

http://www.work.caltech.edu/~htlin/publication/doc/infkernel.pdf?fbclid=IwAR0y3-BpeP3_nwaz9tLFtJZh3-XiGTFv7M9CvkurYlgkj_MQbJx9tWfav54
<https://pouannes.github.io/blog/initialization/?fbclid=IwAR1JtHrBBNL68MnBHFQ4its6vNzb1kKoNsCVG8pu9IVdqq-rFHZY7tG8WPU>
http://d2l.ai/chapter_multilayer-perceptrons/numerical-stability-and-init.html?fbclid=IwAR11jhydSiNDLjKGn3euWHVIRtku39han8rJd9iimJCUbhVhlBoxZChdLn8

◆ date / ◆ page /

1. Conditioned on $X^{(k-1)}$

$$E[S_j^{(k)}] = E\left[\sum_{i=1}^{d(k-1)} w_{ij}^{(k)} X_i^{(k-1)}\right] = \sum_{i=1}^{d(k-1)} E[w_{ij}^{(k)} X_i^{(k-1)}] = \sum_{i=1}^{d(k-1)} E[w_{ij}^{(k)}] X_i^{(k-1)}$$

$$= \sum_{i=1}^{d(k-1)} 0 \cdot X_i^{(k-1)} = 0 \Rightarrow S_j^{(k)} \text{ are zero mean}$$

let F as the distribution function and f as the density function.

let $W_{jk} = (w_{j1}^k, w_{j2}^k, \dots, w_{jd(k-1)}^k)$

$$F_{j_1, j_2, j_3, \dots, j_k}(t_{j_1}, t_{j_2}, \dots, t_{j_k})$$

$$= P(S_{j_1}^k \leq t_{j_1}, S_{j_2}^k \leq t_{j_2}, \dots, S_{j_k}^k \leq t_{j_k})$$

$$= P(\sum_{i=1}^{d(k-1)} w_{ij}^k X_i^{(k-1)} \leq t_{j_1}, \sum_{i=1}^{d(k-1)} w_{ij_2}^k X_i^{(k-1)} \leq t_{j_2}, \dots, \sum_{i=1}^{d(k-1)} w_{ij_k}^k X_i^{(k-1)} \leq t_{j_k})$$

$w_{j1}, w_{j2}, \dots, w_{jk}$ are independent random vector

$$= P(\sum_{i=1}^{d(k-1)} w_{j1}^k X_i^{(k-1)} \leq t_{j_1}) P(\sum_{i=1}^{d(k-1)} w_{j2}^k X_i^{(k-1)} \leq t_{j_2}) \dots P(\sum_{i=1}^{d(k-1)} w_{jk}^k X_i^{(k-1)} \leq t_{j_k})$$

$$= P(S_{j_1}^k \leq t_{j_1}) P(S_{j_2}^k \leq t_{j_2}) \dots P(S_{j_k}^k \leq t_{j_k})$$

$$= F_{j_1}(t_{j_1}) \cdot F_{j_2}(t_{j_2}) \dots F_{j_k}(t_{j_k})$$

$\Rightarrow f_{j_1, j_2, \dots, j_k}(t_{j_1}, t_{j_2}, \dots, t_{j_k}) = f_{j_1}(t_{j_1}) f_{j_2}(t_{j_2}) \dots f_{j_k}(t_{j_k})$

$\Rightarrow S_1^k, S_2^k, S_3^k, \dots, S_d^k$ are independent.

$$2. \text{Var}[X_i^{(t-1)}] = E[(X_i^{(t-1)})^2] - E[X_i^{(t-1)}]^2$$

$$\Rightarrow E[(X_i^{(t-1)})^2] = 6\bar{x} + \bar{x}^2$$

$$\text{Var}[w_{ij}^{(t)}] = E[(w_{ij}^{(t)})^2] - E[w_{ij}^{(t)}]^2$$

$$\Rightarrow E[(w_{ij}^{(t)})^2] = \sigma_w^2$$

$$\text{Var}[s_j^{(t)}] = E[(s_j^{(t)})^2] - E[s_j^{(t)}]^2$$

$$= E[(\sum_{k=1}^{d-1} w_{ij}^{(t)} X_i^{(t-1)})^2]$$

$$= E[\sum_{i=1}^{d-1} \sum_{k=1}^{d-1} w_{ij}^{(t)} w_{kj}^{(t)} X_i^{(t-1)} X_k^{(t-1)}]$$

$$= E[\sum_{i=1}^{d-1} (w_{ij}^{(t)})^2 (X_i^{(t-1)})^2 + \sum_{i \neq k} w_{ij}^{(t)} w_{kj}^{(t)} X_i^{(t-1)} X_k^{(t-1)}]$$

$$= \sum_{i=1}^{d-1} E[(w_{ij}^{(t)})^2 (X_i^{(t-1)})^2] + \sum_{i \neq k} E[w_{ij}^{(t)} w_{kj}^{(t)} X_i^{(t-1)} X_k^{(t-1)}]$$

