

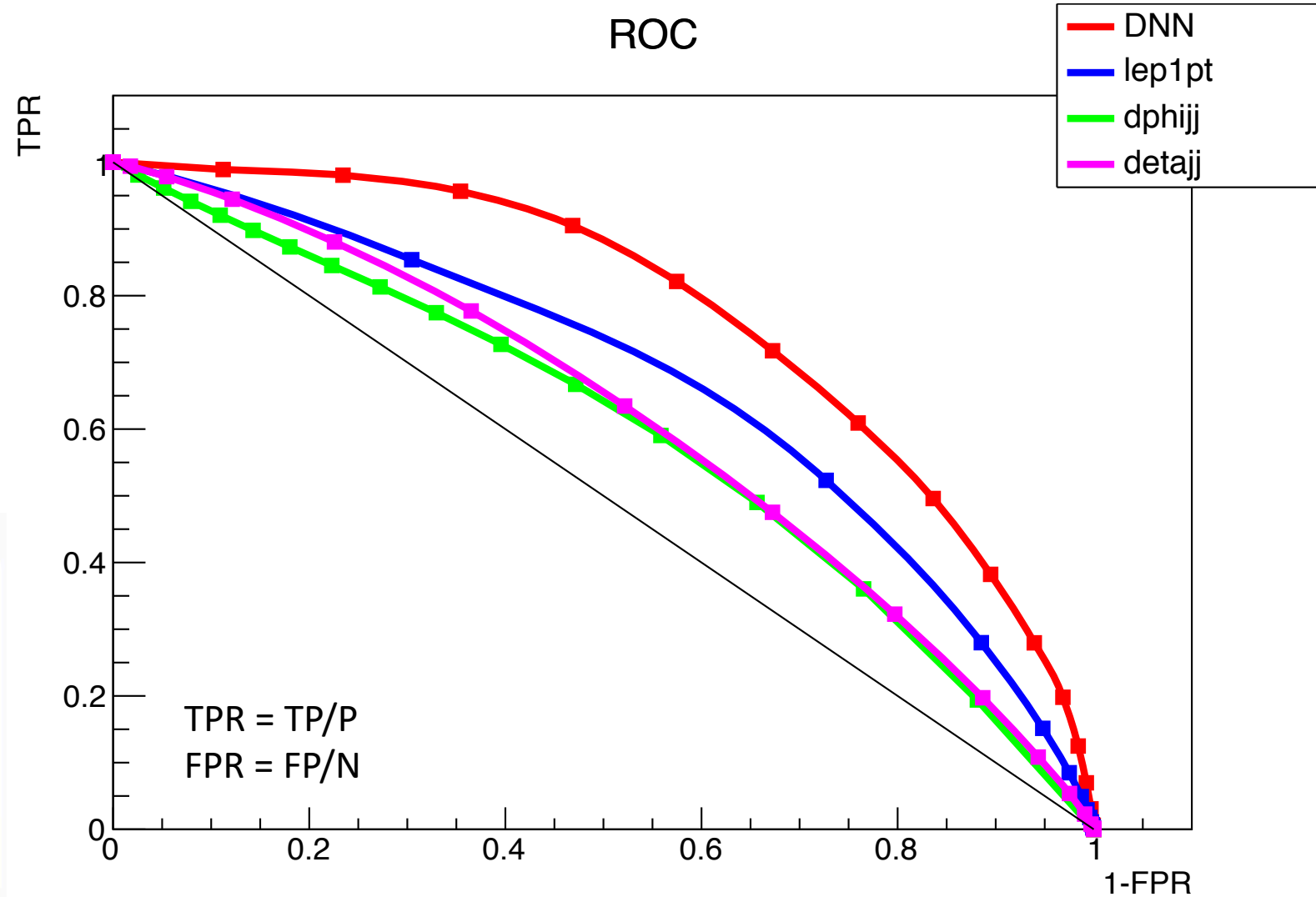
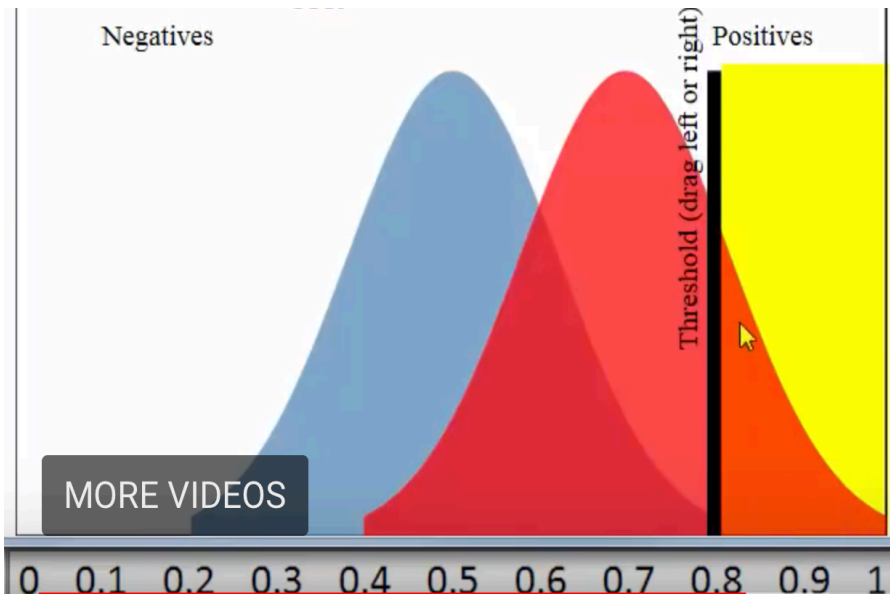
머신러닝 & 통계 기초

2019/11/29

이준호

ROC curve

Classifier	TRUE CLASS	
	p (positive)	n (negative)
Predicted class		
Y	True Positives	False Positives
N	False Negatives	True Negatives
Total	P	N



