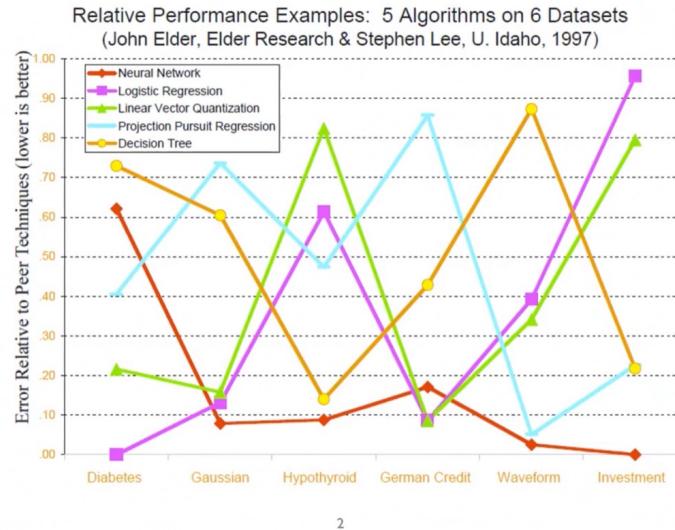


Ensemble Learning.

- Can we have a superior algorithm for all datasets?

✓ Every algorithm scored best or next-to-best on at least two of the six data sets.

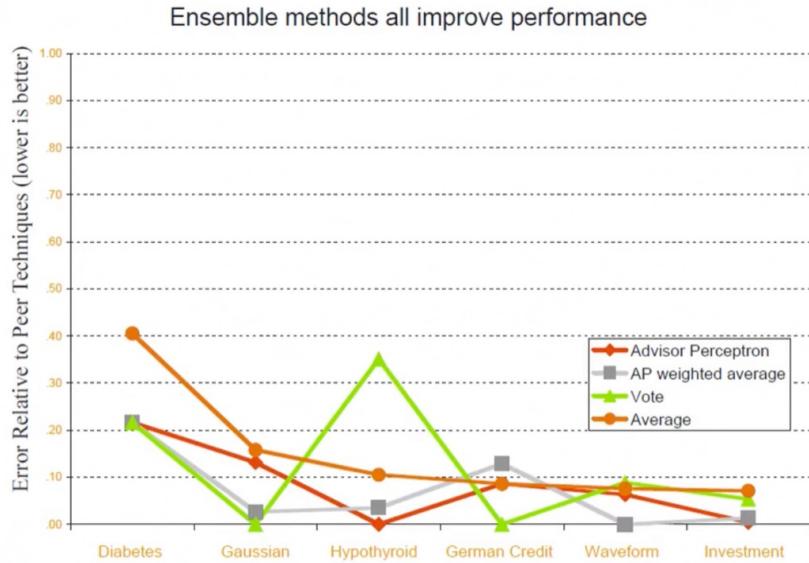


모든 상황에서 우월한 알고리즘? X

⇒ 이것저것 해보는게 훨씬 안전
(여러 알고리즘을 공부하는 이유).

However, if they are properly combined ...

≥ best of individual algorithms.



No Free Lunch Theorem.

179 algorithms
& 121 data sets.

Journal of Machine Learning Research 15 (2014) 3133-3181

Submitted 11/13; Revised 4/14; Published 10/14

Do we Need Hundreds of Classifiers to Solve Real World Classification Problems?

