



Ref sheet

Taylor expansions


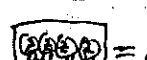
Isserlis theorem

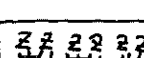
$Z_\alpha \sim \mathcal{N}(0, G^{\alpha\beta})$  = Π

 = $\Pi\Pi + \Pi\Pi + \Pi\Pi$ = $3\Pi\Pi$

 = $\Pi\Pi\Pi + \dots$ (15 terms)

$Z_{i\alpha} \sim \mathcal{N}(0, S_{ij} G^{\alpha\beta})$

 = 2Π |  = $n^2[\Pi\Pi] + n[\Pi\Pi + \Pi\Pi]$

 = $n^3[\Pi\Pi\Pi]$ + $n^2[\dots]$

$\frac{1}{1+x} = 1 - x + x^2 - x^3 \dots$

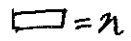
$\ln(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 \dots$

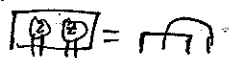

$e^x = 1 + x + \frac{1}{2}x^2 + \frac{1}{6}x^3 \dots$

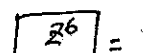
$(1+x)^n = 1 + nx + \frac{n(n-1)}{2}x^2 + \frac{n(n-1)(n-2)}{6}x^3 \dots$

$\sqrt{1+x} = 1 + \frac{1}{2}x - \frac{1}{8}x^2 + \frac{1}{16}x^3 \dots$


Kronecker delta: δ_{ij}

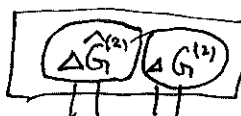
 = n

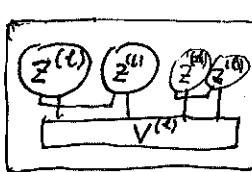
 = n^2 |  = $n^2 + n^2 + n^2$

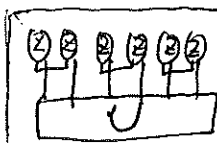
 = $n^3 + n^3 + n^3$ + $n^3 + n^3 + n^3$ + $n^3 + n^3 + n^3$ + $n^3 + n^3 + n^3$ = $8n^3$

field fluctuation

 = 0

 = $\frac{1}{n_i} \langle V^{(2)} \rangle$

 = $n_i^2 \langle V^{(3)} \rangle + 2n_i \langle V^{(2)} \rangle$

 = $n_i^3 \langle V^{(4)} \rangle + 6n_i^2 \langle V^{(3)} \rangle + 8n_i \langle V^{(2)} \rangle$

$\hat{G}_{\alpha\beta}^{(L+1)} = C_b^{(L+1)} + C_w^{(L+1)} \frac{1}{n_c} \sum_i \sigma_{i\alpha}^L \sigma_{i\beta}^L = C_b^{(L+1)} + C_w^{(L+1)} \langle \hat{\sigma}_\alpha^L, \hat{\sigma}_\beta^L \rangle$

$\frac{1}{n_c} V_{\alpha\beta}^{(L+1)} = (C_w^{(L+1)})^2 \text{Cov}[\langle \hat{\sigma}_\alpha^L, \hat{\sigma}_\beta^L \rangle, \langle \hat{\sigma}_\alpha^L, \hat{\sigma}_\beta^L \rangle]$

quartic action

$S^{(2)} = \frac{1}{2} \langle z^{(2)} z^{(2)} \rangle - \frac{1}{8} \langle z^{(4)} z^{(4)} \rangle + \dots$

$\langle e^{\frac{1}{8n_c} z z z z V} \rangle_{G^{(4)}} = 1 + \frac{n_c^2}{8} \langle V^{(4)} \rangle + \dots$

$\langle F \rangle_{S^{(4)}} = \langle F \rangle_{G^{(4)}} + \frac{1}{8} \text{Cov}_{G^{(4)}} [F, z z z z V] + \dots$

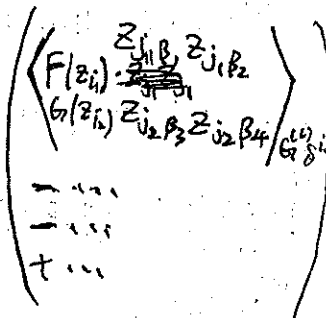
$\text{Cov}_{S^{(4)}} [F, V] = \text{Cov}_{G^{(4)}} [F, V] + \frac{1}{8} [\langle F z z z z V \rangle - \langle F \rangle \langle z z z z V \rangle - \langle G \rangle \langle F z z z V \rangle + 2 \langle F \rangle \langle G \rangle \langle z z z z V \rangle - \langle F G \rangle \langle z z z z V \rangle]$

