

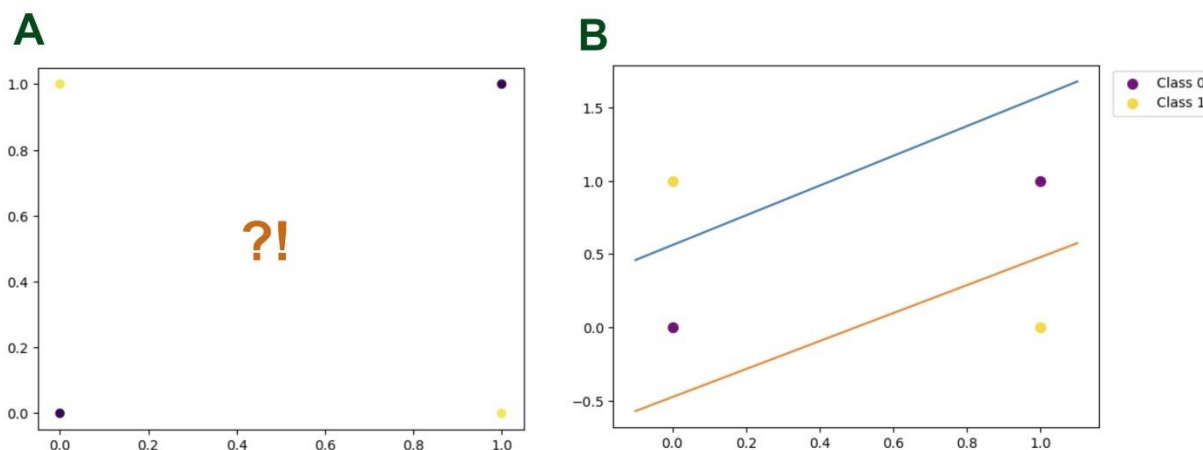
Deep Learning

Homework-1

تمرین 1:

الف) نشان دهید تابع XOR بدون لایه پنهان پرسپترون قابل شبیه سازی نیست.

Marvin Lee Minsky، بنیانگذار آزمایشگاه هوش مصنوعی MIT، در سال 1969 تردید خود را در قالب تحلیل دقیق محدودیت های Perceptron اولیه، منتشر کرد. او نتیجه گرفت که رویکرد Perceptron در هوش مصنوعی یک بن بست است. مهمترین عنصر مورد بحث در این تحلیل، روشن کردن حدود یک Perceptron بدون لایه پنهان است. برای مثال، آنها نمیتوانند یاد بگیرند تا تابع ساده boolean function XOR را پیشبینی کنند، زیرا عملگرهای XOR به صورت خطی جدایی پذیر نیستند و Perceptron منفرد (بدون لایه های پنهان و تابع فعال سازی) تنها میتواند داده های linearly separable را تفکیک کنند. به این ترتیب با اجرای Perceptron بدون لایه پنهان بر روی عملگر XOR (قسمت A تصویر زیر) تابع در یک loop تا بینهایت بار می گردد و هرگز به جواب نمیرسد و جهت شبیه سازی تابع XOR، مطابق قسمت B تصویر، به دو خط برای جداسازی 2 کلاس 0 و 1 آن نیاز خواهد بود که با اعمال شبکه های عصبی پیچیده تر غیر خطی مثل Perceptron چند لایه (مثلا لایه ای با تابع فعال سازی غیرخطی sigmoid) قابل اجرا است.

