

Analytics in India a few trends – field notes from a layman

Hindol Basu

India presents its own unique challenges in leveraging the power of data

2

Focus on quantified impact leading to top and bottom line results



Focus on data quality and data management



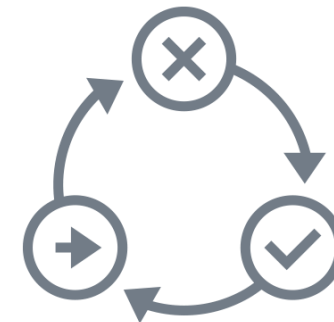
Balance between intuitive feel of models and higher predictive power



Communicating analytics across the organization



Adopt a “test & learn” approach



1. Start with KPIs and understand how they translate to financials and what are the benchmarks
2. Identify how benefits will be measured – test vs. control, pre vs. post, champion vs. challenger
3. Identify the change management requirement and manage the softer aspects of implementation
4. Start with glass box methodologies – so that business users have confidence
5. Move from glass box to black box methods (if needed)

The limitations of a data driven approach is often not understood properly ³

Analytics will solve all my problems

I need the most complex solutions
that I have Googled

Models will give me benefits,
irrespective of any process glitch

I don't need to understand what has
gone into your model, I just want
benefits

Truth lies
somewhere in the
middle

I know have been running this
business for the last 20 years, I
know what to do and how to do

Give me pivot tables and I can find
out the rest

I am fixing basic IT systems, we will
talk about models 3 years from noe

Show me the variables, their trends
and only if I think that all of them
make business sense, will I
implement the model

These gives rise to 4 key challenges

1. Difficulty in implementing relatively complex models (neural network, tree ensemble etc.)
2. Data quality issues, sparse data and lack of data
3. Change management required to push implementation
4. Sound methodology of tracking, attributing and demonstrating financial benefits

Difficulty in implementing complex models

5

Models with no or very limited interactions

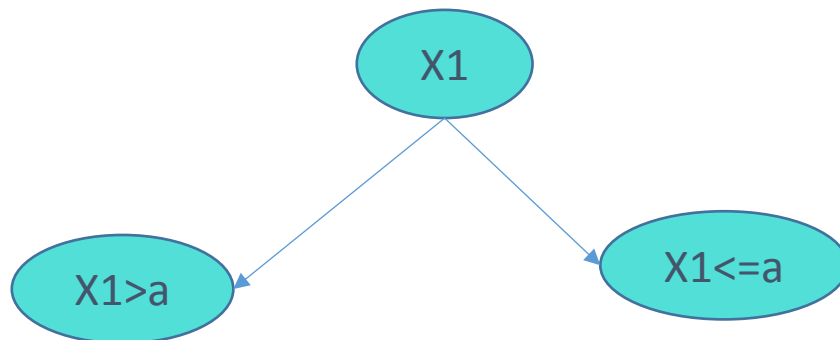
$$\frac{d(\text{Prediction})}{dX_1} = K \text{ (irrespective of } X_2 \dots X_n)$$

$$\frac{d(\text{Prediction})}{dX_1} = KX_2 \text{ (irrespective of } X_3 \dots X_n)$$

Complex interactions

$$\frac{d(\text{Prediction})}{dX_1} = K_{X_2 \dots X_n}$$

Simple decision trees



Complex interactions

- Use intuitive models in cases where current processes requires subject matter expert judgement
- Use complex modes in cases where only human intelligence needs to be automated
- Create sensitivity analysis to help unbox complex models
- Leverage intelligent methods of segmentation, tree interaction variables etc. to build power of interaction within simple models

