# Imputation of the Production Domain

*Michael C. J. Kao*

**Food and Agriculture Organization
of the United Nations**

# Outline

1 Introduction

2 Yield

3 Production and Area Harvested

# Outline for section 1

## Relationships

The relationship of production and its components by definition
can be expressed as:

$$P_t := A_t \times Y_t \quad P_t \geq 0, A_t \geq 0, Y_t > 0 \qquad (1)$$

Where $P_t$ , $A_t$ and $Y_t$ denotes production, area harvested and
yield, respectively, at time t.

## Scope of the project

A total of **169** commodity just in the crop domain.

**228** countries including obsolete classification and territories.

A total of **31,797** time series require imputation.

Percentage of missing value can be as high as **80%** (By commodity).
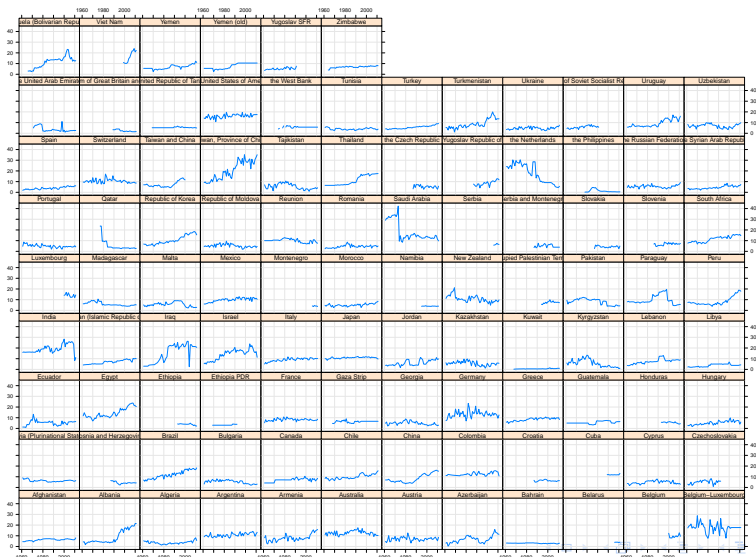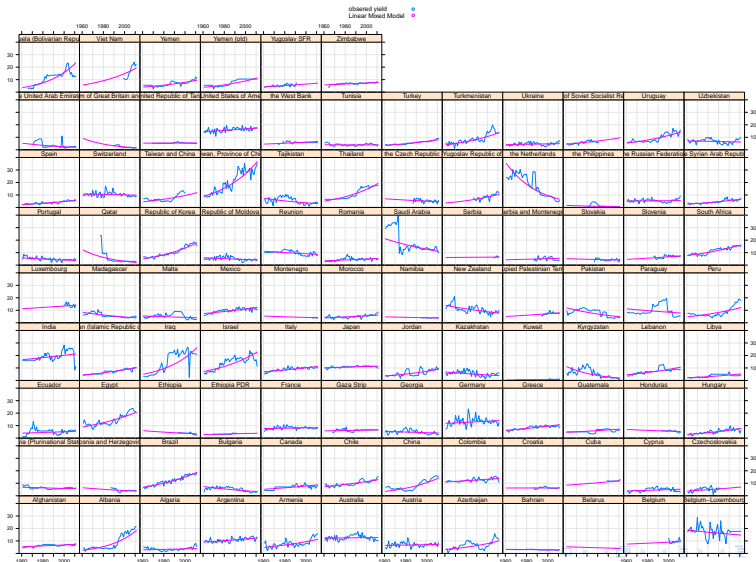
# Outline for section 2

# Linear mixed model

The model implemented in for the imputation of yield is the linear mixed model.

