

Day 12, Oct-14, 2024

## # Setting up Supervised learning Problems - Part I

We, human causally make predictions like estimating and predicting the bus arrival time, raining or sunny without any explicit mathematics.

But our Computer relies on mathematics

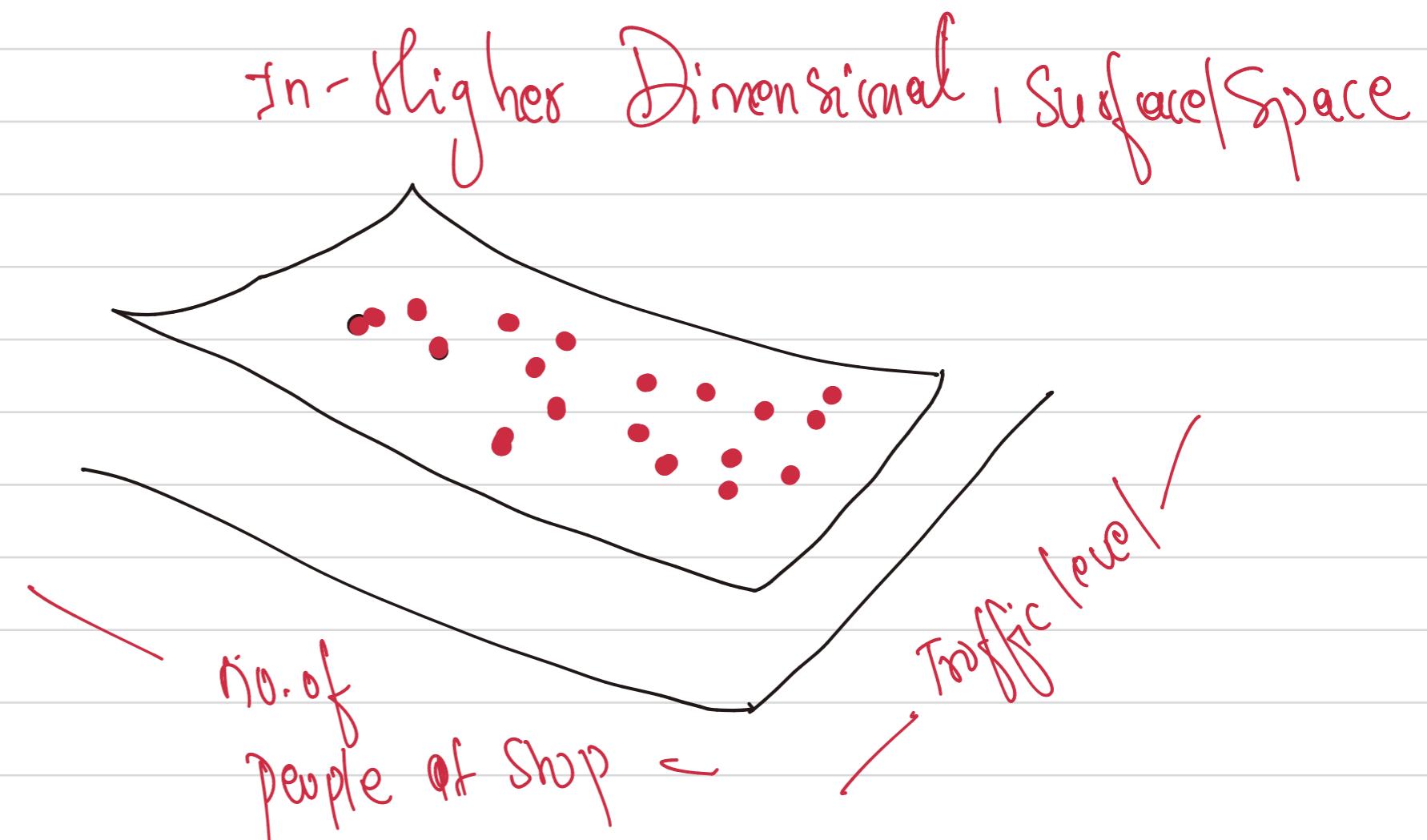
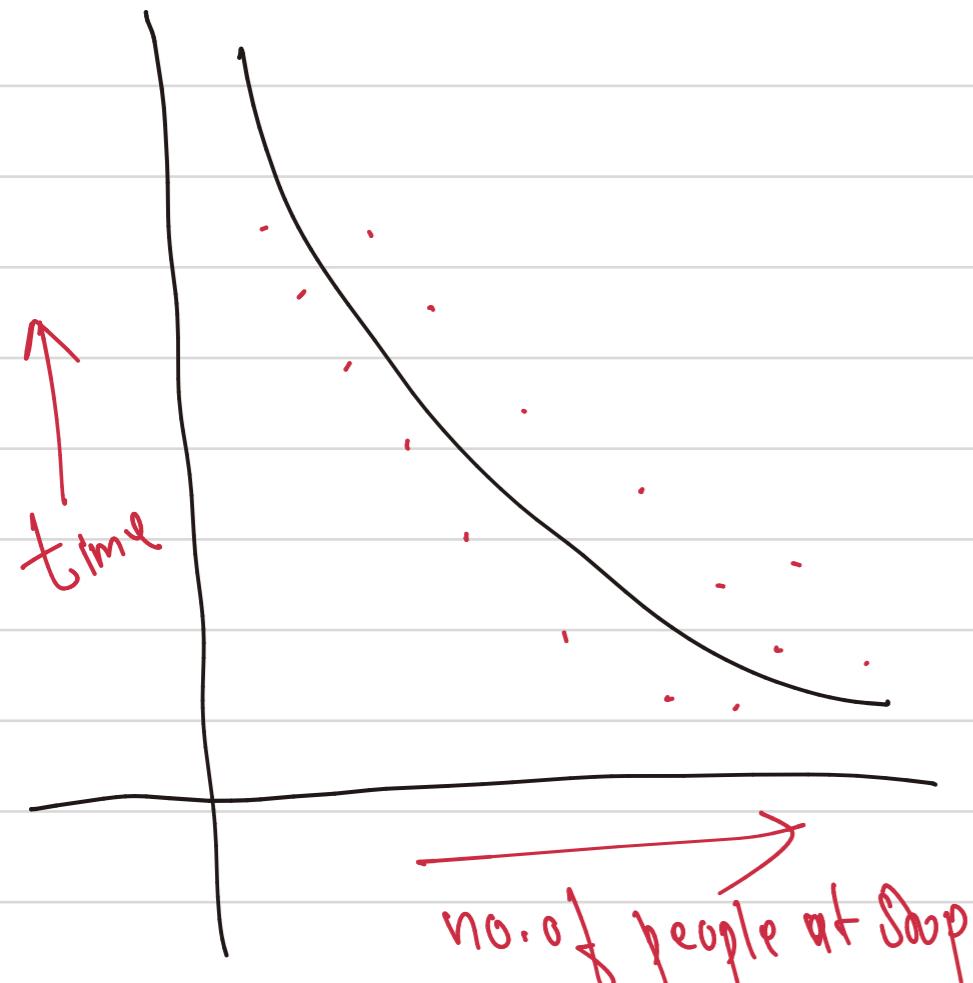
$$x_i \in \mathbb{R}^D$$

