

$$P(y|x; w) = \prod_{i=1}^n P(y_i | x_i, w) \Rightarrow \sum \begin{cases} f(x, w) & y=+1 \\ 1-f(x, w) & y=-1 \end{cases} =$$

(الف) (۱)

این توابع میزنند

$$\frac{1}{2} \log \sum (1+y_n) \log f(x_n; w) + (1-y_n) \log f(x_n; w))$$

این توابع چند مثال می دهند زیرا اگر $y_n = +1$ یا $y_n = -1$ را جایگزین کنیم می شود رابطه بالا.

$$\frac{dg}{dx} = \frac{e^x}{1+e^x} \rightarrow \text{آنچه داریم}$$

$$\frac{e^x(e^x+1) - e^x(e^x+0)}{(e^x+1)^2} =$$

سوال ۱ ب داریم.

$$\frac{e^x(e^x+1) - e^{2x}}{(e^x+1)^2} = \frac{e^x}{(e^x+1)^2}$$

$$e^x \geq 0 \quad \text{سری}$$

$$(e^x+1)^2 \geq 0 \quad \text{سری}$$

$$x > 0 \rightarrow e^x > 0$$

$$x < 0 \rightarrow e^{-x} \rightarrow \frac{1}{e^x} > 0$$

$$g \circ f = g(f(x)) \Rightarrow \frac{dg}{dx} = g'(f(x)) f'(x) = k$$

$$\frac{dk}{dx} = \underbrace{g'(f(x))}_{\text{معمولی}} \underbrace{f''(x)}_{\text{محدب}} + \underbrace{g''(f(x))}_{\text{محدب}} \underbrace{f'(x) f'(x)}_{(f'(x))^2} \geq 0.$$



$$k = \frac{dh}{dx} = \sum_{k=1}^K c_k P'_k \rightarrow \frac{d^2 h}{dx^2} = \sum_{k=1}^K c_k P''_k = \text{معدب}$$

می دانیم $c_k \geq 0$ فرض شود
 $P''_k \geq 0$ معدب

① ②

محدب

① ثابت
 $\log(1 + e^{\omega^T x_n})$ محدب است. با توجه به این در قضیه ۱ ب ثابت برسیم.

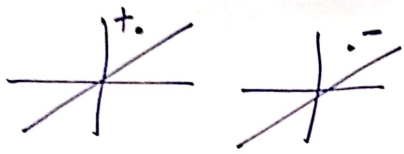
$$\frac{1}{2} \sum \underbrace{(1+y_n) \log(1 + e^{\omega^T x_n})}_{\text{محدب مثبت}} + \underbrace{(1-y_n) \log(1 + e^{\omega^T x_n})}_{\text{محدب مثبت}}$$

محدب است

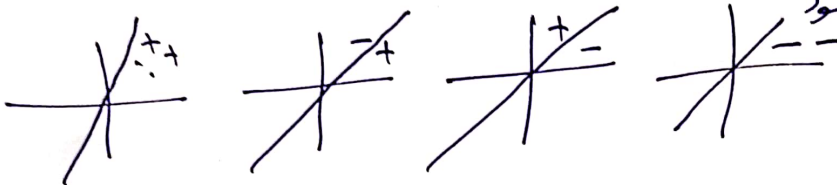
مقدور

زیرا ترکیب خطی هم مقعر هم با مقادیر
 محدب است. و متنی محدب مقعر می شود.

