Essence: Pervasive & Distributed Intelligence

# BIG DATA THINNING: KNOWLEDGE DISCOVERY FROM RELEVANT DATA

MR NAJI SHEHAB

SUPERVISED BY: DR CHRISTOS ANAGNOSTOPOULOS
@ SAWB- F121; 21 MAR 2019

# THE PROBLEM

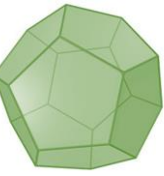| What is the problem? | • Building Predictive Models Over Big Data |
|---|---|
| Who has this problem? | • Data Scientists - Data Analysts – IOT Applications – Edge Centric Services |
| Why should this problem be solved? | • Big Data Is Big – It Can Be Thinned |
| The Solution ? | • Exploit & Learn Only The Most Interested & Important Subspaces |

# HYPOTHESIS 1

*"Upon quantizing a large-scale dataset, each respective cluster can be associated with an individual predictive local model. Given this, the local models should predict more accurately than its counterpart global model over the entire data space."*
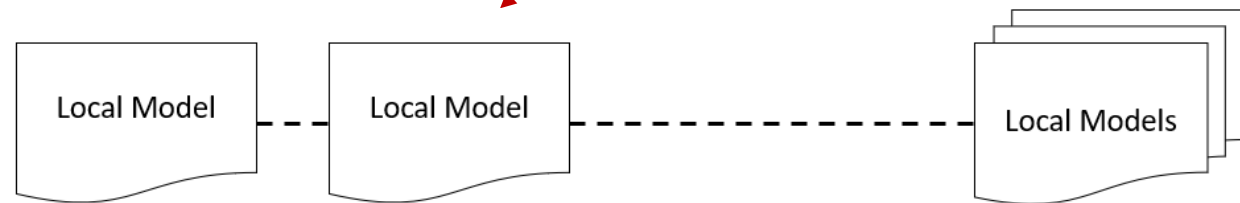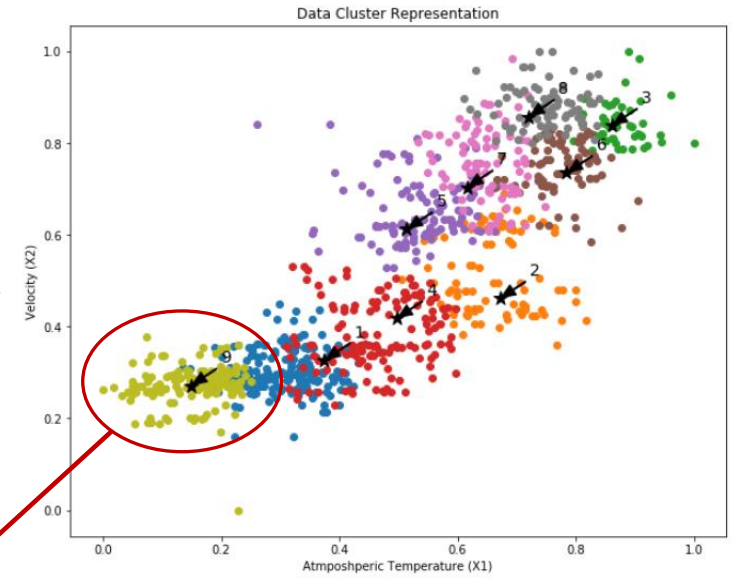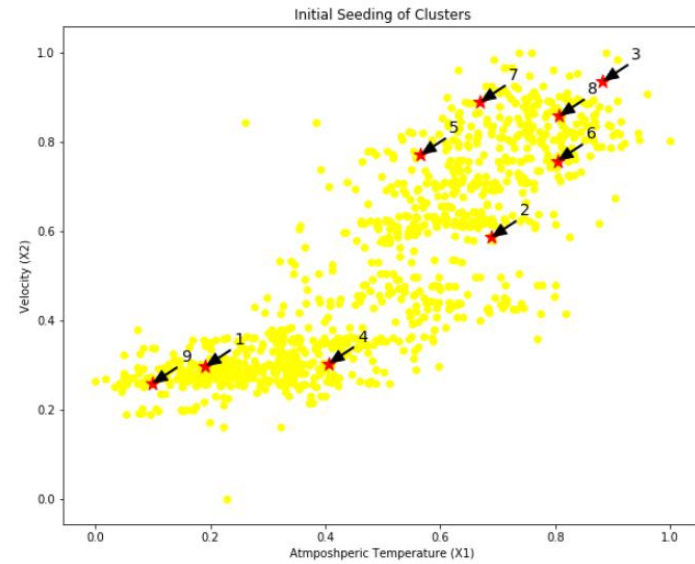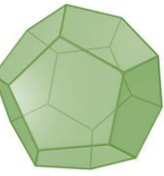
# HYPOTHESIS 2