If $F(z_{i_1}^{(c)})$ and $G(z_{i_2}^{(c)})$

depend on two different neurons $i_1 \neq i_2$

then

$$\text{Cov}_{z^{(c)} \sim S^{(c)}} [F, G] = \frac{1}{8} \left[\langle F G z^{\dagger} \rangle_{G(z_{i_1}^{(c)})} - \langle F \rangle \langle G z^{\dagger} \rangle \right. \\ \left. - \langle G \rangle \langle F z^{\dagger} \rangle + \langle F \rangle \langle G \rangle \langle z^{\dagger} \rangle \right] V^{(c)} = \frac{1}{8} \sum_{\beta_1, \beta_2, \beta_3, \beta_4} V^{(\beta_1, \beta_2)(\beta_3, \beta_4)}$$



By enumeration, the summation terms are zero, except when giving us just 8 terms, which ~~are~~ then gives

$$i_1 = j_1, i_2 = j_2 \quad \text{or} \quad i_1 = j_2, i_2 = j_1$$

$$\text{Cov}_{z^{(c)} \sim S^{(c)}} [F(z_{i_1}^{(c)}), G(z_{i_2}^{(c)})] = \frac{1}{4} \sum_{\beta_1, \beta_2, \beta_3, \beta_4} V^{(\beta_1, \beta_2)(\beta_3, \beta_4)} \left\langle F(z_{A_1}^{(c)}) (z_{\beta_1} z_{\beta_2} - G_{\beta_1, \beta_2}^{(c)}) \right\rangle_{G^{(c)}} \\ \left\langle G(z_{A_2}^{(c)}) (z_{\beta_3} z_{\beta_4} - G_{\beta_3, \beta_4}^{(c)}) \right\rangle_{G^{(c)}}$$

more succinctly,

~~$$\text{Cov}_{S^{(c)}} [F(z_{i_1}^{(c)}), G(z_{i_2}^{(c)})] = \frac{1}{4} V \left\langle F_{A_1}(z) \right\rangle \left\langle G_{A_2}(z) \right\rangle$$~~

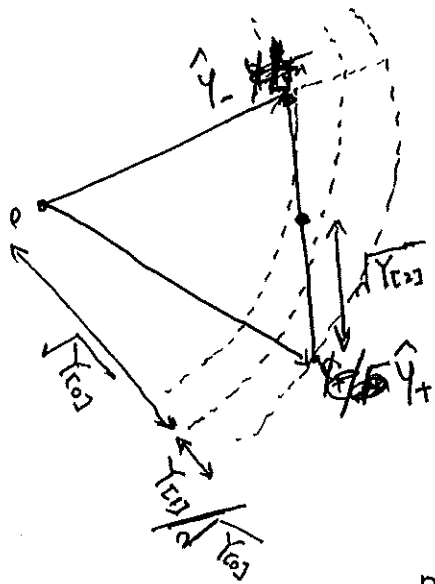
In particular, we have 4-point covariance

$$\text{Cov} [\sigma_{i_1 \alpha}^{(c)}, \sigma_{i_2 \alpha_2}^{(c)}, \sigma_{i_2 \alpha_3}^{(c)}, \sigma_{i_1 \alpha_4}^{(c)}] = \frac{1}{4} \sum_{\beta_1, \beta_2, \beta_3, \beta_4} V^{\beta_1, \beta_2, \beta_3, \beta_4} \left\langle \sigma_{\alpha_1} \sigma_{\alpha_2} (z_{\beta_1} z_{\beta_2} - G_{\beta_1, \beta_2}^{(c)}) \right\rangle_{G^{(c)}} \left\langle \dots \right\rangle$$

