# Generic Pipeline for Production of Official Statistics Using Satellite Data and Machine Learning

InKyung Choi (UNECE)

HLG-MOS ML Project Webinar (17 November 2020)
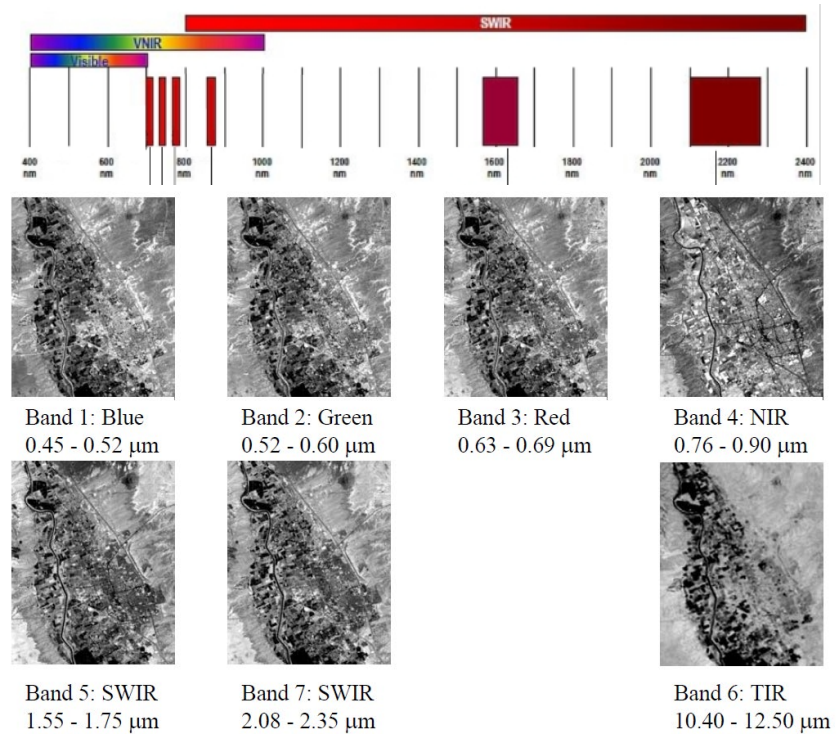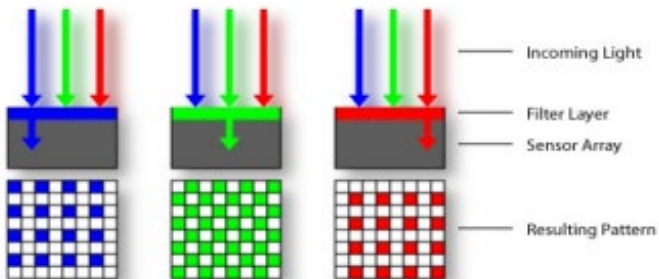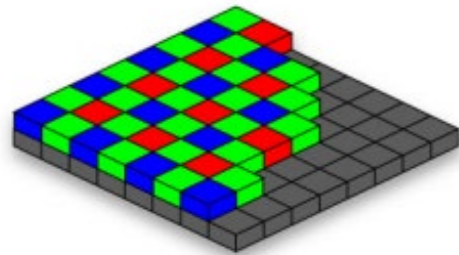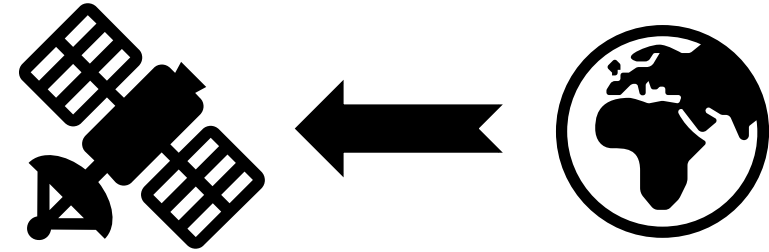
# Satellite data



