



AI 융합전문가
(마스터 과정)

데이터분석

2024.06.22.





강사 소개_ 홍용기 컨설팅학 박사



소속	주요 업무
(주)제타데이터 대표이사	데이터 전략 및 컨설팅 ODA 컨설팅 (Official Development Assistance)
(주)지구파트너스 감사	창업보육, 투자, 기업·기술가치평가
(주)메타로직 컨설팅 수석	ISP 컨설팅 (Information Strategy Planning) ISMP 컨설팅 (Information System Master Plan)
(주)이너피플 이사	데이터분석, 데이터 가치평가 컨설팅, 데이터바우처 사업수행

◎ 자격증

1. 경영지도사31기 (인적자원, 2016)
2. 창업보육매니저 (BI협회, 2018)
3. 기업·기술가치평가사 (기업·기술가치평가협회, 2018)
4. 기업재난관리사 실무과정 (행정안전부, 2019)
5. 데이터분석 준전문가 ADsP (데이터산업진흥원 K-Data, 2021)
6. 빅데이터 분석기사 (과학기술정보통신부 · 통계청, 2021)
7. 국제공인컨설턴트 CMC (ICMCI, 2023)
8. 인공지능(AI) 활용마스터1급 (뉴미디어교육연구소, 2024)



홍용기 컨설팅학 박사

- 데이터분석 및 실전 R코딩 (경영기술지도사회, 빅데이터 분석기사 자격증 취득 과정)
- 데이터분석 Python 심화과정 (서울 여성능력개발원 강동 여성인력개발센터 / 용산 여성인력개발센터)
- AI & ChatGPT 활용 및 데이터분석 컨설팅 방법론 (경영기술지도사회, 국제공인컨설턴트 CMC 양성과정)
- 파이썬 코딩을 통한 크롤링 자동화 인텐시브 과정 (경영지도사 및 컨설턴트)

HOME > ChyGlobal

홍용기 컨설팅학박사, 한국경영기술지도사회 '빅데이터 분석 실무 과정' 성료

김전선 기자 | 송인 2023.03.27 16:54 | 댓글 0

경영컨설팅 과정에서 데이터 분석을 결합한 객관적인 사업계획수립, 마케팅 전략 제시 역량을 갖추고자

중요기사

- 외대교수 시국선언 "니들 알아들"...
- 한일 정상회담 전하는 KBS1 라디오...
- 민경욱 대표 "중앙선관위를 직무..."
- 김기현, 가족들을 청남을 만나 "말..."
- 언론단체들, "대통령은 기소된 한..."

인기뉴스

- 1 "호세킹 관련 선관위에 대한 의..."
- 2 [차재용 합법] 안국지능(AI) 시대, 위협과 가...
- 3 [황교안의 손원지] (359) 도둑이 재 발 저연...
- 4 [전정희 작가 합법] 가정의 달 5월을 맞이하여...
- 5 송도호 서울시의원, 전가차 중전구역 회재만...

(주)제타데이터 홍용기 대표, 강동여성인력개발센터/강동여성새로일하기센터의 '빅데이터 분석 파이썬(Python) 심화 과정' 성료

윤성일 | 송인 2023.06.05 11:53 | 댓글 0

중요기사

- 이동준 아들 친구, "난 박록희가 아니, 당시..."
- "대한민국 외교부 동북아 관련 부서 장신 차원..."
- 이와중에 '선략적 감시' 받았다는 선관위에 자제...

(주)제타데이터 홍용기 대표(컨설팅학박사/경영지도사)는 지난 5월 23일부터 6월 1일까지 서울시 강동여성인력개발센터/강동여성새로일하기센터에서 진행된 '빅데이터 분석에 필요한 파이썬(Python) 심화' 과정을 진행하였다.

(주)제타데이터 홍용기 대표, 용산여성인력개발센터의 '빅데이터 분석을 위한 파이썬(Python) 심화 과정' 성료

윤성일 | 송인 2023.06.10 19:51 | 댓글 0

중요기사

- "북한공공 영웅에게 40억원들..." 박민식 보훈장...
- 새KBS5공투위, "KBS의 지속적인 자사 공립 편향..."
- 김장겸 전 mbc사장, "국가 복력어린..."
- 법사위 살전 "최강육 안목거리지 달라 VS 한동..."
- 김의겸, 고(前) 재로 상행 사건 "우사 기록 갖고..."

(주)제타데이터 홍용기 대표는 지난 7월 13일부터 7월 26일까지 용산 전자랜드에서 빅데이터 분석을 위한 파이썬(Python) 심화과정 강의를 진행했다.

다자합업교육원 & TODAY
교육 (2023.06.10) 19:51

제1기 챗GPT 활용
핵심기초 양성 과정

강의장소: 용산 전자랜드 3층 301호 (용산역 5분 도보, 지하철 5호선 용산역 5분 도보)
강의기간: 2023.11.08 ~ 11.18 (11회)



홍용기 컨설팅학 박사

◎ 정보화전략계획수립(ISP) 컨설팅 수행

20.05~20.08. 창업진흥원
20.10~20.12. 한국연구재단
21.01~21.04. 소상공인시장진흥공단
21.11~22.06. 서울특별시
22.08~22.12. 경찰대학교
23.08~24.07 ODA (요르단 경찰청 PSD)

◎ 2022년 AI학습용데이터 구축사업 평가

1차 08 방송 콘텐츠 대화체 음성인식 데이터
09 방송 콘텐츠 한국어·영어 통번역 데이터
43 갑각류 종자생산 데이터
48 식생 탄소포집량 식별 데이터
2차 74 축산 기자재(소, 돼지) 3D 데이터
75 소(한우, 젃소) 및 돼지 발정행동 데이터
3차 06 인공지능 신기술 선도(자유 공모)

◎ 데이터 가치평가 컨설팅

23.09~23.11 중소벤처기업진흥공단

발급번호:00KH-183K-W6YQ-0A64-CG1Z

소프트웨어기술자 경력증명서

성명	홍용기		생년월일	1964.08.25			
현 근무처	회사명			사업자등록번호			
	전화번호			업종			
	소재지						
근무경력	확인여부	근무기간	회사명	담당업무	부서/직위		
기술자격	종목 및 등급		등록번호	취득일	발급기관		
	빅데이터분석기사		BAE-002000023	2021.07.16	한국데이터산업진흥원		
	ADSP(데이터 분석 준전문가)		ADSP-028000961	2021.04.09	한국데이터산업진흥원		
학력	학교명	학과(전공)	수학기간	학위			
교육	기간	과정	수료번호	교육기관			
상훈	수여일	종류	상훈기관	근거			
기술경력	확인여부	참여사업명	참여기간	발주자	소속사	직위	담당업무
	확인	국민 제강형 치안 안심 플랫폼 구축 정보화전략계획 사업	2022.08.04 ~ 2022.12.31	경찰대학교	(주)메타로 직권선택	프리랜서	IT컨설팅 > 정보기술기획
	확인	서울시 차세대 지방세 징수시스템 통합 구축 변화관리 컨설팅	2021.11.01 ~ 2022.06.30	서울특별시	(주)메타로 직권선택	프리랜서	IT컨설팅 > 정보기술컨설팅
	확인	소상공인지원사업 디지털전환 정보화전략계획(OSP) 용역	2021.01.04 ~ 2021.05.03	소상공인시장진흥공단	(주)메타로 직권선택	수석컨설턴트	IT컨설팅 > 정보기술컨설팅
	확인	한국연구재단 중장기 정보화전략계획(OSP) 수립	2020.10.05 ~ 2021.01.04	한국연구재단	(주)메타로 직권선택	수석컨설턴트	IT컨설팅 > 정보기술컨설팅
	확인	창업기업확인시스템 구축을 위한 정보화전략계획(OSP) 수립	2020.05.18 ~ 2020.08.17	창업진흥원	(주)메타로 직권선택	수석컨설턴트	IT컨설팅 > 정보기술컨설팅

「소프트웨어 진흥법」 제24조제3항 및 같은 법 시행규칙 제13조제3항에 따라 소프트웨어기술자의 경력 사항을 증명합니다.

2023년 01월 25일





홍용기 컨설팅학 박사

2018년 1월

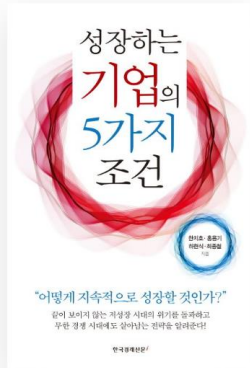
KYOB0 교보문고

통합검색 > 대입을 결정하는 초등 영어 공부법

대학교재 > 어린이 > 베스트 > 상품 > 이벤트 > 사운품 > PICKS > CASTing > 교보ONLY

국내도서 > 경제/경영 > 경영전략 > 경영전략일반

성장하는 기업의 5가지 조건



챗GPT AI로 데이터분석 마스터하기

초보자도 가능한 노코딩 AI 데이터분석



2023년 5월

AI와 데이터분석은 전문가들의 전유물이 아니다.
AI & Data Literacy for Everyone.

