

A photograph of the Princeton University Chapel, a Gothic-style building with multiple towers and spires, set against a warm, orange-hued sky at sunset or sunrise. The building is partially obscured by trees in the foreground.

# **MOOC LAB 5: FINANCIAL MARKET REGIME PREDICTION WITH MACHINE LEARNING**

PRINCETON UNIVERSITY

FEBRUARY 2021



# Outline

---

- Regime-based Models in Finance
- Regime Identification
  - Trend Filtering
- Machine Learning in Finance
  - Overview of Machine Learning Methods
    - Supervised Learning: Regression and Classification
    - Unsupervised Learning
  - Applications in Economics and Finance
  - Classification Methods
    - Logistic Regression, Decision Trees and Ensemble Methods
- Machine Learning Approach in Regime Prediction
  - Dataset description
  - Model framework for regime prediction
  - Results
- References

# Why are Regime-Based Models used in Finance?

---

- Idea of regime change is **natural and intuitive**
- Regimes identified by econometric methods often correspond to **different periods in regulation, policy and other secular changes**
  - interest rate behavior markedly changed from 1979 through 1982
  - in equities, different regimes correspond to periods of high and low volatility, and bull and bear market periods
- Regime switching models can **capture stylized behavior of many financial series**
  - including fat tails
  - persistently occurring periods of turbulence followed by periods of low
  - skewness
  - time-varying correlations
- Regime switching models started to appear in finance literature in 1989
- For a recent overview of the concept
  - Ang and Timmerman (2011) "Regime Changes and Financial Markets"

# Trend Filtering

---

- We are given a time series  $\mathbf{y}_t, t = 1, \dots, n$ , assumed to consist of an underlying slowly varying trend  $\mathbf{x}_t$  and a more rapidly randomly varying component  $\mathbf{z}_t$ , which is equal to  $\mathbf{z}_t = \mathbf{y}_t - \mathbf{x}_t$
- **Goal:** To estimate the underlying trend in time series data,  $\mathbf{x}_t$
- Optimization problem formulation for  $l_1$ - trend filtering for  $y \in \mathbb{R}^n$

$$\hat{\mathbf{x}} = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \lambda \|D\mathbf{x}\|_1$$

where  $\lambda \geq 0$  is the regularization parameter and  $D \in \mathbb{R}^{(n-1) \times n}$  is the first order difference matrix

$$D = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 & 0 \\ 0 & 1 & -1 & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & & \cdot & \cdot \\ \cdot & \cdot & \cdot & & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & -1 \end{bmatrix} \in \mathbb{R}^{(n-1) \times n}$$

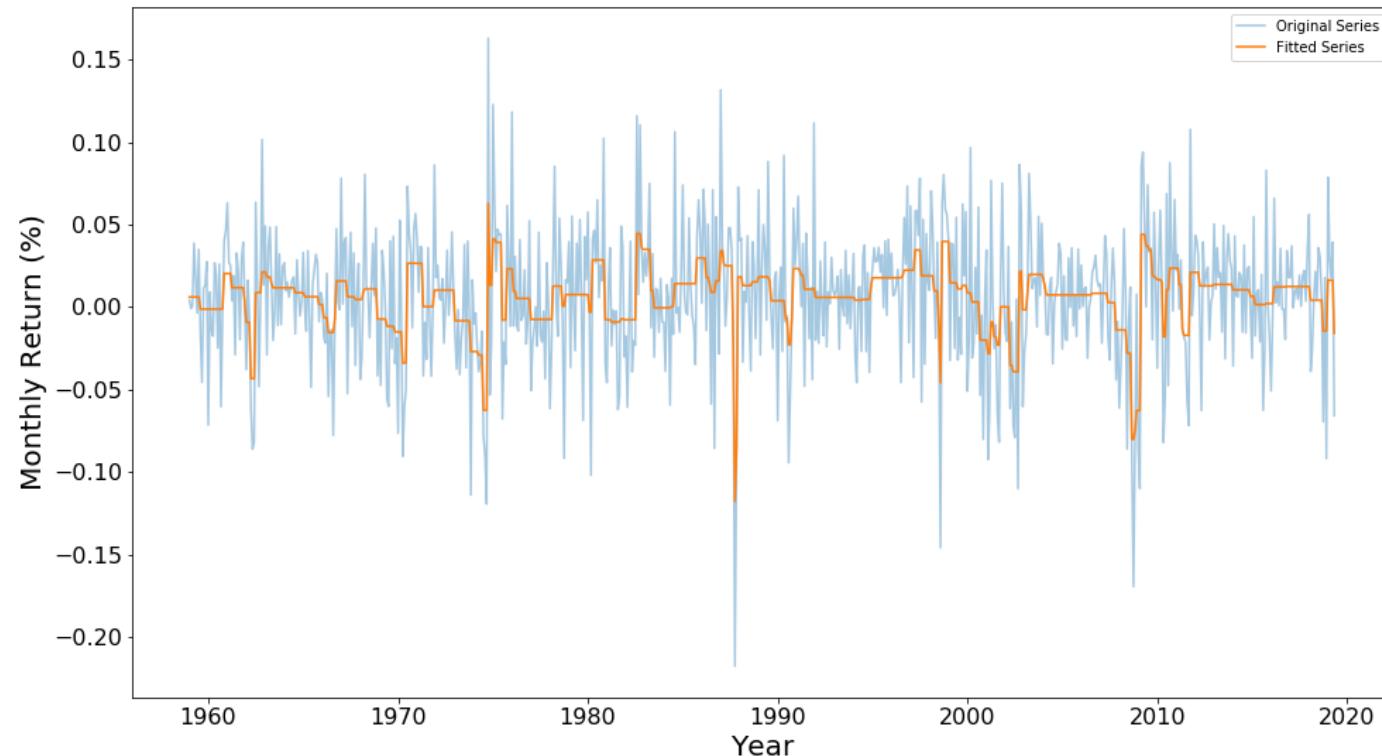
## Properties of $l_1$ - Trend Filtering Estimator

---

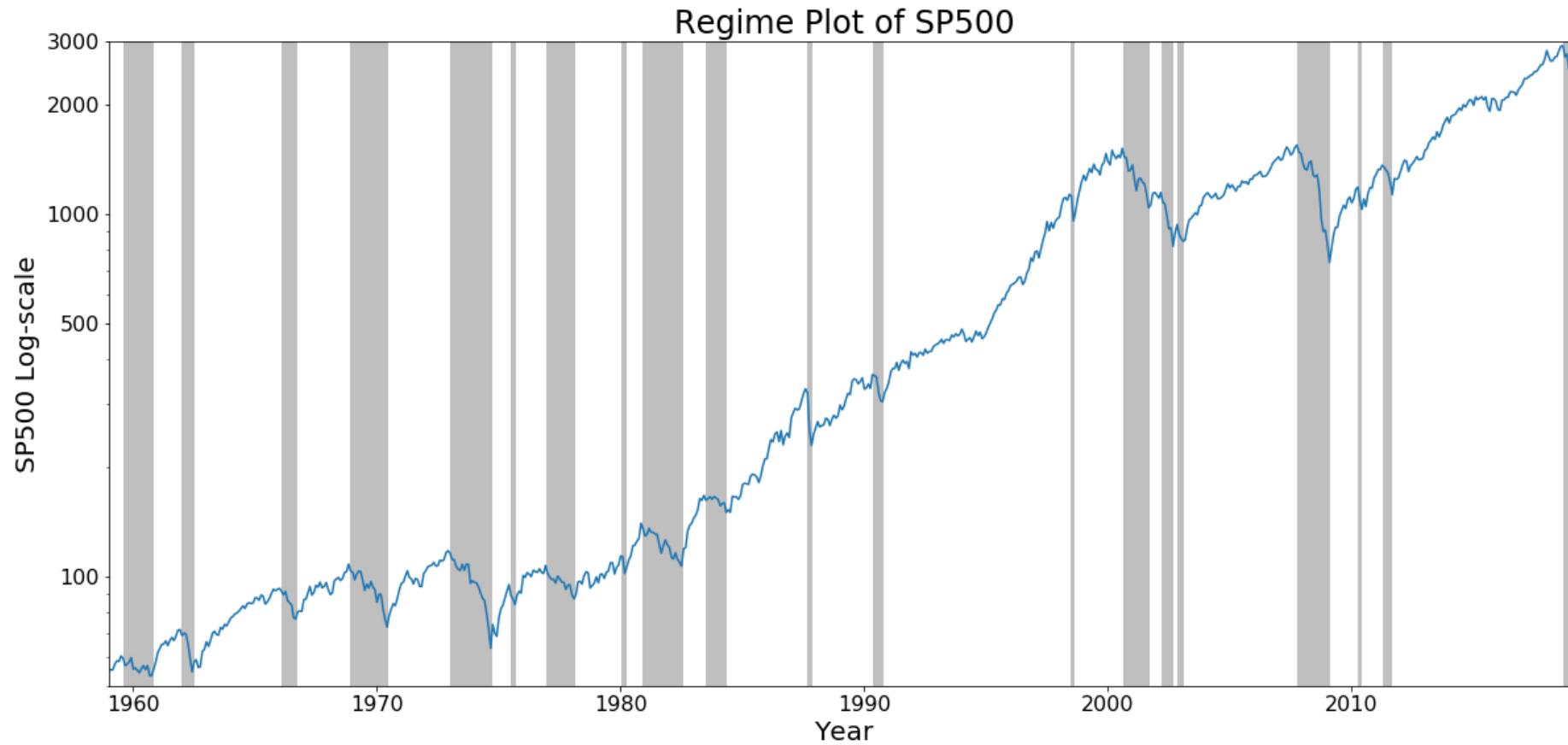
- Convex optimization problem
- The output will be a **piece-wise constant time series** which can be conveniently identified as a regime
- It is a **non-parametric regression model** which means it doesn't have a fixed structure of model and the nature of the parameters can be quite flexible.
- The estimator converges to **original data**  $y$  as  $\lambda \rightarrow 0$ , and converges to the **best constant fit** (i.e. the estimator is constant throughout the whole time) as  $\lambda \rightarrow \infty$