Multivariate analysis

Dataset

Train Dataset

Train Dataset

Obtain Trained DNN model

→ LL and TTTL classifier

Validation Dataset

→ Validate Training

Test Dataset

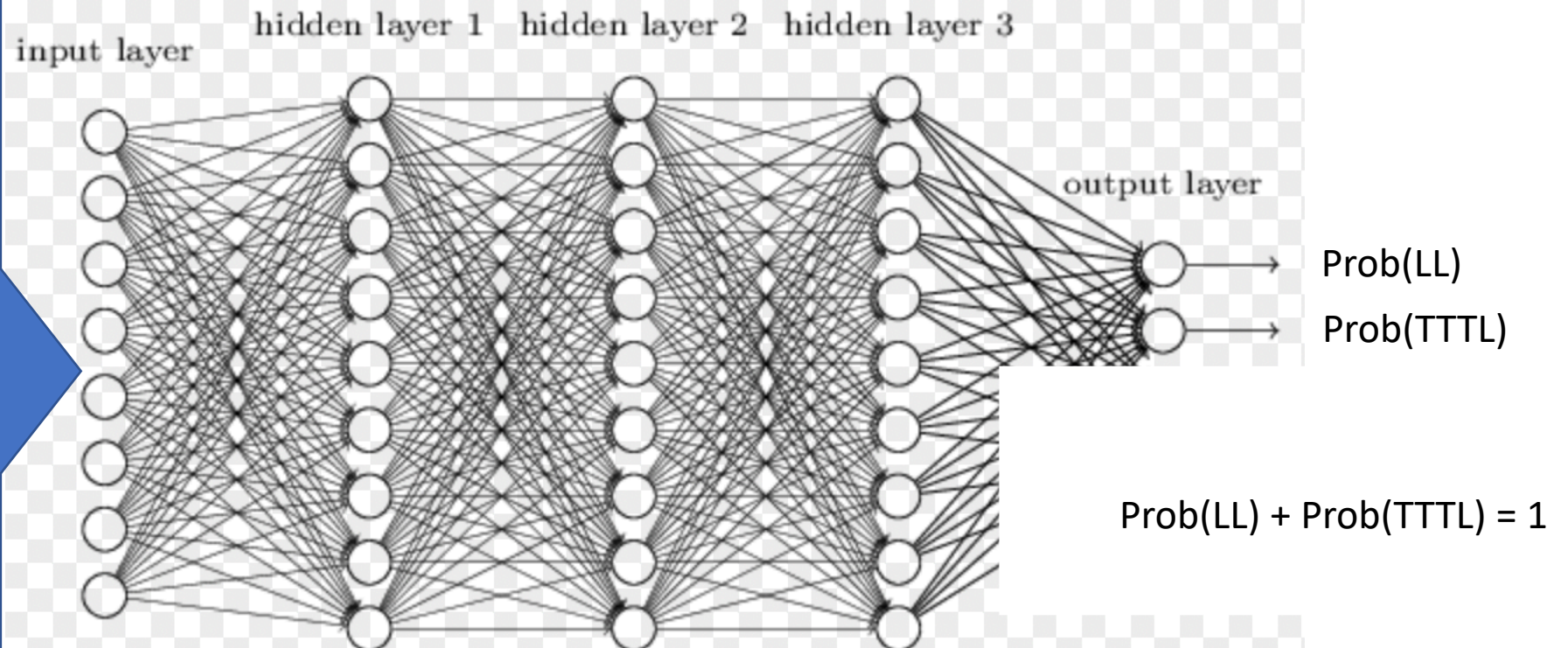
→ Test performance of the model

Low level:

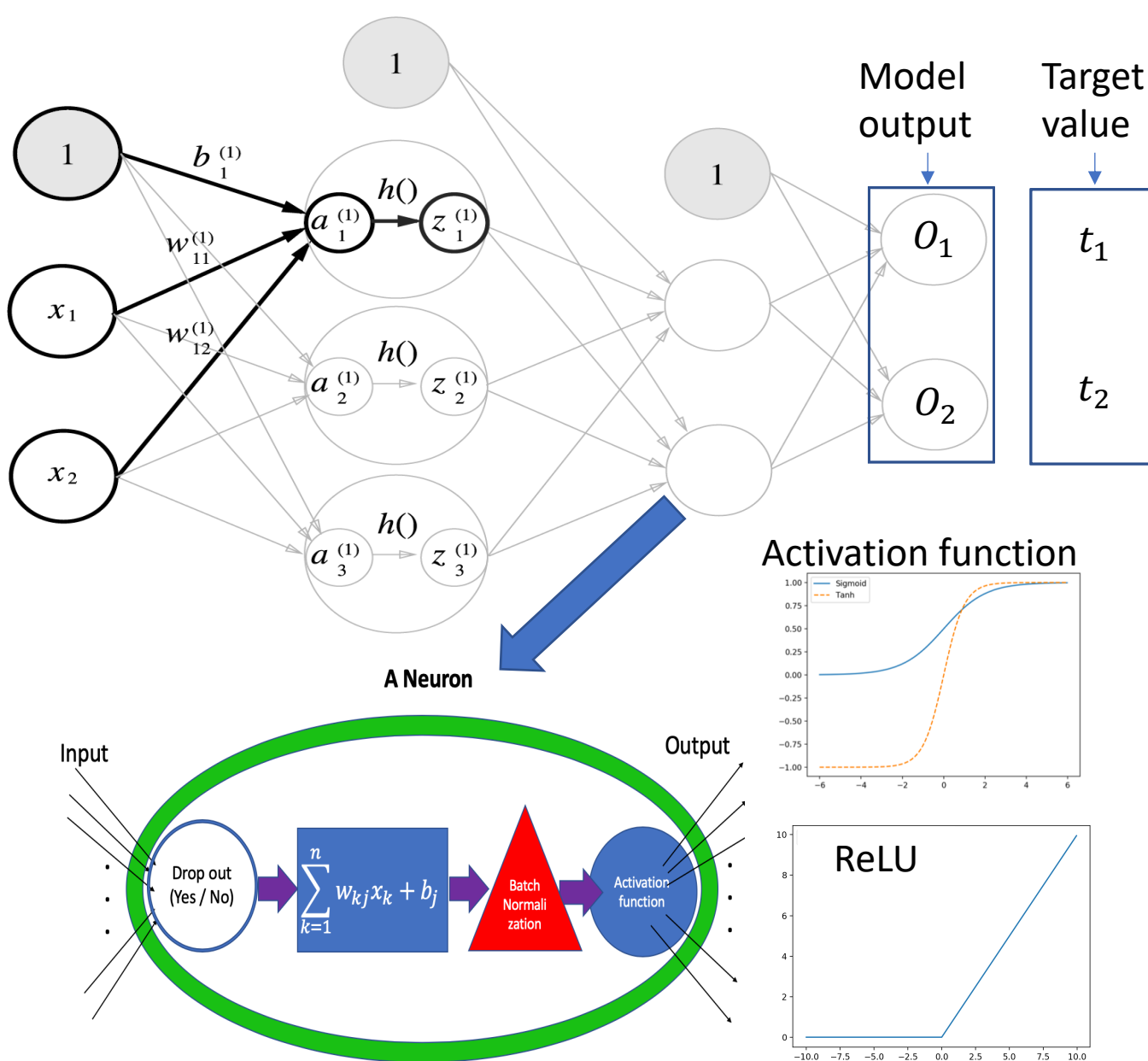
lep1pt,
lep1eta,
lep1phi,
lep2pt,
lep2eta,
lep2phi,
jet1pt,
jet1eta,
jet1phi.,
jet1M
jet2pt,
jet2eta,
Jet2phi,
jet2M,
MET,
METphi

High level:

dr_ll_jj,
dphijj,
detajj,
Zeppen_lep1,
Zeppen_lep2



Neural network : 신경망



- Given Dataset : $[[x_1, x_2], [t_1, t_2]]$
- Initialize weights of a NN model

1. Forward propagation
 - Activation function
 - Until reach to output layer
 - Calculates output of the model
2. Make loss function given output and Target value

$$\text{MSE} : L(\theta) = \frac{1}{n} \sum_{j=1}^J (t_j - O(x_j, \theta))^2$$

$$\text{Cross entropy} : L(\theta) = -\frac{1}{n} \sum_{j=1}^J [t_j \log(O(x_j, \theta)) + (1 - t_j) \log(1 - O(x_j, \theta))]$$

