

Analytics on streaming data from Environmental Sensors

Rajrup Ghosh (M Tech CDS)

Shib Shankar Das (M Tech ECE)

Streams Everywhere

- Internet and Social Networks generate streams of posts, hashtags, videos, etc
 - Twitter, Facebook, Internet packets
- Business
 - Stock market prediction
 - Monetary Transactions
- Cybersecurity
 - Telecom call logs,
 - financial transactions, Malware
- Internet of Things
 - Smart Transport, Power and smart grids
 - Weather forecast and Pollutant monitoring
 - Smart phone, Health appliances, TV

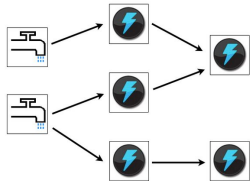
Apache Storm

- Originally developed by Nathan Marz at Backtype/Twitter
- Distributed, fault-tolerant stream-processing platform

[Distributed and fault-tolerant realtime computation](#) [about](#) [documentation](#) [blog](#) [downloads](#) [community](#)






Apache Storm is a [free and open source](#) distributed realtime computation system. Storm makes it easy to reliably process unbounded streams of data, doing for realtime processing what Hadoop did for batch processing. Storm is [simple](#), can be used with [any programming language](#), and is a lot of fun to use!

Storm has many use cases: realtime analytics, online machine learning, continuous computation, distributed RPC, ETL, and more. Storm is fast: a benchmark clocked it at over **a million tuples processed per second per node**. It is [scalable](#), [fault-tolerant](#), [guarantees your data will be processed](#), and is [easy to set up and operate](#).



Storm [integrates](#) with the queueing and database technologies you already use. A Storm topology consumes streams of data and processes those streams in arbitrarily complex ways, repartitioning the streams between each stage of the computation however needed. Read more in [the tutorial](#).

Companies & Projects Using Storm





STORM

<http://storm.apache.org/>

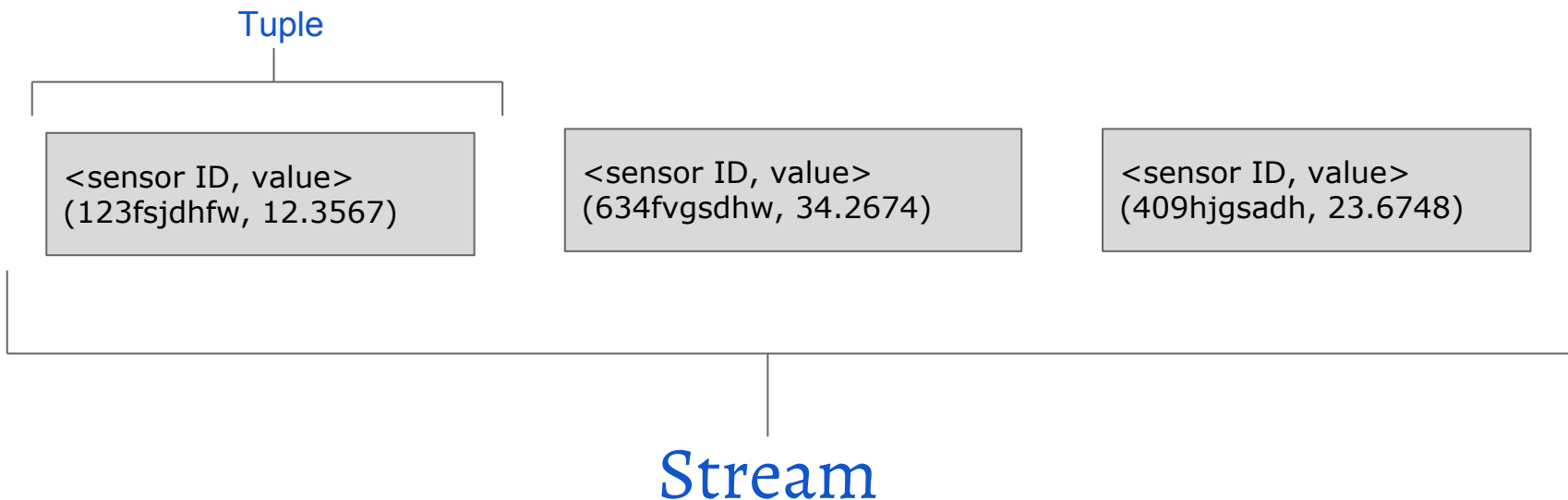
Storm Concepts

- Tuples and Streams
- Spouts, Bolts, Topologies
- Tasks and Workers
- Stream Grouping

