

$$\begin{aligned}
 1. \quad \text{Var}(\hat{P}_n(r)) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \frac{1}{\varphi_n} \varphi\left(\frac{y-r_i}{h}\right)\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}\left(\frac{1}{\varphi_n} \varphi\left(\frac{y-r_i}{h}\right)\right) \stackrel{\text{i.i.d. } r_i}{=} \frac{1}{n} \text{Var}\left(\frac{1}{\varphi_n} \varphi\left(\frac{y-r_1}{h}\right)\right) \\
 &= \frac{1}{n} E\left[\frac{1}{\varphi_n^2} \varphi^2\left(\frac{y-r_1}{h}\right)\right] - \frac{1}{n} E^2\left[\frac{1}{\varphi_n} \varphi\left(\frac{y-r_1}{h}\right)\right] \\
 &= \frac{1}{n} \int \frac{1}{\varphi_n^2} \varphi^2\left(\frac{y-r_1}{h}\right) P(r_1) dr_1 - \frac{1}{n} E^2[\hat{P}(r)] \\
 &\leq \frac{1}{n} \frac{\text{Sup}(\varphi)}{\varphi_n} \underbrace{\int \frac{1}{\varphi_n} \varphi\left(\frac{y-r_1}{h}\right) P(r_1) dr_1}_{E[\hat{P}(r)]} \leq 0
 \end{aligned}$$

$$\rightarrow \text{Var}[\hat{P}(r)] \leq \frac{\text{Sup}(\varphi) E[\hat{P}(r)]}{n \varphi_n}$$

2.

$$\varphi(r) = \begin{cases} e^{-x} & x > 0 \\ 0 & x \leq 0 \end{cases}, \quad \varphi_n = \int \varphi\left(\frac{x}{h_n}\right) dx = \int_0^\infty e^{-\frac{x}{h_n}} dx = -h_n e^{-\frac{x}{h_n}} \Big|_0^\infty = h_n$$

$$P_n(r) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\varphi_n} \varphi\left(\frac{x-x_i}{h_n}\right) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_n} \varphi\left(\frac{x-x_i}{h_n}\right)$$

$$\begin{aligned}
 \rightarrow \bar{P}_n(r) &= E[P_n(r)] = E\left[\frac{1}{n} \sum_{i=1}^n \frac{1}{h_n} \varphi\left(\frac{x-x_i}{h_n}\right)\right] = \frac{1}{n} \sum_{i=1}^n E\left[\frac{1}{h_n} \varphi\left(\frac{x-x_i}{h_n}\right)\right] \\
 &\stackrel{\text{i.i.d. } x_i}{=} E\left[\frac{1}{h_n} \varphi\left(\frac{x-x_1}{h_n}\right)\right] \\
 &= \int \frac{1}{h_n} e^{-\frac{x-x_1}{h_n}} P(x_1) dx_1
 \end{aligned}$$

$$x \leq 0 \rightarrow \varphi\left(\frac{x-x_1}{h_n}\right) = 0; \quad 1 \leq i \leq n \rightarrow E[P_n(x)] = 0$$

$$0 < x < a \rightarrow \varphi\left(\frac{x-x_1}{h_n}\right) = \begin{cases} e^{-\frac{x-x_1}{h_n}} & x < x_i \\ 0 & x > x_i \end{cases} \rightarrow E[P_n(x)] = \int_0^x \frac{1}{h_n} e^{-\frac{x-x_1}{h_n}} \frac{1}{a} dx_1 = \frac{1}{ah_n} e^{-\frac{x}{h_n}} \left(h_n e^{\frac{x_1}{h_n}} \Big|_0^x\right) = \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}}\right)$$

$$x \geq a \rightarrow \varphi\left(\frac{x-x_1}{h_n}\right) = e^{-\frac{x-x_1}{h_n}} \rightarrow E[P_n(x)] = \int_0^a \frac{1}{h_n} e^{-\frac{x-x_1}{h_n}} \frac{1}{a} dx_1 = \frac{1}{a} e^{-\frac{x}{h_n}} \left(e^{\frac{x_1}{h_n}} - 1\right)$$