Data quality – data is our best friend, poor data quality is our worst enemy ⁶

Models with no or very limited interactions

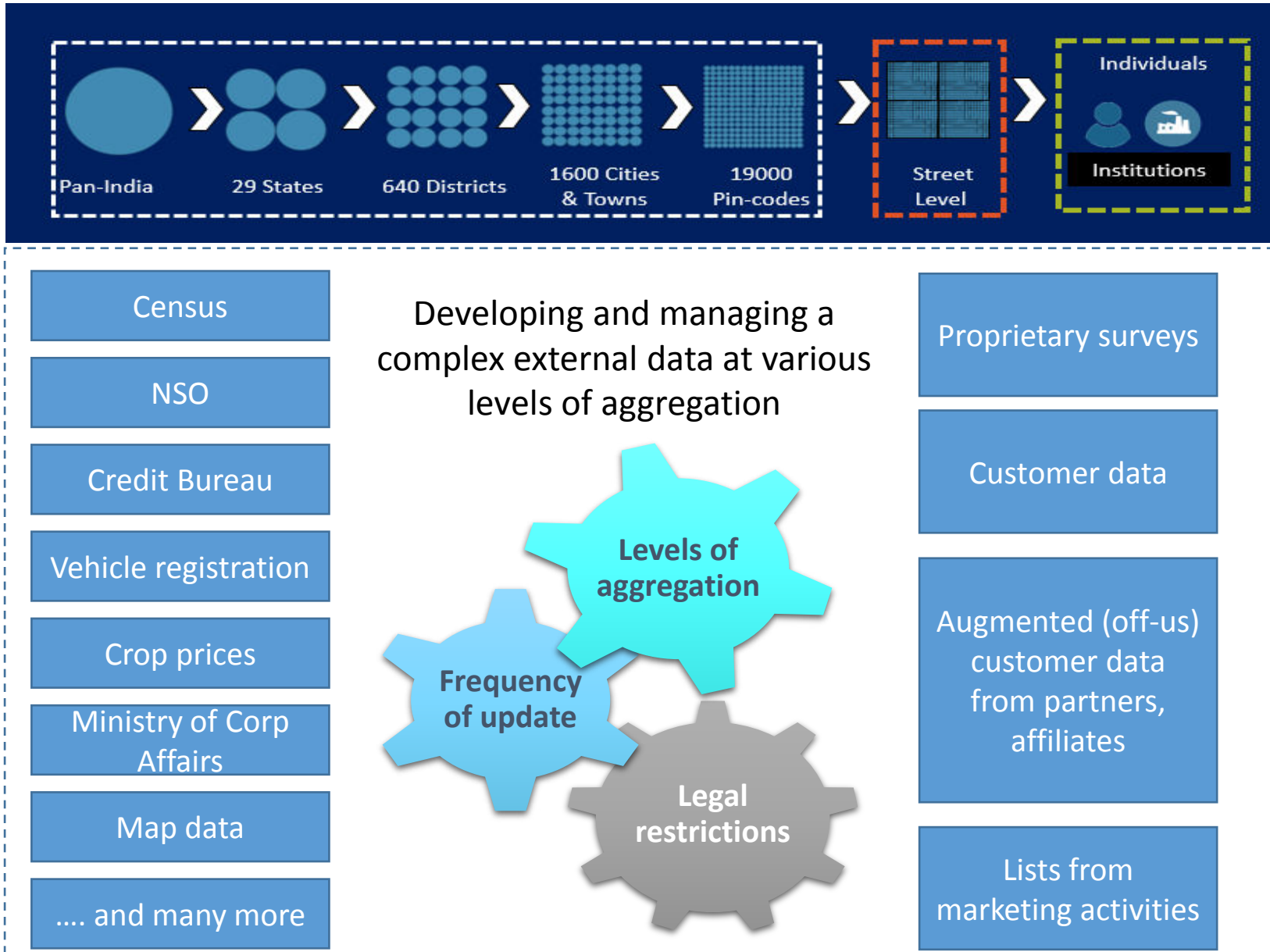
- Missing value
- Erroneous values
- Errors creeping into the data during summarization and storage in warehouse
- Lack of unique ID for the entity of interest (customer, supplier etc.)
- Lack of data capture
- Lack of history
- Issues in human data capture
- Lack of a proper data owner who can help