where  $D$  is the characteristics of envir

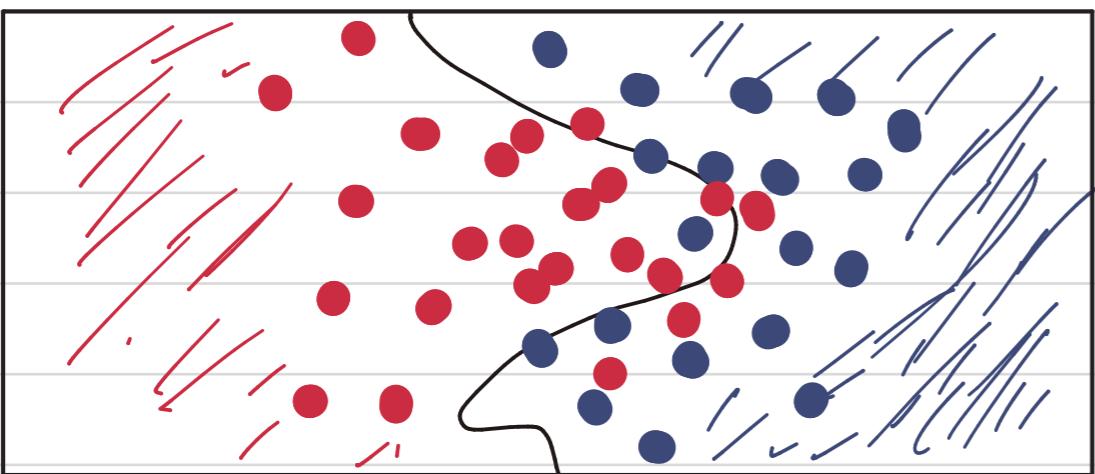
$$\text{So, } \left( \begin{array}{l} x_1 | y_1 \\ x_2 | y_2 \\ \vdots \\ x_n | y_n \end{array} \right)$$

$x_i \Rightarrow$  features like people, weather

- 1 Regression  $\rightarrow$  Response is Continuous
- 2 Classification  $\rightarrow$  Integer/Discrete or Classification



Classification Example Blood Pressure  $\rightarrow$  1 D



→ 2D Example of Blood Pressure  
Prediction (Categorical)

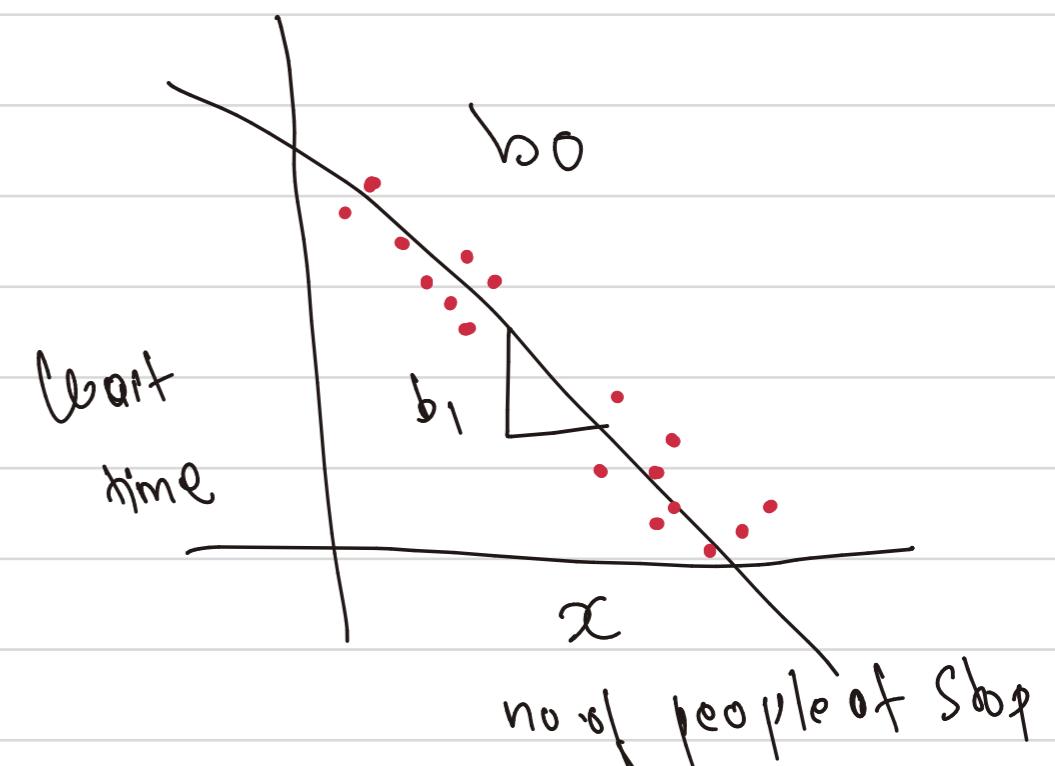
## # linear Models

$$f_b(x) = b_0 + b_1 x$$

Eq for 'n' features.

