# Machine Learning applied to Planetary Sciences

## PTYS 595B/495B

## Leon Palafox

https://leonpalafox.github.io/MLClass/

# Experts can be wrong

# In the news!

## An artificial intelligence system that correctly predicted the last 3 elections says Trump will win

Pamela Engel ✉ 🐦
🕐 Oct. 28, 2016, 8:24 PM    🔥 152,927

The polls have consistently showed Hillary Clinton with a lead over Donald Trump in recent weeks, but an artificial intelligence system has a different prediction for the outcome of the presidential election.

The system, called MogIA, uses 20 million data points from online platforms like Google, YouTube, and Twitter to come up with its predictions, according to CNBC. MogIA correctly predicted the past three presidential elections as well as the

Donald Trump appears at a campaign event in Geneva, Ohio, U.S., October 27 2016. REUTERS/Carlo Allegri

# Application to time series data

# Problem Statement

- Both in astronomy and planetary science time series data is important.

- Even if it's not a time series, we can use these techniques in other datasets.

- Most of the techniques used for time series are rather old.
  - We won't be talking about those.

- Traditional Machine Learning approaches are very effective.

# Datasets

- Sunspots dataset (people just love this one)

- Light curves.

- Terrain Profiles (not strictly time-series, but works!)

# Problems

- Classification:

  - We would like to know whether there is a signal at one particular time.

  - We would like to know if the time series is informative of different effects.

- Regression:

  - We would like to predict the time series in the future.

# Problems

- Time Series data is usually 1D, noisy and incomplete.

# Feature selection

- If we have multiple time series, I can get easier, or harder.

- Time series are remarkably non-stationary.

- Non-Stationarity:

  - The signal does not remain constant over time
  - Single data points are very uninformative
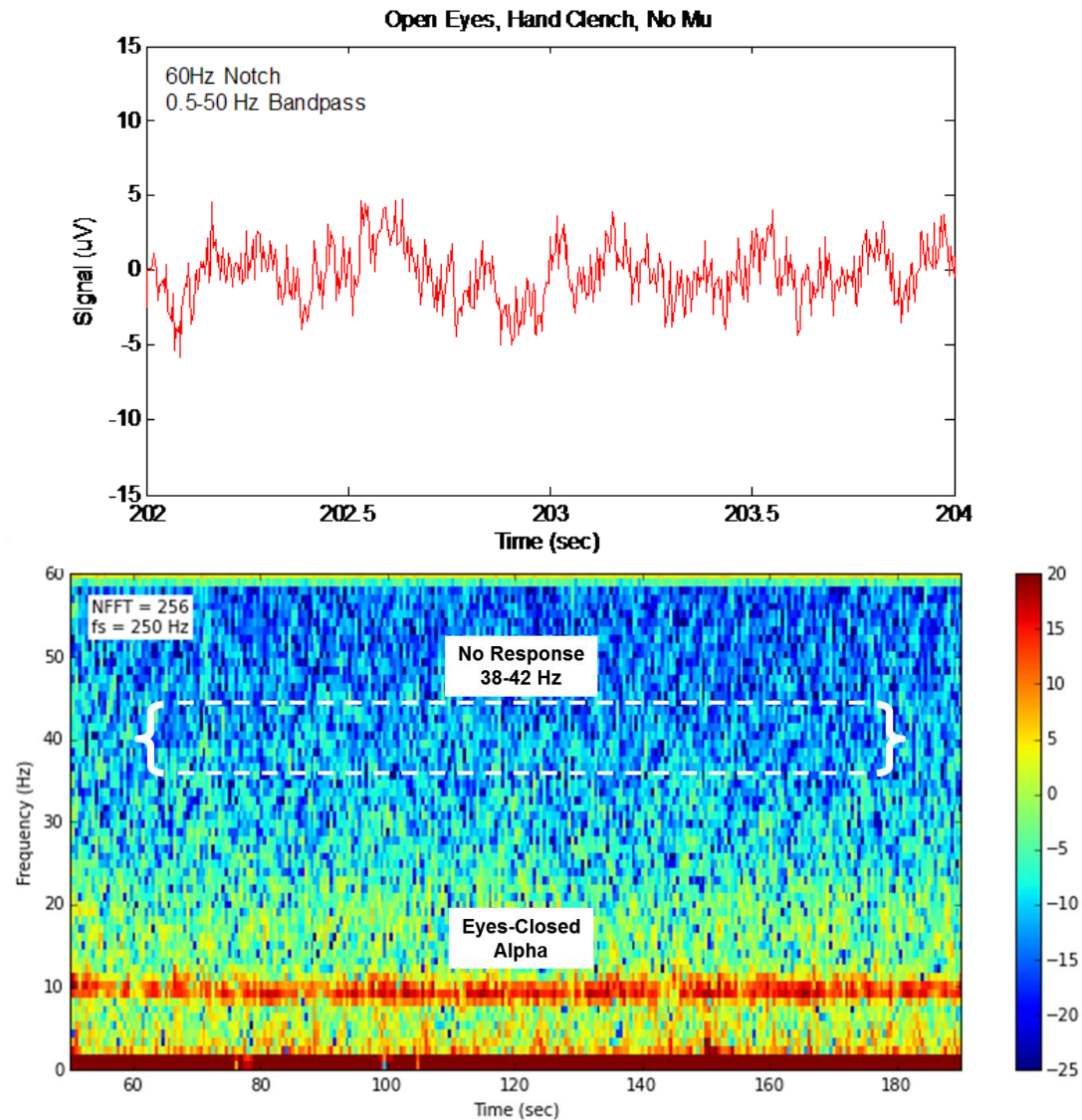
# Stock Market



