

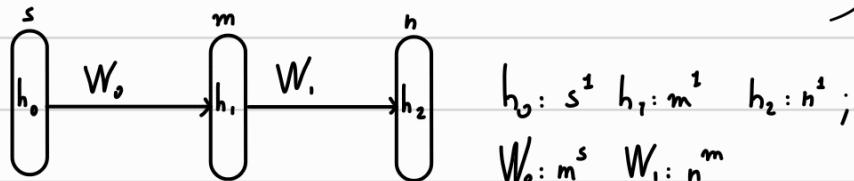
Ideas to solve the instability issues in predictive coding networks

Luca Pinchetti - 22/07/2024

! • GENERAL NOTE I DIDN'T KNOW WHERE TO PUT:

In the following document, I used the same network structure for both bp and pc, this means that the term " h_i " has a different meaning, as in bp it is the activation of a layer, while in pc it is its state. Furthermore, it is quite implicit at the beginning, but in pc we are requiring that $V[h_i] = V[p_i]$

SETUP: Linear network with 2 layers (no bias + no act f_n), gaussian sampling for weights.



$$h_0: s^1 \quad h_1: m^1 \quad h_2: n^1;$$

$$W_0: m^s \quad W_1: n^m$$

$$W: m^{\frac{m}{n}} = n^{\frac{m}{n}}$$

↓
input dim
output dim
dim

NOTATION: a $n \times m$ matrix W is defined as $W: m^{\frac{m}{n}} = n^{\frac{m}{n}}$, such that a vector is $x: m^1$ and $y = Wx: n^{\frac{m}{n}} = n^1$

Backprop: ① $V[h_0] = V[h_1] = V[h_2]$ ② $V\left[\frac{\partial F}{\partial h_0}\right] = V\left[\frac{\partial F}{\partial h_1}\right] = V\left[\frac{\partial F}{\partial h_2}\right] = k$ // my understanding is that $V[x]$ means all the elements of x have the same variance $V[x]$
(Xavier et al., 2010) - since $E[h_i] = \phi$; $E[W_i] = \phi$ and $h_i = W_i h_0 \Rightarrow V[W_i h_0] = V[h_i]$ // ① ②: the intuition is that we want constant "scale of information" through the network, so no exploding or vanishing activations / gradients

$$\text{then } V[W_0] = V[h_1] / sV[h_0] = \frac{1}{s} \Rightarrow \sigma_{W_0} = \frac{1}{\sqrt{s}}$$

1 ⇒ = ○

$$\text{using that } V[Ax] = aV[A]V[x] \text{ with } A: b^a$$

$$\text{and similarly, } V[W_0] = \frac{1}{m}. \text{ Then, we can do a similar analysis for the gradients.}$$

PREDICTIVE CODING

- same analysis holds for forward pass, having that if $V[W_i] = \frac{1}{d_i}$ then $V[h_i] = 1$, with $W: d_i^{\frac{d_i}{n}}$, which satisfies ①

- HOWEVER, this does NOT seem true for the gradients, and ② is not satisfied.

$$\text{ANALYSIS: } F: \frac{1}{2} (W_0 h_0 - h_1)^2 + \frac{1}{2} (W_1 h_1 - h_2)^2 \Rightarrow \frac{\partial F}{\partial h_1} = h_1 - W_0 h_0 + W_1^T (W_1 h_1 - h_2)$$

$$\text{thus } (1) V\left[\frac{\partial F}{\partial h_1}\right] = V[W_0 h_0] + V[(W_1^T W_1 + I) h_1] + V[W_1^T h_2]$$

*: using $V[X+Y] = V[X] + V[Y] + 2\text{Cov}[X,Y]$
// I'm not 100% sure about $\text{Cov}[A,B]$ bc since W_1 appears in both, but the final formula is correct and since we multiply W_1 by h_1 in A and by h_2 in B, they should be completely random.

$$= V[h_1] + (m-1)nV[W_1]V[h_1] + V[h_1](2nV[W_1]^2 + 1 + 2nV[W_1] + n^2V^2[W_1]) + nV[W_1]V[h_2]$$

is true only if we are sampling from a Normal distribution.

