

Previews of TDWI course books offer an opportunity to see the quality of our material and help you to select the courses that best fit your needs. The previews cannot be printed.

TDWI strives to provide course books that are contentrich and that serve as useful reference documents after a class has ended.

This preview shows selected pages that are representative of the entire course book; pages are not consecutive. The page numbers shown at the bottom of each page indicate their actual position in the course book. All table-of-contents pages are included to illustrate all of the topics covered by the course.

This page intentionally left blank.

Machine & Deep Learning: Delivering Insights from Big Data

You Will Learn

- The foundations of machine learning in chaos theory, game theory, and algorithms
- What "machine talk" is and how we architect for it
- How to bring search implementation theory to data with deep learning
- Introduction to TensorFlow, a library of open source algorithms
- Implementing machine and deep learning, including what processing techniques are involved
- Available analytics insights, outcomes, mashups, and KPIs
- How to visualize the results of machine and deep learning



Chaos Theory

- Chaos theory has been defined as the science that determines order in the randomness and unpredictability that exist in the natural and social systems
- Ex: coastlines, mountains, clouds, galaxy clusters, leaves, heart rhythms, etc.
- · Chaos theory underlying principles
 - Non-linearity
 - Determinism
 - Sensitivity to initial conditions
 - Sustained irregularity in the behavior of the system
 - Unpredictable long-term behavior
 - Self-organization
 - Positive feedback loops
 - Qualitative character of systems rather than numerical predictions of future states
 - Unstable
 - Aperiodic no periodic repetition of values

SystemOrderChaosRandomnessParadigmatic ExampleClocks, PlanetsClouds, WeatherSnow on TV ScreenPredictabilityVery HighFinite, Short TermNone, Simple LawsEffect of Small ErrorsVery SmallExplosiveNothing BUT ErrorsSpectrumPureYes!Noisy, BroadDimensionFiniteLowInfinite	Characteristics of Systems				
Paradigmatic ExampleClocks, PlanetsClouds, WeatherSnow on TV ScreenPredictabilityVery HighFinite, Short TermNone, Simple LawsEffect of Small ErrorsVery SmallExplosiveNothing BUT ErrorsSpectrumPureYes!Noisy, BroadDimensionFiniteLowInfinite	System	Order	Chaos	Randomness	
PredictabilityVery HighFinite, Short TermNone, Simple LawsEffect of Small ErrorsVery SmallExplosiveNothing BUT ErrorsSpectrumPureYes!Noisy, BroadDimensionFiniteLowInfinite	Paradigmatic Example	Clocks, Planets	Clouds, Weather	Snow on TV Screen	
Effect of Small ErrorsVery SmallExplosive ErrorsNothing BUT ErrorsSpectrumPureYes!Noisy, BroadDimensionFiniteLowInfinite	Predictability	Very High	Finite, Short Term	None, Simple Laws	
Spectrum Pure Yes! Noisy, Broad Dimension Finite Low Infinite	Effect of Small Errors	Very Small	Explosive	Nothing BUT Errors	
Dimension Finite Low Infinite	Spectrum	Pure	Yes!	Noisy, Broad	
	Dimension	Finite	Low	Infinite	
Control Easy Tricky, Very Poor Effective	Control	Easy	Tricky, Very Effective	Poor	
Attractor Point, Cycle, Torus Strange, Fractal No!	Attractor	Point, Cycle, Torus	Strange, Fractal	No!	

Fractals

Fractals are geometric shapes that are very complex and infinitely detailed. You can zoom in on a section and it will have just as much detail as the whole fractal. They are recursively defined and small sections of them are similar to large ones. One way to think of fractals for a function f(x) is to consider x, f(x), f(f(x)), f(f(f(x))), f(f(f(x))), etc. Fractals are related to chaos because they are complex systems that have definite properties.

Game Theory

Game Theory

- Any situation in which individuals must make strategic choices and in which the final outcome will depend on what each person chooses to do can be viewed as a game.
- Game theory models seek to portray complex strategic situations in a highly simplified setting.
- All games have three basic elements:
 - Players
 - Strategies
 - Payoffs
- Players can make binding agreements in *cooperative* games, but can not in *noncooperative* games, which are studied in this chapter.

Equilibrium Concepts

- In the theory of markets an equilibrium occurred when all parties to the market had no incentive to change his or her behavior.
- When strategies are chosen, an equilibrium would also provide no incentives for the players to alter their behavior further.
- The most frequently used equilibrium concept is a Nash equilibrium.

An Illustrative Advertising Game

- Two firms (A and B) must decide how much to spend on advertising
- Each firm may adopt either a higher (H) budget or a low (L) budget.
- The game is shown in extensive (tree) form.
- A makes the first move by choosing either H or L at the first decision "node."
- Next, B chooses either H or L, but the large oval surrounding B's two decision nodes indicates that B does not know what choice A made.





