

ARTIFICIAL INTELLIGENCE IMPLEMENTATION BOOT CAMP

Classroom Manual

65500PG_2.0 2019

aspetraining.com
877-800-5221



ARTIFICIAL INTELLIGENCE IMPLEMENTATION BOOT CAMP



NOTES:



Part 1: Introduction



The man depicted in this photo is Alan Turing.

Welcome. Did you know...

- AI is more than 40 years old.
- Only 10,000 people in the world can actually “do” serious AI.



NOTES:

Welcome. Did you know...

- Only 23 % of companies have actually deployed AI.
- Only a small percentage (approx.. 4%) are actually performing substantive, state-of-the-art AI operations
- 54% either have no plans to deploy AI applications or have not begun implementing their planned solutions.

— *David Kiron, MIT Sloan Management Review*



NOTES:

Introductions

What is your Name and Job Role?

Your company or team?

Expectations for the class.
Why are you here?

Name something interesting about yourself.



NOTES:

Our Agenda

- **Part 1:** AI – Introduction and working definitions
 - **Part 2:** Big Data and its relationship to AI
 - **Part 3:** Implementing machine learning
 - **Part 4:** Creating concrete value
 - **Part 5:** Machine intelligence and customer experience
 - **Part 6:** Machine intelligence & cybersecurity
 - **Part 7:** Teaming and internal capability
 - **Part 8:** Discussion and charting your course



NOTES:

What to expect from this class

- **Flexibility**
- **Conversations**
- **Literacy and awareness of what's possible with AI right now**
- **General examples of good AI use cases**
- **Focus on history and current landscape of AI and high-level overview of machine learning**
- **Focus on how AI and machine learning can help your organization so you can act on what you learn**



NOTES:

What not to expect from this class

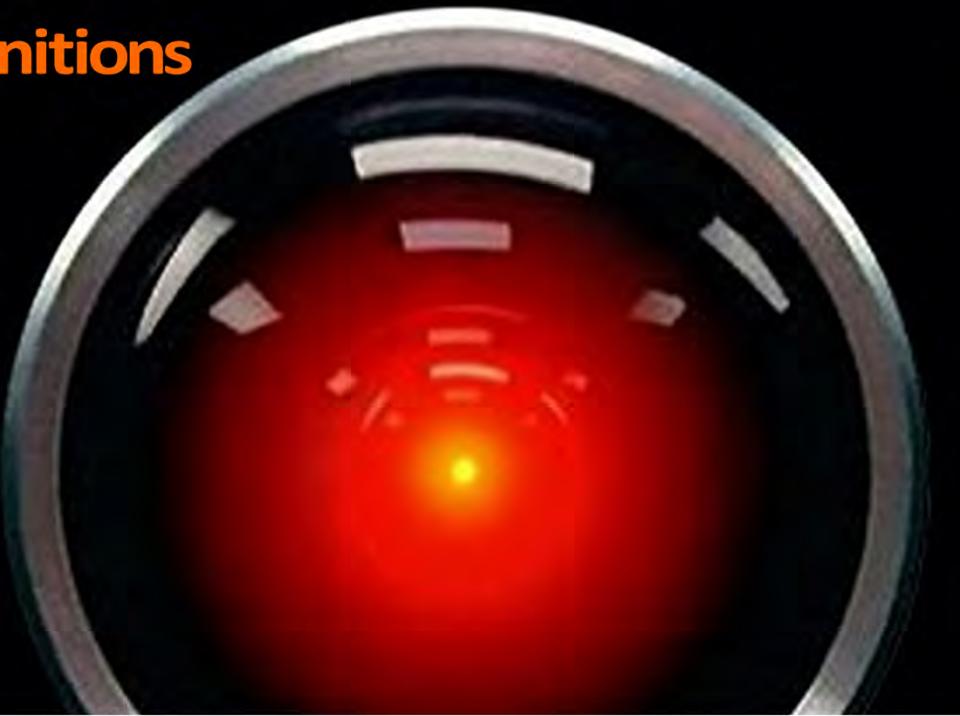
- **Prescriptions, formulas, one-size-fits-all solutions**
- **Perfect and/or mature best practices**
- **Rigid processes or step-by-step instructions**
- **Big overnight transformations**
- **Extended technical discussions or deep focus on any specific technology or tool**



NOTES:

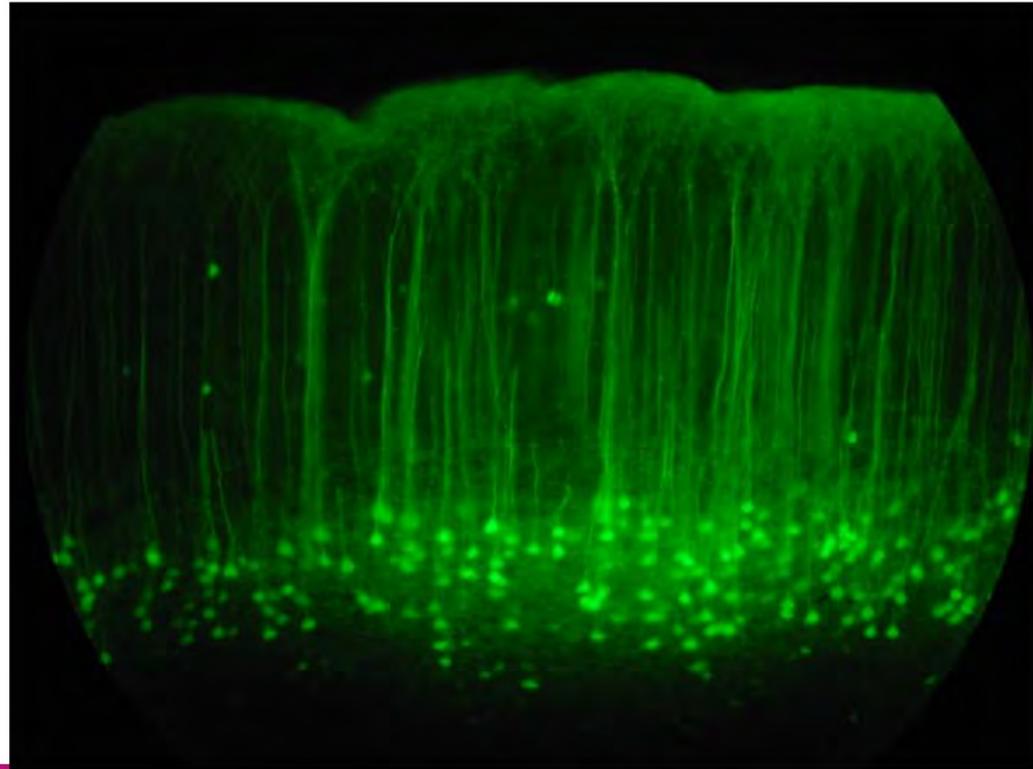
AI: Working Definitions

- AI
- Machine Learning
- Deep Learning
- Data Science
- Data Model
- Big Data



NOTES:

What is intelligence?



Pictured: 20X magnification of the somatosensory cortex of a mouse brain slice.

What are the characteristics of intelligence?

Beware deep metaphysical discussion. We are interested in practical application of machine intelligence. At present, the characteristics of our own intelligence which are able to be mimicked or improved upon are pretty straightforward, defined as:

- Ability to detect
- And act upon
- Patterns

ML is good at the detection, but needs to be given the values and goals that enable action.

AI: Working Definitions

- AI
- Machine Learning
- Deep Learning
- Data Science
- Data Model
- Big Data



NOTES:

Artificial Intelligence - Definition

An intelligent agent that perceives its environment and makes decisions to maximize chances of achieving its goal.

3 Types of AI:

- Artificial Narrow Intelligence
- Artificial General Intelligence (strong AGI)
- Artificial Super Intelligence



NOTES:

Artificial Intelligence - Definition

An intelligent agent that perceives its environment and makes decisions to maximize chances of achieving it's goal.

3 Types of AI:

- Artificial Narrow Intelligence
- Artificial General Intelligence(strong AGI)
- Artificial Super Intelligence

*Let's get these
out of the way
first*



NOTES:

AI Stages	Artificial Narrow Intelligence (ANI)	Artificial General Intelligence (AGI)	Artificial Super Intelligence (ASI)
	Execute specific focused tasks, without ability to self-expand functionality	Perform broad tasks, reason, and improve capabilities comparable to humans	Demonstrate intelligence beyond human capabilities
Timing	Today	About 2040?	Soon after AGI
Implications	Outperform humans in specific repetitive functions, such as driving, medical diagnosis and financial advice	Compete with humans across all endeavors, such as earning university degrees and convincing humans that it is human	Outperform humans, helping to achieve societal objectives or threatening human race



Source – acceleratingbiz.com/proof-point/future-artificial-intelligence/

Is AGI even possible?

This is an important question, but we won't spend much time on AGI in this session beyond the current discussion, because there's little practical application to be gained and AGI. Still, the trends in AI tie to potential developments in the future, and even if AGI is distant or not possible, it's the possibilities between then and now that practical value.

Kevin Hill, PhD, one of the contributors to this class, is one of those skeptics that thinks that maybe strong AGI isn't really that possible. To be 10x as smart as a human, he says, "Might require more computation that is available in the universe." Some good neuro studies show that humans are near the mathematically defined limit for information processing. The human brain runs on just 100 watts, whereas a single modern GPU takes 900 W, and 9 humans are way smarter than one GPU. Another big blocker to AGI is single-trial learning. Right now machines need thousands of examples to learn from.

BUT in repeated contests small advantages snowball, and there are huge advantages to be gained in processing different TYPES of data than humans are optimized for. And many scientists and prominent leaders in the AI and technology community have concerns about the future of AI and what it might mean for humanity.

Myth: Superintelligence by 2100 is inevitable		Fact: It may happen in decades, centuries or never: AI experts disagree & we simply don't know	
Myth: Superintelligence by 2100 is impossible			
Myth: Only Luddites worry about AI		Fact: Many top AI researchers are concerned	
Mythical worry: AI turning evil		Actual worry: AI turning competent, with goals misaligned with ours	
Mythical worry: AI turning conscious			



Source: Max Tegmark, “Life 3.0”

NOTES:

Myth: Robots are the main concern		Fact: Misaligned intelligence is the main concern: it needs no body, only an internet connection	
Myth: AI can't control humans		Fact: Intelligence enables control: we control tigers by being smarter	
Myth: Machines can't have goals		Fact: A heat-seeking missile has a goal	



Source: Max Tegmark, "Life 3.0"

NOTES:



Source: Max Tegmark, "Life 3.0"

NOTES:

AI Today...

Many organizations want to deploy AI but do not have enough, or consistent, or strategically structured enterprise data architecture to do so. Microsoft explained that many orgs simply don't have the historical data in place in order for AI algorithms to learn enough to be useful. For instance – one company wanted to use AI to predict future failures in their system, but they did not have data on past failures in order to train their AI solution.

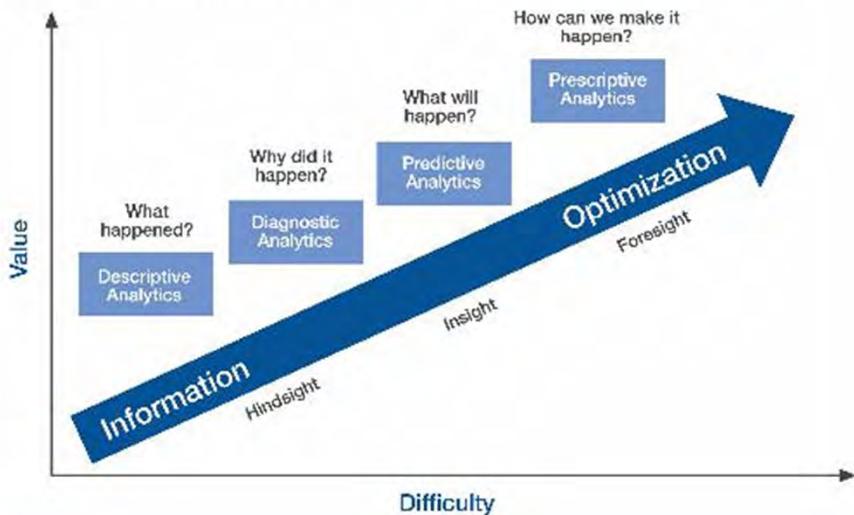


Source – David Kiron, MIT Sloan Management Review

NOTES:

Levels of Algorithmic Sophistication

Predictive Analytics



Source: Gartner

Most companies live at the bottom left. It is risky to skip a step without major new investments in people and systems.

AI Today...

AI leaders are:

- 12X more likely to understand the process of training algorithms
- 10X more likely to understand development costs
- These orgs already have mature analytics practices and already understand how to manage large, unified data engineering practices and large-scale analytic



Source – David Kiron, MIT Sloan Management Review

NOTES:

AI Today...

- AI solutions are currently incomplete. They are like very young children who need to learn in order to be productive and useful.
- If orgs do not understand how to train their AI solutions, AND the organized stores of data upon which to build: these are big hurdles to adoption.
- In one well-known example, Wells Fargo cannot exploit many potential AI applications because the data stores in their organization are so fragmented. So the underlying substrate of data is not in a state that is ready to harness AI.



Source – David Kiron, MIT Sloan Management Review

NOTES:

AI: Working Definitions

- AI
- **Machine Learning**
- Deep Learning
- Data Science
- Data Model
- Big Data



NOTES:

Machine Learning – Definition

Machine learning is a subfield of artificial intelligence. Its goal is to enable computers to learn on their own. A machine's learning algorithm enables it to identify patterns in observed data, build models that explain the world, and predict things without having explicit pre-programmed rules and models.



NOTES:

Machine Learning - Types

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning



NOTES:

Example: Supervised Learning for Classification

With supervision, machine models can be effective at learning how to categorize things like:

- Is this email spam or not
- Is this a fraudulent charge
- Is this photo a dog or cat
- Is this fake news or not



NOTES:

What problem are you trying to solve? (think small)

- No such thing as machine learning fairy dust
- Ask small questions first
- Big problems are a mass of small problems, you can make more progress by tackling small ones first
- Chances of triggering change on a small problem are much greater



NOTES:

Machine Learning - Data

- **Importance of Data**
- **Structured vs. Unstructured**
- **Cleaning Data**
- **Training, Validation, and Test**



NOTES:

Machine Learning – Structured Data

Pre-defined and machine readable. Usually has a relational data model. **Examples of structured data:**

- Library catalogues
- Census records
- Meta-data
- Databases
- XML



NOTES:

Machine Learning – Unstructured Data

There's no predefined model. It's often text. But always some structure.

Examples of data sources:

- Email
- Newspapers
- Health Records
- Books
- PDF Documents



NOTES:

AI: Working Definitions

- AI
- Machine Learning
- **Deep Learning**
- Data Science
- Data Model
- Big Data



NOTES:

Deep Learning – Definition

A subset of machine learning composed of algorithms that permit software to train itself to perform tasks, like speech and image recognition, by exposing multi-layered artificial neural networks to vast amounts of data. The levels in these learned statistical models correspond to different levels of concepts, where higher level concepts are defined from lower level ones, and the same lower level concepts can help to define many higher-level concepts.



NOTES:

Neural Networks – How they learn

How do they learn?

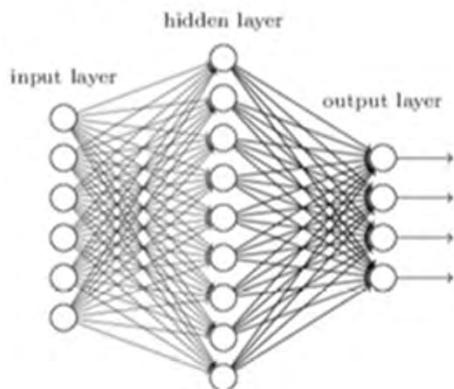
- Setup Network (defined schedule) up so that it takes in inputs and produces output
- If the output does not match desired output network creates an error signal
- Network will pass the error signal backwards through the network. Learning occurs by adjusting the weights to match the inputs and outputs



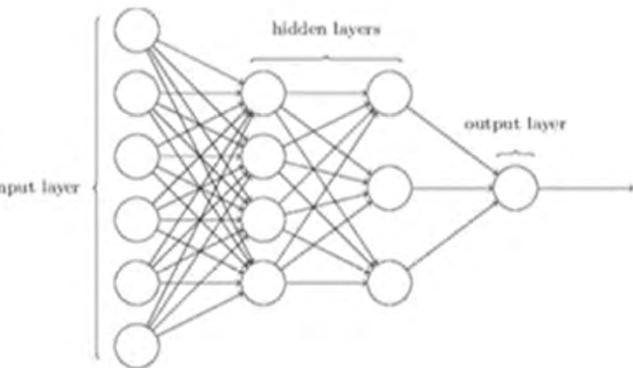
NOTES:

Neural Network – Deep Learning

- Not Deep



Deep Learning



Source: Homstock, Peter. Harvard University 2017

The first OCR system was adopted by the US post office in 1965.