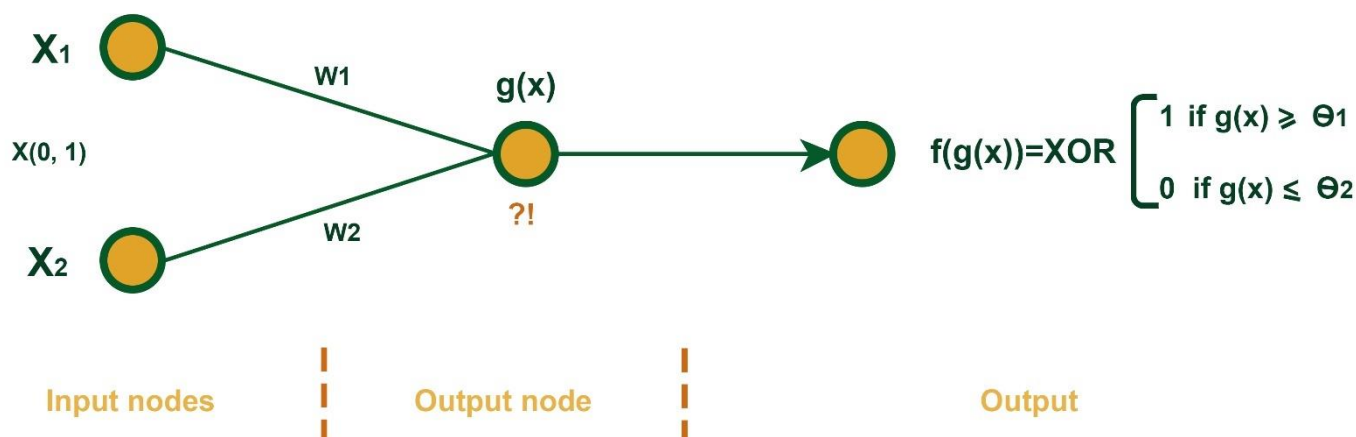


به این ترتیب برای پیدا کردن θ در ضابطه Perceptron برای حل تابع XOR نیاز است تا مطابق جدول رسم شده در قسمت ب سوال، $g(x)$ ی را پیدا کنیم که بتواند بر روی ترکیب خطی X_1 و X_2 با وزن یک اعمال شود و نتیجه $X_1 \text{ XOR } X_2$ را دهد، ولی با استفاده از روش SLP، هیچ θ ی را نمیتوان در نظر گرفت تا این شبکه رو شبیه سازی کند و برای اینکار نیاز است تا $g(x)$ بزرگتر مساوی θ و کوچکتر مساوی θ را بدست بیاوریم که شامل دو خط مجزا می شود:

Single Layer Perceptron (SLP):

Θ = Threshold

$$y = f(g(x)) = \begin{cases} 1 & \text{if } g(x) \geq \Theta \\ 0 & \text{o.w.} \end{cases}$$

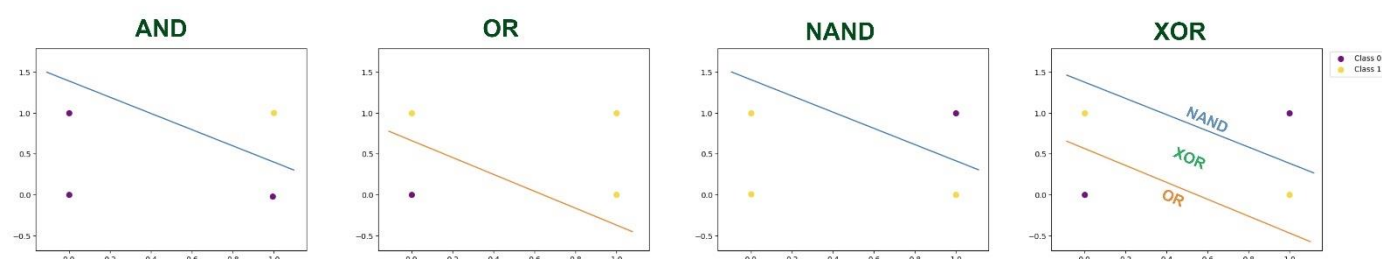


ب) XOR را با یک MLP با کمترین لایه شبیه سازی کنید.

XOR حاصل AND دو عملگر OR و NAND است که جدول حقیقت این boolean function ها به شرح زیر است:

X_1	X_2	$g(x)=X_1, X_2 \text{ Com.}$	$X_1 \text{ AND } X_2$	$X_1 \text{ OR } X_2$	$X_1 \text{ NAND } X_2$	$X_1 \text{ XOR } X_2$
T	T	2	1	1	0	0
T	F	1	0	1	1	1
F	T	1	1	1	1	1
F	F	0	0	0	1	0
T=1 F=0						

این عملگرها در فضای دو بعدی به صورت زیر نمایش داده می شوند که همانطور که مشخص است عملگرهای AND، OR و NAND به صورت خطی جداپذیر هستند ولی همانطور که در قسمت قبل توضیح داده شد XOR linearly separable نیست.