<http://storm.apache.org/releases/current/Tutorial.html>

Tuples and Streams

- Tuple: ordered list of elements
- Stream: unbounded sequence of tuples



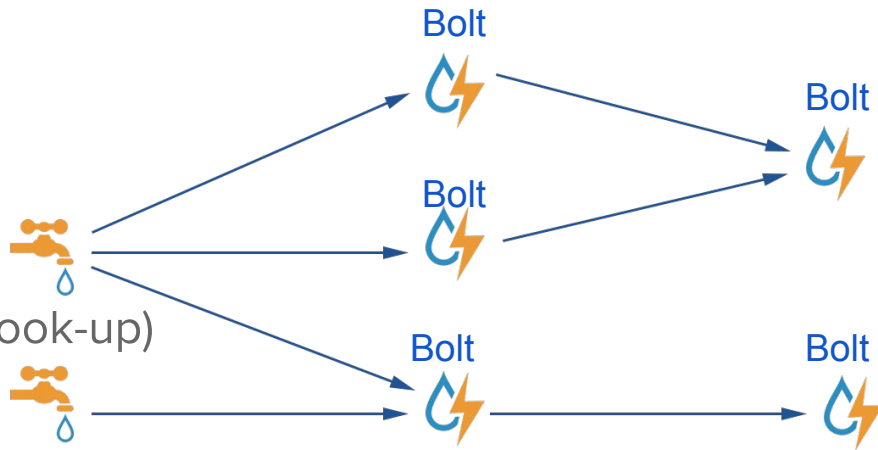
Spouts

- The sources of streams
- Can talk with
 - Queues (Kafka, Kestrel, etc.)
 - Web logs
 - API calls
 - Filesystem (MapReduce-FS / HDFS)



Bolts

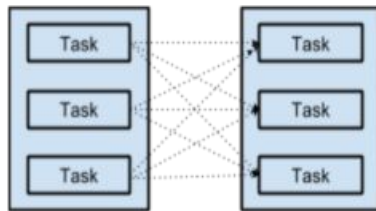
- Process tuples and create new streams
- Implement business logic via ...
 - Transform
 - Filter
 - Aggregate
 - Join
 - Access datastores & DBs
 - Access APIs (e.g., geo location look-up)



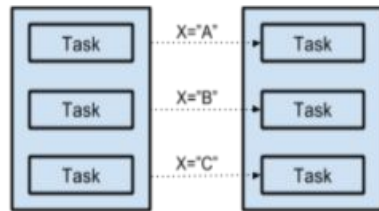
TOPOLOGY - Directed graph of spouts and bolts

Stream Grouping

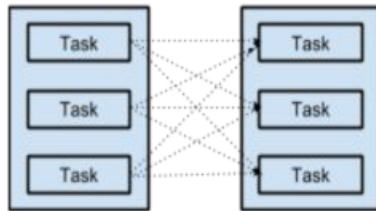
- Shuffle grouping: tuples are randomly distributed across all of the tasks running the bolt
- Fields grouping: groups tuples by specific name field and routes to the same task



