

학과: 응용계산학 이음: 백형민

인공지능 기초를 위한 FAQ

1. 인공지능에서 학습에 해당하는 기능은 무엇인가?

인공지능에서 학습에 해당하는 기능은 인간의 인지 능력을 모방하거나 구현하는 다양한 기능을 의미한다.

<대표적인 학습>

1. 학습(learning)

• 데이터에서 패턴을 찾아내고, 경험을 통해 개선하는 학습 (지도학습, 비지도학습, 강화학습)

2. 추론(reasoning)

• 기존의 지식을 활용하여 새로운 사실을 도출하는 학습 (인역적 추론, 귀납적 추론)

3. 문제 해결(problem solving)

• 주어진 목표를 달성하기 위해 최적의 방법을 찾는 학습 (탐색, 연결, 최적화, 계획)

2. 인공지능의 종류 3가지에 대해서 설명해줘 (지도학습, 비지도학습, 강화학습)

1. 지도학습(supervised learning)

• 정의: 입력 데이터(input)와 정답(label)이 주어진 상태에서, 모델이 입력과 출력간의 관계를 학습하는 방식

• 특징: 데이터에 대한 명확한 정답이 있어야 하며, 학습이 끝난 후 새로운 데이터에 대해 정확한 예측을 수행하는 것이 목표

2. 비지도학습(unsupervised learning)

• 정의: 일부 데이터에 정답(label)이 없고, 대부분의 데이터에는 정답이 없는 상태에서 학습하는 방법

• 특징: 지도학습과 비지도학습의 중간 형태인, 소량의 라벨링된 데이터와 대량의 라벨링되지 않은 데이터를 함께 활용

3. 강화학습(reinforcement learning, RL)

• 정의: 환경과 상호작용하며 시행 착오를 통해 보상을 최적화하는 방식

• 특징: 지도학습처럼 정답이 주어져 양극, 행동에 대한 보상과 벌칙을 통해 최적의 행동을 학습

3. 전통적인 프로그래밍 방법과 인공지능 프로그래밍의 차이점은 무엇인가?

전통적인 프로그래밍

• 코딩 방식: 개발자가 직접 규칙을 정의

• 입력과 출력: 입력 → 명시적인 규칙 적용 → 출력

• 문제 해결 방법: 논리적 규칙과 알고리즘을 직접 설계

• 변화 대응력: 새로운 문제가 생기면 개발자가 직접 수정

인공지능 프로그래밍

• 데이터에서 패턴을 학습

• 입력 → 학습된 모델 적용 → 출력

• 데이터에서 패턴을 찾아 자동으로 규칙을 학습

• 학습을 통해 새로운 패턴 적용 가능

4. 딥러닝과 머신러닝이 차이점은 무엇인가?

- 머신러닝 : 데이터에서 패턴을 학습하여 예측하거나 분류하는 알고리즘을 포함하는 인공지능의 하위 분야
- 딥러닝 : 머신러닝의 한 종류로, 인공신경망(ANN)을 기반으로 여러 층의 뉴런을 활용해 복잡한 패턴을 학습하는 방법
- 딥러닝은 작업의 복잡도를 선택. 딥러닝 모델이 자동으로 순차적 학습
- 머신러닝은 작은 데이터도 학습 가능해 상대적으로 연산량이 적다. 딥러닝은 대량의 데이터가 필요하고 높은 연산량이 요구된다.

5. Classification과 Regression의 주된 차이점은?

- 분류 (Classification) : 데이터를 특정 범주 (category) 또는 클래스 (class)로 구분하는 문제
- 회귀 (Regression) : 연속적인 숫자 값 (Numeric Value)을 예측하는 문제
- 출력 값이 이산적일 때 (분류) 또는 연속적일 때 (회귀)에 따라 적절한 알고리즘을 선택하는 것이 중요하다.

6. 머신러닝에서 차원 저주 (curse of dimensionality)란?

차원의 저주는 데이터의 차원 증가할수록 모델 성능이 저하되는 현상을 의미한다. 즉, 특성 (Feature)의 개수가 많아질수록 데이터가 희소해지고, 계산 복잡도가 증가하며, 일반화 성능이 떨어지는 문제가 발생한다.

