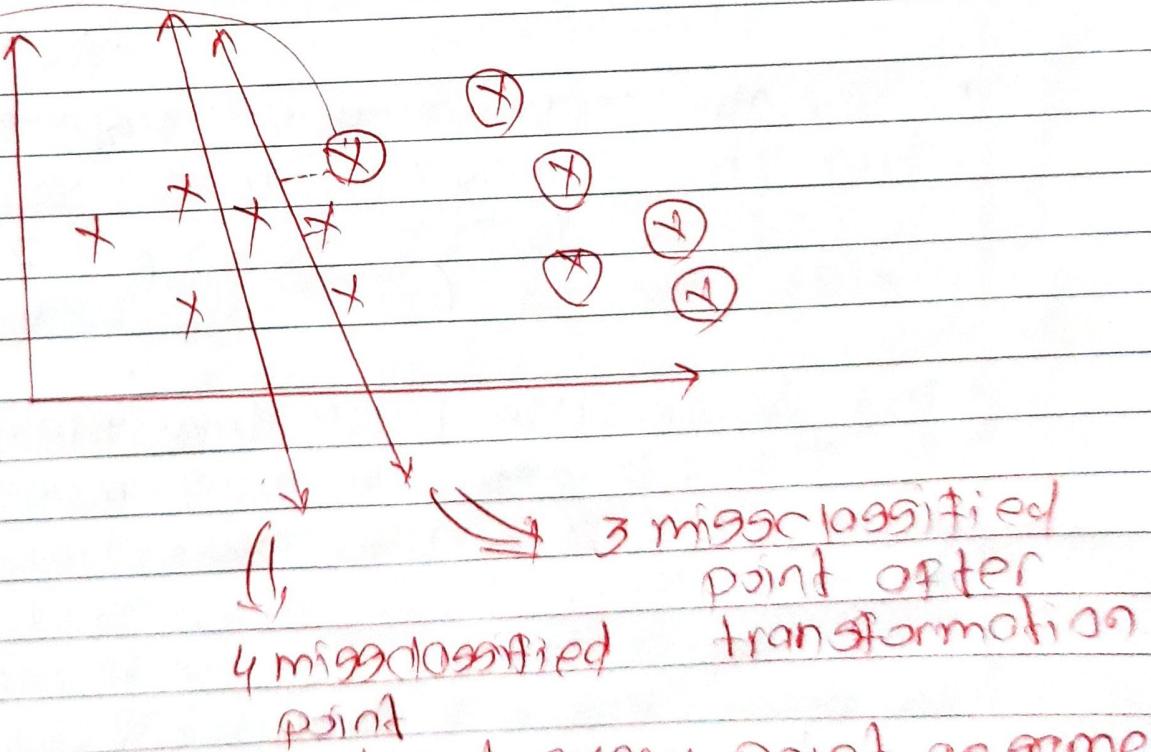


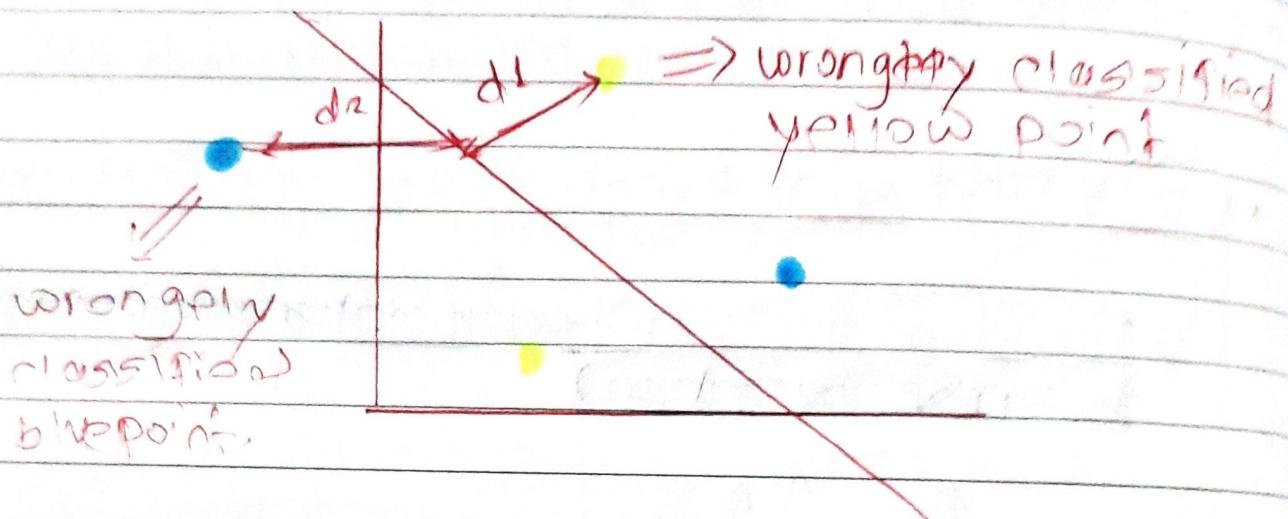
LOSS FUNCTION FOR PERCEPTIONS

- Loss function is function of $\omega_1, \omega_2, \omega_0$ and b .
 - We need such values of weights and bias that our error / loss becomes minimum.
 - Let's find error function for perception,
- i) No. of misclassified points - (could be this is error function).



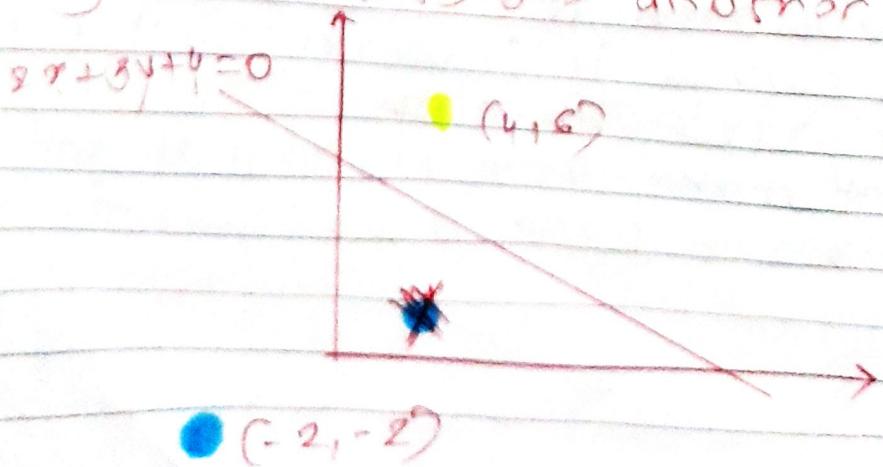
- Above approach treat every point as same, but it cannot tell which point make how many magnitude of error.
- This point make more error than other points. (basically this point misclassified many)

ii) Second approach could be calculate distance of points from our model line
perpendicular



- Calculate perpendicular distance of points from this line for missclassified points
 - $(d_1) \Rightarrow 10$
 - $(d_2) \Rightarrow 13$
- This loss function better than previous loss function, because we bring magnitude of error for each miss classified point.

Also, perceptron uses different loss function than distance of misclassification magnitude. Let's see another approach.



- Both of this 'blue' and yellow point are misclassified point.
- Rather than finding distance, we put point's coordinates in model line equation (i.e. perceptron equation)

$$\begin{cases} \text{Yellow } (4, 6) \Rightarrow 2 \times 4 + 3 \times 6 + 4 = 26 \\ \text{Blue } (-2, -2) \Rightarrow 2 \times -2 + 3 \times -2 + 4 = -6 \end{cases}$$

- Rather than distance we use this, Because calculating distance computationally expensive.
- This is just dot product.

To sklearn what is actual loss function for perceptron?

