

EXPERT INSIGHT

# Machine Learning for Algorithmic Trading

퀀트 투자를 위한 머신러닝·딥러닝  
알고리즘 트레이딩 2/e

파이썬, Pandas, 텐서플로 2.0, Scikit-learn을 활용한  
효과적인 트레이딩 전략

Second Edition



스테판 젠슨 지음 홍창수·이기홍 옮김

## 출판사 책소개

머신러닝과 딥러닝 기술을 이용해 알고리즘 트레이딩의 아이디어에서 실행까지 전반적인 프로세스를 서술하는 좋은 안내서다. 2 판에서는 전략 백테스팅, 오토인코더, 적대적 생성 신경망(GAN), 이미지 형식으로 변환된 시계열에 합성곱 신경망(CNN) 적용과 같은 최신 내용을 추가했다. 또한 금융 시장 분석으로 100 개 이상의 알파 팩터를 설명하고 있으며, 부록으로 기술적 분석에 관한 모든 것을 수록했다. 세계적으로 저명한 AI 콘퍼런스인 NeurIPS 에서 발표된 최신 금융 딥러닝 연구를 소개하고 구현한다.

## 저자소개

### 저자: 스테판 젠슨 Stefan Jansen

어플라이드 에이아이(Applied AI)의 창립자이자 CEO 다. 포춘지 선정 500 대 기업, 투자 기업 및 업계 전반에서 데이터와 AI 전략에 대해 조언하고, 데이터 과학 팀을 구성하고, 광범위한 비즈니스 문제에 대한 엔드투엔드 머신 학습 솔루션을 개발하고 있다.이전에는 국제 투자회사의 파트너이자 상무이사를 역임하며 예측 분석 및 투자 연구 업무를 구축했다. 15 개 시장에서 영업을 하는 글로벌 핀테크 기업의 고위 임원으로 신흥국 중앙은행에 자문, 세계은행과 상담하기도 했다.조지아 공대에서 컴퓨터 공학 석사 학위를, 하버드 및 자유 대학교 베를린에서 경제학 석사 학위를 취득했고 CFA 자격증을 보유하고 있다. 유럽, 아시아, 미주 지역에서 6 개 언어로 근무했으며 제너럴 어셈블리(General Assembly)와 데이터캠프(Datacamp)에서 데이터 과학을 가르쳤다.

### 역자: 홍창수

경북대학교에서 재무 전공으로 경영학 석사 학위를 받고, 한국외국어대학교에서 「장외 개별주식옵션의 내재변동성에 관한 실증분석」으로 경영학 박사 학위를 받았다. 2000 년 외환선물(주) 투자공학 팀 애널리스트로 첫 직장 생활을 시작한 이후로 리딩투자증권 파생상품운용 팀 옵션 딜러, 한국투자증권 리스크관리부 리스크퀀트로 근무했다. 한화투자증권에서 장외 파생상품경영인가 취득에 참여한 이후 금융공학 팀, OTC 파생 팀에서 장외파생상품 담당 차장으로 근무했다. 이후 금융자산 평가회사인 NICE P&I(주) 금융공학연구소 실장으로 근무했으며, 현재 정보사업본부에서 장외 파생상품 및 리스크 관련 금융공학 컨설팅 업무를 수행하고 있다. 지은 책으로는 『장외 파생상품 실무입문』(서울경제경영, 2014)이 있다. SCI 급 저널인 PLOS ONE 을 비롯해 국내 유수의 재무금융 및 경영 학술지인 증권학회지, 재무관리연구, 금융공학연구, 경영연구 등에 재무 및 파생상품 관련 논문을 게재했다. 중앙대학교에서 금융공학 융합전공 강사로 '금융공학 실습' 과목을 강의했으며, 한국금융공학회 산학협력위원으로 활동하고 있다.

### 역자: 이기흥

카네기멜론 대학교에서 석사 학위를 받았고, 피츠버그 대학교 Finance Ph.D, CFA,FRM 이며, 금융, 투자, 경제분석 전문가다. 삼성생명, HSBC, 새마을금고 중앙회, 한국투자공사 등과 같은 국내 유수의 금융 기관, 금융 공기업에서 자산운용 포트폴리오 매니저로 근무했으며, 현재 딥러닝과 강화학습을 금융에 접목시켜 이를 전파하고 저변을 확대하는 것을 보람으로 삼고 있다. 저서로는 『엑셀 VBA 로 쉽게 배우는 금융공학 프로그래밍』(한빛미디어, 2009)이 있으며, 번역서로는 『포트폴리오 성공 운용』(미래에셋투자교육연구소,

2010), 『딥러닝 부트캠프 with 케라스』(길벗, 2017), 『프로그래머를 위한 기초 해석학』(길벗, 2018)과 『실용 최적화 알고리즘』(에이콘, 2020), 『초과 수익을 찾아서 2/e』(에이콘, 2020), 『자산운용을 위한 금융 머신러닝』(에이콘, 2020), 『실전 알고리즘 트레이딩 배우기』(에이콘, 2021), 『존 헐의 비즈니스 금융 머신러닝 2/e』(에이콘, 2021) 등이 있다. 누구나 자유롭게 머신러닝과 딥러닝을 자신의 연구나 업무에 적용해 활용하는 그날이 오기를 바라며 매진하고 있다.