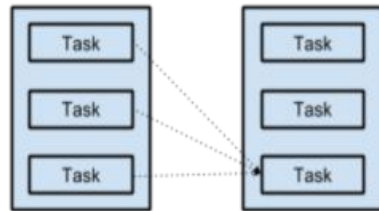
Shuffle Grouping



Fields Grouping



All Grouping

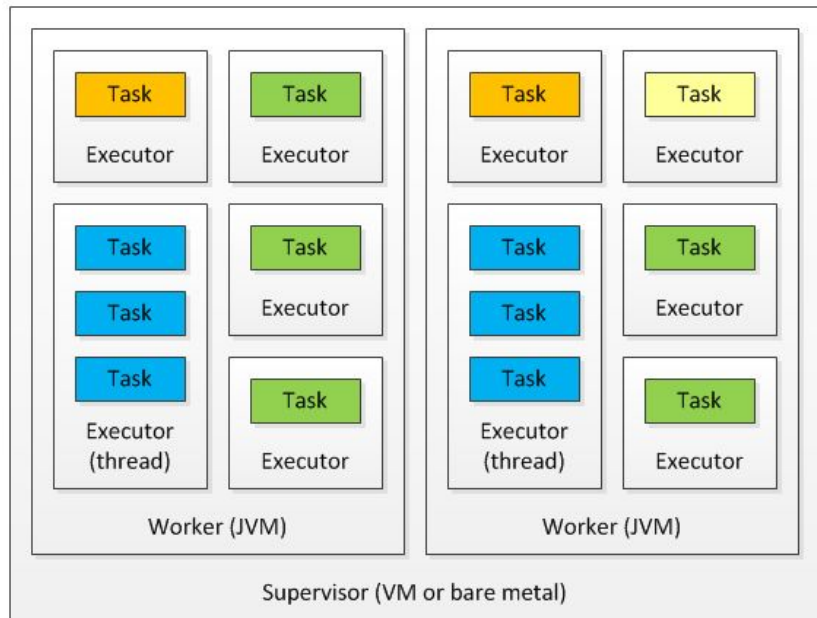


Global Grouping

Storm Groupings

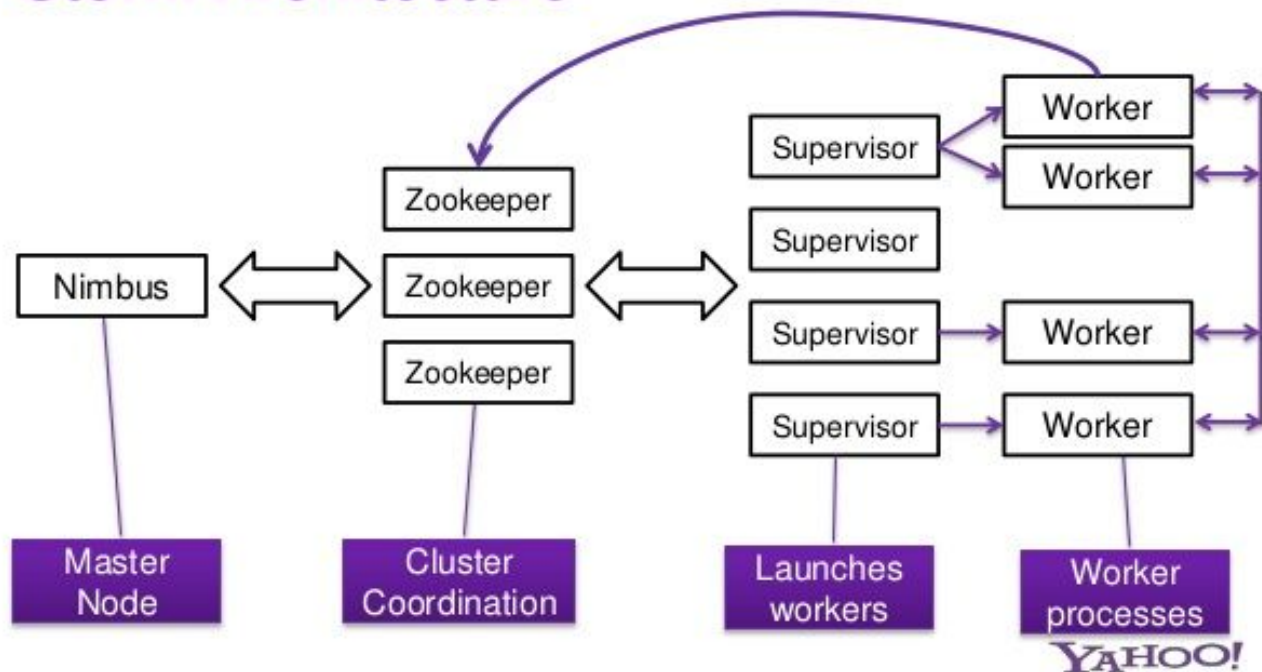
Tasks and Workers

- **Task:** each spout/bolt executes as many threads of execution across the cluster
- **Worker:** a physical JVM that executes a subset of all the tasks for the topology



Storm Architecture

Storm Architecture



Storm Architecture

DataStreamGeneratorSpout

FilterBolt : Kalman Filter

BuildParameterBolt (mean, standard deviation, second moment, third moment)