- Go to SGD page in documentation,

$$E(w, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \alpha R(w)$$

• This is general equation of SGD
error

• $\alpha R(w) \Rightarrow$ regularization parameter

Then we have loss function,

$$L(y_i, f(x_i)) = \max(0, 1 - y_i f(x_i))$$

- SGD Derivation

$$L(y_i, f(x_i)) = \max(0, 1 - y_i f(x_i))$$

- So according SGD general equation our loss function for perceptron becomes,

$$L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + 0$$

(No regularization for now)

$$\frac{\partial R(w_1, w_2)}{\partial} = 0$$

- What is $f(x_i)$?

$$f(x_i) = w_1 x_1 + w_2 x_2 + b = z$$

- So final loss function is,

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i f(x_i))$$

→ Average error of all data points/100%

n = No. of rows in data point.

Q What is L1?

- out, out in data. (Originally 1 or 0)

classmate

Date _____

Page _____

- This loss function depends on three things, (w_1, w_2, b) basically loss functions function of this (w_1, w_2, b)

- We have to make such changes in w_1, w_2 and b such that $L(w_1, w_2, b)$ value becomes minimum.

- Mathematically, We have to use gradient descent.

$$L = \underset{w_1, w_2, b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

Before that what is geometric intention of this loss function?

- This is our loss function,

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

where, $f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$

- We have this dataset,

	x_1	x_2	y
1	x_{11}	x_{12}	y_1
2	x_{21}	x_{22}	y_2
:			

throw

r

- What is $\max(0, -y_i f(x_i))$?

Let's assume $-y_i f(x_i) = x$,

$\max(0, -y_i f(x_i))$ becomes $\max(0, x)$

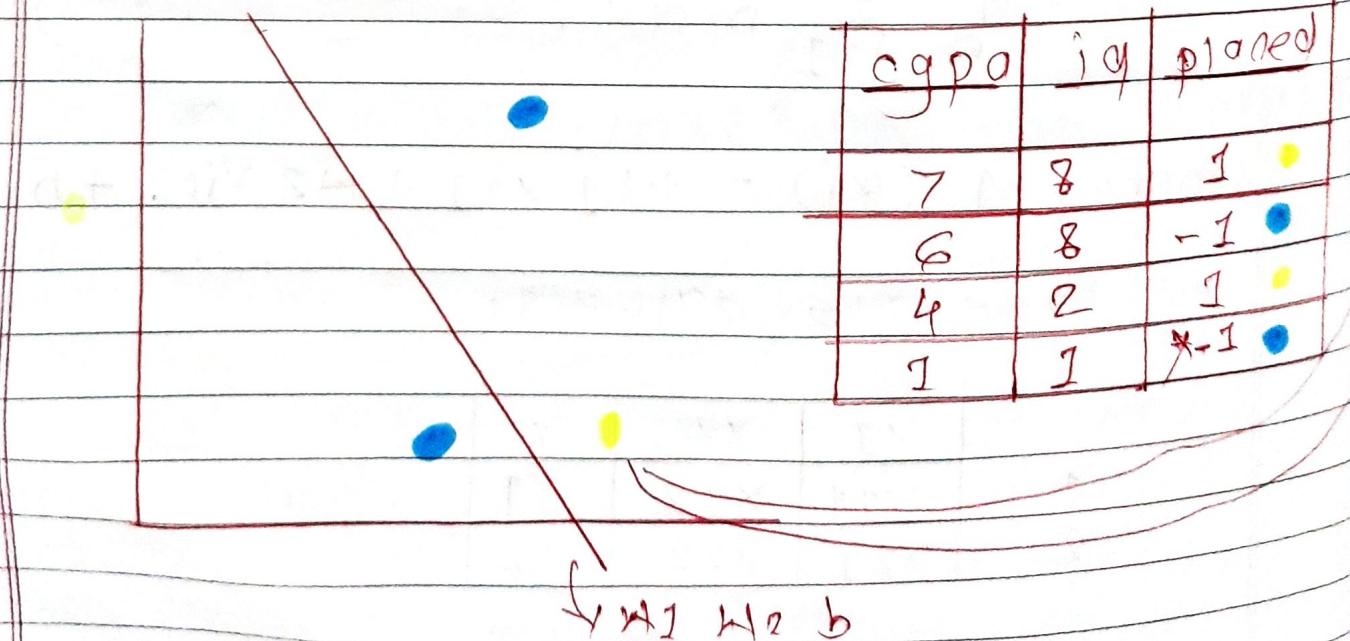
- Consider we have only two points, so break down this loss function.