- Powered by an early neural network
- Among original success stories of NN

Currently:

- 55 million hand-addressed s-mail each day ◦ 83% are automatically processed
- Accuracy > 98%
- Saved \$100 million during first year (1997)

Neural Networks – What are they good at?

- Automatic speech translation
 - Recognize traffic signs > people
 - Google maps read all building #s
 - Baidu: image understanding
 - Sentiment analysis
 - Read Chinese ~ native levels
 - Write sentences as image captions
 - Medicine: cancer research & insights
 - Medical decisions & diagnostics in 15 minutes



NOTES:

Neural Networks

- **Association between human brains and neural networks is inspired, but weak**
 - **Neural nets seem particularly good at human-like tasks such as:**
 - Recognizing images
 - Recognizing speech
 - Recent work over past few years suggest they are good at predicting
-



Neural networks: key facts

- NN are non-linear as a whole
- Can have linear processing units
- Still trying to find a function f to map input x to output y with minimal loss
- Able to learn from examples
- Deep learning is great at finding complex relationships
- Rules for manipulation of input are learned not preloaded
- The network architecture is predefined but changes its architecture based on the outputs that the network creates
- Fault tolerant

Brain vs. Computer?

Connections

100 billion neurons in brain
Core i7 has 1.4 billion transistors
Iphone8 has 4.3 billion transistors

Neuron switching time 0.001 sec

Much slower than transistors (1GHz)

Power:

Brain 20W
Core i7-920 120W

Connections

10K-100K connections per neuron
Few for computer -> bus architecture

AI – Challenges of “Black Box” vs. “White Box”



- Sophisticated AI applications can operate in ways that are opaque to humans, even designers of the technology. For example, deep learning that has been used for self-driving cars is not understood.
- This presents problems for adoption, trust, buy-in, compliance, troubleshooting, and “explainability”



Black box AI and the “explainability” problem

The problem of “black box” AI is an important consideration in today’s landscape. Anyone working with or attempting to understand AI should be familiar with this debate.

“White box” AI can be understood as an AI application which can be traced and explained throughout every step of its operation. The problem is that many of the most interesting applications of AI rely on machine operations that are set in motion by human creators who do not then fully understand how the AI arrives at its outputs.

Even if an AI application is not truly a “black box” it can present challenges for its stakeholders if it is not easy to explain. Consider a basic example: At its simplest, a deep learning algorithm or neural network may be only a series of recursive linear regression analyses. Although linear regression is straightforward and well-understood, nesting many linear regression algorithms together can quickly create an application which produces difficult-to-explain outputs or which is simply difficult to explain to a non-technical audience. Therefore, it may be challenging to integrate into audit requirements, regulator, or senior decision makers. **Here are a few common challenges presented by the question of black box AI:**

- Users may be reluctant to adopt the application if how it works is not well understood
- Regulatory, audit and compliance requirements may be impossible to meet if an application’s behavior can’t be thoroughly explained
- Improvements and troubleshooting can be difficult if an AI is a black box
- Trust and support can be difficult to win within an organization attempting to deploy AI when the way it functions is impossible to explain

From MIT Technology Review:

Already, mathematical models are being used to help determine who makes parole, who's approved for a loan, and who gets hired for a job. If you could get access to these mathematical models, it would be possible to understand their reasoning. But banks, the military, employers, and others are now turning their attention to more complex machine-learning approaches that could make automated decision-making altogether inscrutable. Deep learning, the most common of these approaches, represents a fundamentally different way to program computers. "It is a problem that is already relevant, and it's going to be much more relevant in the future," says Tommi Jaakkola, a professor at MIT who works on applications of machine learning. "Whether it's an investment decision, a medical decision, or maybe a military decision, you don't want to just rely on a 'black box' method."

There's already an argument that being able to interrogate an AI system about how it reached its conclusions is a fundamental legal right. Starting in the summer of 2018, the European Union may require that companies be able to give users an explanation for decisions that automated systems reach. This might be impossible, even for systems that seem relatively simple on the surface, such as the apps and websites that use deep learning to serve ads or recommend songs. The computers that run those services have programmed themselves, and they have done it in ways we cannot understand. Even the engineers who build these apps cannot fully explain their behavior.

This raises mind-boggling questions. As the technology advances, we might soon cross some threshold beyond which using AI requires a leap of faith. Sure, we humans can't always truly explain our thought processes either—but we find ways to intuitively trust and gauge people. Will that also be possible with machines that think and make decisions differently from the way a human would? We've never before built machines that operate in ways their creators don't understand. How well can we expect to communicate—and get along with—intelligent machines that could be unpredictable and inscrutable?"

For further reading:

- <https://www.forbes.com/sites/jasonbloomberg/2018/09/16/dont-trust-artificial-intelligence-time-to-open-the-ai-black-box/#7930a8513b4a>
- <https://www.technologyreview.com/s/604087/the-dark-secret-at-the-heart-of-ai/>

AI: Working Definitions

- AI
- Machine Learning
- Deep Learning
- **Data Science**
- Data Model
- Big Data



NOTES:

Data Science - Definition

Data science is a multi-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data. Data science is the same concept as data mining and big data: "use the most powerful hardware, the most powerful programming systems, and the most efficient algorithms to solve problems".

Data science is a "concept to unify statistics, data analysis, machine learning and their related methods" in order to "understand and analyze actual phenomena" with data. It employs techniques and theories drawn from many fields within the context of mathematics, statistics, information science, and computer science.



In 2012, when Harvard Business Review called it "The Sexiest Job of the 21st Century", the term "data science" became a buzzword. It became conflated with:

- Business analytics
- Business intelligence
- Predictive modeling
- Statistics (*Nate Silver referred to data science as a "sexed up term for statistics."*)

While many university programs now offer a data science degree, there exists no consensus on a definition or suitable curriculum contents.

Many data-science and big-data projects fail to deliver useful results, often as a result of poor management and utilization of resources.

Data Science - Definition

In 2012, when Harvard Business Review called it "The Sexiest Job of the 21st Century", the term "data science" became a buzzword. It became conflated with:

- **Business analytics**
- **Business intelligence**
- **Predictive modeling**
- **Statistics** (*Nate Silver referred to data science as a "sexed up term for statistics."*)

While many university programs now offer a data science degree, there exists no consensus on a definition or suitable curriculum contents.

Many data-science and big-data projects fail to deliver useful results, often as a result of poor management and utilization of resources.



NOTES:

Who is a real data scientist?

Data scientists are highly qualified and can be hard to hire and retain. They are the key professional role behind AI and advanced analytics.



Pictured: DJ Patil, Deputy Chief Technology Officer for Data Policy and Chief Data Scientist in the Office of Science and Technology Policy, The White House. 2016

Programming and databases

- Computer science fundamentals
- Scripting language i.e. Python
- Statistical computing package i.e. R
- Databases: SQL and NoSQL
- Relational algebra
- Parallel databases and parallel query processing
- MapReduce concepts
- Hadoop and Hive/Pig, Spark
- Custom reducers
- Experience with XaaS like AWS

Communication and visualization

- Able to engage with senior management
- Storytelling skills
- Translate data-driven insights into decisions and actions
- Visual design skills
- R packages like ggplot or lattice
- Knowledge of data visualization tools like Flare, D3.js, or Tableau

Math and statistics

- Machine learning
- Statistical modeling
- Experiment design
- Bayesian inference
- Supervised and unsupervised learning
- Optimization: gradient descent and variants

Domain knowledge and soft skills

- Passionate about the business
- Curious about data
- Influence without authority
- Hacker mindset
- Problem solver
- Strategic, proactive, creative, innovative and collaborative

AI: Working Definitions

- AI
- Machine Learning
- Deep Learning
- Data Science
- **Data Model**
- Big Data



NOTES:

Data Model

For the purpose of discussing AI initiatives, a **data model** is the collection of algorithms and digital application behavior which operationalizes an AI application. Data models are mathematical and logical constructs used for:

- **Machine learning**
- **Data analysis**
- **Advanced analytics**
- **AI applications**



Data Models – Different definitions for different contexts

For the purposes of this class, we relate **data models** to their roles as components in AI applications. It's also useful to understand more conventional definitions of a data model, as various interpretations of the term will often relate to and overlap with how AI initiatives fit into the overall data strategies and ecosystem in the organization.

From Wikipedia:

A **data model** is an abstract model that organizes elements of data and standardizes how they relate to one another and to properties of the real world entities. For instance, a data model may specify that the data element representing a car be composed of a number of other elements which, in turn, represent the color and size of the car and define its owner.

The term **data model** is used in two distinct but closely related senses. Sometimes it refers to an abstract formalization of the objects and relationships found in a particular application domain, for example the customers, products, and orders found in a manufacturing organization. At other times it refers to a set of concepts used in defining such formalizations: for example concepts such as entities, attributes, relations, or tables. So the "data model" of a banking application may be defined using the entity-relationship "data model".

Data models: Still created by a human.

```
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
```



NOTES:

AI: Working Definitions

- AI
- Machine Learning
- Deep Learning
- Data Science
- Data Model
- **Big Data – We will spend some time on it in the next section of class**



Big data is a separate but closely related area of practice for any organization that wishes to implement AI solutions.

We will spend the next section discussing a few of the most important fundamentals of big data and how they relate to AI.



Part 2: Big Data



Big data is a closely related area of practice, and in many ways the precursor to today's AI boom. Many organizations have not yet figured out foundational big data practices. Until they do, their readiness for an AI initiative is limited.

- The amount of data generated today is growing exponentially.
- Both the volume and type of much of this data—the type we'll refer to as “big data”—is not suitable for storing in relational databases, yet has significant business value
- Engineering new systems or augmenting existing systems often requires the use of this data and so big data skills are mandatory for data engineers practicing today
- Stores of big data are usually a prerequisite to most interesting AI applications. Machine models need large datasets to learn, and the best AI applications are based on good machine learning
- Big datasets require preparation of the dataset and effective architecture of technical dependencies – for example, storage or streaming infrastructure
- Once stored the data must be able to be processed in a timely manner. Currently, the most available technique for doing so is massively parallel computing on low cost hardware
- The ability to interactively query Big Data sets to determine patterns and trends is a very powerful facility.
- Presenting the results of Big Data processing and querying in a coherent manner is a technical challenge and requires planning and commitment for business benefits to be realized

What is Big Data?

The 3 V's

- **Volume**
Large Quantities (think gigabytes, terabytes of information daily)
 - **Velocity**
Speed of generation and/or need to be processed quickly
 - **Variety**
Data might be structured, unstructured, varying sources
 - **And some additional considerations...**
-



Big data can be described by the following characteristics:

- **Volume** – The quantity of data that is generated is very important in this context. It is the size of the data which determines the value and potential of the data under consideration and whether it can actually be considered Big Data or not. The name 'Big Data' itself contains a term which is related to size and hence the characteristic.
- **Variety** - The next aspect of Big Data is its variety. This means that the category to which Big Data belongs to is also a very essential fact that needs to be known by the data analysts. This helps the people, who are closely analyzing the data and are associated with it, to effectively use the data to their advantage and thus upholding the importance of the Big Data.
- **Velocity** - The term 'velocity' in the context refers to the speed of generation of data or how fast the data is generated and processed to meet the demands and the challenges which lie ahead in the path of growth and development.

Other Considerations

- **Variability**
Is your data consistent? Does it have holes in it?
 - **Veracity**
How accurate are your data sets?
 - **Complexity**
The more complex your data, the “bigger” it is. Do you have multiple sources that must be linked, connected, or correlated in order to grasp its meaning?
-



- **Variability** - This is a factor which can be a problem for those who analyze the data. This refers to the inconsistency which can be shown by the data at times, thus hampering the process of being able to handle and manage the data effectively.
- **Veracity** - The quality of the data being captured can vary greatly. Accuracy of analysis depends on the veracity of the source data.
- **Complexity** - Data management can become a very complex process, especially when large volumes of data come from multiple sources. These data need to be linked, connected and correlated in order to be able to grasp the information that is supposed to be conveyed by these data. This situation, is therefore, termed as the ‘complexity’ of Big Data.

New Data Sources

- **Logfiles**
- **Sensors / RFID**
- **Location / Geospatial**
- **Video / Recognition**
- **Social / Natural Language**

Power Sensors



Driving Sensors

Helmet Sensors



NOTES:

Sensors, Machine-generated data, and IoT



NOTES:

Machine vision



NOTES:

Machine vision



NOTES:

The big data prerequisite

15 years ago...

- Big, exotic supercomputers
(Cray J90 pictured)
- 2.5 million bucks
- Big, important projects
- Wait lists...PhDs



NOTES:

The “Condor Cluster” (U.S. Air Force, 2011)



The Condor Cluster and others like it were early, highly effective computing clusters built from Playstation 3 consoles.

- A single PS3, for a time, was able to perform like a cluster of 30 PCs at the price of only one. (*Marc Stevens, Arjen K. Lenstra, and Benne de Weger, 2007*)
- In Summer 2007, Gaurav Khanna, a professor in the Physics Department of the [University of Massachusetts Dartmouth](#) independently built a message-passing based cluster using 8 PS3s running Fedora Linux, and performed astrophysical simulations of large supermassive black holes capturing smaller compact objects.
- The Condor cluster operated 168 separate graphical processing units and 84 coordinating servers in a parallel array capable of performing 500 trillion [floating-point operations per second](#) (500 TFLOPS). As built the Condor Cluster was the 33rd largest supercomputer in the world, used to analyse high definition satellite imagery.

Compare the cost:

- Cray J90 = 2.5 million bucks
- Playstation 3 = 300 bucks x 1700 = ~ 500k bucks

Virtualization

2006

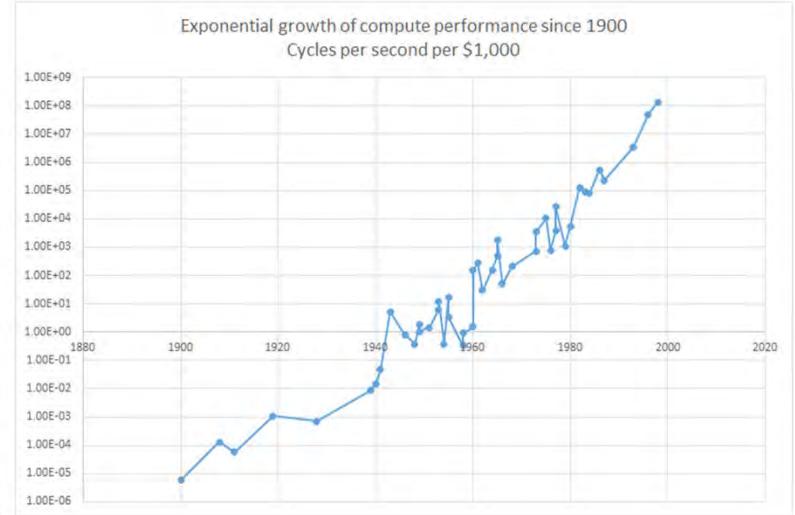
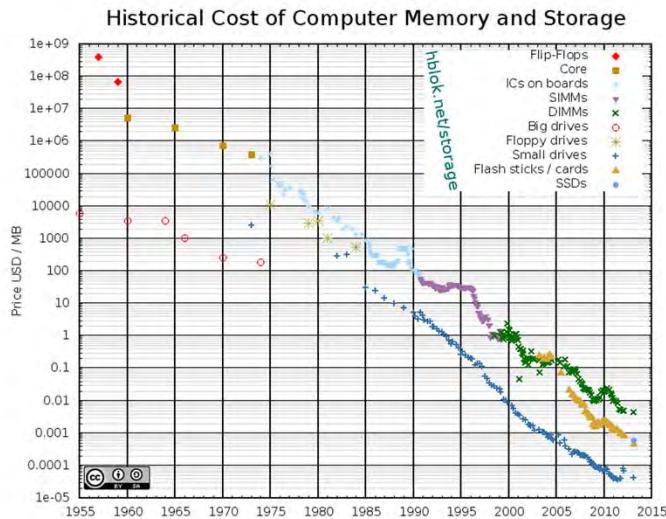


Today



NOTES:

Compute capability, costs, and performance



The class should take a few minutes to review Moore's law and consider technical implications. Discuss for 5-10 minutes. Questions to consider:

1. If the trends shown above can be relied upon to continue, what are the technical implications?
2. What types of applications could be enabled by twice the computing power at half the cost?
3. What kind of timeframe might one reasonably expect to associate with the ideas the class suggested in #2?

Applications

What can we do with it?