Band 1: Blue
0.45 - 0.52 μm

Band 2: Green
0.52 - 0.60 μm

Band 3: Red
0.63 - 0.69 μm

Band 4: NIR
0.76 - 0.90 μm

Band 5: SWIR
1.55 - 1.75 μm

Band 7: SWIR
2.08 - 2.35 μm

Band 6: TIR
10.40 - 12.50 μm

Incoming Light

Filter Layer

Sensor Array

Resulting Pattern

# Satellite data



Image source: Natural Resources Canada - Fundamentals of Remote Sensing
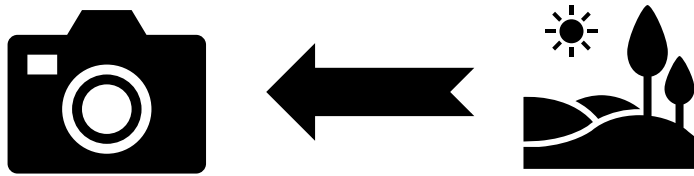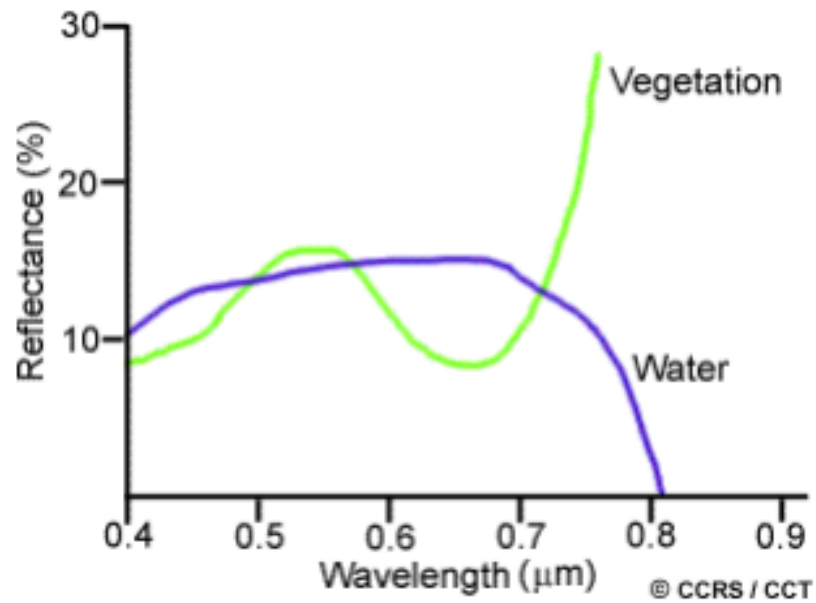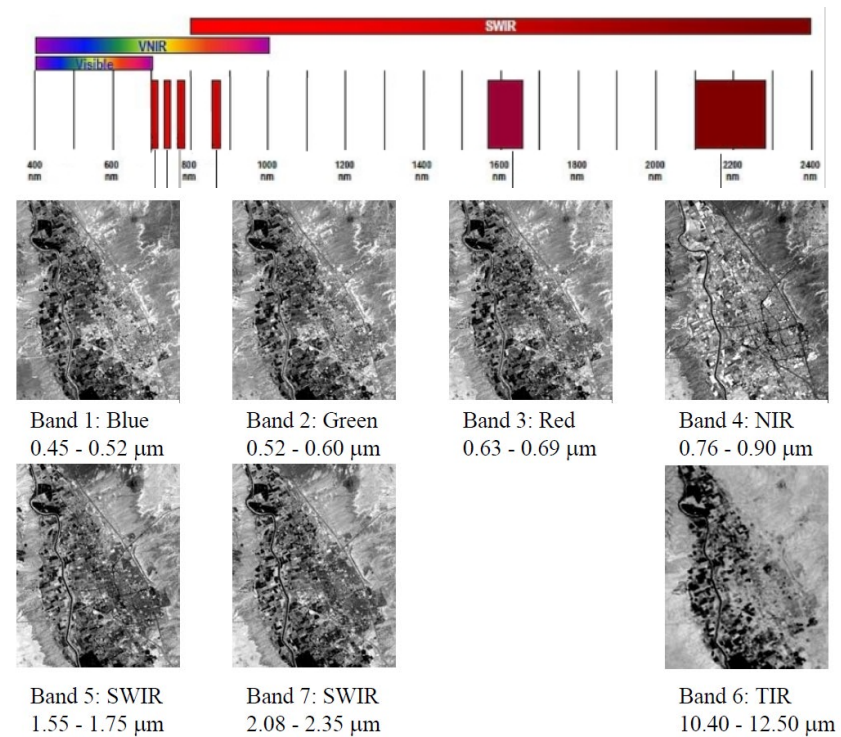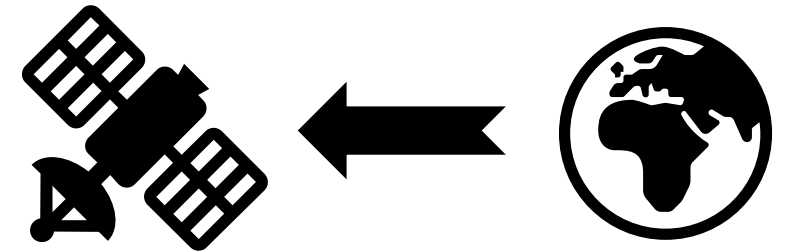
Image source: https://web.nmsu.edu/~aulery/docs/Lab_9_appendix_e.pdf