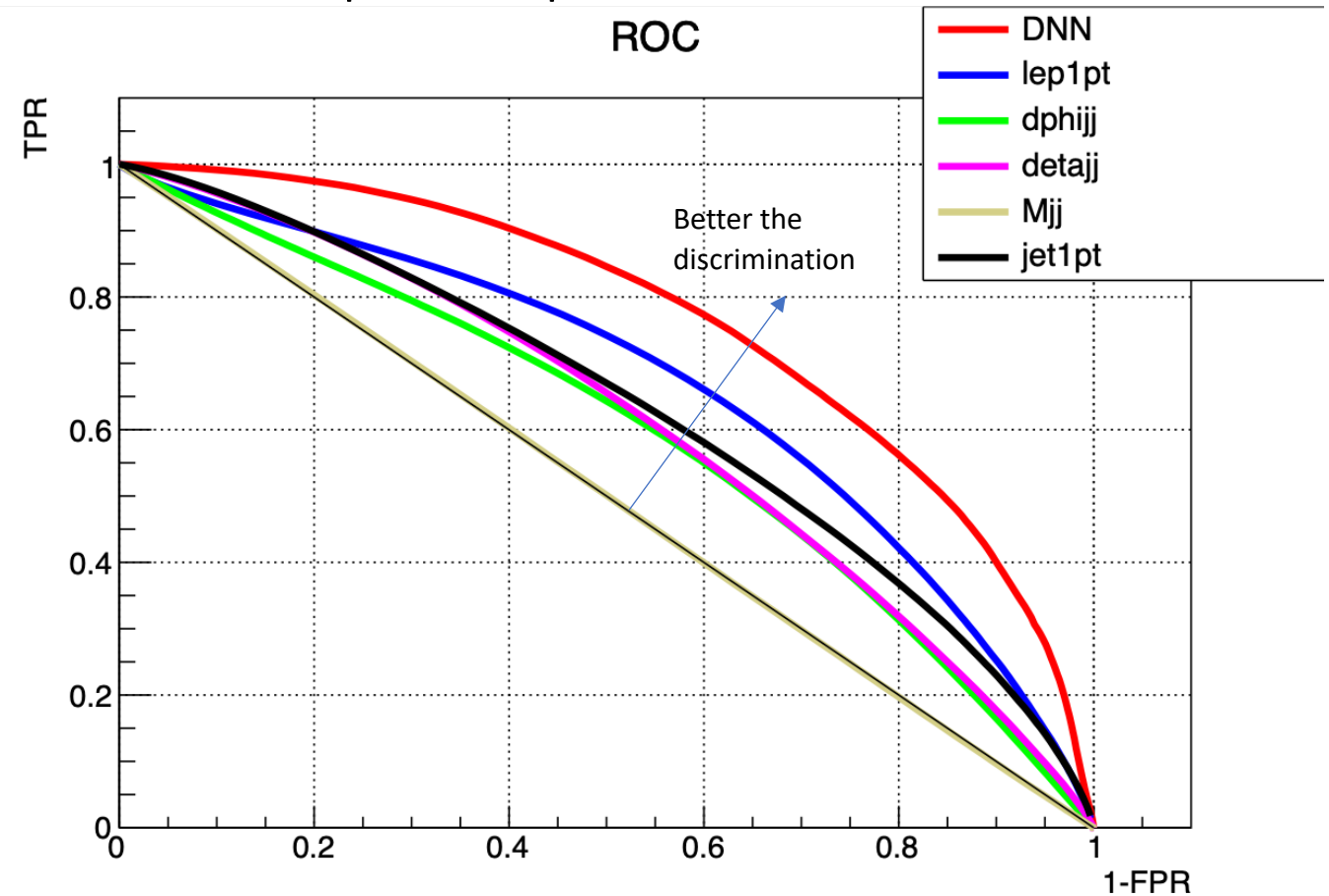
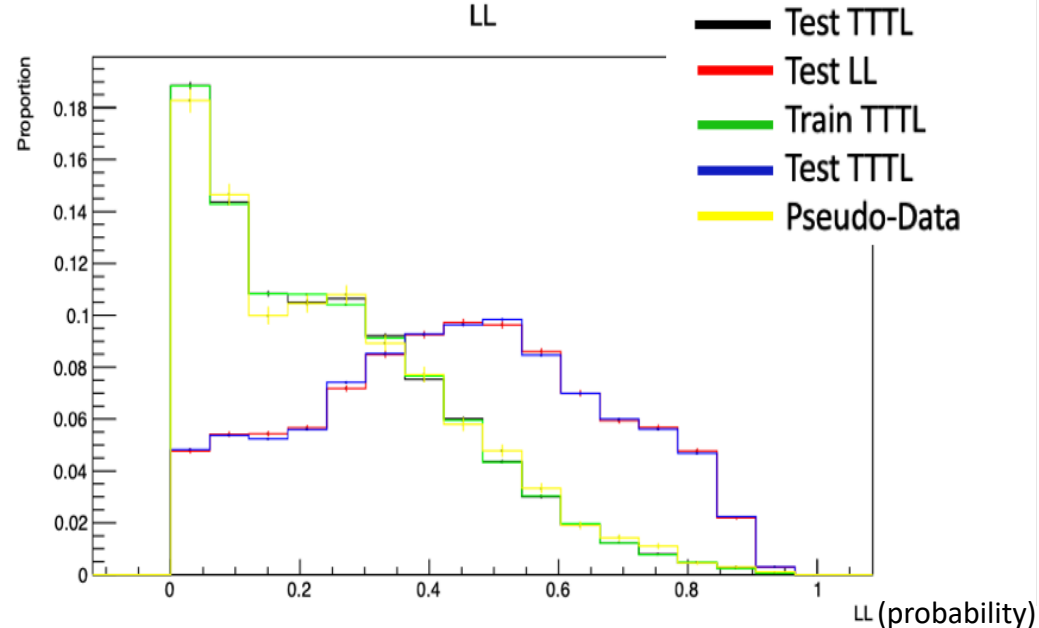
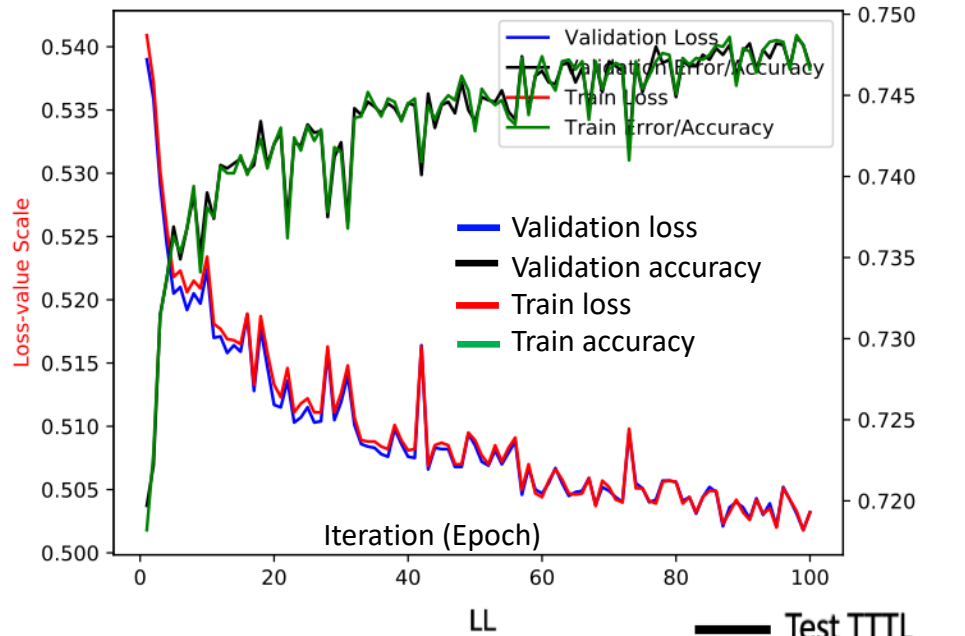
3. Back propagation : adjust weights and biases by gradients given loss function
 - Optimizer (SGD, Momentum, Adam)
 - Learning-rate (η)

$$w_{ij}^{\text{updated}} = w_{ij} + \Delta w_{ij} = w_{ij} - \eta \frac{\partial L_j}{\partial w_{ij}}$$