일반인을 위한
인공지능과 데이터 리터러시

챗GPT-4o(omni) 활용 데이터분석

컨설팅학박사 홍용기 저

일반인을 위한 AI와 데이터의 세계로 안내하는 입문서

인공지능과 데이터분석의 기본개념부터 응용까지

마케팅, 인사, 생산관리, 재무 등 실제 응용 가능한 7가지 분석사례

퍼플

2024년 6월
예정

2024년 4월

이 모든 것을 담은
All-In-One

생성형 AI 활용 전략 전자책

“놀랍도록 똑똑하고 충격적으로 어리석다”

최예진 교수

(미 워싱턴대 컴퓨터과학과 교수, 옥스퍼드 대학교 AI윤리 연구소 선임연구원)



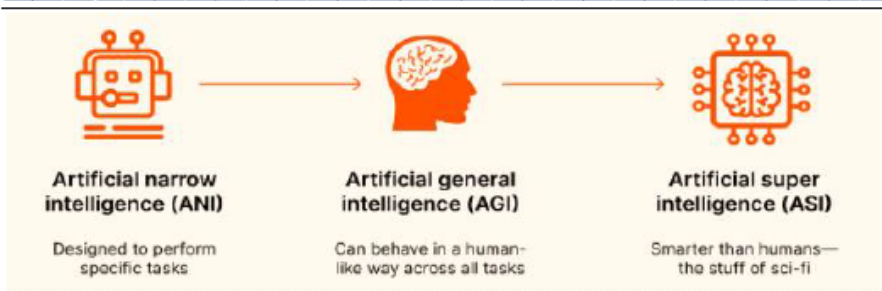
<https://v.daum.net/v/20240612211300429>

Road to AGI (범용인공지능, Artificial General Intelligence)

현재 널리 사용되는 AI는 '좁은 의미의 인공지능(ANI)'이며, AGI, ASI로 진화

- ANI, Artificial Narrow Intelligence, 좁은 의미의 인공지능 언어 번역, 이미지 인식, 게임 등 특정 영역에서만 인간을 능가하는 성능을 보유한 최적화된 인공지능 의미(유아 수준)
 - 얼굴 인식, 체스 게임, 언어 번역과 질병 진단 등 인간이 할 수 있는 특정 작업을 수행할 수 있지만 해당 영역을 넘어서 학습하거나 일반화할 수는 없음
 - ANI는 좁은 의미의 인공지능(때때로 약한 AI로 간주되지만 AI 연구의 의미에 실제로 적용)은 특정 작업을 수행하도록 설계된 AI 시스템으로 ChatGPT 및 Bard와 같은 기타 AI 챗봇은 여전히 좁은 의미의 AI
- AGI, Artificial General Intelligence, 범용인공지능 특정 분야에만 특화되어 있는 것이 아닌, 기본적으로 기본적인 이해 능력, 추론, 문제 해결, 창의적 사고 등을 갖춘 인간과 유사한 또는 높은 지능 수준을 가진 인공지능을 의미(성인 수준)
 - 경험을 통해 배우고, 추론하고, 이해하고, 소통하고, 일반 지식과 상식이 필요한 문제를 해결
 - AGI는 그것이 무엇을 의미하든 일반적으로 지능적(generally intelligent)이어야 함(G가 중요)
 - AGI, 즉 강력한 AI는 인간과 유사한 지능(또는 "일반적으로 인간보다 똑똑함")을 나타내는 AI
- ASI, Artificial Super Intelligence, 초인공지능 범용인공지능보다 한 단계 발전한 개념으로 모든 면에서 인간의 지능과 능력을 뛰어넘는 혹은 인간보다 뛰어난 지능 수준(신과 동등한 수준*)
 - AGI는 특정 조건에서만 쓰이는 현재 AI 기술과 달리 모든 상황에 두루 적용할 수 있는 차세대 AI 모델이며, 물질과 에너지 제어 등 인류와 우주의 미래에 심오한 영향을 미치게 됨
 - The Terminator의 Skynet 또는 Avengers: Endgame의 Jarvis와 같은 일부 가상 캐릭터
 - ASI는 인간의 지능을 훨씬 능가하는 AI 시스템으로, 지각 있는 공상과학 슈퍼컴퓨터

그림7 ANI vs AGI vs ASI



자료 : Zapier

멀티모달 AI 경쟁과 다가오는 AGI (정보통신기획평가원 IITP, 2024.05)

인간과 같은 수준의 AGI(범용인공지능)는 언제 나올 수 있는가

- 레이 커즈와일(Ray Kurzweil, 미래학자)은 '특이점이 온다(The Singularity is Near)'에서 '기술적 특이점' 즉, 싱귤래리티 도래를 주장(2005년)
 - 앞으로 30년 후인 2045년 경 인류는 인간보다 뛰어난 기계가 출현하는 특이점에 도달
 - 최근, 레이 커즈와일은 인간의 두뇌를 닮은 인공지능 AGI가 5년 후, 2029년까지 현실이 될 것이고 주장(SXSW 2024)

그림8 특이점(Singularity)



자료 : 언론 보도자료 정리 등

- 젠슨 황(Jensen Huang, 엔비디아 CEO)은 인간과 같은 수준의 범용인공지능(AGI)이 5년 이내 등장할 것이라고 전망(GTC 2024)
 - '인간이 치르는 모든 종류의 시험을 통과할 수 있는 AI'를 AGI라고 전제
- 일론 머스크(Elon Musk, 테슬라 CEO)는 AGI를 가장 똑똑한 인간보다 똑똑한 AI로 정의한다면 내년 혹은 후년에 가능하고, 2029년에는 모든 인간지능을 합친 것보다 더 똑똑해질 것임(2024년, 노르웨이 국부펀드 니콜라이 탕겐 CEO 인터뷰)
 - 머스크는 데이터센터 장비와 전력망 공급이 관건이라고 지적, 작년의 칩 제약에서 올해는 변압기 공급정지로 이동
- 샘 알트먼(Sam Altman, OpenAI CEO), AGI는 인류가 발명한 가장 강력한 기술이 될 것이라고 생각(2023년)하며, 오픈AI 역시 '인간과 같은 추론 능력'의 AGI 개발이 핵심
 - 향후 4~5년 이내에 AGI가 구축될 것으로 전망
 - AGI는 매우 다른 세상이고 그것은 공상과학이 오랫동안 우리에게 약속해 온 세상이고, 처음으로 그것이 어떤 모습일지 볼 수 있게 된 것 같음
 - 오픈AI는 '23년 말 '인간처럼 생각하는' 능력으로 정답이 확실한 수학 문제 등을 풀어내는 'Q*(큐스타)'라는 모델을 개발, AGI 개발로 가는 돌파구를 찾는 것으로 알려졌다
- 데미스 하사비스(Demis Hassabis, 딥마인드 CEO)는 AGI가 10년 내 실현하더라도 놀라지 않을 거다. 하지만 꼭 그렇게 될 거라는 말은 아니다. AGI가 10년내 실현될 확률은 50%이며, 이 시간표는 딥마인드 설립 이래 결코 변하지 않았다고 말함(파이낸셜타임즈 인터뷰)
 - 이를 위해 '성능과 규모가 기하급수적으로 증가하는 문제 해결'과 '발열과 전기' 해결이 선행조건



8 NEWS



"방금 왜 사과를 준 거야?"
신기함 넘어 섬뜩한 대답



Like every great presentation, I've divided my talk into three subjects. Steve Jobs -

I .

**Data Analysis
Basic Theory**

II .

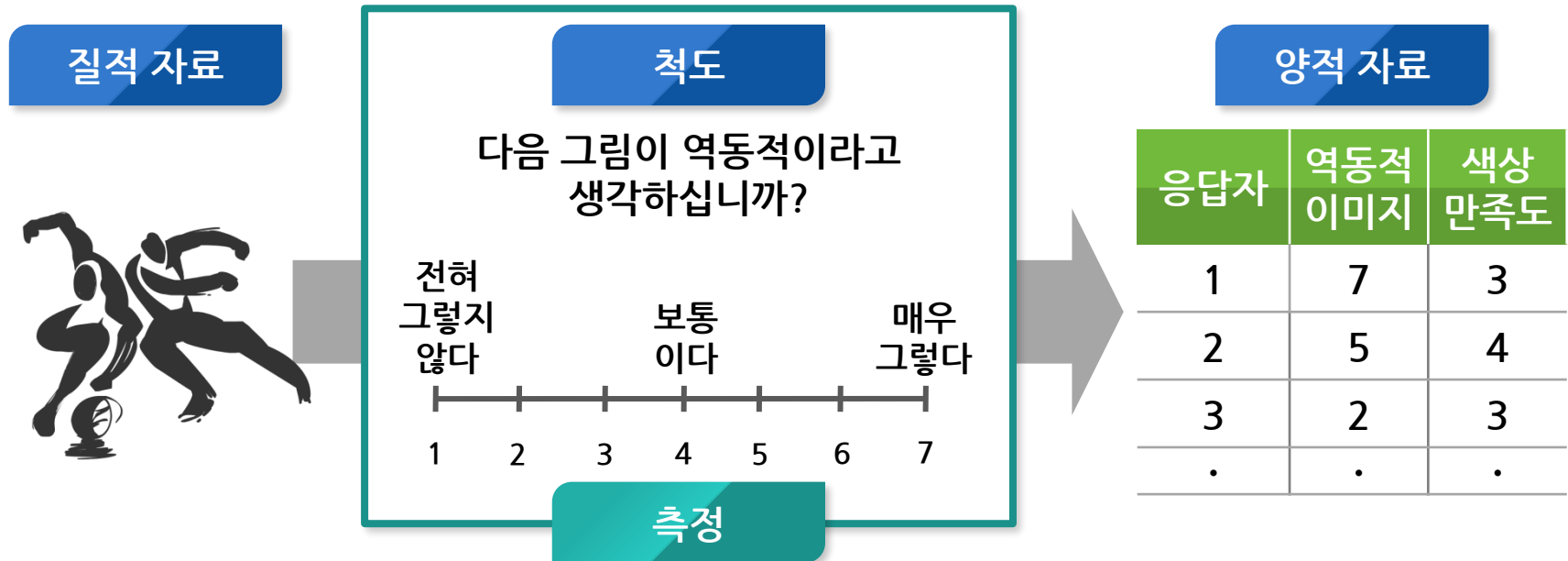
**Classification
Analysis**

III .

