



기계학습을 이용한 우리나라 여름철 장기기온예측

2022년 10월 28일
APCC 기후정보서비스 사용자워크숍

예측기술과
이진영 선임연구원

오늘의 학습

I. 기계학습이란

- 인공지능이란
- 기계학습이란
- 기계학습의 종류

II. 기계학습을 이용한 장기예측

- 장기예측을 위한 준비
- 다양한 기계학습을 이용한 기온예측 – Toy Examples

I. 기계학습이란

- 인공지능이란
- 기계학습이란
- 기계학습의 종류

I. 기계학습이란

인공지능이란

“인간의 지능을 필요로 하는 일들을 컴퓨터나 컴퓨터에 의해 제어되는 로봇이 수행하는 능력”

<http://movie.naver.com/movie/bi/mi/basic.nhn?code=10200>



터미네이터 2 1984~

미래, 인류와 기계의 전쟁은 계속 되는 가운데 스카이넷은 인류 저항군 사령관 존 코너를 없애기 위해 ...

movie.naver.com <http://movie.naver.com/movie/bi/mi/basic.nhn?code=24452>



매트릭스 1999

2199년, 인공 두뇌를 가진 컴퓨터(AI: Artificial Intelligence)가 지배하는 세계. 인간들은 태어나자마자...

movie.naver.com

<http://movie.naver.com/movie/bi/mi/basic.nhn?code=38420>



아이, 로봇 2004

(법칙 1. 로봇은 인간을 다치게 해선 안되며, 행동하지 않음으로써 인간이 다치지도록 방관해서도 안된다...)

movie.naver.com

출처: 네이버 영화

<https://movie.naver.com/movie/bi/mi/basic.nhn?code=69105>

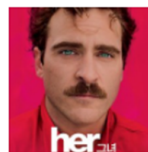


월-E 2008

땅 빈 지구에 홀로 남은 스페이스카 인간은 어떻게 ...

movie.naver.com

<http://movie.naver.com/movie/bi/mi/basic.nhn?code=101950>



그녀 2013

'테오도르'(조아킨 피닉스)는 다른 사람들의 편지를 대신 써주는 대필 작가로, 아내(루니 마라)와 별...

movie.naver.com

<http://movie.naver.com/movie/bi/mi/basic.nhn?code=118361>

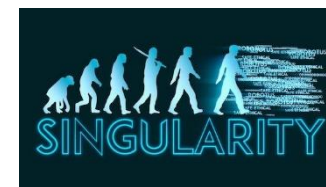


엑스 마키나 2014

유능한 프로그래머 '칼랩'(동늘 글리슨)은 치열한 경쟁 끝에 인공지능 분야의 천재 개발자 '네이든'...

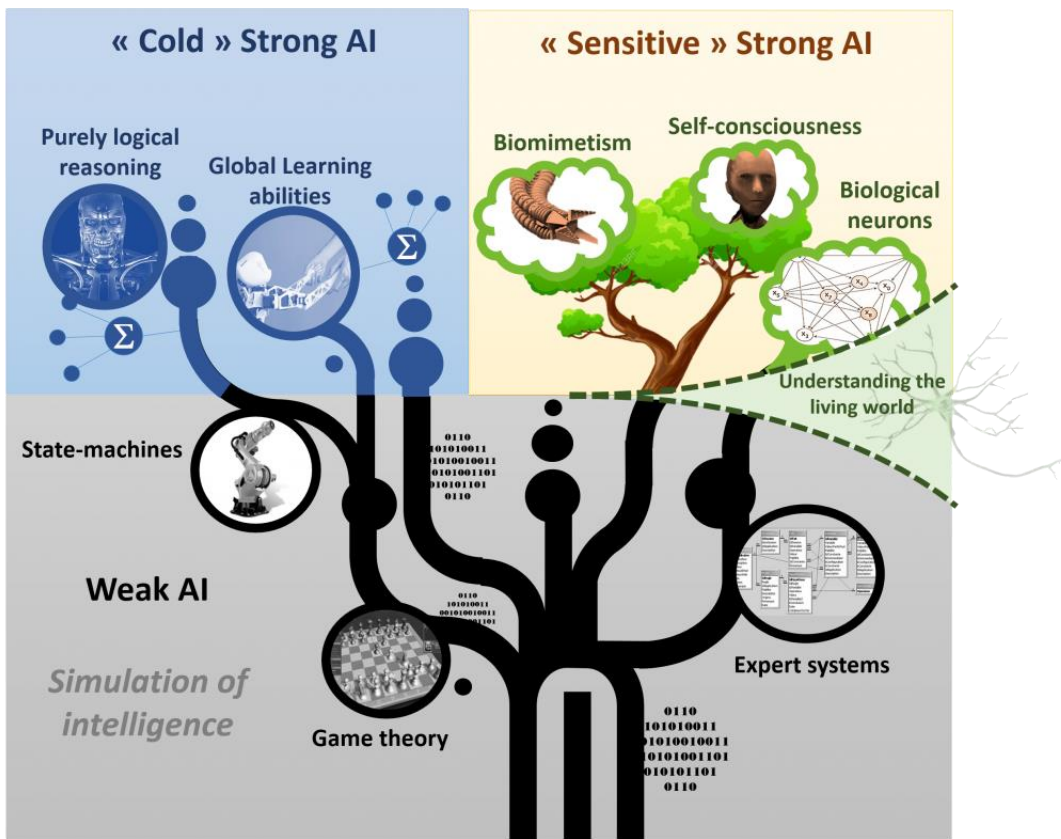
movie.naver.com

출처: 네이버 영화

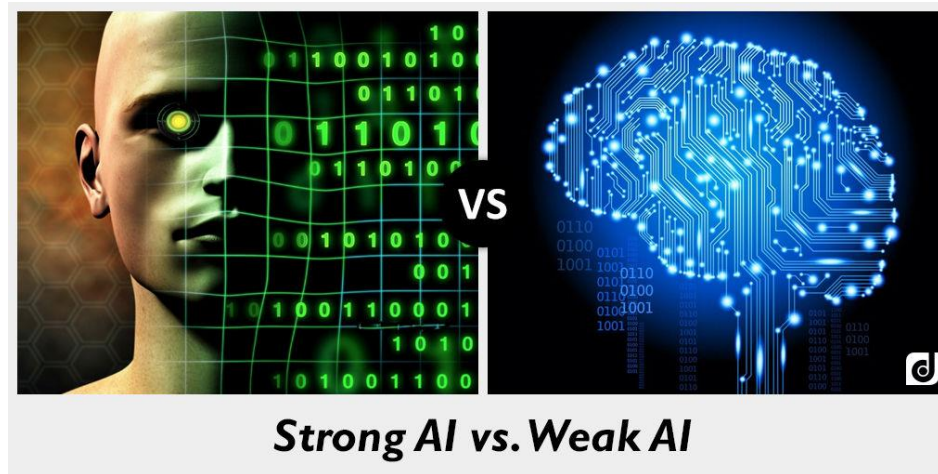


I. 기계학습이란

인공지능이란



이미지 출처: Théophile Gonos



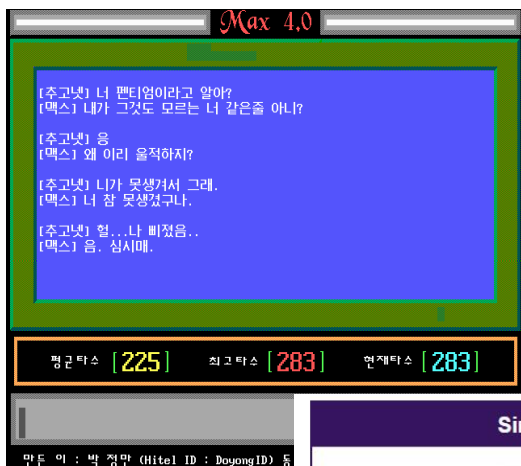
이미지 출처: Difference Between

Strong AI는 인간을 모방하는 것 뿐 아니라 인간과 유사한 지능을 갖추고 인간처럼 사고하고 행동

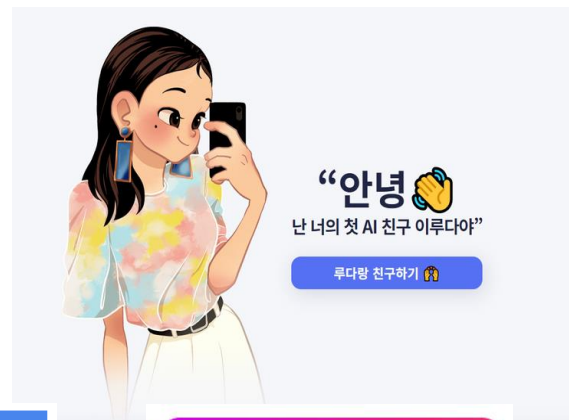
Weak AI는 인간이 미리 계획한 일을 능숙히 수행하기 위한 기술을 개발하기 위해서 사고

I. 기계학습이란

인공지능이란



심심이	
국가	대한민국
개발사	심심이(주)
대표자	김유진
출시일	2002년 ^[1]
지원 운영체제	Android, iOS
지원 언어	81개 언어 ^[2]

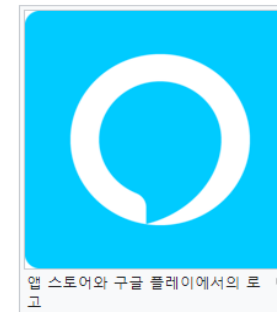


Siri	
개발	Apple
출시	2011년 10월 12일
OS	iOS 5 이상, iPhone 4s 이상의 iPhone 및 5세대 이상의 iPod touch iPadOS 모든 버전 macOS 10.12 Sierra 이상 watchOS 모든 버전 tvOS 모든 버전

구글 어시스턴트 Google Assistant	
개발	구글
발표	2016년 5월
유형	크로스 플랫폼형
지원 대상	Android, iOS, 스마트 TV
지원 언어	문단 참고
홈페이지	🏠 🗨️ 도움말 (한국어)
도움말	🏠

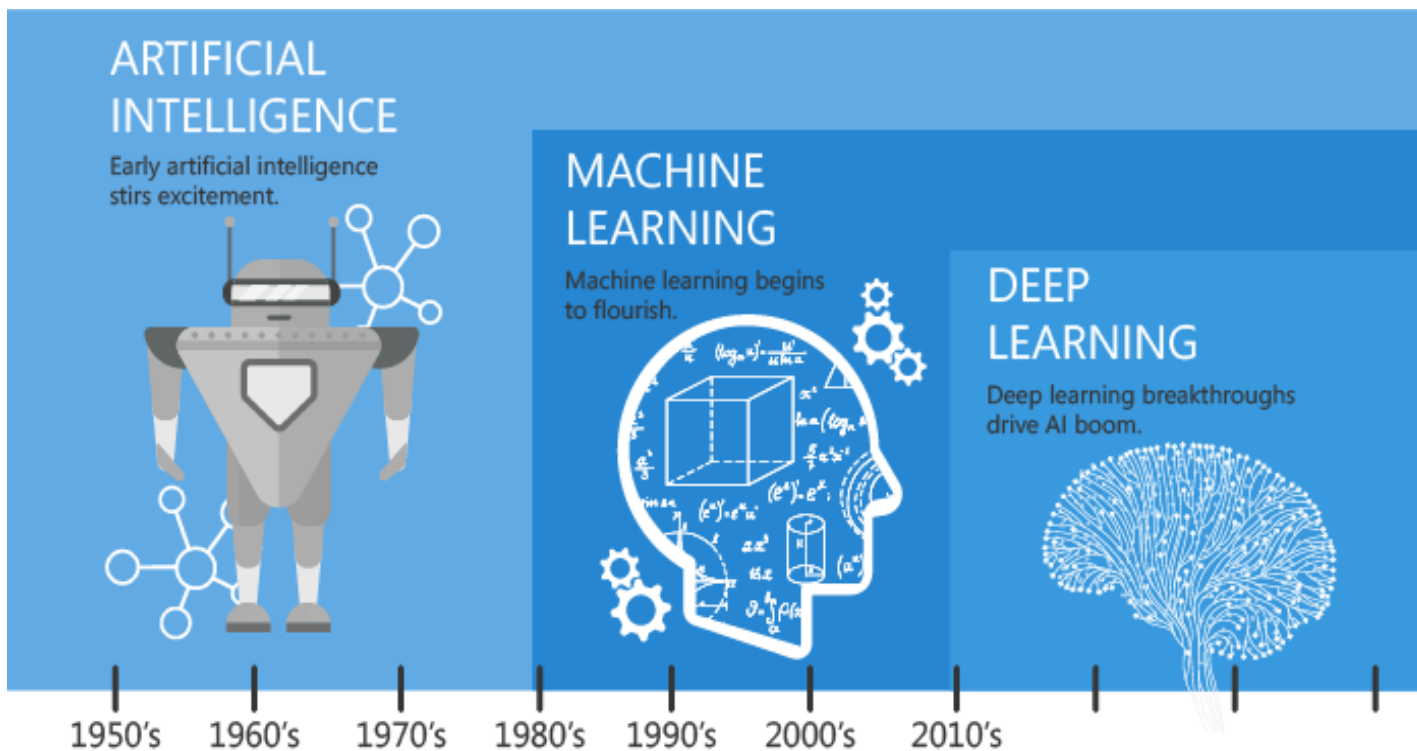
빅스비 Bixby	
개발 및 유통	SAMSUNG
플랫폼	안드로이드 7.0 이상 ^[1] , 기기에 따라 다릅니다.
타이젠	4.0 이상

알렉사 amazon alexa	
개발자	아마존
발표일	2014년 11월 (6년 전)
운영 체제	iOS 8.0 이상 안드로이드 4.4 이상
크기	4.6 MB (iOS) 38.76 MB (안드로이드)
언어	영어(US), 영어(UK), 독일어(DE), 영어(India), 영어(CA)
웹사이트	alexa.amazon.com/📄🔗



I. 기계학습이란

기계학습이란



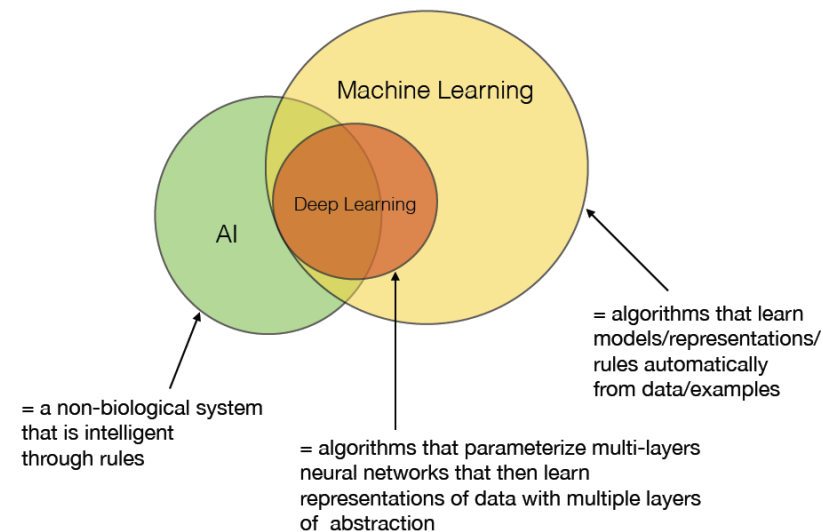
Since an early flush of optimism in the 1950's, smaller subsets of artificial intelligence - first machine learning, then deep learning, a subset of machine learning - have created ever larger disruptions.

