



# Essentials of XGBoost Machine Learning Algorithm

Yunfan Li PhD

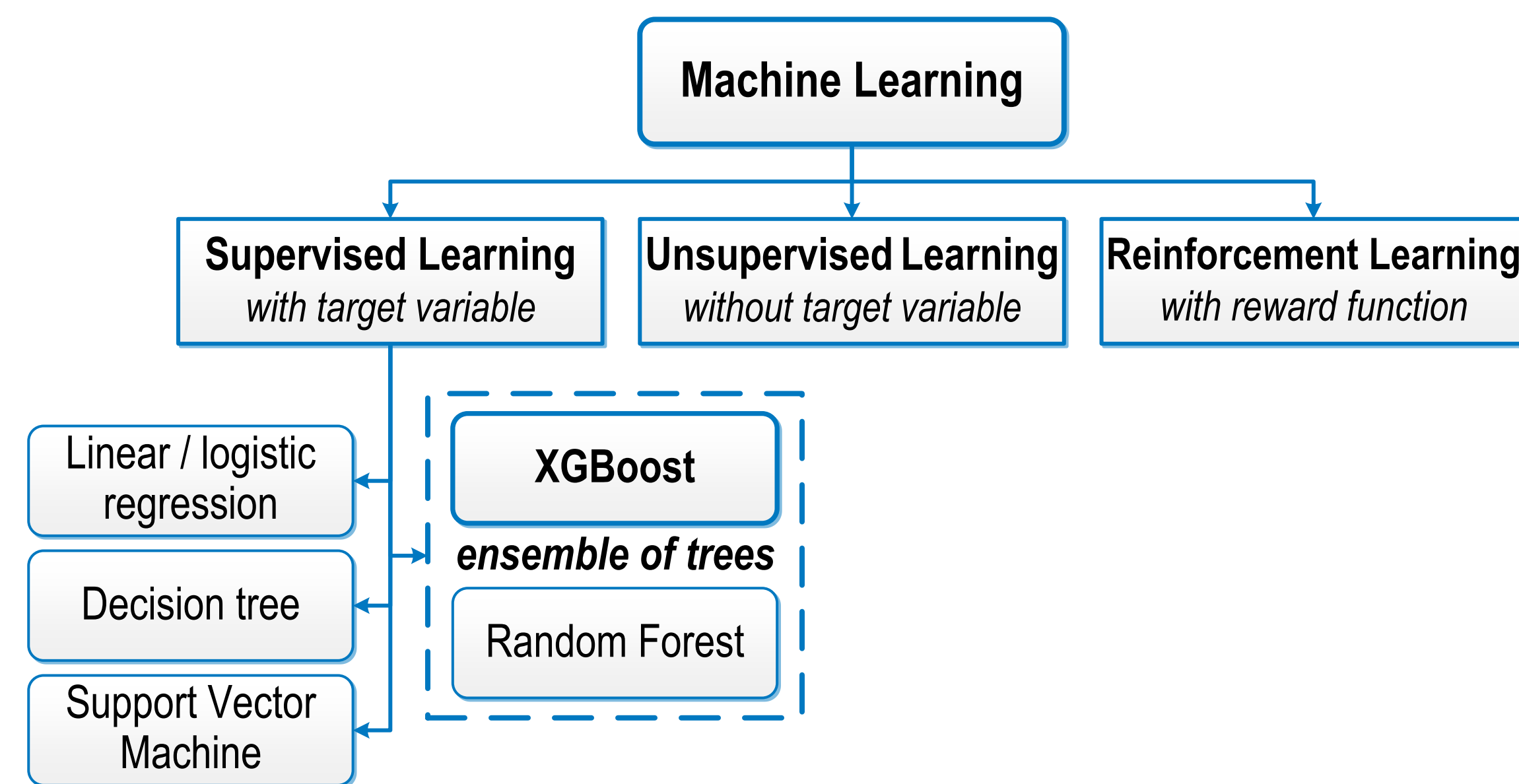
Tsinghua University, University of California at Riverside

## About This Poster

This poster is about the XGBoost, a supervised machine learning algorithm, that builds a sequence of decision trees to predict continuous and target variables. Background, model, training steps and mathematics are presented.

