
Hidden Market Regimes in Equity Returns via Hidden Markov Models

Charlie Kushelevsky, Parth Paliwal, Max Zhang, Cole Carter

1 **1 Milestone 1: Project Plan**

2 **1.1 Problem Description**

3 Financial markets appear noisy from day to day, yet empirical evidence suggests that price movements
4 are driven by a small number of persistent latent “regimes,” such as low-volatility growth periods,
5 high-volatility downturns, or transitional sideways phases (e.g. bull trends, bear markets, high
6 volatility). These regimes are not directly observable, but they affect return distributions, volatility
7 clustering, and risk.

8 The goal of this project is to use a Hidden Markov Model (HMM) to uncover latent market regimes
9 from historical equity return data. We treat the daily market return as the observed variable and the
10 underlying regime as a hidden state that evolves slowly over time. Using probabilistic reasoning,
11 expectation maximization (Baum-Welch), and sequence inference (Viterbi), we want to characterize
12 how many regimes best describe the behavior of the market and how persistent these regimes are.
13 This problem uses the course topics of latent variable modeling, EM-based learning, and probabilistic
14 inference.

15 We seek to answer questions such as: (1) How many distinct regimes best capture the distributional
16 structure of daily market returns? (2) What statistical properties (mean, variance) define each regime?
17 (3) How persistent are different regimes, and how frequently do transitions occur? (4) Do inferred
18 high-volatility states align with known market stress periods, such as the 2008 crisis or the 2020
19 COVID crash?

20 **1.2 Dataset**

21 We will use daily adjusted closing prices of the SPDR S&P 500 ETF (ticker: SPY), obtained from
22 publicly available historical data through the yfinance Python API. The dataset will span 1958
23 through 2025, covering multiple market cycles including the dot com aftermath, the 2008 global
24 financial crisis, the long 2010-2019 expansion, the COVID crash, and the 2022-2023 inflationary
25 period.

26 Preprocessing consists of: (1) restricting the dataset to valid trading days; (2) computing daily log
27 returns,

$$r_t = \log\left(\frac{P_t}{P_{t-1}}\right),$$

28 where P_t is the adjusted closing price; (3) dropping missing values from holidays and non-trading
29 days; and (4) optionally standardizing returns for numerical stability. No additional feature engineer-
30 ing is required for the baseline HMM.

31 This dataset is well suited for sequence models because it provides a long, continuous, high quality
32 return time series, enabling reliable estimation of persistent hidden dynamics.

33 **1.3 Methodology**

34 We will model the return sequence $\{r_t\}$ using a K -state Hidden Markov Model with Gaussian
35 emissions. Each hidden state $S_t \in \{1, \dots, K\}$ represents an unobserved market regime, while the

36 emission distribution $r_t \mid S_t = k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ captures the characteristic return behavior of that
37 regime. The transition matrix encodes the persistence and switching behavior between regimes,
38 which we expect to be highly skewed toward self-transitions.

39 Model parameters (initial state distribution, transition probabilities, and emission parameters) will be
40 learned with the Baum-Welch algorithm (EM). We will experiment with different values of K (e.g.,
41 $K = 2, 3, 4$) to compare fit, interpretability, and stability. After learning, we will apply the Viterbi
42 algorithm to infer the most likely regime sequence over time and examine whether high-volatility
43 regimes correspond to known historical events.

44 Evaluation will include log-likelihood, qualitative inspection of decoded regimes, and analysis of
45 transition probabilities and expected regime durations. If time permits, we may extend the model to
46 incorporate additional observable variables such as realized volatility or trading volume.