

Title *

Tianyu Du [†]

Sunday 2nd February, 2020

Contents

1	Introduction	3
2	Data	3
2.1	The West Texas Intermediate Crude Oil Dataset	3
2.2	Missing Data in Crude Oil Dataset	3
2.3	Day of the Week Effect in Crude Oil Dataset	5
2.3.1	Difference in Returns across the Week	5
2.3.2	Kolmogorov-Smirnov test for Distributional Similarities	7
2.4	News and Sentiment Datasets	8
2.5	Classifying News Type	9
2.6	Case Studies of Events	9
2.6.1	Positive Spike on November 30, 2016	9
2.6.2	Negative Spike on December 6, 2018	9
2.6.3	Positive Spike on June. 12 ~ 13, 2019	9
3	Models	9
4	Experiments	9

*Compile Date: 21:10 Sunday 2nd February, 2020

[†]tianyu.du@mail.utoronto.ca

References	9
5 Appendix	10

1 Introduction

2 Data

This study involves three major datasets: i) the daily spot price of crude oil at the West Texas Intermediate (WTI), ii) news sentiment dataset from Ravenpack News Analytics (RPNA), and iii) macroeconomic indicators proxying the overall economic state.

2.1 The West Texas Intermediate Crude Oil Dataset

West Texas Intermediate market spot price (coded as DCOILWTICO in the St. Louis Federal Reserve Economic Data) has been the most commonly used price for crude oil in current literature. The dataset retrieved from the Federal Reserve Bank of St. Louis is measured on a daily level and spans from January 1986 up to the present day (citation: data series). Because of the limited availability of the RavenPack dataset, this paper focuses only on crude oil prices after January 1, 2000. Analysis of the crude oil market (citation: 40-year paper) shows the spot price is highly responsive to news and other macroeconomic shocks, which is exactly the tricky part of forecasting financial time series. If the proposed forecasting pipeline performs well on the crude dataset, such a pipeline is conceivably promising on other datasets as well.

Almost all financial time series suffers from missing data problem, so is the crude oil dataset. Much existing research studying stock market datasets simply drop missing values. Instead, this paper uses an autoregressive integrated moving average (ARIMA) model to interpolate and fill missing data, so that the time gap between two consecutive observations is exactly one trading day.

2.2 Missing Data in Crude Oil Dataset

Given the focus of this paper is on forecasting crude oil returns, which captures the difference between two consecutive prices, The missing data problem can be crucial here.

This paper calculates crude oil returns on one particular day t by taking the difference in logged

prices at t and the previous trading day:

$$r_t := \ln(p_t) - \ln(p_{t-\Delta}) \quad (2.1)$$

where $t - \Delta$ is the last trading day before day t . As mentioned before, the time gap between two observed prices are not even. For instance, the return on a Monday can be computed by taking difference between the log close price on Monday Friday (if available). In this case, $\Delta = 3$. If the previous Friday was a holiday without valid price data, r_t will be $\ln(p_{\text{Mon}}) - \ln(p_{\text{Prev Thu}})$, and $\Delta = 4$. According to [the table below](#), 33 days are in this case.

Day of the week	Num. Days.	Num. Trading Days	$\Delta=1$	2	3	4	5
Monday	1031	927	0	0	883	33	11
Tuesday	1030	1018	921	0	0	97	0
Wednesday	1030	1022	1011	5	0	0	6
Thursday	1030	1002	994	8	0	0	0
Friday	1030	986	969	17	0	0	0
Saturday	1030	0	0	0	0	0	0
Sunday	1030	0	0	0	0	0	0
Total	7211	4955	3895	30	883	130	17

Table 1: The values of Δ used to calculate returns. This table only include trading days, but the first day with price observation in this dataset was dropped because it did not have a previous trading day, so return on this day cannot be computed using our definition.

As mentioned before, the oil price dataset does not have any prices over weekends. [The table below](#) reports dates that are most frequently associated with a missing data over the span of 20 years. The pool of days with missing data is pretty consistent overtime, the market is always closed on January 1, July 4 (Independence Day) and December 25 (Christmas). The group of dates in late November are responsible for missing data on Thanksgiving holiday.