What techniques we will see

- kNN algorithm
- Winnow algorithm
- Naïve Bayes classifier
- Decision trees
- Reinforcement learning (Rocchio algorithm)
- Genetic algorithm







- The unsupervised weight adapting algorithms are usually based on some form of global competition between the neurons.
- Applications of self-organizing networks are:
 - **clustering:** the input data may be grouped in `clusters' and the data processing system has to find these inherent clusters in the input data.
 - vector quantisation: this problem occurs when a continuous space has to be discretised. The input of the system is the n-dimensional vector *x*, the output is a discrete representation of the input space. The system has to find optimal discretisation of the input space.
 - **dimensionality reduction**: the input data are grouped in a subspace which has lower dimensionality than the dimensionality of the data. The system has to learn an "optimal" mapping.
 - **feature extraction:** the system has to extract features from the input signal. This often means a dimensionality reduction as described above.





Types of Search Algorithm (Google Driven)

- Page-Rank Algorithm
- Penguin Algorithm
- Panda Algorithm
- Humming-Bird Algorithm



















Readout Layer

• Finally, we add a softmax layer, just like for the one layer softmax regression.

```
W_fc2 = weight_variable([1024, 10])
b_fc2 = bias_variable([10])
```

```
y_conv=tf.nn.softmax(tf.matmul(h_fc1_drop, W_fc2) + b_fc2)
```













Watson – the computer system we developed to play Jeopardy! is based on the DeepQA softate archtiecture. Here is a look at the **DeepQA** architecture. This is like looking inside the brain of the Watson system from about 30,000 feet high.

Remember, the intended meaning of natural language is ambiguous, tacit and highly contextual. The computer needs to consider many possible meanings, attempting to find the evidence and inference paths that are most confidently supported by the data.

So, the primary computational principle supported by the DeepQA architecture is to assume and pursue multiple interpretations of the question, to generate many plausible answers or **hypotheses** and to collect and evaluate many different competing evidence paths that might support or refute those hypotheses.

Each component in the system adds assumptions about what the question might means or what the content means or what the answer might be or why it might be correct.

DeepQA is implemented as an extensible architecture and was designed at the outset to support interoperability.

<UIMA Mention>

For this reason it was implemented using UIMA, a framework and OASIS standard for interoperable text and multi-modal analysis contributed by IBM to the open-source community.

Over 100 different algorithms, implemented as UIMA components, were integrated into this architecture to build *Watson*.

Microsoft

- Computational Neuroscience
- Hebbian learning
- Hopfield Net
- Bolzmann machines
- Memory models
- Bio-inspired Al
- RNN
- Computer vision NLP
- IR



Tools for Data Preparation/Feature Engineering

- Languages/Environments
 - PigLatin
 - HiveQL
 - Need to deal with mismatch between offline/online feature generation
- Java/Scala APIs
 - Crunch (Cloudera)
 - <u>Scoobi</u> (NICTA)
 - Cascading (Concurrent)
 - Jaql (IBM)

Apache Mahout

- The starting place for MapReduce-based machine learning algorithms
 - Not machine-learning-in-a-box
 - Custom tweaks/modifications are the rule
- A disparate collection of algorithms for:
 - Recommendations
 - Clustering
 - Classification
 - Frequent Itemset Mining











Interactive performance of execution engine Code generation for operators (similarly to Impala) Data is pipelined MPP-style Runs at Facebook scale *Capable of querying other non-HDFS data stores as well*



Notes from https://www.facebook.com/notes/facebook-engineering/presto-interacting-with-petabytes-of-data-at-facebook/10151786197628920

The diagram below shows the simplified system architecture of Presto. The client sends SQL to the Presto coordinator. The coordinator parses, analyzes, and plans the query execution. The scheduler wires together the execution pipeline, assigns work to nodes closest to the data, and monitors progress. The client pulls data from output stage, which in turn pulls data from underlying stages.

The execution model of Presto is fundamentally different from Hive/MapReduce. Hive translates queries into multiple stages of MapReduce tasks that execute one after another. Each task reads inputs from disk and writes intermediate output back to disk. In contrast, the Presto engine does not use MapReduce. It employs a custom query and execution engine with operators designed to support SQL semantics. In addition to improved scheduling, all processing is in memory and pipelined across the network between stages. This avoids unnecessary I/O and associated latency overhead. The pipelined execution model runs multiple stages at once, and streams data from one stage to the next as it becomes available. This significantly reduces end-to-end latency for many types of queries.

The Presto system is implemented in Java because it's fast to develop, has a great ecosystem, and is easy to integrate with the rest of the data infrastructure components at Facebook that are primarily built in Java. Presto dynamically compiles certain portions of the query plan down to byte code which lets the JVM optimize and generate native machine code. Through careful use of memory and data structures, Presto avoids typical issues of Java code related to memory allocation and garbage collection.



The client sends SQL to the Presto coordinator node. The coordinator in this case determines there are operations needed from more than just Hive data.

•A coordinator (a master daemon) uses connectors to get metadata (such as table schema) that is needed to build a query plan. Workers use connectors to get actual data that will be processed by them.

Presto supports pluggable connectors that provide metadata and data for queries. The Hive connector supports Text, SequenceFile, RCFile, ORC and Parquet (?) formats.

Presto does NOT access the Hive server or Hive. It accesses Hive tables in HDFS. The graphic is kept simple but the details are not so simple.



Hive tables and HCatalog

Apache Cassandra

Apache Kafka

Kafka topics = Presto tables,

messages = rows

MySQL

Single node access only -- no sharding

Postgres

Single node access only