4. Minimize loss function by iterating 1~3.

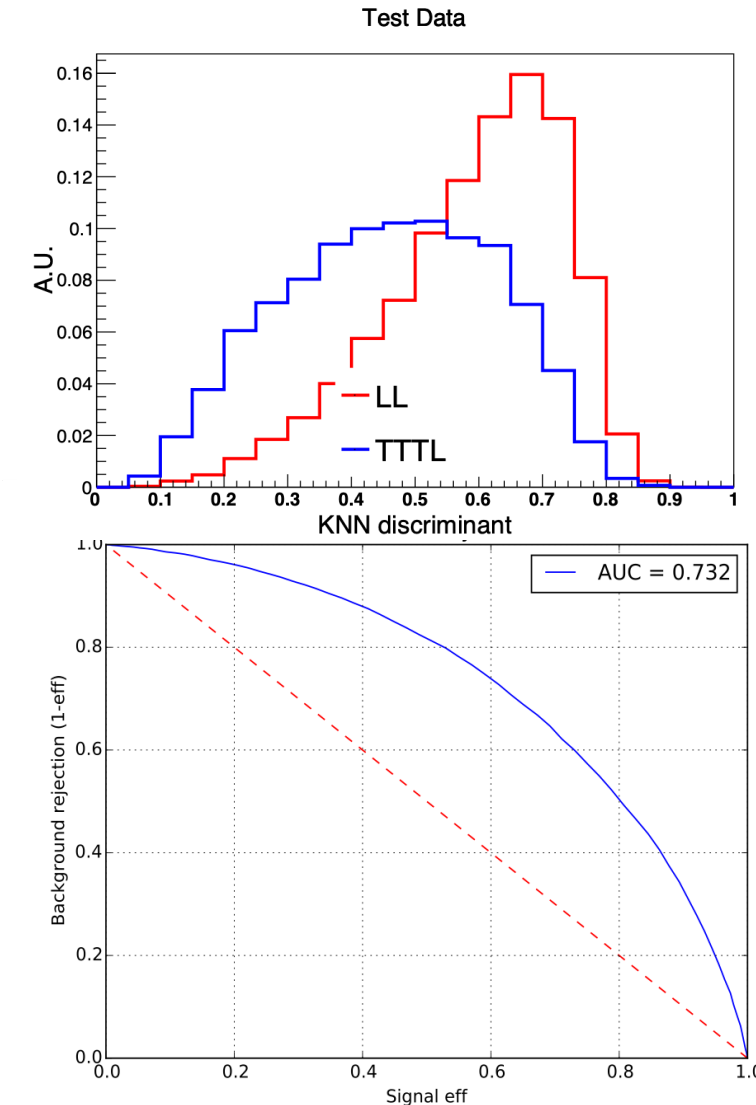
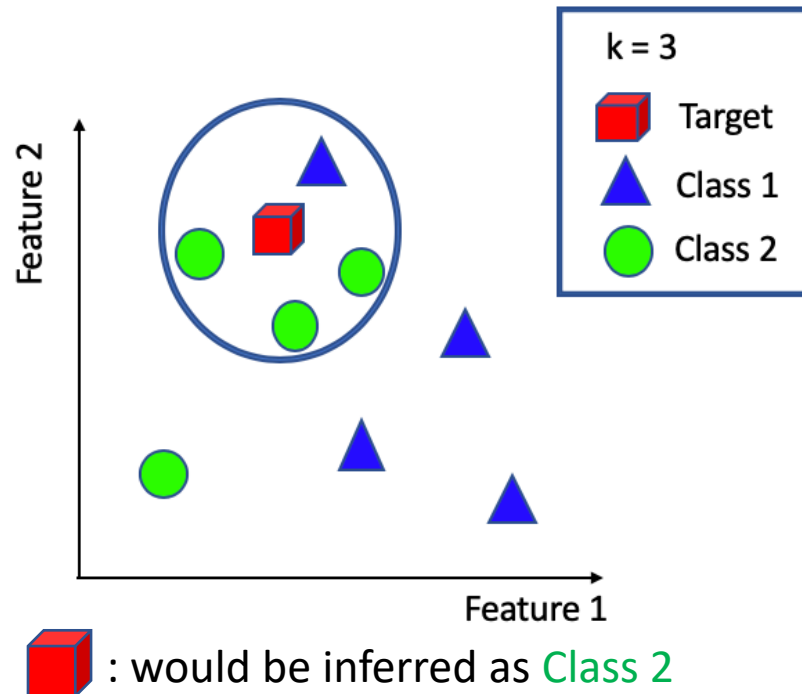
Multivariate classification (1) : Dense neural-network

- 10 hidden-layered dense-NN
 - The best among 1~15 layered dense-NN
 - 150 neurons for each layer, ReLU activation function
 - 50% dropout and 50 epoch for early stopping
- Discrimination power improvement found from dense-NN

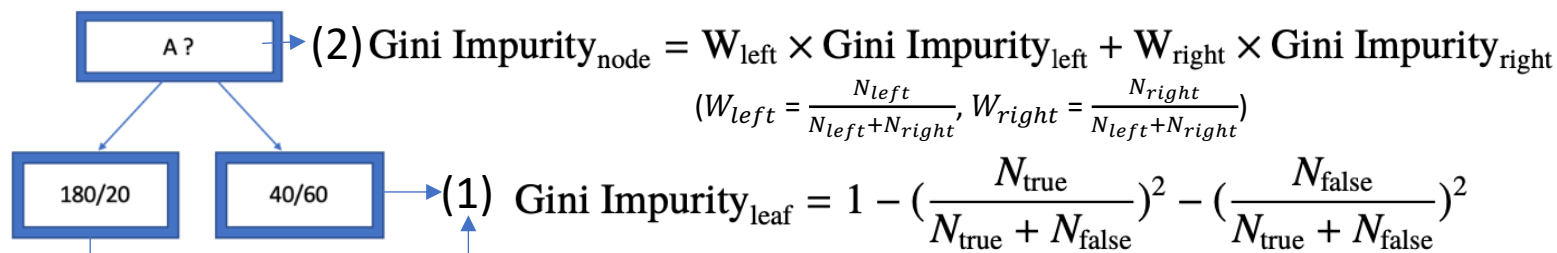
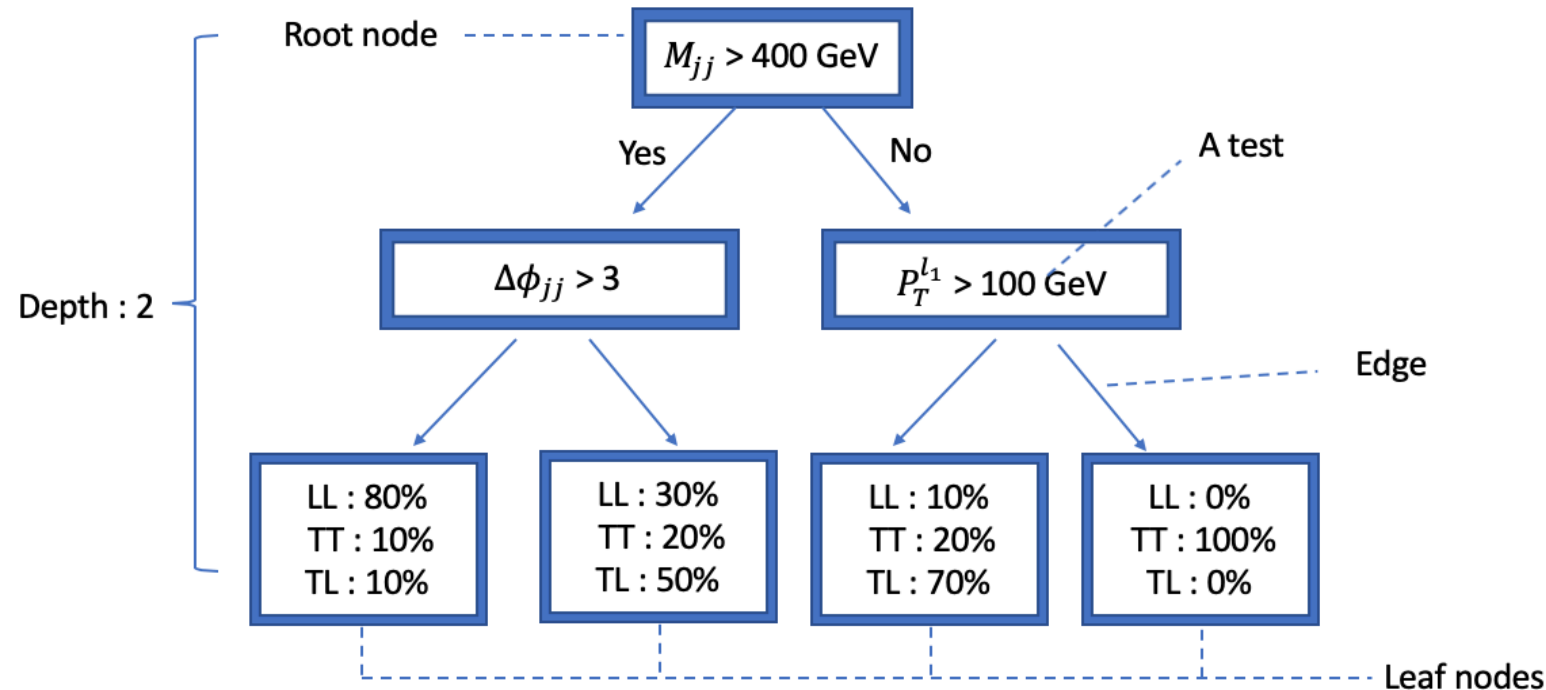


