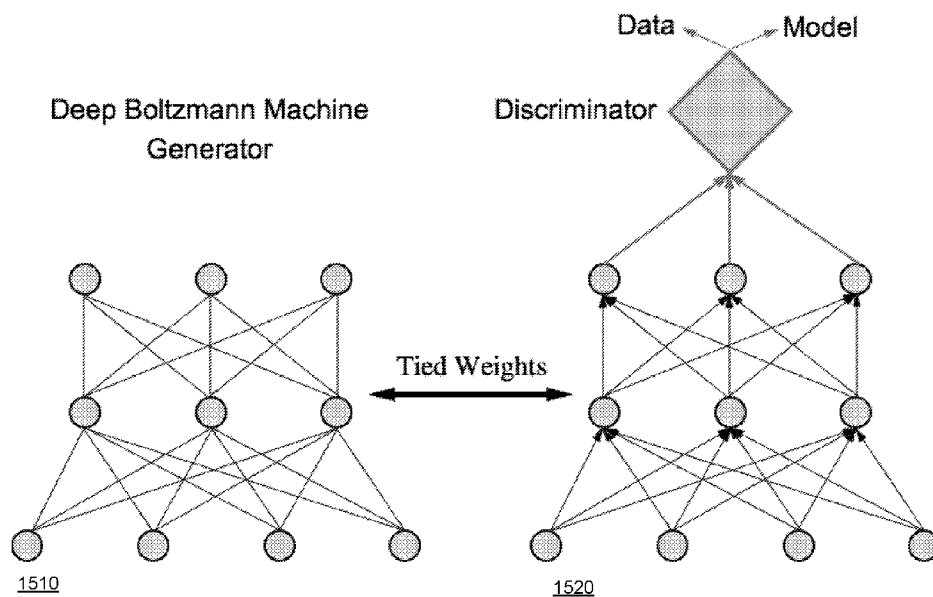


FIG. 15

(57) Abstract: Systems and methods for modeling complex probability distributions are described, One embodiment includes a method for training a restricted Boltzmann machine (RBM), wherein the method includes generating, from a first set of visible values, a set of hidden values in a hidden layer of a RBM and generating a second set of visible values in a visible layer of the RBM based on the generated set of hidden values. The method includes computing a set of likelihood gradients based on the first set of visible values and the generated set of visible values, computing a set of adversarial gradients using an adversarial model based on at least one of the set of hidden values and the set of visible values, computing a set of compound gradients based on the set of likelihood gradients and the set of adversarial gradients, and updating the RBM based on the set of compound gradients.



SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN,
TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

- (84) Designated States** (*unless otherwise indicated, for every kind of regional protection available*): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

Published:

- *with international search report (Art. 21(3))*

Systems and Methods for Modeling Probability Distributions

CROSS-REFERENCE TO RELATED APPLICATIONS

[0001] The present application claims the benefit of and priority to U.S. Provisional Patent Application No. 62/618,440 entitled 'Systems and Methods for Modeling Probability Distributions', filed January 17, 2018, and U.S. Provisional Patent Application No. 62/792,648 entitled 'Simulating Biological and Health Systems with Restricted Boltzmann Machines' filed January 15, 2019. The disclosure of U.S. Provisional Patent Application Serial Nos. 62/618,440 and 62/792,648 are herein incorporated by reference in their entirety.

FIELD OF THE INVENTION

[0002] The present invention generally relates to modeling probability distributions and more specifically relates to training and implementing a Boltzmann machine to accurately model complex probability distributions.

BACKGROUND

[0003] In a world of uncertainty, it is difficult to properly model probability distributions across multiple dimensions based on diverse and heterogeneous sets of data. For example, in the health industry, individual health outcomes are never certain. The condition of one patient with a disease may deteriorate rapidly, while another patient quickly recovers. The inherent stochasticity of individual health outcomes implies that health informatics must aim to predict health risks rather than deterministic outcomes. The ability to quantify and predict health risks has important implications for business models that depend on the health of a population.

SUMMARY OF THE INVENTION

[0004] Systems and methods for modeling complex probability distributions in accordance with

embodiments of the invention are illustrated. One embodiment includes a method for training a restricted Boltzmann machine (RBM), wherein the method includes generating, from a first set of visible values, a set of hidden values in a hidden layer of a RBM and generating a second set of visible values in a visible layer of the RBM based on the generated set of hidden values. The method also includes computing a set of likelihood gradients based on at least one of the first set of visible values and the generated set of visible values, computing a set of adversarial gradients using an adversarial model based on at least one of the set of hidden values and the set of visible values and computing a set of compound gradients based on the set of likelihood gradients and the set of adversarial gradients. The method includes updating the RBM based on the set of compound gradients.

[0005] In a further embodiment, the visible layer of the RBM includes a composite layer composed of a plurality of sub-layers for different data types.

[0006] In still another embodiment, the plurality of sub-layers includes at least one of a Bernoulli layer, an Ising layer, a one-hot layer, a von Mises-Fisher layer, a Gaussian layer, a ReLU layer, a clipped ReLU layer, a student-t layer, an ordinal layer, an exponential layer, and a composite layer.

[0007] In a still further embodiment, the RBM is a deep Boltzmann machine (DBM), wherein the hidden layer is one of a plurality of hidden layers.

[0008] In yet another embodiment, the RBM is a first RBM and the hidden layer is a first hidden layer of the plurality of hidden layers. The method further includes sampling the hidden layer from the first RBM, stacking the visible layer and the hidden layer from the first RBM into a vector, training a second RBM, and generating the DBM by copying weights from the first and second RBMs to the DBM. The vector is a visible layer of the second RBM.

[0009] In a yet further embodiment, the method further includes steps for receiving a phenotype vector for a patient, using the RBM to generate a time progression of a disease, and treating the patient based on the generated time progression.

[0010] In another additional embodiment, the visible layer and the hidden layer are for a first time instance, wherein the hidden layer is further connected to a second hidden layer that incorpo-

rates data from a different second time instance.

[0011] In a further additional embodiment, the visible layer is a composite layer includes data for a plurality of different time instances.

[0012] In another embodiment again, computing the set of likelihood gradients includes performing Gibbs sampling.

[0013] In a further embodiment again, the set of compound gradients are weighted averages of the set of likelihood gradients and the set of adversarial gradients.

[0014] In still yet another embodiment, the method further includes steps for training the adversarial model by drawing data samples based on authentic data, drawing fantasy samples based from the RBM, and training the adversarial model based on the adversarial model's ability to distinguish between the data samples and the fantasy samples.

[0015] In a still yet further embodiment, training the adversarial model includes measuring a probability that a particular sample is drawn from either the authentic data or the RBM.

[0016] In still another additional embodiment, the adversarial model is one of a fully-connected classifier, a logistic regression model, a nearest neighbor classifier, and a random forest.

[0017] In a still further additional embodiment, the method further includes steps for using the RBM to generate a set of samples of a target population.

[0018] In still another embodiment again, computing a set of likelihood gradients includes computing a convex combination of a Monte Carlo estimate and a mean field estimate.

[0019] In a still further embodiment again, computing a set of likelihood gradients includes initializing a plurality of samples and initializing an inverse temperature for each sample of the plurality of samples. For each sample of the plurality of samples, computing a set of likelihood gradients further includes updating the inverse temperature by sampling from an autocorrelated Gamma distribution, and updating the sample using Gibbs sampling.

[0020] Additional embodiments and features are set forth in part in the description that follows, and in part will become apparent to those skilled in the art upon examination of the specification or may be learned by the practice of the invention. A further understanding of the nature and

advantages of the present invention may be realized by reference to the remaining portions of the specification and the drawings, which forms a part of this disclosure.

BRIEF DESCRIPTION OF THE DRAWINGS

[0021] The description and claims will be more fully understood with reference to the following figures and data graphs, which are presented as exemplary embodiments of the invention and should not be construed as a complete recitation of the scope of the invention.

[0022] Figure 1 illustrates a system that provides for the gathering and distribution of data for modeling probability distributions in accordance with some embodiments of the invention.

[0023] Figure 2 illustrates a data processing element for training and utilizing a stochastic model.

[0024] Figure 3 illustrates a data processing application for training and utilizing a stochastic model.

[0025] Figure 4 conceptually illustrates a process for preparing data for analysis.

[0026] Figure 5 illustrates data structures for implementing a generalized Boltzmann Machine in accordance with certain embodiments of the invention.

[0027] Figure 6 illustrates a bimodal distribution and a smoothed, spread distribution that is learned by a RBM distribution in accordance with several embodiments of the invention.

[0028] Figure 7 illustrates an architecture for a generalized Restricted Boltzmann Machine in accordance with some embodiments of the invention.

[0029] Figure 8 illustrates a schema for implementing a generalized Boltzmann Machine in accordance with certain embodiments of the invention.

[0030] Figure 9 illustrates an architecture for a generalized Deep Boltzmann Machine in accordance with certain embodiments of the invention.

[0031] Figure 10 conceptually illustrates a process for reverse layerwise training in accordance with an embodiment of the invention.

[0032] Figure 11 illustrates an architecture for a generalized Deep Temporal Boltzmann Machine

in accordance with many embodiments of the invention.

[0033] Figure 12 conceptually illustrates a process for training a Boltzmann Encoded Adversarial Machine in accordance with some embodiments of the invention.

[0034] Figure 13 illustrates resulting samples drawn from RBMs trained to maximize log likelihood and from RBMs trained as BEAMs.

[0035] Figure 14 illustrates results of training a BEAM on a 2D mixture of Gaussians in accordance with a number of embodiments of the invention.

[0036] Figure 15 illustrates an architecture for implementing a Boltzmann Encoded Adversarial Machine in accordance with a number of embodiments of the invention.

[0037] Figure 16 illustrates a comparison between samples drawn from a Boltzmann machine with regular Gibbs sampling to those drawn using Temperature Driven Sampling.

[0038] Figure 17 illustrates a a comparison between fantasy particles generated by GRBMs trained on the MNIST dataset using regular Gibbs sampling to those using TDS.

DETAILED DESCRIPTION

[0039] Machine learning is one potential approach to modeling complex probability distributions. In the following description, many examples are described with reference to medical applications, but one skilled in the art will recognize that techniques described herein can be readily applied in a variety of different fields including (but not limited to) health informatics, image/audio processing, marketing, sociology, and lab research. One of the most pressing problems is that one often has little, or no, labeled data that directly addresses a particular question of interest. Consider the task of predicting how a patient will respond to an investigational therapeutic in a clinical trial. In a supervised learning setting, one would give the therapeutic to many patients and observe how each patient responds. Then, one would use this data to build a model that predicts how a new patient will respond to the therapeutic. For example, a nearest neighbor classifier would look through the pool of previously treated patients to find a patient that is most similar to the new patient,

then it would predict the new patient's response based on the previously treated patient's response. However, supervised learning requires significant amounts of labeled data and, particularly where sample sizes are small or labeled data is not readily available, unsupervised learning is critical to the successful application of machine learning.

[0040] Many machine learning applications, such as computer vision, require the use of homogeneous information (e.g., images of the same shape and resolution), which must be pre-processed or otherwise manipulated to normalize the input and training data. However, in many applications it is desirable to combine data of various types (e.g., images, numbers, categories, ranges, text samples, etc.) from many sources. For example, medical data can include a variety of different types of information from a variety of different sources, including (but not limited to) demographic information (e.g., a patient's age, ethnicity, etc.), diagnoses (e.g., binary codes that describe whether or not a patient has a particular disease), laboratory values (e.g., results from laboratory tests, such as blood tests), doctor's notes (e.g., hand written notes taken by a physician or entered into a medical records system), images (e.g., x-rays, CT scans, MRIs, etc.), and 'omics data (e.g., data from DNA sequencing studies that describe a patient's genetic background, the expression of his/her genes, etc.). Some of these data are binary, some are continuous, and some are categorical. Integrating all of these different types and sources of data is critical, but treating a variety of data types with traditional approaches to machine learning is quite challenging. Typically, the data have to be heavily pre-processed so that all of the features used for machine learning are of the same type. Data pre-processing steps can take up a large portion of an analyst's time in training and implementing a machine learning model.

[0041] In addition to processing many different types of data, the data used for an analysis is often incomplete or irregular. In the example of medical data, physicians often do not run the same set of tests on every patient (though, clinical trials are an important exception). Instead, a doctor will order a test if he/she has a specific concern about the patient. Therefore, medical records contain many fields with missing observations. But, these observations may not be missing at random. Handling these missing observations is an important part of any application of machine learning in

health care.

[0042] There are two implications of missing data for machine learning in healthcare. First, any algorithm needs to be able to learn from data where there are missing observations in the training set. Second, the algorithm needs to be able to make predictions even when it is only presented with a subset of input observations. That is, one needs to be able to express any conditional relationship from the joint probability distribution.

[0043] One approach that has recently gained a lot of popularity is the use of Generative Adversarial Networks (GANs). GANs, in their traditional formulation, use a generator that transforms random Gaussian noise into a visible vector through a feed-forward neural network. Models with this formulation can be trained using the standard back-propagation process. However, GAN training tends to be unstable – requiring a careful balance between training of the generator and the discriminator (or critic). Moreover, it is not possible to generate samples from arbitrary conditional distributions with GANs, and it can be very difficult to apply GANs to problems involving heterogeneous datasets with different data types and missing observations.

[0044] Many embodiments of the invention provide novel and innovative systems and methods for the use of heterogeneous, irregular, and unlabeled data to train and implement stochastic, unsupervised machine learning models of complex probability distributions.

System for Modeling Probability Distributions

[0045] Turning now to the drawings, a system that provides for the gathering and distribution of data for modeling probability distributions in accordance with some embodiments of the invention is shown in Figure 1. Network 100 includes a communications network 160. The communications network 160 is a network such as the Internet that allows devices connected to the network 160 to communicate with other connected devices. Server systems 110, 140, and 170 are connected to the network 160. Each of the server systems 110, 140, and 170 is a group of one or more servers communicatively connected to one another via internal networks that execute processes that provide cloud services to users over the network 160. For purposes of this discussion, cloud

services are one or more applications that are executed by one or more server systems to provide data and/or executable applications to devices over a network. The server systems 110, 140, and 170 are shown each having three servers in the internal network. However, the server systems 110, 140 and 170 may include any number of servers and any additional number of server systems may be connected to the network 160 to provide cloud services. In accordance with various embodiments of this invention, a network that uses systems and methods that model complex probability distributions in accordance with an embodiment of the invention may be provided by a process (or a set of processes) being executed on a single server system and/or a group of server systems communicating over network 160.

[0046] Users may use personal devices 180 and 120 that connect to the network 160 to perform processes for providing and/or interaction with a network that uses systems and methods that model complex probability distributions in accordance with various embodiments of the invention. In the shown embodiment, the personal devices 180 are shown as desktop computers that are connected via a conventional “wired” connection to the network 160. However, the personal device 180 may be a desktop computer, a laptop computer, a smart television, an entertainment gaming console, or any other device that connects to the network 160 via a “wired” connection. The mobile device 120 connects to network 160 using a wireless connection. A wireless connection is a connection that uses Radio Frequency (RF) signals, Infrared signals, or any other form of wireless signaling to connect to the network 160. In Figure 1, the mobile device 120 is a mobile telephone. However, mobile device 120 may be a mobile phone, Personal Digital Assistant (PDA), a tablet, a smart-phone, or any other type of device that connects to network 160 via wireless connection without departing from this invention.

[0047] A data processing element for training and utilizing a stochastic model in accordance with a number of embodiments is illustrated in Figure 2. In various embodiments, data processing element 200 is one or more of a server system and/or personal devices within a networked system similar to the system described with reference to Figure 1. Data processing element 200 includes a processor (or set of processors) 210, network interface 225, and memory 230. The network inter-

face 225 is capable of sending and receiving data across a network over a network connection. In a number of embodiments, the network interface 225 is in communication with the memory 230. In several embodiments, memory 230 is any form of storage configured to store a variety of data, including, but not limited to, a data processing application 232, data files 234, and model parameters 236. Data processing application 232 in accordance with some embodiments of the invention directs the processor 210 to perform a variety of processes, such as (but not limited to) using data from data files 234 to update model parameters 236 in order to model complex probability distributions.

[0048] A data processing application in accordance with a number of embodiments of the invention is illustrated in Figure 3. In this example, data processing element 300 includes a data gathering engine 310, database 320, a model trainer 330, a generative model 340, a discriminator model 350, and a simulator engine 345. Model trainer 330 includes a schema processor 332 and a sampling engine 334. Data processing applications in accordance with many embodiments of the invention process data to train stochastic models that can be used to model complex probability distributions.

[0049] Data gathering engines in accordance with many embodiments of the invention gather data from various sources in various formats. The gathered data in accordance with many embodiments of the invention include data that may be heterogeneous (e.g., data with various types, ranges, and constraints) and/or incomplete. One skilled in the art will recognize that various types and amounts of data can be utilized as appropriate to the requirements of specific applications in accordance with embodiments of the invention. In some embodiments, data gathering engines are further for pre-processing the data to facilitate the training of the model. However, unlike pre-processing performed in other methods, pre-processing in accordance with some embodiments of the invention is automatically performed based on a datatype and/or a schema associated with each data input. For example, in certain embodiments, bodies of unstructured text (e.g., typed medical notes, diagnoses, free-form questionnaire responses, etc.) are processed in a variety of ways, such as (but not limited to) vectorization (e.g., using word2vec), summarization, sentiment

analysis, and/or keyword analysis. Other pre-processing steps can include (but are not limited to) normalization, smoothing, filtering, and aggregation. In some embodiments, the pre-processing is performed using various machine learning techniques, including (but not limited to) Restricted Boltzmann machines, support vector machines, recurrent neural networks, and convolutional neural networks.

[0050] Databases in accordance with various embodiments of the invention store data for use by data processing applications, including (but not limited to) input data, pre-processed data, model parameters, schemas, output data, and simulated data. In some embodiments, databases are located on separate machines (e.g., in cloud storage, server farms, networked databases, etc.) from a data processing application.

[0051] Model trainers in accordance with a number of embodiments of the invention are used to train generative and/or discriminator models. In many embodiments, model trainers utilize schema processors to build the generator and/or discriminator models based on schemas that are defined for the various data available to the system. Schema processors in accordance with some embodiments of the invention build composite layers for a generative model (e.g., restricted Boltzmann machine) that are made up of several different layers for handling different types of data in different ways. In some embodiments, model trainers train the generative and discriminator models by optimizing a compound objective function based on a log-likelihood and adversarial objectives. Training generative models in accordance with certain embodiments of the invention utilizes sampling engines to draw samples from the models to measure the probability distributions of the data and/or the models. Various methods for sampling from such models to train and/or draw generated samples from a model are described in greater detail below.

[0052] In many embodiments, generative models are trained to model complex probability distributions, which can be used to generate predictions/simulations of various probability distributions. Discriminator models discriminate between data-based samples and model-generated samples based on the visible and/or hidden states.

[0053] Simulator engines in accordance with several embodiments of the invention are used to

generate simulations of complex probability distributions. In some embodiments, simulator engines are used to simulate patient populations, disease progressions, and/or predicted responses to various treatments. Simulator engines in accordance with several embodiments of the invention use a sampling engine for drawing samples from the generative models that simulate the probability distribution of the data.