$s_j^{(t)}, w_{ij}^{(t)}, X_i^{(t-1)}$'s independence

$$= \sum_{i=1}^{d-1} E[(w_{ij}^{(t)})^2] E[(X_i^{(t-1)})^2] + \sum_{i \neq k} E[w_{ij}^{(t)}] E[w_{kj}^{(t)}] E[X_i^{(t-1)}] E[X_k^{(t-1)}]$$

$$= d-1 \sigma_w^2 (6\bar{x} + \bar{x}^2) + \sum_{i \neq k} 0 \cdot 0 \cdot \bar{x} \cdot \bar{x}$$

$$= d-1 \sigma_w^2 (6\bar{x} + \bar{x}^2)$$

$$3. \text{ReLU}(x) \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} = \max(x, 0)$$

$$E[(x_i^{(t-1)})^2] = E[(\max(s_i^{(t-1)}, 0))^2]$$

let F as the distribution function and f as the density function

$$F_{s_i^{(t-1)}}(t) = P(s_i^{(t-1)} \leq t) = \begin{cases} P(s_i^{(t-1)} \leq t) = F_{s_i^{(t-1)}}(t) & \text{if } t \geq 0 \\ 0 & \text{if } t < 0 \end{cases}$$

$$\Rightarrow f_{s_i^{(t-1)}}(t) = F'_{s_i^{(t-1)}}(t) = \begin{cases} F'_{s_i^{(t-1)}}(t) = f_{s_i^{(t-1)}}(t) & \text{if } t \geq 0 \\ 0 & \text{if } t < 0 \end{cases}$$

$$\Rightarrow E[(x_i^{(t-1)})^2] = \int_{-\infty}^{\infty} t^2 f_{s_i^{(t-1)}}(t) dt = \int_0^{\infty} t^2 f_{s_i^{(t-1)}}(t) dt + \int_0^{\infty} t^2 f_{s_i^{(t-1)}}(t) dt$$

$$= \int_0^{\infty} t^2 f_{s_i^{(t-1)}}(t) dt \quad ; s_i^{(t-1)} \text{ is symmetric.}$$

$$= \frac{1}{2} \int_0^{\infty} t^2 f_{s_i^{(t-1)}}(t) dt$$

$$= \frac{1}{2} E[(s_i^{(t-1)})^2]$$

4. By problem 2, 3, we have:

$$\mathbb{E}[(S_j^{(1)})^2] = 2\mathbb{E}[(X_i^{(1)})^2] = 2(6x^2 + \bar{x}^2)$$

$$\Rightarrow \text{Var}[S_j^{(1)}] = \mathbb{E}[(S_j^{(1)})^2] - \mathbb{E}[S_j^{(1)}]^2 \\ = 2(6x^2 + \bar{x}^2) - 0^2$$

$$\begin{aligned} \text{Var}[S_j^{(1)}] &= \frac{1}{d-1} G_W(6x^2 + \bar{x}^2) \\ &= \frac{1}{d-1} G_W \cdot 2(6x^2 + \bar{x}^2) \\ &= \frac{1}{d-1} G_W \cdot \text{Var}[S_j^{(1)}] \end{aligned}$$

5. Leaky ReLU(x) $\begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases} = \max(x, \alpha x)$

let F as the distribution function and f as the density function.