$$f_b(x) = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_D x_D$$

$$\boxed{f_b(x) \Rightarrow b^T x}$$



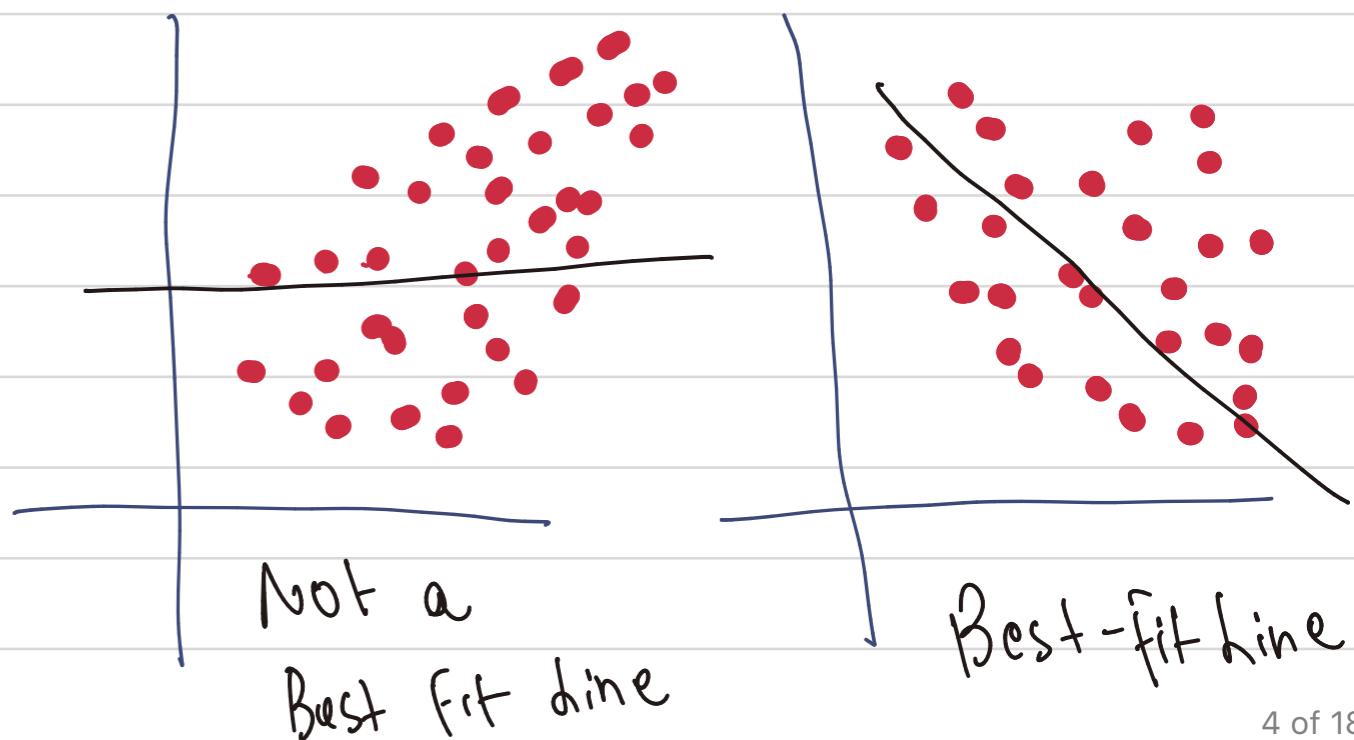
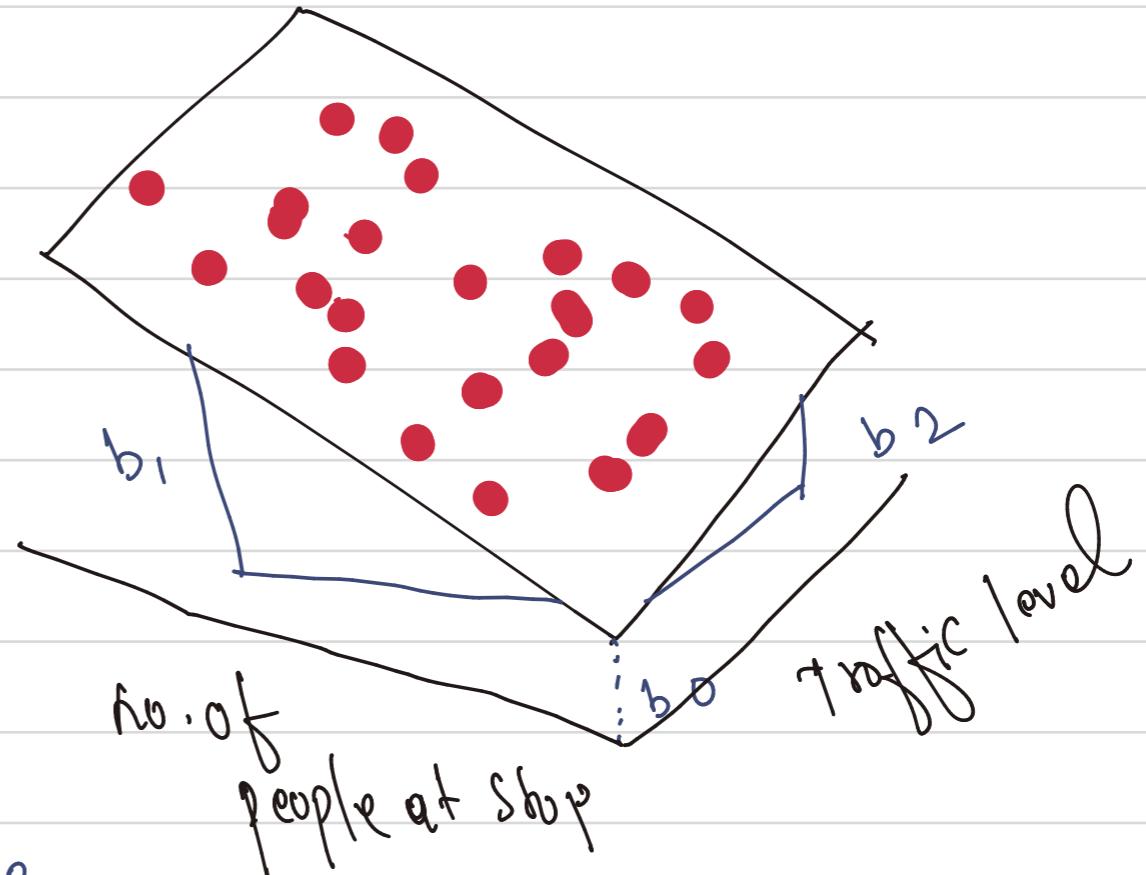
Since waiting time decreases with Counter