# $l_1$ - Trend Filtering on S&P500 Returns

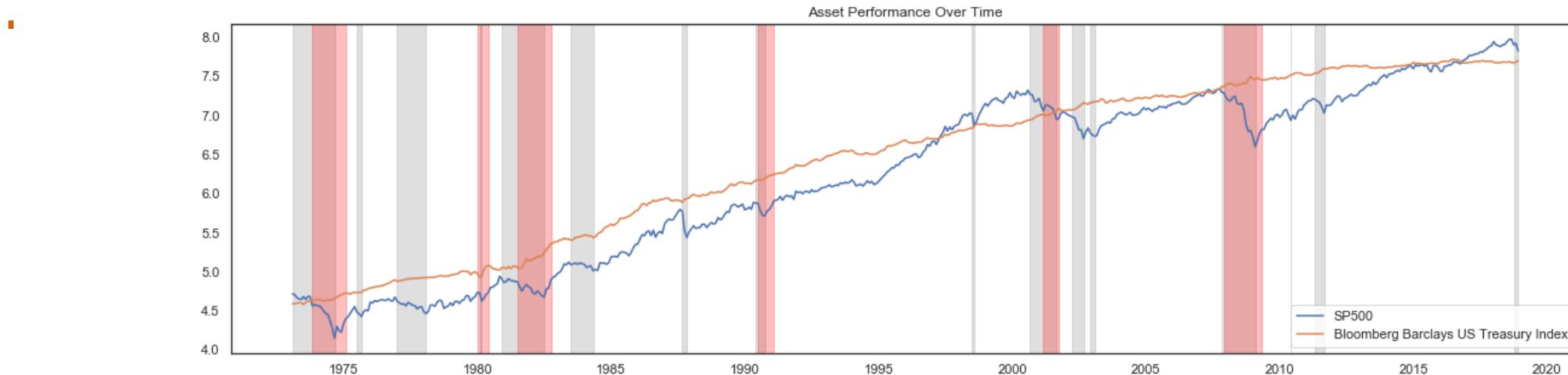
- Input: monthly S&P500 return data from 1959 until 2020



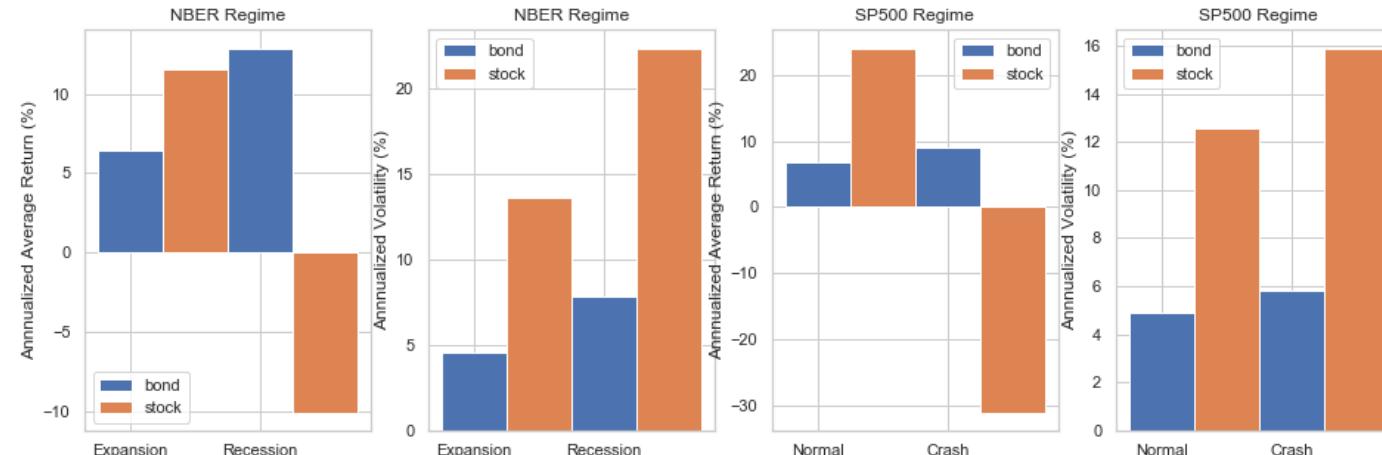
# $l_1$ - Trend Filtering on S&P500 Returns



# Asset Performances over Different Regimes



Red areas denote **recession periods** and gray areas show **stock market crashes** identified by trend filtering algorithm



# Machine Learning - Overview

## Supervised Learning

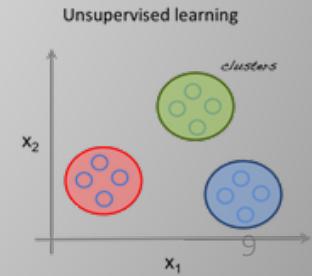
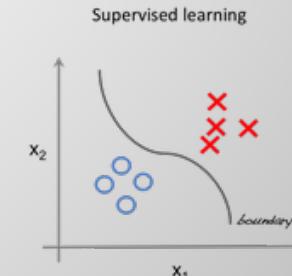
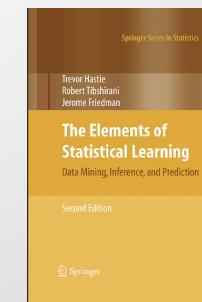
use an algorithm to learn the mapping function from the input ( $x$ ) to the output ( $y$ )

- Regression – continuous output
  - Linear Regression and it's cousins (LASSO, Ridge, ElasticNet, Best-subset)
  - Tree based models (Random Forest, Boosting Trees)
- Classification – categorical output
  - Logistic Regression
  - Tree based models (Random Forest, Boosting Trees)
  - Support Vector Machines

## Unsupervised Learning

model the underlying structure or distribution in the data ( $x$ )

- Clustering – want to discover inherent groupings in the data
  - k-means
- Anomaly detection – fraud detection
- Principal Component Analysis



# Machine Learning Applications in Finance and Economics

- Recent papers and books on machine learning in econ and finance

*Journal of Economic Perspectives—Volume 28, Number 2—Spring 2014—Pages 3–28*

## Big Data: New Tricks for Econometrics<sup>1</sup>

Hal R. Varian

Computers are now involved in many economic transactions and can capture and analyze associated data sets that are too large to be handled by traditional statistical and econometric techniques such as regression often used well, but there are issues unique to big datasets that may require different tools.

In this essay, I will describe three ways that machine learning can help analyze big data. I believe that these methods have a lot to offer and should be more widely known and used by economists. In fact, my standard advice to graduate students these days is to go to the computer science department and take a class in machine learning. There have been very fruitful collaborations between computer scientists and statisticians in the last decade or so, and I expect collaborations between computer scientists and economists will also be productive in the future.

<sup>1</sup> Hal Varian is Chief Economist, Google Inc., Mountain View, California, and Emeritus Professor of Economics, University of California, Berkeley, California. His email address is hal@econ.berkeley.edu.

<sup>2</sup> To access the Appendix and disclosure statement, visit <http://dx.doi.org/10.1257/jep.28.2.5>

*Journal of Economic Perspectives—Volume 31, Number 2—Spring 2017—Pages 87–106*

## Machine Learning: An Applied Econometric Approach

Sendhil Mullainathan and Jann Spiess

Machines are increasingly doing “intelligent” things: Facebook recognizes faces in photos, Siri understands voices, and Google translates websites. The fundamental insight behind these breakthroughs is as much statistical as it is computational. Machine learning has become popular because it stopped approaching intelligence tasks procedurally and began tackling them empirically. Face recognition algorithms, for example, do not follow a hardwired rule to scan for certain pixel combinations. Instead, they learn an underlying model from a large dataset of labeled faces. These algorithms use a large dataset of photos labeled as having a face or not to estimate a function  $f(x)$  that predicts the presence of a face from pixels  $x$ . This similarity to econometrics raises questions: Are these machine learning tools just glorified regression? What does it mean to “learn”? If there are fundamentally new empirical tools, how do they fit with what we know?