$$L = \frac{1}{2} \left[\max(0, -y_1 f(x_1)) + \max(0, -y_2 f(x_2)) \right]$$

$$f(x_1) = w_1 x_1 + w_2 x_2 + b \quad (\text{first row})$$

$$f(x_2) = w_1 x_2 + w_2 x_2 + b \quad (\text{second row})$$

- Geometric intuition



- Line context of four point's separation,

y_i	\hat{y}_i
1	1
1	-1
-1	1
-1	-1

- We are going to find meaning of this point,
for this, $\max(0, -y_i \hat{f}(x_i))$
- i) condition one, $y_i=1$ and $\hat{y}_i=1$.
→ For this point, this point is on positive side of line and not classified after putting this point co-ordinate in line equation.

$w_1 x_1 + w_2 x_2 + b$ if will become positive.

$$w_1 x_1 + w_2 x_2 + b \geq 0 = \hat{f}(x_i)$$

$$\max(0, -\frac{1}{w_1}(-w_1 y_i))$$

$$\boxed{\hat{y}_i = 1}$$

$$\max(0, -\text{negative value})$$

max becomes 0.

- ii) condition 2, $y_i = -1$ and $\hat{y}_i = -1$

$$+\hat{y}_i = -1$$

$$w_1 x_1 + w_2 x_2 + b < 0$$

$-y_p x - y_p = +ve$, so,

$$\max(0, -y_i (-y_p \text{ value}))$$

$$\max(0, -(-1 (-y_p \text{ value})))$$

$$\max(0, -(+y_p \text{ value}))$$

~~$$\max(0, -)$$~~

so, max becomes zero, i.e., no error.

- iii) condition 3, $y_i = 1$ and $\hat{y}_i = -1$.
(misclassified point)

$$+\hat{y}_i = +1$$

$f(x)$ becomes negative.

$$-(+y_p x - (\text{negative } f(x)))$$

$$\max(0, +y_p)$$

+y_p value

If point is misclassified then it has some contribution in loss function.

CLASSESMATE

Date _____

Page _____

If point is correctly classified it will have zero contribution.

i) condition $y_i = -1 \quad g=1$

$$-1(-1 \times f(x) + y_i)$$

(misclassified point)

$$\max(0, +y_i)$$

+y_i value.

ii) predicted if used form 1 calculation and actual value put in loss equation y_i

$$\max(0, y_i (w_1 x_1 + w_2 x_2 + b))$$

ii) Gradient Descent

$$L = \underset{w_1 w_2 b}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

$$\text{where } f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$$

for i in epochs:

$$w_1 = w_1 + \eta \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 + \eta \frac{\partial L}{\partial w_2}$$

$$b = b + \eta \frac{\partial L}{\partial b}$$

- For w_1 ,

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \times \frac{\partial f(x_i)}{\partial w_1}$$

$$\frac{\partial L}{\partial f(x_i)} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial f(x_i)}{\partial w_1} = x_{i1}$$

~~combinations~~

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \times \frac{\partial f(x_i)}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i1} & \text{if } y_i f(x_i) < 0 \end{cases}$$