*"Upon observing user queries analytics, we can learn and make use of the principles of the RPCL and Learning Automata methodologies in order to extract the distributions and associations of the queries over the large-scale data space. Given this and Hypothesis 1; we should be able to create respective local models over the past issued query spaces providing high quality analytics."*
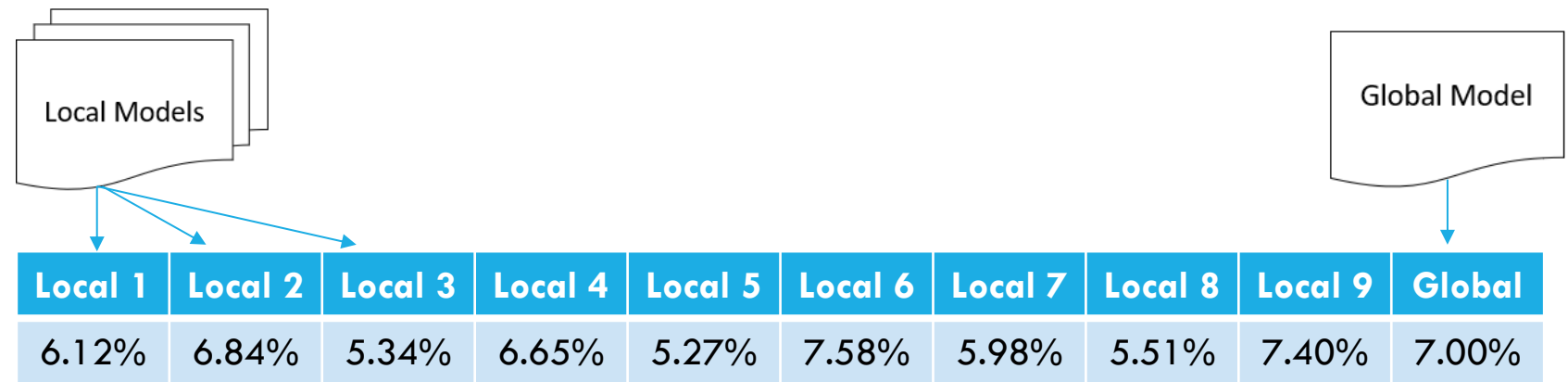
# FORMING LOCAL MODELS