## More Geometric Interpretation:

Finding Slope that fits the data well so we can use to optimize a "loss" function.

For linear Regression, a good choice is a Squared Error loss

$$\therefore L(b) = \sum_{i=1}^N (y_i - b^\top x)^2$$



## # For Logistic Regression

$$f_b(x) \Rightarrow \frac{1}{1 + \exp(b^T x)}$$

$$\Rightarrow \frac{1}{1 + e^{-z}}$$

where  $z = b_0 + b_1 x$

$$\text{for } z_n = b^T x$$

(Threshold = 0.5)

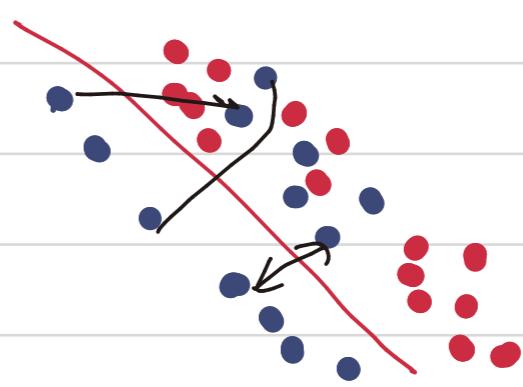
Threshold < Outcome (1)

Threshold > Outcome (0)

## # Loss function in Logistic Regression

↳ is perpendicular to predicted error.

$$\left[ \sum_{i=1}^N y_i \log(f_b(x_i)) + (1-y_i) \log(1 - f_b(x_i)) \right]$$

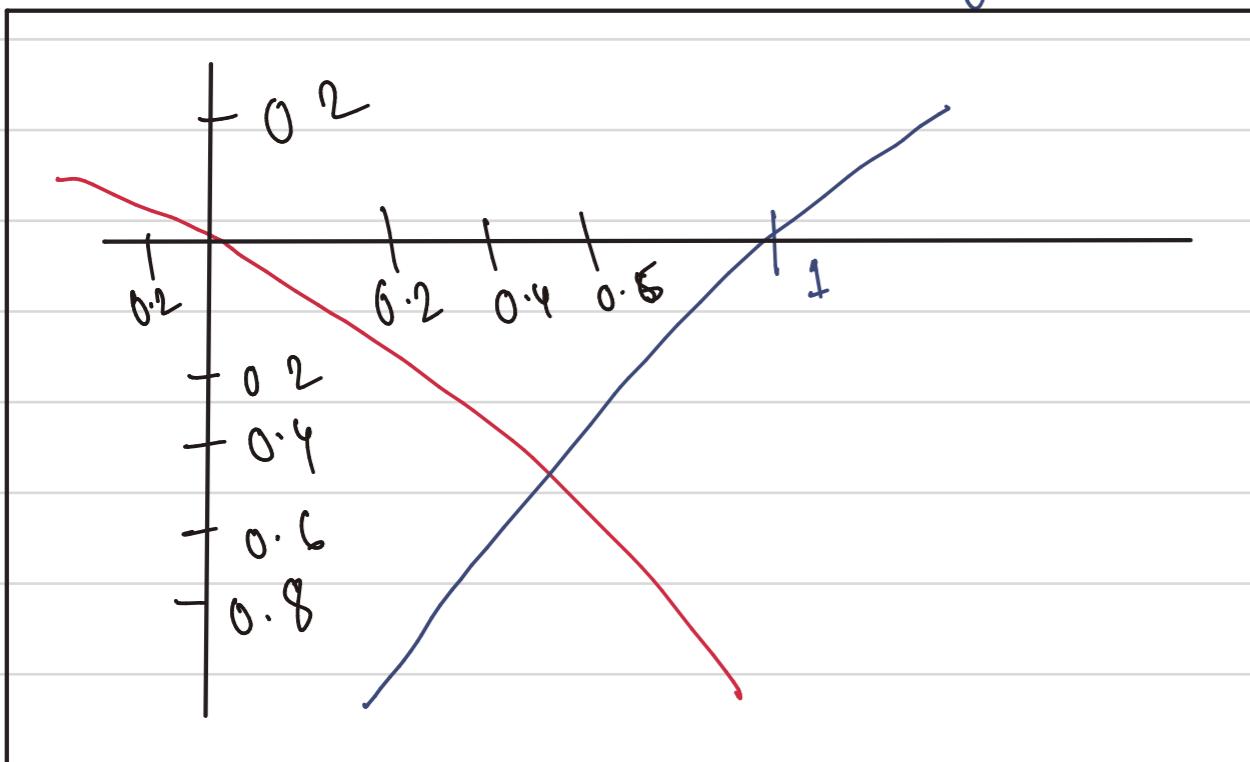


So the above loss function for logistic Regression

$$\left[ \sum_{i=1}^n y_i \log(f_b(x_i)) + (1 - y_i) \log(1 - f_b(x_i)) \right]$$

When  $y_i = 1$  red point on the graph  
 $y_i = 0$  blue point so it's a (binary Cross Entropy loss)

Graph for  $\log(x), \log(1-x)$



So for  $y_i = 0$

$$\begin{aligned} & \Rightarrow \sum_{i=1}^n 0 + (1-0) \log(1 - f_b(x_i)) \\ & \Rightarrow (1-0) \log(1 - f_b(x_i)) \end{aligned}$$

# Interpretation: if 1 is predicted with probability 1 then there is no loss  
and if 0 is predicted with probability 1 then there is loss.

