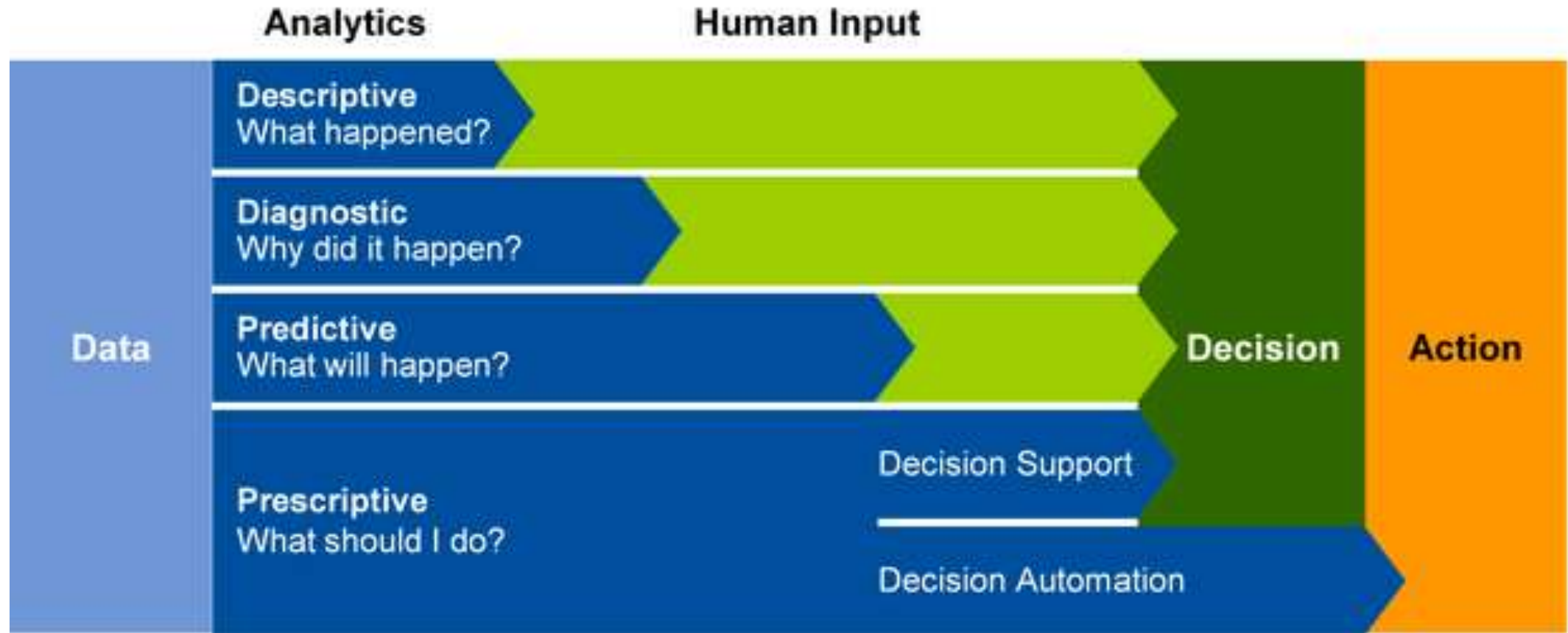# MATLAB EXPO 2017
## KOREA

4월 27일, 서울

등록 하기 matlabexpo.co.kr

# 빅데이터 처리 및 머신 러닝 기법

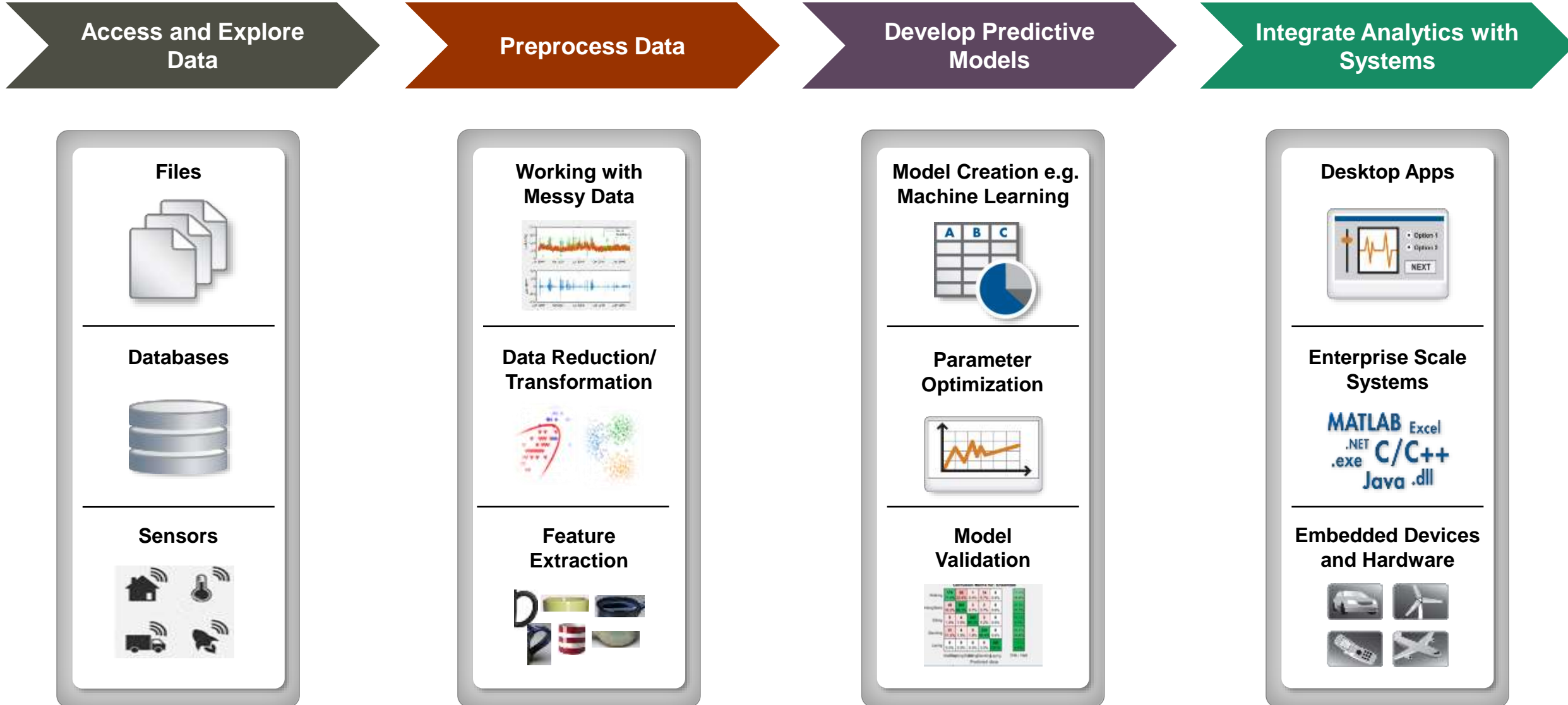Application Engineer

엄준상 과장

# Data Analytics



*Turn large volumes of complex data into actionable information*
*source: Gartner*

# Data Analytics Workflow

# Example: Working with Big Data in MATLAB

- **Objective:** Create a model to predict the cost of a taxi ride in New York City

- **Inputs:**
  - Monthly taxi ride log files
  - The local data set is **small** (~20 MB)
  - The full data set is **big** (~25 GB)



- **Approach:**
  - Acecss Data
  - Preprocess and explore data
  - Develop and validate predictive model (linear fit)