[0054] As described above, as a part of the data gathering process, the data in accordance with several embodiments of the invention is pre-processed in order to simplify the data. Unlike other pre-processing which is often highly manual and specific to the data, this can be performed automatically based on the type of data, without additional input from another person.

[0055] A process for preparing data for analysis in accordance with some embodiments of the invention is conceptually illustrated in Figure 4. The process 400 processes (405) unstructured data. Unstructured data in accordance with many embodiments of the invention can include various types of data that can be pre-processed in order to speed up processing and/or to reduce the memory requirements for storing the relevant data. Examples of such data can include (but are not limited to) bodies of text, signal processing data, audio data, and image data. Processing unstructured data in accordance with many embodiments of the invention can include (but is not limited to) feature identification, summarization, keyword detection, sentiment analysis, and signal analysis.

[0056] The process 400 reorders (410) the data based on a schema. In certain embodiments, processes reorder the data based on the different data types defined in schemas by grouping similar data types to allow for efficient processing of the data types. The process 400 in accordance with some embodiments of the invention rescales (415) the data to prevent the overrepresentation of certain data elements based purely on the scale of the measurements. Process 400 then routes (420) the pre-processed data to the sublayers of a Boltzmann machine that are structured based on data types identified in the schema. Examples of Boltzmann machine structures and architectures are described in greater detail below. In some embodiments, the data is pre-processed into temporally sequenced data structures for inputs to a deep temporal Boltzmann machine. Deep temporal

Boltzmann machines are described in further detail below.

[0057] Temporal data structures for inputs to a Boltzmann machine in accordance with a number of embodiments of the invention are illustrated in Figure 5. The example of Figure 5 shows three data structures 510, 520, and 530. Each of the data structures represents a set of the data values captured at a particular time (i.e., times t_0 , t_1 , and t_n). In this example, certain traits (e.g., gender, ethnicity, birthdate, etc.) do not usually change over time, while other characteristics (e.g., test results, medical scans, etc.) do change over time. The example further shows that certain data may be missing for some fields for certain times for certain individuals. In this example, each individual is assigned a separate identification number in order to maintain patient confidential information.

Boltzmann Encoded Adversarial Machines

[0058] Models trained to minimize forward KL divergence, $D_{KL}(p_{data}||p_{\theta})$, tend to spread the model distribution out to cover the support of the data distribution. An example of a spread distribution is illustrated in Figure 6. Specifically, Figure 6 illustrates a bimodal distribution 610 and the pretty good, smoothed, spread distribution that is learned by a RBM distribution 620. While RBMs are able to generate such good approximations, they can struggle when faced with finer, more complex distributions.

[0059] To overcome the problems with traditional Boltzmann machines, several embodiments of the invention implement a framework for training Boltzmann machines against an adversary, referred to herein as a Boltzmann Encoded Adversarial Machine (BEAM). A BEAM minimizes a loss function that is a combination of the negative log-likelihood and an adversarial loss. The adversarial component ensures that BEAM training performs a simultaneous minimization of both the forward and reverse KL divergences, which prevents the oversmoothing problem observed with regular RBMs.

Boltzmann Machine Architectures

[0060] With many traditional machine learning techniques, supervised learning is used to train

a model on a large set of labeled data to make predictions and classifications. However, in many cases, it is not feasible or possible to gather such large samples of labeled data. In many cases, the data cannot be readily labeled or there are simply not enough samples of an event to meaningfully train a supervised learning model. For example, clinical trials often face difficulties in gathering such labeled data. A clinical trial typically proceeds through three main phases. In phase I, the therapeutic is given to healthy volunteers to assess its safety. In phase II, the therapeutic is given to approximately 100 patients to obtain initial estimates for safety and efficacy. Finally, in phase III, the therapeutic is given to a few hundred to a few thousand patients to rigorously investigate the efficacy of the drug. Before phase II, there is no in-human data on the effect of the investigational drug for the desired indication, making supervised learning impossible. After phase II, there is some in-human data on the effect of the investigational drug, but the sample size is quite limited, rendering supervised learning techniques ineffective. For comparison, a phase II clinical trial may have 100-200 patients, whereas a typical application of machine learning in computer vision may use millions of labeled images. As with many situations with limited data, the lack of large labeled datasets for many important problems implies that health informatics must heavily rely on methods for unsupervised learning.

Restricted Boltzmann Machines (RBMs)

[0061] One machine learning model (or method) that uses unsupervised learning is a Restricted Boltzmann Machine (RBM). RBMs are bidirectional neural networks, where the neurons (also called units) are divided into two layers, a visible layer and a hidden layer. The visible layer \mathbf{v} describes the observed data. The hidden layer \mathbf{h} consists of a set of unobserved latent variables that capture the interactions between the visible units. The model describes the joint probability distribution of \mathbf{v} and \mathbf{h} using an exponential form,

$$p(\mathbf{v}, \mathbf{h}) = Z^{-1} e^{-E(\mathbf{v}, \mathbf{h})}. \quad (1)$$

Here, $E(\mathbf{v}, \mathbf{h})$ is called the energy function, and $Z = \int d\mathbf{v}d\mathbf{h}e^{-E(\mathbf{v}, \mathbf{h})}$ is called the partition function. In many embodiments, processes use the integral operator, $\int dx$, to denote both standard integration or a sum over all of the elements in a discrete set.

[0062] In a traditional RBM, both the visible and hidden units are binary. Each can only take on the values 0 or 1. The energy function can be written as,

$$E(\mathbf{v}, \mathbf{h}) = -\sum_i a_i v_i - \sum_\mu b_\mu h_\mu - \sum_{i\mu} W_{i\mu} v_i h_\mu \quad (2)$$

or, in vector notation, $E(\mathbf{v}, \mathbf{h}) = -\mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$. Notice that visible units interact with the hidden units through the weights, W . However, there are no visible-visible or hidden-hidden interactions.

[0063] A key feature of an RBM is that it is easy to compute the conditional probabilities,

$$p(\mathbf{v}|\mathbf{h}) = \prod_i \frac{e^{(a_i + \sum_\mu W_{i\mu} h_\mu) v_i}}{1 + e^{a_i + \sum_\mu W_{i\mu} h_\mu}} \quad (3)$$

and,

$$p(\mathbf{h}|\mathbf{v}) = \prod_\mu \frac{e^{(b_\mu + \sum_i W_{i\mu} v_i) h_\mu}}{1 + e^{b_\mu + \sum_i W_{i\mu} v_i}}. \quad (4)$$

Similarly, it is easy to compute the conditional moments,

$$\langle \mathbf{v} \rangle_{p(\mathbf{v}|\mathbf{h})} = \frac{1}{1 + e^{-(\mathbf{a} + \mathbf{W} \mathbf{h})}} \quad (5)$$

and,

$$\langle \mathbf{h} \rangle_{p(\mathbf{h}|\mathbf{v})} = \frac{1}{1 + e^{-(\mathbf{b} + \mathbf{W}^T \mathbf{v})}}. \quad (6)$$

However, it is generally very difficult to compute statistics from the joint distribution. As a result, statistics from the joint distribution have to be estimated using random sampling processes such as Markov Chain Monte Carlo (MCMC).

[0064] RBMs can be trained by maximizing the log-likelihood $\mathcal{L} := \langle \log p(\mathbf{v}) \rangle_{data} = \langle \log \int d\mathbf{h} p(\mathbf{v}, \mathbf{h}) \rangle_{data}$.

Here, $\langle \cdot \rangle_{data}$ denotes an average over all of the observed samples. The derivative of the log-likelihood with respect to some parameter of the model θ is:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \theta} &= \left\langle \frac{\partial}{\partial \theta} \log \int d\mathbf{h} p(\mathbf{v}, \mathbf{h}) \right\rangle_{data} \\
&= \left\langle \frac{\partial}{\partial \theta} \log \int d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})} \right\rangle_{data} - \frac{\partial}{\partial \theta} \log Z \\
&= \left\langle \frac{\int d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})} \left(-\frac{\partial E}{\partial \theta}\right)}{\int d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})}} \right\rangle_{data} - \frac{\int d\mathbf{v} d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})} \left(-\frac{\partial E}{\partial \theta}\right)}{\int d\mathbf{v} d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})}} \\
&= \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{p(\mathbf{v}, \mathbf{h})} - \left\langle \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{p(\mathbf{h}|\mathbf{v})} \right\rangle_{data} \tag{7}
\end{aligned}$$

In the standard formulation of an RBM, there are three parameters a , b , and W . The derivatives are:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial a} &= \langle \mathbf{v} \rangle_{p(\mathbf{v}, \mathbf{h})} - \langle \mathbf{v} \rangle_{data} \\
\frac{\partial \mathcal{L}}{\partial b} &= \langle \mathbf{h} \rangle_{p(\mathbf{v}, \mathbf{h})} - \left\langle \langle \mathbf{h} \rangle_{p(\mathbf{h}|\mathbf{v})} \right\rangle_{data} \\
\frac{\partial \mathcal{L}}{\partial W} &= \langle \mathbf{v} \mathbf{h}^T \rangle_{p(\mathbf{v}, \mathbf{h})} - \left\langle \langle \mathbf{v} \mathbf{h}^T \rangle_{p(\mathbf{h}|\mathbf{v})} \right\rangle_{data} \tag{8}
\end{aligned}$$

[0065] Computing expectations from the joint distribution is generally computationally intractable. Therefore, the derivatives have to be computed using samples from the model drawn with an MCMC process. Samples can be drawn from an RBM using alternating Gibbs sampling.

Input: Initial configuration (\mathbf{v}, \mathbf{h}) .
 A number of Monte Carlo steps, k .
 An RBM.

Output: A new configuration $(\mathbf{v}', \mathbf{h}')$.

```

set  $\mathbf{v}_0 = \mathbf{v}, \mathbf{h}_0 = \mathbf{h}$ ;
for  $i = 1, \dots, k$  do
  | draw  $\mathbf{h}_i \sim p(\mathbf{h}|\mathbf{v}_{i-1})$ ;
  | draw  $\mathbf{v}_i \sim p(\mathbf{v}|\mathbf{h}_i)$ ;
end
return  $(\mathbf{v}_k, \mathbf{h}_k)$ 

```

[0066] In theory, Gibbs sampling produces uncorrelated random samples from $p(\mathbf{v}, \mathbf{h})$ in the limit that $n \rightarrow \infty$. Of course, infinity is a long time. Therefore, the derivatives of the log-likelihood of an RBM are usually approximated using one of two processes: Contrastive Divergence (CD), or Persistent Contrastive Divergence (PCD). K-step CD is very simple: Grab a batch of data. Compute an approximate batch of samples from the model by running k-steps of Gibbs sampling starting from the data. Compute the gradients of the log-likelihood and update the model parameters. Importantly, the samples from the model are re-initialized using the batch of observed data for each gradient update. K-step PCD is similar: First, samples from the model are initialized using a batch of data. The samples are updated for k steps, the gradients are computed, and the parameters are updated. In contrast to CD, the samples from the model are never re-initialized. Many architectures of Boltzmann machines in accordance with several embodiments of the invention utilize sampling to compute derivatives for training the Boltzmann machines. Various methods for sampling in accordance with several embodiments of the invention are described in greater detail below.

Generalized RBMs

[0067] One challenge that arises in the use of traditional Boltzmann machines is that many RBMs use binary units, while much of the data that is to be processed can come in a variety of different

forms. To overcome this limitation, some embodiments of the invention use a generalized RBM. A generalized RBM in accordance with a number of embodiments of the invention is illustrated in Figure 7. The example of Figure 7 shows a generalized RBM 700 with a visible layer 710 and a hidden layer 720. The visible layer 710 is a composite layer comprised of several nodes of various types (i.e., continuous, categorical, and binary). The nodes of visible layer 710 are connected to nodes of hidden layer 720. Hidden layers of generalized RBMs in accordance with several embodiments of the invention operate as a low dimensional representation of individuals (e.g., patients in a clinical trial) based on the compiled inputs to a composite visible layer.

[0068] Generalized RBMs in accordance with a number of embodiments of the invention are trained with an energy function,

$$E(\mathbf{v}, \mathbf{h}) = -a(\mathbf{v}) - b(\mathbf{h}) - \mathbf{v}^T \frac{W}{(\sigma\varepsilon^T)^2} \mathbf{h} \quad (9)$$

where $a(\cdot)$ and $b(\cdot)$ are arbitrary functions, and $\sigma > 0$ and $\varepsilon > 0$ are scale parameters of the visible and hidden layers, respectively. Different functions (called layer types) are used to represent different types of data. Examples of layer types used for modeling various types of data are described below.

[0069] Bernoulli Layer: A Bernoulli layer is used to represent binary data $v_i \in \{0, 1\}$. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters are set to $\sigma_i = 1$.

[0070] Ising Layer: An Ising layer is a symmetrized Bernoulli layer for visible units $v_i \in \{-1, +1\}$. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters are set to $\sigma_i = 1$.

[0071] One-hot Layer: A one-hot layer represents data where $v_i \in \{0, 1\}$ and $\sum_i v_i = 1$. That is, one of the units is turned on and all of the other units are turned off. One-hot layers are commonly used to represent categorical variables. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters are set to $\sigma_i = 1$.

[0072] von Mises-Fisher Layer: A von Mises-Fisher layer represents data where $v_i \in [0, 1]$ and $\sum_i v_i^2 = 1$. That is, the units are confined to the surface of an n-dimensional sphere. This layer is

particularly useful for modeling fractional data where $x_i \in [0, 1]$ and $\sum_i x_i = 1$ because $v_i = \sqrt{x_i}$ satisfies the spherical property. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters are set to $\sigma_i = 1$.

[0073] Gaussian Layer: A Gaussian layer represents data where $v_i \in \mathbb{R}$. The bias function is $a(\mathbf{v}) = -\sum_i \frac{(v_i - \bar{v}_i)^2}{2\sigma_i^2}$. Both the location, \bar{v}_i , and scale, σ_i , parameters of the layer are generally trainable. In practice, it helps to parameterize the model in terms of $\log \sigma_i$ to ensure that the scale parameter stays positive.

[0074] ReLU Layer: A Rectified Linear Unit (ReLU) layer represents data where $v_i \in \mathbb{R}$ with $v_i \geq v_i^{low}$. In the context of a Boltzmann machine, a ReLU layer is essentially a one-sided truncated Gaussian layer. The bias function is $a(\mathbf{v}) = -\sum_i \frac{(v_i - \bar{v}_i)^2}{2\sigma_i^2}$ over the domain $v_i \geq v_i^{low}$. Both the location, \bar{v}_i , and scale, σ_i , parameters of the layer are generally trainable whereas v_i^{low} is typically specified before training. In practice, it helps to parameterize the model in terms of $\log \sigma_i$ to ensure that the scale parameter stays positive.

[0075] Clipped Relu Layer: A Clipped Rectified Linear Unit (ReLU) layer represents data where $v_i \in \mathbb{R}$ with $v_i^{high} \leq v_i \leq v_i^{low}$. In the context of a Boltzmann machine, a Clipped ReLU layer is essentially a two-sided truncated Gaussian layer. The bias function is $a(\mathbf{v}) = -\sum_i \frac{(v_i - \bar{v}_i)^2}{2\sigma_i^2}$ over the domain $v_i^{high} \leq v_i \leq v_i^{low}$. Both the location, \bar{v}_i , and scale, σ_i , parameters of the layer are generally trainable whereas v_i^{high} and v_i^{low} are typically specified before training. In practice, it helps to parameterize the model in terms of $\log \sigma_i$ to ensure that the scale parameter stays positive.

[0076] Student-t Layer: A Student-t distribution is similar to a Gaussian distribution, but has fatter tails. In a variety of embodiments, implementation of a Student-t layer is implicit. The layer has three parameters, a location parameter \bar{v}_i that controls the mean, a scale parameter v_i that controls the variance, and a degrees of freedom parameter d_i that controls the thickness of the tails. The layer is defined by drawing a variance $\sigma_i^2 \sim \text{InverseGamma}(\frac{d_i}{2}, \frac{d_i}{2v_i})$ and then taking the energy as $a(\mathbf{v}) = -\sum_i \frac{(v_i - \bar{v}_i)^2}{2\sigma_i^2}$.

[0077] Ordinal Layer: An Ordinal layer is a generalization of a Bernoulli layer that is used to represent integer valued data $v_i \in \{0, N_i\}$. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters

are set to $\sigma_i = 1$. The upper value N_i is specified ahead of time.

[0078] Gaussian-Ordinal Layer: A Gaussian-ordinal layer is a generalization of an ordinal layer that is used to represent integer valued data $v_i \in \{0, N_i\}$ with a more flexible distribution. The bias function is $a(\mathbf{v}) = -\sum_i \frac{(v_i - \bar{v}_i)^2}{2\sigma_i^2}$. The upper value N_i is specified ahead of time.

[0079] Exponential Layer: An exponential layer represents data where $v_i \in \mathbb{R}_+$. The bias function is $a(\mathbf{v}) = a^T \mathbf{v}$ and the scale parameters are set to $\sigma_i = 1$. Note, exponential layers have some constraints because $a_i + \sum_i W_{i\mu} h_\mu > 0$ for all values of the connected hidden units. Typically, this limits the types of layers that can be connected to an exponential layer, and requires ensuring that all of the weights are positive.

[0080] Composite Layer: A composite layer is not a mathematical object *per se* as was the case for the previously described layer types. Instead, a composite layer is a software implementation for combining multiple sub-layers of different types to create a meta-layer that can model heterogeneous data.

[0081] Specific examples of layers for modeling data in accordance with embodiments of the invention are described above; however, one skilled in the art will recognize that any number of processes can be utilized as appropriate to the requirements of specific applications in accordance with embodiments of the invention.

Schema

[0082] A schema in accordance with several embodiments of the invention is conceptually illustrated in Figure 8. A schema with descriptions of different layers of a generalized RBM is illustrated in Figure 8. A schema allows for a model to be tuned to handle particular types of data, without requiring burdensome pre-processing by a person. The different layers allow for heterogeneous data of different types that may be incomplete and/or irregular.

[0083] Specific examples of a schema for building models in accordance with embodiments of the invention are described above; however, one skilled in the art will recognize that any number of processes can be utilized as appropriate to the requirements of specific applications in accordance

with embodiments of the invention.

Generalized Deep Boltzmann Machines (DBMs)

[0084] Deep learning refers to an approach to machine learning where the model processes the data through a series of transformations. The goal is to enable the model to learn to construct appropriate features rather than requiring the researcher to craft features using prior knowledge.

[0085] A generalized Deep Boltzmann Machine (DBM) is essentially a stack of RBMs. A generalized DBM in accordance with some embodiments of the invention is illustrated in Figure 9. The generalized DBM 900 shows a visible layer 910 connected to a hidden layer 920. Hidden layer 920 is further connected to another hidden layer 930. The visible layer 910 is encoded to hidden layer 920, which then operates like a visible layer for the next hidden layer 930.