## 작가의 말

이 글을 읽고 있다면, 여러분은 아마 투자 산업을 포함한 많은 산업에서 머신러닝(ML)이 전략적인 역량이 됐다는 것을 알고 있을 것이다. ML의 부상과 밀접한 관련이 있는 디지털 데이터의 폭발은 특히 투자에 강력한 영향을 미치고 있으며, 이미 정교한 모델을 사용해 정보를 처리해 온 오랜 역사를 갖고 있다. 이러한 추세는 계량 투자를 새로운 방식으로 접근할 수 있게 만들었으며, 재량적 거래 전략과 알고리즘 거래 전략 모두에 대한 데이터 과학의 수요를 증가시키고 있다. 자산클래스 간 거래 범위는 주식과 국채에서 상품과 부동산에 이르기까지 광범위하다. 이는 매우 광범위한 새로운 대체 데이터 소스가 시장 위와 그 밖의 데이터, 과거에 대부분의 분석 노력의 중심에 있었던 기초 데이터와 관련될 수 있음을 시사한다. 머신러닝이나 데이터 과학을 성공적으로 적용하려면 개인 또는 팀 차원의 통계 지식, 컴퓨터 기술과 도메인 전문 지식이 통합돼야 한다. 다시 말해 올바른 질문을 하고, 답을 제공할 수 있는 데이터를 식별하고 이해하며, 결과를 얻기 위한 광범위한 도구를 배포하고, 올바른 결정을 내리는 방식으로 이를 해석하는 것이 필수다. 따라서 이 책은 머신러닝의 투자 및 트레이딩 영역에 대한 통합적 관점을 제공한다.

## 목차

### 1 부. 데이터, 알파 팩터, 포트폴리오

#### 1 장. 트레이딩용 머신러닝: 아이디어에서 주문 집행까지

- \_\_ 투자업계에서 머신러닝의 부상
- \_\_ 전자 거래에서 고빈도 매매(HFT)까지
- \_\_ 팩터 투자와 스마트 베타 펀드
- \_\_ 알고리즘 개척자는 인간보다 우위에 있다
- \_\_ 머신러닝과 대체 데이터
- \_\_ 클라우드 소싱 거래 알고리즘
- \_\_ 머신러닝 기반 전략의 설계와 실행
- \_\_ 데이터의 소싱과 관리
- \_\_ 알파 팩터 리서치에서 포트폴리오 관리까지
- \_\_ 전략 백테스팅
- \_\_ 트레이딩을 위한 머신러닝: 전략과 사용 사례
- \_\_ 알고리즘 트레이딩의 진화
- \_\_ 거래를 위한 머신러닝 사용 사례
- \_\_ 요약

## 2 장. 시장 데이터와 기본 데이터: 소스와 기법

- 시장 데이터는 시장 환경을 반영한다
- 시장 미시 구조: 기본 지식
- 트레이딩 방법: 여러 종류의 주문
- 트레이딩 장소: 거래소에서 다크풀까지
- 고빈도 데이터로 작업
- 나스닥 주문 호가창으로 작업
- FIX 프로토콜을 활용한 거래의 통신
- 나스닥 토탈뷰-잇치 데이터 피드
- 틱에서 바로: 시장 데이터 정규화 방법
- AlgoSeek 분 바: 주식 호가와 거래 데이터
- 시장 데이터에 대한 API 접근
- 판다스를 이용한 원거리 데이터 접근
- yfinance: 데이터를 야후! 파이낸스로부터 스크레이핑
- 퀀토피안
- 집라인
- 퀀들
- 기타 시장 데이터 공급자
- 기본 데이터 작업 방법
- 재무제표 데이터
- 기타 기본 데이터 소스
- 판다스를 이용한 효율적 데이터 저장
- 요약

## 3 장. 금융을 위한 대체 데이터: 범주와 사용 사례

- 대체 데이터 혁명
- 대체 데이터의 원천
- 개인
- 비즈니스 프로세스
- 센서
- 대체 데이터 평가를 위한 기준
- 신호 내용의 질
- 데이터의 질
- 기술적 측면
- 대체 데이터 시장
- 데이터 제공업체와 사용 사례
- 대체 데이터로 작업
- 오픈테이블 데이터 스크래핑
- 어닝 콜 트랜스크립트 스크래핑과 파싱
- 요약

#### 4 장. 알파 팩터 리서치

- 알파 팩터 실무: 데이터부터 시그널까지
  - 알파 팩터 공학
    - 수십 년간의 팩터 연구를 바탕으로 구축
    - 모멘텀과 심리: 추세는 당신의 친구다
    - 가치 팩터
    - 변동성과 규모 이례 현상
  - 우량주 팩터
    - 수익률을 예측하는 알파 팩터 공학
    - 판다스와 넘파이를 사용한 팩터 공학
    - 칼만 필터를 사용한 알파 팩터의 잡음 제거
    - 웨이블릿을 이용한 잡음 데이터의 전처리 방법
  - 시그널에서 트레이딩으로: 백테스트를 위한 집라인
    - 단일 팩터 전략을 백테스트하는 방법
    - 다양한 데이터 소스의 팩터 결합
  - 알파렌즈를 사용한 시그널과 잡음의 분리
  - 미래 수익률과 팩터 분위수의 생성
  - 팩터 5 분위수에 의한 예측 성과
  - 정보 계수
  - 팩터 회전을
- 알파 팩터 리소스
- 기타 알고리즘 트레이딩 라이브러리
  - 요약

#### 5 장. 포트폴리오 최적화와 성과 평가

- 포트폴리오 성과를 측정하는 방법
  - 단일 숫자로 위험 수익 트레이드오프 포착
  - 적극적 운용의 근본 법칙
- 포트폴리오 위험과 수익률 관리하는 방법
  - 현대 포트폴리오 관리의 진화
  - 평균 분산 최적화
  - 평균 분산 최적화의 대안
  - 리스크 패리티
  - 리스크 팩터 투자
  - 계층적 리스크 패리티
- 집라인으로 트레이딩과 포트폴리오 관리
  - 신호 생성과 거래 실행 계획 수립
  - 평균 분산 포트폴리오 최적화 구현
- 파이폴리오를 이용한 백테스팅 성과 측정
  - 수익률과 벤치마크 입력 생성
  - 표본 외 수익률로 전진 분석