- Data warehousing
- Business intelligence
- Analytics
- Predictive statistics
- Data science



NOTES:

In Summary

- **There are many, many sources of data**
- **The speed of data creation is increasing**
- **We need to process it**
- **We need to make sense of it**
- **We can use it to establish advantage**



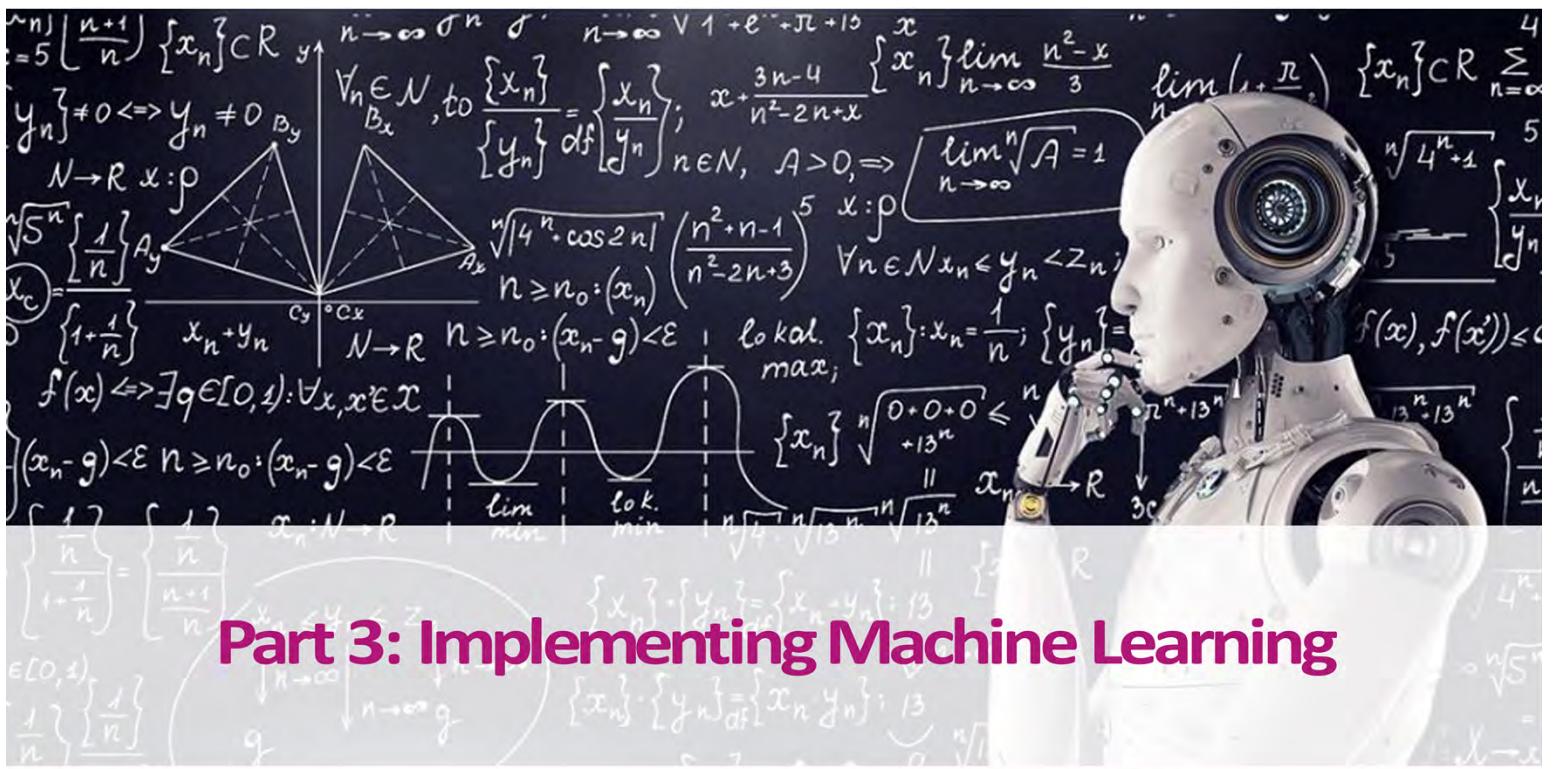
NOTES:

Big data and AI: the main takeaway

- In general, effective AI is usually powered by machine learning
- Effective machine learning requires existing datasets to learn
- While there are some interesting things AI can do with local machines and small datasets, in general **the most exciting applications rely on very large datasets**



NOTES:



Part 3: Implementing Machine Learning



Image CC license: <https://www.flickr.com/photos/mikemacmarketing/30212411048>

Pillars of a Successful AI Team

- **Algorithms**
 - **Modeling**
 - **Business Case & Business Integration**
 - **Domain Expertise**
 - **Data**
 - **Automation**
 - **Scalability**
-



You need all of these to be able to expect success.

As we will continue to see throughout this session, AI is about more than just data, algorithms, or computing. AI is also about integrating advanced technology into your already-complex organization. Teams, human factors, and business processes are just as critical to a successful AI solution as the technology itself.

Cross Industry Standard Process for Data Mining (CRISP-DM)



Of course, you need a way to sequence and plan across all these pillars, which is where business process models such as CRISP-DM come in. First published 1999, but still a great and simple model to follow. IBM is the current champion of CRISP-DM, and the more complicated ASUM-DM.

Models like CRISP-DM are independent of project management tools. Think of it as a way to sequence multiple sub-projects.

State of the Tools



pythonTM



We will discuss tools more in section 7.

NOTES:

Machine Learning: Two main types

- Supervised
- Unsupervised



NOTES:

Supervised Learning – Components

- Create a function f that can map input X to an output Y with minimal loss.
- Features(input or X) – independent inputs of some specific variable
- Labels(output or Y) – dependent variables you are predicting from a feature or set of features.



NOTES:

Supervised Learning – Types

- **Linear Regression** – predicting a continuous value from a feature or set of features. Predict income based on years of education.
- **Classification** – assign a discrete value from a feature or set of features. Predict whether income is higher or lower than 50k based on years of education.



NOTES:

Unsupervised Learning

- Initially no labels, only features
- AI finds the underlying structure of data
- AI groups data into new groups



NOTES:

Unsupervised Learning - Types

- Clustering
- Dimensionality Reduction

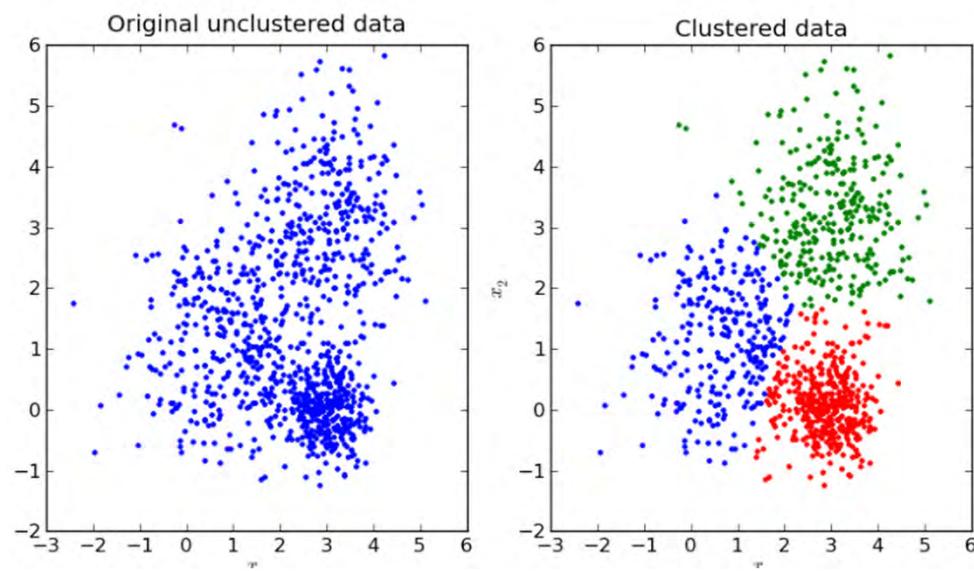


NOTES:

Unsupervised Learning - Clustering

K-Means Clustering

- Common type of clustering
- Create data points such that points in different clusters are dissimilar while points within a cluster are similar.
- Cluster data into K groups
- Larger K creates smaller groups with more granularity
- Lower K creates larger groups with no granularity



Unsupervised Learning - Clustering

Kmeans Steps

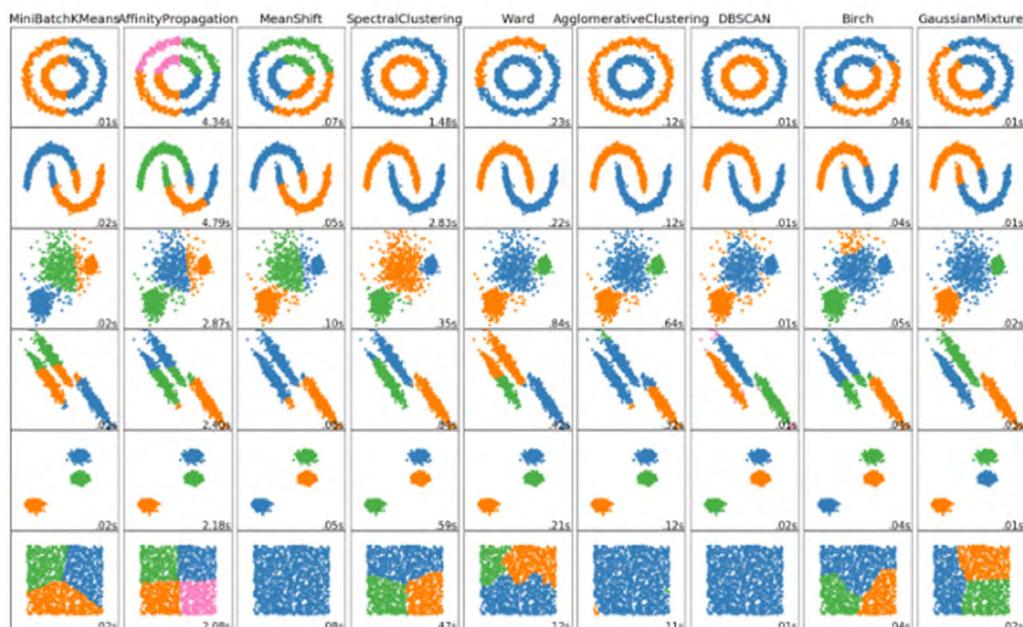
- Step 1: K centroids are randomly selected
- Step 2: Assign data point to closest centroid's cluster (most commonly Euclidean distance)
- Step 3: Distance between the points and the closest centroid is minimized with each iterations of the algorithm by moving the centroid.
- Step 4: Move the centroids to the center of clusters
- Step 5: Repeat steps 3 and 4 until centroid stops moving a lot. This is called converging.

Demos:

- <https://www.youtube.com/watch?v=jxqvBeJCLPA>
- <https://www.youtube.com/watch?v=BVFG7fd1H30>



There are dozens of types. In fact, most real world data does not fit into a simple model.



Data models:

Let's walk through the most common types of models.

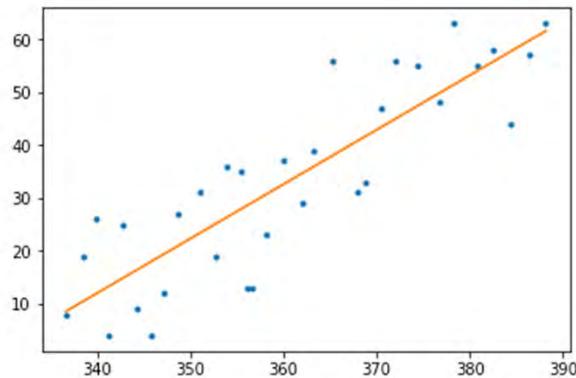
```
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
```



NOTES:

Linear Regression

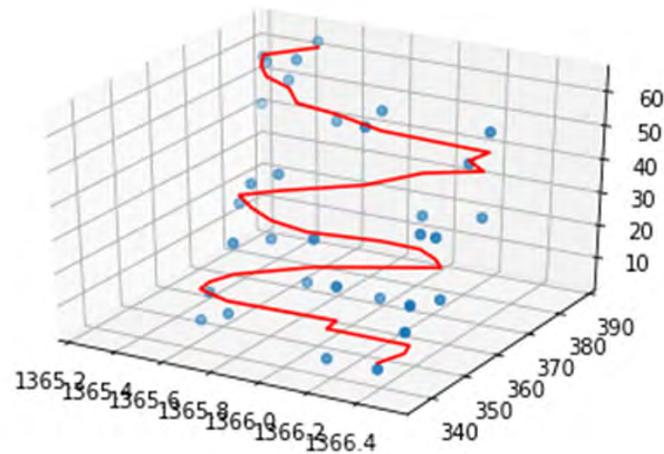
	Year	CO2	Temp
0	1979	336.67	8
1	1980	338.57	19
2	1981	339.92	26
3	1982	341.30	4
4	1983	342.71	25



Linear regression is probably familiar to everyone, and is one of the most straightforward types of data analysis.

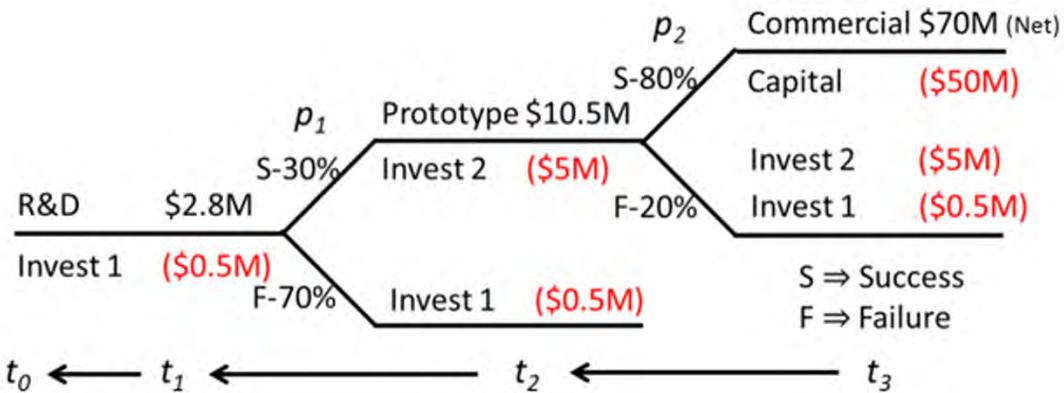
Multiple Linear Regression

	Temp	CO2	Solar
0	8	336.67	1366.43
1	19	338.57	1366.51
2	26	339.92	1366.51
3	4	341.30	1366.16
4	25	342.71	1366.18



No need to go into details, just the fact that we generalize the simple model to multiple inputs.

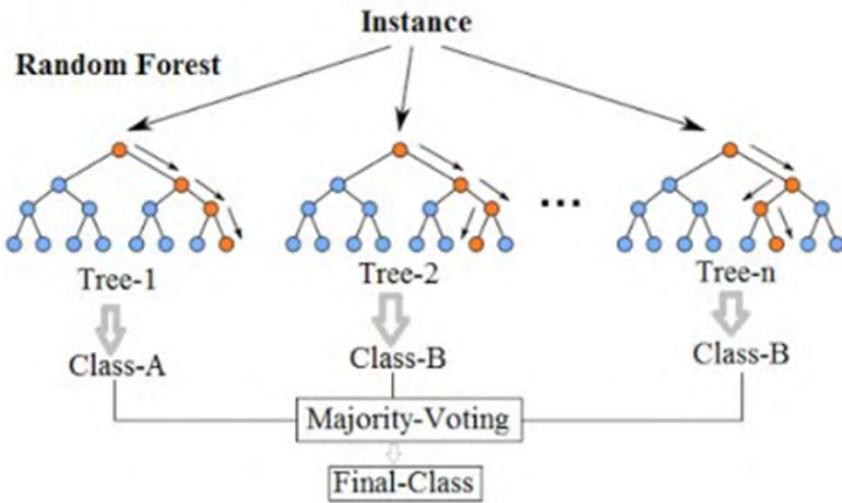
Decision trees



This conceptual tool should be familiar to anyone who studies business process models

NOTES:

Random Forest



AKA “Decision Forest”. There are often a few fast and accurate models for a wide variety of business needs. Often the default ‘first pass’ model

An important consideration here is to contrast with decision trees. A single tree is very like a business process, but here we have many competing trees, all trained on different subsets of the data. A 2nd level system pools the outputs of the 1st level models. Very common to see these “ensemble” structures in production.

