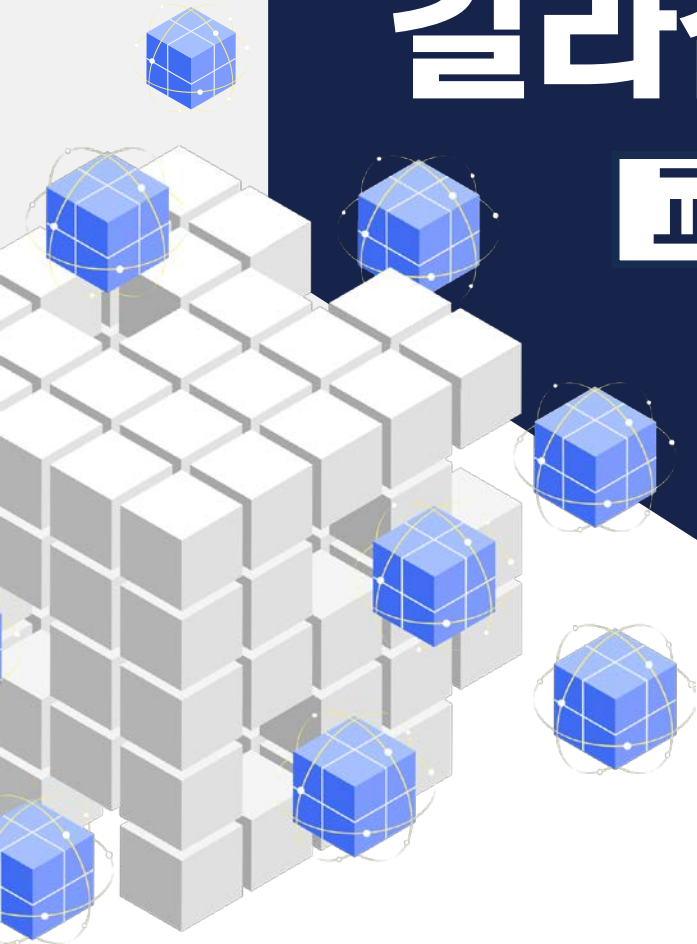


데이터 과학 길라잡이

교사용



사이트 이용 안내

원활한 학습을 위한 학습 전 유의 사항 및 사이트 이용 안내입니다.
아래 내용을 반드시 숙지하신 후 학습 참여 부탁드립니다.

학습 방법

● 수강 완료 후 이수증 출력이 가능한 [학습 도장 모으기] 학습하기

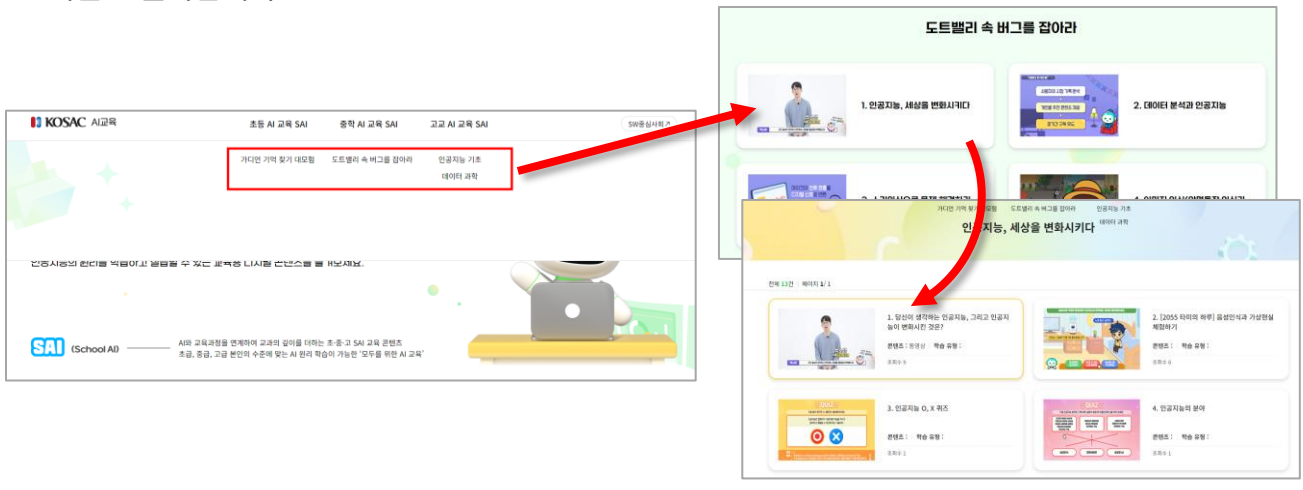
[학습 도장 모으기] 버튼을 클릭하여 학습 시 수강 후 수강생의 이름이 적힌 이수증 출력이 가능합니다.
수강 방법은 홈페이지 메인 중앙의 학년별 배너를 클릭 후 [학습 도장 모으기] 버튼을 클릭하여 학습을 진행하셔야 이수증 출력이 가능합니다.



- * [학습 도장 모으기] 버튼 클릭 외 다른 경로로 콘텐츠 학습 시 이수증 출력이 불가할 수 있으니 이수증 출력이 필요한 경우 반드시 해당 경로로 학습을 진행해 주시기 바랍니다.
- * [학습 도장 모으기]으로 학습 중 학습 미완료 상태로 학습창을 종료할 경우 처음부터 재학습하셔야 하니 유의하셔서 학습에 참여해 주시기 바랍니다.
- * 이수증은 학습페이지 마지막 단계에서 진행되며 이수증 출력을 위해 수강생의 이름을 정확히 입력해주세요.
- * 학습창을 종료한 콘텐츠의 이수증 재발급은 불가하며 발급이 필요한 경우 처음부터 학습을 진행해야 하므로 이수증 출력 시 PDF 파일로도 저장하여 보관해주시기 바랍니다.
- * [학습 도장 모으기]으로 학습 시 학습창 하단의 영상 재생바 조작(재생바 이동)이 불가하므로 정배속으로 순차 학습하셔야 합니다.

● 학습 제한없이 자유롭게 학습하기

이수증 출력 없이 영상만 시청을 원할 경우 홈페이지 상단의 위치한 메뉴에서 각 학년별 과정명을 클릭하여 학습하시면 영상 재생바 이동 등 단순 시청이 가능합니다. 다만 해당 경로로 학습할 경우 학습을 모두 완료하였어도 이수증 출력은 불가하며 영상을 재학습하여도 이전 학습한 부분부터 이어서 학습은 불가합니다.



원활한 학습을 위한 학습 전 유의 사항 및 사이트 이용 안내 입니다.
아래 내용을 반드시 숙지하신 후 학습 참여 부탁드립니다.

학습 시작 전 유의 사항

● 본 사이트는 회원가입 및 로그인 없이 바로 콘텐츠 학습이 가능합니다.

해당 사이트는 별도의 회원가입 절차가 없고 바로 콘텐츠 학습이 가능합니다.
다만 로그인이 없기 때문에 학습에 대한 기록이 남지 않아 학습 중단 후 학습창 이탈 시 이전에 학습한 영상을 이어서 하거나 이수증 재출력이 불가하오니 종료 전 반드시 확인 후 종료해 주세요.

● 학습창을 종료하고 재접속 시 이어서 학습은 불가합니다.

개인의 학습 이력 관리가 없기 때문에 학습창 이탈 후 재접속 시 이전 학습에서 이어서 학습이 불가합니다. 이미 수강이 완료된 콘텐츠도 재접속 시 처음부터 새롭게 학습이 진행됩니다.
[학습 도장 모으기]의 콘텐츠를 학습할 경우 영상의 재생바 조작이 불가하며 이미 학습을 완료한 차시도 재생바 조작이 불가합니다.

● 이수증 출력을 원하시는 경우 반드시 [학습 도장 모으기] 버튼을 클릭하여 학습해주세요.

해당 사이트에서는 이수증을 출력할 수 있는 학습 방법과 단순 영상만 시청하는 학습 방법이 있습니다.
이수증 출력을 원하시는 경우 홈페이지 메인 화면에서 학년별 배너를 클릭 후 [학습 도장 모으기] 버튼을 클릭하여 학습하셔야 학습 완료 후 이수증 출력이 가능합니다.
다만 로그인 없이 학습하기 때문에 학습 이탈 시 이어서 학습하기가 불가하여 처음부터 학습을 진행해야 하니 [학습 도장 모으기]으로 학습하는 경우 반드시 끝까지 학습을 완료하여 이수증을 출력하고 종료해 주시기 바랍니다.

* 이수증 재발급이 필요한 경우 처음부터 학습을 재시작하셔야 합니다.

* 이수증 인쇄 전 이수증 내 입력한 이름 정보가 틀리지 않았는지 반드시 확인 후 출력해주세요.

● 학습은 데스크톱, 노트북, 태블릿PC 기기에서 학습해주세요.

본 콘텐츠는 화면을 클릭하여 진행해야 하는 다양한 상호작용이 있는 콘텐츠로 모바일 학습 시 원활한 학습이 불가합니다. 데스크톱, 노트북, 태블릿 PC를 통해 학습해 주시기 바랍니다.

무엇을 도와 드릴까요?

대표전화 1522-6841

문의메일 ai4school@kosac.re.kr



*평일 9시~18시 | 점심시간 12시~13시 (일요일/공휴일 휴무)



목차

1과	데이터과학으로 여는 세계	3
2과	데이터의 속성	4-10
3과	데이터셋과 데이터베이스	11-14
4과	데이터 수집의 중요성	15-17
5과	데이터 전처리	18-21
6과	데이터 시각화	22-25
7과	데이터 상관관계	26-30
8과	데이터 분석방법	31-39
9과	데이터 모델링 및 데이터 분석 도구 탐색하기	40
10과	단순 선형 회귀, 다중 선형 회귀	41-42
11과	회귀 분석, 결정 계수	43-45
12과	군집화 개념 이해하기	46-48
13과	군집분석	49-51
14과	지지도, 신뢰도, 향상도	52-55
15과	장바구니 분석	56-58



1. 전체 콘텐츠 명

데이터과학

2. 콘텐츠 개요

◆ 과정 목차

모듈	순서	차시명	주제명	유형	소요 시간	
데이터 과학의 이해	1	데이터과학으로 여는 세계	데이터 기반 의사 결정	영상형	8분	
	2	데이터의 속성	다양한 데이터의 속성	영상형	6분	
			인공지능(AI) 학습을 위한 데이터 속성 분류하기	실습형	10분	
			정형데이터와 비정형 데이터 구분하기	실습형	10분	
	3	데이터셋과 데이터베이스	데이터셋과 데이터베이스 기초	영상형	6분	
			데이터셋과 데이터베이스	실습형	10분	
데이터 준비와 분석	4	데이터 수집의 중요성	데이터 수집방법 알아보기	영상형	7분	
			데이터 수집의 중요성	실습형	10분	
	5	데이터 전처리	결측치, 이상치, 정규화	영상형	6분	
			데이터 전처리	실습형	10분	
	6	데이터 시각화	데이터 시각화 특징과 그래프 종류	영상형	7분	
			데이터 시각화	실습형	10분	
	7	데이터 상관관계	데이터 시각화를 통한 상관관계 이해하기	영상형	4분	
			데이터 상관관계	실습형	10분	
	8	데이터 분석방법	데이터 분석방법 알아보기	영상형	4분	
			건강데이터 살펴보기	실습형	15분	
			건강데이터로 비만 분포 확인하기	실습형	15분	
	데이터 모델링과 평가	9	데이터 모델링 및 데이터 분석 도구 탐색하기	다양한 분석 도구 살펴보기	영상형	6분
		10	단순 선형 회귀, 다중 선형 회귀	선형 회귀의 기본 이해	영상형	6분
				회귀분석으로 예측하기	실습형	10분
		11	회귀 분석, 결정 계수	R^2 값 알아보기	영상형	6분
회귀분석과 결정계수				실습형	10분	
12		군집화 개념 이해하기	보로노이 다이어그램과 센트로이드	영상형	5분	
			SNS 사용 데이터를 활용한 군집화 알아보기	실습형	10분	
13		군집분석	포켓몬빵 사례로 알아보는 군집분석	영상형	7분	
	군집분석		실습형	10분		
14	지지도, 신뢰도, 향상도	연관분석 기초	영상형	6분		
		지지도, 신뢰도, 향상도	실습형	10분		
15	장바구니 분석	장바구니 분석 이해하기	영상형	4분		
		장바구니 분석	실습형	5분		



◆ 과정 구성

영상형 콘텐츠 15개, 실습형 콘텐츠 15개 (총 30개 콘텐츠)

1) 영상형

데이터과학의 이해 모듈(3개), 데이터 준비와 분석 모듈(5개),
데이터 모델링과 평가 모듈(7개)

❖ 구성: 지식영상 & 캐릭터 애니메이션 & 인터뷰 영상

2) 실습형

데이터과학의 이해 모듈(3개), 데이터 준비와 분석 모듈(6개),
데이터 모델링과 평가 모듈(6개)

◆ 추천 학습 환경

실습형 콘텐츠 오류 발생 시 대처 방법

- 오류 발생 이유: 실습환경은, 사용하시는 기기의 메모리를 활용하여 진행할 수 있도록 구성되어 있습니다. 실습 진행 시, 기기의 메모리를 초과하는 실습을 진행하거나(다수의 버튼을 클릭하여 메모리 실행 초과 등) 기기의 메모리가 현저히 낮다면, 메모리 관련 오류가 발생할 수 있습니다.
- 오류 발생 해결 방법: 실습 진행 시, 기기의 메모리를 소모할 수 있는 다른 프로그램을 종료하여 최적의 실습 환경으로 구성해 주시기를 바랍니다. 이러한 조치에도 오류가 발생하였다면, 새로고침 혹은 이전/다음 페이지를 눌러 실습 환경을 초기화 후 다시 시도해 보시기 바랍니다.

실습형 실행 환경

- 해상도: 실습 콘텐츠는 해상도 1680 * 1050을 기준으로 제작되어 있습니다. 필요시 해상도를 조정해 주세요.
- 실행 환경: 실습 콘텐츠는 PC 환경 및 Chrome, Edge 브라우저 환경에 최적화되어 있습니다.



3. 세부 콘텐츠

◆ 개요

1차시. 데이터과학으로 여는 세계

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과아-01] 데이터과학의 개념을 이해하고, 문제 해결 사례를 데이터 기반 의사 결정 상황에 적용한다. [12데과아-04] 데이터로 인한 사회 변화를 인식하고, 진로 및 직업과 관련한 데이터 기반 문제 해결 사례를 분석한다.
학습 목표	<ul style="list-style-type: none"> • 데이터과학의 개념을 이해하고 데이터 기반 의사 결정의 사례를 설명할 수 있다. • 데이터로 인한 사회 변화를 이해하고 진로 및 직업과 관련한 데이터 기반 문제 해결 사례를 설명할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	1. (강의실) 동영상 콘텐츠

◆ 구성

강의실

- ❖ **주제명:** 데이터 기반 의사 결정
- ❖ **세부 주제**
 - 1) 데이터과학의 개념
 - 2) 데이터 기반 의사 결정을 함으로써 문제를 해결하는 사례와 사회적 변화 안내
- ❖ **콘텐츠 개발 목적:** 현직자 인터뷰를 통해 데이터과학의 필요성을 전달한다.
- ❖ **선생님을 위한 팁!**
추가적인 데이터 기반 의사 결정 사례 예시를 다루어 볼 수 있다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

02차시. 데이터의 속성

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과아-02] 정형 데이터와 비정형 데이터를 구분하고, 데이터 속성에서 데이터의 잠재적 가치를 파악한다.
학습 목표	<ul style="list-style-type: none"> 정형 데이터와 비정형 데이터의 개념을 이해하고 특징을 설명할 수 있다. 데이터 속성이 데이터과학 분야에서 갖는 잠재적 가치를 설명할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

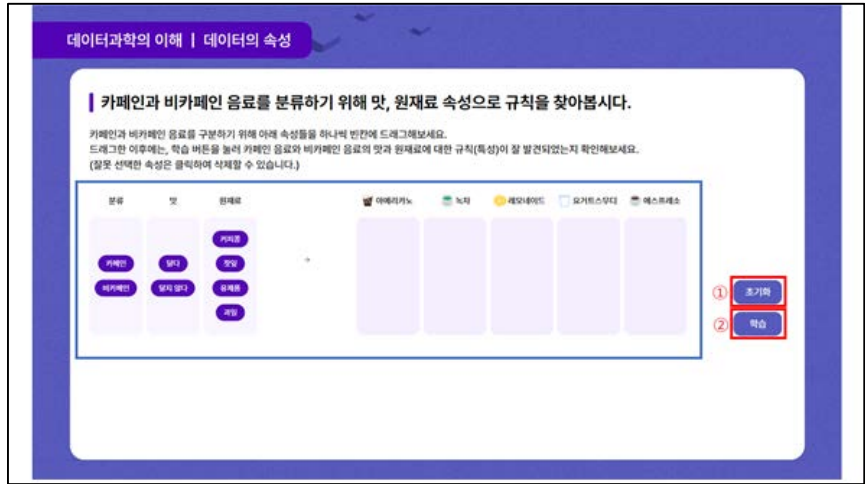
- ❖ **주제명:** 다양한 데이터의 속성
- ❖ **세부 주제**
 - 1) 정형 데이터와 비정형 데이터
 - 2) 범주형 데이터와 수치형 데이터
- ❖ **콘텐츠 개발 목적:** 같은 데이터라도 목적에 따라 속성들의 중요성이 달라질 수 있음을 이해할 수 있다.
- ❖ **선생님을 위한 팁!**
정형 데이터와 비정형 데이터의 구조에 대해 추가 설명이 필요할 수 있다.