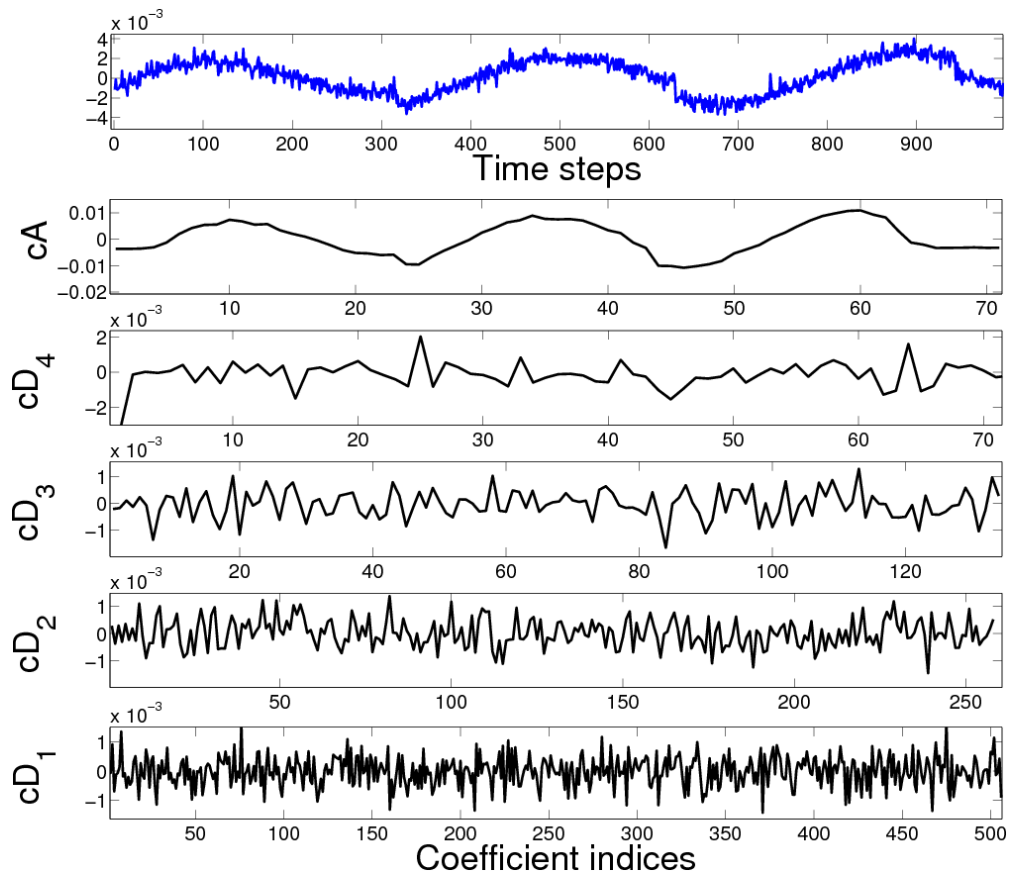You are hired to create a classifier that will sell or buy.