출처: NVIDIA 블로그

처리과정을 학습하는 컴퓨터 모델링:
컴퓨터로 하여금 귀납이나 연역과 같은 특정 추론 전략을 사용하여 현존하는 자료나 이론으로부터 지식을 획득하게 함

출처: Jensen, 2016

Examples From The Three Related "Areas"

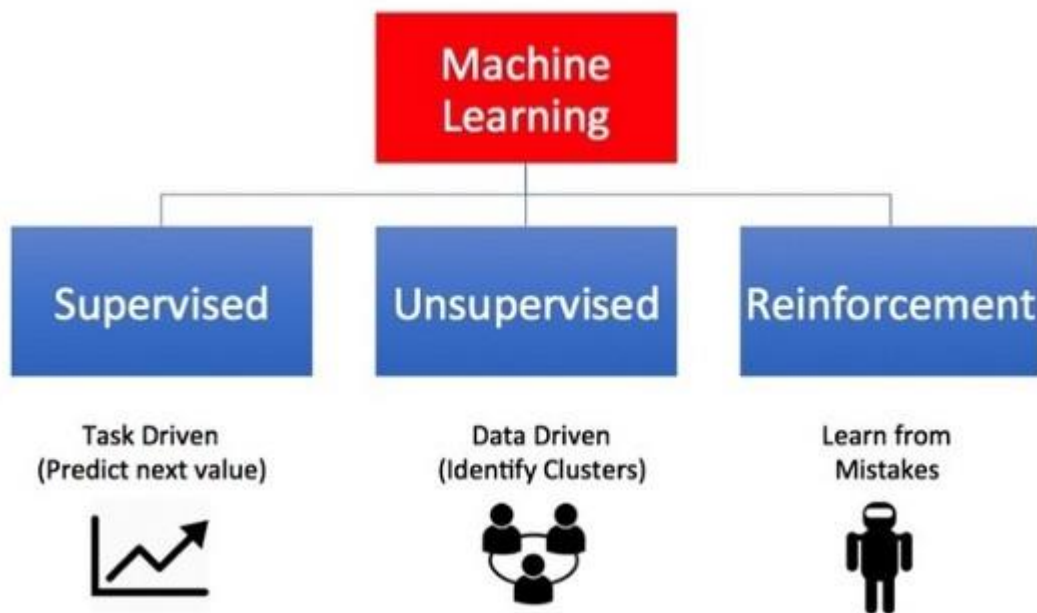


I. 기계학습이란

기계학습의 종류

지도학습 비지도학습 강화학습

Types of Machine Learning



Supervised Learning	<ul style="list-style-type: none"> > Labeled data > Direct feedback > Predict outcome/future
Unsupervised Learning	<ul style="list-style-type: none"> > No labels/targets > No feedback > Find hidden structure in data
Reinforcement Learning	<ul style="list-style-type: none"> > Decision process > Reward system > Learn series of actions

출처: Raschka and Mirjalily, 2019

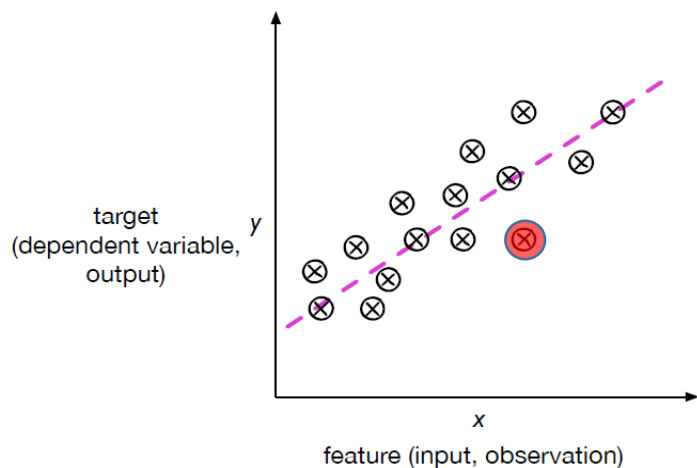
출처: towardsdatascience

I. 기계학습이란

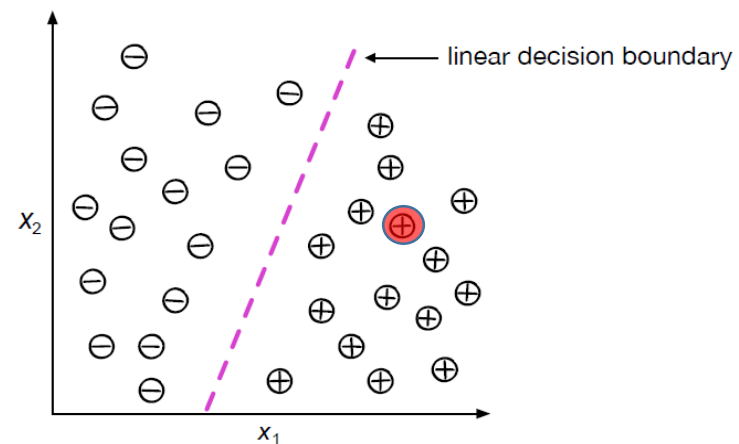
기계학습의 종류

지도학습

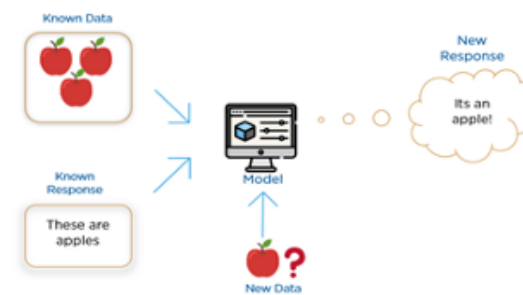
회귀



분류 / 패턴인식



출처: Raschka and Mirjalily, 2019



출처: towardsdatascience

I. 기계학습이란

기계학습의 종류

비지도학습

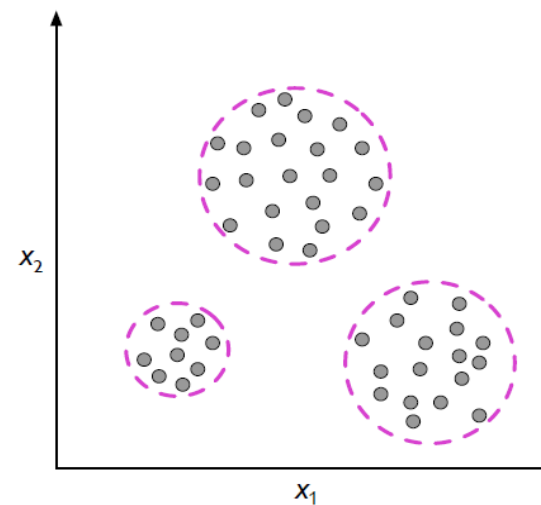
차원 감소



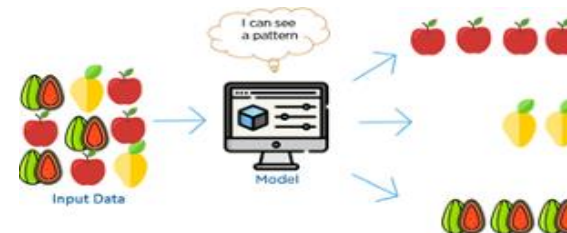
<그림 6> k개의 eigenface만을 이용한 데이터 복원(reconstruction)

출처: darkpgmr 블로그

군집



출처: Raschka, 2020

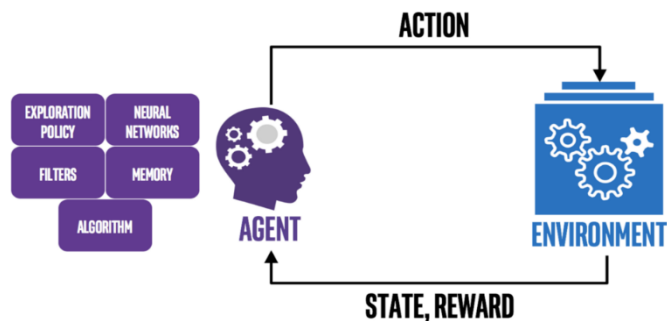


출처: towardsdatascience

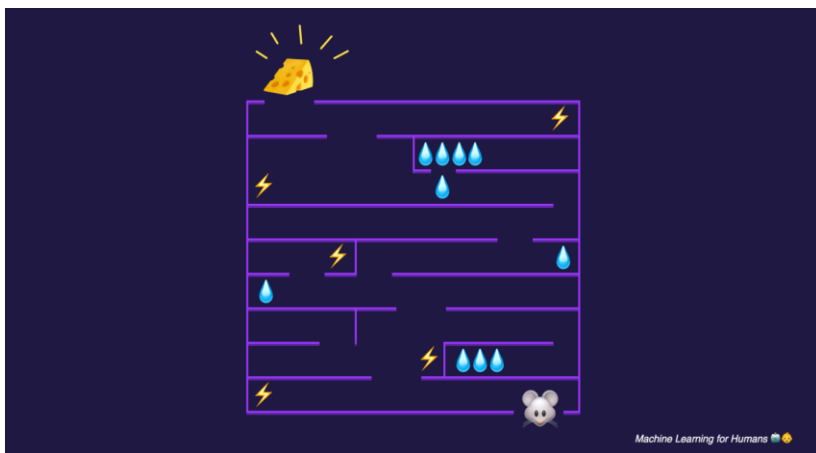
I. 기계학습이란

기계학습의 종류

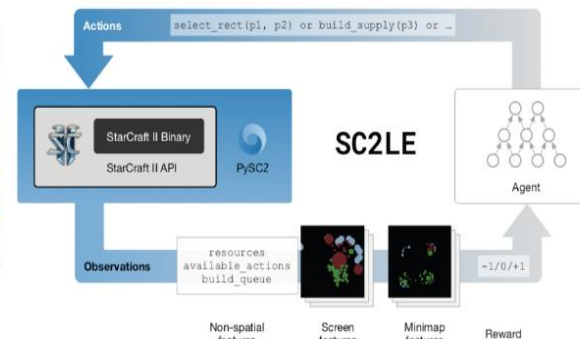
강화학습



출처: medium.com



Starcraft II



Vinyals, Oriol, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani et al. "Starcraft II: A new challenge for reinforcement learning." *arXiv preprint arXiv:1708.04782* (2017).



어떤 환경 안에서 정의된 에이전트(행위자)가 현재의 상태를 인식하여 선택 가능한 행동들 중 보상을 최대화하는 행동을 선택하는 방법

II. 기계학습을 이용한 장기예측

- 장기예측을 위한 준비
- 다양한 기계학습을 이용한 기온예측 - Toy Examples

II. 기계학습을 이용한 장기예측

장기예측을 위한 준비

지도학습 활용 의 단계

풀고 싶은 문제를
정의하자

자료를
수집하자

활용할 기법을
선택하자

모델 학습을 위한
최적화 방법을
정하자

모델 평가를 위한
기준을 정하자

II. 기계학습을 이용한 장기예측

풀고 싶은 문제를 정의하자: 단정예측

지도학습 활용 의 단계

풀고 싶은 문제를
정의하자

자료를
수집하자

활용할 기법을
선택하자

모델 학습을 위한
최적화 방법을
정하자

모델 평가를 위한
기준을 정하자

기온 값 자체를 예측
→ 단정예측

5월(ex. 5/15)에 예측한 6, 7, 8월 평균 기온 예측

삼분위 카테고리 분류
→ 각 카테고리별 확률
→ 확률예측

확률분포 예측(베이지안)
→ 각 카테고리별 확률

5월(ex. 5/15)에 예측한 6, 7, 8월 평균 기온
확률 예측

II. 기계학습을 이용한 장기예측

자료를 수집하자

지도학습 활용
의 단계

풀고 싶은 문제를
정의하자

자료를
수집하자

활용할 기법을
선정하자

모델 학습을 위한
최적화 방법을
정하자

모델 평가를 위한
기준을 정하자

입력자료 X

관측자료 y

- 훈련/검증/테스트 기간의 6, 7, 8월 평균기온
→ 전국 62개 ASOS 지점 자료 평균으로 산정

II. 기계학습을 이용한 장기예측

자료를 수집하자

지도학습 활용
의 단계

풀고 싶은 문제를
정의하자

자료를
수집하자

활용할 기법을
선정하자

모델 학습을 위한
최적화 방법을
정하자

모델 평가를 위한
기준을 정하자

입력자료 X

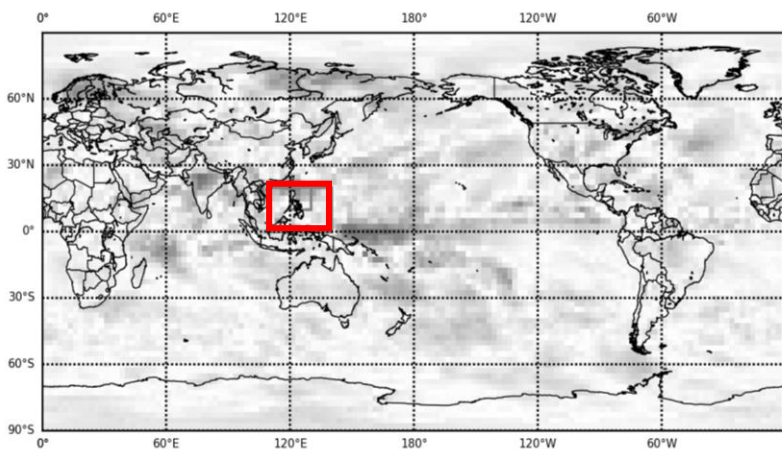
관측자료 y

- Outgoing Long wave Radiation (OLR)
- Sea Surface Temperature (SST)
- Precipitation (PRCP)
- Snow Cover Extent (SCE)
- 500 hPa & 850 hPa Geopotential Height (Z500, Z850)
- Sea Ice Area (SIA)