$i_1 \neq i_2$
 $i_1 = i_2$

We also have

$$\left\langle F(z_{i_1}^{(c)}) G(z_{i_1}^{(c)}) \right\rangle_{S^{(c)}} = \left\langle F_{A_1} G_{A_2} \right\rangle_{G^{(c)}} +$$



$$Y_{1/n} = \left\| \frac{1}{2} (y_+ + y_-) / \frac{1}{\sqrt{n}} \right\|^2 = \left\| \frac{1}{2} (\hat{y}_+ + \hat{y}_-) \right\|^2$$

$$Y_{1/n} = \left\| \frac{1}{2} y_+ / \frac{1}{\sqrt{n}} \right\|^2 - \left\| \frac{1}{2} y_- / \frac{1}{\sqrt{n}} \right\|^2 = \left\| \frac{1}{2} \hat{y}_+ \right\|^2 - \left\| \frac{1}{2} \hat{y}_- \right\|^2$$

$$Y_{\perp} = \left\| \frac{1}{2} (y_+ - y_-) / \frac{1}{\sqrt{n}} \right\|^2 = \left\| \frac{1}{2} (\hat{y}_+ - \hat{y}_-) \right\|^2$$

$$K = \cancel{K_{1/n} Y_{1/n} + K_{\perp} Y_{\perp} + K_{1/n} Y_{1/n}} K_{1/n} [1; 1] + K_{\perp} [1; -1] + K_{1/n} [1; 1]$$

and $y \sim \mathcal{N}(0, \delta_{ij} K_{\alpha\beta})$

then $Y_{1/n}, Y_{\perp}, Y_{1/n}$ has joint PDF

$$P(y_+, y_- | H) = (4\pi^2 (4K_0 K_2 - K_1^2))^{-\frac{n}{2}} \exp\left(-\frac{\phi n}{4K_0 K_2 - K_1^2} (2K_0 Y_2 + 2K_2 Y_0 - K_1 Y_1)\right)$$

Let $\hat{K}_{1/n} = \frac{1}{n} K_{1/n}$, $\hat{y}_+ = \frac{1}{\sqrt{n}} y_+$, then (equation 6.48)

$$P(\hat{y}_+, \hat{y}_- | H) = (4\pi^2 (4\hat{K}_0 \hat{K}_2 - \hat{K}_1^2))^{-\frac{n}{2}} \exp\left(-\frac{2\hat{K}_0 Y_2 + 2\hat{K}_2 Y_0 - \hat{K}_1 Y_1}{2\hat{K}_0 \hat{K}_2 - \hat{K}_1^2}\right)$$

$$\hat{\sigma}_\alpha^L = \sigma_\alpha^L / \sqrt{n}$$

$$\hat{\sigma}_{\alpha\beta}^{(L)} = C_b^{(L)} + C_w^{(L)} (\hat{\sigma}_\alpha^L \cdot \hat{\sigma}_\beta^L)$$

$$K_{\alpha\beta}^{(L)} = \langle \hat{\sigma}_{\alpha\beta}^{(L)} \rangle = C_b^{(L)} + C_w^{(L)} \langle \hat{\sigma}_\alpha^L \cdot \hat{\sigma}_\beta^L \rangle = \langle \hat{z}_\alpha^{(L)} \hat{z}_\beta^{(L)} \rangle$$

$$V_{\alpha_1 \alpha_2 \beta_1 \beta_2}^{(L)} = n_c (C_w^{(L)})^2 \text{Cov}[\hat{\sigma}_{\alpha_1}^L \cdot \hat{\sigma}_{\alpha_2}^L, \hat{\sigma}_{\beta_1}^L \cdot \hat{\sigma}_{\beta_2}^L]$$

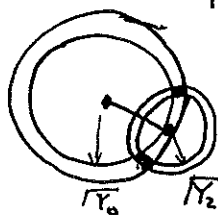
$$= \frac{(C_w^{(L)})^2 (C_w^{(L)})^2}{n_c} \text{Cov}[\sigma_{\alpha_1}^L \cdot \sigma_{\alpha_2}^L, \sigma_{\beta_1}^L \cdot \sigma_{\beta_2}^L]$$

$$G^{(L)} = \frac{G^{(0)(L)}}{n_c} + \frac{1}{n_c} G^{(1)(L)} + \frac{1}{n_c^2} G^{(2)(L)} + \dots = \langle \hat{z}_{\alpha}^{(L)} \hat{z}_{\beta}^{(L)} \rangle = \langle \hat{z}_\alpha^{(L)} \cdot \hat{z}_\beta^{(L)} \rangle$$

$$V^{(L)} = V^{(0)(L)} + \frac{1}{n_c} V^{(1)(L)} + \frac{1}{n_c^2} V^{(2)(L)} + \dots$$