7. Dimensionality Reduction는 왜 필요인가?

차원 축소 (Dimensionality Reduction)은 데이터의 특성 (Feature) 개수를 줄여 모델의 성능을 향상시키는 기법이다. 차원이 증가할수록 데이터가 희소해지고 연산 비용이 증가하며, 과적합 위험이 커지기 때문에 차원 축소가 필요하다.

8. Ridge와 Lasso의 공통점과 차이점은? (Regularization, 규제, scaling)

- 공통점
 - 둘 다 회귀 모델의 과적합 (Overfitting)을 방지하기 위한 정규화 (Regularization) 기법으로 사용된다.
 - 손실 함수 (Loss Function)에 페널티 (Penalty)를 추가하여 모델이 과도하게 복잡해지는 것을 방지한다.
 - 특성 (Feature) 스케일링이 필요하다 (standard scaler or min max scaler 사용)
- 차이점
 - Ridge는 L_2 정규화, Lasso는 L_1 정규화 방식 사용
 - Ridge는 가중치를 0에 차감하지 않는다, Lasso는 불필요한 정적분을 완전히 없애서 특성 선택 효과
 - Ridge는 모든 특성이 중요할 때 사용, Lasso는 일부 특성만 중요한 경우 사용

9. Overfitting vs. Underfitting

1. Overfitting (과적합)

정의: 모델이 훈련 데이터에 너무 과적합 되어 테스트 데이터에서 성능이 낮아지는 현상

원인: 모델이 너무 복잡함, 특징이 너무 많음, 학습 데이터에 과하게 적합

해결방법: Regularization (L1/L2, Dropout), 불필요한 특징 제거, 데이터 증가, 훈련 데이터 크기 증가

2. Underfitting (과소적합)

정의: 모델이 너무 단순하여 데이터의 패턴을 제대로 학습하지 못하는 현상

원인: 모델이 너무 단순함, 학습 데이터 부족, 학습 시간이 너무 짧음

해결방법: 더 복잡한 모델 사용, 더 많은 feature 추가, 훈련 데이터 증가, 학습 시간 증가

10. Feature Engineering 과 Feature selection의 차이점은?

1. Feature Engineering (특성 공학)

정의: 새로운 특징을 생성하거나 기존 특징을 변형하여 모델의 성능을 향상시키는 과정, 원본 데이터에서 더 유용한 정보를 추출하는 것이 목표

방법: 데이터 병합, 조합, 변환 (다항식 변환, 로그 변환, 원-핫 인코딩)

결과: 데이터의 표현력을 향상

2. Feature selection (특성 선택)

정의: 모델 성능에 영향을 주는 특징을 제거하여 학습 속도를 높이고 과적합을 방지하는 과정, 불필요한 변수를 제거하여 모델을 단순화하고 일반화 성능을 향상시키는 것이 목표

방법: 통계적 방법, 모델 기반 선택 (Lasso, RFE, PCA, 카계법 검증)

결과: 계산량 감소, 과적합 방지

11. 전처리 (Preprocessing)의 목적과 방법? (노이즈, 이상치, 결측치)

전처리 (Preprocessing)은 이상치인 모델을 학습시키기 전에 원시 데이터를 정제하고 변환하여 모델이 데이터를 잘 학습할 수 있도록 만드는 과정이다. 데이터 품질을 높이고 모델 성능을 제고하는 데 중요한 역할을 한다.

· 노이즈 처리

노이즈는 데이터에 포함된 무의미한 변동에 불규칙적인 값이다.

- 방법: 스무딩: 데이터를 평활화하여 불필요한 변동 제거 (이동 평균, 로컬 RFI)

필터링: 주어진 기준에 맞지 않는 데이터를 제거하거나 수정

· 이상치 처리

이상치는 데이터에서 다른 값들과 현저하게 다른 값을 의미한다.

- 방법: IQR 사용, 박스 플롯을 이용해 식별 → 삭제, 변경 (평균, 아 중앙값)

· 결측치 처리

결측치는 일부 데이터가 누락된 상태를 의미한다.

- 방법: 결측치를 포함한 행이나 열 삭제, 평균/중앙값/모평균으로 대체

12. EDA (Exploratory Data Analysis)란? 데이터의 특성 파악 (분포, 상관관계)

정의: 데이터의 특성을 탐색하고 이해하는 과정. 데이터를 분석하기 전에 데이터의 구조, 분포, 상관관계 등을 시각적으로 파악하는 중요한 첫번째 단계이다.

