

Time Series Forecasting With a Learning Algorithm: An Approximate Dynamic Programming Approach

Ricardo Collado¹ Germán Creamer¹

¹School of Business
Stevens Institute of Technology
Hoboken, New Jersey



Time Series Drivers:

- Historical data:
 - Incorporated in classical time series forecast methods
 - Works best when the underlying model is fix
- Exogenous processes:
 - Not included in historical observations
 - Difficult to incorporate via classical methods
 - Could indicate changes in the underlying model

Techniques to handle changes due to external forces:

- Jump Diffusion Models
- Regime Switching Methods
- System of Weighted Experts
- Others . . .

These methods do not directly integrate alternative data sources available to us

Alternative data sources:

- **Text & News Analysis**
- Social Networks Data
- Sentiment Analysis



automated data mining survey
responses com... ter transcripts
qualatativ... root cause
classificati... insights
ad-hoc and... is product
reviews ser... it vor... of the
customer dashboards consumer
trends ad-hoc analysis early warning

Alternative data sources:

- Text & News Analysis
- **Social Networks Data**
- Sentiment Analysis



We study time series forecast methods that are:

- Dynamic
- Context-Based
- Capable of Integrating Social, Text, and Sentiment Data

In this presentation we develop:

- Stochastic dynamic programming model for time series forecast
 - Rely on an “external forecast” for future values
 - External forecast allows to incorporate alternative data sources