Background and a demo model are presented firstly. Model building, including flowchart, underlying mathematics and a demo training is illustrated. Missing value processing, feature importance and advantage summary are also presented.

## Background



**Ensemble method** A machine learning technique that combines multiple base models to obtain better performance.

**XGBoost and Random Forest** Two tree-based ensemble learning algorithms that builds trees sequentially and in parallel. Final outputs of both are weighted sums of individual tree outputs.

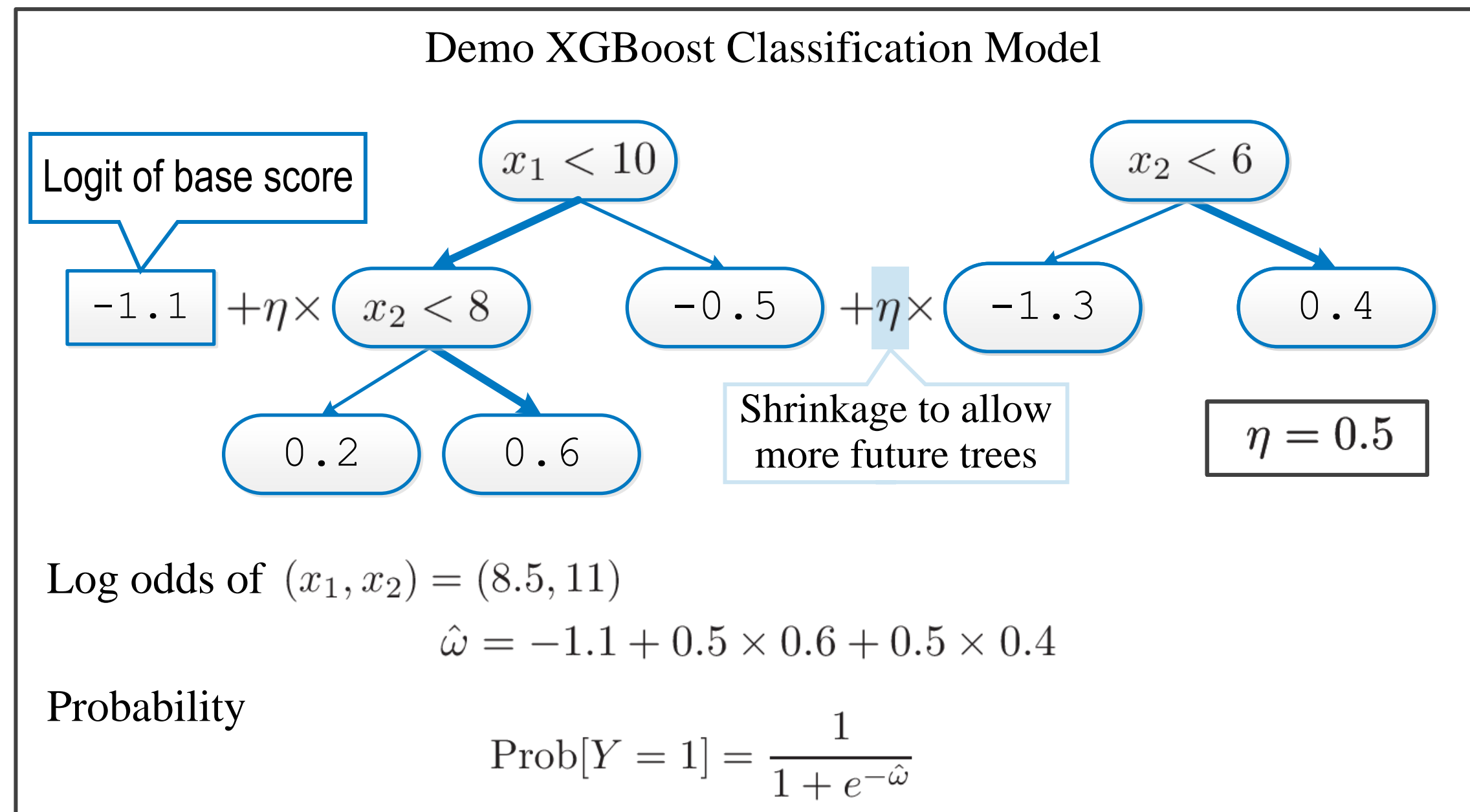
**XGBoost** stands for **eXtreme Gradient Boosting**. Its core algorithm is gradient boosting that builds a sequence of decision trees based on Gradients of cost function. The sum of scaled outputs of the trees is the final output.