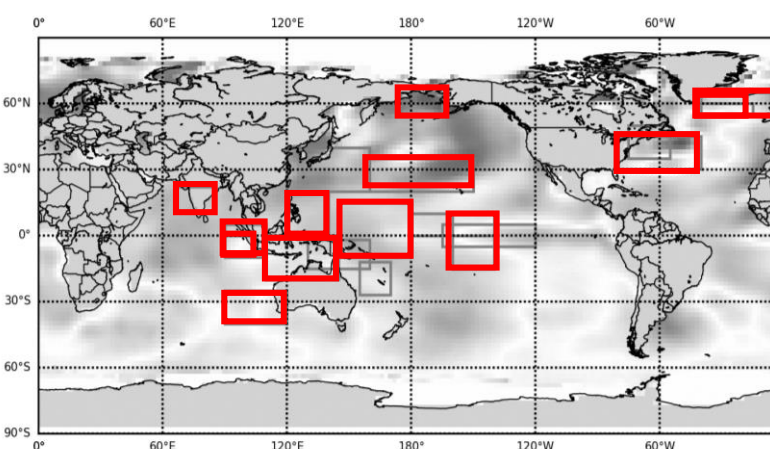
II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 기후감시요소의 수가 적으므로 장마기간 강수, 6~8월 기온에 영향을 미치는 요소 모두 사용
 - 쌍극, 삼극 패턴의 경우 개별 및 차이 모두 사용, 전년도 12월 ~ 당해 년도 4월 자료 모두 사용
- >> 약 30여개의 기후감시요소 도출



5월 북서태평양 지역 OLR 편차

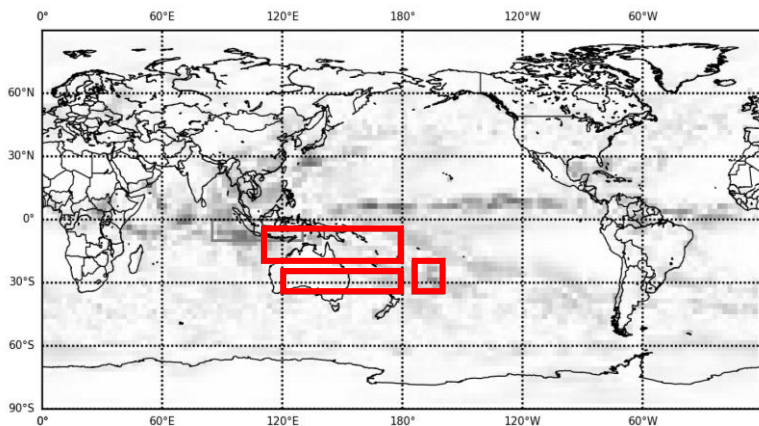


- 4월 북서태평양 SST 편차 지역차
- 4월 북태평양 지역 SST 편차(시기차)
- 4월 열대중태평양 지역 SST 편차
- 4월 북인도양 지역 SST 편차
- 3월 베링해 SST 편차
- 1월 열대서태평양 SST 편차
- 4월 동인도양~필리핀동쪽 SST 편차 삼극
- 2월 호주북쪽 지역 SST 편차
- 4월 열대인도양 지역 SST 편차
- 4월 북대서양 SST 편차

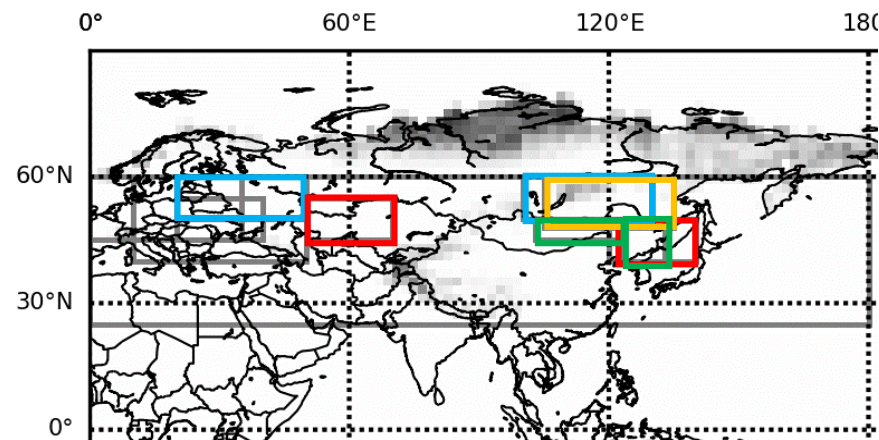
II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 기후감시요소의 수가 적으므로 장마기간 강수, 6~8월 기온에 영향을 미치는 요소 모두 사용
 - 쌍극, 삼극 패턴의 경우 개별 및 차이 모두 사용, 전년도 12월 ~ 당해 년도 4월 자료 모두 사용
- >> 약 30여개의 기후감시요소 도출



4월 호주부근 PRCP 편차 지역차
3월 호주북쪽 지역 PRCP 편차



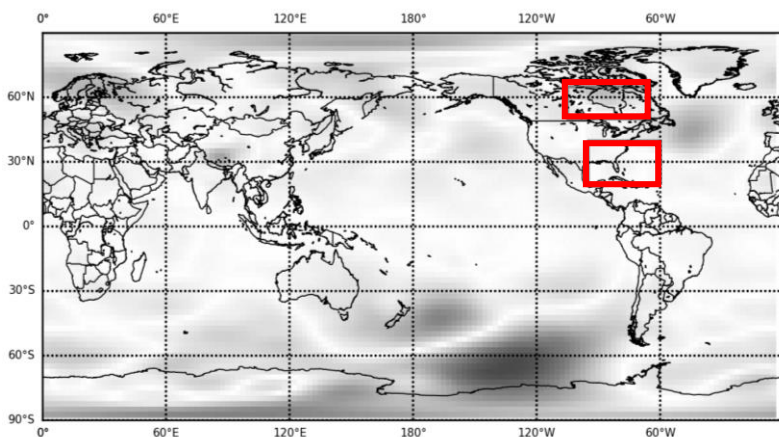
3~5월 유라시아 눈덮임 편차 지역차
3월 중국 북동부 지역 눈덮임 편차
4월 동아시아 / 서아시아 눈덮임 편차 지역차
4월 만주 눈덮임 편차

II. 기계학습을 이용한 장기예측

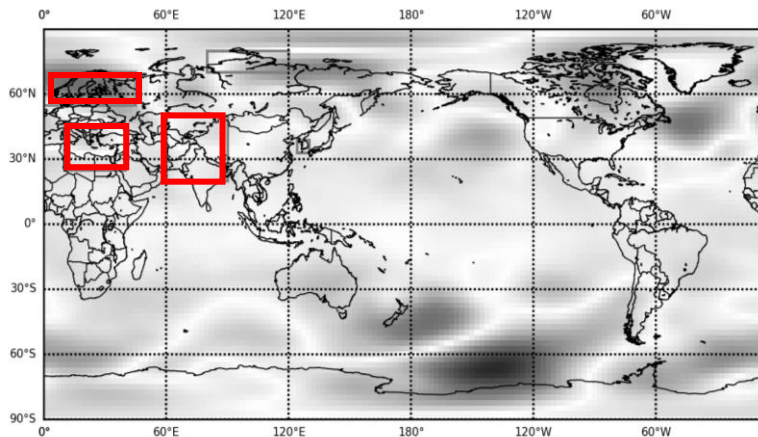
자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 기후감시요소의 수가 적으므로 장마기간 강수, 6~8월 기온에 영향을 미치는 요소 모두 사용
 - 쌍극, 삼극 패턴의 경우 개별 및 차이 모두 사용, 전년도 12월 ~ 당해 년도 4월 자료 모두 사용
- >> 약 30여개의 기후감시요소 도출

3~4월 베링해 해빙 편차



4월 대서양 지역 850 hPa 지위고도 편차 쌍극



3월 유라시아 지역 500 hPa 지위고도 편차 쌍-삼극

II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 각 기후감시요소에 대하여 시계열 자료 산정

자료 수집

NOAA NCEP CPC GLOBAL monthly olr Data Files

This dataset has bytes (2.3630976E07 22.536255MB) of data in it, which should give you a rough idea of the size of any file that you ask for.

Download Data To Specific Software

ingrid	The Postscript-based software on which the Data Library is built.
CPT	Climate Predictability Tool More information
ferret	Interactive computer visualization and analysis software. More information
GRADS	Grid Analysis and Display System More information
matlab	Data analysis and visualization software. More information
NCL	NCAR Command Language More information
WinDisp	A public domain software package for the display and analysis of satellite images, maps and associated databases, with an emphasis on early warning for food security. More information

Other Available File Formats

Full Information Formats	
These files contain all of the available metadata.	
OPeNDAP	A system which downloads data directly to software, such as matlab, Ferret, GRADS, etc. Specific instructions are available in the table above. Note: OPeNDAP was for System). More information
netCDF (network Common Data Form)	A commonly supported self-describing data format. More information

Home » ERSST_V3n

On this page: [Temporal Coverage](#) | [Spatial Coverage](#) | [Levels](#) | [Update Schedule](#) | [Download/Plot Data](#) | [Analysis Tools](#) | [Restrictions](#) | [Details](#) | [Caveats](#) | [File Naming](#) | [Citation](#) | [References](#) | [Original Source](#) | [Contact](#)

NOAA Extended Reconstructed Sea Surface Temperature (SST) V5

Values from 2008-2018 have chanced at the source. We are changing our data updates to update older values instead of just appending values each month (as of 2020/03/23).

Brief Description:

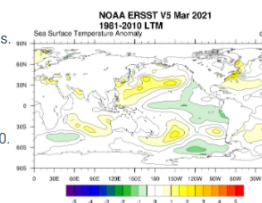
- A global monthly SST analysis from 1854 to the present derived from ICOADS data with missing data filled in by statistical methods. [More Details...](#)

Temporal Coverage:

- Monthly values for 1854/01 - present.
- Long term monthly means, derived from data for years 1981 - 2010.

Spatial Coverage:

- 2.0 degree latitude x 2.0 degree longitude global grid (89x180).
- 88.0N - 88.0S, 0.0E - 358.0E.



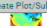

Levels:

- Sea Surface

Update Schedule:

- Variable

Download/Plot Data: ([Download Issues](#))

Variable	Statistic	Level	Units	Download File	Create Plot/Subset
Sea Surface Temperature	Monthly Mean	Surface	degC	est.monmean.nc	
Sea Surface Temperature	Monthly Long Term Mean	Surface	degC	est.mon.ltm.1981-2010.nc	

II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 각 기후감시요소에 대하여 시계열 자료 산정

자료 수집

Home > PSL Home > GPCP V2 Precipitation

On this page: [Temporal Coverage](#) | [Spatial Coverage](#) | [Levels](#) | [Update Schedule](#) | [Download/Plot Data](#) | [Analysis Tools](#) | [Restrictions](#) | [Details](#) | [Caveats](#) | [File Naming](#) | [Citation](#) | [References](#) | [Original Source](#) | [Contact](#)

GPCP Version 2.3 Combined Precipitation Data Set

Note: This dataset has been updated to version 2.3 and will be updated regularly. 10/09/2020: Grids from years since 2015 have been replaced. See the NCEI webpage for more information

Brief Description:

- Global Precipitation Climatology Project monthly precipitation dataset from 1979-present combines observations and satellite precipitation data into 2.5°x2.5° global grids.

Temporal Coverage:

- Monthly values 1979/01 through Feb 2021 (some months are interim).
- Long term monthly means, derived from years 1981 - 2010.

Spatial Coverage:

- 2.5 degree latitude x 2.5 degree longitude global grid (144x72)
- 88.75N - 88.75S, 1.25E - 358.75E

Levels:

- N/A

Update Schedule:

- Monthly

Latest available data: [Click to Enlarge](#)

Formerly the National Climatic Data Center (NCDC)... [more about NCEI](#)

Home > Climate Information > Data Access > Customer Support > Contact > About

Home > Climate Data Record Program > Terrestrial > Snow Cover Extent (Northern Hemisphere)

Snow Cover Extent (Northern Hemisphere)

NOAA CDR Snow Cover Extent (Northern Hemisphere)

0.0 0.2 0.4 0.6 0.8 1.0

Data Access

- Download
- THREDDS

Documentation

- CDR Flyer
- Use Agreement
- Algorithm Description
- Data Flow Diagram
- Maturity Matrix
- Source Code

Information

- Contacts
- Registration (optional)

II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드선스 III(기상청, 2018) 활용
 - 각 기후감시요소에 대하여 시계열 자료 산정

자료 수집

Home > Gridded Climate Data > NCEP Reanalysis Products: Pressure level variables

On this page: [Temporal Coverage](#) | [Spatial Coverage](#) | [Levels](#) | [Update Schedule](#) | [Download/Plot Data](#) | [Analysis Tools](#) | [Restrictions](#) | [Details](#) | [Caveats](#) | [File Naming](#) | [Citation](#) | [References](#) | [Original Source](#) | [Contact](#)

NCEP/NCAR Reanalysis 1: Pressure

We have transitioned the data files from netCDF3 to netCDF4-classic format on Monday Oct 20th, 2014.