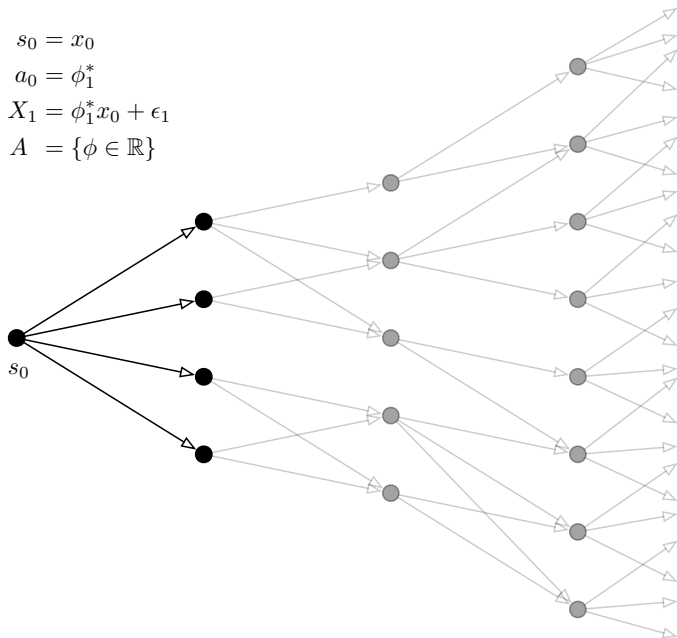
Time Series Fitting Process

$$s_0 = x_0$$

$$a_0 = \phi_1^*$$

$$X_1 = \phi_1^* x_0 + \epsilon_1$$

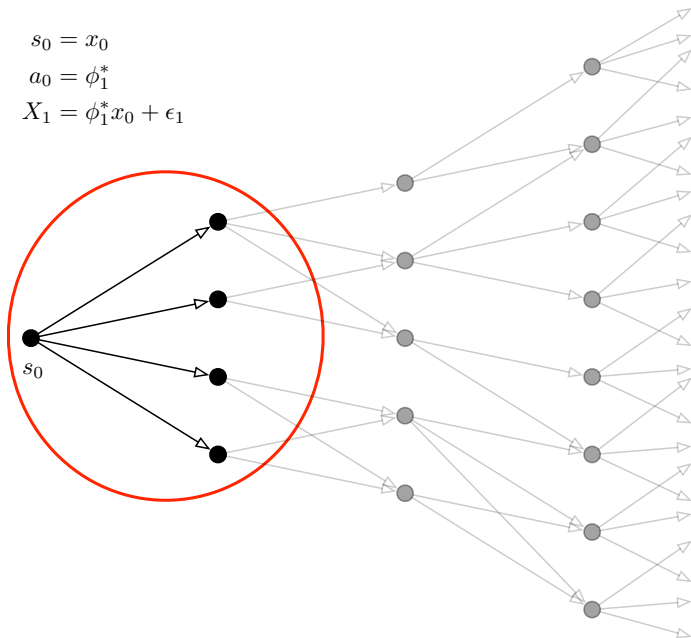
$$A = \{\phi \in \mathbb{R}\}$$



$$s_0 = x_0$$

$$a_0 = \phi_1^*$$

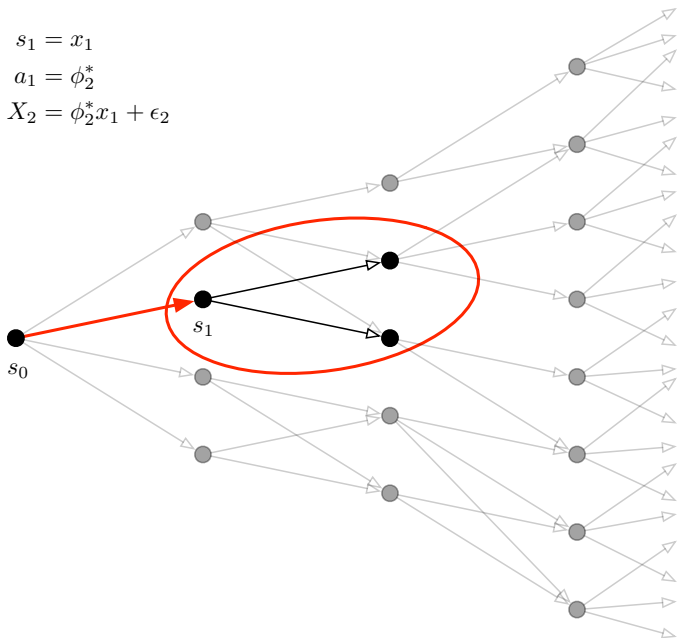
$$X_1 = \phi_1^* x_0 + \epsilon_1$$



$$s_1 = x_1$$

$$a_1 = \phi_2^*$$

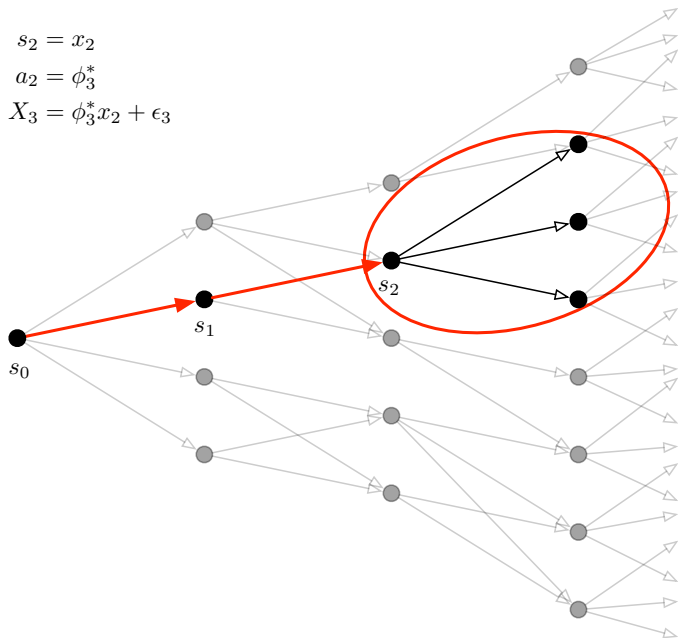
$$X_2 = \phi_2^* x_1 + \epsilon_2$$



$$s_2 = x_2$$

$$a_2 = \phi_3^*$$

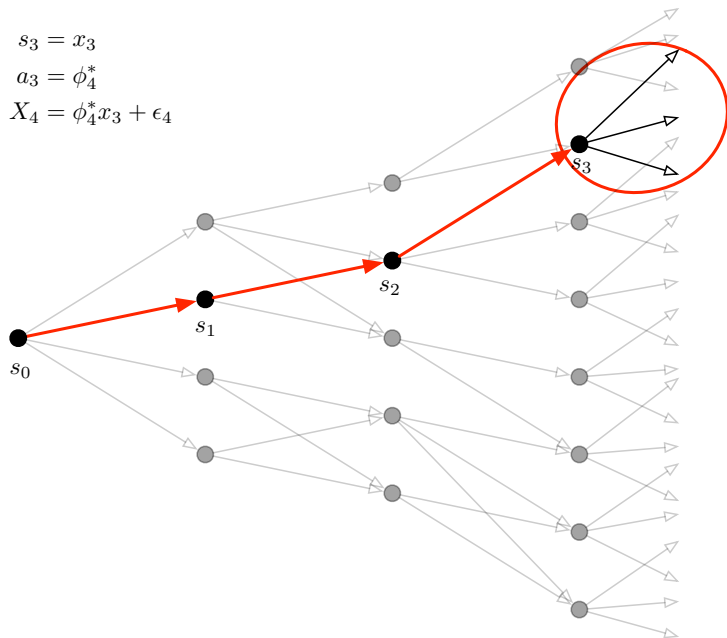
$$X_3 = \phi_3^* x_2 + \epsilon_3$$



$$s_3 = x_3$$

$$a_3 = \phi_4^*$$

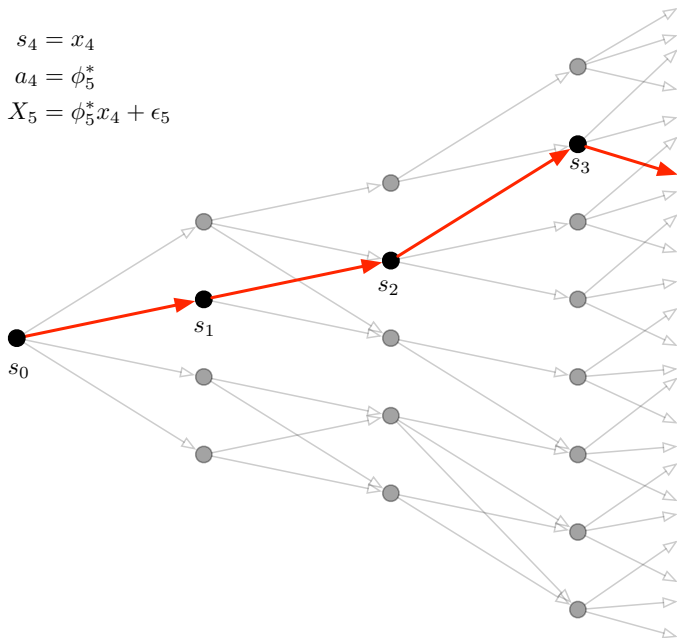
$$X_4 = \phi_4^* x_3 + \epsilon_4$$



$$s_4 = x_4$$

$$a_4 = \phi_5^*$$

$$X_5 = \phi_5^* x_4 + \epsilon_5$$



Main Problem:

$$\min_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=1}^T c_t(s_t, a_t) \right],$$

where $a_t = \pi_t(x_1, \dots, x_t)$ is an admissible fitting policy.

- The time series model is parametrized by $\Theta \subseteq \mathbb{R}^d$
- $A(s) = \Theta$ for all states s
- $c_t(s, a)$ is the result of a goodness of fit test for the observations $s = (x_0, \dots, x_t)$ and model selection $a = \theta$
- Solution via Bellman's optimality equations

Conditions to guarantee optimality:

- The set of actions $A(s)$ is compact
- The cost functions $c_t(s, \cdot)$ are lower semicontinuous
- For every measurable selection $a_t(\cdot) \in A_t(\cdot)$, the functions $s \mapsto c_t(s, a_t(s))$ and $c_T(\cdot)$ are elements of $\mathcal{L}_1(\mathcal{S}, \mathcal{B}_S, P_0)$
- The DP stochastic kernel function $Q_t(s, \cdot)$ is continuous

Value Functions: Value functions $v_t : \mathcal{S} \rightarrow \mathbb{R}$, $t = 1, \dots, T$, given recursively by:

$$v_T(s) = c_T(s)$$

$$v_t(s) = \min_{a \in A_t(s)} \{c_t(s, a) + \mathbb{E}[v_{t+1} | s, a]\},$$

for all $s \in \mathcal{S}$ and $t = T - 1, \dots, 0$.

Bellman's Optimality Equations:

Then an optimal Markov policy $\pi^* = \{\pi_0^*, \dots, \pi_{T-1}^*\}$ exists and satisfies the equations:

$$\pi_t^*(s) \in \arg \min_{a \in A_t(s)} \{c_t(s, a) + \mathbb{E}[v_{t+1} | s, a]\}, \quad s \in \mathcal{S}, \quad t = T - 1, \dots, 0.$$

Conversely, any measurable solution of these is an optimal Markov policy π^* .

Notice that:

- Our choice of model does not affect future observations and cost.
- So, $\mathbb{E}[v | s, a] = \mathbb{E}[v | s, a']$, for any $(s, a), (s, a') \in \text{graph}(\mathcal{A})$.
- Therefore we can rewrite the optimal policy as:

$$\pi_t^*(s) \in \arg \min_{a \in A_t(s)} \{c_t(s, a)\}, \quad s \in \mathcal{S}, \quad t = T - 1, \dots, 0.$$

Notice that:

- Our choice of model does not affect future observations and cost.
- So, $\mathbb{E}[v | s, a] = \mathbb{E}[v | s, a']$, for any $(s, a), (s, a') \in \text{graph}(\mathcal{A})$.
- Therefore we can rewrite the optimal policy as:

$$\pi_t^*(s) \in \arg \min_{a \in A_t(s)} \{c_t(s, a)\}, \quad s \in \mathcal{S}, \quad t = T - 1, \dots, 0.$$

- **The optimal policy π^* is purely myopic**

Q: How to break with the myopic policy?

A: Introduce a new Markov model

New Markov Model:

- Given stochastic process $\{X_t \mid t = 0, \dots, T\}$, s.t. $X_0 = \{\phi_0\}$
- Time series model parameterized by $\Theta \subseteq \mathbb{R}^d$
- State space:

$$\mathcal{S}_t = \left\{ (x_t, h_{t-1}, \theta_{t-1}) \left| \begin{array}{l} x_t \text{ observation from } X_t, \\ h_t = x_0, \dots, x_{t-1} \text{ sample sequence,} \\ \theta_{t-1} = (\phi_1, \dots, \phi_p) \in \Theta \end{array} \right. \right\}$$

- Action space: $A(s) = \Theta$ for all states $s \in \mathcal{S}_t$

Cost function:

$$c_t(s, \theta_t) = \gamma(s_t, \theta_t) + r \delta(s_t, \theta_{t-1}, \theta_t)$$

- γ : Goodness of fit test
- δ : Penalty on changes from previous model selection
- $r \geq 0$: Scaling factor used to balance fit and penalty

Example:

$$\chi^2(\theta_t | h_{t-1}, x_t) + r \left(1 - \exp \left\{ -\lambda \left| \mathbb{E} [P_{\theta_t} | h_{t-1}, x_t] - \mathbb{E} [P_{\theta_{t-1}} | h_{t-1}, x_t] \right| \right\} \right),$$

where $r, \lambda \geq 0$.

- $\gamma(s_t, \theta_t) := \chi^2(\theta_t | h_{t-1}, x_t)$
- $\delta(s_t, \theta_{t-1}, \theta_t) := 1 - \exp \left\{ -\lambda \left| \mathbb{E} [P_{\theta_t} | h_{t-1}, x_t] - \mathbb{E} [P_{\theta_{t-1}} | h_{t-1}, x_t] \right| \right\}$

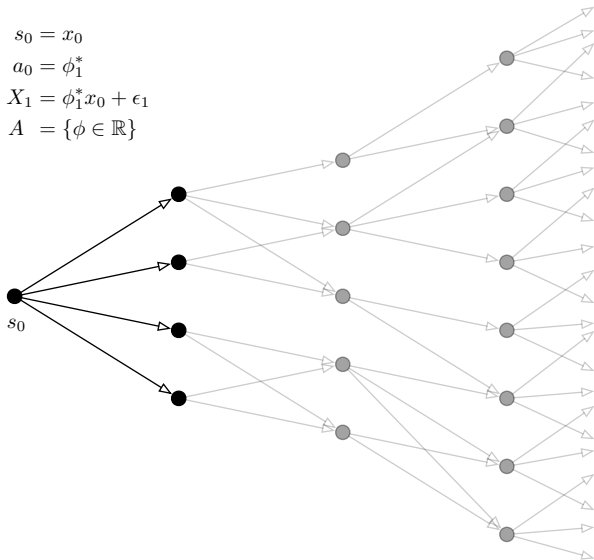
Dynamic Optimization Model:

$$s_0 = x_0$$

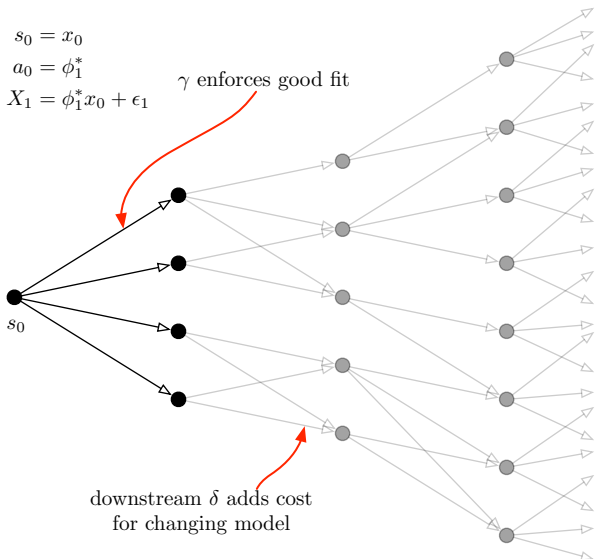
$$a_0 = \phi_1^*$$

$$X_1 = \phi_1^* x_0 + \epsilon_1$$

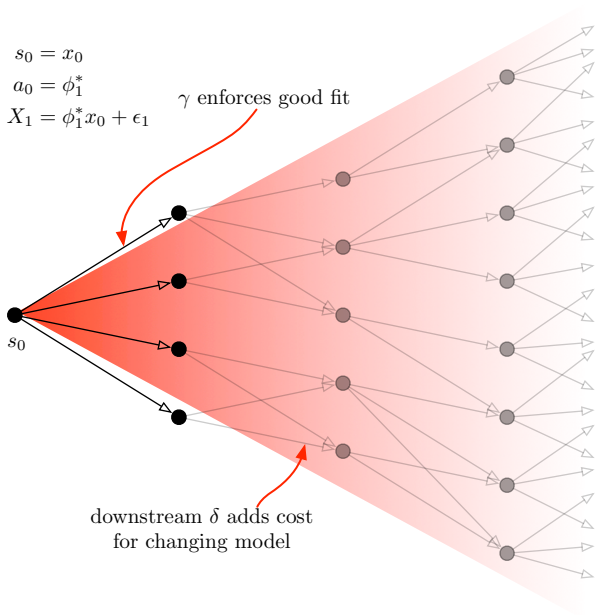
$$A = \{\phi \in \mathbb{R}\}$$



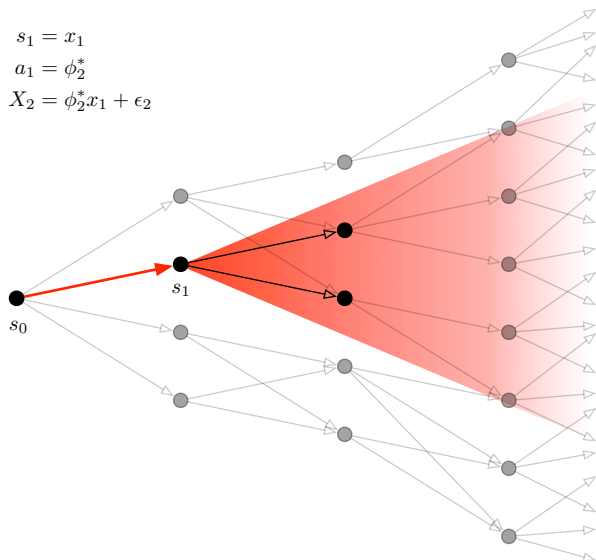
Dynamic Optimization Model:



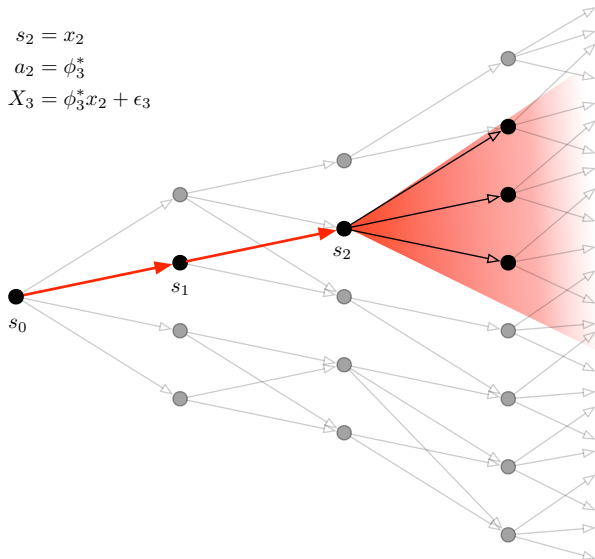
Dynamic Optimization Model:



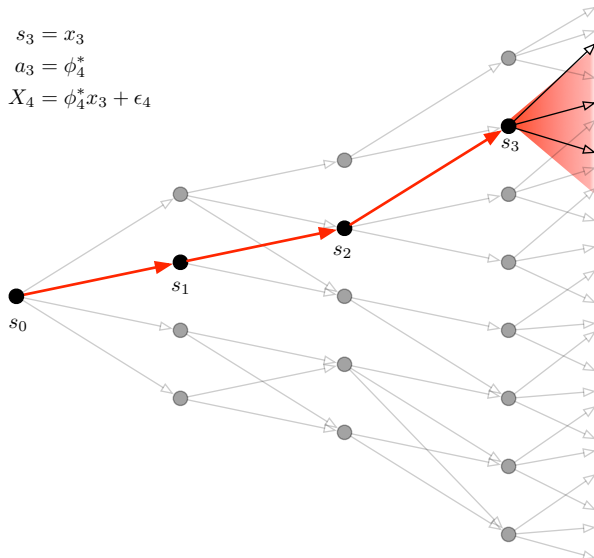
Dynamic Optimization Model:



Dynamic Optimization Model:



Dynamic Optimization Model:



Introducing Exogenous Information Via Lookahead Methods

Lookahead Methods:

- At time t we have access to forecast random variables:

$$\widehat{X}_{t+1}^t, \dots, \widehat{X}_{t+h}^t$$

- We assume these are discrete approximations to X_{t+1}, \dots, X_{t+h}
- Each forecast induces a discrete probability measure on \mathcal{S}

We expect the forecasts to be finite random variables with few atoms

This has the effect of simplifying calculations, but at the expense of rough approximations.

Basic lookahead method:

Step 0 Initialization:

Step 0a Initialize $\bar{v}_t(s)$ for all time periods t and all $s \in \mathcal{S}$.

Step 0b Choose an initial state s_0^1 .

Step 1 For $t = 0, \dots, T$ do:

Step 1a Update the state variable: observe a value x_t of the stochastic process, let s_t be obtained from x_t and model selection at $t - 1$, and get forecasts $\hat{X}_{t+1}, \dots, \hat{X}_{t+h}$.

Step 1b Solve

$$\hat{v}_t = \min_{a_t \in A_t(s_t)} \{c_t(s_t, a_t) + \mathbb{E}[\bar{v}_{t+1} | s_t, a_t]\},$$

by solving a stochastic optimization problem. Let a_t be the obtained optimal solution to the minimization problem.

Step 1c Update the value function approximation \bar{v}_t :

$$\bar{v}_t(s) = \begin{cases} \hat{v}_t, & \text{if } s = s_t, \\ \bar{v}_t(s), & \text{otherwise.} \end{cases}$$

Step 2 Return the value functions $\{\bar{v}_t | t = 0, 1, \dots, T\}$.

