

벡터화 처리 분석

1. 텍스트 마이닝 ; 자연어로 구성된 비정형데이터에서 패턴 또는 관계를 추출하여 의미있는 정보를 얻는 작업

→ HLP (언어 처리 분석을 컴퓨터가 같은 기능을 수행하여 보다 쉽게 분석 가능)

즉, 텍스트 마이닝은 데이터 마이닝과 다르게 비정형 데이터를 다루고, NLP를 적용하여 텍스트 분석을 진행.

→ 웹 크롤링, 도표 분석, 음성 분석 + 조사 상관관계 분석, 인과 분석, 도표 모델링 등

→ 킴 쿼리 ; 빈도가 높은 키워드를 시각적으로 표시.

II. 감성 분석, sentiment analysis ; 텍스트를 작성한 사람의 태도, 의견, 성향과 같은 주관적 태도 분석.

감정 (Emotion) vs 감성 (Sentiment) ; 감성은 행복·공포 등 복잡한 심리 상태.

감성은 감정에 비해 단순한 정서적 태도.

→ 감성을 기반으로 한 텍스트 분류로 표현.

◦ 감성 점수 계산 ; 감성 점수 = 긍정 단어 수 - 부정 단어 수

but. 역설 (Irony) 과 풍자 (sarcasm)와 같은 문장은 다뤄 보류

III. 소셜 네트워크 분석.

◦ 소셜 네트워크란? 인터넷상에 사회적으로 연결된 많은 노드와 노드 사이에 도메인들.

→ 노드 또는 객체들이 노드와 노드를 기반으로 상호작용을 통해 네트워크를 형성하고 있는 개념을
개인의 행동을 통해 온라인에서 발생하는 다양한 사회적 현상을 설명하는데 있어
기본적 개념과 다른 새로운 노드를 적용한다.

◦ 노드 ; 노드 개체 (Node)들의 집합 V와 개체들 연결되는 Edge의 집합 E의 쌍.

- 종류 : 방향성 · 무방향 · 가중 (보통 노드 연결) · 가중 노드

- 특징 : degree (노드) - indegree, outdegree

hub (허브) - 상대적으로 노드가 높은 노드

- 노드의 정도 : $P(k) = \frac{H(k)}{H}$, $H(k)$: 노드 k를 가지는 노드의 개수, H : 전체 노드의 개수

- 노드의 밀도 : $\frac{\text{전체 노드의 정도}}{\text{전체 노드의 개수}}$

- 노드의 중심성 ; 개체가 가지는 영향력을 분석하는데 사용.

노드 중심성 $D_c(i) = \text{노드 i와 연결된 노드 수} / \text{노드 i와 연결된 노드 수}$

중심 중심성 $C_c(i) = (H-1) / \sum_{j \neq i} D_c(i, j)$, $D_c(i, j)$: 노드 i와 j 사이의 최단 경로 길이

매개 중심성 $B_c(i) = \sum_{j, k} s_{jk}(i) / s_{jk}$, s_{jk} : 노드 j와 k 사이의 최단 경로 길이, $s_{jk}(i)$: 노드 i를 포함하는 최단 경로 길이

⇒ 링크 관계도, 현상도, 유량 관계도, 흐름도 등 다양한 활용 사례.

IV. 텍스트 클러스터링

◦ 텍스트 클러스터링? 텍스트의 개체들 간의 유사도를 계산한 후 이를 바탕으로 유사한 개체들을 묶어주는 작업.

◦ 텍스트 데이터 수회 표현.

i) TF (Term Frequency) ; Term-document 행렬. 같은 단어, 같은 문서로 발생빈도 행렬.

→ 단문으로 별로 중요하지 않은 단어가 빈도가 높아 영향을 미침.

ii) TF-IDF (TF-inverse document frequency) ; \log_2 처리한 것은 outlier 영향 감소.

$IDF = \log_2 \left(\frac{H}{DF} \right)$, DF : 문서가 나타난 문서 수. → 단어의 중요도 평가 기준.

$TF-IDF = TF(t, d) \times IDF(t)$ but. 문서의 길이가 길수록 클러스터링에 영향.

∴ 문서 길이 또한 고려해야 함

n_d : 문서 d의 길이

→ TF의 표준화 $\tilde{TF}_{td} \leftarrow \frac{TF_{td} - \min_j TF_{td}}{\max_j TF_{td} - \min_j TF_{td}}$ / TF의 정규화 $\tilde{TF}_{td} \leftarrow \frac{1 + \log_2 TF_{td}}{n_d}$

◦ 유사도 측정.

i) 유클리드 거리 ; 점에 대한 (공공화·점과 점의), 최소 거리로 간격 측정.

ii) 코사인 유사도 ; '각도'를 반영한 유사도.

◦ 클러스터링 방법 → #개별항 처리 분석.

V. 도표 모델링

◦ 도표 모델링이란? 문서에서 추출한 도표들 (구체) 중 나타내는 기술 일련의 관계 분석.

→ 도표 ; 문서에서 각 주제 등장하는 단어들이 문서 전체 도표를 이루고 있음. 문서와 같은 구조 분석.

즉, 큰 문서에서 각 주제 등장하는 단어들을 주제별로 분류하여 도표 모델링 가능.

◦ LDA (Latent Dirichlet Allocation)

step 1. 각 주제 k에 대한 단어들의 $\phi_k \in R^V$ 는 Dirichlet 분포를 따름.

i) 도표 개수 k 지정.

step 2. D개의 문서의 각각은 K개 주제 도표들의 혼합으로 구성되어 있다고 가정. 도표들 $\theta_d \in R^K$ 는 Dirichlet 분포를 따름.

ii) 각 주제별로 단어의 분포 측정.

step 3. 문서 d에 있는 n번째 단어에 대하여 도표들 중 선택.

D : 전체 문서 수	K : 전체 도표 수
θ_d : 문서 d의 주제 분포	ϕ_k : 주제 k의 도표 분포
$z_{d,n}$: 문서 d에서 n번째 주제	
$w_{d,n}$: 문서 d에서 n번째 단어	

도표 $z_{d,n} \sim \text{multinomial}(\theta_d)$

도표 (각 주제별로)를 가지고 $(\phi_k = \phi_{z_{d,n}})$ 단어를 선택.

$w_{d,n} \sim \text{Multinomial}(\phi_{z_{d,n}})$

→ Dirichlet 분포?

× Dirichlet 분포.

$f(x_1, \dots, x_K | p_1, \dots, p_K) = \frac{1}{B(p)} \prod_{k=1}^K x_k^{p_k-1}$

$p(x_1, \dots, x_K | d_1, \dots, d_K) = \frac{n!}{x_1! \dots x_K!} \prod_{k=1}^K p_k^{x_k}$

$B(p) = (\pi^K \Gamma(p_1) \dots \Gamma(p_K)) / \Gamma(\sum p_k)$ 는 다변량 베타 함수.

◦ LSA ; DMM을 차원 축소하여 근접점들을 도표로 묶음.