실험실

- ❖ **주제명:** 데이터의 속성
- ❖ **세부 주제**
 - 1) 데이터 속성 분류의 중요성
 - 2) 인공지능(AI)을 학습시키기 위한 데이터 속성 분류 체험
- ❖ **콘텐츠 개발 목적:** 인공지능(AI)을 만들 때 데이터 속성을 분류하는 과정을 단계별로 실습하며 중요성을 인식



❖ 실습 콘텐츠 안내:

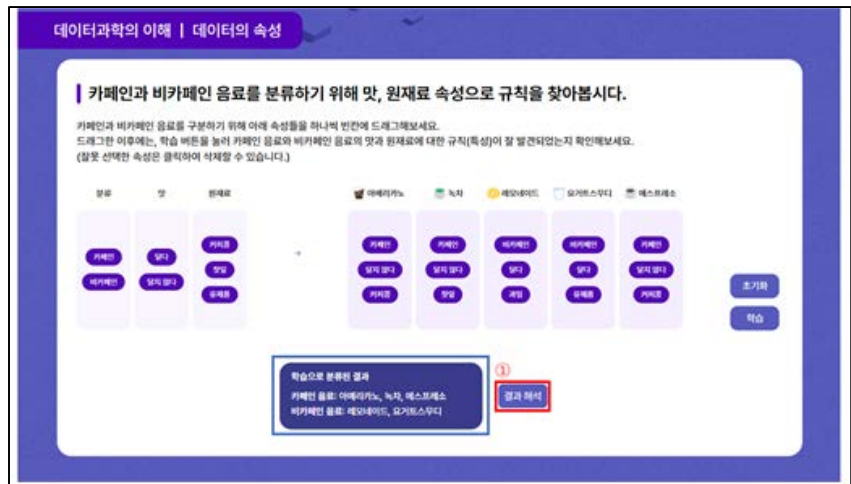


Step 1

분류, 맛, 원재료 속성을 각 음료 영역으로 드래그하여 카페인/비카페인 음료 구분을 위한 속성을 인공지능에게 학습시키기

기능 사용법

- ① 초기화 버튼: 각 속성을 드래그 앤 드롭으로 분류하는 과정에서 초기화 버튼을 누르면 분류 전 처음 상태로 돌아갈 수 있습니다. (전체 초기화를 하지 않고 특정 속성만 지우고 싶을 경우, 클릭하여 지울 수 있습니다.)
- ② 학습 버튼: 모든 분류가 끝나면, 분류한 결과로 인공지능(AI)을 학습시킬 수 있는 버튼입니다.



Step 2

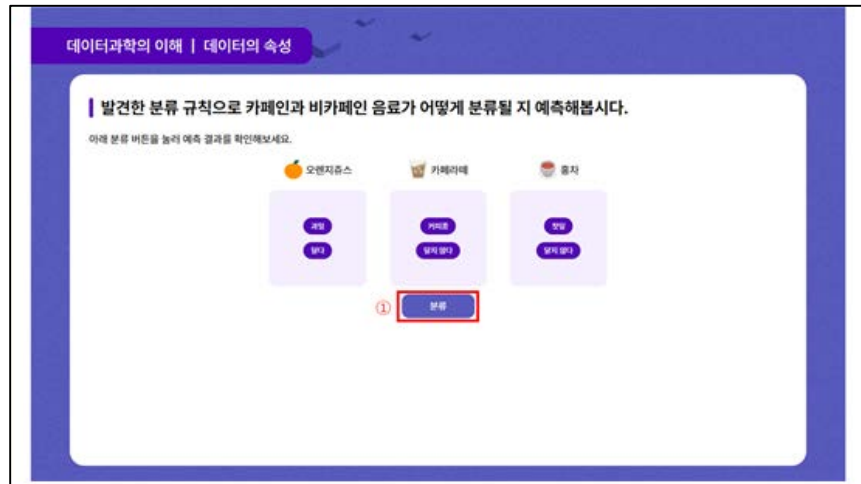
학습된 인공지능(AI)의 음료 분류 결과를 확인

기능 사용법

- ① 결과 해석 버튼: 인공지능(AI)이 학습한 결과를 토대로 분류된 결과에 대한 설명을 볼 수 있습니다.



❖ 실습 콘텐츠 안내:



Step 3

기존에 분류해 둔 속성들을 기준으로, 다른 음료들의 분류 결과가 어떻게 나올지 예측해 본 뒤 분류 버튼을 눌러 확인하기

기능 사용법

① 분류 버튼: 인공지능(AI)을 학습시킨 분류 규칙으로 음료를 분류하는 버튼입니다. 버튼을 누르면 분류 결과와 설명을 확인할 수 있습니다.



Step 4

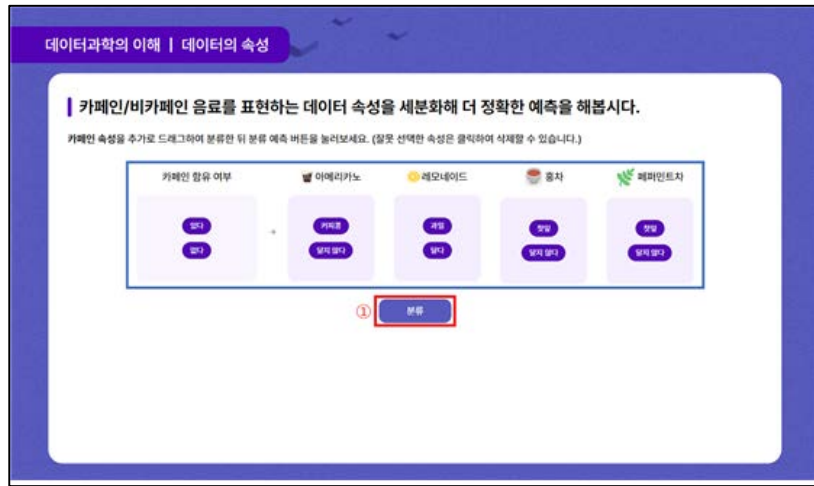
기존에 학습된 결과로 분류할 수 없는 예외 케이스를 분류 버튼을 눌러 확인하기

기능 사용법

① 분류 버튼: 인공지능(AI)을 학습시킨 분류 규칙으로 음료를 분류하는 버튼입니다. 버튼을 누르면 분류 결과와 설명을 확인할 수 있습니다.



❖ 실습 콘텐츠 안내:



Step 5

예외 케이스까지 포함할 수 있는 속성을 드래그 앤 드롭으로 추가 분류 하고, 분류 버튼을 눌러 결과 확인하기

① 분류 버튼: 인공지능(AI)을 학습시킨 분류 규칙으로 음료를 분류하는 버튼입니다. 버튼을 누르면 분류 결과와 설명을 확인할 수 있습니다.

주의

분류 실습 시, 아이콘을 너무 빠르게 움직이면 아이콘이 없어지는 경우가 있으므로 천천히 움직여주세요.

❖ 선생님을 위한 팁!

인공지능에 학습시킬만한 주제, 그에 따른 속성들을 추가적으로 탐구해보는 시간을 가지면 원리 학습에 더욱 도움이 됩니다.





3. 세부 콘텐츠

◆ 개요

◆ 구성

02차시. 데이터의 속성

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과아1-02] 정형 데이터와 비정형 데이터를 구분하고, 데이터 속성에서 데이터의 잠재적 가치를 파악한다.
학습 목표	<ul style="list-style-type: none"> 정형 데이터와 비정형 데이터의 개념을 이해하고 특징을 설명할 수 있다. 데이터 속성이 데이터과학 분야에서 갖는 잠재적 가치를 설명할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

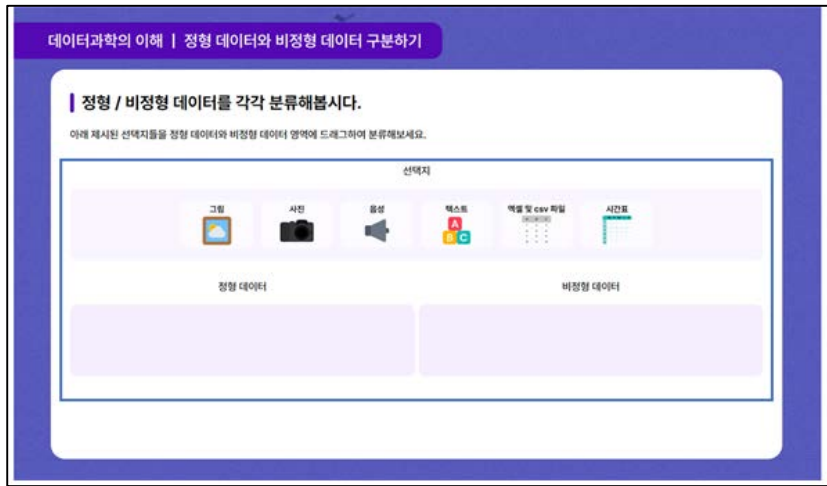
- ❖ **주제명:** 다양한 데이터의 속성
- ❖ **세부 주제**
 - 1) 정형 데이터와 비정형 데이터
 - 2) 범주형 데이터와 수치형 데이터
- ❖ **콘텐츠 개발 목적:** 같은 데이터라도 목적에 따라 속성들의 중요성이 달라질 수 있음을 이해할 수 있다.
- ❖ **선생님을 위한 팁!**
정형 데이터와 비정형 데이터의 구조에 대해 추가 설명이 필요할 수 있다.

실험실

- ❖ **주제명:** 정형 데이터와 비정형 데이터 구분하기
- ❖ **세부 주제**
 - 1) 정형/비정형 데이터를 구분
 - 2) 비정형 데이터를 정형 데이터로 바꾸어 표현
 - 3) 정형 데이터 내에서 문제를 해결하기 위해 어떤 데이터 속성이 필요한지 결정
- ❖ **콘텐츠 개발 목적:** 분석을 위한 데이터의 유형 및 특성 이해

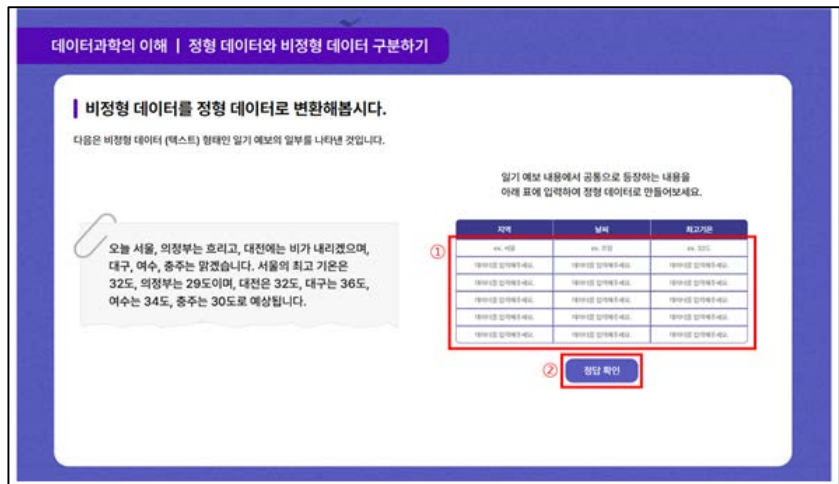


❖ 실습 콘텐츠 안내:



Step 1

선택지에 있는 데이터 유형을 정형/비정형 데이터로 각각 분류



Step 2

비정형 데이터 (텍스트) 형태인 왼쪽의 일기예보를 오른쪽 정형 데이터 틀에 옮겨 적어 정형 데이터로 만들고, 정답 확인 버튼을 눌러 올바르게 변환했는지 확인하기

기능 사용법

① 데이터 입력란: 왼쪽 일기예보를 보고 첫번째 행에 ex로 표시된 정형데이터 형식과 같이 작성합니다.

② 정답 확인 버튼: ①번 입력란을 전부 기입한 뒤, 정답 확인 버튼을 누르면 일기예보를 정형 데이터로 변환했을 때의 정답과 입력한 내용을 비교해 확인할 수 있습니다.



❖ 실습 콘텐츠 안내:

데이터과학의 이해 | 정형 데이터와 비정형 데이터 구분하기

1시간 전 기상 상황에 따른 자전거 대여량을 예측하기 위한 데이터 속성을 골라봅시다.

아래 데이터를 보고 필요한 속성명을 하단 우측 영역으로 드래그해주세요. (잘못 선택한 속성은 클릭하여 삭제할 수 있습니다.)

고유 번호	대여 시간	1시간 전 기온	1시간 전 비	1시간 전 풍속	1시간 전 습도	1시간 전 가시성	1시간 전 오존	1시간 전 미세먼지(pm10)	1시간 전 초미세먼지(pm2.5)	대여 개수
3	20	16.3	1	1.5	89	576	0.027	76	33	49
6	13	20.1	0	1.8	48	956	0.042	73	40	159
7	6	13.9	0	5.7	79	1382	0.033	32	19	26
8	23	8.1	0	2.7	54	946	0.04	75	64	57
9	18	29.5	0	4.8	7	2000	0.067	27	11	431

데이터 속성

고유 번호, 대여 시간, 1시간 전 기온, 1시간 전 비, 1시간 전 풍속, 1시간 전 습도, 1시간 전 가시성, 1시간 전 오존, 1시간 전 미세먼지(pm10), 1시간 전 초미세먼지(pm2.5), 대여 개수

선택된 데이터 속성

1시간 전 기온, 1시간 전 미세먼지(pm10), 1시간 전 초미세먼지(pm2.5), 대여 개수

① 제출

Step 3

예측 조건이 주어졌을 때, 예측에 필요한 데이터 속성을 ‘선택된 데이터 속성’ 영역에 드래그하고 제출 버튼을 눌러 제대로 분류 했는지 확인하기 가능 사용법

① 제출 버튼: 데이터 속성 영역에서, 조건에 맞는 속성을 선택된 데이터 속성으로 드래그 앤 드롭한 후 제출 버튼을 누르면 정/오답 여부를 판별할 수 있습니다.

주의

일기 예보를 정형 데이터로 변환하는 실습에서, 날씨를 글로 작성할 때는 예시와 동일한 형태로 작성해주세요. 예를 들어, 날씨를 흐림/비/맑음 세 가지로만 작성할 수 있으며, 흐리다/비가온다/맑다 등으로 작성하면 오답 처리 될 수 있습니다.

❖ 선생님을 위한 팁!

Step 3에 해당하는 실습에서, 미세먼지 수치에 따른 대여량은 어떤 속성을 가지고 예측할 수 있을지 학생들에게 추가 질문을 해보세요.

정형 데이터

엑셀 및 csv 파일

시간표

비정형 데이터

그림

사진

음성

텍스트

선택된 데이터 속성

1시간 전 기온

1시간 전 미세먼지(pm10)

1시간 전 초미세먼지(pm2.5)

대여 개수

1시간 전 오존

1시간 전 가시성

1시간 전 습도

1시간 전 풍속

1시간 전 비



3. 세부 콘텐츠

◆ 개요

◆ 구성

03차시. 데이터셋과 데이터베이스

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과아-03] 데이터셋의 집합인 데이터베이스를 이해하고, 서로 다른 데이터셋의 데이터를 분석이 가능한 형태로 통합하는 것의 의미를 파악한다.
학습 목표	<ul style="list-style-type: none"> • 데이터셋과 데이터베이스를 이해하고 설명할 수 있다. • 서로 다른 데이터셋의 데이터를 분석 가능한 형태로 통합하고 의미를 설명할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> 1. (강의실) 동영상 콘텐츠 2. (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 데이터셋과 데이터베이스 기초
- ❖ **세부 주제**
 - 1) 데이터셋과 데이터베이스의 개념
 - 2) 데이터베이스를 사용하는 것의 장점
 - 3) 데이터베이스를 사용하는 사례
- ❖ **콘텐츠 개발 목적:** 데이터셋과 데이터베이스의 개념을 이해하고 실생활에 적용할 수 있다.
- ❖ **선생님을 위한 팁!**
데이터베이스 실생활 예시를 추가 제공할 수 있다.

실험실

- ❖ **주제명:** 데이터셋과 데이터베이스
- ❖ **세부 주제**
 - 1) 파일 시스템 관리 방법 이해
 - 2) 데이터셋과 DBMS 이해
 - 3) DBMS로 데이터 관리하는 방법, 장점 이해
- ❖ **콘텐츠 개발 목적:** 데이터셋을 관리하는 방법과, DBMS의 장점을 실습을 통해 이해



❖ 실습 콘텐츠 안내:

데이터과학의 이해 | 데이터 셋과 데이터베이스

필수가 수강과목을 변경했을때 데이터를 수정해봅시다.