[0086] Consider a DBM with L hidden layers \mathbf{h}_l for $l = 1, \dots, L$. The energy function of the DBM is:

$$E(\mathbf{v}, \mathbf{h}_1, \dots, \mathbf{h}_L) = -a(\mathbf{v}) - \sum_{l=1}^{l=L} b_l(\mathbf{h}_l) - \mathbf{v}^T \frac{W}{(\sigma \epsilon_1^T)^2} \mathbf{h}_1 - \sum_{l=1}^{l=L-1} \mathbf{h}_l^T \frac{W_l}{(\epsilon_l \epsilon_{l+1}^T)^2} \mathbf{h}_{l+1} \quad (10)$$

[0087] A DBM can, in principle, be trained in the same way as an RBM. However, in practice, DBMs are often trained using a greedy layer-wise process. Examples of greedy layer-wise process are described in R. Salakhutdinov and G. Hinton, in *Artificial Intelligence and Statistics* (2009) pp. 448-455, which is incorporated by reference herein. In essence, forward layerwise training of a DBM proceeds by training a sequence of RBMs with energy functions:

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}_1) &= -a(\mathbf{v}) - b_1(\mathbf{h}_1) - \mathbf{v}^T \frac{W}{(\sigma \epsilon_1^T)^2} \mathbf{h}_1 \\ E(\mathbf{h}_1, \mathbf{h}_2) &= -b_1(\mathbf{h}_1) - b_2(\mathbf{h}_2) - \mathbf{h}_1^T \frac{W_1}{(\epsilon_1 \epsilon_2^T)^2} \mathbf{h}_2 \\ &\vdots \\ E(\mathbf{h}_{L-1}, \mathbf{h}_L) &= -b_{L-1}(\mathbf{h}_{L-1}) - b_L(\mathbf{h}_L) - \mathbf{h}_{L-1}^T \frac{W_{L-1}}{(\epsilon_{L-1} \epsilon_L^T)^2} \mathbf{h}_L \end{aligned}$$

where the outputs of the previous RBM are used as the inputs of the next RBM. It can be difficult to get information from the data distribution to propagate into the deep layers of the model when training a DBM in this forward layerwise way. As a result, it is generally difficult to train DBMs with more than a couple of hidden layers.

[0088] To overcome the limitations with forward layerwise training of DBMs, methods in accordance with many embodiments of the invention train DBMs in reverse – starting with the deepest hidden layer \mathbf{h}_L and working backwards towards \mathbf{v} . This ensures that the deepest hidden layer must contain as much information about the visible layer as possible. The reverse layerwise training procedure makes use of the fact that a three layer DBM with connectivity $\mathbf{v} - \mathbf{h}_1 - \mathbf{h}_2$ is the same as a two layer RBM with connectivity $[\mathbf{v}, \mathbf{h}_2] - \mathbf{h}_1$, allowing RBMs with Composite Layers to talk backwards down the connectivity graph of the DBM.

[0089] A process for reverse layerwise training in accordance with an embodiment of the invention is conceptually illustrated in Figure 10. Process 1000 trains (1005) a first RBM with connectivity $\mathbf{v} - \mathbf{h}_L$. Process 1000 samples (1010) $\mathbf{h}_L \sim p(\mathbf{h}_L|\mathbf{v})$ from the trained RBM. The process then stacks (1015) \mathbf{v} and \mathbf{h}_L into a vector $[\mathbf{v}, \mathbf{h}_L]$ and trains (1020) a second RBM with connectivity $[\mathbf{v}, \mathbf{h}_L] - \mathbf{h}_{L-1}$. Process 1000 then determines (1025) whether $[\mathbf{v}, \mathbf{h}_2] - \mathbf{h}_1$ has been reached. When it has not been reached, process 1000 returns to step 1005. When process 1100 determines that $[\mathbf{v}, \mathbf{h}_2] - \mathbf{h}_1$ has been reached, the process copies (1030) the weights from each of these intermediate RBMs into their respective positions in the DBM. In some embodiments, DBMs can then be fine-tuned by regular end-to-end training.

Boltzmann Machines for Time Series

[0090] Many problems (e.g., modeling patient trajectories) require the ability to generate time series. That is, to generate a sequence of states $\{\mathbf{v}(t)\}_{t=0}^T$. Two approaches in accordance with numerous embodiments of the invention are described below.

[0091] An Autoregressive Boltzmann Machine (ADBMs) is a DBM where the hidden layers have undirected edges connecting neighboring time points. As a result, an ADBM relates nodes to

their previous timepoints. A generalized ADBM in accordance with some embodiments of the invention is illustrated in Figure 11. The generalized ADBM 1100 shows a visible layer 1110 at time t connected to a hidden layer 1120, also at time t . Hidden layer 1120 is further connected to another hidden layer 1130 that incorporates data that is offset from time t by τ .

[0092] As a result, an ADBM is a model for entire sequences that describes the joint probability distribution $p(\mathbf{v}(0), \dots, \mathbf{v}(\tau))$. Specifically, let $\mathbf{x}(t) = [\mathbf{v}(t), \mathbf{h}_1(t), \dots, \mathbf{h}_L(t)]$ denote the state of all of the layers at time t . Moreover, let $E_{DBM}(\mathbf{x}(t))$ be the energy of a DBM given by

$$E(\mathbf{v}, \mathbf{h}_1, \dots, \mathbf{h}_L) = -a(\mathbf{v}) - \sum_{l=1}^{l=L} b_l(\mathbf{h}_l) - \mathbf{v}^T \frac{W}{(\sigma \epsilon_l^T)^2} \mathbf{h}_1 - \sum_{l=1}^{l=L-1} \mathbf{h}_l^T \frac{W_l}{(\epsilon_l \epsilon_{l+1}^T)^2} \mathbf{h}_{l+1} \quad (11)$$

The energy function of the ADBM is:

$$E(\{\mathbf{x}(t)\}_{t=0}^{\tau}) = \sum_{t=0}^{\tau} E_{DBM}(\mathbf{x}(t)) - \sum_{l=1}^L \mathbf{h}_L^T(t) \frac{\Omega}{(\epsilon_L \epsilon_L^T)^2} \mathbf{h}_L(t-1) \quad (12)$$

For simplicity, this has been illustrated with a single autoregressive connection connecting the last hidden layer with its previous value. However, one skilled in the art will recognize that this model can be extended to include multiple time delays or inter-temporal connections between layers.

[0093] ADBMs, as described in the previous section, are able to capture correlations through time, but they are often unable to represent non-stationary distributions or distributions with drift. For example, most patients with a degenerative disease will tend to worsen over time – an effect that the ADBM cannot capture. To capture this effect, many embodiments of the invention implement a Generalized Conditional Boltzmann Machine (GCBM). Consider a time series of visible units $\{\mathbf{v}(t)\}_{t=0}^{\tau}$. The joint probability distribution can be factorized into a product $p(\mathbf{v}(t), \dots, \mathbf{v}(\tau)) = p_0(\mathbf{v}(t)) \prod_{t=1}^{\tau} p(\mathbf{v}(t) | \mathbf{v}(t-1))$. In several embodiments, this model can be constructed from two DBMs. First, a *non-time dependent* DBM, p_0 , can be trained on all of the data. Next, a *time dependent* DBM can be trained on a Composite Layer created by joining all of the neighboring time points $[\mathbf{v}(t), \mathbf{v}(t-1)]$. In this example, the second DBM describes the joint distribution $p(\mathbf{v}(t), \mathbf{v}(t-1))$, which makes it possible to compute both $p(\mathbf{v}(t) | \mathbf{v}(t-1))$ and

$p(\mathbf{v}(t-1)|\mathbf{v}(t))$ allowing for both forward and backwards prediction.

[0094] Although this example is described using a single time lag, one skilled in the art will recognize that processes in accordance with many embodiments of the invention can be adjusted to consider longer and/or multiple time lags. For example, the second DBM can be trained on a Composite Layer that can be readily extended to include multiple time lags, e.g., $[\mathbf{v}(t), \mathbf{v}(t-1), \dots, \mathbf{v}(t-n)]$.

Training RBMs

[0095] There are multiple pathways for improving the performance of RBMs. These include new approaches to regularization, novel optimization algorithms, alternative objective functions, and improved gradient estimators. Systems and methods in accordance with several embodiments of the invention implement alternative objective functions and improved gradient estimators.

Adversarial objectives for RBMs

[0096] A machine learning model is generative if it learns to draw new samples from an unknown probability distribution. Generative models can be used to learn useful representations of data and/or to enable simulations of systems with unknown, or very complicated, mechanistic laws. A generative model defined by some model parameters θ describes the probability of observing some variable \mathbf{v} . Therefore, training a generative model involves minimizing a distance between the distribution of the data, $p_d(\mathbf{v})$, and the distribution defined by the model, $p_\theta(\mathbf{v})$. The traditional method for training a Boltzmann machine maximizes the log-likelihood, which is equivalent to minimizing the forward Kullback-Liebler (KL) divergence:

$$D_{\text{KL}}(p_d \| p_\theta) = \int d\mathbf{v} p_d(\mathbf{v}) \log \left(\frac{p_d(\mathbf{v})}{p_\theta(\mathbf{v})} \right). \quad (13)$$

[0097] The forward KL divergence, $D_{\text{KL}}(p_d \| p_\theta)$, accumulates differences between the data and model distributions weighted by the probability under the data distribution. The reverse KL diver-

gence, $D_{\text{KL}}(p_{\theta} \parallel p_d)$, accumulates differences between the data and model distributions weighted by the probability under the model distribution. As a result, the forward KL divergence strongly punishes models that underestimate the probability of the data, whereas the reverse KL divergence strongly punishes models that overestimate the probability of the data.

[0098] There are a variety of sources of stochasticity that enter into the training of an RBM. The stochasticity implies that different models may become statistically indistinguishable if the differences in their log-likelihoods are smaller than the errors in estimating them. This creates an entropic force because there will be many more models with a small $D_{\text{KL}}(p_d \parallel p_{\theta})$ than there are models with both a small $D_{\text{KL}}(p_d \parallel p_{\theta})$ and $D_{\text{KL}}(p_{\theta} \parallel p_d)$. As a result, training an RBM using a standard approach with PCD decreases $D_{\text{KL}}(p_d \parallel p_{\theta})$ (as it should) but tends to increase $D_{\text{KL}}(p_{\theta} \parallel p_d)$. This leads to distributions with spurious modes and/or to distributions that are over-smoothed.

[0099] One can imagine overcoming the limitations of maximum likelihood training of RBMs by minimizing a combination of the forward and reverse KL divergences. Unfortunately, computing the reverse KL divergence requires knowledge of p_d , which is unknown. In many embodiments, rather than the reverse KL divergence, RBMs can be trained using a novel type of f-divergence as a discriminator divergence:

$$D_D(p_d \parallel p_{\theta}) := - \int d\mathbf{v} p_{\theta}(\mathbf{v}) \log \left(\frac{2p_d(\mathbf{v})}{p_d(\mathbf{v}) + p_{\theta}(\mathbf{v})} \right), \quad (14)$$

[0100] Notice that the optimal discriminator between p_d and p_{θ} will assign a posterior probability

$$p(\text{data}|\mathbf{v}) = \frac{p_d(\mathbf{v})}{p_d(\mathbf{v}) + p_{\theta}(\mathbf{v})} \quad (15)$$

that the sample \mathbf{v} was drawn from the data distribution. Therefore, the discriminator divergence can be written as

$$D_D(p_d \parallel p_{\theta}) = -\log 2 - \int d\mathbf{v} p_{\theta}(\mathbf{v}) \log(p(\text{data}|\mathbf{v})) \quad (16)$$

to show that it measures the probability that the optimal discriminator will incorrectly classify a

sample drawn from the model distribution as coming from the data distribution.

[0101] The discriminator divergence belongs to the class of f -divergences defined as $D_f(p||q) := \int dx q(x) f(p(x)/q(x))$. The function that defines the discriminator divergence is

$$f(t) = \log\left(\frac{t+1}{2t}\right) \quad (17)$$

which is convex with $f(1) = 0$, as required. It can be shown that the discriminator divergence upper bounds the reverse KL divergence:

$$\begin{aligned} \log 2 + D_D(p_d || p_\theta) &= \int d\mathbf{v} p_\theta(\mathbf{v}) \log\left(1 + \frac{p_\theta(\mathbf{v})}{p_d(\mathbf{v})}\right) \\ &\geq D_{\text{KL}}(p_\theta || p_d). \end{aligned}$$

[0102] It is often difficult to access $p_d(\mathbf{v})$ directly or to compute the reverse KL divergence. However, methods in accordance with numerous embodiments of the invention can train a discriminator to approximate Equation 15 and, therefore, can approximate the discriminator divergence.

[0103] A generator that is able to trick the discriminator so that $p(\text{data}|\mathbf{v}) \approx 1$ for all samples drawn from p_θ will have a low discriminator divergence. The discriminator divergence closely mirrors the reverse KL divergence and strongly punishes models that overestimate the probability of the data.

[0104] Methods in accordance with numerous embodiments of the invention implement a Boltzmann Encoded Adversarial Machine (BEAM) for training an RBM against an adversary. A BEAM in accordance with a number of embodiments of the invention minimizes a loss function that is a combination of the negative log-likelihood and an adversarial loss. The adversarial component ensures that BEAM training performs a simultaneous minimization of both the forward and reverse KL divergences, which prevents the oversmoothing problem observed with regular RBMs.

[0105] A method for training a BEAM in accordance with many embodiments of the invention is described below:

Input:

n = number of epochs;
 m = number of fantasy particles;
 k = number of Gibbs sampling steps;
 α = weight of the likelihood and adversarial gradients

Initialize:

sample $F \sim p_{\theta}(\mathbf{v})$ using k -steps of Gibbs sampling;

for $epoch = 1, \dots, n$ **do**

while *True* **do**

$V \leftarrow \text{minibatch}$;

if $\text{len}(V) == 0$ **then**

break;

end

 sample $F \sim p_{\theta}(\mathbf{v})$ using k -steps of Gibbs sampling;

 compute the log-likelihood gradient $g_{\mathcal{L}}(V, F, \theta)$;

 encode $\tilde{V} = \{E_{p_{\theta}(\mathbf{h}|\mathbf{v})}[\mathbf{h}]\}_{\mathbf{v} \in V}$ and $\tilde{F} = \{E_{p_{\theta}(\mathbf{h}|\mathbf{v})}[\mathbf{h}]\}_{\mathbf{v} \in F}$;

 train discriminator on \tilde{V} and \tilde{F} ;

 compute the adversarial gradient $g_V(\tilde{F}, \theta)$;

 compute the full gradient $g = \alpha g_{\mathcal{L}} + (1 - \alpha) g_V$;

 update the model parameters using the gradient;

end

end

[0106] A process for training an adversarial model in accordance with some embodiments of the invention is conceptually illustrated in Figure 12. The process 1200 draws (1205) samples from a model, such as (but not limited to) Boltzmann machines such as those described above. Samples can be drawn from a model according to a variety of methods, including (but not limited to) k -steps Gibbs sampling and TDS. The process 1200 then computes (1210) gradients based on the drawn samples. Process 1200 trains (1215) a discriminator based on the drawn samples and computes

an adversarial gradient based on the classification of the samples, as either drawn from the model or drawn from the data. In many embodiments, the process 1200 then computes (1220) a full compound gradient and updates (1225) the model parameters using the full gradient.

[0107] Figure 13 presents some comparisons between Boltzmann machines trained to maximize log likelihood and those trained as BEAMs. The examples of this figure illustrate three multimodal data distributions: a bimodal mixture of Gaussians in 1-dimension (1310), a mixture of 8 Gaussians arranged in a circle in 2-dimensions (1320), and a mixture of 25 Gaussians arranged in a grid in 2-dimensions (1330). Problems similar to the 2-dimensional mixture of Gaussians examples are commonly used for testing GANs. In each case, the regular Boltzmann machine learns a model with a pretty good likelihood by spreading the probability over the support of the data distribution. In contrast, the Boltzmann machines trained using as BEAMs learn to reproduce the data distributions very accurately.

[0108] An example of results of training a BEAM on a 2D mixture of Gaussians is illustrated in Figure 14. The first panel 1405 illustrates estimates of the forward KL divergence, $D_{KL}(p_d||p_\theta)$, and the reverse KL divergence, $D_{KL}(p_\theta||p_d)$, per training epoch. The first panel 1405 illustrates that training an RBM as a BEAM decreases both the forward and reverse KL divergences. The second panel 1410 illustrates distributions of fantasy particles at various epochs during training. In the early stages of training, the BEAM fantasy particles are spread out across the support of the data distribution capturing the modes near the edge of the grid. These early epochs resemble the distributions obtained with GANs, which also concentrate density in the modes near the edge of the grid. As training progresses, the BEAM progressively learns to capture the modes near the center of the grid.

[0109] An architecture of a Boltzmann Encoded Adversarial Machine (BEAM) in accordance with some embodiments of the invention is illustrated in Figure 15. The illustrated example shows two steps of the BEAM architecture. In the first stage 1510, a generator (e.g., an RBM) with a visible layer (circles) and a hidden layer (diamonds). Generators in accordance with a number of embodiments of the invention are trained to encode input data by passing the input data through

the visible layer to be encoded in a set of nodes of a hidden layer. Generators in accordance with several embodiments of the invention are trained with an objective to generate realistic samples from a complex distribution. In many embodiments, objective functions for training generators can include a contribution from an adversarial loss generated by a critic (or discriminator).

[0110] In the second stage 1520, the hidden layer of the generator feeds into a discriminator (or critic) that evaluates the hidden layers to distinguish samples drawn from the data from samples drawn from the model using tied weights learned by the generator. The discriminator (or adversary) is constructed by encoding the visible units using a single forward pass through the layers of the generator and then applying a classifier (e.g., logistic regression, nearest neighbor classifiers, and random forest) trained to discriminate between samples from the data and samples from the model. By refining the discriminator, processes in accordance with many embodiments of the invention allow for an improved model of complex probability distributions. Although shown in separate stages, the BEAM in accordance with many embodiments of the invention is trained with a compound objective that trains both the critic and the generator simultaneously. In certain embodiments, the discriminator is a simple classifier that requires very little training.

[0111] The objective function in accordance with a number of embodiments of the invention is

$$\mathcal{C} = -\gamma\mathcal{L} - (1 - \gamma)\mathcal{A}, \quad (18)$$

which includes a contribution from adversarial term, \mathcal{A} , from a critic. Adversarial terms in accordance with a number of embodiments of the invention can be defined as

$$\mathcal{A} := \int d\mathbf{v}d\mathbf{h} p_{\theta}(\mathbf{v}, \mathbf{h}) T(\mathbf{v}, \mathbf{h}). \quad (19)$$

where $T(\mathbf{v}, \mathbf{h})$ is a critic function. In some embodiments, the adversary uses the same architecture and weights as the RBM, and encodes visible units into hidden unit activations. These hidden unit activations, computed for both the data and fantasy particles sampled from the RBM, are used by a critic to estimate the distance between the data and model distributions.