In empirical economics, how can we use them?<sup>2</sup>

We present a way of thinking about machine learning that gives it its own place in the econometric toolbox. Central to our understanding is that machine learning

<sup>1</sup> In this journal, Varian (2014) provides an excellent introduction to many of the more novel tools and “tricks” from machine learning, such as dimensionality reduction and regularization. Emai and Levin (2016) offer a good introduction to machine learning in economics. Belloni et al. (2015) provide a good introduction to machine learning in economics, including an introduction on how LASSO (and close cousins) can be used for inference in high-dimensional data. Athey (2015) provides a brief overview of how machine learning relates to causal inference.

<sup>2</sup> Sendhil Mullainathan is the Robert C. Waggoner Professor of Economics and Jean Spiess is a PhD candidate in Economics, both at Harvard University, Cambridge, Massachusetts. Their email addresses are mullainathan@harvard.edu and jspies@fas.harvard.edu.

<sup>3</sup> For additional material such as appendices, tables, and author disclosure statements, see the article page at <https://doi.org/10.1257/jep.31.2.87>



## The Impact of Machine Learning on Economics

Susan Athey

Chapter in NBER book *The Economics of Artificial Intelligence: An Agenda* (2019), Ajay Agrawal, Joshua Gans, and Avi Goldfarb, editors (p. 507 - 547)

Conference held September 13-14, 2017

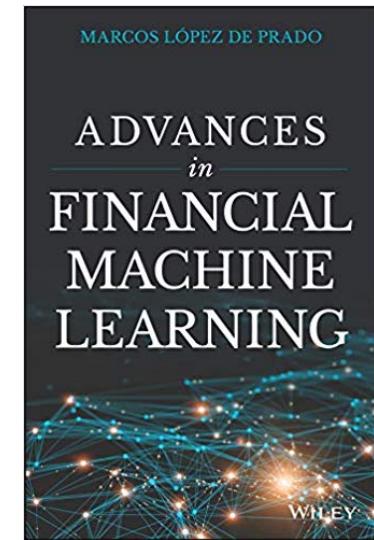
Published in May 2019 by University of Chicago Press

© 2019 by the National Bureau of Economic Research

This paper provides an assessment of the early contributions of machine learning to economics, as well as predictions about its future contributions. It begins by briefly overviews some themes from the literature on machine learning, and then draws some contrasts with traditional approaches to estimating the impact of counterfactual policies in economics. Next, we review some of the initial “off-the-shelf” applications of machine learning to economics, including applications in analyzing text and images. We then describe new types of questions that have been posed surrounding the application of machine learning to policy problems, including “prediction policy problems,” as well as considerations of fairness and manipulability. We present some highlights from the emerging econometric literature combining machine learning and causal inference. Finally, we overview a set of broader predictions about the future impact of machine learning on economics, including its impacts on the nature of collaboration, funding, research tools, and research questions.

This chapter is no longer available for free download, since the book has been published. To obtain a copy, you must buy the book.

Order from Amazon.com

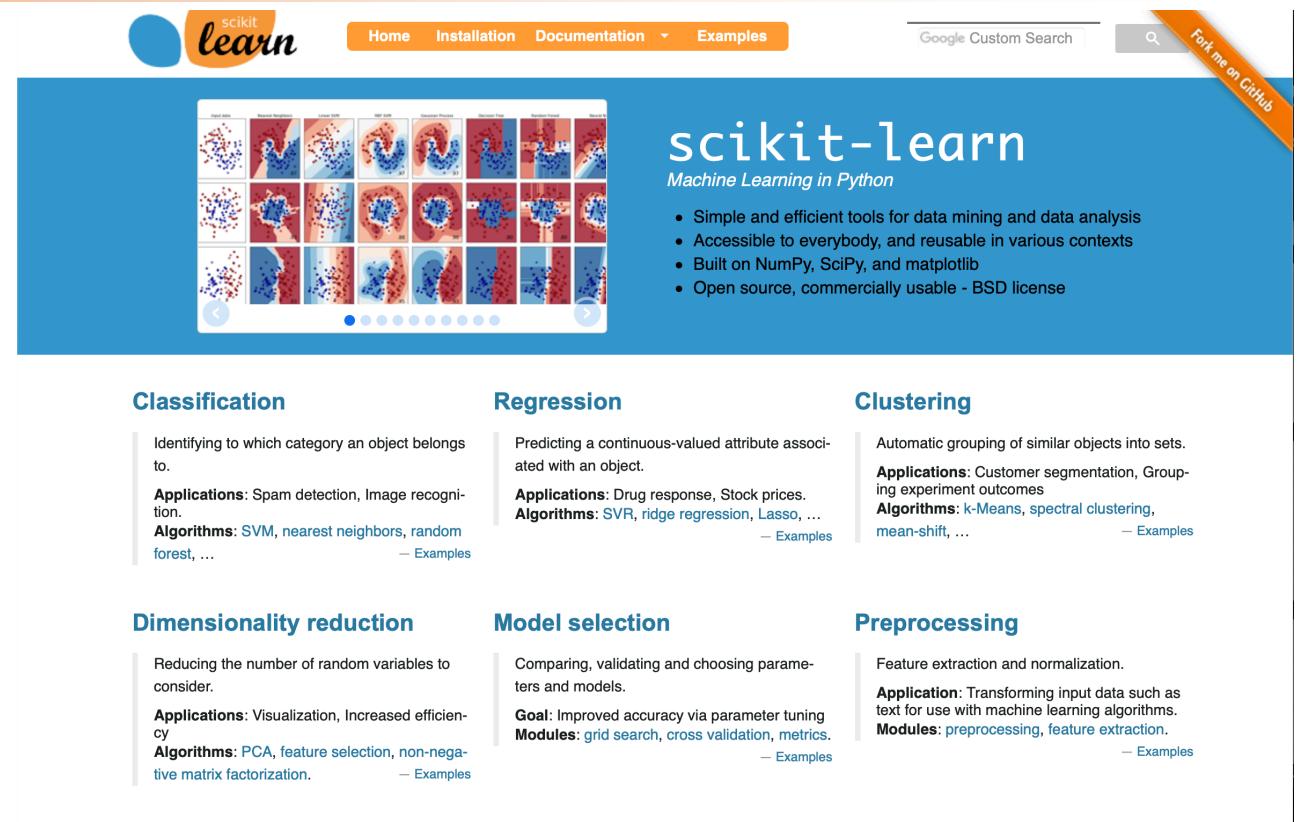


- Various applications in finance

- factor models, asset pricing, regime prediction, credit default prediction, news sentiment analysis

# Machine Learning – Software Package in Python

- [scikit-learn](#) is the machine learning library in Python
- Includes many regression and classification algorithms
- [sklean.datasets](#) package contains some small toy datasets



# Example: Linear Regression in Scikit-Learn

- [Diabetes dataset](#)
- n=442 diabetes patients
- 10 features: age, sex, body mass index, average blood pressure, and six blood serum measurements
- Target: a quantitative measure of disease progression

```
In [4]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn import datasets, linear_model
from sklearn.metrics import mean_squared_error, r2_score
```

```
In [49]: # Load the diabetes dataset
diabetes = datasets.load_diabetes() #returns a dictionary-like object
type(diabetes)
```

```
Out[49]: sklearn.utils.Bunch
```

```
In [50]: X=diabetes.data
y=diabetes.target
df=pd.DataFrame(diabetes.data,columns=diabetes.feature_names)
df['Target']=y
df.head()
```

	age	sex	bmi	bp	s1	s2	s3	s4	s5	s6	Target
0	0.038076	0.050680	0.061696	0.021872	-0.044223	-0.034821	-0.043401	-0.002592	0.019908	-0.017646	151.0
1	-0.001882	-0.044642	-0.051474	-0.026328	-0.008449	-0.019163	0.074412	-0.039493	-0.068330	-0.092204	75.0
2	0.085299	0.050680	0.044451	-0.005671	-0.045599	-0.034194	-0.032356	-0.002592	0.002864	-0.025930	141.0
3	-0.089063	-0.044642	-0.011595	-0.036656	0.012191	0.024991	-0.036038	0.034309	0.022692	-0.009362	206.0
4	0.005383	-0.044642	-0.036385	0.021872	0.003935	0.015596	0.008142	-0.002592	-0.031991	-0.046641	135.0

In

```
[76]: # Use only one feature - bmi  
diabetes_X = diabetes.data[:, np.newaxis, 2]  
# Split the data into training/testing sets  
diabetes_X_train = diabetes_X[:-20]  
diabetes_X_test = diabetes_X[-20:]  
  
# Split the targets into training/testing sets  
diabetes_y_train = diabetes.target[:-20]  
diabetes_y_test = diabetes.target[-20:]
```

In

```
[80]: print('Size of Training Set: ' +str(diabetes_X_train.shape[0]))  
print('Size of Test Set: ' +str(diabetes_X_test.shape[0]))
```