1. 필수가 수강과목을 인증하는 기초에서 데이터 과학으로 바꾸었습니다.
2. 필수이 데이터 과학 과목 성적은 95, 등급은 A입니다.

학급 학생 데이터				데이터 과학 교과 수업생 학생 데이터				재학생 데이터			
학번ID	이름	연락처	수강과목	학번ID	이름	성적	등급	학번ID	이름	연락처	재학 여부
01-01	홍길동	010-1234-5678	데이터 과학	01-01	홍길동	90	A	01-01	홍길동	010-1234-5678	재학
01-02	김영희	010-4567-8901	데이터 과학	01-02	김영희	85	B	01-02	김영희	010-8765-4321	재학
01-03	박정수	010-1234-1234	인공지능 기초					01-03	박정수	010-1234-1234	재학

데이터를 직접 수정하고 제출버튼을 눌러주세요.

제출

Step 1

파일 시스템 환경에서, 데이터 변경/추가로 인한 관리 방법을 실습 (초록색으로 표시된 부분에 필요한 내용을 입력하여 수정, 추가) 기능 사용법

- ① 데이터 입력란: 지시문을 따라 데이터를 추가, 수정하여 입력합니다.
- ② 제출 버튼: 데이터를 수정/추가한 후 해당 버튼을 누르면 정/오답 여부를 확인할 수 있습니다.

데이터과학의 이해 | 데이터 셋과 데이터베이스

파일 시스템에서 관리하던 데이터를, 데이터베이스 관리 시스템(DBMS)에 옮겨봅시다.

아래 파일들을 데이터베이스 관리 시스템(DBMS)으로 드래그해보세요.

학급 학생 데이터 파일				데이터과학 교과 수업생 데이터 파일				재학생 데이터 파일			
학번ID	이름	연락처	수강과목	학번ID	이름	성적	등급	학번ID	이름	연락처	재학 여부

데이터베이스 관리 시스템 (DBMS)

Step 2

파일 시스템으로 관리하던 파일들을 데이터 베이스 관리 시스템 영역으로 드래그하여 옮기기



❖ 실습 콘텐츠 안내:

데이터과학의 이해 | 데이터 셋과 데이터베이스

각 데이터 셋에 공통적으로 존재하는 속성을 식별하고, 중복을 제거해봅시다.

데이터셋의 속성 중, 각 데이터를 다른 데이터와 구별해주는 최소한의 속성 집합을 '후보키'라고 부릅니다. 이 후보키들 중에서 가장 중요하고 핵심적인 내용을 포함하는 키를 '기본키'로 지정할 수 있습니다.

(1) 오른쪽 학급 학생 데이터셋에서 기본키가 될 수 있는 속성을 선택해 보세요.
기본키는 데이터를 간결하고 명확하게 식별하는 데 사용됩니다.
데이터의 무결성을 보장하기 위해 변경할 수 없으며, 중복되거나 공백의 값을 가질 수 없습니다.

학생ID	이름	연락처	수강과목
01-01	홍길동	010-1234-5678	데이터 과학
01-02	김영희	010-8765-4321	데이터 과학
01-03	박철수	010-1234-1234	데이터 과학
01-04	최민수	010-3333-4444	데이터 과학

Step 3

DBMS로 옮긴 데이터셋에 공통적으로 존재하는 속성을 식별하고, 중복을 제거하기

기능 사용법

- ① 속성 선택: 지시문에 해당하는 속성을 클릭하여 선택합니다. 열(학생ID, 이름, 연락처, 수강과목)을 기준으로 데이터를 선택할 수 있습니다. 선택 시 해당 열이 초록색으로 표시되며 정오답 여부를 확인할 수 있습니다. 올바른 속성을 선택할 경우 (2)번 문제가 화면에 표시됩니다.

(2) 1번에서 선택한 기본키는 다른 데이터셋에도 존재할 수 있습니다. 이를 기반으로 데이터과학 교과 학생 데이터와 재학생 데이터에서 식별해도 되는 속성을 선택해 보세요. (3개 선택)

기본키는 데이터 셋들 간의 연결성을 제공하여 중복되는 속성들을 제거할 수 있습니다. 이로써 데이터의 정규화를 이루고 데이터 관리의 효율성을 높일 수 있습니다.

학생ID	이름	성적	등급
01-01	홍길동	90	A
01-02	김영희	75	B
01-03	박철수	95	A
01-04	최민수		

학생ID	이름	연락처	성명
01-01	홍길동	010-1234-5678	재학생
01-02	김영희	010-8765-4321	재학생
01-03	박철수	010-1234-1234	재학생
01-04	최민수	010-3333-4444	재학생

Step 4

기본키로 연결된 데이터셋에서, 중복되는 속성들을 제거하기

기능 사용법

- ① 속성 선택: 지시문에 해당하는 속성을 클릭하여 선택합니다. 열(학생ID, 이름, 연락처, 수강과목)을 기준으로 데이터를 선택할 수 있으며, 선택 시 해당 열이 초록색으로 표시되며 정오답 여부를 확인할 수 있습니다.

주의

3, 4 페이지에서 초록색으로 표시된 영역 외의 부분을 수정할 경우 정답 처리에 영향을 줄 수 있습니다.

❖ 선생님을 위한 팁!

DBMS로 데이터를 관리하려면, 관리할 데이터가 정형 데이터여야 합니다. 이전 차시의 실습에서 배웠던 비정형 데이터를 정형데이터로 변환하는 실습을 다시 언급하며, DBMS로 관리하고 싶은 데이터가 정형데이터가 아닐 경우, 정형화 과정이 필요하다는 점을 알려주시면 더욱 좋습니다.



❖ 실습 콘텐츠 안내:

❖ 철수가 수강과목을 변경했을때 데이터를 수정해봅시다.

1. 철수가 수강과목을 인공지능 기초에서 데이터 과학으로 바꾸었습니다.
2. 철수의 데이터 과학 과목 성적은 95, 등급은 A입니다.

학급 학생 데이터				데이터 과학 교과 수업 학생 데이터				재학생 데이터			
학생ID	이름	연락처	수강과목	학생ID	이름	성적	등급	학생ID	이름	연락처	재학 여부
01-01	홍길동	010-1234-5678	데이터 과학	01-01	홍길동	90	A	01-01	홍길동	010-1234-5678	재학
01-02	김영희	010-8765-4321	데이터 과학	01-02	김영희	75	B	01-02	김영희	010-8765-4321	재학
01-03	박철수	010-1234-1234	데이터 과학	01-03	박철수	95	A	01-03	박철수	010-1234-1234	재학

❖ 최민수 학생이 전학을 왔을때 데이터를 관리해봅시다.

1. 최민수라는 학생이 "전학"을 왔습니다.
2. 최민수 학생의 학생 ID는 01-04, 연락처는 010-3333-4444, 수강과목은 데이터 과학입니다.

학급 학생 데이터				데이터 과학 교과 수업 학생 데이터				재학생 데이터			
학생ID	이름	연락처	수강과목	학생ID	이름	성적	등급	학생ID	이름	연락처	재학 여부
01-01	홍길동	010-1234-5678	데이터 과학	01-01	홍길동	90	A	01-01	홍길동	010-1234-5678	재학
01-02	김영희	010-8765-4321	데이터 과학	01-02	김영희	75	B	01-02	김영희	010-8765-4321	재학
01-03	박철수	010-1234-1234	데이터 과학	01-03	박철수	95	A	01-03	박철수	010-1234-1234	재학
01-04	최민수	010-3333-4444	데이터 과학	01-04	최민수			01-04	최민수	010-3333-4444	재학

(1) 오른쪽 학급 학생 데이터셋에서 기본키가 될 수 있는 속성을 선택해보세요.

기본키는 데이터를 간결하고 명확하게 식별하는 데 사용됩니다.

데이터의 무결성을 보장하기 위해 변경될 수 없으며, 중복되거나 공백의 값을 가질 수 없습니다.

학급 학생 데이터셋

학생ID	이름	연락처	수강과목
01-01	홍길동	010-1234-5678	데이터 과학
01-02	김영희	010-8765-4321	데이터 과학
01-03	박철수	010-1234-1234	데이터 과학
01-04	최민수	010-3333-4444	데이터 과학

(2) 1번에서 선택한 기본키는 다른 데이터셋에도 존재할 수 있습니다.

이를 기반으로 데이터과학 교과 학생 데이터와 재학생 데이터에서 삭제해도 되는 속성을 선택해보세요. (3개 선택)

기본키는 데이터 셋들 간의 연결점을 제공하여 중복되는 속성들을 제거할 수 있습니다. 이로써 데이터의 정규화를 이루고 데이터 관리의 효율성을 높일 수 있습니다.

데이터 과학 교과 수업 학생 데이터셋

학생ID	이름	성적	등급
01-01	홍길동	90	A
01-02	김영희	75	B
01-03	박철수	95	A
01-04	최민수		

재학생 데이터셋

학생ID	이름	연락처	재학 여부
01-01	홍길동	010-1234-5678	재학
01-02	김영희	010-8765-4321	재학
01-03	박철수	010-1234-1234	재학
01-04	최민수	010-3333-4444	재학



3. 세부 콘텐츠

◆ 개요

◆ 구성

04차시. 데이터 수집의 중요성

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-01] 데이터를 편향되지 않도록 수집하고, 수집된 데이터의 특성을 분석한다.
학습 목표	<ul style="list-style-type: none"> • 데이터를 수집하는 방법을 이해하고 목적에 맞게 데이터를 수집할 수 있다. • 데이터의 편향성에 유의하며 데이터를 수집할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> 1. (강의실) 동영상 콘텐츠 2. (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 데이터 수집방법 알아보기
- ❖ **세부 주제**
 - 1) 데이터셋 직접 구성
 - 2) 공유 데이터셋 사용
- ❖ **콘텐츠 개발 목적:** 데이터 분석의 기초가 되는 데이터 수집에 대해 이해할 수 있다.
- ❖ **선생님을 위한 팁!**
좋은 데이터를 선별하는 활동을 추가 구성할 수 있다.

실험실

- ❖ **주제명:** 데이터 수집의 중요성
- ❖ **세부 주제**
 - 1) 데이터 수집 방법 학습
 - 2) 수집한 데이터 특성 파악
- ❖ **콘텐츠 개발 목적:** 데이터를 수집하고, 수집한 데이터의 특성을 파악



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 수집의 중요성

| 도서관 이용 현황을 분석하기 위한 질문에 응답해봅시다.

응답 후 제출 버튼을 누르고, 설문 내용이 정형 데이터로 만들어지는 형태를 관찰해 보세요.

①

학년 일주일에 몇번 도서관을 방문하시나요?

성별

아닌 종류의 책을 좋아하나요?

③ 초기화

학년	성별	주당 도서관 방문횟수	선호 종류
3학년	여	2번	역사
4학년	남	3번	과학
4학년	남	1번	가사
3학년	여	4번	문학
2학년	여	4번	예술
5학년	여	3번	예술
3학년	여	30번	가사

② 제출

Step 1

설문지 형태로 데이터를 기입하고, 정형 데이터로 변환되는 과정 관찰 기능 사용법

① 데이터 입력란: 학년 / 도서관 방문 횟수는 숫자를 입력하고, 성별과 좋아하는 책 종류는 아이콘을 선택하여 데이터를 기입합니다.

② 제출 버튼: 해당 버튼을 누르면 ① 데이터 입력란에 입력된 내용을 기반으로, 오른쪽 정형 데이터 형식의 표에 데이터 행이 추가됩니다.

(제출은 여러번 할 수 있으며, 응답을 자유롭게 추가하고 제출하여 나만의 도서관 이용 현황 데이터를 만들어보세요.)

③ 초기화 버튼: 초기 데이터 외에, 추가한 데이터를 삭제하고 싶을 경우 해당 버튼을 눌러 한번에 삭제할 수 있습니다.



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 수집의 중요성

기존에 구성된 데이터셋을 확인하고, 데이터가 적절한 데이터인지 점검해봅시다.

이전 페이지에서 만들었던 완성된 형태의 데이터셋을 살펴보고, 아래 문제를 풀어주세요.

① 데이터의 개수는 총 몇 개인가요?
ex. 4

학년 속성은 어떤 형태의 데이터인가요?
 범주형 수치형

주당 도서관 이용일수는 어떤 형태의 데이터인가요?
 범주형 수치형

②

학년	성별	주당 도서관 이용일수	선호 종류
3학년	여	2번	역사
4학년	남	3번	과학
1학년	남	1번	가타
2학년	여	4번	문학
2학년	여	4번	예술
5학년	여	10번	예술
3학년	여	30번	가타

Step 2

Step 1 (3 페이지)에서 생성한 데이터의 특성을 파악하여 답안을 작성하고, 제출 버튼을 눌러 정답 확인

기능 사용법

① 데이터 입력란: 데이터 개수에는 숫자 형식, 학년 속성 및 주당 도서관 이용 일수의 형태는 체크 형식으로 답안을 작성해 주세요.

② 정답 확인 버튼: ①번 입력란을 전부 기입한 뒤, 제출 버튼을 누르면 해당 답안에 대한 정/오답 여부를 확인할 수 있습니다.

주의

3,4 페이지 실습에서 사용되는 데이터에서, 처음에 주어지는 4개의 데이터는 초기데이터로 고정되어 있으며, 초기화 버튼을 눌러도 수정할 수 없습니다.

❖ 선생님을 위한 팁!

데이터를 추가하는 3페이지 실습에서, 성별은 여성만 혹은 남성만 선택하는 등의 행위는 데이터 편향이 발생할 수 있음을 함께 안내 해주시면 더욱 좋습니다.

기존에 구성된 데이터셋을 확인하고, 데이터가 적절한 데이터인지 점검해봅시다.

이전 페이지에서 만들었던 완성된 형태의 데이터셋을 살펴보고, 아래 문제를 풀어주세요.

데이터의 개수는 총 몇 개인가요?

학년 속성은 어떤 형태의 데이터인가요?
 범주형 수치형

주당 도서관 이용일수는 어떤 형태의 데이터인가요?
 범주형 수치형

학년	성별	주당 도서관 이용일수	선호 종류
3학년	여	2번	역사
4학년	남	3번	과학
1학년	남	1번	가타
2학년	여	4번	문학
5학년	여	10번	예술



3. 세부 콘텐츠

◆ 개요

◆ 구성

05차시. 데이터 전처리

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-02] 이상치와 결측치 탐색 및 정규화를 통해 전처리하여 오류 가능성을 최소화하고, 데이터 분석을 위해 시각화한다.
학습 목표	<ul style="list-style-type: none"> • 이상치와 결측치의 개념과 탐색 방법을 이해한다. • 이상치와 결측치의 처리 방법을 이해하고, 적절한 처리 방법을 선택할 수 있다. • 정규화의 개념과 방법을 이해하고, 데이터의 분포를 일정하게 조절할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> 1. (강의실) 동영상 콘텐츠 2. (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 결측치, 이상치, 정규화
- ❖ **세부 주제**
 - 1) 결측치 처리방법
 - 2) 이상치 탐색방법
 - 3) 정규화 방법
- ❖ **콘텐츠 개발 목적:** 데이터 전처리의 중요성을 이해하고, 데이터 분석의 전 단계임을 인지한다.
- ❖ **선생님을 위한 팁!**
데이터 전처리의 기본 개념에 대해 이해하지 못하는 학생이 있다면 보충 설명을 먼저 진행해야 한다.

실험실

- ❖ **주제명:** 데이터 전처리
- ❖ **세부 주제**
 - 1) 데이터의 결측치를 확인하고, 제거
 - 2) 데이터 결측치와 이상치를 박스 플롯으로 확인
- ❖ **콘텐츠 개발 목적:** 데이터의 결측치, 이상치를 제거하거나 시각화하는 과정으로 데이터 전처리의 중요성 인식




❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 전처리

실습에 앞서, 공공 자전거 대여량 데이터를 먼저 확인해 봅시다.

데이터를 살펴보는 것은 데이터 분석 및 의사 결정 과정에서 매우 중요한 단계입니다.

공공 자전거 대여량은 CSV 파일로 제공됩니다. CSV는 정형 데이터의 일종으로, 열과 행의 구조로 데이터가 표현되어 있습니다. 열로 구분 되는 것은 속성, 행으로 구분되는 것은 각 개별 데이터입니다.



```


속성
├── 대여일시
├── 대여인수
├── 평균기온
├── 일강수량
├── 대여여부 (quantity)
├── 일강수율
└── ...
    
```

필요한 열 선택

```

필요한 열 선택
├── 대여일시
├── 대여인수
├── 평균기온
├── 일강수량
├── 대여여부 (quantity)
├── 일강수율
└── ...
    
```

필요한 열 선택



Step 1


특정 데이터를 살펴보고, 구조를 이해

데이터 준비와 분석 | 데이터 전처리

실습에 앞서, 공공 자전거 대여량 데이터를 먼저 확인해 봅시다.