$$F_{X_i^{(1)}}(t) = P(S_i^{(1)} \leq t) = \begin{cases} P(S_i^{(1)} \leq t) & \text{if } t \geq 0 \\ P(S_i^{(1)} \leq \frac{t}{\alpha}) & \text{if } t < 0 \end{cases}$$

$$\Rightarrow f_{X_i^{(1)}}(t) = F'_{X_i^{(1)}}(t) dt = \begin{cases} F'_{S_i^{(1)}}(t) & \text{if } t \geq 0 \\ F'_{S_i^{(1)}}(\frac{t}{\alpha}) = \frac{1}{\alpha} f_{S_i^{(1)}}(\frac{t}{\alpha}) & \text{if } t < 0 \end{cases}$$

$$\begin{aligned} \Rightarrow \mathbb{E}[X_i^{(1)}]^2 &= \int_0^\infty t^2 f_{X_i^{(1)}}(t) dt + \int_0^\infty t^2 f_{S_i^{(1)}}(t) dt \\ &= \int_0^\infty \frac{t^2}{\alpha} f_{S_i^{(1)}}(\frac{t}{\alpha}) dt + \int_0^\infty t^2 f_{S_i^{(1)}}(t) dt \\ &= \alpha^2 \int_0^\infty t^2 f_{S_i^{(1)}}(t) dt + \int_0^\infty t^2 f_{S_i^{(1)}}(t) dt \end{aligned}$$

$S_i^{(1)}$ is symmetric

$$\begin{aligned} &= \frac{\alpha^2}{2} \int_0^\infty t^2 f_{S_i^{(1)}}(t) dt \\ &= \frac{\alpha^2}{2} \mathbb{E}[(S_i^{(1)})^2] \end{aligned}$$

$$\text{Also by problem 2. } \mathbb{E}[(S_i^{(1)})^2] = \frac{2}{\alpha^2+1} \mathbb{E}[X_i^{(1)}]^2 = \frac{2}{\alpha^2+1} (6x^2 + \bar{x}^2)$$

$$\Rightarrow \text{Var}[S_i^{(1)}] = \mathbb{E}[(S_i^{(1)})^2] - \mathbb{E}[S_i^{(1)}]^2 = \frac{2}{\alpha^2+1} (6x^2 + \bar{x}^2)$$

$$\text{By problem 2 we have. } \text{Var}[S_i^{(1)}] = \frac{1}{d-1} G_W^2 (6x^2 + \bar{x}^2) = \frac{1}{d-1} G_W^2 (\alpha^2 \mathbb{E}[(S_i^{(1)})^2]) \text{Var}[S_i^{(1)}]$$

$$\Rightarrow \text{取 } G_W^2 = \frac{2}{\alpha^2+1}, \text{ and all the receive from p1-p4,}$$

$$\text{we have } \text{Var}[S_i^{(1)}] = \text{Var}[S_j^{(1)}] \neq 1$$

6

$$V_1 = (1-\beta)\Delta t$$

$$V_2 = \beta(1-\beta)\Delta_1 + (1-\beta)\Delta_2$$

$$V_3 = \beta^2(1-\beta)\Delta_1 + \beta(1-\beta)\Delta_2 + (1-\beta)\Delta_3$$

$$\text{guess } V_k = \sum_{i=1}^k \beta^{k-i}(1-\beta)\Delta_i$$

when $k=1, 2, 3$, assumption holds

Assume that when $k=n$, assumption also hold

$$V_n = \sum_{i=1}^n \beta^{n-i}(1-\beta)\Delta_i$$

What $k=n+1$

$$V_{n+1} = \beta V_n + (1-\beta)\Delta_{n+1}$$

$$= \sum_{i=1}^n \beta^{n-i+1}(1-\beta)\Delta_i + (1-\beta)\Delta_{n+1}$$

$$= \sum_{i=1}^{n+1} \beta^{n+1-i}(1-\beta)\Delta_i, \text{ assumption holds!}$$

By mathematical induction,

$$\text{VTE}, V_T = \sum_{i=1}^T \beta^{T-i}(1-\beta)\Delta_i$$

$$\Rightarrow \alpha_t = \beta^{T-t}(1-\beta) \#$$

$$7. \quad 0 < \beta < 1$$

$$\log_2 \beta < 0$$

$$\alpha_1 \leq \frac{1}{2}$$

$$\Rightarrow \beta^{T-1}(1-\beta) \leq \frac{1}{2}$$

$$\Rightarrow \log_2(\beta^{T-1}(1-\beta)) \leq \log_2 \frac{1}{2}$$

$$\Rightarrow (T-1)\log_2 \beta + \log_2(1-\beta) \leq -1$$

$$\Rightarrow (T-1) \geq \frac{-1 - \log_2(1-\beta)}{\log_2 \beta}$$

$$\Rightarrow T \geq \frac{-1 - \log_2(1-\beta)}{\log_2 \beta} - 1$$

$$\text{Ans : } T = \left\lceil \frac{-1 - \log_2(1-\beta)}{\log_2 \beta} - 1 \right\rceil$$