Size of Training Set: 422

Size of Test Set: 20

- Split dataset into training and test sets
- Training set is used for model training
- Test set is used for out-of-sample prediction to evaluate model performance

In [89]:

```
# Create linear regression object
regr = linear_model.LinearRegression()

# Train the model using the training sets
regr.fit(diabetes_X_train, diabetes_y_train)

# Make predictions using the testing set
diabetes_y_pred = regr.predict(diabetes_X_test)

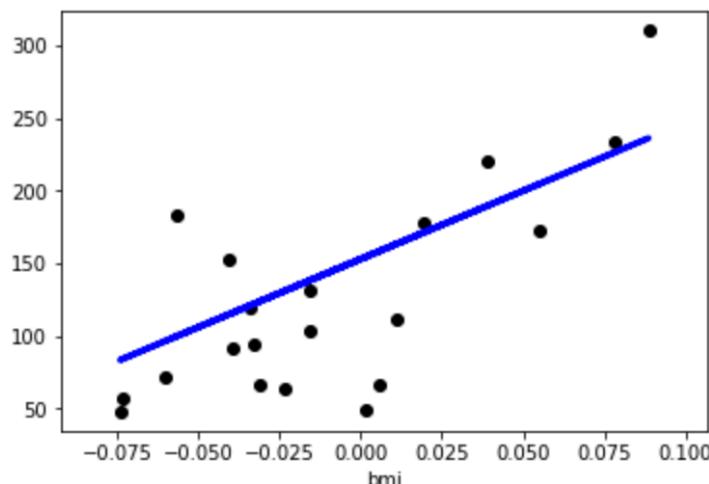
# The coefficients
print('Coefficients: \n', regr.coef_)

# the mean-squared error
print("Mean squared error: %.2f"
      % mean_squared_error(diabetes_y_test, diabetes_y_pred))
# Explained variance score: 1 is perfect prediction
print('Variance score: %.2f' % r2_score(diabetes_y_test, diabetes_y_pred))

# Plot outputs
plt.scatter(diabetes_X_test, diabetes_y_test, color='black')
plt.plot(diabetes_X_test, diabetes_y_pred, color='blue', linewidth=3)
plt.xlabel('bmi')
plt.show()
```

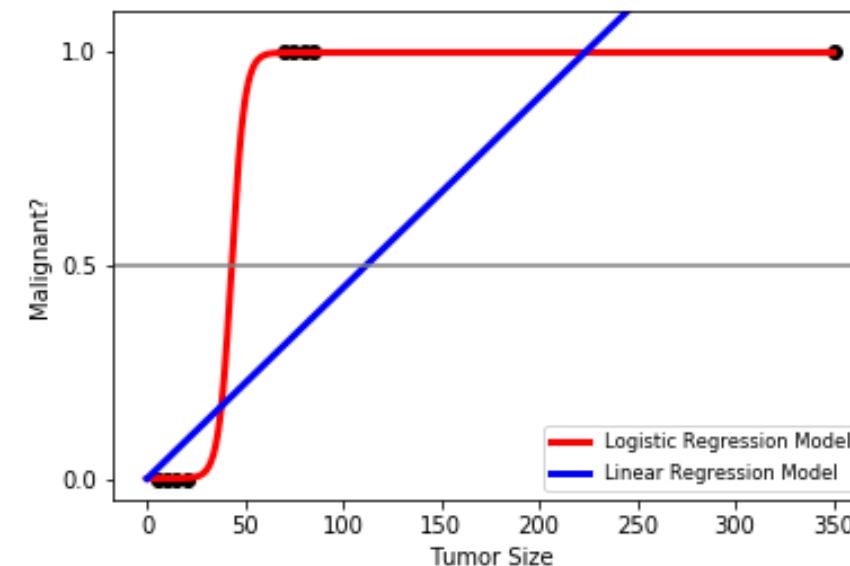
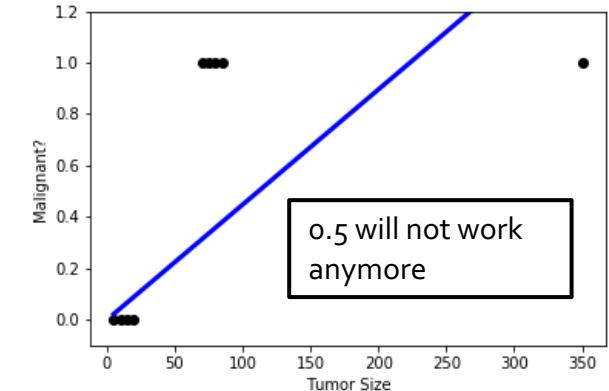
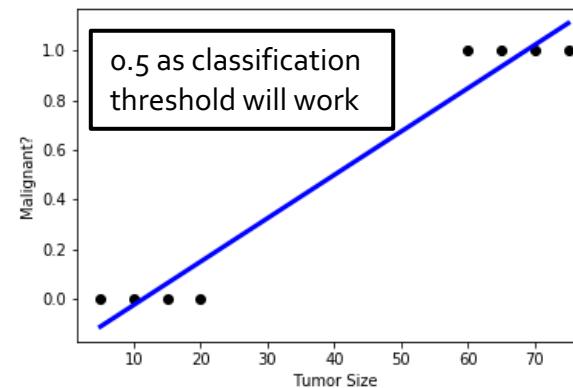
- `sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=None)`
- `.fit` is used for training model the model on training set
- `.predict` is used for predicting on test set
- Blue line shows fitted regression line and black points denote data points from test set

Coefficients:  
[938.23786125]  
Mean squared error: 2548.07  
Variance score: 0.47



# What is Wrong with using Linear Regression for Classification?

- You can label one of the classes with 0 and other with 1 and use linear regression but there are few problems
  - It simply interpolate between points so does not output probabilities
  - Does not extend to classification problems with multiple classes
- A solution for classification is **logistic regression**



# Logistic Regression

- Logistic regression model uses the logistic function to squeeze the output of a linear equation between 0 and 1.

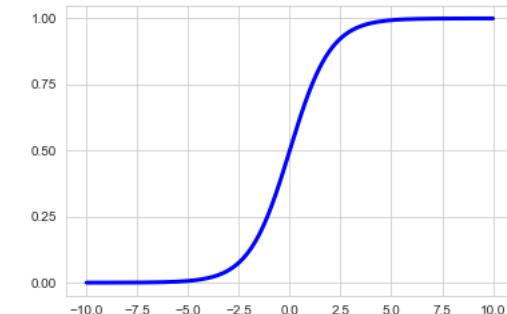
$$\text{logistic}(x) = \frac{1}{1+e^{-x}}$$

- From Linear Regression

$$y^{(i)} = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_p x_p^{(i)}$$

- to Logistic Regression

$$P(y^{(i)} = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_p x_p^{(i)})}}$$



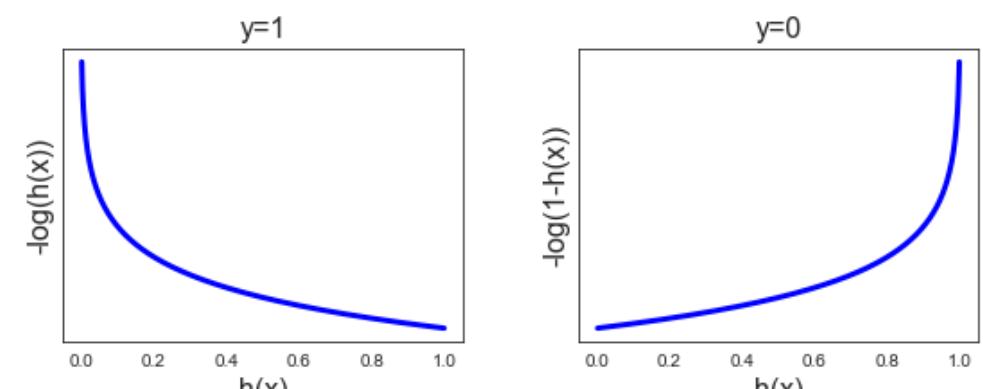
# Logistic Regression Learning

- Logistic Regression for binary classification

$$P(y = 1|x; \beta) = h_\beta(x) = \frac{1}{1 + e^{(-\beta^T x)}} \Rightarrow \begin{cases} y = 1 & \text{if } h(x) \geq 0.5 \\ y = 0 & \text{if } h(x) < 0.5 \end{cases}$$

- Cost function – logistic loss function

$$cost(h_\beta(x), y) = \begin{cases} -\log(h_\beta(x)) & \text{if } y = 1 \\ -\log(1 - h_\beta(x)) & \text{if } y = 0 \end{cases}$$



- Set up optimization problem as cost minimization

$$\min_{\beta \in R^p} J(\beta) = \frac{1}{n} \sum_i cost(h_\beta(x^{(i)}), y^{(i)})$$

$$= \frac{1}{n} \sum_{i=1}^n -y^{(i)} \log(h_\beta(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_\beta(x^{(i)}))$$