Explanation for ...: analysis of $W_1^T W_1 + I$

$$\begin{matrix} n \\ m \end{matrix} \begin{matrix} m \\ n \end{matrix} + \begin{matrix} \frac{1}{n} & \dots \\ \dots & \frac{1}{n} \end{matrix} = \begin{matrix} \bullet & \bullet \\ \bullet & \bullet \end{matrix} \quad \text{with } E[\bullet] = \phi, V[\bullet] = nV[W_1]^2$$

$$E[\bullet] = nV[W_1] + 1, V[\bullet] = 2nV[W_1]^2 \quad \text{using } \square = D+1$$

1) Each off-diagonal element is obtained as $\bullet = \sum_{i=1}^n x_i y_i$ where $x_i \in W_1^T$ and $y_i \in W_1$ are different elements of W_1 for all i , so they are independent

$$\Rightarrow V[\bullet] = nV[W_1]^2, E[\bullet] = \phi$$

*: using $V[x]V[y]V[z]V[w] = V[x]V[y]V[z]V[w] + V[x]V[y]V[w]V[z]$
and $E[xy] = E[x]E[y]$ if x, y independent

2) Each diagonal element is obtained as

$$\square = \sum_{i=1}^n x_i^2 \quad \text{since we pick the same element}$$

in both W_1^T and W_1 . Since x_i is Gaussian r.v.

$$\text{we have } E[x_i^2] = V[x] + E[x]^2 = V[W_1] \Rightarrow$$

$$E[\square] = nV[W_1] \text{ and } V[\square] = 2nV[W_1]^2$$

check online for this

thus we compute $V[A h_i]$ as ... Using *, with $A = W_i^T W_i + I$.

Having $V\left[\frac{\partial F}{\partial h_i}\right]$, let's try to satisfy $V[h_i] = 1$ and $V\left[\frac{\partial F}{\partial h_i}\right] = 1$

formally I required all of them to be equal to constants j and k, but let's just use 1.

$V[h_i] = 1 \Rightarrow V[W_i] = \frac{1}{d_i}$ as shown initially. Let's substitute:

$$V\left[\frac{\partial F}{\partial h_i}\right] = 1 + (m-1) \underbrace{\frac{1}{m^2}}_{n} \cdot 1 + 2n \underbrace{\frac{1}{m^2}}_{n} + 1 + n^2 \underbrace{\frac{1}{m^2}}_{n} + 2n \cdot \underbrace{\frac{1}{m}}_{n} + n \cdot \underbrace{\frac{1}{m}}_{n} \cdot 1 = 2 + \frac{(m-1)n}{m^2} + \frac{2n}{m^2} + \frac{n^2}{m^2} + 3 \frac{n}{m}$$

$$= 2 + \frac{n(m+1)}{m^2} + \frac{3n}{m} + \frac{n^2}{m^2} \neq 1 \quad // \text{This formula has been verified numerically.}$$

- $E[\square]$ and $V[\square]$ play a major role in the formula above. If somehow we assume we manage to "break the symmetry" of $W_i^T W_i$, we have that $E[\square] = 1$ and $V[\square] = n V[W_i]^2$. Then $V\left[\frac{\partial F}{\partial h_i}\right]$ becomes:

$$\tilde{V}\left[\frac{\partial F}{\partial h_i}\right] = 1 + (m-1) \underbrace{\frac{n}{m^2}}_{n} + \underbrace{\frac{n}{m^2}}_{n} + 1 + \underbrace{\frac{n}{m}}_{n} = 2 \left(1 + \frac{n}{m}\right) \neq 1$$

which looks better than the above.

and the original (w/out substitution) becomes:

$$\tilde{V}\left[\frac{\partial F}{\partial h_i}\right] = V[h_i] + (m-1) n V[W_i] V[h_i] + V[h_i] (n V[W_i]^2 + 1) \quad) + n V[W_i] V[h_i]$$

- QUESTION: can an activation function mitigate such symmetry?

Now we don't have any, using one should lead to something like $W_i^T \sigma(W_i h_i)$ instead of $W_i^T W_i h_i$. Let's briefly consider having

$$F = \frac{1}{2} (W_0 h_0 - h_i)^2 + \frac{1}{2} \left[\sigma(W_i h_i) - h_2 \right]^2 \Rightarrow \frac{\partial F}{\partial h_i} = h_i - W_0 h_0 + W_i^T \left[(\sigma(W_i h_i) - h_2) \odot \sigma'(W_i h_i) \right]$$

$$= h_i W_0 h_0 + W_i^T \sigma(W_i h_i) \sigma'(W_i h_i) - W_i^T h_2 \odot \sigma'(W_i h_i)$$

QUESTION: is this enough for independence?

LET'S GO BACK TO THE LINEAR CASE.

