

به نام خدا

یادگیری ماشین تمرین سوم

محمد ناصری

۸۱۰۱۰۰۴۸۶

بهار ۱۴۰۱

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

۱- سوال اول (۲۰ نمره)

توزیع نرمال $p(x) \sim N(\mu, \sigma^2)$ و تابع پنجره پارزن $\varphi(x) \sim N(0, 1)$ را در نظر بگیرید. نشان دهید که تخمین پنجره پارزن

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

برای h_n های کوچک دارای ویژگی های زیر است:

- $\tilde{p}_n(x) \sim N(\mu, h_n^2 + \sigma^2)$
- $p_n(x) - \tilde{p}_n(x) \cong \frac{1}{2} \left(\frac{h_n}{\sigma}\right)^2 \left[1 - \left(\frac{x-\mu}{\sigma}\right)^2\right] p(x)$
- $\text{var}[p_n(x)] \cong \frac{1}{2nh_n\sqrt{\pi}} p(x)$

الف

$$\rightarrow \bar{p}_n(x) = E[p_n(x)] = \frac{1}{nh_n} \sum_{i=1}^n E\left[\varphi\left(\frac{x - x_i}{h_n}\right)\right]$$

$$= \frac{1}{h_n} \int_{-\infty}^{\infty} \varphi\left(\frac{x-u}{h_n}\right) p(u) du$$

$$= \frac{1}{h_n} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\left(\frac{x-u}{h_n}\right)^2}{2}\right) \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\left(\frac{u-\mu}{\sigma}\right)^2}{2}\right) du$$

$$= \frac{1}{2\pi h_n \sigma} \exp\left(-\frac{\left(\frac{x^2}{h_n^2} + \frac{\mu^2}{\sigma^2}\right)}{2}\right) \int_{-\infty}^{\infty} \exp\left(-\frac{u^2}{2} \left(\frac{1}{h_n^2} + \frac{1}{\sigma^2}\right) - 2u\left(\frac{x}{h_n^2} + \frac{\mu}{\sigma^2}\right)\right) du$$

$$\rightarrow \begin{cases} \theta^2 = \frac{\sigma^2 h_n^2}{\sigma^2 + h_n^2} \\ \alpha = \theta^2 \left(\frac{x}{h_n^2} + \frac{\mu}{\sigma^2}\right) \end{cases} \Rightarrow \bar{p}_n(x) = \frac{\sqrt{2\pi}\theta}{2\pi h_n \sigma} \exp\left[\frac{-1}{2} \left(\frac{x^2}{h_n^2} + \frac{\mu^2}{\sigma^2}\right) + \frac{1}{2} \frac{\alpha^2}{\theta^2}\right]$$

$$\rightarrow \bar{p}_n(x) = \frac{1}{\sqrt{2\pi} h_n \sigma} \times \frac{h_n \sigma}{\sqrt{h_n^2 + \sigma^2}} \exp\left[\frac{-1}{2} \left(\frac{x^2}{h_n^2} + \frac{\mu^2}{\sigma^2} - \frac{\alpha^2}{\theta^2}\right)\right]$$

$$= \frac{(x h_n)^2}{(h_n^2 + \sigma^2) h_n^2} + \frac{(\mu \sigma)^2}{(h_n^2 + \sigma^2) \sigma^2} - \frac{2x\mu}{h_n^2 + \sigma^2} = \frac{(x-\mu)^2}{h_n^2 + \sigma^2}$$

$$\bar{p}_n(x) \sim N\left(\mu, h_n^2 + \sigma^2\right)$$

$$\rightarrow \bar{p}_n(x) = \frac{1}{\sqrt{2\pi} \sqrt{h_n^2 + \sigma^2}} \exp\left[\frac{-1}{2} \frac{(x-\mu)^2}{h_n^2 + \sigma^2}\right]$$

۲- سوال دوم (۱۰ نمره)

متریک فاصله اقلیدسی را در d بعد در نظر بگیرید:

$$D(a, b) = \sqrt{\sum_{k=1}^d (a_k - b_k)^2}$$

فرض کنید عناصر هر بعد را در یک مقدار حقیقی غیر صفر ضرب میکنیم. یعنی $k = 1, 2, \dots, d$ داریم:

$$x_k = \alpha_k x_k$$

ثابت کنید پس از ضرب نیز این متریک فاصله همچنان یک فاصله استاندارد است یعنی ویژگی‌های گفته شده برای یک فاصله استاندارد را دارا می‌باشد. در مورد تاثیر این امر بر طبقه بند knn بحث کنید.

① $D(a, b) \geq 0$ ← جمع مربعات مثبت است ✓

② $D(a, b) = 0 \Leftrightarrow a = b$ ← زیر رادیکال بخواند و فرستد باید $a = b$ ✓

③ $D(a, b) = D(b, a)$ ← $(a_k - b_k)^2 = (b_k - a_k)^2$ ✓

④ $D(a, b) + D(b, c) \geq D(a, c)$ ←

$\text{diag}[\alpha_1, \dots, \alpha_d]$
 \uparrow
 $D(a, b) = \|A(a-b)\|_2$

$\begin{matrix} x = a - c \\ y = c - b \end{matrix} \rightarrow D(a, b) + D(b, c) \geq D(a, c) \Rightarrow \|Ax\|_2 + \|Ay\|_2 \geq \|A(x+y)\|_2$

$\begin{matrix} w = Ax \\ z = Ay \end{matrix} \rightarrow \|w\|_2 + \|z\|_2 \geq \|w+z\|_2 \xrightarrow{^2} \|w\|_2^2 + \|z\|_2^2 + 2\|w\|_2 \|z\|_2 \geq \|w+z\|_2^2$
 $\geq \|w\|_2^2 + \|z\|_2^2 + 2 \sum_{i=1}^d w_i z_i$
 $\Rightarrow \|w\|_2 \|z\|_2 \geq \sum_{i=1}^d w_i z_i \quad \checkmark$

فشار متریک knn همیشه یک نمونه اولیه ذخیره شده P با جرم α فاصله از n^* بعد خواصیات
 $\delta = \min_x D(P, x) : n \rightarrow \infty \Rightarrow D(P, x_n) \rightarrow \delta$

$D(P, x^*) = \delta \Leftrightarrow \begin{cases} D(P, x^*) \leq \delta \\ D(P, x^*) \geq \delta \end{cases} \leftarrow \begin{matrix} D(P, x_n) + D(x_n, x^*) \geq D(P, x^*) \Leftrightarrow x_n \rightarrow x^* \end{matrix}$

$\rightarrow P(w_i | P) \simeq P(w_i | x^*)$

۳- سوال سوم (۲۰ نمره)

یک مسئله طبقه بندی با روش knn را در نظر بگیرید. مجموعه داده دو کلاسه D را نیز به صورت $D = \{x^q, \omega_i^q\}$, $q=1, \dots, Q$ داریم. این داده ها نتایج یک نظر سنجی بوده و دیتاپوینت ها به صورت پرچسب خورده، مستقل از هم هستند و فرض می کنیم تعداد داده های دو کلاسه یکسان است. برای هر سَمپل تست نزدیک ترین k دیتاپوینت را به صورت $\{x_i\}_{i=1, \dots, k}$ نمایش می دهیم. هر دوی $p(x|1)$ و $p(x|2)$ توزیع یکنواخت بر روی یک کره به شعاع واحد دارند و مرکز دو ابر کره نیز از هم ۱۰ واحد فاصله دارند.

الف) نشان دهید اگر k فرد باشد متوسط احتمال خطا از رابطه زیر به دست می آید:

$$p_Q(e) = \frac{1}{2^Q} \sum_{j=0}^{\frac{k-1}{2}} \binom{Q}{j}$$

ب) با توجه به بخش قبل نشان دهید که در این حالت خطای طبقه بند نزدیک ترین همسایه کمتر از حالت $k \geq 2$ است و دلیل مشاهده این موضوع را توضیح دهید.

توزیع بلنواخت

$$\left. \begin{array}{l} \left\{ \begin{array}{l} |x_1| < 1 \\ 0 < \omega \end{array} \right\} \frac{1}{2} P(\omega_1|x) \\ \left\{ \begin{array}{l} |x_2| < 1 \\ 0 < \omega \end{array} \right\} \frac{1}{2} P(\omega_2|x) \end{array} \right\}$$

ج) نشان دهید: $\lim_{Q \rightarrow \infty} p_Q(e) = 0$

الف

$$\begin{aligned} P_Q(e) &= \Pr[e \in \omega_1, \omega_2 \text{ is frequent}] + \Pr[e \in \omega_2, \omega_1 \text{ is frequent}] \\ &= 2 \Pr[e \in \omega_1, \omega_2 \text{ is frequent}] = 2 P(\omega_1) \Pr(\text{count}(\omega_2) < \frac{k-1}{2}) \\ &= \sum_{j=0}^{(k-1)/2} \binom{Q}{j} \frac{1}{2^Q} \times \frac{1}{2^{(Q-j)}} = \frac{1}{2^Q} \sum_{j=0}^{(k-1)/2} \binom{Q}{j} \end{aligned}$$

ب

$$k=1 \rightarrow P_n(e) > \frac{1}{2^Q} \rightarrow \frac{1}{2^Q} \sum_{j=0}^{(k-1)/2} \binom{Q}{j} \rightarrow k > 1$$

ج

$$\begin{aligned} P_n(e) &= \Pr[\text{Bin}(n, \frac{1}{2}) \leq (k-1)/2] = \Pr[Y_1 + \dots + Y_n \leq \frac{k-1}{2}] \\ Q \rightarrow \infty \Rightarrow P_n \rightarrow 0 &\leftarrow \frac{k-1}{2} < \frac{\frac{\alpha}{\sqrt{Q}} - 1}{2} < 0 \leftarrow \end{aligned}$$

(الف)

در ابتدا برای این سوال لازم است که مفهوم Bias variance trade off در یادگیری ماشین بیان گردد. در ابتدا مفهوم بایاس یا سوگیری بیان میشود که این مفهوم ناشی از فرضیات غلط و اشتباه در یادگیری است به عبارتی یک مدل با مقدار بایاس بالا توجه بسیار کمی به داده‌های آموزشی دارد و مدل را بیش از حد ساده در نظر میگیرد. اما واریانس به این معناست که یک خطا از حساسیت به نوسانات کوچک در مجموعه داده‌ی آموزش است، به این معنا که واریانس بالا سبب میشود که یک الگوریتم به جای خروجی‌های مورد نظر نویزهای تصادفی را در داده‌های آموزشی مدل کند. به عبارت دیگر یک مدل با میزان واریانس بالا توجه زیادی به داده‌های آموزشی دارد و قابل تعمیم به داده‌های دیده نشده نیست. یک عبارت ریاضی برای این مفاهیم به شکل زیر است:

$$\text{Bias}(X) = E[\hat{f}(X)] - f(X)$$

$$\text{Var}(X) = E\left[(\hat{f}(X) - E[\hat{f}(X)])^2\right]$$

بایاس در واقع تفاوت بین برچسب واقعی و مقدار برچسب پیش‌بینی شده می‌باشد اما واریانس در آمار تعریف میشود و امید ریاضی مربع انحراف یک متغیر تصادفی از میانگین آن است که در اینجا f نشان دهنده‌ی مدل ما در دنیای واقعی است. همچنین باید توجه کرد که ما دارای یک نویز تصادفی هستیم، که از آن نمیتوانیم اجتناب کنیم و آن را با اپسیلون نشان میدهند همچنین توجه شود که true label به شکل زیر نشان میدهند

$$y = f(X) + \epsilon.$$

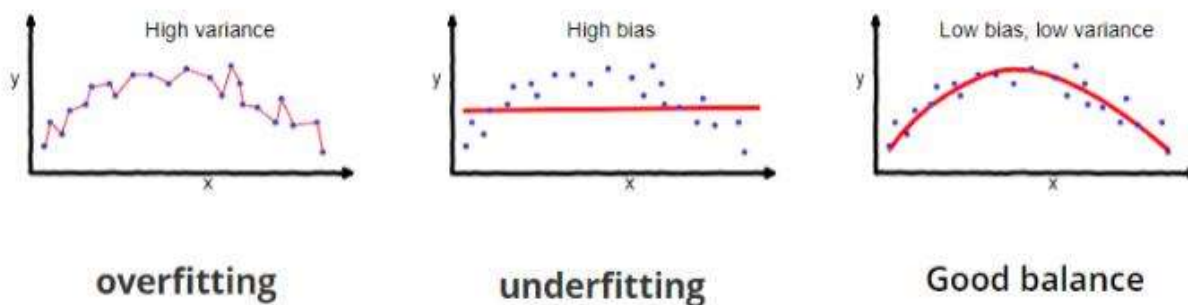
و ما میتوانیم خطا را به شکل زیر محاسبه کنیم

$$\begin{aligned} \text{Err}(X) &= E\left[(y - \hat{f}(X))^2\right] \\ &= E\left[(f(X) + \epsilon - \hat{f}(X))^2\right] \\ &= (E[\hat{f}(X)] - f(X))^2 + E\left[(\hat{f}(X) - E[\hat{f}(X)])^2\right] + \sigma_\epsilon^2 \\ &= \text{Bias}^2 + \text{Variance} + \text{Random Error}. \end{aligned}$$

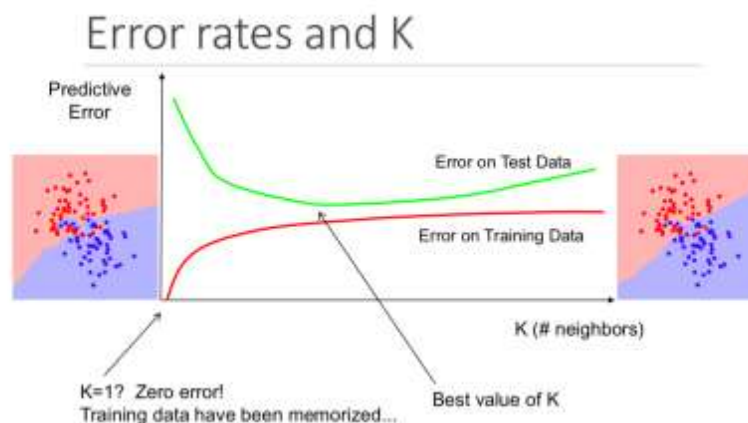
ضمناً توجه شود که دو اصطلاح دیگر در مورد بایاس و واریانس وجود دارد به نامهای **Overfitting** و **Underfitting**. منظور از **Overfitting** به این معناست که مدل داده های آموزش را به خوبی پیش بینی میکند اما در مورد داده های دیده نشده عملکرد خوبی ندارد

منظور از **Underfitting**: این است که مدل با داده های آموزش به خوبی مطابقت ندارد و به عبارت دیگر، داده های آموزش را به خوبی پیش بینی نمیکند.

به طور معمول **Underfitting** نشان دهنده ی بایاس زیاد و واریانس کم و **Overfitting** نشان دهنده ی بایاس کم و واریانس زیاد است



با توجه به میزان K باید توجه داشت که **KNN** در ابتدا میزان یادگیری پایینی دارد، در $K=1$ داده های یادگیری به خوبی پیش بینی میشوند و به این معناست که مقدار بایاس باید برابر با صفر باشد، اما در $K=1$ زمانی که داده ی جدید دیده نشده داریم (داده ی تست) شانس خطا بالاتر میرود به این معنا که واریانس بالایی داریم. اما زمانی که میزان K بیشتر میشود، خطای یادگیری (افزایش بایاس) افزایش می یابد اما خطای تست ممکن است با کاهش روبرو شود (کاهش واریانس). ما میتوانیم فکر کنیم که با افزایش K تعداد همسایه ها بیشتر میشوند و مدل پیچیده تر میگردد، با مدل زیر به راحتی میتوان میزان بایاس و واریانس را با افزایش سایز k بیان کنیم.



همانطور که در شکل بالا مشخص است در $K=1$ میزان خطای داده های تست بسیار بالا و میزان خطای یادگیری در حدود صفر است و مطالب گفته شده مطابق نمودار بالا نشان داده شده است

برای پنجره پارزن ما با مفهومی به نام Undersmoothing و Oversmoothing روبرو هستیم زمانی که h کوچک باشد، مقدار واریانس زیاد است و مقدار بایاس کوچک است. زمانی که h بزرگ است، مقدار واریانس کوچک است اما بایاس زیاد است.

(ب)

در ابتدا لازم است که توضیح کوتاهی در مورد روشهای پارامتریک بیان گردد، مفروضات در مورد شکل یک تابع میتواند روند یادگیری را ساده نماید و مدلهای پارامتریک با ساده سازی تابع به شکل شناخته شده مشخص میشود و مدل پارامتریک یادگیرندهای است که داده‌ها را از طریق مجموعه‌ای از پارامترها خلاصه میکند. پارامترها دارای اندازه‌ی ثابتی هستند یعنی اینکه مدل از قبل تعداد پارامترهای مورد نیاز خود را بدون توجه به داده‌های آنها میداند. در واقع پارامترها نیز مستقل از تعداد نمونه‌های آموزشی هستند. در مدلهای پارامتریک، دو مرحله وجود دارد، اولین مرحله انتخاب فرم تابع است و مرحله‌ی دوم یادگیری ضرایب تابع از داده‌های آموزشی است. نمونه‌های از مدلهای پارامتری شامل Logistic regression و SVM خطی است.

مدلهای غیرپارامتریک

الگوریتمهایی که مفروضات خاصی در مورد نوع تابع نگاشت ندارند به عنوان الگوریتمهای غیرپارامتریک شناخته میشوند. این روشها آزادی را دارند که هر فرم را از داده‌های آموزشی انتخاب کنند. در نتیجه مدل‌های پارامتریک برای تخمین تابع نگاشت به داده‌های بیشتری نسبت به مدل‌های پارامتریک نیاز دارند. ممکن است تصور شود که روشهای غیرپارامتریک به این معنی است که هیچ پارامتری وجود ندارد. اما این موضوع صحیح نیست بلکه به این معنی است که پارامترها نه تنها قابل تنظیم هستند بلکه میتوانند تغییر کنند. این منجر به تمایز کلیدی بین الگوریتمهای پارامتریک و غیرپارامتریک است، همانطور که پیش از این گفتیم الگوریتمهای پارامتریک بدون توجه به میزان داده‌های آموزشی تعداد پارامترهای ثابتی دارند. با این حال در مدل‌های غیرپارامتریک تعداد پارامترها به مقدار داده‌های آموزشی بستگی دارد و هر چه داده‌های آموزشی بیشتر باشد تعداد پارامترها بیشتر میشود و نتیجه این موضوع این است که در الگوریتمهای غیرپارامتریک ممکن است فرایند آموزش بسیار طول بکشد. به عنوان یک مثال الگوریتم KNN نمونه‌ای از الگوریتم غیرپارامتریک است. توجه شود که هیچ فرضی در مورد شکل تابع نگاشت به جز یک فرض وجود ندارد. فرض بر این است که مشابه‌ترین الگوهای آموزشی احتمال بیشتری برای تولید خروجی مشابه دارند. در این قسمت در مورد فواید و محدودیت‌های هر کدام از این مدلها صحبت میکنیم

سادگی: روش‌های الگوریتمهای پارامتریک برای درک آسانتر است. تفسیر پذیری نتایج نیز در مقایسه با مدل‌های غیرپارامتریک آسانتر است.

داده‌های آموزشی: مدل‌های پارامتریک به نسبت مدل‌های غیرپارامتریک به داده‌های کمتری نیاز دارند.

سرعت آموزش: روش‌های پارامتریک از نظر محاسباتی سریعتر از روش‌های غیرپارامتریک هستند، آنها را میتوان سریعتر از روش‌های غیرپارامتریک آموزش داد زیرا معمولاً پارامترهای کمتری برای آموزش نیاز دارند

در مورد مدلهای غیر پارامتریک:

کارایی: مدل های غیر پارامتریک ممکن است، پیش بینی های دقیقتری ارائه دهند، زیرا تناسب بهتری با داده ها نسبت به مدل های پارامتریک ارائه میدهند.

انعطاف پذیری: این الگوریتم های تناسب خوبی برای داده ها فراهم میکنند و میتوانند بسیاری از اشکال یک تابع را جا بدهند.

فرضیات کمی دارند: در مقایسه با الگوریتم های پارامتریک الگوریتم های غیر پارامتریک بیشتر از داده ها یاد میگیرند. این به این دلیل است که یادگیری الگوریتم های پارامتری ممکن است با مفروضاتی که آنها ایجاد میکنند محدود شوند

محدودیتها:

در مدل های پارامتریک:

محدودیت فرم: محدودیت مدل های پارامتریک در تعیین فرم تابعی است.

محدودیت در Fit: این روش ها بهترین تناسب را با داده ها ارائه نمیدهند.

پیچیدگی: الگوریتم های پارامتریک پیچیدگی محدودی ارائه میدهند، به این معنی است که آنها برای مسئله های با مقدار کمتری پیچیدگی مناسب هستند.

مدلهای پارامتریک:

بیش برآزش:

به همان اندازه که این الگوریتم ها تمایل دارند، داده ها را بهتر از الگوریتم های پارامتریک برآزش دهند بیشتر مستعد برآزش بیش از حد هستند

داده های آموزش: این روش ها به داده های بسیار بیشتری نسبت به الگوریتم های پارامتریک نیاز دارد.

سرعت: الگوریتم های غیر پارامتریک آهسته تر آموزش داده میشوند زیرا معمولاً پارامترهای بیشتر برای آموزش در نظر میگیرند

(ج)

تعدادی از معایب و مشکلات این روش ها را بیان میکنیم:

اولین مشکل این روش ها غیر قابل درک بودن آنهاست: درک آن چیزی که روشهای کرنل آموزش دیده اند دشوار است. دومین مشکل، انتخاب کردن یک کرنل مناسب است. امکان تعمیم این مسائل به چند کلاس وجود ندارد، در مقایسه با روش هایی مثل درخت تصمیم که میتواند به راحتی با چند کلاس سازگار گردد اما تعمیم به مسائل چند کلاسه با روش های مبتنی بر کرنل سخت است. مشکل بعدی این است که این روش ها به خوبی روش های یادگیری عمیق عمل نمیکند. یک مشکل این است که هزینه محاسباتی روشهای کرنل معقول است اما ناچیز نیست و زمانی که تعداد ابعاد

زیاد باشد ولی تعداد نمونه‌ها نسبتاً کم باشد (کمتر از یک میلیون) روش‌های مبتنی بر کرنل از نظر محاسباتی مناسب هستند.

(د)

مقدار $P(X)$ طریق رابطه $P(X) = \frac{K}{nV}$ محاسبه میشود که در دو روش پارزن و KNN تفاوت در ثابت نگه داشتن مقادیر K یا V میباشد. در روش پارزن مقدار حجم V ثابت نگه داشته میشود و مشخص میشود چه مقدار K درون آن حجم قرار میگیرد که این مقدار حجم دارای رابطه $V_n = 1/\sqrt{n}$ با تعداد نمونه‌ها دارد. اما در روش KNN پارامتر K ثابت نگه داشته میشود و حجم مربوط از داده‌ها که این مقدار K را در خود جا میدهد مشخص میشود. جزو شرایط این روش‌های این است که مقدار V_n در بی نهایت به صفر میل نماید و همچنین یعنی اگر چه V_n در بینهایت به صفر میل میکند اما سرعت رشد آن باید از سرعت رشد داده‌ها کمتر باشد.

سوال ۴-۲

(ج)

طبقه بندی KNN دارای این فرض است که نقاط مشابه یکدیگر دارای برچسب مشابه هم هستند. اما در فضاهای با ابعاد بالا نقاطی که براساس یک توزیع احتمالی هستند، به یکدیگر نزدیک نخواهند بود. به طور کلی طبقه بندیهایی مثل KNN که به فاصله‌ی زوجی بین نقاط متکی هستند به شدت تحت تاثیر مشکل نحسی ابعاد قرار میگیرند. در مسئله‌ی KNN ما یک نقطه تست داریم و میخواهیم K تا از نزدیکترین همسایه‌ها به این نقطه‌ی تست را بیابیم. به وضوح میتوانیم مشاهده کنیم با افزایش بعد مسئله، در فضاهای با تعداد بعد بالاتر داده‌های آموزشی چگالی کمتری خواهند داشت و برای یافتن همسایگان نقطه‌ی تست باید حجم زیادی از فضا جستجو شود. پس فاصله بین جفت نقاط با افزایش بعدها افزایش مییابد و در این صورت همسایه‌ها ممکن است از نقطه‌ی تست آنقدر دور باشند که دیگر اشتراک زیادی با آن نقطه تست نداشته باشند. به طور کلی طول کوچکترین Hyper Cube که شامل همه K نزدیکترین همسایه یک نقطه تست میباشد برابر است با $(k/N)^{1/d}$ که در اینجا N تعداد نمونه‌ها و d هم اندازه ابعاد ماست.

از این عبارت میتوان به این نتیجه رسید که با افزایش خطی تعداد ابعاد، تعداد نمونه‌های آموزشی باید به صورت تصاعدی افزایش یابد تا با نحسی ابعاد مقابله کند. همچنین روش دیگر برای مقابله با این موضوع این است که بعد را کاهش بدهیم یا داده‌ها را به ابعادی با فضایی پایینتر تبدیل کنیم.