Other Grid Types: [NCEP Reanalysis Main Page](#) | [Pressure Level Data](#) | [Surface Data](#) | [Surface Flux Data](#) | [Other Flux Data](#) | [Spectral Coefficients Data](#) | [Tropopause Data](#)

Brief Description:

- NCEP/NCAR Reanalysis 1

Temporal Coverage:

- 4-times daily, daily and monthly values for 1948/01/01 to present
- Long term monthly means, derived from data for years 1981 - 2010
- Values are instantaneous at the time indicated in the files

Spatial Coverage:

- 2.5 degree x 2.5 degree global grids (144x73)
- 0.0E to 357.5E, 90.0N to 90.0S

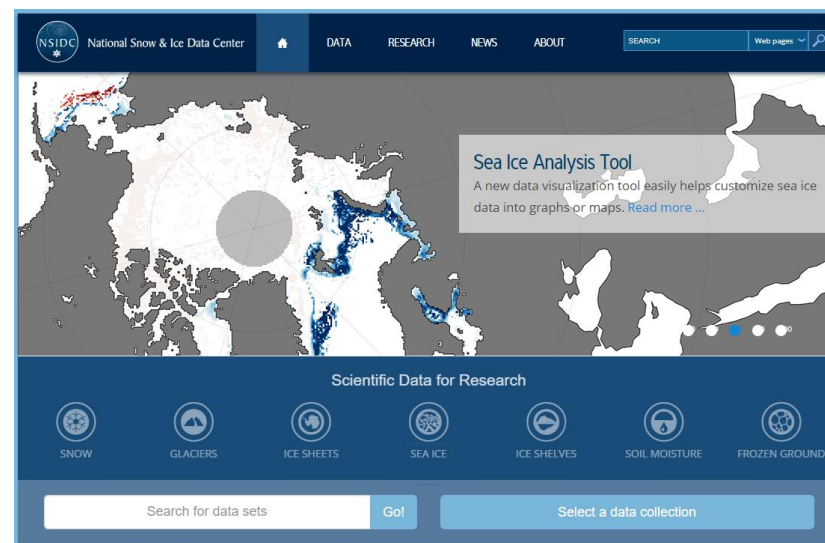
Levels:

- 17 Pressure levels (mb): 1000,925,850,700,600,500,400,300,250,200,150,100,70,50,30,20,10
- Some variables have less: omega (to 100mb) and Humidities (to 300mb).

Update Schedule:

- Daily

[Download/Plot Data: \(Download Issues\)](#)



II. 기계학습을 이용한 장기예측

자료를 수집하자

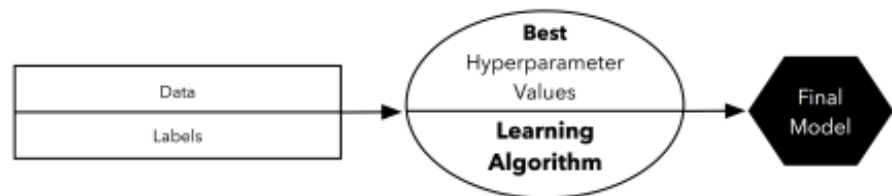
모델의 훈련/검증 및 테스트를 위한 자료 분리

- 본 예시에서는 테스트 수행하지 않고 Leave One-year Out Cross Validation

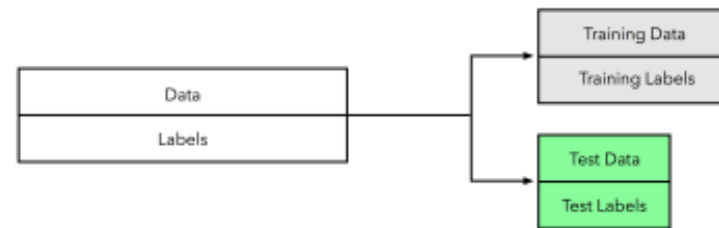
4



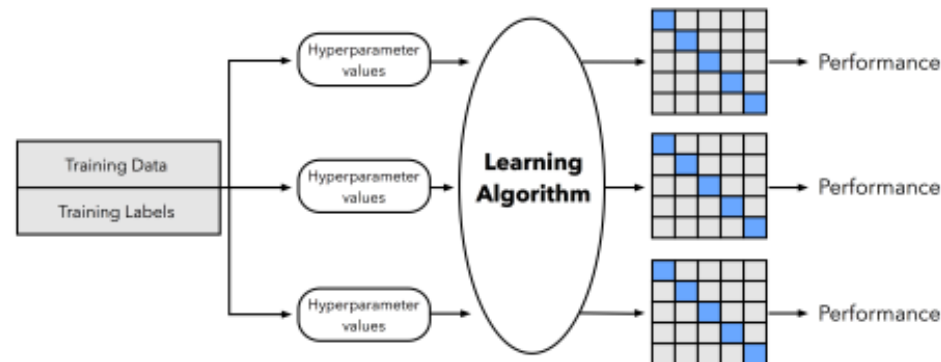
5



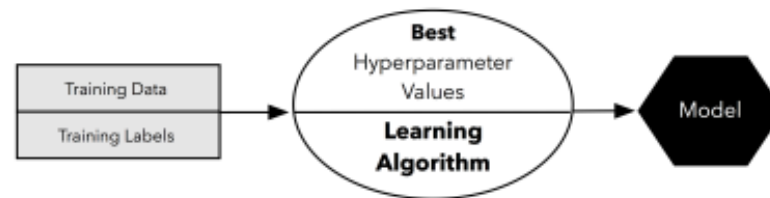
1



2



3



II. 기계학습을 이용한 장기예측

자료를 수집하자

- 모델에서 사용할 입력변수들: 장기예보 가이드스 III(기상청, 2018) 활용
 - 각 기후감시요소에 대하여 시계열 자료 산정

자료 전처리 및 입력변수(Feature) 선정

- NetCDF 등의 자료 읽기
- 기후감시요소 지역에 대해 영역평균
- 각 지역에 대해 전년도 12월 ~ 당해 년도 4월 중 실제값과 상관관계 높은 월 사용(상관관계 산정 시 검증 년도 미포함)
- 본 예시에서는 온난화 트렌드 고려하지 않음

변수	지역번호	기후감시시점(월)	r	p	n
sst	rg1_sub_2	1	-0.29	0.070	40
sst	rg2	3	0.29	0.068	40
sst	rg3	4	0.27	0.089	40
sst	rg5	4	0.48	0.002	40
sst	rg27	1	0.35	0.028	40
sst	rg26	4	0.31	0.054	40
sst	rg30	4	0.30	0.060	40
sst	rg14	4	0.53	0.000	40
sst	rg20	1	0.33	0.038	40
sst	rg19	4	0.46	0.003	40
prcp	rg12	1	0.31	0.054	40
prcp	rg12_sub_13	1	0.33	0.040	40
snow	rg8	12	-0.31	0.052	40
snow	rg8_sub_9	12	-0.34	0.034	39
snow	rg10_add_11	3	0.51	0.001	40
snow	rg9	3	0.32	0.047	40
snow	rg10	3	0.40	0.010	40
snow	rg24	3	0.39	0.012	40
snow	rg28	3	0.55	0.000	40
snow	rg11	3	0.54	0.000	40
snow	rg24_sub_25	3	0.36	0.022	40
500hpa	rg16	4	0.45	0.004	40
500hpa	rg15_add_16_sub_17	12	0.33	0.036	40
500hpa	rg15	3	0.29	0.067	40
seaice	laptev	4	-0.43	0.006	40
seaice	bering	1	-0.32	0.046	39
seaice	kara	3	-0.31	0.052	40
seaice	barents	3	-0.32	0.047	40

1980년 6월 평균온도 교차검증 예측을 위한 (1981-2020년 자료 이용) 입력변수 목록

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

지도학습, 회귀모델

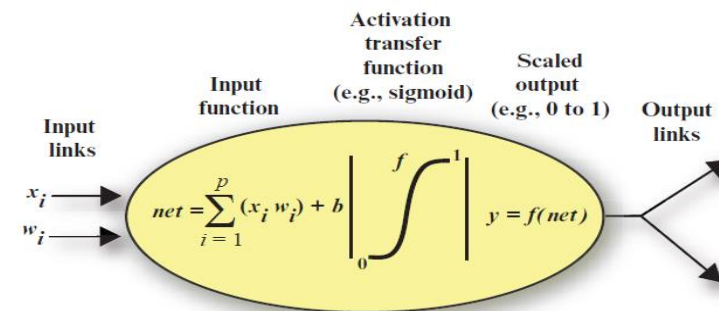
- 인공신경망
- 트리를 이용한 기계학습
 - 앙상블 / 부스팅 기법
- 서포트 벡터 머신

II. 기계학습을 이용한 장기예측

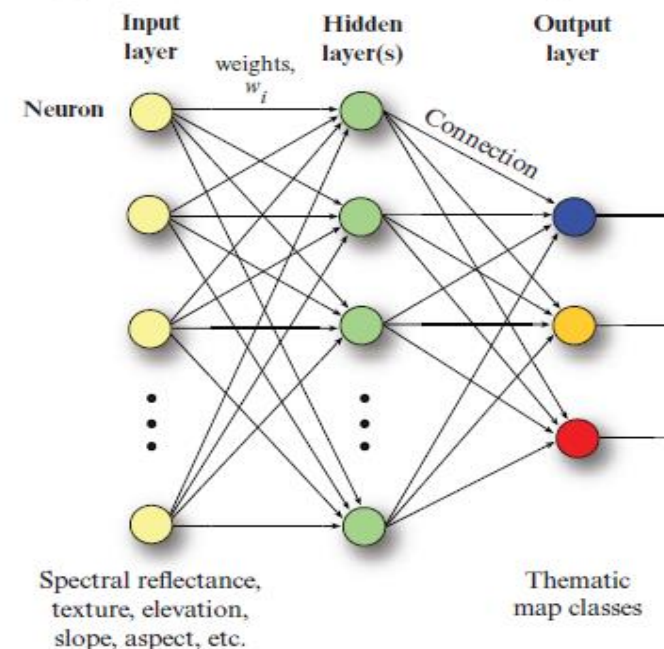
활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

- 뉴런의 출력값은 y 로 입력값(x_i)과 각각의 가중치(w_i)를 곱한 것의 합에 편향(b)을 더하여 활성 함수(activation function)를 적용한 값
- 신경망은 **입력층, 은닉층, 출력층**으로 구성되며 각 층은 노드를 통해 서로 연결됨
- 상호 연결을 통해 정보는 다중 방향으로 흐르게 되고 신경망이 훈련됨
- **상호 연결의 세기(혹은 가중치)**는 신경망에 의해 학습되고 저장되며 검증 단계에서 사용됨



Typical Artificial Neural Network Components



II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

지도학습

- **Multilayer Feedforward Backpropagation Network**
- Radial Basis Function (RBF) Network
- Convolutional Neural Network (CNN)
- Recurrent Neural Network (RNN) 등

비지도학습

- Self Organizing Map (SOM)
- Autoencoder
- Generative Adversarial Networks (GAN) 등

강화학습

- Finding the balance between exploration (new territory) and exploitation (known areas/knowledge)

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

성과와 정확도에 영향을 미치는 요소들

- 은닉층과 뉴런의 개수
- 활성화 함수
- 학습 알고리즘
- 학습률 등과 같은 매개변수 등

```
#ANN
from sklearn.neural_network import MLPRegressor

mymodel = MLPRegressor(hidden_layer_sizes = (100), \
                        activation = 'relu', \
                        solver = 'sgd', \
                        alpha = 0.1, \
                        learning_rate = 'constant', \
                        learning_rate_init = 0.001, \
                        tol = 1e-2, \
                        max_iter = 1000, \
                        momentum = 0.9, \
                        random_state = 42, \
                        )
```

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

성과와 정확도에 영향을 미치는 요소들

- 은닉층과 뉴런의 개수
- **활성 함수**
- 학습 알고리즘
- 학습률 등과 같은 매개변수 등

```
#ANN
from sklearn.neural_network import MLPRegressor

mymodel = MLPRegressor(hidden_layer_sizes = (100), \
                        activation = 'relu', \
                        solver = 'sgd', \
                        alpha = 0.1, \
                        learning_rate = 'constant', \
                        learning_rate_init = 0.001, \
                        tol = 1e-2, \
                        max_iter = 1000, \
                        momentum = 0.9, \
                        random_state = 42, \
                        )
```

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

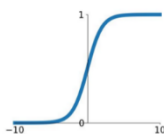
활성 함수

자료의 비선형성을 다루기 위해
선형 자료를 비선형 자료로 전환

임계점 이상이면 활성화

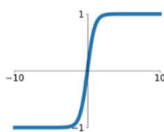
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



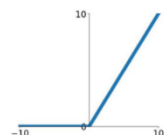
tanh

$$\tanh(x)$$



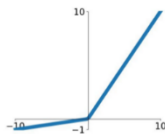
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

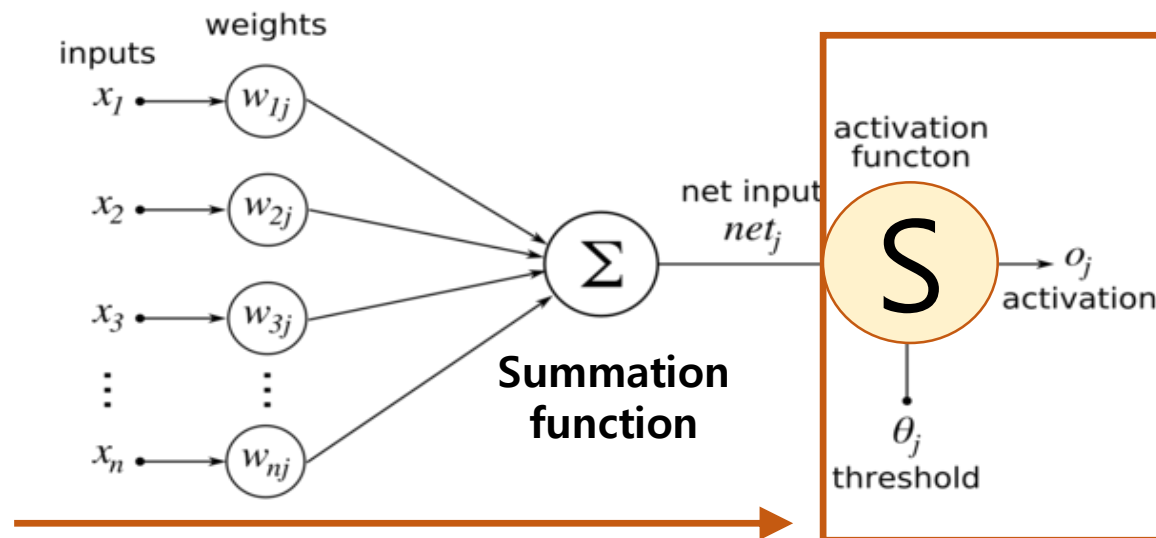
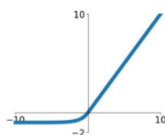