HYPOTHESIS 1



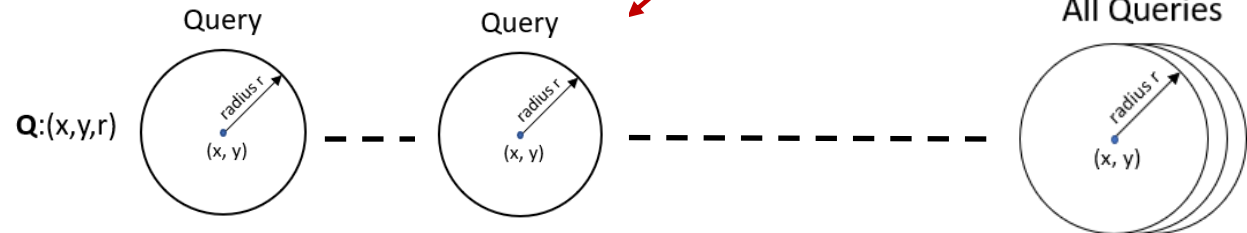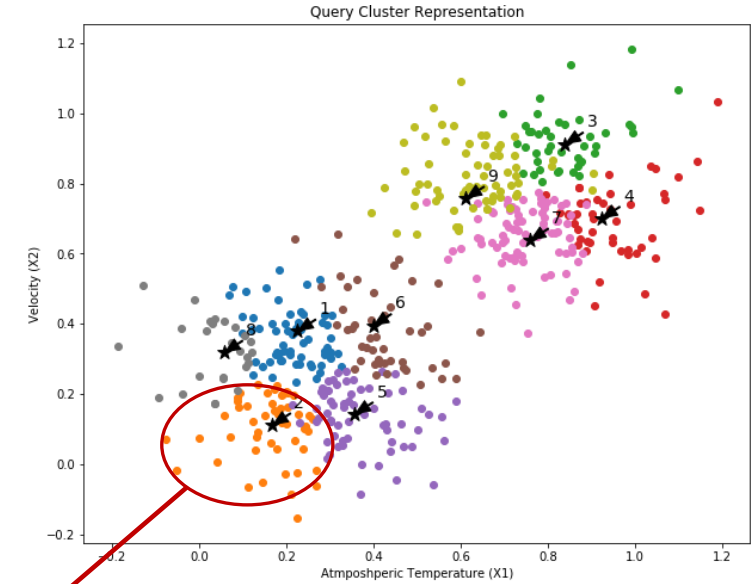Local Model - - - Local Model - - - Local Models

# RMSE: LOCAL VS GLOBAL

## HYPOTHESIS 1

Local Models

Global Model

| Local 1 | Local 2 | Local 3 | Local 4 | Local 5 | Local 6 | Local 7 | Local 8 | Local 9 | Global |
|---------|---------|---------|---------|---------|---------|---------|---------|---------|--------|
| 6.12% | 6.84% | 5.34% | 6.65% | 5.27% | 7.58% | 5.98% | 5.51% | 7.40% | 7.00% |

Each local model produces on average a **10%** increase in predictive performance with the use of less than **15%** of the entire dataset
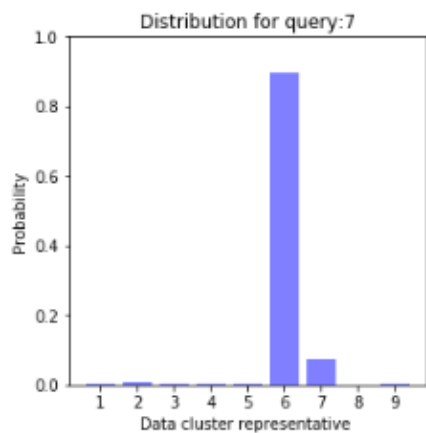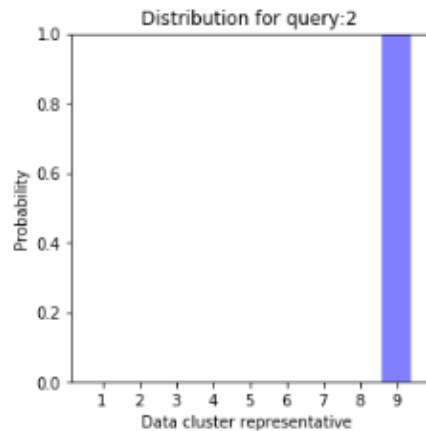
# ANALYSING USER QUERIES