2. bias =  $P(x) - E[\hat{P}_n(x)]$

$$\rightarrow \text{bias} = \frac{1}{a} - \frac{1}{a} \left(1 - e^{-\frac{x}{h_n}}\right) < 0.01$$

$$\rightarrow e^{-\frac{x}{h_n}} < \frac{0.01}{100} \xrightarrow{\ln} -\frac{x}{h_n} < \ln \frac{0.01}{100}$$

$$\rightarrow h_n < x \ln \frac{100}{0.01}$$

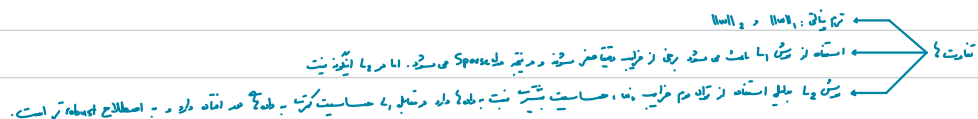
$$\begin{aligned} \rightarrow \hat{\beta} &= \arg \min_{\beta} \underbrace{\|Y - A\beta\|_2^2}_{J(\beta)} + \lambda \underbrace{\|\beta\|_2^2}_{\text{penalty}} = \\ &= \frac{(Y - A\beta)^T (Y - A\beta)}{Y^T Y - 2\beta^T A^T Y + \beta^T A^T A \beta} \\ &= Y^T Y - 2\beta^T A^T Y + \beta^T A^T A \beta \end{aligned}$$

$$\rightarrow \nabla_{\beta} J(\beta) = -2A^T Y + 2A^T A \beta + 2\lambda \beta = 0 \rightarrow (A^T A + \lambda I) \beta = A^T Y$$

$$\rightarrow \hat{\beta} = (A^T A + \lambda I)^{-1} A^T Y$$

2. در حالت کلی regularization روشی را منظور از overfitting می باشد. این روش با افزودن یک تابع خط (Loss Function) انجام می شود. در  $L_1$  Regularization که به نام LASSO هم شناخته

می‌شود، قدر مطلق ضرایب؛ به بیان نرم‌یافتی - Loss Function افزوده می‌شود. در  $L_2$  Regularization که به نام Ridge Regularization هم شناخته می‌شود، نرم  $L_2$  ضرایب به تابع خطا افزوده می‌شود.



W1 $\Delta$		W2 $\times$		W3 $\bullet$	
X1	X2	X1	X2	X1	X2
1	10	1	5	1	2
2	0	2	0	2	-5
3	5	3	5	3	10
	-2		5		-4

$$w_1 = \begin{bmatrix} 3 \\ -4 \end{bmatrix}$$

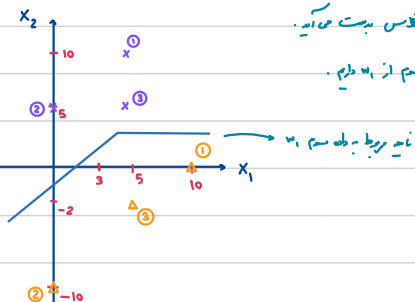
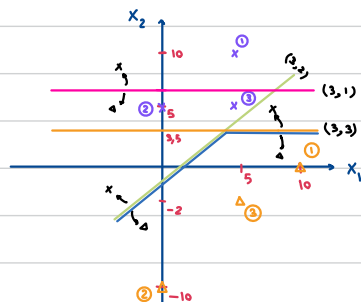
$$w_2 = \begin{bmatrix} 3.33 \\ 6.66 \end{bmatrix}$$