    - Work with subset of data for prototyping
    - Scale to full data set on a cluster

# Example: Working with Big Data in MATLAB

# Data Access and Pre-processing – Challenges

**Challenges**

- Data aggregation
  - Different sources (files, web, etc.)
  - Different types (images, text, audio, etc.)

- Data clean up
  - Poorly formatted files
  - Irregularly sampled data
  - Redundant data, outliers, missing data etc.

- Data specific processing
  - Signals: Smoothing, resampling, denoising, Wavelet transforms, etc.
  - Images: Image registration, morphological filtering, deblurring, etc.
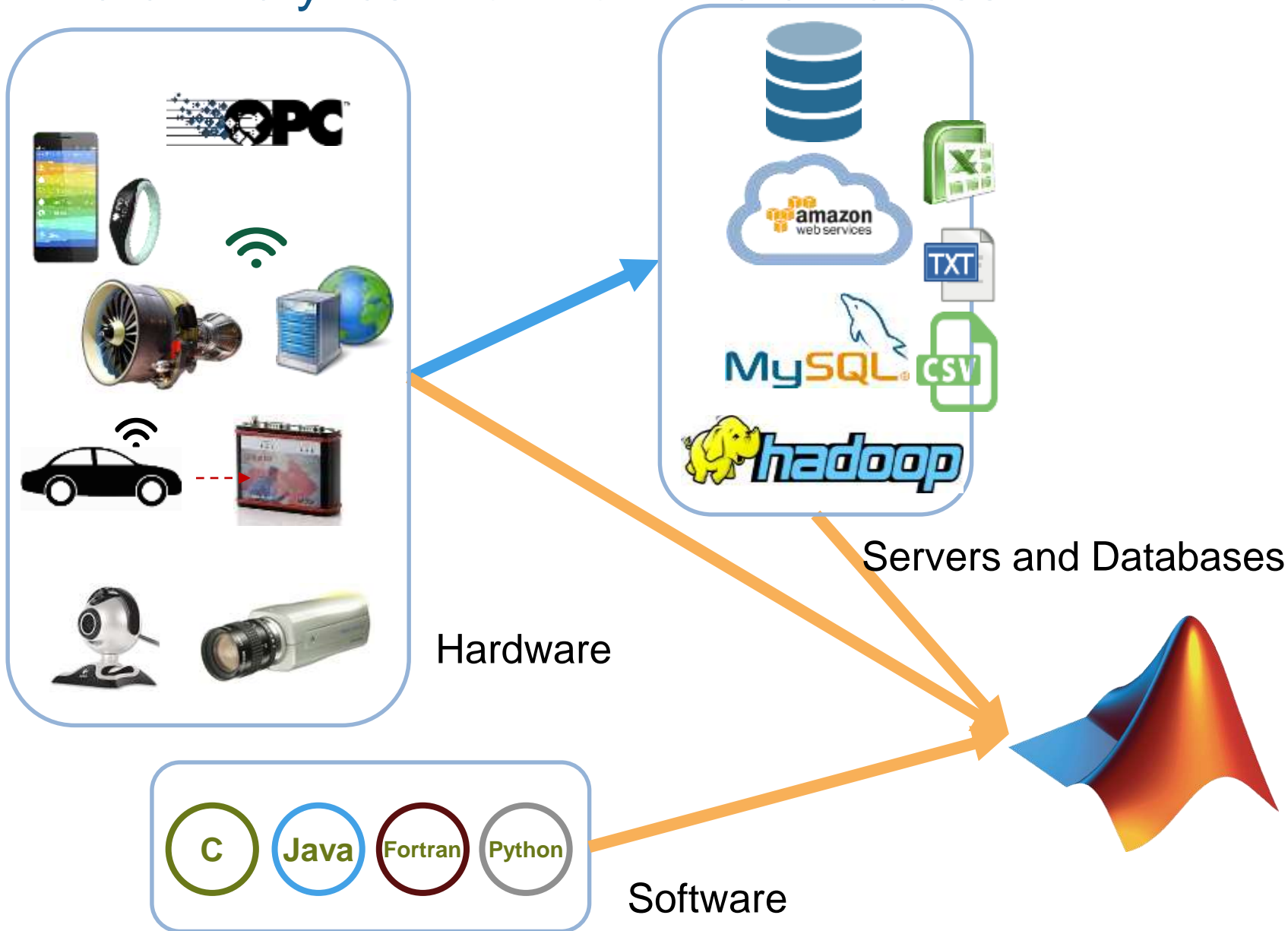
- Dealing with out of memory data (big data)