k-nearest neighborhood : 분류

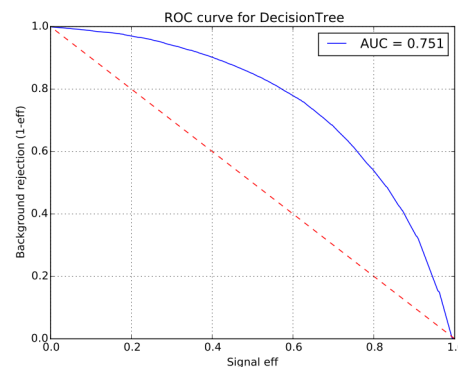
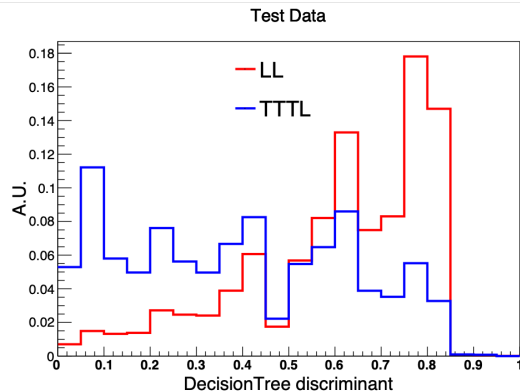
- Data itself constructs a k-NN model
- The model outputs probability on a certain data point to be classified as certain class
 - Feature space
 - With hyper-parameter 'k'
 - Euclidean distance (optional)