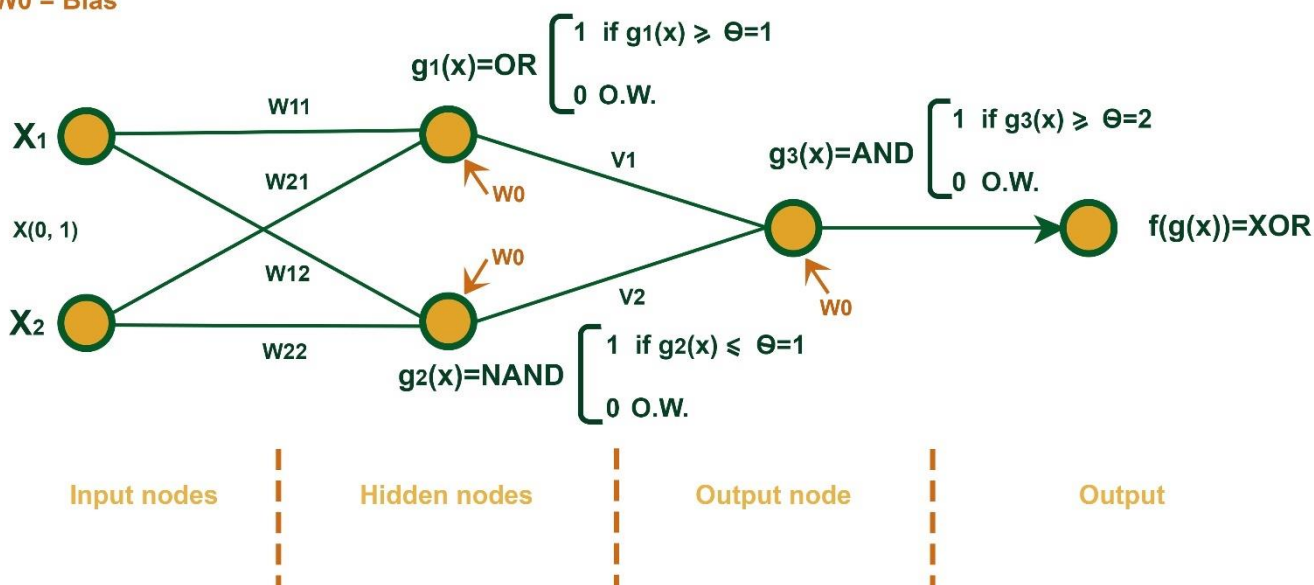


برای شبیه سازی تابع XOR می توان از شبکه رسم شده در زیر استفاده نمود که حاصل ترکیب دو عملگر OR و NAND بصورت AND است که خروجی XOR را ایجاد می کند:

Multi Layer Perceptron (MLP)

Θ = Threshold

W_0 = Bias



در این شبکه $g_1(x)$ ، $g_2(x)$ و $g_3(x)$ توابع فعال ساز یا Active function های ما هستند که به ترتیب نمایانگر توابع OR، NAND و AND که به عنوان hidden layer شبکه MLP ما بکار می روند. این لایه ها برای عبور داده ها مجموع Weight ها (W) و Bias ها (W_0) را در هر پرسپترون ایجاد می کنند. با در نظر گرفتن θ صحیح شبکه تابع ما را محاسبه خواهد کرد.

در بخش پیاده سازی تمرین، از **Active function** سیگموئید به عنوان **hidden layer** استفاده گردید. برای آموزش مدل، نمونه های آموزشی زیادی به شبکه منتقل می شود. شبکه با **Weight** ها و **Bias** های تصادفی راه اندازی می شود. پس از هر آموزشی که از طریق شبکه عبور داده می شود، مقداری را خروجی می دهد که پیش بینی آن برای خروجی (y) با توجه به ورودی (x) است. این پیش بینی با مقدار واقعی 1 یا 0 مقایسه می شود. مقایسه بین خروجی پیش بینی شده و واقعی در تابع هزینه انجام می شود. صرفاً جهت مقایسه و تمرین در پیاده سازی از دو تابع هزینه **MSE** و **BCE** به طور جداگانه برای این کار استفاده گردید. پس از محاسبه **Cost / Loss function**، **Weight** ها و **Bias** ها با استفاده از گرادینت (SGD) تابع هزینه به روز می شوند تا به **local** یا **global minimum** نزدیک تر شوند. هدف از بهینه سازی به حداقل رساندن تابع هزینه است. وقتی تابع هزینه صفر باشد، به این معنی است که تمام مقادیر پیش بینی شده با مقادیر واقعی یکسان هستند و هر چه مقدار آن به صفر نزدیکتر باشد به این معنی است که شبکه ما بهتر عمل کرده است.