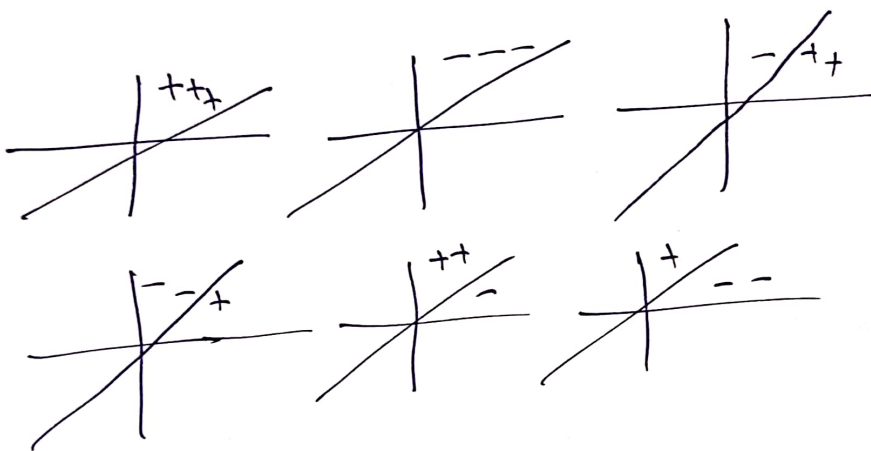
الف) خط با یک نقطه که ۲ حالت می شود



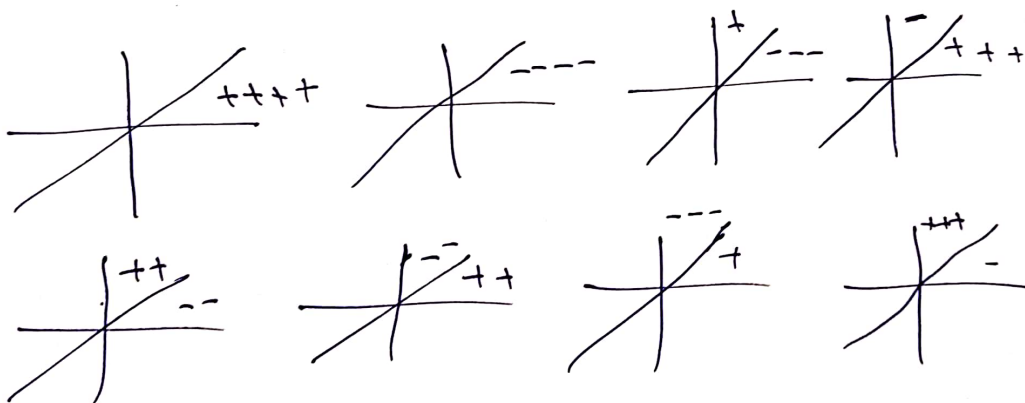
دو نقطه را ۴ حالت می شود



۳ نقطه می شود ۶ حالت



۴ نقطه ۸ حالت است



پس می شود 2^N تابع دارد

(۲) (ب) زیرا قوانین ثابت داریم ~~هرچند~~ هرچند باشد هر نوعی در داشته باشد
 به هر حال ۲ تا افزایش یافته اضافه نمی ترانه بکنند که نهایتاً اگر بشود ۲ تا افزایش می کنیم.

با یک نقطه که دو حالت بسته نمی ترانه بشود

با دو نقطه که ۴ حالت بسته نمی ترانه بشود

با ۳ نقطه ۶ حالت بوده و ۲ حالت اضافه می شود که می شود ۸

با ۴ نقطه ۸ حالت بوده و ۲ حالت اضافه می شود که می شود ۱۰ حالت.

(۲) (ب) پس ~~بسیار~~ بعد از این سؤال می شود ۳ و بزرگترین ۴ می شود $1 \leq 4$
 پس داریم.

$$4m_H(2N)e^{-\frac{1}{8} \epsilon^2 N} \leq 8$$

$$4(2N)^3 e^{-\frac{1}{8} \times \frac{1}{2} N} = 4 \times 8 N^3 e^{-\frac{1}{16} N}$$

$$1 - 32 N^3 e^{-\frac{N}{16}}$$

و احتمال های شود

② ③ در واقع همان سؤال ۲ است معنی وقتی مشتق بگیریم داریم.

$$\frac{\partial |h - y_n|}{\partial y_n} = \begin{cases} -1 & h < y_n \\ \infty & h = y_n \\ 1 & h > y_n \end{cases}$$

$$\frac{\partial \sum |h - y_n|}{\partial y_n} = \cancel{\text{sgn}(h - y_n)} \sum \text{sgn}(h - y_n)$$

که در واقع دارد میانه را حساب می‌کند.

۱ الف ۳) زیرا مجموع مربعات در واقع فاصله را اندازه می گیرد و در مسئله کلاس بندی فاصله مهم نیست مثلا

۱+ و ۵+ هر دو در یک کلاس دسته بندی می شوند و در واقع فاصله مثلا ۱+ و ۱- هیچ بایه نیست از ۱+ و ۱-+ باشد که مجموع مربعات خطا این را خوب انجام نمی دهد.