If $K_{1/n} = 0$, ~~symmetric~~ and $Y_{1/n} = 0$, then

$$P(Y_{1/n}, Y_{\perp} | Y_{1/n} = 0) \text{ approximately } \propto Y_{1/n}^{\frac{n}{2}-1} Y_{\perp}^{\frac{n}{2}-1} \exp\left(-\frac{Y_{1/n} \hat{K}_{1/n} - Y_{\perp} \hat{K}_{\perp}}{2\hat{K}_0 \hat{K}_2}\right)$$



which has maximum

$$Y_{1/n} = (n-2) \hat{K}_{1/n} \approx K_{1/n}$$

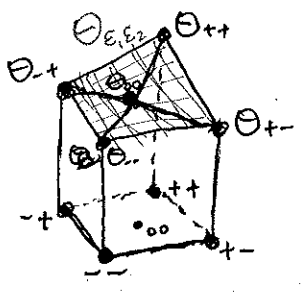
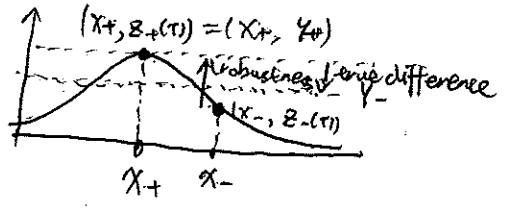
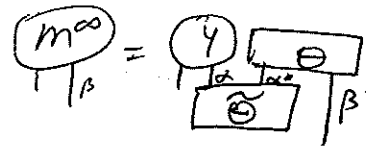
$$Y_{\perp} = (n-2) \hat{K}_{\perp} \approx K_{\perp}$$

$$= \text{Cov}\left[\hat{z}_{\alpha_1}^L \cdot \hat{z}_{\alpha_2}^L, \hat{z}_{\beta_1}^L \cdot \hat{z}_{\beta_2}^L\right] = \text{Cov}\left[\hat{G}_{\alpha_1 \alpha_2}^{(L)}, \hat{G}_{\beta_1 \beta_2}^{(L)}\right]$$

Bayesian correction.

$$\langle \hat{z}_{\alpha\beta} | y_A \rangle = \hat{m} + \frac{1}{2} \left(\text{Eq (6.88)} \right)$$

Kernel learning I



Example of using Feynmann diagram: evaluating (6.88).

$$2 \langle\langle W_{i;\beta} Q(w) \rangle\rangle_G = \text{diagram 1} + \text{diagram 2}$$

Diagram 1: A box labeled v with three Φ nodes on top. The rightmost Φ node is connected to a W node, which is connected to another W node. A dashed box encloses the two W nodes. A label $i;\beta$ is below the first W node. A G label is at the top right.

Diagram 2: A box labeled v with three Φ nodes on top. The rightmost Φ node is connected to a W node, which is connected to another W node. A dashed box encloses the two W nodes. A G label is at the top right.

$$= \text{diagram 3} + \text{diagram 4} + \text{diagram 5} + \text{diagram 6}$$

Diagram 3: Box v with three Φ nodes. The rightmost Φ node has a loop.

Diagram 4: Box v with three Φ nodes. The rightmost Φ node is connected to a W node, which has a loop.

Diagram 5: Box v with three Φ nodes. The rightmost Φ node is connected to a W node, which is connected to another W node. The second W node has a loop.

Diagram 6: Box v with three Φ nodes. The rightmost Φ node is connected to a W node, which is connected to another W node. The first W node has a loop.

A large bracket on the right side of diagrams 3-6 is labeled "Wick".

$$= \text{diagram 7} + 2 \text{diagram 8} + n_L \text{diagram 9}$$

Diagram 7: Box v with three Φ nodes. The rightmost Φ node has a loop.

Diagram 8: Box v with three Φ nodes. The rightmost Φ node is connected to a W node, which has a loop.

Diagram 9: Box v with three Φ nodes. The rightmost Φ node is connected to a W node, which is connected to another W node. The second W node has a loop.