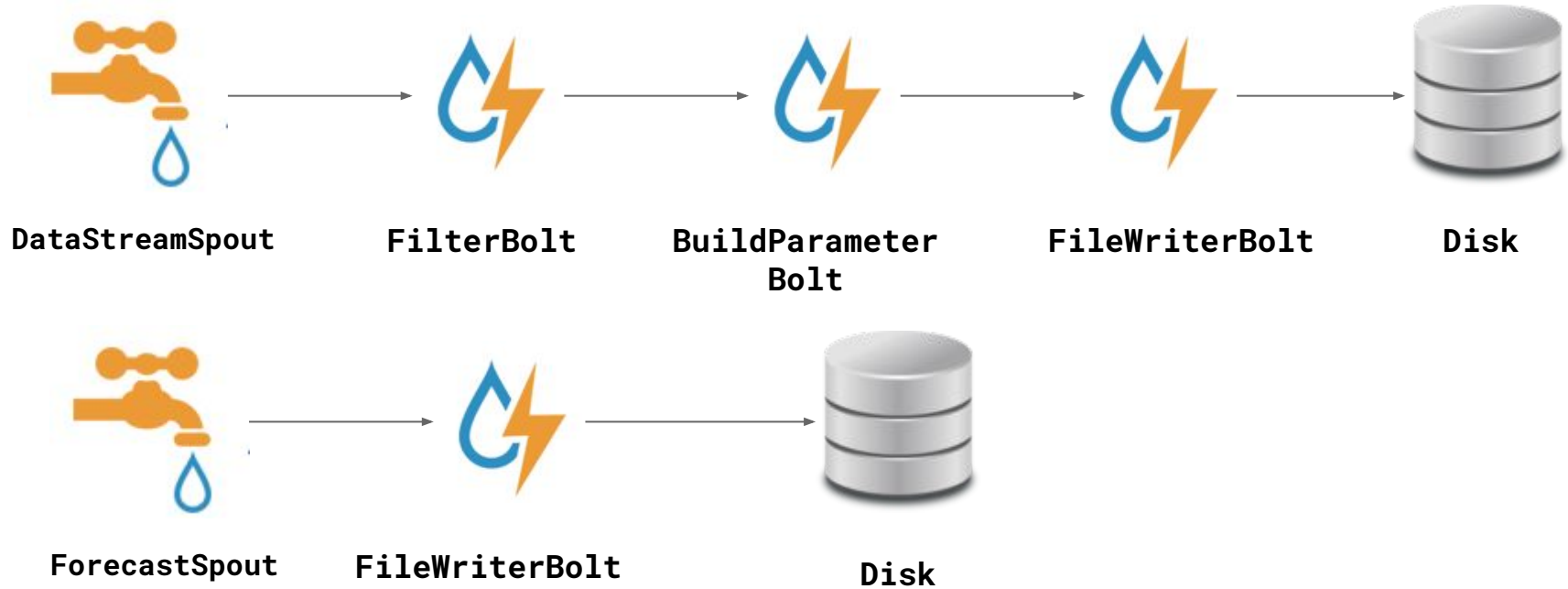
Used Alon-Matias-Szegedy Algorithm for moment estimation

FileWriterBolt

DataEstimatorSpout

DataWriterBolt

Topology



Alon-Matias-Szegedy Algorithm

Pick random element

Store the element with a count 1

If existing element found in further data increase count

Estimated Second moment = $n \cdot (2 \cdot X.\text{value} - 1)$

Estimated Third moment = $n \cdot (3 \cdot v^2 - 3 \cdot V + 1)$

Forecasting

- Forecasting is a process of estimation of an unknown event/parameter such as demand for a product, amount of rainfall, etc.
- Forecasting is commonly used to refer timeseries data.
- Time Series is a sequence of data points measured at successive time intervals.
- Time series analysis helps to identify and explain:
 - Any systematic variation in the series of data which is due to seasonality.
 - Cyclical pattern that repeat.
 - Trends in the data.
 - Growth rates in the trends.

Additive and Multiplicative Models

1. Additive Forecasting Model


$$Y_t = T_t + S_t + C_t + R_t$$


Diagram illustrating the Additive Forecasting Model components: Trend, Seasonality, Cyclical, and Random, connected by lines to the equation $Y_t = T_t + S_t + C_t + R_t$.

2. Multiplicative Forecasting Model


$$Y_t = T_t \times S_t \times C_t \times R_t$$


Diagram illustrating the Multiplicative Forecasting Model components: Trend, Seasonality, Cyclical, and Random, connected by lines to the equation $Y_t = T_t \times S_t \times C_t \times R_t$.

Time Series Techniques

- Moving Average.
 - Simple moving average
 - Weighted moving average

$$F_{t+1} = \frac{1}{n} \sum_{i=t-n+1}^t Y_i$$

F_{t+1} = Forecast for period $t + 1$

- Exponential Smoothing.