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

Data preparation accounts for about **80%** of the work of data scientists - Forbes

# Data Analytics Workflow: Data Access



**Servers and Databases**

**Hardware**

**Software**

**Business and Transactional Data**

- Repositories – SQL, NoSQL, etc.
- File I/O – Text, Spreadsheet, etc.
- Web Sources – RESTful, JSON, etc.

**Engineering, Scientific and Field Data**

- Real-Time Sources – Sensors, GPS, etc.
- File I/O – Image, Audio, etc.
- Communication Protocols – OPC (OLE for Process Control), CAN (Controller Area Network), etc.

# Data Analytics Workflow: Big Data Access and Pre-processing

www.nyc.gov/html/tlc/html/about/trip_record_data.shtml

| | Data Analytics - Hom | | Discover MATLAB & | | CRE - Home | | MATLAB | | Fleet Data Analysis |

**2016**

**2015**

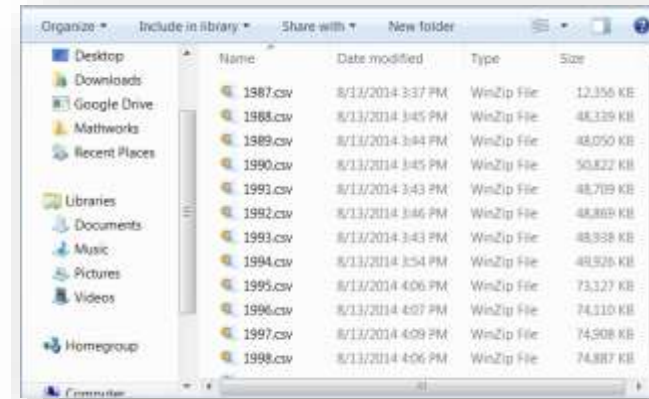| January | Yellow | Green | FHV |
| February | Yellow | Green | FHV |
| March | Yellow | Green | FHV |
| April | Yellow | Green | FHV |
| May | Yellow | Green | FHV |
| June | Yellow | Green | FHV |
| July | Yellow | Green | FHV |
| August | Yellow | Green | FHV |
| September | Yellow | Green | FHV |
| October | Yellow | Green | FHV |
| November | Yellow | Green | FHV |
| December | Yellow | Green | FHV |

**2014**

Download 2015 Taxi Data from Web using *'websave' in parallel*

```
parfor i=1:12
    fileName = ['taxiData2015_', num2str(i)]
    url      = ['https://s3.amazonaws.com/nyc-tlc/trip+data/yellow_tripdata_2016-0',num2str(i), '.csv']
    websave(fileName, url)
end
```

# Big Data in Recent Releases

- **datastore**
  - Tabular text files
  - Images
  - Excel spreadsheets
  - (SQL) Databases
  - HDFS (Hadoop)
  - S3 (Amazon Web Services)
- **MATLAB MapReduce**
  - Scales from Desktop to Hadoop



```
>> preview(ds)
ans =
    Year    Month    DayofMonth    DayOfWeek
    ____    _____    _____    _____
    1987     10          21            3
    1987     10          26            1
    1987     10          23            5
    1987     10          23            5
```

```
airdata = datastore('*.csv');
airdata.SelectedVariables = {'Distance', 'ArrDelay'};

data = read(airdata);
```

# Data Analytics Workflow: Big Data Access and Pre-processing

www.nyc.gov/html/tlc/html/about/trip_record_data.shtml

| | | | |
|---|---|---|---|
| January | Yellow | Green | FHV |
| February | Yellow | Green | FHV |
| March | Yellow | Green | FHV |
| April | Yellow | Green | FHV |
| May | Yellow | Green | FHV |
| June | Yellow | Green | FHV |
| July | Yellow | Green | FHV |
| August | Yellow | Green | FHV |

## Create a `datastore` to represent the data

A `datastore` is a repository for data and allows you to read part of the data, memory.

```
fileLoc = fullfile('taxiData','*.csv');
ds = datastore(fileLoc);
preview(ds)
```

Select variables of interest and give them more intuitive labels.

```
vars = [2:3,5,12:13,16,19];
ds.VariableNames(vars) = {'Pickup','Dropoff','TripDistance',.
    'PaymentType','Fare','Tip','Total'};
ds.SelectedVariableNames = ds.VariableNames(vars);
```

## Connect to the database application

```
conn = database('taxiDemo', 'root', 'matlab', ...
    'Vendor', 'MYSQL', ...
    'Server', 'localhost', ...
    'PortNumber', 3306);
```