Date	Counts (all)	Counts (excl. weekends)
July 4	20	16
January 1	20	14
December 25	19	14
July 3	10	5
November 23	10	5
November 24	10	4
November 25	10	3
November 22	9	4
November 26	9	3

Table 2: Dates most frequently associates with missing data. Data on January 1, July 4, and December 25 are missing ever year. Because the entire dataset ranges from January 3, 2000 to September 30, 2019, missing data problems on December 25 are only reported 19 times.

2.3 Day of the Week Effect in Crude Oil Dataset

2.3.1 Difference in Returns across the Week

Gibbons and Hess' work examined returns on stocks from S&P 500, Dow Jones 30, and Treasury Bills. They found strong negative mean returns on Monday compared with other weekdays. The seasonality persisted even after market adjustment measures, such as using mean-adjusted returns instead, were taken (Gibbons & Hess, 1981). Analysis in this paper suggests a similar daily seasonality presents in crude oil returns as well. **Panels in the figure below** demonstrate the empirical distributions of returns on each day of the week. We can see that Monday and Wednesday have relatively larger variances, which again matches Gibbons and Hess' observations.

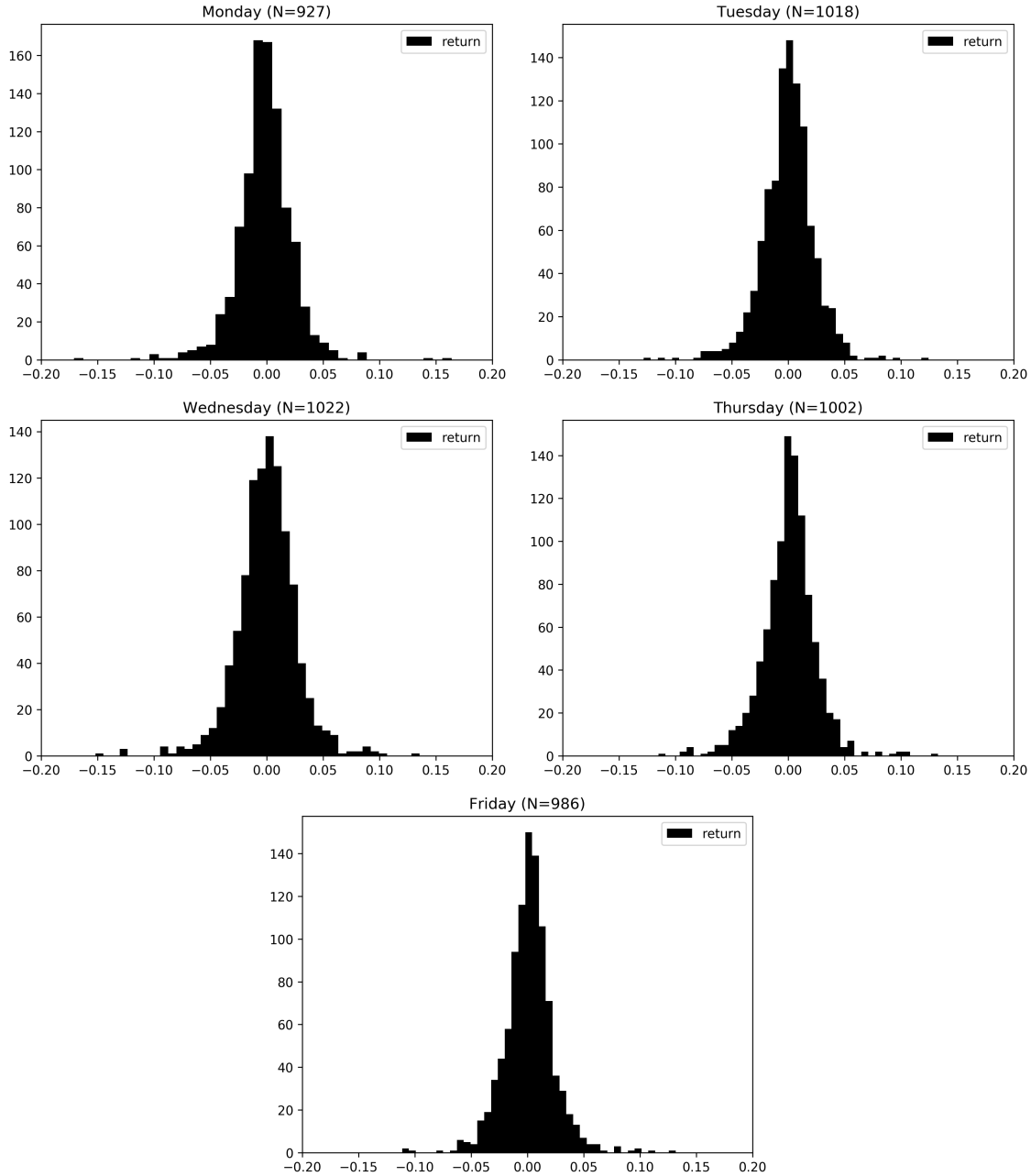


Figure 1: Crude oil returns on each weekday. Weekend data are not available in the daily dataset provided by EIA. The range of y-axis in all five histograms are from -0.2 to 0.2. N s in parentheses denote the number of observations. See appendix for distributions of crude oil prices.

The **two tables** below provide summary statistics for prices and returns on each day. It turns out that Monday is the only weekday with a mean return significantly less than zero.

Day of the week	Num. Obs.	Mean	Std.	3 rd Moment
Monday	927	62.072	26.493	7081.163
Tuesday	1019	61.828	26.317	6895.638
Wednesday	1022	61.810	26.398	7049.810
Thursday	1002	62.005	26.431	6955.555
Friday	986	62.079	26.247	6676.566
Total	4956			

Table 3: Summary statistics of crude oil prices on each day of week

Day of the week	Num. Obs.	Mean (P -Value)	Std.	3 rd Moment
Monday	927	-0.002 (0.049)	0.025	-0.0000019
Tuesday	1018	-0.000 (0.900)	0.023	-0.0000031
Wednesday	1022	0.000 (0.884)	0.027	-0.0000054