Manuel Fernández-Delgado

MANUEL.FERNANDEZ.DELGADO@USC.ES

Eva Cernadas

EVA.CERNADAS@USC.ES

Senén Barro

SELEN.BARRO@USC.ES

CITIUS: Centro de Investigación en Tecnologías da Información da USC

University of Santiago de Compostela

Campus Vida, 15872, Santiago de Compostela, Spain

Dinani Amorim

DINANIAMORIM@GMAIL.COM

Departamento de Tecnologia e Ciências Sociais- DTCS

Universidade do Estado da Bahia

Av. Edgard Chastinet S/N - São Geraldo - Juazeiro-BA, CEP: 48.305-680, Brasil

Empirical Evidence

- Empirical study 3: 179 algorithms on 121 datasets

Rank	Acc.	κ	Classifier	Rank	Acc.	κ	Classifier
32.9	82.0	63.5	parRF.t (RF)	67.3	77.7	55.6	pda.t (DA)
31.3	82.3	63.6	rL1t (RF)	67.6	78.7	55.2	elm.m (NNET)
36.8	81.8	62.2	svm.C (SVM)	67.6	77.8	54.2	SimpleLogistic.w (LMR)
38.0	81.2	60.1	svmPoly.t (SVM)	69.2	78.3	57.4	MAB.J48.w (BST)
39.4	81.9	62.5	rforest.R (RF)	69.8	78.8	56.7	BG.REPTree.w (BAG)
39.6	82.0	62.0	elm.kernel.m (NNET)	69.8	78.1	55.4	SMO.w (SVM)
40.3	81.4	61.1	svmRadialCost.t (SVM)	70.6	78.3	58.0	MLP.w (NNET)
42.5	81.0	60.0	svmRadial.t (SVM)	71.0	78.8	58.2	BG.RandomTree.w (BAG)
42.9	80.6	61.0	C5.0.t (BST)	71.0	77.1	55.1	mlmr.R (GLM)
44.1	79.4	60.5	avNNTen.t (NNET)	71.0	77.8	56.2	BG.J48.w (BAG)
45.5	79.5	61.0	nnet.t (NNET)	72.0	75.7	52.6	rbf.t (NNET)
47.0	78.7	59.4	pcANNet.t (NNET)	72.1	77.1	54.8	fda.t (DA)
47.1	80.8	53.0	BG.LibSVM.w (BAG)	72.4	77.0	54.7	kda.t (DA)
47.3	80.3	62.0	mlpt.t (NNET)	72.4	79.1	55.6	svmlight.C (NNET)
47.6	80.6	60.0	RotationForest.t (RF)	72.6	78.4	57.9	AdaBoostML.J48.w (BST)
50.1	80.9	61.6	RRF.t (RF)	72.7	78.4	56.2	BG.LibSVM.w (BAG)
51.6	80.7	61.4	RRFglobal.t (RF)	72.9	77.1	54.6	lbaBag.R (BAG)
52.5	80.6	58.0	MAB.LibSVM.w (BST)	73.2	78.3	56.2	BG.LibWLS.w (BAG)
52.6	79.9	56.9	MAB.LibSVM.w (SVM)	73.7	77.9	56.0	MAB.REPTree.w (BST)
57.6	79.1	59.3	adaBoost.R (BST)	74.0	77.4	52.6	RandomSubSpace.w (DT)
58.5	79.7	57.2	pmn.m (NNET)	74.4	76.9	54.2	lds2.t (DA)
58.9	78.5	54.7	cforest.t (RF)	74.6	74.1	51.8	svmlag.Bag.R (BAG)
59.9	79.7	42.6	dkp.C (NNET)	74.6	77.5	55.2	LibLINEAR.w (SVM)
60.4	80.1	55.8	gaussprad.RLLR (OM)	75.9	77.2	55.6	rbfDDAD.t (NNET)
60.5	80.0	57.4	RandomForest.w (RF)	76.5	76.9	53.8	sda.t (DA)
62.1	78.7	56.0	svmLinear.t (SVM)	76.6	78.1	56.5	END.W (OEN)
62.5	78.4	57.5	fda.t (DA)	76.6	77.3	54.8	LogitBoost.w (BST)
62.6	78.6	56.0	knn.t (NN)	76.6	78.2	57.3	MAB.RandomForest.w (BST)
62.8	78.5	58.1	mlp.C (NNET)	77.1	78.4	54.0	BG.RandomForest.w (BAG)
63.0	79.9	59.4	RandomCommittee.w (OEN)	78.5	76.5	53.7	Logistic.w (LMR)
63.4	78.7	58.4	Deconome.w (OEN)	78.7	76.6	50.5	ctreeBag.R (BAG)
63.6	76.9	56.0	mipWeightDecay.t (NNET)	79.0	76.8	53.5	BG.Logistic.w (BAG)
63.8	78.7	56.7	rda.t (DA)	79.1	77.4	53.0	lvq.t (NNET)
64.0	79.0	58.6	MAB.MLP.w (BST)	79.1	74.4	50.7	pls.t (PLSR)
64.1	79.9	56.9	MAB.RandomForest.w (BST)	79.8	76.9	54.7	hdda.R (DA)
65.0	79.0	56.8	knn.R (NN)	80.6	75.9	53.3	MCC.Cm (OEN)
65.2	77.9	56.2	multinom.t (LMR)	80.9	76.9	54.5	mrd.R (DA)
65.5	77.4	56.6	gevEarth.t (MARS)	81.4	76.7	55.2	C5.0Rules.t (RL)
65.5	77.8	55.7	glimnet.R (GLM)	81.6	78.3	55.8	lssvmRadial.t (SVM)
65.6	78.6	58.4	MAB.PART.w (BST)	81.7	75.6	50.9	JRip.t (RL)
66.0	78.5	56.5	CVRL.w (OM)	82.0	76.1	53.3	MAB.Logistic.w (BST)
66.4	79.2	58.9	treebag.t (BAG)	84.2	75.8	53.9	C5.0Tree.t (DT)
66.6	78.2	56.8	BG.PART.w (BAG)	84.6	75.7	50.8	BG.DecisionTable.w (BAG)
66.7	77.5	55.2	mda.t (DA)	84.9	76.5	53.4	NBTtree.w (DT)

Rank	Acc.	κ	Classifier
32.9	82.0	63.5	parRF.t (RF)
33.1	82.3	63.6	rf.t (RF)
36.8	81.8	62.2	svm_C (SVM)
38.0	81.2	60.1	svmPoly.t (SVM)
39.4	81.9	62.5	rforest._R (RF)

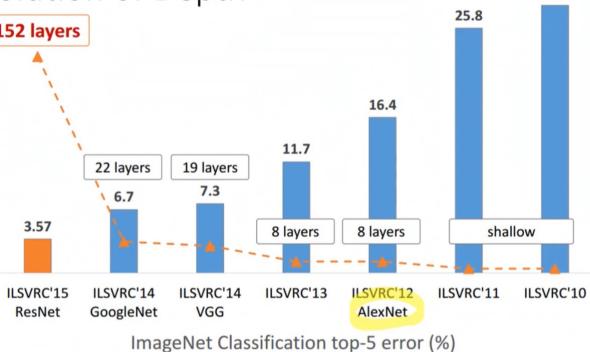
Algorithm \ Dataset 1 2 . .

26

179 3

- Large Scale Visual Recognition Challenge (~ILSVRC2015)

Revolution of Depth



✓ 2016



양상률의 시대



✓ 2017

Object detection (DET)^{final}

Task 1a: Object detection with provided training data

Ordered by number of categories won

Team name	Entry description	Number of object categories won	mean AP
CLImage	Ensemble of 6 models using provided data	109	0.66703
Hikvision	Ensemble A of 3 RPN and 6 FCNN models, mAP is 67 on val2	30	0.652704
Hikvision	[Ensemble B of 3 RPN and 5 FCNN models, mean AP is 69.9, median AP is 69.3 on val2]	18	0.652003

Object localization (LOC)^{final}

Task 2a: Classification+localization with provided training data

Ordered by localization error

Team name	Entry description	Localization error	Classification error
Trinps-Southern	Ensemble 3	0.077087	0.02991
Trinps-Southern	Ensemble 4	0.077429	0.02991
Trinps-Southern	Ensemble 2	0.077988	0.02991
Trinps-Southern	Ensemble 1	0.079008	0.03144

Object detection (DET)^{final}

Task 1a: Object detection with provided training data

Ordered by number of categories won

Team name	Entry description	Number of object categories won	mean AP
KISTIA	Submission1	80	0.731392
KISTIA	Submission2	80	0.729229
KISTIA	Submission3	80	0.7237172
DeepView(ETRI)	Ensemble_A	10	0.593084
NUS		0	0.588001
Choo, DPNs	Ensemble of DPN models	9	0.656932
(DET)		0	0.656932
KAISTNIA_ETRI	Ensemble Model1	0	0.660402
KAISTNIA_ETRI	Ensemble Model2	0	0.663299
KAISTNIA_ETRI	Ensemble Model1	0	0.660376
KAISTNIA_ETRI	Ensemble Model3	0	0.660376