[0112] To compute the derivatives for training the generator, methods in accordance with some embodiments of the invention use the stochastic derivative trick:

$$\begin{aligned}
\partial_{\theta} \mathcal{A} &= \frac{\partial}{\partial \theta} \int d\mathbf{v} d\mathbf{h} p(\mathbf{v}, \mathbf{h}) T(\mathbf{v}, \mathbf{h}) \\
&= \int d\mathbf{v} d\mathbf{h} T(\mathbf{v}, \mathbf{h}) \frac{\partial}{\partial \theta} p(\mathbf{v}, \mathbf{h}) \\
&= \int d\mathbf{v} d\mathbf{h} T(\mathbf{v}, \mathbf{h}) \frac{p(\mathbf{v}, \mathbf{h})}{p(\mathbf{v}, \mathbf{h})} \frac{\partial}{\partial \theta} p(\mathbf{v}, \mathbf{h}) \\
&= \int d\mathbf{v} d\mathbf{h} T(\mathbf{v}, \mathbf{h}) p(\mathbf{v}, \mathbf{h}) \partial_{\theta} \log p(\mathbf{v}, \mathbf{h}) \\
&= \langle T(\mathbf{v}, \mathbf{h}) \partial_{\theta} \log p(\mathbf{v}, \mathbf{h}) \rangle_{p(\mathbf{v}, \mathbf{h})} \\
&= -\langle T(\mathbf{v}, \mathbf{h}) \rangle_{p_{\theta}(\mathbf{v}, \mathbf{h})} \langle -\partial_{\theta} E_{\theta}(\mathbf{v}, \mathbf{h}) \rangle_{p_{\theta}(\mathbf{v}, \mathbf{h})} + \langle T(\mathbf{v}, \mathbf{h}) (-\partial_{\theta} E_{\theta}(\mathbf{v}, \mathbf{h})) \rangle_{p_{\theta}(\mathbf{v}, \mathbf{h})} \\
&= \text{Cov}_{p_{\theta}(\mathbf{v}, \mathbf{h})} [T(\mathbf{v}, \mathbf{h}), -\partial_{\theta} E_{\theta}(\mathbf{v}, \mathbf{h})]. \tag{20}
\end{aligned}$$

where $\partial_{\theta} \log p_{\theta}(\mathbf{v}, \mathbf{h}) = -\langle -\partial_{\theta} E_{\theta}(\mathbf{v}, \mathbf{h}) \rangle_{p_{\theta}(\mathbf{v}, \mathbf{h})} - \partial_{\theta} E_{\theta}(\mathbf{v}, \mathbf{h})$ is used for an RBM.

[0113] In principle, the critic can be any function of the visible and hidden units. However, based on the discriminator divergence, methods in accordance with several embodiments of the invention use a critic that is monotonically related to $p(\text{data}|\mathbf{v})$. Although the discriminator divergence suggests that one could use $\log p(\text{data}|\mathbf{v})$, methods in accordance with certain embodiments of the invention use a linear function $T(\mathbf{v}) = 2 * p(\text{data}|\mathbf{v}) - 1$. Typically, the optimal discriminator can be approximated as a function of the hidden units activations $p(\text{data}|\mathbf{v}) \approx g(\langle \mathbf{h} \rangle_{p_{\theta}(\mathbf{h}|\mathbf{v})})$. The function $g(\cdot)$ could be implemented by a neural network, as in most GANs, or using a simpler algorithm such as a random forest or nearest neighbor classifier. In a number of embodiments, a simple approximation to the optimal discriminator can be sufficient because the classifier can operate on the hidden unit activities of the RBM generator rather than the visible units. Therefore, the optimal critic can be approximated using nearest neighbor methods.

[0114] Suppose $X = \{x_1, \dots, x_N\}$ are identically and independently distributed samples from an unknown probability distribution with pdf $p(\mathbf{x})$ in \mathbb{R}^n . In a variety of embodiments, $p(\mathbf{x})$ is estimated at an arbitrary point \mathbf{x} based on a k -nearest-neighbor estimate. Specifically, methods in

accordance with some embodiments of the invention fix some positive integer k and compute the k nearest neighbors to \mathbf{x} in X . Then, d_k is defined to be the distance between x and the furthest of the nearest-neighbors and the density $p(\mathbf{x})$ is estimated to be the density of the uniform distribution on a ball of radius d_k . That is,

$$p(\mathbf{x}) \approx k \left(\frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2} + 1)} d_k^n \right)^{-1}. \quad (21)$$

[0115] Now denote by $p_\theta(\mathbf{v})$ and $p_d(\mathbf{v})$ the unknown pdfs of the model and data distributions, respectively, and define the distance between two vectors \mathbf{v} and \mathbf{v}' as the Euclidean distance between their hidden unit activations, $d(\mathbf{v}, \mathbf{v}') = \|\langle h \rangle_{p_\theta(\mathbf{h}|\mathbf{v})} - \langle h \rangle_{p_\theta(\mathbf{h}|\mathbf{v}')}\|$. This distance may no longer satisfy all of the properties of a proper metric. Let $X = \{\mathbf{v}_1, \dots, \mathbf{v}_{2N}\}$ be a collection of samples where exactly half are drawn from p_θ and half from p_d . Fix some k and compute the k nearest neighbors in X , denoting by d_k the distance to the furthest. Then the denominator is estimated as described above. Let j be the number of nearest neighbors which come from p_d as opposed to p_θ . The numerator then can be estimated as uniform on the same size ball with only j/k of the density of the denominator, allowing the nearest-neighbor critic to be defined $T_{NN}(\mathbf{v}) := j/k$. In many embodiments, the nearest neighbors can be computed from a cached minibatch of samples from the model combined with a minibatch of samples from the training dataset.

[0116] The distance-weighted nearest-neighbor critic is a generalization which adds some continuity to the nearest-neighbor critic by applying an inverse distance weighting to the ratio count. Specifically, let $\{d_0, \dots, d_k\}$ be the distances of the k -nearest neighbors, with $\{d_0, \dots, d_j\}$ the distances for the neighbors originating from the data samples and $\{d_{j+1}, \dots, d_k\}$ the distances for the neighbors originating from the model samples. In many embodiments, the distance-weighted nearest-neighbor critic can be defined as:

$$T_{DNN}(\mathbf{v}) := \frac{\sum_{i=1}^j \frac{1}{d_i + \epsilon}}{\sum_{i=1}^k \frac{1}{d_i + \epsilon}}, \quad (22)$$

where ϵ is a small parameter that regularizes the inverse distance.

[0117] In the context of most formulations of GANs, which use feed-forward neural networks

for both the generator and the discriminator, one could say that BEAMs use the RBM as both the generator and as a feature extractor for the adversary. In various embodiments, this double-usage allows the reuse of a single set of fantasy particles for multiple steps of the training algorithm. Specifically, a single set of M persistent fantasy particles are updated k times per gradient evaluation. In many embodiments, the same set of fantasy particles are used to compute the log-likelihood derivative and the adversarial derivative. Then, these fantasy particles can replace the fantasy particles from the previous gradient evaluation in the nearest neighbor estimates of the critic value. Reusing the fantasy particles for each step means that BEAM training has roughly the same computational cost as training an RBM with PCD.

Improved gradient estimates

[0118] The gradients of the log-likelihood and the adversarial term both involve expectation values with respect to the model distribution. Unfortunately, these expectation values cannot be computed exactly. As a result, the expectation values can be approximated using Monte Carlo methods or other approximations. The accuracy of these approximate gradients can have a significant effect on the utility of the resulting model. Different approaches to improving the accuracy of the approximate gradients in accordance with certain embodiments of the invention are described below.

Mean-field approximations and shrinkage estimates

[0119] Monte Carlo estimates of the gradients have the advantage of being unbiased. That is, $\frac{1}{N} \sum_k f(\mathbf{v}_k, \mathbf{h}_k) \rightarrow \langle f(\mathbf{v}, \mathbf{h}) \rangle_{p_{\theta}(\mathbf{v}, \mathbf{h})}$ as $N \rightarrow \infty$. However, the estimates may have a high variance when N is small. On the other hand, mean field estimates such as those derived from the Thouless-Andersen-Palmer (TAP) expansion are analytic and have zero variance, but have a bias that can be difficult to control. Let $f(\omega) = \omega f_{MC} + (1 - \omega) f_{MF}$ be an estimate created from a convex combination of a Monte Carlo estimate f_{MC} and a mean field estimate f_{MF} . It is easy to show that $\text{Bias}^2[f] = (1 - \omega)^2 \text{Bias}^2[f_{MF}]$ and $\text{Var}[f] = \omega^2 \text{Var}[f_{MC}]$ so that the mean squared error of f

is $\text{MSE}[f] = \text{Bias}^2[f] + \text{Var}[f] = (1 - \omega)^2 \text{Bias}^2[f_{\text{MF}}] + \omega^2 \text{Var}[f_{\text{MC}}]$. Therefore, one can generally choose a value of ω to minimize the mean squared error of the combined estimator.

Tempered sampling

[0120] Drawing samples from a probability distribution is an important component of many processes for training models in accordance with many embodiments of the invention. This can often be done with a simple function call for many 1-dimensional distributions. However, random sampling from Boltzmann machines is much more complicated.

[0121] Sampling from a Boltzmann machine is usually performed using Gibbs sampling. Gibbs sampling is a local sampling process, which means that successive samples are correlated. Drawing uncorrelated samples requires one to make many Gibbs sampling steps for each successive sample. As a result, drawing a batch of uncorrelated random samples from a Boltzmann machine can take a long time. A batch of random samples is required for each gradient update – if it takes a long time to generate each batch, it can make training a Boltzmann machine take such a long time that it becomes impractical. Therefore, methods that decrease the correlation between successive samples from a Boltzmann machine can greatly accelerate the learning process.

[0122] Many methods for accelerated sampling from Boltzmann machines rely on an analogy with temperature from statistical physics. To do this, methods in accordance with a number of embodiments of the invention introduce a fictional inverse temperature β into a Boltzmann machine by defining the probability distribution as:

$$p_{\beta}(\mathbf{v}, \mathbf{h}) = Z_{\beta}^{-1} e^{-\beta E(\mathbf{v}, \mathbf{h})} \quad (23)$$

The original distribution of the Boltzmann machine is recovered by setting $\beta = 1$.

[0123] The fictional temperature is useful because raising the temperature (i.e., decreasing β) decreases the autocorrelation between samples. Consider a situation with starting configuration (\mathbf{v}, \mathbf{h}) and ending at configuration $(\mathbf{v}', \mathbf{h}')$. The initial energy is $E(\mathbf{v}, \mathbf{h})$. As one moves from the

initial to the final configuration, the intermediate configurations will have varying energies. If the maximal energy from these intermediate configurations is E_{max} then the time to travel from (\mathbf{v}, \mathbf{h}) to $(\mathbf{v}', \mathbf{h}')$ roughly scales as:

$$\tau \sim e^{\beta(E_{max} - E(\mathbf{v}, \mathbf{h}))} \quad (24)$$

Therefore, decreasing β will decrease the number of Gibbs sampling steps required to move between distant configurations.

[0124] Although raising the temperature will decrease the mixing time, it also changes the resulting probability distribution. Therefore, simply sampling from a model with a $\beta \ll 1$ during training will not allow a model to learn correctly. Processes in accordance with certain embodiments of the invention use a process called parallel tempering (in the machine learning and statistics literature) or replica exchange (in the physics community). In parallel tempering in accordance with a variety of embodiments of the invention, multiple Gibbs sampling chains are run in parallel, each at a different temperature. Periodically, one attempts to swap the configurations of two chains. In several embodiments, the swap can be accepted or rejected based on a criterion (e.g., the Metropolis criterion) to ensure that entire system stays at equilibrium. After a long time, a configuration that started out at $\beta = 1$ will travel to a chain with a lower temperature (where it can cross energy barriers more easily) and back to the chain running at $\beta = 1$. This ensures that the chain running at $\beta = 1$ has a faster mixing time while still sampling from the correct probability distribution. There is a computational cost, however, because many Gibbs sampling chains have to be run in parallel.

[0125] In some embodiments of the invention, the process uses Temperature Driven Sampling (TDS), which greatly improves the ability to train Boltzmann machines without incurring significant additional computational cost. TDS is a variant of a sequential Monte Carlo sampler. A collection of m samples are evolved independently using Gibbs sampling updates from the model. Note that this is not the same as running multiple chains for a parallel tempering process because each of the m samples in the sequential Monte Carlo sampler will be used compute statistics, as opposed to just the samples from the $\beta = 1$ chain during parallel tempering. Each of these samples has an inverse temperature that is drawn from a distribution with mean $\langle \beta \rangle = 1$ and a variance

$\text{Var}[\beta] < 1$. In several embodiments, the inverse temperatures of each sample can be independently updated once for every Gibbs sampling iteration of the model. In a variety of embodiments, the updates are autocorrelated across time so that the inverse temperatures are slowly varying. As a result, the collection of samples are drawn from a distribution that is close to the model distribution, but with fatter tails. This allows for much faster mixing, while ensuring that the model averages (computed over the collection of m samples) remain close approximations to averages computed from the model with $\beta = 1$. An example of sampling from an autocorrelated Gamma distribution is described below.

Input:

Autocorrelation coefficient $0 \leq \phi < 1$.

Variance of the distribution $\text{Var}[\beta] < 1$.

Current value of β .

Set: $v = 1/\text{Var}[\beta]$ and $c = (1 - \phi)\text{Var}[\beta]$.

Draw $z \sim \text{Poisson}(\beta * \phi / c)$.

Draw $\beta' \sim \text{Gamma}(v + z, c)$.

return β'

[0126] TDS includes a standard Gibbs sampling based sequential Monte Carlo sampler in the limit that $\text{Var}[\beta] \rightarrow 0$. The samples drawn with TDS are *not* samples from the equilibrium distribution of the Boltzmann machine. In certain embodiments, the drawn samples are re-weighted to correct for the bias due to the varying temperature.

Input:

Number of samples m .

Number of update steps k .

Autocorrelation coefficient for the inverse temperature $0 \leq \phi < 1$.

Variance of the inverse temperature $\text{Var}[\beta] < 1$.

Initialize:

Randomly initialize m samples $\{(\mathbf{v}_i, \mathbf{h}_i)\}_{i=1}^m$.

Randomly initialize m inverse temperatures $\beta_i \sim \text{Gamma}(1/\text{Var}[\beta], \text{Var}[\beta])$.

for $t = 1, \dots, k$ **do**

for $i = 1, \dots, m$ **do**

 Update β_i using a driven gamma sampler.

 Update $(\mathbf{v}_i, \mathbf{h}_i)$ using Gibbs sampling.

end

end

[0127] Temperature Driven Sampling (TDS) improves sampling from a Boltzmann machine. A direct comparison between samples drawn from a Boltzmann machine with regular Gibbs sampling to those drawn using TDS is illustrated in Figure 16. GMM (gray) refers to samples from a Gaussian mixture model. GRBM (blue) refers to samples from the equivalent Boltzmann machine drawn using 10 steps of Gibbs sampling. TDS (red) refers to samples from the equivalent Boltzmann machine drawn using TDS with 10 steps of Gibbs sampling. This example shows a Gaussian mixture model with three modes at $(-1, 0, +1)$ with various standard deviations and using a simple construction to create an equivalent Boltzmann machine with a Gaussian visible layer and a One-hot hidden layer with 3 hidden units. The autocorrelation coefficient and the standard deviation of the inverse temperature were set to 0.9 and 0.95, respectively. All starting samples were initialized from the middle mode. Starting from the middle mode, regular Gibbs sampling is unable to sample from the neighboring modes after 10 steps when the modes are well separated. TDS, by contrast, has fatter tails allowing for better sampling of the neighboring modes.

[0128] Using TDS at train time can have a pretty dramatic effect on the resulting model. In Fig-

ure 17, two identical Gaussian-Bernoulli RBMs were trained on grayscale images of handwritten digits from the MNIST dataset. Images are from models with identical architectures trained with identical hyperparameters, except that one used regular Gibbs sampling (1710) whereas the other used TDS (1720), or (a) is trained with $\text{Var}[\beta] = 0$ and (b) is trained with $\text{Var}[\beta] = 0.9$. Both models are Gaussian-Bernoulli RBMs with 256 hidden units, trained for 100 epochs of persistent contrastive divergence using the ADAM optimizer with a learning rate of 0.0005 and batch size of 100. Temperature Driven Sampling (TDS) improves learning for a model of the MNIST handwritten digits (grayscale). Both models achieve a low reconstruction error (data not shown), but the GRBM trained with the regular Gibbs sampler fails to generate realistic fantasy particles. The GRBM trained with TDS, by contrast, generates fantasy particles that look like realistic handwritten digits.

[0129] Specific processes for drawing samples from a probability distribution in accordance with embodiments of the invention are described above; however, one skilled in the art will recognize that any number of processes can be utilized as appropriate to the requirements of specific applications in accordance with embodiments of the invention.

Applications

[0130] That is, even though it may only be possible to predict the probability of a health outcome for an individual patient, this ability makes it possible to precisely predict the number of patients with that health outcome in a large population. For example, predicting health risks makes it possible to accurately estimate the cost of insuring a population. Similarly, predicting the likelihood that a patient will respond to a particular therapeutic makes it possible to estimate the probability of a positive outcome in a clinical trial.

Simulating Patient Trajectories

[0131] Developing the ability to accurately predict patients' prognoses is a necessary step towards precision medicine. A patient can be represented as a collection of information that de-

scribes their symptoms, their genetic information, results from diagnostic tests, any medical treatments they are receiving, and other information that may be relevant for characterizing their health. A vector containing this information about a patient is sometimes called a phenotype vector. A method for prognostic prediction in accordance with many embodiments of the invention uses past and current health information about a patient to predict a health outcome at a future time.

[0132] A patient trajectory refers to a time series that describes a patient's detailed health status (e.g., a patient's phenotype vector) at various points in time. In several embodiments, prognostic prediction takes in a patient's trajectory (i.e., their past and current health information) and makes a prediction about a specific future health outcome (e.g., the likelihood they will have a heart attack within the next 2 years). By contrast, predicting a patient's future trajectory involves predicting all of the information that characterizes the state of their health at all future times.

[0133] To frame this mathematically, let $\mathbf{v}(t)$ be a phenotype vector containing all of the information characterizing the health of a patient at time t . Therefore, a patient trajectory is a set $\{\mathbf{v}(t)\}_{t=0}^T$. Many of the examples are described with discrete time steps (e.g., one month), but one skilled in the art will recognize that this is not necessary and that various other time steps can be employed in accordance with various embodiments of the invention. In some embodiments of the invention, models for simulating patient trajectories use discrete time steps (e.g., one month). The length of the time step in accordance with a number of embodiments of the invention will be selected to approximately match the frequency of treatment. A model for patient trajectories in accordance with many embodiments of the invention describes the joint probability distribution of all points along the trajectory, $p(\mathbf{v}_0, \dots, \mathbf{v}_T)$. Such a model can be used for prediction by sampling from the conditional probability distribution $p(\mathbf{v}_\tau, \dots, \mathbf{v}_T | \mathbf{v}_0, \dots, \mathbf{v}_{\tau-1})$. In many embodiments, the model is a Boltzmann machine, as they make it easy to express conditional distributions and can be adapted to heterogeneous datasets, but one skilled in the art will recognize that many of the processes described herein can be applied to other architectures as well.