$$w_3 = \begin{bmatrix} 5.25 \\ 2 \end{bmatrix}$$

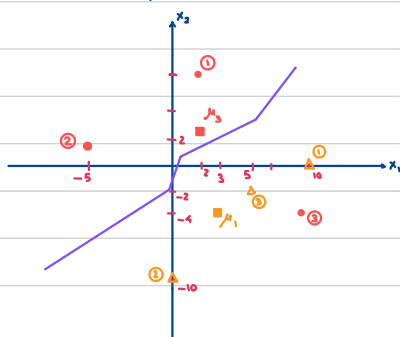
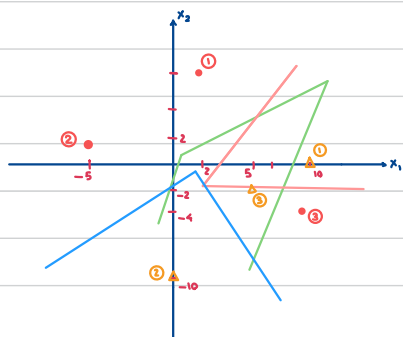
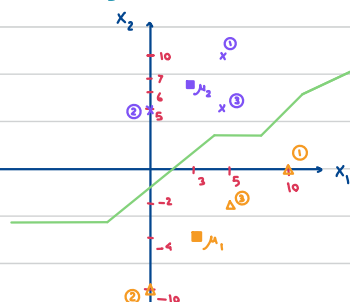
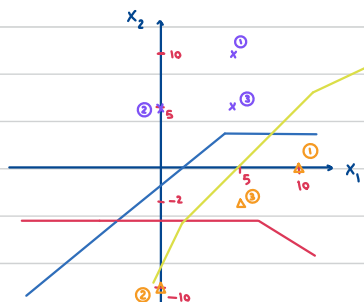
1. نسل اول حرکت از داده  $w_1$  کلاس  $w_1$  را با تمام داده  $w_2$  میل می کند و در نهایت  $w_1$  را رسم می کند و مشخص با اشتراک  $w_1$  و  $w_2$  به هر داده وابسته می کند. در نهایت با اجماع گرفتن می توانی.

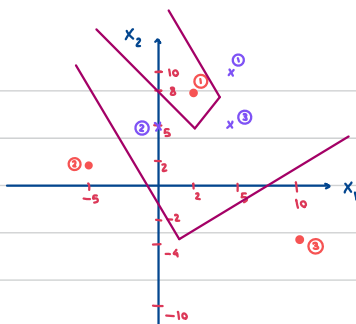
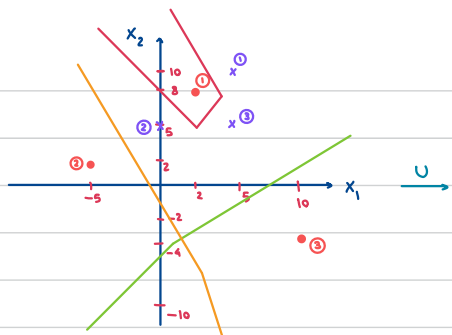
در نهایت در کلاس  $w_1$  به دست می آید.

تجربه داده  $w_1$  از  $w_2$  داریم.

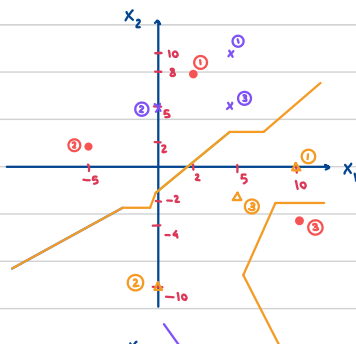
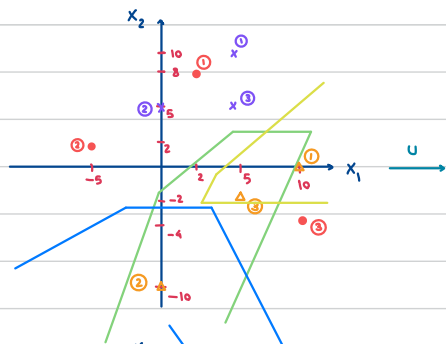


حال نسل دوم به داده  $w_2$  را میل می کند و در نهایت  $w_2$  را رسم می کند و مشخص با اشتراک  $w_1$  و  $w_2$  به هر داده وابسته می کند. در نهایت با اجماع گرفتن می توانی.