## Create a database datastore and import data of interest

```
sqlquery = ['select pickuptime, dropofftime, trip_distance,'...
    'payment_type, fare_amount from taxiData'];
ds = databaseDatastore(conn,sqlquery, 'ReadSize',100000);
```
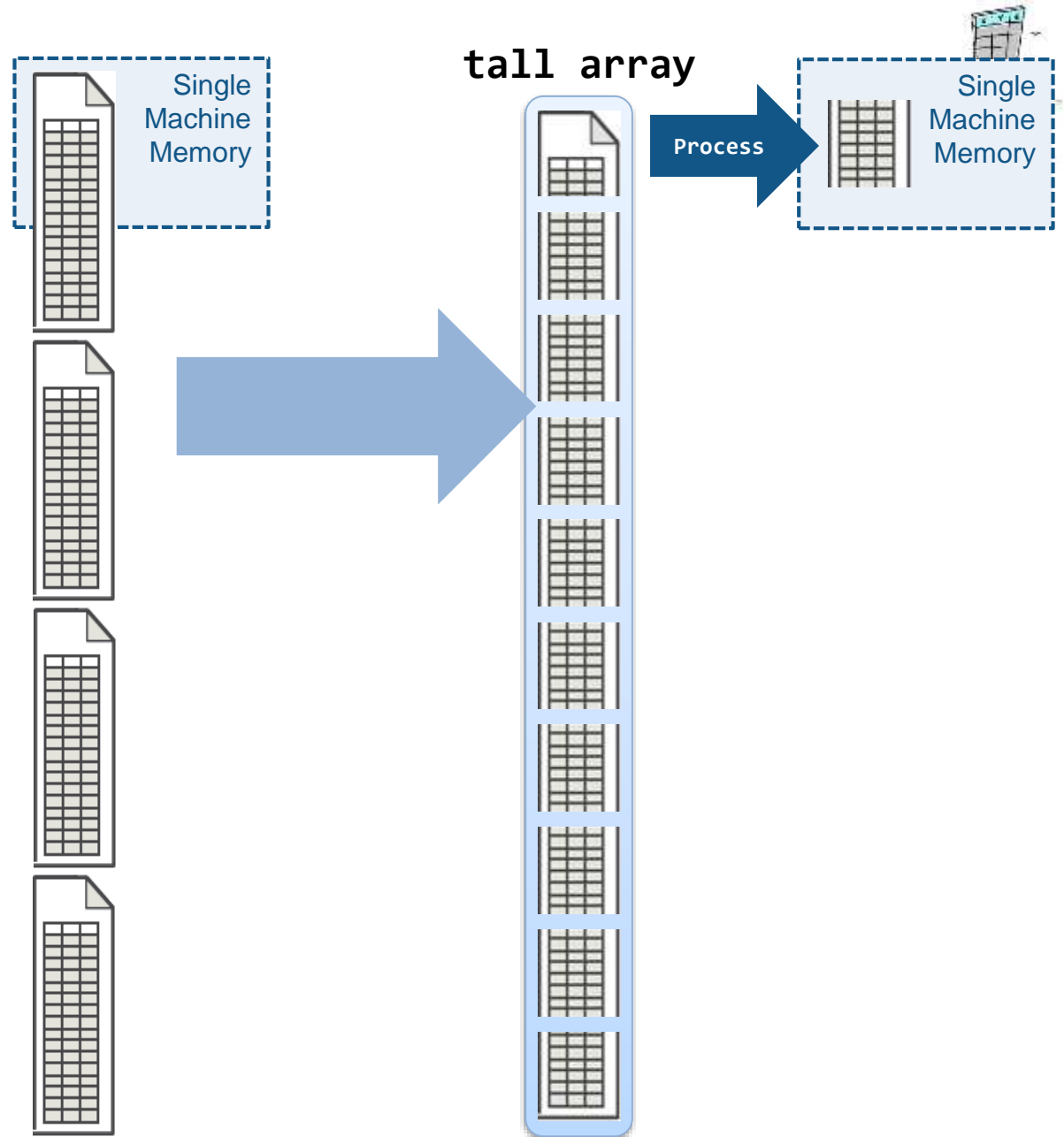
# `tall` arrays in **R**2016**b**

- New data type designed for data that doesn't fit into memory

- Lots of observations (hence "tall")