Clinical Decision Support Systems

[0134] Clinical decision support systems provide information to patients, physicians, or other caregivers to help guide choices about patient care. Simulated patient trajectories provide insights into a patient's future health that can inform choices of care. For example, consider a patient with mild cognitive impairment. A physician or caregiver would benefit from knowing the risks that the patient's condition progresses to Alzheimer's disease, or that he or she begins to exhibit other cognitive or psychological systems. In certain embodiments, systems based on simulated patient trajectories can forecast these risks to guide care choices. Aggregating such predictions over a population of patients can also help estimate population level risks, enabling long-term planning by organizations, such as elder care facilities, that act as caregivers to large groups of patients.

[0135] In some embodiments, a set of patient trajectories is collected from electronic medical records (also known as real world data), from natural history databases, or clinical trials. The patient trajectories in accordance with many embodiments of the invention can be normalized and used to train a time-dependent Boltzmann machine. To use the model, the medical history for a patient can be input in the form of a trajectory $\{\mathbf{v}(t)\}_{t=0}^{t_0}$ where t_0 is the current time and use the Boltzmann machine to simulate trajectories from the probability distribution $p(\mathbf{v}_{t_0+1}, \dots, \mathbf{v}_T | \mathbf{v}_0, \dots, \mathbf{v}_{t_0})$. Then, these simulated trajectories can be analyzed to understand the risks associated with specific outcomes (e.g., Alzheimer's diagnosis) at various future times. In some cases, models that are trained on data with treatment information would contain variables that describe treatment choices. Such a model could be used to assess how different treatment choices would change the patient's future risks by comparing simulated outcome risks conditioned on different treatments. In many embodiments, a caretaker or physician can treat a patient based on the treatment choices and/or the simulated trajectories.

Simulating Control Arms for Clinical Trials

[0136] Randomized Clinical Trials (RCTs) are the gold-standard for evidence in assessing therapeutic efficacy. In an RCT, each patient is randomly assigned to one of two study arms: a treatment arm where the patients are treated with an experimental therapy, and a placebo arm where the pa-

tients receive a dummy treatment and/or the current standard of care. At the end of the trial, a statistical analysis is performed to determine if patients in the treatment arm were more likely to respond positively to the new therapy than patients in the placebo arm were to respond to the dummy therapy.

[0137] In order to have enough statistical power to accurately assess the efficacy of the experimental therapy, RCTs need to include a large number of patients. For example, it is not uncommon for Phase III clinical trials to include thousands of patients. Recruiting the large number of patients necessary to achieve sufficient power is challenging, and many clinical trials never meet their recruitment goals. Although there is, almost by definition, little-to-no data about an experimental therapy there is likely a lot of data about the efficacy of the current standard of care. Therefore, one way to reduce the number of patients needed for clinical trials is to replace the control arm with a synthetic control arm that contains virtual patients simulated from a Boltzmann machine trained to model the current standard of care.

[0138] Methods in accordance with several embodiments of the invention use simulations to create a synthetic, or virtual, control arm for a clinical trial by training a Boltzmann machine using data from the control arms of previous clinical trials. In many embodiments, data sets can be constructed by aggregating data from the control arms of multiple clinical trials for a chosen disease. Then, Boltzmann machines can be trained to simulate patients with that disease under the current standard of care. This model can then be used to simulate a population of patients with particular characteristics (e.g., age, ethnicity, medical history) to create a cohort of simulated patients that match the inclusion criteria of new trial. In some embodiments, each patient in the experimental arm can be matched to a simulated patient with the same baseline measurements by simulating from the appropriate conditional distribution of the Boltzmann machine. This can provide a type of counterfactual (i.e., what would have happened to this patient if they had been given a placebo rather than the experimental therapy). In either case, data from simulated patients can be used to supplement, or in place of, data from a concurrent placebo arm using standard statistical methods in accordance with many embodiments of the invention.

Simulating Head-to-Head Clinical Trials

[0139] Traditionally, health care in the United States has been provided on a fee-for-service basis. However, there is an ongoing shift towards value based care. In the context of pharmaceuticals, value based care means that the cost of a drug will be based on how effective it is, rather than a simple cost per pill. As a result, governments and other payers need to be able to compare the effectiveness of alternative therapies.

[0140] Consider two drugs A and B with the same indication. There are two standard ways to compare the efficacy of A and B. First, one can use electronic health records and insurance claims data to observe how well the drugs are working in the context of real world clinical practice. Alternatively, one can run an RCT to perform a head-to-head comparison of the drugs. Both of these methods take years of additional observation and/or experimentation to arrive at a conclusion about the comparative effectiveness of A and B.

[0141] Simulations in accordance with many embodiments of the invention provide an alternative approach for performing head-to-head trials. In some embodiments, detailed individual level data from clinical trials of each drug can be included in the training data for a Boltzmann machine. In some embodiments, samples generated with a Boltzmann machine, such as a BEAM, can be used to simulate a head-to-head clinical trial between A and B. However, individual level data are not usually released for the experimental arms of clinical trials. In the absence of these data, aggregate level data from the experimental arms in accordance with a number of embodiments of the invention can be used to adjust a model that was trained on control arm data.

Learning Unsupervised Genomic Features

[0142] The human genome encodes for more than 20 thousands genes that engage in an incredibly complex network of interactions. This network of genetic interactions is so complex that it is intractable to develop a mechanistic model linking genotype to phenotype. Therefore, studies that aim to predict a phenotype from genomic information have to use machine learning methods.

[0143] A common goal of a genomic study in the clinical setting is predicting whether or not a patient will respond to a given therapeutic. For example, data describing gene expression (e.g., from messenger RNA sequencing experiments) may be collected at the beginning of a phase-II clinical trial. The response of each patient to the therapeutic is recorded at the end of the trial, and a mathematical model (e.g., linear or logistic regression) is trained to predict the response of each patient from their baseline gene expression data. Successful prediction of patient response would enable the sponsor of the clinical trial to use a genomic test to narrow the study population to a subset of patients where the drug is most likely to be successful. This improves the likelihood of success in a subsequent phase-III trial, while also improving patient outcomes through precision medicine.

[0144] Unfortunately, phase-II clinical trials tend to be small (200 people). Moreover, sequencing experiments used to measure gene expression are still fairly expensive. As a result, even non-clinical gene expression studies are limited in size. Therefore, the standard task involves training a regression model with up to 20 thousand features (i.e., the expression of the genes) using less than 200 measurements. In general, a linear regression model is underdetermined if the number of features is greater than the number of measurements. Although there are techniques to mitigate this problem, the situation in most 'omics studies is so lopsided that standard approaches fail.

[0145] In many embodiments, raw gene expression values are combined into a smaller number of composite features. For example, individual genes interact as parts of biochemical pathways, so one approach is to use known biochemical information to derive scores that describe the activation of pathways. Then, pathway activation scores can be used as features instead of raw expression values. However, due to the complexity of biochemical networks, it can be unclear how to construct pathway activation scores in the first place.

[0146] In certain embodiments, Deep Boltzmann Machines (DBMs) are implemented as a tool for unsupervised feature learning that may be useful for 'omics studies. Let \mathbf{v} be a vector containing gene expression values determined from an experiment. A DBM describes the distribution of gene expression vectors using a probability distribution $p(\mathbf{v}) = \int d\mathbf{h}_1 \cdots d\mathbf{h}_L p(\mathbf{v}, \mathbf{h}_1, \dots, \mathbf{h}_L)$ where the

layers of hidden units \mathbf{h}_l describe progressive transformations of the gene expression values into higher level features. The model in accordance with many embodiments of the invention can be trained without labels; therefore, in some embodiments, a large data set can be compiled by combining many different studies. In a number of embodiments, the pre-trained DBM can be used to transform a vector of raw gene expression values into a lower dimensional vector of features by computing $\langle \mathbf{h}_L \rangle_{\mathbf{v}} = \int d\mathbf{h}_1 \cdots d\mathbf{h}_L \mathbf{h}_L p(\mathbf{h}_1, \dots, \mathbf{h}_L | \mathbf{v})$. These lower dimensional features in accordance with certain embodiments of the invention can then be used as input to a simpler supervised learning algorithm to construct a predictor of drug response for a given therapeutic.

Predicting Transcriptomic Responses

[0147] Predicting the effect that a change in the activity, or expression, of a gene will have in-human is important for both drug design and drug development. For example, if one could predict the effect that a compound will have in-human then one could perform high-throughput computational screens for drug discovery. Similarly, if one could predict the effect that an investigational drug will have on different types of patients then one could optimize patient selection for phase II clinical trials even though there is no direct data on the action of the drug in-human.

[0148] There isn't an obvious way to use supervised learning methods to develop a predictor of transcriptomic response. In many embodiments, transcriptomic responses are predicted using a *generative model* of gene expression. Let \mathbf{v} be a vector of raw gene expression values and let $p_{\theta}(\mathbf{v})$ be a model of the distribution of gene expression values that is parameterized by θ . Moreover, suppose that the model is parameterized such that θ_i is related to the mean value of v_i , such that increasing (or decreasing) θ_i leads to an increase (or decrease) in $\langle v_i \rangle$. In many embodiments, the effect of a drug that decreases the activity of gene i is simulated by decreasing θ_i and computing the change in $\langle \mathbf{v} \rangle$. In a number of embodiments, when the change is small, then this involves computing the derivative $\partial_{\theta_i} \langle \mathbf{v} \rangle = \partial_{\theta_i} \int d\mathbf{v} \mathbf{v} p_{\theta}(\mathbf{v})$.

[0149] The utility of generative models in accordance with several embodiments of the invention relies on the ability of the model to *implicitly* learn interactions between gene expression values.

That is, the model must know that decreasing the activity of gene i using a therapeutic will – via a complex network of interactions – lead to a decrease in the expression of some other gene j . In numerous embodiments, DBMs as described in previous sections of this application are used as a generative model that implicitly (i.e., without trying to construct a mechanistic understanding of biochemical pathways or other methods of direct gene interaction) learns interaction between genes.

[0150] In many embodiments, DBMs trained on gene expression data in a fully unsupervised manner do not have a notion of an *individual* patient. Instead, the vector of observations \mathbf{v} can be broken into two pieces: the vector of gene expression values \mathbf{x} and a vector of metadata \mathbf{y} . The metadata in accordance with some embodiments of the invention may describe characteristics of the sample such as (but not limited to) which tissue it came from, the health status of the patient, or other information. Then, in a number of embodiments, predictions can be made from the conditional distributions $\partial_{\theta_i} \langle \mathbf{x} \rangle_{\mathbf{y}} = \partial_{\theta_i} \int d\mathbf{x} \mathbf{x} p_{\theta}(\mathbf{x}|\mathbf{y})$.

[0151] Finally, predictions for individual patients in accordance with several embodiments of the invention can use a notion of locality in gene expression space. Let $\mathcal{F}_{\theta}(\mathbf{x}|\mathbf{y}) := -\log p_{\theta}(\mathbf{x}|\mathbf{y})$ define the energy \mathbf{x} given \mathbf{y} . In a DBM, this also involves integrating over all of the hidden layers. In certain embodiments, local measures of gene interactions can be computed from the derivatives of \mathcal{F} evaluated at \mathbf{x} .

[0152] Although the present invention has been described in certain specific aspects, many additional modifications and variations would be apparent to those skilled in the art. It is therefore to be understood that the present invention may be practiced otherwise than specifically described. Thus, embodiments of the present invention should be considered in all respects as illustrative and not restrictive.

What is claimed is:

- 1.** A method for training a restricted Boltzmann machine (RBM), wherein the method comprises:

 - generating, from a first set of visible values, a set of hidden values in a hidden layer of a RBM;
 - generating a second set of visible values in a visible layer of the RBM based on the generated set of hidden values;
 - computing a set of likelihood gradients based on at least one of the first set of visible values and the generated set of visible values;
 - computing a set of adversarial gradients using an adversarial model based on at least one of the set of hidden values and the set of visible values;
 - computing a set of compound gradients based on the set of likelihood gradients and the set of adversarial gradients; and
 - updating the RBM based on the set of compound gradients.
- 2.** The method of claim **1**, wherein the visible layer of the RBM comprises a composite layer composed of a plurality of sub-layers for different data types.
- 3.** The method of claim **1**, wherein the plurality of sub-layers comprises at least one of a Bernoulli layer, an Ising layer, a one-hot layer, a von Mises-Fisher layer, a Gaussian layer, a ReLU layer, a clipped ReLU layer, a student-t layer, an ordinal layer, an exponential layer, and a composite layer.
- 4.** The method of claim **1**, wherein the RBM is a deep Boltzmann machine (DBM), wherein the hidden layer is one of a plurality of hidden layers.
- 5.** The method of claim **4**, wherein the RBM is a first RBM and the hidden layer is a first hidden layer of the plurality of hidden layers, wherein the method further comprises:

 - sampling the hidden layer from the first RBM;

stacking the visible layer and the hidden layer from the first RBM into a vector;
training a second RBM, wherein the vector is a visible layer of the second RBM; and
generating the DBM by copying weights from the first and second RBMs to the DBM.

6. The method of claim 1 further comprising:
 - receiving a phenotype vector for a patient;
 - using the RBM to generate a time progression of a disease; and
 - treating the patient based on the generated time progression.
7. The method of claim 1, wherein the visible layer and the hidden layer are for a first time instance, wherein the hidden layer is further connected to a second hidden layer that incorporates data from a different second time instance.
8. The method of claim 1, wherein the visible layer is a composite layer comprising data for a plurality of different time instances.
9. The method of claim 1, wherein computing the set of likelihood gradients comprises performing Gibbs sampling.
10. The method of claim 1, wherein the set of compound gradients are weighted averages of the set of likelihood gradients and the set of adversarial gradients.
11. The method of claim 1 further comprising training the adversarial model by:
 - drawing data samples based on authentic data;
 - drawing fantasy samples based from the RBM; and
 - training the adversarial model based on the adversarial model's ability to distinguish between the data samples and the fantasy samples.
12. The method of claim 1, wherein training the adversarial model comprises measuring a probability that a particular sample is drawn from either the authentic data or the RBM.

13. The method of claim 1, wherein the adversarial model is one of a fully-connected classifier, a logistic regression model, a nearest neighbor classifier, and a random forest.

14. The method of claim 1 further comprising using the RBM to generate a set of samples of a target population.

15. The method of claim 1, wherein computing a set of likelihood gradients comprises computing a convex combination of a Monte Carlo estimate and a mean field estimate.

16. The method of claim 1, wherein computing a set of likelihood gradients comprises:
initializing a plurality of samples;
initializing an inverse temperature for each sample of the plurality of samples;
for each sample of the plurality of samples:
 updating the inverse temperature by sampling from an autocorrelated Gamma distribution; and
 updating the sample using Gibbs sampling.

17. A non-transitory machine readable medium containing processor instructions for training a restricted Boltzmann machine (RBM), wherein execution of the instructions by a processor causes the processor to perform a process that comprises:

 generating, from a first set of visible values, a set of hidden values in a hidden layer of a RBM;

 generating a second set of visible values in a visible layer of the RBM based on the generated set of hidden values;

 computing a set of likelihood gradients based on at least one of the first set of visible values and the generated set of visible values;

 computing a set of adversarial gradients using an adversarial model based on at least one of the set of hidden values and the set of visible values;

computing a set of compound gradients based on the set of likelihood gradients and the set of adversarial gradients; and

updating the RBM based on the set of compound gradients.

18. The non-transitory machine readable medium of claim **17**, wherein the visible layer of the RBM comprises a composite layer composed of a plurality of sub-layers for different data types.

19. The non-transitory machine readable medium of claim **17**, wherein the RBM is a deep Boltzmann machine (DBM), wherein the hidden layer is one of a plurality of hidden layers.

20. The non-transitory machine readable medium of claim **19**, wherein the RBM is a first RBM and the hidden layer is a first hidden layer of the plurality of hidden layers, wherein the process further comprises:

sampling the hidden layer from the first RBM;

stacking the visible layer and the hidden layer from the first RBM into a vector;

training a second RBM, wherein the vector is a visible layer of the second RBM; and

generating the DBM by copying weights from the first and second RBMs to the DBM.

