
S

Synthetic Data Sets

Sargur N. Srihari
Department of Computer Science and
Engineering, University at Buffalo, The State
University of New York, Buffalo, NY, USA

Synonyms

[Anonymized Data](#); [Augmented Data](#)

Glossary

Machine Learning	Method to perform a task from examples
Probabilistic Graphical Model (PGM)	A graphical representation of joint probability distributions where nodes represent variables and edges represent influences.
Regularization	Method used to generalize to previously unobserved evidence

Definition

Synthetic data, as used in computing, consists of symbols and values relevant to a given situation that are not directly observed but are artificially constructed.

Introduction

Synthetic data is created to design or improve performance of information processing systems. Principal uses of synthetic data are in designing machine learning systems to improve their performance and in the design of privacy-preserving algorithms that need to filter information to preserve confidentiality.

Key Points

The design of many information processing systems rely on representative samples of data. Such data may not be readily available due to the cost of acquiring the data or due to privacy concerns. This situation can be alleviated by constructing representative samples from available data sets. This can be done by applying appropriate transformations, e.g., injecting noise, that preserve semantic relationships, e.g., class labels.

Historical Background

Synthetic data has long been used to demonstrate design principles in pedagogy. In the early days of machine learning research in the 1950s statistically representative datasets were hard to obtain. Thus augmenting data sets by using appropriately transformed data was a work-around. Even today, state-of-the-art methods such as deep learning

improve their performance when provided with synthetic data. In information systems that deal with human information, synthetic data is useful to anticipate unforeseen situations in which data may be compromised.

Synthetic Data Sets

A prerequisite for the construction of many algorithms with real-world applications is the availability of representative data sets. Artificial intelligence algorithms based on machine learning, e.g., to perform speech recognition or recognize objects in scenes, require data to determine necessary parameters. Supervised learning algorithms, such as neural networks, need examples of input and output, commonly referred to as *ground truth* to determine weights of connections. Even after models are constructed, sample data may be needed to perform inference with the models. *Probabilistic graphical models* (Koller and Friedman 2009; Srihari 2017), which are widely used for artificial intelligence applications, are notoriously intractable when it is necessary to respond to queries. The most popular approximate inference methods depend on generating samples (also referred to as *particles*) from the model.

Ground-truthed data sets that are large enough to make statistically significant conclusions are rarely available. This may be due to privacy concerns, such as with social network data. In such cases, automatic construction of synthetic data sets with characteristics similar to real-world data is necessitated. Many machine learning algorithms, including those based on *deep learning* improve their performance when synthetic data sets are used.

Data Set Augmentation for Machine Learning

The central problem of machine learning is to design algorithms that will perform well not just on training data but on new inputs as well. Regularization refers to any modification we make to a learning algorithm to reduce its generalization error but not its training error. The goal is to reduce the test error even at the expense of increasing training error.

Many forms of regularization available. Principal goals of regularization are to encode prior knowledge and to express preference for a simpler model. Some methods put extra constraints on the objective function that is optimized by the learning algorithm – they are equivalent to a soft constraint on parameter values, resulting in improved performance on a test set. Some strategies for regularization are: parameter norm penalties (L_2 - and L_1 - regularization), norm penalties as constrained optimization, data set augmentation, early stopping, sparse representations, bagging and other ensemble methods, dropout, and tangent methods. Data augmentation is one of the most effective methods of regularization, since a model generalizes better when it is trained on more data.

Data Sets to Overcome Privacy

Privacy issues are paramount when dealing with domains such as social networks. Examples of social network analytics are: attribute prediction given user social connections and phone usage log. It is of most interest to service providers and mobile equipment manufacturers to know particular user attributes and behavior patterns in order to provide better service and appropriate hardware technologies. This makes prediction of user attributes and behavior patterns one of the most actual research issues. Another is predicting future states of the social network. This would be useful for planning and marketing purposes.

The social network data set should suit most of the needs of the proposed analytic algorithms, e.g., it contains data for a large number of users for a long period of time (more than a year), has a unique combination of social interactions data and contextual attributes and basic user attributes (gender, age, marital status, etc.), and personal data in combination with location data and phone status information makes it possible to reveal novel patterns and relationships between users semantic location and its phone usage, etc.

Methods for Generating Synthetic Data

Data Augmentation

Data augmentation is easiest for classification, A classifier takes high-dimensional input \mathbf{x} and summarizes it with a single category identity y . Since the main task of the classifier is to be invariant to a wide variety of transformations, we can generate new samples (\mathbf{x}, y) just by transforming inputs. This approach is not easily generalized to other problems, e.g., for the density estimation problem it is not possible generate new data without solving density estimation.