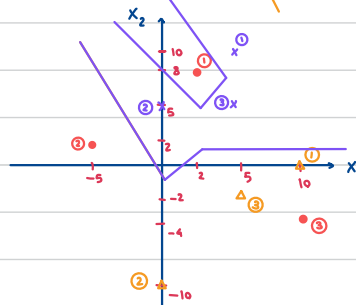
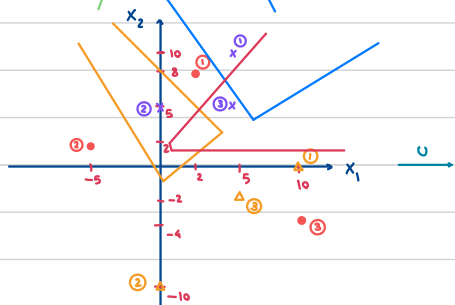




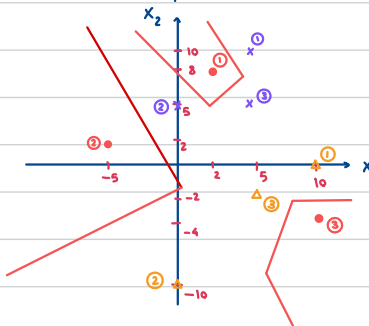
3.



4. - نام ربط به  $w_1$



نام ربط به  $w_2$



نام ربط به  $w_3$

(الف)

A graph of a bivariate probability density function  $P(x_1, x_2)$ . The horizontal axis is  $x_1$  and the vertical axis is  $x_2$ . The density is 1 for  $x_1 \in [0, 1]$  and  $x_2 \in [0, 1]$ , and 0 otherwise. A yellow shaded region is shown for  $x_1 \in [1/3, 2/3]$  and  $x_2 \in [0, 1]$ . The region is labeled  $R_1 \cup R_2$ .

$$P(c) = P(y=u_1, \hat{y}=u_2) + P(y=u_2, \hat{y}=u_1)$$

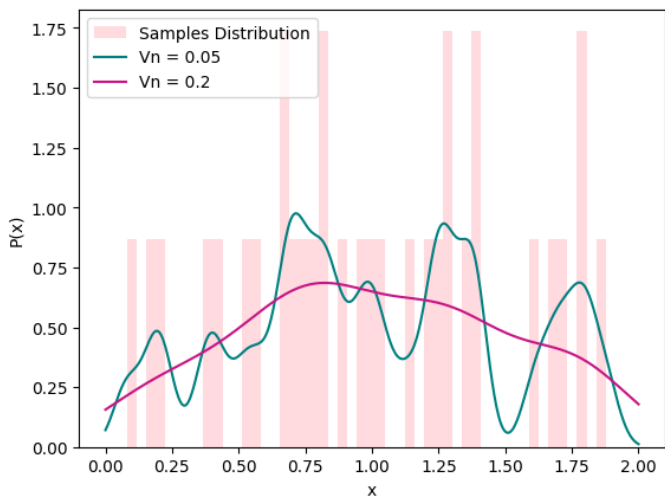
$$= \int_{R_2} P(x|w_1)P(w_1) dx + \int_{R_1} P(x|w_2)P(w_2) dx$$

از به  $\frac{2}{3}$  و  $\frac{1}{3}$  در  $\frac{1}{2}$  است.