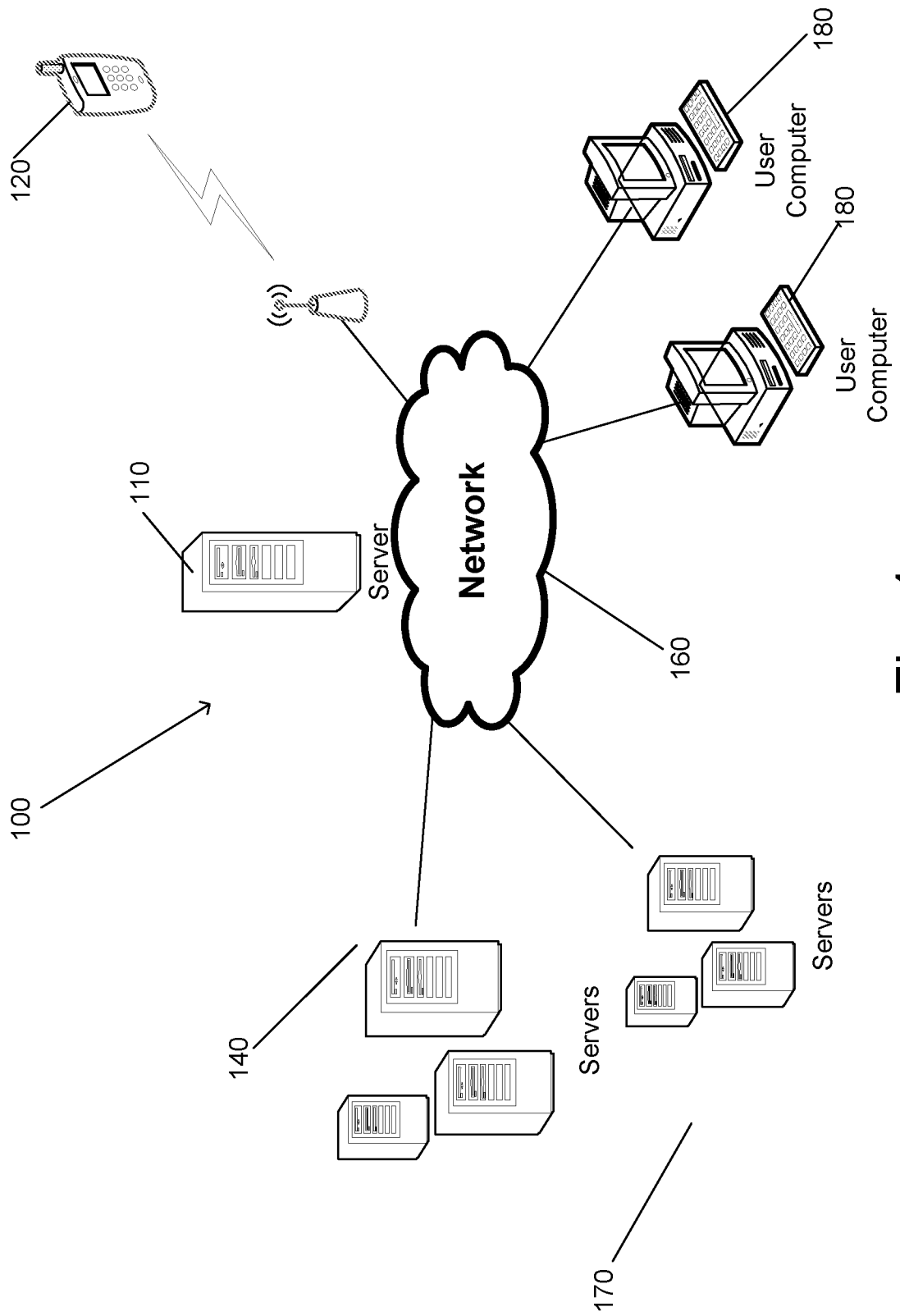


Fig. 1

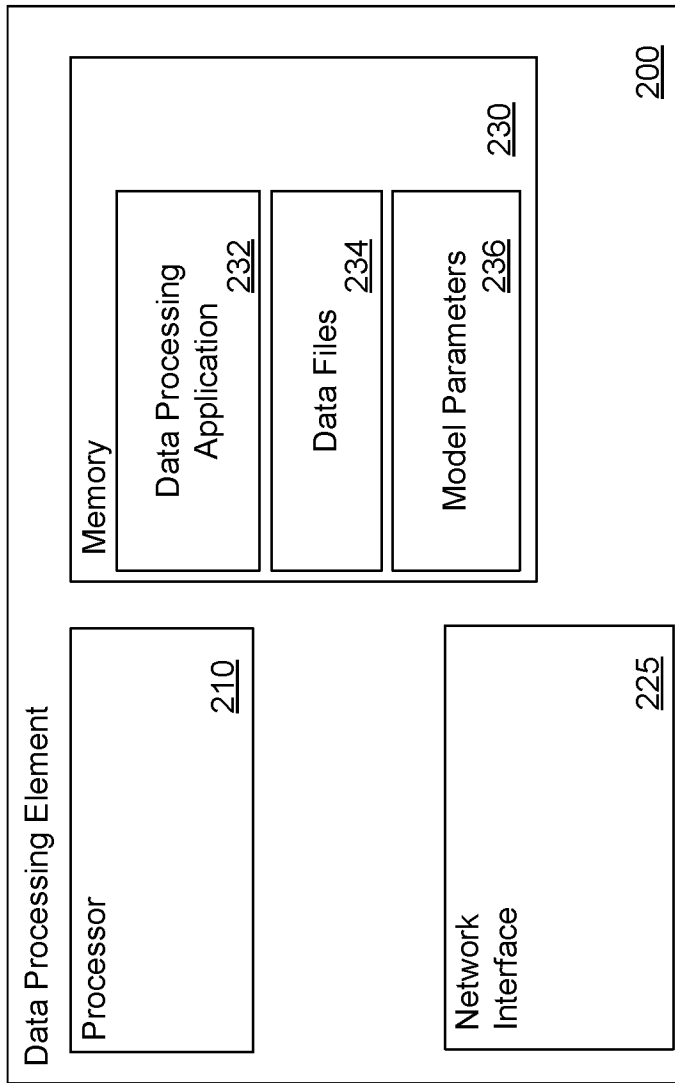


FIG. 2

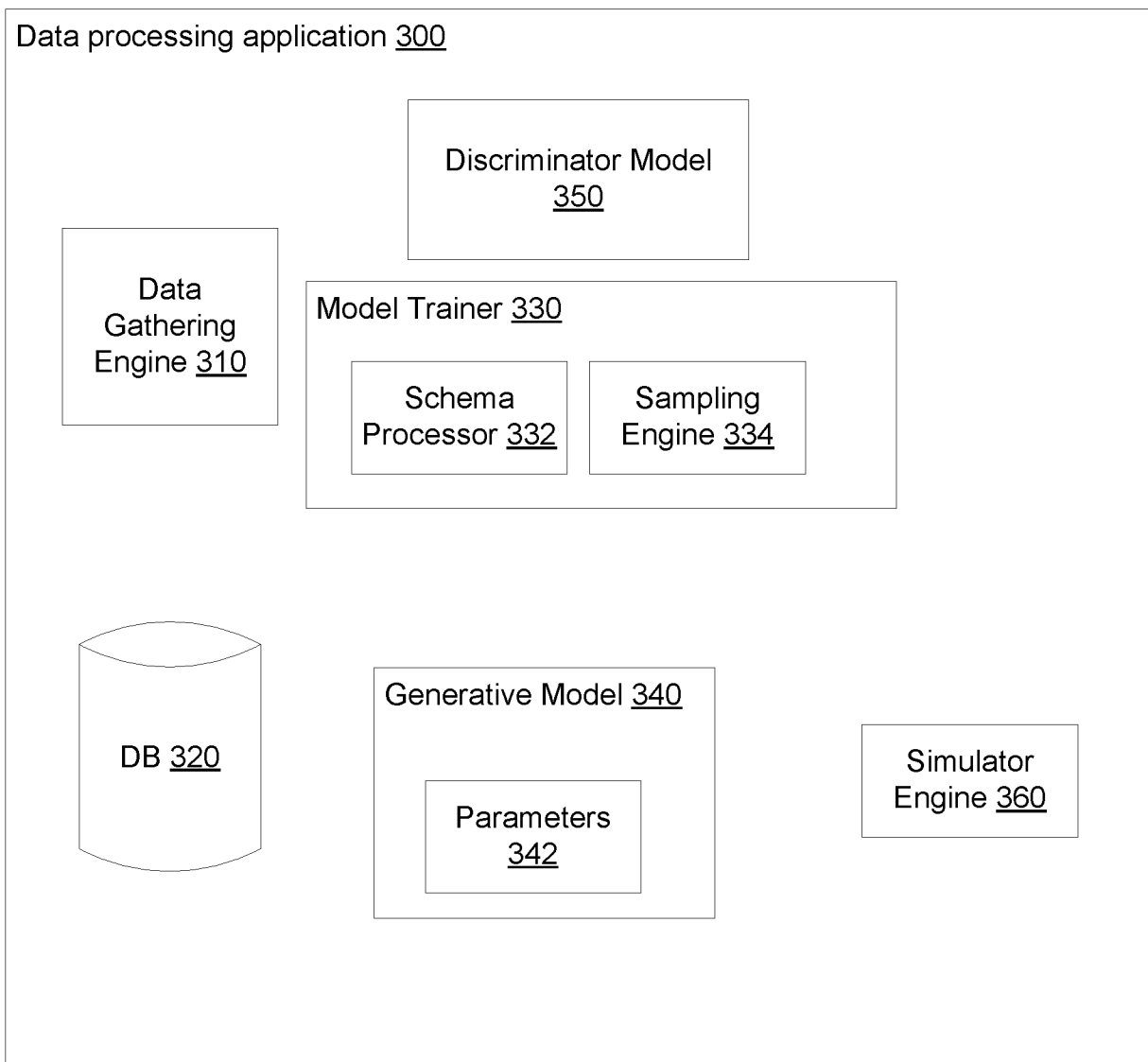


FIG. 3

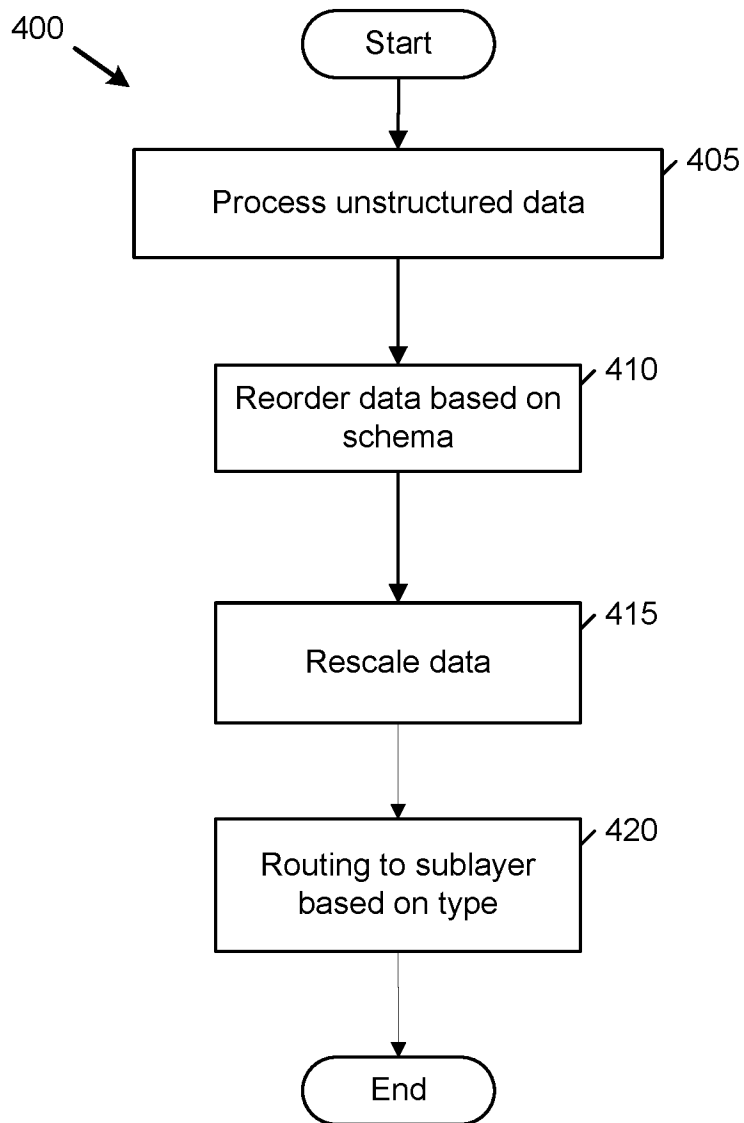


FIG. 4

t0

ID	Test 1	Image 1	Test 2	Gender
1	TRUE	V1	-1.5	M
2	FALSE	V2	1.5	M
.				
.				
.				
n	TRUE	Vn	4.2	F

510

t1

ID	Test 1	Image 1	Test 2	Gender
1	TRUE	V1		M
2	FALSE	V2		M
.				
.				
.				
n	TRUE	Vn	3.5	F

520

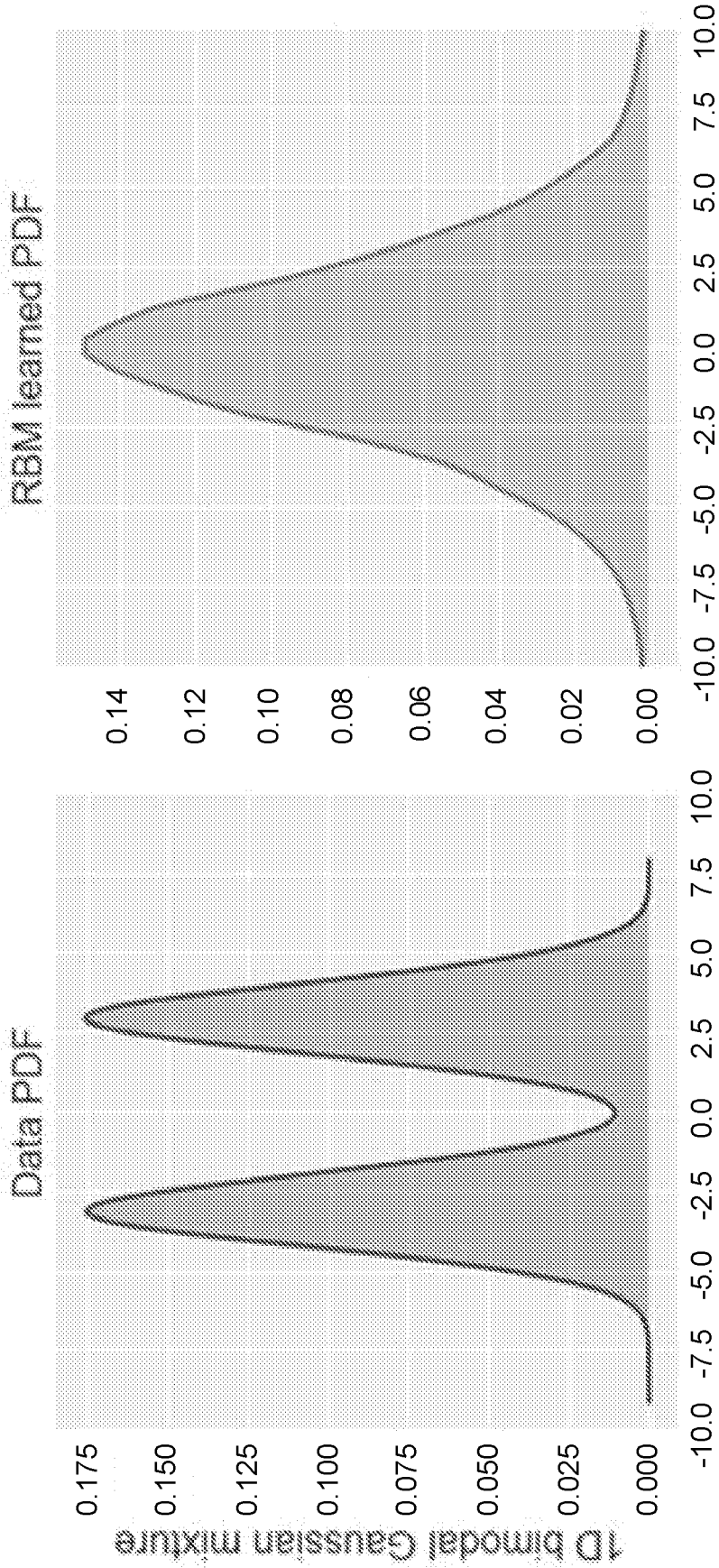
-
-
-

tn

ID	Test 1	Image 1	Test 2	Gender
1	TRUE	V1		M
2	FALSE	V2	1.7	M
.				
.				
.				
n	TRUE	Vn	2.8	F

530

FIG. 5



610

620

FIG. 6

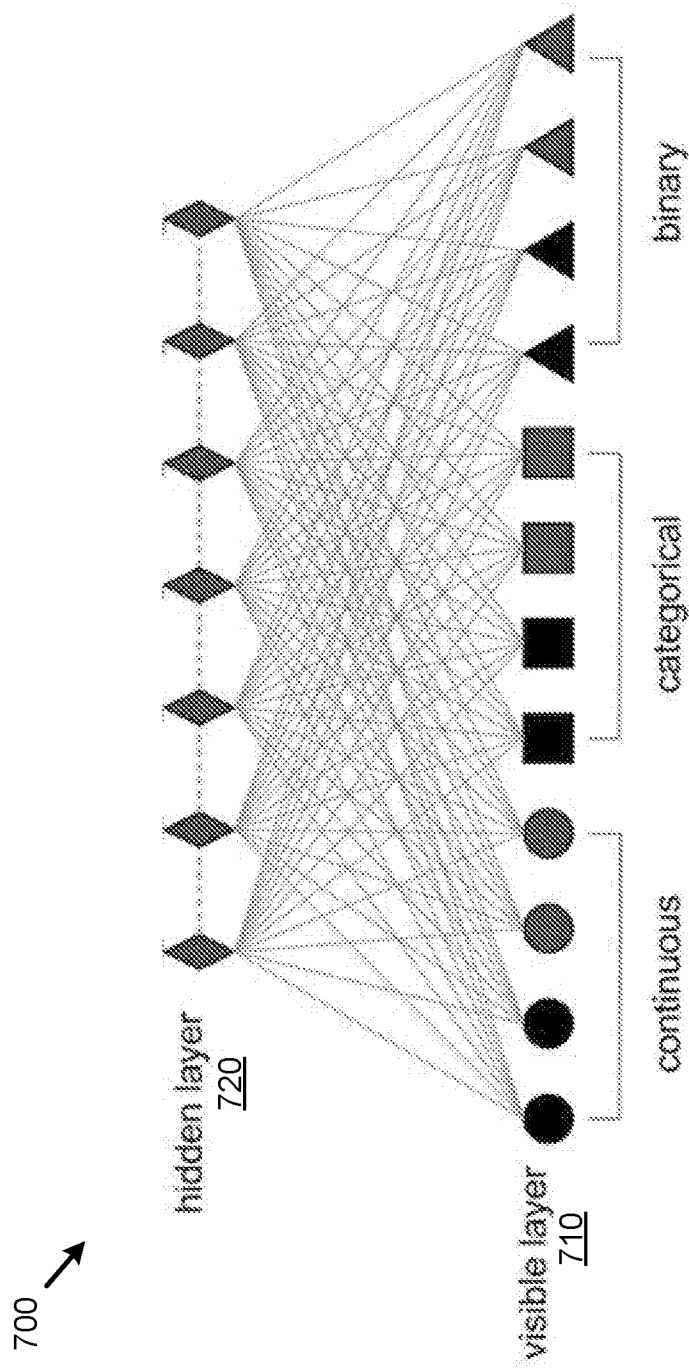


FIG. 7

800



```
model:
  type: BoltzmannMachine
  layers:
    -
      layer_type: BernoulliLayer
      num_units: 8
      center: False
      time_dependent: False
    -
      layer_type: GaussianLayer
      num_units: 5
      center: True
      time_dependent: False
```