Legend: $\square = n_L = \sum_{i,j} \delta_{ij} \delta_{i2}$



$\gamma^{(0)} \quad \gamma^{(1)} \quad \gamma^{(2)}$

++
++

0-
0-

+ -
- +

$$K_{00}^{(0)} = \langle \frac{1}{\sqrt{N}} \sum_{i=1}^N z_i^{(0)} | z^{(0)} \rangle^2$$

$$R^{(0)} = \langle \frac{1}{\sqrt{N}} \sum_{i=1}^N z_i^{(0)} | z^{(0)} \rangle - \frac{1}{\sqrt{N}} \langle z^{(0)} | z^{(0)} \rangle$$

$$D^{(0)} = \langle \frac{1}{\sqrt{N}} \sum_{i=1}^N z_i^{(0)} | z^{(0)} \rangle^2 - \frac{1}{N} \langle z^{(0)} | z^{(0)} \rangle^2$$

$$\begin{bmatrix} K_{00} \\ K_{01} \\ K_{02} \\ R \\ D \end{bmatrix} = \begin{bmatrix} \gamma_0 & \gamma_1 & \gamma_2 \\ 2\gamma_0 & 2\gamma_1 & 0 \\ \gamma_0 & \gamma_1 & -2\gamma_2 \\ 1 & -1 & 0 \\ 1 & 1 & -2 \end{bmatrix} \begin{bmatrix} K_{++} \\ K_{--} \\ K_{+-} \end{bmatrix}$$

$$K = K_{00} \gamma^{(0)} + K_{01} \gamma^{(1)} + K_{02} \gamma^{(2)}$$

coincidence-limit of 2 inputs (coincident fixed-point)

$$K_{00}^* [1,1], R=D=0$$

is a fixedpoint of $K_{\alpha\beta}^{(L+1)} \leftarrow C_0 + C_W \langle \sigma(z_{\alpha}) \sigma(z_{\beta}) \rangle_{K^{(L)}}$

$$K_{\alpha\beta}^{(L)} := \left\langle \frac{1}{N} \sum_{i=1}^N z_{i,\alpha}^{(L)} z_{i,\beta}^{(L)} \right\rangle \quad (\text{kernel recursion formula})$$

3 perturbations $\left\{ \begin{array}{l} \Delta K_{00} \\ \Delta R = \Delta(K_{++} - K_{--}) \\ \Delta D = \Delta(K_{++} + K_{--} - 2K_{+-}) \end{array} \right\}$ single input something new

input midpoint perturbation.

let X_0 be a whatever input, ("midpoint input")

let SX be a whatever input, ("perturbation")

Define $X_{\pm} := X_0 \pm SX$ (infinitesimally-separated input points)

Define $K_{00}^{(L)}$ to be the $K_{00}^{(L)}$ if we use X_0 as input.

$$K_{00}^{(L)} := \left\langle \frac{1}{N} \sum_{i=1}^N [z_{i,\alpha}(X_0)]^2 \right\rangle$$

"midpoint kernel"

strategy:

① Let $X_0 = \frac{1}{2}(X_+ + X_-)$, $X_{\pm} = X_0 \pm SX$

② analyze midpoint kernel $K_{00}^{(L)}$, $\Delta K_{00}^{(L)}$ to order δ^2

③ analyze $\Delta R^{(L)} = \Delta(K_{++}^{(L)} - K_{--}^{(L)})$ to order δ^2

④ analyze ΔD to order δ^2

$$g(K) := \frac{1}{\sqrt{2\pi K}} \int dz e^{-\frac{z^2}{2K}} \sigma(z) \sigma(z)$$

$$= \mathbb{E}_{z \sim \mathcal{N}(0, K)} [\sigma(z)^2]$$

fixed-point of 1 input.

$$K_{00}^{(L+1)} \leftarrow C_0 + C_W g(K_{00}^{(L)})$$

$$K_{00}^{*(L+1)} = K_{00}^{*(L)} = K_{00}^*$$

parallel sus. $K_{00}^{(L+1)} = K_{00}^* + \Delta K_{00}^{(L+1)}$

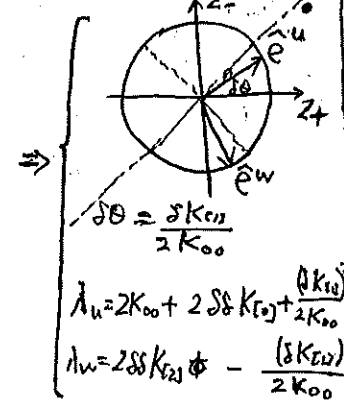
$$\Delta K_{00}^{(L+1)} = \chi_{11}(K_{00}^*) \cdot \Delta K_{00}^{(L)} + O(\delta^4)$$

$$\chi_{11}(K) = \frac{C_W}{2K^2} \langle \sigma(z) \sigma(z) (z^2 - K) \rangle_{z \sim \mathcal{N}(0, K)}$$



$$K = (K_{00} + \delta^2 K_{02}) \gamma_0^{(2)} + (\delta K_{01}) \gamma_0^{(1)} + (\delta^2 K_{02}) \gamma_0^{(2)}$$