**Regression
Analysis**

• 척도의 개념

- 특정 속성을 측정하여 그 정도를 숫자로 나타내는 규칙
- 질적 자료를 양적 자료로 전환시켜 주는 도구



- 척도의 종류
- 어떤 척도를 사용하는지에 따라 측정된 숫자에 내재된 정보량이 달라지며, 적용 가능한 통계분석 기법이 달라짐

척도의 종류	내용
명목척도	응답보기들을 상호 배타적으로 구분 하기 위해 임의의 숫자를 부여하는 척도
서열척도	응답보기들을 구분 하고, 구분한 응답보기들의 순서 까지 측정하는 척도
등간척도	서열 척도에 포함된 정보(분류, 서열정보)외에 거리(간격)정보까지 가지는 척도
비율척도	절대 영점을 가지고 있어서 속성의 상대적 크기 뿐만 아니라, 절대적 크기의 비교도 가능한 척도

1. 명목 척도(Nominal scale)

- 응답보기들을 상호 **배타적으로 구분**하기 위해 임의의 숫자를 부여하는 척도
- 선택한 응답을 기준으로 응답자들을 특정 집단으로 분류하기 위해 사용(=**분류정보**)



귀하는 다음 중 어떤 훈련과정에 입학을 원하십니까?

1) A과정 2) B과정 3) C 과정 4) D 과정 5) 기타

- 숫자는 '크기'의 의미가 없는 명칭에 해당하기 때문에 **사칙연산은 무의미함**
- 대표치는 **최빈치(Mode)** : 응답보기 중 가장 많이 선택된 응답보기의 선택된 수
- 4가지 척도 중 정보량이 가장 적은 척도 : **분류 정보만 보유**

2. 서열 척도(Ordinal scale)

- 응답보기들을 **구분**하고, 구분한 응답보기들의 **순서**까지 측정하는 척도
- 응답보기들의 속성을 서열로 나타내는 척도(=**분류정보** + **순서정보**)
- 응답보기 간의 **간격은 측정하지 않고** 순서만 측정함
 - 응답 보기들 간의 순위만 나타낼 뿐, 얼마나 더 선호되는지는 측정이 불가능함



다음 교육과정 중 귀사에서 가장 중요하다고 생각하는 대로 순서를 기입해 주십시오.
A과정 (), B 과정 (), C 과정 (), D 과정 ()

- 사칙연산은 무의미
 - 순위 간 간격이 서로 달라 숫자 차이에 절대적 의미가 없기 때문
 - 1, 2순위의 차이보다 3, 6순위의 차이가 3배 크다고 할 수 없음
- 대표치로서 중앙값(Median)를 사용함
- 명목 척도 다음으로 적은 정보를 보유함 : 분류 정보 + 순서 정보

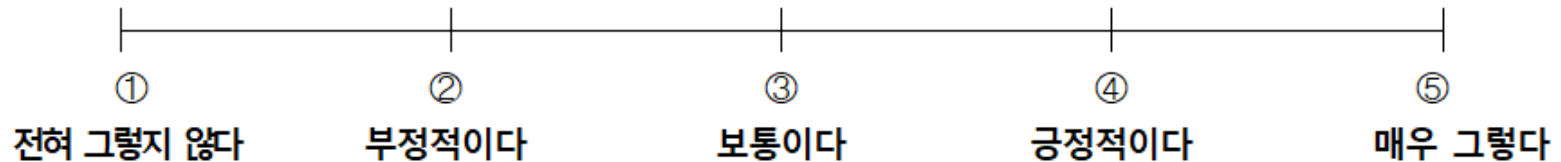
3. 등간 척도(Interval scale)

- 서열 척도에 포함된 정보(분류, 서열정보)외에 거리(간격)정보까지 가지는 척도
- **간격이 동일한 서열척도**
- 속성의 **상대적 크기**를 측정하기 위해 균일한 간격으로 분할한 길이를 이용하여 측정
 - 예) 온도계, IQ 등
 - 온도계로 측정한 1도와 2도 간의 차이는 2도와 3도 간의 차이와 동일함
- 간격 척도의 **숫자 자체**는 절대적 의미를 가지지 않음
- 절대 영점이 없기 때문에 **숫자 간 비율개념 없음**
- 간격 척도에서 **숫자 간의 차이**는 절대적 의미를 가짐(**차이 값 간 비율개념있음**)
- 대표치로서 산술평균을 사용
- 정보량 : 분류 정보 + 순서 정보 + 상대적 크기 정보

3. 등간 척도(Interval scale)

예

지난 6개월간 참여하신 교육과정이 취업역량 확보에 도움이 되셨습니까?



- 5점 응답자와 3점 응답자의 만족도 차이가 5점과 4점 응답자의 만족도 차이 보다 2배 크다고 할 수 없음(응답보기(척도점) 간 간격이 동일하다고 볼 수 없기 때문)
- 따라서 **간격 척도라기 보다 서열척도에 가까움**
- 하지만, 사회과학연구의 특성을 고려하여 척도점 간 간격이 동일하고, 각 척도점의 의미를 응답자들이 동일하게 이해한다는 전제 하에 간격 척도로 인정함

4. 비율 척도(Ratio scale)

- 절대 영점을 가지고 있어서 속성의 상대적 크기 뿐만 아니라, 절대적 크기의 비교도 가능한 척도

예

나이 ()세, 근무기간()년, 연봉()원

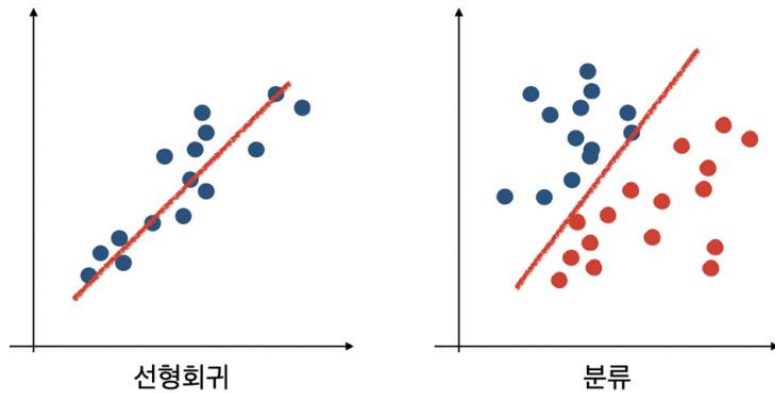
예

하루에 보통 몇시간 정도 훈련이 가능하십니까? ()시간

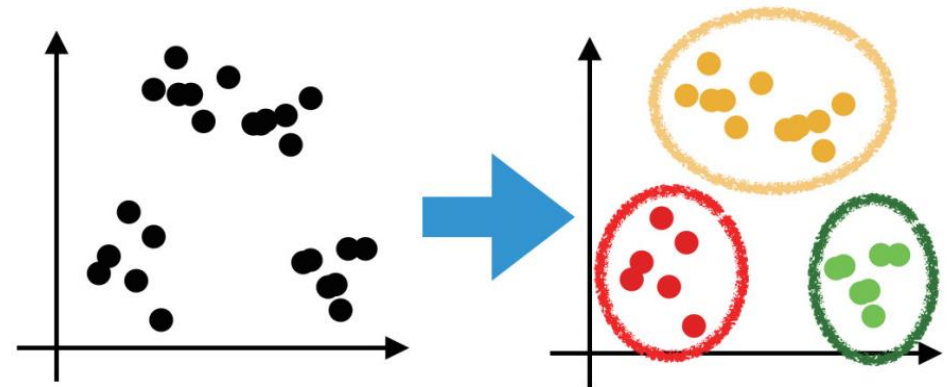
- 만족도, 선호도, 인지도 등 절대 영점이 존재하기 어려운 소비자의 사고나 인지수준에 대한 측정은 한계가 있음
- 직접 관찰할 수 있는 물리적 사건이나 현상을 측정하는데 주로 사용함
- 사칙연산이 가능하며, 대표치는 평균값
- 4가지 척도 중 가장 정보량이 많은 척도 : 분류정보 + 순서정보 + 상대적 크기 정보 + 절대적 크기 정보

Supervised learning VS Unsupervised learning

지도 학습 supervised learning



비지도 학습 unsupervised learning



※출처: 으뜸 데이터 분석과 머신러닝(박동규·강영민, 2021)

Data split

Training data

[illegible]

Test data

The diagram illustrates the relationship between different data sets in a machine learning workflow. It features a horizontal bar divided into two main sections. The left section, labeled 'Test data' in dark blue, contains a grid of 12 light blue squares. The right section, labeled 'y_predict' in red, contains a single light blue square. Below the 'Test data' section, the label 'X_test' is placed over a 2x2 grid of squares. To the right of the 'y_predict' section, the label 'y_test' is placed over a single square. The overall layout suggests that 'Test data' is used to generate 'X_test' and 'y_test', which are then used to produce 'y_predict'.