- Going by some logic, for w_2

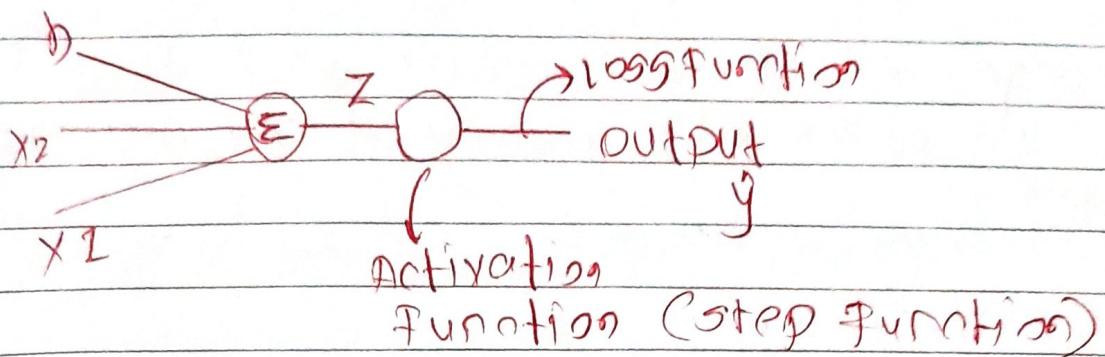
$$\frac{\partial L}{\partial w_2} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i2} & \text{if } y_i f(x_i) < 0 \end{cases}$$

- ALSO for, b

$$\frac{\partial L}{\partial b} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

- Refer notebook from video.

More Loss Function



- To get output in probabilities we can replace with sigmoid function as a activation function.

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

- Binary class entropy as loss function and sigmoid as activation function then our perceptron equals to logistic regression.

- For multiclass classification,

Activation function \Rightarrow softmax \Rightarrow this is softmax regression.

Categorical cross entropy (loss function) \Rightarrow

- For regression

Activation function \Rightarrow Linear \rightarrow Non-activating function
Loss function \Rightarrow mse \rightarrow loss function

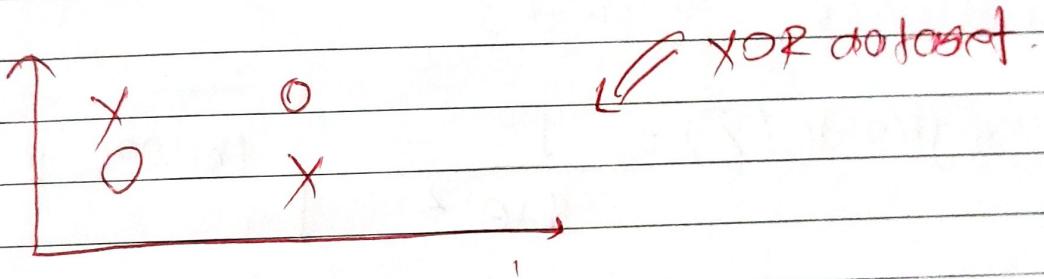
- Change activation function and loss function we get parameters every time by using SGD.

- In summary,

LOSS FUNCTION	ACTIVATION	OUTPUT
Hinge loss	Step	perceptron-binary classifier
Log-loss (Binary cross entropy)	Sigmoid	Logistic regression \rightarrow binary classifier
Categorical cross entropy	softmax	softmax regression (Multiclass classifier) output is probability
MSE	Linear	Linear regression