It is an extension of the simple linear regression, but enables the user to take advantage of any correlation structure exist.
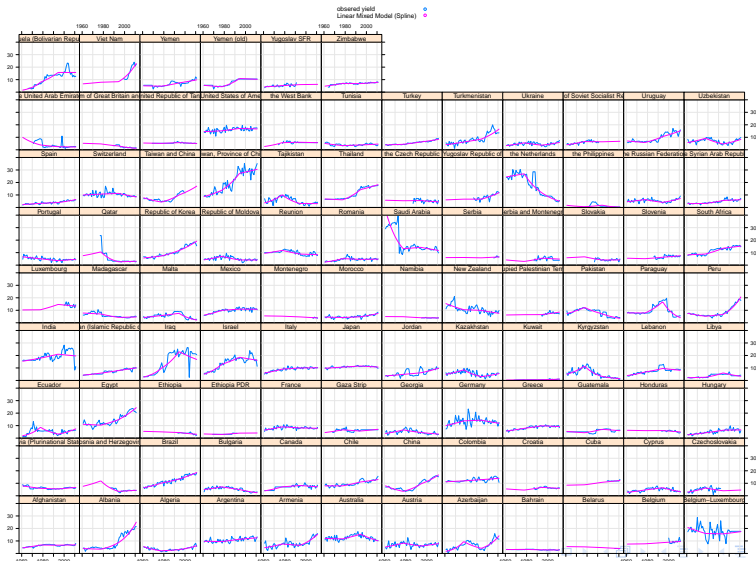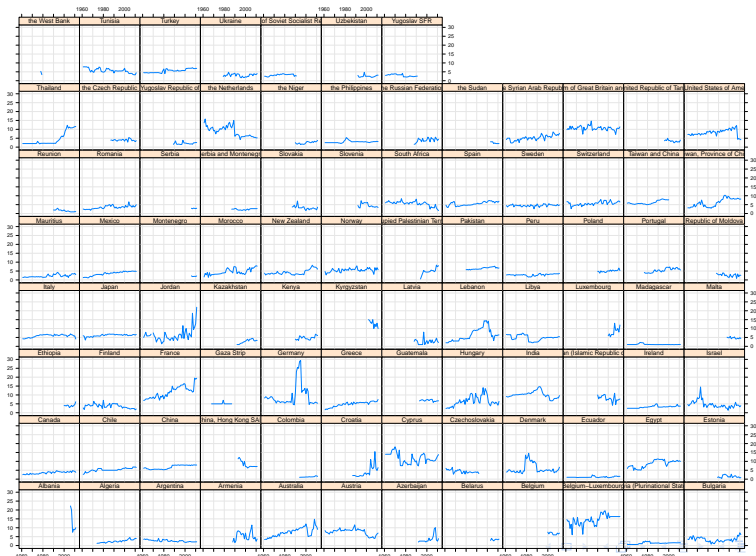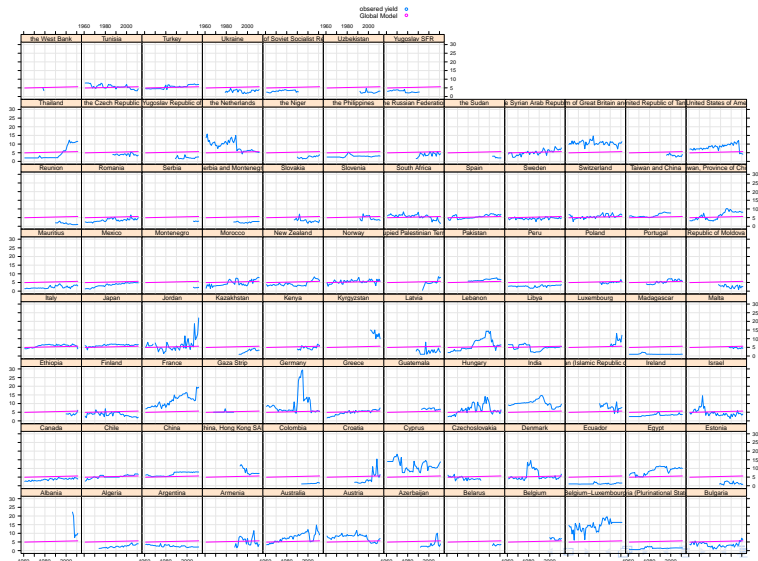
# Yield of grape