$\log(x)$  and  $\log(1-x)$  for  $x=0$

$\log(0)$  and  $\log(1-0)$

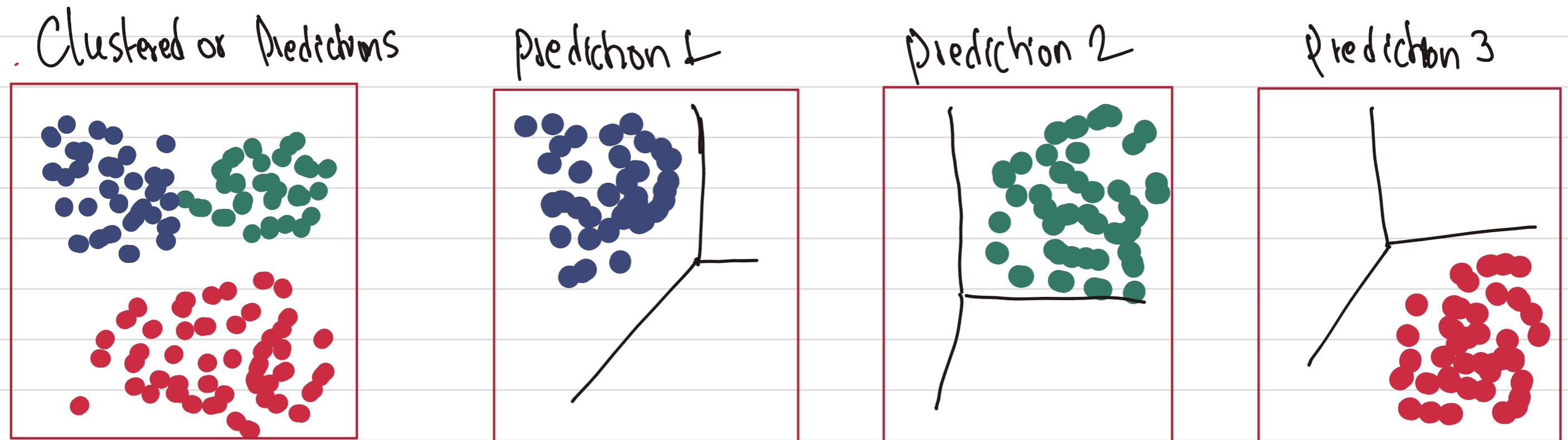
$\log(0)$  and  $\log(1)$

# When model gets confident and totally wrong it gets penalty.

# When model not gets confident and is correct (also the model gets penalty (it could be L1 and L2)).

## # Coding Part See the Code.

⇒ We do prediction using the features data so the features are  $x_1, x_2, x_3, \dots, x_n$  and finding its slope  $m$  or  $b_1$  OR Coefficient  $b_1, b_2, \dots, b_n$  and we have y-intercept or bias'  $b_0 \rightarrow$  this is a single variable (error).



for linear regression Decision Boundary is linear (But Not for Complex Process)

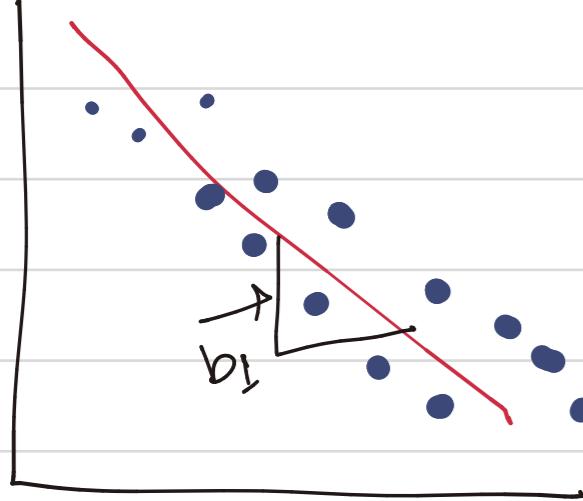
## # Sparse linear Models:

A model that knows how to ignore the irrelevant features will always do better than that tries to use them all. This is the main idea behind using sparsity in linear regression to fit the model because we can have 'n-' features' and it is likely that our model will learn noise too.

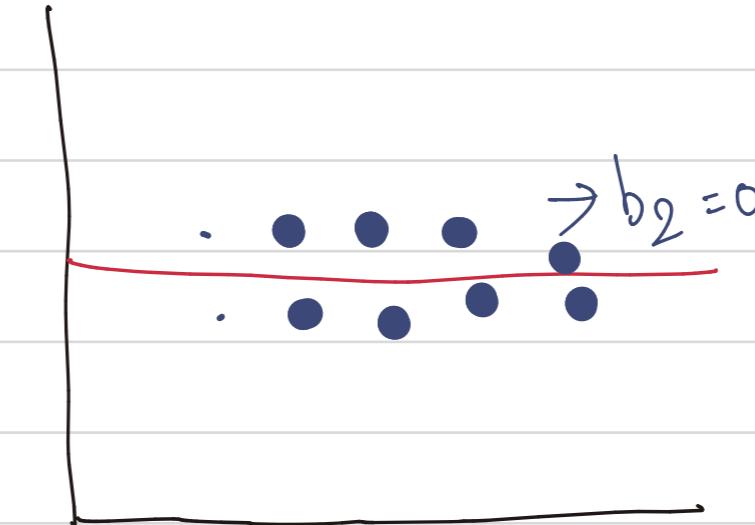
$$f_b(x) = b_0 + b_1x_1 + \dots + b_nx_n \Rightarrow b^T x$$

but what we make the assumption that many of the  $b_j$  are exactly 0. Generally, we imagine that response don't change at all we