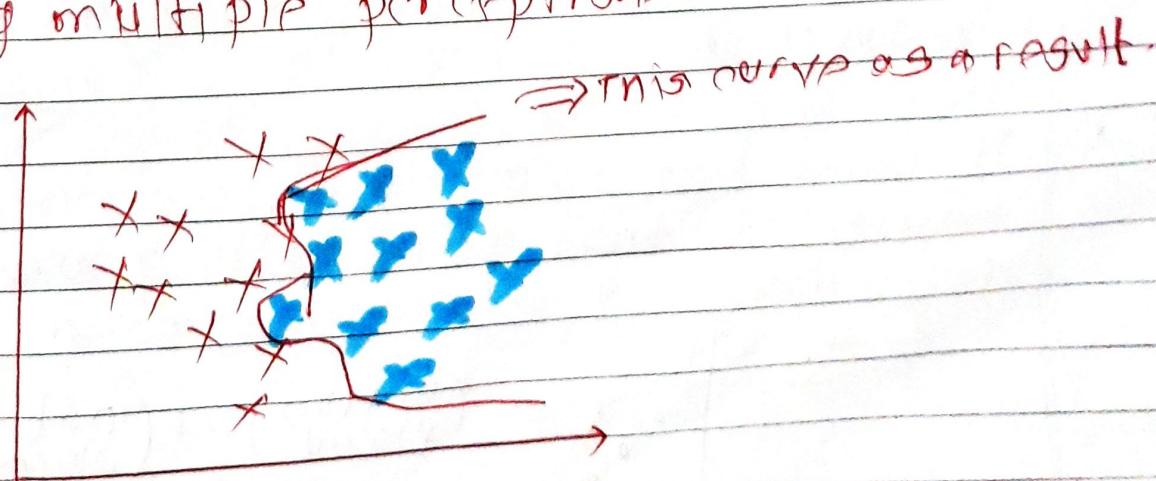
Problem with perceptron

- It not work with ^{non-} linear data.
- perceptron work with OR dataset, AND dataset but not work XOR dataset (it cannot even able create boundary/decision tree between point).
- You can check this with tensorflow playground tool.



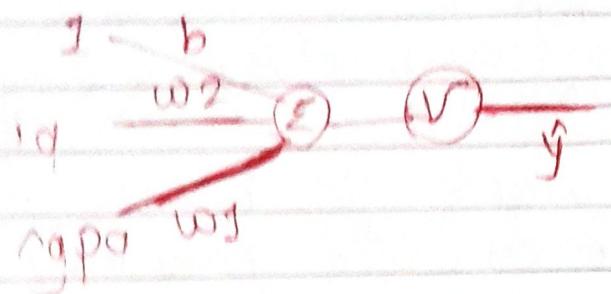
The problem

- We want to make algorithm that going to ~~make~~ capture non-linearity in data by using multiple perceptron.



- Before we have to know perceptron with sigmoid.

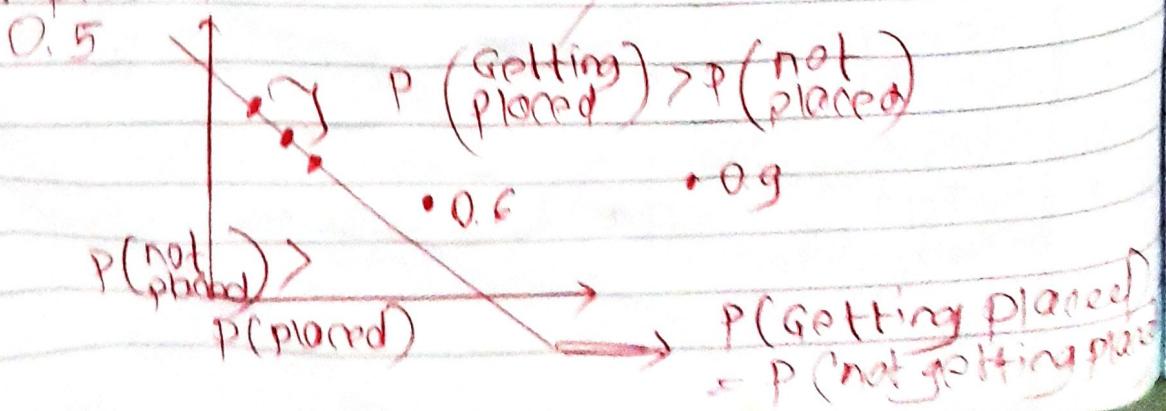
* Perceptron with Sigmoid.



- And loss function is log-loss (logistic regression)
- Sigmoid activation function gives probability between 0 and 1

$$\text{sigmoid}(z) = \frac{1}{1+e^{-z}} \quad \} \text{ Answer in } 0-1.$$

- The equation is after calculating weights and bias
 $w_1 \text{ cgpa} + w_2 \text{ iq} + b = \theta \cdot z.$
- And we put this 'z' value in sigmoid activation function.
- If point lies exactly on this line, then probability of that point getting placed and not placed is exactly similar.



जैसी ऐसी हम line की दूर अल्पों की तरफ साझी probabilities होती है, जिनमें $P(\text{Getting placed}) > P(\text{not getting placed})$ region and opposite side probability of getting placed reduce होती है.