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

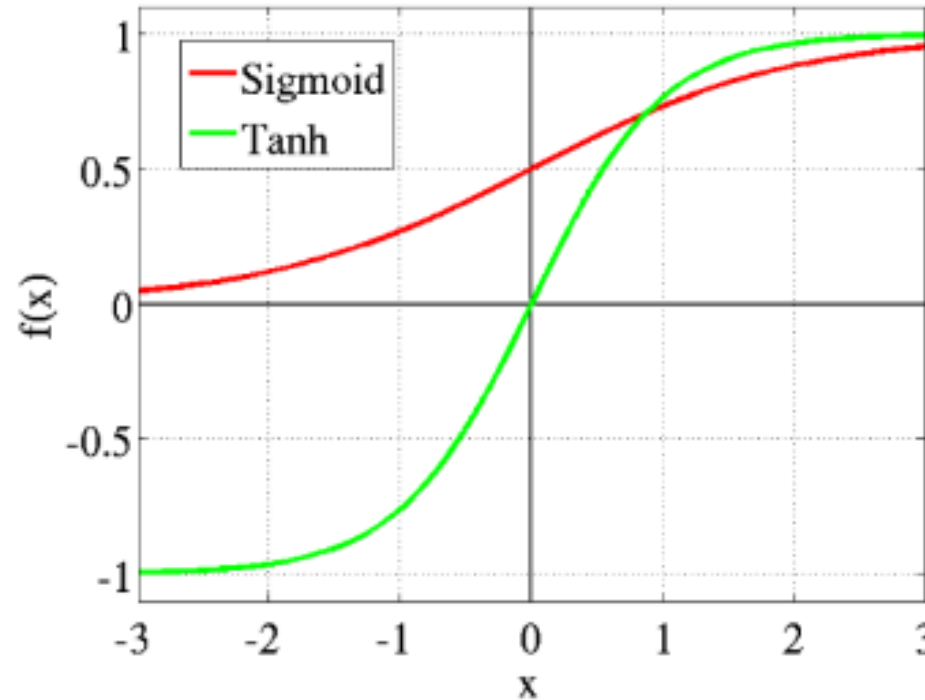


출처: UNIST IRIS Lab

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

활성 함수



Sigmoid (Logistic) / tanh (Hyperbolic tangent):
기울기 소실 문제 - x값이 크거나 작은 경우 기울기가 0에 수렴

$$h(x) = \frac{1}{1 + \exp(-x)}$$

$$h(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- 분류, 회귀에 모두 사용됨: 회귀에서는 은닉층에서만 사용되며 마지막 출력층은 활성화 함수 사용하지 않거나 활성화 함수 결과에 가중치 곱하여 사용
- 각 층마다 다른 활성화 함수 활용할 수 있음

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

성과와 정확도에 영향을 미치는 요소들

- 은닉층과 뉴런의 개수
- 활성화 함수
- **학습 알고리즘**
- 학습률 등과 같은 매개변수 등

비용 함수를 최소화 하는 가중치와 편향 찾기

- 비용함수는 일반적으로 실제값과 예측값의 제곱오차
- 전역 최대/최소를 찾기 위해 반복적으로 업데이트
- 어느 방향으로 가중치와 편향을 업데이트 해야 할지, 얼마나 움직여야 할지 중요

```
#ANN
from sklearn.neural_network import MLPRegressor

mymodel = MLPRegressor(hidden_layer_sizes = (100), \
                        activation = 'relu', \
                        solver = 'sgd', \
                        alpha = 0.1, \
                        learning_rate = 'constant', \
                        learning_rate_init = 0.001, \
                        tol = 1e-2, \
                        max_iter = 1000, \
                        momentum = 0.9, \
                        random_state = 42, \
                        )
```

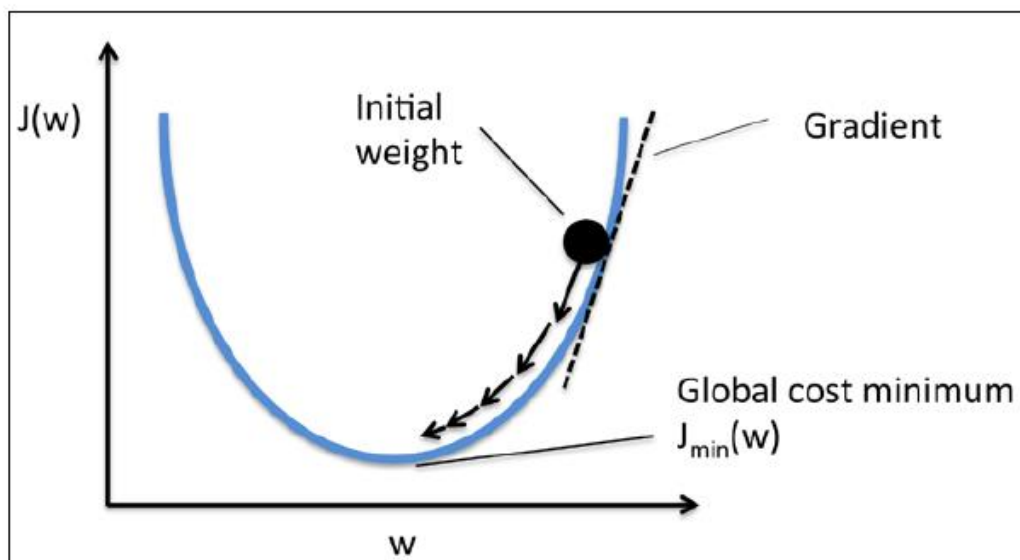
- 경사 하강법(Gradient descent)
- 다양한 경사 하강법의 변형들
 - SGD: Stochastic gradient descent (확률적 -)
 - Momentum: SGD with momentum
 - NAG: Nesterov Accelerated Gradient
 - Adagrad: Adaptive Gradient
 - Adadelat: Adaptive Delta
 - Rmsprop
- Quasi-Newton
- Levenberg-Marquardt 등

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

- 경사하강법(Gradient Descent)
 - 산을 내려가려는데 안개가 끼어 바로 주변만 보일때... 자기 위치에서 가장 큰 경사가 쪽으로 내려가야 함 → 이런 방식으로 주어진 곡선의 최소값을 찾아가는 방식이 경사하강법



비용함수 예시

$$J(w) = \frac{1}{2} \sum_i \left(y - \sum_j w_j x_j \right)^2$$

업데이트된 가중치

$$w_j = w_j + \Delta w_j = w_j - \eta \Delta J(w_j)$$

$$w_j = w_j + \eta \sum_i (y^{(i)} - \hat{y}^{(i)}) x_j^{(i)}$$

경사 즉 기울기

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

인공신경망

성과와 정확도에 영향을 미치는 요소들

- 은닉층과 뉴런의 개수
- 활성화 함수
- 학습 알고리즘
- 학습률 등과 같은 매개변수 등

```
#ANN
from sklearn.neural_network import MLPRegressor

mymodel = MLPRegressor(hidden_layer_sizes = (100), \
                        activation = 'relu', \
                        solver = 'sgd', \
                        alpha = 0.1, \
                        learning_rate = 'constant', \
                        learning_rate_init = 0.001, \
                        tol = 1e-2, \
                        max_iter = 1000, \
                        momentum = 0.9, \
                        random_state = 42, \
                        )
```

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

트리를 이용한 기계학습 - 결정트리(Decision Trees)

- 자기 발견적인 지식기반 전문가 시스템(Expert Systems)을 이용한 문제해결 방식

전문가의 머릿속에 저장되어 있는 도메인 지식은 가설(문제), 규칙, 조건들로 구성된 지식기반의 형태로 추출



사용자 인터페이스와 추론 엔진은 지식기반 규칙들을 암호화

온라인 데이터베이스로부터 필요한 정보를 추출하여 문제를 해결하는데 사용

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

트리를 이용한 기계학습 - 결정트리(Decision Trees)

- 인간 도메인 전문가와 지식공학자에 기반한 전문가 시스템

“그 지형은 태양에너지를 최대한 이용하는 주거지 개발에 적합하다, 즉 지붕에 태양패널을 설치할 수 있을 것이다”

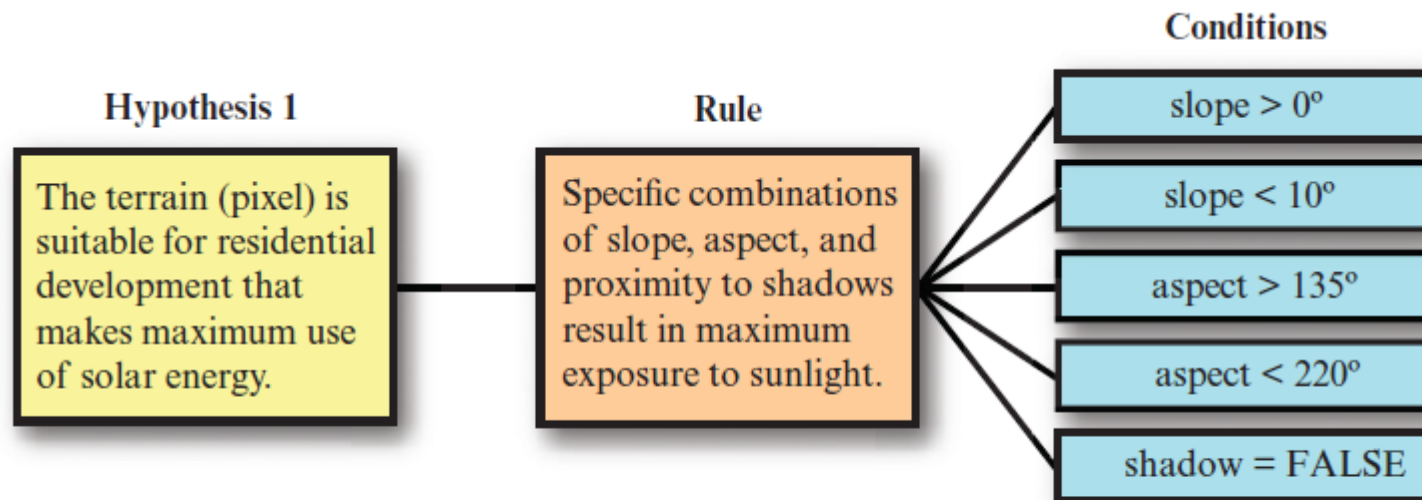


FIGURE 10-4 A human-derived decision-tree expert system with a rule and conditions to be investigated by an inference engine to test hypothesis 1.

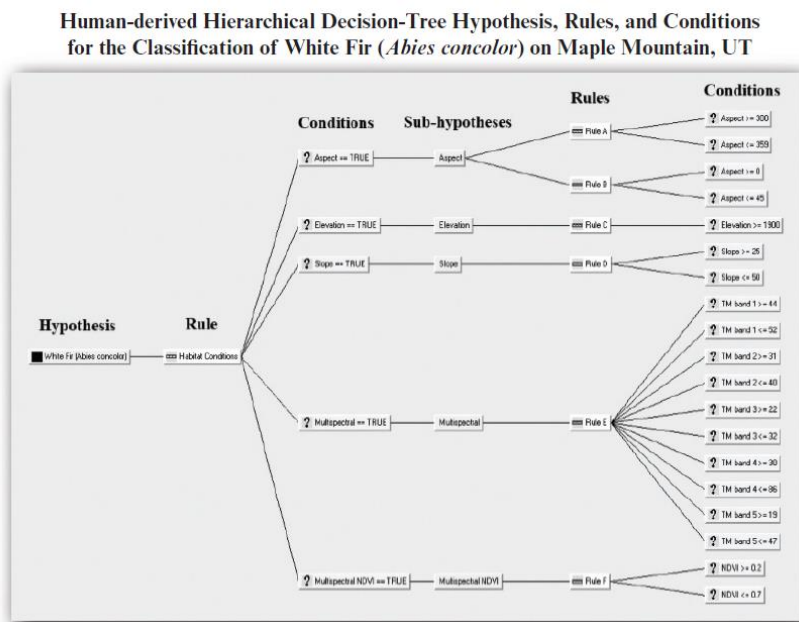
II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

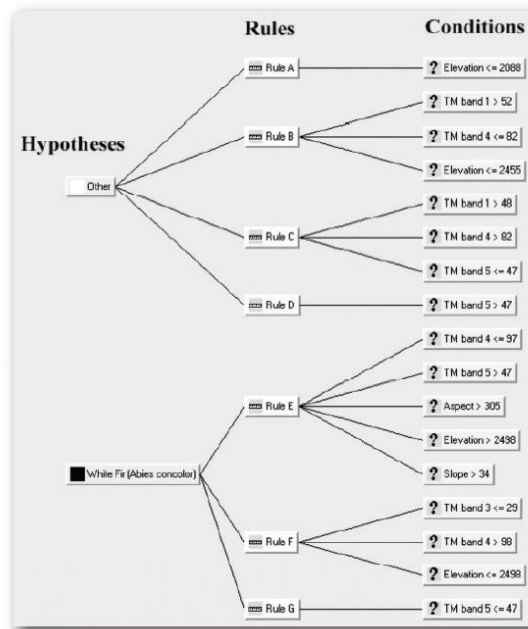
지식기반 구축의 자동화: 기계학습

- 기계학습에서는 귀납적 학습과정을 통해 전문가 시스템을 위한 지식기반을 구축
- “좋은 훈련자료가 있다면 인간 전문가로부터 일반 이론을 명확히 추출하는 것 보다 훨씬 쉽다.”