FIG. 8

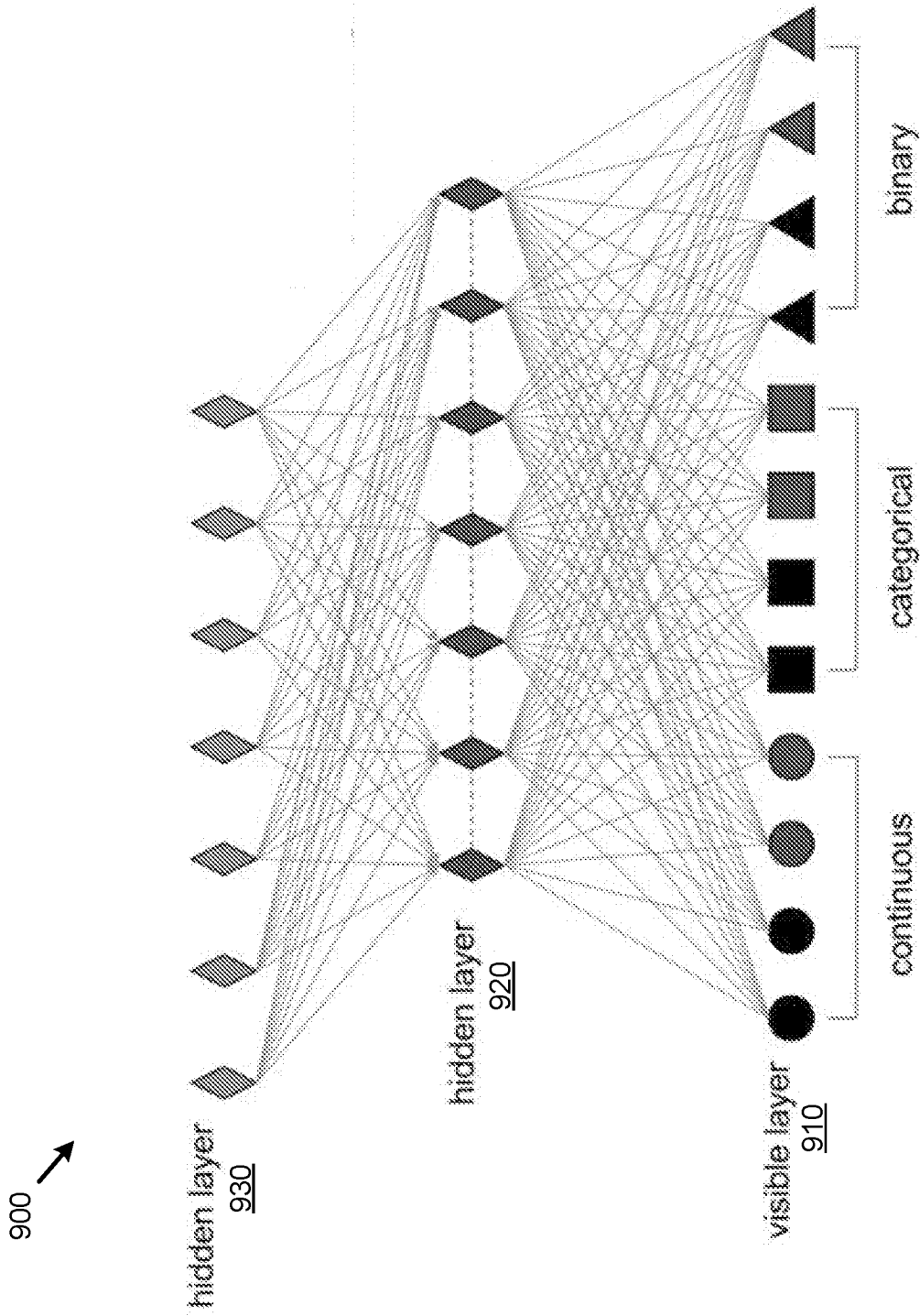


FIG. 9

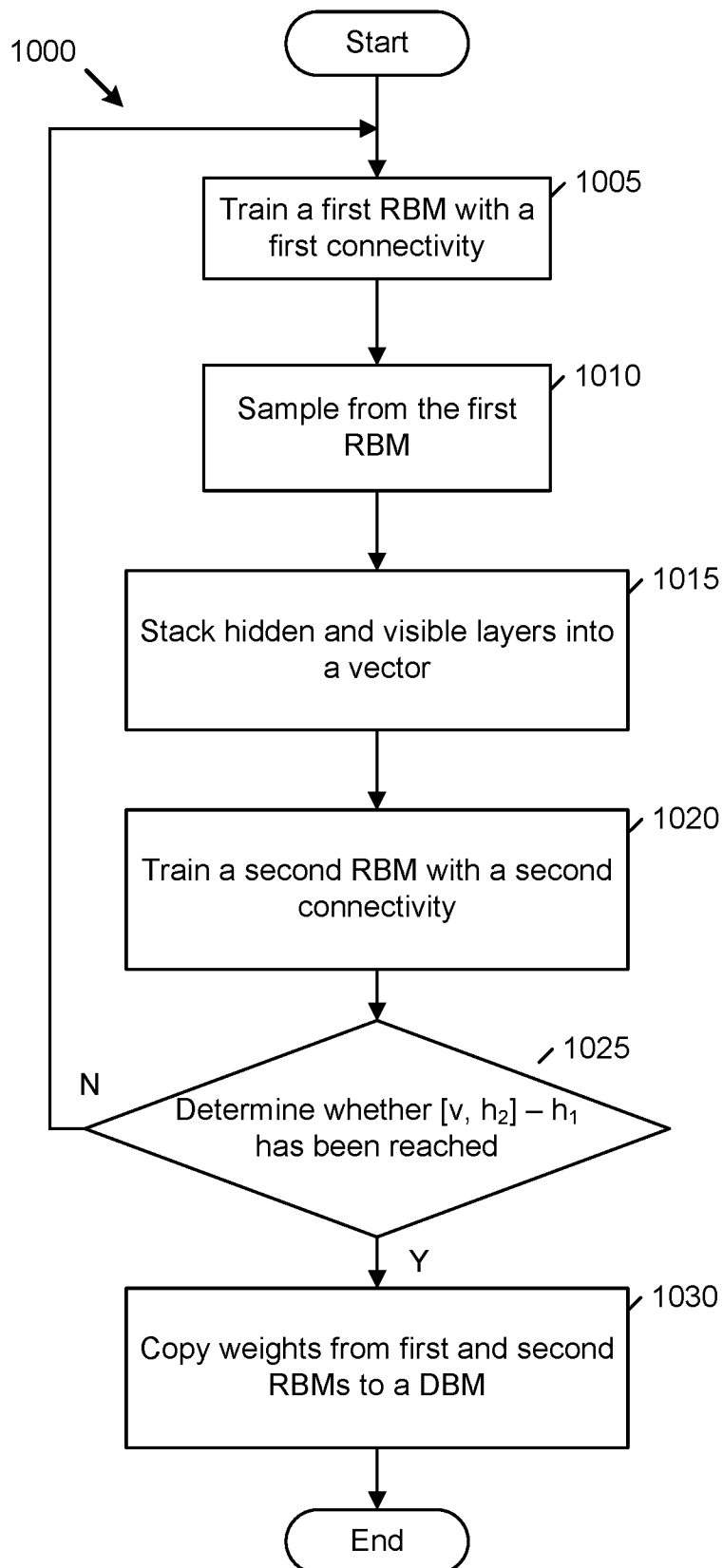


FIG. 10

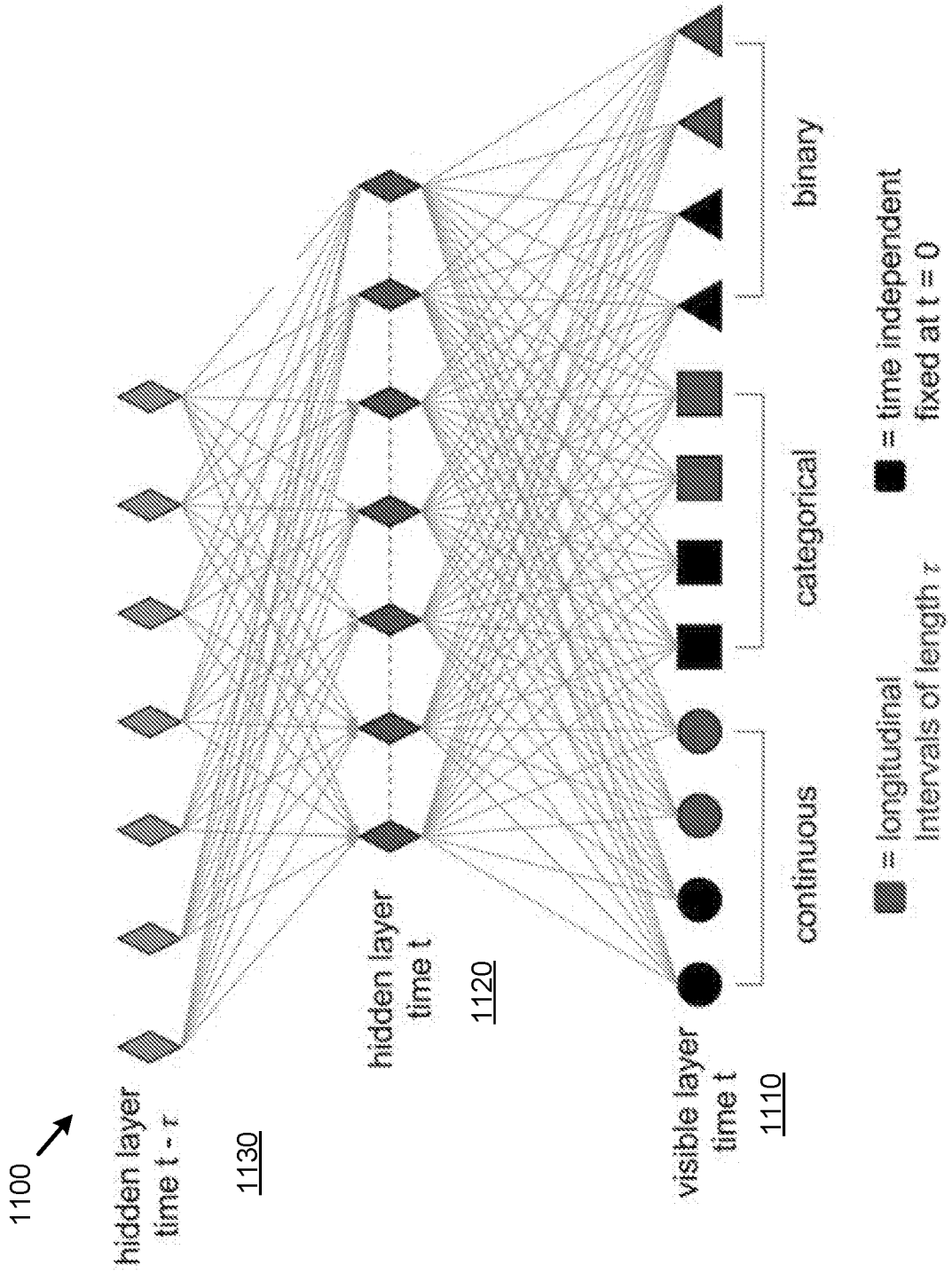
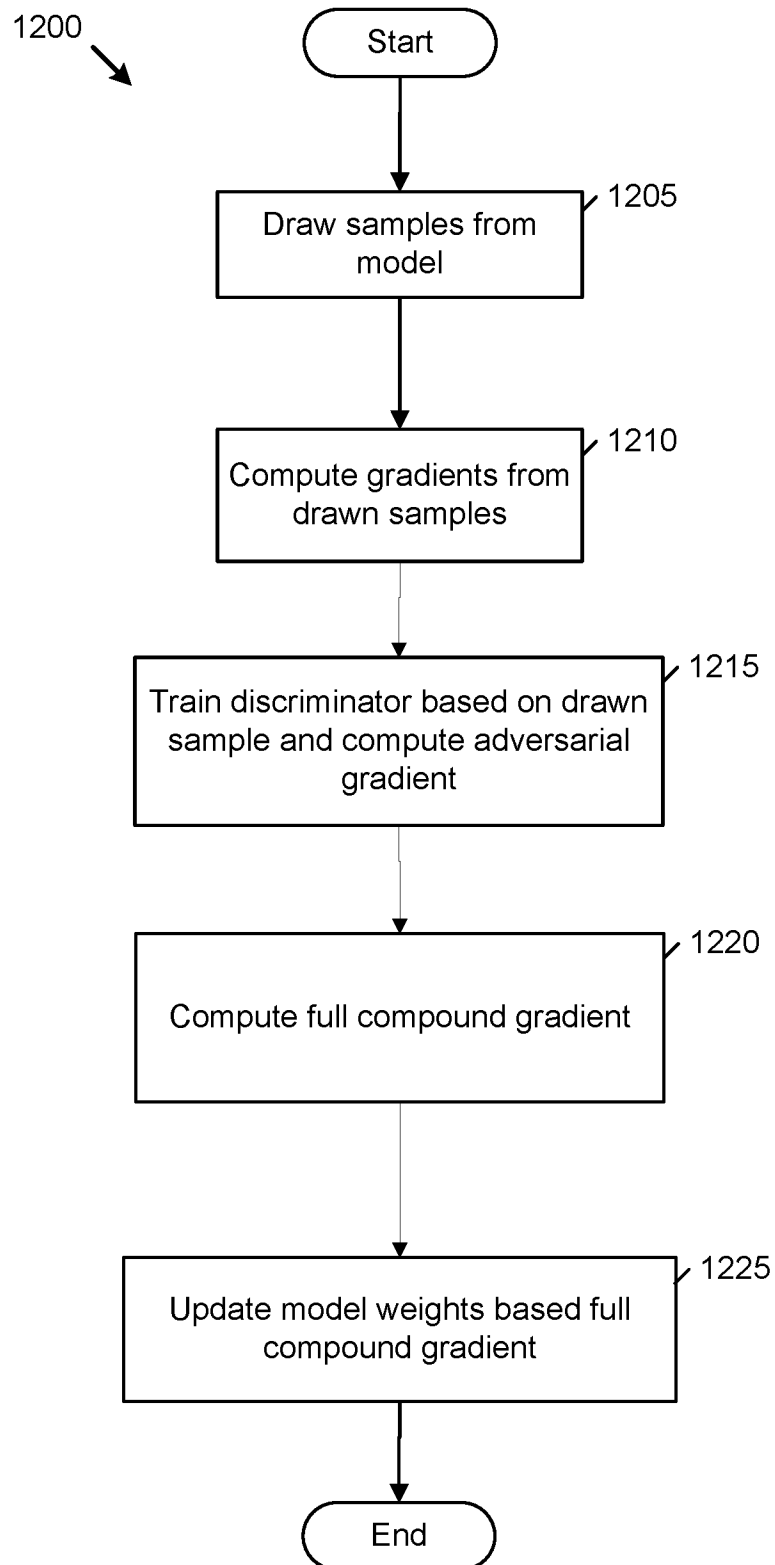
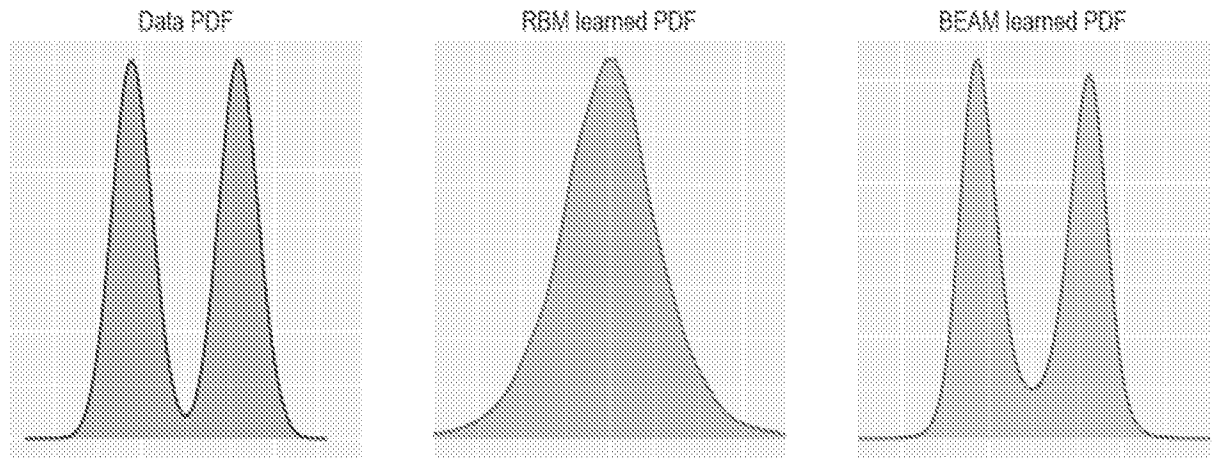


FIG. 11

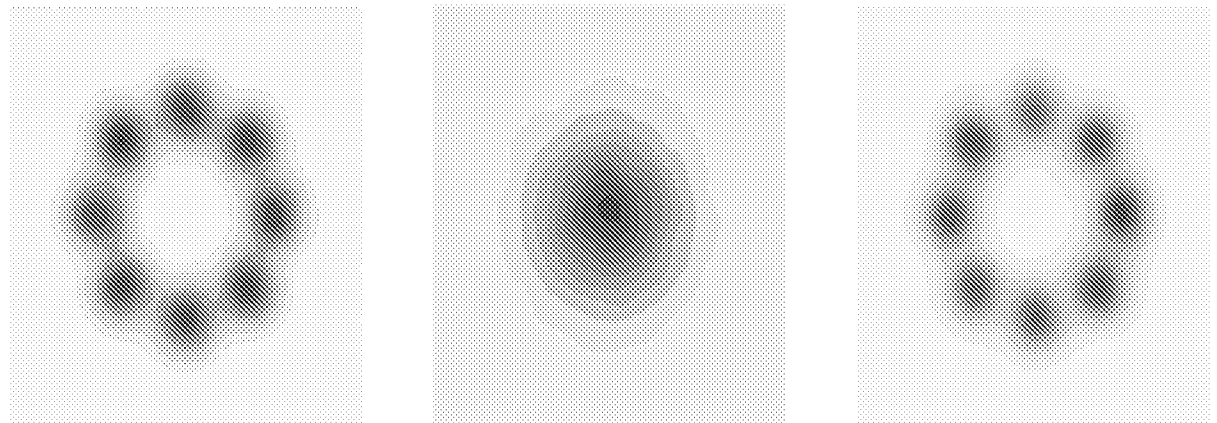
12/17

**FIG. 12**

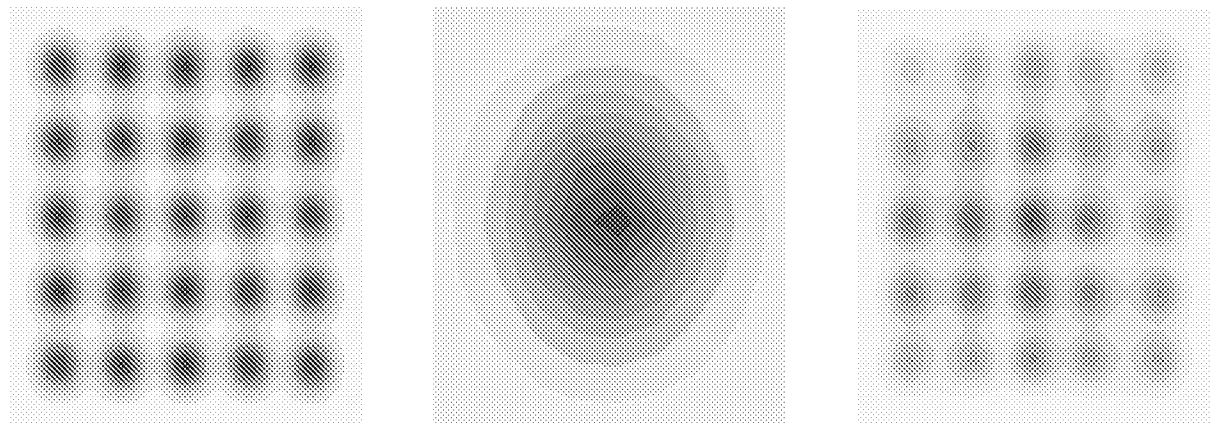
Comparison between RBMs and BEAMs on multimodal data