Naïve Bayes

- **Built on Bayes Theorem**
- **Assumes observations are independent (naïve)**
- **Scales well to very large datasets**
- **You also get “generative” models for free**

$$p(C_k|x) = \frac{p(C_k)p(x|C_k)}{p(x)}$$

$$p(C_k|x) \propto p(C_k) \prod_{i=1}^n p(x_i|C_k)$$

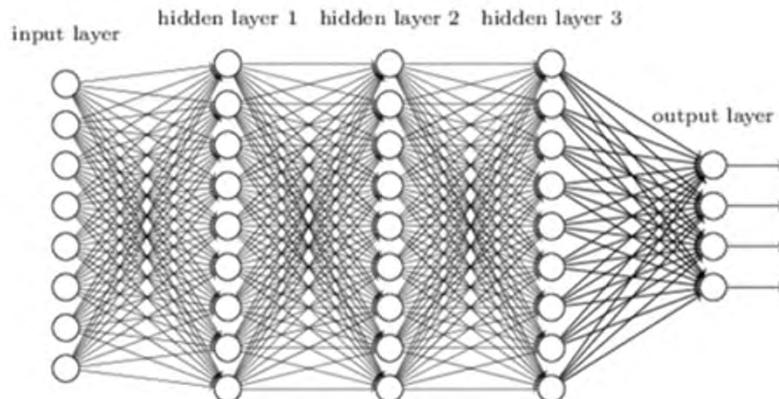


One of the models that is hard to delve deep into without some math background, but useful to know about because of the idea of generative models.

Generative models mean that you can create semi-realistic fake data given an assumed class output. This can be very useful for a lot of indirect business challenges: eg how do you test a data system without increasing exposure of sensitive customer data (GDPR / HIPAA compliance)

Neural Network – Deep Learning

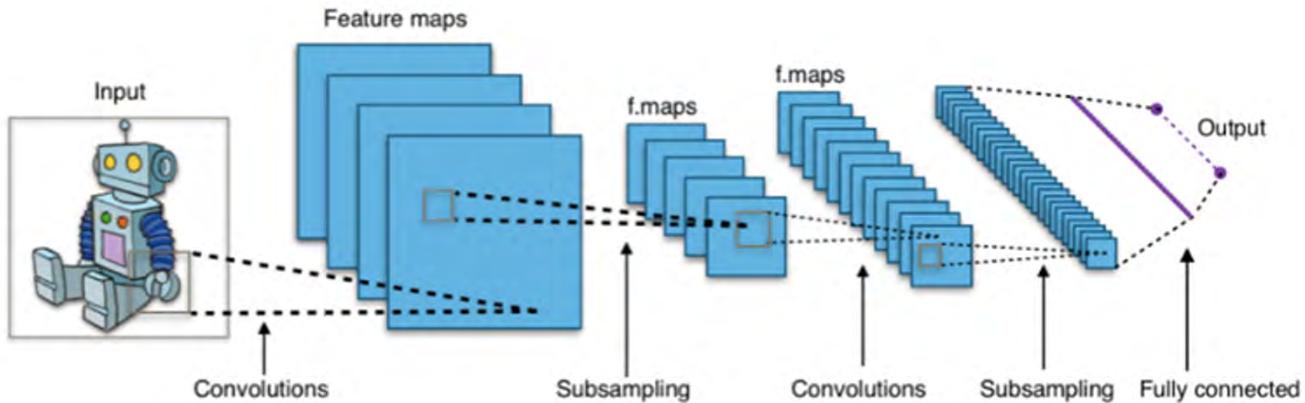
Deep NN has
More than one
hidden layer



Homstock, Peter Harvard University 2017

Example of a fully-connected network: All nodes are linked to all nodes of the next layer. A fully connected network can, in theory, represent any possible model. But you might need millions of layers and millions of years to train such a generic model. Many of the connections would be at or near 0, representing wasted computation.

Convolutional Neural Networks

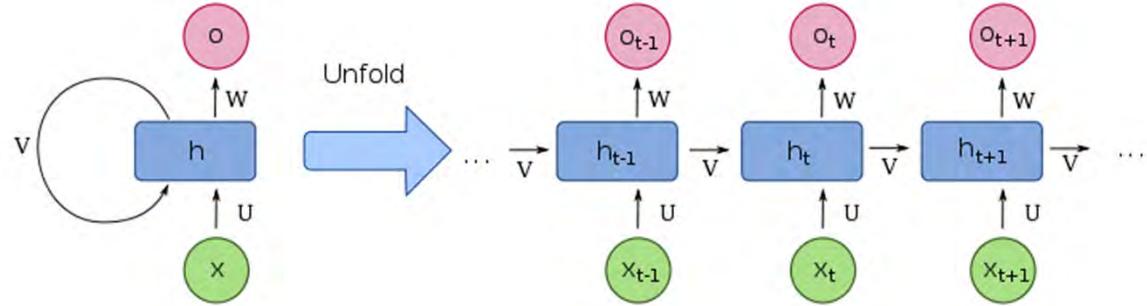


Deep learning networks should be thought of as a design space. Find a simpler, non-fully connected model that solves your specific problem, not any possible problem. These design skills are a subset of general data science skills.

The first clear example of this principle was Convolutional Neural Networks (CNN) in 1989. It was specialized for processing images, and based on the way the visual system of the brain processes images. CNNs take advantage of the real-world fact that the pixels you need to understand an image are spatially arranged.

They usually have a fully connected layer near the end of the network to map image features into a specific business problem space.

Recurrent Neural Network



Another example of a Neural Network design component. Recurrent Neural Networks (RNNs) take not just one input, but a sequence of inputs. The output of the previous input is used to help process the subsequent inputs. Often used in time-series analysis and language processing. Eg “I went fishing on the bank” where the previous word “fishing” helps distinguish between “a bank of a river” and “a bank that stores money”.

Model Overfitting

- **Extreme overfitting = memorizing every data in the training set**
- **High accuracy on training data**
- **Low accuracy on new data!**



More accuracy doesn't always mean better.

NOTES:

Model Hyper Parameters

- **Values or settings that impact how the model trains, and the potential for overfitting**
- **Need to train model many times with different hyperparameters to find ideal setup**



NOTES:

Train, Validate, Test

- **Split data into 3 groups randomly**
 - True random can be hard
 - **80/10/10 is a good rule-of-thumb**
 - **Validate set makes sure you don't overfit and sets hyperparameters**
 - **Test set is a final sanity check on your hyperparameters**



NOTES:

Accelerating Training

- **The power of GPUs**
 - XGBoost for trees
 - Deep learning uses the same linear algebra as video games
- **Bootstrap**
 - Like decision tree -> random forest
 - Faster training but slower inference



NOTES:

Encoding Domain Expertise

- **Feature engineering - simple algorithms to combine two or more datapoints**
- **Great way to build collaboration between business units**
- **Often related to KPIs**
- **Can also uncover hidden business assumptions**



Regarding hidden business assumptions, one of our favorite examples is a sales pipeline prediction algorithm that identified client location as a predictive measure. It turns out that sales people only put in an address for a potential client when they booked an in-person meeting!

Most businesses and organizations are dense with these types of hidden assumptions and missed connections between departments and teams. This is critical to understand because it is most often these types of human and organizational factors that will sabotage success with an AI initiative – not the difficulty of mathematics required or complexity of the required technology. Both are formidable requirements, so it is all the more important not to let our own organizational dysfunction be the downfall.

Cross-functional collaboration between siloes and business units is one of the most important prerequisites to success with any kind of AI or advanced analytics practice. Therefore, the exercise of forming appropriate teams and examining flows of information pave the way not only for effective AI, but probably other improvements also.

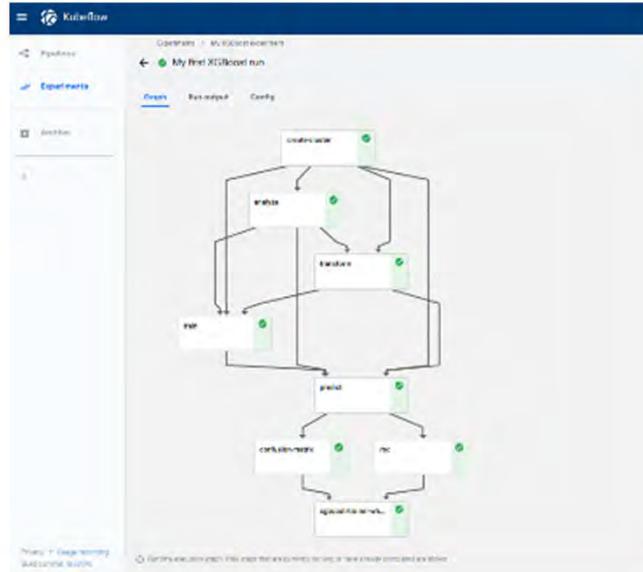
Machine Learning Model Management

- **Every day you get new data and new ideas**
- **Models should be constantly retrained**
- **Need to record which model made which prediction**
- **Might have multiple models in production (A/B testing)**



NOTES:

Automating Data Science



A platform to help you manage the data science workflow, and experiment with different algorithms.

Kubeflow is the new fancy kid on the block. Based on Google's internal Machine Learning automation pipelines, and open source. Excellent choice if you already have Kubernetes running in your company.

Other options:

- Airflow (Open source from AirBnB, very flexible, can do core ETL and ML pipelines, but with flexibility comes a lot of manual setup)
- Tibco's Product Line (Spotfire etc. The whole line has python, but built to use R powered by Tibco's Enterprise Runtime for R)
- Alteryx
- Google Cloud – AutoML
- AWS – SageMaker
- Azure – Machine Learning Services (yeah, not a very creative name)



Tool discussion: Tensorflow



You can't have a conversation about AI and machine learning without considering a tool like Tensorflow. As a very capable and accessible machine learning technology, understanding how this tool is used can give you a better mindset for how AI can have practical application.

Before we continue, let's discuss:

- Has anyone NOT heard of Tensorflow?
- Has anyone encountered any scenarios in which it was used? If so, please share.

Tensorflow

Tensorflow offers AI capability for:

- **Classification**
 - **Perception**
 - **Understanding**
 - **Discovering**
 - **Prediction**
 - **Creation**
-



Main Use Cases of TensorFlow

Voice and sound Recognition - With the right data, neural networks are capable of understanding audio signals. These can be:

- Voice recognition – For IoT, automotive, security and UX/UI features
- Voice search – numerous telecommunication application
- Sentiment analysis – often associated with CRM
- Flaw detection (i.e. mechanical noise, machinery) – for automotive and aviation uses

Understanding language is another use case for voice recognition. Speech-to-text can be used to determine samples of sound in larger audio files and transcribe the spoken word as text.

Sound-based applications also can be used in customer service. With the right design, Tensorflow algorithms might augment customer service agents, routing customers to the information they need faster than the agents.

Text Based Applications

Text-based applications are used for sentimental analysis on top of CRMs or social media. Threat detection (in social media or government intelligence) and fraud detection (insurance and finance) are other popular applications.

Language Detection

Language detection is a popular use of text based applications. Google Translate, which supports over 100 languages is now well known. Evolved versions can be used for applications like translating legal jargon into plain language.

Text Summarization

Google found out that for shorter texts, summarization can be used to produce headlines for news articles with a technique called sequence-to-sequence learning.

Image Recognition

Facial recognition, image search, motion detection, machine vision and photo clustering can be used also in automotive, aviation and healthcare industries. Image recognition seeks to recognize and identify people and objects in images as well as understanding the content and context.

TensorFlow object recognition algorithms classify and identify arbitrary objects within larger images. This is usually used in engineering applications to identify shapes for modeling purposes (3D extrapolation from 2D images) and by social networks for photo tagging. For instance, by analyzing thousands of photos of trees, the technology can learn to identify a tree it has never seen before.

Image recognition is starting to expand in healthcare, where TensorFlow algorithms can process more information and spot more patterns than their human counterparts. As we've already seen in this course, computers are now able to review scans and spot more illnesses than humans.

4. Time Series

TensorFlow time series algorithms are used for analyzing time series data in order to extract meaningful statistics. They allow forecasting non-specific time periods in addition to generate alternative versions of the time series.

The most common use case for time series is **recommendation**. You've probably used Netflix, which analyzes customer activity, compares it to millions of other users, and determines what the customer might watch. Recommendations are getting even smarter, for example, they can now suggest gifts that your family members might like.

The other uses of TensorFlow Time Series algorithms are useful in finance, accounting, government, security and IoT with risk detections, predictive analysis and enterprise resource planning.

5. Video Detection

TensorFlow neural networks also work on video data. Applications include image detection, real-time threat detection in security, and UX/UI fields. Recently, universities are working on large-scale =video classification datasets like YouTube-8M aiming to accelerate research on large-scale video understanding, representation learning, noisy data modeling, transfer learning, and domain adaptation approaches for video.

Tensorflow – who can use it?

The Google developer community provides these expectations:

- **Mastery of intro-level algebra.**
 - **Proficiency in programming basics, and some experience coding in Python.** Although Tensorflow now supports many environments, in general one should feel comfortable reading and writing Python code that contains basic programming constructs.
-

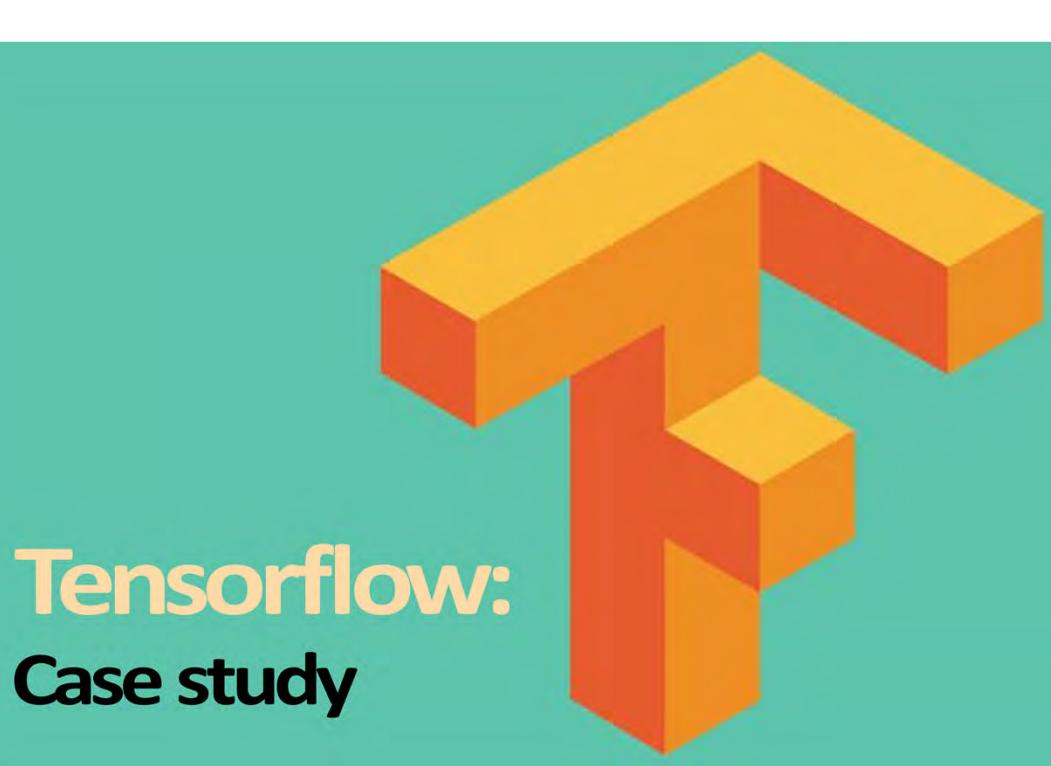


Source: developers.google.com

Tensorflow: how has it changed the accessibility of AI?

Tensorflow and tools like it may seem intimidating to a non technical audience, because there are still technical prerequisites required to use them. However, it should also be understood that Tensorflow is a tool available to anyone which removed technical requirements which – prior to its creation – were an order of magnitude more difficult in terms of technical expertise. While proficiency with Tensorflow depends on some higher math and application engineering skills, before it was available machine learning applications in production required much more advanced mathematics coupled with hand-built computing and application capability. *From the Google developer community, prerequisites to getting started with Tensorflow*

- **Mastery of intro-level algebra.** You should be comfortable with variables and coefficients, linear equations, graphs of functions, and histograms. (Familiarity with more advanced math concepts such as logarithms and derivatives is helpful, but not required.)
- **Proficiency in programming basics, and some experience coding in Python.** Programming exercises in Machine Learning Crash Course are coded in [Python](#) using [TensorFlow](#). No prior experience with TensorFlow is required, but you should feel comfortable reading and writing Python code that contains basic programming constructs, such as function definitions/invocations, lists and dicts, loops, and conditional expressions.



Tensorflow: Case study



AirBnB

Read about this case study in detail at

- <https://medium.com/airbnb-engineering/categorizing-listing-photos-at-airbnb-f9483f3ab7e3>



Part 4: Creating Concrete Value



NOTES:

Let's think about value for a moment.

What constitutes value?

- Better products
 - Path to innovation
 - New business opportunities
 - Improved customer experience
 - Higher revenues
 - Reduced production times
 - *(and many more)*
 - **The business case should be directly tied to value, and finance should be the ultimate scorekeepers.**



NOTES:

Where is AI valuable?

