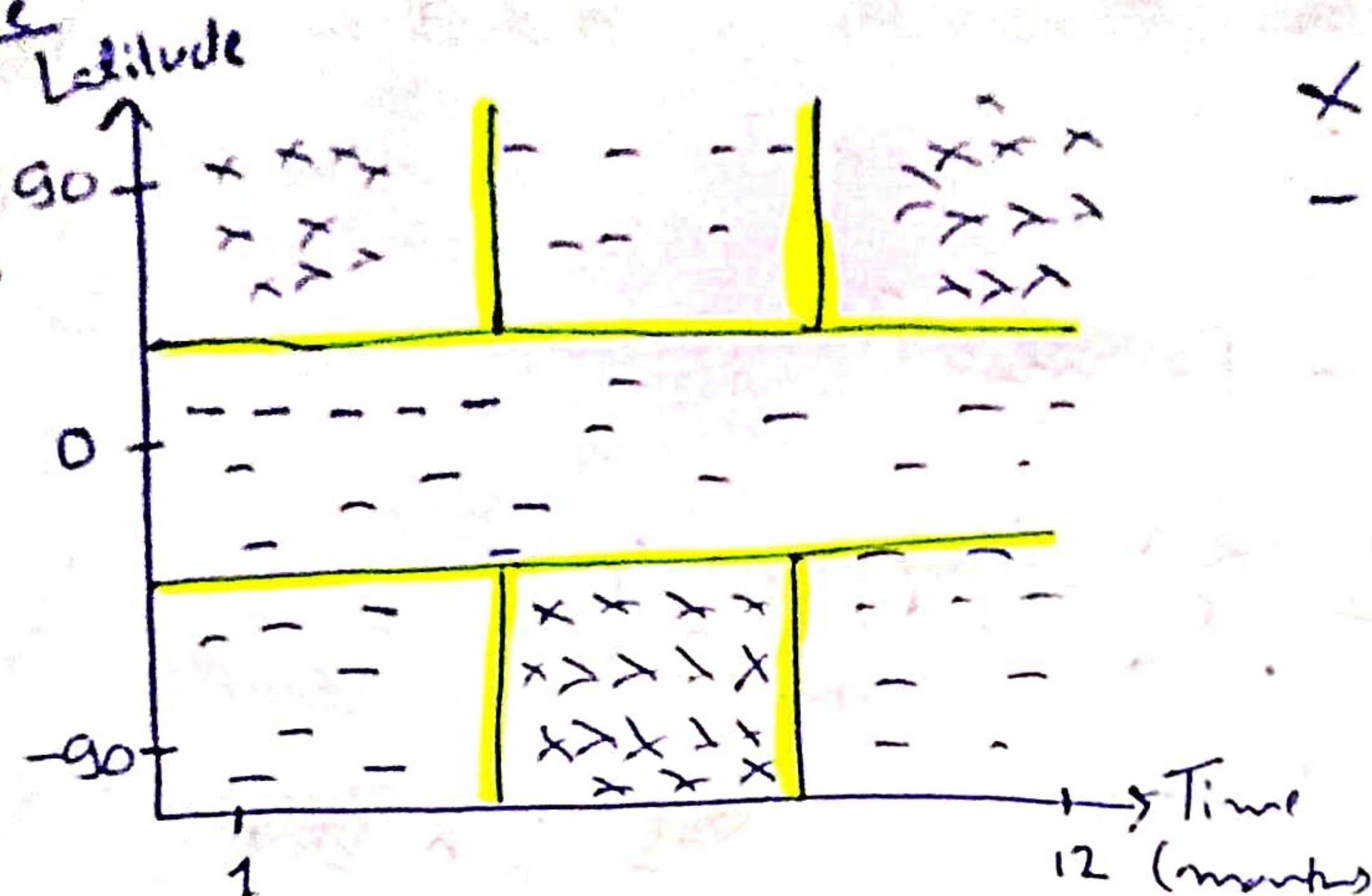


(10)