- Looks like a normal MATLAB array
  - Supports numeric types, tables, datetimes, strings, etc…
  - Supports several hundred functions for basic math, stats, indexing, etc.
  - **Statistics and Machine Learning Toolbox** support
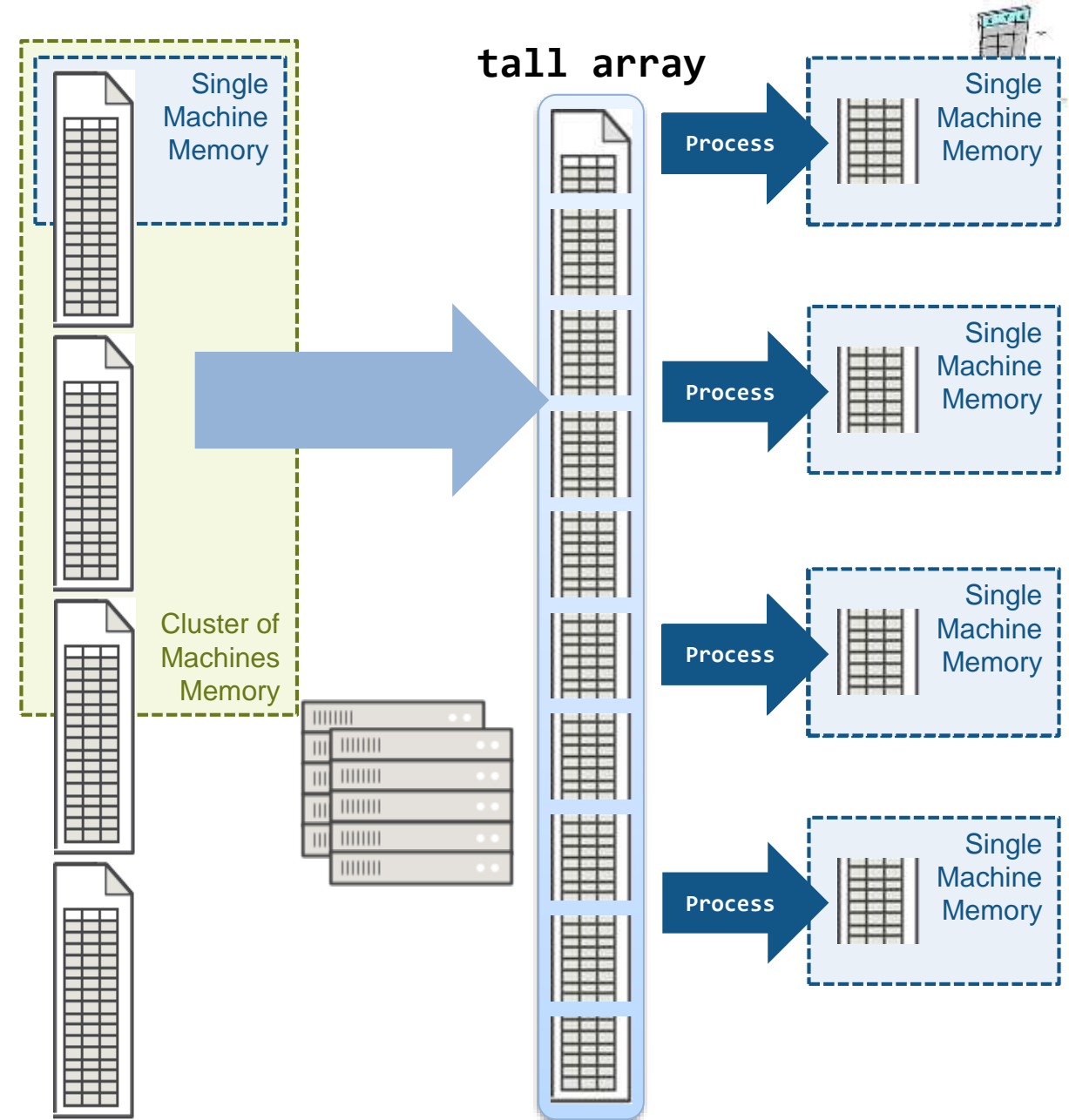    (clustering, classification, etc.)

# tall arrays R2016b

- Automatically breaks data up into small "chunks" that fit in memory

- Tall arrays scan through the dataset one "chunk" at a time

- Processing code for tall arrays is the same as ordinary arrays



Single Machine Memory

**tall array**

Process

Single Machine Memory

# tall arrays R2016b

- With Parallel Computing Toolbox, process several "chunks" at once

- Can scale up to clusters with MATLAB Distributed Computing Server



Single Machine Memory

Cluster of Machines Memory

**tall array**

Process → Single Machine Memory

Process → Single Machine Memory

Process → Single Machine Memory

Process → Single Machine Memory

# Demo: Working with Tall Arrays

# Data Access and pre-processing – challenges and solution

**Challenges**

- Data aggregation
  - Different sources (files, web, etc.)
  - Different types (images, text, audio, etc.)

- Data clean up
  - Poorly formatted files
  - Irregularly sampled data
  - Redundant data, outliers, missing data etc.

- Data specific processing
  - Signals: Smoothing, resampling, denoising, Wavelet transforms, etc.
  - Images: Image registration, morphological filtering, deblurring, etc.

- Dealing with out of memory data (big data)

Files

Signals

Databases

Images

- Point and click tools to access variety of data sources

- High-performance environment for **big data**

- Built-in algorithms for data preprocessing including sensor, image, audio, video and other real-time data

# Machine Learning with Big Data

## R2016b

- Descriptive statistics (skewness, tabulate, crosstab, cov, grpstats, …)

- K-means clustering (kmeans)

- Visualization (ksdensity, binScatterPlot; histogram, histogram2)

- Dimensionality reduction (pca, pcacov, factoran)

- Linear and generalized linear regression (fitlm, fitglm)

- Discriminant analysis (fitcdiscr)

## R2017a

- Linear classification methods for SVM and logistic regression (fitclinear)

- Random forest ensembles of classification trees (TreeBagger)

- Naïve Bayes classification (fitcnb)

- Regularized regression (lasso)