The #1 barrier to AI is often not the tech itself – it's how to figure out how to use AI for the applications it is good at performing.

So what do many high-profile uses of AI have in common?

Unfortunately, they suck. Not from a technological perspective, as technical capabilities are growing and exciting. Google's self-driving car initiative has been underway for more than 10 years now.

But consider end users: is there any demonstrable value? Do we go any faster compared to a human-driven Uber?



NOTES:

Common examples of mediocre AI

- **Extreme time to value** – don't meet expectations, and consistently overpromise and underdeliver. We have such high expectations from working with humans who do these jobs.
- **Virtual agents** – it's been 6 years since Apple launched Siri. There is now a rush to build virtual service agents (chatbots), and now many attempts at virtual sales folks.
- **Care providers, consumer medical recommendations** – More than 50% of user reviews are only 1 star on such services.



Virtual agents – How far along are we in interacting with a virtual service agent which does more than simply accepting an inquiry and passing us along to a human agent? How long will it be until there are robust, end-to-end sales or service representatives which people are willing to adopt?

Care providers – Medical advice and recommendations seem a good use case. You can input symptoms and they give back medical advice. But these systems have overwhelmingly not demonstrated value to end users.

Where (and why) does AI suck?

AI offers “subhuman automation”

- Usually can't handle long tail situations
- **Context** – can't fill in gaps in information as a human, because humans can empathize
- **Expertise** – can't reason or explain as a human could
- **Trust** – can't instill confidence



NOTES:

What is AI good at?

Jobs you would “feel bad asking a human to do.”

- Password resets
- Answering repetitive questions
- Data entry
- Moving things around a warehouse
- Connecting you with the right person
- Cost savings and profit opportunities are real, but those are a long way from producing “4th industrial revolution.”



NOTES:

What does AI do well?

There is however good progress with transformative potential. Treat AI as a tool for “superhuman insight.” When applied correctly AI has the ability to outperform any individual or team:

- **Data** – Can handle what’s beyond the scope of a person to digest
- **Analysis** – ability to compute in ways humans cannot
- **Complexity** – Can detect patterns too complex for humans to envision
- **Speed** – it’s a no brainer...this is where AI dominates over humans

Example: Self-driving cars? Consider instead driver assistance and navigation

Example: Medical recommendations? Consider instead imaging: AI is better than human doctors at identifying benign vs. malignant skin tumors



NOTES:

A few successful examples

Rethink AI for “superhuman insight.” Consider some successful examples:

- Starsky robotics for semi-autonomous trucks
- IBM Watson for oncology – (physician recommendation systems)
- Better results come from leveraging AI and humans effectively
- Encourage adoption



- **Starsky robotics for semi-autonomous trucks** – The Starsky platform allows remote central control of trucks which are largely, but not completely, autonomous. In February 2018, Starsky Robotics completed a 7-mile fully driverless trip in Florida without a single human in the truck. Starsky is the first company to publicly test an empty cabin for autonomous trucks.
- **IBM Watson for oncology** – provides both diagnostic and treatment guidance to a real doctor (physician recommendation systems). Instead of attempting to replace or automate a doctor, the AI automates portions of the oncologist’s job which are repetitive and error-prone, and uses this capability to further enable the doctor.

The fundamental advantage: These use cases avoid the long tail of different, unexpected, individual situations that don’t have a lot of precedent to learn from.

Better results: The most effective use cases leverage AI and humans effectively in a collaborative setting.

Faster adoption: More adoption means less change management. It may be hard to get a doctor to adopt a new treatment tool, but it’s vastly easier than adopting a “virtual doctor.”

Takeaways?

- Beware AI for subhuman automation
 - Leverage AI for superhuman insight
 - Rethink AI use cases accordingly if your goal to value looks like subhuman automation.
 - Start small, solve a specific problem, demonstrate success and value, and use quick wins to build adoption and support
 - Understand the prerequisites
 - Have a plan for human factors (more in section 7)



NOTES:

AI and the jobs landscape



NOTES:

AI and the jobs landscape

- AI and machine automation won't just involve job losses – also job gains and transformation
- Transformation will involve different ways of working alongside computers
- Human processes CAN be offloaded to machines
- Job losses will also lead to job creation
- “Robots are taking your jobs” fear in media is a little misleading and overblown. The scariest numbers attach the least firm dates and metrics.
- Job loss numbers make better headlines than job creation estimates.
- Few analysts focus on job gains also, and indeed they are poorly equipped to do so because data and prognostication around job losses
- Consider whether jobs will be displaced, created, and automated



NOTES:

AI and the Jobs Landscape

Consider Wordsmith, a service which uses AI to write convincing news stories: it can produce 4000 automated press releases but really hasn't replaced any journalists jobs.

- All the same, the robot revolution will reshape jobs and the economy
- It won't just involve job losses – also job gains and job transformation
- 67% of automation professionals see fear of job losses creating negative attitudes toward automation.
- Jobs are not immune – 17% (1 in 6) jobs will be lost to automation
- 10% equivalent jobs will be created by the automation economy



Source – JP Gownder, VP Principal Analyst, Forrester

For more on Wordsmith and how the application produces high-quality written press releases, check out <https://www.theverge.com/2015/1/29/7939067/ap-journalism-automation-robots-financial-reporting>

What will changes mean for humans who intend to remain working?

- **Upskill:** higher-value, more strategic and sophisticated decision making jobs
- **Job losses:** construction and mining, production, office and admin staff, sales and related roles (*i.e. customer service on lower end*)
- **Job gains:** next-generation management, business, financial
(*i.e. auditing ...how do you value a company that has AI in place?*)
- **New categories of professional services for literacy and AI collaboration**
- **Human-machine resources may disrupt HR.** Creative and artistic roles are required (*i.e. hiring novelists to write language for chatbots because software doesn't know how to write natural dialog*)



NOTES:

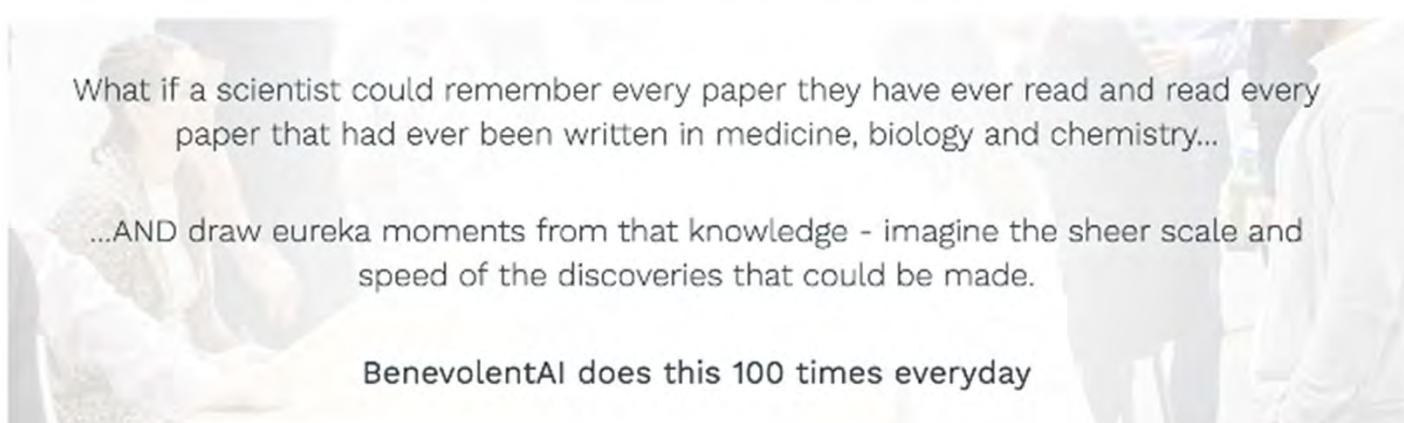
AI and the Jobs Landscape

- Prepare for digital workers (e.g. RPA bots, various flavors of AI, AI-infused robots)
- Learn about human-machine teaming. Coined by military for battlefield context
- Invest in change management. E.g. Baxter robot – he doesn't get programmed by C++, he learns from a human who demonstrates how to do the job on a factory line that Baxter needs to do
- Build in time for experimentation and failure. *Algorithmic learning and machine solutions can take a LOT more time than you might think or want.*



NOTES:

Use case breakout: Scoring data for AI consumption



What if a scientist could remember every paper they have ever read and read every paper that had ever been written in medicine, biology and chemistry...

...AND draw eureka moments from that knowledge - imagine the sheer scale and speed of the discoveries that could be made.

BenevolentAI does this 100 times everyday

A bold proclamation by London-based BenevolentAI (screenshot from About Us page, August 2017).



NOTES:

Scoring data in two medical use cases

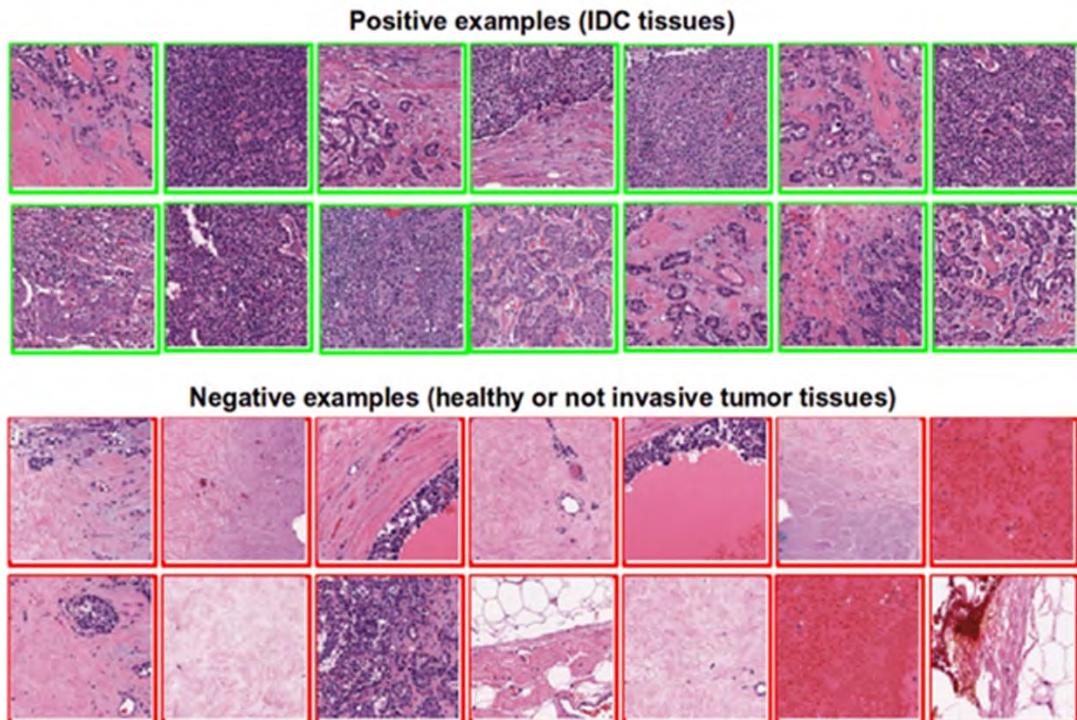
Two healthcare examples provide an opportunity to grade medical applications for AI. Grading criteria:

- **Quantity of data**
- **Quality of data**
- **Machine learning techniques**



NOTES:

Scorecard: Example 1



Scorecard 1 – Medical Imaging

There are many use cases for feeding the machine learning application images as datasets, then asking it to predict: (for example)

- 1) Breast cancer
- 2) Metastatic foci
- 3) Melanoma
- 4) Blood flow

B –

Quantity of data:

There is a lot of data, often with annotations and good quality.
However, it's sometimes hard to get to or takes some curation.

A

Quality of data:

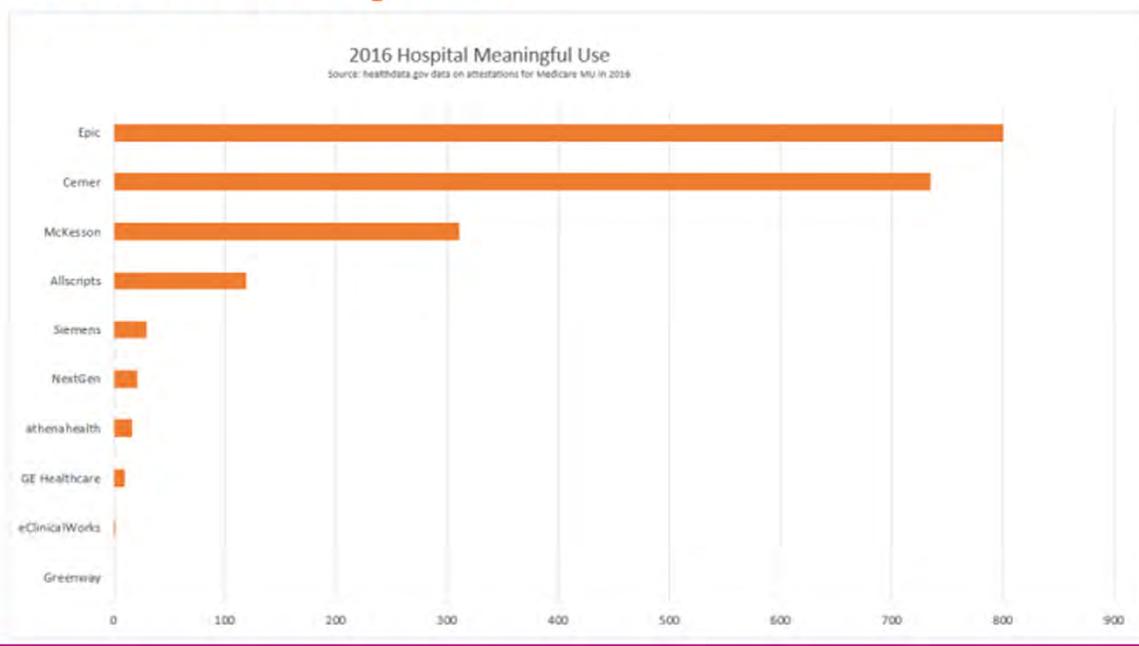
Medical imaging is excellent. Machine vision turned a corner in 2012.
Visual data is a ripe area for analysis.

A

Machine Learning Techniques:

Deep learning models are practical and accessible.
Machines are now almost as good as humans at identifying objects.

Scorecard example 2



Scorecard 2 – Electronic Medical Records.

A

Quantity of data:

Quantity is high, often with annotations

C

Quality of data:

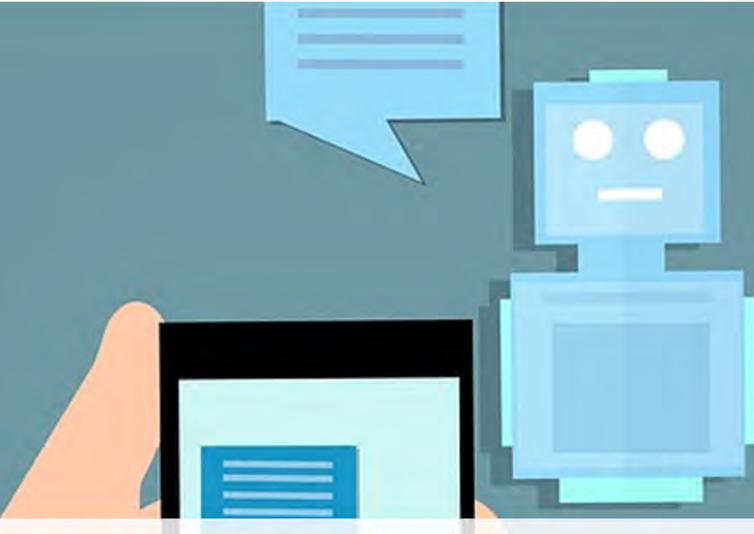
Quality, often low. It's rich but noisy...text is unstructured, full of abbreviations, acronyms, inconsistent, ad hoc, full of errors. **The dataset is often challenging for current AI capabilities.**

B –

Machine Learning Techniques:

- Natural language processing (NLP)
- Topic models
- Neural Networks

Models work but success is very focal given the specificity of the domain – looking at AI solutions that are more general or broad, they do a pretty good job at very surface, mainstream type questions (weather, sports). So success is limited, and ML works well with mostly more structured problems. The amount of specificity and nuance associated with the data and what it means can be difficult to extract at scale in a way that's well-suited to machines.