\_\_요약

## 2 부. 트레이딩을 위한 ML: 기본

### 6 장. 머신러닝 프로세스

\_\_데이터로부터 머신러닝이 작동하는 방법

\_\_도전 과제: 알고리즘을 작업에 매칭

\_\_지도학습: 예제에 의한 학습

\_\_비지도학습: 유용한 패턴의 발견

\_\_강화학습

\_\_머신러닝 워크플로

\_\_기본 설명: k-최근접 이웃

\_\_문제의 구성: 목적과 성과 측정

\_\_데이터의 수집과 준비

\_\_특성 탐험, 추출, 특성 공학

\_\_ML 알고리즘 선택

\_\_모델 설계와 조정

\_\_모델 선택을 위한 교차 검증의 활용

\_\_파이썬으로 교차 검증을 구현하는 방법

\_\_금융에서 교차 검증의 문제

\_\_사이킷런을 이용한 파라미터 조정과 엘로우브리크

\_\_요약

### 7 장. 선형 모델: 리스크 팩터에서 수익률 예측까지

\_\_추론에서 예측까지

\_\_베이스라인 모델: 다중 선형 회귀 모델

\_\_모델을 만드는 방법

\_\_모델을 훈련하는 방법

\_\_가우스-마르코프 정리

\_\_통계적 추론을 수행하는 방법

\_\_문제를 진단하고 해결하는 방법

\_\_실제로 선형 회귀를 실행하는 방법

\_\_statsmodels 를 활용한 최소 자승법

\_\_sklearn 을 활용한 확률적 경사 하강법(SGD)

\_\_선형 요인 모델을 구축하는 방법

\_\_CAPM 에서 파마-프렌치 5 요인 모델까지

\_\_위험 요인 입수

\_\_파마-맥베스 회귀 분석

\_\_수축 방법을 사용한 선형 회귀 규제화

\_\_과적합을 방지하는 방법

\_\_리지 회귀 분석의 작동 방식

- \_\_\_ 라쏘 회귀 분석의 작동 방식
- \_\_\_ 선형 회귀로 수익률을 예측하는 방법
- \_\_\_ 모델 특성과 선도 수익률 준비
- \_\_\_ statsmodels 를 사용한 선형 OLS 회귀 분석
- \_\_\_ 사이킷런을 이용한 선형 OLS 회귀
- \_\_\_ 사이킷런을 이용한 리지 회귀 분석
- \_\_\_ sklearn 을 이용한 라쏘 회귀 분석
- \_\_\_ 예측 신호의 품질 비교
- \_\_\_ 선형 분류
- \_\_\_ 로지스틱 회귀 모델
- \_\_\_ statsmodels 로 추론을 수행하는 방법
- \_\_\_ 가격 움직임을 로지스틱 회귀 분석으로 예측
- \_\_\_ 요약

## 8 장. ML4T 작업 흐름: 모델에서 전략 백테스트까지

- \_\_\_ ML 기반 전략의 백테스트 방법
- \_\_\_ 백테스트의 단점과 이를 회피하는 방법
- \_\_\_ 데이터를 올바르게 얻기
- \_\_\_ 시뮬레이션을 올바르게 수행
- \_\_\_ 통계량을 올바르게 얻기
- \_\_\_ 백테스트 엔진 작동법
- \_\_\_ 벡터화 대 이벤트 기반 백테스트
- \_\_\_ 주요 구현 측면
- \_\_\_ 백트레이더: 로컬 백테스트용 유연한 도구
- \_\_\_ 백트레이더의 세레브로 구조의 주요 개념
- \_\_\_ 실제에서 백트레이더를 사용하는 방법
- \_\_\_ 백트레이더 요약과 다음 단계
- \_\_\_ 집라인: 퀀토피안이 만든 확장 가능한 백테스트
- \_\_\_ 강건한 시뮬레이션을 위한 캘린더와 파이프라인
- \_\_\_ 분 데이터를 사용해 자체 번들 인제스트
- \_\_\_ 파이프라인 API: ML 시그널 백테스트
- \_\_\_ 백테스트 중 모델을 훈련하는 방법
- \_\_\_ 노트북 사용법
- \_\_\_ 요약

## 9 장. 시계열 모델

- \_\_\_ 진단과 특성 추출을 위한 분석 도구
- \_\_\_ 시계열 패턴의 분해
- \_\_\_ 롤링 윈도우 통계량과 이동 평균
- \_\_\_ 자기 상관계수의 측정
- \_\_\_ 정상성 진단과 회복

- \_\_\_ 정상성을 얻기 위한 시계열 변환
- \_\_\_ 단위근의 처리 방법
- \_\_\_ 시계열 변환의 실제 적용
- \_\_\_ 일변량 시계열 모델
- \_\_\_ 자기 회귀 모델의 구축
- \_\_\_ 이동 평균 모델의 구축
- \_\_\_ ARIMA 모델 구축과 확장
- \_\_\_ 매크로 펀더멘털 예측
- \_\_\_ 시계열 모델을 활용한 변동성 예측
- \_\_\_ 다변량 시계열 모델
- \_\_\_ 방정식 체계
- \_\_\_ 벡터 자기 회귀(VAR) 모델
- \_\_\_ VAR 모델을 활용한 매크로 펀더멘털 예측
- \_\_\_ 공적분: 공통 추세를 가진 시계열
- \_\_\_ 앵글-그레인저 2 단계 방법
- \_\_\_ 요한센 우도 비율 테스트
- \_\_\_ 공적분을 활용한 통계적 차이 거래
- \_\_\_ 같이 움직이는 자산 쌍을 선택하는 방법
- \_\_\_ 페어 트레이딩 실제 구현
- \_\_\_ 전략 백테스트 준비
- \_\_\_ 백트레이더를 이용한 전략 백테스트
- \_\_\_ 확장: 개선 방법
- \_\_\_ 요약

## 10 장. 베이지안 머신러닝: 동적 샵프 비율과 페어 트레이딩

- \_\_\_ 베이지안 머신러닝 학습법
- \_\_\_ 경험적 증거에서 가정을 업데이트하는 방법
- \_\_\_ 정확한 추론: 최대 사후 확률 추정
- \_\_\_ 결정론적 및 확률론적 근사 추론
- \_\_\_ PyMC3 를 사용한 확률 프로그래밍
- \_\_\_ 테아노를 이용한 베이지안 머신러닝
- \_\_\_ PyMC3 워크플로: 경기 침체 예측
- \_\_\_ 트레이딩을 위한 베이지안 머신러닝
- \_\_\_ 성과 비교를 위한 베이지안 샵프 비율
- \_\_\_ 페어 트레이딩을 위한 베이지안 롤링 회귀 분석
- \_\_\_ 확률 변동성 모델
- \_\_\_ 요약