و البته وقتی کنیم در سوال اشتباه نایی وجود دارد و بایه ~~مجموع مربعات~~ در نظر گرفته شود.
بدی ~~مجموع مربعات~~ این است که

۳ ب) در واقع این روش یک $soft\ threshold$ در نظریه گیری که خیلی بهتر از $hard\ threshold$ یا همان $sgn(x)$ است علت این موضوع آن است که در $cross\ entropy$ ما نیاز به احوال رعلق به صدیک از دسته ها داریم و روش های

$sgn(x)$ رفته صرفا کلاس را برمی گردانده و بصورت دلی بهتر است از $cross\ entropy$ استفاده کنیم.

$cross\ entropy$ استفاده کنیم.

البته مشتق پذیری یا $smooth$ بودن تابع هم مهم است.

$$\frac{-2}{N} \sum_{n=1}^N (x_n - \theta) = 0 \Rightarrow \frac{-2}{N} \sum_{n=1}^N x_n + \frac{2N\theta}{N} = 0$$

$$\theta^* = \frac{\sum_{n=1}^N x_n}{N}$$

باستفاده از این سه مورد



۱. ۲. ۳. ~~۴.~~ شرف

$$\frac{\partial |x_n - \theta|}{\partial \theta} = \begin{cases} 1 & x_n > \theta \\ \infty & x_n = \theta \\ -1 & x_n < \theta \end{cases} = \text{sgn}(x_n - \theta)$$

$$\frac{\partial L_{\text{MAR}}(\theta)}{\partial \theta} = \frac{1}{N} \sum \text{sgn}(x_n - \theta)$$

در واقع داریم به دنبال مقدری می گردیم که دیناست ما دقت
کوچک پس همان میانگین θ^* ما خواهد بود

③ ② ① روش MSE ← در واقع روش خوبی است برای این که مدل آموزش داده شده

داده outlier نداشته باشد زیرا می توان آن فیلتر کرد.

برای آن نیز آن است که اگر یک سری فیلتر داشته باشیم آن توان ۲ خطا را حذف می کند

MAE دقیقاً همان مقادیر MSE را می دهد یعنی خطا را به اندازه خود خطا تأثیر می دهد.

و بنابراین MSE را هم ندارد یعنی اگر داده outlier داشته باشیم آن را حذف نمی کند.

MLE، رای خاص

$$L(w) = \prod_{i=1}^N P(y_i | x_i, w) \Rightarrow \log L(w) = \sum_{i=1}^N \log P(y_i | x_i, w)$$

کدام سوال در واقع می‌توان گفت ادامه سوال ۱ الف است که بایه مشتق بگیریم و به دست
صفر قرار دهیم که در وقت شب!

④ الف) اگر تعداد هیا کم باشد بهتر است لایروزش های ساده تر (مدل های ساده تر) استفاده کنیم و اگر

تعداد هیا زیاد باشد می توان با مدل های پیچیده تر نیز جواب های خوب گرفت اما باید که عدد نسبتاً بزرگ باشد.

④ ب)

~~$$x_1 w_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} -1 \\ 5 \end{bmatrix} \quad x_2 w_2 = \begin{bmatrix} -2 \\ 3 \end{bmatrix} \begin{bmatrix} -2 & 2 \\ 3 & -1 \end{bmatrix} =$$~~

$$w_1 x_1 = \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ 5 \end{bmatrix} \quad w_2 x_2 = \begin{bmatrix} -2 & 2 \\ 3 & -1 \end{bmatrix} \begin{bmatrix} -2 \\ 3 \end{bmatrix} = \begin{bmatrix} 10 \\ -9 \end{bmatrix}$$

$$\max(w_1 x_1, 0) = \begin{bmatrix} 0 \\ 5 \end{bmatrix} \quad \max(w_2 x_2, 0) = \begin{bmatrix} 10 \\ 0 \end{bmatrix}$$

$$\tanh\left(\begin{bmatrix} 0 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 10 \\ 0 \end{bmatrix}\right) = \tanh\left(\begin{bmatrix} 10 \\ 0.5 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ 0.46 \end{bmatrix}$$