데이터를 살펴보는 것은 데이터 분석 및 의사 결정 과정에서 매우 중요한 단계입니다.

공공 자전거 대여량은 CSV 파일로 제공됩니다. CSV는 정형 데이터의 일종으로, 열과 행의 구조로 데이터가 표현되어 있습니다. 열로 구분 되는 것은 속성, 행으로 구분되는 것은 각 개별 데이터입니다.



```

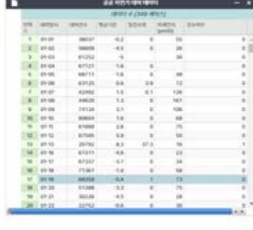
속성
├── 대여일시
├── 대여인수
├── 평균기온
├── 일강수량
├── 대여여부 (quantity)
├── 일강수율
└── ...
    
```

필요한 열 선택

```

필요한 열 선택
├── 대여일시
├── 대여인수
├── 평균기온
├── 일강수량
├── 대여여부 (quantity)
├── 일강수율
└── ...
    
```

필요한 열 선택



Step 2

Step 1 (3 페이지)의 데이터를 요약한 정보를 통해 전반적인 데이터 분포를 확인



❖ 실습 콘텐츠 안내:


데이터 준비와 분석 | 데이터 전처리

| 결측치를 찾아봅시다.

데이터를 분석하기 전, 정확한 분석 결과를 도출하기 위해 누락된 값, 즉 결측치가 있는지 확인해보겠습니다. 결측치는 값이 비어있거나 NA 또는 NaN으로 표시되어 있습니다. 확인된 결측치는 제거해주는 전처리를 수행합니다. 이 과정으로 더 정확하고 신뢰성 있게 데이터 분석을 진행할 수 있습니다.

다음의 단계를 따라 결측치를 가진 데이터를 전부 삭제해보세요!

1. 결측치가 있는 데이터 선택하세요.
2. 휴지통 아이콘 클릭 후 '선택한 케이스 삭제' 버튼 클릭하세요.
3. 결과 확인으로 버튼으로 실습 내용과 전처리 결과 비교해보세요.



주요 지표에 대한 요약 통계

지표	최소값	최대값	평균값	표준편차	중앙값
1. 결측치	0	0	0	0	0
2. 결측치	0	0	0	0	0
3. 결측치	0	0	0	0	0
4. 결측치	0	0	0	0	0
5. 결측치	0	0	0	0	0
6. 결측치	0	0	0	0	0
7. 결측치	0	0	0	0	0
8. 결측치	0	0	0	0	0
9. 결측치	0	0	0	0	0
10. 결측치	0	0	0	0	0

결과 확인

Step 3

결측치(비어있는 칸)가 있는 데이터 행을 삭제하여 전처리 과정을 실습해 보고, '결과 확인' 버튼을 눌러 제대로 수행했는지 확인

기능 사용법

① 데이터 조작 영역: 데이터를 조작할 수 있는 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

② 결과 확인 버튼: 말풍선의 지시문에 따라 실습을 수행한 후, 해당 버튼을 클릭해 정답을 확인하고 작성한 답안과 일치하는지 비교해 보세요.

데이터 준비와 분석 | 데이터 전처리


| 이상치는 어떻게 판별하는 걸까요?

이전 실습에서, 결측치를 제거하는 전처리를 수행했습니다. 결측치와 더불어, 또다른 전처리 대상이 있습니다. 그것은 바로 **이상치**입니다. 이상치는 데이터 수집 과정에서 데이터 범위에서 너무 많이 벗어난 아주 작은 값이나 큰 값을 의미합니다. 이상치는 박스플롯이라는 그래프를 통해 쉽게 시각화가 가능합니다. 이상치를 눈으로 살펴보기 위해, 지전거 대머수 값의 박스플롯을 살펴봅시다.

마우스를 박스플롯 위에 올리보세요. 데이터의 사용처는 박스플롯에서, 이상치는 박스 모양과 직선을 벗어나는 부분에서 탐색할 수 있습니다.

상위 10% 이상치

위 그림을 예시로, 보라색으로 표시된 부분이 이상치에 해당하는 값이 위치합니다.



1차분할(1Q)
데이터를 작은 순서대로 나열했을 때, 하위 25%에 해당하는 값입니다.

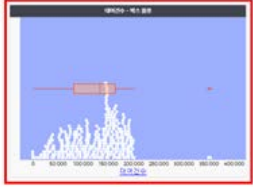
중위값(Median)
데이터를 정렬 후 정렬된 데이터의 중간 값입니다.

3차분할(3Q)
데이터를 작은 순서대로 나열했을 때, 상위 75%에 해당하는 값입니다.

하위값(Min) = 데이터(정렬된) - 1.5 * IQR
이보다 작은 값은 박스플롯에서 제외된다고 간주합니다.

상위값(Max) = 데이터(정렬된) + 1.5 * IQR
이보다 큰 값은 박스플롯에서 제외된다고 간주합니다.

상위 10% 이상치



Step 4

설명을 참고하여 박스플롯 그래프에 표시된 지표 확인하면서, 이상치 판별

기준을 확인

기능 사용법

① 그래프 조작 영역: 말풍선과 말풍선 밑 설명을 토대로 그래프에서 지표를 탐색하는 영역입니다.

20



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 전처리

| 자전거 대여수 값의 이상치를 박스플롯으로 탐색해봅시다.

박스플롯을 기반으로, 이상치를 확인해보았다면 이번에는 직접 자전거 대여수 값을 박스플롯으로 시각화하는 과정을 실습해봅시다.
또한, 박스플롯으로 시각화 한 테이블 데이터를 수정하여 직접 이상치를 만들고, 시각화 한 결과를 확인해 봅시다.

1. 대여건수 속성을 우측 그래프의 x축, 또는 가운데로 드래그하세요.

2. 그래프를 클릭해 우측 메뉴를 눌러보세요.

- 상자그림, 이상치 보기 옵션을 체크하면 박스플롯이 시각화됩니다.

2. 인덱스가 1인 데이터의 대여건수를 400000으로 수정해보세요.

- 수정 후, 박스플롯을 확인해 보면 오른쪽 끝에 이상치를 시각화하는 점(+)이 추가된 것을 확인할 수 있습니다.

Step 5

데이터를 시각화하고, 특정 행의 대여건수 값을 매우 큰 값으로 수정하고 박스플롯의 변화를 관찰

기능 사용법

① 데이터 조작 영역: 데이터를 조작할 수 있는 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

주의

실습 공간이 한정적이므로, 테이블의 크기, 위치 등을 임의로 조정할 경우 실습 진행에 방해가 될 수 있습니다. 가급적으로 주어진 실습 공간의 배치를 움직이지 말아주세요!

❖ 선생님을 위한 팁!

결측치, 이상치를 삭제하는 방법 외에 다른 전처리 방법 (결측, 이상치 전치, 데이터 정규화, 구간화 등)을 추가 설명하면 더욱 좋습니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

06차시. 데이터 시각화

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-02] 이상치와 결측치 탐색 및 정규화를 통해 전처리하여 오류 가능성을 최소화하고, 데이터 분석을 위해 시각화한다.
학습 목표	<ul style="list-style-type: none"> • 데이터 시각화의 개념과 필요성을 이해한다. • 시각화의 종류와 특징을 알고, 적절한 시각화를 선택할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> 1. (강의실) 동영상 콘텐츠 2. (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 데이터 시각화 특징과 그래프 종류
- ❖ **세부 주제**
 - 1) 데이터 시각화 특징
 - 2) 시각화 그래프 종류와 활용 예시
- ❖ **콘텐츠 개발 목적:** 그래프 예시를 통해 데이터 시각화를 알기 쉽게 설명한다.
- ❖ **선생님을 위한 팁!**
필요시, 서울시 공공 자전거 데이터를 활용할 수 있습니다.

실험실

- ❖ **주제명:** 데이터 시각화
- ❖ **세부 주제**
 - 1) 데이터 시각화의 개념과 필요성 인식
 - 2) 데이터 시각화 실습
- ❖ **콘텐츠 개발 목적:** 데이터 분석에서 중요한 단계인 시각화의 개념과 필요성을 인식하고, 배운 이론을 활용해 직접 실습 진행



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 시각화

산점도 그래프를 통해 공공 자전거 대여 데이터를 시각화해봅시다.

산점도는 두 변수간의 관계를 시각적으로 표현하는 사용되는 그래프입니다. 각 데이터 포인트는 x, y 축에서 좌표를 가지며 두 변수간의 관계를 이해할 수 있습니다. 미세먼지(pm10)에 따른 자전거 대여 상황을 산점도 그래프로 시각화해 미세먼지가 많고 적을 때의 자전거 대여 건수 분포를 살펴봅시다. 오른쪽 상단의 산점도 보기 버튼을 클릭해 실습을 진행해보세요.

1. '미세먼지(pm10)' 속성을 그래프의 x축으로 드래그하세요.

2. '대여건수' 속성을 그래프의 y축으로 드래그하세요.

① 산점도 보기

일	대여일시	대여건수	평균기온	일강수량	미세먼지 (pm10)	강수여부
1	01-01	38037	-0.2	0	55	0
2	01-02	56609	-4.5	0	26	0
3	01-05	68711	-1.6	0	48	0
4	01-06	63125	0.6	3.9	72	0
5	01-07	42492	1.5	0.1	126	0
6	01-08	44620	1.3	0	161	0
7	01-09	74134	3.1	0	106	0
8	01-10	80604	1.6	0	68	0
9	01-11	81988	2.6	0	75	0
10	01-12	87595	5.9	0	50	0
11	01-13	20792	8.3	37.3	16	1
12	01-16	61311	-4.6	0	23	0
13	01-17	67337	-3.7	0	34	0
14	01-18	71361	-1.9	0	58	0
15	01-19	68359	-0.4	1	73	0
16	01-20	51309	-3.9	0	76	0

Step 1

미세먼지, 대여 건수 속성들로 구성된 산점도를 시각화 기능 사용법

① 산점도 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

속성

인덱스	대여일시	대여건수	평균기온	일강수량	미세먼지 (pm10)	강수여부
1	01-01	38037	-0.2	0	55	0
2	01-02	56609	-4.5	0	26	0
3	01-05	68711	-1.6	0	48	0
4	01-06	63125	0.6	3.9	72	0
5	01-07	42492	1.5	0.1	126	0
6	01-08	44620	1.3	0	161	0
7	01-09	74134	3.1	0	106	0
8	01-10	80604	1.6	0	68	0
9	01-11	81988	2.6	0	75	0
10	01-12	87595	5.9	0	50	0
11	01-13	20792	8.3	37.3	16	1
12	01-16	61311	-4.6	0	23	0
13	01-17	67337	-3.7	0	34	0
14	01-18	71361	-1.9	0	58	0
15	01-19	68359	-0.4	1	73	0
16	01-20	51309	-3.9	0	76	0

그래프



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 시각화

데이터를 그룹화해 시간대 별 공공 자전거 이용 현황을 파악해봅시다.

각 시간대별로 자전거 이용 건수의 합계를 시각화하여 한눈에 파악하기 위해, 먼저 데이터를 시간대별로 그룹화하는 작업을 수행해야 합니다. 이 단계는 시간대별 이용 현황을 시각화하는 것에 중요한 기초 자료를 제공합니다. 아래의 설명을 따라 실습을 진행해보세요.

1. "데이터" 속성을 테이블 왼쪽의 영역에 드래그하세요. (필요 없을 때는 데이터가 그룹화 됩니다.)
2. 시간 속성에서 그룹 선택하세요. 그룹을 두 개로 나누고 선택하여 그룹 범위(이름)를 선택합니다.
3. "데이터" 열에 새로운 열을 추가하고, 열 이름은 "세 속성"을 입력해 "자전거 이용건수 합"으로 필드를 추가하세요.
4. "자전거 이용건수 합" 필드를 클릭하여 필드 목록에서 "이름" 필드를 사용하여 필드 이름을 추가하세요. 사용 되는 수식은 sum("이름")입니다.

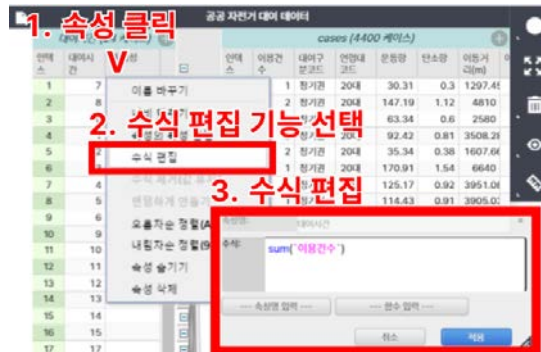
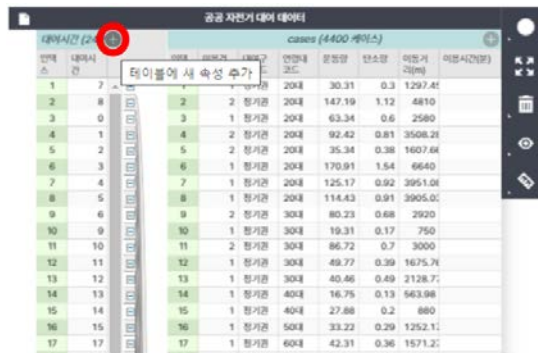
주의: 필드명에는 영문, 숫자, 밑줄(_)만 사용할 수 있습니다. 필드명에는 공백을 사용하지 않습니다.

연도	월	일	시간대	이용건수	평균속도	평균속도	이동거리	이동시간(분)
2018	1	1	20대	30.31	0.3	1297.4t		
2018	1	2	20대	147.19	1.12	4810		
2018	1	3	10대	63.34	0.6	2580		
2018	1	4	20대	92.42	0.81	3508.2t		
2018	1	5	20대	35.34	0.38	1607.6t		
2018	1	6	10대	170.91	1.54	6640		
2018	1	7	10대	125.17	0.92	3951.0t		
2018	1	8	10대	114.43	0.91	3905.0t		
2018	1	9	20대	80.23	0.68	2920		
2018	1	10	10대	19.31	0.17	750		
2018	1	11	20대	86.72	0.7	3000		
2018	1	12	10대	49.77	0.39	1675.7t		
2018	1	13	10대	40.46	0.49	2128.7t		
2018	1	14	10대	16.75	0.13	563.98		
2018	1	15	10대	27.88	0.2	880		
2018	1	16	10대	33.22	0.29	1252.1t		
2018	1	17	10대	42.31	0.36	1571.2t		

Step 2

시각화를 위한 데이터 그룹화 작업을 실습 기능 사용법

① 데이터 조작 영역: 데이터 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.





❖ 실습 콘텐츠 안내:



Step 3

그룹화 한 데이터를 선 그래프 버튼을 눌러 시각화하기
기능 사용법

① 선 그래프 버튼: 버튼을 누르면 그룹화한 데이터를 선그래프로
시각화한 결과와, 시각화에 대한 설명을 확인할 수 있습니다.

주의

실습 공간이 한정적이므로, 테이블의 크기, 위치 등을 임의로
조정할 경우 실습 진행에 방해가 될 수 있습니다. 가급적으로
주어진 실습 공간의 배치를 움직이지 말아주세요!

❖ 선생님을 위한 팁!

말풍선은 한정적인 가이드라인을 표시하지만, 다른 조건으로 시각화 하는
것은 자유입니다. 실습 시 콘텐츠에서 제공한 조건 외, 여러 조건을
제시해서 시각화하는 수업을 다양하게 구성하시면 더욱 좋습니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

07차시. 데이터 상관관계

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-03] 데이터를 분석하기 위해 데이터 속성 간의 관계를 파악하고 통합한다.
학습 목표	<ul style="list-style-type: none"> 상관관계의 개념을 설명할 수 있다. 상관관계를 파악하기 위한 데이터의 시각화와 통계적 분석 방법을 이해할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 데이터 시각화를 통한 상관관계 이해하기
- ❖ **세부 주제**
 - 1) 데이터 시각화
 - 2) 통계적 분석
- ❖ **콘텐츠 개발 목적:** 상관관계의 개념을 이해하고 해석할 수 있다.
- ❖ **선생님을 위한 팁!**
학습 전, 학생들이 데이터 시각화 차시 내용을 복습할 필요가 있습니다.

실험실

- ❖ **주제명:** 데이터 상관관계
- ❖ **세부 주제**
 - 1) 상관관계 이해
 - 2) 데이터 속성 간의 상관관계 시각화
- ❖ **콘텐츠 개발 목적:** 데이터 상관관계를 시각화하는 실습으로 분석에서 상관관계의 역할을 확인



❖ 실습 콘텐츠 안내:


데이터 준비와 분석 | 데이터 상관관계


① 산점도 보기

시각화를 통해 날씨와 자전거 대여 건수의 상관관계를 파악해 봅시다.

상관관계는 두 변수 간 관계의 정도를 나타내는 것으로, 한 변수가 변할때 다른 변수가 어떻게 변하는지에 대한 정보를 제공합니다. 자전거 대여 건수와 평균 기온 간에 어떠한 상관관계가 있는지 산점도로 시각화해봅시다. 오른쪽 상단의 산점도 보기 버튼을 클릭해 실습을 진행해보세요.