# Global Model

# Country Model

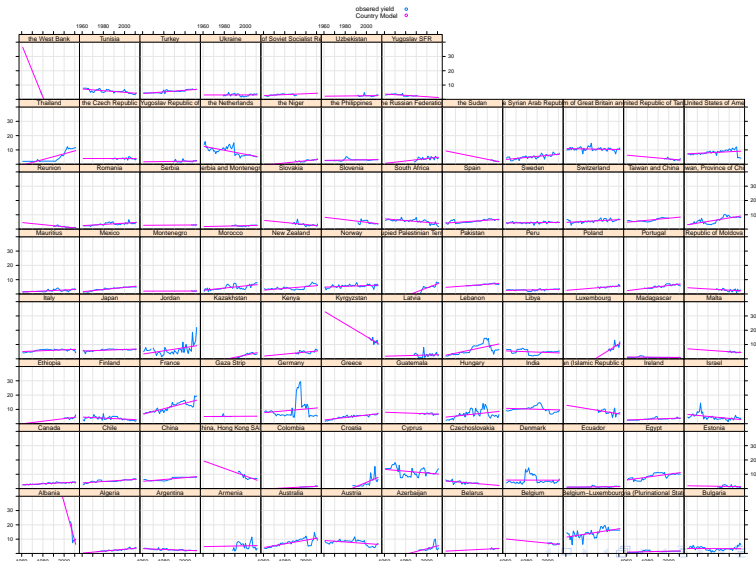# Linear Mixed Model

# Linear Mixed Model with Splines
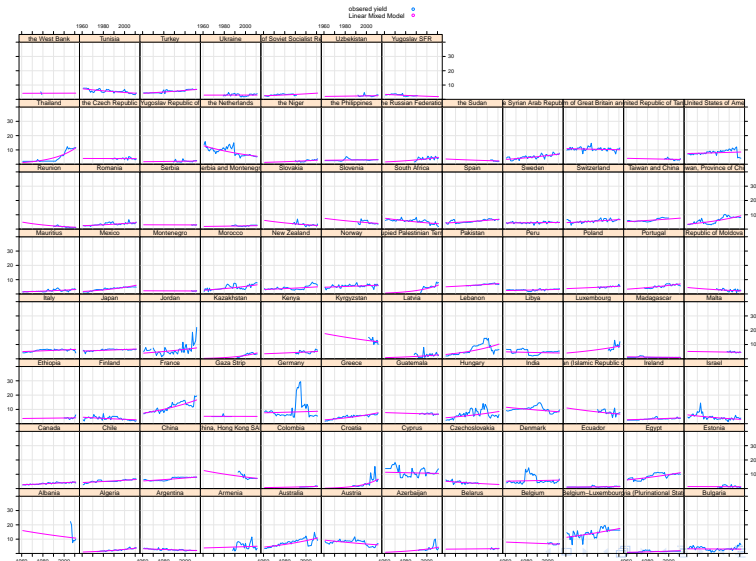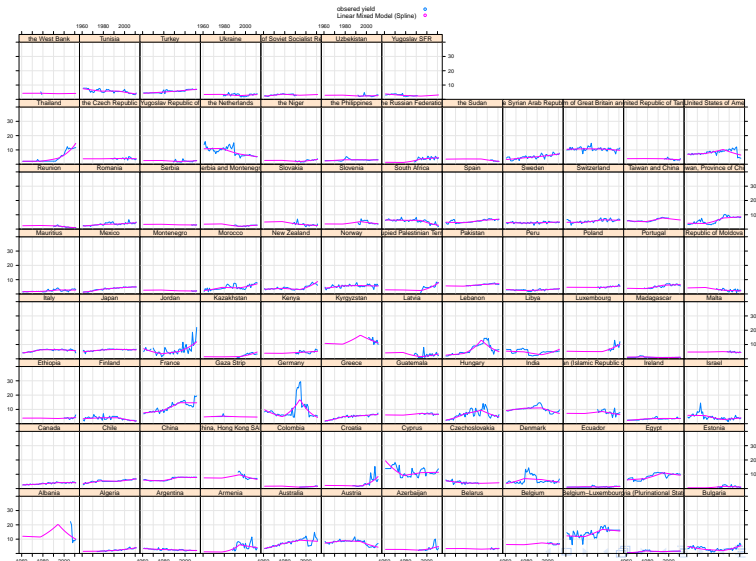
# Yield of green peas

# Global Model

# Country Model

# Linear Mixed Model

# Linear Mixed Model with Splines

There are several reasons why linear mixed model resolve many problems which are associated with fitting a global or country level model.

- It utlize cross-country information if it exists.
- It captures country specific pattern, but conforms to global constraints.
- This allow us to fit more flexible models without over-fitting.

It makes a much more reasonable assumption to our data, countries are different but similar.

# Outline for section 3

1 Introduction

2 Yield

3 Production and Area Harvested

After the imputation of the yield, we can balance the area harvested or production provided the counter-part exists.

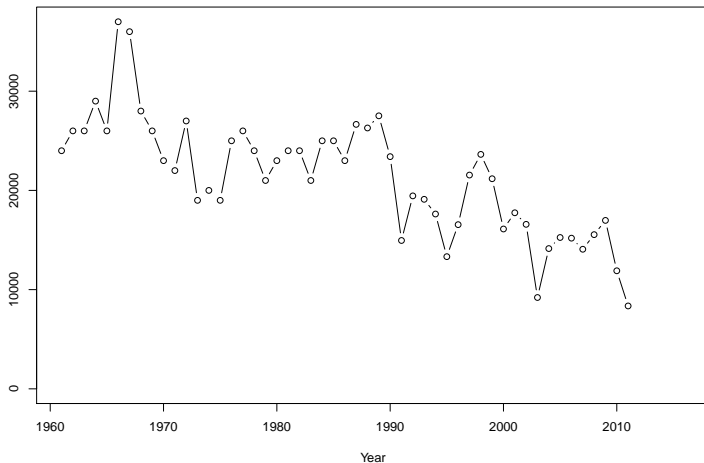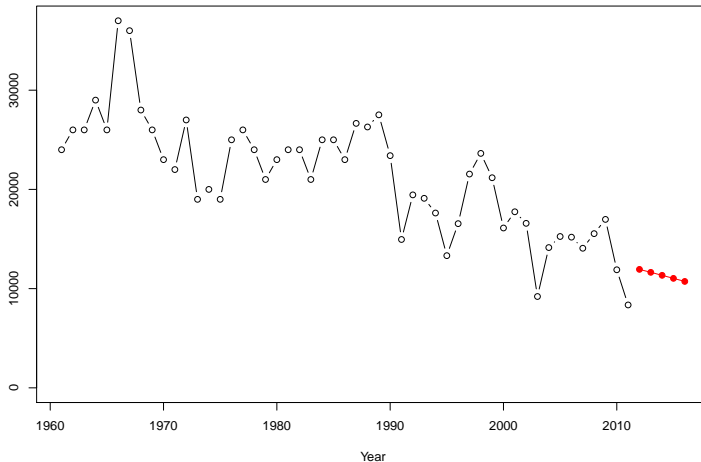However, the problem is when both of them are missing.