Part 5: Machine Intelligence as Part of Customer Experience



NOTES:

AI and customer experience

Key usage examples:

- Chatbots
- Common inquiries
- Sentiment analysis

Not good for:

- Pretending to be a human
- Anything beyond short tail scenarios



NOTES:

AI and customer experience

Gartner predicts that in the future the vast majority of customer interactions with a company will happen without human customer service rep. They also posit that a well-designed chatbot can handle 80% of customer inquiries. Consider just a few tasks that are fully serviceable by an autonomous agent:

- **Book flights**
 - **Buy tickets**
 - **Order food**
 - **Schedule rides**
 - **Purchase products**
 - **Send order updates**
 - **Book appointments**
-



Source: Ilya Gelfenbeyn, Product Lead, Google

Brands increasingly interact with their customers with active engagement via social media, messaging apps and smart devices. The increase in smart devices and the shift in brand engagement is leading to rapid adoption and fast growth for conversational applications.

AI and customer experience

As a result, customer demands are growing. Customers get used to this level of interaction and they expect you to be reachable anywhere, available any time, know their likes and needs, and to respond instantly.



Source: Ilya Gelfenbeyn, Product Lead, Google

Conversational experiences will be ubiquitous in the future, so it's important to develop your strategy now. AI-powered chatbots can handle many simple customer service tasks, preventing customers from having to spend time on hold waiting for the call center, and frees up the capacity of the call center to focus its time on more valuable customer interactions – such as upselling or resolving more difficult inquiries.

AI and customer experience

- **How can mobile site content, platforms, and Siri questions work together to present a desired experience?**
- **Starbucks is doing this now.**
- **Create digital experiences, not apps.**
- **Design still matters.**
- **A common question is, “How do you build a chatbot?” The correct answer is, “Why?”**



Michael Facemire, VP Principal Analyst, Forrester

What is Starbucks doing?

Members of the rewards program and mobile app authorize Starbucks to gather a lot of info about their coffee-buying habits from their preferred drinks to what time of day they're usually ordering. So, even when people visit a “new” Starbucks location, that store’s point-of-sale system is able to identify the customer through their smartphone and give the barista their preferred order. In addition, based on ordering preferences, the app will suggest new products (and treats) customers might be interested in trying. (*Forbes*)

The design process matters more than ever.

A common question is, “How do you build a chatbot?” The correct answer is, “why?” Sometimes the answer is, “It’s cool,” or “Somebody else has done it.” These aren’t really good reasons. But if they say, “It solves a cross-channel platform problem,” then they are in a much better starting place.

Design always matters! What does the brand look like, or better yet...what does it sound like to an end customer? When someone’s mom, son, or grandmother engages with your brand, will their side of the conversation reflect that your features are hitting the target of the experience your brand wants to deliver?

AI and customer experience

- Got a customer who wants to tie their e-commerce to a chatbot? Slow down!
- Intelligence bonds channel experiences.
- It can be a bad idea to just throw AI out there.
- It's no longer about platforms devices languages or sites.
- Intent is the tie that binds.



Michael Facemire, VP Principal Analyst, Forrester

Remember to ask, “What is AI good for?”

Got a customer who wants to tie their e-commerce to a chatbot? Slow down! Customers should think about what success looks like. It should be: conversation. If you treat an AI agent like an answering machine, maybe it's OK. Right now (2019) current data shows that with every exchange between the system and the person, the success rate goes down.

Intelligence should be applied to bond **channel experiences**. Humans behind the AI initiative should understand what people and customers want. It's best to build very small pieces and talk to people, because otherwise risk is high what gets built isn't what people want. Intent is the tie that binds.

Some folks want to just throw AI out there. Consider the horrible miss from Google's image recognition system several years ago, in which their image recognition service returned the label “gorilla” when shown the faces of people of color.

AI and Customer Experience

- **AI is improving fast but there is a lot more work to be done**
- **It's about more than technology – how do we as humans and society make decisions?**
- **What we call AI is not about replicating the human mind**
- **Context matters, and humans are still far better at understanding a myriad little details, intentions, and context**
- **Computation can help us understand things that we may not see or may not understand.**
- **AI is about expanding and amplifying our own mental abilities.**
- **Cognitive systems amplify human cognition. The IBM team thinks of Watson as “augmented intelligence,” not “artificial intelligence.”**



Source – Rob High, IBM Fellow, VP – CTO IBM Watson

Starting with fundamentals, what we call AI is not about replicating the human mind. If you look through history at Turing, Minksy, Von Neumann, and others, there is a pattern of thinking about making computers behave as a human so we can get answers that we may seek at a human level. However, these guys were mathematicians and they thought of computing primarily in terms of mathematics. However, it is very hard when we think about all the subtleties and innuendo that go into our everyday experience to come up with some sort of a model that would accurately describe what we experience.

Context matters, and what we as humans are good at is taking in and forming an understanding of what another human speaker on stage is saying right now – a myriad little details, intentions, and context that inform our understanding of what is being said.

What we can do however, is gather enough knowledge and use of computation to help us understand things that we may not see or we may not understand. Thus, the IBM Watson team thinks of AI not as artificial intelligence, but augmented intelligence.

We aren't very good at consuming and assimilating vast amounts of information, reading millions of pages, and discovering the relationships and insights in that much information. So what computing and AI allow us to do is amplify our own ability. Ever since tool making began in human history, the tools that have worked the best and lasted the longest have been about somehow amplifying our own strength, our own reach, our own human capabilities. AI is about expanding and amplifying our own mental abilities.

Cognitive systems amplify human cognition.

IBM Watson is delivered as a set of APIs in the cloud. (ibm.com/bluemix). It is a product which is all about provide this amplification to third parties in a broad range and huge number of applications. They have boiled down the capability of Watson's AI to these two higher reasoning skills: conversation and discovery.

The goal is to inspire and help understand how to think about problems in a different way, and ask questions we may not otherwise think about. If all we do is ask and answer the same question again and again, we won't evolve, grow or make progress. So higher-order ways of reasoning should help us do these things.

Conversation – it should:

- Engage the user
- Focus on the user's broader concern
- Build on an idea
- Leave the user inspired and satisfied

For the most part, right now, machine powered conversations really only focus on one term or one question. For instance, "What's my account balance?" may be the immediate question you need an answer to, but there is a deeper intent there. You may be trying to figure out whether you can afford to buy something. You may be planning to save for some specific goal. Those are the types of deeper concerns that IBM wants to understand with Watson in order to help a user with the real concern or the real goal.

Most of applied AI around 'conversation' are chatbots. Today's "chatbots" are really not serving deeper needs like this. They also rely on typing, which you can't do in a car. Chatbots have the highest level of successful interactions with only one exchange of interaction or question. At two exchanges, success as defined by user satisfaction drops precipitously, and by three interactions the measure of success is incredibly low. A full toolbox of language, speech, vision and empathy is really what's needed for conversational interaction with a machine.



Part 6: Machine Intelligence & Cybersecurity



NOTES:

How can ML help with security

- Fully automated systems are fast, but error prone
- Augmented / hybrid / extended intelligence gaining the most traction
- Most ML systems use row-based data, most threats are best represented as a graph



The buzzwords are in flux around augmented, hybrid or extended intelligence, but the core concepts are the same: let machines do what they do best (read log lines) and let humans do what they do best (infer motives, goals, and adapt to new situations)

AI in Security – how and what?

FOUR QUESTIONS to start the process:

- **Where in IT operations is there already high automation?**
- **Where in a DevOps-style or CI/CD pipeline for application development or IT services is work already heavily automated?**
- **Where is high-volume, high noise-to-signal machine data being produced?**
- **What is the structure of that type of data?**



Evaluate your AI prerequisites to find security improvement opportunities

This isn't just an exercise in applying advanced analytics and AI to your security practices. It's also an opportunity to use what we have already learned about AI in general to begin thinking about inter-domain use cases. Domains have specific considerations, but many AI capabilities transcend them. To wield AI effectively, one should find lessons proven in one domain and look for novel ways to apply them to other domains.

Before we proceed: Class exercise

Let's begin describing a user story for a security-related application using AI.

Remembering the previous "Scoring data for AI consumption" exercise, start with the **FOUR QUESTIONS** in the previous slide.



Please answer:

Where in IT operations is there already high automation?

Where in a DevOps-style or CI/CD pipeline for application development or IT services is work already heavily automated?

Where is high-volume, high noise-to-signal machine data being produced?

What is the structure of that type of data?

Class exercise (continued)

Now let's grade the three components of an AI use case in the same style demonstrated in the earlier healthcare example. This time we're thinking about security.

- **Quantity of data**
- **Quality of data**
- **Machine learning models available**



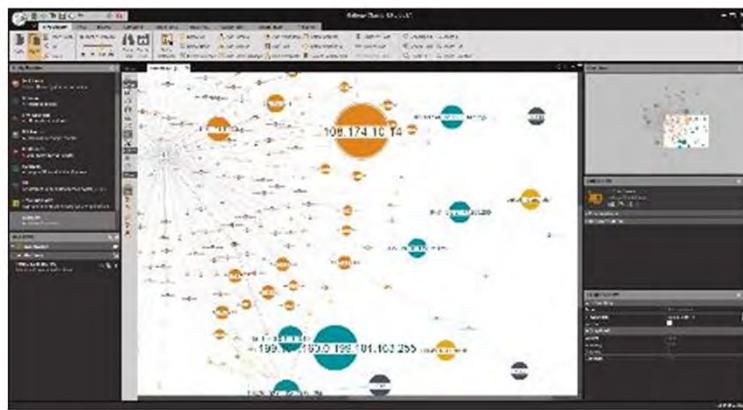
Quantity of data

Quality of data

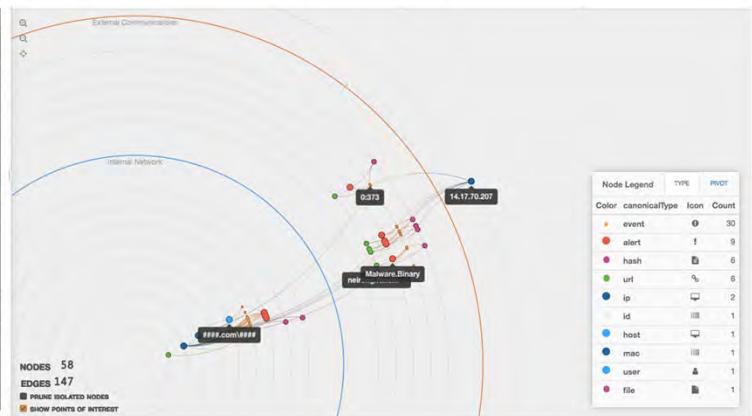
Machine learning models available

Graph Data for Cybersecurity

Maltego



Graphistry



Maltego is an established player in the space, Graphistry is the new upstart, but both involve manual graph exploration.

There is a host of graph based data stores such as Neo4J, and a large less-well-known series of graph algorithms that may need specialized skills, but have large untapped potential.

How are attackers leveraging ML and AI

- **AI as the tool**
 - **Identifying vulnerabilities**
 - **Automated attack escalation**
- **AI as the target**
 - **Reverse engineering the algorithm to steal IP**
 - **Adversarial input**



On tools:

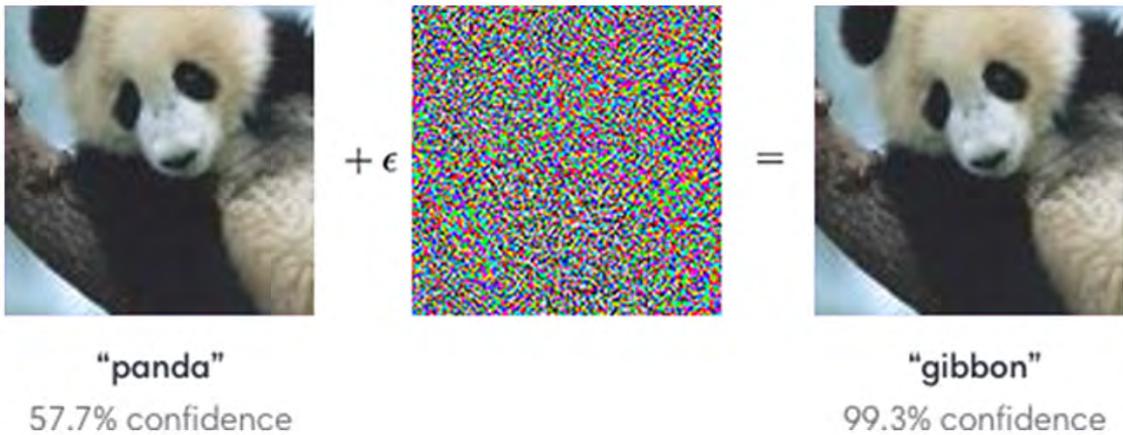
You want to stay ahead of the curve. If you are automating the detection of vulnerabilities, there is a good chance your automation system is as good or better than (non-state) attackers.

On reverse engineering:

Don't allow ungated access directly to the ML system. If attacker can give arbitrary input and receive direct output, reverse engineering is trivial https://regmedia.co.uk/2016/09/30/sec16_paper_tramer.pdf

Adversarial input looked at next slide.

Adversarial Input



From open AI blog: <https://openai.com/blog/adversarial-example-research/>

NOTES:

Automated application monitoring

- **Subset of the “anomaly detection” problem**
- **You might need a separate model for each type of application**



NOTES:

Identifying Vulnerabilities

- **Too many systems and too many ports to test everything with nmap**
- **Prioritize, learn, reprioritize**



The “prioritize, learn, reprioritize” flow can’t be a simple model, and requires fairly sophisticated data pipelines.

Automated Red Team

- You want generative models: “If I was an attacker in this situation what is the chance I would do X”
- Sandboxing is hard to get right.



“Generative models” ties back to Naïve Bayes type of ML discussed in the “Implementing Machine Learning” section.

Sandboxing: Too small and your test case is too artificial, too loose and you risk damaging real systems.

Modeling previous breaches

“It has been said critically that there is a tendency in many armies to spend the peace time studying how to fight the last war.”

- *Lt. Col. J. L. Schley, 1929*



This quote is from *The Military Engineer*, “Some Notes on the World War,” one of the earliest examples of this idea.

The reality is, it's a tough balance to get right. We want to learn from past mistakes, but the cybersecurity space is rapidly evolving. Archived data from 6 years ago might not be good training data. Conversely, this presents opportunities for AI to be a good tool, since the ability to adapt and evolve is fundamental to an AI application. We are not looking for a pattern of “learn and repeat,” but of “inspect and anticipate.”

For further reading:

https://www.barrypopik.com/index.php/new_york_city/entry/general_always_fight_the_last_war

Automating Incident Response

- **Alerting vs automated action**
- **Need to be very careful with the algorithm's reward function**
 - **No network traffic = no attack traffic**



NOTES:

Detection of APTs

- Need global scale data beyond one organization
- Can help build better models because more variance across APTs than within a single APT



NOTES:

Using Natural Language Processing (NLP)

- **Near-hit domain names**
 - www-ebay.com, weightloss4tmz[.]com
- **Phishing email detection**
- **Bleeding edge: hacker message board monitoring**



NOTES:

Fraud Detection

- **Incorporating timeseries data often breaks independence assumptions**
 - Fewer assumptions = more expensive and slower projects
- **Secondary benefit of better understanding customer habits.**



NOTES:

Reducing compliance cost

- Dimension reduction techniques also partially anonymize data
- Can you train on aggregate data?
- Remember generative models for testing



NOTES:

Building trust in automated security decisions

- **Always start with humans in the loop**
- **Data provenance as a first-class feature**
- **Great idea to hire outside experts to audit**
- **Remember to visualize. A picture is worth 1000 words**



NOTES:



Part 7: Teaming and Internal Capabilities



NOTES:

Teaming and Internal Capabilities

PART 1 – “Human Architecture”



 ASPE

NOTES:

Data analysis roles

Champions

Professionals

Semiprofessionals

Amateurs

- **Four categories of analytical people**
- **Broad categories to “describe the challenges of managing analytical talent”¹**

¹ *Analytics at Work: Smarter Decisions, Better Results* by Robert Morison; Jeanne Harris; Thomas Davenport, HBR Press, 2010



The roles in a successful analytics organization

We break down into four categories the actors in the data driven organization. As is usual, different sources have different granularity, but these four provide a nice subdivision of skills one can expect to encounter. These are not jobs, but classifications of jobs. The goal is to provide an approach to the challenges of managing talent in a data driven organization.

Data analysis roles

Champions

Professionals

Semiprofessionals

Amateurs

- Executive decision makers.
- Rely on analytics to understand business process and decisions.
- Support and advocate business initiatives aligning to metric-driven methods.
- Advocate analytics across groups in the organization.



Data Champions

These are the characteristics of the champion in an organization. The champion faces in two directions – towards the stakeholder/customer and towards the analytics team. Champions advocate for analytics in the org AND use analytics in the org. It is not enough to simply talk about using data and analysis to make decisions. You have to put it into practice and give it the place of prominence. In other words, champions lead by example.

A few examples of the Champion Role...

For example, this might be a person whose title or role is: higher up in org, pushing people to be more data-driven. Leading the charge – CFO/CIO/CTO. Could be technical or non-technical person. Anyone with clout who can campaign for data-driven actions.