## HYPOTHESIS 2
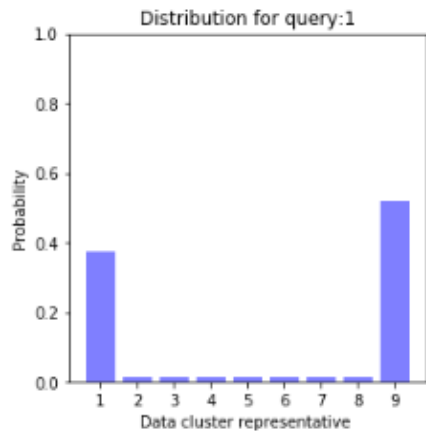
MAPPING USER QUERIES ONTO DATA
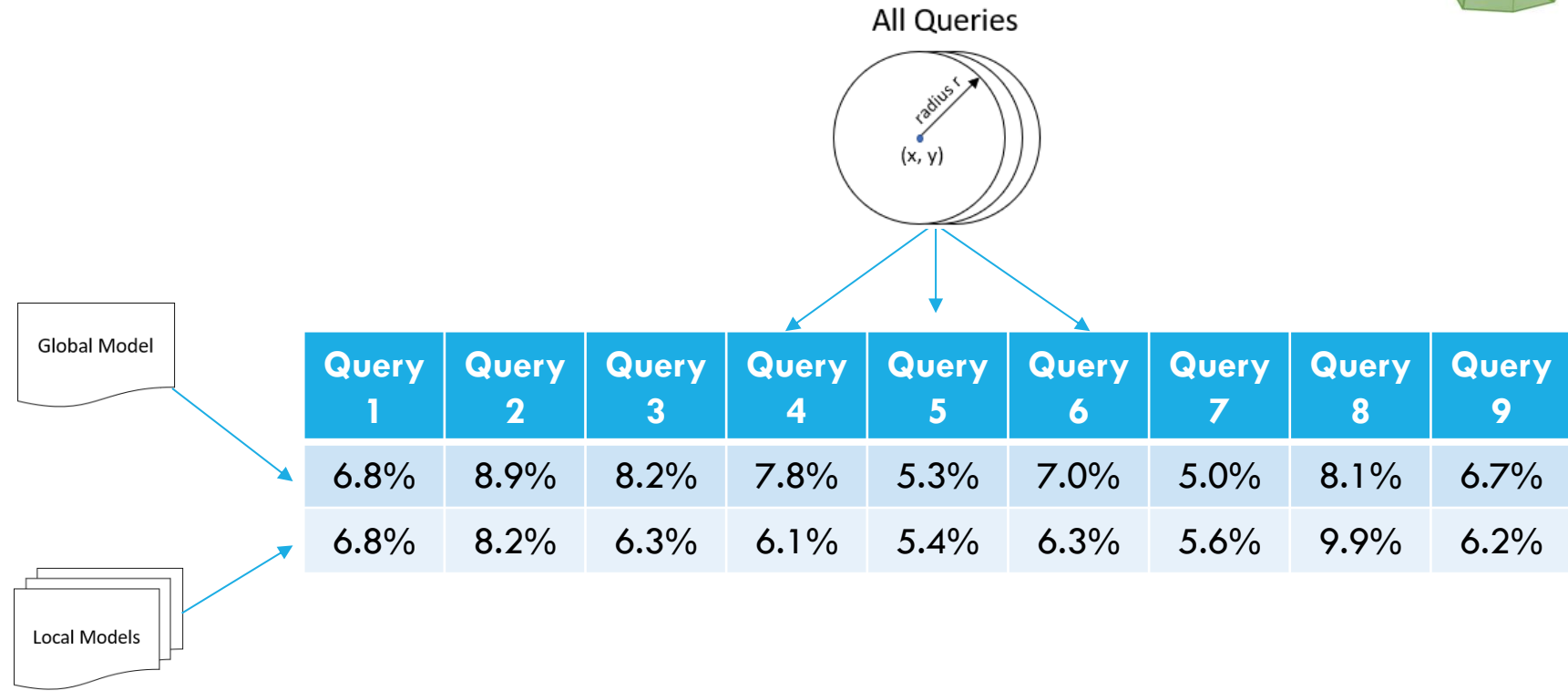
HYPOTHESIS 2

# CONCLUSION

- We were able to exploit important subspaces from query analytics

- Map each significant subspace into small predictive local models that:
  - Yielded higher predictive accuracy.
  - Made use of less data.

- Thus thinning the data space, and making use of less training data instances.