

<1장>

페이지	행	기존	변경
33	박스	사실상 데이터 마이닝을 전제 조건으로 한다.	사실상 데이터 마이닝의 필요조건이다.
33~34	마지막 행	데이터에 대한 판매자는 파격적인 할인가로 타겟팅한 예비 임신부 기저귀, 젖병, 아기 장난감들을 계속 구매해 단골 고객이 되길 바란다.	판매자는 타겟팅한 예비 임신부가 기저귀, 젖병, 아기 장난감을 계속 구매해 단골 고객이 되길 바란다.
34	5행	홍보 메일에 이런 데이터를 사용한 후 자신이 10대 딸이 출산 물품에 대한 쿠폰을 받은 이유를 판매자에게 듣고자 화가 난 딸의 아버지가 연락해 왔다.	홍보 메일에 이런 데이터를 사용한 후, 자신의 10대 딸이 왜 출산 물품에 관한 쿠폰을 받았는지 알고자 하는 화가 난 아버지에게서 연락이 왔다.
36	아래에서 5행	관찰, 기억 장치를 활용하고 추론하기 위한 사실적 근거를 제공할 수 있다.	미래의 추론을 위한 사실적 근거를 제공하기 위해 관찰, 기억 공간을 활용한다.
37	7-9행	일반적으로 사용하는 개념도와 개요를 만드는 학습 전략은 기계가 지식 추상화를 하는 방법과 유사하다.	개념화하고 개요를 만드는 일반적인 학습 전략은 기계가 지식을 추상화하는 방법과 유사하다.
37	12-16행	학습 전략과 시험 결과가 ~ 좋은 선생님은 이점에 꽤 도움이 된다.	시험지를 채점하거나 학습 전략이 적절했는지를 확인할 때는 항상 긴장이 된다. 선생님이나 교수님이 선택한 문제에 학습전략이 일반화되었는지 알 수 있다. 일반화는 많은 추상화된 데이터뿐만 아니라 보지 못한 주제에 이러한 지식을 적용하는 방법에 대한 고수준 이해가 필요하다. 좋은 선생님은 이점에 꽤 도움이 된다.
37	17행	학습 과정을 3가지 단계로 설명했더라도 구체적인 목적을 위해 이 방법이 이뤄짐을 기억하자.	내용 설명을 위해 이렇게 학습 과정을 3 단계로 나누었다.
38	아래에서 3행	지식 표현성의 과정 동안 ~ 모델로 요약한다.	지식 표현성의 처리 과정 동안, 컴퓨터는 원시 입력을 데이터 간의 구조적 패턴의 명시적 기술인 모델(model)로 요약한다.
39	아래에서 4행	모델이 추상적인 ~ ~ ~ 쉽지 않다.	모델이 추상적인 항(떨어지는 물체를 설명하는 모델의 g 항)으로 중력을 알려주기 전까지 이 개념을 인식하기 쉽지 않다.
40	13행	추상화된 모델들(즉, 이론들)을 검색으로 생각할 수 있다.	추상화된 전체 모델(즉, 이론)에서 찾는 것으로 생각할 수 있다.
41	1-3행	휴리스틱은 사람이 ~ ~ 널리 사용된다. 상황을 평가하기 ~ ~ 있는 것이다.	사람은 경험을 새로운 상황에 대해 빠르게 일반화하기 위해 휴리스틱을 사용한다. 주어진 상황을 충분히 살피기 전에 결정하려고 직관을 활용한다면, 여러분은 자신도 모르게 정신적 휴리스틱을 사용하고 있는 것이다.
43	10행	기계 학습 태스크는 일련의 관리할 수 있는 단계로	기계 학습 태스크를 관리할 수 있는 일련의 단계로
45	4행	인스턴스는	예제는