Y_i = Data corresponding to time period i

$$F_t = \alpha Y_{t-1} + \alpha(1-\alpha)Y_{t-2} + \alpha(1-\alpha)^2 Y_{t-3} + \dots$$

- Auto-regression Models (AR Models)
- ARIMA (Auto-regressive Integrated Moving Average) Models

AR(p), MA(q) and ARMA(p, q)

- Auto-Regressive Process: AR(p) process models each future observation as a function “p” previous observations. If $\{Y_t\}$ is purely random with mean zero and constant standard deviation σ (White Noise). Then the autoregressive AR(p) process:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t$$

- Moving Average Process: MA(q) models each future observation as a function of “q” previous errors. If $\{Y_t\}$ is a moving average process of order q (MA(q)) if for some constants $\beta_0, \beta_1, \dots, \beta_q$

$$Y_t = \beta_0 + \beta_1 \varepsilon_t + \beta_2 \varepsilon_{t-2} + \dots + \beta_q \varepsilon_{t-q} + \varepsilon_t$$

- ARMA(p,q) model is a combination of AR(p) and MA(q) process, given by -

$$Y_t = \alpha_0 + \alpha_1 Y_{t-1} + \alpha_2 Y_{t-2} + \dots + \alpha_p Y_{t-p} + \\ + \beta_0 + \beta_1 \varepsilon_{t-1} + \beta_2 \varepsilon_{t-2} + \dots + \beta_q \varepsilon_{t-q} + \varepsilon_t$$

ARIMA (p, d, q)

ARIMA has the following three components:

- Auto-regressive component (p): Function of past values of the time series.
- Integration Component (d): Differencing the time series to make it a stationary process.
- Moving Average Component (q): Function of past error values.
- The q and p values are identified using auto-correlation function (ACF) and Partial auto-correlation function (PACF) respectively. The value d identifies the level of differencing.
- Box-Jenkins Methodology is used to identify the model parameters.

Measures of aggregate error

Mean absolute error MAE	$MAE = \frac{1}{n} \sum_{t=1}^n E_t $
Mean absolute percentage error MAPE	$MAPE = \frac{1}{n} \sum_{t=1}^n \left \frac{E_t}{Y_t} \right $
Mean squared error MSE	$MSE = \frac{1}{n} \sum_{t=1}^n E_t^2$
Root mean squared error RMSE	$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n E_t^2}$

$$E_t = F_t - Y_t$$

What we have used?

- R forecast package has been used (developed by Hyndman & Khandakar (JSS, 2008). We have used JRI for R function calls from Java.
- We have used `auto.arima()` function to estimate ARIMA(p, d, q) parameters, with learning seasonal parameters as well
- We have trained a part of the dataset, and used multi-step forecast with re-estimation for the test data in streaming environment.

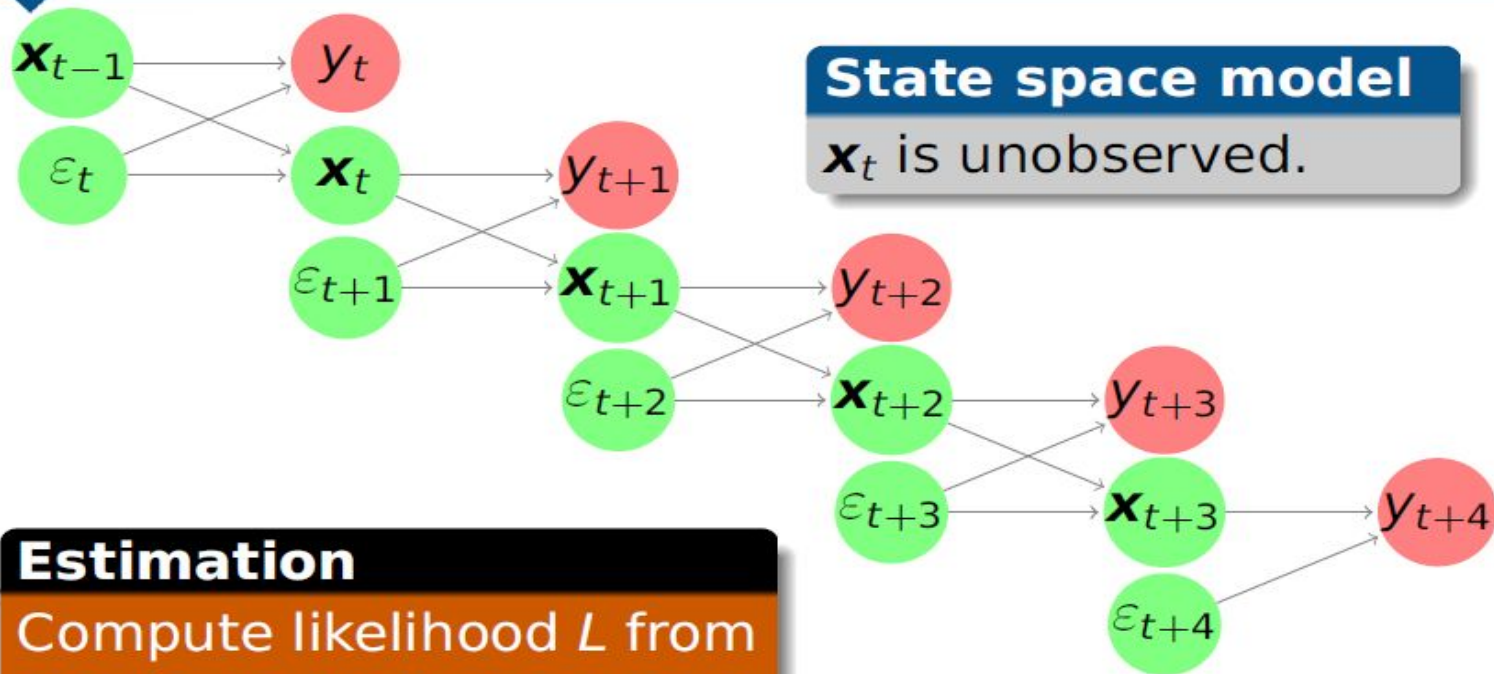
Test Code

```
t <- ts(temperature, start= start_time, deltat = 1)
train <- window(t, end=end_train)
test <- window(t, start=start_test)
fit <- auto.arima(train)
order <- arimaorder(fit)
```