Thursday	1002	0.001 (0.361)	0.024	-0.0000006
Friday	986	0.002 (0.0311)	0.023	0.0000021
Total	4955			

Table 4: Summary statistics of crude oil returns on each day of week. The first day (January 1, 2000) of the oil price dataset was Saturday, and the observation on the following Monday (January 3) was missing. Hence, the return on Tuesday (January 4) could not be computed because it was the first trading day in this dataset, and there are only 1018 Tuesday in the dataset of returns. A value of -0.000 indicates a negative value with magnitude less than 0.0005. P -values are calculated in a two-tailed t -test with $\mu_0 = 0$. Bold fonts indicate statistically significance at level $\alpha = 0.05$.

2.3.2 Kolmogorov-Smirnov test for Distributional Similarities

Smirnov developed a non-parametric method of testing the equality between two continuous distributions, with CDFs $F(x)$ and $G(x)$ respectively, (Smirnov, 1939). Refer to Hodges' work for a detailed review on the Kolmogorov-Smirnov test (Hodges, 1958). I am using the two-tailed version of Kolmogorov-Smirnov test to check whether distributions of two different days are similar. Given two datasets, take returns on Monday and Tuesday for example, the null hypothesis says those two datasets are drawn from the same distribution, and the alternative says they are from different distributions¹. Firstly, the Kolmogorov-Smirnov test constructs the empirical CDFs $F_{Mon,927}(x)$ and $F_{Tue,1018}(x)$ from the dataset. Then, the Kolmogorov-Smirnov statistic measures the maximum

¹Different alternative hypotheses can be used in Kolmogorov-Smirnov test: i) $H_1 : F(x) \geq G(x)$, ii) $H_1 : F(x) \leq G(x)$, and iii) $H_1 : F(x) \neq G(x)$. This paper is using the third (two-tailed) alternative hypothesis.

discrepancy between two distribution functions, which is

$$D := \sup_x |F_{Mon,927}(x) - F_{Tue,1018}(x)| \in [0, 1] \quad (2.2)$$

A smaller D -statistic implies stronger distributional similarity between two distributions. For instance, when $F_{Mon,927}(x)$ and $F_{Tue,1018}(x)$ are exactly the same, the D -statistic is zero. In contrast, let $X = 0$ and $Y = 1$ be two deterministic random variables, in this case, $D_{X,Y} = 1$.