- Prediction applied to tall arrays

# Regression Learner

# Demo: Training a Machine Learning Model

# Demo: Training a Machine Learning Model

# Regression Learner

App to apply advanced regression methods to your data

- Added to Statistics and Machine Learning Toolbox in R2017a

- Point and click interface – no coding required

- Quickly evaluate, compare and select regression models

- Export and share MATLAB code or trained models

# Classification Learner

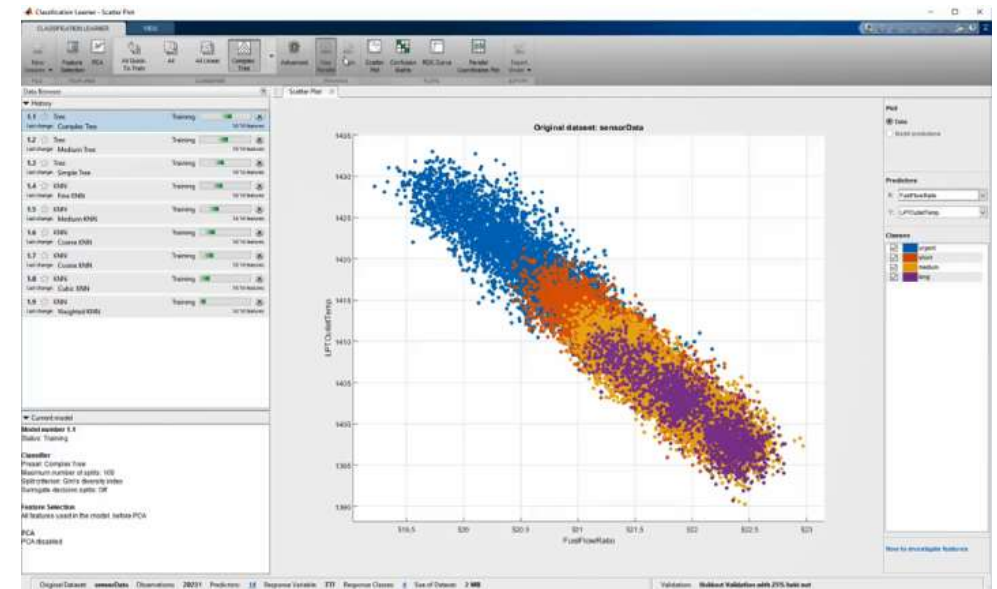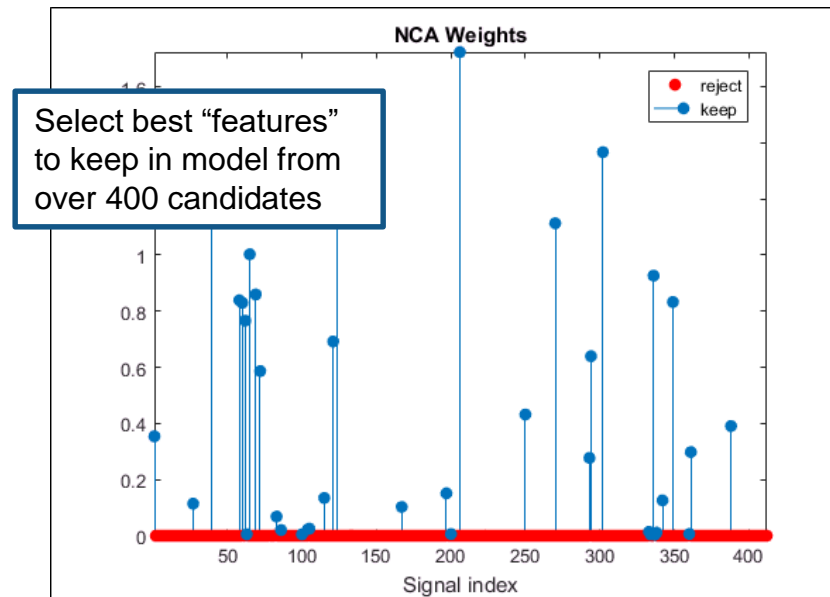App to apply advanced classification methods to your data

- Added to Statistics and Machine Learning Toolbox in R2014a

- Point and click interface – no coding required

- Quickly evaluate, compare and select classification models

- Export and share MATLAB code or trained models

# Tuning Machine Learning Models
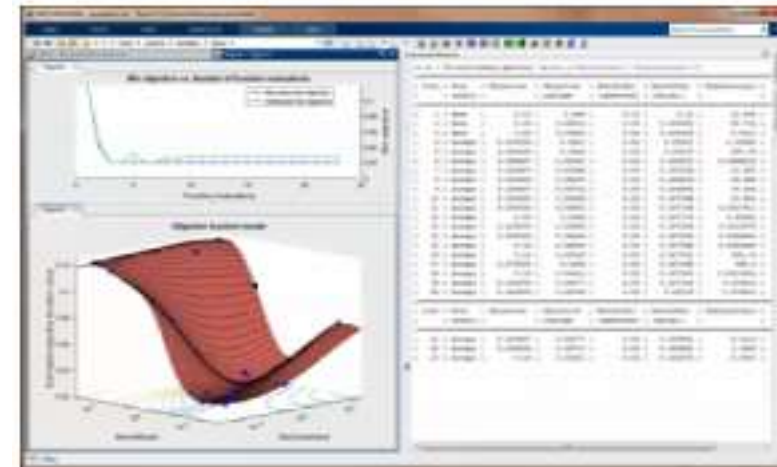Get more accurate models in less time

**Automatically** select best machine leaning "features"



R2016b

NCA:  Neighborhood Component Analysis

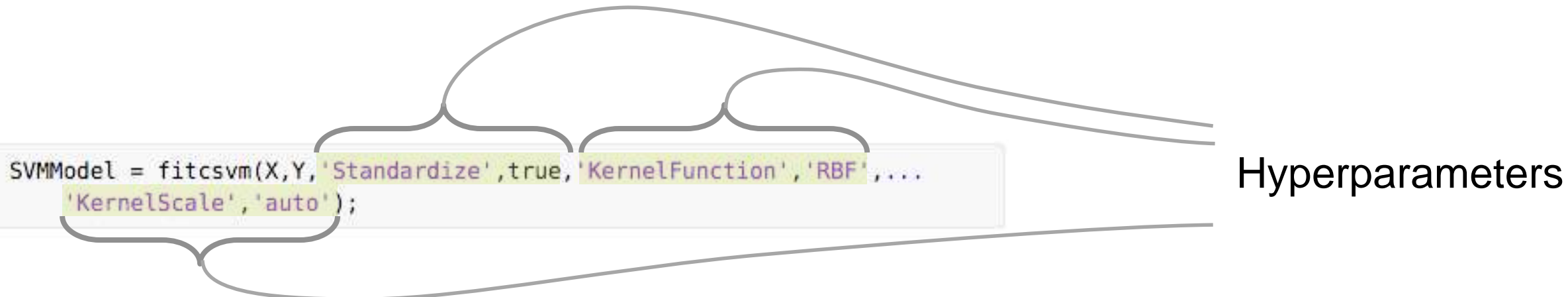**Automatically** fine-tune machine learning parameters