$$P_1 = x_1 w_1$$

$$P_2 = x_2 w_2$$

$$P_1 = w_1 x_1$$

$$P_2 = w_2 x_2$$

$$\max P_1 = \max(P_1, 0) \rightarrow \text{Fraction} = \frac{1}{2} \max P_1$$

$$\max P_2 = \max(P_2, 0)$$



$$P = \text{Sum} = \text{Fraction} + \max P_1, \text{ tanResult} = \tanh(\text{sum})$$

$$\frac{\partial f}{\partial \text{Result}} = 1 - \tanh^2 h = \begin{bmatrix} 0 \\ 0.75 \end{bmatrix}$$

$$\frac{\partial f}{\partial \text{Fraction}} = \frac{\partial f}{\partial \text{sum}} \times \frac{\partial \text{sum}}{\partial \text{Fraction}} = 1 - \tanh^2 x^2 \times \frac{\partial f}{\partial P_1} = \frac{\partial f}{\partial \max P_1} \times \frac{\partial \max P_1}{\partial P_1} = \alpha \begin{cases} 0 & P_1 \leq 0 \\ 1 & P_1 > 0 \end{cases}$$

$$\frac{\partial f}{\partial \max P_1} = \frac{\partial f}{\partial \text{sum}} \times \frac{\partial \text{sum}}{\partial \max P_1} = 1 - \tanh^2(x) = \beta$$

$$\frac{\partial f}{\partial P_1} = \frac{\partial f}{\partial \max P_1} \times \frac{\partial \max P_1}{\partial P_1} = \beta \begin{cases} 0 & P_1 \leq 0 \\ 1 & P_1 > 0 \end{cases}$$

$$\frac{\partial f}{\partial \max P_1} = \frac{\partial f}{\partial \text{Fraction}} \times \frac{\partial \text{Fraction}}{\partial \max P_1} = (1 - \tanh^2 x) \times 1 = \alpha$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial P_1} \times \frac{\partial P_1}{\partial x_1} = \frac{\partial f}{\partial P_1} \times w_1 \quad \frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial P_1} \times \frac{\partial P_1}{\partial w_1} = \frac{\partial f}{\partial P_1} \times x_1$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial P_2} \times \frac{\partial P_2}{\partial w_2} = \frac{\partial f}{\partial P_2} \times x_2$$

$$\frac{\partial f}{\partial x_2} = \frac{\partial f}{\partial P_2} \times \frac{\partial P_2}{\partial x_2} = \frac{\partial f}{\partial P_2} \times w_2$$

$$\frac{\partial f}{\partial \max P_1} = \begin{bmatrix} 0 \\ 0.75 \end{bmatrix} \quad \frac{\partial f}{\partial \max P_1} = \begin{bmatrix} 0 \\ 0.75 \end{bmatrix}$$

$$\frac{\partial f}{\partial P_1} = \begin{bmatrix} 0.075 \\ 0 \end{bmatrix} \quad \frac{\partial f}{\partial P_1} = \begin{bmatrix} 0.75 \\ 0 \end{bmatrix}$$

$$\frac{\partial f}{\partial x_1} = \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 0.075 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.22 \\ 0.075 \end{bmatrix}$$

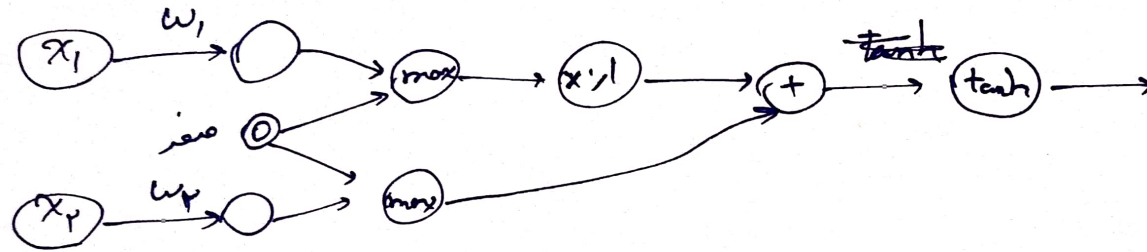
~~$$\frac{\partial f}{\partial x_1} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 0.075 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.075 \\ 0 \end{bmatrix} \begin{bmatrix} 0.75 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 0.75 & 0.15 \\ 0 & 0 \end{bmatrix}$$~~

$$\frac{\partial f}{\partial x_1} = \begin{bmatrix} -2 & 2 \\ 3 & -1 \end{bmatrix} \begin{bmatrix} 0.75 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.15 \\ 0 \end{bmatrix}$$