## XGBoost Classification and Regression

XGBoost classification for binary target and regression for continuous target are under the same framework that sequentially builds decision trees and sum of outputs of the trees is the final output.

**Classification model** outputs the log odds of binary target  $\ln \frac{\Pr[y=1]}{\Pr[y=0]}$ .

**Regression model** outputs the expectation of continuous target  $\hat{y}$ .



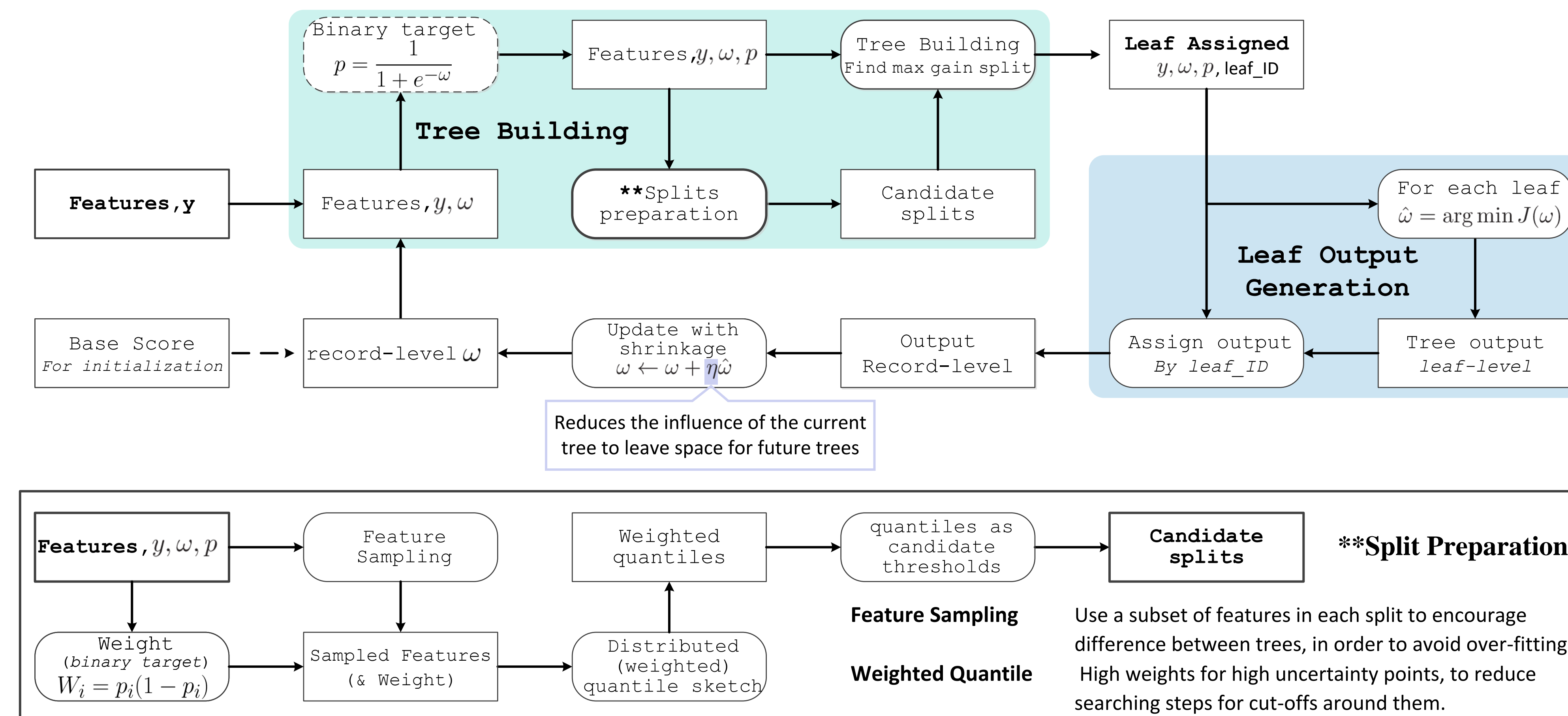
## Model Building Steps

**Initialization** Set a initial score to all points.  
Regression  $\bar{y}$   
Classification  $\ln \frac{\bar{y}}{1-\bar{y}}$

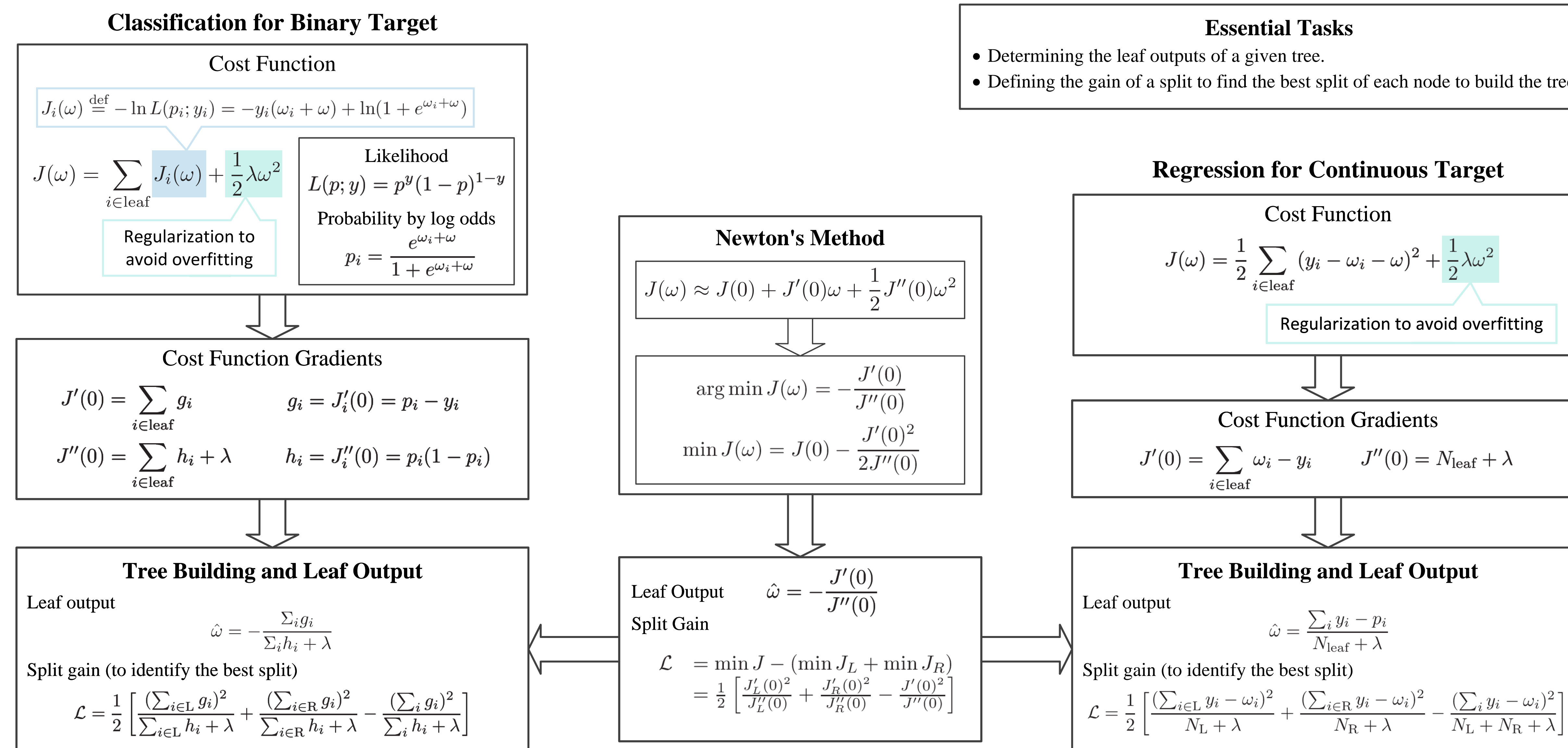
**Iteration** Generate tree output based on previous score.