• QUESTION: can we find different initialisations/energy functions that satisfy ① & ②?

my ANSWER: ① $\Rightarrow V[h_i] = V[h_{i+1}] = V[W; h_i] \Rightarrow V[h_i] = d_i V[W] V[h_i]$, with $W_i := d_i$

$\Rightarrow 1 = d_i V[W] \cdot 1 \Rightarrow V[W] = \frac{1}{d_i}$ (which we know) // is this correct? - what about havin an act fn?

\Rightarrow The initialisation is forced

\Rightarrow we can try to act on the energy:

Let's assume $\frac{\partial F}{\partial h_i} = \alpha(h_i - W_0 h_0) + \beta(W_i^T (W_i h_i - h_0))$,

as it is the only reasonable modification of the energy fn I could think of, with $\alpha = f(m)$ and $\beta = f(n)$. Then:

$$V\left[\frac{\partial F}{\partial h_i}\right] = V[\alpha W_0 h_0] + V\left[\beta(W_i^T W_i + 1)h_i - W_i^T h_0\right] = \alpha^2 V[h_i] + \beta^2 n V[W_i] V[h_i] + \beta^2 (m-1) n V[W_i]^2 V[h_i] + \beta^2 \left(2n V[W_i]^2 + \alpha^2 + n^2 V[W_i]^2 + 2\alpha n V[W_i]\right) V[h_i]$$

$E[\square] = n V[W_i] + \alpha$

// indicates terms that disappear

when we "break the symmetry".

• if $V[h_i] = 1 \Rightarrow V[W_i] = \frac{1}{m}$

$$\Rightarrow V\left[\frac{\partial F}{\partial h_i}\right] = \alpha^2 + \beta^2 \frac{n}{m} + \beta^2 \frac{(m-1)}{m} \frac{n}{m} + \beta^2 \frac{2n}{m^2} + \beta^2 \frac{\alpha^2}{m^2} + \beta^2 \frac{n^2}{m^2} + 2\alpha \beta^2 \frac{n}{m}$$

--- terms imply that α and β must be constants \Rightarrow also this doesn't work.

so: we CAN'T satisfy both ① and ② // are my conclusions correct?

NOTES: $\alpha = \frac{1}{m}$ and $\beta = \frac{1}{n}$ would guarantee $V\left[\frac{\partial F}{\partial h_i}\right] < k$ (so no exploding grads, but still vanishing)

• if we "break the symmetry", with $E[\square] = \alpha$ and $V[\square] = n V[W_i]^2$, then

$$V\left[\frac{\partial F}{\partial h_i}\right] = \alpha^2 + \alpha^2 \beta^2 + \beta^2 \frac{n}{m} + \beta^2 \frac{n}{m} = \alpha^2 \left(1 + \beta^2\right) + 2\beta^2 \frac{n}{m}$$

\Rightarrow still no solution (as α, β must be constant)

SO WE NEED TO FIND DIFFERENT CONSTRAINTS THAT GUARANTEE STABLE TRAINING!!!

• QUESTION: can we rewrite the constraints as:

$$(1) \quad V[h_{i+1}] = V[W_i h_i] \quad // \text{activation and state have the same } V$$

$$(2) \quad V\left[\frac{\partial F}{\partial h_i}\right] = V[h_i] \quad // \text{state and its error have the same } V.$$

$$(3) \quad d_i = d_j \Rightarrow a < V[h_i] = V[h_j] < b, \quad a, b \in \mathbb{R}^+$$

// if two layers have the same size, they must have the same variance, this is to write that the variance is a property of the layer and it doesn't depend on its position in the network (maybe there are better ways to write this).?

• do they achieve analogous training stability?

INTUITION: we don't need the variance to be GLOBALLY the same, but just to be LOCALLY stable: (1) guarantees that the incoming activation V matches the state V , such that they both have the same "scale"; (2) ensures that the same is valid for the gradient (so that the error is proportional to state, so it can't vanish or explode); finally (3) ensures V of a layer is bounded and non-vanishing.

Are these constraint correct? Am I missing something from Xavier 2010?