Data analysis roles

Champions

Professionals

Semiprofessionals

Amateurs

- Knowledgeable in advanced analytics methods techniques, statistics, machine learning, math.
- Create the models and algorithms to be used in organization.
- Typically hold advanced degrees and want to continue growing.
- Develop best strategies for achieving goals.



Data Professionals

Professionals are the people on the team who are performing advanced analytics. This category needs the support of the champions because what they do can be difficult to understand for people in the other categories - remember, literacy does not mean expertise. Especially when working to create new algorithms and models, success is not always imminent, or even guaranteed, so champions have to hold the line and advocate for the practice.

When planning for building the team, decision makers have to accommodate for the cost it takes to find the professionals as well. Since they often hold advanced degrees and have specialized expertise, they do not come cheap. However, their usefulness becomes apparent, especially when working towards a more predictive path in the organization.

A few examples of the Professional role...

Data architect, engineer, scientist, DBA expert, statisticians

Data analysis roles

Champions

Professionals

Semiprofessionals

Amateurs

- **Can apply the models and algorithms.**
- **Understand the business through its data and use.**
- **Link analysis to business insights.**
- **Skilled at analytical applications and visualization tools.**



Data semiprofessionals

Semiprofessionals are people who can apply data models to the business. These are people who have domain expertise, instead of theoretical expertise in computer science or statistics, and through the domain expertise, apply the models and analytical techniques to the business in which the org works.

A few examples of the Semiprofessional role...

Person who is using the results, connecting the pipes by using tools, programming dashboards: Data Analyst, Data Visualization Specialist, Business Analyst, Marketing Analyst

Data analysis roles

Champions

Professionals

Semiprofessionals

Amateurs

- **Basic understanding in order to perform job tasks.**
- **Consumers of analyses.**
- **Can apply insights to their siloed job, and need support to gain greater picture.**
- **Often report summaries, run the same analyses repeatedly.**



Data Amateurs

The Amateur is the ultimate consumer of most analytics output. This does mean that the person is amateur at all things, just that the level of expertise in analytics is as the user of the final product developed by the professional (new models and techniques) and the semiprofessional (building applications of the techniques to the business). This person is the person that loads data into the tools, steps through the predefined processes, and then makes reports with data, visuals, and job relevant insights for decision makers. Decision makers and leaders themselves can be data Amateurs, and often are. That doesn't mean they aren't excellent leaders...especially if they recognize their need and dependence on the right analytics team who can inform their decisions and empower their leadership. For leaders to receive this benefit, they must in turn also build, empower, encourage and protect the team.

A few examples of the Amateur role...

Who in the organization connects results and generate reports?

- Marketing professionals
- Sales professionals.
- “Champions without power”
- Those who consume data but don’t know how it’s actually produced
- Inventory Manager
- Call center employee
- Business Manager

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

A typical organization will have a variety of analysts and engineers that fall into the categories previously described. The next few slides will walk through a typical set.



From Creating a Data Driven Organization by Carl Anderson, O'Reilly Media, 2015

Specific roles in an advanced analytics practice

This is a further subdivision of the data analysis enterprise as seen from a different source. Instead of the larger granularity of the previous slides where we provide a few sample roles in each, here we break the analysts into more specific roles. To provide some context, these fall under the previous descriptors except for Champion. We will discuss some of the responsibilities of each of these roles in the organization, touching on some more technical details.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- **Wrangle the data:**
 - Internal sources
 - External sources
- **Specialist in ETL:**
 - Extract
 - Transform
 - Load
- **Get the data to the analysts**



Data engineers

The Data Engineer is someone who lives with the data day to day and provides the analysts with access to the data they require. Data for an organizational problem needs to be acquired then “wrangled” into the organization. While identification is a team effort, the technical skills of this role are catered to acquisition. The data engineer is strong on the database technologies and the computer science of Extract-Transform-Load (ETL) which will be discussed in detail in later sections. At this time, it is enough to understand that data is in the world in many forms, but the org needs it specific formats. The tools needed to transform and store the data in the correct form, in sufficient time, and with sufficient space is the responsibility of the data eng. The analysts and reporting specialists need data provided them in forms they can manage and this is also the area of the data eng.

To use a relational database example, a data eng imports data into the appropriate normalized form for optimal storage and querying. When another role requires some portion of the data, the data eng develops the proper access description (not to limit access in principle, but to protect data against corruption) then builds views of the data with SQL and informs the consumer of the names of said views and how to gain access.

There is more to this person’s role, but for understanding how a mature analytics team underpins success with AI, this is enough for now.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- **Liaison between the tech and the business stakeholders.**
- **Identify process improvement.**



Business Analysts

The Business Analyst is a consumer of data. The role provides the connection between the specialty of business and the specialty of stats/comp science. Between data and analyses, the business analysts focuses on the business process. For instance, a specialist in Just-In-Time supply chain will use analytics to improve aspects of the JIT process to minimize the product inventory on hand during a 24 hour period, or work with raw material providers to lower the raw material inventories by improving material flow.

Each aspect of a given business has internal business processes with room for improvement and the business analyst is the domain expert that brings process and analysis together.

This person's role is to define the metric inventory for the organization and align to those goals. They also work on data quality and process inefficiencies, which often increases their value to the organization and gives them greater purpose.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- “**Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.**”

—Josh Wills

- **Typically have advanced degrees in an area of business use and know statistics.**
- **Commonly report spending 80 percent of their time collecting and preparing data before analysis.**



Data Scientists

The Data Scientist is the Professional from the previous category discussion. This is the role that knows stats/math and computer science. In a lot of ways, when dealing with the high level data scientist, this is birthplace of machine learning and data mining algorithms. They are not the people who design the next-generation computer or the people who create mathematical theories, but combine the current tech and math to develop new ways of finding information from our data.

This is the person that build the models and get the Data Engineers/Architect to align their efforts to support the machine learning processes that they're going to put in place.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- **Typically hold advanced degree in statistics and has deep appreciation for statistics.**
- **Focus on modeling and algorithms.**
- **Ability to make inferences in the statistical process.**
- **Work with the Business Analyst to build presentations and data storytelling process.**
- **"One quarter of statisticians in the US work for federal, state, and local government."¹**



Statisticians

Another in the Professional category, the Statistician is the role that is developing the analytics needed. Now it is important to note that this does not mean creating new statistics, or writing papers on new mathematical theorems. Rather, it is the selection of methods in such a way as to build an analytics set. For a very simple example, the statistician determines that the mean, median, and standard deviation interpreted on a normal curve are the statistical elements that are relevant to the data available for the business problem.

A statistician is focusing on models and methods that help provide the spectrum of analysis – historical, descriptive, predictive, and inferential. It is no accident that government employs so many. Predictive models of health and economics (among many) set the forecasts of budgetary requirements for years.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- **Focus on internal financials.**
- **Develop models for market development.**
- **Primarily work on reporting.**
- **Use models for prediction.**



Accountants and financial analysts

Accountants and Financial Analysts are in the Semiprofessional category as well. A financial analyst looks at the financial picture of a company (past present and future) in the context of market trends and business goals. Accountants look more at the day-to-day details of the finances. Both are proficient with “the numbers” of the business operations. Both can consume analytical tools. To be more precise, the statisticians and the data scientists work to deliver models that may apply to the financial forecasting for the business. This set of roles can both test and (when/if proven) use the models when working to achieve the financial health desired for the business.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

- **Strong design skills (graphic designer).**
- **Good tech skills (programmer)**
- **Often build in programming tools such as D3, Rshiny, Sparklyer, PowerBI, Qlik, Tableau, Cognos, MSTR.**
- **Can generate reports, but works on interactive tools (dashboards, balanced scorecards).**



Data Visualization Specialists

Data Visualization Specialists fall in to the category of amateur. But don't let that lead you to believe they are new hires and without skills necessary to the enterprise. Lots of reporting is assumed to be a PDF with a few charts that are handed to relevant consumers of the information, but this is a subset. Proactive organizations use tools that attach to data in "real time" (caveat for data considerations that will be covered in more detail later) and can benefit from novel ways of viewing data. The art of reporting, so to speak, is making the data speak to the specific audience in the correct way. The question always being asked is "how do we make the argument and tell the story?" The days of synopsis with bar charts are fading. Today, the art of drawing one's attention to the important conclusions from the data with the right style and chart is visualization. Connecting the art to the live data makes it interactive and immediate. These are important skills for this particular consumer of the analyses.

Data analysis roles

Data engineers

Business analysts

Data scientists

Statisticians

Accountants and financial analysts

Data visualization specialists

Domain experts

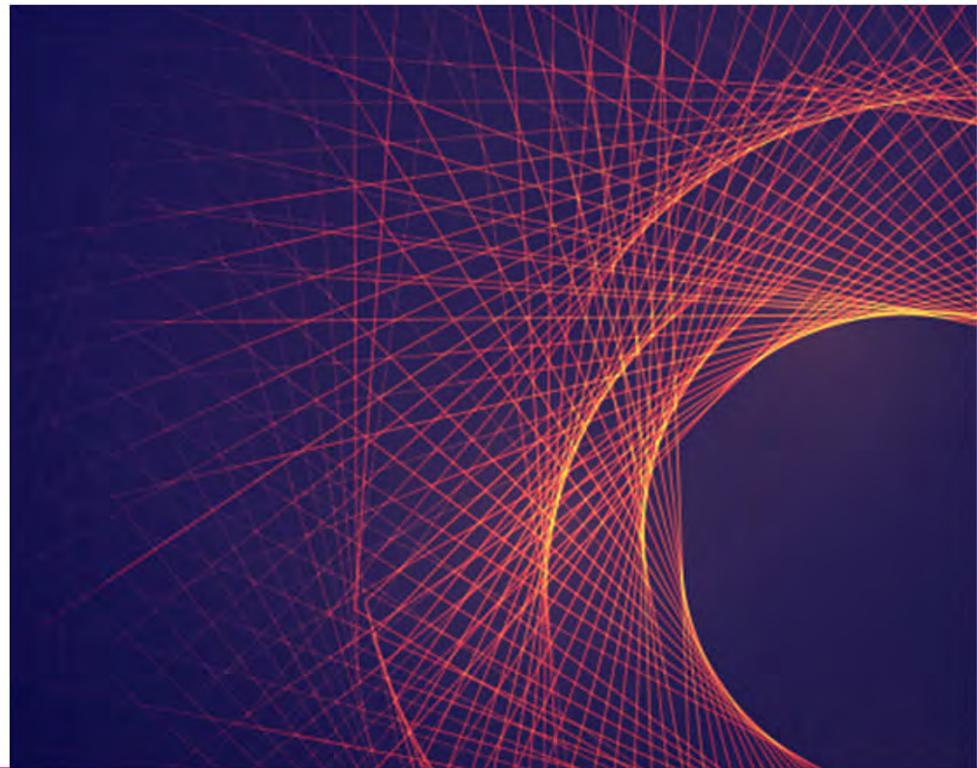
Not a role specific to data, but is usually a necessary collaborator in order to inform and define business requirements so engineers and data scientists can build and implement a useful model.



NOTES:

Teaming and Internal Capabilities

PART 2 – The Technology Ecosystem



NOTES:

Technology ecosystems are evolving from the relational databases to big data and AI systems. Where are we today?

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

- **Structured data management**
- **Often the combination of:**
 - ETL/discovery/transformation/data prep tools
 - Analytical tools for statistics and modeling
 - Operational processing, governance, data security and recovery tools
- **Available in open source and license packages**



Relational data

Starting from the databases, let's break down a typical data ecosystem. The first we will cover is relational databases. They are used for structured data, even though one can dump unstructured data into it by setting the least stringent restrictions possible, they excel when they are enforcing type and relationship constraints. The technology is meant to provide high flexibility in access to data with optimized memory. Most products, from Oracle to SQL Server, come with additional tools to help with ETL, performance analysis and recovery tools. These days most products also have community editions which are free to the user.

Technology ecosystems: NoSQL

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

- **Structured and unstructured data**
- **Heavy on the unstructured**
- **Compromise consistency (from relational) for easier scaling and speed**



NoSQL

NoSQL databases refer to “not only SQL” and generally mean non relational databases. These are databases which usually use unstructured data and provide simpler scalability in clusters. Often, consistency is compromised in order to have availability in clusters. What this means in a practical sense is that, when clustered, speed of writing data may dictate that some nodes do not get notified of new data until much later than the actual write on the original node. Therefore queries may not reflect the “up to the minute” data. These are not necessarily problems, just considerations for selecting the right tool for the job.

Always identify the business need. This is key for structures that have lots of empty data or redundant.

Technology ecosystems: NoSQL, graph databases

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

- **Graph databases can be queried in many different ways.**
- **Provide great flexibility in traversing the graph.**
- **Adding new relationships to a graph data store is easy.**



NoSQL (*continued*)

Graph databases can be queried in many different ways

- For example, get all nodes of people that like Sony products
- Known as traversing the graph

Provide great flexibility in traversing the graph

- Also very high performance

Adding new relationships to a graph data store is easy

- Relational database requires schema changes and data movement

Lastly, though graph databases are hard to imagine and we have provided only some small definition, here are few points about them. We can query graph database from a number of different perspectives (along the edges) and provide flexibility in traversal along the edges from node to node. Storing the relationships with the nodes provides a way to add new relationships as needed without changing the structure of the database (just write the data and the edges). Remember though, if the relationships are largely static, this is just a relational database in another form.

Technology ecosystems: Big data

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

- **Parallelize the processing to trillions of rows and many variable types of data and data types. Lots of use cases that NoSQL and SQL databases were not designed to.**
- **Often merge and have relational techniques that are brought into the mix to be able to get the value of both tools working together.**



Big Data

NoSQL databases play a large part in the Big Data ecosystem. In general terms, with large datasets, scalability becomes a primary function. Tools like Hadoop bring in the ability to parallelize functionality across clusters of computers. Without being too technical, imagine a scenario where you are looking to provide a similar purchases recommendation tool for a large company (on the order of Amazon or Facebook). **The solution would go something like this:**

Choose some criterion, like genre purchased by a user, and begin to search all purchases for similar purchases in that genre. Now the data structure the company has is vast catalogs of users and what they purchased as child properties. the data is largely unstructured except for that requirement. To provide scalable uptime for the product, all the data is stored in easy scalable NoSQL databases where users with last name starting with A are on one machine, B on another and so and so on. We have to find the genres that people like and are similar to our current user. We can parallelize the problem by searching all A's on one machine, B's on another because one search result does not affect another. When we have done that, we can gather results on one machine, perhaps a relational database with a more structured data schema so that it is on call for when the user logs in. The “Big Data” access is parallelized, the NoSQL allows variable structures (some people may not have ever bought books, but instead washing machines) and the relational database allows faster recall of the specific information in a structured manner.

This is the type of activity and interplay between components in the ecosystem. The example glosses over many things and the analysis is a trivial comparison, but gets to the heart of what we might try to do with “Big Data” system.

Technology ecosystems: Big data

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

Volume, variety, velocity:

- Hadoop handles all the data and will show results... eventually.
- Velocity issues are met with "eventual consistency."
- Storm can intercept the data, process it, and give instant results.

Use cases:

- Log data can be checked for DDOS or hacking.
- Patterns in sales data can drive real time marketing.
- Twitter and other social media comments can be mined.
- Sensor data analysis can predict problems before failure occurs.



Storm Use Cases by industry

source: <https://hortonworks.com/apache/storm/>

Prevention Use Cases:

- Financial services – Securities fraud, Operational risks & compliance violations
- Telecom – Security breaches, Network outages
- Retail – Shrinkage, Stock outs
- Manufacturing – Preventative maintenance, Quality assurance
- Transportation – Driver monitoring, Predictive maintenance
- Web – Application failures, Operational issues

Optimization Use Cases:

- Financial services – Order routing, Pricing
- Telecom – Bandwidth allocation, Customer service
- Retail – Offers, Pricing
- Manufacturing – Supply chain optimization, Reduced plant downtime
- Transportation – Routes, Pricing
- Web – Personalized content

Technology ecosystems: statistical tools

Relational databases

NoSQL databases

Big data tools

Statistical tools

Machine learning

Visualization and reporting tools

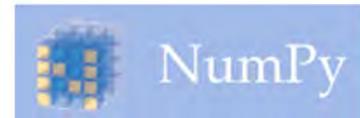
- **Analysis tools for modelling and calculation**
- **Sometimes libraries for general use languages (NumPy, SciPy for Python)**
- **Sometimes programming languages all their own (R, SAS, SPSS, Spark)**



NOTES:

Technology ecosystems – statistical tools

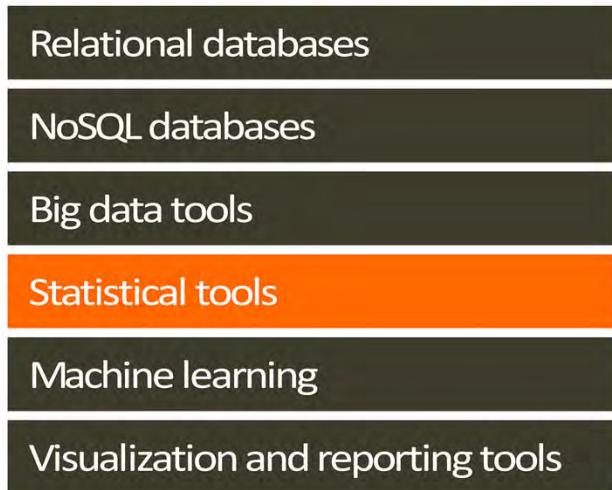
- Relational databases
- NoSQL databases
- Big data tools
- Statistical tools
- Machine learning
- Visualization and reporting tools



Statistical tools - Python

Python is a general purpose programming language which has additional libraries available through a variety of open source organizations that enhance its mathematical capabilities. These include statistics, plotting, and scientific packages.