Training data

[illegible]

Validation data

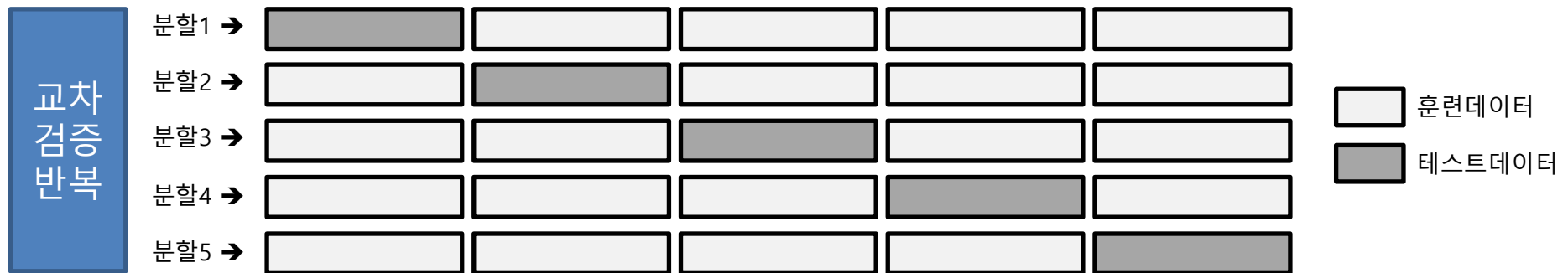
The diagram illustrates the split of validation data. A large rectangle labeled "Validation data" is divided into two parts. The left part is labeled "X_val" and the right part is labeled "y_val". The "y_val" label is in red, matching the "y_predict" label in the diagram above.

Test data

A diagram illustrating the relationship between X_{test} and y_{test} . X_{test} is represented by a large blue rectangle on the left, and y_{test} is represented by a smaller blue rectangle on the right. A white arrow points from X_{test} to y_{test} .

교차 검증의 의의

- 교차 검증은 일반화 성능을 재기 위해 훈련 세트와 검증 세트로 한 번 나누는 것보다 더 안정적이고 뛰어난 통계적 평가 방법
- 교차 검증에서는 데이터를 여러 번 반복해서 나누고 여러 모델을 학습
- 가장 널리 사용되는 교차 검증 방법은 k-겹 교차 검증(k-fold CV)으로 보통 5 또는 10을 사용
- 5-겹 교차 검증을 하면 데이터를 비슷한 크기의 부분 집합(5개의 폴드)으로 나누고, 일련의 모델을 만들어 훈련과 테스트를 반복

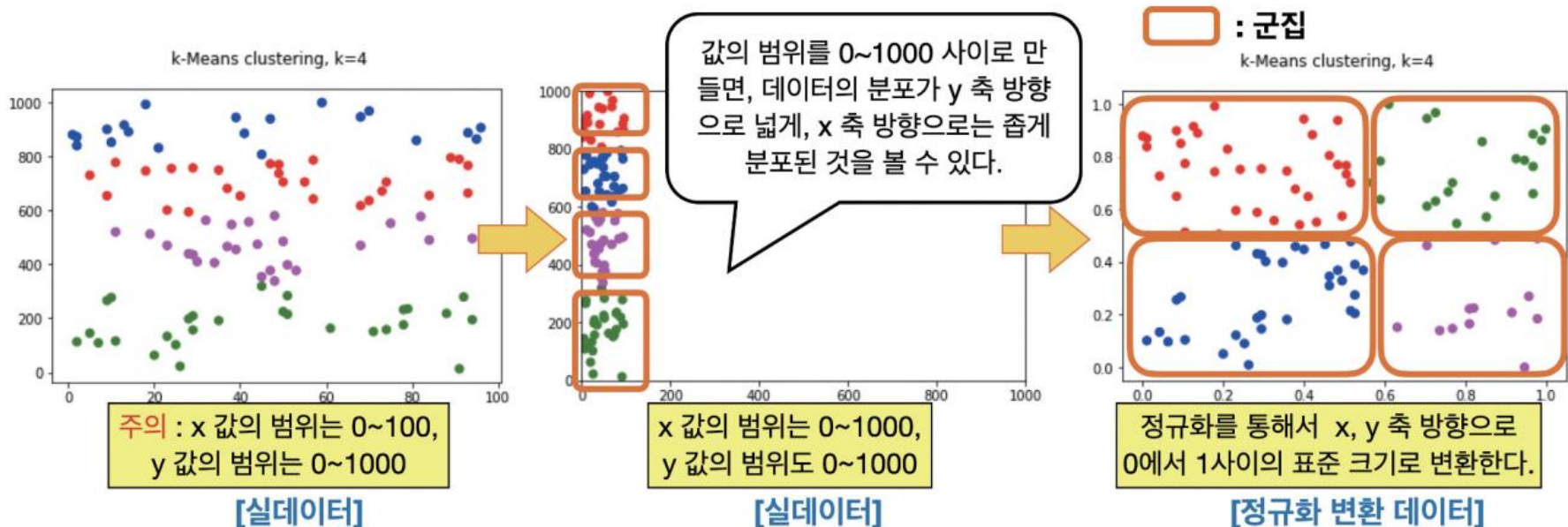


※출처: 파이썬 라이브러리를 활용한 머신러닝(번역개정판, 안드레아스뮐러 & 세라가이드, 2021.7)

Missing value : NA, NAN, Null

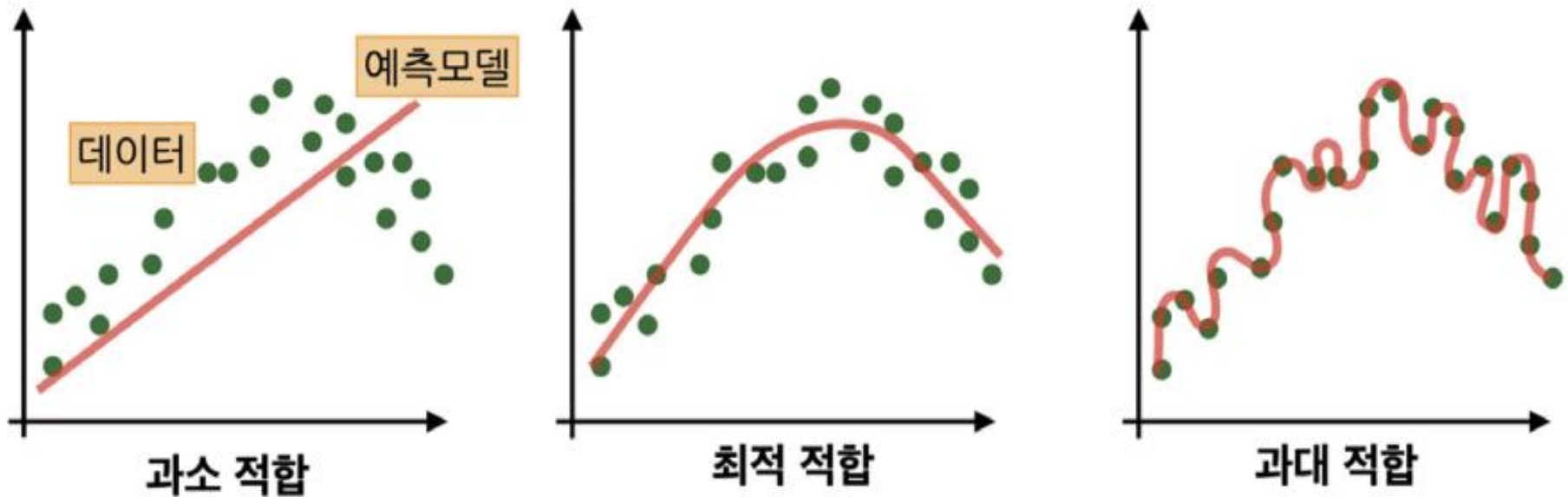
➔ 결측치 제거 또는 대체 (평균, 중위수, 최빈값)

정규화normalization(min-max scale), **표준화**standardization(평균=0, 분산=1로 만듦)



※출처: 으뜸 데이터 분석과 머신러닝(박동규·강영민, 2021)

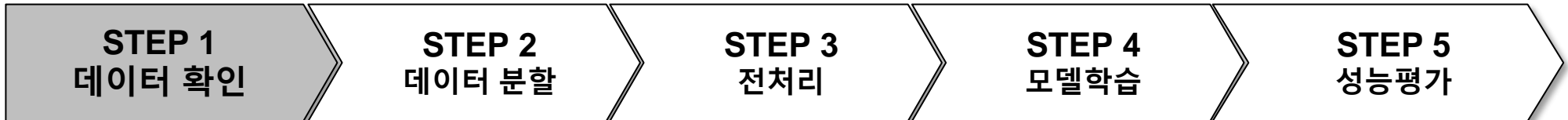
Underfitting, Optimal fitting, Overfitting



※출처: 으뜸 데이터 분석과 머신러닝(박동규·강영민, 2021)

1단계 : 데이터 확인

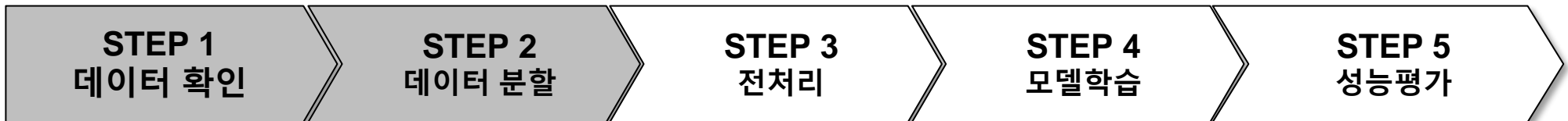
- 분석할 데이터의 특성을 확인하는 단계
- 변수의 특성(독립변수/입력변수)과 타겟(종속변수/반응변수)의 존재 여부 파악
- 적용가능한 분석모델 확인(ex. 타겟 연속된 수치형이라면 회귀분석, 범주형이라면 분류분석)
- 타겟이 없는 데이터라면 비지도학습 적용