Train Code

```
fc <- ts(numeric(length(test)), start=start_test, deltat = 1)
for(wind in 1:num_retrain)
{
  x <- window(temperature, end=end_train + wind_size*(wind-1))
  refit <- Arima(x, model = fit)
  temp_forecast <- forecast(refit, h=10)$mean
  for(i in 1:wind_size)
  {
    fc[wind_size*(wind-1)+i] <- temp_forecast[i]
  }
}
```

Time series forecasting



Estimation

Compute likelihood L from

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T$.

Use optimization

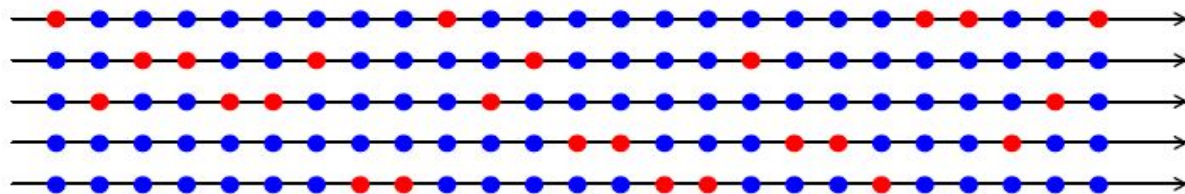
algorithm to maximize L .