Decision tree : 분류

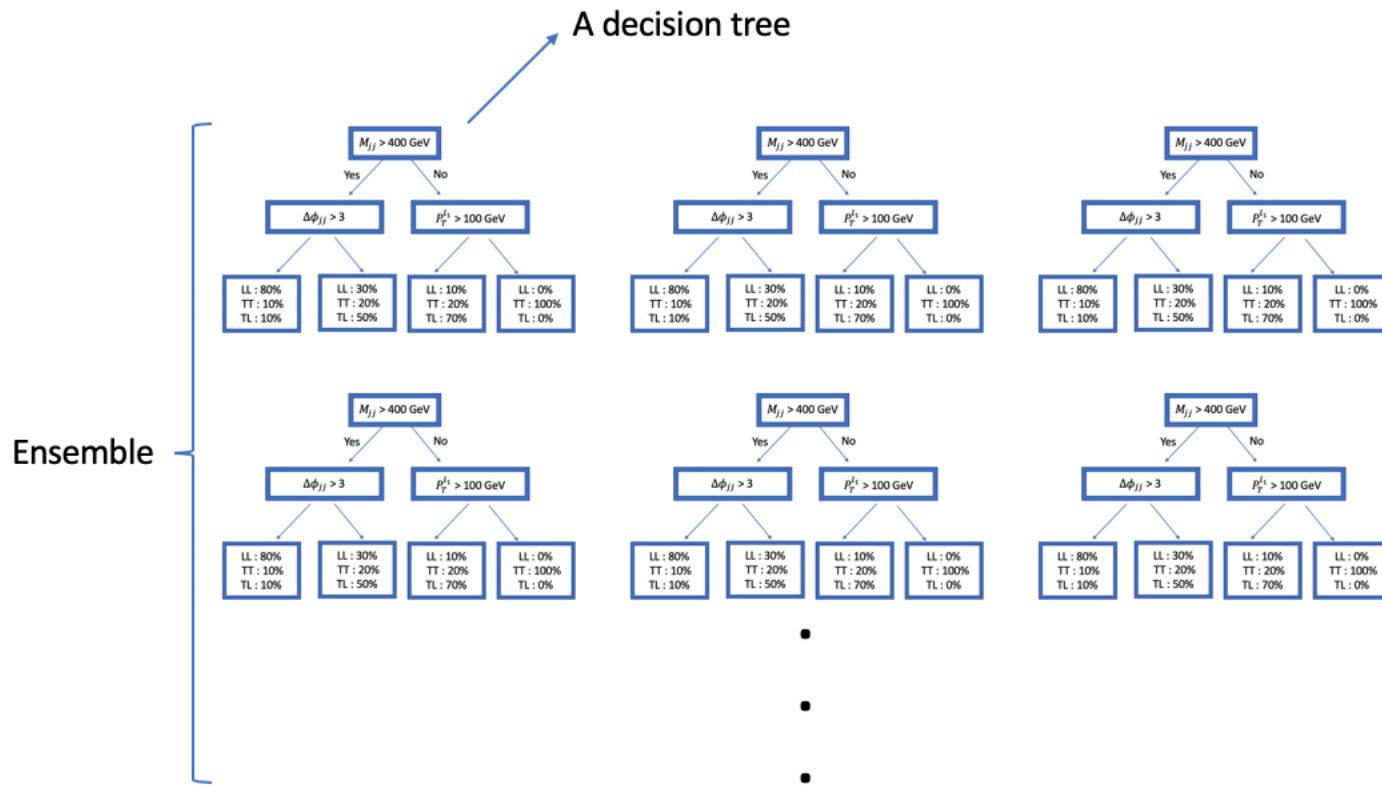


Only slightly better
than k-NN



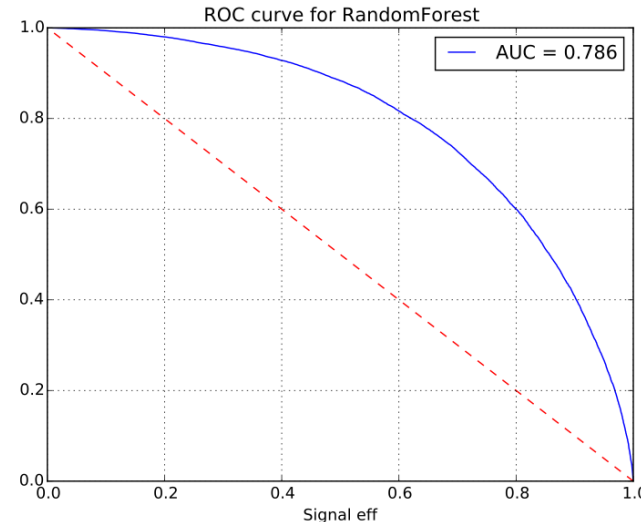
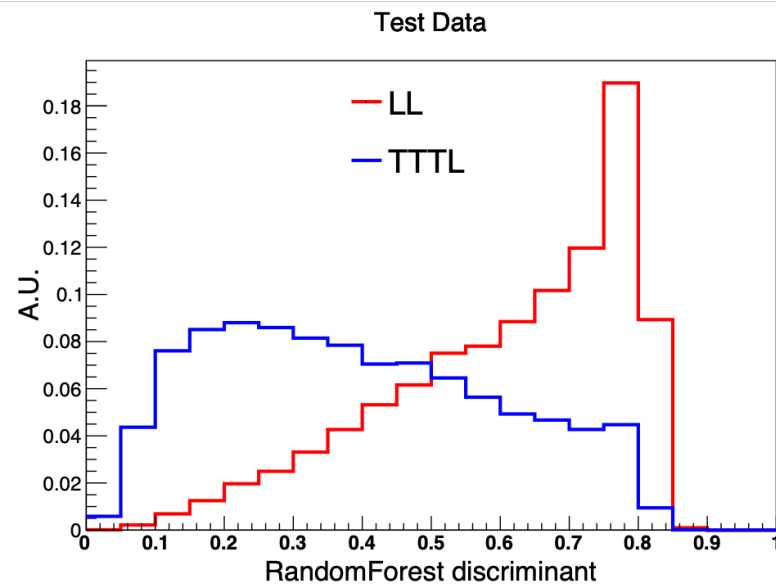
- Given Dataset : $[[x_1, x_2, \dots, x_N], [t_1, t_2]]$
- Find the best test which will be used for a node
 - Gini impurity leaf
 - Gini impurity node
 - If node itself has lower impurity score than any of its leaves' score, than the node **become a leaf node**
 - Stop separating if depth of tree reaches $N_{\text{max_depth}}$, which is a hyper-parameter for avoiding over-fitting.
 - Separate** if not satisfies '2' and '3'.
 - Iterate 1~5.

Random forest



- Given Dataset : $[[x_1, x_2, \dots, x_N], [t_1, t_2]]$

- Create a **bootstrapped** dataset.
 - Some data point might not included
- Create a **decision tree**, while only taking n features, which is given hyper-parameter, when making a node.
 - Usually, $n = \sqrt{N_{feature}}$
- Iterate 1~2 until number of trees reaches given hyper-parameter N_{tree} .



Better than k-NN, Logistic Regression, and Decision tree.

Linear regression : 선형회귀

Odds & Odds ratio

happened



not happened

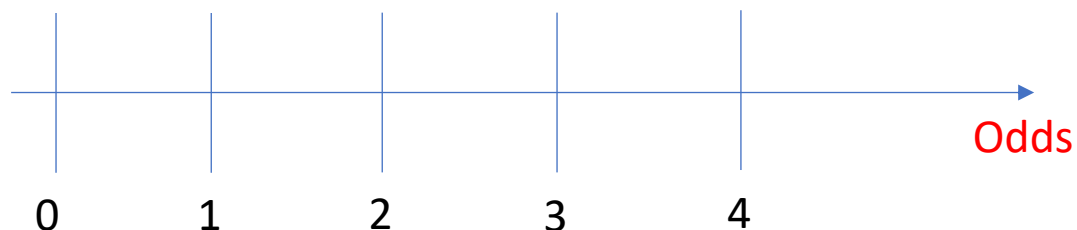


Odds = $2/3 \rightarrow$ '1' a boundary of a favor

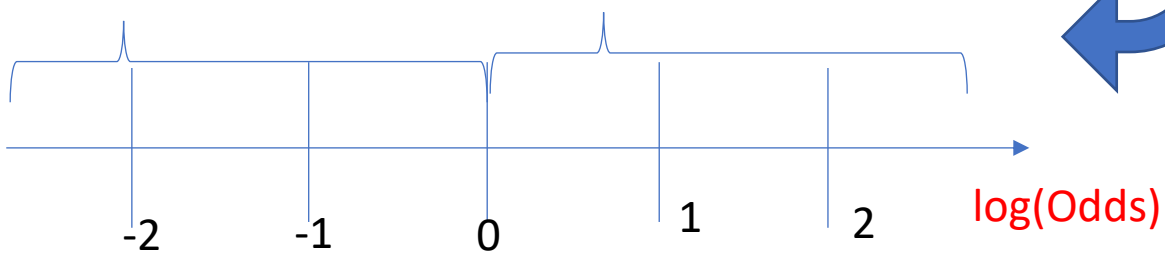
Probability (p) = $2/5$

Odds = $p / (1-p)$

($\log(\text{Odds}) = \log(p/(1-p)) \rightarrow$ logit function)



Taking $\log(\text{Odds})$
for **symmetric**



Odds ratio & $\log(\text{odds ratio})$

		Has Cancer		
		Yes	No	
Has the mutated gene	Yes	23	117	$\rightarrow 23/117$
	No	6	210	$\rightarrow 6/210$

Use odds ratio to determine relationship between mutated gene and cancer

- Odds ratio = $\frac{23/117}{6/210} = 6.88$
- $\log(\text{Odds ratio}) = 1.93$

Higher 'odds ratio' : better predictor (mutated gene)

Observed		Has Cancer	
		Yes	No
Has the mutated gene	Yes	23	117
	No	6	210

Odds (Wald test)

- Test if the odds ratio is statistically significant

1. Fisher's exact test

m&m example : knowing expected distribution, how special the observation is, given by p-value

2. Chi-square test → p-value

Compare observed values to expected values assume there is no relation between mutated gene and cancer

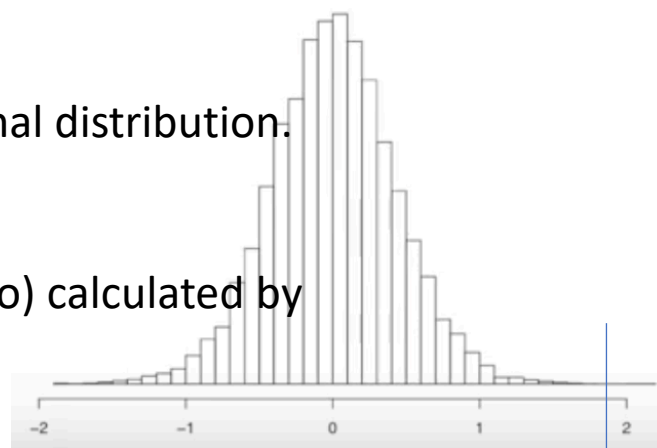
3. The Wald test → confidence interval

- Takes advantage of $\log(\text{odds ratio})$ is normally distributed (as well as $\log(\text{odds})$)

- Compute standard deviation of the normal distribution. (e.g. 0.43)

- Compute p-value based on $\log(\text{Odds ratio})$ calculated by observed data

- $1.93/0.43 = 4.11 \sigma$



$\log(\text{Odds ratio}) = 1.93$

- $p(\text{has cancer}) = 29/356 = 0.08$ (expected)
- if no relationship, for those has mutated genes : $140 * 0.08 = 11.2$
- No cancer for **128.8** people (expeted)