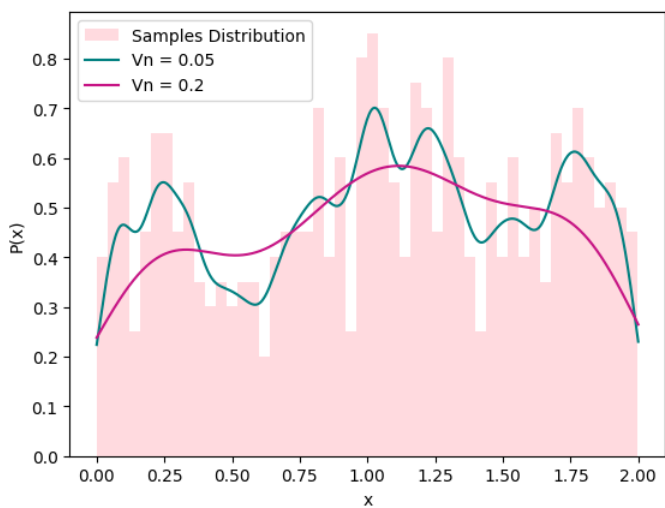
(ب) احتمال خلا برابر است با احتمال اینکه نقطه متوجه به ۳ قرار گیرد باشد، در نتیجه ما خود مشترک ۳ و ۳ داریم، احتمال خلا ما مشترک می باشد.

$$P(c_1) = \int_{\kappa_1} P(x|w_2) = \int_{\frac{1}{3}}^{\frac{2}{3}} \frac{3}{2} dx = \frac{1}{2}$$

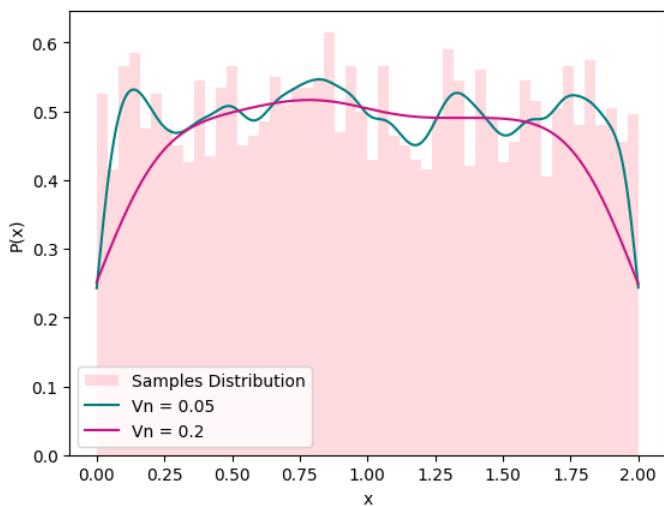
۶. نتایج را به شکل زیر خواهیم داشت:



شکل ۱: نتایج برای ۳۲ داده آموزشی



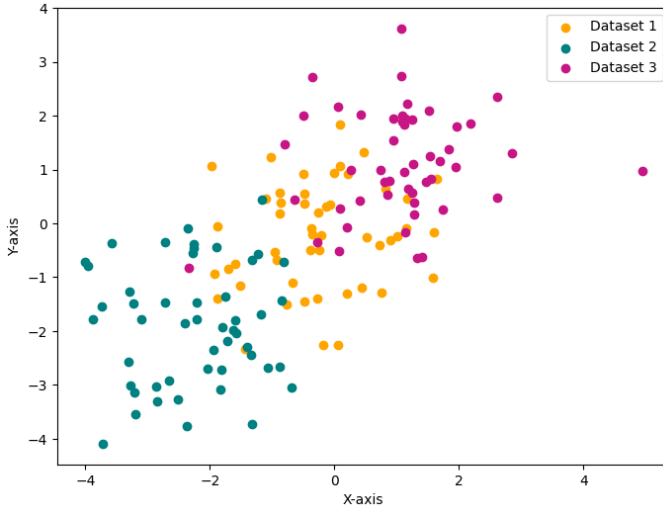
شکل ۲: نتایج برای ۲۵۶ داده آموزشی



شکل ۳. نتایج برای ۵۰۰۰ داده آموزشی

با توجه به نمودارها مشاهده می‌کنیم به ازای  $N$  یکسان برای  $V_n$  بیشتر تقریب نرم‌تری خواهیم داشت. در واقع به نسبت نویز کمتری داریم. همچنین با افزایش تعداد داده‌های آموزشی به ازای  $V_n$  یکسان، تقریب به نمودار واقعی نزدیک‌تر می‌شود.