Time Series Cross-validation

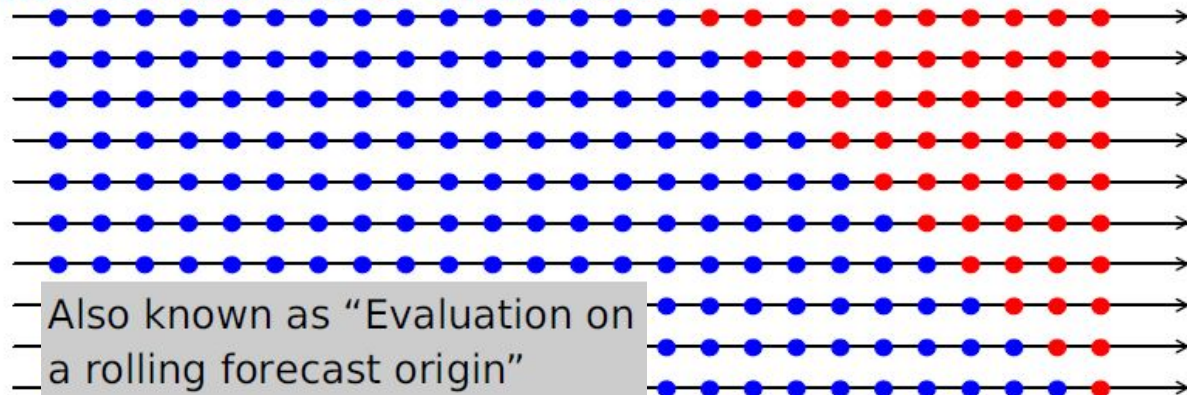
Traditional evaluation



Standard cross-validation



Time series cross-validation



Results

Comparison between Actual Temperature and Predicted Temperature



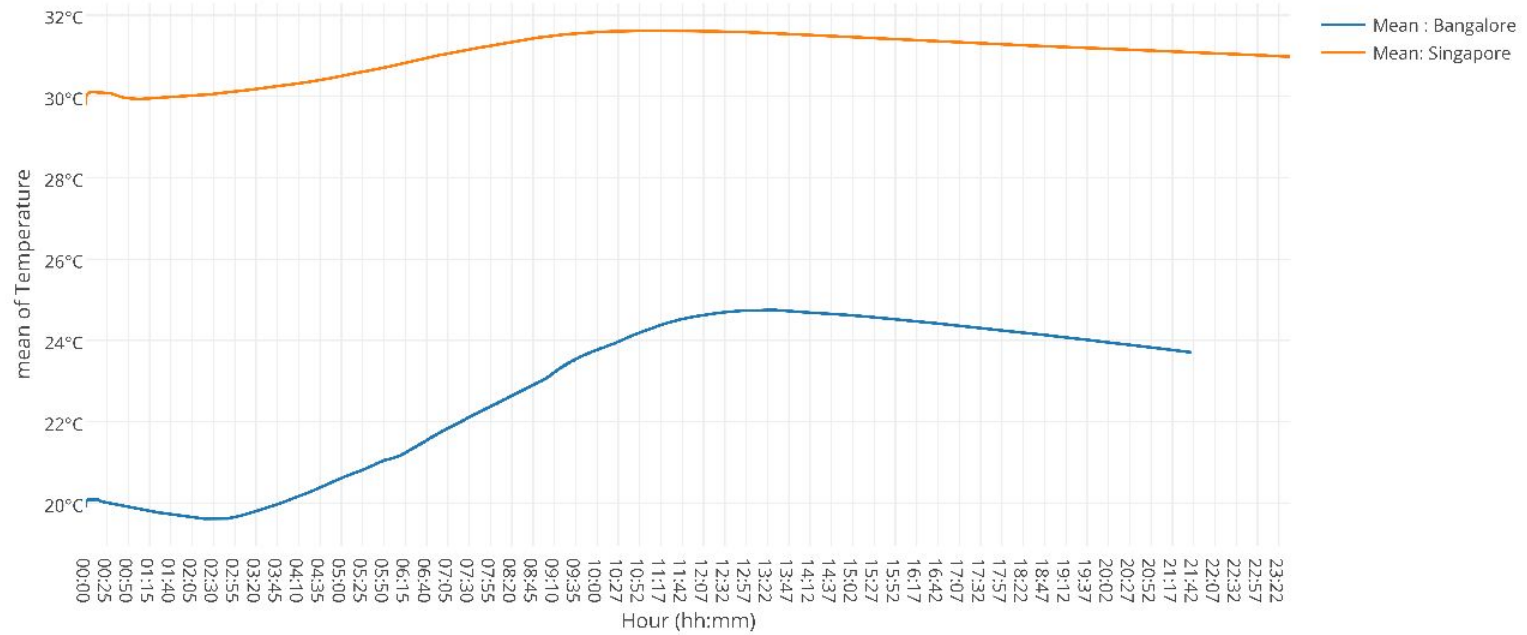
RMSE: 0.124

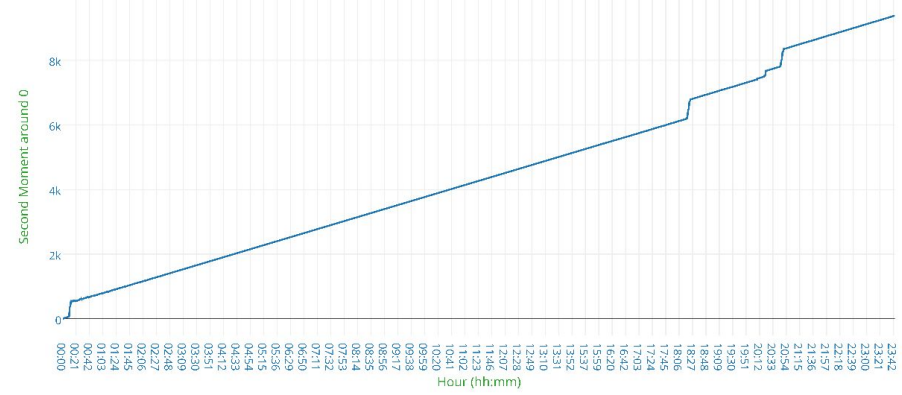
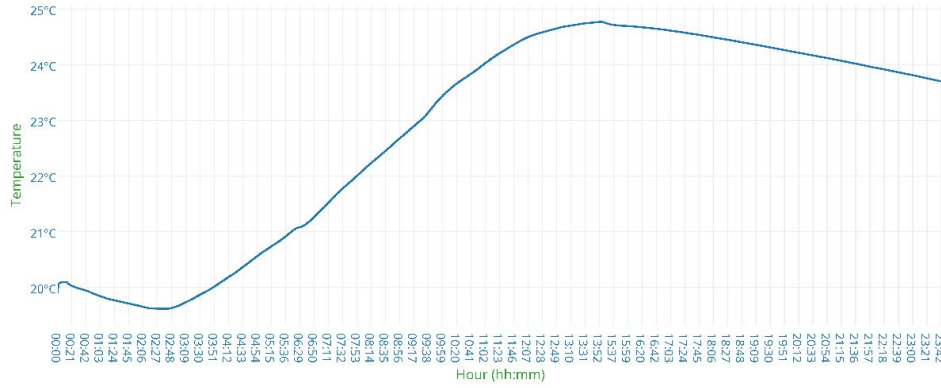
MAPE: 0.310