Object localization (LOC)^{final}

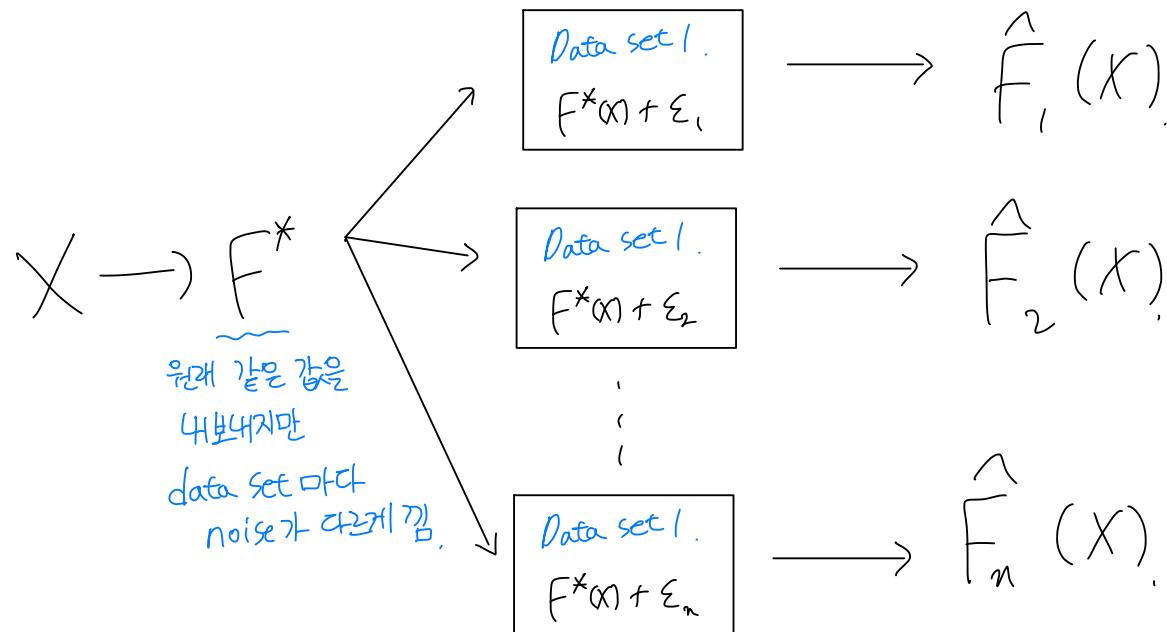
Task 2a: Classification+localization with provided training data

Ordered by localization error

Team name	Entry description	Localization error	Classification error
Alpha-DPN	2D3LOC: Dual Path Networks + Basic Ensemble	0.062283	0.03413
Trinps-Southern	Result_3	0.064681	0.02481
Trinps-Southern	Result_2	0.065235	0.02481
Trinps-Southern	Result_4	0.065261	0.02481
Trinps-Southern	Result_5	0.065302	0.02481
Trinps-Southern	Result_1	0.067698	0.02481

Bias - Variance Decomposition.

$$y = F^*(x) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$



$$\bar{F} = E[\hat{F}_0(x)]$$

$$\begin{aligned} \text{Err}(x_0) &= E[\gamma - \hat{F}(x) \mid X=x_0] \\ (\text{MSE}) &= E\left[F^*(x_0) + \varepsilon - \hat{F}(x_0)\right]^2 \quad (\gamma = F^*_0 + \varepsilon). \end{aligned}$$

$$= E\left[F^*(x_0) - \hat{F}(x_0)\right]^2 + \sigma^2$$

$$= E\left[\underbrace{F^*(x_0) - \bar{F}(x_0)}_A + \underbrace{\bar{F}(x_0) - \hat{F}(x_0)}_B\right]^2 + \sigma^2.$$

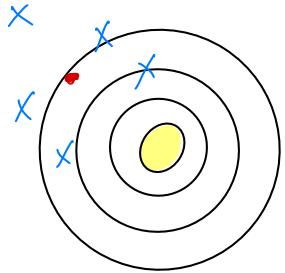
$$(2E[(F^*(x_0) - \bar{F}(x_0))(\bar{F}(x_0) - \hat{F}(x_0))]) = 0.$$

$$= E \left[F^*(x_0) - \bar{F}(x_0) + \bar{F}(x_0) - \hat{F}(x_0) \right]^2 + \sigma^2$$

$$= E \left[F^*(x_0) - \bar{F}(x_0) \right]^2 + E \left[\bar{F}(x_0) - \hat{F}(x_0) \right]^2 + \sigma^2$$

$$= \left[F^*(x_0) - \bar{F}(x_0) \right]^2 + E \left[\bar{F}(x_0) - \hat{F}(x_0) \right]^2 + \sigma^2$$

$$= \text{Bias}^2(\hat{F}(x_0)) + \text{Var}(\hat{F}(x_0)) + \sigma^2$$

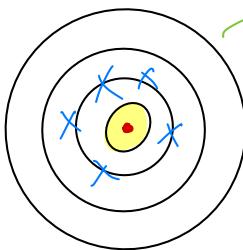


Bias

High

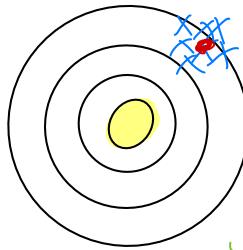
Variance.

High



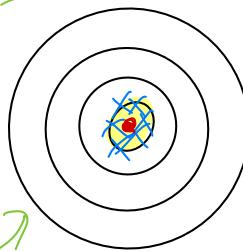
Low

High



High

Low



Boosting

Low

Low.

Bagging

Complexity ↑
model.

Purpose of Ensemble.

Reduce the error through constructing multiple learners.

Var ↓: Bagging, Random Forests, (\neq Bagging)

Bias ↓: Boosting

ex. 카드 대회.

A팀

76점
+ 76점 29명.

↓
diversity.

B팀

28
+ 선.박사 29명

C팀

28점 30명.

↓ good performance individually.

↳ Key fact of the ensemble Construction.

Average Error vs Ensemble Error.

$$E_{\text{Avg}} = \frac{1}{M} \sum_{m=1}^M E[\varepsilon_m(x)^2]$$

양상을의 측정: 개별 모델들의 평균.

$$E_{\text{Ensemble}} = E \left[\left\{ \frac{1}{M} \sum_{m=1}^M y_m(x) - f(x) \right\}^2 \right]$$

$$= E \left[\left\{ \frac{1}{M} \cdot \left(\sum_{m=1}^M y_m(x) - M \cdot f(x) \right) \right\}^2 \right]$$

$$= E \left[\left\{ \frac{1}{M} \cdot \sum_{m=1}^M (y_m(x) - f(x)) \right\}^2 \right]$$

$$= E \left[\left\{ \frac{1}{M} \cdot \sum_{m=1}^M \varepsilon_m(x) \right\}^2 \right]$$

$$= E \left[\left\{ \frac{1}{M} \cdot \sum_{m=1}^M \varepsilon_m(x) \right\}^2 \right]$$