## 11 장. 랜덤 포레스트: 일본 주식 롱/숏 전략

- \_\_\_ 결정 트리: 데이터로부터의 규칙 학습
- \_\_\_ 트리의 결정 규칙 학습과 적용



- \_\_\_ 결정 트리의 실용적 사용법
- \_\_\_ 과대적합과 규제화
- \_\_\_ 하이퍼파라미터의 조정
- \_\_\_ 랜덤 포레스트: 트리의 신뢰성을 높인다
- \_\_\_ 앙상블 모델의 성과가 더 좋은 이유
- \_\_\_ 배경
- \_\_\_ 랜덤 포레스트를 구축하는 방법
- \_\_\_ 랜덤 포레스트의 학습과 조정
- \_\_\_ 랜덤 포레스트를 위한 특성 중요도
- \_\_\_ 아웃오브백 테스트
- \_\_\_ 랜덤 포레스트의 장점과 단점
- \_\_\_ 일본 주식에 대한 롱/숏 시그널
- \_\_\_ 데이터: 일본 주식
- \_\_\_ LightGBM 을 사용한 ML4T
- \_\_\_ 전략: 집라인을 이용한 백테스트
- \_\_\_ 요약

## 12 장. 거래 전략 강화

- \_\_\_ 시작하기: 적응형 부스팅
- \_\_\_ 에이다부스트 알고리즘
- \_\_\_ 에이다부스트를 사용해 월별 가격 변동 예측
- \_\_\_ 그래디언트 부스팅: 대부분의 작업을 위한 앙상블
- \_\_\_ GBM 모델을 훈련하고 조정하는 방법
- \_\_\_ 사이킷런으로 그래디언트 부스팅을 사용하는 방법
- \_\_\_ XGBoost, LightGBM, CatBoost 사용
- \_\_\_ 알고리즘 혁신이 성능을 향상시키는 방법
- \_\_\_ 부스팅을 사용한 롱/숏 거래 전략
- \_\_\_ LightGBM 과 CatBoost 로 시그널 생성
- \_\_\_ 블랙박스의 내부: GBM 결과를 해석하는 방법
- \_\_\_ 부스팅 앙상블에 기반을 둔 전략 백테스팅
- \_\_\_ 학습한 교훈과 다음 단계
- \_\_\_ 일중 전략을 위한 부스팅
- \_\_\_ 고빈도 데이터를 위한 특성 공학
- \_\_\_ LightGBM 으로 분 단위 빈도 시그널
- \_\_\_ 거래 시그널 품질 평가
- \_\_\_ 요약

## 13 장. 비지도학습을 활용한 데이터 기반 리스크 팩터와 자산 배분

- \_\_\_ 차원 축소
- \_\_\_ 차원의 저주
- \_\_\_ 선형 차원 축소

- \_\_\_ 매니폴드 학습: 선형 차원 축소
- \_\_\_ 트레이딩을 위한 PCA
- \_\_\_ 데이터 주도형 리스크 팩터
- \_\_\_ 고유 포트폴리오
- \_\_\_ 군집화
- \_\_\_ k-평균 군집화
- \_\_\_ 계층적 군집화
- \_\_\_ 밀도 기반 군집화
- \_\_\_ 가우시안 혼합 모델
- \_\_\_ 최적 포트폴리오를 위한 계층적 군집화
- \_\_\_ 계층적 리스크 패리티의 작동 원리
- \_\_\_ ML 트레이딩 전략을 사용한 HRP 백테스트
- \_\_\_ 요약

### 3 부. 자연어 처리

- 14 장. 트레이딩을 위한 텍스트 데이터: 감성 분석
- \_\_\_ 텍스트를 이용한 머신러닝: 언어에서부터 특성까지
  - \_\_\_ 텍스트 데이터 작업의 주요 과제
  - \_\_\_ 자연어 처리 워크플로
  - \_\_\_ 응용 사례
  - \_\_\_ 텍스트에서 토큰으로: NLP 파이프라인
  - \_\_\_ spaCy 와 textacy 를 이용한 NLP 파이프라인
  - \_\_\_ TextBlob 을 사용한 자연어 처리
  - \_\_\_ 토큰 계산: 문서 단어 행렬
  - \_\_\_ 단어 주머니 모델
  - \_\_\_ sklearn 을 사용한 문서 단어 행렬(DTM)
  - \_\_\_ 핵심적인 교훈
  - \_\_\_ 트레이딩을 위한 자연어 처리
  - \_\_\_ 나이브 베이즈 분류기
  - \_\_\_ 뉴스 기사 분류
  - \_\_\_ 트위터와 옐프 데이터를 사용한 감성 분석
  - \_\_\_ 요약

15 장. 토픽 모델링

- \_\_\_ 잠재 토픽 학습: 목적과 접근법
- \_\_\_ 잠재 의미 인덱싱
- \_\_\_ sklearn 을 이용한 LSI 의 구현 방법
- \_\_\_ 장점과 제약점
- \_\_\_ pLSA
- \_\_\_ sklearn 을 이용한 pLSA 의 구현 방법

