

[8~11] 다음 글을 읽고 물음에 답하시오.

데이터를 처리할 때 데이터의 정확성은 매우 중요하다. 그런데 데이터에 결측치와 이상치가 포함되면 데이터의 특징을 제대로 ㉠ 나타내기 어렵다.

결측치는 데이터 값이 ㉡ 빠져 있는 것이다. 결측치를 처리하는 방법 중 하나인 대체는 다른 값으로 결측치를 채우는 것인데, 대체하는 값으로는 평균, 중앙값, 최빈값을 많이 사용한다. 중앙값은 데이터를 크기순으로 정렬했을 때 중앙에 위치한 값이다. 크기가 같은 값이 복수일 경우에도 순위를 매겨 중앙값을 찾고, 데이터의 개수가 짝수이면 중앙에 있는 두 값의 평균이 중앙값이다. 또 최빈값은 데이터에 가장 많이 나타나는 값을 이른다. 일반적으로 데이터 값이 연속적인 수치이면 평균으로, 석차처럼 순위가 있는 값에는 중앙값으로, 직업과 같이 문자인 경우에는 최빈값으로 결측치를 대체한다.

이상치는 데이터의 다른 값에 비해 유달리 크거나 작은 값으로, 데이터를 수집할 때 측정 오류 등에 의해 주로 ㉢ 생긴다. 그러나 정상적인 데이터라도 데이터의 특징을 왜곡하는 데이터 값이 있을 수 있다. 예를 들어, 데이터가 어떤 프로 선수들의 연봉이고 그중 한 명의 연봉이 유달리 많다면, 이상치가 포함된 데이터에 해당한다. 이런 데이터의 특징을 하나의 수치로 나타내려는 경우 ㉣ 대푯값으로 평균보다 중앙값을 주로 사용한다.

평면상에 있는 점들의 위치를 나타내는 데이터에서도 이상치를 발견할 수 있다. 대부분의 점들이 가상의 직선 주위에 모여 있다면 이 직선은 데이터의 특징을 잘 나타낸다고 할 수 있다. 이 직선을 직선 L 이라고 하자. 그런데 직선 L 로부터 멀리 떨어진 위치에도 몇 개의 점이 있다. 이 점들이 이상치이다.

㉤ 이상치를 포함하는 데이터에서 직선 L 을 찾는다고 하자. 이때 사용할 수 있는 기법의 하나인 A기법은 두 점을 무작위로 골라 정상치 집합으로 가정하고, 이 두 점을 ㉥ 지나는 후보 직선을 그어 나머지 점들과 후보 직선 사이의 거리를 구한다. 이 거리가 허용 범위 이내인 점들을 정상치 집합에 추가한다. 정상치 집합의 점의 개수가 미리 정해 둔 기준, 즉 문턱값보다 많으면 후보 직선을 최종 후보군에 넣는다. 반대로 점의 개수가 문턱값보다 적으면 후보 직선을 버린다. 만약 처음에 고른 점이 이상치이면, 대부분의 점들은 해당 후보 직선과의 거리가 너무 ㉦ 멀어 이 직선은 최종 후보군에서 제외되는 것이다. 이 과정을 반복하여 최종 후보군을 구하고, 최종 후보군에 포함된 직선 중에서 정상치 집합의 데이터 개수가 최대인 직선을 직선 L 로 선택한다. 이 기법은 이상치가 있어도 직선 L 을 찾을 가능성이 높다.

8. 윗글을 이해한 내용으로 적절하지 않은 것은?
- ① 데이터가 수치로 구성되지 않아도 최빈값을 구할 수 있다.
 - ② 데이터의 특징이 언제나 하나의 수치로 나타나는 것은 아니다.
 - ③ 데이터가 정상적으로 수집되었다면 이상치가 존재하지 않는다.
 - ④ 데이터에 동일한 수치가 여러 개 있어도 중앙값으로 결측치를 대체할 수 있다.
 - ⑤ 데이터를 수집하는 과정에서 측정 오류가 발생한 값이라도 이상치가 아닐 수 있다.

9. 윗글을 참고할 때, ㉦의 이유로 가장 적절한 것은?
- ① 중앙값은 극단에 있는 이상치의 영향을 덜 받기 때문이다.
 - ② 중앙값을 찾기 위해 데이터를 나열할 때 이상치는 제외되기 때문이다.
 - ③ 데이터의 개수가 많아질수록 이상치도 많아지고 평균을 구하기 어렵기 때문이다.
 - ④ 이상치가 포함되면 평균을 구하는 것이 중앙값을 찾는 것보다 복잡하기 때문이다.
 - ⑤ 이상치가 포함되면 평균은 데이터에 포함되지 않는 값일 가능성이 큰 반면 중앙값은 항상 데이터에 포함된 값이기 때문이다.

10. ㉠과 관련하여 윗글의 A기법과 <보기>의 B기법을 설명한 내용으로 가장 적절한 것은? [3점]

—<보 기>—

다음과 같은 방법으로 직선 L 을 찾는 B기법을 가정해 보자. 후보 직선을 임의로 여러 개 가정한 뒤에 모든 점에서 각 후보 직선들과의 거리를 구하여 점들과 가장 가까운 직선을 선택한다. 그러나 이렇게 찾은 직선은 직선 L 로 적합한 직선이 아니다. 이상치를 포함해서 찾다 보니 대부분 최적의 직선과 이상치 사이에 위치한 직선을 선택하게 된다.

- ① A기법과 B기법 모두 최적의 직선을 찾기 위해 최대한 많은 점을 지나는 후보 직선을 가정한다.
 ② A기법은 이상치를 제외하고 후보 직선을 가정하지만 B기법은 이상치를 제외하는 과정이 없다.
 ③ A기법에서 최종적으로 선택한 직선은 이상치를 지나지 않지만 B기법에서 선택한 직선은 이상치를 지난다.
 ④ A기법은 이상치의 개수가 문턱값보다 적으면 후보 직선을 버리지만 B기법은 선택한 직선이 이상치를 포함할 수 있다.
 ⑤ A기법에서 후보 직선의 정상치 집합에는 이상치가 포함될 수 있고 B기법에서 후보 직선은 이상치를 지날 수 있다.

11. 문맥상 ㉠~㉥와 바꿔 쓰기에 가장 적절한 것은?

- ① ㉠: 형성(形成)하기
 ② ㉡: 누락(漏落)되어
 ③ ㉢: 도래(到來)한다
 ④ ㉣: 투과(透過)하는
 ⑤ ㉤: 소원(疏遠)하여