# What is the solution?

Oppose to yield, we have no clear cross-country information which we can utilize.

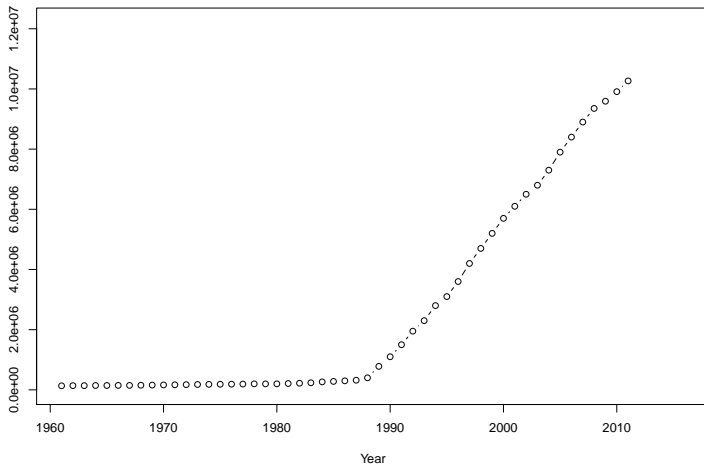Furthermore, we actually have no idea about the trend, or even the shape of production!

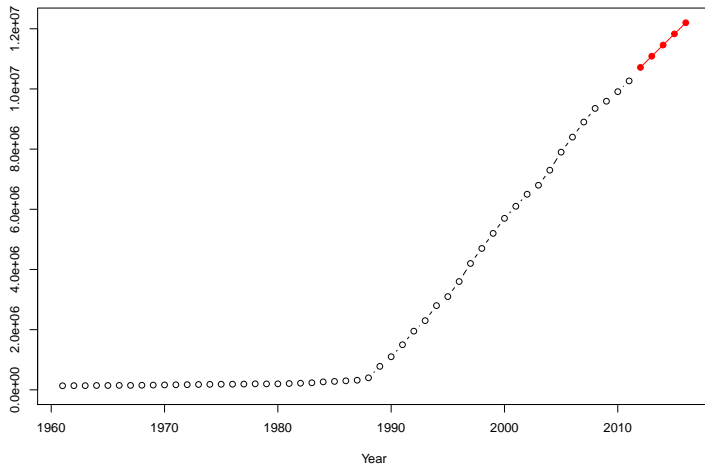Production of green peas in South Africa
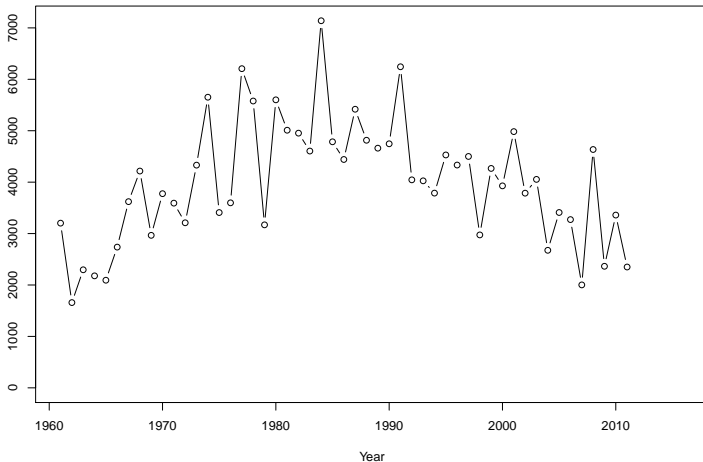
Production of green peas in South Africa
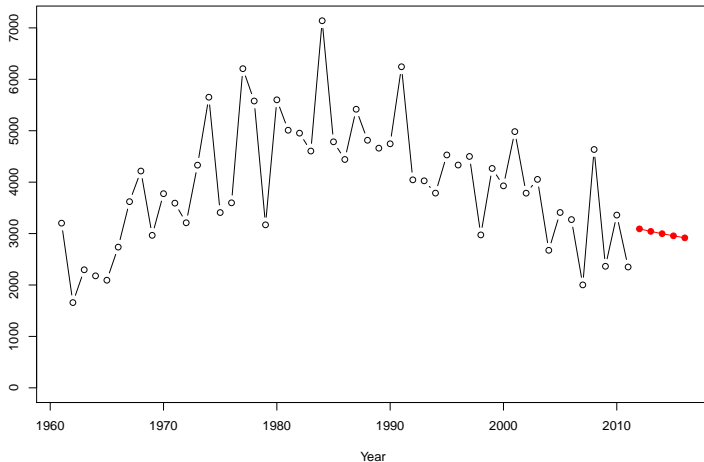
Production of green peas in China
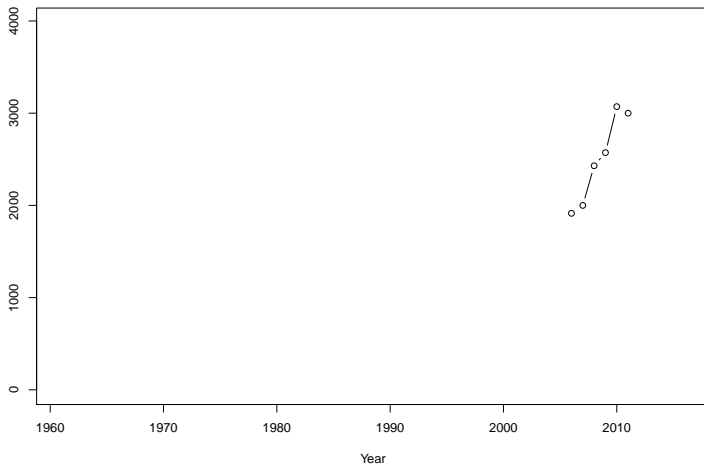
**Production of green peas in China**