کد مرتبط با تمرین 1:

• DL-HW1-Q1-Mahdavi.ipynb

تمرین 2:

روش کلاس بندی لاجستیک رگرسیون برای حالت چندکلاسه را با یک شبکه نمایش دهید و جزئیات آنرا شرح دهید.

لاجستیک رگرسیون چند کلاسه به صورت زیر شرح داده می شود:

$$\max \text{loss}(y, \hat{y}) = \prod_j \hat{y}_j^{y_j} = \sum_j y_j \ln \hat{y}_j$$

$$\max \sum_i \text{loss}(y_i, \hat{y}_i) = \sum_i \sum_j y_{ij} \ln \hat{y}_{ij}$$

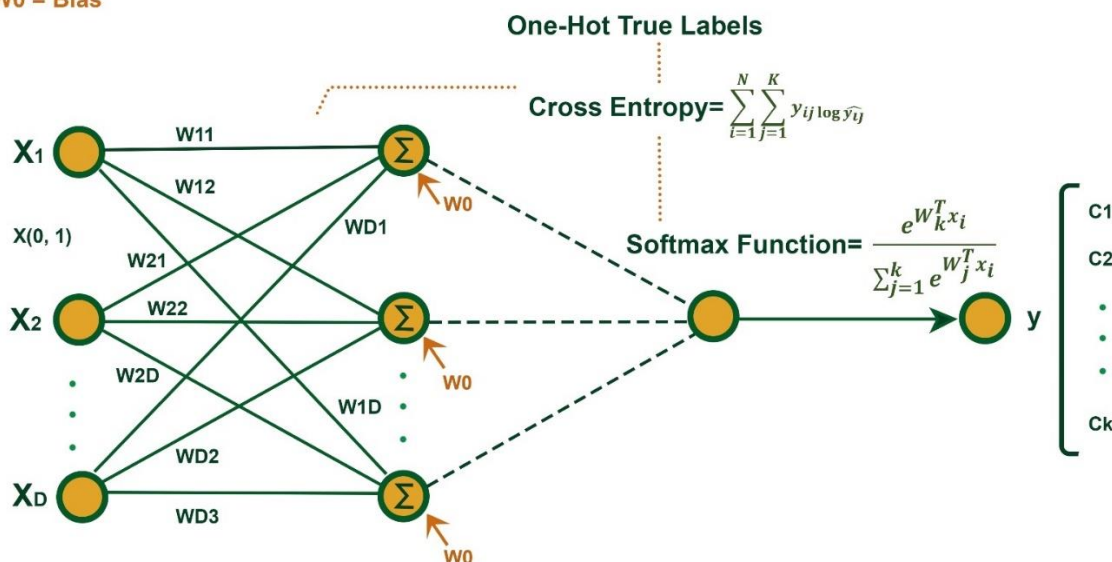
که معادل **Cross Entropy** است:

$$\text{Cross Entropy} = \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log \hat{y}_{ij}$$

اگر لاجستیک رگرسیون دو کلاسه داشتیم میتوانستیم از تابع فعال سازی **sigmoid** استفاده کنیم ولی **sigmoid** برای لاجستیک رگرسیون چند کلاسه جواب خوبی نمیدهد به این ترتیب برای این چنین مسائل از تابع فعال سازی **Softmax** استفاده میگردد که در آن مقادیر تابع هدف با کمک **One-Hot True Labels** به شکل برداری در می آیند و با یک **loss function Cross Entropy** آموزش مبینند و یک نوع رابطه غیر خطی از لاجستیک رگرسیون چند جمله ای را ارائه میدهند:

Softmax Regression (SMR)

W0 = Bias



Softmax regression (یا لاجستیک رگرسیون چند جمله‌ای) تعمیم لاجستیک رگرسیون به مواردی است که می‌خواهیم k کلاس را در حالت چند بعدی مدیریت کنیم و از آن به عنوان آخرین تابع فعال‌سازی (Activation Function) برای شبکه عصبی مان جهت نرمال سازی خروجی شبکه و تبدیل آن به توزیع احتمال بهره می‌بریم. نرمال سازی در این حالت نسبت به کلاس‌های خروجی پیش بینی شده، صورت می‌گیرد. پس قبل از استفاده از تابع Softmax، بعضی از اجزای برداری ممکن است منفی یا بیشتر از یک باشند، اما بعد از استفاده از این تابع، هر مولفه در بازه صفر و یک قرار می‌گیرد بطوری که مجموع آن‌ها برابر با یک باشد.

