

빅데이터

대량히 Big data 가 아닌 Big Value from data (분석)

- ① 빅데이터 시대가 도래하면서 누구나 빅데이터를 분석할 수 있어야 하는 능력이 요구되었다
- ② 누구나 분석하기 위해서는 플랫폼이 필요했고, 사용하기 쉬워야 했다
- ③ 그렇게 하둡이라는 플랫폼이 탄생했다

하둡 에코시스템

빅데이터 기술 + 분산구조(R)

빅데이터 플랫폼 → 하둡

데이터 처리과정

데이터 수집 → 저장 → 분석 → 시각화 → 서비스

빅데이터의 활용 ① 예측 ② 추천
예) 구매(행위), 이벤트(행위), 행렬(행위)

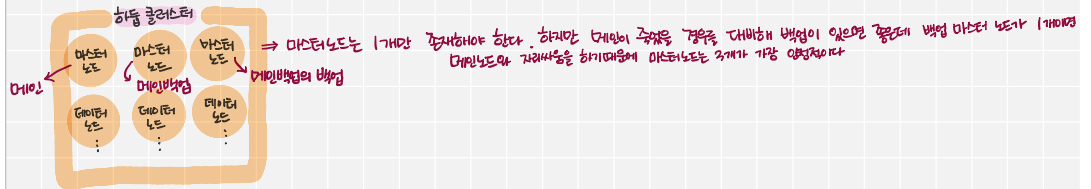
하둡 시스템 과정

- ① 예측
- ② 분류
- ③ 클러스터링 (군집분석)
- ④ 연관규칙 (연관성)
- ⑤ 협업 필터링

하둡 에코시스템

2006년 즈음부터 하둡은 기업적인 환경에 정착했다 (영양분 ↑)
한 대의 컴퓨터로 여러 영병병과 과부하 때문에 cluster 방식 (N대의 컴퓨터영역) 탄생

슈퍼컴퓨터 ① 1000 노드 이상 ② 속도 Top 500 (대규모 병렬)
하둡 클러스터 : 분산 파일 시스템 + 분산병렬처리시스템



2008년에 하둡이 슈퍼컴퓨터보다 빨랐다. 하둡은 다른 컴퓨터의 일을 신경쓰지 않고 그저 연결되어있는 (클러스터) 1개의 컴퓨터만 **각자** 잘 돌아가게 해준다
슈퍼컴퓨터는 모든 처리과정이 잘 진행되든지 끊임없이 서로 확인하기 때문에 요청과 응답을 계속해서 속도가 하둡에 비해 느림

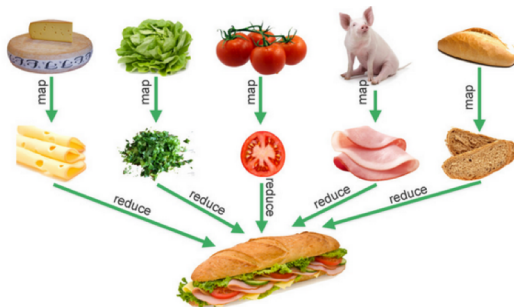
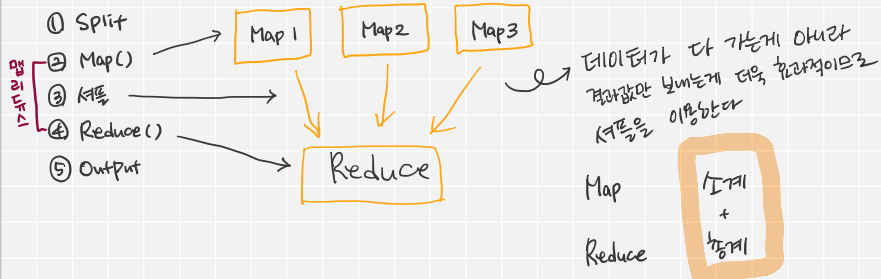
하둡 분산파일 시스템 HDFS

분산시스템의 유일한 단점은 어디에 뭐가 저장되어있는지 모른다는 것이다 '이'!
저장하려는 파일이 128mb 이상이면 128mb씩 자른다
128mb 이하이면 원본 그대로 유지한다



네임노드가 데이터노드의 상태를 보고 어디다 저장할 것인지 결정한다

Map Reduce



맵리듀스의 원리

	정렬시간	Read/Write 시간
백만개	1분	0.1분
천만개	100분	1분
억개	10,000분	10분

Map $\frac{1}{2}$ Reduce
 100분 10분 10분 → 120분

* 셔플을 하게되면 꼭 오버헤드(20%~30%)가 발생한다

Shuffle (모양별로 수집) Reduce (모양별로 Merge Sort)

예)

Map 개수	Reduce 개수	Map 시간	Shuffle 시간	Reduce 시간
10	1	100분	10분	10분
10	4	$100/4$ 분	$10/4$ 분	$10/4$ 분

컴퓨터 10대를 1000배의 속도를 낼 수 있다