# Satellite data - benefits

**LandSat 8 Program**
Temporal resolution: Every 16 days
Spatial resolution: 30m

**Sentinel Program**
Temporal resolution: Every 5 days
Spatial resolution: 10-60m

- Global coverage with
- High (temporal) resolution
- Barrier is getting lower
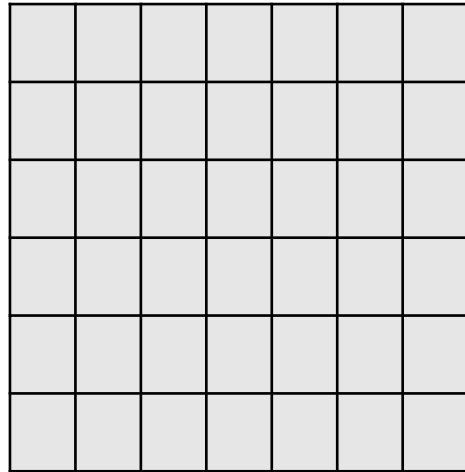
# Satellite data - benefits



- Global coverage with
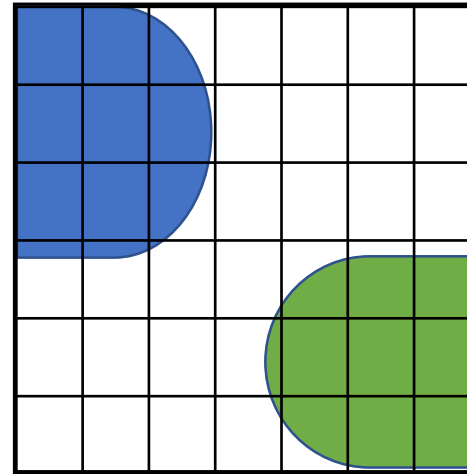- High (temporal) resolution
- Barrier is getting lower

Need of National Statistics Offices (NSOs) to produce statistics at more (spatial) disaggregated at higher (temporal) frequency with lower cost

# Use case - prediction

Satellite data
in pixels

Ground-
truth data in
*some* shapes

Urban

Unknown

Rural

# Use case - prediction

**Satellite data in pixels**

**Ground-truth data in *some* shapes**

Urban

Unknown

Rural

? ? ? ?
? ? ? ?
? ? ? ?
? ? ?
? ? ? ?
? ? ? ?