R2016b

Hyperparameter Tuning

# Machine Learning Hyperparameters

```
SVMModel = fitcsvm(X,Y,'Standardize',true,'KernelFunction','RBF',...
    'KernelScale','auto');
```
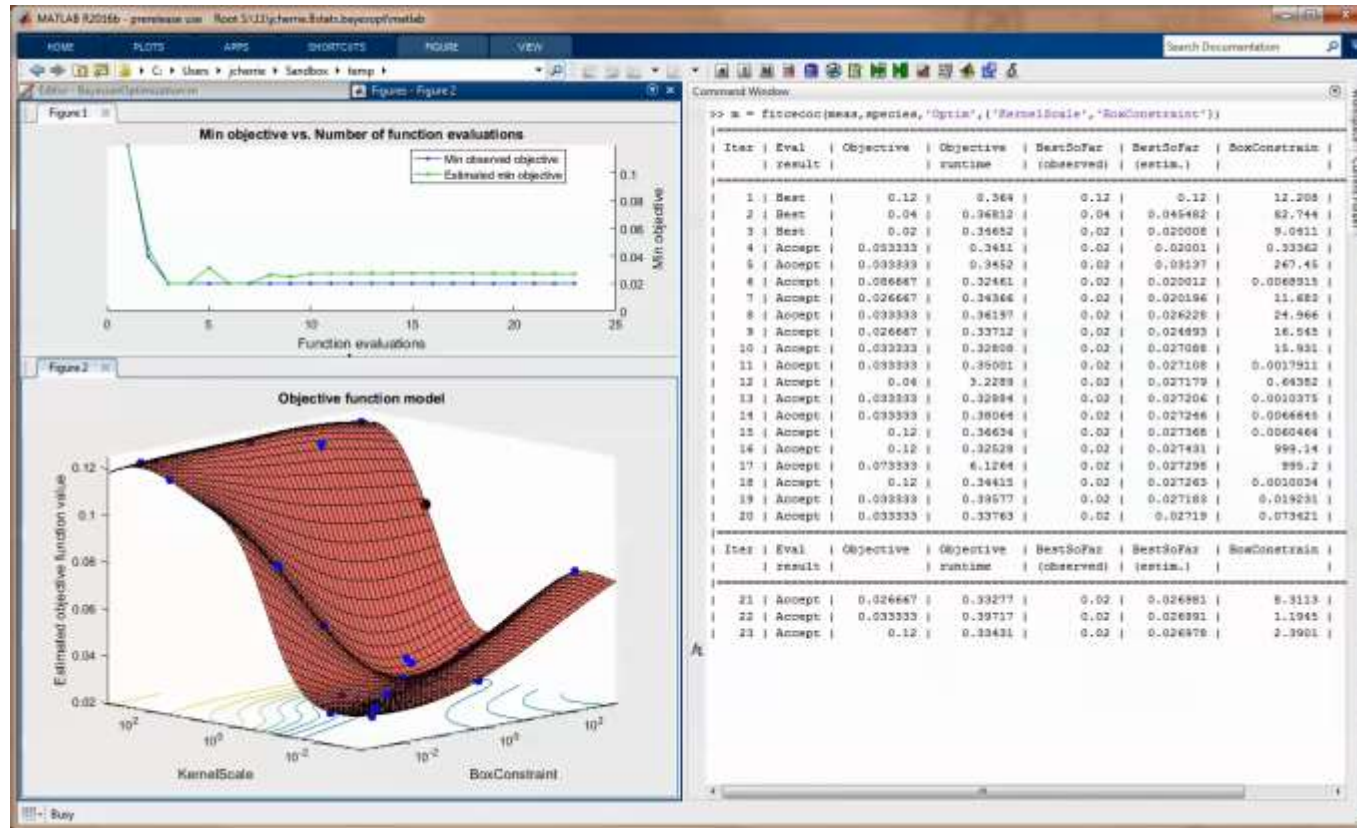
**Hyperparameters**

```
SVMModel = fitcsvm(X,Y,'OptimizeHyperparameters','auto');
```