MAE: 0.050

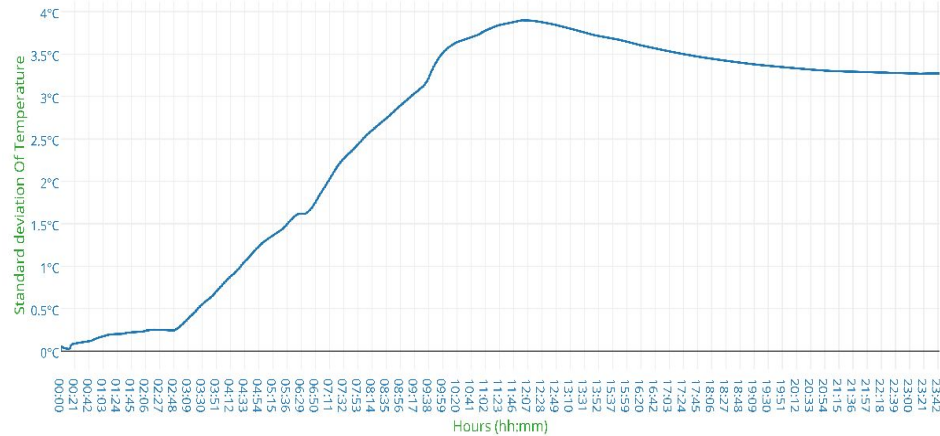
Dataset: Data
Canvas Sensor
Data [2]

Mean of Bangalore and Singapore





Comparison of Mean, Standard deviation, second moment about 0 (in degree Celsius) for temperature in Bangalore



Conclusion

- In this project we have learned real-time analytics over high velocity streaming data
- Apache Storm scaled nicely with different input stream rates for performing statistical analysis on the sensor data.
- It can be observed from the results obtained that forecasting for weather data worked very well using `auto.arima()` with periodic retraining.
- All these study has a huge application in weather stations and pollution monitoring centers for dealing with numerous sensor data and is of a growing importance in developing future “SMART CITIES”.

Future Work

- Due to lack of JAVA based implementation of ARIMA, we used R packages which have some overhead on end to end latency for streaming applications.
- ARIMA and `auto.arima()` can be implemented in JAVA so the stream application can better scale with forecasting, decreasing the end to end latency as a whole.
- Another important study may involve comparison with different forecasting methods and how they scale with streaming applications.
- Deploy a statistical package (maybe query based) for real-time analytics of streaming data coming from weather stations and pollution monitoring centers.

References

1. A. Toshniwal, S. Taneja, A. Shukla, K. Ramasamy, J. M. Patel, S. Kulkarni, J. Jackson, K. Gade, M. Fu, J. Donham, N. Bhagat, S. Mittal, and D. Ryaboy, “Storm@twitter,” in Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data
2. D. Canvas, “Data Canvas Dataset,” url: <http://datacanvas.org/sense-your-city/>, 2015..
3. A. Rajaraman and J. D. Ullman, Mining of Massive Datasets. New York, NY, USA: Cambridge University Press, 2011.
4. G. E. P. Box and G. Jenkins, Time Series Analysis, Forecasting and Control. Holden-Day, Incorporated, 1990.
5. R. J. Hyndman and Y. Khandakar, “Automatic time series forecasting: the forecast package for R,” Journal of Statistical Software, vol. 26, no. 3, pp. 1–22, 2008.
6. Y. Sakamoto, M. Ishiguro, and G. Kitagawa, Akaike information criterion statistics, ser. Mathematics and its applications. Japanese series. Tokyo: KTK Dordrecht, 1986, includes index.