Data set augmentation is very effective for the classification problem of object recognition. Images are high-dimensional and include a variety of variations, and may be easily simulated. The transformations include translation, rotation, scaling, and compression. Translating the images a few pixels can greatly improve performance. This is observed even when a deep learning classifier has been designed to be invariant using convolution and pooling. Rotating and scaling are also effective.

A precaution in data augmentation for classification is that we should not apply transformations that would change the class. For instance, in the task of recognizing handwritten characters, e.g., b versus d and 6 versus 9 , horizontal flips and 180 degree rotations are not appropriate transformations. Some transformations are not easy to perform, e.g., out-of-plane rotation cannot be implemented as a simple geometric operation on pixels.

Injecting noise into the input of a neural network can be seen as data augmentation. Neural networks are not robust to noise. To improve robustness, train them with random noise applied to their inputs. Part of some unsupervised learning, such as denoising autoencoder. Noise can also be applied to hidden units. Dropout, a powerful regularization strategy, can be viewed as constructing new inputs by multiplying by noise.

Hand-designed data set augmentation can dramatically improve performance. When comparing algorithms A and B , the same data set augmentation technique should be used for both. If A performs poorly with no dataset augmentation

and B performs well with synthetic transformations of the input, reason may be the data set rather than algorithm. Adding Gaussian noise is considered part of learning while cropping input images is not.

Generative Models and Sampling

Data can be synthesized from probabilistic generative models. Such models are typically in the form of directed or undirected graphs where nodes represent variables and edges represent influences. The model can be constructed from available data and any number of samples can be synthesized from the resulting model.

Synthesized data plays a role not only in constructing models but also in performing inference with models. A typical inference task is to determine the conditional probability of a variable when value assignments to evidence variables are given. A naive approach is to eliminate the irrelevant variables by summing over all possible values that they can take. But this is intractable since the number of terms increases exponentially with the number of variables.

In the *particle-based* approach to obtain an approximate answer, the data are generated appropriate to the query. Particle-based inference methods, also known as *Monte Carlo* methods, synthesize samples from a probability distribution. Sampling provides a flexible way to approximate many sums and integrals at reduced cost. Sometimes for speedup of a costly but tractable sum, e.g., subsample training cost with minibatches. In other cases, learning algorithms require us to approximate an intractable sum or integral, e.g., gradient of the log partition function of an undirected model.

If the model is in the form of a directed acyclic graph (also called a Bayesian network), samples can be obtained by proceeding through the variables in such a way that the parent nodes are sampled before the descendant nodes. This is known as ancestral or forward sampling. In the case where the query is conditional – when evidence variables are observed – it becomes necessary to modify forward sampling by using a method such as importance sampling.

When the model is an undirected graph, a *Markov chain Monte Carlo (MCMC)* method is used. In one such method known as *Gibbs sampling*, one of the variables is changed while holding the other variables constant. It is the conceptually simplest approach for drawing samples from an undirected graph. It works as follows. Suppose we have a graphical model over an n -dimensional vector of random variables \mathbf{x} . We iteratively visit each variable x_i and draw a sample conditioned on all the other variables, i.e., from $p(x_i | \mathbf{x}_{\setminus i})$. Due to the separation properties of the graphical model, we can equivalently condition on only the neighbors of x_i .

Key Applications

Some uses of synthetic data are: (i) providing machine learning systems with more data than what is readily available so that they can improve their generalization performance, (ii) performing tractable inference from probabilistic models, and (iii) in designing filters for information systems so that confidentiality can be preserved.

Future Directions

An area of current research in deep learning is disentangling factors of variation in data (Goodfellow et al. 2016). Once such factors are discovered they can be used within generative models, which can then be sampled to synthesize realistic data sets.

Cross-References

- ▶ [Gibbs Sampling](#)
- ▶ [Markov Monte Carlo Model](#)
- ▶ [Models of Social Networks](#)
- ▶ [Probabilistic Graphical Models](#)

References

- Bishop C (2006) Pattern recognition and machine learning. Springer, New York, NY
- Goodfellow I, Bengio Y, Courville A (2016) Deep learning. MIT Press, Cambridge, MA
- Koller D, Friedmann N (2009) Probabilistic graphical models: principles and techniques. MIT Press, Cambridge, MA
- Srihari S (2017) Lecture slides and videos on machine learning and PGMs. <http://www.cedar.buffalo.edu/~srihari/CSE674>