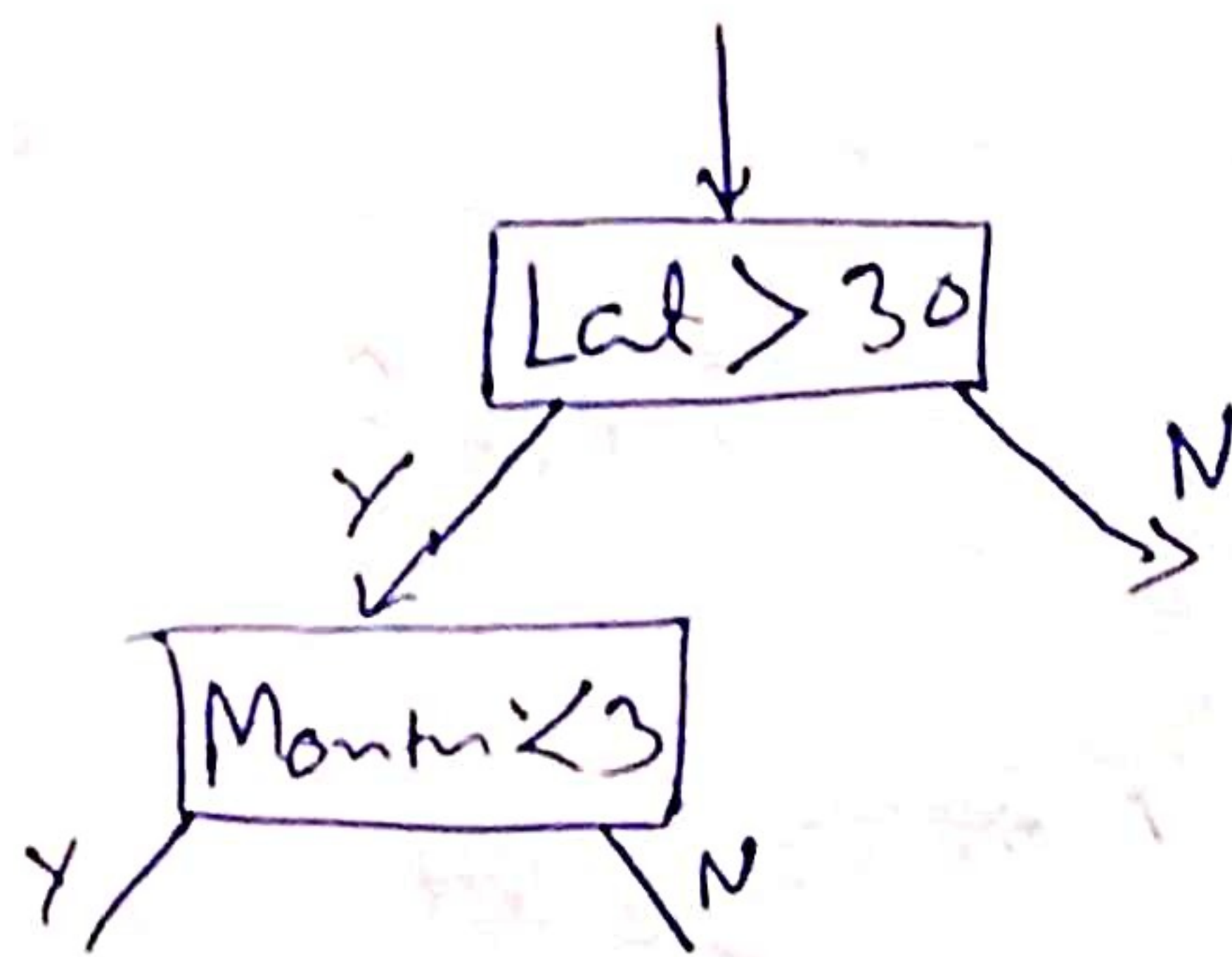
Decision Trees & Ensemble Methods

* Decision Trees

Example



⇒ Greedy, Top-Down, Recursive Partitioning.



⇒ Region R

⇒ Looking for a split S_p

$$S_p(j, t) = \left(\{x | x_j < t, x \in R_p\}, \right.$$

Feature number Threshold

$$\{x | x_j \geq t, x \in R_p\} \Bigg)$$

→ R_1

→ R_2

⇒ How to choose splits?

⇒ Define $L(R)$: loss on R

⇒ Given C classes, define \hat{P}_c to be the proportion of examples in R that are of class c .

