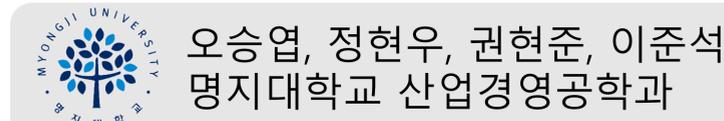


변수들을 통한 데이터 분석 기법으로 미래 금값 예측 Forecast Future Gold Price with Data Analysis Techniques



1. 개요

연구 배경

금융위기를 겪으며 금리가 낮아지고 안전자산을 선택하는 경향이 높아짐

이에 대해서 안전자산인 금값을 예측

연구 목표

현실세계에서 금값에 미치는 변수들을 탐색, 발굴하여 결과적으로 금값에 가장 큰 영향을 미치는 변수와 미래의 금값을 예측하는 것이다.

2. 이론적 배경

관련 이론

-데이터 마이닝
: 일반적인 의미는 대규모로 저장된 데이터 안에서 체계적이고 자동적으로 통계적 규칙이나 패턴을 찾아 내는 것

-선형 회귀 분석
: 회귀 분석은 관찰된 변수들에 대해 변수 사이의 모형을 구한 뒤 적합도를 측정해 내는 분석 방법

-주성분 분석
: 고차원의 데이터를 저차원의 데이터로 환원시키는 기법

-의사결정 나무
: 의사결정 규칙과 그 결과들을 트리 구조로 도식화한 의사 결정 지원 도구의 일종

3. 변수 선정

선정 이유

1. 금값의 변수 선정에 있어서 가장 중요한 요인은 금(Gold)은 수요와 공급에 따라 값이 결정되는 시장경제가 적용되는 것이다. 그렇기 때문에 대표적인 금융시장의 지표들을 변수로 선정

2. 금값은 '값'이라는 특성상 돈의 가치와 밀접한 관련이 있다고 판단하였다. 특히 기축통화인 달러를 선두로 하여 엔화, 유로화, 위안화 등등 국가들의 돈의 가치를 또한 변수로 추가

3. 금은 무형의 재산이 아닌 유형의 재산이기 때문에 금값에는 어느 정도의 물류비가 차지하고 있는 영역도 있을 것이라고 판단

4. 그렇기 때문에 국제 유가, 한국 유가, 두바이 유가 등을 변수로 선정

5. 마지막으로 경기 순환 이론에 의한 각종 비철금속 및 에너지 자원과 대표적 원자재의 가격을 변수로 추가

4. 설계 방법

실험 방법

실험은 데이터 분석 프로그램 중 XLMINER를 이용한다. 선정된 변수들을 입력한 뒤 위에서 언급한 기법들을 실행한다. 구체적인 실험 방법은 515개의 데이터를 Train : Validation : test의 데이터로 나누고 각각의 비율을 5 : 3 : 2 로 나누어서 분석을 진행한다.

Data Set

연도	1월	2월	3월	4월	5월	6월	7월	8월	9월	10월	11월	12월
2008	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2009	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2010	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2011	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2012	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2013	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2014	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2015	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2016	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2017	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2018	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2019	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2020	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2021	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2022	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2023	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2024	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2025	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2026	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2027	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2028	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2029	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00
2030	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00	1240.00

Data Set 분석

- 다중회귀분석
Y = aX1+bX2+cX3+...+C
처럼 진행되고 이때의 Y는 X로부터 1달 뒤의 예측 금값이며 X들은 국제 유가, 국내 휘발유, 백금, 아연, ... 등이 되며 C는 상수가 됨

- 주성분 분석
주성분으로 변경되는 것은 X에 해당하는 독립변수들이며 Y에 해당하는 종속변수는 변화하지 않고 주성분으로 변경된 X'와 Y를 다시 a. 경우와 같은 방식으로 다중회귀분석을 실시

- 의사결정 나무
X는 a.와 마찬가지로 독립변수이며 Y는 1달 후의 금값으로 지정하여 실시한다. 의사결정 나무의 맨 마지막 Node는 결과값 Y의 각각의 평균값으로 결정
이때 오분류율이 증가하는 시점에서 적절한 가지치기를 수행한다. 이를 통해 가장 상위의 변수와 2위, 3위의 변수 등을 확인