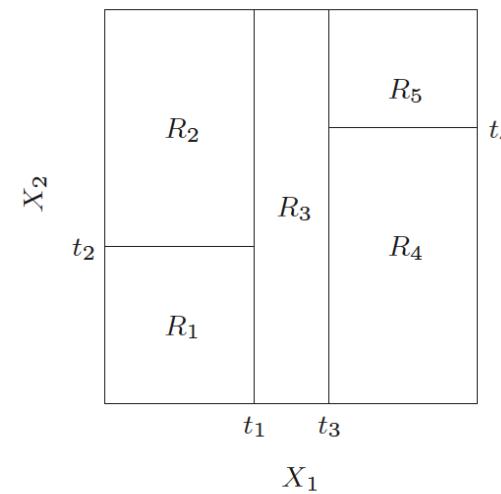
# Regularized Logistic Regression

- Regularization terms can be added
  - in high dimensional problems to avoid overfitting
  - if you have prior assumptions on data – e.g. sparsity
- For  $l_1$  regularized logistic regression  $\min_{\beta \in R^p, \lambda \geq 0} J(\beta) + \lambda ||\beta||_1$
- For  $l_2$  regularized logistic regression  $\min_{\beta \in R^p, \lambda \geq 0} J(\beta) + \lambda ||\beta||_2^2$

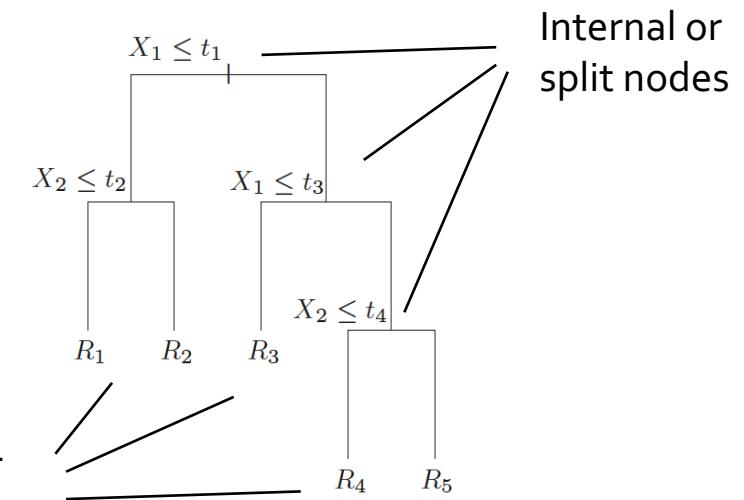
# Decision Trees

- Linear and Logistic Regression models fail in situations where
  - the relationship between features and outcome is **nonlinear**
  - **features interact** with each other
- Tree based models **split the data multiple times** according to certain cutoff values in the features.
- To predict the outcome in each leaf node, the **average outcome of the training data in this node** is used. Trees can be used for **classification and regression**.

Partition of a two-dimensional feature space by recursive binary splitting



Terminal or  
leaf nodes



# Decision Tree Learning

- The classification and regression trees (**CART**) algorithm is probably the most popular algorithm for tree induction
- The goal is to find boxes  $R_1, \dots, R_J$  that minimizes residual sum of squares, given by

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

where  $\hat{y}_{R_j}$  is the mean response for the training data in  $j^{th}$  box

- Decision trees are **easy to use, interpretable, flexible** but a single decision tree is quite **unstable** and usually does not have good prediction accuracy
- However, by **aggregating many decision trees**, the predictive performance of trees can be substantially improved

# Ensemble Methods

---



## ■ Bagging

- Generate  $B$  different bootstrapped training datasets.
- Train the learner on the  $b^{th}$  bootstrapped training set and get  $\widehat{f}^b(x)$ , the prediction at point a point  $x$
- Then average all predictions to obtain  $\widehat{f}_{bag}(x) = \frac{1}{B} \sum_b \widehat{f}^b(x)$

## ■ Random Forest

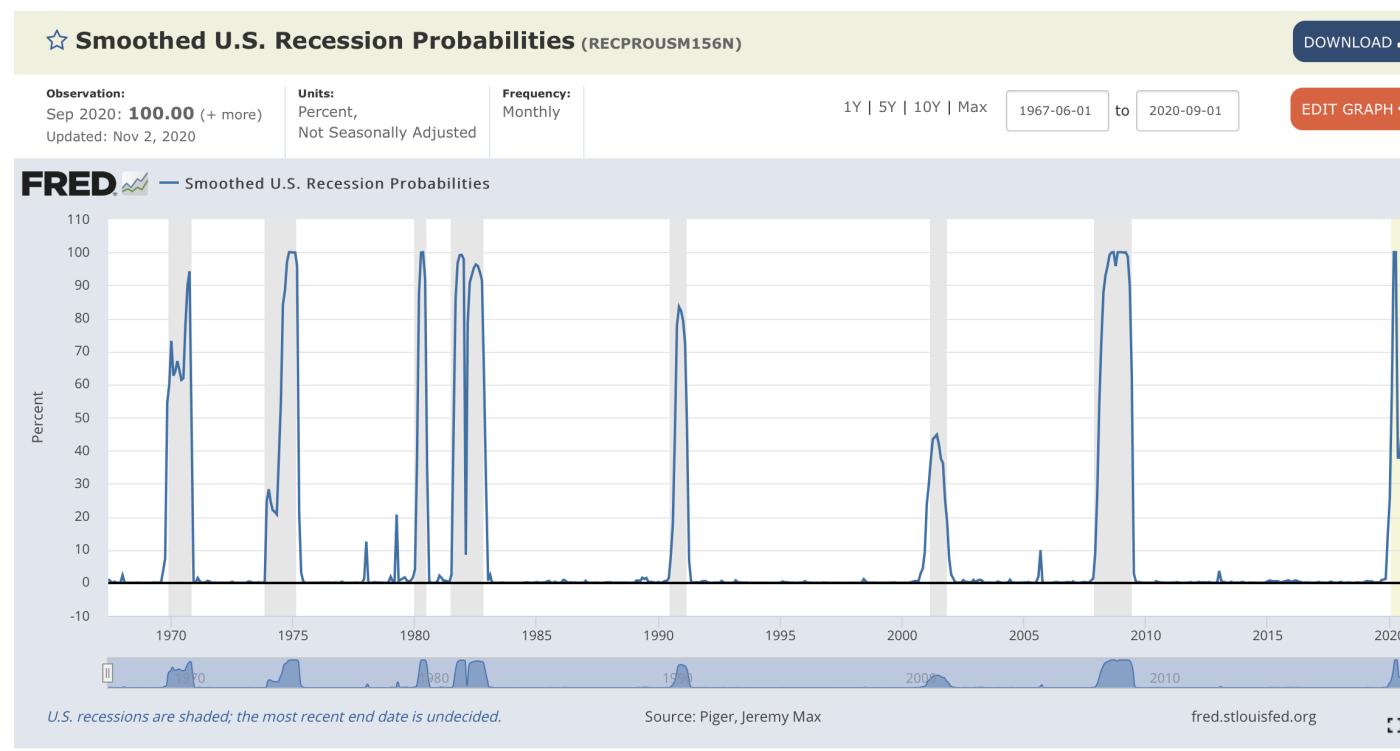
- Provides an improvement over bagged trees by decorrelating trees
- This reduces the variance when we average the trees

## ■ Boosting

- “Boosting is one of the most powerful learning ideas introduced in the last twenty years” – Hastie et.al (2009)
- Trees are grown sequentially: each tree is grown using information from previously grown trees

# Business Cycle Prediction

- Months and quarters of macroeconomic time series are separated into periods of **recession and expansion**
- NBER provides a chronology of business cycle dates but its methodology is **not explicitly formalized**,

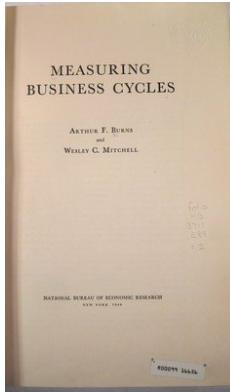


## Business Cycle Dating Committee Announcement June 8, 2020



# Literature Overview

## Initial work



James H. Stock  
KENNEDY SCHOOL OF GOVERNMENT  
HARVARD UNIVERSITY  
CAMBRIDGE, MASSACHUSETTS

Mark W. Watson  
DEPARTMENT OF ECONOMICS  
NORTHWESTERN UNIVERSITY  
EVANSTON, ILLINOIS

### New Indexes of Coincident and Leading Economic Indicators

#### 1. Introduction

During six weeks in late 1937, Wesley Mitchell, Arthur Burns, and their colleagues at the National Bureau of Economic Research developed a list of leading, coincident, and lagging indicators of economic activity in the United States as part of the NBER research program on business cycles. Since their development, these indicators, in particular the leading and coincident indexes constructed from these indicators, have played an important role in summarizing and forecasting the state of macroeconomic activity.