$$L_{\text{misclass}} = 1 - \max_c \hat{P}_c$$

$$\max_{j, t} L(R_p) - (L(R_1) + L(R_2))$$

Parent loss

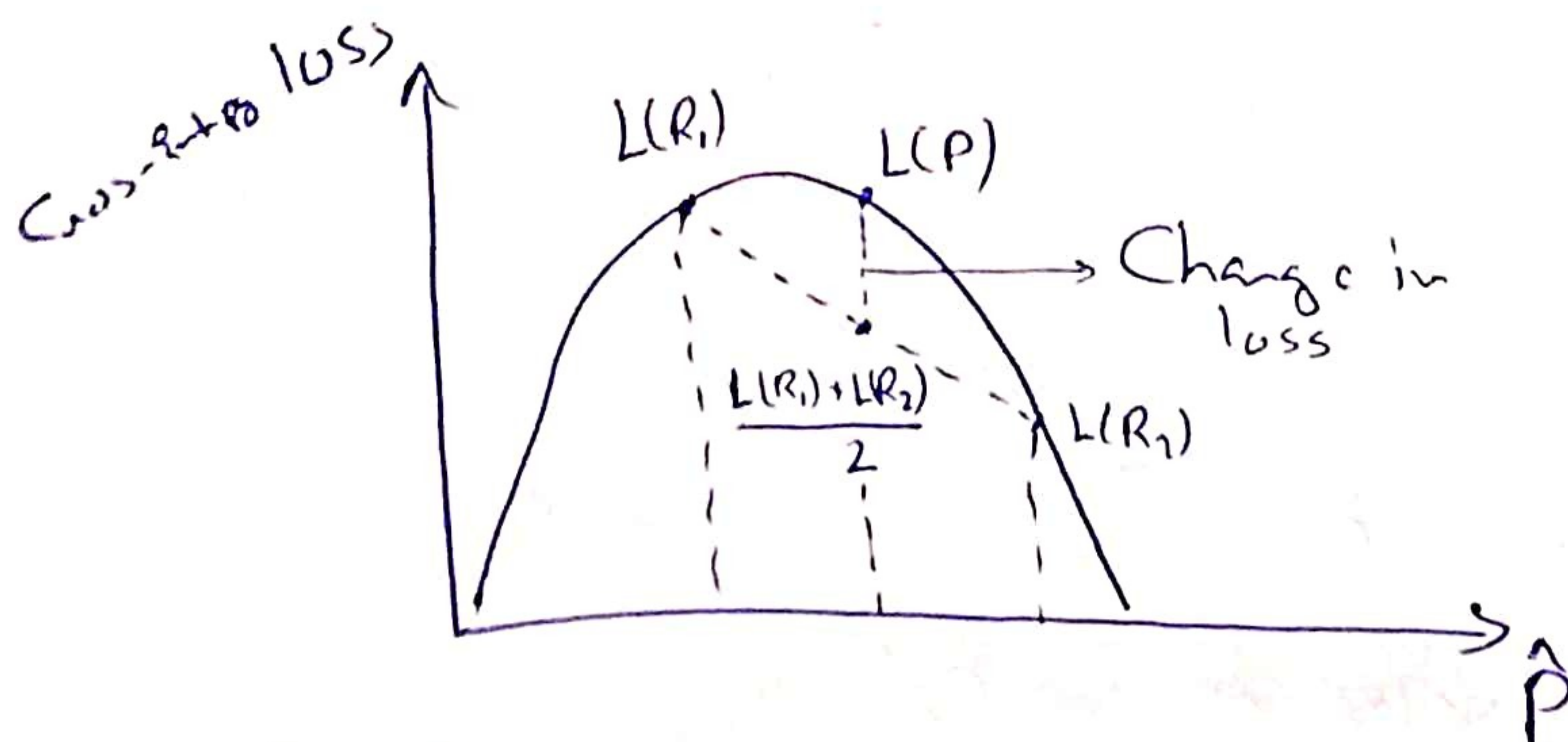
Children loss

Misclassification loss

⇒ Misclassification Loss has issues!!