북미 전나무 분류:
인간 전문가



Machine-learning Hierarchical Decision-Tree Classification of White Fir (*Abies concolor*) on Maple Mountain, UT



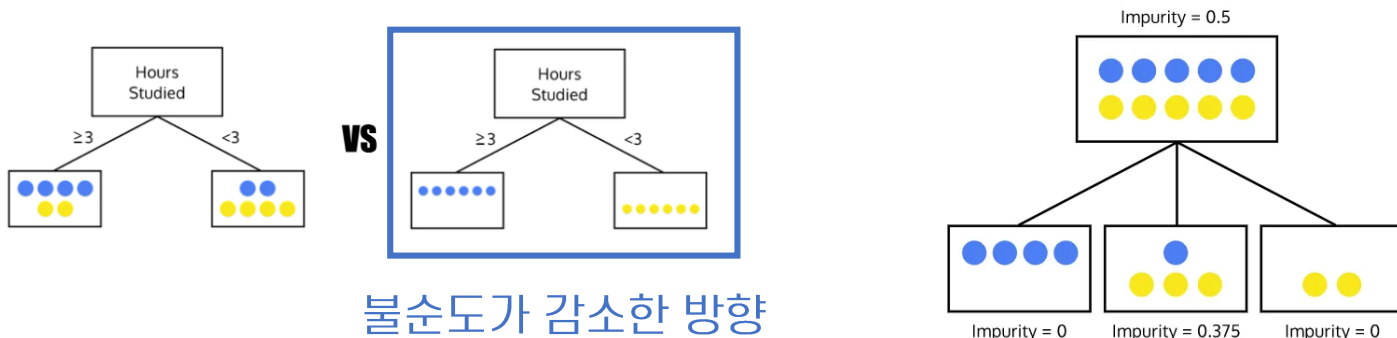
북미 전나무 분류:
결정트리

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

결정트리(Decision Trees)의 학습 원리

- 순도(homogeneity)가 증가하고 불순도(impurity) 혹은 불확실성(uncertainty)이 최대한 감소하는 방향으로 학습을 진행
- 순도가 증가 / 불확실성이 감소 = 정보획득(information gain)



$$\text{Information Gain} = \text{Mother} - \text{Children} \\ : 0.5 - (0 + 0.375 + 0) = 0.125$$

정보획득량(Information Gain)
이 커지는 방향으로 가지치기

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

트리를 이용한 기계학습 - 결정트리(Decision Trees)

전문가 시스템을 위한 지식기반을 자동적으로 구축하기 위해 귀납적 학습 기법을 적용

성과와 정확도에 영향을 미치는 요소들

- 트리의 최대 깊이(None인 경우 최대로 성장)
- 가지가 나누어질 수 있는 잎의 수
- 가지를 나눌 때 사용할 입력변수의 최대 개수 ('sqrt' 인 경우 전체 입력변수 수의 제곱근)

```
#DT
from sklearn.tree import DecisionTreeRegressor

my_model = DecisionTreeRegressor(criterion = 'mse', \
                                 max_depth = None, \
                                 min_samples_split = 2, \
                                 max_features = None, \
                                 random_state = 42, \
                                 )
```


II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

트리를 이용한 기계학습 - 결정트리(Decision Trees)

결정트리의 장점

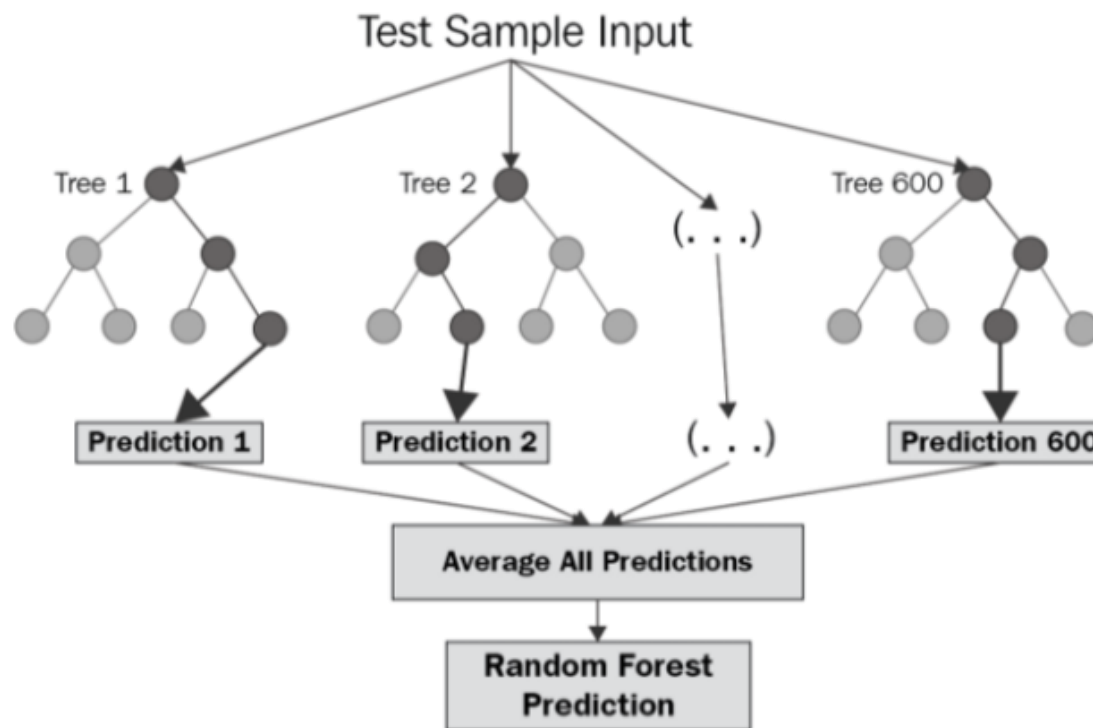
- 훈련 자료로부터 직접 규칙과 조건들을 생성하여 전문가 시스템을 훈련
- 전문가 시스템의 결과를 평가할 수 있고 결론에 어떻게 도달했는지 역으로 확인할 수 있음
- 입력변수의 분포에 대한 어떠한 가정도 필요 없음
- 입력변수 사이의 비선형적 및 계층적 관계를 밝힐 수 있음

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

앙상블 기법 - 랜덤포레스트(Random Forests)

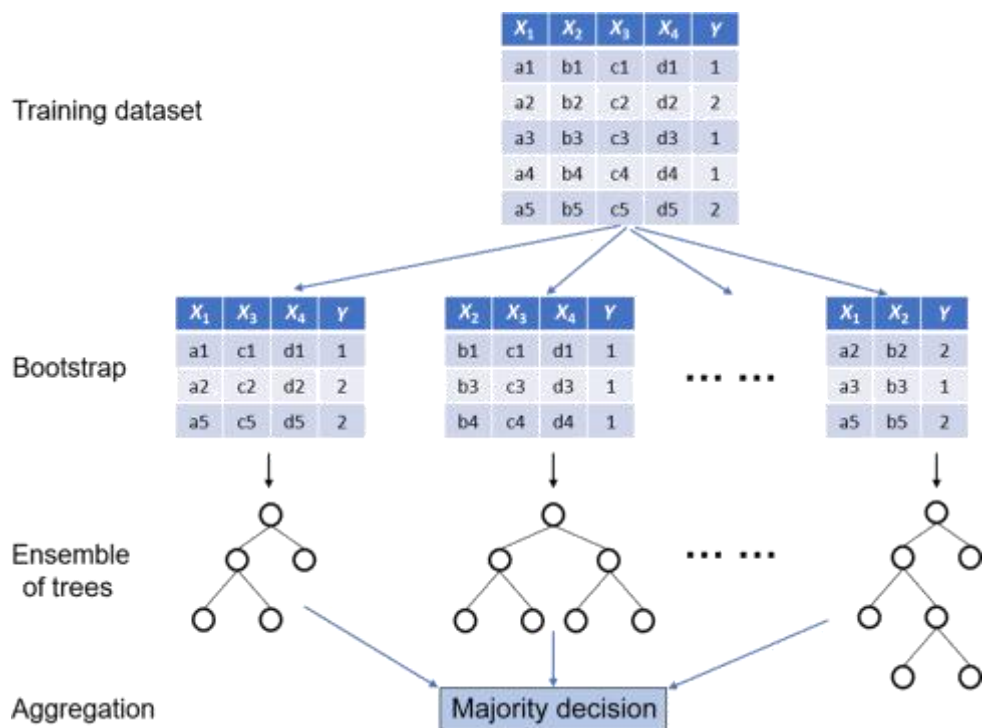
여러 결정트리를 취합하여 학습 성능을 높이는 앙상블 기법



II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

앙상블 기법 - 랜덤포레스트(Random Forests)



- 자료의 부트스트랩(Bootstrap) 랜덤 샘플링
- 각 트리에 대해 사용할 자료를 N개 무작위 복원 추출
 - 예를 들어 전체 자료의 대략 70% 사용; 나머지 30%는 OOB (out-of-bag)

- 입력변수의 랜덤 샘플링
- 각 트리에서 가지를 나눌 때마다 입력변수 중 m개를 무작위 복원 추출
 - 전체 입력변수가 M개라면 예를 들어 M의 제곱근을 사용

다수의 트리에 대해 수행하여 결과 취합 (voting, 평균)

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

앙상블 기법 - 랜덤포레스트(Random Forests)

성과와 정확도에 영향을 미치는 요소들

- Base Estimator는 결정트리
 - 트리의 최대 깊이(None인 경우 최대로 성장)
 - 가지가 나누어질 수 있는 잎의 수
 - 가지를 나눌 때 사용할 입력변수의 최대 개수 ('sqrt' 인 경우 전체 입력변수 수의 제곱근)
- 트리의 개수

```
#RF
from sklearn.ensemble import RandomForestRegressor

mymodel = RandomForestRegressor(criterion = 'mse', \
                                n_estimators = 100, \
                                max_depth = None, \
                                min_samples_split = 2, \
                                max_features = None, \
                                bootstrap = True, \
                                random_state = 42
                                )
```

II. 기계학습을 이용한 장기예측

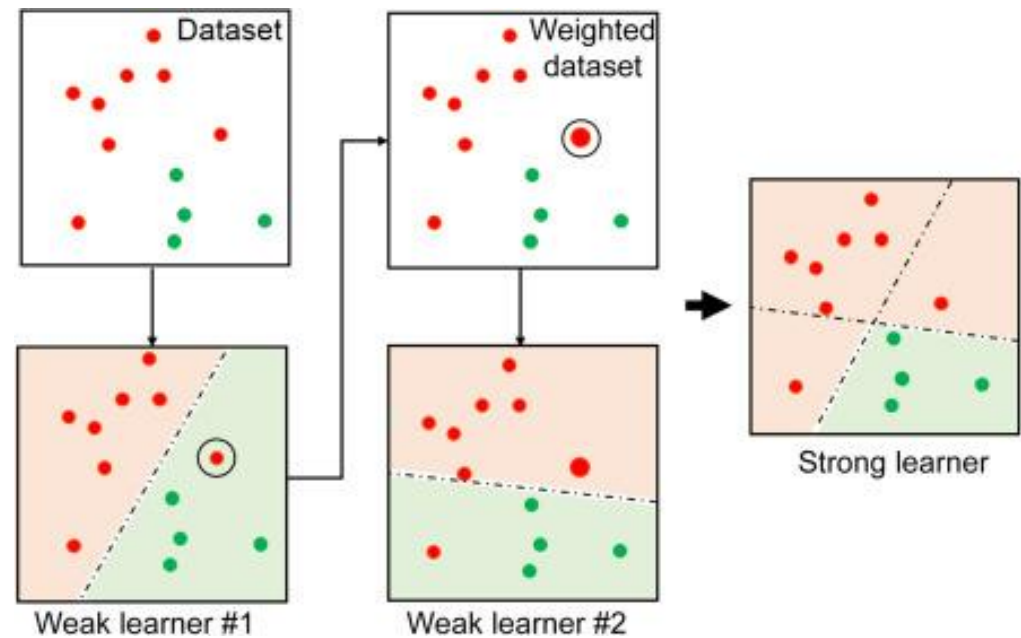
활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

부스팅 기법 - 에이다부스트(Adaptive Boosting; AdaBoost)

여러 학습 알고리즘의 결과물에 가중치를 주어 학습 성능을 높이는 가속화(부스팅) 기법

메타 알고리즘

에이다부스트는 약한 학습기(weak learner)의 결과물에 가중치를 주어 성능을 향상시킴



출처: Misra and Li, 2020

II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

부스팅 기법 - 에이다부스트(Adaptive Boosting; AdaBoost)

성과와 정확도에 영향을 미치는 요소들

- Base Estimator는 결정트리
 - 트리의 최대 깊이(None인 경우 최대로 성장)
 - 가지가 나누어질 수 있는 잎의 수
 - 가지를 나눌 때 사용할 입력변수의 최대 개수('sqrt' 인 경우 전체 입력변수 수의 제곱근)
- 트리의 개수
- 학습률: 새로운 학습자가 얼마나 기여하도록 할지 정함(0~1)
- 손실 함수: 가중치 업데이트할 때 사용

```
#Adaboost
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import AdaBoostRegressor

mydt = DecisionTreeRegressor(criterion = 'mse', \
                             max_depth = 5, \
                             min_samples_split = 2, \
                             max_features = None, \
                             random_state = 42, \
                             )
mymodel = AdaBoostRegressor(base_estimator = mydt, \
                             n_estimators = 500, \
                             learning_rate = 1, \
                             loss = 'square', \
                             random_state = 42
                             )
```

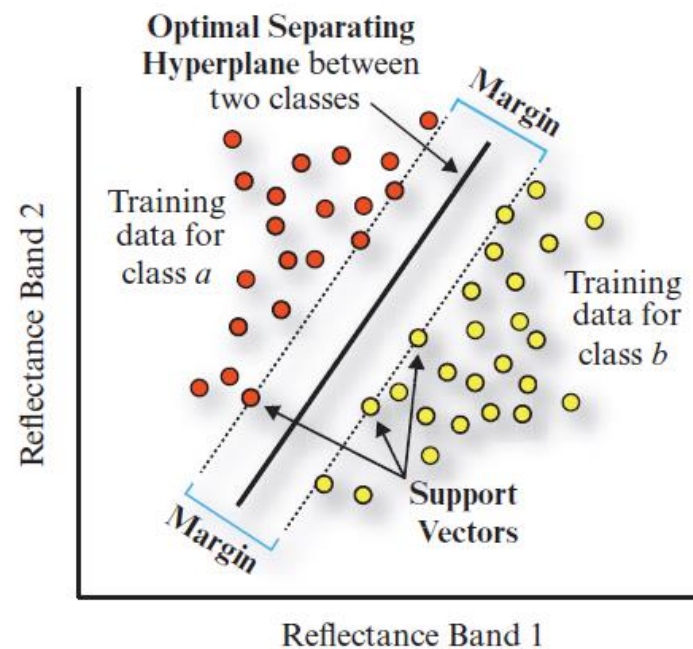
II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

서포트 벡터 머신(Support Vector Machine)

- 훈련 자료를 이용해 클래스 간을 분리하는 최적의 초평면(hyperplane)을 찾음: 클래스의 가장 가까운 훈련 샘플들 사이의 마진을 최대화
- 경계에 놓여 있는 포인트들이 **서포트 벡터**라고 불리며 그 마진의 중간이 클래스를 분리하는 최적의 초평면
- 잘못 분리된 훈련 자료들에 음의 가중치를 주어 영향을 줄임