Basic Monte Carlo Algorithm:

Step 0 Initialization:**Step 0a** Initialize $\bar{v}_t^0(s)$ for all time periods t and all $s \in \mathcal{S}$.**Step 0b** Choose an initial state s_0^1 .**Step 0c** Solve the problem the previous method and denote the resulting value function approximations \bar{v}_t^0 , $t = 0, \dots, T$.**Step 0d** Set $n = 1$.**Step 1** For $t = 0, \dots, T$ do:**Step 1a** Update the state variable: Observe a value x_t^n of the stochastic process, let s_t^n be obtained from x_t^n and model selection at $t - 1$, and get forecasts $\hat{X}_{t+1}^n, \dots, \hat{X}_{t+h}^n$.**Step 1b** Solve

$$\hat{v}_t^n = \min_{a_t \in A_t(s_t^n)} \{c_t(s_t^n, a_t) + \mathbb{E}[\bar{v}_{t+1} | s_t^n, a_t]\},$$

as before. Let a_t^n be the obtained optimal solution to the minimization problem.**Step 1c** Update the value function approximation \bar{v}_t^{n-1} :

$$\bar{v}_t^n(s) = \begin{cases} (1 - \alpha_{n-1})\bar{v}_t^{n-1}(s) + \alpha_{n-1}\hat{v}_t^n, & \text{if } s = s_t^n, \\ \bar{v}_t^{n-1}(s), & \text{otherwise.} \end{cases}$$

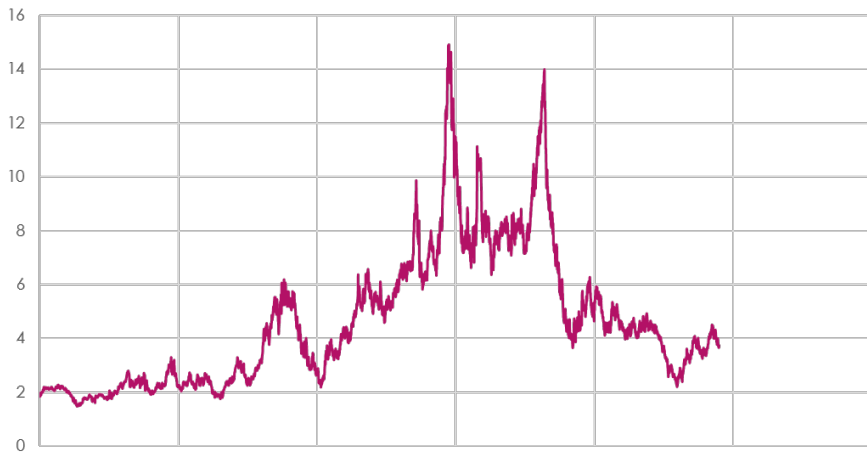
Step 2 Let $n = n + 1$. If $n < N$, go to **Step 1**.**Step 3** Return the value functions $\{\bar{v}_t^N \mid t = 0, 1, \dots, T\}$.

Numerical Results I

Estimation Techniques I:

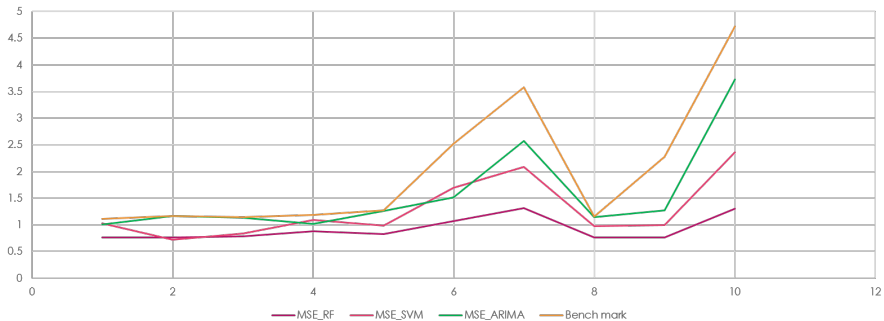
- Initial test: Natural Gas Futures Contract 4 - times series (RNGC4), 1999 – 2013.
- Test approximate dynamic programming method (MSE₋) over 1 to 10 days forecast window, comparing the calibration of machine learning algorithms used for regression and classification:
 - Support vector machine (SVM)
 - Random forests (RF)
 - ARIMA
 - Logistic regression (LR) – Benchmark
- 2nd test: Crude Oil Spot (WTI), Aug. 2015 – Jan. 2016.
- We made two main tests by using approximations to “future” WTI Spot on different time windows:
 - Actual values + random white noise with increasing variance
 - Actual values + white noise + increasing bias

RNGC4

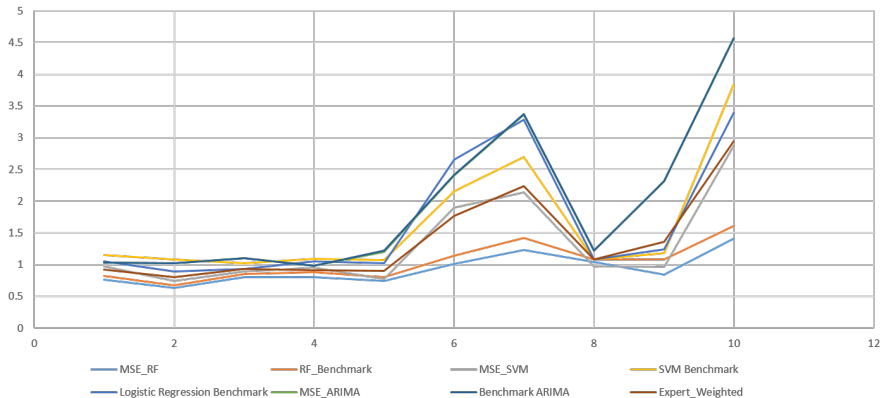


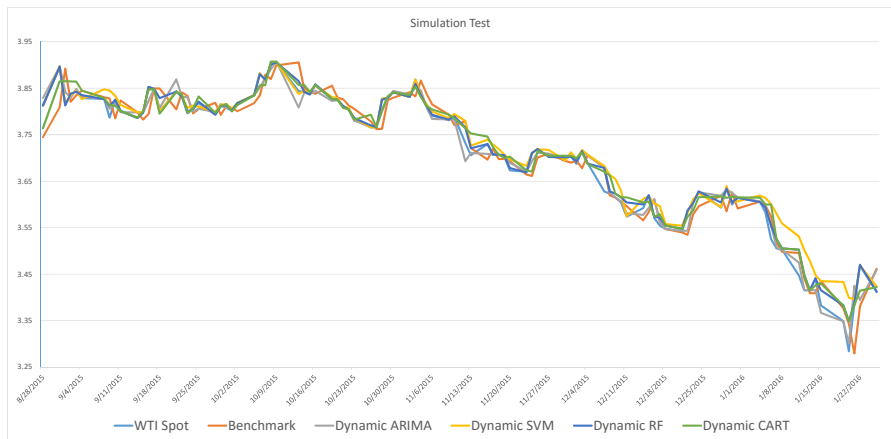
Energy Time Series

Performance of Three Algorithms vs Benchmark



Learners vs LR Benchmark





Model	Model RMSE					
	Noise Variance					
	0.001	0.01	0.05	0.1	0.5	1
ARIMA	0.0003227	0.000389	0.001943	0.004214	0.008405	0.007514
SVM	0.000776849	0.000818	0.002475	0.008127	0.099774	0.16513
RF	0.000175392	0.00023	0.001843	0.003496	0.007223	0.008425
CART	0.000483638	0.000504	0.002376	0.006527	0.023588	0.024053

Benchmark	0.001116097
-----------	-------------

Unbiased Random Noise External Forecast Test

ARIMA Baseline	490.6124
RF Baseline	502.5576
CART Baseline	967.2576

	0.001	0.01	0.03	0.05	0.1	0.5	1	2
ARIMA	384.0882	380.7184		379.9081	373.224	380.9279	386.4873	473.5217
RF	28.16006	30.8183	36.77366	45.00305	56.97973	195.9457	282.2603	389.0374
CART	6.680402	29.53345	46.15347	64.68325	147.4693	497.2315	725.8215	863.787

ARIMA Baseline	627.327
RF Baseline	804.8075
CART Baseline	1588.515

	0.001	0.01	0.03	0.05	0.1	0.5	1	2
ARIMA	586.0408	586.1985		586.4885	587.7807	590.063	624.0278	705.4027
RF	71.76113	87.90271	117.1728	146.3925	199.0221	424.4266	528.3886	618.359
CART	12.98331	65.70607	135.9031	240.7952	369.7038	913.2118	1188.266	1422.569

Biased Random Noise External Simulation: 10 and 50 Days Windows

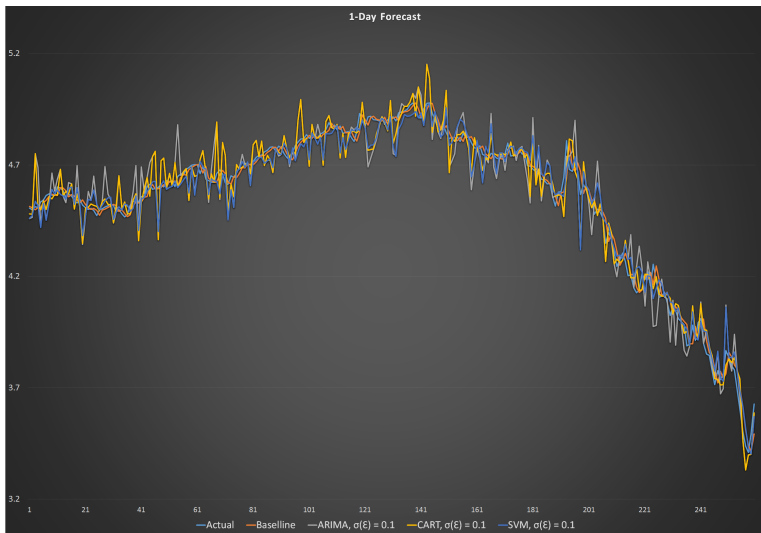
Numerical Results 2

Estimation Techniques II:

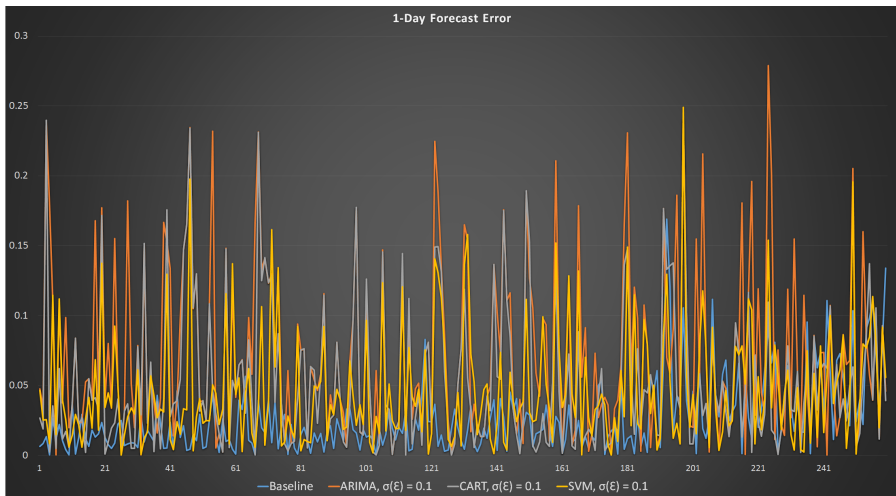
- 3rd test: Crude Oil Spot (WTI). Test approximate dynamic programming method over 1, 5, and 5 days forecast window, using different ML dynamic techniques with actual + white noise as “future” data:
 - Support vector machine (SVM)
 - Classification and regression trees (CART)
 - ARIMA
 - Logistic regression (LR) – Benchmark
- 4th test:SVM futures WTI and SVM sparse “future” data on WTI spot.