Assume that zero mean & uncorrelated.

$$E[\varepsilon_m(x)] = 0 \quad E[\varepsilon_m(x) \cdot \varepsilon_\ell(x)] = 0 \quad (m \neq \ell)$$

$$\Rightarrow E_{\text{ensemble}} = \frac{1}{M} \cdot E_{\text{Avg.}} \quad \left(\text{But in reality errors are correlated.} \right)$$

In reality, by Cauchy inequality

$$(a^2+b^2)(x^2+y^2) \geq (ax+by)^2.$$

$$\left\{ 1 \cdot \varepsilon_1(x) + 1 \cdot \varepsilon_2(x) + \cdots + 1 \cdot \varepsilon_M(x) \right\}^2 \leq \underbrace{(1+1+\cdots+1)}_{M} (\varepsilon_1(x)^2 + \varepsilon_2(x)^2 + \cdots + \varepsilon_M(x)^2)$$
$$\hookrightarrow \left(\sum_{m=1}^M \varepsilon_m(x) \right)^2 = M \cdot \sum_{m=1}^M \varepsilon_m(x)^2$$

$$\Rightarrow \left\{ \frac{1}{M} \sum_{m=1}^M \varepsilon_m(x) \right\}^2 \leq \frac{1}{M} \cdot \sum_{m=1}^M \varepsilon_m(x)^2.$$

$$\Rightarrow E_{\text{ensemble}} \leq E_{\text{Avg}}$$

← 양상률의 핵심은
모델들의 "diversity"

Bagging.

diversity 주는 방법.

{ Data, $\xrightarrow{\text{red}}$ Bagging.
Model, $\xrightarrow{\text{blue}}$ Random Forests.

i) k-fold.

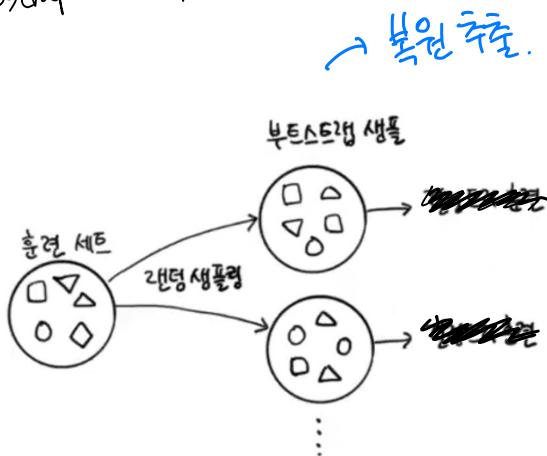
- Sampling with replacement.



iii) Bagging (Bootstrap Aggregating).

: Choosing N uniformly random with replacement.

Bootstrap Sampling

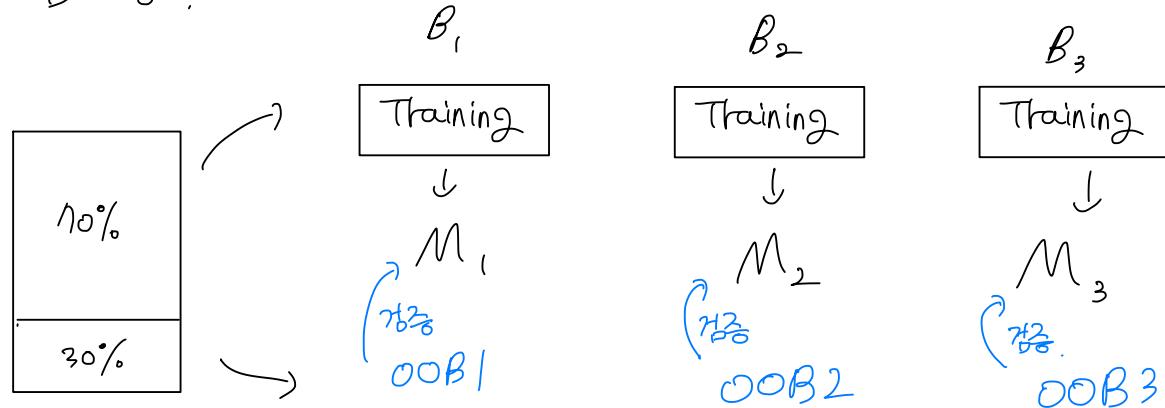


이론적으로 한 번도 선택되지 않은 확률.

$$\left(1 - \frac{1}{N}\right)^N \Rightarrow \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^N = e^{-1} = 0.368.$$



OOB 평가.



개별적인 샘플들이

{ 어떤 경우에는 학습데이터
어떤 경우에는 검증데이터.

Result Aggregating

Training Accuracy	Ensemble population	$P(y=1)$ for a test instance	Predicted class label
0.80	Model 1	0.90	1
0.75	Model 2	0.92	1
0.88	Model 3	0.87	1
0.91	Model 4	0.34	0
0.77	Model 5	0.41	0
0.65	Model 6	0.84	1
0.95	Model 7	0.14	0
0.82	Model 8	0.32	0
0.78	Model 9	0.98	1
0.83	Model 10	0.57	1

Seung Bum Kim Courtesy of Pilsung Kang's lecture slides

각 모델의 비율.