- 독립변수, 종속변수 확인
- 연속형 vs 범주형 확인
- 범주형 독립변수 여부확인
- 적용가능한 분석모델 확인
(회귀, 분류, 비지도 학습)

2단계 : 데이터 분할

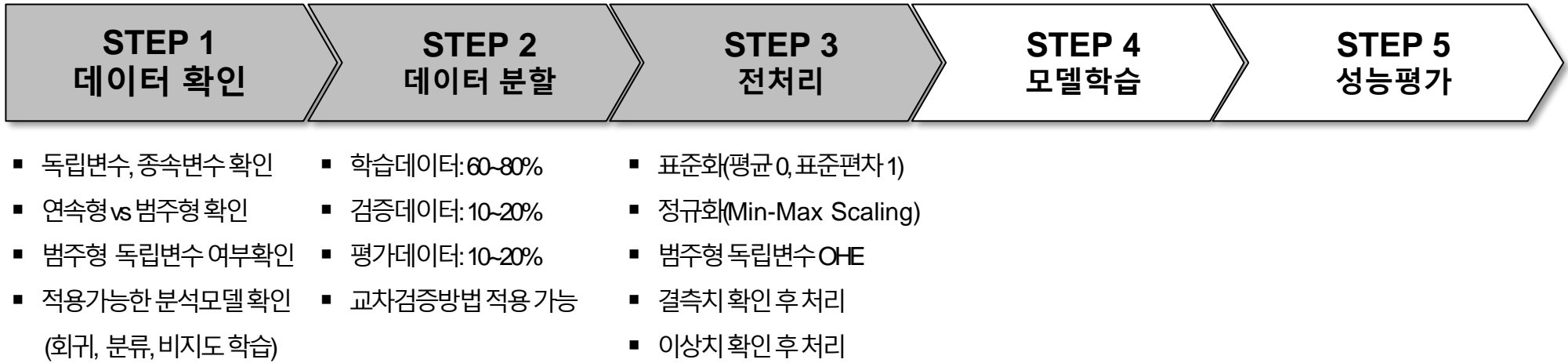
- 학습용 데이터와 평가용 데이터를 분할하는 단계
- 데이터는 학습데이터(60~80%), 검증데이터(10~20%), 평가데이터(10~20%)로 분할
- 예측을 수행하는 데이터 세트는 학습용 데이터 세트가 아니라 평가 전용 데이터세트여야 함
- 단순 학습데이터 + 복잡한 평가데이터의 경우 평가데이터의 특징을 반영하지 못할 수 있음
- 데이터 크기가 작은 경우나, 검증 결과를 일반화하기 위해 교차검증방법을 적용



- 독립변수, 종속변수 확인
- 연속형 vs 범주형 확인
- 범주형 독립변수 여부확인
- 적용가능한 분석모델 확인
(회귀, 분류, 비지도 학습)
- 학습데이터: 60~80%
- 검증데이터: 10~20%
- 평가데이터: 10~20%
- 교차검증방법 적용 가능

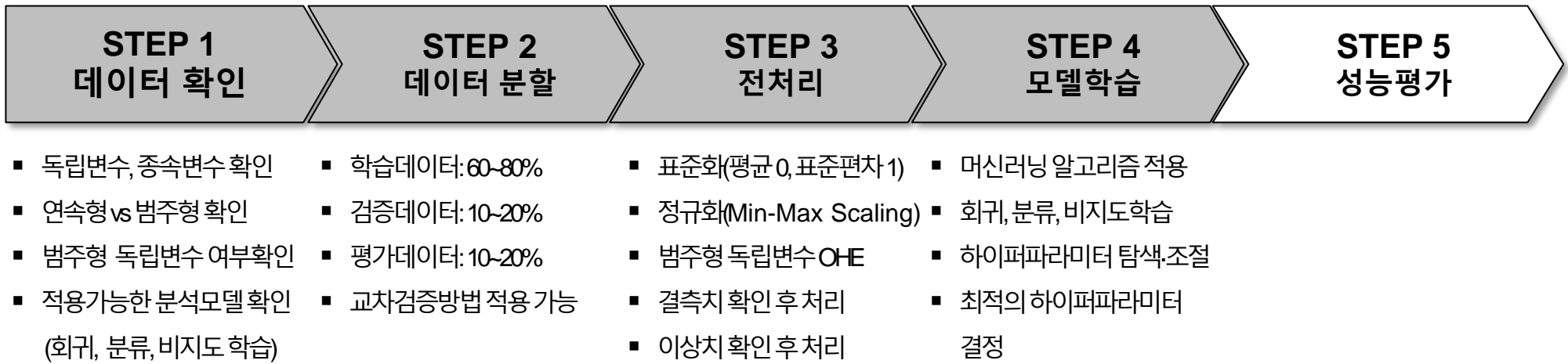
3단계 : 전처리

- 데이터의 특성에 따라 분석이 가능한 형태로 변형하는 단계
- 독립변수에 범주형 변수가 있을 경우 데이터 분할 전 One-hot Encoding으로 데이터를 변형
- 변수마다 단위 특성에 차이가 클 때 분석결과에 영향을 줄 수 있으므로, 정규화나 표준화 실시
- 결측치와 이상치는 분석가의 판단과 도메인 상황에 따라 적절한 방법으로 처리



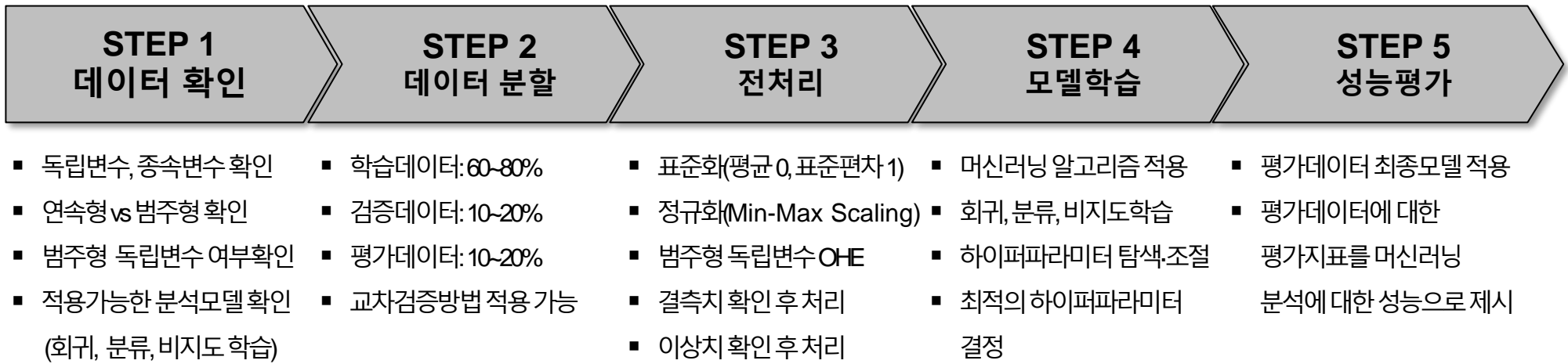
4단계 : 모델학습

- 머신러닝 알고리즘을 학습데이터에 적용하는 단계
- 1단계에서 파악한 분석방법에 따라 적합한 라이브러리를 사용해 머신러닝 알고리즘을 적용
- 머신러닝 분석방법은 지도학습과 비지도학습으로 구분되며, 지도학습은 회귀와 분류로 나뉨
- 학습데이터로 학습을 수행, 검증데이터로 학습결과 확인 후 하이퍼파라미터 탐색 및 조절



5단계 : 성능평가

- 최적의 하이퍼파라미터 및 최종모델 결정 단계
- 최종모델에 평가데이터를 적용하여 머신러닝 알고리즘의 예측성능을 평가
- 평가데이터는 반드시 학습 과정이나 검증 과정에서 사용되지 않은 데이터로 사용해야 함
- 평가데이터에 대한 평가지표를 머신러닝 분석에 대한 최종성능으로 제시





Like every great presentation, I've divided my talk into three subjects. Steve Jobs -

I .

Data Analysis
Basic Theory

II .