The test rejects H_0 at a significance level of α if

$$D > \sqrt{-\frac{1}{2} \ln \frac{\alpha}{2}} \sqrt{\frac{n+m}{nm}} \quad (2.3)$$

where m and n denote sizes of two datasets.

D -Statistic (P -Value)	Monday	Tuesday	Wednesday	Thursday	Friday
Monday	0.000 (1.000)	0.061 (0.048)	0.065 (0.030)	0.092 (0.001)	0.092 (0.001)
Tuesday		0.000 (1.000)	0.044 (0.260)	0.036 (0.505)	0.044 (0.264)
Wednesday			0.000 (1.000)	0.053 (0.114)	0.073 (0.009)
Thursday				0.000 (1.000)	0.025 (0.900)
Friday					0.000 (1.000)

Table 5: The Kolmogorov-Smirnov D -Statistic for all pairs of distributions. Bold font indicates the null hypothesis is rejected at a significance level of 0.05, which implies discrepancy in distributions.

The table above presents the Kolmogorov-Smirnov D -Statistic for distributions of every pairs of days. At a significance level of 0.05, we can see that Monday follows a distribution significantly different from distributions other days follow. Because the dataset does not contain weekend data, the return on Monday is always computed using the difference between log prices on Monday and the previous Friday (Thursday if Friday is not a trading day and so on). Therefore, returns associated with Mondays pick the weekend effect. In fact, the distribution of returns on Monday (over weekend) is the only one with negative mean among distributions of all five days.

2.4 News and Sentiment Datasets

The event sentiment dataset from RavenPack News Analytics (RPNA) tracks and analyzes all information of companies, organizations, countries, commodities, and currencies from Dow Jones

Newswires, Wall Street Journal, Barron's and MarketWatch ranges from January 1, 2000, to September 30, 2019. RavenPack records the exact date and time (measured using Coordinated Universal Time, UTC) when each news is published. Because WTI crude oils are traded New York Mercantile Exchange, the UTC time is converted to Eastern Standard Time before further processing. Moreover, RPNA categorizes each event following the RavenPack taxonomy (figure: Ravenpack Event Taxonomy), which assigns a sequence of attributes describing each piece of news. In addition, for each event entry, using an algorithm combines results from surveying financial experts and pattern matching, RPNA assigns an Event Sentiment Score (ESS) between 0 and 100 to each event, measuring the short-term positive (100) and negative (0) financial or economic impact of this particular event. Refer to Appendix I for a complete list of attributes includes in this dataset.

Focusing our attention on events about crude oil only, there are 106, 960 entries from the raw dataset left (around 15 events per day). In the figure below, panel A presents a distribution of ESS for all news related to crude oil in the time span of 20 years and panel B shows all distributions of events within each year.

2.5 Classifying News Type

2.6 Case Studies of Events

2.6.1 Positive Spike on November 30, 2016

2.6.2 Negative Spike on December 6, 2018

2.6.3 Positive Spike on June. 12 ~ 13, 2019

3 Models

4 Experiments

References

Gibbons, M. R., & Hess, P. (1981). Day of the Week Effects and Asset Returns. *The Journal of Business*, 54(4), 579. doi: 10.1086/296147

- Hodges, J. L. (1958). The significance probability of the smirnov two-sample test. *Arkiv för Matematik*, 3(5), 469–486. doi: 10.1007/bf02589501
- Smirnov, N. (1939). On the estimation of the discrepancy between empirical curves of distribution for two independent samples. *Bulletin Moscow University*, 2, 3–16.

5 Appendix