1 - 6개 \rightarrow 1 예측.
0 - 4개

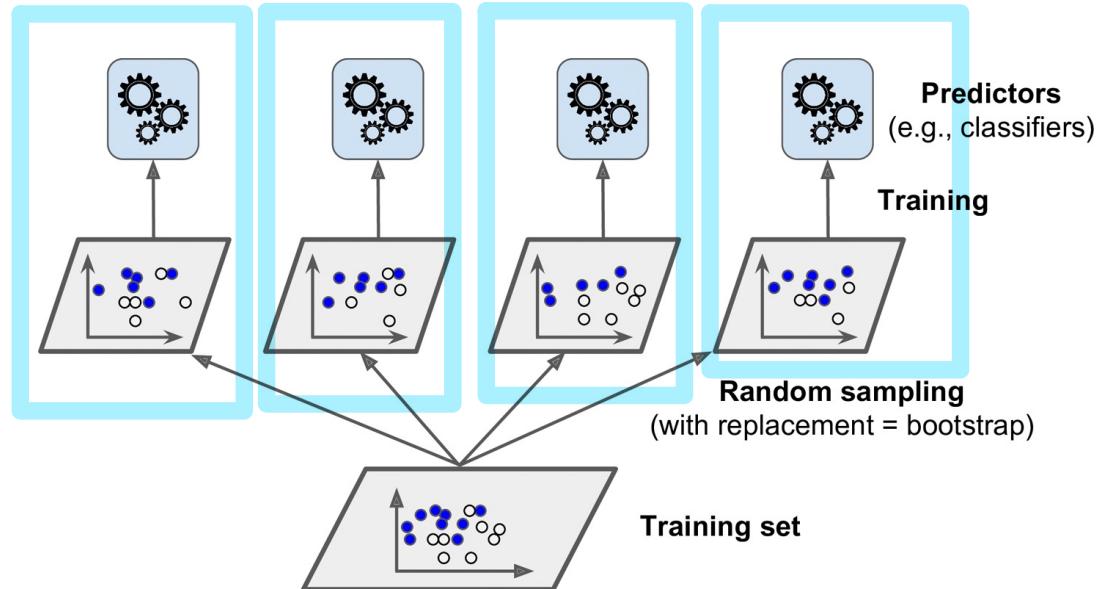
정확도 가중치 적용

$1 \times 1 + 1 \times 1 + \dots + 10$.

학습의 평균 적용

⋮

동시에 훈련 가능 !

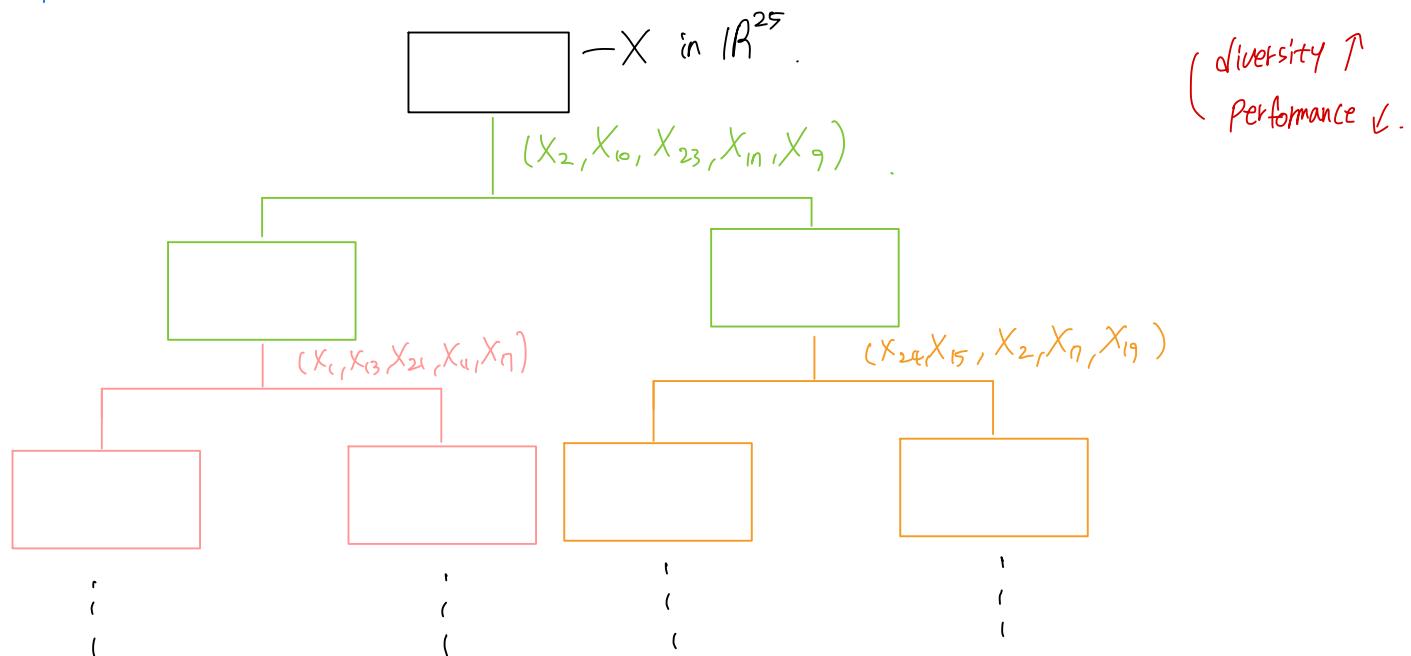


Random Forests.

Bagging + Decision Tree (with random subspace).

↖
Data를 이용한
diversity.

↗
Model을 이용한
diversity.



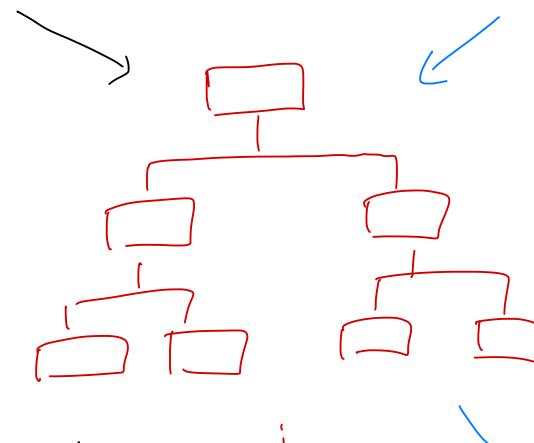
Feature Importance.

Original oob Data.

	X_i	
	1	
	2	
	3	
	4	
	5	

(after permutation X_i) oob Data.

	X_i	
	2	
	5	
	3	
	1	
	4	



oob Error of the
original Data e_i .

oob Error of the
permuted Data P_i .

변수의 중요도가 높다면?

- Permutation 전후 oob Error 차이가 커야 함.
- 그 차이의 평차가 적어야 함.

Let $d_i^m = p_i^m - e_i^m$ (m 번째 tree, X_i 에 대한 oob error 차이).

$$\bar{d}_i = \frac{1}{M} \sum_{m=1}^M d_i^m, \quad s_i^2 = \frac{1}{M-1} \sum_{m=1}^M (d_i^m - \bar{d}_i)^2$$

$$\Rightarrow \text{Feature importance of } X_n = \frac{\bar{d}_i}{s_i^2}$$

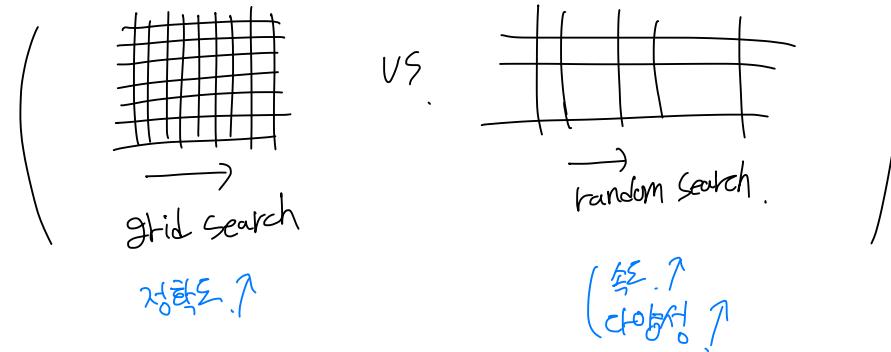
Extra Thees.