Change some of the inputs, all else held equal



Variable 1  
With the relevant feature



Variable 2 (Slope=0)  
With the irrelevant features has no contribution

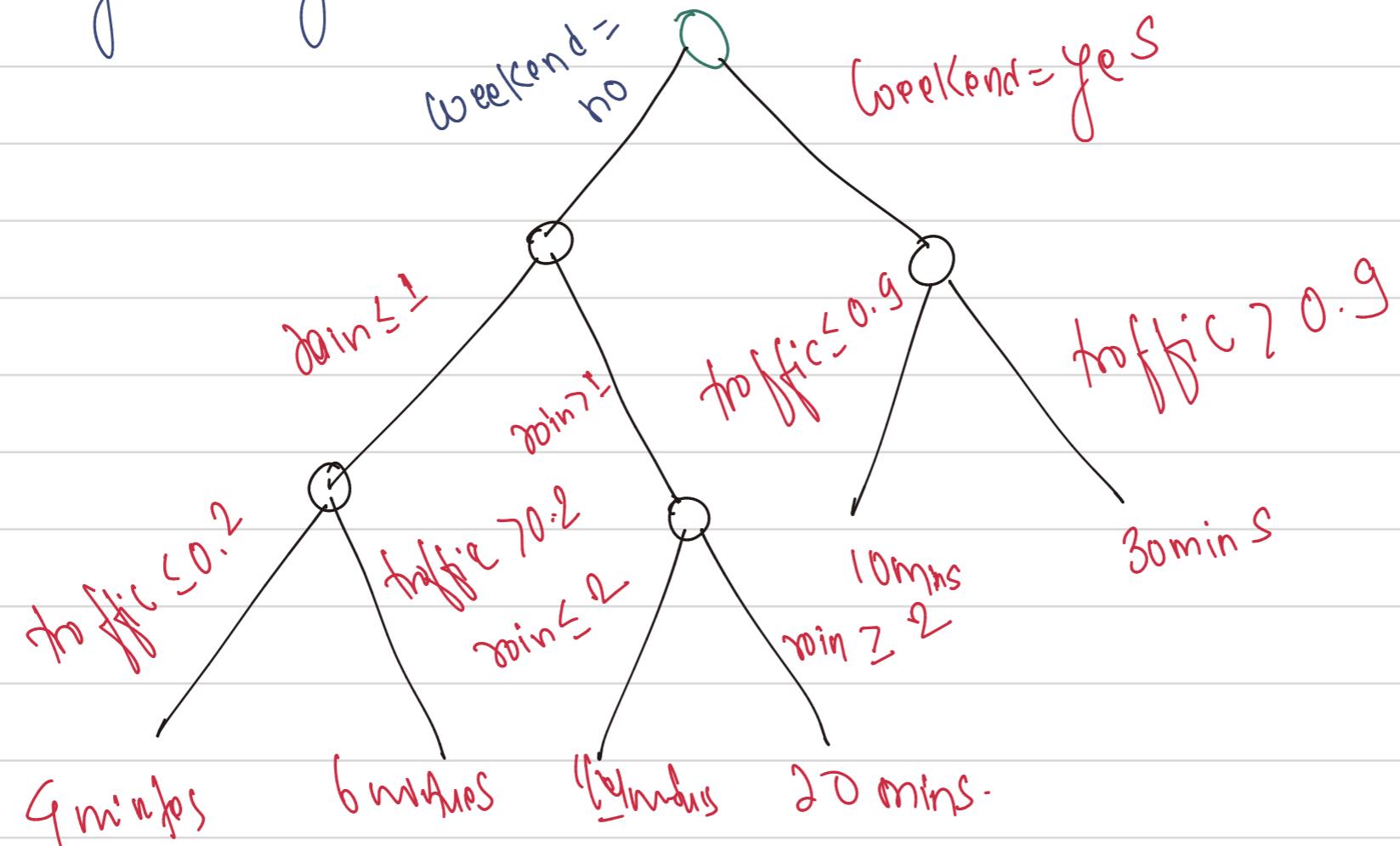
# Example Decision Boundary in Above Current Blood Pressure does  
the some level of Sparsity.

# Elastic Net in the Sparsity, So the linear-model. Elastic Net ( $\alpha = 1e-1$ )  
use shot ↑

## # Tree-Based Models

→ Tree-Based models fit a different class of curves. To motivate them,

Consider making a prediction for the bus time arrival problem using the following diagram:



We can use some logic to do either regression or classification,  
for Regression each "leaf" at the bottom of the tree is a continuous  
prediction. For Classification, we associate leaves with probabilities of  
different classes

It turns out that we can train these models using Squared Error &

Cross-Entropy ( $\times$ ) losses before though the details are beyond the  
scope of these notes.

# So, we can infer that these rules are equivalent to drawing curves  
that are piecewise constant over subsets of the input space. Let's  
convince ourselves using some pictures. First, notice that a tree with a single

Split is exactly a "curve" that takes two values depending upon the Split point.