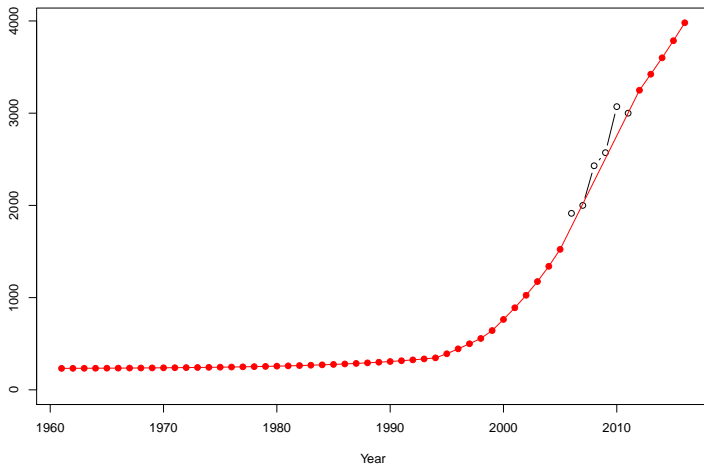Production of green peas in Norway
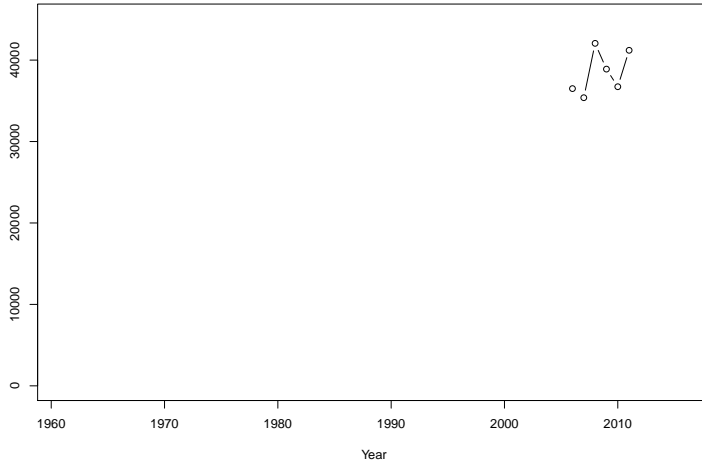
Production of green peas in Norway