\_\_장점과 제한점

\_LDA

\_\_LDA의 원리

\_\_LDA 토픽 평가

\_\_sklearn을 이용한 LDA의 구현 방법

\_\_pyLDAvis를 이용한 LDA 결과의 시각화

\_\_Gensim을 사용한 LDA의 구현 방법

\_실적 발표 콜을 위한 토픽 모델링

\_\_데이터 전처리

\_\_모델 훈련과 평가

\_\_실험 실행

\_금융 뉴스를 위한 토픽 모델링

\_요약

16 장. 어닝 콜과 SEC 공시 보고서를 위한 단어 임베딩

\_단어 임베딩이 의미를 인코딩하는 방법

\_\_신경 언어 모델이 문맥에서 사용법을 배우는 방법

\_\_word2vec: 확장 가능한 단어와 구분 임베딩

\_\_의미론적 산술을 활용한 임베딩 평가

\_사전 훈련된 단어 벡터를 사용하는 방법

\_\_GloVe: 단어 표현을 위한 글로벌 벡터

\_금융 뉴스를 위한 사용자정의 임베딩

\_\_전처리: 문장 인식과 앵그램

\_\_텐서플로 2를 이용한 스킵그램 아키텍처

\_\_텐서보드를 사용해 임베딩 시각화

\_\_Gensim으로 임베딩을 더 빠르게 훈련하는 방법

\_SEC 공시 문서를 이용한 트레이딩을 위한 word2vec

\_\_전처리: 문장 인식과 앵그램

\_\_모델 훈련

\_doc2vec 임베딩을 사용한 감성 분석

\_\_엘프 감성 데이터에서 doc2vec 입력 생성

\_\_doc2vec 모델 훈련

\_\_문서 벡터로 분류기 훈련

\_\_학습한 교훈과 다음 단계

\_새로운 개척자: 사전 훈련된 트랜스포머 모델

\_\_어텐션은 당신이 필요한 모든 것

\_\_BERT: 좀 더 보편적인 언어 모델을 지향

\_\_텍스트 데이터에 대한 트레이딩: 교훈과 다음 단계

\_요약

4 부. 딥러닝과 강화학습

## 17 장. 딥러닝

- 딥러닝: 무엇이 새롭고 왜 중요한가
- 계층적 구조가 고차원 데이터의 문제를 완화한다
- 표현 학습으로서의 DL
- DL의 ML과 AL와의 연관성
- 신경망의 설계
  - 단순 순전파 신경망 구조
  - 주요 설계 선택
  - 심층 신경망의 규제화
  - 더 빠른 훈련: DL을 위한 최적화
  - 요약: 주요 하이퍼파라미터를 조정하는 방법
- 파이썬을 이용해 처음부터 신경망을 구축하는 방법
  - 입력층
  - 은닉층
  - 출력층
  - 순전파
  - 교차 엔트로피 비용 함수
  - 파이썬을 사용해 역전파를 구현하는 방법
- 인기 있는 딥러닝 라이브러리
  - GPU 가속화의 활용
  - 텐서플로 2를 사용하는 방법
  - 텐서보드 사용법
  - 파이토치 1.4 사용법
  - 다른 옵션
- 룭/숏 전략을 위한 NN의 최적화
  - 일 수익률을 예측하기 위한 특성 공학
  - NN 구조 프레임워크 정의
  - NN을 조정하한 설계 옵션 교차 검증
  - 예측 성능 평가
  - 앙상블 시그널을 기반으로 하는 전략의 백테스트
  - 결과를 더욱 개선하는 방법
- 요약

## 18 장. 금융 시계열과 인공위성 이미지를 위한 CNN

- CNN이 그리드 같은 데이터를 모델링하는 방법
- 수작업 코딩부터 데이터의 학습 필터까지
- 합성곱층의 요소 연산 방법
- CNN 아키텍처의 진화: 주요 혁신
- 인공위성 이미지와 물체 인식을 위한 CNN
- LeNet5: 산업용 애플리케이션을 갖춘 최초의 CNN

- \_\_\_ 알렉스넷: 딥러닝 연구의 재점화
- \_\_\_ 전이학습: 적은 데이터로 더 빠른 훈련
- \_\_\_ 객체 탐지와 분할
- \_\_\_ 실제 객체 탐지
- \_\_\_ 시계열 데이터용 CNN: 수익률 예측
- \_\_\_ 1 차원 합성곱이 있는 자기 회귀 CNN
- \_\_\_ CNN-TA-2D 형식의 시계열 군집화
- \_\_\_ 요약

## 19 장. 순환 신경망

- \_\_\_ RNN 의 작동 원리
- \_\_\_ 순환 주기로 계산 그래프 펼치기
- \_\_\_ 시간을 통한 역전파
- \_\_\_ 대안적 순환 신경망 구조
- \_\_\_ 심층 RNN 의 설계
- \_\_\_ 장기 의존성 학습 시 문제점
- \_\_\_ GRU
- \_\_\_ 텐서플로 2 를 이용한 시계열을 위한 RNN
- \_\_\_ 일변량 회귀: S&P 500 예측
- \_\_\_ 시계열 데이터를 RNN 용 형태로 변환
- \_\_\_ 적층 LSTM: 주가 움직임과 수익률 예측
- \_\_\_ 매크로 데이터에 대한 다변량 시계열 회귀
- \_\_\_ 텍스트 데이터를 위한 RNN
- \_\_\_ 감성 분류를 위한 LSTM 과 단어 임베딩
- \_\_\_ 사전 학습된 단어 벡터를 이용한 감성 분석
- \_\_\_ SEC 보고서 임베딩으로부터 수익률 예측
- \_\_\_ 요약

## 20 장. 조건부 위험 요인과 자산 가격 결정을 위한 오토인코더

- \_\_\_ 비선형 특성 추출을 위한 오토인코더
- \_\_\_ 선형 차원 축소 일반화
- \_\_\_ 이미지 압축을 위한 합성곱 오토인코더
- \_\_\_ 규제화된 오토인코더로 과적합 관리
- \_\_\_ 노이즈 제거 오토인코더로 손상된 데이터 수정
- \_\_\_ 시계열 특성을 위한 Seq2seq 오토인코더
- \_\_\_ 변분 오토인코더로 생성 모델링
- \_\_\_ 텐서플로 2 로 오토인코더 구현
- \_\_\_ 데이터를 준비하는 방법
- \_\_\_ 단층 순전파 오토인코더
- \_\_\_ 희소성 제약 조건이 있는 순전파 오토인코더
- \_\_\_ 심층 순전파 오토인코더