데이터 분포 확인: 각 특성의 분포를 확인하기 위해 히스토그램, 박스플롯 등을 활용. 연속형 변수는 히스토그램, 아인도그램, 변곡점 분석, 막대 그래프 활용.

상관관계: 두 변수 간의 관계를 나타내는 지표. 한 변수의 변화가 다른 변수에 어떻게 영향을 미치는지를 나타낸다. -1에서 1 사이의 상관계수를 갖고 1에 가까울수록 양의 상관관계, 0에 가까울수록 상관관계가 있음을 나타내고 -1에 가까울수록 음의 상관관계를 나타낸다.

13. 회귀에서 절편과 기울기가 의미하는 바는? 슬러핑과 어떻게 연관되는가?

선형 회귀식: $y = mx + b$

y 는 종속변수 (예측값), x 는 독립변수 (영입 값), m 는 기울기, b 는 절편이다.

기울기 (m): 독립변수가 1 증가할 때 종속변수가 얼마나 변하는지를 나타낸다.

절편 (b): $x=0$ 일때 y 값. 데이터의 기본적인 위치를 결정하는 역할을 한다.

슬러핑과 연관성: 슬러핑에서는 선형 회귀의 개념이 늘어나. 가중치와 편향과 직접 연결된다.

뉴런: 신경망에서 데이터를 처리하는 기본 단위.

가중치: 입력 데이터에 곱해지는 값 \rightarrow 기울기 (m) 역할 (입력 데이터가 모델에 출력에 미치는 영향을 조절)

편향: 모델이 조정할 수 있는 추가적인 선 값 \rightarrow 절편 (b) 역할 (뉴런의 활성화 여부를 조정하여 학습의 유연성을 증가시킴)

28. 결정 트리에서 불순도 (Impurity) - 지니 계수 (Gini Index)란 무엇인가?

1. 불순도: 한 노드에 있는 데이터가 서로 다른 클래스로 섞여 있는 정도를 나타내는 값.

- 불순도가 높다 \rightarrow 클래스가 다양해서 섞여 있을 확률이 높음 (예측)

- 불순도가 낮다 \rightarrow 대부분의 데이터가 같은 클래스에 속함 (분류가 쉬움)

2. 지니 계수: 결정 트리에서 불순도를 측정하는 대표적인 방법 중 하나이다.

(1) 공식

$$Gini = 1 - \sum_{i=1}^C p_i^2$$

C 는 클래스 개수
 p_i 는 특정 클래스 i 에 속하는 데이터 샘플의 비율
($0 \leq Gini \leq 0.5$)

• $Gini = 0$ (불순도 0%) \rightarrow 한 노드에 한 가지 클래스만 존재. (완전 순서, 최적의 상태)

• $Gini = 0.5$ (불순도 50%) \rightarrow 두 개의 클래스가 동일한 비율로 섞여 있음 (최악의 상태)

29. 앙상블이란 무엇인가?

정의: 여러 개의 다중학습 모델을 조합하여 더 나은 성능을 얻는 기법이다.

- 개별 모델이 가진 한계를 극복하고 일반화 성능을 향상시킬 수 있다.
- 단일 모델보다 과적합(Overfitting)을 방지하고 예측력을 높일 수 있다.
- 여러 모델을 조합하는 방식에 따라 Bagging, Boosting, Stacking 등의 기법이 있다.

30. 부트스트래핑(bootstrapping)이란 무엇인가?

정의: 데이터 샘플링 기법으로, 주어진 데이터에서 복원 추출을 통해 새로운 데이터셋을 여러개 만든 방법이다.
즉, 하나의 데이터셋에서 여러개의 샘플을 생성하여 통계적 추정이나 다중학습 모델 학습에 활용하는 기법이다.

31. 배깅(Bagging)이란 무엇인가?

정의: 부트스트래핑 기법을 활용하여 여러 개의 모델을 학습한 후, 결과를 평균 또는 투표 방식으로 결합하는 앙상블 기법이다. 즉, 하나의 모델이 아니라 여러 개의 모델을 독립적으로 학습하여 편향은 유지하면서 분산을 줄이는 방식이다.

32. 주성분 분석(PCA)이란 무엇인가?

정의: 고차원의 데이터를 저차원으로 변환하는 차원 축소 기법이다. 즉, 데이터의 중요한 정보는 최대한 유지하면서 차원을 줄이는 방법이다.