این تابع با یک تحلیل آماری و به صورت زیر بدست می‌آید:

طبق قضیه Bayes برای x_i اگر عضو کلاس c_i باشد، شانس این کلاس به شرط رخداد آن به این صورت است:

$$P(c_i | x_i) = \frac{P(c_i).P(x_i | c_i)}{P(x_i)} = \frac{P(c_i).P(x_i | c_i)}{\sum_{i=1}^k P(c_i).P(x_i | c_i)}$$

برای این که مولفه‌ها در بازه صفر و یک قرار گیرند خواهیم داشت:

$$P(c_i | x_i) = \frac{e^{\ln P(c_i).P(x_i | c_i)}}{\sum_{i=1}^k e^{\ln P(c_i).P(x_i | c_i)}}$$

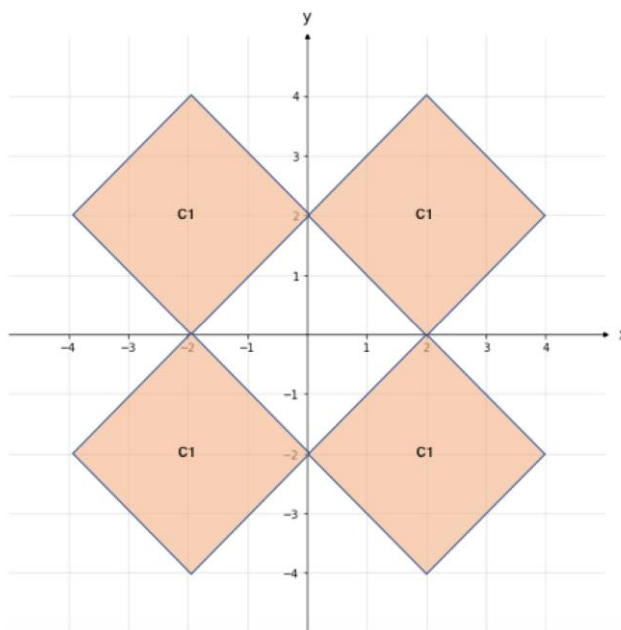
با تقریب زدن $\ln P(c_i).P(x_i | c_i)$ با $W_k^T x_i$ تابع Softmax به صورت زیر تعریف می‌شود:

$$\text{softmax}(x)_i = \frac{e^{W_k^T x_i}}{\sum_{j=1}^k e^{W_j^T x_i}}$$

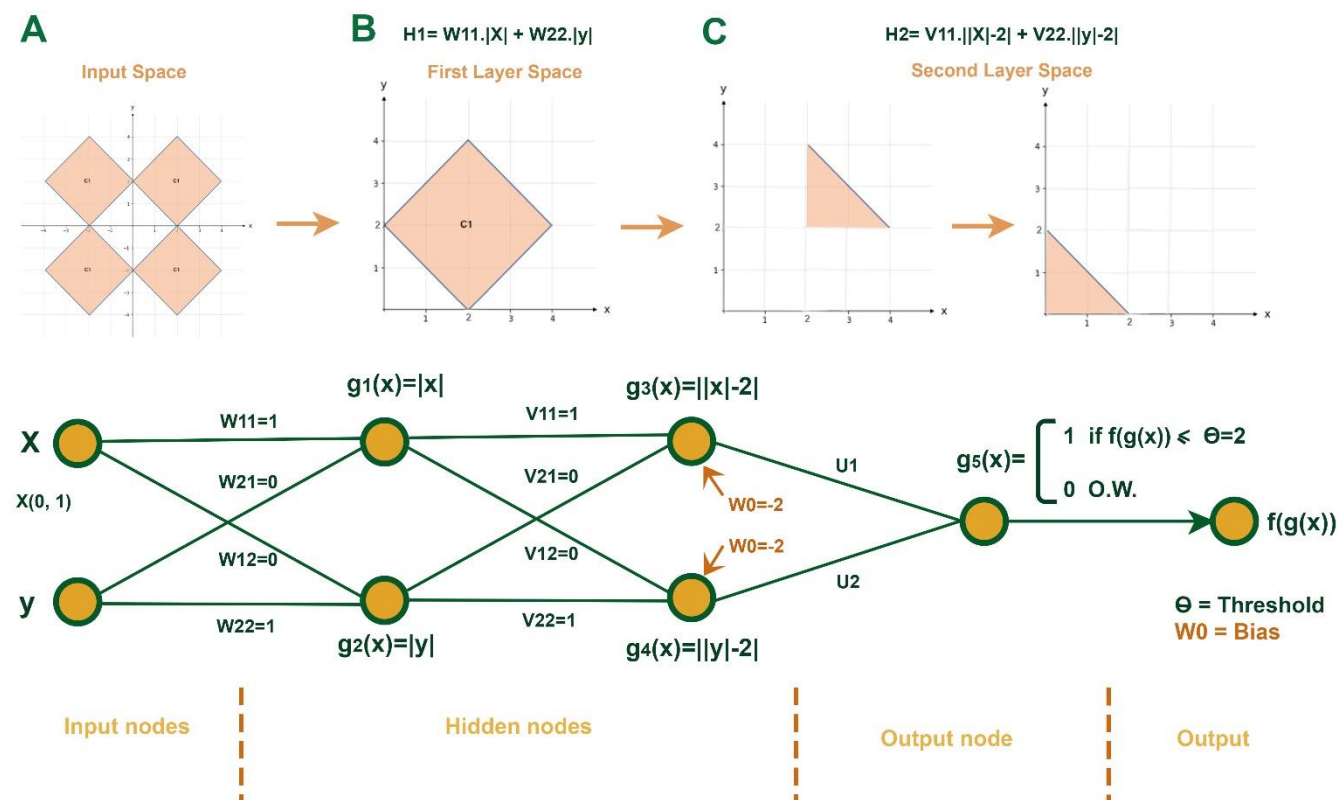
ایده پشت این روش به این صورت است که این مدل با دریافت یک نمونه، اول برای همه کلاس‌ها یک امتیاز رو محاسبه می‌کند و بعد با دادن این امتیازها به تابع Softmax احتمال بودن در هر کلاس رو حدس می‌زند. بعد از محاسبه امتیاز هر کلاس برای نمونه X_i ، میتوان احتمال تعلق X_i به کلاس k رو به دست آورد. پس این تابع اول \exp همه امتیازها رو به دست میاره بعد با تقسیم این مقدار به جمع همه \exp ها، آن را نرمال می‌کند.

تمرین 3:

در شکل زیر کلاس $C1$ مشخص شده است. خارج از آن کلاس $C2$ است. با یک شبکه با کمترین لایه این دو کلاس را جدا کنید.



Absolute Value Active Function



برای جدا کردن دو کلاس $C1$ و $C2$ باید رویکردی را در نظر گرفت تا این تصویر را به تصویر ساده تری تبدیل کرد تا بتوان آن را بصورت خطی جداپذیر کرد. برای این منظور در ابتدا و در لایه پنهان اول، صفحه یک بار حول محور x و بار دیگر حول محور y تا میخورد تا نتیجتاً یک ربع از تصویر که تصویر یک لوزی-مربع است بدست آید (قسمت B). این کار با استفاده از Activation Function قدر مطلق $|x|$ و $|y|$ انجام میگردد که جزئیات نودها و وزن های ورودی در شکل نمایش داده شده است. در مرحله بعد و لایه پنهان دوم دوباره تصویر حول محورهای فرضی x و y مجدداً با استفاده از توابع فعال سازی قدر مطلق $||x|-2|$ و $||y|-2|$ تا میخورد تا تصویر مثلت نمودار اول بخش C بدست آید و در نهایت با یک شیفت به مرکز محور x و y از طریق Bias منفی دو به فضای دوبعدی ای خواهیم رسید که به راحتی با یک خط می توان دو کلاس $C1$ (ناحیه رنگی) و $C2$ (خارج از ناحیه رنگی) را کلاسه بندی و از هم تفکیک کرد.