- \_\_\_ 합성곱 오토인코더
- \_\_\_ 오토인코더 노이즈 제거
- \_\_\_ 거래를 위한 조건부 오토인코더
- \_\_\_ 주식 가격과 메타데이터 정보 습득
- \_\_\_ 예측 자산 특징 계산
- \_\_\_ 조건부 오토인코더 아키텍처 만들기
- \_\_\_ 교훈과 다음 단계
- \_\_\_ 요약

## 21 장. 합성 시계열 데이터를 위한 적대적 생성 네트워크

- \_\_\_ GAN 으로 합성 데이터 생성
- \_\_\_ 생성 모델과 판별 모델 비교
- \_\_\_ 적대적 훈련: 사기의 제로섬 게임
- \_\_\_ GAN 아키텍처 ZOO 의 급속한 진화
- \_\_\_ 이미지를 위한 GAN 애플리케이션과 시계열 데이터
- \_\_\_ 텐서플로 2 를 사용해 GAN 을 빌드하는 방법
- \_\_\_ 생성자 네트워크 빌드
- \_\_\_ 판별자 네트워크 생성
- \_\_\_ 적대적 훈련 프로세스 설정
- \_\_\_ 평가 결과
- \_\_\_ 합성 금융 데이터를 위한 TimeGAN
- \_\_\_ 특성과 시간에 걸쳐 데이터를 생성하는 방법 학습
- \_\_\_ 텐서플로 2 를 사용한 TimeGAN 구현
- \_\_\_ 합성 시계열 데이터의 품질 평가
- \_\_\_ 교훈과 다음 단계
- \_\_\_ 요약

## 22 장. 심층 강화학습: 트레이딩 에이전트의 구축

- \_\_\_ RL 시스템의 구성 요소
- \_\_\_ 정책: 상태에서 행동으로 전이
- \_\_\_ 보상: 행동으로부터 학습
- \_\_\_ 가치 함수: 장기적으로 최적 선택
- \_\_\_ 모델 기반 대 모델 프리: 도약 전에 보기?
- \_\_\_ RL 문제의 해법
- \_\_\_ RL 문제 풀이의 주요 문제점
- \_\_\_ RL 문제를 푸는 근본적인 방법
- \_\_\_ 동적 계획 문제 풀기
- \_\_\_ 유한 MDP
- \_\_\_ 정책 반복
- \_\_\_ 가치 반복

- \_\_\_ 일반화된 정책 반복
- \_\_\_ 파이썬을 이용한 동적 계획
- \_\_\_ Q-러닝: 바둑에서 최적 정책 찾기
- \_\_\_ 탐험 대 활용 트레이드오프:  $\epsilon$ -탐욕 정책
- \_\_\_ Q-러닝 알고리즘
- \_\_\_ 파이썬을 이용한 Q-러닝 에이전트 훈련법
- \_\_\_ OpenAI Gym 을 사용한 심층 강화학습
- \_\_\_ 신경망을 이용한 가치 함수의 근사
- \_\_\_ 딥 Q-러닝 알고리즘과 확장
- \_\_\_ OpenAI GYM 소개
- \_\_\_ 텐서플로 2 를 이용한 DDQN 의 구현
- \_\_\_ 간단한 트레이딩 에이전트 작성
- \_\_\_ 사용자 정의 OpenAI 트레이딩 환경을 설계하는 방법
- \_\_\_ Deep Q-learning on the stock market
- \_\_\_ 얻은 교훈
- \_\_\_ 요약

## 23 장. 결론과 다음 단계

- \_\_\_ 학습에서 얻은 주요 요점과 교훈
- \_\_\_ 데이터는 가장 중요한 단일 재료
- \_\_\_ 도메인 전문 지식: 잡음으로부터 시그널의 구별
- \_\_\_ ML 은 데이터로 문제를 풀기 위한 도구 모음이다
- \_\_\_ 백테스트의 과적합을 주의하라
- \_\_\_ 블랙박스 모델에서 통찰력을 얻는 방법
- \_\_\_ 실전 트레이딩을 위한 ML
- \_\_\_ 데이터 관리 테크놀로지
- \_\_\_ ML 도구
- \_\_\_ 온라인 트레이딩 플랫폼
- \_\_\_ 결론

## 부록. 알파 팩터 라이브러리

- \_\_\_ TA-Lib 에서 구현된 공통 요인 알파
- \_\_\_ 핵심 구성 요소: 이동 평균
- \_\_\_ 중첩 연구: 가격과 변동성 추세
- \_\_\_ 모멘텀 지표
- \_\_\_ 거래량과 유동성 지표
- \_\_\_ 변동성 지표
- \_\_\_ 기본 위험 요인
- \_\_\_ 정형화된 알파에 대한 월드퀀트의 탐구
- \_\_\_ 횡단면 및 시계열 함수
- \_\_\_ 정형화된 알파 표현식

\_\_이변량 및 다변량 요인 평가

\_\_정보 계수와 상호 정보량

\_\_특성 중요도와 SHAP 값

\_\_비교: 각 지표에 대한 Top 25 특성

\_\_금융 성과: 알파렌즈

## 출판사 서평

◆ 이 책에서 다루는 내용 :

- ◆ 투자와 트레이딩 문제를 푸는 머신러닝 기법 구현
- ◆ 시장, 기본, 대체 데이터를 활용한 알파 팩터 연구
- ◆ 지도학습, 비지도학습, 강화학습 모델 설계
- ◆ 판다스, 넘파이와 Scikit-learn 을 이용한 포트폴리오 위험과 성과 최적화
- ◆ 머신러닝 모델을 Quantopian 플랫폼 위에서 실전 트레이딩 전략으로 통합
- ◆ 신뢰성 있는 시계열 백테스트 방법으로 전략 평가
- ◆ 케라스, 파이토치와 텐서플로를 사용한 딥신경망 설계와 평가
- ◆ 강화학습을 이용한 OpenAI Gym 에서의 트레이딩 전략 개발