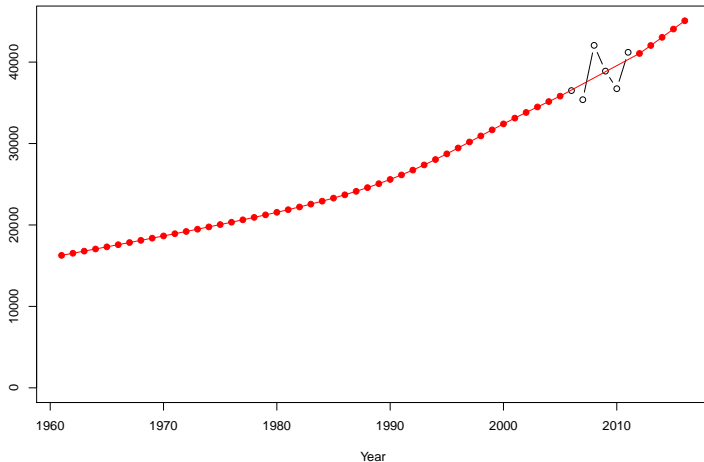Production of green peas in Albania
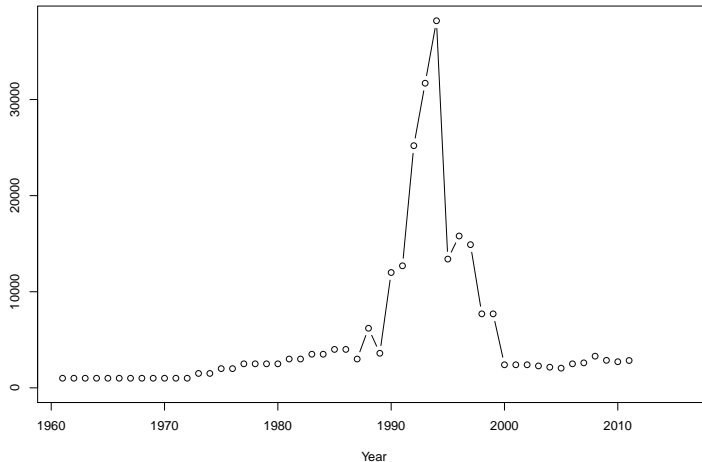
Production of green peas in Albania
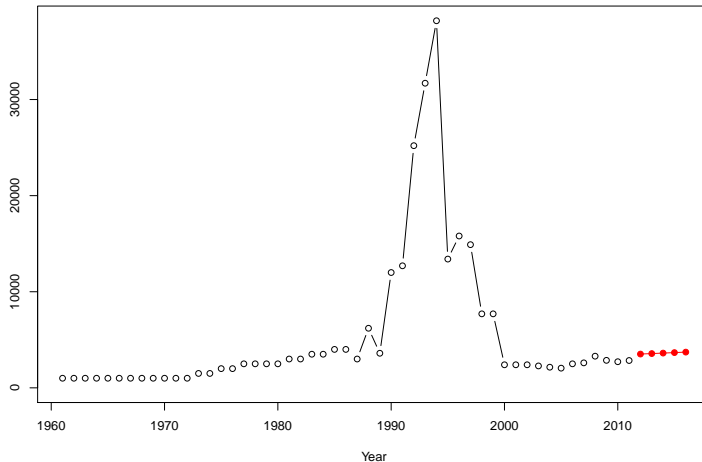
Production of green peas in Serbia
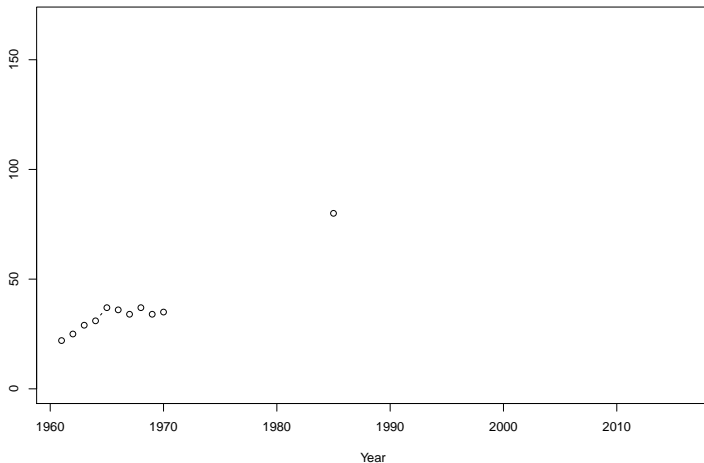
**Production of green peas in Serbia**

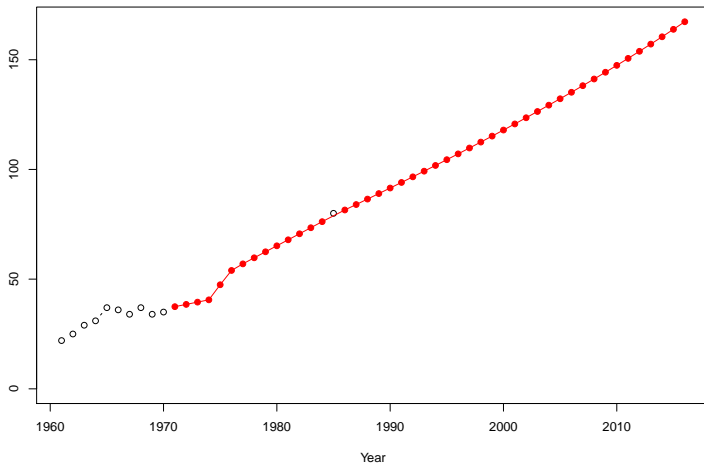**Production of green peas in Ireland**

Production of green peas in Ireland

**Production of carrots in Trinidad and Tobago**

Production of carrots in Trinidad and Tobago

## What is Ensemble Learning?

Ensemble learning in its simplest sense, is to build multiple models/learners and combine them to obtain the final model or prediction.

It consist of two components:

1. Building multiple models or learners.
2. Combine the models and predictions.
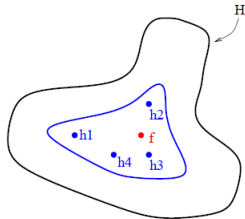
# Why Ensemble Learning?