ML to learn relationship between category and satellite signal based on "labelled" pixels

# Use case - prediction

Satellite data in pixels
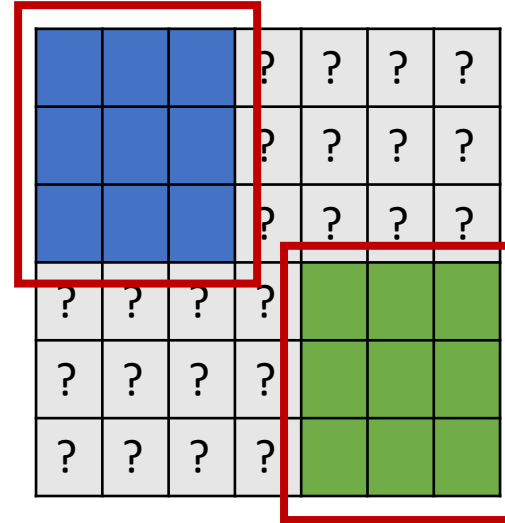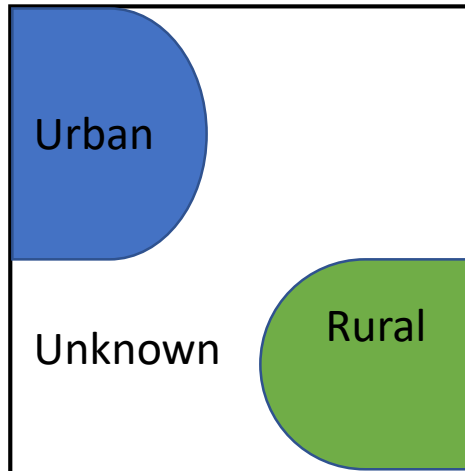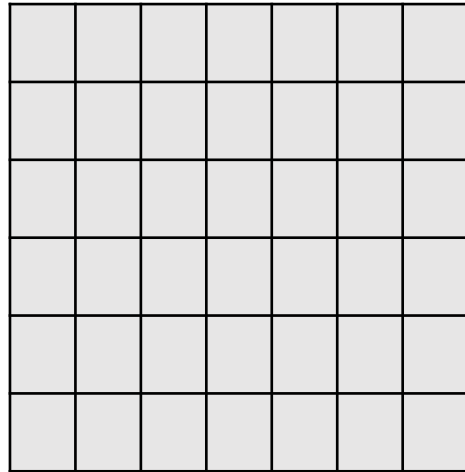
Ground-truth data in *some* shapes

Urban

Unknown

Rural

Can make prediction on the "unlabelled" pixels

# Satellite data - issues

CES (2019) In-depth review of satellite imagery / earth observation technology in official statistics

Conference of European Statisticians(CES) 67th Plenary Session (2019) identified issues such as

- Lack expertise handling satellite data

- Satellite data are big

- Satellite data alone cannot produce statistics

- Institutional commitment to satellite data integration

# Generic Pipeline – why?



Generic Statistical Business Process Model (GSBPM)

# Generic Pipeline - aims

- To improve understanding about workflow needed to use satellite data and ML for statistical production

- To clarify scope and boundary of works

- To serve as common reference points to link increasing body of works

# Generic Pipeline - overview

- What activities should be carried out
- Who should play the leading role?
- What resources are available?

6 Stages

3 Roles



| Business Understanding | Data Collection and Preparation | Modeling (ML) | Prediction | Dissemination | Evaluation |

Thematic Exeprts

Start

Establish problem to be solved

Obtain ground-truth data
- Obtain layers of geo. info
- Build polygons
- Assign labels

Obtain satellite imagery data

EO scientist

Identify satellite data to address the problem

Data scientists, statisticians and computer scientists

Translate the problem into statistical problems (e.g. classification, regression)

Integrate ground-truth data with satellite data

Define features and generate a dataset to be used for analysis

Decide ML algorithm (e.g. NN, Random Forest) and validation method (e.g. CV, LOO)

Train a model with training dataset and test with test dataset. Repeat for each ML and split of dataset