# Features

- We need to create features that will be stationary.

- Spectrogram
  - Use Fourier transforms at certain interval windows.
  - The plot becomes one of Time vs Frequency
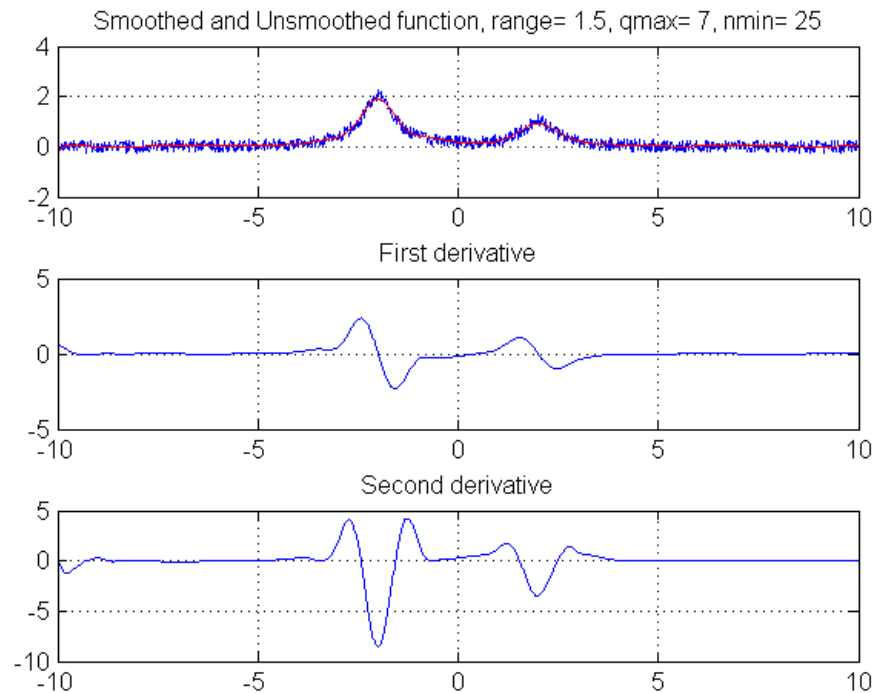  - Single command in Matlab and in Python