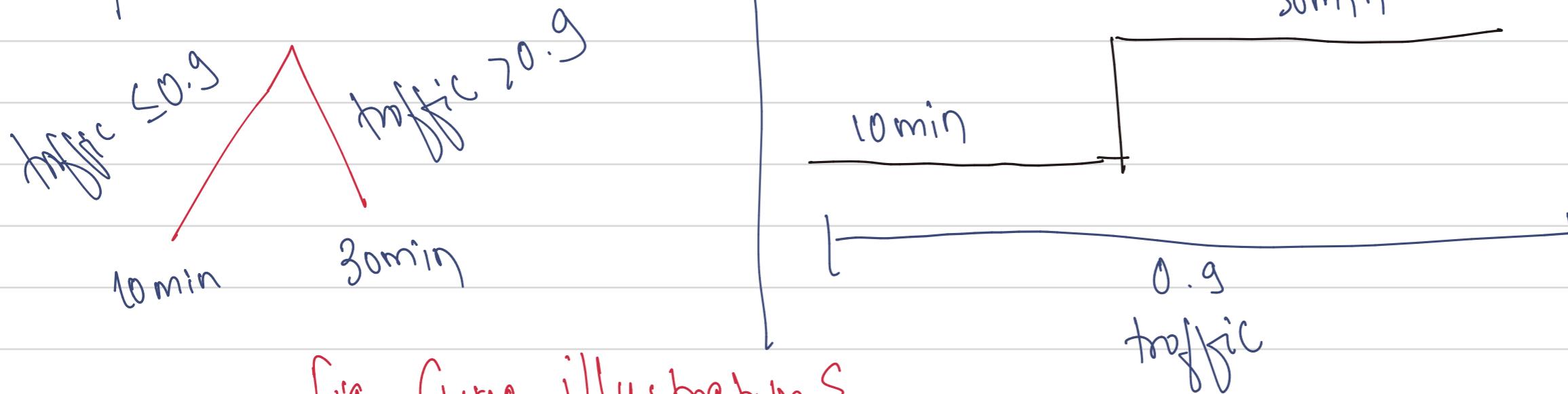
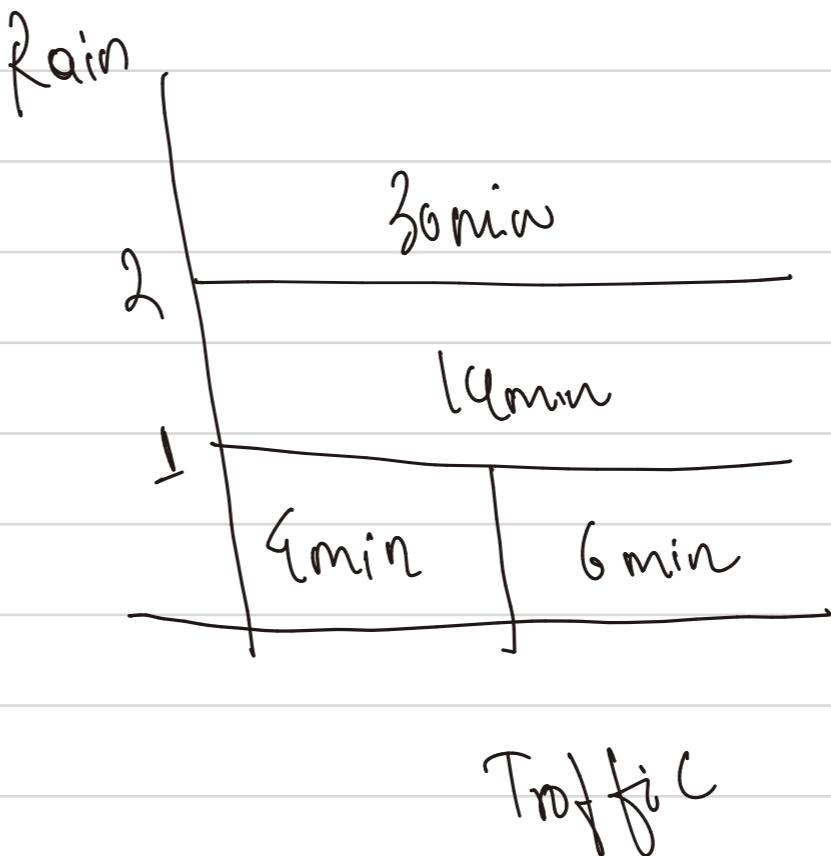
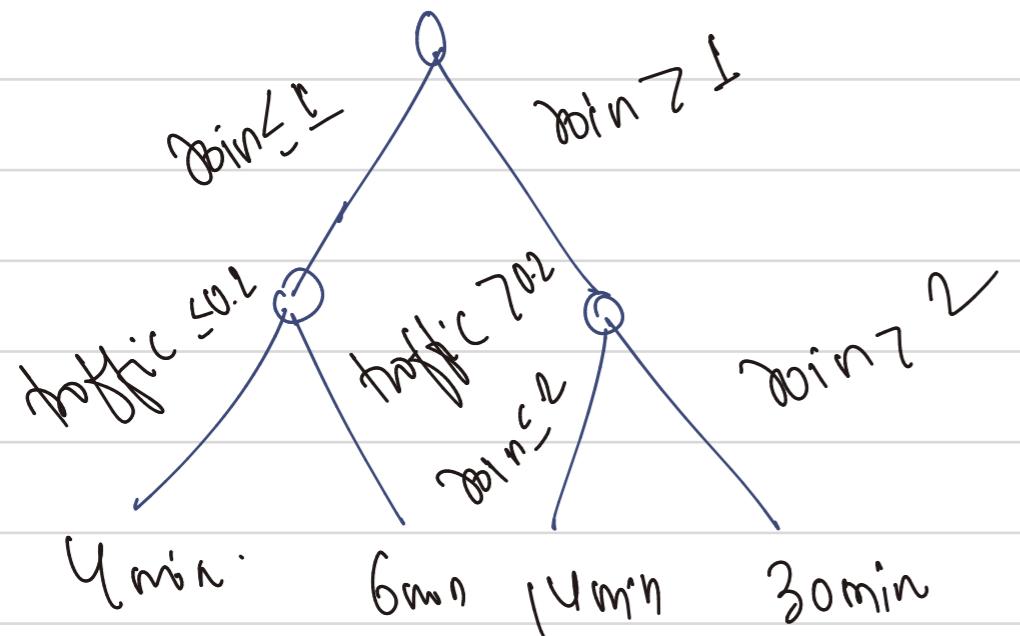


Fig. Curve illustrations

if we split the same variable deeper, it creates more steps.