<2장>

페이지	행	기존	변경
60	4행	위해 또는	위해
64	아래에서 6행	추가하거나 삭제했을 때	추가하거나 삭제해
66	7행	행은 예다.	행은 예제다.
73	3행	데이터를 저장하려면	데이터가 저장되어 있다면,
80	1행(제목)	퍼짐 측정: 사분위수와 5개 수의 요약	퍼짐 측정: 사분위수와 다섯수치요약
80	8행	5개 수의 요약은	다섯수치요약(five-number summary)은
81	아래에서 5행	quantile() 함수는 5개 수 요약을 반환한다.	quantile() 함수는 다섯수치요약을 반환한다.
82	박스	사분위수를 계산할 때 중앙값이 없는 데이터와 값 ~~ 명시할 수 있게 한다.	사분위수를 계산할 때, 중앙(middle) 값이 없는 데이터셋이나 같은 값들을 다루는 여러 가지 기법이 있다. quantile() 함수는 매 개변수로 9개의 다른 알고리즘을 명시할 수 있다.
82	아래에서 10행	5개 수 요약에	다섯수치요약에
83	7행	5개 요약의 일반적인	다섯수치요약의 일반적인
85	3행	라벨을 명시할 수 있는	라벨을 명시할 수 있는

<3장>

페이지	행	기존	변경
105	표	Nuts, Orange, Protein	nuts, orange, protein
105	아래에서 7행	단백질, 과일 증 하나로	단백질, 과일 증 하나로
112	6행	최소약(즉, 최대)	최소(즉, 최대)

<4장>

페이지	행	기존	변경
128	주석1	베이지스(Bayse)를 베이즈로	베이지스(Bayes)를 베이즈로
130	아래에서 6행	표기법은 $P(\sim\text{spam}) = 0.80$ 과 같은	표기법은 $P(\sim\text{spam}) = 0.80$ 과 같이
131	1행(제목)	조건부 확률	결합 확률
133	12행(제목)	베이지스 이론과 조건 확률	베이지스 이론과 조건 부 확률
133	아래에서 10행	독립적이기 때문에 조건 확률(conditional probability)로	독립적이기 때문에 조건 부 확률(conditional probability)로
136	7행	대부분의 경우 이 가정을 어길 때	대부분 경우 이 가정을 어기지만,
137	아래에서 7행	독립적이라는 범주 조건 독립	독립적이라는 범주 조건 부 독립
137	아래에서 7행	조건 독립의	조건 부 독립의
139	10행(식)	$0 / (0 + 0.0005) = 0$	$0 / (0 + \mathbf{0.00005}) = 0$
139	12행(식)	$0.00005 / (0 + 0.00005) = 1$	$0.00005 / (0 + \mathbf{0.00005}) = 1$
139	아래에서 1행	주어진 식료품 용어는 무력화하며, 다른 증거까지 모두 영향을 미친다.	주어진 식료품이라는 용어는 다른 증거까지 모두 무효로 만든다.
142	2행	너무 적은 빈은	너무 작은 구간은
142	3행	많은 구간	큰 구간
142	주석2	확률 분포의 ~~ 중의 하나	* 삭제
143	아래에서 8행	그 녀석은 며칠 안에 ~~ 시작했어.	그 녀석이 새 일을 찾는다면 금방 찾을 거야.
145	2행	로 저장한다.	로 저장한다. * 주석번호 2 추가 주석 2번 추가 : 파일 인코딩에 문제가 있으면 다음을 실행한다. sms_raw\$text <- iconv(enc2utf8(sms_raw\$text),sub="byte") - 옮긴이
147	4행	Vectrosouce() 를 명시하고,	Vectorsource() 를 명시하고,
150	4행	주어지 메시지에	주어진 메시지에
150	6행	주어지 tm	주어진 tm
152	2번째 박스	R이 그림에 모든 단어를 적합하게 할 수 없다는 경고 메시지를 받을 수 있다.	R이 그림에 모든 단어를 적합하게 나타낼 수 없다는 메시지를 출력할 수 있다.
156	13행(코드)	x <- factor(x, levels = c(0, 1), labels = c("No", "Yes"))	x <- factor(x, levels = c(0, 1), labels = c("No", "Yes"))
156	아래에서 2행	이전 문장에서 ~~ 호출한다.	이전 명령에도 비슷하게 사용했다. 그 함수는 apply()이다.
157	7행	이 결과는 행으로 ~~~ 두 매트릭스다.	그 결과, 행은 메시지, 열은 각 단어에 대해 Yes, No 로 명시한 팩터인 두 매트릭스가 된다.