◆ 이 책의 대상 독자 ◆

금융시장에 대한 이해와 트레이딩 전략에 관심이 있는 애널리스트, 데이터 과학자, ML 엔지니어라면 이 책이 도움이 될 것이다. 또한 ML 을 활용해 더 나은 의사결정을 내리고자 하는 투자 전문가로서 가치도 찾을 수 있다. 소프트웨어와 ML 에 대한 배경지식이 있는 경우 이 영역에 대한 소개자료를 생략할 수 있다. 마찬가지로 전문 지식이 투자에 관한 것이라면 이 책이 제공하는 다양한 배경의 금융 맥락 일부 또는 전부를 잘 알고 있을 것이다. 이 책은 독자가 매우 역동적인 이 영역을 계속 배우고 싶어 한다고 가정한다. 이를 위해 깃허브(GitHub) 저장소 내의 각 장에 대해 README 파일에 연결된 수많은 학술 참고 자료와 추가 자료를 책의 끝에 제공한다. 파이썬(Python) 3 과 넘파이, 판다스, SciPy 와 같은 과학 컴퓨팅 라이브러리를 사용하는 것이 편할 것이며, 계속해서 다른 많은 자료를 얻을 수 있을 것을 기대한다. ML 과 사이킷런에 대한 경험이 도움이 되겠지만 기본 워크플로를 간략하게 다루고 격차를 메꾸거나 더 깊이 파고들려면 다양한 리소스를 참조하길 바란다. 마찬가지로 금융과 투자에 기본 지식이 있다면 용어를 더 쉽게 사용할 수 있을 것이다.

◆ 이 책의 구성 ◆

ML 이 트레이딩 전략의 설계와 실행에 가치를 더할 수 있는 방법을 종합적으로 소개한다. 데이터 소싱과 전략 개발 프로세스의 다양한 측면과 ML 과제에 대한 다양한 솔루션을 다루는 4 개의 부로 구성돼 있다.

1 부에서는 머신러닝을 활용하는 트레이딩 전략 전반에 걸쳐 관련된 기본적인 측면을 다룬다.

1 장, '트레이딩용 머신러닝: 아이디어에서 주문 집행까지'에서는 ML 이 거래에서 중요한 이유와 방법을 요약하고 투자 프로세스를 설명하며, ML 이 가치를 높일 수 있는 방법을 간략히 설명한다.

2 장, '시장 데이터와 기본 데이터: 소스와 기법'에서는 거래소 제공 틱 데이터와 보고된 재무 정보를 포함해 시장 데이터를 소싱하고 활용하는 방법을 다룬다.

3 장, '금융을 위한 대체 데이터: 범주와 사용 사례'에서는 폭발적으로 증가하는 수의 출처와 공급자를 평가하기 위한 범주와 기준을 설명한다.



4 장, '금융 특성 공학: 알파 팩터 리서치'에서는 예측 시그널을 포착하는 데이터 변환을 생성하고 평가 프로세스와 팩터 성과를 측정하는 방법을 설명한다.

5 장, '포트폴리오 최적화와 성과 평가'에서는 전략의 실행에 따른 포트폴리오를 관리, 최적화, 평가하는 방법을 소개한다.

2 부에서는 근본적인 지도학습과 비지도학습 알고리즘이 엔드투엔드 워크플로의 맥락에서 트레이딩 전략에 어떻게 정보를 제공할 수 있는지 설명한다.

6 장, '머신러닝 프로세스'에서는 ML 모델의 예측 성과를 체계적으로 수립, 훈련, 튜닝, 평가하는 방법을 개략적으로 설명함으로써 단계를 설정한다.

7 장, '선형 모델: 리스크 팩터에서 수익률 예측까지'에서는 선형 회귀와 로지스틱 회귀를 추론과 예측에 사용하는 방법과 규제화를 사용해 과적합 위험을 관리하는 방법을 설명한다.

8 장, 'ML4T 작업 흐름: 모델에서 전략 백테스팅까지'에서는 지금까지 별도로 다룬 ML4T 워크플로의 다양한 구성 요소를 통합한다.

9 장, '시계열 모델'에서는 벡터 자기 회귀 모델뿐 아니라 변동성 예측에 대한 ARCH/GARCH 모델을 비롯한 일변량 및 다변량 시계열 진단과 모델을 다룬다.

10 장, '베이저안 머신러닝: 동적 샵 비율과 페어 트레이딩'에서는 확률적 모델과 마르코프 체인 몬테카를로(MCMC) 샘플링 및 변분 베이스 방법이 대략적인 추론을 어떻게 용이하게 하는지 제시한다.

11 장, '랜덤 포레스트: 일본 주식 롱/숏 전략'에서는 통찰력과 예측을 위한 비선형 트리 기반 모델을 구축, 훈련, 조정하는 방법을 설명한다.

12 장, '거래 전략 강화'에서는 그래디언트 부스트(Gradient Boost)를 소개하고 고성능 훈련과 예측을 위해 XGBoost, LightBGM, CatBoost 라이브러리를 사용하는 방법을 보여준다.

13 장, '비지도학습을 활용한 데이터 기반 리스크 팩터와 자산 배분'에서는 알고리즘 거래를 위해 차원 축소와 군집화를 사용하는 방법을 설명한다.

3 부는 텍스트 데이터에 초점을 맞추고 대체 데이터의 핵심 출처에서 고품질 시그널을 추출하기 위한 최첨단 비지도학습 기법을 도입한다.

14 장, '트레이딩을 위한 텍스트 데이터: 감성 분석'에서는 텍스트 데이터를 수치형으로 변환하는 방법을 보여주고 감성 분석을 위한 2 부의 분류 알고리즘을 대규모 데이터 세트에 적용한다.

15 장, '토픽 모델링'에서는 비지도학습을 사용해 많은 문서를 요약한 주제를 추출하고 분류 모델의 특성으로 텍스트 데이터를 탐색하거나 주제를 사용하는 좀 더 효과적인 방법을 제공한다.