Ensemble as described by Dietterich (2000) can mitigate the following three issues.

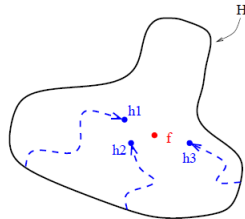Statistical: Lack of data to identify an unique solution.

Computational: Optimization
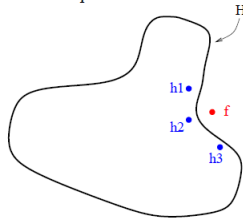
Representational: Complex model

## Implementation

The details of the ensemble implemented is describe here:
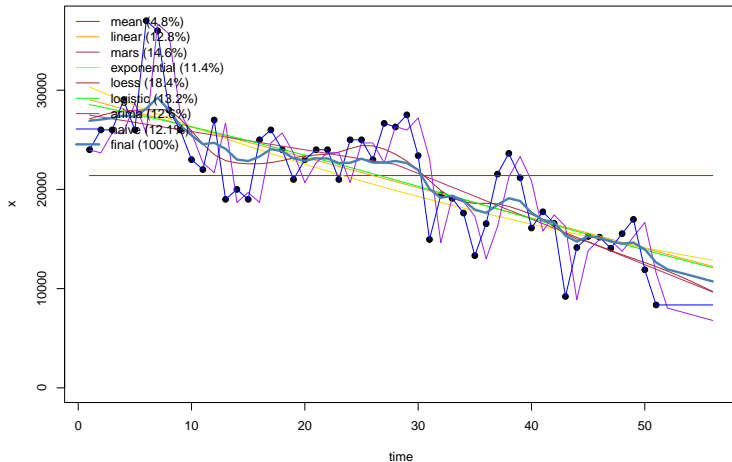
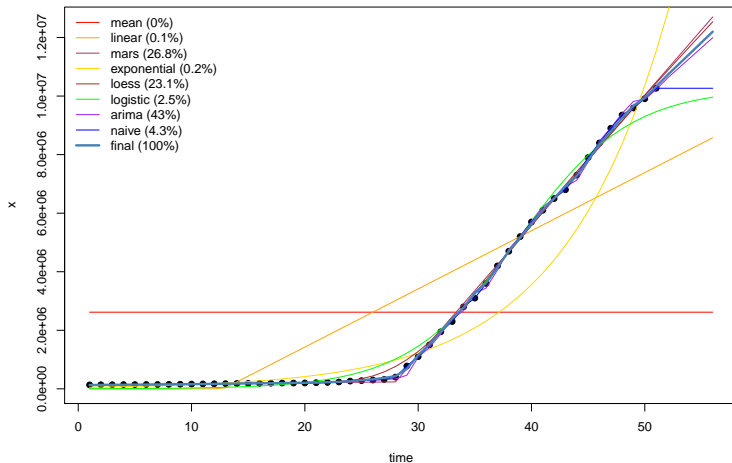- **Base learners:**
  - mean
  - linear
  - MARS
  - exponential
  - locally smooth linear
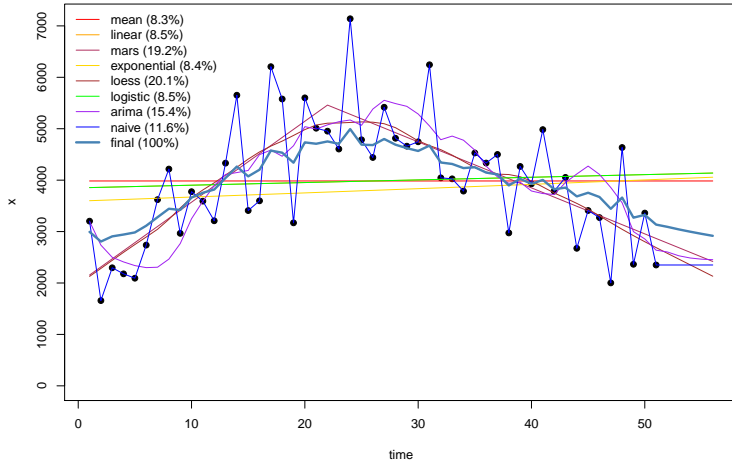  - logistic
  - ARIMA
  - naive

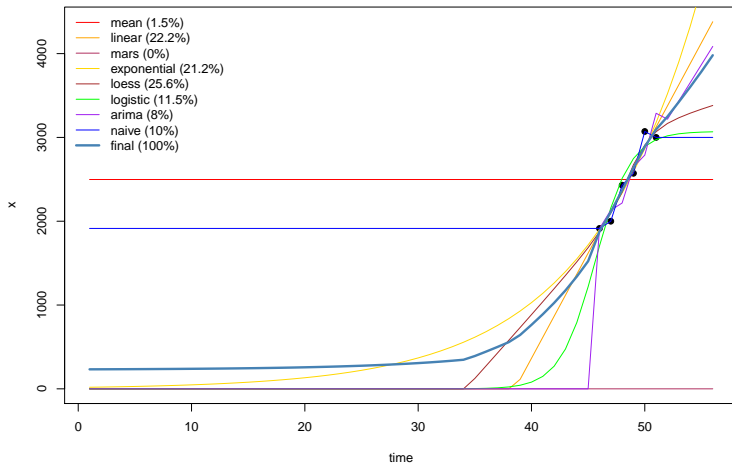- **Combiner:** non-trainable algebraic combiner - Weighted sum rule