What if we had two variables? Depending on the order of the splits, we create different axis-aligned partitions.

Let's take more complex -



If these tree-based methods arises the decision-tree and give rise to a new method we call it a random forest.

If it does not work well as the complexity (decision-tree) so the collection of optimized decision-tree are often called Random Forest.

The Random Forest and Gradient Boosted Decision Trees.

## Random forest

vs

## Gradient Boosted Decision Trees

→ Creates multiple decision

→ builds trees Sequentially

trees independently focus

→ Each tree corrects the error

on reducing variance (

Made by Previous one (boosting)

Avoids overfitting by averaging)

→ focusing on reducing bias

→ used for both classification

(Improves model accuracy

and Regression

Progressively)

So, the tree-base methods have non-linear decision boundary  
as compared to linear Models. Let's understand and  
summarize.

## Baselines

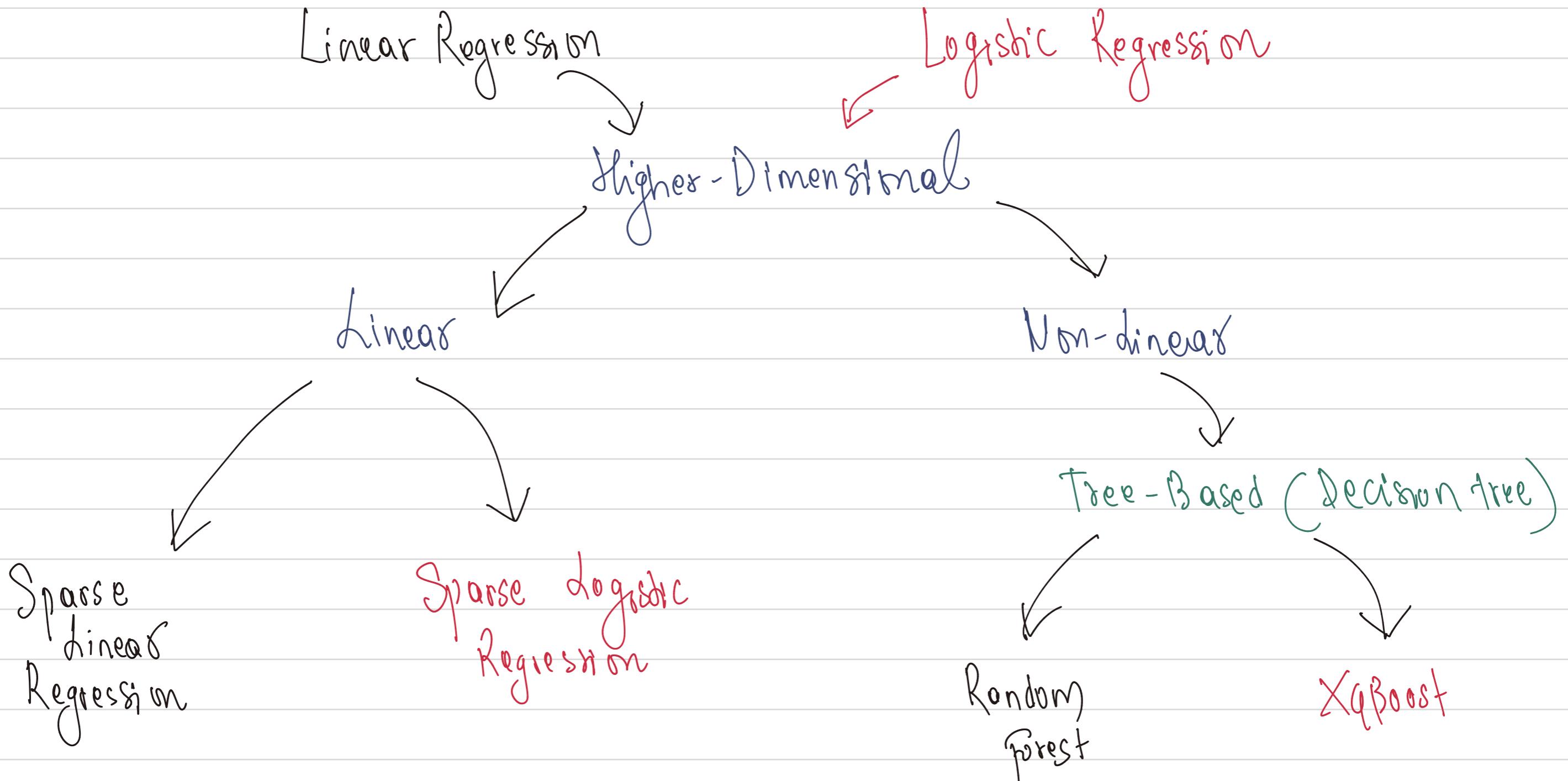


Fig. Hierarchical Tree showing the linear and Non-linear Model for Regression & classification

- Quick and less features use logistic and linear regression
- Above figure Black is Regression and Red is Classification
- Tree-Method Based Method are slower to train but handle complex decisions and can be difficult to interpret.

	Strengths	Weakness
Linear / Logistic Regression	<ul style="list-style-type: none"><li>- No tuning parameters</li><li>→ very fast to train</li></ul>	<ul style="list-style-type: none"><li>→ Unstable when many features to pick from</li><li>→ Can only fit linear curves</li></ul>

Sparse linear/  
logistic Regression

→ Stable when  
every/many features to  
pick up

→ Can only fit  
linear & curves

Tree-based  
Classification &  
Regression

→ Can fit nonlinear  
functions of inputs

→ Hard to interpret