Support Vector Machine Classification (Linear)

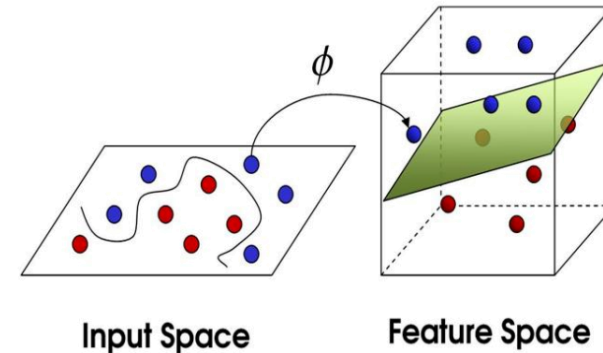


II. 기계학습을 이용한 장기예측

활용할 기법을 선정하고 모델 학습을 위한 최적화 방법을 정하자

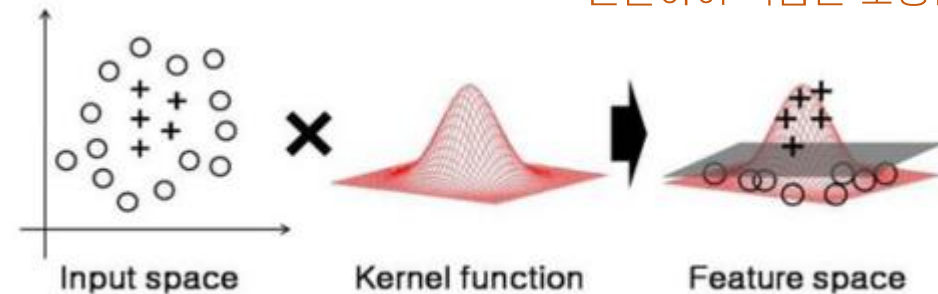
서포트 벡터 머신(Support Vector Machine)

- 선형의 초평면을 찾는 것이 불가능할 경우 커널 함수를 이용하여 고차원 공간으로 변환하여 적합한 초평면을 찾음
- 널리 이용되는 SVM 커널 함수는 선형, 다항식, 가우시안(RBF), 시그모이드 등이 있음
- 클래스의 경계에서 선택된 적은 훈련 자료만으로도 꽤 정확한 분류를 수행함



출처: UNIST IRIS Lab

RBF Kernel을 이용하여 3차원 공간으로 변환하여 적합한 초평면을 찾음



출처: bskyvision 블로그

II. 기계학습을 이용한 장기예측

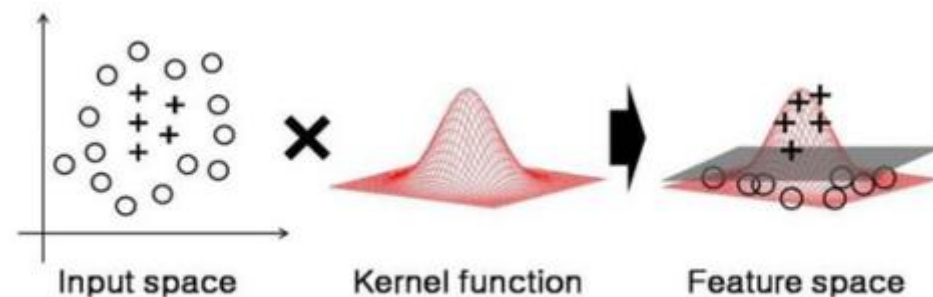
활용할 기법을 선정하고
 — 모델 학습을 위한 최적화 방법을 정하자

서포트 벡터 회귀(Support Vector Regression)

- Kernel: 알고리즘에서 사용하는 커널 타입
- Gamma: 커널의 계수
 - 하나의 자료 샘플이 영향력을 행사하는 거리를 결정
 - RBF의 경우 표준편차와 반비례하는 계수로 gamma가 크면 영향력을 행사하는 거리가 짧아지고 gamma가 낮으면 커짐
 - 'scale' 로 지정하면 $1 / (\text{입력 변수의 개수} * X \text{의 분산})$ 이용
- C: 정규화 매개변수
 - 정규화 강도는 C에 반비례
 - C가 너무 낮으면 과소적합, 너무 크면 과대적합 가능성
- Epsilon: 실제값과 예측값 차이 허용 범위

```
#SVR
from sklearn.svm import SVR

myModel = SVR(kernel = 'rbf', \
               gamma = 'scale', \
               C = 1.0, \
               epsilon = 0.1, \
               )
```



출처: bskyvision 블로그

II. 기계학습을 이용한 장기예측

모델 평가를 위한 기준을 정하자

단정예측을 위한 모델 평가

$$r_{XY} = \frac{\frac{\sum_i^n (x_i - \bar{X})(y_i - \bar{Y})}{n-1}}{\sqrt{\frac{\sum_i^n (x_i - \bar{X})^2}{n-1}} \sqrt{\frac{\sum_i^n (y_i - \bar{Y})^2}{n-1}}}$$

- 피어슨 상관계수
- 평균절대오차(Mean Absolute Error; MAE)

$$MAE = \frac{\sum_{i=1}^n |y_i - x_i|}{n} = \frac{\sum_{i=1}^n |e_i|}{n}$$

삼분위로 나타냈을 때

- Heidke Skill Score (HSS)
- Total Accuracy = Proportion Correct
- Hit Rate = Probability of Detection
- False Alarm Rate = Probability of False Detection

TABLE B1. Contingency table for tercile forecasts.

Predicted category	Observed category			
	AN	NN	BN	
AN	<i>a</i>	<i>b</i>	<i>c</i>	PRED _{AN}
NN	<i>d</i>	<i>e</i>	<i>f</i>	PRED _{NN}
BN	<i>g</i>	<i>h</i>	<i>i</i>	PRED _{BN}
	OBS _{AN}	OBS _{NN}	OBS _{BN}	Total, <i>n</i>

$$HSS = \frac{\sum_k p(\text{PRED}_k, \text{OBS}_k) - \sum_k p(\text{PRED}_k)p(\text{OBS}_k)}{1 - \sum_k p(\text{PRED}_k)p(\text{OBS}_k)}$$

$$PC = \sum_k p(\text{PRED}_k, \text{OBS}_k)$$

II. 기계학습을 이용한 장기예측

Toy Example 단정예측 결과

인공신경망

8월 평균기온 예측결과(단정예측을 삼분위에 적용)

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	0.40	0.61	0.46	0.45	0.50	0.56	0.73
Jul	0.23	0.49	0.42	0.85	0.33	0.54	0.56
Aug	0.16	0.44	0.33	0.91	0.43	0.54	0.36

YEAR	TRUE	PRED
2016	AN	AN
2017	NN	NN
2018	AN	NN
2019	AN	BN
2020	AN	AN

트리를 이용한 기계학습 - 결정트리(Decision Trees)

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	-0.01	0.34	0.27	0.63	0.00	0.38	0.53
Jul	0.19	0.46	0.28	1.18	0.50	0.38	0.50
Aug	0.27	0.51	0.23	1.07	0.43	0.62	0.50

YEAR	TRUE	PRED
2016	AN	NN
2017	NN	NN
2018	AN	AN
2019	AN	AN
2020	AN	BN

II. 기계학습을 이용한 장기예측

Toy Example 단정예측 결과

8월 평균기온 예측결과(단정예측을 삼분위에 적용)

앙상블 기법 - 랜덤포레스트(Random Forests)

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	0.09	0.44	0.40	0.47	0.00	0.81	0.33
Jul	0.30	0.54	0.41	0.85	0.33	0.69	0.56
Aug	0.24	0.49	0.51	0.81	0.29	0.77	0.43

YEAR	TRUE	PRED
2016	AN	AN
2017	NN	NN
2018	AN	NN
2019	AN	NN
2020	AN	NN

부스팅 기법 - 에이다부스트(Adaptive Boosting; AdaBoost)

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	0.10	0.44	0.45	0.48	0.00	0.69	0.47
Jul	0.27	0.51	0.37	0.89	0.42	0.54	0.56
Aug	0.20	0.46	0.45	0.84	0.29	0.69	0.43

YEAR	TRUE	PRED
2016	AN	AN
2017	NN	NN
2018	AN	NN
2019	AN	BN
2020	AN	NN

서포트 벡터 회귀(Support Vector Regression)

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	0.11	0.44	0.49	0.43	0.10	0.63	0.47
Jul	0.27	0.51	0.42	0.80	0.33	0.69	0.50
Aug	0.20	0.46	0.49	0.81	0.21	0.77	0.43

YEAR	TRUE	PRED
2016	AN	NN
2017	NN	NN
2018	AN	NN
2019	AN	NN
2020	AN	AN

II. 기계학습을 이용한 장기예측

풀고 싶은 문제를 정의하자: 확률예측

지도학습 활용 의 단계

풀고 싶은 문제를
정의하자

자료를
수집하자
(실제값)

활용할 기법을
선택하자

모델 학습을 위한
최적화 방법을
정하자

모델 평가를 위한
기준을 정하자

확률분포 예측(베이지안)
→ 각 카테고리별 확률

5월(ex. 5/15)에 예측한 6, 7, 8월 평균 기온 확률 예측

기온 값 자체를 예측
→ 단정예측

5월(ex. 5/15)에 예측한 6, 7, 8월 평균 기온 예측

삼분위 카테고리로 분류
→ 각 카테고리별 확률
→ 확률예측

II. 기계학습을 이용한 장기예측

확률과 불확실성

경험적 확률(빈도주의) (여러 번 되풀이 할 수 있는) 사건의 빈도

동전던지기
앞면이 나올 확률은?

$$P(\text{앞면}) = \frac{\text{앞면이 나온 수}}{\text{지금까지 던진 수}}$$

베이저안 확률 (여러 번 되풀이 할 수 없는 경우도 포함하여) 사건에 대한 믿음

백두산 화산 폭발의 확률은?
확률 = 믿음

$$p(w|D) = \frac{p(D|w)p(w)}{p(D)}$$

업데이트된 믿음 (posterior probability) 사건 w가 있을 때 단서 D가 발생할 Likelihood 특정 사건 w에 대해 기존에 가지고 있던 믿음 (prior probability)

새로운 단서 (evidence)

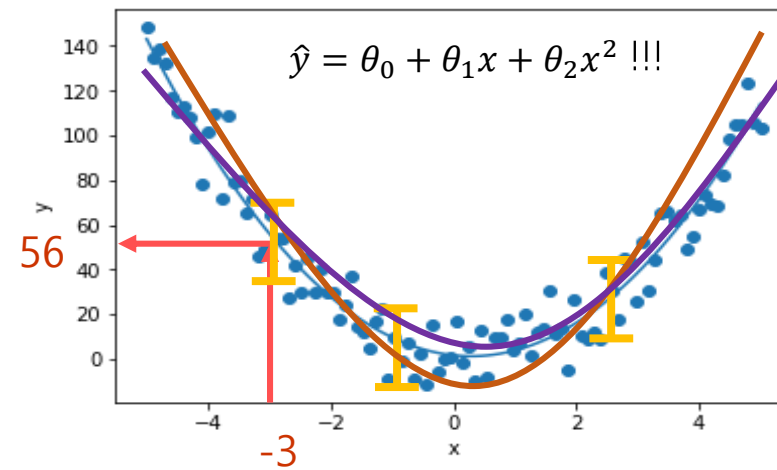
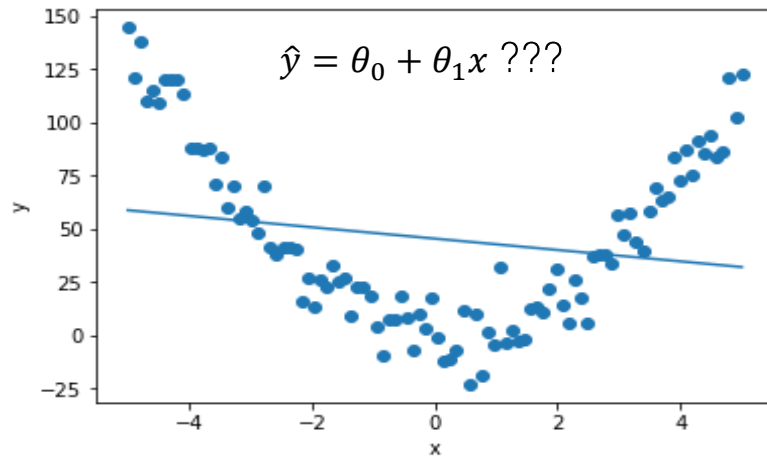


사건과 관련된 여러 확률을 이용해 새로 일어날 사건의 확률을 추정
이러한 단서들이 추가되면 사건에 대한 불확실성을 정량화

II. 기계학습을 이용한 장기예측

베이지안 추론

전통적 추론(Classical Inference)

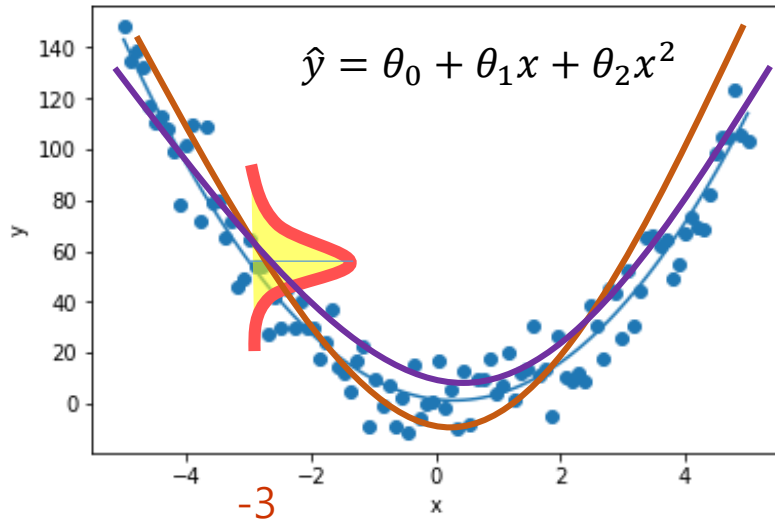


- 우리의 목적: 새로운 x_* 에 대한 y 값 추정하도록 → 모집단의 특성을 파악(즉 적합한 함수 도출)
- 전통적 추론: 관측된 자료(훈련 자료)에 관한 함수를 사용하여 그 함수가 가지는 미지의 모수를 추론
- 하지만 전통적 추론에 의한 결과는 자료에 대한 불확실성 정보를 알려주지 않음!