• QUESTION: can we satisfy these constraints?

my ANSWER: be $V[h_i] = \frac{1}{\sqrt{d_i}}$ // satisfies (3)

$$\stackrel{\text{then}}{\Rightarrow} d_i V[W_i] \cdot V[h_i] = V[h_{i+1}] \Rightarrow V[W_i] = \frac{\sqrt{d_i}}{\sqrt{d_{i+1} \cdot d_i}} = \frac{1}{\sqrt{d_{i+1} \cdot d_i}} \quad // \text{satisfies (1)}$$

$$\text{we have that: } V\left[\frac{\partial F}{\partial h_i}\right] = \underbrace{\alpha^2}_{\textcolor{purple}{m}} \cdot \underbrace{\frac{1}{\sqrt{m}}} + \underbrace{\beta^2}_{\textcolor{blue}{n}} \cdot \underbrace{\frac{1}{\sqrt{n}}} \cdot \underbrace{\frac{1}{\sqrt{m \cdot n}}} + \underbrace{\beta^2}_{\textcolor{red}{(m-1)}} \cdot \underbrace{n}_{\textcolor{brown}{m \cdot m}} \cdot \underbrace{\frac{1}{\sqrt{m}}} + \underbrace{2\beta^2}_{\textcolor{yellow}{n}} \cdot \underbrace{\frac{1}{\sqrt{m}}} \cdot \underbrace{\frac{1}{\sqrt{m}}} + \underbrace{\beta^2}_{\textcolor{orange}{m}} \cdot \underbrace{\frac{\alpha^2}{\sqrt{m}}} + \underbrace{\beta^2}_{\textcolor{green}{n}} \cdot \underbrace{\frac{n^2}{\sqrt{m \cdot n}}} + \underbrace{2\alpha\beta^2}_{\textcolor{red}{m}} \cdot \underbrace{\frac{n}{\sqrt{m \cdot n}}}$$

$$= \frac{1}{\sqrt{m}} \left(\alpha^2 + \beta^2 + \beta^2 \frac{(m-1)}{m} + \beta^2 \alpha^2 + \beta^2 \frac{n^2}{m} + 2\alpha\beta^2 \frac{n}{\sqrt{m \cdot n}} \right) \neq V[h_i] = \frac{1}{\sqrt{m}}$$

Thus it is NOT a solution.

HOWEVER, if we "break the symmetry": (remove — from above)

$$\tilde{V}\left[\frac{\partial F}{\partial h_1}\right] = \frac{1}{\sqrt{m}} \left(\alpha^2 + 2\beta^2 + \beta^2 \alpha^2 \right) = V[h_1] \text{ if } \alpha^2 + 2\beta^2 + \beta^2 \alpha^2 = 1 !!$$

// note that there may be an error in the derivation, since in the numerical solution I have that $\alpha = \beta = \frac{1}{2}$ works (however the numerical solution uses a different setup,) so the above could be right.

so we (ALMOST) FOUND A SOLUTION!

↓ since for example if $d_i = d_{i+1} \Rightarrow \alpha = \beta \Rightarrow 3\alpha^2 + \alpha^4 = 1$

$$\Rightarrow \alpha^2 = \frac{-3 \pm \sqrt{9+4}}{2} \text{ is clearly not a realistic solution}$$

\Rightarrow there's likely a small error somewhere.

(but let's ignore it for a second)

QUESTION: How do we "break the symmetry"? 3 possibilities:

- using an activation function? (to be tested)
- using a separate set of weights E for the backward pass instead of W^T (see Neural Generative Coding by Ororbia)?
- Using complex numbers!!

Apparently if $x \in \mathcal{N}(\phi, k+j)$ not sure how to write "sampled from a normal distribution with ϕ mean k variance for both real and imaginary part"

$\Rightarrow E[x^2] = \phi$ which we can use to have $E[\square] = \phi$
so to "break the symmetry".

In particular (numerically tested), using the following satisfies (1), (2) and (3): check py notebook

- $V[h_i] = \frac{1}{\sqrt{d_i}}, h_i: d_i^1$

..... : factors probably due to the complex numbers

- $V[W_i] = \frac{1}{2\sqrt{d_i \cdot d_{i+1}}}, W_i: d_i^{d_{i+1}}$

// All of these was done assuming normal distributions, slightly different results may be true for Uniform distributions.

- $\alpha = \beta = \frac{1}{2}$

// and of course this holds only for linear networks.

DISCLAIMER: everything related to complex numbers was found via trial and error, so not sure the explanation I gave are correct.

QUESTION: is this solution unique? Does it apply to actual networks?
can we get rid of complex numbers?

DISCUSSION:

or a constant, but it shouldn't matter.