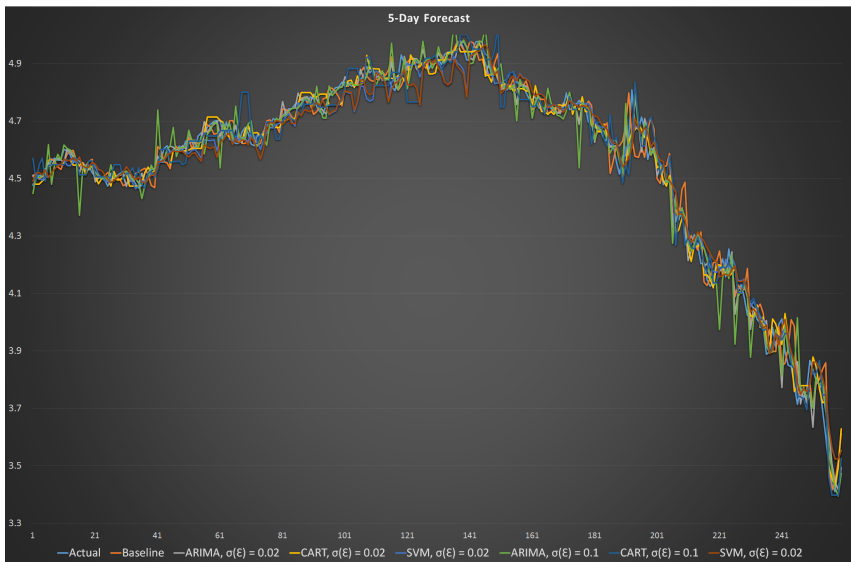
WTI Spot



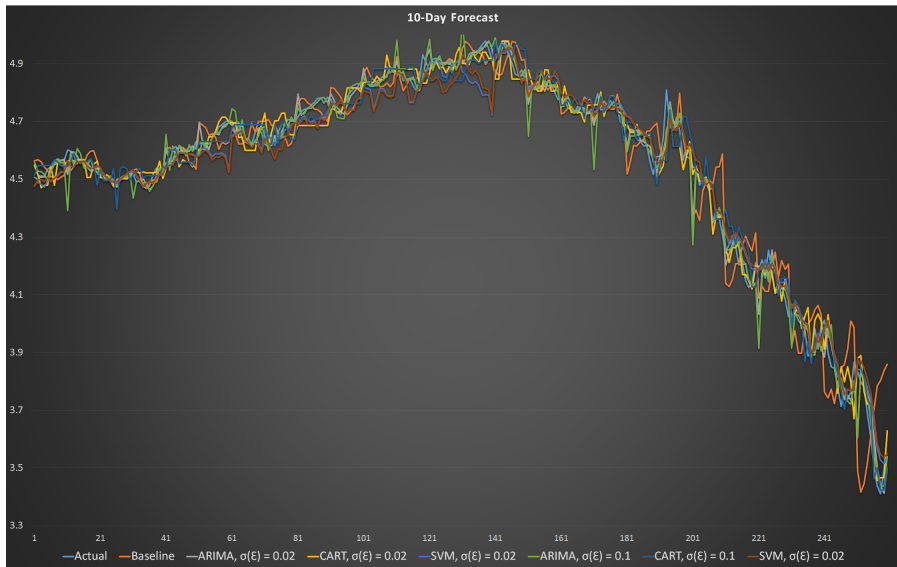
1-Day Forecast



1-Day Forecast Error



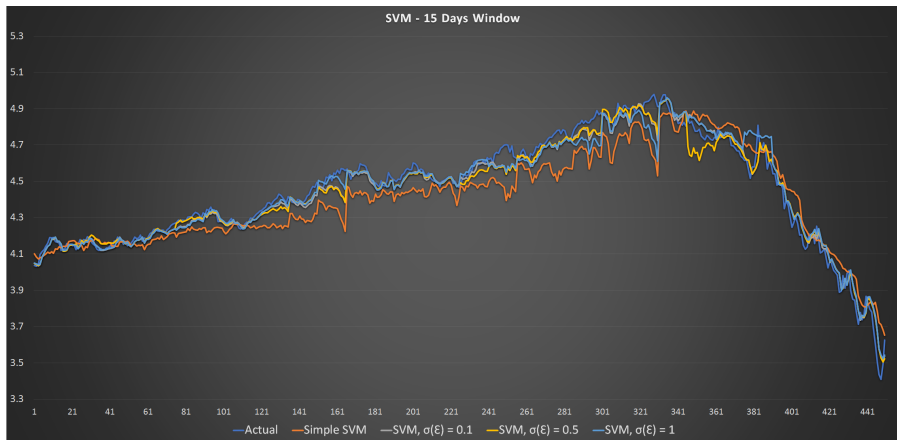
5-Day Forecast



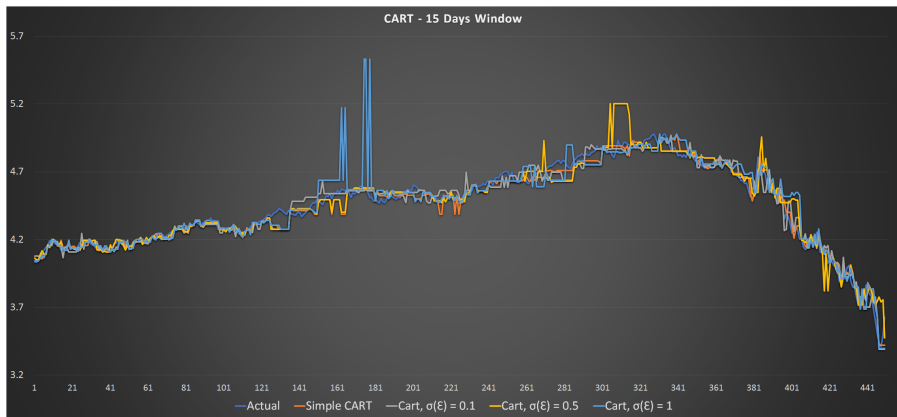
10-Day Forecast

Numerical Results 2.2

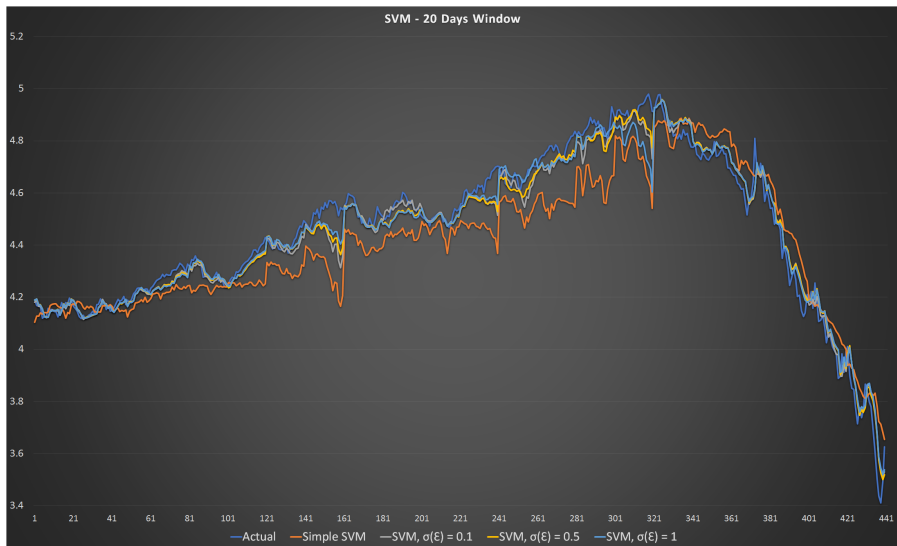
SVM – CART Analysis



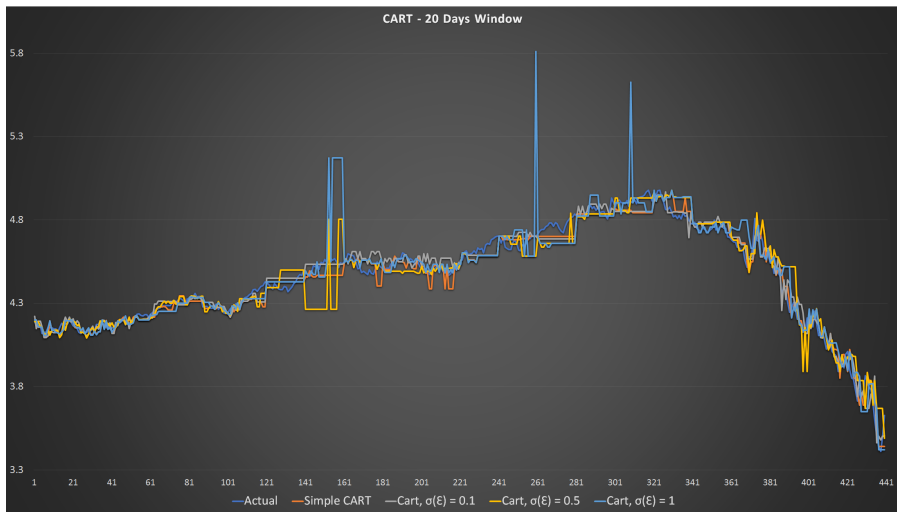
SVM 15 Days Window



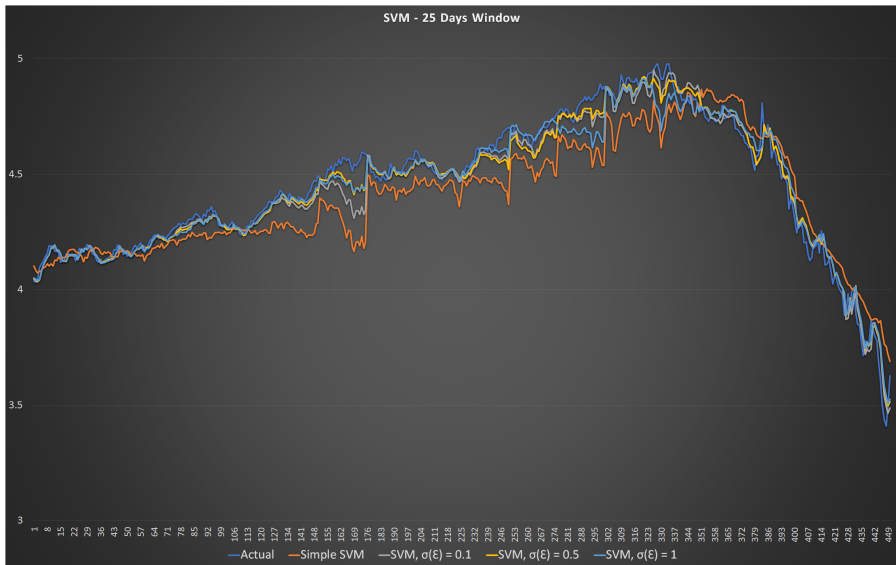
CART 15 Days Window



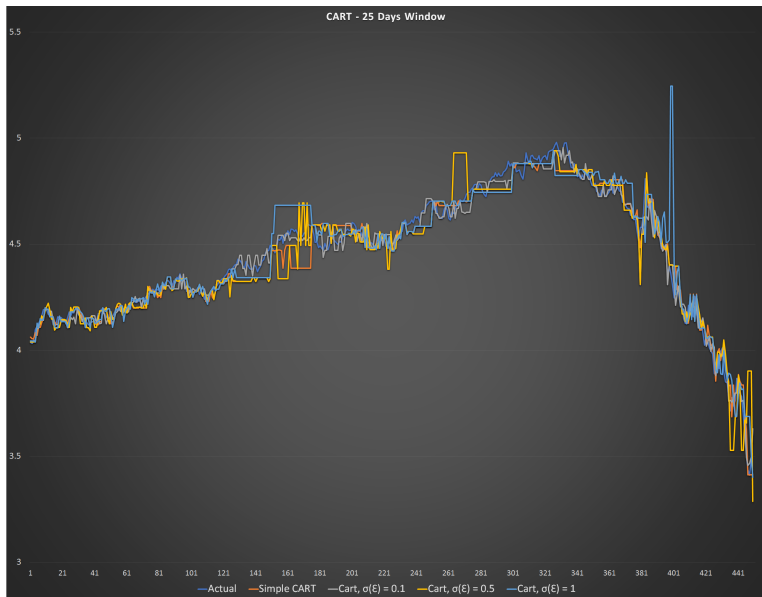
SVM 20 Days Window



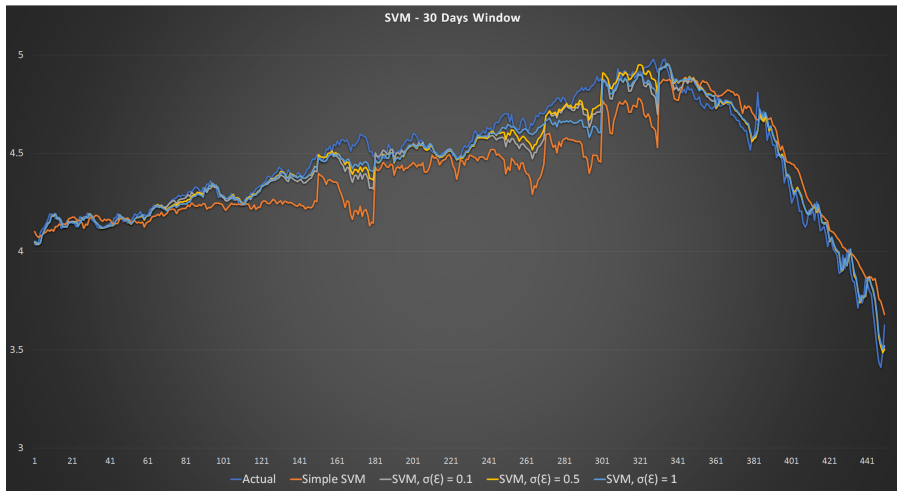
CART 20 Days Window



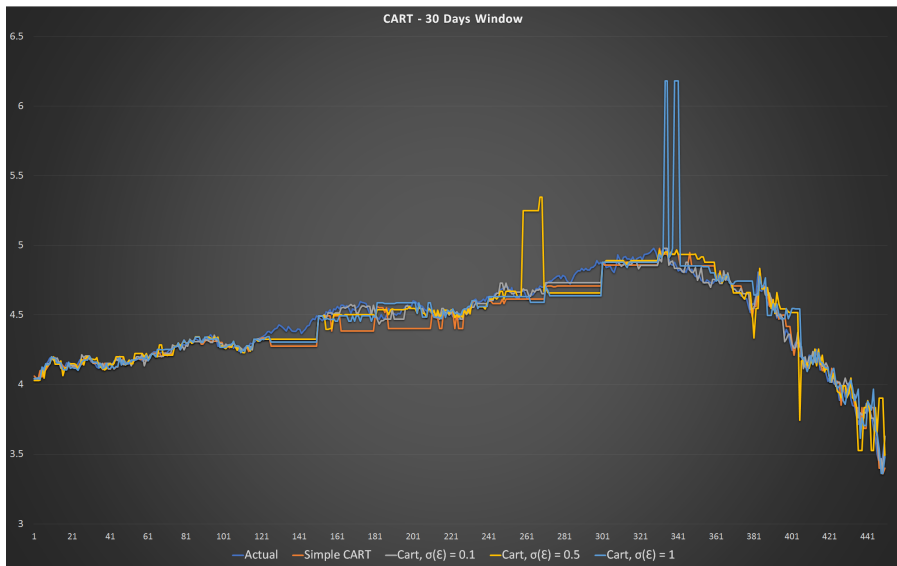
SVM 25 Days Window



CART 25 Days Window



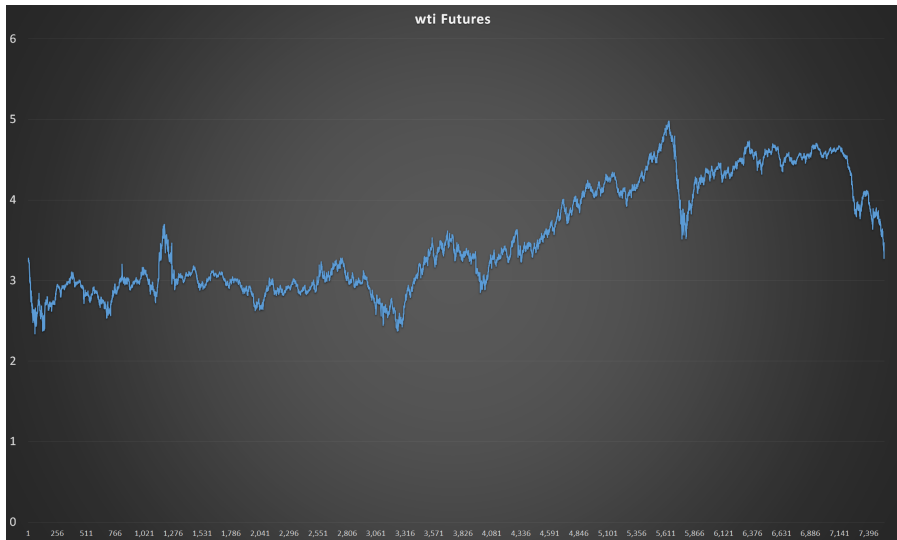
SVM 30 Days Window



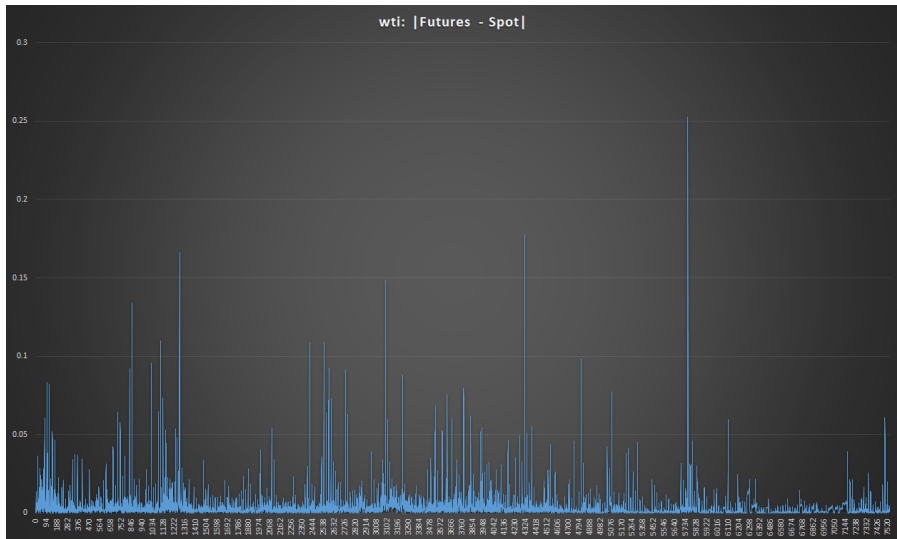
CART 30 Days Window

Numerical Results 2.3

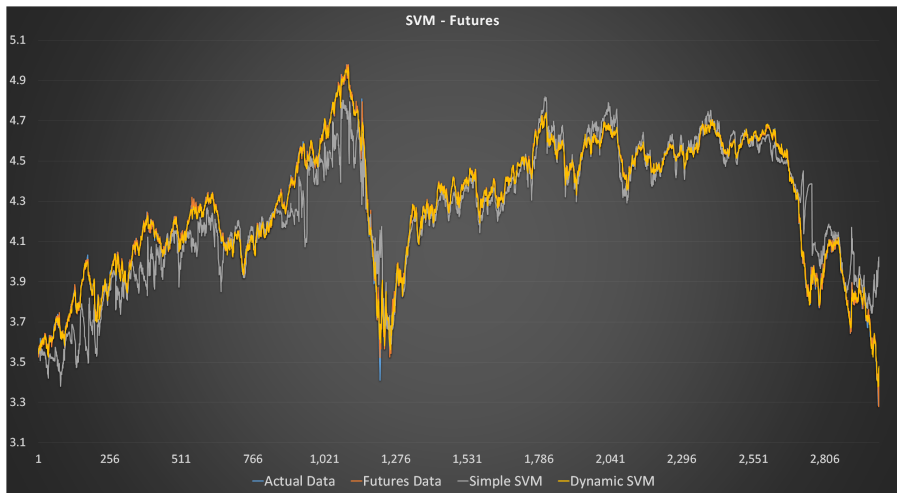
WTI Futures – Sparse SVM Analysis



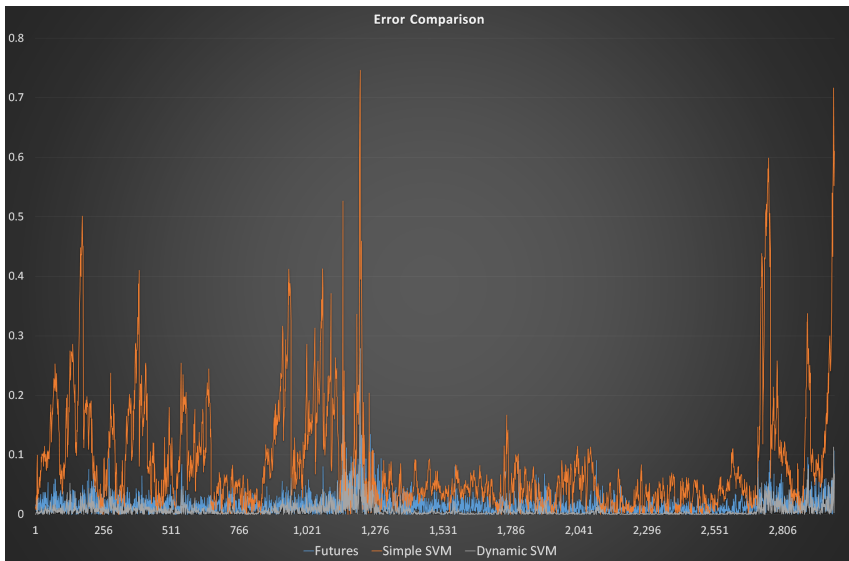
WTI Futures



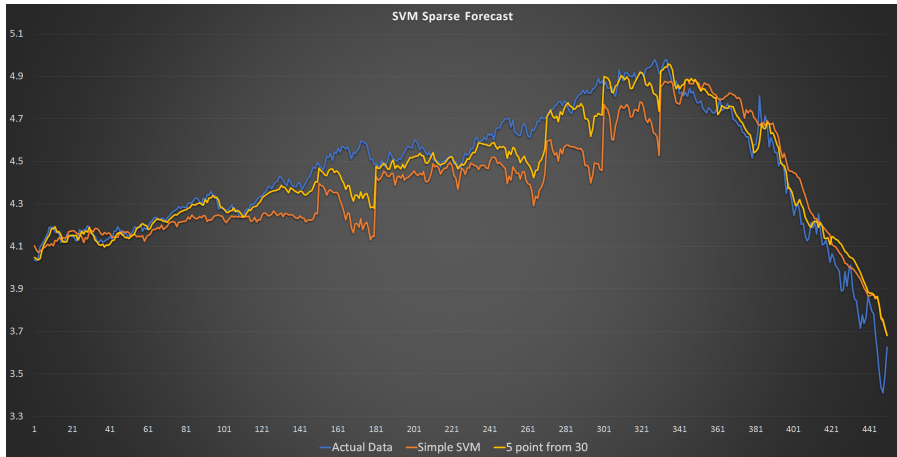
WTI: |Futures – Spot|



SVM Futures



SVM Futures Error Comparison

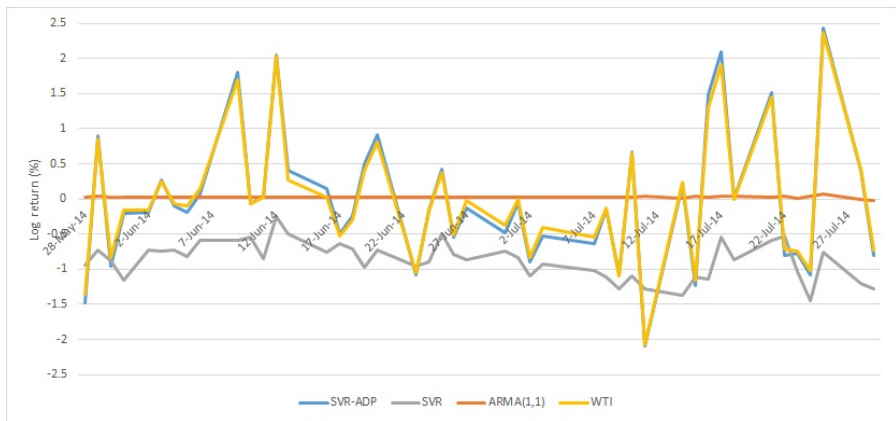


SVM Sparse Forecast

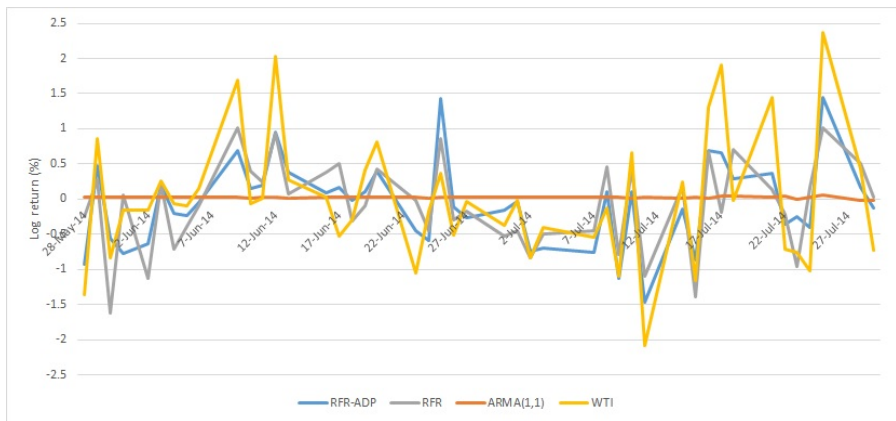
Numerical Results 3