Evaluate ML algorithms and decide one (consider also ensemble)

Apply the model to un-labelled satellite dataset

Evaluate the results using external dataset or against domain knowlege

Publish prediction results with quality criteria for publishing estimates

Evaluate the work process

End

# Generic Pipeline – Business understanding



- Establish the problem to be resolved
- Identify satellite data and ground-truth data to address the problem. Factors to consider:
  - Temporal resolution
  - Spectral resolution
  - Spatial resolution
  - Sustainability
  - Easy of use
  *\* UN Global Working Group on Big Data (2017) Satellite Imagery and Geospatial Data Task Team Report Chapter 2. Data Sources*
- Translate the problem into statistical problem, typically
  - Regression
  - Classification

# Generic Pipeline – Data collection and preparation



- Obtain ground-truth data with attention to geo-referencing

- Obtain satellite data

- Integrate ground-truth data with satellite imagery data

*\* HLG-MOS Data Integration Project (2016-17)*

- Define variables and generate a dataset to be used for analysis (e.g. Normalised Difference Vegetation Index (NDVI))
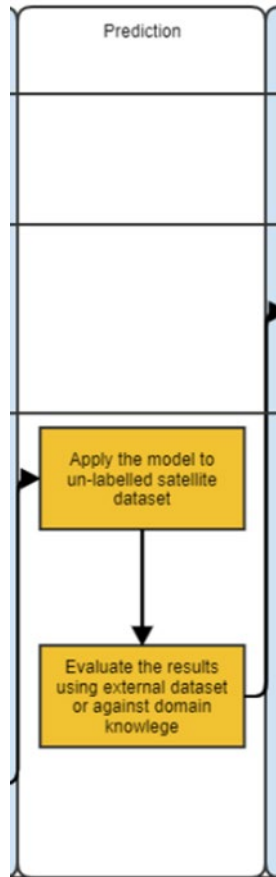
# Generic Pipeline – Modeling (ML)



Modeling (ML)

Decide ML algorithm (e.g. NN, Random Forest) and validation method (e.g. CV, LOO)

Train a model with training dataset and test with test dataset. Repeat for each ML and split of dataset

Evaluate ML algorithms and decide one (consider also ensemble)

- Decide ML algorithms (e.g. SVM, Random Forest, NN), measure of performance (e.g. precision, F1) and validation method

*\* Methodological Approaches for Utilising Satellite Imagery to Estimate Official Crop Area Statistics (2014; ABS)*
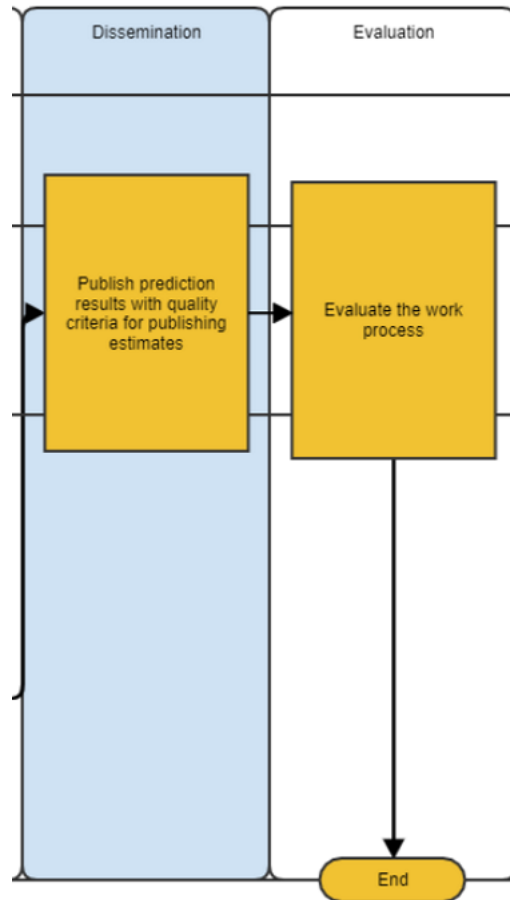
- Train a model with training dataset and test with test dataset. Repeat for each ML and split of dataset