Models with no or very limited interactions

- Unique entity creation (unique customer, household, supplier etc.)
 - Particularly critical for non-financial services (in the absence of Aadhar, PAN etc.)
 - Fuzzy matching of name, address
 - Issues of multiple phone, email
- Value correlation analysis
- Identify if missings are part of data generation process
- Analytical methodologies for sparse data handling
 - Expectation maximisation
 - MCMC
- Non-human captured data: clickstreams, sensors etc. tend to have far lesser issues
- Plan for future data capture
- Takes about 75%-80% of project timeline

Key Initiative-1: Aggregate multiple sources of public external data

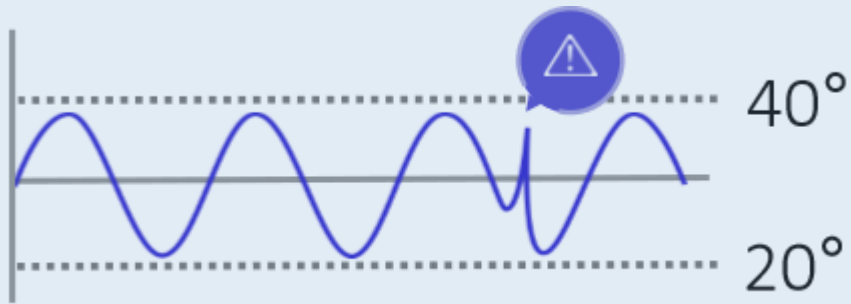
7



- Most organizations view external data only for the purpose of solving a specific problem
 - Store/branch location
 - Demand forecasting at micro geo level
 - Distribution planning
 - And others
- There may be a lack of concerted efforts in creating a comprehensive economic data repository by companies
- Indicus and Nielsen may be a good example

Key Initiative-2: Sensor data for manufacturing

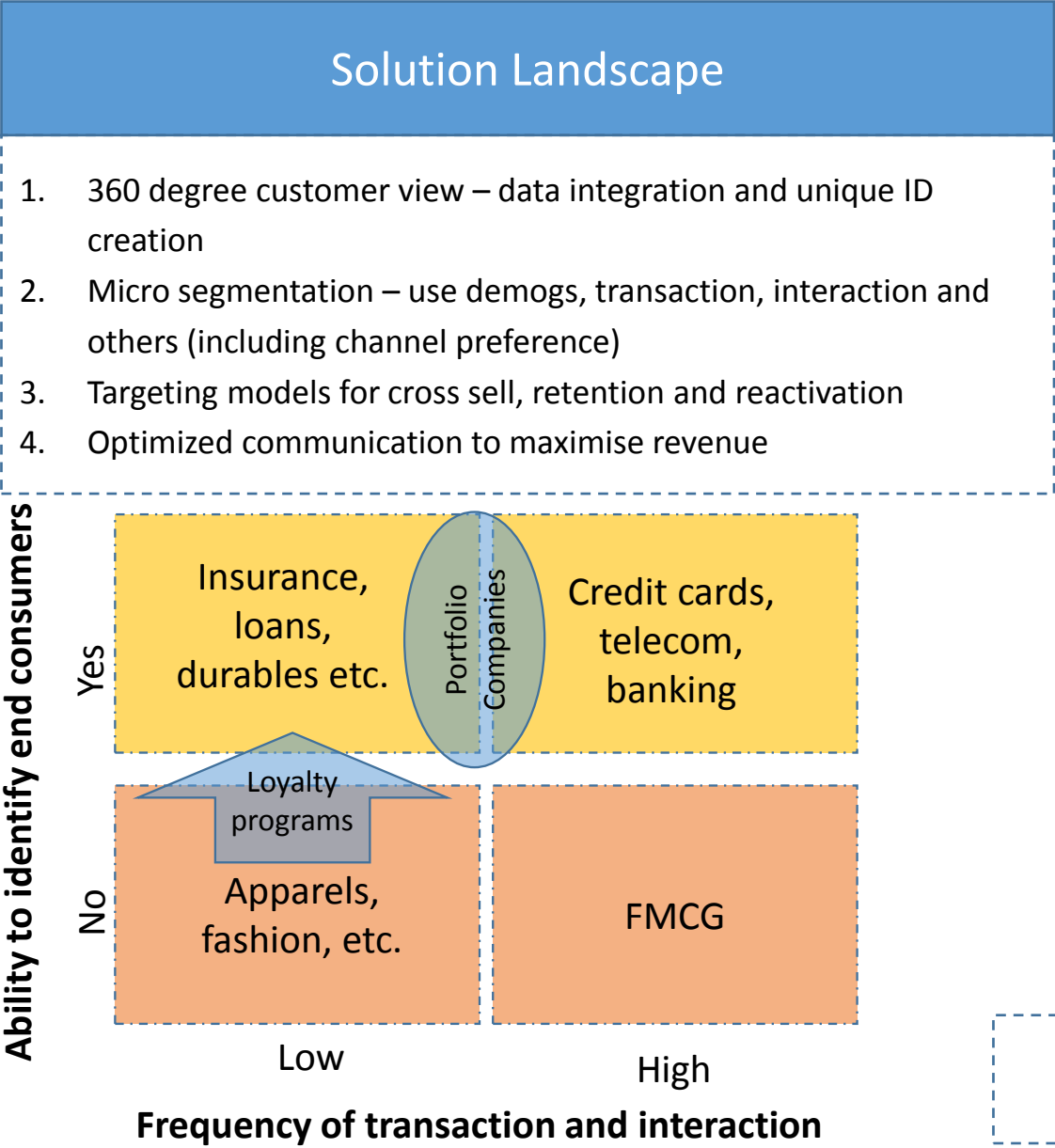
Operations managers have been using traditional control charts for ages