۷. پراکندگی نمونه‌ها به شکل زیر خواهند بود:



شکل ۱: نمودار پراکندگی داده‌ها

۴ نقطه به شکل زیر در نظر می‌گیریم:

$$P_1 = \begin{pmatrix} -4 \\ -4 \end{pmatrix}, \quad P_2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \quad P_3 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \quad P_4 = \begin{pmatrix} 1.5 \\ 0 \end{pmatrix}$$

با توجه به نمودار پراکندگی داده‌ها انتظار داریم برچسب داده‌ها به ترتیب زیر باشد:

$$\hat{y}_1 = 2, \quad \hat{y}_2 = 1, \quad \hat{y}_3 = 3, \quad \hat{y}_4 = 3$$

برای بررسی دو حالت خواهیم داشت:

$$N = 50, \quad V_n = 1 \quad (\text{آ})$$

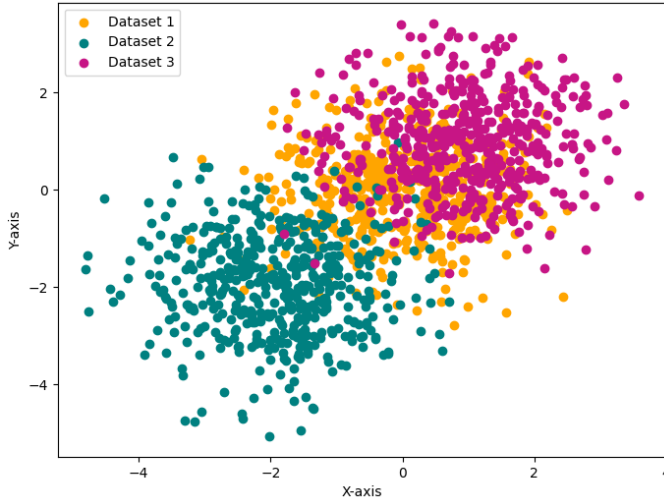
$$N = 50, \quad V_n = 0.1 \quad (\text{ب})$$

در حالت اول حجم همسایگی بیشتر است و تعداد بیشتری از نمونه‌ها داخل همسایگی قرار می‌گیرند. در نتیجه دقت تصمیم‌گیری در کل بالاتر است.

با تست نقاط، حالت (آ) در بیشتر مواقع پیش‌بینی درست است اما در حالت (ب) به دلیل کمتر بودن نمونه‌ها در یک همسایگی (به دلیل کاهش  $V_n$ ) جواب بسته به پراکندگی‌های داده‌های آموزشی متفاوت خواهد بود. به خصوص نقطه  $P_4$  که برچسب آن بین کلاس ۲ و ۳ تغییر می‌کند.



با افزایش تعداد داده‌های هر کلاس خواهیم داشت:



شکل ۲: نمودار پراکندگی داده‌ها برای  $N = 500$

مشابه قبل دو حالت خواهیم داشت:

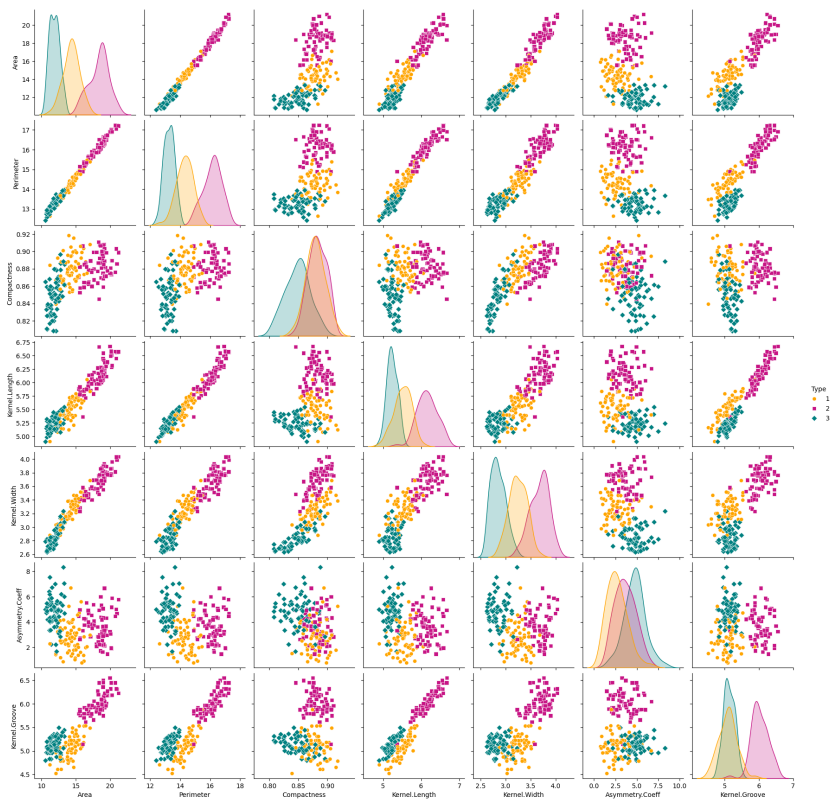
$$N = 500, \quad V_n = 1 \quad (\bar{A})$$

$$N = 500, \quad V_n = 0.1 \quad (B)$$

اینبار تعداد داده‌ها بیشتر شده و در یک حجم یکسان تعداد بیشتری نمونه داخل همسایگی قرار می‌گیرند.

برای حالت (آ) نتایج دقت بالاتری خواهند داشت. در حالت (ب) در مقایسه با حالت قبل برچسب نقطه  $P_4$  به درستی تعیین می‌شود اما در برخی تست‌ها برچسب نقطه  $P_2$  بین کلاس ۱ و ۲ تغییر می‌کند.

در مجموع هر قدر  $V_n, N$  بیشتر باشند دقت بالاتری خواهیم داشت.