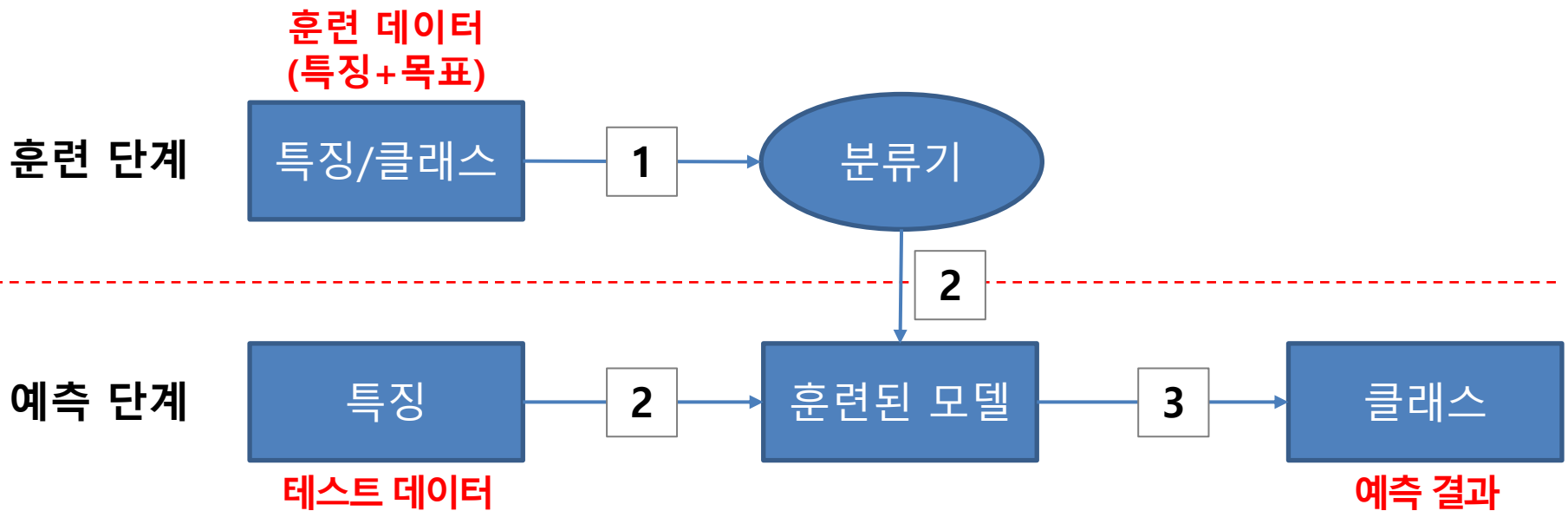
Classification
Analysis

III .

Regression
Analysis

분류의 개념

- 분류(classification)는 지도학습의 하나로 관측값과 해당 관측값에 대한 범주형 출력을 포함하는 훈련데이터셋이 주어졌을 때 관측값을 목표 범주에 올바르게 매핑하는 규칙을 학습하는 것
- 관측값(observation)은 특징(feature) 또는 예측변수라고도 하며, 목표 범주(category)는 레이블(label), 클래스(class) 또는 타겟(target)이라고도 한다.



분류의 종류와 클래스

- 일반적으로 분류는 두개의 클래스로 분류하는 이진 분류(binary classification)와 셋 이상의 클래스로 분류하는 다중 분류(multiclass classification)로 나눌 수 있음
- 이진 분류에서 한 클래스를 양성(positive) 클래스, 다른 하나를 음성(negative) 클래스라 함
- 양성 클래스라고 해서 좋은 값이나 장점을 나타내는 것이 아니고 학습하고자 하는 대상을 의미
- 일반 메일에서 스팸 메일을 골라내는 분석의 경우 스팸메일이 양성 클래스가 되고, 양성 종양과 악성 종양을 분별하는 분석에서는 악성 종양이 양성 클래스가 됨

[일반화, 과대적합, 과소적합]

- ✓ 지도학습에서는 훈련데이터로 학습한 모델이 훈련데이터와 특성이 같다면 새로운 데이터가 주어져도 정확히 예측할 거라 기대
- ✓ 모델이 처음 보는 데이터에 대해 정확하게 예측할 수 있으면 이를 "훈련세트에서 데이터 세트로 일반화" 되었다고 함
- ✓ 과대적합은 모델이 훈련세트의 각 데이터에 너무 맞춰져서 새로운 데이터에 일반화되기 어려움
- ✓ 과소적합은 모델이 너무 간단하여 데이터의 면면과 다양성을 잡아내지 못하고 훈련세트에도 잘 맞지 않음

단순회귀분석 (Simple Linear Regression)

- **하나의 독립변수와 하나의 종속변수 간의** 선형적인 관계를 모델링하는 회귀분석 방법
- 독립변수와 종속변수 간의 관계를 파악하고, 독립변수의 값을 통해 종속변수 값을 예측하거나 설명

다중회귀분석(Multiple Linear Regression)

- **둘 이상의 독립변수와 하나의 종속변수 간의** 선형적인 관계를 모델링하는 회귀분석 방법

로지스틱 회귀분석(Logistic Regression)

- **종속변수가 이항형(binary)**일 때, 독립변수와 종속변수 간의 선형적인 관계를 모델링하는 회귀분석 방법
- 로지스틱 회귀분석은 이진 분류(binary classification, 0과 1로 분류)에 널리 사용되며, 예측하려는 결과가 두 가지 중 하나인 경우에 사용

- 결과변수는 범주형 범주로서 1(사건발생), 0(사건 미 발생)의 값을 갖기 때문에 결과변수의 기대 값은 항상 0과 1 사이의 값을 가짐
- 결과변수 예측 값은 사건이 발생할 확률을 나타냄
 - 특정 고객의 카드연체 가능성 예측 모델 결과값이 0.78이라면, 카드를 연체할 확률이 78%라는 의미임
- 분류 문제에서는 0.5를 기준으로 0과 1을 분류함

예1 카드회사에서 신규 카드 발급 시 고객정보를 기반으로 연체 가능성 예측

예2 환자의 연령, 성별, 혈액 검사 결과를 기반으로 암 진단

Algorithms → 분류와 회귀 모두 가능한 알고리즘이 많이 있음

● 의사결정 나무 (Decision Tree)

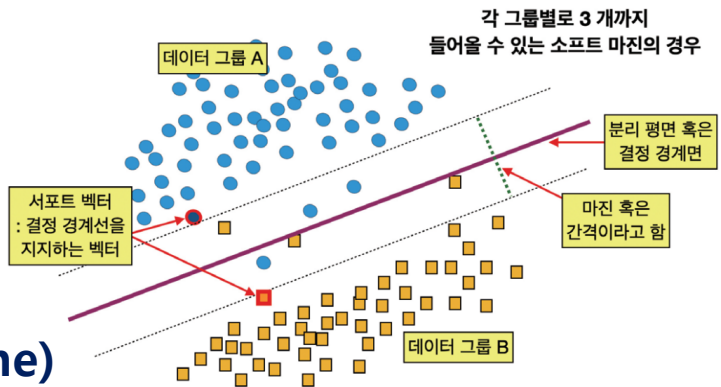
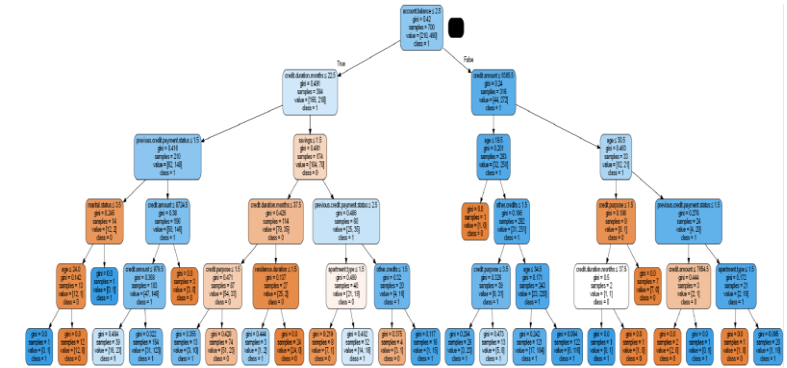
● 앙상블 모형 (Ensemble)

1. Bagging

2. Boosting

- AdaBoost (Adaptive Boosting)
- GBM (Gradient Boosting Machine)
- XGBoost
- LightGBM
- CatBoost

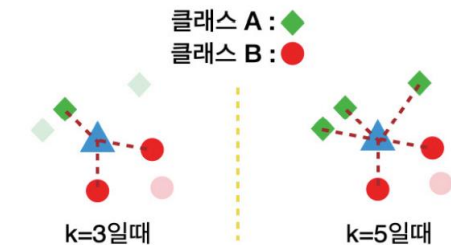
3. Random Forest



● 서포트 벡터 머신 (SVM; Support Vector Machine)

● K 최근접 이웃 (K-Nearest Neighbor)

● 소프트맥스 (Softmax) 회귀 → 다항 로지스틱 회귀라고도 함



		예측(prediction)	
		양성(Positive)	음성(Negative)
		Positive	Negative

Confusion matrix

		예측(prediction)	
		양성(Positive)	음성(Negative)
실제 (real)	양성	Positive	Negative
	음성	Positive	Negative

		예측(prediction)	
		양성(Positive)	음성(Negative)
실제 (real)	양성	True Positive	False Negative
	음성	False Positive	True Negative

Confusion matrix

		예측(prediction)	
		양성(Positive)	음성(Negative)
실제 (real)	양성	TP (True Positive)	FN (False Negative)
	음성	FP (False Positive)	TN (True Negative)

		예측(prediction)	
		양성(Positive)	음성(Negative)
실제 (real)	양성	TP	FN
	음성	FP	TN

- 정확도(Accuracy) = (제대로 예측)/(전체) = $(TP+TN)/(TP+FN+FP+TN)$
- 정밀도(Precision) = (실제 양성)/(양성으로 예측) = $TP/(TP+FP)$
- 재현률(Recall) = (양성으로 예측)/(실제 양성) = $TP/(TP+FN)$ = 민감도(Sensitivity)
- 특이도(Specificity) = (음성으로 예측)/(실제 음성) = $TN/(TN+FP)$
- 거짓양성율(FPR) = $1 - \text{특이도}$
- F1 score = $2 \times \text{정밀도} \times \text{재현률} / (\text{정밀도} + \text{재현률})$