1. 평균 기온' 속성을 그래프의 y축으로 드래그하세요.
2. '대여건수' 속성을 그래프의 x축으로 드래그하세요.





산점도를 통해, 평균기온이 높아질수록 자전거 대여건수가 증가하는 양의 상관 관계를 확인할 수 있습니다. 하지만 이로 인한 상관 정도를 정확히 수치로 파악하기는 어렵습니다.

Step 1

날씨와 자전거 대여건수의 상관관계를 그래프로 시각화해 확인
기능 사용법

① 산점도 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

공공 자전거 대여 데이터

속성

인덱스	대여일시	대여건수	평균기온	일강수량	미세먼지 (pm10)	강수여부
1	01-01	38037	-0.2	0	55	0
2	01-02	56609	-4.5	0	26	0
3	01-05	68711	-1.6	0	48	0
4	01-06	63125	0.6	3.9	72	0
5	01-07	42492	1.5	0.1	126	0
6	01-08	44620	1.3	0	161	0
7	01-09	74134	3.1	0	106	0
8	01-10	80604	1.6	0	68	0
9	01-11	81988	2.6	0	75	0
10	01-12	87595	5.9	0	50	0
11	01-13	20792	8.3	37.3	16	1
12	01-16	61311	-4.6	0	23	0
13	01-17	67337	-3.7	0	34	0
14	01-18	71361	-1.9	0	58	0
15	01-19	68359	-0.4	1	73	0
16	01-20	51398	-3.2	0	76	0

그래프

y축

x축

여기서 클릭하거나 속성을 드래그하세요.



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 시각화

데이터를 그룹화해 시간대 별 공공 자전거 이용 현황을 파악해봅시다.

각 시간대별로 자전거 이용 건수의 합계를 시각화하여 한눈에 파악하기 위해, 먼저 데이터를 시간대별로 그룹화하는 작업을 수행해야 합니다. 이 단계는 시간대별 이용 현황을 시각화하는 것에 중요한 기초 자료를 제공합니다. 아래의 설명을 따라 실습을 진행해보세요.

1. '대역시간' 속성을 대역으로 선택하여 왼쪽 공백에 드래그하세요. (왼쪽 공백 가장 상단에서부터 그룹화 합니다.)

2. '시간' 속성에서 '월' 선택해봅시다. 그룹화 창 상단 주 메뉴를 클릭하여 그룹 범위 (대역)에서 '연월일' 선택합니다.

3. '대역시간' 옆에 새로운 열을 추가하고, 열 이름은 '내 속성'을 클릭해 '자전거 이용건수' 열 선택을 수행해줍니다.

4. '자전거 이용건수' 열 '합계'를 클릭하여 해당 열의 '합계' 기능을 사용해 '내 속성' 열을 추가해줍니다. 사용 시는 수식 'sum('이동건수')' 입니다.

수익 집계범위 설정

1. sum(이동건수) 입력

2. 대역시간으로 범위를 나누고, 해당 범위를 사용

공공 자전거 대여 데이터

연월일	대역시간	이동건수	연월일	대역시간	이동건수	연월일	대역시간	이동건수
2018-01-01	00:00-01:00	30.31	2018-01-01	01:00-02:00	147.19	2018-01-01	02:00-03:00	63.34
2018-01-01	03:00-04:00	92.42	2018-01-01	04:00-05:00	35.34	2018-01-01	05:00-06:00	170.91
2018-01-01	06:00-07:00	125.17	2018-01-01	07:00-08:00	114.43	2018-01-01	08:00-09:00	80.23
2018-01-01	09:00-10:00	19.31	2018-01-01	10:00-11:00	86.72	2018-01-01	11:00-12:00	49.77
2018-01-01	12:00-13:00	40.46	2018-01-01	13:00-14:00	16.75	2018-01-01	14:00-15:00	27.88
2018-01-01	15:00-16:00	33.22	2018-01-01	16:00-17:00	42.31	2018-01-01	17:00-18:00	1571.21

Step 2

산점도로 속성 간 상관관계를 시각화해 보고, 상관관계에 대한 지표인 상관계수를 확인한 뒤 지표를 토대로 상관관계 해석 기능 사용법

① 데이터 조작 영역: 데이터 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

공공 자전거 대여 데이터

연월일	대역시간	이동건수	연월일	대역시간	이동건수	연월일	대역시간	이동건수
2018-01-01	00:00-01:00	30.31	2018-01-01	01:00-02:00	147.19	2018-01-01	02:00-03:00	63.34
2018-01-01	03:00-04:00	92.42	2018-01-01	04:00-05:00	35.34	2018-01-01	05:00-06:00	170.91
2018-01-01	06:00-07:00	125.17	2018-01-01	07:00-08:00	114.43	2018-01-01	08:00-09:00	80.23
2018-01-01	09:00-10:00	19.31	2018-01-01	10:00-11:00	86.72	2018-01-01	11:00-12:00	49.77
2018-01-01	12:00-13:00	40.46	2018-01-01	13:00-14:00	16.75	2018-01-01	14:00-15:00	27.88
2018-01-01	15:00-16:00	33.22	2018-01-01	16:00-17:00	42.31	2018-01-01	17:00-18:00	1571.21

1. 속성 클릭

2. 수식 편집 기능 선택

3. 수식 편집

공공 자전거 대여 데이터

연월일	대역시간	이동건수	연월일	대역시간	이동건수	연월일	대역시간	이동건수
2018-01-01	00:00-01:00	30.31	2018-01-01	01:00-02:00	147.19	2018-01-01	02:00-03:00	63.34
2018-01-01	03:00-04:00	92.42	2018-01-01	04:00-05:00	35.34	2018-01-01	05:00-06:00	170.91
2018-01-01	06:00-07:00	125.17	2018-01-01	07:00-08:00	114.43	2018-01-01	08:00-09:00	80.23
2018-01-01	09:00-10:00	19.31	2018-01-01	10:00-11:00	86.72	2018-01-01	11:00-12:00	49.77
2018-01-01	12:00-13:00	40.46	2018-01-01	13:00-14:00	16.75	2018-01-01	14:00-15:00	27.88
2018-01-01	15:00-16:00	33.22	2018-01-01	16:00-17:00	42.31	2018-01-01	17:00-18:00	1571.21



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 상관관계

상관 정도를 정확히 수치로 파악하기 위해, 상관 계수를 살펴봅시다.

1. 데이터셋의 상관 관계를 파악하고자 하는 속성을 y축으로 선택해 드래그해주세요.

2. 시각화된 그래프를 확인하고 아래 실습을 진행해주세요.

상관 계수란 상관관계의 정도를 나타내는 수치로, 이는 -1 ~ +1 사이의 값을 가집니다.

y축으로 옮겨 놓은 속성을 아래에서 선택해 자전거 대역 간수와의 상관 계수를 확인해보고, 오른쪽 표를 토대로 상관계수를 해석해봅시다.

평균기온
평균수령
미세먼지(pm10)
강수량

대역 간수와 _____ 의 상관 계수

상관계수 크기	의미적인 범위
상관계수 = +1.0	완전 양의 상관관계
상관계수 = +0.8	강한 양의 상관관계
상관계수 = +0.6	중간 정도의 상관관계
상관계수 = +0.4	약한 양의 상관관계
상관계수 = +0.2	매우 약한 양의 상관관계

Step 3

그룹화한 데이터를 선 그래프 버튼을 눌러 시각화하기
기능 사용법

- ① 데이터, 그래프 조작 영역: 데이터 및 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.
- ② 속성 선택 버튼: 각 속성들을 선택하면 아래 영역에 상관계수를 표시하는 버튼입니다.



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 상관관계

1. 데이터간의 상관 계수 값을 히트맵으로 시각화해 봅시다.

1. 데이터수치의 상관 관계를 파악하고자하는 속성을 y축으로 지정하고 선택하세요.
2. 시각화된 그래프를 확인하고 아래 실습을 진행해보세요.

y축으로 옮겨 놓은 속성을 아래에서 선택해, 자전거 데이터수치의 상관계수 및 일부 히트맵으로 시각화 된 결과를 관찰해봅시다. 관찰이 끝난 후에는 전체 상관계수 히트맵 보기 버튼을 눌러 히트맵에 관한 설명을 확인해보도록 합니다.

2

속성: 필수기간 일당수령
미세먼지(pm10) 강수여부

데이터 간수와의 상관 계수

3

상관 계수 히트맵

데이터 간수

Step 4

상관관계를 산점도, 히트맵을 시각화하여 상관관계 시각화의 다각적 측면 이해하기

기능 사용법

① 데이터, 그래프 조작 영역: 데이터 및 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

② 속성 선택 버튼: 각 속성들을 선택하면 아래에 상관계수를 표시하고, 오른쪽 영역에 상관계수를 기반으로 시각화한 히트맵의 일부를 표시하는 버튼입니다.

③ 전체 상관계수 히트맵 확인하기 버튼: ②번 버튼은 히트맵으로 시각화된 일부만 확인할 수 있던 반면, 해당 버튼은 데이터 속성들에 대한 전체적인 히트맵을 확인하고 그에 대한 설명을 확인할 수 있습니다.

주의

실습 공간이 한정적이므로, 테이블의 크기, 위치 등을 임의로 조정할 경우 실습 진행에 방해가 될 수 있습니다. 가급적으로 주어진 실습 공간의 배치를 움직이지 말아주세요!

❖ 선생님을 위한 팁!

상관관계를 이용해 분석할 만한 데이터, 속성이 무엇이 있을지 탐구해 보는 시간을 가지면 상관관계 원리 이해에 더욱 도움이 됩니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

08차시. 데이터 분석방법

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-04] 동일한 데이터를 서로 다른 분석방법을 적용하여 분석 결과를 비교한다.
학습 목표	<ul style="list-style-type: none"> 탐색적 데이터 분석과 확증적 데이터 분석을 이해한다. 동일한 데이터로 서로 다른 분석방법을 적용해 분석결과를 비교한다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

- ❖ 주제명: 데이터 분석방법 알아보기
- ❖ 세부 주제
 - 1) 탐색적 데이터 분석
 - 2) 확증적 데이터 분석
- ❖ 콘텐츠 개발 목적: 분석방법에 따른 분석결과의 차이를 이해할 수 있다.
- ❖ 선생님을 위한 팁!

다양한 응용활동이 가능합니다.

실험실

- ❖ 주제명: 데이터 분석(1) - 건강데이터 살펴보기
- ❖ 세부 주제
 - 1) 학생 건강 검진 데이터를 이해하고 해석
 - 2) 특정 데이터로 데이터 분석을 하는 과정 실습
- ❖ 콘텐츠 개발 목적: 건강 검진 데이터를 사례로, 데이터를 이해하고, 시각화하면서 해석하는 과정을 통해 데이터 분석의 기초를 실습할 수 있도록 구성



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(1) - 건강데이터 살펴보기

학생 건강 검진 데이터를 직접 살펴봅시다.

학생들의 건강 검진 데이터를 오른쪽에서 확인하고, 왼쪽 입력란에 파악된 정보를 작성해 주세요.

① 데이터 속성은 총 몇개인가요? 건강 검진 데이터는 총 몇개인가요?
 ex. 5 ex. 1000

학생들이 받은 건강관련 항목은 무엇일까요?

성별 키 몸무게

병주령 속성은 무엇일까요?

성별 키 몸무게

② 제출

학생 건강검진 데이터			
성별	키	몸무게	병주령
남	170.0	65.0	1
남	175.0	70.0	1
남	180.0	75.0	1
남	185.0	80.0	1
남	190.0	85.0	1
남	195.0	90.0	1
남	200.0	95.0	1
남	205.0	100.0	1
남	210.0	105.0	1
남	215.0	110.0	1
남	220.0	115.0	1
남	225.0	120.0	1
남	230.0	125.0	1
남	235.0	130.0	1
남	240.0	135.0	1
남	245.0	140.0	1
남	250.0	145.0	1
남	255.0	150.0	1
남	260.0	155.0	1
남	265.0	160.0	1
남	270.0	165.0	1
남	275.0	170.0	1
남	280.0	175.0	1
남	285.0	180.0	1
남	290.0	185.0	1
남	295.0	190.0	1
남	300.0	195.0	1
남	305.0	200.0	1
남	310.0	205.0	1
남	315.0	210.0	1
남	320.0	215.0	1
남	325.0	220.0	1
남	330.0	225.0	1
남	335.0	230.0	1
남	340.0	235.0	1
남	345.0	240.0	1
남	350.0	245.0	1
남	355.0	250.0	1
남	360.0	255.0	1
남	365.0	260.0	1
남	370.0	265.0	1
남	375.0	270.0	1
남	380.0	275.0	1
남	385.0	280.0	1
남	390.0	285.0	1
남	395.0	290.0	1
남	400.0	295.0	1
남	405.0	300.0	1
남	410.0	305.0	1
남	415.0	310.0	1
남	420.0	315.0	1
남	425.0	320.0	1
남	430.0	325.0	1
남	435.0	330.0	1
남	440.0	335.0	1
남	445.0	340.0	1
남	450.0	345.0	1
남	455.0	350.0	1
남	460.0	355.0	1
남	465.0	360.0	1
남	470.0	365.0	1
남	475.0	370.0	1
남	480.0	375.0	1
남	485.0	380.0	1
남	490.0	385.0	1
남	495.0	390.0	1
남	500.0	395.0	1
남	505.0	400.0	1
남	510.0	405.0	1
남	515.0	410.0	1
남	520.0	415.0	1
남	525.0	420.0	1
남	530.0	425.0	1
남	535.0	430.0	1
남	540.0	435.0	1
남	545.0	440.0	1
남	550.0	445.0	1
남	555.0	450.0	1
남	560.0	455.0	1
남	565.0	460.0	1
남	570.0	465.0	1
남	575.0	470.0	1
남	580.0	475.0	1
남	585.0	480.0	1
남	590.0	485.0	1
남	595.0	490.0	1
남	600.0	495.0	1
남	605.0	500.0	1
남	610.0	505.0	1
남	615.0	510.0	1
남	620.0	515.0	1
남	625.0	520.0	1
남	630.0	525.0	1
남	635.0	530.0	1
남	640.0	535.0	1
남	645.0	540.0	1
남	650.0	545.0	1
남	655.0	550.0	1
남	660.0	555.0	1
남	665.0	560.0	1
남	670.0	565.0	1
남	675.0	570.0	1
남	680.0	575.0	1
남	685.0	580.0	1
남	690.0	585.0	1
남	695.0	590.0	1
남	700.0	595.0	1
남	705.0	600.0	1
남	710.0	605.0	1
남	715.0	610.0	1
남	720.0	615.0	1
남	725.0	620.0	1
남	730.0	625.0	1
남	735.0	630.0	1
남	740.0	635.0	1
남	745.0	640.0	1
남	750.0	645.0	1
남	755.0	650.0	1
남	760.0	655.0	1
남	765.0	660.0	1
남	770.0	665.0	1
남	775.0	670.0	1
남	780.0	675.0	1
남	785.0	680.0	1
남	790.0	685.0	1
남	795.0	690.0	1
남	800.0	695.0	1
남	805.0	700.0	1
남	810.0	705.0	1
남	815.0	710.0	1
남	820.0	715.0	1
남	825.0	720.0	1
남	830.0	725.0	1
남	835.0	730.0	1
남	840.0	735.0	1
남	845.0	740.0	1
남	850.0	745.0	1
남	855.0	750.0	1
남	860.0	755.0	1
남	865.0	760.0	1
남	870.0	765.0	1
남	875.0	770.0	1
남	880.0	775.0	1
남	885.0	780.0	1
남	890.0	785.0	1
남	895.0	790.0	1
남	900.0	795.0	1
남	905.0	800.0	1
남	910.0	805.0	1
남	915.0	810.0	1
남	920.0	815.0	1
남	925.0	820.0	1
남	930.0	825.0	1
남	935.0	830.0	1
남	940.0	835.0	1
남	945.0	840.0	1
남	950.0	845.0	1
남	955.0	850.0	1
남	960.0	855.0	1
남	965.0	860.0	1
남	970.0	865.0	1
남	975.0	870.0	1
남	980.0	875.0	1
남	985.0	880.0	1
남	990.0	885.0	1
남	995.0	890.0	1
남	1000.0	895.0	1
남	1005.0	900.0	1
남	1010.0	905.0	1
남	1015.0	910.0	1
남	1020.0	915.0	1
남	1025.0	920.0	1
남	1030.0	925.0	1
남	1035.0	930.0	1
남	1040.0	935.0	1
남	1045.0	940.0	1
남	1050.0	945.0	1
남	1055.0	950.0	1
남	1060.0	955.0	1
남	1065.0	960.0	1
남	1070.0	965.0	1
남	1075.0	970.0	1
남	1080.0	975.0	1
남	1085.0	980.0	1
남	1090.0	985.0	1
남	1095.0	990.0	1
남	1100.0	995.0	1
남	1105.0	1000.0	1
남	1110.0	1005.0	1
남	1115.0	1010.0	1
남	1120.0	1015.0	1
남	1125.0	1020.0	1
남	1130.0	1025.0	1
남	1135.0	1030.0	1
남	1140.0	1035.0	1
남	1145.0	1040.0	1
남	1150.0	1045.0	1
남	1155.0	1050.0	1
남	1160.0	1055.0	1
남	1165.0	1060.0	1
남	1170.0	1065.0	1
남	1175.0	1070.0	1
남	1180.0	1075.0	1
남	1185.0	1080.0	1
남	1190.0	1085.0	1
남	1195.0	1090.0	1
남	1200.0	1095.0	1
남	1205.0	1100.0	1
남	1210.0	1105.0	1
남	1215.0	1110.0	1
남	1220.0	1115.0	1
남	1225.0	1120.0	1
남	1230.0	1125.0	1
남	1235.0	1130.0	1
남	1240.0	1135.0	1
남	1245.0	1140.0	1
남	1250.0	1145.0	1
남	1255.0	1150.0	1
남	1260.0	1155.0	1
남	1265.0	1160.0	1
남	1270.0	1165.0	1
남	1275.0	1170.0	1
남	1280.0	1175.0	1
남	1285.0	1180.0	1
남	1290.0	1185.0	1
남	1295.0	1190.0	1
남	1300.0	1195.0	1
남	1305.0	1200.0	1
남	1310.0	1205.0	1
남	1315.0	1210.0	1
남	1320.0	1215.0	1
남	1325.0	1220.0	1
남	1330.0	1225.0	1
남	1335.0	1230.0	1
남	1340.0	1235.0	1
남	1345.0	1240.0	1
남	1350.0	1245.0	1
남	1355.0	1250.0	1
남	1360.0	1255.0	1
남	1365.0	1260.0	1
남	1370.0	1265.0	1
남	1375.0	1270.0	1
남	1380.0	1275.0	1
남	1385.0	1280.0	1
남	1390.0	1285.0	1
남	1395.0	1290.0	1
남	1400.0	1295.0	1
남	1405.0	1300.0	1
남	1410.0	1305.0	1
남	1415.0	1310.0	1
남	1420.0	1315.0	1
남	1425.0	1320.0	1
남	1430.0	1325.0	1
남	1435.0	1330.0	1
남	1440.0	1335.0	1
남	1445.0	1340.0	1
남	1450.0	1345.0	1
남	1455.0	1350.0	1
남	1460.0	1355.0	1
남	1465.0	1360.0	1
남	1470.0	1365.0	1
남	1475.0	1370.0	1
남	1480.0	1375.0	1
남	1485.0	1380.0	1
남	1490.0	1385.0	1
남	1495.0	1390.0	1
남	1500.0	1395.0	1
남	1505.0	1400.0	1
남	1510.0	1405.0	1
남	1515.0	1410.0	1
남	1520.0	1415.0	1
남	1525.0	1420.0	1
남	1530.0	1425.0	1
남	1535.0	1430.0	1
남	1540.0	1435.0	1
남	1545.0	1440.0	1
남	1550.0	1445.0	1
남	1555.0	1450.0	1
남	1560.0	1455.0	1
남	1565.0	1460.0	1
남	1570.0	1465.0	1
남	1575.0	1470.0	1
남	1580.0	1475.0	1
남	1585.0	1480.0	1
남	1590.0	1485.0	1
남	1595.0	1490.0	1
남	1600.0	1495.0	1
남	1605.0	1500.0	1
남	1610.0	1505.0	1
남	1615.0	1510.0	1
남	1620.0	1515.0	1
남	1625.0	1520.0	1
남	1630.0	1525.0	1
남	1635.0	1530.0	1
남	1640.0	1535.0	1
남	1645.0	1540.0	1
남	1650.0	1545.0	1
남	1655.0	1550.0	1
남	1660.0	1555.0	1
남	1665.0	1560.0	1
남	1670.0	1565.0	1
남	1675.0	1570.0	1
남	1680.0	1575.0	1
남	1685.0	1580.0	1
남	1690.0	1585.0	1
남	1695.0	1590.0	1
남	1700.0	1595.0	1
남	1705.0	1600.0	1
남	1710.0	1605.0	1
남	1715.0	1610.0	1
남	1720.0	1615.0	1
남	1725.0	1620.0	1




❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(1) - 건강데이터 살펴보기


건강검진을 받은 남학생과 여학생의 인원수를 살펴봅시다.

수집된 건강검진 데이터의 성별 비율이 차이가 크다면 분석할 때 편향이 발생할 수 있습니다. 데이터 편향이 존재하는 지 확인하기 위해 데이터에서 남학생과 여학생의 빈도수 분포를 막대 그래프로 살펴보겠습니다. 오른쪽 상단의 막대 그래프 보기 버튼을 클릭해 실습 결과에 따른 빈도수 분포를 확인하고, 편향 여부를 생각해보십시오.

1. '성별' 속성을 그래프로 드래그 하세요.
2. 그래프를 클릭해 우측 메뉴를 눌러 주세요.
 - '장을 막대로 변환' 옵션을 체크하세요.



②



Step 2

‘성별’ 속성을 막대그래프로 시각화하여 분포를 살펴보고, 데이터 편향 여부 검토하기

기능 사용법

① 막대그래프 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(1) - 건강데이터 살펴보기

건강검진을 받은 학생들의 키 데이터를 살펴봅시다.

데이터를 박스 플롯으로 시각화하면 분포와 중앙값, 이상치 등을 한눈에 파악할 수 있습니다. 키 데이터를 박스 플롯으로 살펴보면 학생들의 성장 패턴을 이해하고, 비정상적인 키 변화가 있는 학생들을 조기 식별하는 데 유용합니다. 오른쪽 상단의 박스 플롯 보기 버튼을 클릭해 박스 플롯을 확인해보세요.

1. '키' 속성을 그래프의 x축으로 드래그하세요.
2. 그래프를 클릭해 우측 메뉴를 불러 주세요.
 - '상자 그림'과 '이상치 보기' 옵션을 체크해주세요.
3. 박스 플롯에 마우스를 올려 아래 값을 찾아보세요.

- 1번째(Q1): 160.3
- 중앙값(Median): 162.4
- 3번째(Q3): 172.4
- 최후값(Max): 187.5
- 최전값(Min): 155.2

박스 플롯 보기

키 속성으로 만든 박스 플롯 그래프에서 이상치를 찾고, 그 값을 적어주세요.

ex. 160 cm ex. 160 cm 제출

Step 3

데이터의 '키' 속성을 박스플롯으로 시각화한 뒤 분포, 중앙값, 이상치 확인하기

기능 사용법

① 박스플롯 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

③ 데이터 입력란 및 제출 버튼: ②번 기능에서 수행한 내용을 토대로 답안을 숫자로 기입한 뒤, 해당 버튼을 누르면 정/오답 여부를 확인할 수 있습니다.

34



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(1) - 건강데이터 살펴보기

남학생과 여학생의 키 분포를 히스토그램으로 시각화하여 확인해봅시다. 1. 히스토그램 보기

히스토그램을 사용하면 각 성별의 키가 어떻게 분포되어 있는지 한눈에 볼 수 있으며, 빈도수와 경향성을 파악할 수 있습니다. 오른쪽 상단의 히스토그램 보기 버튼을 클릭해 실습을 진행해보세요.


1. '키' 속성을 그래프의 x축으로 드래그하세요.

2. '성별' 속성을 그래프의 y축으로 드래그하세요.

3. 그래프를 클릭해 우측 메뉴를 눌러주세요.

4. '히스토그램으로 변환' 옵션을 체크하세요.

2



3

여학생 중 가장 많이 분포된 키 범위는? 남학생 중 가장 많이 분포된 키 범위는?

E30-140

130-140

제출

Step 4

‘키’ 분포를 ‘성별에 따라 히스토그램으로 시각화하여 빈도수와 경향성 파악하기

기능 사용법

① 히스토그램 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

③ 데이터 입력란 및 제출 버튼: ②번 기능에서 수행한 내용을 토대로 답안을 선택한 뒤, 해당 버튼을 누르면 정/오답 여부를 확인할 수 있습니다.




❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(1) - 건강데이터 살펴보기

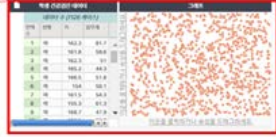
키와 몸무게의 상관 관계를 산점도로 시각화하여 확인해봅시다. ① 산점도 보기

산점도는 두 변수 간의 관계를 시각적으로 표현하는 데 유용한 도구입니다. 키와 몸무게 데이터를 산점도로 표현하면 각 데이터가 어떻게 분포하는지 볼 수 있으며, 키와 몸무게 사이에 어떤 상관관계가 있는지 파악할 수 있습니다. 오른쪽 상단의 산점도 보기 버튼을 클릭해 실습을 진행해보세요.

1. '키' 속성을 그래프의 x축으로 드래그하세요.
2. '몸무게' 속성을 그래프의 y축으로 드래그하세요.



②



③ 키와 몸무게의 상관관계는 어떤 특성을 가지고 있나요?

음의 상관관계
양의 상관관계
상관관계 거의 없음
제출

Step 5

‘키’와 ‘몸무게’ 속성을 산점도로 시각화하여 상관관계 파악하기
기능 사용법

① 산점도 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해보세요.

③ 데이터 입력란 및 제출 버튼: ②번 기능에서 수행한 내용을 토대로 답안을 체크한 뒤, 해당 버튼을 누르면 정/오답 여부를 확인할 수 있습니다.

주의

답안은 ex로 주어진 형식으로 작성해 주세요. 만약 'ex. 160'으로 표시된 입력란에 '160cm'라고 적을 경우 오답 처리될 수 있습니다.

❖ 선생님을 위한 팁!

실습 콘텐츠에서 언급하지 않은 다른 시각화 그래프 혹은 이전 차시에 배웠던 시각화의 중요성 등에 대해 추가 부연 설명을 하면서 실습을 진행하면 더욱 좋습니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

08차시. 데이터 분석방법

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과02-04] 동일한 데이터를 서로 다른 분석방법을 적용하여 분석 결과를 비교한다.
학습 목표	<ul style="list-style-type: none"> 탐색적 데이터 분석과 확증적 데이터 분석을 이해한다. 동일한 데이터로 서로 다른 분석방법을 적용해 분석결과를 비교한다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

- ❖ **주제명:** 데이터 분석방법 알아보기
- ❖ **세부 주제**
 - 1) 탐색적 데이터 분석
 - 2) 확증적 데이터 분석
- ❖ **콘텐츠 개발 목적:** 분석방법에 따른 분석결과의 차이를 이해할 수 있다.
- ❖ **선생님을 위한 팁!**
다양한 응용활동이 가능합니다.

실험실

- ❖ **주제명:** 데이터 분석(2) - 건강데이터로 비만 분포 확인하기
- ❖ **세부 주제**
 - 1) 데이터를 조합해 분석에 필요한 새로운 지표를 생성하는 방법 습득
 - 2) 분석에 필요한 데이터를 조합하여 시각화
- ❖ **콘텐츠 개발 목적:** 분석에 필요한 데이터를 조합하여 새로운 지표를 생성해 보고, 시각화하는 역량을 키우기 위한 실습



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(2) - 건강데이터로 비만 분포 확인하기

학생들의 건강 검진 데이터를 이용하여 비만 정도를 확인해봅시다.

우리나라 학생들의 비만도 분포 및 비만인 학생과 과체중인 학생의 비율을 알기 위해, 키와 몸무게 데이터를 이용해 비만도를 나타내는 체질량지수(BMI)를 데이터 속성으로 추가해보겠습니다.

1. 데이터 테이블을 클릭해 테이블의 오른쪽의 **비**으로 속성을 추가하고, '체속성'을 'BMI'로 수정하세요.
2. 'BMI' 속성명을 클릭한 후 '수식 편집' 메뉴를 누르세요.
3. 몸무게, 키 속성명은 **수식 편집** 버튼을 눌러 상단에 입력의 값이 수식을 작성해주세요.
수식: 몸무게 / ((키 / 100)² * (100))
* 수식만 입력하면 unit은 100으로 나타내지 않습니다.

BMI 구하는 공식 = 몸무게(kg) / 키(m)²

Step 1

키와 몸무게 속성값으로 체질량 지수 속성 추가하기
기능 사용법

① 데이터, 조작 영역: 데이터 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

데이터 준비와 분석 | 데이터 분석(2) - 건강데이터로 비만 분포 확인하기

체질량 지수(BMI)를 이용해 비만도 결과를 구해봅시다.

체질량 지수(BMI)는 다음과 같은 범위로 비만도를 분류합니다.

키, 몸무게 데이터로 구성된 BMI 값에 따라 비만도를 측정하는 속성을 데이터 테이블에 추가해 보겠습니다.

1. 데이터 테이블을 클릭한 뒤, 테이블 오른쪽의 **비**으로 속성을 추가하고 속성명을 클릭해 '비만도' 결과로 수정하세요.
2. '비만도 결과' 속성명을 클릭해 '수식 편집' 메뉴를 이용해, BMI값에 따라 비만도를 분류한 결과를 나타내는 수식을 작성해주세요.
수식: 키(BMI<18.5, "저체중", BMI<23, "정상", BMI<25, "과체중", "비만")

수식 쉽게 작성하는 방법
BMI 속성을 수식에 직접 입력 대신: **수식 편집** 버튼을 누르세요.
BMI 속성 선택에 단축키를 제공합니다.

Step 2

산점도로 속성 간 상관관계를 시각화해 보고, 상관관계에 대한 지표인 상관계수를 확인한 뒤 지표를 토대로 상관관계 해석

기능 사용법

① 데이터 조작 영역: 데이터 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.



❖ 실습 콘텐츠 안내:

데이터 준비와 분석 | 데이터 분석(2) - 건강데이터로 비만 분포 확인하기

비만도 결과를 시각화해봅시다.

막대그래프로 비만도 결과를 시각화하면 정상과 과체중 및 비만 학생의 비율을 직관적으로 비교할 수 있고, 각 범주의 분포를 명확하게 한눈에 파악할 수 있습니다. 오른쪽 상단의 **막대 그래프 보기** 버튼을 클릭해 실습을 진행해 보세요.

1. '비만도 결과' 속성값을 그래프로 드래그합니다.

2. 그래프를 클릭해 우측 메뉴를 눌러 주세요.

- '범주 막대로 변환' 옵션에 체크하세요.
- 1. 그래프를 클릭해 우측 메뉴를 눌러 주세요.
- 각 막대의 빈도수, 백분율 옵션을 체크하세요.

전체 학생들 중 과체중, 비만의 비율을 확인하고 답안을 제출해주세요.

③ 과제점 % 비만 % 제출

① 막대 그래프 보기

② 데이터, 그래프 조작 영역

Step 3

그룹화한 데이터를 선 그래프 버튼을 눌러 시각화하기
기능 사용법

① 막대 그래프 보기 버튼: 상단의 설명을 토대로 진행할 실습 공간을 여는 버튼입니다.

② 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

③ 데이터 입력란 및 제출 버튼: ②번 기능에서 수행한 내용을 토대로 답안을 체크한 뒤, 해당 버튼을 누르면 정/오답 여부를 확인할 수 있습니다.

주의

답안은 ex로 주어진 형식으로 작성해 주세요. 만약 'ex. 10'으로 표시된 입력란에 '10%'라고 적을 경우 오답 처리될 수 있습니다.

❖ 선생님을 위한 팁!

새롭게 생성한 지표를 가지고 다른 분석 주제를 선정할 수 있을지, 가능하다면 어떤 주제를 선정하여 어떤 분석 목표를 이룰 수 있을 지에 대한 토론을 실시하면 더욱 좋습니다.



3. 세부 콘텐츠

◆ 개요

09차시. 데이터 모델링 및 데이터 분석 도구 탐색하기

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-01] 데이터 모델 개념을 이해하고 데이터 분석에 활용할 수 있는 도구를 탐색한다.
학습 목표	<ul style="list-style-type: none"> • 데이터 모델 개념 이해 및 데이터간 관계 분석 및 상호 연관성을 파악한다. • 데이터 속성에 대한 유사성 측정, 분석의 원리 이해할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	1. (강의실) 동영상 콘텐츠

◆ 구성

강의실

- ❖ 주제명: 다양한 분석 도구 살펴보기
- ❖ 세부 주제
 - 1) 데이터 모델링의 특징
 - 2) 데이터 분석 도구
- ❖ 콘텐츠 개발 목적: 다양한 데이터 분석 도구를 소개한다.
- ❖ 선생님을 위한 팁!

데이터 분석 도구 예시로 설명된 사이트들을 추가적으로 실습해 볼 수 있다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

10차시. 단순 선형 회귀, 다중 선형 회귀

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-02] 동일한 데이터를 통계적 회귀모델과 기계학습을 통한 회귀모델로 분석하여 결과 해석 내용을 비교한다.
학습 목표	<ul style="list-style-type: none"> 회귀 분석의 개념에 대해 설명할 수 있다. 단순 선형 회귀, 다중 선형 회귀 분석의 개념을 설명할 수 있다. 예에 따른 회귀 분석 방법을 선택할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

❖ **주제명:** 선형 회귀의 기본 이해

❖ **세부 주제**

1) 단순 선형 회귀

2) 다중 선형 회귀

❖ **콘텐츠 개발 목적:** 회귀 분석의 개념을 알고 실행할 수 있다.