<5장>

페이지	행	기존	변경
166	아래에서 2행	영화는 스타 영화배우가 나오는 영화와 나오지 않는 영화로 나뉜다.	영화는 영화배우가 많이 나오는 영화와 많이 나오지 않는 영화로 나뉜다.
166	그림 y축	영화배우 리스트	영화배우 수
167	그림 y축	영화배우 리스트	영화배우 수
167	아래에서 3행	마지막 그룹은 적은 스타 배우와	마지막 그룹은 출연한 스타 영화배우 수가 적은 한편,
198	아래에서 10행	버섯이 외간상 서로	버섯이 외관상 서로
199	3	식별하기 위해 카네기 멜론 대학(Carnegie Mellon University) UCI 기계 학습 저장소의 ~~ 사용한다.	식별하기 위해 카네기 멜론 대학(Carnegie Mellon University)의 제프 쉐리머(Jeff Schlimmer)가 기부한 UCI 기계 학습 저장소의 버섯 데이터셋을 사용한다.
201	아래에서 2행	그리고 자바를 설치해야 한다.	그리고 자바가 설치되어 있어야 한다.
205	5행	1R은 실제로 안전하게 역할을	1R은 실제로 안전한 역할을
207	9-11행(항목)	냄새가 foul(악취)라면 버섯은 독성이 있다. 주름의 크기가 좁고 색이 담황색이라면 버섯은 독성이 있다. 주름의 크기가 좁고 냄새가 pungent(자극적)이라면 버섯은 독성이 있다.	냄새(odor)가 악취(foul)라면 버섯은 독성이 있다. 주름의 크기(gill_size)가 좁고(narrow) 색(gill_color)이 담황색(buff)이라면 버섯은 독성이 있다. 주름의 크기(gill_size)가 좁고(narrow) 냄새(odor)가 자극적(pungent)이라면 버섯은 독성이 있다.

<6장>

페이지	행	기존	변경
221	아래에서 1행	해석하는 데 여러 가지 다양한 경험	해석하는 데 다양한 경험
222	박스 안 아래에서 2	의사가 모두~~ 경향이 많다.	의사가 영화를 좀 더 보라고 추천하기 전에, 우리는 다른 설명을 배제할 필요가 있다: 노인들은 더 적은 영화를 보고 더 빨리 돌아가실 가능성이 높다.
223	4행	독립 변수 y는 i번째 속성 x와 추정된 b 값에~~ 나타낸다.	독립 변수 y는 절편 항의 총합과 i번째 속성에 대한 x 값과 추정된 β 값의 곱을 다 더한 것으로 나타낸다.
223	11행	bi만큼 y는 변화한다.	β_i 만큼 y는 변화한다.
223	11행	절편이 예상된 값이다.	절편은 y의 예상 값이다.
223	12행	b0으로	β_0 으로
224	아래에서 8행	최적의 해답 찾기는 선형 대수가 필요하다. 결과 전개는 ~~ 더 많다.	최적의 해답을 찾으려면 선형 대수가 필요하다. 전개는 이 책의 범위를 넘는다.
225	14-15행(항목)	solve() 매트릭스의 역을 갖는다. t()는 매트릭스를 전치로 사용한다.	solve()는 매트릭스의 역을 구한다. t()는 매트릭스의 전치를 구한다.
226	7행	temperature는 ~~ 수 있다.	launch의 3번째 열인 temperature와 reg() 함수를 실행한다.
230	9행	female로 나눈다. ~~ 나눈다.	female로 나 뉐 다. ~~ 나 뉐 다.
231	15행	때문에 완벽 대각선은	때문에 대각선은
231	아래에서 10행	매트릭스의 상관관계가 없음은 강한 것을 의미한다.	매트릭스에서 뚜렷한 상관관계는 없다.
231	아래에서 4행	속성 간을 관계를	속성 간 의 관계를
231	아래에서 3행	많은 속성의 수는 어울리지 않을 수 있다.	많은 속성의 수는 오히려 복잡하게 만들 수 있다.
235	표 아래 3행	독립 변수	종속 변수
237	아래에서 10행	선형 회귀 모델의 ~~ 있게 한다.	선형 회귀 모델의 결과를 논리적으로 이해할 수 있다.
238	아래에서 5행	중요한 레벨은 ~~~ 0.05 레벨을 사용한다.	중요 레벨은(목록 아래 Signif. 코드) 측정값을 고려해 얼마나 실제 계수와 유사한지를 제공한다. 별 3개 표시는 종속 변수와 관련이 매우 높은 중요 레벨 0을 나타낸다. 일반적인 사례는 통계적으로 중요한 변수를 나타내기 위해 0.05인 중요 레벨을 사용한다.
239	2행	논리적인 방법으로 ~~ 보인다.	이들 변수만이 논리적으로 결과와 관련있어 보인다.
239	4행	종속 변수의 ~~ 제공한다.	종속 변수 값에 대해 모델이 얼마나 설명하는지를 측정하여 제공한다.
243	아래에서 4행	하는데 사용하는 일련의	하는 데 사용 되 는 일련의
243	아래에서 4행	이런 트리는 ~~ 수 있다.	트리를 커지게 하는 알고리즘을 약간 변경하여 수치를 예측하는 데 트리를 사용할 수 있다.
247	아래에서 2행	발견하기 위해 사용한다.	발견하기 위해 사용 된 다.
257	12행	한편, ~~ 상기하자.	한편, 대부분 와인은 매우 좋거나 매우 나쁜지는 않다는 점을 상기하자.
258	5행	1992년 ~~ 알고리즘이다.	왕(Wang)과 위튼(Witten)이 제안한 최신 모델 트리인 M5'(M5-prime) 알고리즘은 1992년 퀴란(Quinlan)이 제안한 본래의 M5 모델 트리 알고리즘을 향상한 것이다.