~~$$\frac{\partial f}{\partial x_1} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 0.75 \\ 0 \end{bmatrix} = \begin{bmatrix} -1.5 \\ 0 \end{bmatrix}$$~~

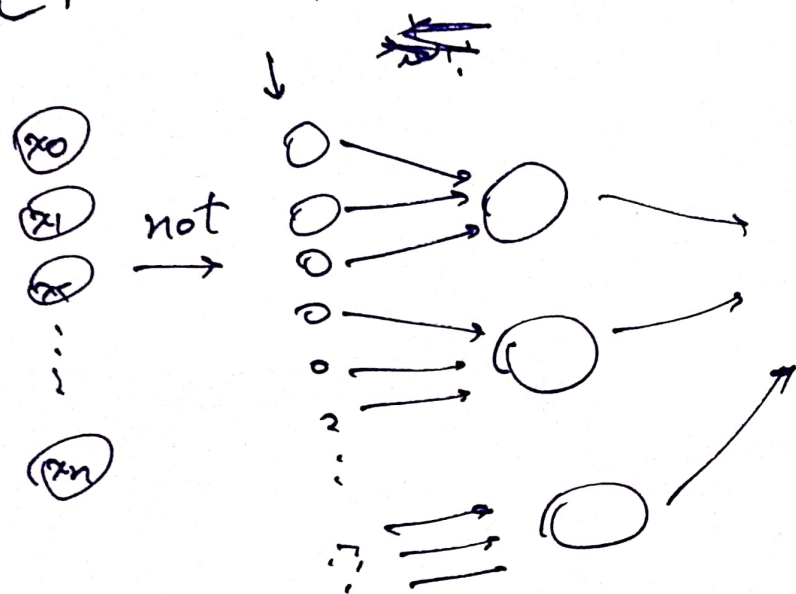
$$\frac{\partial f}{\partial w_1} = \begin{bmatrix} 0 \\ 0.75 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0.75 & 0.15 \end{bmatrix}$$

$$\frac{\partial f}{\partial w_2} = \begin{bmatrix} 0 \\ 0.75 \end{bmatrix} \begin{bmatrix} -2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0.15 & 2.25 \end{bmatrix}$$



② ب. رأف

۴) پ با یک رایه Not هارا ایجاد می کنیم با رایه بعد در وزن ها ضرب می کنیم در رایه بعد هم جمع -
 انجام می دهیم.



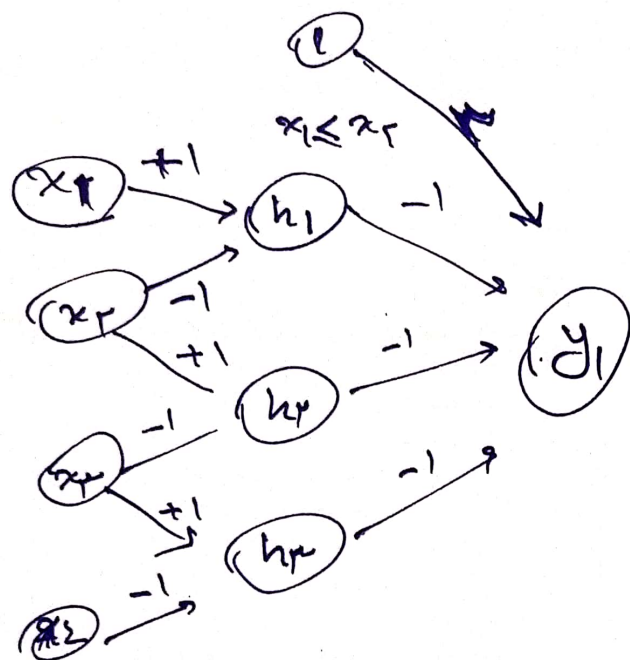
رایه اول Not هارا داریم

در رایه دوم ضرب ها انجام می شود

و در رایه سوم می توان جمع هارا انجام داد.