❖ **선생님을 위한 팁!**

상황에 따라 단순 선형 회귀와 다중 선형 회귀 중 어떤 것을 적용해야 하는지 설명할 수 있습니다.

실험실

❖ **주제명:** 회귀분석으로 예측하기

❖ **세부 주제**

1) 회귀분석에 사용되는 독립변수, 종속변수 이해

2) 회귀 모델 생성 과정 이해 및 예측 오차 최소화 해보기

❖ **콘텐츠 개발 목적:** 단순 선형 회귀 방식으로 키에 기반한 몸무게를 예측하며 회귀 분석의 원리 이해

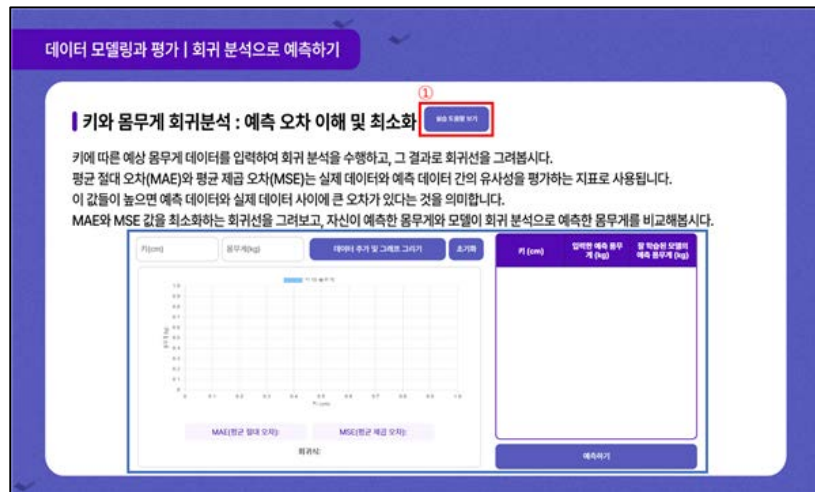


❖ 실습 콘텐츠 안내:



Step 1

드래그 앤 드롭으로 회귀 분석을 위한 종속 변수/독립변수 분류



Step 2

회귀 분석으로 키에 따른 몸무게 예측해보기

기능 사용법

① 실습 도움말 보기 버튼: 실습을 진행하기 위한 가이드를 차례대로 제공하는 버튼입니다. 실습을 진행하기 전 해당 버튼을 눌러 실습 과정을 이해해보세요.

주의

4 페이지에서, 키와 예측 몸무게를 입력할 때 키는 100(cm) 이상, 몸무게는 10(kg) 이상 입력해야 합니다.

❖ 선생님을 위한 팁!

회귀 분석에 사용되는 회귀식에 대한 추가 설명을 해주신다면 학생들의 회귀분석 이해에 더욱 도움이 됩니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

11차시. 회귀 분석, 결정 계수

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-02] 동일한 데이터를 통계적 회귀모델과 기계학습을 통한 회귀모델로 분석하여 결과 해석 내용을 비교한다.
학습 목표	<ul style="list-style-type: none"> 성능 지표(평균제곱오차, 결정계수)의 개념을 알고 값을 구할 수 있다. SI 모델의 적합도를 확인할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

- ❖ 주제명: R^2 값 알아보기
- ❖ 세부 주제
 - 1) 회귀 분석
 - 2) 결정 계수
- ❖ 콘텐츠 개발 목적: 동일한 데이터셋을 활용하더라도 데이터 전처리 과정, 데이터 선택, 분석방법에 따라 모델링 결과가 다를 수 있음을 인지한다.
- ❖ 선생님을 위한 팁!

R^2 값에 대한 부연설명이 필요할 수 있다.

실험실

- ❖ 주제명: 회귀분석으로 예측하기
- ❖ 세부 주제
 - 1) 회귀분석에 사용되는 변수 선정하기
 - 2) 선정한 변수에 따라 달라지는 결정계수 이해하기
- ❖ 콘텐츠 개발 목적: 더 정확하고 정교한 회귀 분석 인공지능(AI) 모델을 만들기 위해 속성을 선택하는 기준들을 학습



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 회귀분석과 결정계수

주책 중위 가격 예측 AI 모델학습을 위해, 선형 회귀분석에 사용할 변수를 골라봅시다.

아래 자료는 캘리포니아 주택 데이터셋을 히스토그램으로 시각화한 그래프와 상관계수입니다. 상관계수가 -1 또는 1에 가까울수록 변수 간에 강한 상관관계 및 선형 관계가 있다는 것을 의미 합니다. 이를 토대로 AI 모델 학습을 위해 선형 회귀 분석에 사용할 변수 2가지를 선택해보세요.

②

<p>면적</p> <p>상관계수: -0.64967</p>	<p>방문객수</p> <p>상관계수: -0.14490</p>	<p>중위 소득</p> <p>상관계수: 0.88875</p>	<p>방문객 개수</p> <p>상관계수: -0.15413</p>	<p>방문객 수</p> <p>상관계수: 0.64968</p>
<p>방문객 수</p> <p>상관계수: -0.02000</p>	<p>방문객 개수</p> <p>상관계수: 0.00741</p>	<p>건축물 용량 면적</p> <p>상관계수: 0.70573</p>	<p>주책 중위 가격</p> <p>상관계수: 1.00000</p>	

기타의 다양한 다른 변수는, 주책 중위 가격과 가 수식에 대한 상관 관계를 나타내지 않습니다. 주책 중위 가격과 높은 상관 관계를 가지는 유일한, 예측 AI 모델 학습에 중요한 변수가 될 수 있습니다. 이를 고려하여 변수를 선택해 보세요.

Step 1

드래그 앤 드롭으로 회귀 분석을 위한 종속 변수/독립변수 분류
기능 사용법

- ① 도움말 보기 버튼: 변수들에 대한 설명을 확인할 수 있는 버튼입니다.
- ② 아이콘 선택 기능: 설명을 따라 올바른 변수들을 두 가지 선택합니다.
선택한 변수 아이콘은 밝은색으로 표시되며, 두 가지를 모두 선택할 경우
정/오답 여부가 판별됩니다.



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 회귀분석과 결정계수

선택된 변수에 따라 달라지는 회귀 분석의 결정 계수를 확인해봅시다.

결정계수는 우리가 만든 예측 모델이 실제로 관측된 데이터를 얼마나 잘 설명하는지를 나타내는 수치입니다. 결정계수의 값이 1에 가까우면, 모델이 데이터를 매우 잘 반영하고 있음을 의미하며, 0에 가까우면 모델의 설명력이 매우 낮음을 의미합니다.

1) 독립 변수를 선택해 보세요.

①
건축물 중위 연령
중위 소득
면적 당 가구수
방의 총 개수

2) 아래 '학습하기' 버튼을 눌러 학습에 따른 결정 계수를 확인해 보세요.

②
학습 하기

결정계수 :

TIP

만약 결정계수가 1에 가까우면, 모델이 데이터를 매우 잘 설명한다는 뜻입니다. 하지만, 이 숫자가 높다고 해서 항상 좋은 것은 아닙니다. 가끔 모델이 학습 데이터에만 너무 잘 맞춰져 있어서 새로운 데이터로 다른 상황에서는 예측력이 떨어지는 경우가 있습니다. 이를 '과적합'이라고 합니다.

Step 2

변수를 선택하고 달라지는 결정계수 확인하기
기능 사용법

- ① 아이콘 선택 기능: 회귀 분석에 사용할 독립 변수를 선택합니다. 최소 한 가지 ~ 최대 네 가지의 변수를 선택할 수 있습니다. (다중 선택이 가능)
- ② 학습하기 버튼: 해당 버튼을 누르면 ①번 기능으로 선택한 변수들에 대한 결정계수를 확인할 수 있습니다.

주의

3페이지에서, 요소를 선택할 때 TIP이 꺼지지 않은 상태일 경우 선택이 잘 되지 않을 수 있습니다. 선택이 잘 되지 않는다면 TIP을 닫아주세요.

❖ 선생님을 위한 팁!

이전 실습에서 학습한 MAE, MSE 지표와 결정계수에 관계를 추가 설명해주면 회귀분석 이해에 더욱 도움을 줄 수 있습니다.





3. 세부 콘텐츠

◆ 개요

12차시. 군집화 개념 이해하기

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-03] 데이터의 속성에 대한 유사성을 측정하고 분석하여 군집을 형성하고, 군집 분석 결과의 의미를 해석한다. [12데과03-04] 데이터 간의 관계를 분석하고 상호 연관성을 파악하여 결과의 의미를 해석한다.
학습 목표	<ul style="list-style-type: none"> • 데이터 속성에 따라 유사도가 높은 데이터를 묶을 수 있다. • 데이터의 속성에 따라 다수의 군집으로 나누고 군집 내 유사성과 군집 간 상이성을 이해한다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> 1. (강의실) 동영상 콘텐츠 2. (실험실) 실습형 콘텐츠

◆ 구성

강의실

- ❖ **주제명:** 보로노이 다이어그램과 센트로이드
- ❖ **세부 주제**
 - 1) 보로노이 다이어그램 활용 예시
 - 2) 센트로이드 설정 과정
- ❖ **콘텐츠 개발 목적:** 센트로이드 개념을 설명하기 위해 보로노이 다이어그램으로 흥미를 유발한다.
- ❖ **선생님을 위한 팁!**
보로노이 다이어그램 실습을 위해서는 하얀 그릇, 색상별 초콜릿 물을 준비하셔야 합니다.

실험실

- ❖ **주제명:** SNS 사용 데이터를 활용한 군집화 알아보기
- ❖ **세부 주제**
 - 1) 데이터 그룹화 방법 이해
 - 2) 데이터 그룹별 시각화 방법 이해
- ❖ **콘텐츠 개발 목적:** 데이터를 원하는 목적에 따라 그룹으로 분류하고, 시각화하여 군집 유형별 의미 유추해 보기



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | SNS사용 데이터를 활용한 군집화 알아보기

히스토그램으로 시각화한 데이터를, 나이 속성별로 그룹화해서 확인해보시다.

SNS를 통한 정보 습득을 목적으로 하는 사용자들의 SNS 사용 시간을 나이에별로 그룹화해 살펴봅시다.
아래 설명을 따라 실습을 진행해보세요.

1. 테이블에서 '나이' 속성을 그래프에 드래그하여 추가하고 변화를 확인해 보세요.
2. 그래프를 클릭해 우측 메뉴를 눌러보세요.

- 범례의 색상을 바꾸는 옵션을 선택하여 더욱 눈에 띄게 시각화 해보세요.

Step 1

히스토그램으로 시각화한 데이터를 나이 속성별로 그룹화
기능 사용법

① 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다.
말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.

데이터 모델링과 평가 | SNS사용 데이터를 활용한 군집화 알아보기

소셜 네트워크 서비스 (SNS) 사용 데이터를 먼저 살펴봅시다.

아래는 SNS 사용에 관한 설문 조사를 통해 수집된 데이터입니다. 데이터를 살펴보기 위해, 아래 설명을 따라 실습을 진행해보세요.

1. 아래 테이블에서 SNS 사용 시간 속성명을 클릭하고, 오른쪽순 정렬을 선택 하세요.
2. 그래프 x축에 'SNS 사용 시간' 속성을 드래그 하세요.
3. 그래프 y축에 'p1_정보습득' 속성을 그래프로 드래그하세요.
4. 그래프를 한번 클릭하여 우측 메뉴를 눌러주세요.

- '히스토그램으로 변환' 옵션을 선택해 시각화 해보세요.
- 4번 그래프 메뉴를 다시 눌러보면, 글자(침사이 간격)을 조절해 시각화 단위를 변경할 수 있습니다.

Step 2

정보 습득을 위해 sns를 사용하는 사용자들의 사용 시간을 그룹화
기능 사용법

① 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다.
말풍선을 따라 해당 영역을 조작하여 실습을 수행해 보세요.



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | SNS사용 데이터를 활용한 군집화 알아보기

히스토그램으로 시각화한 데이터를, 중독 속성별로 그룹화해서 확인해보십시오.

이번에는 SNS를 통한 정보 습득을 목적으로 하는 사용자들의 SNS 사용 시간을 중독 속성을 기준으로 그룹화한 데이터를 살펴봅시다. 아래 설명을 따라 실습을 진행해보십시오.

1. 테이블에서 's3' 중독 속성을 그래프 중간에 드래그하여 추가하고 변화를 확인해 보세요.
2. 그래프를 클릭해 우측 메뉴를 눌러주세요.
 - 메뉴의 색상을 바꾸는 옵션을 선택하여 더욱 눈에 띄게 시각화해보세요.
 - 메뉴바에서, 자 오(인)출(정)을 누르면 반드시 수 배분율도 계산할 수 있습니다.

번호	성명	성별	나이	SNS 사용 시간	카카오톡	인스타그램	네이버	블로그	유튜브	스카이프	이메일	스카이프	이메일	스카이프	이메일	스카이프	이메일	스카이프	이메일
1	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
2	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
3	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
4	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
5	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
6	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
7	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
8	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
9	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2
10	김민준	남자	17	30.0%	1	4	2	2	2	2	2	2	2	2	2	2	2	2	2

결과 예시 및 설명 확인

Step 3

Step 2에서 그룹화 했던 데이터를 중독 속성별로 시각화 기능 사용법

- ① 데이터, 그래프 조작 영역: 데이터, 그래프 조작이 가능한 영역입니다. 말풍선을 따라 해당 영역을 조작하여 실습을 수행해보세요.
- ② 결과 예시 및 설명 확인 버튼: 지금까지 단계별로 시각화 해왔던 실습에 대한 결과 예시 및 그에 대한 설명을 확인할 수 있는 버튼입니다.

주의 실습 공간이 한정적이므로, 테이블의 크기, 위치 등을 임의로 조정할 경우 실습 진행에 방해가 될 수 있습니다. 가급적으로 주어진 실습 공간의 배치를 움직이지 말아주세요!

❖ 선생님을 위한 팁!

해당 데이터를 그룹화 하는 이유와, 그룹화 해서 확인하면 알 수 있는 점을 군집 분석과 연계하여 설명하면 더욱 좋습니다.



3. 세부 콘텐츠

◆ 개요

13차시. 군집분석

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-03] 데이터의 속성에 대한 유사성을 측정하고 분석하여 군집을 형성하고, 군집 분석 결과의 의미를 해석한다. [12데과03-04] 데이터 간의 관계를 분석하고 상호 연관성을 파악하여 결과의 의미를 해석한다.
학습 목표	<ul style="list-style-type: none"> 데이터 간 관계를 분석하고 데이터 분석에 활용할 수 있는 도구를 탐색한다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

◆ 구성

강의실

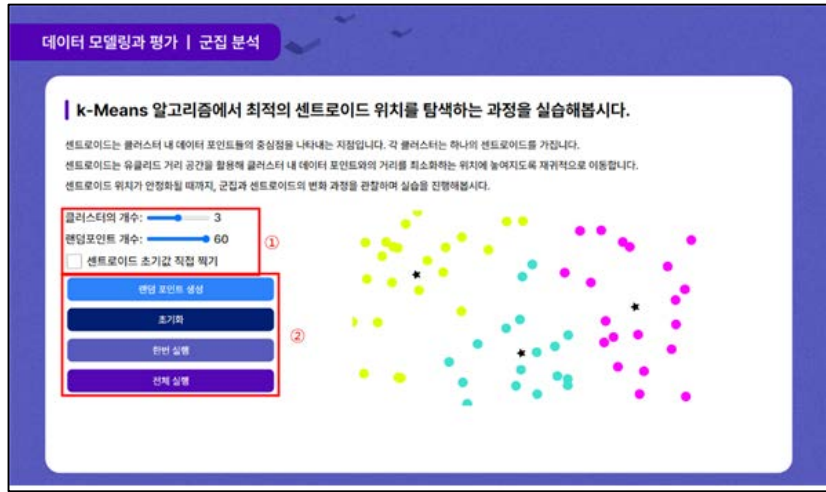
- ❖ 주제명: 포켓몬빵 사례로 알아보는 군집분석
- ❖ 세부 주제
 - 1) 군집분석, 군집간 유사성
 - 2) 군집간 상이성 결과 해석 및 설명
- ❖ 콘텐츠 개발 목적: 군집분석 사례를 통해 데이터 분석에 대한 이해를 돕는다.
- ❖ 선생님을 위한 팁!
포켓몬빵 데이터를 활용하여 활동을 구성할 수 있습니다.