$$8. \hat{X}'_t = \frac{\alpha_t}{\sum_{t=1}^T \alpha_t}$$

$$= \frac{\beta^{T-t}(1-\beta)}{\sum_{t=1}^T \beta^{T-t}(1-\beta)}$$

$$= \frac{\beta^{T-t}(1-\beta)}{(1-\beta) \sum_{t=1}^T \beta^{T-t}}$$

$$= \beta^{T-t} \frac{1-\beta}{1-\beta^T} \neq 1$$

$$\text{q: } 0 < \beta < 1$$

$$\therefore \log \beta < 0$$

$$\alpha_t \leq \frac{1}{2}$$

$$\Rightarrow \beta^{T-1} \frac{1-\beta}{1-\beta^T} \leq \frac{1}{2}$$

$$\Rightarrow \frac{1-\beta^{-1}}{1-\beta^{-T}} \leq \frac{1}{2}$$

$$\Rightarrow 2(1-\beta^{-1}) \leq 1-\beta^{-T}$$

$$\Rightarrow \beta^{-T} \leq 2\beta^{-1} - 1$$

$$\Rightarrow T \ln \beta^{-1} \leq \ln(2\beta^{-1} - 1)$$

$$\Rightarrow T \geq \frac{\ln(2\beta^{-1} - 1)}{\ln \beta^{-1}}$$

$$\text{Ans: } T = \left\lceil \frac{\ln(2\beta^{-1} - 1)}{\ln \beta^{-1}} \right\rceil$$

10.

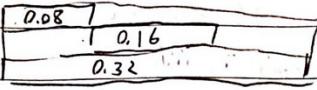
11 We need at least 2 classifier to correctly classify
then G will classify the data wrongly.

Therefore, the lower bound would be 0, as
the three classifier didn't overlap.

The upper bound would be .

With the max overlap,

As a result, $0 \leq E_{\text{out}}(G) \leq 0.24$



12. Because K is odd, we need at least $\frac{K+1}{2}$

classifier to wrongly classify, the G will wrongly classify.

Let the data size as $1, N$:

then we will have $\sum_{n=1}^N$ data wrongly classify by g_K .

Thus, we will have $\sum_{n=1}^K e_n N$ data wrongly classify by g_1, g_2, \dots, g_K

When the datas are correctly classify by G are also correctly

classify by g_1, g_2, \dots, g_K and the datas are wrongly classify by G

are only wrongly classify by $\frac{K+1}{2}$ of the classifier, g_1, g_2, \dots, g_K ,

at this situation the upper bound of the data wrongly classify

would be $\frac{\sum_{n=1}^K e_n N}{\frac{K+1}{2}} = \frac{2}{K+1} \sum_{n=1}^K e_n N$.

\Rightarrow one of the upper bound of $E_{\text{out}}(G)$ is $\frac{2}{K+1} \sum_{n=1}^K e_n N = \frac{2}{K+1} \sum_{n=1}^K e_n$

13 If the results of every sample are independent;

the possibility that at sample $N \geq pN$, the data hadn't been sample once is $(\frac{N-t}{N})^{pN} = (1 - \frac{1}{N})^{pN} \approx e^{-p}$ if N is large enough.
Then, the possibility for the data to get sample at least once is $(1 - e^{-p})^N$. This is the approx results by Watson: $(1 - \frac{1}{N})^{pN} \approx e^{-ap}$.

Let X_i represents that whether the i th data got sample, 1 means it got sample at least once, 0 means it didn't get sample.

We now consider k random variable, where there are s times

succesed, and t times failed, $s+t=k$. Also let $x_{i1}, x_{i2}, \dots, x_{ik} \in \{0, 1\}$

$$\Rightarrow P(X_{i1}=x_{i1})P(X_{i2}=x_{i2}) \dots P(X_{ik}=x_{ik}) \approx (1-e^{-p})^s (e^{-p})^t.$$

$$\Rightarrow P(X_{i1}=x_{i1}, X_{i2}=x_{i2}, \dots, X_{ik}=x_{ik})$$

$$\begin{aligned} &= \left(\frac{N-t}{N}\right)^{pN} - C_1^s \left(\frac{N-(t+1)}{N}\right)^{pN} + C_2^s \left(\frac{N-(t+2)}{N}\right)^{pN} - \dots \\ &+ (-1)^s C_s^s \left(\frac{N-(t+s)}{N}\right)^{pN} \\ &\approx e^{-tp} - C_1^s e^{-c(t+1)p} + C_2^s e^{-c(t+2)p} - \dots + (-1)^s C_s^s e^{-c(t+s)p} \\ &= e^{-tp} (1 - C_1^s e^{-p} + C_2^s e^{-2p} - \dots + (-1)^s C_s^s e^{-sp}) \\ &= e^{-tp} (1 - e^{-ps}) \end{aligned}$$

$$= P(X_{i1}=x_{i1})P(X_{i2}=x_{i2}) \dots P(X_{ik}=x_{ik})$$

Which means that they are independent random variable.