259	마지막 행~260p, 2	각 수는 ~~~~ 증가됨이 예측된다.	각 수는 예측된 와인 품질에 대한 연관 속성의 순수 효과(net effect)다. fixed acidity에 대한 0.266 계수는 acidity 단위가 1씩 증가할 때 와인 품질도 0.266씩 증가함을 나타낸다.
260	본문 3행~ 본문 끝	예측된 효과는 ~~ 기억하자.	예측된 효과는 단지 이 노드에 해당하는 와인 표본에만 적용한다는 점이 중요하다. 10개 다른 속성과 fixed acidity의 다른 영향 추정으로 총 36개 선형 모델이 이 모델 트리에서 만들어진다. 모델이 훈련 데이터를 얼마나 잘 적합화했는지에 대한 통계를 알기 위해, summary() 함수를 MSP 모델에 적용시킨다. 그러나 이런 통계는 훈련 데이터에 제한되어 있기 때문에, 대략적으로 진단하는 데 사용한다는 것을 기억하자.
262	5행	후자에서 두 가지 형태를~~ 회귀 모델을 만드는 모델 트리다.	후자에서 두 가지 형태, 즉 수치 예측을 하는 잎 노드에서 예제의 평균값을 사용하는 회귀 트리, 회귀와 트리 기법의 장점을 취하는 하이브리드 접근법으로 각 잎 노드에서 회귀 모델을 만드는 모델 트리를 사용했다.
262	9행	변수 사이의 관계는	변수 사이의 관계를
262	11행	흡연자와 비만인 사람으로	흡연자이자 비만인 사람으로
262	12행	회귀 트리에서 모델 트리는 ~~ 사용했다.	회귀 트리와 모델 트리를 측정된 특성으로부터 와인의 주관적인 품질을 모델화하는 데 사용했다.

<7장>

페이지	행	기존	변경
265	아래에서 4행(항목)	사무 건물의 ~~ 자동화	사무 건물의 환경 조절기, 스스로 운전하는 자동차, 스스로 운행하는 드론 같은 스마트 디바이스의 자동화
268	아래에서 2행	계단 활성화 함수	단위 계단 활성화 함수
269	2행	생화학의 제한으로	생화학의 제한에서
269	5행	이는 유사한 계단 모양이거나 경계 활성화 함수와 S형이지만	이는 계단 모양 혹은 S 모양으로 경계 활성화 함수와 유사하지만
269	아래에서 3행	이는 입력의 전체 ~~ 의미다.	입력 전체 범위에 걸쳐 도함수를 계산할 수 있다.
270	1행	시그모이드는 활성화 함수 가장 많이 사용되지만, 기본 형태로 사용한다.	시그모이드는 기본 형태인 활성화 함수가 가장 많이 사용된다.
284	1행	\$net.results다.	\$net.result다.
289	2행	MMH와 가장 가까운	최대 마진 초평면과 가장 가까운
289	4행	MMH를 정의할 수 있다.	최대 마진 초평면을 정의할 수 있다.
290	3행	MMH는	최대 마진 초평면은
290	5행	MMH는	최대 마진 초평면은
293	10행	추가적인 여백 변수에 추가해	여백 변수를 추가해
300	1-3행	마찬가지로 ~~ ~~ 함수를 제공한다.	마찬가지로, SVMlight 알고리즘을 활용한 도르트문트 기술 대학(Dortmund, Dortmund University of Technology) 통계학과의 klaR 패키지는 SVM 함수를 제공한다.
304	2행	정확하게 식별했다.	정확하게 식별했다는 점을 알 수 있다.