- Key sensors have been part of most equipment
- Univariate approach using only max and min
- Multi variate anomaly detection
- Availability of low cost of sensors

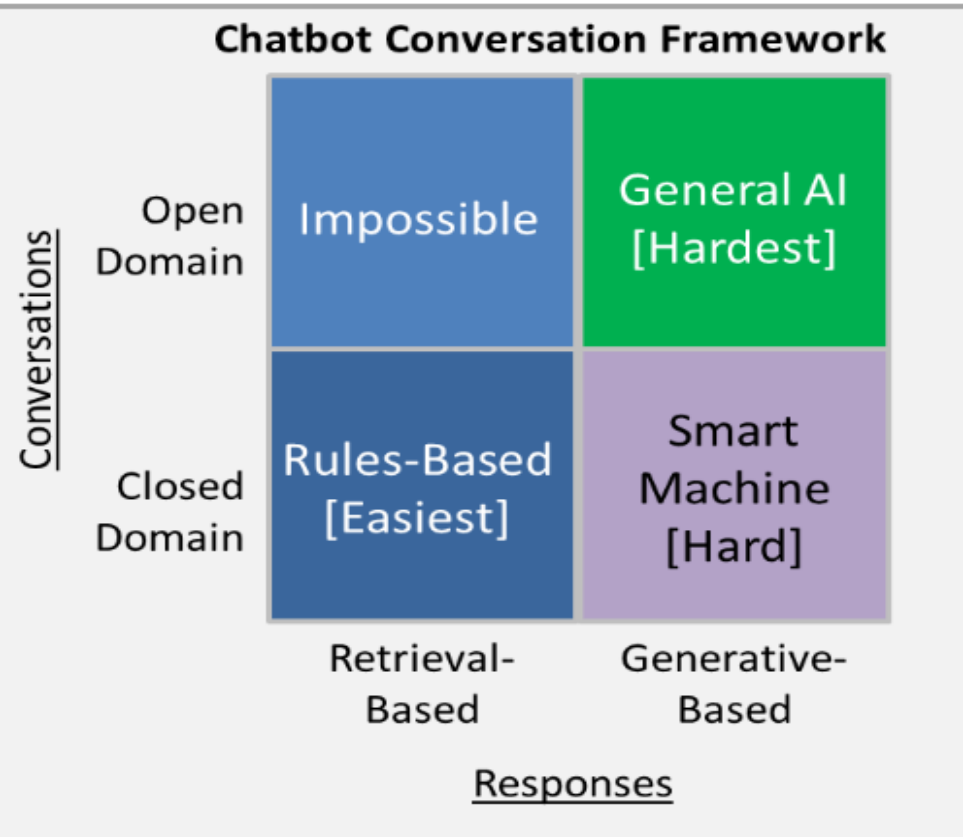
- Key use cases
 - Predictive Asset Management
 - Predictive Quality Management
- Unsupervised learning for anomaly detection
- Supervised learning for identifying quality issues, equipment failure
- Thresholding is another key problem
- Type-1 and Type-2 errors have different costs associated. Hence, meta cost analysis is important for thresholding
- The data is structured but of high frequency in update; hence, specific data engineering is important
- Linking operations technology to ERP, CRM etc.