II. 기계학습을 이용한 장기예측

베이지안 추론

베이지안 추론(Bayesian Inference)



- 우리가 추론하고자 하는 모수는 불확실하며 이 불확실성의 정도를 확률모델로 표현하고 싶음
- 모수 각각에 대한 확률모델과 관측된 자료를 사용
- 확률예측을 통해 불확실성 정보를 제공

관측된 자료 y 에 대한 확률모델

$$\text{사후확률 } p(\theta|y) = \frac{p(y|\theta) p(\theta)}{p(y)} \text{ 사전확률}$$

- 모수에 관한 과거의 경험이나 사전 지식 같은 주관적 견해를 수량화한 모수의 특성(사전확률분포)을
- 자료로부터 얻은 모수에 관한 정보(관측된 자료)와 결합하여
- 사후확률분포를 얻음(확률예측결과를 제공)

$\theta_0, \theta_1, \theta_2$ 각각의 사전확률분포
입력 x 대해 관측된 자료 y



$\theta_0, \theta_1, \theta_2$ 각각의 사후확률분포

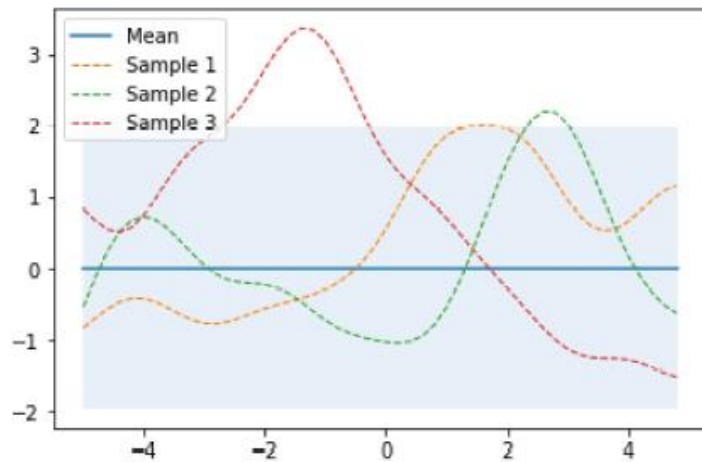


새로운 x_* 에 대한 예측 값 \hat{y} 의 분포

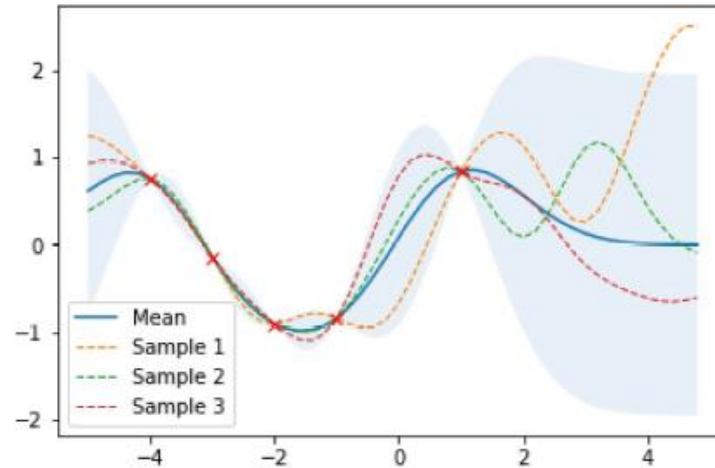
II. 기계학습을 이용한 장기예측

가우시안 프로세스: 비모수 베이지안 모델

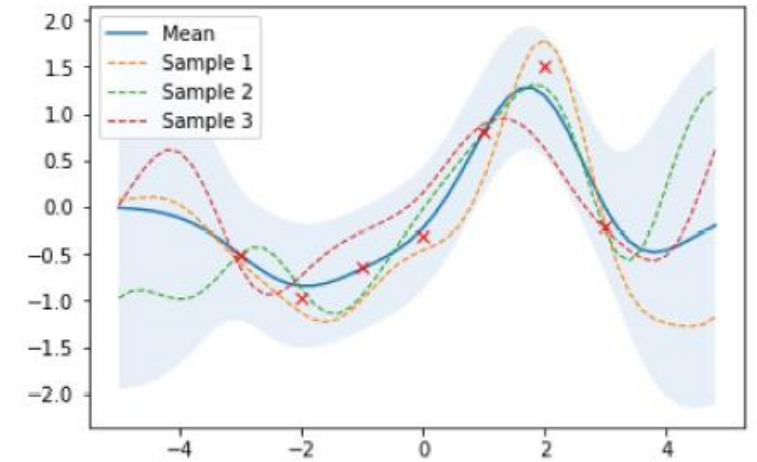
- 앞서 본 예시처럼 어떤 함수를 사용할지(모수가 몇 개인지, 모수가 어떤 분포를 따를지) 규정하는 대신
- “모든 가능한 함수”에 사전확률을 주고 더 그럴듯한 함수에 높은 사후확률을 주면 어떨까?
- 각각의 함수를 특정한 입력 x 에서 함수값 $f(x)$ 를 가지는 굉장히 긴 벡터로 생각하고,
- 무한한 모든 함수를 고려할 필요 없이 관측된 유한한 지점에서의 함수값의 특성만을 보면 됨



Three samples from prior



Three samples from posterior
Noise-free data

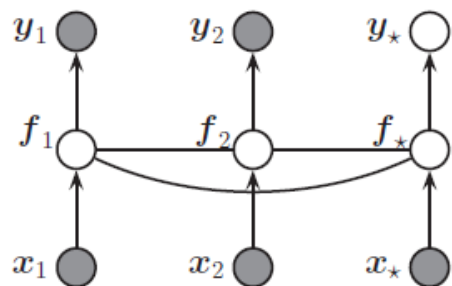


Three samples from posterior
Noisy data

II. 기계학습을 이용한 장기예측

가우시안 프로세스 모델

Gaussian Process는 확률과정(stochastic process)으로 임의의 점 $x \in R^d$ 에서 확률변수 $f(x)$ 를 가지며 유한한 확률변수들의 결합확률분포 $p(f(x_1), \dots, f(x_N))$ 는 Gaussian



$$p(f|X) = N(f|\mu, K)$$

$$\mu(x) \quad \text{평균 함수}$$

$$K_{ij} = \kappa(x_i, x_j) \quad \text{공분산함수(커널)}$$

x_i 와 x_j 가 가까우면 $f(x_i)$ 와 $f(x_j)$ 도 가까울 것이라 가정



공분산함수(커널)에 의해 결정

II. 기계학습을 이용한 장기예측

가우시안 프로세스 모델

이미 관측한 자료 y 와 예측 f_* 의 결합분포 $\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim N\left[0, \begin{pmatrix} \mathbf{K}_y & \mathbf{K}_* \\ \mathbf{K}_*^T & \mathbf{K}_{**} \end{pmatrix}\right]$

GP 사전분포 $p(\mathbf{f}|\mathbf{X}) = N(\mathbf{f}|\boldsymbol{\mu}, \mathbf{K})$ GP 사후분포 $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$

$$\mathbf{K}_y = \kappa(\mathbf{X}, \mathbf{X}) = \mathbf{K}$$

사후 예측분포 $p(\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int p(\mathbf{f}_*|\mathbf{f}, \mathbf{X}_*)p(\mathbf{f}|\mathbf{X}, \mathbf{y})d\mathbf{f}$.

$$\mathbf{K}_* = \kappa(\mathbf{X}, \mathbf{X}_*)$$

$$\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$$

RBF 커널 및 매개변수 $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left[-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M}(\mathbf{x}_i - \mathbf{x}_j)\right] + \sigma_y^2 \delta_{ij}$

새로운 x_* 에 대한 예측 f_*

$$p(\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}) = N(\mathbf{f}_*|\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\boldsymbol{\mu}_* = \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_*^T \mathbf{K}_y^{-1} \mathbf{K}_*$$

↑
Vertical variation parameter

↑
Length scale parameter

↑
Noise parameter

커널의 매개변수는 LML (로그 주변 우도함수) 최대화하도록 결정

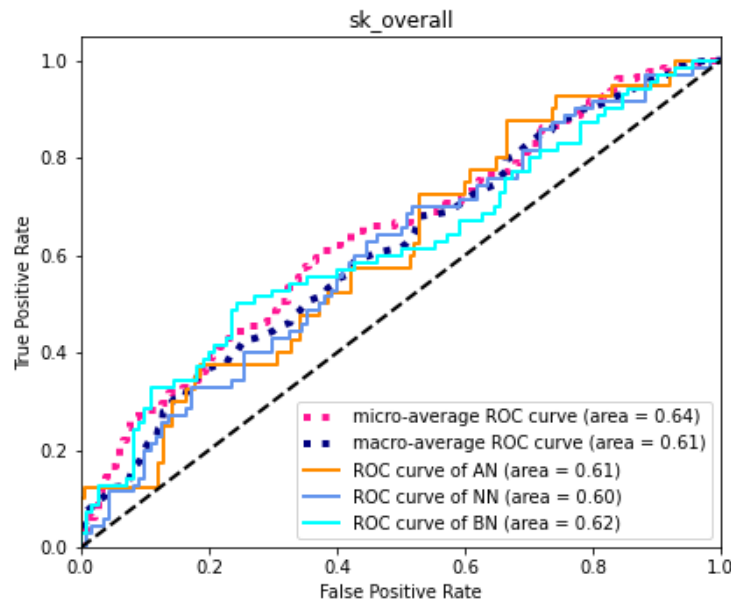
$$\log p(\mathbf{y}|\mathbf{X}) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_y) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{N}{2} \log(2\pi)$$

II. 기계학습을 이용한 장기예측

모델 평가를 위한 기준을 정하자

확률예측을 위한 모델 평가

- AUC (Area Under the Receiver Operating Characteristic Curve)
 - False positive rate을 x축에, True positive rate을 y축에 둔 ROC 곡선 아래 면적
 - Perfect Forecast: 1.0
 - Random Forecast: 0.5
- RPSS (Rank Probability Skill Score)



$$RPS = \sum_k (Y_k - O_k)^2,$$

$$RPS_{Ref} = \sum_k (P_k - O_k)^2,$$

$$RPSS = 1 - \frac{RPS}{RPS_{Ref}},$$

where k refers to categories of AN, NN, and BN; Y_k , O_k , and P_k denote cumulative forecast probability, cumulative observation probability, and cumulative climatological probability for category k , respectively. A perfect set of forecasts would be scored as 1.0, and a set of random forecasts would be scored as zero.

II. 기계학습을 이용한 장기예측

Toy Example 확률예측 결과

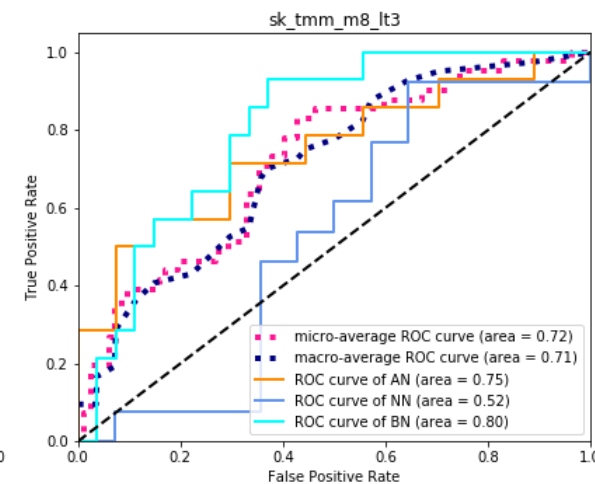
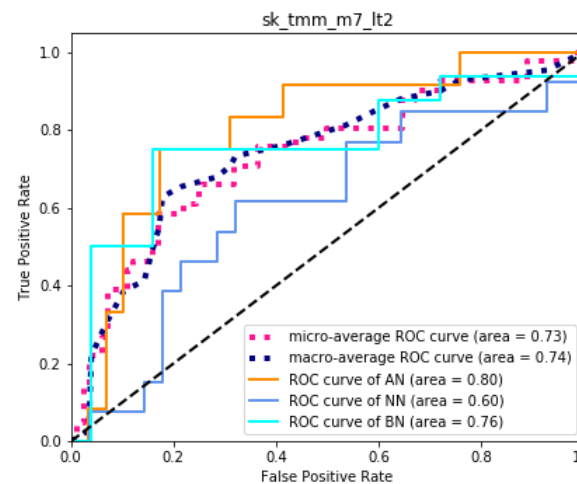
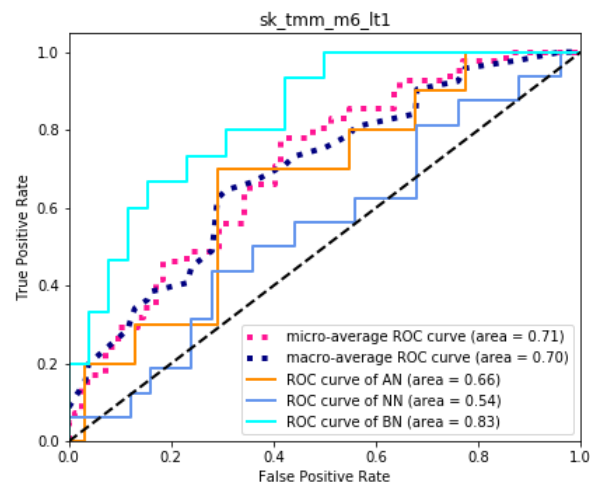
가우시안 프로세스

Month	HSS	Total Accuracy	r	MAE	HR (AN)	HR (NN)	HR (BN)
Jun	0.25	0.50	0.50	0.45	0.48	0.44	0.59
Jul	0.33	0.56	0.42	0.82	0.49	0.49	0.68
Aug	0.28	0.52	0.46	0.84	0.57	0.31	0.67

Month	AUC	RPSS
Jun	0.70	0.10
Jul	0.74	0.13
Aug	0.71	0.13

8월 평균기온 예측결과 (확률예측을 삼분위에 적용)

YEAR	TRUE	PRED
2016	AN	AN
2017	NN	NN
2018	AN	NN
2019	AN	BN
2020	AN	AN



오늘의 학습

I. 기계학습이란

- 인공지능이란
- 기계학습이란
- 기계학습의 종류

II. 기계학습을 이용한 장기예측

- 장기예측을 위한 준비
- 다양한 기계학습을 이용한 기온예측 – Toy Examples



감사합니다.