Tune a typical set of
hyperparameters for this model

```
SVMModel = fitcsvm(X,Y,'OptimizeHyperparameters','all');
```
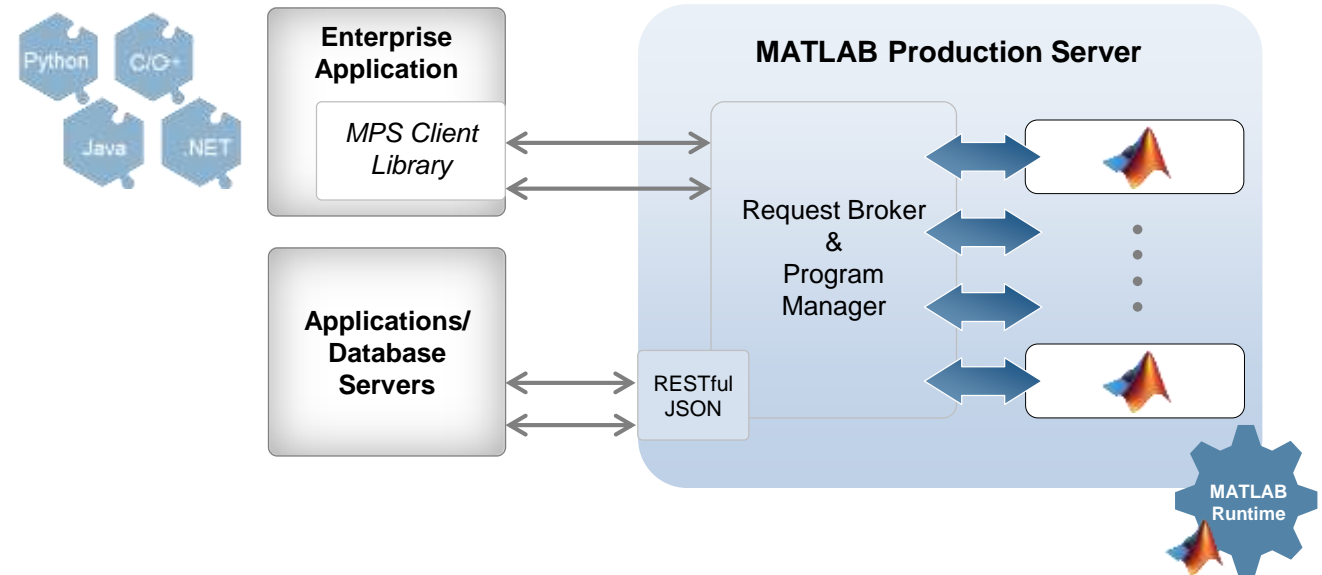
Tune all
hyperparameters for this model
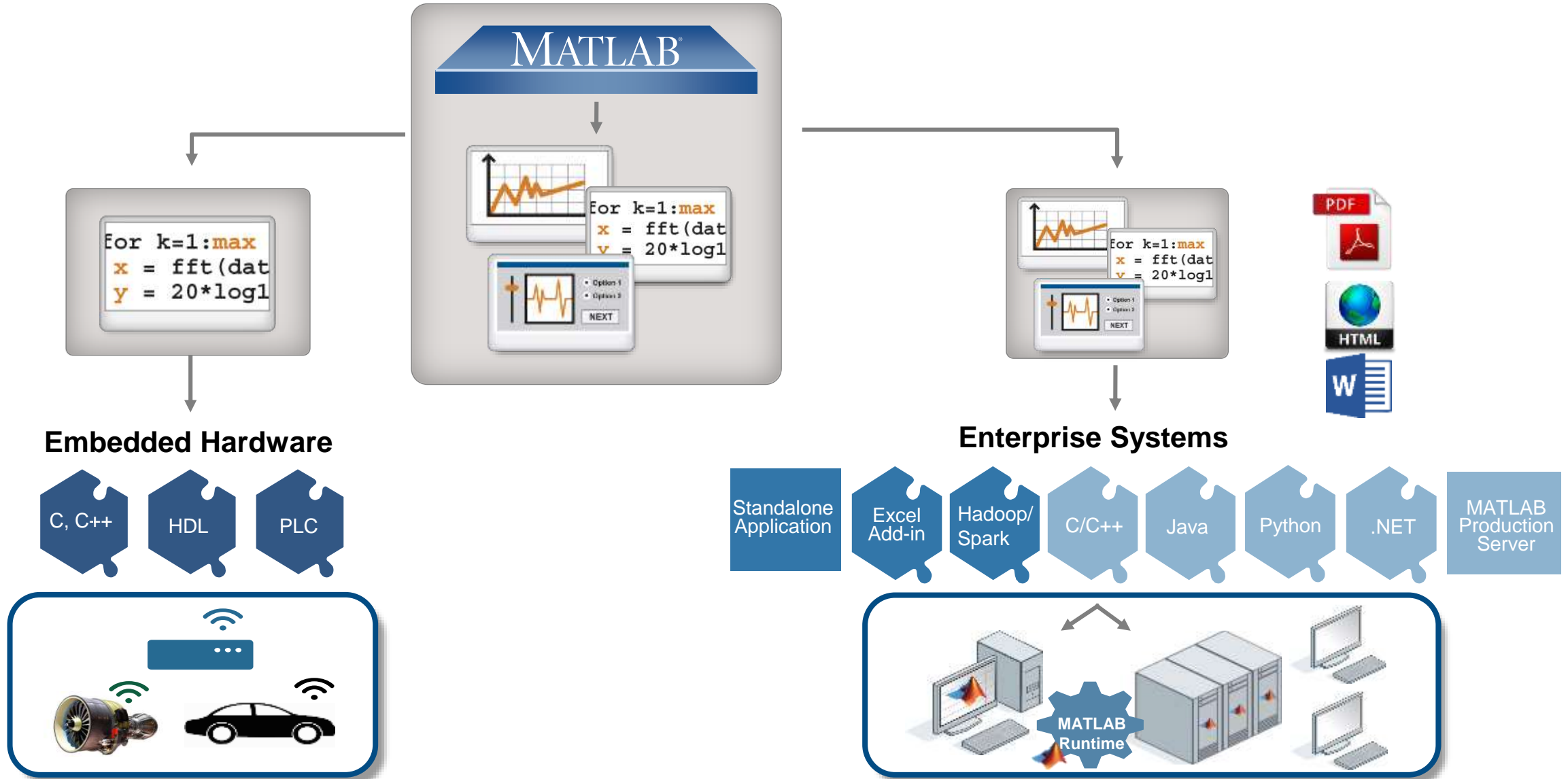
# Bayesian Optimization in Action

# MATLAB Production Server

- ## Server software
  - Manages packaged MATLAB programs and worker pool

- ## MATLAB Runtime libraries
  - Single server can use runtimes from different releases

- ## RESTful JSON interface

- ## Lightweight client libraries
  - C/C++, .NET, Python, and Java

# Integrate analytics with systems



**Embedded Hardware**

C, C++   HDL   PLC

**Enterprise Systems**

Standalone Application | Excel Add-in | Hadoop/Spark | C/C++ | Java | Python | .NET | MATLAB Production Server

MATLAB Runtime

# Key Takeaways

**MATLAB Analytics work with business and engineering data**  ①

▪ Utilize all of your data.

**MATLAB enables domain experts to do Data Science**  ②

▪ Apply advanced analytics techniques.

**MATLAB Analytics run anywhere**  ③

▪ Operationalize analytics to enterprise systems and embedded devices.

# Resources to learn and get started

mathworks.com/machine-learning

mathworks.com/big-data

eBook