Key Initiative-3: consumer companies trying to bring efficiency of marketing⁹



Managing the curse of false positives in churn models

Key Initiative-4: Automated customer interaction



- **Retrieval bots:** Retrieval bots do not generate any response on their own, they pick up one of the pre-defined responses from the thousands of available responses
- **Conversational bots:** are a sub set of retrieval bots, that uses business rules derived from expected flow of conversations
- **Generative bots:** these bots actually generate a response given the question
- **Importance of the context:** Specific retrieval rules (models) should be created for specific context so as to identify the right response
- **Usage of paid frameworks (e.g. Microsoft bot framework):** Tends to be very expensive in the long run. Loss of IP

Key Initiative-5: Explore new data sources

11

Images

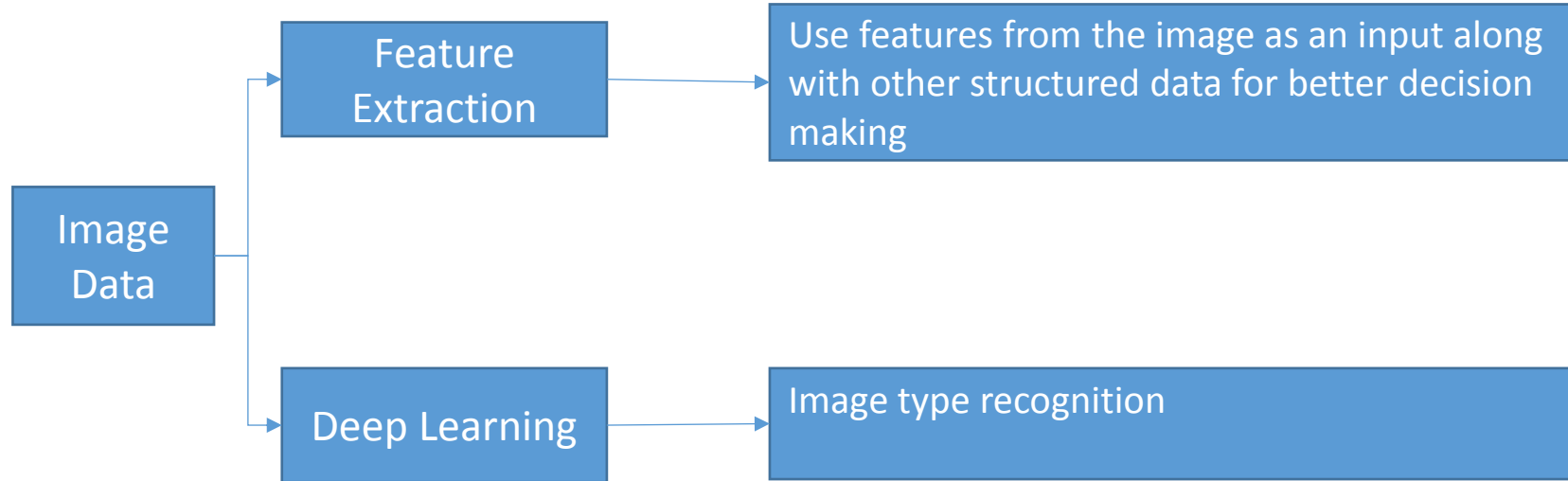
Scanned documents

Map data

Satellite image data

Video data

Wifi Latching data



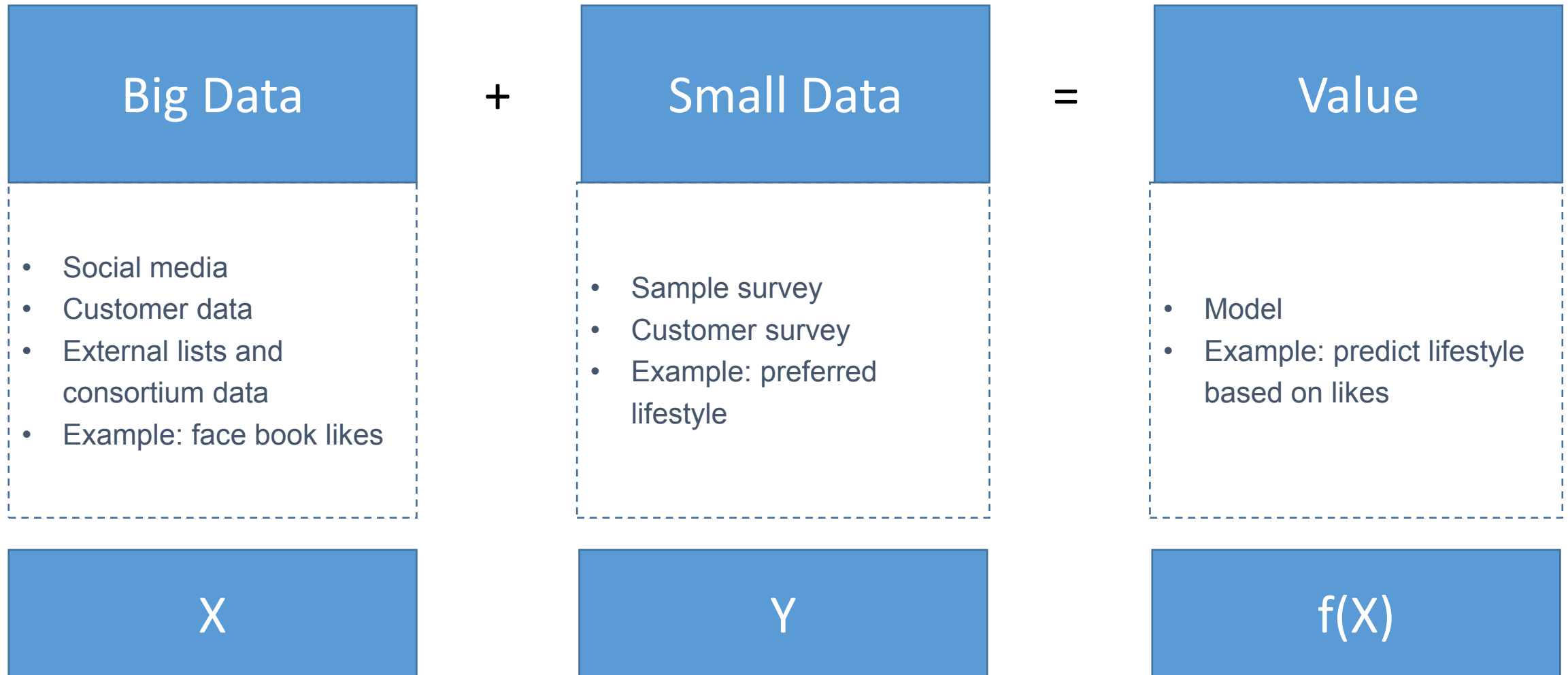
- Applications of satellite images for agricultural applications