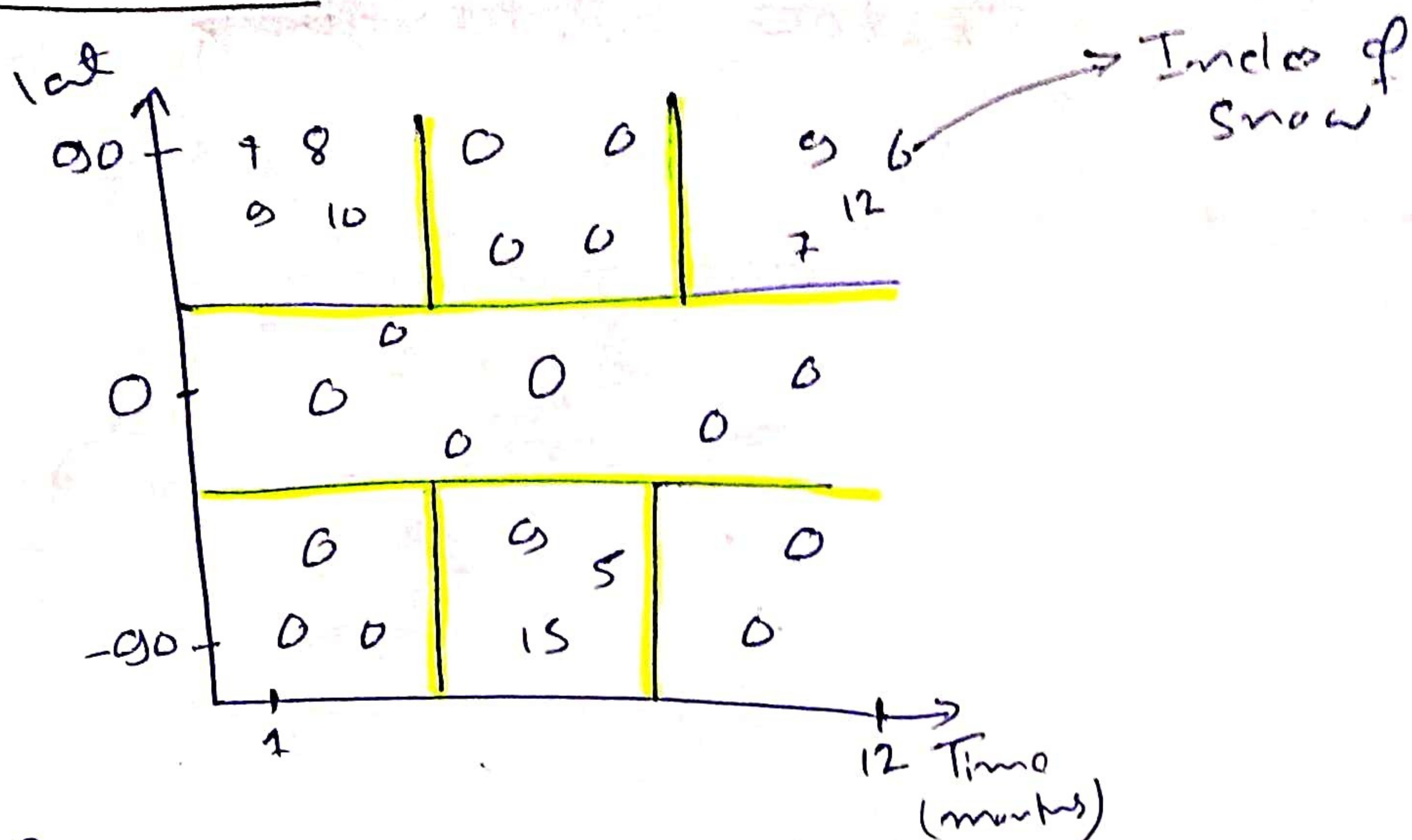
⇒ Instead, define cross-entropy loss

$$L_{\text{cross}} = \sum_c \hat{P}_c \log_2 \hat{P}_c$$



$$\text{Gini Loss} = \sum_c \hat{P}_c (1 - \hat{P}_c)$$

* Regression Trees



Predict $\hat{y}_m = \frac{\sum_{i \in R_m} y_i}{|R_m|}$

$$L_{\text{Squared}} = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$$

* Regularization of DTs {Heuristics}

- ① min leaf size
- ② max depth
- ③ max number of nodes
- ④ min decrease in loss
- ⑤ Pruning (misclassification with validation set)

⇒ Down side of decision trees:

→ No additive structure

→ High Variance model

→ Generally Low predictive accuracy

* Ensembling

⇒ Take X_i 's which are RV and are IID.

$$\text{Var}(X_i) = \sigma^2$$

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum X_i\right) = \frac{\sigma^2}{n}$$

⇒ Drop the independence assumption

So now X_i 's are id

X_i 's correlated by ρ

$$\text{Var}(\bar{X}) = \rho\sigma^2 + \frac{1-\rho}{n}\sigma^2$$

⇒ Ways to ensemble:-

- 1) different algorithms
- 2) different training set
- 3) Bagging (Random Forest)
- 4) Boosting (Adaboost, Xgboost)

⊕ Bagging (Bootstrap Aggregation)

- Have a true population P
- Training Set $S \sim P$
- Assume $P=S$
- Bootstrap sample $Z \sim S$

"We will take a bunch of bootstrap samples, train separate models on each and then average their output"

Bootstrap Samples Z_1, \dots, Z_M

Train model G_m on Z_m

$$G(m) = \frac{\sum_{m=1}^M G_m(x)}{M}$$

Random Forests:

↳ At each split, Consider only a fraction of your total features.

* Boosting

↳ Decreasing bias
↳ Additive

⇒ It increases the weights on the mistakes of previous model.