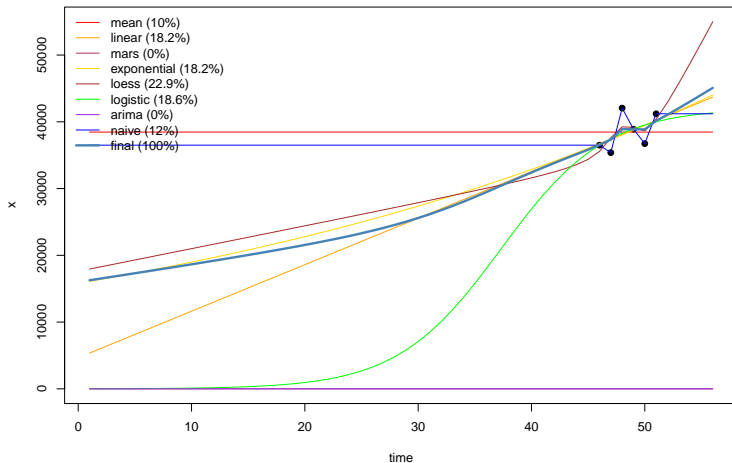$$u_j(x) = \sum_{i=1}^{N} w_i d_{n,j}(x)$$

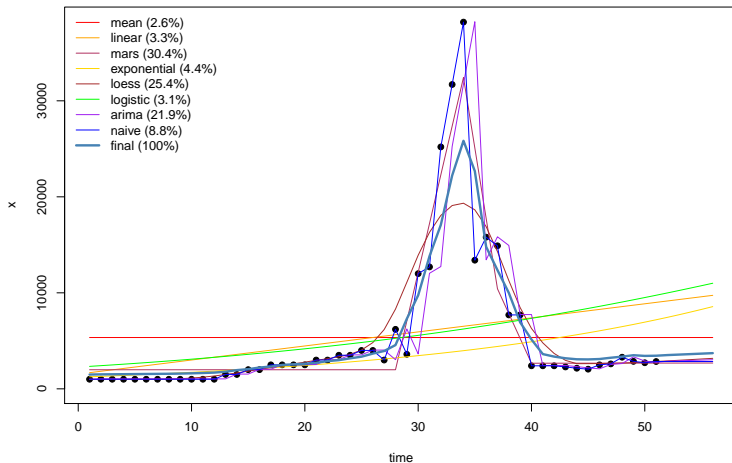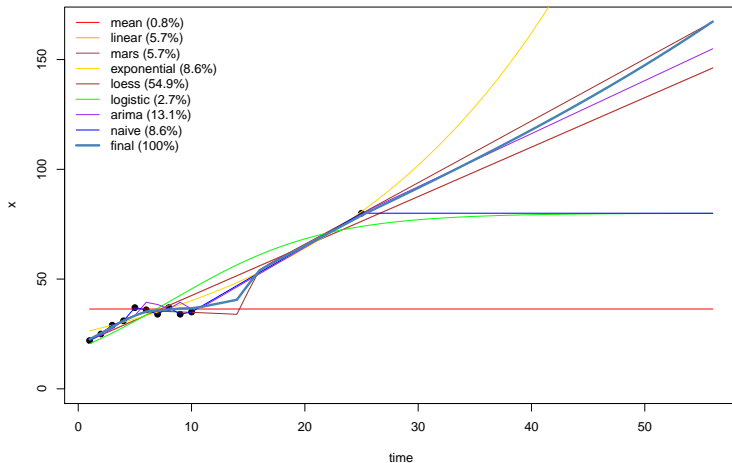Where the weights depends on the fit on the available data.

## Practical benefits of adopting the ensemble model

From the examples, we can see that we have an extremely flexible model with almost no possibility of over-fitting since none of the model are itself complex.

We can continue to add in more models in which we consider appropriate.

The reason it is called ensemble learning is becaues it is a model that can learn. If a production that has been growing linearly in the past but suddenly at an exponential rate, the model will learn and shift the weight of the linear model to model which are more flexible.

## Explanatory variables?

Why we don't use other explanatory variables:

- Can contain as many missing value or more than our production domain.

- Too many of them, hard to choose from the set. A simple set would contain precipitation and temperature, should we also include number of sunny days? and their interactions? What about $CO_2$/phosphorus concentraction? Number of bee hives? Fertilizer, pesticide consumption?

- Difficult to maintain and integrate into the system.