Step 1: build a tree based on previous score, target and features.  
Step 2: generate leaf outputs based on previous score and target.

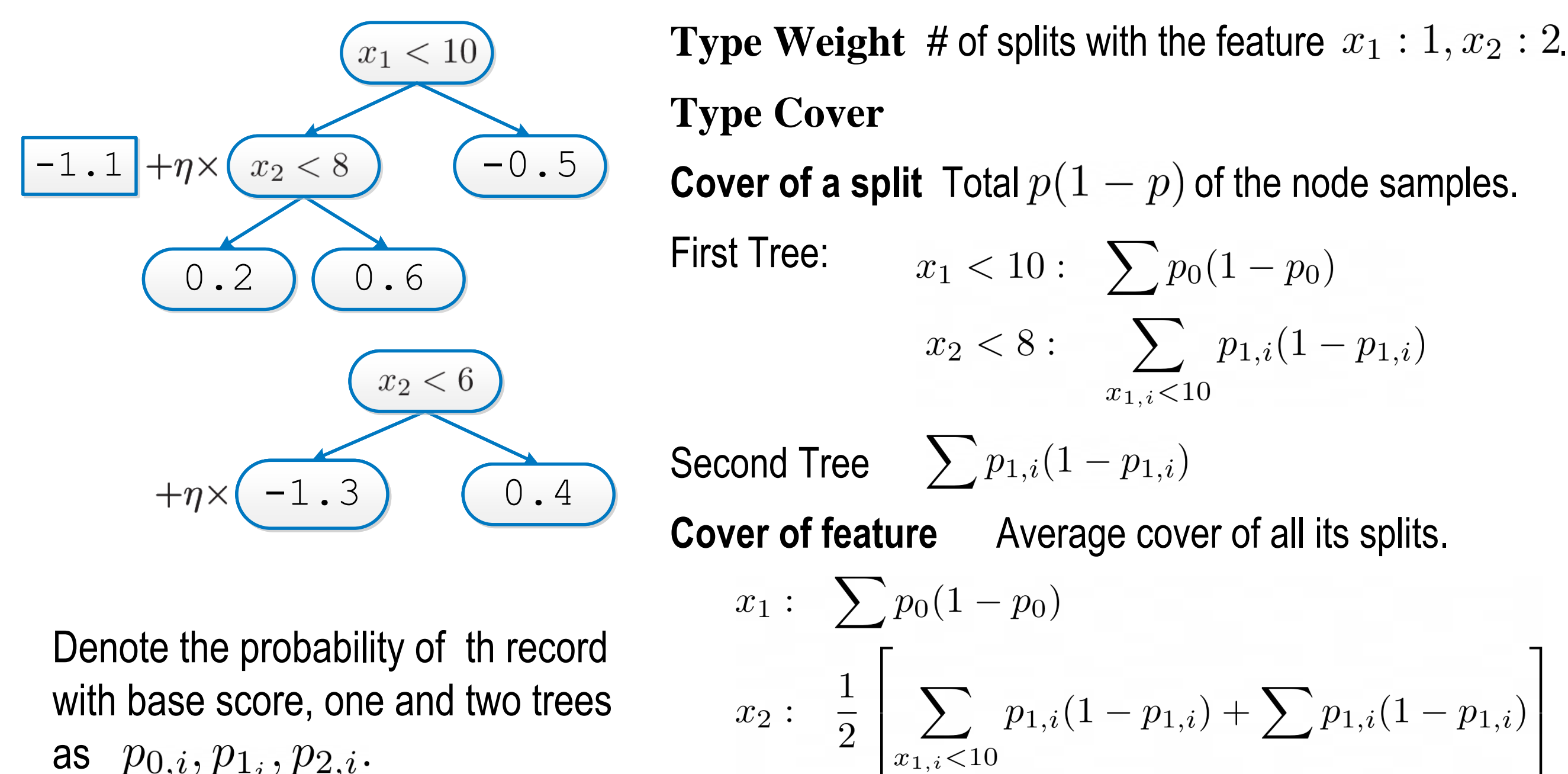
## Model Building Flowchart



## Mathematics of XGBoost Classification and Regression



## Three Types of Feature Importance



## XGBoost in Python — model building

```
df=pd.DataFrame({'y':[0,0,0,0,1,1,1,0,1],
                  'x1':['H','H','H','H','H','H','L','L',np.nan,'H'],
                  'x2':[2,8,3,2,2,6,3,2,7,4,5,1,3,4,2],
                  'x3':[np.nan,2,4,2,2,6,3,2,7,6,2,2,4] })
X=df[['x1','x2','x3']].replace(r'^s+$',np.nan,regex=True)
cat_cols,num_cols,y=['x1','x3','x4'],df['y']
```

### Feature preprocessor

- Preprocessor: impute missings, one-hot encode.

```
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
cat_transformer=Pipeline(steps=[
    (["imputer",SimpleImputer(strategy='constant', fill_value='NA')],
     ("enc",OneHotEncoder(sparse=False, handle_unknown='ignore')))]))
preprocessor = ColumnTransformer([("numerical", "passthrough", num_cols),
                                  ("categorical", cat_transformer,cat_cols)])
```