Model diversity ↑

random subspace

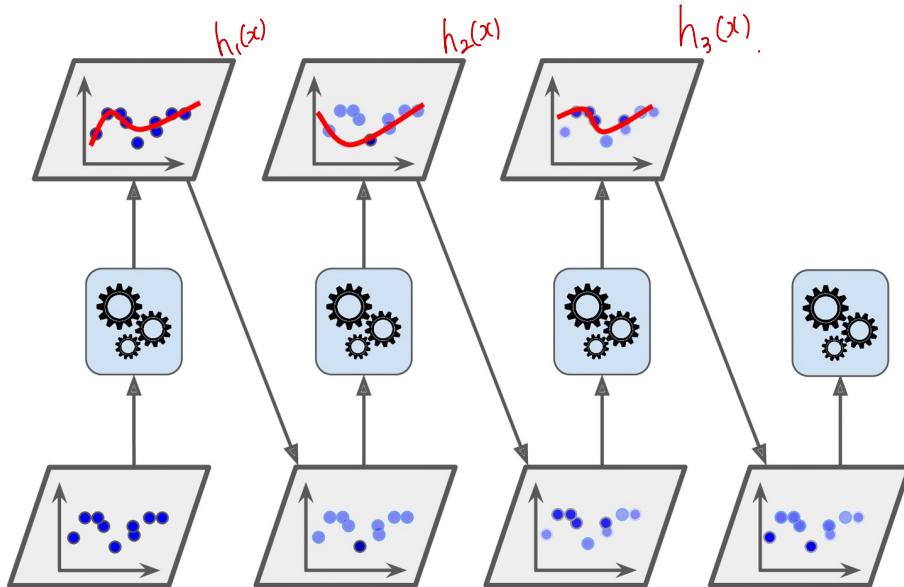
+

random search.

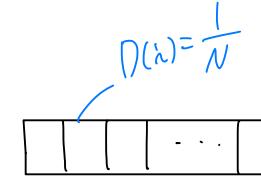


+ bootstrap sampling X

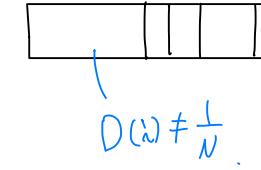
Boosting.



Bagging



Boosting



$$H(x) = \text{Sign} (\lambda_1 h_1(x) + \lambda_2 \cdot h_2(x) + \dots + \lambda_k \cdot h_k(x)).$$

Ada Boosting.

Training Data $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, where $y_n = (-1, 1)$.

Define uniform distribution $\tilde{D}_t(i)$ over elements of S .
첫번째 Data set에
example i 가 선택될 확률.

i) Train h_t using distribution D_t .

ii). calculate $\varepsilon_t = P_{D_t}(h_t(x) \neq y)$ 오른쪽을. ($\varepsilon_t > 0.5$: stop).

$$\text{set } \lambda_t = \frac{1}{2} \cdot \ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

$$\lambda_t \begin{cases} \varepsilon_t = 0.5 \rightarrow \frac{1}{2} \ln \left(\frac{0.5}{0.5} \right) = 0 \\ \varepsilon_t = 0 \rightarrow \frac{1}{2} \ln \left(\frac{1}{0} \right) = \infty \end{cases}$$

iii) Update $D_{t+1}(i) = \frac{D_t(i) \cdot \exp(-\beta_t \cdot Y_t \cdot h_t(x_i))}{Z_t}$

↑ or ↓

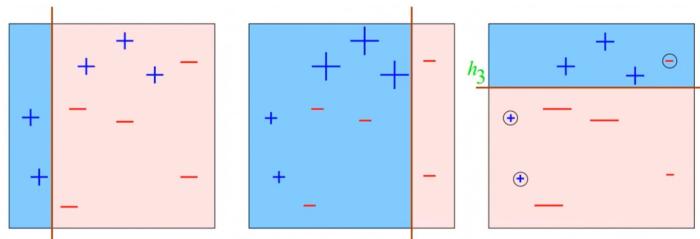
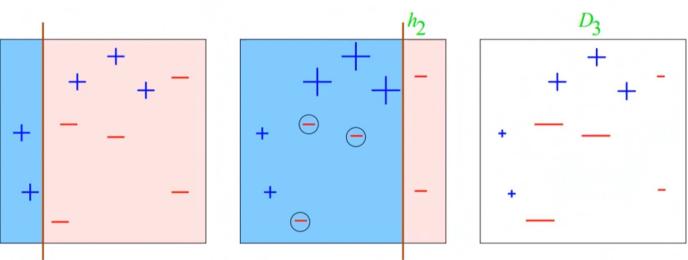
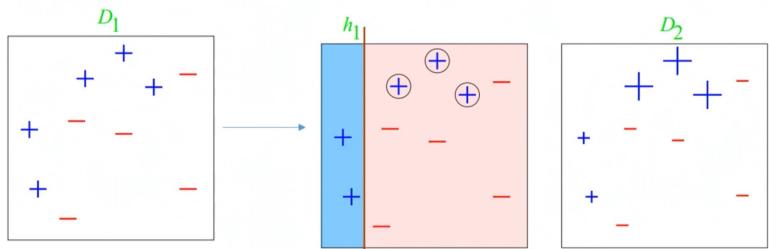
$$\frac{1}{Z_t} \times D_t(i) \times \exp(-\beta_t \cdot Y_t \cdot h_t(x_i)).$$

normalize factor
t 시점 선택 확률이 기준.

h_t 가 정답하면
Sample rate 의
병동성 ↑

$$\Rightarrow \begin{cases} h_t \text{가 } i \text{의 정답} \\ D_{t+1}(i) \downarrow (e^{-\beta_t}) \end{cases} .$$

$$\begin{cases} h_t \text{가 } i \text{의 정답} \\ D_{t+1}(i) \uparrow (e^{\beta_t}) \end{cases} .$$