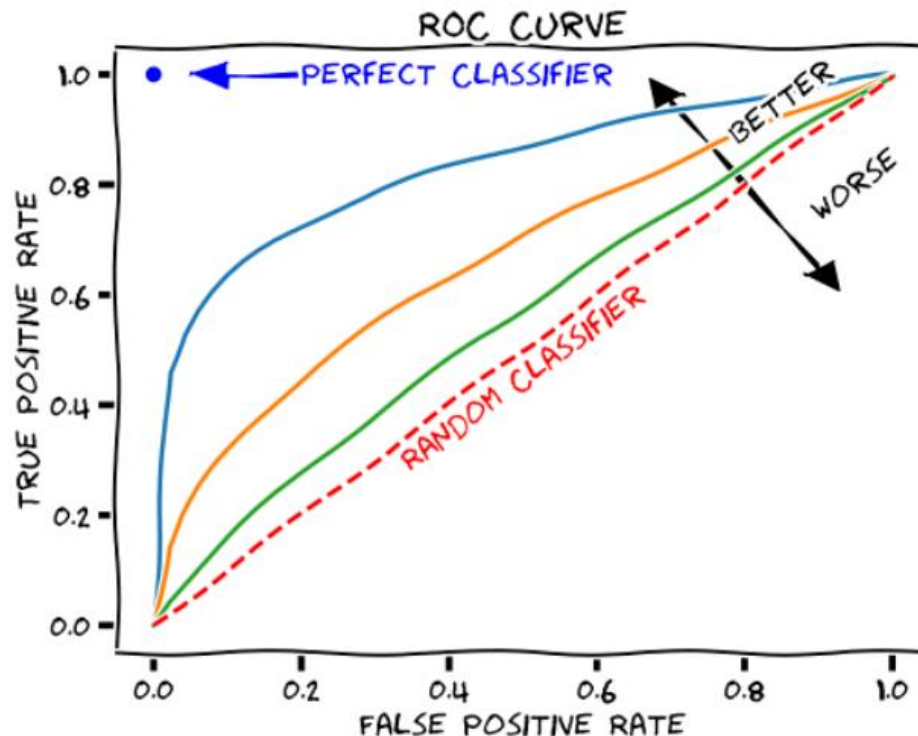
분류 모델의 성능 지표 : ROC 곡선, AUC

- ROC(Receiver Operating Characteristic) 곡선

분류 모델의 임계값(threshold)을 변화시켰을 때, 모델의 TPR(True Positive Rate)과 FPR(False Positive Rate)이 어떻게 변화하는지를 나타내는 그래프

- AUC(Area Under the Curve)

ROC 커브 아래 면적을 나타내는 지표로 1에 가까울 수록 성능이 우수한 것으로 판단



Classification analysis practice using ChatGPT

Data : heart_disease.csv (미국 심장질환 데이터셋)

✓ 303행, 14개의 변수 (Target = target)

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 303 entries, 0 to 302
```

```
Data columns (total 14 columns):
```

```
#   Column   Non-Null Count  Dtype
```

---	-----	-----	-----	
0	age	303 non-null	int64	
1	sex	303 non-null	int64	성별 (1=남성, 0=여성)
2	cp	303 non-null	int64	가슴 통증 (1=안정형 협심증, 2=불안정형 협심증, 3=협심증 이외 통증, 4=무증상)
3	trestbps	303 non-null	int64	휴식 시 혈압
4	chol	303 non-null	int64	콜레스테롤 수치
5	fbs	303 non-null	int64	공복 혈당
6	restecg	303 non-null	int64	휴식 상태의 심전도
7	thalach	303 non-null	int64	최대 심장 박동수
8	exang	303 non-null	int64	운동 유발 협심증
9	oldpeak	303 non-null	float64	운동에 의한 상대적 휴식 시 ST 하강
10	slope	303 non-null	int64	최대 운동 ST 세그먼트의 기울기
11	ca	303 non-null	int64	형광 투시로 착색된 주요 혈관 수
12	thal	303 non-null	int64	탈라세미아 유형
13	target	303 non-null	int64	심장질환의 존재 여부 (1=예, 0=아니오)