- Physical click stream
- Movement within store
- Store fraud

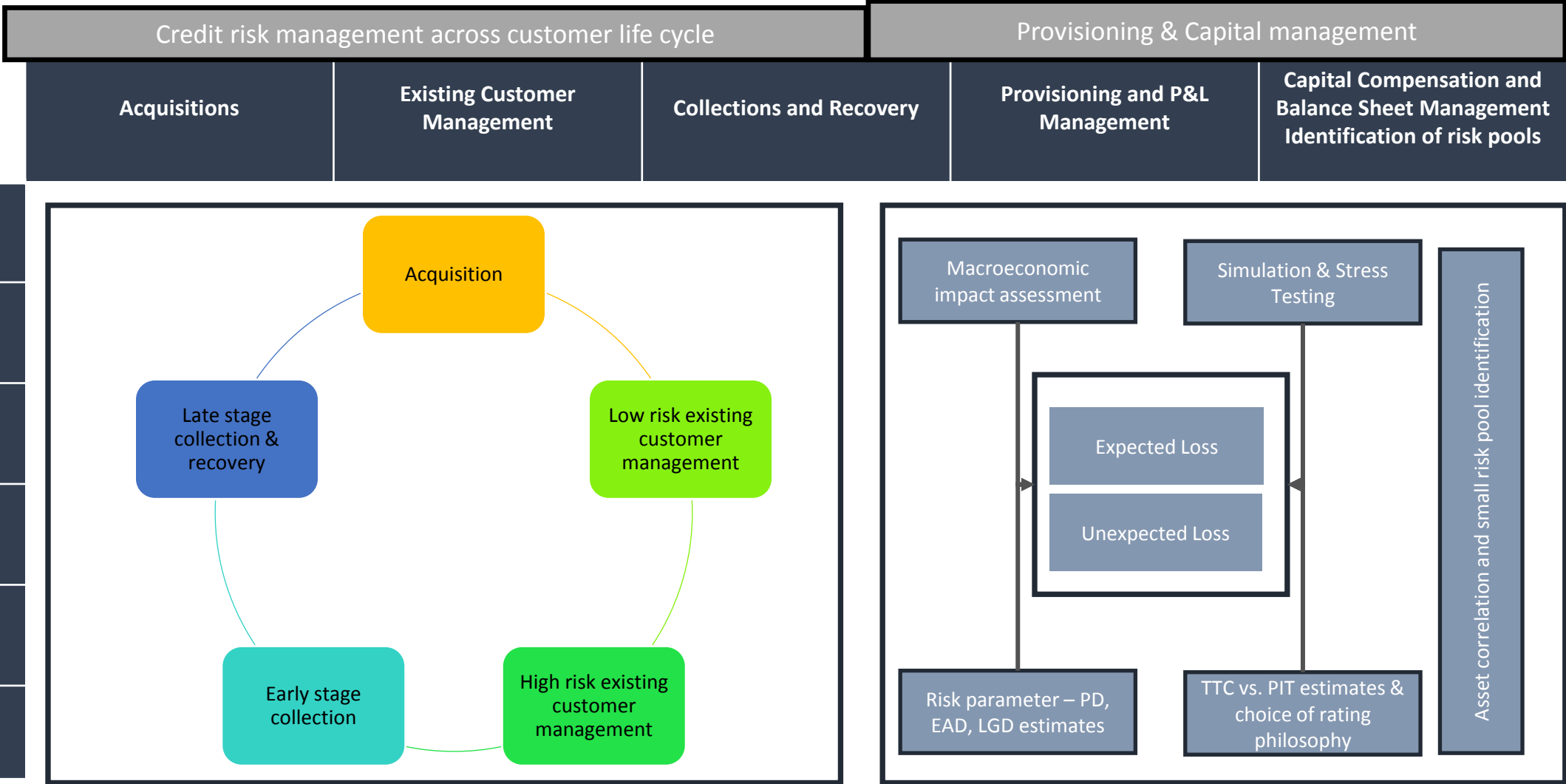
Key Initiative-6: GPS data for logistics management

- Availability of real time GPS data has very interesting applications for India
 - Avoiding toll roads
 - Pilferage
 - Idle time
 - And many more
- Stochastic nature of lead times
- Semi real time re-routing
- Driver behaviour monitoring
- Plugging in telematics for failure prediction