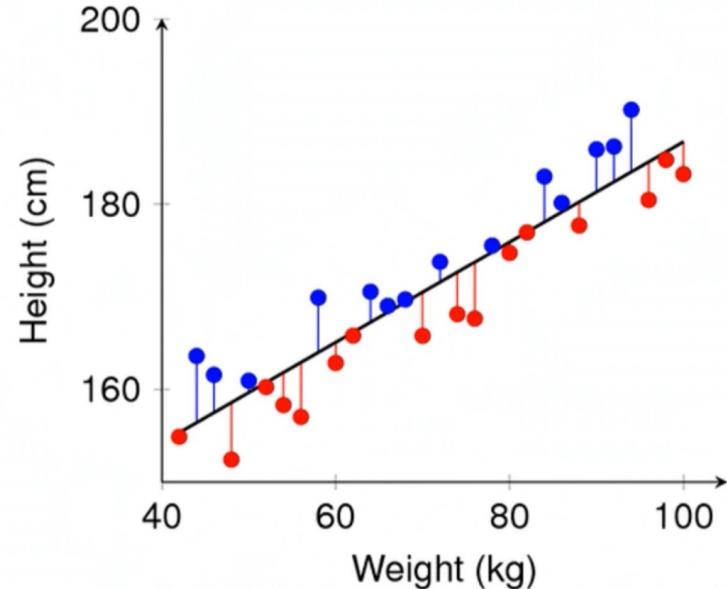


$H_{\text{final}} = \text{sign} \left(0.42 \begin{array}{|c|c|} \hline \text{blue} & \text{pink} \\ \hline \end{array} + 0.65 \begin{array}{|c|c|} \hline \text{blue} & \text{pink} \\ \hline \end{array} + 0.92 \begin{array}{|c|c|} \hline \text{blue} & \text{pink} \\ \hline \end{array} \right)$

=

https://www.cse.buffalo.edu/~jcorso/CSE455/fileslecture_booleans.pdf 9

Gradient Boosting.



$$y = f_1(x).$$

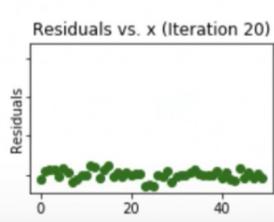
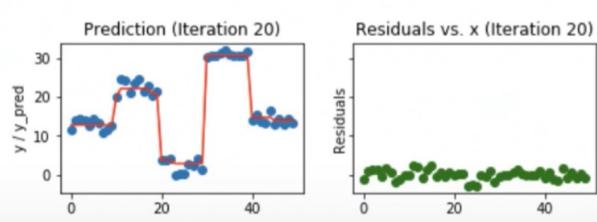
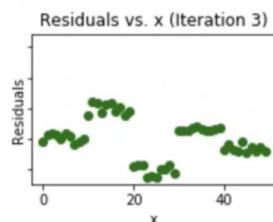
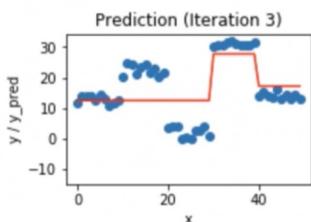
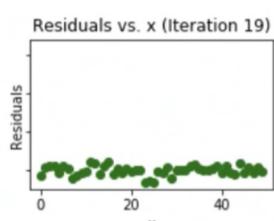
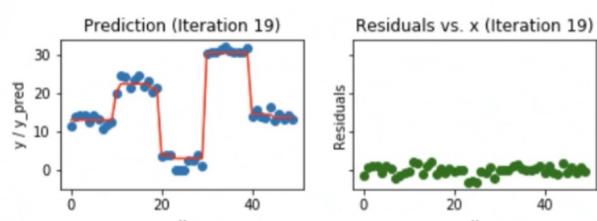
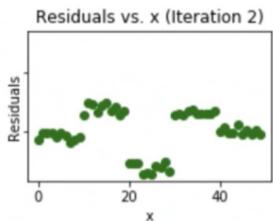
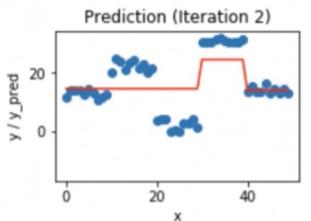
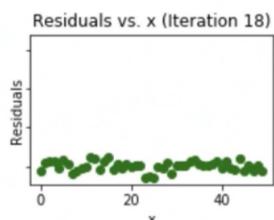
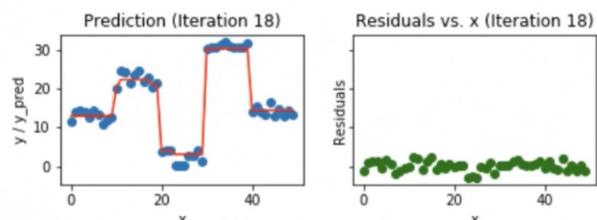
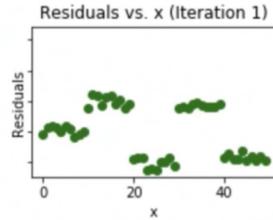
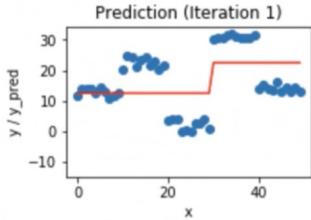
$$\underbrace{y - f_1(x)}_{\text{residual}} = f_2(x).$$

$$\Rightarrow y = f_1(x) + f_2(x).$$

(Ada Boosting 2t 쓰기) Data Sampling X)