1310



1320



1330

FIG. 13

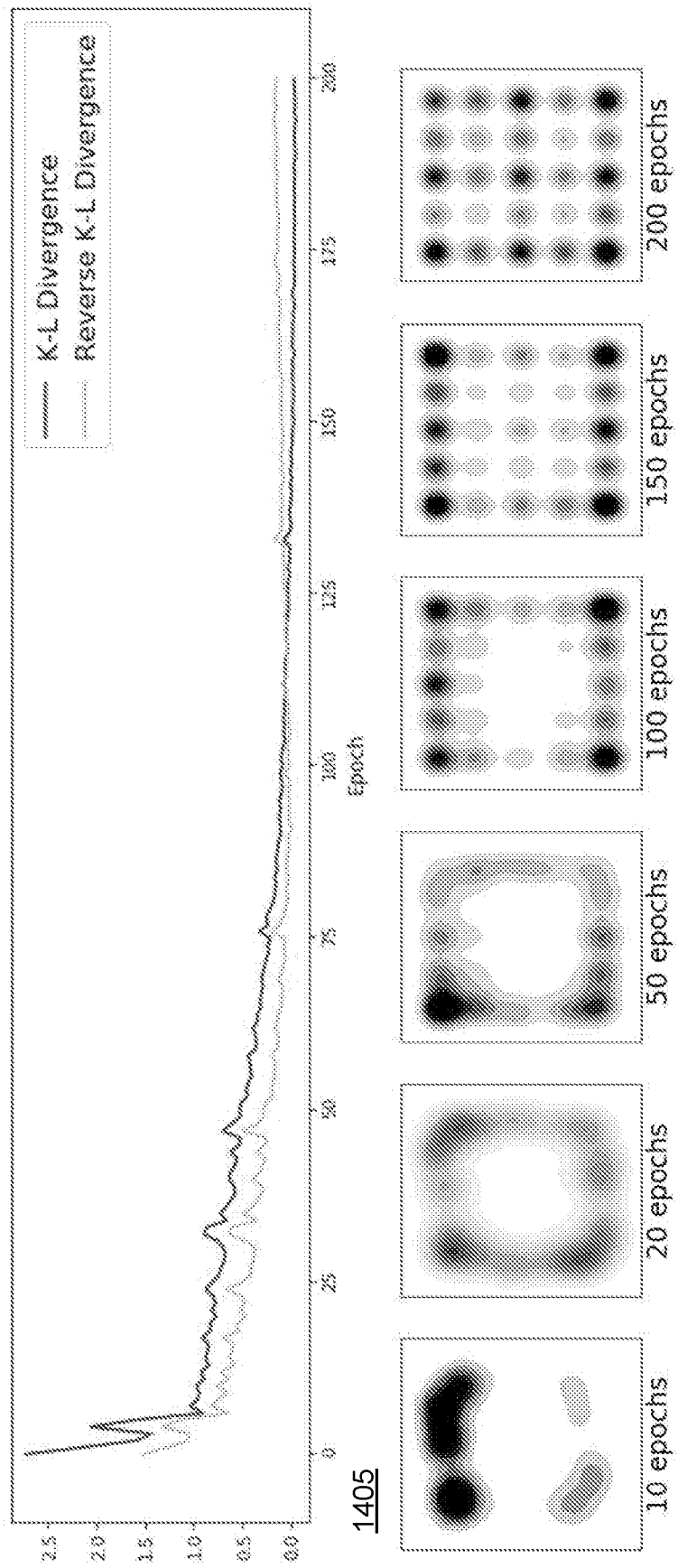


FIG. 14

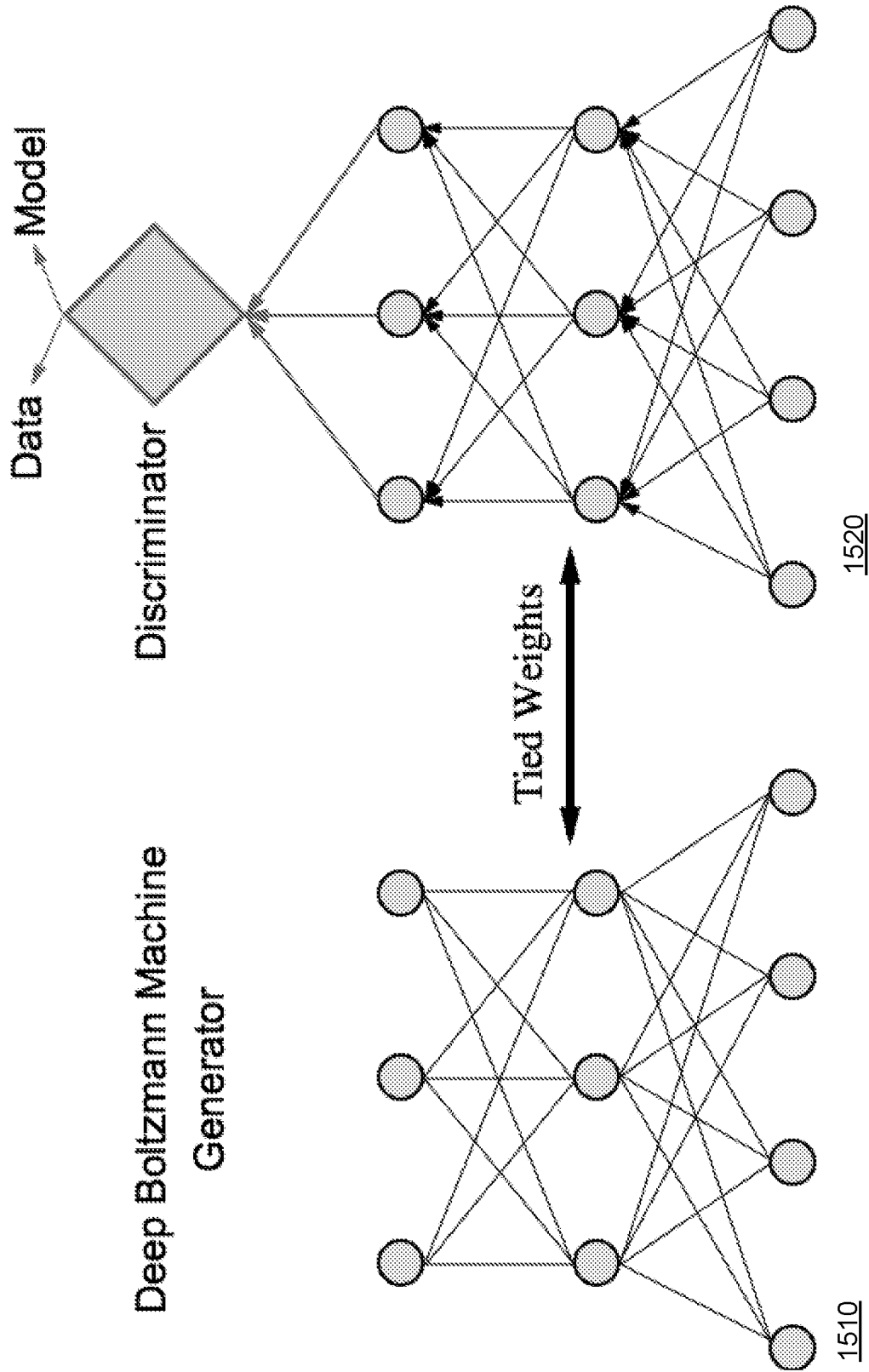


FIG. 15

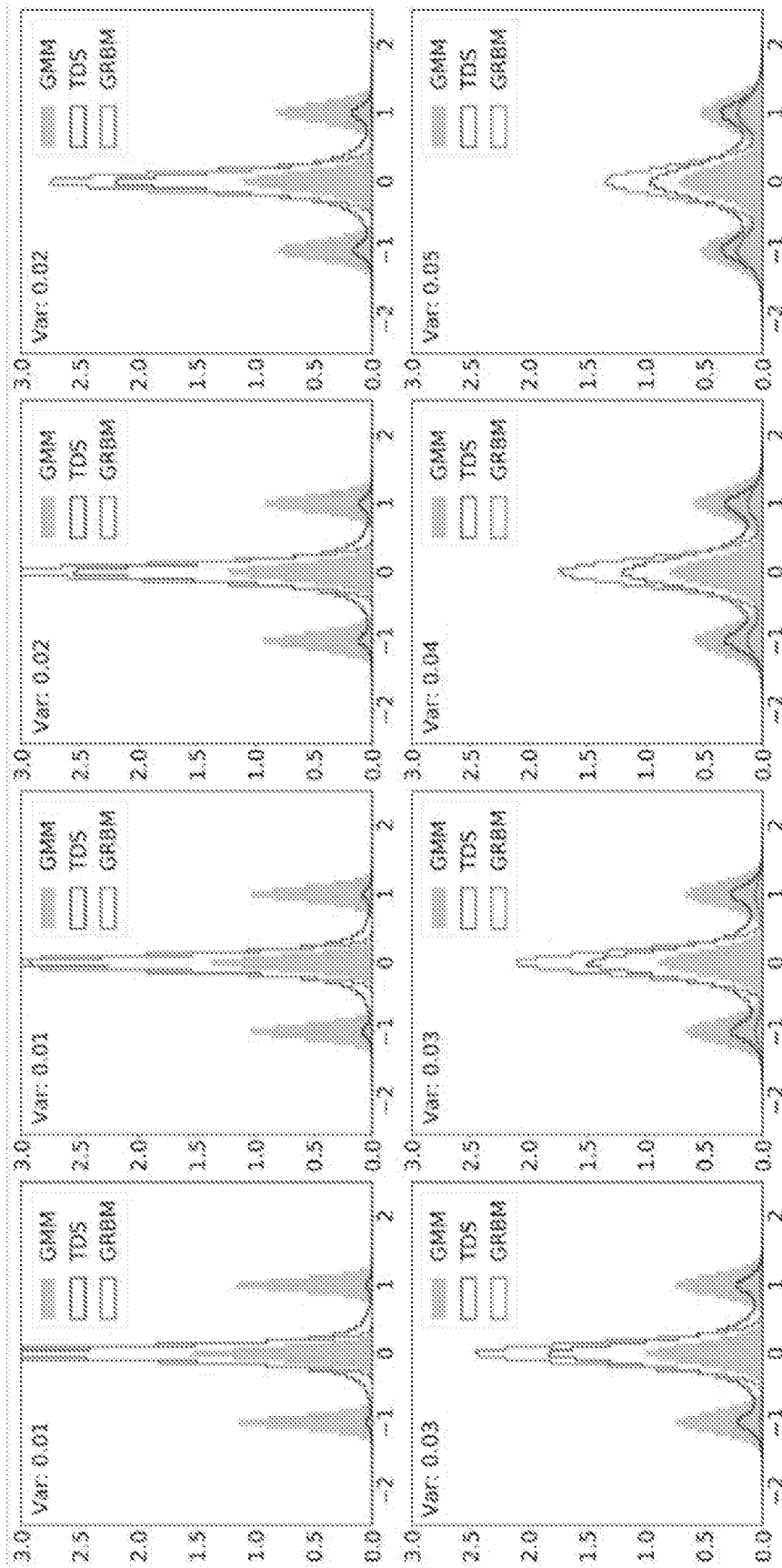
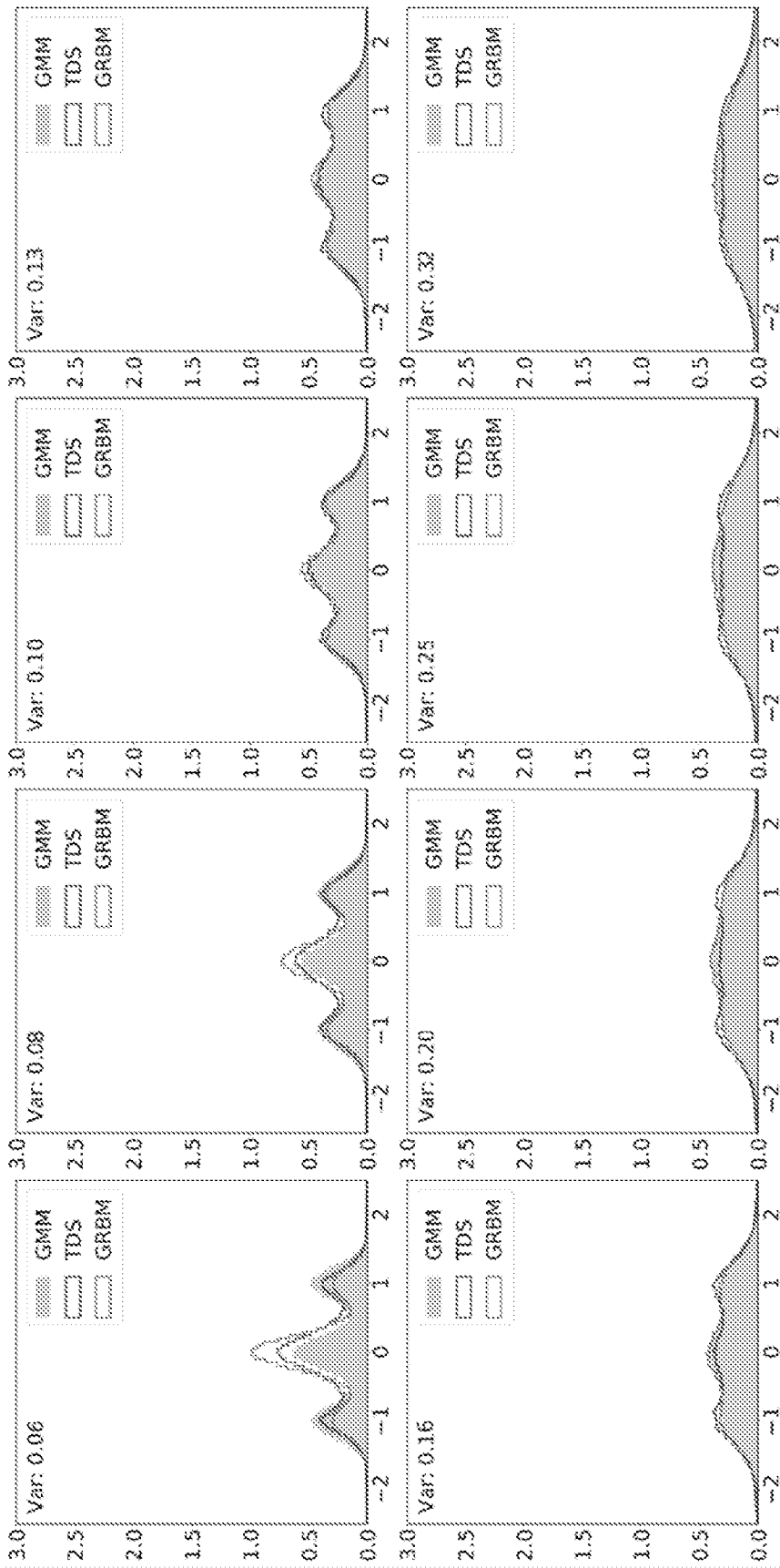


FIG. 16A



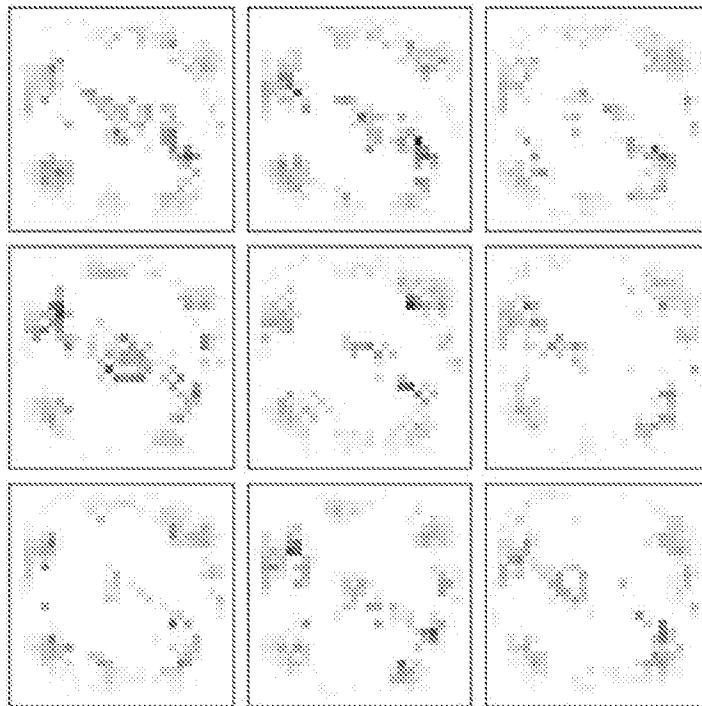
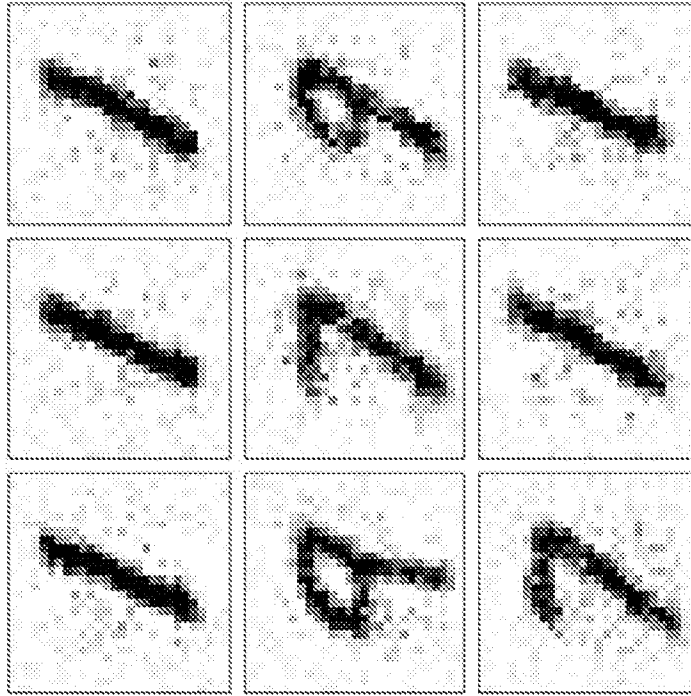


FIG. 17

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US19/13870

A. CLASSIFICATION OF SUBJECT MATTER
IPC - G06N 3/08, 3/02 (2019.01)
CPC - G06N 3/0445, 3/0454, 3/0472, 3/08, 3/02

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
See Search History document

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
See Search History document

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
See Search History document

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	WO 2016/145379 A1 (WILLIAM MARSH RICE UNIVERSITY) 15 September 2016; entire document.	1-20
A	US 2017/0372193 A1 (SIEMENS HEALTHCARE GMBH) 28 December 2017; entire document.	1-20
A	AKHTAR, S et al. "Improving the Robustness of Neural Networks Using K-Support Norm Based Adversarial Training"; IEEE Access; Publication [online]. 28 December 2016 [retrieved 17 March 2019]. Retrieved from the Internet: <URL: https://www.researchgate.net/publication/311879529_Improving_the_Robustness_of_Neural_Networks_Using_K-Support_Norm_Based_Adversarial_Training >; pp 1-9	1-20

Further documents are listed in the continuation of Box C. See patent family annex.

* Special categories of cited documents:	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
"A" document defining the general state of the art which is not considered to be of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"E" earlier application or patent but published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&" document member of the same patent family
"O" document referring to an oral disclosure, use, exhibition or other means	
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search
18 March 2019 (18.03.2019)

Date of mailing of the international search report
27 MAR 2019

Name and mailing address of the ISA/
Mail Stop PCT, Attn: ISA/US, Commissioner for Patents
P.O. Box 1450, Alexandria, Virginia 22313-1450
Facsimile No. 571-273-8300

Authorized officer
Shane Thomas
PCT Helpdesk: 571-272-4300
PCT OSP: 571-272-7774