<8장>

페이지	행	기존	변경
309	8행	살펴본 분류와	살펴본 분류 나
309	15행	질적인 유용성을 위한 학습기 ~~ 점이다.	실질적으로 학습기 유용성을 평가하기도 어렵다는 점이다.
310	아래에서 2행(본문)	있음을 유추한다.	있음이 유추 된 다.
312	12행	명백하거나	명백한 사실,
314	14행	반복 i의 모든 아이템셋은 ~~ 결합한다.	반복 i+1 평가에서 사용할 후보 아이템셋을 생성하기 위해 반복 i의 모든 아이템셋을 결합한다.
314	15행	그러나 아프리오리 원칙은 ~~~ 수 있다.	그러나 다음 반복이 시작하기 전에, 아프리오리 원칙은 아이템의 일부를 제거할 수 있다.
314	아래에서 8행	{A, C}는 빈번하지 않음을~~사실 필요 없다.	{A, C}가 빈번하지 않다면, 반복 3에서 {A, B, C}에 대한 지지도 평가를 하려하지만 평가는 일어나지 않는다.
317	아래에서 6행	상단에 큰 제품 수가 가장 많은 거래를	상단에 제품 수가 가장 많은 거래를
319	아래에서 11행	매트릭스의 0이	매트릭스에서 0이
328	박스	세 개의 제품을 의미한다.	네 개의 제품을 의미한다.
330	1-6행	lift 값은 ~~ 지지도는 아니다).	lift 값은 potted plants 를 구매한 후 평균 소비자와 비교해, 얼마나 더 whole milk 을 함께 구매하는 소비자가 있는지 알려준다. whole milk 를 구매한 약 25.4%의 소비자(지지도)와 40%의 potted plant 을 사고 whole milk 를 산(신뢰도) 소비자를 알기 때문에, $0.40/0.256 = 1.56$ 으로 값을 계산할 수 있다. 이 값은 주어진 표의 값과 일치한다(support 열은 lhs나 rhs에 대한 지지도가 아닌 규칙에 대한 지지도다.).
331	아래에서 2행	arules 패키지는 가장 ~~ 함수를 포함한다.	arules 패키지에는 낮거나 높은 순으로 규칙을 정렬하는 sort() 함수가 있다.
332	아래에서 7행	whipped cream~~~~ 때문일까?	berries(딸기) 를 구매한 소비자는 일반 소비자보다 3배 더 많이 whipped cream 을 구매한다(디저트와 관련되었을까?).
333	10행	berriessms yogurt	berries는 yogurt