Technology ecosystems: statistical tools

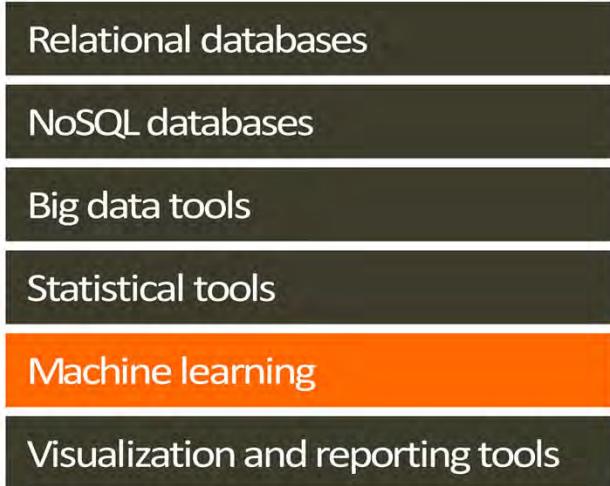


Statistical tools – R

The R programming language is a language designed with analysis in mind. Data is loaded into data frames and can be cast into collections of various types where mathematical operations and functions can operate on the entire set at once. The R product space has packages to make available web controls for manipulation of data, and other visualization tools.

R is an open sourced statistical tool.

Technology ecosystems: machine learning



- **Blending statistics with computer science**
- **From pattern recognition to artificial intelligence**
- **Make data-driven predictions and decisions**



Machine learning – expanding the definition:

“Machine Learning at its most basic is the practice of using algorithms to parse data, learn from it, and then make a determination or prediction about something in the world.”

– [Nvidia](#)

“Machine learning is the science of getting computers to act without being explicitly programmed.”

– [Stanford](#)

“Machine learning is based on algorithms that can learn from data without relying on rules-based programming.”

– [McKinsey & Co.](#)

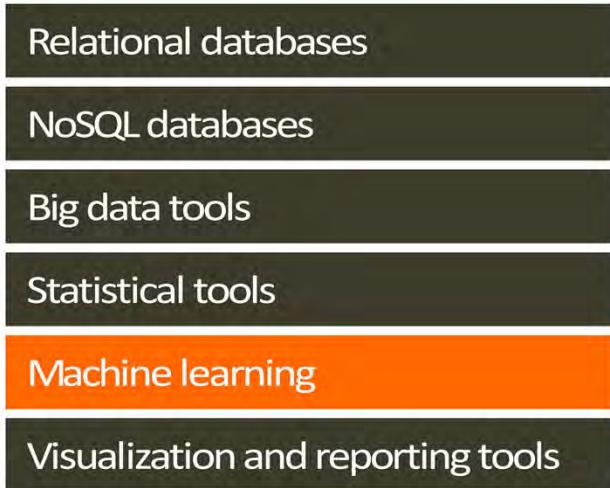
“Machine learning algorithms can figure out how to perform important tasks by generalizing from examples.”

– [University of Washington](#)

“The field of Machine Learning seeks to answer the question “How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?”

– [Carnegie Mellon University](#)

Technology ecosystems: machine learning



- **Blending statistics with computer science**
- **From pattern recognition to artificial intelligence**
- **Make data-driven predictions and decisions**



In the “Implementing Machine Learning” section of this class we covered the pillars necessary for successful machine learning:

- **Algorithms**
- **Modeling**
- **Business Case & Business Integration**
- **Domain Expertise**
- **Data**
- **Automation**
- **Scalability**

Who's using it:

Financial services – Banks and other businesses in the financial industry use machine learning technology for two key purposes: to identify important insights in data, and prevent fraud. The insights can identify investment opportunities or help investors know when to trade. Data mining can also identify clients with high-risk profiles, or use cybersurveillance to pinpoint warning signs of fraud.

Government – Government agencies such as public safety and utilities have a particular need for machine learning since they have multiple sources of data that can be mined for insights. Analyzing sensor data, for example, identifies ways to increase efficiency and save money. Machine learning can also help detect fraud and minimize identity theft.

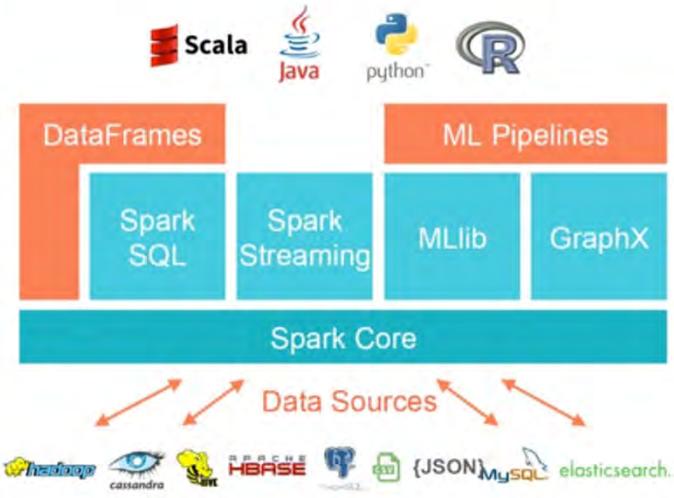
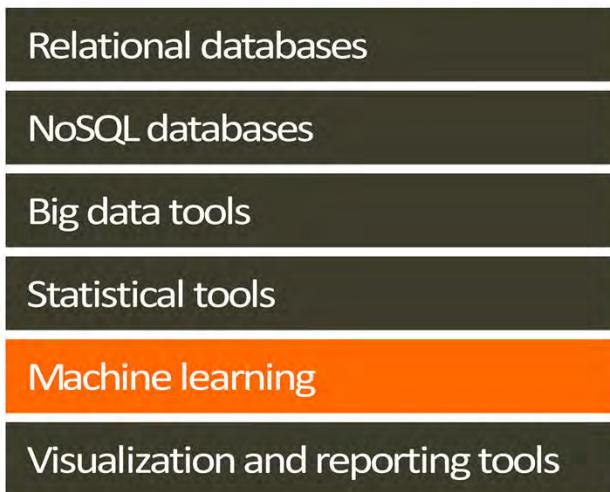
Healthcare – Machine learning is a fast-growing trend in the health care industry, thanks to the advent of wearable devices and sensors that can use data to assess a patient's health in real time. The technology can also help medical experts analyze data to identify trends or red flags that may lead to improved diagnoses and treatment.

Marketing and sales – Websites recommending items you might like based on previous purchases are using machine learning to analyze your buying history – and promote other items you'd be interested in. This ability to capture data, analyze it and use it to personalize a shopping experience (or implement a marketing campaign) is the future of retail.

Oil and gas – Finding new energy sources. Analyzing minerals in the ground. Predicting refinery sensor failure. Streamlining oil distribution to make it more efficient and cost-effective. The number of machine learning use cases for this industry is vast – and still expanding.

Transportation – Analyzing data to identify patterns and trends is key to the transportation industry, which relies on making routes more efficient and predicting potential problems to increase profitability. The data analysis and modeling aspects of machine learning are important tools to delivery companies, public transportation and other transportation organizations.

Technology ecosystems: machine learning



Machine Learning (*continued*)

Machine learning algorithms are implemented as libraries and functions in the Spark, R and Python environments mentioned before. For example, the scikit-learn package in Python has Tree Classifier algorithms and R has kMeans and kNN built in. Both languages come well stocked with useful machine learning algorithms built in, these are just a couple.

<http://scikit-learn.org/stable/> is the link to the Python library

Technology ecosystems: visualization and reporting

Relational databases	<ul style="list-style-type: none">• Allow easy and/or advanced visualization of data sets and results.
NoSQL databases	<ul style="list-style-type: none">• Quickly extract and present useful information.
Big data tools	<ul style="list-style-type: none">• Examples:<ul style="list-style-type: none">— Excel— Tableau— D3— PowerBI
Statistical tools	
Machine learning	
Visualization and reporting tools	



The importance of data visualization

Data visualizations come in many forms. Everything from simple excel spreadsheets to full BI packages that allow you to build complex dashboards or provide reporting on AI projects require effective visualization.

It is critical that you define your organizations needs before trying to decide on a tool. This decision can be affected by monetary considerations, technical capabilities, infrastructure issues, and more.

In the end you are trying to identify the quickest, cheapest, most compelling and effective way of communicating your analysis.

Technology ecosystems: visualization and reporting

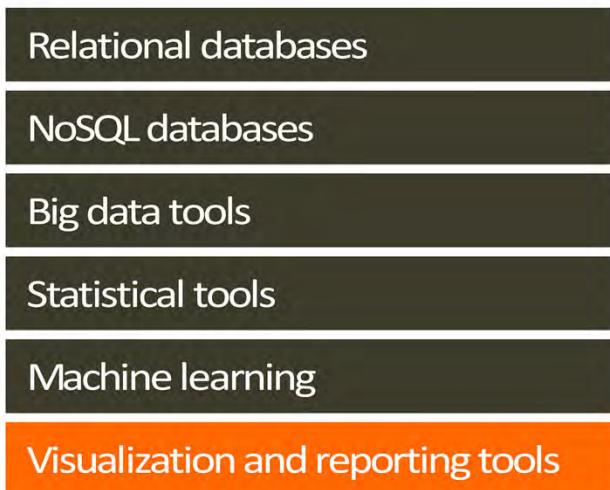
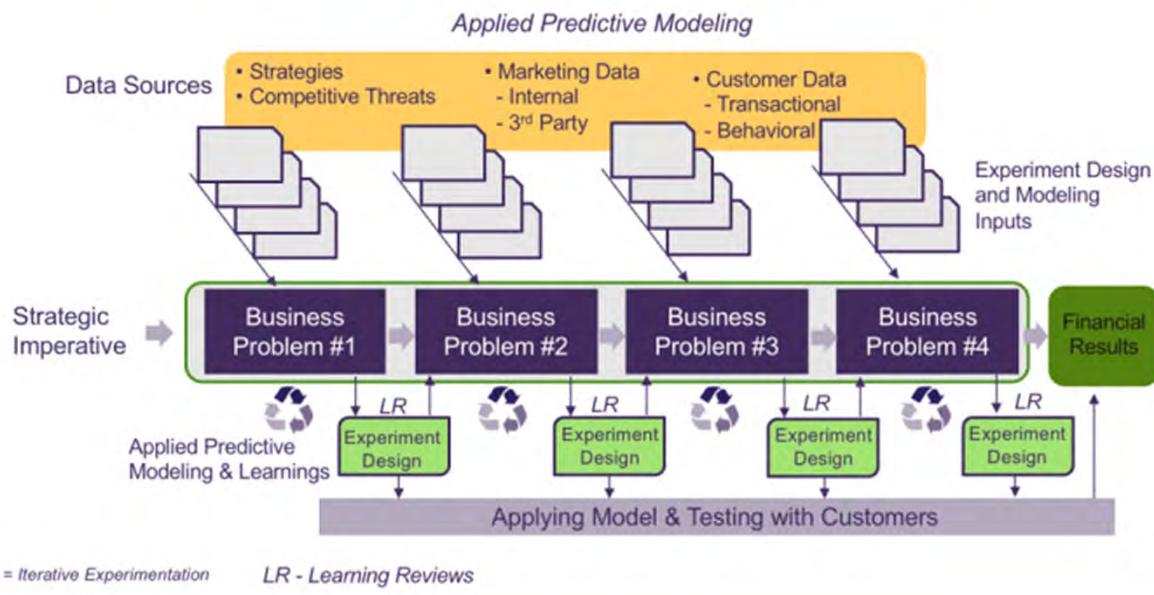


Tableau for data visualization

Tableau is an application that helps people analyze and visualize data. It is available for license at www.tableau.com/

This is an example of a product that provides a platform for building reports and sharing them with colleagues. It does have a public (free) version as well as commercial licenses allowing new data analysts to learn the core skills without paying enterprise prices right from the beginning.

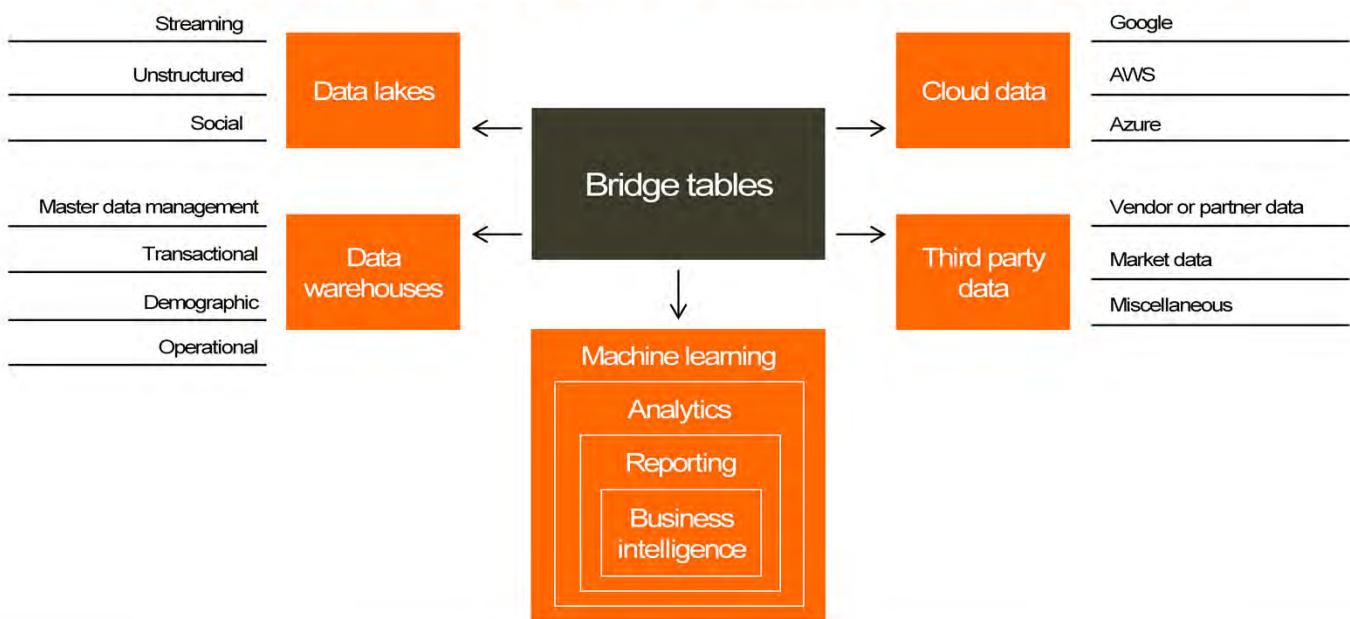
Technology ecosystems: Putting it all together



Source:
ASPE

Source: Jordan Martz, DataMartz, LCC

Technology ecosystems: Putting it all together



Solutions must be integrative.

In the relational DB section we discussed the concept of bridge tables with regards to traditional databases. This is not the only place where a bridge table may be needed in the modern data landscape.

In order to provide full analysis it is important to bring together many diverse potential data sources. In a successful analytics you will need to consider how to tie these sources together.

- What commonalities exist?
- Are they tied to similar metrics and/or scales?
- How can/should these data sets be used together?
- Which components add meaningful context to your analysis?

It is not enough to simply build a pretty graph... That visualization is the culmination of the full ecosystem at work.



Part 8: Conclusion and Charting Your Course



NOTES:

Class discussion: charting your course

We'll conclude class by taking a few minutes to answer the three questions in your classroom manual. When we've finished, we'll go around the room and share the answers.



Please answer:

What have you learned about AI that you didn't know before?

What ideas do you have for possible use cases of AI?

What are some actions you can take now to begin advancing AI initiatives in your own work?



Thank you for your time!

Please take a moment to give us your feedback on the class using the participant evaluation link.

We hope you will turn to ASPE for any future learning needs on how to put emergent technologies to practical use!



NOTES:
