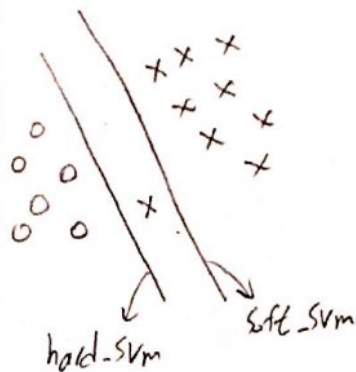


مسئله (۱)

الف :



یک) خیر جواب سخت دترم کزدمایگی نیست. مثلاً داده‌های "۰" را در نظر بگیرید.
 در این داده‌ها مشخص است که خط مناسب همان است که توسط Soft-SVM مشخص شده و در واقع آن یک \times به عنوان نویز باید در نظر گرفته شود.
 اما در Hard-SVM به همان یک نویز هم خیلی بهای می‌دهد و محدوده‌ها را یک \times را به عنوان support-vector در نظر می‌گیرد.

اما اگر در Soft-SVM پارامتر C را به بله داریم، می‌توانیم مدل را مجبور کنیم که به misclassified کردن یا حتی تراز کردن داده در margin خیلی بهای دهد و اینگونه محدوده‌ها مانند همان hard-margin محال می‌کند.

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

اگر به ∞ میل کند، مدل معادل Hard-SVM می‌شود.

(دو) برای Soft-SVM داریم:

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$y^n (w^T x^n + w_0) \geq 1 - \xi_i$$

$\xi_i > 0$: داده‌ی i ام به اشتباه دسته‌بندی شده است. (misclassified)

$\xi_i < 0$: داده‌ی i ام به درستی دسته‌بندی شده اما داخل margin قرار می‌گیرد.

$\xi_i = 0$: داده‌ی i ام با margin مناسب به درستی دسته‌بندی شده است.

در واقع با داده‌ی i ام می‌توانیم margin است یعنی خطی بردار پشتیبان است.

با این کار کلاً inactive است. در هر صورت خارج margin (یا می‌تواند margin) ۱ در دسترس است.

$z_i = 1$: این ویژگی خاص است. به این مفاسد داده نام دقیقاً این صفحه جدا کننده قرار دارد.
 به عبارتی $(w^T x_i + w_0) \geq 0$ شده است. و از آن جا که برای تابع sgn داریم:

$$\text{sgn}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

مثلاً داده نام به دسته خاصی تعلق می‌گیرد. مگر این که زیر فرض افشانه کمتر مثلاً $\text{sgn}(x) \geq 0$ را به دسته ۱ اختصاص دهیم.

Primal problem: $p^* = \min_x \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \mathcal{L}(x, \alpha, \lambda).$

Dual problem: $d^* = \max_{\{\alpha_i \geq 0\}, \{\lambda_i\}} \min_x \mathcal{L}(x, \alpha, \lambda)$

in general: $d^* \leq p^* \xrightarrow[n \text{ is affine}]{f \text{ and } g \text{ convex}} d^* = p^*$

Primal: $\min_{w, w_0} \max_{\{\alpha_n \geq 0\}} \left\{ \frac{1}{2} \|w\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)}(w^T x^{(n)} + w_0)) \right\}$ در حالت خاص برای SVM داریم:

dual: $\max_{\{\alpha_n \geq 0\}} \min_{w, w_0} \left\{ \frac{1}{2} \|w\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)}(w^T x^{(n)} + w_0)) \right\}$

- استفاده از لاگرانژین مرصع است زیرا:
- معمولاً حل آن راحت‌تر است.
 - دانش و بینش (insights) بیشتری در رابطه با زیر صفحه جدا کننده بهینه در اختیار ما می‌گذارد.
 - این امکان را به ما می‌دهد که از روش‌های مبتنی بر Kernel بهره ببریم.

Perceptron : $\max \{ -y^{(n)} h_{\theta}(x^n), 0 \}$

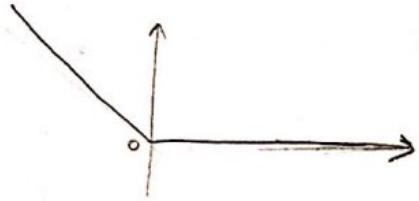
(1,68)

Logistic regression : $-\{ y^{(n)} \log(h_{\theta}(x^n)) + (1-y^{(n)}) \log(1-h_{\theta}(x^n)) \} \rightarrow$ log-loss or cross entropy

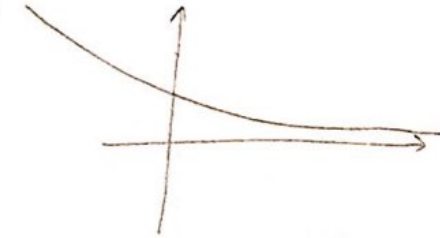
hard SVM : $\frac{1}{p} \|W\|^p$

soft SVM : $\frac{1}{p} \|W\|^p + C \sum_{n=1}^N \underbrace{\max(0, 1 - y^{(n)} (W^T x^n + w_0))}_{\text{hinge loss}}$

Perceptron Loss :



Logistic regression \rightarrow cross entropy loss :



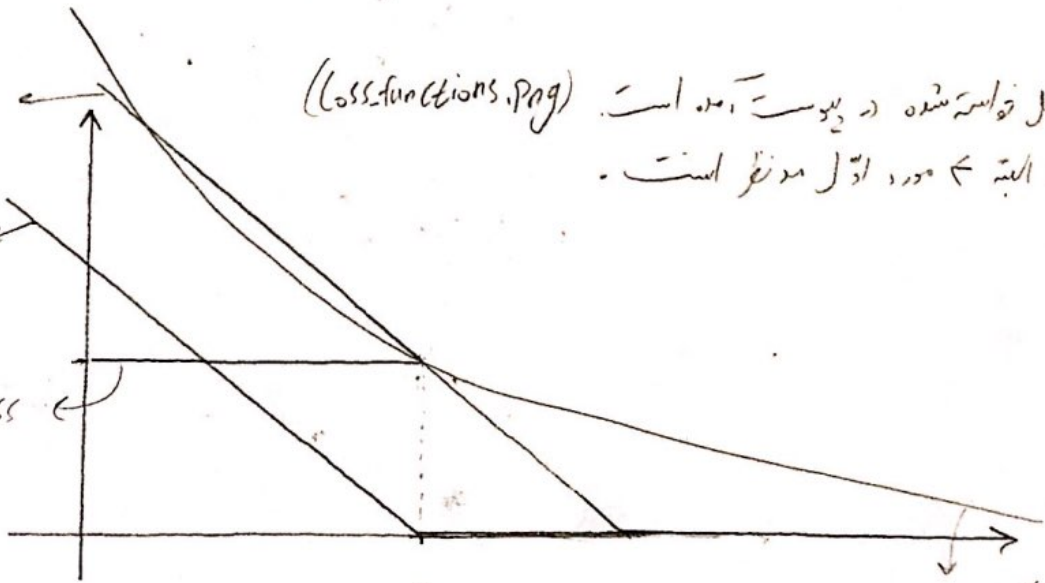
Soft-SVM \rightarrow hinge loss :



hinge loss

Perceptron loss

Zero-one loss



نکته: شکل فاصله شده در پیوسته است. (loss functions, pag)
 که البته \leftarrow مورد اول مد نظر است.

log loss = cross entropy loss

hard-SVM \rightarrow در تابع تابع قرار ندارد و یک Regularization $\|W\|^2$ اضافه می‌شود.

نکته: شکل فاصله شده در پیوسته است.

$$\xi_i = \max \{0, 1 - y^{(n)}(w^T x^{(n)} + w_0)\}$$

ب) تحلیل تابع گرانژ:

$$\min \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$

بنابراین تنها تفاوت در تابع گرانژ $\mathcal{L}(w, w_0, \alpha, \beta)$ در عبارت $\sum_{n=1}^N \xi_n$ خواهد بود.
چرا که شرایط مسئله یعنی $\xi_n \geq 0$ و $y^{(n)}(w^T x^{(n)} + w_0) \geq 1 - \xi_n$ مانده حالت قبل است. پس داریم:

$$\mathcal{L}(w, w_0, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - y^{(n)}(w^T x^{(n)} + w_0)) - \sum_{n=1}^N \beta_n \xi_n$$

پس مقادیر $\nabla_w \mathcal{L}$ و $\frac{\partial \mathcal{L}}{\partial w_0}$ مانده حالت اصلی شود.

اما برای مقدار $\frac{\partial \mathcal{L}}{\partial \xi_n}$ خواهیم داشت: (برای حالت اصلی)

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \alpha_n - \beta_n = 0 \Rightarrow \alpha_n = C - \beta_n, \quad \alpha_n \geq 0, \beta_n \geq 0 \Rightarrow 0 \leq \alpha_n \leq C$$

برای حالت جدید:

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 2C\xi_n - \alpha_n - \beta_n = 0 \Rightarrow \alpha_n = 2C\xi_n - \beta_n$$

$$\left. \begin{array}{l} \alpha_n \geq 0 \\ \beta_n \geq 0 \end{array} \right\} \Rightarrow 0 \leq \alpha_n \leq 2C\xi_n \times C \Rightarrow C \geq 2C \rightarrow 0 \leq \frac{\alpha_n}{\xi_n} \leq C$$

تحلیل تابع هزینه:

در این تابع جدید، بجای ξ از ξ^2 استفاده شده است. بنابراین به داده‌های نویز حساسیت بیشتری داریم.
نشان می‌دهد، چرا که برای داده‌های نویز، یعنی داده‌هایی که استیلا به دست می‌دهند، $\xi < 1$ است در نتیجه $\xi^2 < \xi$ و عملاً داده‌های نویز را دارد به نسبت از حالت عادی جریمه می‌کند پس باید انتظار داشته باشیم داده‌های نویز کمتر شوند.

از طرفی داده‌های بدون margin را کمتر جریمه می‌کند چرا که برای آن‌ها $1 - \xi < 0$ پس $\xi^2 < \xi$ خواهد بود.
پس اگر این ۲ حرف را کنار هم بگذاریم، این تابع جدید علاوه بر این که ξ^2 دارد که margin را کوچک می‌کند تا داده‌های کمتری نویز شوند و ایراد کمتری بینند که margin کوچک شود و داده‌ها داخل margin بیفتند.
یک جوی می‌توان گفت شبیه SVM-soft margin در حالت عادی عمل می‌کند با این زن که C آن را زیاد کرد تا margin را بیشتر کند.

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{n=1}^N \alpha_n (1 - y^{(n)} (w^T x^{(n)} + w_0))$$

باید

$$\max_{\{\alpha_n \geq 0\}} \min_{w, w_0} L(w, w_0, \alpha) \rightarrow \text{دلفی بهینه سازی}$$

بنابراین $N=5$ و x_1 تا x_5 هم که در صورت سوال مشخص است. y نامطلوب هم مشخص است در صورت سوال ($y=1$ یا -1 و بقیه $-$). پارامترهای w , w_0 نیز از رابطی زیر بدست می آیند:

$$\nabla_w L(w, w_0, \alpha) = 0 \Rightarrow w = \sum_{n=1}^N \alpha_n y^{(n)} x^{(n)}$$

$$\frac{\partial L(w, w_0, \alpha)}{\partial w_0} = 0 \Rightarrow \sum_{n=1}^N \alpha_n y^{(n)} = 0$$

$$L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y^{(n)} y^{(m)} x^{(n)T} x^{(m)}$$

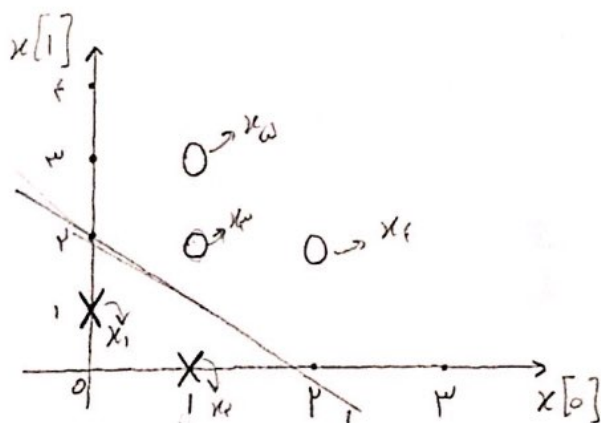
با جایگذاری خواهیم داشت:

که به ازای رابطی $L(\alpha) = \max_{\alpha} \alpha$ بدست می آید.

پس از آن به ازای یک نقطه که $\alpha_n > 0$ دارد (یعنی Support Vector است) از رابطی زیر برای بدست آوردن w_0 استفاده می کنیم:

$$y^{(s)} (w^T x^{(s)} + w_0) = 1 \Rightarrow w_0 = y^{(s)} - w^T x^{(s)}$$

پس مشکل بدست آمدن α_n ها هست.



$$y^{(n)} = 1 \rightarrow x$$

$$y^{(n)} = -1 \rightarrow 0$$

(د)

بنابراین x_4 و x_5 نمی توانند بردار پشتیبان باشند.

(نکته: در ادامه خط جدا کننده x_2 و x_3 است اما این جا تقریبی رسم شده است.)

داده و ضرایب برای $L(\alpha)$ را از قسمت یک به دست می آوریم:

$$L(\alpha) = \sum_{n=1}^3 \alpha_n - \frac{1}{r} \sum_{n=1}^3 \sum_{m=1}^3 \alpha_n \alpha_m y^{(n)} y^{(m)} x^{(n)T} x^{(m)}$$

$$= (\alpha_1 + \alpha_2 + \alpha_3) - \frac{1}{r} \left\{ \alpha_1^2 x_1 + \alpha_1 \alpha_2 x_0 + \alpha_1 \alpha_3 x_{-1} + \alpha_2 \alpha_1 x_0 + \alpha_2^2 x_1 + \alpha_2 \alpha_3 x_{-1} + \alpha_3 \alpha_1 x_{-1} + \alpha_3 \alpha_2 x_{-1} + \alpha_3^2 x_{-1} \right\}$$

$$= \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{r} \left\{ \alpha_1^2 - 2\alpha_1 \alpha_3 + \alpha_2^2 - 2\alpha_2 \alpha_3 + \alpha_3^2 \right\}$$

$$\sum_{n=1}^3 \alpha_n y^{(n)} = 0 \Rightarrow \boxed{\alpha_1 + \alpha_2 = \alpha_3}$$

از طرفی داریم:

با جایگذاری در رابطه $L(\alpha)$ خواهیم داشت:

$$L(\alpha) = \alpha_1 + \alpha_2 - \frac{1}{r} \left\{ \alpha_1^2 - 2\alpha_1^2 - 2\alpha_1 \alpha_2 + \alpha_2^2 - 2\alpha_2^2 + 2\alpha_1 \alpha_2 + \alpha_2^2 + \alpha_1^2 + \alpha_1 \alpha_2 + \alpha_2 \alpha_1 + \alpha_1^2 \right\}$$

$$= 2\alpha_1 + 2\alpha_2 - \alpha_1^2 - 2\alpha_1 \alpha_2 - \alpha_2^2 = L(\alpha_1, \alpha_2)$$

برای یافتن $\arg \max_{\alpha_1, \alpha_2} L(\alpha_1, \alpha_2)$ به دست می آوریم:

$$\left. \begin{aligned} \frac{\partial L(\alpha_1, \alpha_2)}{\partial \alpha_1} &= 2 - 2\alpha_1 - \alpha_2 = 0 \\ \frac{\partial L(\alpha_1, \alpha_2)}{\partial \alpha_2} &= 2 - \alpha_1 - 2\alpha_2 = 0 \end{aligned} \right\} \begin{aligned} \alpha_1 &= 1 \\ \alpha_2 &= 0 \end{aligned} \right\} \alpha_2 = 1$$

$$W = \sum_{n=1}^3 \alpha_n y^{(n)} x^{(n)} \xrightarrow{\alpha_1=1, \alpha_2=0, \alpha_3=0} \alpha_1 y^{(1)} x^{(1)} + \alpha_2 y^{(2)} x^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \quad (\text{درج})$$

$$W_0 = y^{(3)} - W x^{(3)} = -1 - \begin{bmatrix} -1 \\ 0 \end{bmatrix}^T \begin{bmatrix} 1 \\ 1 \end{bmatrix} = -1 - (-1) = 0$$

پس رابطه به شکل $W^T x + W_0$ که $W = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ و $W_0 = 2$.

$$\Rightarrow \begin{bmatrix} -1 \\ -1 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + 2 = \boxed{-x_1 - x_2 + 2} \rightarrow \text{رابطی مرز}$$

$$\Rightarrow y^{(n)} = \text{sgn} \left\{ -x_1^{(n)} - x_2^{(n)} + 2 \right\}$$

مسئله ۲) هسته

الف) K_1 و K_2 هسته معبر هسته. پس یک فضای برداری f_1 و f_2 وجود دارد:

$$K_1(x, x') = f_1(x)^T f_1(x')$$

$$K_2(x, x') = f_2(x)^T f_2(x')$$

تک: کافی است فضای برداری f_1 را از کنار هم قرار دادن f_1 و f_2 تشکیل دهیم.

$$f(x) = [f_1(x)^T, f_2(x)^T]^T$$

$$\Rightarrow K(x, x') = f(x)^T f(x') = [f_1(x)^T, f_2(x)^T] \begin{bmatrix} f_1(x') \\ f_2(x') \end{bmatrix} = f_1(x)^T f_1(x') + f_2(x)^T f_2(x')$$

پس فضای f_1 وجود دارد که K_1 در آن ضرب داخلی باشد. پس K_2 نیز یک هسته معبر است.

$$= K_1(x, x') + K_2(x, x') \rightarrow$$

$$K_f(x, x') = K_1(x, x') + K_2(x, x') = f_1(x)^T f_1(x') + f_2(x)^T f_2(x')$$

$$f_1(x) = [f_{11}(x), f_{12}(x), \dots, f_{1n}(x)]^T, f_2(x) = [f_{21}(x), f_{22}(x), \dots, f_{2m}(x)]^T$$

$$\Rightarrow K_f(x, x') = \sum_{i=1}^n f_{1i}(x) f_{1i}(x') + \sum_{j=1}^m f_{2j}(x) f_{2j}(x') = \sum_{i=1}^n \sum_{j=1}^m (f_{1i}(x) f_{2j}(x)) (f_{1i}(x') f_{2j}(x'))$$

$$= \sum_{i=1}^n \sum_{j=1}^m A_{ij}(x) A_{ij}(x'), A_{ij}(x) = f_{1i}(x) f_{2j}(x) \Rightarrow \text{یک فضای برداری مانند } A \text{ وجود دارد}$$

که $K_f(x, x')$ در واقع ضرب داخلی A در آن به فضای A باشد.

$$a \geq 0 \Rightarrow a = (\sqrt{a})^2$$

$$K_1(x, x') = a f_1(x)^T f_1(x') = \sqrt{a} f_1(x)^T \sqrt{a} f_1(x') = (\sqrt{a} f_1(x))^T (\sqrt{a} f_1(x')) \\ = f_0(x)^T f_0(x') = K_0(x, x'), \quad f_0(x) = \sqrt{a} f_1(x)$$

$$K_0 = \exp(K_1)$$

با استفاده از بسط تیلور داریم:

$$\exp(K_1) = \exp(0) + \exp(0) K_1 + \frac{\exp(0)}{2!} K_1^2 + \frac{\exp(0)}{3!} K_1^3 + \dots$$

$$\Rightarrow \exp(K_1) = 1 + K_1 + \frac{1}{2} K_1^2 + \frac{1}{6} K_1^3 + \dots$$

$$\Rightarrow \exp(K_1) = \sum_{j=0}^{\infty} \exp(0) \times K_1^j \times \frac{1}{j!}$$

پس K_0 از جمع تعدادی عبارت به دست آمده است.

هر عبارت از ضرب یک بهر مثبت در حاصل ضرب تعدادی K_1 به دست آمده است.

طبق بحثی در دوره اثبات کردیم که ضرب غریب مثبت در کرنل، باز هم یک کرنل می دهد و ضرب ۲ کرنل نیز یک کرنل معبر است.

پس هر عبارت یک کرنل معبر است. پس جمع آن ها نیز همین قسمت یک کرنل معبر است.

پس K_0 یک کرنل معبر است.

نکته: عبارت اول $\exp(0) = 1$ است که خودش یک کرنل است.

آیا $K_0(x, x') = 1$ را میتوانیم فضای f_0 وجود دارد:

$$f_0(x) = [1]$$

$$\Rightarrow K_0(x, x') = f_0(x)^T f_0(x') = 1 \times 1 = 1 \Rightarrow K_0(x, x') = 1 \rightarrow \text{کرنل معبر آ.}$$

ب) برای این که $K(A, B)$ یک عدد باشد، باید فضای دیگری را به گونه ای تعیین کنیم که در آن

$$|A \cap B| = Q(A)^T Q(B)$$

تابع Q_U را اینگونه تعریف می کنیم:

$$Q_U(X) = \begin{cases} 1 & \text{if } U \subseteq X \\ 0 & \text{o.w.} \end{cases}$$

بنابراین ضرب داخلی را اینگونه تعریف می کنیم:

$$Q(A)^T Q(B) = \sum_{U \subseteq S} Q_U(A) Q_U(B)$$

سنگینای بالا روی فضای زیرمجموعه های S (که مجموعه مرجع است) اعمال می شود.

به ازای هر U ، حاصل در صورتی 1 می شود که هم $Q_U(A) = 1$ هم $Q_U(B) = 1$.

پس حاصل زمانی 1 می شود که U هم زیرمجموعه A باشد هم زیرمجموعه B .

پس سنگینای بالا به ازای $U \subseteq A \cap B$ برابر با 1 است، در غیر این صورت 0 است.

بنابراین حاصل عبارت بالا برابر تعداد زیرمجموعه های $A \cap B$ است که می شود $|A \cap B|$.

پس داریم:

$$Q(A)^T Q(B) = |A \cap B| = K(A, B)$$

یک: $x, x' \in \mathbb{R}^2$ باشند قاعده است:

$$\begin{aligned} \kappa(x, x') &= (C + x^T x')^2 = (C + x_1 x'_1 + x_2 x'_2)^2 \\ &= C^2 + 2C x_1 x'_1 + 2C x_2 x'_2 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 x_2 x'_2 \end{aligned}$$

$$\Rightarrow Q(x) = \left[C, \sqrt{C} x_1, \sqrt{C} x_2, x_1^2, x_2^2, \sqrt{2} x_1 x_2 \right]^T$$

حال برای $x, x' \in \mathbb{R}^d$ داریم: (مستطاب)

$$\kappa(x, x') = (C + x^T x')^2 = (C + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^2$$

$$\begin{aligned} &= C^2 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + \dots + x_d^2 x'^2_d + 2C x_1 x'_1 + 2C x_2 x'_2 + \dots + 2C x_d x'_d + 2x_1 x'_1 x_2 x'_2 + \dots \\ &\quad + 2x_1 x'_1 x_d x'_d + 2x_2 x'_2 x_d x'_d + \dots + 2x_{d-1} x'_{d-1} x_d x'_d \end{aligned}$$

$$\Rightarrow Q(x) = \left[C, x_1^2, x_2^2, \dots, x_d^2, \sqrt{C} x_1, \sqrt{C} x_2, \dots, \sqrt{C} x_d, \sqrt{2} x_1 x_2, \sqrt{2} x_1 x_3, \dots, \sqrt{2} x_1 x_d, \dots, \sqrt{2} x_2 x_3, \dots, \sqrt{2} x_2 x_d, \dots, \sqrt{2} x_{d-1} x_d \right]^T$$

و: عملاً d بعد آن صفر (یا به عبارت دیگر صفر) می شود.

$$\kappa(x, x') = (x^T x')^2 = (x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^2$$

$$\Rightarrow Q(x) = \left[x_1^2, x_2^2, \dots, x_d^2, \sqrt{2} x_1 x_2, \sqrt{2} x_1 x_3, \dots, \sqrt{2} x_1 x_d, \sqrt{2} x_2 x_3, \dots, \sqrt{2} x_2 x_d, \dots, \sqrt{2} x_{d-1} x_d \right]^T$$

و جنبه دیگری bias ندارد.

$$K(x, x') = (c + x^T x')^m = (c + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^m$$

این عبارت تعدادی چندجمله ای می شود که هر کدام به تعداد مساوی x_i و x'_i دارند. بنابراین می توان به دیگری اشاره کرد.

پس محلاً سؤال این است که عبارت $(x_1 x'_1 + \dots + x_d x'_d)^m$ چند عبارت می شود.

$$\left(\sum_{i=1}^d x_i x'_i \right)^m$$

عبارت را نیز به شکل $x_i x'_i$ می بینیم. پس محلاً خواهیم داشت:

برای سادگی $x_i x'_i$ را معادل a_i در نظر می گیریم. پس داریم:

$$(a_0 + a_1 + a_2 + \dots + a_d)^m$$

هر عبارت به این شکل خواهد بود:

$$a_0^{b_0} a_1^{b_1} \dots a_d^{b_d}$$

$$0 \leq b_i \leq m$$

$$\sum_{i=0}^d b_i = m$$

پس مسئله تبدیل شد به این که a_i متغیر b_i داریم که می خواهیم محضاً m شود. مقادیر صحیح بین 0 تا m می گیرند. حالات مختلف این b_i ها جواب مسئله است.

$$\binom{m+d}{d} = \binom{m+d}{m}$$

که برابر است با: پس تعداد اشیاء برابر است انتخاب m از $m+d$.

$$\forall b_i : 0 \leq b_i, \sum_{i=1}^k b_i = m$$

روش حل: k متغیر داریم که به ازای هر کدام:

$$\forall b_i : 0 \leq b_i, \sum_{i=1}^k b_i = m+k$$

می توان به هر متغیر یک واحد اضافه کرد. پس مسئله تبدیل می شود به: پس محلاً مانند این است که $m+k$ مهره داشته باشیم. پس $m+k-1$ جایگاه داریم بین مهره ها، کافی است $k-1$ جایگاه را انتخاب کنیم. که می شود $\binom{m+k-1}{k-1}$.

$$\binom{m+d}{m} = \binom{m+d}{d}$$

که این مثال $k=d$ (چون k صفر داریم). پس جواب نهایی می شود

(الف)

information gain \rightarrow IG

$$IG(X_d, Y) = H(Y) - H(Y|X_d)$$

$$H(Y) = - \sum_{y_j \in Y} P(y_j) \log P(y_j)$$

$$H(Y|X_d) = - \sum_i \sum_j P(X_d=i, Y=j) \log P(Y=j|X_d=i)$$

$$P(Y=j) = \sum_{i \in X_d} P(Y=j, X_d=i) \Rightarrow H(Y) = - \sum_i \sum_j P(X_d=i, Y=j) \log P(Y=j)$$

$$\Rightarrow IG(X_d, Y) = - \sum_i \sum_j \left\{ P(X_d=i, Y=j) \times [\log P(Y=j) - \log P(Y=j|X_d=i)] \right\}$$

$$\log P(Y=j) - \log P(Y=j|X_d=i) = \log \frac{P(Y=j)}{P(Y=j|X_d=i)} = \log \frac{P(Y=j)P(X_d=i)}{P(Y=j, X_d=i)}$$

$$\Rightarrow IG(X_d, Y) = - \sum_i \sum_j P(X_d=i, Y=j) \log \frac{P(Y=j)P(X_d=i)}{P(Y=j, X_d=i)}$$

اگر X_d و Y با هم مستقل باشند، $P(Y=j)P(X_d=i) = P(Y=j, X_d=i)$ در این صورت عبارت داخل لگاریتم برابر با ۱ می شود پس عبارت داخل سلیک با ۰ می شود و $\log 1 = 0$ ضرب می شود و $IG(X_d, Y) = 0$ خواهد بود.

از آن جا که IG نمی تواند منفی باشد، پس X_d در واقع کمترین gain را بین ویژگی ها دارد. پس نمی تواند به عنوان ریشه انتخاب شود چرا که ID_3 ویژگی با بیشترین IG را در مرحله انتخاب می کند. (مگر این که همگی ویژگی ها مستقل از هم باشند)

(ب)

۲. دیدگاه کلی در رابطه با نحوه برخورد با شماره پدیده وجود دارد.

اگر به چشم یک متغیر پیوسته به آن نگاه کنیم و روی آن یک $threshold$ قرار دهیم که مثلاً شماره های کمتر از K در دسته ۱ و مابقی در دسته ۲ قرار گیرند، در این صورت چون شماره پدیده مستقل از نوع بیماری است، به عنوان رتبه انتخاب خواهد شد. مگر این که بهمانسان سیاست خاصی در تعیین شماره پدیده انداخته باشد. مثلاً یک چکاپ ادبی انجام دهد و به بیماران بدضم عدد کوچکی نسبت به بیماران خوش خیم تخصیص دهد. در این صورت اگر K مناسب انتخاب شود، می تواند به عنوان رتبه انتخاب شود و خیلی هم خوب داده ها را دسته بندی کند و می داده های نسبت هم خوب عمل خواهد کرد و $overfitting$ نداریم. مثلاً فرض کنید بهمانسان به بیماران بدضم شماره ۵ تا ۱۰۰ و به بیماران خوش خیم شماره ۱۰۱ تا ۲۰۰ بدهد. اگر $K = ۱۵۰$ انتخاب شود، مدل بسیار مناسبی خواهد داشت.

دیدگاه دوم این است که نحوه استفاده از شماره پدیده به شکل یک متغیر $categorical$ باشد. یعنی مثلاً در ازای هر یک حالت، به ازای هر مقدار موجود برای شماره پدیده، یک دسته «نظر بگیرم» یعنی در دست تصمیم دهنی از مجید از این مدل، به تعداد بیماران زنند خواصم داشت که چون هر شماره منحصر به فرد است، در هر دسته دقیقاً یک بیمار قرار می گیرد و دست می داده های $train$ ۱۰۰ داده می شود. که در این صورت قطعاً شماره پدیده به عنوان رتبه انتخاب خواهد شد. و شب این شکل به است چرا که مدل دست ۱۰۰ داده می دهد اما در حقیقت هیچ چیزی یاد نگرفته است و کاملاً $overfit$ شده است. در نقطه این که در تحلیل متغیر شماره پدیده مناسب نیست. پس به جای این که هر شماره پدیده را یک زنند «نظر بگیرم»، می توانیم یک $threshold$ قرار دهیم و یک باره را به عنوان زنند «نظر بگیرم» در نهایت به نحوی تعداد زنان را محدود کنیم.

با این که دست را کامل بسازیم و پس از آن از پایین به بالا به $tune$ کردن آن بپردازیم. اما در حالت کلی اگر شماره پدیده قبل از تشخیص بیماری به بیمار داده شود (یعنی عملاً مستقل باشد) و تعداد زنان و بیماران شماره پدیده را محدود «نظر بگیرم» قاعده IG خاصی نمی دهد و به عنوان رتبه انتخاب نخواهد شد. پس نکته هم این است که به ازای هر شماره پدیده یک زنند بسازیم و تعداد دسته ها را محدود در نظر بگیریم.

ب) یک:

برای این کار کافی است IP_3 را یک بار تا انتها اجرا کنیم. اثباتی خود را با شرایط گفته شده در صورت سوال، دست نوشته خطای آموزشی منفر خواهد داشت.

یک بزرگ را در نظر بگیریم. اگر مقادیر این بزرگ یکسان باشند مشکلی نداریم. چون با همگی مقادیر دست را دارند، یا همگی یک مقدار غلط دارند که کافی است بنده آخرین درستی را برای این دسته به مقدار صحت تغییر دهیم.

پس مشکل زمانی رخ می دهد که یک بزرگ مقادیر یکسان نداشته باشند. چون داریم یک بزرگ را بررسی می کنیم، پس IP_3 متوقف شده و دیگر آن بزرگ را دسته بندی نکرده است (چون اگر زیاد بزرگ بود).

شرط توقف IP_3 این است که همه مقادیر یک دسته یکی باشند، یا این که از تمام درستی ها استفاده کرده باشند.

چون مقدار این بزرگ یکسان نیستند و IP_3 متوقف شده، پس بنده می گیریم که احمق این دسته همگی از تمام درستی ها بهره گرفته. چون دست بیان نگار است، از هر بزرگ به ریشه تنها یک مسیر وجود دارد. پس تمام احمق یک دسته به ازای هر درستی، دقیقاً در یک دسته یکسان اشتباه شده اند. از آن جا که هر درستی، داده ها را به دسته تقسیم می کند، پس در آن گروهی که به این محاسبات که داده ها درستی له ام یکسانی داشته اند.

پس تمام درستی های داده های که در یک دسته قرار دارند یکی است. پس محدوده ۲ یکسان داریم که ۳ متفاوت دارند. اگر این را متناقص در نظر بگیریم، بنده می شود که دست محدوده، یکی وجود ندارد که مقادیر غیر یکسان داشته باشند پس خطای آموزشی صفر است. اگر وجود ۲ داده یکسان که ۳ متفاوت دارند را متناقص در نظر بگیریم، مطلقاً به خطای آموزشی منفر نخواهیم رسید چرا که تحت هیچ شرایطی آن ۲ داده یکسان از هم قابل تشخیص نخواهند بود.

پس کافی است از IP_3 استفاده کنیم.

یعنی به ازای درستی له ام، به ازای هر مقدار ممکن یعنی $\{1, 2, \dots, 10\}$ یک زیر شاخه تشکیل دهیم و تا زمانی که مقادیر موجود در یک دسته همگی یکسان نشده اند یا این که درستی ها تمام نشده است به کار ادامه دهیم. طبق اثبات دست حاصل خطای آموزشی صفر خواهد داشت.

ضریح
 مثال نفی: داده ها را در نظر بگیرید به صورتی که تنها یک ویژگی پیوسته وزن دارند. هدف تعیین جنسیت ازاد است. مجموعه داده ما به این شکل است:

$$D = \{ (پسر, ۵), (دختر, ۵), (پسر, ۵) \}$$

تنها یک ویژگی داریم و هدف وزن کنار است. پس در همان تفلیک اول کار تمام است.
 فرض کنید آساند را به نحوی انتخاب کنیم که خطای آموزش صفر شود. یعنی برای وزن های کمتر از K پسر باشد و بیشتر از K دختر، پسر تشخیص داده شود.

$$K < x \rightarrow \text{پسر} \Rightarrow x < K \text{ و } x < K \Rightarrow \text{پسر}$$

$$K < x \rightarrow \text{دختر} \Rightarrow x < K \text{ و } x < K \Rightarrow \text{دختر}$$

آر وزن های کمتر از K را دختر بگیریم و

پسر، دختر تشخیص داده شود.

پس چنین K ای وجود ندارد. پس همین هدف ضمیمی نمی توان ساخت.

سه :

وقتی در حالت گسسته بر اساس یک ویژگی تقسیم بندی را انجام می دهیم، مقدار دقیق آن ویژگی را می دانیم پس استفاده مجدد از آن ویژگی هیچ منفی برای ما ندارد. پس استفاده مجدد از آن گامی در افزایش دقت و ادله دست بندی به ما می گذرد.

اما در حالت پیوسته، وقتی بر اساس یک ویژگی تقسیم را انجام می دهیم، در دایره مقدار دقیق آن را نمی دانیم و صرفاً برای مقدار آن یک باره در نظر گرفته ایم. پس استفاده مجدد از آن ویژگی می تواند سودمند باشد و به ما کمک کند که مقدار امی آن ویژگی را برای هر داده به صورت دقیق تر تخمین بزنیم. در نتیجه می توانیم با مقادیر دقیق تر، داده ها را بهتر تفلیک کنیم و دست بندی را ادامه دهیم و در نهایت خطای آموزش را کاهش دهیم.

املاً برای همین است که قسمت "گام" جواب داشت چون با استفاده از هر ویژگی معطایک بار، در حالت گسسه تمام اطلاعات موجود را به دست می آوریم. اما در حالت پیوسته این چنین نیست.

(الف)

$$H_M(x) = \frac{1}{M} \sum_{m=1}^M h_m(x)$$

$$\Rightarrow (H_M(x) - h(x))^2 = \left(\frac{1}{M} \sum_{m=1}^M h_m(x) - h(x) \right)^2 = \frac{1}{M^2} \left(\sum_{m=1}^M h_m(x) \right)^2 + h(x)^2 - \frac{2h(x)}{M} \sum_{m=1}^M h_m(x)$$

$$\Rightarrow E_{\text{Cor}} = E_x \left[\frac{1}{M^2} \left(\sum_{m=1}^M h_m(x) \right)^2 + h(x)^2 - \frac{2h(x)}{M} \sum_{m=1}^M h_m(x) \right]$$

$$= \frac{1}{M^2} E_x \left[\left(\sum_{m=1}^M h_m(x) \right)^2 \right] + E_x [h(x)^2] - \frac{2}{M} E_x \left[h(x) \sum_{m=1}^M h_m(x) \right] \rightarrow \text{نزاره ①}$$

$$E_{\text{avg}} = \frac{1}{M} \sum_{m=1}^M E_x [(h_m(x) - h(x))^2] = \frac{1}{M} E_x \left[\sum_{m=1}^M (h_m(x) - h(x))^2 \right]$$

$$= \frac{1}{M} E_x \left[\sum_{m=1}^M (h_m(x)^2 + h(x)^2 - 2h(x)h_m(x)) \right] = \frac{1}{M} E_x \left[\sum_{m=1}^M h_m(x)^2 + \sum_{m=1}^M h(x)^2 - 2h(x) \sum_{m=1}^M h_m(x) \right]$$

$$E_x \left[\frac{1}{M} \sum_{m=1}^M h_m(x)^2 \right] = E_x [h(x)^2]$$

$$= \frac{1}{M} E_x \left[\sum_{m=1}^M h_m(x)^2 \right] + \frac{1}{M} E_x \left[\sum_{m=1}^M h(x)^2 \right] - \frac{2}{M} E_x \left[h(x) \sum_{m=1}^M h_m(x) \right] \rightarrow \text{نزاره ②}$$

$$E_{\text{Cor}} - E_{\text{avg}} = \frac{1}{M^2} E_x \left[\left(\sum_{m=1}^M h_m(x) \right)^2 \right] - \frac{1}{M} E_x \left[\sum_{m=1}^M h_m(x)^2 \right] : \text{بقیه می شود که : ①، ②، ③}$$

پس باید اثبات کنیم که این عبارت کوچکتر مساوی ۰ است.

یلم کوئی حوالہ:

$$\left(\sum_{i=1}^n u_i v_i \right)^2 \leq \left(\sum_{i=1}^n u_i^2 \right) \left(\sum_{i=1}^n v_i^2 \right)$$

بالا دن $u_i = 1$, $n = M$ و $v_i = h_i(x)$ ظاہر دست:

$$\left(\sum_{i=1}^M 1 \times h_i(x) \right)^2 \leq \sum_{i=1}^M 1^2 \times \sum_{i=1}^M h_i^2(x)$$

$$\Rightarrow \left(\sum_{m=1}^M h_m(x) \right)^2 \leq M \sum_{m=1}^M h_m^2(x) \rightarrow \text{طرف دایہ } \frac{1}{M^2} \text{ ضرب کیے۔}$$

$$\Rightarrow \frac{1}{M^2} \left(\sum_{m=1}^M h_m(x) \right)^2 \leq \frac{1}{M} \sum_{m=1}^M h_m^2(x) \rightarrow \text{از د طرف امید ریاضی لی۔}$$

$$\Rightarrow \frac{1}{M^2} E_x \left[\left(\sum_{m=1}^M h_m(x) \right)^2 \right] \leq \frac{1}{M} E_x \left[\sum_{m=1}^M h_m^2(x) \right]$$

$$\Rightarrow \frac{1}{M^2} E_x \left[\left(\sum_{m=1}^M h_m(x) \right)^2 \right] - \frac{1}{M} E_x \left[\sum_{m=1}^M h_m^2(x) \right] \leq 0$$

$$\Rightarrow E_{\text{com}} - E_{\text{avg}} \leq 0 \Rightarrow E_{\text{com}} \leq E_{\text{avg}}$$

$$E_{\text{com}} = E_x \left[\left(\frac{1}{M} \sum_{n=1}^M h_n(n) - h(n) \right)^2 \right] = E_x \left[\left(\frac{1}{M} \sum_{n=1}^M h_n(n) - h(n) \right) \left(\frac{1}{M} \sum_{n=1}^M h_n(n) - h(n) \right) \right] \quad (ب)$$

$$= \frac{1}{M^2} \sum_{m=1}^M \sum_{l=1}^M E_x \left[(h_m(n) - h(n)) (h_l(n) - h(n)) \right]$$

حاصل امید ریاضی بالا به ازای $m \neq l$ برابر با صفر است. (طبق فرض مسئله).

$$\frac{1}{M^2} \sum_{m=1}^M E_x \left[(h_m(n) - h(n)) (h_m(n) - h(n)) \right] \quad \text{پس عبارت بالا برابر است با:}$$

$$= \frac{1}{M^2} \sum_{m=1}^M E_x \left[(h_m(n) - h(n))^2 \right] = \frac{1}{M^2} \left\{ \frac{1}{M} \sum_{m=1}^M E_x \left[(h_m(n) - h(n))^2 \right] \right\}$$

$$= \frac{1}{M} E_{\text{avg}} \Rightarrow \boxed{E_{\text{com}} = \frac{1}{M} E_{\text{avg}}}$$

الف) فرض می‌کنیم $h_t = h_{t+1}$ برای مرتبه t داریم:

$$\epsilon_t = \frac{\sum_{i=1}^n w_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))}{\sum_{i=1}^n w_t^{(i)}}$$

برای راضی، $w_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))$ را با $w_t^{(i)}$ نشان می‌دهیم.برای راضی، $w_t^{(i)} \times I(y^{(i)} = h_t(x^{(i)}))$ را با $w_t^{(i)}$ نشان می‌دهیم.پس w_t وزن داده‌هایی است که h_t نتوانسته تشخیص بدهد و w_t وزن داده‌هایی است که h_t به درستی تشخیص داده است. پس داریم:

$$\epsilon_t = \frac{\sum_{y^i \neq h_t(x^i)} w_t^{(i)}}{\sum_{y^i \neq h_t(x^i)} w_t^{(i)} + \sum_{y^i = h_t(x^i)} w_t^{(j)}} < \frac{1}{2}$$

در AdaBoost فقط وزن داده‌هایی که توسط h_t اشتباه تشخیص داده شده اند کاهش می‌شود. پس وزن $w_{t+1}^{(i)} = w_t^{(i)}$ اما وزن داده‌هایی که به درستی تشخیص داده شده اند:

$$w_{t+1}^{(i)} = w_t^{(i)} e^{\alpha_t}, \quad \alpha_t = \ln\left(\frac{1-\epsilon_t}{\epsilon_t}\right) \Rightarrow \begin{cases} w_{t+1}^{(i)} = w_t^{(i)} \times \frac{1-\epsilon_t}{\epsilon_t} \\ w_{t+1}^{(i)} = w_t^{(i)} \end{cases}$$

دقت شود که داده‌هایی که با h_t درست تشخیص داده می‌شوند، دقیقاً همان داده‌هایی هستند که توسط h_{t+1} درست تشخیص می‌شوند چرا که فرض کردیم $h_{t+1} = h_t$.حال خطای ϵ_{t+1} را حساب می‌کنیم:

$$\epsilon_{t+1} = \frac{\sum w_{t+1}^{(i)}}{\sum w_{t+1}^{(i)} + \sum w_{t+1}^{(j)}} = \frac{\frac{1-\epsilon_t}{\epsilon_t} \times \sum w_t^{(i)}}{\frac{1-\epsilon_t}{\epsilon_t} \sum w_t^{(i)} + \sum w_t^{(j)}}$$

$$\frac{1-\epsilon_t}{\epsilon_t} = \frac{1}{\epsilon_t} - 1 = \frac{\sum w_t^{(i)} + \sum w_t^{(j)}}{\sum w_t^{(i)}} - 1 = \frac{\sum w_t^{(j)}}{\sum w_t^{(i)}}$$

$$\Rightarrow \epsilon_{t+1} = \frac{\frac{\sum w_t^{(j)}}{\sum w_t^{(i)}} \times \sum w_t^{(i)}}{\frac{\sum w_t^{(j)}}{\sum w_t^{(i)}} \times \sum w_t^{(i)} + \sum w_t^{(j)}} = \frac{\sum w_t^{(j)}}{2 \sum w_t^{(j)}} = \frac{1}{2}$$

پس اگر $h_{t+1} = h_t$ باشد، ϵ_{t+1} دقیقاً برابر با ϵ_t می شود. پس خطای یادگیری کاهش می یابد.
این خلاف فرض است که خطای یادگیری های ضعیف از بار اول کمتر است.

پس به تناقض رسیدیم. پس $h_{t+1} \neq h_t$

(ب) الگوریتم AdaBoost از تابع خطای نهایی استفاده می کند. فرض کنید h_t یادگیر ساده به H_t افاده شده باشد. یادگیر $t+1$ ام، h_{t+1} نامی خواهد داشت که کمتر از h_t خطای دارد.

$$H_t(x) = \frac{1}{\rho} (\alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_t h_t(x)) = H_{t-1}(x) + \frac{1}{\rho} \alpha_t h_t(x)$$

$$E = \sum_{i=1}^N e^{-y^{(i)} H_t(x^{(i)})} = \sum_{i=1}^N e^{-y^{(i)} [H_{t-1}(x^{(i)}) + \frac{1}{\rho} \alpha_t h_t(x^{(i)})]}$$

$$= \sum_{i=1}^N \underbrace{e^{-y^{(i)} H_{t-1}(x^{(i)})}}_{\text{این عبارت باید کمینه شود}} e^{-\frac{1}{\rho} \alpha_t y^{(i)} h_t(x^{(i)})} = \sum_{i=1}^N w_t^i e^{-\frac{1}{\rho} \alpha_t y^{(i)} h_t(x^{(i)})}$$

در معادله t ثابت است و برابر با $w_t^{(i)}$ می باشد.

برای به دست آوردن α_t و h_t .

برای به دست آوردن h_t خواهیم داشت:

$$E = \sum_{i=1}^N w_t^{(i)} e^{\frac{1}{2} \alpha_t y^{(i)} h_t(x^{(i)})} = e^{\frac{-\alpha_t}{2}} \sum_{y^{(i)} = h_t(x^{(i)})} w_t^{(i)} + e^{\frac{\alpha_t}{2}} \sum_{y^{(i)} \neq h_t(x^{(i)})} w_t^{(i)}$$

$$= \underbrace{\left(e^{\frac{\alpha_t}{2}} - e^{\frac{-\alpha_t}{2}} \right)}_{\substack{\text{به مقدار } h_t \text{ ربطی ندارد} \\ \text{و در ضمن مقدرای مثبت است}}} \sum_{y^{(i)} \neq h_t(x^{(i)})} w_t^{(i)} + \underbrace{e^{\frac{-\alpha_t}{2}} \sum_{i=1}^N w_t^{(i)}}_{\text{مستقل از } h_t(x) \text{ است.}}$$

به مقدار h_t ربطی ندارد.
و در ضمن مقدرای مثبت است.

مستقل از $h_t(x)$ است.

پس محلاً تابعی که باید بهینه شود به این شکل است:

$$\sum_{\substack{y^{(i)} \\ y^{(i)} \neq h_t(x^{(i)})}} w_t^{(i)} \Rightarrow J_t = \sum_{i=1}^N w_t^{(i)} \times I(y^{(i)} \neq h_t(x^{(i)}))$$

پس محلاً تعیین h_t و نتیجه آنگونه فقط به وزن داده‌های جدا دارد که به انتخاب دسته بندی شده اند.
حال به سراغ به روز رسانی وزن‌های می‌رویم.

طبق روش AdaBoost می‌دانیم به روز رسانی اینچنین اتفاق می‌افتد:

$$\left. \begin{aligned} w_{t+1}^{(i)} &= w_t^{(i)} e^{\frac{1}{2} \alpha_t y^{(i)} h_t(x^{(i)})} \\ y^{(i)} h_t(x^{(i)}) &= 1 - 2 I(y^{(i)} \neq h_t(x^{(i)})) \end{aligned} \right\} \begin{aligned} w_{t+1}^{(i)} &= w_t^{(i)} e^{\frac{1}{2} \alpha_t I(y^{(i)} \neq h_t(x^{(i)}))} \\ \text{این قسمت مستقل از } z \text{ هست و محلاً} \\ \text{در همی وزن‌ها ضرب می‌شود پس تأثیری ندارد.} \end{aligned}$$

$$\Rightarrow w_{t+1}^{(i)} = w_t^{(i)} e^{\alpha_t I(y^{(i)} \neq h_t(x^{(i)}))}$$

پس داده‌هایی که به درستی تشخیص داده شده اند، I مقدار صفر خردی می‌دهد و $1 = e^0$ پس $w_{t+1}^{(i)} = w_t^{(i)}$
پس فقط داده‌هایی که انتخاب دسته بندی شده اند به روز رسانی می‌شوند.
وزن