<9장>

페이지	행	기존	변경
339	아래에서 8행(3번째 항목)	대용량의 데이터셋 ~~~ 단순화	유사한 값을 가진 많은 속성을 그룹화해, 대량 데이터셋을 몇 개의 균일한 범주로 단순화
339	아래에서 6행	군집화는 복잡성 ~~ 만든다.	그 결과, 군집화는 복잡성을 줄이고 관계 패턴에 대한 통찰력을 제공하는 데이터 내의 실행 가능한 구조를 만든다.
343	4행	k 평균 알고리즘은 ~~ 뜻한다.	k 평균(k-means) 알고리즘은 각 n개의 예제를 실행하기 전 명시하는 군집의 개수인 k 군집 중 하나에 지정한다.
343	아래에서 5행	먼저 최초 k개 ~~ 변경한다. 군집 적합화가 ~~ 일어난다. 이쯤에서 ~~ 끝난다.	먼저, 최초 k개 군집 중 하나로 예제를 지정한다. 다음, 현재 군집에 속한 예제에 따라 군집 경계를 조정 후 다시 예제를 새로운 군집에 지정한다. 군집 적합화가 더 이상 변경이 없을 때까지, 군집을 조정하고 예제를 다시 지정하는 과정을 몇 번 실시한다. 이후, 이 과정은 멈추고 군집은 결정된다.
344	박스	K 평균의 ~~ 문제일 수 있다.	K 평균의 휴리스틱한 특성 때문에 초기 조건의 미세한 변경에도 다른 결과를 얻을 수 있다.
344	박스	이런 이유로 결과의 ~~ 편이 좋다.	이런 이유로 결과의 견고성을 위해 여러 번 군집 분석을 하는 편이 좋다.
352	13행	teen 데이터	teens 데이터
361	4행	본래의 인구로 돌아가	본래의 모집단으로 돌아가
362	아래에서 13행	전체적인 군집의 인구학적 특성을	전체적인 군집의 인구통계학적 특성을
363	아래에서 5행	성별에 따라 ~~~ 놀랄 만하다.	알고리즘 입력으로 친구 수를 사용하지 않은 점을 고려하면 이런 발견은 놀랄 만하다.
364	7행	이런 작업에 휴리스틱과 ~~ 대안물이 있다.	이런 작업에 대한 독특한 속성과 휴리스틱을 사용한 많은 대안이 있다.

<10장>

페이지	행	기존	변경
399	마지막 행	특별한 경우 부트스트랩	특별한 형태의 부트스트랩

<11장>

페이지	행	기존	변경
404	1-3행	최적의 k 값을 ~~ 옵션을 사용했다.	k 최근접 이웃 모델을 조절하기 위해 최적의 k 값을 찾고자 했고, 신경망이나 서포트 벡터 머신에 대한 노드 수, 은닉 층, 다양한 커널 함수 같은 옵션을 조절했다.
405	2행	하나를 선택하는 것은 연관된다.	하나를 선택하는 것과 연관된다.
406	아래에서 5행	리샘플링 전략을 선택한 10장에서 이 기법을 사용했다.	리샘플링 전략의 선택 같은 주제를 10장에서 살펴보았다.
406	아래에서 3행	성능 통계는 caret이 지원한다	성능 통계를 caret은 지원한다
407	2-4행	실제로 모델의 ~~ 함수를 제공한다.	이러한 실행은 모델 복잡성을 매우 증가시키는 것으로 눈에 띄게 성능 향상이 된 모델을 선택하기도 한다.
409	8~10행	train() 함수는 ~~ 사용할 필요가 없다.	train() 함수는 전체 입력 데이터에 대해 모델을 생성하기 위해 m\$finalModel과 같이 m에 저장된 최적의 조절된 매개변수를 사용한다. 대부분 finalModel의 하위 객체와 직접 연결할 필요는 없다.
415	아래에서 12행	진화 생물학과 자기 변형과 학습 테스트에 적용하는 유전자에서 차용한 개념을	진화 생물학, 자기 변형의 유전학, 적응 학습 테스트에서 차용한 개념을
416	그림	혼합 함수	조합 함수
416	아래에서 4행	생성된 모델은 일부 방법으로 처리가 필요한 예측을 내놓는다.	생성된 모델은 일부 방법으로 처리가 필요한 예측을 내놓는다. 조합 함수 (combination function)는 예측 간의 불일치를 어떻게 조합하지를 결정한다.
420	10-11행	이런 기능을 작성하기보다 ~~ 제공한다.	이런 기능을 작성하기 보다 caret 패키지의 기본 svmBag 리스트 객체가 제공하는 3개 함수를 사용한다.
421	2행	svmBag 리스트에 3가지	svmBag 리스트의 3가지
427	아래에서 11행	out-of-bag error rate	out-of-bag 오차 비율
427	아래에서 5행	out-of-bag 비율	out-of-bag 오차 비율
427	아래에서 2행	out-of-bag error 비율	out-of-bag 오차 비율

<12장>

페이지	행	기존	변경
434	아래에서 3행	문서 형태로,	평문 텍스트 형태로,
436	3행	중심을 한다.	중심으로 한다.