Estimation Techniques III:

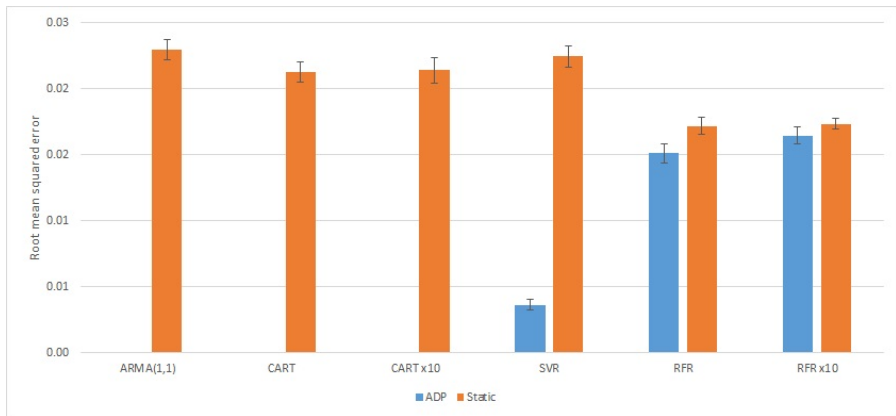
- Test our approximate programming dynamic method comparing the calibration of two machine learning algorithms used for regression and classification: support vector machine (SVM) and random forests (RF).
- For regression analysis:
 - Support vector regression (SVR): minimizes a loss function using only the most relevant values.
 - Random forests for regression (RFR): explores a large search space based on random selection of its features and samples.
- External forecast function for daily data of WTI log returns: 1 month WTI futures price at the end of every month.
- Training & test dataset: 15 lags of the dependent variable and technical indicators:
 - Price indicators: Simple moving averages with 10 and 20 days.
 - Momentum indicators: Relative strength index with 10 and 20 days, and the moving average convergence divergence.
- The test sample includes the forecast function.



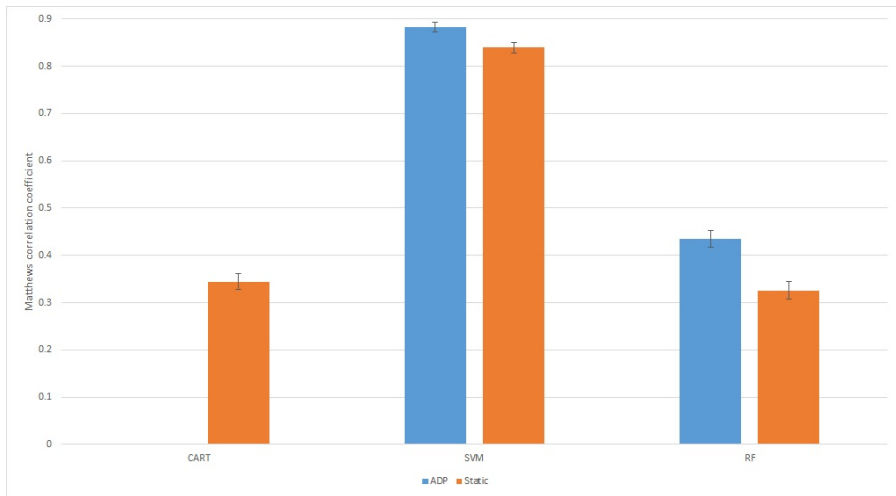
WTI log return and forecasts by ARMA(1,1), support vector regression (SVR), and the approximate dynamic programming SVR (SVR-ADP) for June and July 2014.



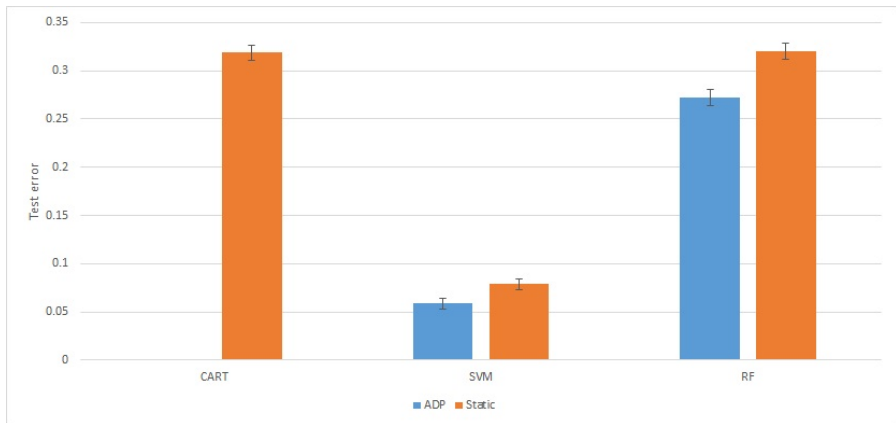
WTI log return and forecasts by ARMA(1,1), random forest for regression (RFR), and the approximate dynamic programming RFR (RFR-ADP) for June and July 2014.



Root mean squared error (RMSE) of ADP and static (benchmark) methods. SVR, RFR, and x10 stand for support vector regression, random forest for regression, and 10 times 40 folds cross-validation respectively. The error bars represent standard error. RMSE mean differences of SVR-ADP and RFR-ADP in relation to their static versions and ARMA(1,1) and CART are statistically significant with 99% confidence level .



Matthews correlation coefficient of ADP and static (benchmark) methods. SVM and RF stand for support vector machine and random forest respectively. The error bars represent standard error. MCC mean differences of SVM-ADP and RF-ADP in relation to their static versions and CART are statistically significant with 95% conf.level.



Test error of ADP and static (benchmark) methods. SVM and RF stand for support vector machine and random forest respectively. The error bars represent standard error. Test error mean differences of SVM-ADP and RF-ADP in relation to their static versions and CART are statistically significant with 95% confidence level.