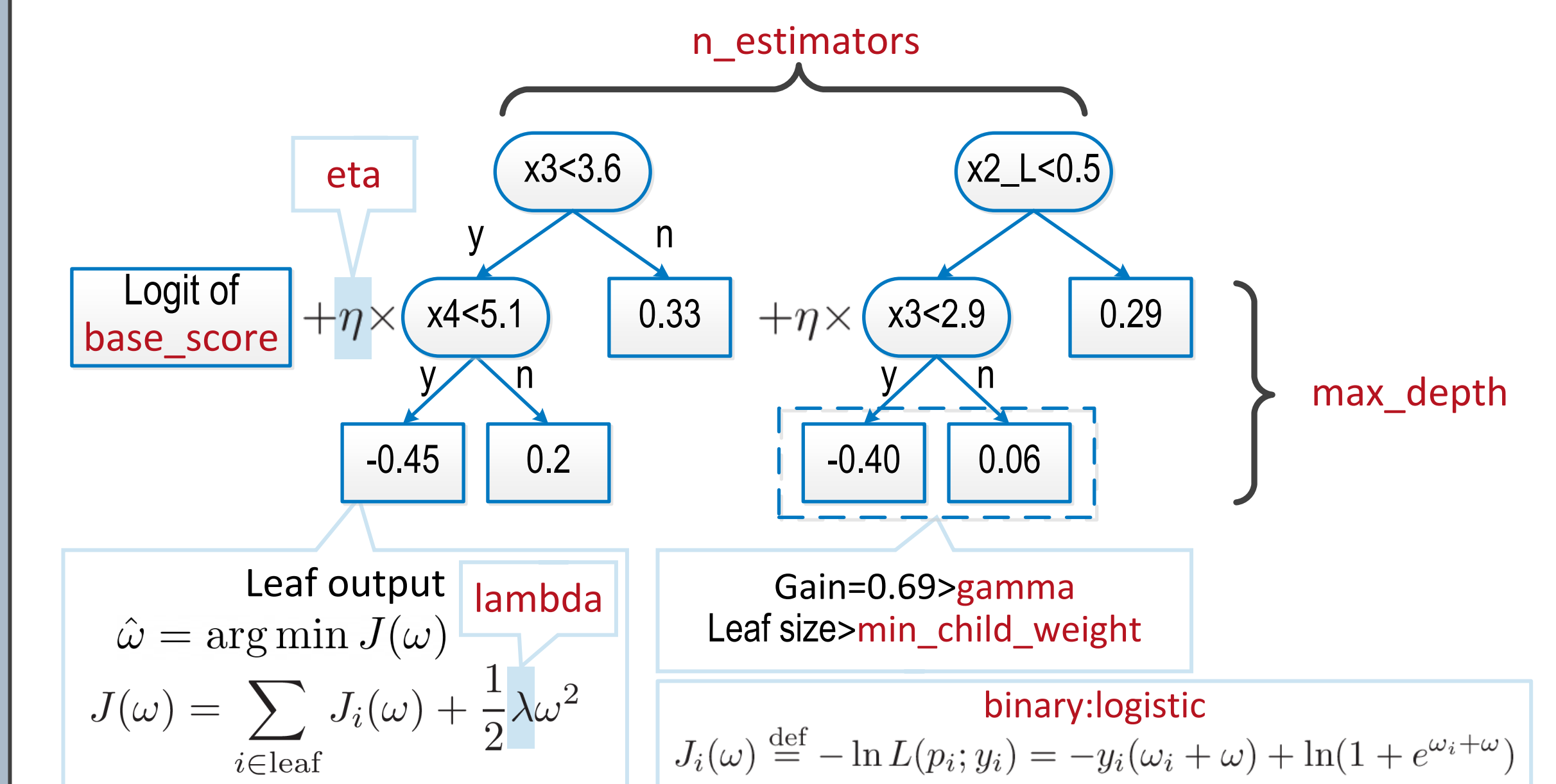
### Demo

```
preprocessor.fit(X)
imputer=preprocessor.named_transformers_["categorical"].named_steps["imputer"]
enc=preprocessor.named_transformers_["categorical"].named_steps["enc"]
X_imputed=imputer.transform(X[["x1"]])
```

enc.categories_	X_imputed	enc.transform(X_imputed)
[ array(['H', 'L', 'NA'], dtype=object)]	[[ 'H'] [ 'H'] [ 'H'] [ 'H'] [ 'H'] [ 'H'] [ 'L'] [ 'L'] [ 'NA'] [ 'H']]	[[1. 0. 0.] [0. 0. 1.] [0. 1. 0.] [0. 1. 0.] [0. 1. 0.] [1. 0. 0.] [1. 0. 0.] [1. 0. 0.] [1. 0. 0.] [0. 1. 0.]

### XGBoost Model

```
param = { 'objective':'binary:logistic', 'base_score':0.5, 'n_estimators':2, 'eta':0.5,
          'max_depth':2, 'gamma':0.6, 'min_child_weight':0, 'lambda':1}
model = xgb.XGBClassifier(n_jobs=-1, **param)
```



### Create Pipeline and Fit Model

```
xgb_model = Pipeline([("preprocessor", preprocessor), ("model", model)])
xgb_model.fit(X,y)
booster = xgb_model.named_steps["model"].get_booster()
preprocessor = xgb_model.named_steps["preprocessor"]
```

## XGBoost in Python — modeling results

**Model Definition** trees=booster.get\_dump()

**Model Output** proba=xgb\_model.predict\_proba(X)

### Feature Importance

```
cat_ohe_names = preprocessor._transformers_[1][1]['enc'].get_feature_names(['x1'])
f_map = {'f'+str(i):name for i,name in enumerate(num_cols+list(cat_ohe_names))}
scores = booster.get_score(importance_type='weight') # or 'total_gain','total_cover'
{f_map[f]:scores.items() for f in scores.keys() if f in f_map}
```