실험실

- ❖ 주제명: 군집분석
- ❖ 세부 주제
 - 1) k-Means 알고리즘으로 군집 분석 이해하기
 - 2) k-Means 알고리즘과 보로노이 다이어그램의 유사성 확인하기
- ❖ 콘텐츠 개발 목적: k-Means 알고리즘을 기반으로 군집 분석 원리 및 과정 이해하기



❖ 실습 콘텐츠 안내:



Step 1

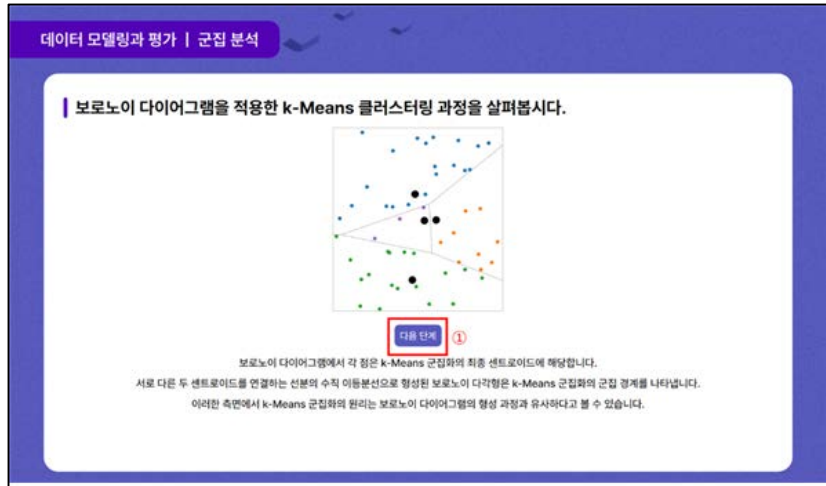
데이터 생성 및 센트로이드 탐색 관련 설정 후 알고리즘 원리 파악 기능 사용법

① 센트로이드 관련 초기 설정: 클러스터의 개수 및 랜덤포인트의 개수, 센트로이드 초기값 설정 여부 등을 확인해주세요. 센트로이드 초기값 직접 찍기 모드를 활성화할 경우 알고리즘이 랜덤하게 센트로이드 위치를 설정하는 것이 아닌, 클릭한 위치를 기반으로 k-means 알고리즘이 실행됩니다.

② 센트로이드 탐색 관련 버튼: 랜덤 포인트 생성 버튼으로 분류할 데이터를 생성할 수 있습니다. 초기화 버튼으로 보드에 생성된 데이터를 초기화할 수 있습니다. 한 번 실행 버튼으로 k-means 알고리즘을 단계별로, 전체 실행 버튼으로 센트로이드 값이 결정될 때까지 k-means 알고리즘을 실행할 수 있습니다.



❖ 실습 콘텐츠 안내:



Step 2

보로노이 다이어그램과 k-means 클러스터링 비교
기능 사용법

① 다음 단계 버튼: 해당 버튼을 클릭하면 보로노이 다이어그램이 형성되는 과정을 단계별로 확인할 수 있습니다.

주의

센트로이드 초기값을 직접 짚을 경우, 짚을 수 있는 센트로이드의 개수는 '클러스터의 개수' 슬라이더를 통해 설정한 값과 동일합니다.

❖ 선생님을 위한 팁!

상기 실습을 반복하며 센트로이드의 초기 위치를 다양한 지점에 설정한 뒤 k-means 알고리즘을 적용하여 적절한 초기값의 조건을 고려할 수 있도록 유도합니다.



3. 세부 콘텐츠

◆ 개요

14차시. 지지도, 신뢰도, 향상도

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-04] 데이터 간의 관계를 분석하고 상호 연관성을 파악하여 결과의 의미를 해석한다.
학습 목표	<ul style="list-style-type: none"> 연관 분석의 개념을 알고 설명할 수 있다. 지지도, 신뢰도, 향상도를 알고 설명할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

◆ 구성

강의실

- ❖ 주제명: 연관분석 기초
- ❖ 세부 주제
 - 1) 지지도
 - 2) 신뢰도
 - 3) 향상도
- ❖ 콘텐츠 개발 목적: 기저귀와 맥주 판매 추이 예시를 통해 연관분석을 설명한다.
- ❖ 선생님을 위한 팁!
공유데이터를 활용하여 연관분석을 실시해보는 활동을 할 수 있습니다.

실험실

- ❖ 주제명: 지지도, 신뢰도, 향상도
- ❖ 세부 주제
 - 1) 지지도, 신뢰도, 향상도의 의미 및 구하는 방법 습득
 - 2) 연관 분석 지표 해석
- ❖ 콘텐츠 개발 목적: 지지도, 신뢰도, 향상도를 실습을 기반으로 연관 분석 수행 방법 습득



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 지지도(support), 신뢰도(confidence), 향상도(lift)

매점 판매 데이터를 기반으로, 연필과 사과 주스의 지지도를 구해보세요.

지지도(Support)는 특정 상품을 동시에 구매하는 비율이 전체 거래에서 차지하는 비율을 나타냅니다.
공식을 기반으로, 연필과 사과 주스가 동시에 포함된 거래 데이터 행의 '선택' 옵션을 체크하고 제출해보세요.

지지도 공식: $Support(A \rightarrow B) = \frac{A와 B가 함께 구매한 거래 수}{전체 거래 수}$

번호	연필	사과주스	선택	입력
1	선택	연필	<input type="checkbox"/>	<input type="checkbox"/>
2	연필	연필	<input type="checkbox"/>	<input type="checkbox"/>
3	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
4	사과주스	연필	<input type="checkbox"/>	<input type="checkbox"/>
5	사과주스	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
6	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
7	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
8	사과주스	연필	<input type="checkbox"/>	<input type="checkbox"/>
9	연필	연필	<input type="checkbox"/>	<input type="checkbox"/>
10	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
11	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
12	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
13	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
14	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>
15	연필	사과주스	<input type="checkbox"/>	<input type="checkbox"/>

지지도 =

전체 거래수 (15)

제출

Step 1

연필과 사과 주스의 지지도 계산하기

기능 사용법

① 행 선택 기능: 조건에 해당하는 행을 선택하는 기능입니다. 표의 맨 오른쪽부분의 박스를 체크하면 선택이 이루어지고, 오른쪽 분자 영역에 데이터 행이 추가 됩니다.

② 제출 버튼: 분자 영역에 추가된 데이터에 대한, 정/오답 여부를 확인하는 버튼입니다. 정답일 경우, 해당 조건에 대한 지지도를 함께 안내합니다.



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 지지도(support), 신뢰도(confidence), 향상도(lift)

매점 판매 데이터를 기반으로, 연필 -> 사과 주스 구매에 대한 신뢰도를 구해보세요.

신뢰도(Confidence)는 원인이 발생할 때 결과가 발생할 가능성을 나타내는 지표입니다. 이 경우에는 연필을 구매했을 때 사과주스를 추가로 구매할 확률을 의미합니다. 공식을 기반으로, 연필이 포함된 거래 데이터 행의 '선택' 옵션을 체크하고 제출해보세요.

신뢰도 공식: $Confidence(A \rightarrow B) = \frac{A와 B가 함께 구매한 거래 수}{A가 구매한 거래 수}$

번호	연필	사과 주스	선택	체크
1	연필	향신료	<input type="checkbox"/>	<input type="checkbox"/>
2	연필	향신료	<input type="checkbox"/>	<input type="checkbox"/>
3	연필	향신료	사과주스	<input type="checkbox"/>
4	향신료	향신료	사과주스	<input type="checkbox"/>
5	향신료	향신료	향신료	<input type="checkbox"/>
6	연필	향신료	<input type="checkbox"/>	<input type="checkbox"/>
7	연필	향신료	향신료	<input type="checkbox"/>
8	향신료	향신료	향신료	<input type="checkbox"/>
9	연필	향신료	향신료	<input type="checkbox"/>
10	연필	향신료	향신료	<input type="checkbox"/>
11	연필	향신료	향신료	<input type="checkbox"/>
12	연필	향신료	향신료	<input type="checkbox"/>
13	연필	향신료	향신료	<input type="checkbox"/>
14	연필	향신료	향신료	<input type="checkbox"/>
15	연필	향신료	향신료	<input type="checkbox"/>

①

연필, 사과주스 동시 포함 거래 수 (3)

신뢰도 =

분자

분모

제출

②

Step 2

연필 -> 사과 주스 구매에 대한 신뢰도 계산하기
기능 사용법

① 행 선택 기능: 조건에 해당하는 행을 선택하는 기능입니다. 표의 맨 오른쪽부분의 박스를 체크하면 선택이 이루어지고, 오른쪽 분모 영역에 데이터 행이 추가 됩니다.

② 제출 버튼: 분모 영역에 추가된 데이터에 대한, 정/오답 여부를 확인하는 버튼입니다. 정답일 경우, 해당 조건에 대한 신뢰도를 함께 안내합니다.



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 지지도(support), 신뢰도(confidence), 향상도(lift)

매점 판매 데이터를 기반으로, 연필 구매와 사과주스 추가 구매 간 향상도를 구해보세요.

향상도(Lift)는 거래 데이터에서 A를 구매했을 때 B를 구매할 가능성에 관한 지표로, 향상도가 1이려면 A와 B가 확률적으로 독립에 가까움을 의미합니다. 공식을 기반으로, 연필을 포함한 거래 및 사과주스를 포함한 거래 데이터 행의 '선택' 옵션을 체크하고 제출해보세요.

향상도 공식: $Lift(A, B) = \frac{\text{전체 거래 수} \times A \text{와 } B \text{가 함께 구매한 거래 수}}{A \text{가 구매한 거래 수} \times B \text{가 구매한 거래 수}}$

행	연필	사과주스	선택	연필	선택
1	갑자입	탄산음료			<input type="checkbox"/>
2	피자왕	탄산음료			<input type="checkbox"/>
3	연필	초콜릿	사과주스	물기유유	<input type="checkbox"/>
4	초코유유	피자왕	사과주스		<input type="checkbox"/>
5	피우개	초코유유	피자왕	탄산음료	<input type="checkbox"/>
6	연필	피우개			<input type="checkbox"/>
7	세탁용	갑자입	초콜릿		<input type="checkbox"/>
8	사과주스				<input type="checkbox"/>
9	세탁용	연필	사과주스		<input type="checkbox"/>
10	갑자입	사과주스			<input type="checkbox"/>
11	연필왕자	물기유유			<input type="checkbox"/>
12	갑자입	물기유유			<input type="checkbox"/>
13	연필왕자	초코유유	탄산음료		<input type="checkbox"/>
14	연필	연필왕자	초콜릿		<input type="checkbox"/>
15	연필	사과주스			<input type="checkbox"/>

①

전체 거래 수 (15) x 연필, 사과주스 동시 포함 거래 수 (3)

연필	사과주스	분자
	x	

② 제출

Step 3

연필 구매와 사과주스 추가 구매 간 향상도 계산하기
기능 사용법

① 행 선택 기능: 조건에 해당하는 행을 선택하는 기능입니다. 표의 맨 오른쪽부분의 박스를 체크하면 선택이 이루어지고, 오른쪽 분모 영역에 각각 데이터 행이 추가 됩니다.

② 제출 버튼: 분모 영역에 추가된 데이터에 대한, 정/오답 여부를 확인하는 버튼입니다. 정답일 경우, 해당 조건에 대한 향상도를 함께 안내합니다.

주의

향상도를 구하는 실습에서 두가지 조건이 있으나, 선택은 나누어서 하지 않습니다. 두 조건에 전부 해당하는 것은 양쪽에 추가가 되고, 조건 하나에만 해당하는 것은 각 조건의 영역에 개별로 추가됩니다.

❖ 선생님을 위한 팁!

다른 조건을 가정하여 지지도, 신뢰도, 향상도 식을 기반으로 값을 도출하는 실습을 진행해도 좋습니다.



3. 세부 콘텐츠

◆ 개요

◆ 구성

15차시. 장바구니 분석

성취기준, 학습목표, 학습내용, 개발유형

성취기준	[12데과03-04] 데이터 간의 관계를 분석하고 상호 연관성을 파악하여 결과의 의미를 해석한다.
학습 목표	<ul style="list-style-type: none"> 장바구니 분석을 해볼 수 있다. 연관 분석의 지표들을 해석하여 적용할 수 있다.
대상 학년	
연계교육과정	정보교과
세부콘텐츠	<ol style="list-style-type: none"> (강의실) 동영상 콘텐츠 (실험실) 실습형 콘텐츠

강의실

- ❖ 주제명: 장바구니 분석 이해하기
- ❖ 세부 주제
 - 1) 아프리오리(Apriori) 알고리즘
 - 2) 데이터 분석 - 최소 지지도 선정
- ❖ 콘텐츠 개발 목적: 연관 분석 지표들을 바탕으로 상품을 진열할 수 있다.
- ❖ 선생님을 위한 팁!

데이터 수집의 중요성 차시와 연계하여 활용가능한 데이터를 찾을 수 있는 사이트를 추가적으로 안내할 수 있다.

실험실

- ❖ 주제명: 장바구니 분석
- ❖ 세부 주제
 - 1) 신뢰도 및 향상도에 대한 해석 실습
 - 2) 아프리오리 알고리즘 적용 데이터 분석
- ❖ 콘텐츠 개발 목적: 데이터를 기반으로 실생활에 유의미한 인사이트를 도출하는 과정 실습



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 장바구니 분석

| 구매 데이터가 10,000건 있을 때, 최소 100건 이상 거래된 구매 패턴을 골라봅시다.

거래 건수를 직접 세지 않아도, 전체 거래 중에서 특정 구매 패턴이 나타나는 비율을 '지지도'를 통해 확인할 수 있습니다. 10,000건의 구매 행위 중 100건의 특정 거래는 1%에 해당합니다. 따라서 최소 100건의 거래가 발생한 구매 패턴을 선정하는 것은 분석에 필요한 최소 지지도가 0.01임을 의미합니다. 아래 표에서 지지도가 0.01 이상인 구매 패턴을 분석 대상으로 선택해봅시다.

선택	상품 A	상품 B	신뢰도	향상도	지지도
<input type="checkbox"/>	식량	화나스	1.13	1.923	0.013
<input type="checkbox"/>	식량	건강바람	5.1	1.631	0.015
<input type="checkbox"/>	화나스	배채도	1.21	1.003	0.012
<input type="checkbox"/>	배채도	화나스	1.11	1.002	0.012
<input type="checkbox"/>	식량	배채도	1.20	1.630	0.008
<input type="checkbox"/>	건강바람	건강바람	1.14	1.235	0.007
<input type="checkbox"/>	화나스	건강바람	1.23	1.006	0.006
<input type="checkbox"/>	화나스	건강바람	1.25	1.234	0.011
<input type="checkbox"/>	건강바람	배채도	1.32	1.135	0.018
<input type="checkbox"/>	배채도	식량	1.08	1.217	0.009

계속

Step 1

지지도 기반 분석 대상 구매 패턴 추측

기능 사용법

① 데이터 선택: '선택' 열에 있는 체크박스를 클릭하여, 발문의 조건에 해당하는 구매 패턴을 분석 대상으로 선택할 수 있습니다.

② 제출 버튼: 체크박스에 선택된 항목을 토대로 정/오답 여부를 판별하는 버튼입니다.



❖ 실습 콘텐츠 안내:

데이터 모델링과 평가 | 장바구니 분석

신뢰도, 향상도를 직접 해석하여 빵집 매대를 진열해봅시다.

아래 표는 빵집 손님들의 구매 데이터에서 최소 지지도 기준을 만족하는 구매 패턴을 바탕으로 산출된 신뢰도와 향상도를 보여줍니다. 손님들이 여러 개의 빵을 동시에 구매할 수 있는 최적의 순서로 빵들을 진열하고, 진열 완료 버튼을 눌러주세요.

A	B	신뢰도	향상도
식빵	도넛	1.13	1.923
식빵	크로와상	1.1	1.631
도넛	크로와상	1.25	1.234
크로와상	바게트	1.32	1.135
도넛	바게트	1.21	1.003
바게트	도넛	1.11	1.002

Tip
 A→B의 신뢰도 : A품목 구매자들 사이에서 B품목의 인기를 나타내는 척도입니다.
 A→B의 향상도 : A품목과 B품목의 구매가 얼마나 서로 의존적인지를 나타내는 척도입니다.

Step 2

신뢰도와 향상도를 고려하여 빵 진열 순서 정하기

기능 사용법

- ① 빵 기본 위치: 진열해야 할 빵이 있는 기본 위치입니다.
- ② 빵 진열 위치: ①에서 드래그하여 빵을 순서대로 진열해야 할 위치입니다. 진열을 잘못했을 경우, 진열대에서 각 빵의 위치를 교체할 수 있습니다.
- ③ 진열 완료 버튼: 진열 완료 시 해당 버튼을 클릭하여 적절한 순서로 진열이 이루어졌는지 확인할 수 있습니다.

주의

빵을 진열하는 순서를 정할 때는 신뢰도와 향상도 두 가지를 모두 고려해야 합니다.

❖ 선생님을 위한 팁!

실습 콘텐츠에서 제시된 경우의 수 외에 학생들이 다른 경우의 수를 스스로 생각해 보고, 해당 경우의 수에 대해서는 어떤 식으로 진열해야 할지 자유롭게 토론할 수 있도록 유도합니다.