## A Comparison of the Real-Time Performance of Business Cycle Dating Methods

Marcelle CHAUVET  
Department of Economics, University of California, Riverside, CA 92521 ([chauvet@ucr.edu](mailto:chauvet@ucr.edu))  
Jeremy PIGER  
Department of Economics, University of Oregon, Eugene, OR 97403 ([j piger@uoregon.edu](mailto:j piger@uoregon.edu))

We evaluate the ability of formal rules to establish U.S. business cycle turning point dates in real time. We consider two approaches, a nonparametric algorithm and a parametric Markov-switching dynamic factor model. Using a new "real-time" dataset of coincident monthly variables, we find that both approaches would have accurately identified the NBER business cycle chronology had they been in use over the past 30 years, with the Markov-switching model most closely matching the NBER dates. Further, both approaches, and particularly the Markov-switching model, yielded significant improvement over the NBER in the speed with which business cycle troughs were identified.

KEY WORDS: Dynamic-factor model; Markov-switching; Recessions; Turning point; Vintage data.

#### 1. INTRODUCTION

There is a long tradition in business cycle analysis of separating periods in which there is broad economic growth, called expansions, from periods of broad economic contraction, called recessions. Understanding these phases and the transitions between them has been the focus of much macroeconomic research over the past century. In the United States, the National Bureau of Economic Research (NBER) establishes a chronology

that the NBER established the correct turning point dates in real time, thus making the NBER chronology the standard for accuracy.

Why are we interested in the speed with which business cycle turning points can be identified? The NBER is likely more concerned with establishing the correct turning point dates than establishing these dates quickly, which breeds additional caution. This caution comes at a low cost if the primary objective is to

## Recent work with machine learning methods

### BOOSTING RECESSIONS

Serena Ng

Department of Economics  
Columbia University\*

August 2013

#### Abstract

This paper explores the effectiveness of boosting, often regarded as the state of the art classification tool, in giving warning signals of recessions three, six and twelve months ahead. Boosting is used to screen as many as 1500 potentially relevant predictors consisting of 132 real and financial time series and their lags. Estimation over the full sample 1961:1–2011:12 finds that there are fewer than ten important predictors and the identity of these variables change with the forecast horizon. There is a distinct difference in the size and composition of the relevant predictors across horizons. The results show that the most important predictors for each type of the term and default spreads are recession specific. The Ann spread is the most robust predictor of recessions three and six months ahead, while the risky bond and 5yr spreads are important for twelve months ahead predictions. Certain employment variables have predictive power for the two most recent recessions when the interest rate spreads were uninformative. Warning signals for the post 1990 recessions have been sporadic and easy to miss. The results underscore the challenge that changing characteristics of business cycles pose for predicting recessions.

Keywords: business cycle chronology, recession probabilities, classification.

JEL Classification: C5,C6, C25, C35.

International Journal of Forecasting 35 (2019) 848–867  
Contents lists available at ScienceDirect  
International Journal of Forecasting  
journal homepage: [www.elsevier.com/locate/ijforecast](http://www.elsevier.com/locate/ijforecast)

### Recession forecasting using Bayesian classification

Troy Davig, Aaron Smalter Hall\*

Research Department, Federal Reserve Bank of Kansas City, Kansas City, MO 64198, USA

#### ABSTRACT

We demonstrate the use of a Naïve Bayes model as a recession forecasting tool. The approach is closely connected with Markov-switching models and logistic regression, but also has important differences. In contrast to Markov-switching models, our Naïve Bayes model treats National Bureau of Economic Research business cycle turning points as data, rather than as hidden states to be inferred by the model. Although Naïve Bayes and logistic regression are asymptotically equivalent under certain distributional assumptions, the assumptions do not hold for business cycle data. As a result, Naïve Bayes has a larger asymptotic error rate, but converges to the error rate more quickly than logistic regression, resulting in more accurate recession forecasts with limited data. We show that Naïve Bayes outperforms competing models and the Survey of Professional Forecasters consistently for real-time recession forecasting up to 12 months in advance. These results hold under standard error measures, and also under a novel measure that varies the penalty on false signals, depending on when they occur within a cycle; for example, a false signal in the middle of an expansion is penalized more heavily than one that occurs close to a turning point.

© 2018 International Institute of Forecasters. Published by Elsevier B.V. All rights reserved.

Received 10 November 2017; revised 14 December 2018; accepted 24 January 2019; available online 27 February 2019  
Peer review handled by Alexander James, Marcelle Chauvet, and Xiaodong Qiao

ALEXANDER JAMES is a co-chair of research at Promontech Technologies Inc., New York, NY, USA. [alex@promontech.com](mailto:alex@promontech.com)  
YASER S. ABU-MOSTAFA is a professor of electrical engineering and computer science at the California Institute of Technology in Pasadena, CA.  
[yaser@caltech.edu](mailto:yaser@caltech.edu)  
XIAODONG QIAO is a co-chair of research at Promontech Technologies Inc., New York, NY, USA. [xiaodong@promontech.com](mailto:xiaodong@promontech.com)

**ABSTRACT:** The authors introduce a novel application of support vector machines (SVM), an important machine learning algorithm, to determine the beginning and end of recessions in real time. Note that the NBER's official dating of recessions is not available until later, key for recessions, which are only determined months after the fact. The authors propose a SVM-based approach for this task, capturing all six recessions from 1973 to 2018 and providing the signal with minimal delay. The authors take advantage of the individual components of the model to improve the discrimination between bonds and bonds. A dynamic risk budgeting approach using SVM outputs appears superior to an equal-weighted contributions portfolio, especially average returns by 85 per annum without increased tail risk.

**TOPICS:** Big data/machine learning, financial crises and financial market history, portfolio construction, tail risks\*

**D**eal-time business cycle dating is of central importance in

investors with enough resources to use this information in their investment process may change their portfolios as the economy turns from growth to contraction. The National Bureau of Economic Research (NBER) provides the official dating of expansions and recessions. The NBER's Business Cycle Dating Committee periodically assesses the latest available data on the macroeconomy and determines if the economy is in an expansion or a recession.<sup>1,2</sup> The committee releases announcements about the date of the turning points. Because of the reliance on macroeconomic data, which may be revised or released with a lag, and the committee's conservative approach to turning points, the NBER has historically announced turning points with a delay of 4 to 21 months (Gusto and Piger 2017).

\*

Is it possible to identify business cycle turning points in a more timely manner? This is the focus of a large body of literature going back to Burns and Mitchell (1946) and greatly expanded in a series of papers

# Dataset Description

## ■ Output

- NBER business cycle dates used for recession periods

## ■ Features: 135 macroeconomic time series

- output and income,
- labor market
- housing
- consumption, orders and inventories,
- money and credit,
- interest and exchange rates,
- prices
- stock market

## ■ Sample period: 1959-2020



Michael W. McCracken  
Assistant Vice-President

Economic Data  
[FRED-MD and FRED-QD: Monthly and Quarterly Databases for Macroeconomic Research](#)

FRED-MD and FRED-QD: Monthly and Quarterly Databases for Macroeconomic Research

FRED-MD and FRED-QD are large macroeconomic databases designed for the empirical analysis of "big data." The datasets of monthly and quarterly observations mimic the coverage of datasets already used in the literature, but they add three appealing features. They are updated in real-time through the FRED database. They are publicly accessible, facilitating the replication of empirical work. And they relieve the researcher of the task of incorporating data changes and revisions (a task accomplished by the data desk at the Federal Reserve Bank of St. Louis).

The accompanying paper shows that factors extracted from the FRED-MD dataset share the same predictive content as those based on the various vintages of the so-called Stock-Watson data. In addition, it suggests that diffusion indexes constructed as the partial sum of the factor estimates can potentially be useful for the study of business cycle chronology.

[Working Paper](#) [FRED-MD Updated Appendix](#)

The monthly and quarterly databases are listed below by release date. MATLAB code is also available to analyze the data.

[FRED-Databases code \(.zip\)](#)  
[Historical Vintages of FRED-MD 1999-08 to 2014-12 \(.zip 45 MB\)](#)  
[All Post-2015 Vintages of FRED-MD \(.zip\)](#)

FRED-QD is a quarterly frequency companion to FRED-MD. It is designed to emulate the dataset used in "Disentangling the Channels of the 2007–2009 Recessions" by Stock and Watson (2012, NBER WP No. 18094) but also contains several additional series. Comments or suggestions are welcome.

[FRED-QD Updated Appendix](#)

Monthly Data	Changes to FRED-MD	Quarterly Data	Changes to FRED-QD
<a href="#">current.csv (current)</a>		<a href="#">current.csv</a>	
<a href="#">2019-10.csv</a>		<a href="#">2019-10.csv</a>	
<a href="#">2019-09.csv</a>		<a href="#">2019-09.csv</a>	
<a href="#">2019-08.csv</a>		<a href="#">2019-08.csv</a>	
<a href="#">2019-07.csv</a>		<a href="#">2019-07.csv</a>	

**Appendix**

The column TCODE denotes the following data transformation for a series  $x_t$ : (1) no transformation; (2)  $\Delta x_t$ ; (3)  $\Delta^2 x_t$ ; (4)  $\log(x_t)$ ; (5)  $\Delta \log(x_t)$ ; (6)  $\Delta^2 \log(x_t)$ ; (7)  $\Delta(x_t/x_{t-1} - 1.0)$ . The FRED column gives mnemonics in FRED followed by a short description. The comparable series in Global Insight is given in the column GSI.

