# A Conversation with …
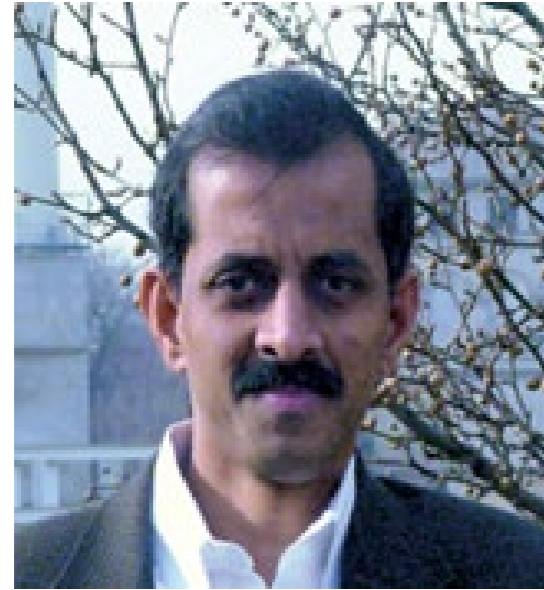## Ganesh Mani

# *Managing the Data Supply Chain*

Ganesh Mani Adj. Faculty Carnegie Mellon

Mehrzad Mahdavi, Executive Director, FDP Institute

Kathy Wilkens, Senior Advisor, FDPI Curriculum

Mirjam Dekker, Project Manager, FDP Institute

www.fdpinstitute.org
April 7, 2020

# Agenda



Ganesh Mani
Adj. Faculty
Carnegie
Mellon

Katherine
Wilkens
Sr. Curriculum
Advisor FDPI

Mehrzad Mahdavi
Executive Director
FDPI

Mirjam
Dekker
PM
FDPI

Download the thought leadership paper here

# By Region View (as of late April 4th)



New deaths in Lombardy and Madrid are flattening off, but the death toll is climbing ever faster in New York and London

Daily deaths with coronavirus (7-day rolling average), by number of days since 3 daily deaths first recorded

FINANCIAL TIMES

Sources: NHS; Covid Tracking Project; Providencialdata19; Santé Publique France; Berliner Morgenpost; OpenZH; Stockholm University; Leuven University. Data updated April 04, 19:00 GMT
FT graphic: John Burn-Murdoch / @jburnmurdoch
© FT

**Source: Financial Times**

# US Demand: Travel vs. Grocery
(Late Jan – Early Mar)



**Source: Exabel**

# Social Distancing Compliance



Where people were still traveling

Percent change in average travel for the week of March 23, compared with travel before the coronavirus outbreak.

No travel — Half of normal — Normal travel
Closer to normal travel →

**Source: NYTimes & CubeIq**

# Health



CDC Hospitalization Data By Underlying Health Conditions

Note: From April 3 update. Data through March 38.
Study based on 7,162 cases with completed information.
Source: Centers for Disease Control and Prevention

**Source: Barron's**



Chinese fatality data (by age)

**Source: (Mar 30, 2020) Lancet, Verity et al.**

Carnegie Mellon University

# Ganesh Mani, PhD, MBA
# Adjunct Faculty

AI / ML:
Useful for Analysis as well as for Data scrubbing / repair.

Ganesh has designed and developed many innovative alpha-generation frameworks based on disparate data (both traditional and alternative) and AI/ML techniques. Clients include leading asset management firms (incl. hedge funds) and plan sponsors.

Alternative

Economy

Natasha

Industry

Company

4

# Reduce Complexity

AI / ML provides a rich toolbox

**Insights**
Actionable for PMs and traders

**Rules or other summaries**
Interpretability trade-offs

**Intermediate Representations**
Features (e.g., Temporal patterns, Cross-asset relationships, Snippets w/ affect)

**Data**
Large amounts often w/ noise

Data
Volume, Velocity & Variety
Provenance, Vintage & Noise-level
Traditional vs. Alternative

AI

Search                    Inference

Planning        Machine
                Learning
Augmented Intelligence
Man with bicycle analogy
                                        Heuristic
                    Deep            Simulation
                    Learning

Machine
Translation
                                Game Theory

Opinion vs. Fact
AI / ML provides strong opinions or hints        Optimization

Models
Choice should be dictated by data
characteristics and goals

Beware of overfitting
Careful use of cross-validation and out-of-
sample data

Diving for insights
Making information actionable

**Validate**

v

Carefully
Time-series data may pose
additional challenges

**Drift?**

Y

Test and Deploy

**Retrain**

Monitor errors and type!
Is new data similar to training data?

r

7

How often?
Many factors (e.g., lookback) go into
this decision

# Multiple models / solutions possible

**_Try to be creative!_**

Whether you use a simpler technique or a more complex one, try to expose insights.

**B** Behaves as expected
If not overfit / regime change!

**E** Explainable
Simpler may be better!

**S** Spurious correlations
More common than you'd think!

**T** Temporal aspects
May need special processing

# Q & A





www.fdpinstitute.org/webinars

# In Closing

➢New curriculum available at the end of April

➢Registration for the October 26 – November 8th exam opens May 10th

➢For a recent candidate webinar go to www.fdpinstitute.org/webinars

## Learn more about the FDP Institute at
### www.fdpinstitute.org