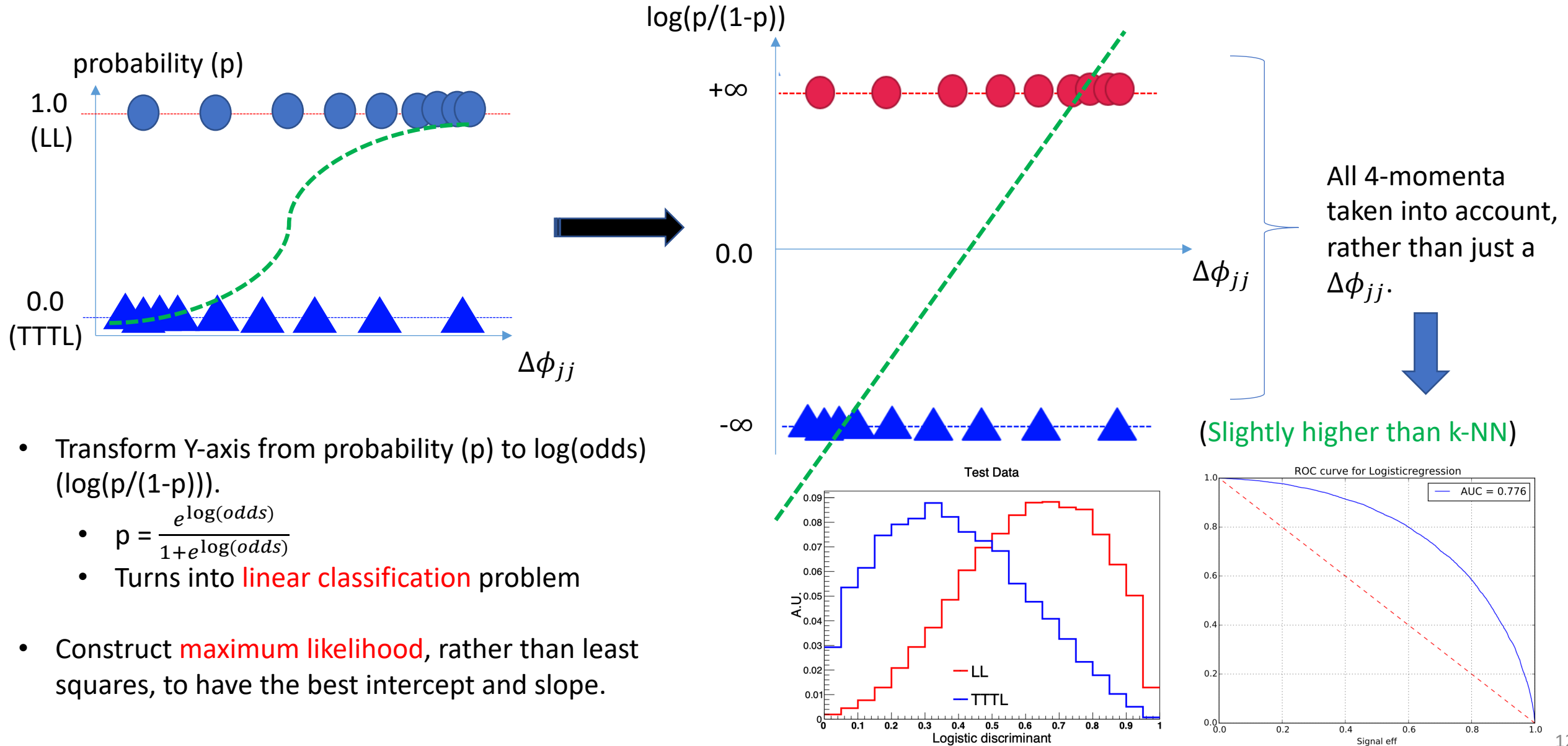


Expected Values

		Has Cancer	
		Yes	No
Yes	11.2	128.8	
No	17.3	198.7	

p-value : 0.00001

Logistic regression



Boosted decision tree

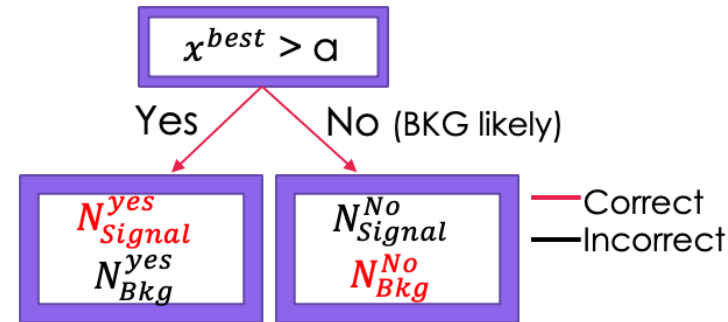
[1] **Adaptive boost** (a.k.a. AdaBoost)

- Forest of stumps
- Diversity of weight for each stump
- Each stump is made by taking the previous stump's mistakes into account

• Given Dataset : $[[x_1, x_2, \dots, x_N], [t_1, t_2]]$

1. Give weights on each data point (w_i^k)
 - Uniform distribution if no preference for initial weighting

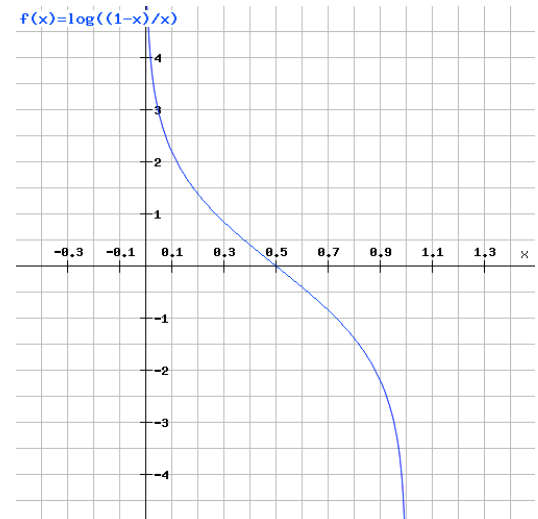
2. Find the best classifying variables among given X with weights, denote as x^{best}
 - Taken as a test of a stump



3. Compute misclassification rate : ϵ_k

$$\epsilon_k = \frac{\sum_{i=1}^N w_i^k \times \text{isMisclassified}_k(i)}{\sum_{i=1}^N w_i^k}$$

4. Weight this stump as $\alpha_k = \beta \times \ln((1 - \epsilon_k)/\epsilon_k)$

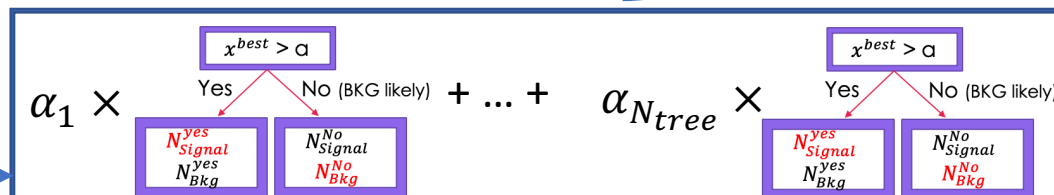


5. Reweight Misclassification data points, while keeping the correctly classified ones

$$w_i^k \rightarrow w_i^{k+1} = w_i^k \times e^{\alpha_k}$$

6. Iterate 2~5.

7. Trained model



$$\frac{1}{\sum_{k=1}^{N_{tree}} \alpha_k} \sum_{k=1}^{N_{tree}} \alpha_k T_k(i)$$

Boosted decision tree

[2] Gradient boost (On Regression)

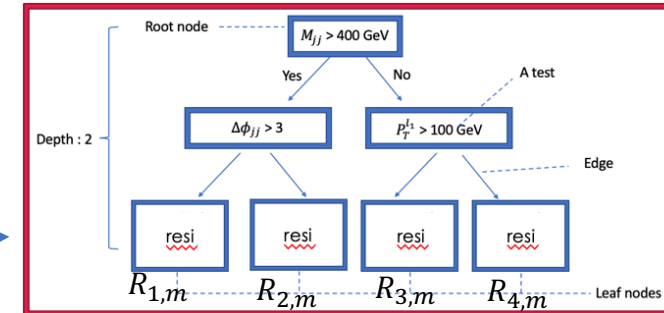
- Given Dataset : $[[x_1, x_2, \dots, x_N], [y_1, y_2, \dots, y_N]]$

- Define a differentiable loss function (ex: MSE)
- Initialized model with a constant value : $F_0(x) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$

$$\longrightarrow \text{MSE: } L(y_i, F(x)) = \frac{1}{2} \sum_{i=1}^n (y_i - F(x))^2$$