Some series require adjustments to the raw data available in FRED. We tag these variables with an asterisk to indicate that they been adjusted and thus differ from the series from the source. A summary of the adjustments is detailed in the paper <https://research.stlouisfed.org/wp/2015/2015-012.pdf>.

Group 1: Output and income

id	tcodes	fred	description	gsi	gsi:description
1	1	5	RPI	M_14386177	PI
2	2	5	W875RX1	M_14526755	PI less transfers
3	6	5	INDPRO	M_116460980	IP: total
4	7	5	IPFPNSS	M_116460981	IP: products
5	8	5	IPFINAL	M_116461268	IP: final prod
6	9	5	IPCONGD	M_116460982	IP: cons gds
7	10	5	IPDCONGD	M_116460983	IP: cons dble
8	11	5	IPNCONGD	M_116460988	IP: cons nondble
9	12	5	IPBUSEQ	M_116460995	IP: bus eqpt
10	13	5	IPMAT	M_116461002	IP: matls
11	14	5	IPDMAT	M_116461004	IP: dble matls
12	15	5	IPNMAT	M_116461008	IP: nondble matls
13	16	5	IPMANSICS	M_116461013	IP: mfg
14	17	5	IPBS122s	M_116461276	IP: res util
15	18	5	IPFUELS	M_116461275	IP: fuels
16	19	1	NAPMPI	M_110157212	NAPM prodn
17	20	2	CUMFNS	M_116461602	Cap util

24

# Model Set-up

---

## Notations

- $y_{t+h}$ : regime at period  $t + h$  and  $y_{t+h} \in \{0,1\}$  (1 for *crash* and 0 *normal* regime)
- Model results include both nowcasting ( $h = 0$ ) and forecasting ( $h = 1, 3, 6, 12$ )
- $X_t$ : input variables at period  $t$  which includes both current values of variables as well as lags
- $\hat{y}_{t+h} = P(y_{t+h} = 1 | X_t)$  crash regime prediction probabilities
- 0.5 threshold can be used for binary classification
- **Linear classifiers** (Logistic Regression with various penalty terms) and **non-linear classifiers** (tree-based ensembles) are used in predictive model framework

## Classification Tree in Regime Prediction

- Recursive splits partition the sample space into  $M$  non-overlapping regions
  - $A_m^*$  denotes region  $m$  for  $m = 1, \dots, M$
  - $T_m^*$  training sample observations in each region  $m$  for  $m = 1, \dots, M$
- Prediction for  $y_{t+h}^c$  is equal to within region proportion of class  $c$

$$P_{A_m^*}^c = \frac{1}{T_m^*} \sum_{X_t \in A_m^*} I(y_{t+h} = c)$$

- The CART classifier is then:

$$\hat{y}_{t+h}^c(X_t) = \sum_{m=1}^M P_{A_m^*}^c I(X_t \in A_m^*)$$

## How is the recursive portioning implemented to arrive at the regions $A_m^*$ ?

- Splits are based on binary condition of the form  $X_{j,t} < \tau$  and  $X_{j,t} \geq \tau$ , where  $j$  and  $\tau$  can differ across splits
- From unsplit node  $A$ , for a given  $j$  and  $\tau^j$  group the data into two regions

$$A_L = \{X_t | X_{j,t} < \tau^j, X_t \in A\} \text{ and } A_R = \{X_t | X_{j,t} \geq \tau^j, X_t \in A\}$$

- To determine splitting variable,  $j$ , and the split threshold,  $\tau^j$  we want to **maximize a measure of homogeneity of class outcomes** in  $A_L$  and  $A_R$
- A common choice Gini Impurity

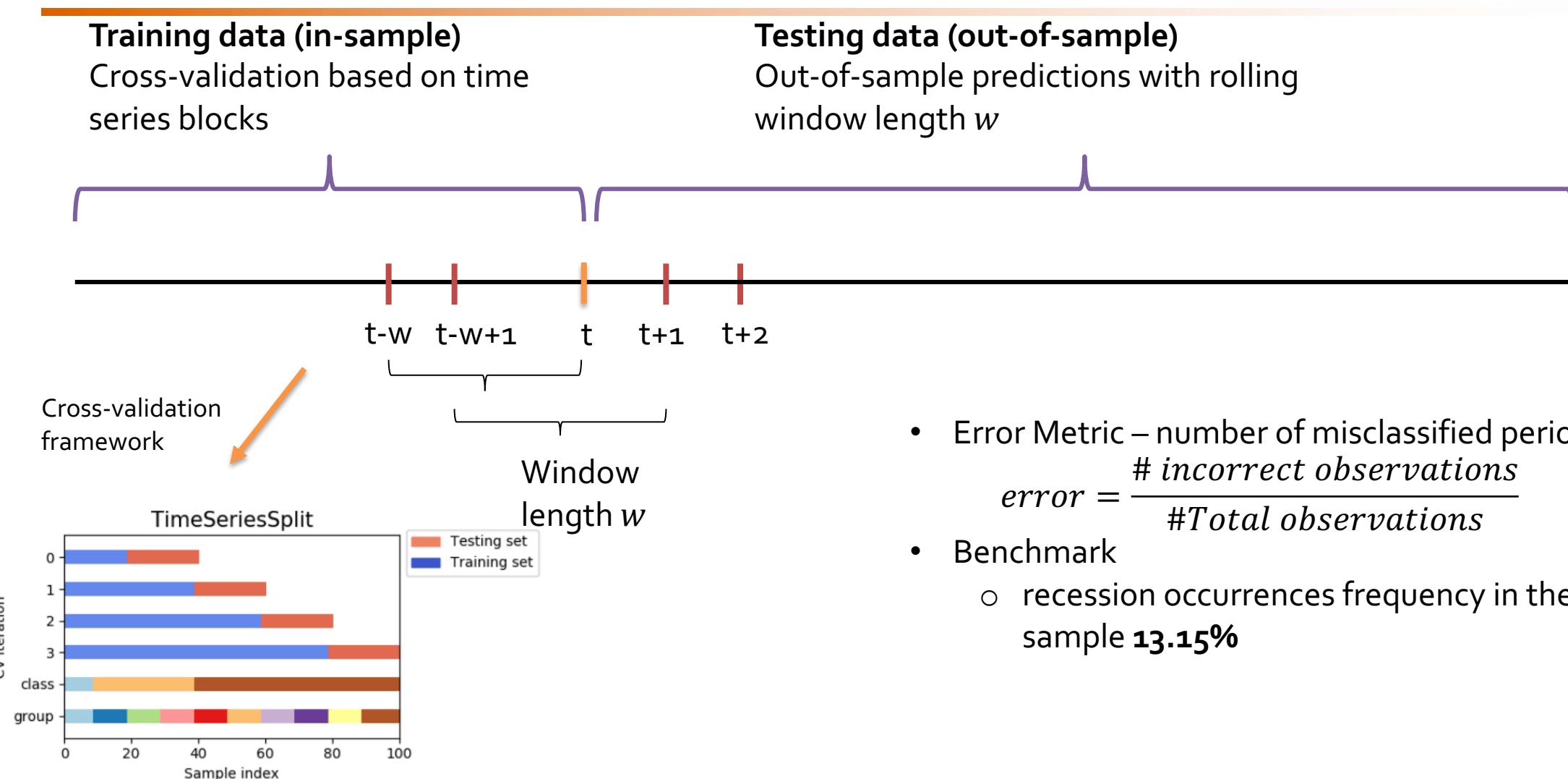
## Gini Impurity

---

$$G_L = \sum_{c=1}^C P_{A_L}^c (1 - P_{A_L}^c) \text{ and } G_R = \sum_{c=1}^C P_{A_R}^c (1 - P_{A_R}^c)$$

- Gini impurity
  - is bounded between 0 and 1
  - Value 0 indicates a pure region where only one class is present
  - Higher values of Gini impurity indicates greater class diversity
- The average Gini impurity for the new proposed regions is
$$\bar{G} = \frac{T^L}{T^L + T^R} G_L + \frac{T^R}{T^L + T^R} G_R$$
- $j$  and  $\tau^j$  are chosen to create regions  $A_L$  and  $A_R$  that minimize the average Gini impurity

# Predictive Model Framework



## Error Metrics

---

- Quadratic Probability Score (QPS)
  - evaluates regime predictions in terms of probabilities
  - In perfect prediction setting QPS will be 0
  - Lower QPS -> better model performance
- Classification Accuracy (ACC) for a given threshold  $c$ 
  - Natural classification metric
  - ACC should be at least better than benchmarks for each regime

$$QPS = \frac{1}{T} \sum_{t=1}^T (\hat{y}_t - y_t)^2$$

$$ACC = \frac{1}{T} \sum_{t=1}^T [\mathbb{1}_{\{\hat{y}_t \geq c\}} y_t + (1 - \mathbb{1}_{\{\hat{y}_t \geq c\}})(1 - y_t)]$$