• Xavier 2010 requires that $V[h_i]$ (and thus $V[\frac{\partial F}{\partial h_i}]$) = 1, this implies that the "energy" of a layer, defined as $E = \sum_{j=1}^{d_i} (h_i - W_{i,j} \cdot h_{i-1})_j^2$, scales with the number of neurons in the last layer.

In particular, $V[E] = d_i V[h_i]^2$, which explains (finally!!) why we use MEAN SQUARE ERROR in the last layer as loss $\mathcal{L} = \frac{1}{2d_i} \sum_{j=1}^{d_i} (h_i - W_{i,j} \cdot h_{i-1})_j^2$. (instead of just square error), since

$V[\mathcal{L}] = \frac{1}{d_i} V[E] = V[h_i]^2$ i.e. the total loss doesn't depend on the size of the layer.

- Having $V[h_i] = \frac{1}{\sqrt{d_i}}$ (and I guess $\alpha = \beta = \frac{1}{2+\frac{1}{\delta}}$ for complex numbers)

achieves the same result, since: $V[E] = d_i V^2[h_i] = d_i \cdot \left(\frac{1}{\sqrt{d_i}}\right)^2 = 1$

and we don't need the "mean" of the square error, justifying the difference between the loss function of a bp and a pc trained network

(i.e., if we assume $V[h_i] = \frac{1}{\sqrt{d_i}}$, then both bp and pc use SQUARE ERROR loss, without any MEAN). \hookrightarrow finally!!

- We are effectively scaling the energy function, since we are going from $V[E_i] = d_i$ to $V[E_i] = 1$. Is this the solution to the "we need a custom optimizer to rescale the energy function" problem we all have been talking about for ~1 year?

- This does intuitively make sense, as each layer now has a fixed energy magnitude (i.e., variance) independently from the number of neurons; that is, we do not generate new energy from thin air by simply adding new neurons to a layer, but, instead, existing neurons have to share the fixed total energy. This seems reasonable if we assume a system with a finite amount of energy like the brain, or, maybe, analog hardware.

- IMPORTANT: using the conditions above ($V[h_i] = \frac{1}{\sqrt{d_i}}$, $V[W_i] = \frac{1}{2\sqrt{d_i d_{i+1}}}$, $\alpha = \beta = \frac{1}{2}$)

also satisfies a fourth constraint that we may find desirable:

- (4) pre and post synaptic errors have the same variance, namely, given $e_{1<} = (W_0 h_0 - h_1)^2$ and $e_{1>} = (W_1 h_1 - h_2)^2$, we have that $V[e_{1<}] = V[e_{1>}]$.
 (see latest sections of the py notebook)

* UPDATE 23/07/2024

- Using a different matrix for the error connections (E_i instead of W_i^T) seems to solve all the problems since it sets $V[\square] = V[\star]$ and $E[\square] = E[\bullet]$. Furthermore it works even with $V[h_i] = 1$ and $V[W_i] = \frac{1}{d_i}$, using $V[E_i] = \frac{1}{d_{i+1}}$ which is the original Xavier method (but also with the new defined d_i variances).

=> the question becomes: "complex numbers or asymmetric weights"? They both require doubling the amount of parameters, so the question is which one trains faster, or other advantages?

=> I found this overall very interesting as PC was proposed as a bio inspired learning algorithm, and here we proved that it cannot work with symmetric backward weights, which is perfectly inline with the asymmetries found in the human brain.

- Using $V[h_i] = \frac{1}{k\sqrt{d_i}}$, k doesn't seem to be important, since it can always be balanced by a coefficient in the energy function. Is there any difference between different K values? How to choose?

* UPDATE 30/07/2024

• FURTHER ANALYSIS OF REC-LRA

- as mentioned, in rec-LRA (and NGC which can be seen as a derivation of it), we introduce an error matrix E_i to backpropagate the error. Further analysis of the algorithm suggests that the goal of E_i is simply to behave as W^T , since $\nabla E = \nabla W^T$. The question then becomes how impactful is the

Luca Poldi
22/07/2024

different initialisation of the two matrixes (W and E). Experimentally, it seems like the norm of the two matrices significantly grows with training, suggesting that the initial point is quite irrelevant. However, the gradient analysis proposed in this document refers to the initial training steps, during which E and W^T are actually different (but also random, so not very good at training the model perhaps).