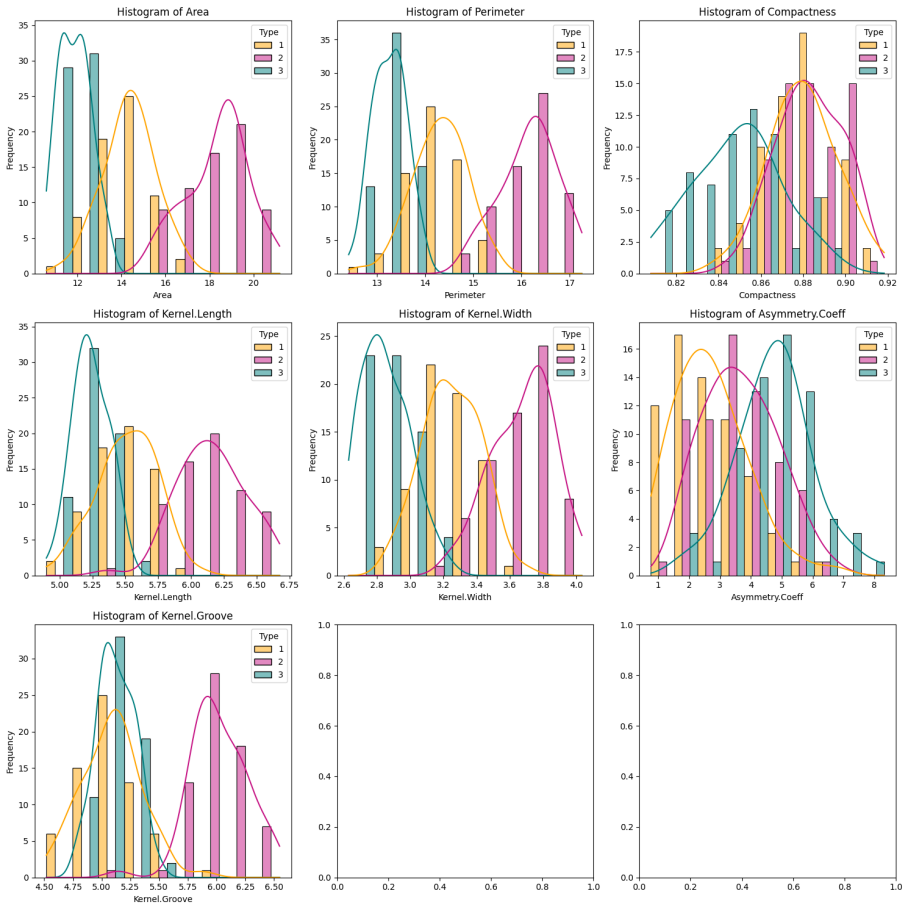


شکل ۱: نمودار ویژگی‌ها

با توجه به شکل ویژگی‌هایی که بهتر سه کلاس را جدا کردند:

$(Area/Kernel:Groove), (Perimeter/Kernel:Groove), (Kernel:Length/Kernel:Groove)$

حال هیستوگرام ویژگی‌ها را رسم می‌کنیم:



شکل ۲: نمودار هیستوگرام کلاس‌های هر ویژگی

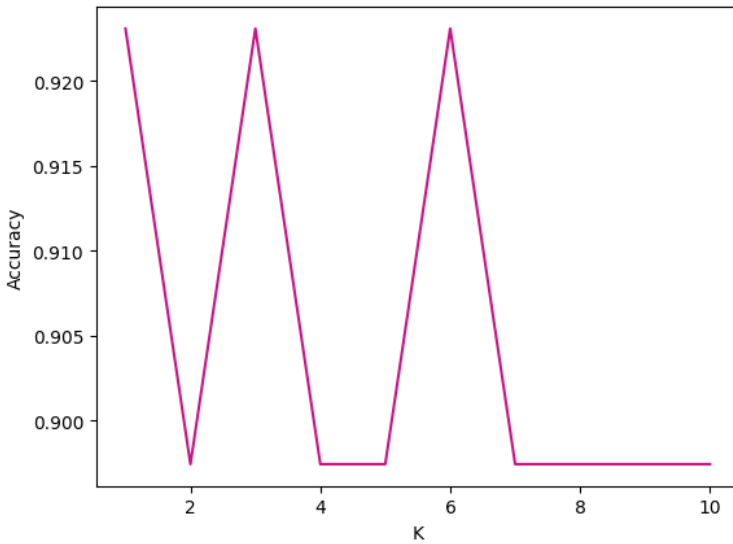
(ب) برای پیش‌پردازش داده‌ها از Z-SCORE استفاده می‌کنیم.

$$Z_{score} = \frac{x - \mu}{\sigma}$$

هر قدر این مقدار بیشتر باشد یعنی داده نویزی‌تر است و بهتر است حذف شود.  
برای نرمال سازی داده‌ها برای هر ویژگی رابطه زیر را اعمال می‌کنیم:

$$x_{Normalized} = \frac{x - x_{min}}{x_{max} - x}$$

(د) با پیاده‌سازی الگوریتم KNN به نتیجه زیر خواهیم رسید:



شکل ۳: نمودار دقت به ازای مقادیر مختلف  $k$

با توجه به نمودار بهترین مقدار  $k$  می‌تواند ۱، ۳ و ۶ باشد.