- Evaluate ML models and decide one (consider also ensemble)
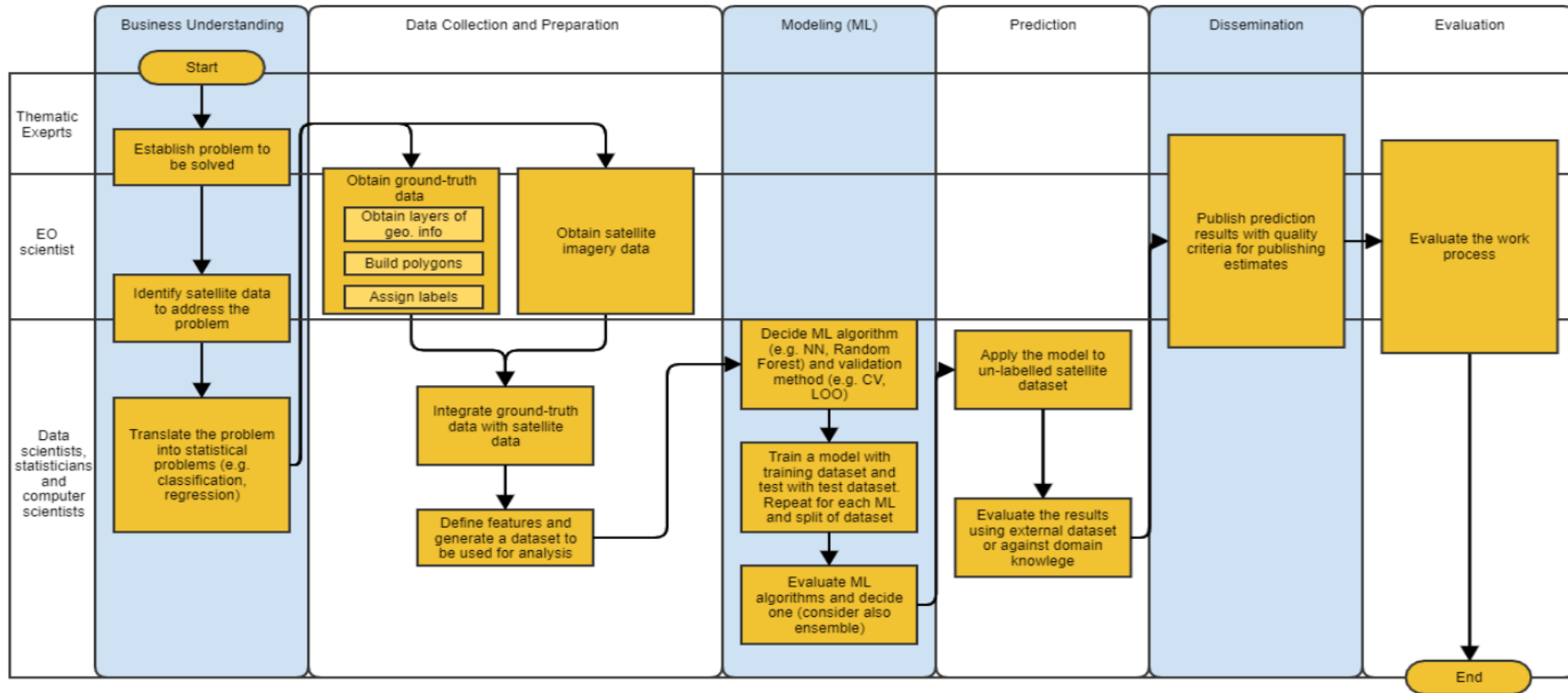
# Generic Pipeline – Prediction



- Apply the model to un-labelled satellite dataset
- Evaluate the prediction results using external dataset or against domain knowledge

# Generic Pipeline – Dissemination and evaluation



- Establish quality criteria for publishing estimates and publish data
- Evaluate the work process

Thank you!