# EEG Example
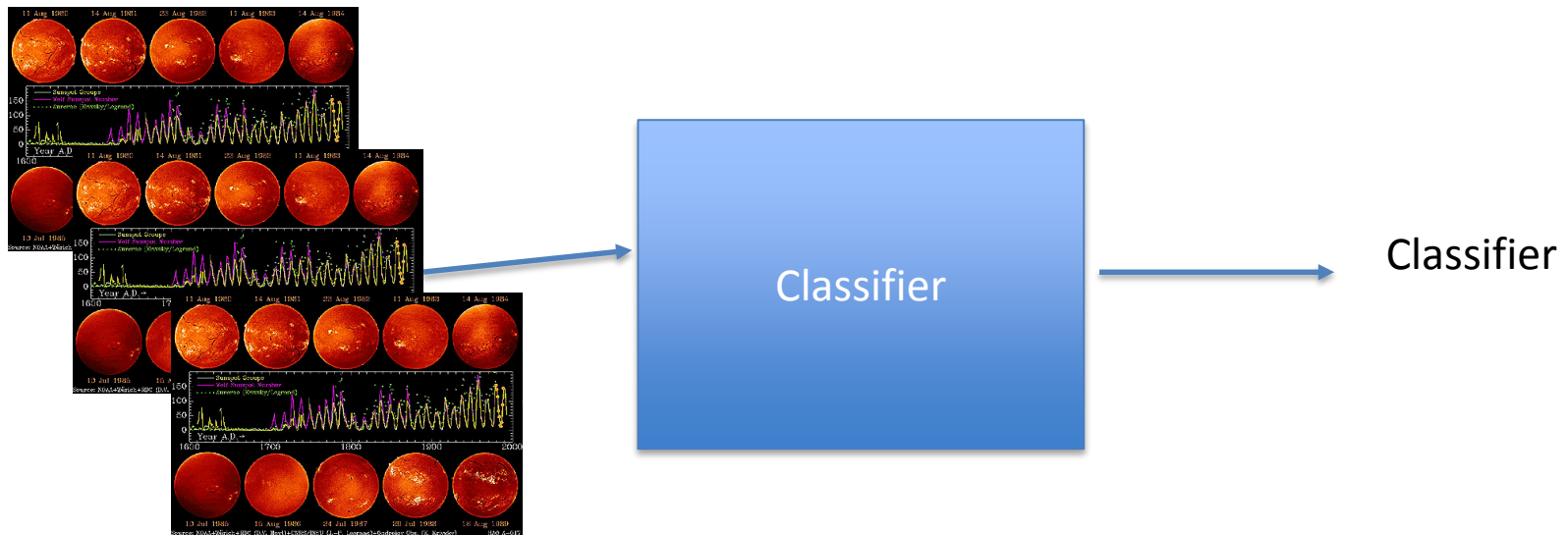
# Features

- ## Wavelet Transform

# Features

- Trend extraction
  - Differencing of time series.



Smoothed and Unsmoothed function, range= 1.5, qmax= 7, nmin= 25
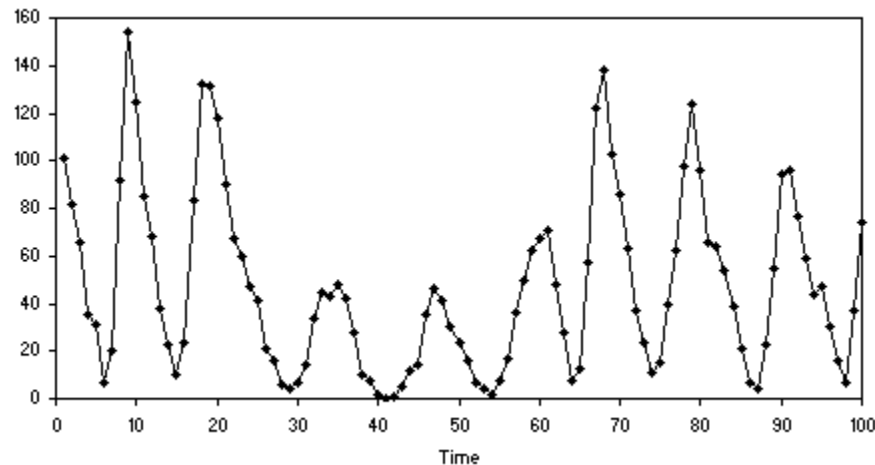
First derivative

Second derivative

# Features

- Full data extraction
  - Only if we have small, multiple time-series.

# Testing-Training

- How do we select the training and testing data?

  - Just using our traditional approach of random sampling won't do anything for us.

# There is no one approach

- Divide the time series in N consecutive folds to do CV.
  - Do not take random points.
- You can train with only half the image and test with the other half.
  - 2-Fold CV
- Testing and Training have to be roughly the same size.
  - You can and will overfit time series data. Regardless of what you do.
    - Course of non-stationarity

# Real time data series analysis

- What is needed:
  - If the problem is regression in real time:
    - Gradient descent, and training as we get new data.
    - We can't have an obscene number of features.
- Which algorithms can do this:
  - SVMs are surprisingly bad to do things in real-time.
  - Neural Nets and Deep Learning are getting there, but they are also not so hot.

# Recurrent Neural Networks