# Class Imbalance Sensitive Metrics

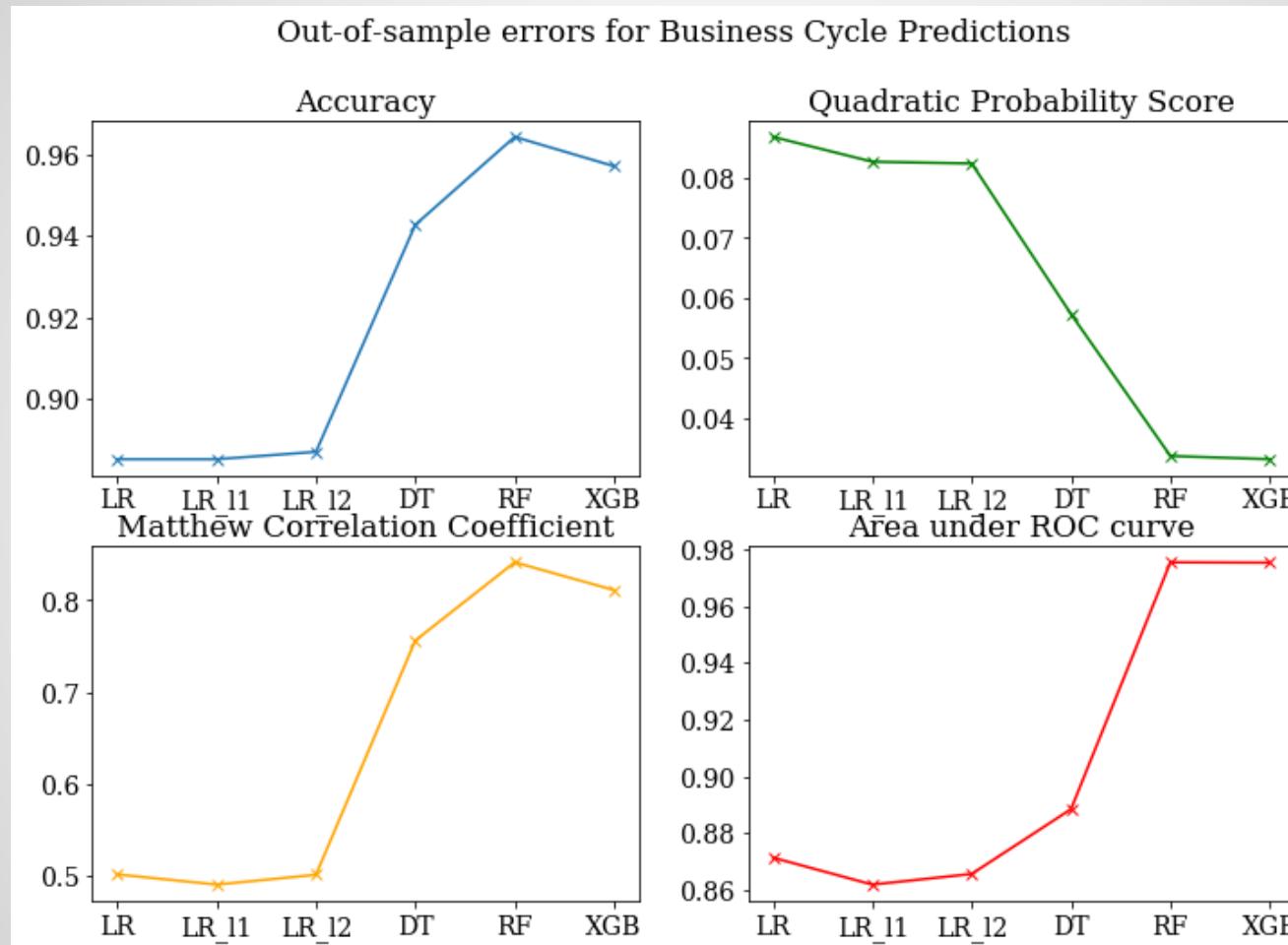
To avoid inflated error metrics, we will consider 2 other error metrics

- Receiver Operating Characteristic (ROC)
  - Common metric for imbalanced classification problems
  - ROC curve plots true positive rate against false positive rate
  - Area under the curve (AUC) is used to generate summary statistics for ROC metric
  - Higher AUC -> better model performance
- Matthew's Correlation Coefficient (MCC)
  - Used widely in computational biology
  - Can be calculated from confusion matrix
  - Higher MCC -> better model performance

Actual/Prediction	P	N
P	TP	FN
N	FP	TN

$$MCC = \frac{tp \times tn - fp \times fn}{\sqrt{(tp + fp)(tp + fn)(tn + fp)(tn + fn)}}$$

# Regime Prediction Error Metrics (1975-2020)



LR: Logistic Regression

LR\_I1: Logistic Regression with l1 penalty

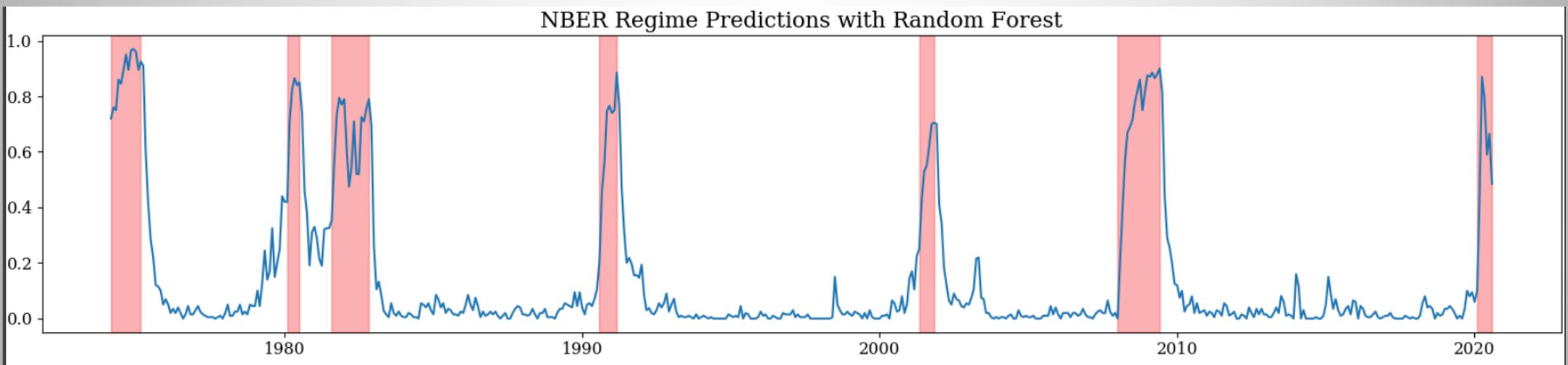
LR\_I2: Logistic Regression with l2 penalty

DT: Decision Trees

RF: Random Forest

XGB: Gradient Boosting Trees

# Recession Prediction Results (1975-2020)



## References

---

- Ang, A. and Timmerman, A. (2011) *Regime Changes and Financial Markets*
- Mulvey, J.M. and Liu, H. (2016) *Identifying Economic Regimes: Reducing Downside Risks for University Endowments and Foundations*. Journal of Portfolio Management
- Hastie, T., Tibshirani, R. and Friedman, J. (2009) *The Elements of Statistical Learning*
- Varian, H. (2014) *Big Data: New Tricks for Econometrics*. Journal of Economic Perspectives
- Mullainathan, S. and Spiess, J. (2017) *Machine Learning: An Applied Econometric Approach*. Journal of Economic Perspectives
- Athey, S. (2018) *Impact of Machine Learning in Economics* NBER
- Lopez De Prado, M. (2018)<https://arxiv.org/abs/1903.10075>) *Advances in Financial Machine Learning*
- Piger, J. (2019) *Macroeconomic Forecasting in the Era of Big Data – Chapter 18*

## References

---

- Burns, A. F., and W. C. Mitchell. (1946) *Measuring Business Cycles*. NBER
- Stock, J. H., and M. W. Watson. (1989). *New Indexes of Coincident and Leading Economic Indicators*. NBER Macroeconomics Annual 4: 351–394.
- Chauvet, M., and J. Piger. (2008) *A Comparison of the Real-Time Performance of Business Cycle Dating Methods*. Journal of Business and Economic Statistics 26: 42–49.
- Ng, S. (2013) *Boosting Recessions*. Canadian Journal of Economics
- Davig, T. and Hall, A.S. (2019) *Recession Forecasting using Bayesian Classification*. International Journal of Forecasting
- James, A., Abu-Mostafa, Y.S. and Qiao, X. (2019) *Machine Learning for Recession Prediction and Dynamic Asset Allocation*. The Journal of Financial Data Science
- Molnar, C. (2019) *Interpretable Machine Learning*
- McCracken,M.W. and Ng, S. (2015) *FRED-MD: A Monthly Database for Macroeconomic Research*. Working Paper