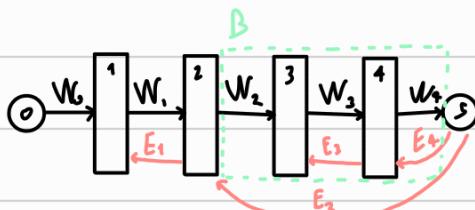
However, there (probably) are alternative ways we can define $\frac{\partial f}{\partial h}$ such that the error is propagated by an auxiliary error matrix E , such that $\frac{\partial F}{\partial E}$ is well defined:

$$1) \tilde{F}_i = \frac{1}{2} \left(W_i h_{i+1} - h_i \right)^2 \Rightarrow \frac{\partial \tilde{F}_i}{\partial h_i} = (h_i - p_i) + (h_i - E_i h_{i+1})$$

$$2) \frac{\partial \tilde{F}_i}{\partial h_i} = (h_i - p_i) + E_{i+1} (p_{i+1} - h_{i+1}) \Rightarrow \frac{\partial \tilde{F}}{\partial E_{i+1}} ? \quad E \text{ metalearning?}$$

- Analysis of ∇E .

In rectra, we have some "Error Skip Connections" as following:



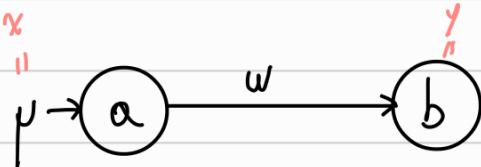
As just described, $\nabla E_i \cdot \nabla W_i^T$ which basically implies that E_1, E_3, E_5 are simply W_1^T, W_3^T, W_5^T and we may not even need them. But what about E_2 ? Clearly E_2 doesn't correspond to W_2^T . My intuition is that we can approximate B as a linear transformation represented by \tilde{W}_2 and assume that $p_5 = \tilde{W}_2 h_2$ and thus $\nabla E_2 = \nabla \tilde{W}_2^T$.

This is effectively a error approximation method. Can then refinement be achieved with inference? What about more complex approximations? (High order or MLPs?)
or CNNs or other

However, this interpretation leads to many different implementations:

- is the error carried by E_2 complementary or alternative to the error carried from h_3 to h_2 ? (Could E_2 be an "initialisation value" for the error node?)
- is B actually representing \tilde{W}_2 , or do we need a forward skip connection as well?

$$x = -1 \quad \begin{cases} \text{Task: invert the input!} \\ y = 1 \end{cases}$$



• For bp is trivial!

OBVIOUS SOLUTION: $w = -1$, HOWEVER, LET'S SET $w = 1$

• for pc it doesn't work at all

\Rightarrow UNLESS: inference is not performed to convergence, but stopped way before!!!

It explains: why a low T works better !!

If convergence is achieved (assuming a minimum of noise)

(but only in this particular example noise is necessary)

w will keep increasing

until it diverges !! \Rightarrow exactly what we experienced.

("summer vacation" $\uparrow \downarrow$)

\Rightarrow might explain the success of iPC: trains giving more emphasis to the forward info.

\Rightarrow backward inference? very likely affected, might not be solution tho.

Brief Analysis:

$$2F = (a-x)^2 + (b-wa)^2, \text{ if } x = -1, y = 1$$

$$\Rightarrow 2F = (a+1)^2 + (1-wa)^2$$

$$\Rightarrow \frac{\partial F}{\partial a} = a+1 - w(1-wa) \cdot a+1 - w + w^2 a$$

$$\frac{\partial F}{\partial a^*} = \Rightarrow a^* = \frac{w-1}{w^2+1} \quad \left\{ \begin{array}{l} \text{if } w \geq 1, \text{ then} \\ a^* > 0 \end{array} \right.$$

$$\Rightarrow \frac{\partial f}{\partial w} = -\vec{a}^*(1-w\vec{a}) = -\vec{a}^* \cdot \vec{e}^* \quad \vec{e}^* = 1-w\vec{a}^* = 1-\underline{w(w+1)} \frac{w^2+1}{w^2+1}$$

$$= -\frac{w-1}{w^2+1} \cdot \frac{w+1}{w^2+1} = -\frac{w^2-1}{w^2+1} \left\{ \begin{array}{l} \text{if } w > 1, \text{ then} \\ \frac{\partial f}{\partial w} < 0 \end{array} \right.$$

w will keep increasing from the value of 1 for ever, instead of approaching the solution $w = -1$.

$$\frac{e^{i\pi w} - 1}{e^{2i\pi w} + 1}$$

$$(a+1)^2 + (1-e^{i\pi w} a)^2$$

$$x = \boxed{\frac{w-1}{w^2+1}}$$

⑤