```
dtypes: float64(1), int64(13)
```

```
memory usage: 33.3 KB
```




Like every great presentation, I've divided my talk into three subjects. Steve Jobs -

I .

Data Analysis
Basic Theory

II .

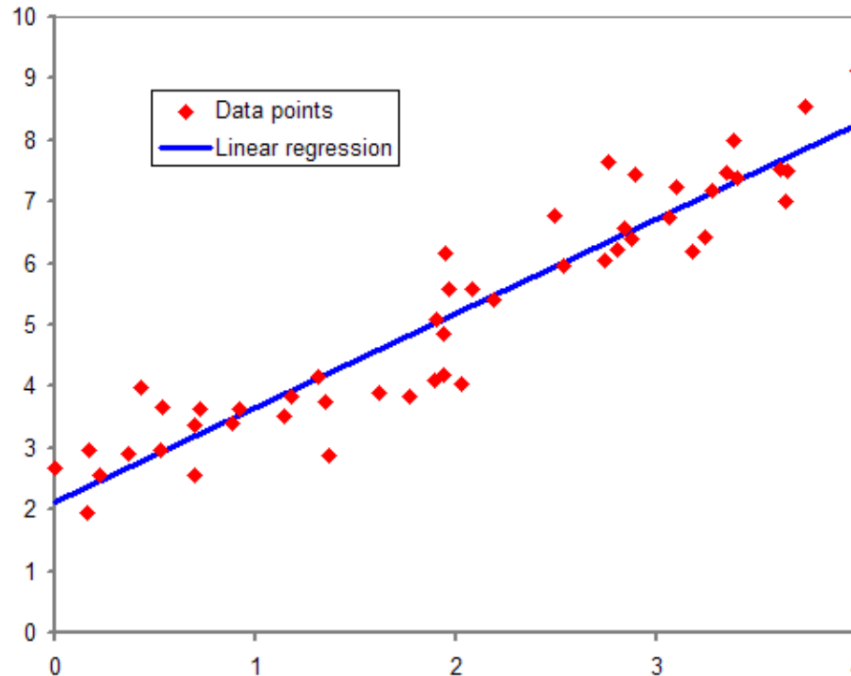
Classification
Analysis

III .

Regression
Analysis

위키백과 : '회귀분석'

- 회귀(regress)의 원래 의미는 옛날 상태로 돌아가는 것을 의미. 영국의 유전학자 프랜시스 골턴은 부모의 키와 아이들의 키 사이의 연관관계를 연구하면서 부모와 자녀의 키 사이에는 선형적인 관계가 있고, 키가 커지거나 작아지는 것보다는 전체 키 평균으로 돌아가려는 경향이 있다는 가설을 세웠으며, 이를 분석하는 방법을 '회귀분석'이라고 함
- 이후 칼 피어슨은 아버지와 아들의 키를 조사한 결과를 바탕으로 함수 관계를 도출하여 회귀분석 이론을 수학적으로 정립



Regression line for 50 random points in a [en:Gaussian distribution](#) around the line $y=1.5x+2$ (not shown). The regression line (shown) that best fits these points is actually $y=1.533858x+2.129333$.

- 단순회귀분석 목적

1

하나의 변수(독립변수, 예측변수)를 이용해서 다른 변수(종속변수, 결과변수)를
예측함

예

영업사원의 수나 판촉행사 횟수, 매장의 면적 등 어떤 특정한 하나의 변수를 이용해서 매출액을 예측함

2

하나의 변수(독립변수, 설명변수)를 이용해서 다른 변수(종속변수, 결과변수)를
설명함

예

가격만족도, 품질만족도 등 어떤 특정한 하나의 변수를 이용해서 전반적인 만족도를 설명함

- 단순회귀분석 회귀식

$$Y = \beta_0 + \beta_1 \cdot X$$

Y : 종속변수 X : 독립변수 β_1 : 회귀계수 β_0 : 상수

예 우리회사 내년도 매출액 규모(Y)를 영업사원 수(X)로 예측

➔ 매출액 = $\beta_0 + \beta_1 \cdot (\text{영업사원 수})$

- 다중회귀분석 목적

1

2개 이상 변수(독립변수, 예측변수)를 이용해서 다른 변수(종속변수, 결과변수)를 **예측**함

예

영업사원의 수, 판촉행사 횟수, 매장의 면적 등 3가지 변수를 이용해서 매출액을 예측함

2

2개 이상 변수(독립변수, 예측변수)를 이용해서 다른 변수(종속변수, 결과변수)를 **설명**함

예

가격만족도, 품질만족도, 디자인만족도, 무게만족도 등 4가지 변수를 이용해서 전반적인 만족도를 설명함

Multiple linear regression

- 다중회귀분석 회귀식

$$Y = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \beta_3 \cdot X_3 + \dots + \beta_i \cdot X_i$$

Y : 종속변수

β_i : X_i 의 회귀계수

X_i : 독립변수

β_0 : 상수

예

우리회사 내년도 매출액 규모(Y)를 '영업사원 수(X_1), 프로모션 횟수(X_2), 광고비 규모(X_3)'를 이용해 예측하는 다중 회귀식

➡ 매출액 = $\beta_0 + \beta_1$ (영업사원 수) + β_2 (프로모션 횟수) + β_3 (광고비)

Regression

OLS Regression Results						
Dep. Variable:	불량률	R-squared:	0.281			
Model:	OLS	Adj. R-squared:	0.228			
Method:	Least Squares	F-statistic:	5.273			
Date:	Tue, 12 Sep 2023	Prob (F-statistic):	0.00117			
Time:	20:12:00	Log-Likelihood:	-102.95			
No. Observations:	59	AIC:	215.9			
Df Residuals:	54	BIC:	226.3			
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	28.1944	6.858	4.111	0.000	14.446	41.943
함수율	-0.4248	0.118	-3.587	0.001	-0.662	-0.187
온도	-0.0912	0.038	-2.430	0.018	-0.166	-0.016
습도	0.0053	0.019	0.280	0.781	-0.033	0.043
미세먼지	-0.0145	0.014	-1.032	0.307	-0.043	0.014
Omnibus:	3.587	Durbin-Watson:	1.887			
Prob(Omnibus):	0.166	Jarque-Bera (JB):	2.774			
Skew:	0.512	Prob(JB):	0.25			
Kurtosis:	3.284	Cond. No.	3.09E+03			

회귀식의 설명력 R^2

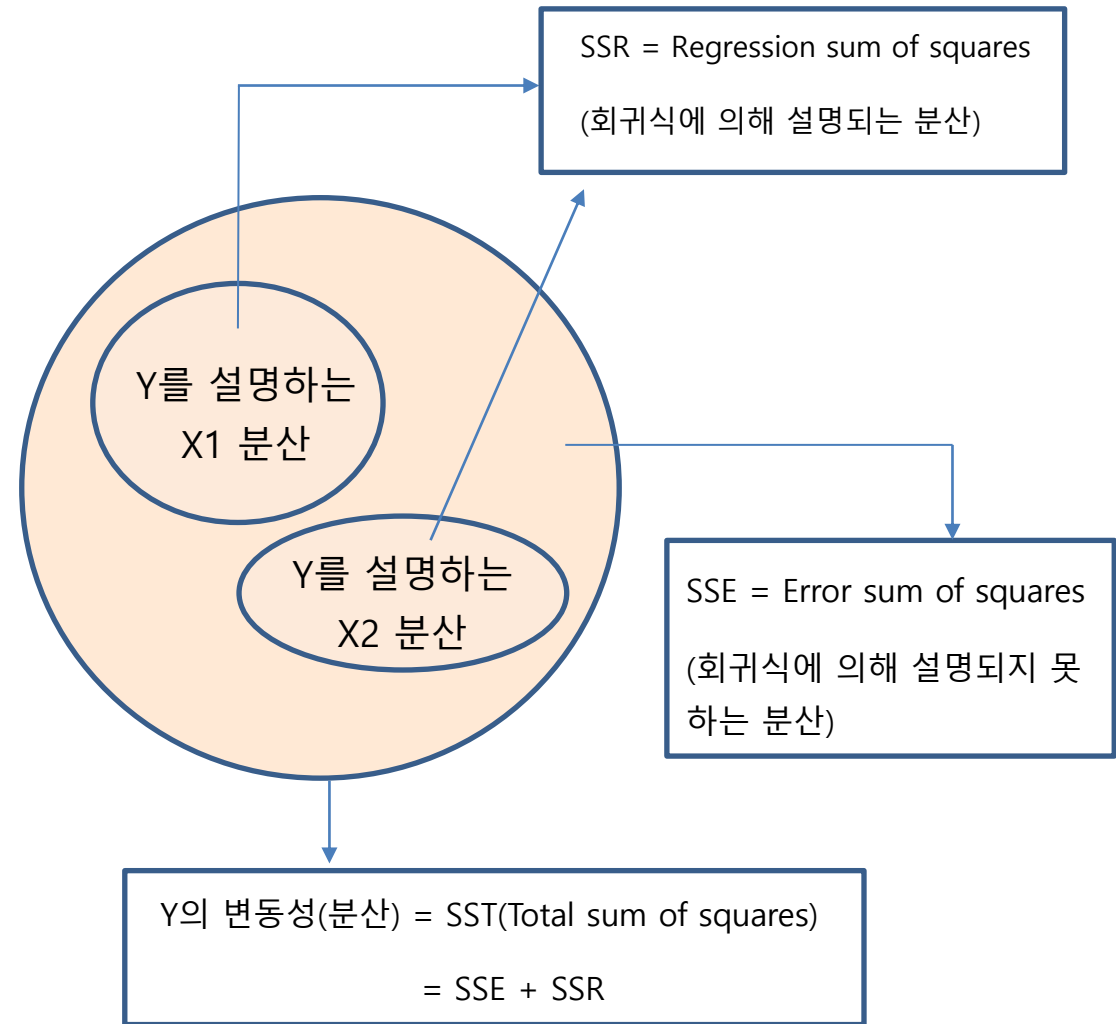
- 회귀식이 종속변수를 설명하고 예측하는데 유용한가를 판단
- 판단지표 : $R^2 = (\text{결정계수, 기여율, 설명력}), 0 < R^2 < 1$
- R^2 은 **종속변수의 분산** 중 독립변수에 의해 **설명되는 비율**을 의미

예 $R^2=0.76$ 이라는 것은 종속변수가 가지는 정보 중에서 76%를 독립변수가 설명할 수 있다는 의미

회귀식의 설명력 R^2

$$R^2 = \frac{SSR}{SST}$$

- 그러나, 변수의 수가 증가하면 SSR이 증가하면서 R^2 도 증가하는 하는 문제가 있음
- R^2 에 변수의 수 만큼 penalty를 주는 지표인 *adjusted* R^2 를 주로 활용



회귀 분석 결과 예시

모형		비표준화 계수		표준화 계수	t	유의 확률	공선성 통계량	
		B	표준오차	베타			공차	VIF
1	(상수)	-.631	.519		-1.215	.235		
	가격만족도	.744	.114	.668	6.528	.000	.298	3.356
	구매 횟수	.302	.094	.331	3.223	.003	.295	3.387
	연령	.011	.011	.054	.962	.345	.983	1.017

a. 종속변수 : 소비자만족도

$$Y = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \beta_3 \cdot X_3$$

$$\text{소비자만족도} = -0.631 + 0.744 \cdot \text{가격만족도} + 0.302 \cdot \text{구매 횟수} + 0.011 \cdot \text{연령}$$

회귀 분석 결과 예시

모형		비표준화 계수		표준화 계수	t	유의 확률	공선성 통계량	
		B	표준오차	베타			공차	VIF
1	(상수)	-.631	.519		-1.215	.235		
	가격만족도	.744	.114	.668	6.528	.000	.298	3.356
	구매 횟수	.302	.094	.331	3.223	.003	.295	3.387
	연령	.011	.011	.054	.962	.345	.983	1.017

a. 종속변수 : 소비자만족도

- 가격만족도와 구매횟수의 유의확률이 유의수준보다 작으므로($p\text{-value} < 0.05$), 통계적으로 유의미한 변수로 판단
- 연령은 유의확률이 유의수준보다 크므로 ($p\text{-value} > 0.05$), 통계적으로 유의하지 않으며 소비자만족도에는 영향을 미치지 않는 변수로 판단

- 회귀모델의 성능 지표

구 분	개 요	수식
평균절대오차 MAE (Mean Absolute Error)	실제 값과 예측한 값의 차이를 절댓값으로 변환해 평균한 값	$MAE = \frac{1}{n} \sum_{i=1}^n Y_i - \hat{Y}_i $
평균제곱오차 MSE (Mean Squared Error)	실제 값과 예측한 값의 차이를 제곱한 후 평균한 값	$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$
평균제곱근오차 RMSE (Root Mean Squared Error)	실제 값과 예측한 값의 차이를 제곱한 후 평균한 값의 제곱근	$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$
평균절대비율오차 MAPE (Mean Absolute Percentage Error)	실제 값과 예측한 값의 차이를 백분율로 표현	$MAPE = \frac{100}{n} \sum_{i=1}^n \left \frac{Y_i - \hat{Y}_i}{Y_i} \right $

Data : insurance.csv (미국 건강보험료 데이터셋)

✓ 1338행, 7개의 변수 (Target = charges)

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1338 entries, 0 to 1337
```

```
Data columns (total 7 columns):
```

#	Column	Non-Null Count	Dtype	
0	age	1338 non-null	int64	연령
1	sex	1338 non-null	object	성별 (male or female)
2	bmi	1338 non-null	float64	체질량 지수(body mass index)
3	children	1338 non-null	int64	자녀의 수(number of children)
4	smoker	1338 non-null	object	흡연 여부(yes or no)
5	region	1338 non-null	object	사는 지역(northeast, southeast, northwest, southwest)
6	charges	1338 non-null	float64	건강보험에서 지출되는 개인별 의료비

dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB

☰

✎

ChatGPT 4o ▾

To

ChatGPT

데이터분석 제타봇

실파

나무위키 GPT

GPT 탐색

오늘

IT 보안 조사 그룹핑.

어제

서울 날씨 오늘

User query, ChatGPT response

Memory Refresh: Recall Request

지난 7일

파일 확인 및 요약

영어 공부 방법

AI와 데이터 리터러시: 이해와 활용

데이터 분석 책 비교

Data 분석 박사의 책

인적자원의 중요성

RFM 분석 결과 업데이트

RFM 분석 기법

새시랑 함께

Team 워크스페이스 추가

Team 플랜에서 함께 작업하세요

🌀

친구에게 결혼식 동반 부탁 메시지 보내기

📷

카메라에 잘 받는 옷 고르기

💡

생산성 극대화를 위한 모닝 루틴

📄

현지인처럼 서울 체험하기

📄 메시지 ChatGPT

🌀

⬆

ChatGPT는 실수를 할 수 있습니다. 중요한 정보를 확인하세요.

?

※출처: ChatGPT 홈페이지(<https://chat.openai.com/> 2024.06.02. 캡처)