16 장, '어닝 콜과 SEC 공시 보고서를 위한 단어 임베딩'에서는 신경망을 사용해 기존 텍스트 특성보다 의미론적 콘텍스트를 훨씬 더 잘 포착하고 텍스트 데이터에서 거래 시그널을 추출하는 매우 유망한 방법을 나타내는 단어 벡터 형태의 최신 언어 특성을 학습한다.

4 부에서는 딥러닝과 강화학습을 소개한다.

17 장, '딥러닝'에서는 4 부에서 사용할 가장 인기 있는 딥러닝 프레임워크인 텐서플로 2와 파이토치(PyTorch)를 소개한다.

18 장, '금융 시계열과 인공지능 이미지를 위한 CNN'에서는 대규모의 비정형 데이터를 활용하는 분류 작업을 위해 매우 강력한 합성곱 신경망(CNN)을 다룬다.

19 장, '순환 신경망'에서는 순환 신경망(RNN)이 어떻게 시퀀스 대 시퀀스 모델링에 유용한지 보여준다. 이는 예측할 일변량 시계열과 다변량 시계열을 포함한다.

20 장, '조건부 위험 요인과 자산 가격 결정을 위한 오토인코더'에서는 고차원 데이터의 비선형 압축을 위한 오토인코더를 다룬다.

21 장, '합성 시계열 데이터를 위한 적대적 생성 네트워크'에서는 딥러닝에서 흥미로운 발전 중 하나를 소개한다.

22 장, '심층 강화학습: 트레이딩 에이전트의 구축'에서는 강화학습이 환경에 대응해 의사결정을 최적화하는 방법을 배우는 에이전트의 설계와 훈련을 어떻게 허용하는지 설명한다.

23 장, '결론과 다음 단계'에서는 학습한 교훈을 요약하고, 계속해서 학습하고 트레이딩 전략을 수립하고자 취할 수 있는 몇 가지 단계를 개략적으로 설명한다.

부록, '알파 팩터 라이브러리'에서는 거의 200 개의 인기 있는 금융 특성을 나열하고 근거를 설명하며, 이러한 특성을 계산하는 방법을 보여준다.

#### ◆ 옮긴이의 말 ◆

금융시장의 흐름을 이해하려면 시장 데이터, 공시, 뉴스, SNS 등 수많은 요인을 분석해야 합니다. 데이터를 효율적으로 관리하기 위해 금융시장에서도 점차 머신러닝과 딥러닝 기술이 사용되고 있으며, 트레이딩 분야에도 활용도가 높아진 상황입니다. 이 책은 머신러닝과 딥러닝 기술을 이용해 알고리즘 트레이딩의 아이디어에서 실행까지 전반적인 프로세스에 초점을 맞춰 서술하고 있는 이 분야의 바이블이 될 수 있는 좋은 안내서입니다. 2 판의 새로운 내용은 저자가 말했듯 더 많은 알파 팩터 사례를 소개하고 있으며, 최신 금융 딥러닝 논문에서 발표된 중요한 내용을 꼼꼼히 구현했다는 점이 특징입니다. 실제로 각 장의 내용을 상당 부분 고쳐가며 개정판 수준으로 설명해 초기에 생각했던 번역 기간보다 8 개월 정도가 더 소요됐습니다. 1 판보다 번역 품질을 높이려고 했고, 의역으로도 이해가 쉽지 않은 부분은 각주를 많이 달아 독자의 이해를 높이려고 노력했습니다. - 홍창수

최근에 일부 계량 금융학자들이 서클을 이루면서 AI 와 금융을 융합하려고 노력하고 있다. 이 책의 저자 스테판 켄슨도 이러한 노력에 상당한 공헌을 하고 있으며, 새로운 AI 의 흐름을 적극 반영하고자 1 판을 개정한 2 판을 출간했다. 이 책은 머신러닝의 전통 이론에서 최신 이론뿐 아니라 최신 딥러닝을 금융에 응용하는 부분도 다루고 있어 계량 금융과 머신러닝의 금융 응용에 관심 있는 사람들에게 특히 도움이 될 것이다. 1 판에서의 실험적인 내용들이 2 판에서는 많이 정립돼 독자들이 더욱 체계적으로 접근할 수 있을 것이다. 특히 2020 년, 일각에서 상당한 공헌을 하던 퀀토피안(Quantopian)이 문을 닫으면서 계량 금융 커뮤니티에 커다란 충격을 줬는데, 많은 학자와 실무 경험자가 모여 혁신적이면서도 실전적인 연구와 작품들을 만들어내고 있던 상황이라 충격은 더욱 컸다. 특히 이들이 만든 라이브러리인 집라인(zipline)과 파이폴리오(pyfolio)는 실제 선진 대형 자산운용사에서 사용하는 수준을 과시하고 있었으므로 이들이 사장되는가 하는 아쉬움이 많았다. 심지어 이 책의 초판도 이들을 광범위하게 사용하고 있어 개정판에서는 어떻게 전개할지 궁금하기도 했다. 놀랍게도 저자는 자체적으로 집라인을 개선한 집라인 리로드(zipline reloaded)를 만들어내 모든 내용을 살렸을 뿐 아니라 최신 파이썬 버전에서도 잘 작동하게 만들어 다시금 계량 금융 커뮤니티에 집라인 열풍을 일으키는 큰 공헌을 했다(참고로 요즘 핫한 인터넷 증권사인 알파카(Alpaca)와 인터랙티브 브로커(InteractiveBroker)와 같은 증권사가 거래 API 에 최근 집라인을 사용할 수 있도록 하고 있다). 저자의 이런 노력에 찬사를 보내며, 지속적인 개발이 이뤄지기를 바란다. 마지막으로 저자도 언급하듯이 방대한 내용인 만큼 모든 내용을 자세히 다루기는 힘들어 핵심 요점 위주로 서술돼 있으며 오히려 저자의 깃허브에 있는 노트북에 많은 내용을 실었다는 것을 다시 한번 상기시키고 싶다. 이 책을 통해 이론과 실전이 통합된 최첨단의 고급 계량 금융의 여정을 즐기기를 바란다. - 이기홍