\Rightarrow The expected value that data got sample at least once

$$\text{is } (1 - e^{-p})N = N - (N \cdot e^{-p})$$

14. First consider two special case: all classified as 1 or -1, i.e.
 $\theta \leq L$, $\theta > R$.

Then let us consider cases where $L \leq \theta \leq R$.

Notice that these decision stumps are same in a specific range: $K \leq \theta \leq K+1$, $K = L+1, L+2, \dots, R$.

This means on a special dimension, we have

$2 \cdot (R-L)$ different decision stumps (2 is because $\{1, -1\}$)

For d dimension, we have $2 \cdot (R-L) \cdot d + 2$ different ds.

This is the 2 special case mentioned at the beginning

$$\Rightarrow 2 \cdot (5-0) \cdot 4 + 2 = 42$$

15. First, notice that $K_d(X, X') = |G| - 2 \cdot \text{diff}(X, X')$, where $\text{diff}(X, X')$ means the number of different decision stumps that $g(x) \neq g(x')$, $g \in G$. Consider a dimension i, then

$$\text{diff}(x_i, x'_i) \begin{cases} x_i = x'_i, 0 \\ x_i \neq x'_i, 2 \cdot |x_i - x'_i| \end{cases} \Rightarrow \text{diff}(X, X') = \sum_{i=1}^d 2 \cdot |x_i - x'_i|$$

$$\begin{aligned} \Rightarrow K_d(X, X') &= 2 \cdot (R-L) \cdot d + 2 - \sum_{i=1}^d 2 \cdot |x_i - x'_i| \\ &= 2(R-L)d + 2 - 4 \sum_{i=1}^d |x_i - x'_i| \end{aligned}$$

$$16. \int_L^R g_{\text{bias}}(x) g_{s,\theta}(x') d\theta$$

$$= \int_L^R \text{sign}(x_i - \theta) \text{sign}(x'_i - \theta) d\theta$$

let $\text{sign}(x_i - \theta) = S_i(\theta)$, $\text{sign}(x'_i - \theta) = S'_i(\theta)$

$$= \int_L^R S_i(\theta) S'_i(\theta) d\theta + \int_{\min(x_i, x'_i)}^{\max(x_i, x'_i)} S_i(\theta) S'_i(\theta) d\theta + \int_{\max(x_i, x'_i)}^R S_i(\theta) S'_i(\theta) d\theta$$

when $\theta \in [\min, \max]$ $S_i(\theta) S'_i(\theta) < 0$

else $S_i(\theta) S'_i(\theta) > 0$

$$= \int_L^{\min} (+1) d\theta + \int_{\min}^{\max} (-1) d\theta + \int_{\max}^R (+1) d\theta$$

$$= (\min - L) + (\min - \max) + (R - \max)$$

$$= (R - L) - 2(\max - \min)$$

$$= (R - L) - 2|x_i - x'_i|$$

now, we summation over i ,

$$\sum_{i=1}^d ((R - L) - 2|x_i - x'_i|)$$

$$= d(R - L) - 2 \sum_{i=1}^d |x_i - x'_i|$$

because $s \in \{+1, -1\}$

\Rightarrow we need to d times ?

$$2d(R - L) - 4 \sum_{i=1}^d |x_i - x'_i| \neq$$

...

17. Actually, I like all of them. Prof. Lin, uses the simplest way to let us understand the hardest concepts. If I have to choose one, I would have to say the newly added course "Activation In Deep Learning". It gave me a better understand on the other ML courses I taken. Prof Lee's ML course is more on the application aspect, where as Prof Lin is more theoretical. With the combination of the two courses, I think I reached a higher level on ML.

18. It would also be the newly added course of Deep Learning. This is because the videos often have background noises, which made learning even harder. Also, personally, I like to review/preview with courses' powerpoint, I'm not used to handwritten notes.