- Make pseudo-residuals for each data point by
- Establish m^{th} regression tree, where denote each leaf node as $R_{j,m}$ (j^{th} leaf of m^{th} tree)

$$r_{i,m} = -\left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}\right]_{F(x)=F_{m-1}(x)}$$

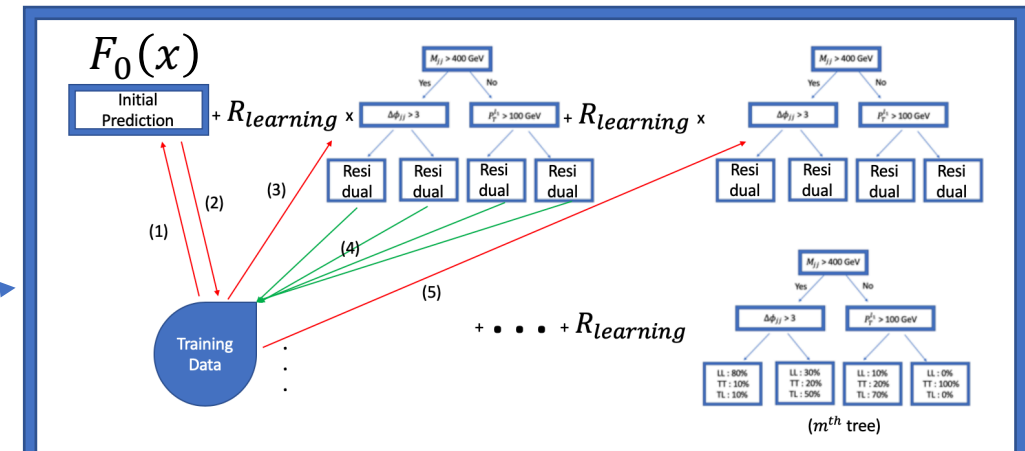


- For all $R_{j,m}$, compute $\gamma_{j,m} = \operatorname{argmin}_{\gamma} \sum_{x_i \in R_{i,j}} L(y_i, F_{m-1}(x_i) + \gamma)$

- Update $F_m(x) = F_{m-1}(x) + \nu \sum_{j=1} \gamma_{j,m} I(x \in R_{j,m})$ (ν is learning rate)
- Iterate 3~6.

- Trained model $F_M(x)$

- Construct regression trees to minimize residuals
 - Make $F_0(x)$ as representative values of over all given data, and reduce residuals



Multivariate classification : Boosted decision tree

[2] Gradient boost (On Classification)

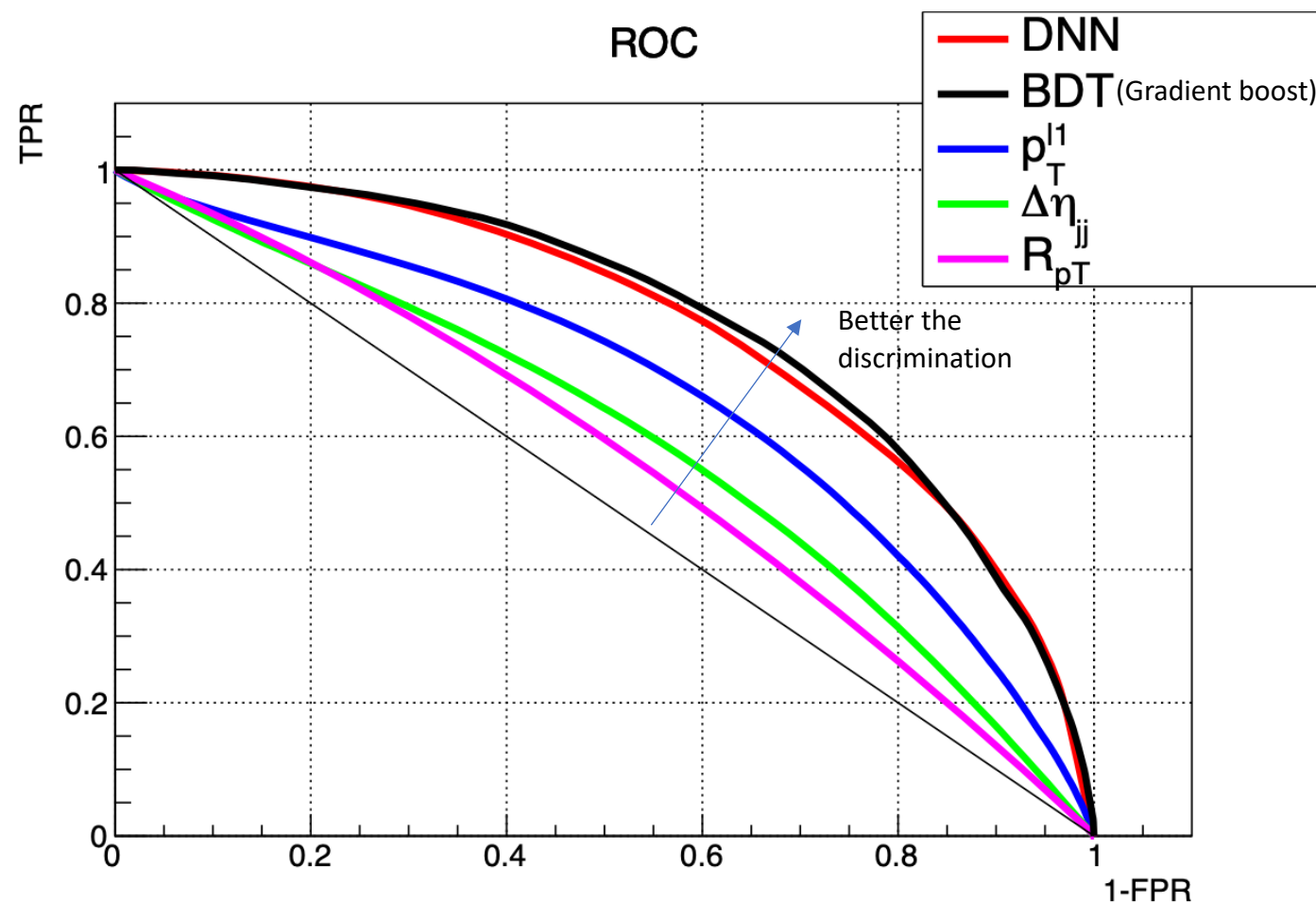
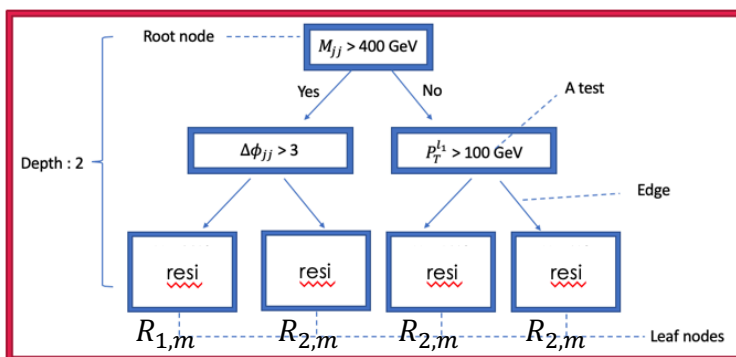
- Follow same steps with regression, usually takes minus log-likelihood as loss function :

$$L(y_i, \gamma) = - \sum_{k=1}^N y_i \times \log(p) + (1 - y_i) \times \log(1 - p)$$

where $\gamma = \log\left(\frac{p}{1-p}\right) = \log(odds)$

- The trained model could be delivered as:

$$\gamma = F_m(x) = F_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{j,m} I(x \in R_{j,m})$$



1000 trees with 5 maximum depth applied for the BDT model.