البته می توان فیوژن لایه ها را حذف کرد و فقط به تعداد لایه های آن لایه ها را حذف کرد و فقط به تعداد لایه های آن لایه ها را حذف کرد.

در لایه مختلف ضرب کرد و جمع هارا اندر لایه های مختلف انجام داد.



(۵) (۳) وقتی ضریب λ یا α زیاد شود $overfitting$ زیاد می شود و $underfitting$ اگر وجود دارد کم می شود که نتیجه می شود که خطای داده های آموزش کم شود یعنی در آموزش بهتر شود و خطای داده های تست زیاد می شود.

(۵) (ب) زیرا در هر مسئله منظم ساز آن فرقی می کند و ما منظم ساز بهینه اصلانه داریم و همچنین یک منظم ساز ساده دیگر مسئله خوب باشد و دیگر مسئله بد باشد و در بهترین حالت نتیجه می شود اگر از منظم ساز به استفاده کنیم λ آن را می نذاریم.

(۵) (۲)

$$L_1 = \sum x_i = 7 \quad L_1(x_2) = 7 \quad L_1(x_3) = 5$$

$$L_2(x_1) = \sqrt{1+4+9} = \sqrt{14} \approx 3.74 \quad L_2(x_2) = \sqrt{6+9} = 5 \quad L_3(x_3) = \sqrt{25} = 5$$

L_1 این طرا انجام می دهد زیرا نرم λ برای بردار x_3 کمتر است و در واقع از آن بردار بیشتر استفاده می شود و به جمع رگولاریزاسیون L_1 این است که x_3 بیشتر از x_1 و x_2 حضور داشته باشد.

نرم λ_2 برای بردار x_1 کمتر است و در واقع منظم ساز L_2 بیشتر از بردار x_1 استفاده می کند.

به صورت کلی هم L_1 متحد تعداد صفرها را بیشتر می کند.

$$\epsilon \sim N(0, B^{-1}) \quad B = \frac{1}{\sigma^2}$$

$$P(w; \alpha) = N(0, \alpha^{-1} I)$$

Posterior \propto likelihood \times Prior

ت و

$$w \sim N(0, \alpha^{-1} I) \quad f(y_k | w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_k - x^T w)^2\right)$$

$$L(w) = \prod_{k=1}^N f(y_k | w) = \prod_{k=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_k - x^T w)^2\right) *$$

رابطه بین ما به این شکل است

$$P(w|D) = \frac{P(D|w) P(w)}{P(D)}$$

$$\Rightarrow \log P(w|D) = \log P(D|w) + \log P(w) - \log P(D)$$

حل: به این معنی می بینیم.

$$w^* = \arg\max_w \log P(D|w) + \log P(w)$$

$$\log P(D|w) + \log P(w)$$

توزیع داده ها توزیع پیش فرض

$$w \sim N(0, \alpha^{-1} I) \Rightarrow f(w) = \frac{\alpha^{D/2}}{(2\pi)^{D/2}} \exp\left(-\frac{\alpha}{2} w^T w\right)$$

$$\log f(w) = \log \alpha^{D/2} - \log(2\pi)^{D/2} - \frac{\alpha}{2} w^T w$$

$$\log P(D|\omega) =$$

حل باید دایمی خود را \max کنیم که در ω به دست اومد

$$\sum \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_k - \bar{x}^T \omega)^2\right) =$$

$$\sum_{k=1}^D \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{k=1}^N (y_k - \bar{x}^T \omega)^2$$

درشتی ω به دست می آید

$$\hat{\omega} = \operatorname{argmax}_{\omega} \left(-\frac{1}{2\sigma^2} \sum (y_k - \bar{x}^T \omega)^2 - \frac{\alpha}{2} \omega^T \omega \right)$$

$$\hat{\omega} = \operatorname{argmin}_{\omega} = -\frac{1}{2\sigma^2} \sum (y_k - \bar{x}^T \omega)^2 + \frac{\alpha}{2} \omega^T \omega$$

$$= -\frac{1}{2} B \sum (y_k - \bar{x}^T \omega)^2 + \frac{\alpha}{2} \omega^T \omega$$

پس داریم