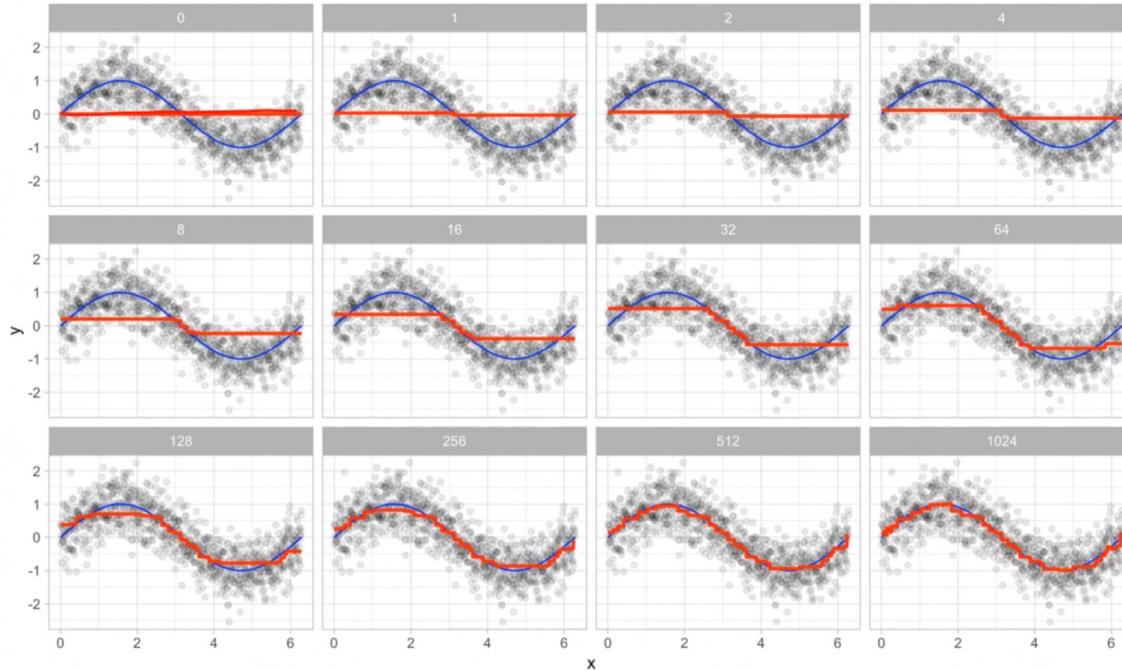
Gradient Boosting Machine: GBM

- **GBM Regression Example I**



Gradient Boosting Machine: GBM

- GBM Regression Example 3



<https://docs.paperspace.com/machine-learning/wiki/gradient-boosting>

Over fitting problem in GBM.

$$f_1(x) = y + \varepsilon$$

$$f_2(x) = (y - f_1(x)) + \varepsilon$$

$$f_3(x) = (y - f_1(x) - f_2(x)) + \varepsilon$$

:



ε 비중 \nearrow (ε 까지 학습).

=> Regularization.

Subsampling



10%

10%

10%

Shrinkage

$$y = f_1(x) + 0.9 f_2(x) + 0.9^2 f_3(x) + \dots$$

Early Stopping.



- Variable Importance in Tree-based Gradient Boosting

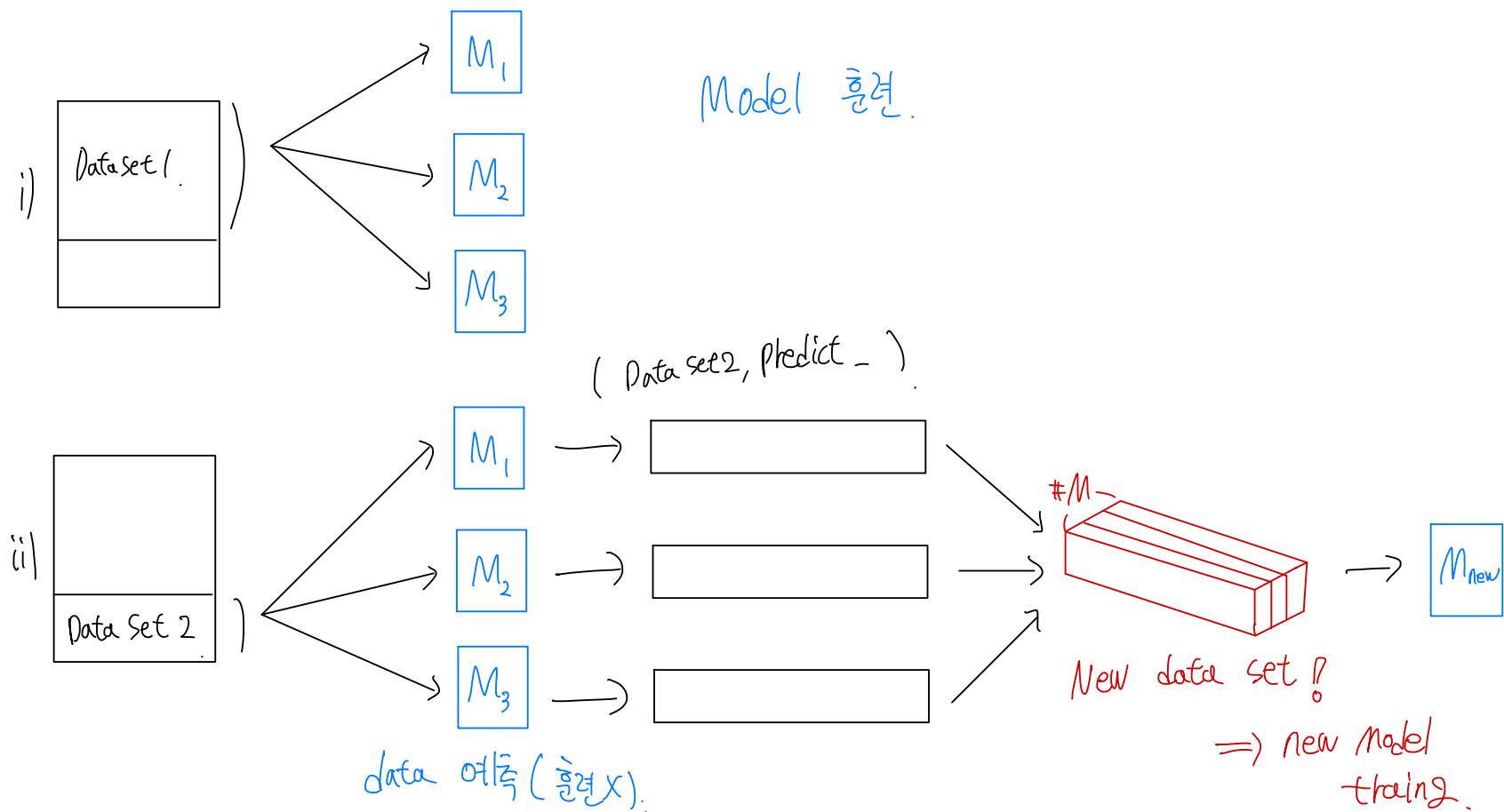
- ✓ $Influence_j(T)$: importance of the variable j in a single tree T .
- ✓ Assume that there are L terminal nodes $\rightarrow L - 1$ splits.

$$Influence_j(T) = \sum_{i=1}^{L-1} (IG_i \times \mathbf{1}(S_i = j))$$

- ✓ Variable importance of Gradient boosting

$$Influence_j = \frac{1}{M} \sum_{k=1}^M Influence_j(T_k)$$

Stacking (양상률의 양상률)



Q & A.

유튜브 강의 & 사진 참고

Pilsung Kang
School of Industrial Management Engineering
Korea University