5. 결과 / 결론

실험 결과

- 다중회귀분석 결과

Residual DF	224
R ²	0.9906963
Adjusted R ²	0.9893257
Std. Error Estimate	436.38469
RSS	42656679

$$R^2 = 0.99$$

평균오차 - 436.38

금값을 예측하는 회귀식

Input Variables	Coefficient	Output
Intercept	-36958.26	1240.00
WTI	-28.91743	1240.00
KR_Oil(won)	9.7446006	1240.00
Platinum	5.0108297	1240.00
Coal	-39.56231	1240.00
Uranium	150.92284	1240.00
Palladium	1.7402597	1240.00
Zinc	-1.055788	1240.00
Lead	0.2334721	1240.00
Nickel	0.3140876	1240.00
USD	-129.8972	1240.00
JPY	23.963357	1240.00
EUR	10.69828	1240.00

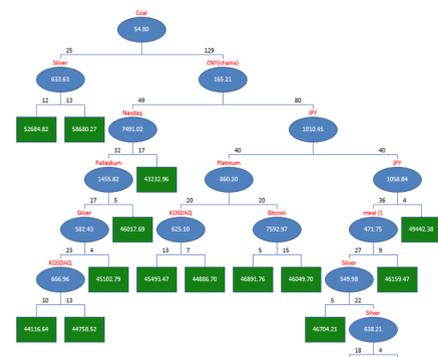
- 주성분 분석 결과

Variances	1	2	3	4	5	6	7
Variance	9479559.359	2915810.92	1110925.9	829934.5	92533.668	32827.83951	8028.4326
Variance Per	65.46371509	20.13593755	7.6718055	5.7313419	0.6390168	0.226701712	0.0554426
Cumulative	65.46371509	85.59965264	93.271458	99.0028	99.641817	99.86851857	99.923961

Residual DF	223
R ²	0.99707041
Adjusted R ²	0.9893284
Std. Error Estimate	436.87012
RSS	42560778

다음과 같은 결과를 보였는데 이는 다중회귀분석과 매우 근소한 차이를 보이며 이는 Data를 랜덤으로 5:3:2로 선정하여 모델화 하였을 때 보이는 오차범위로 보여진다. 결론적으로 다중회귀분석과 주성분 회귀분석의 차이는 없는 것으로 보인다.

- 의사결정 나무 결과



Test Data scoring - Summary Report (Using Best Pruned Tree)

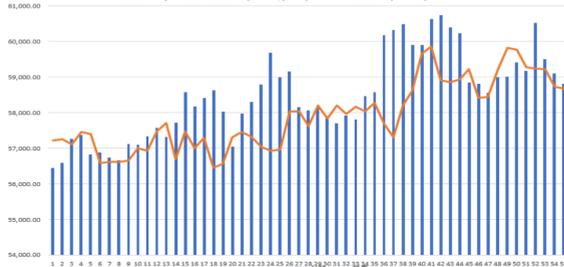
Total sum of squared errors	41557183
RMS Error	635.19115
Average Error	-24.97585229

다음과 같은 결과를 보였는데 금값이 예측치가 약 635원 정도의 평균적 오차를 보임

다중회귀분석과 비교하였을 때 별도로 생각할 것도 없이 다중회귀분석이 더 우수한 것으로 보인다.

결론

실제 금값과 예측 모델의 비교



*정확도 측정기준

정확한 금값이 아닌 금값이 전반에 비해 상승인지 하락인지 맞출 것(O) 예측 금값이 상승이고, 실제 금값이 상승한 경우 상승량은 관계가 없음(O) 예측 금값이 하락이고, 실제 금값이 하락한 경우 하락량은 관계가 없음(O) 예측 금값이 상승인데, 실제 금값이 하락한 경우 카운트(X) 예측 금값이 하락인데, 실제 금값이 상승한 경우 카운트(X)

- 정확도 측정결과 - 55개의 값 중 8번의 카운트
- 약 85%의 정확도

한계성

예측 모델의 명확한 한계점은 바로 연속형의 종속변수를 예측 할 때, 정확한 값을 가져도 실제 값과는 큰 차이가 있다. 현실의 많은 고려 사항들을 변수화 시키기 복잡하기 때문에 이를 예측모델에 담지 못함.