Key Initiative-7: Small data



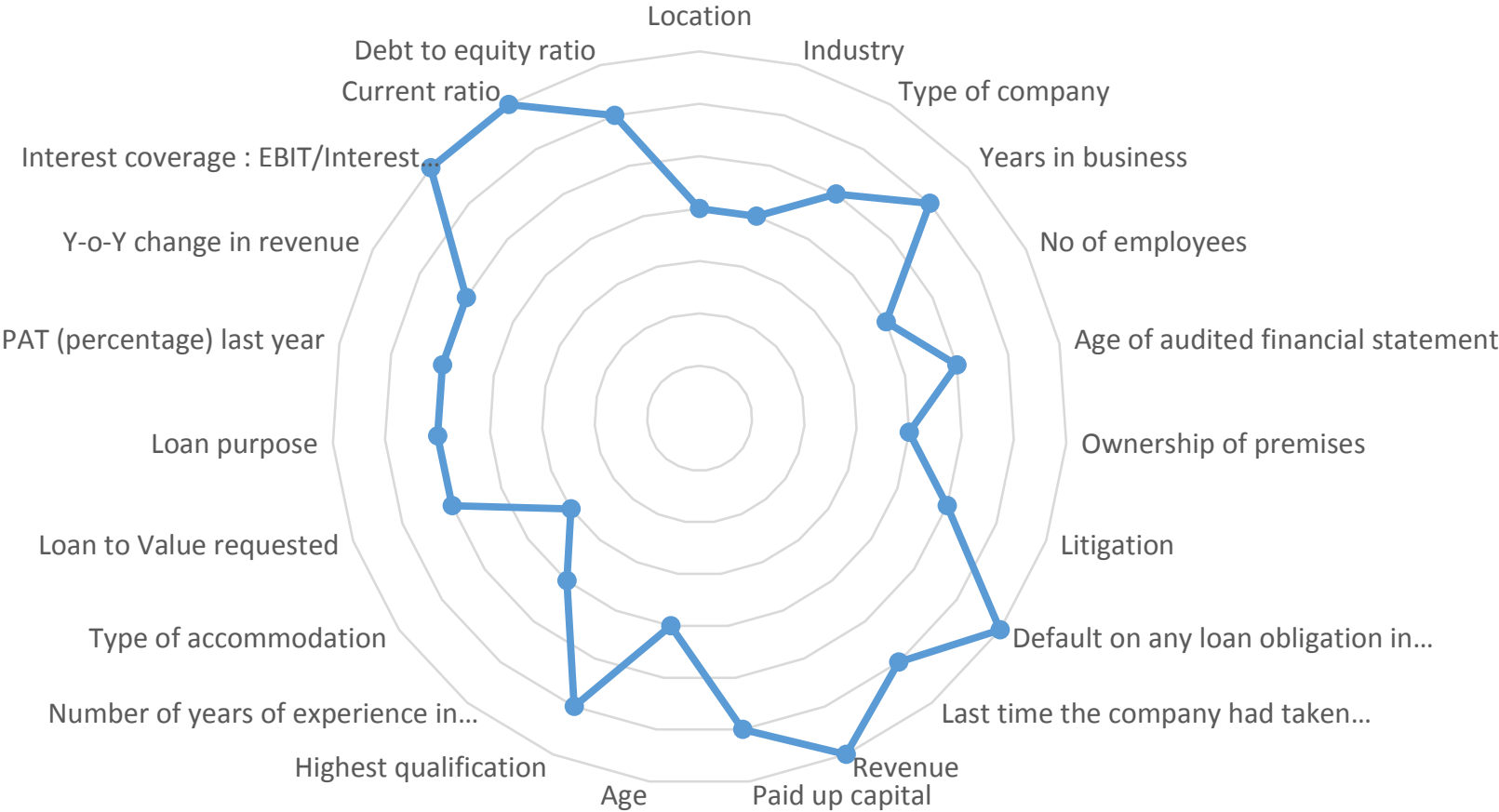
Key Initiative-8: Risk estimation



Key Initiative-8: Risk estimation (consumer and SMEs)

Risk distribution - SMEs

Risk Distribution Across All Risk Drivers



- Lack of financial statement
- Performance definition – grouping facilities, entities, customer group to account for risk spreading across groups
- Alternate source of data – EPF submission, bill payment, etc.