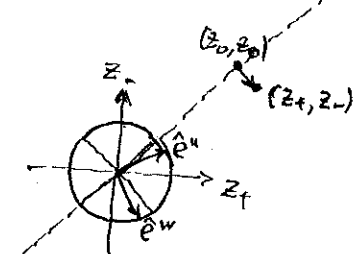
$$K e^u = \lambda_u e^u, K e^w = \lambda_w e^w$$



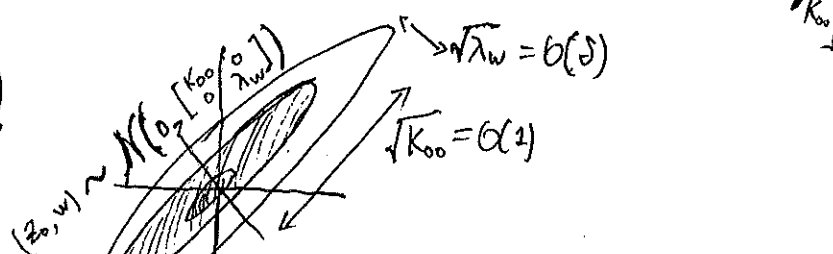
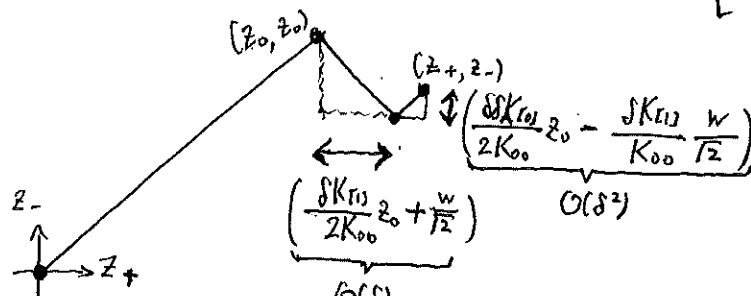
$$\delta \theta = \frac{\delta K_{01}}{2K_{00}}$$

$$\lambda_u = 2K_{00} + 2\delta \delta K_{02} + 2K_{00}$$

$$\lambda_w = 2\delta \delta K_{02} - \frac{(\delta K_{02})^2}{2K_{00}}$$



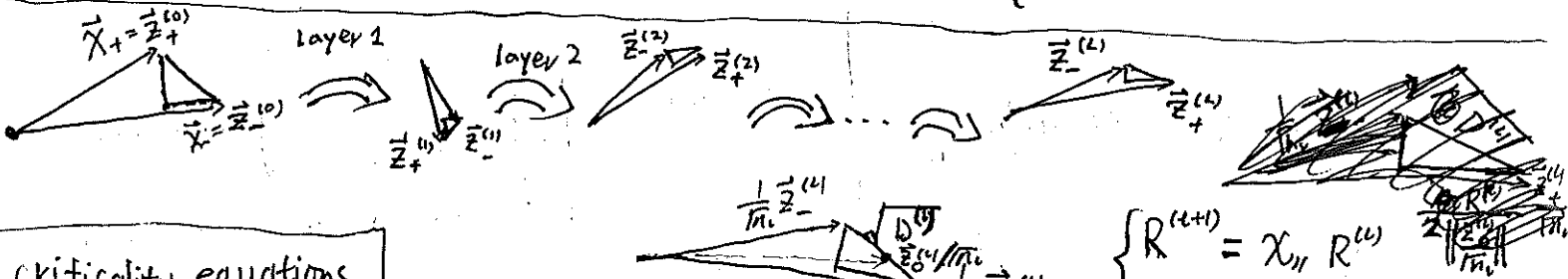
$$\langle F(z_+, z_-) \rangle_{z \sim K} = \langle F(z_0, z_0) \rangle_{K_{00}} + \left[\frac{\delta K_{02}}{2K_{00}} \langle z_0(z_+ - z_-) F \rangle_{K_{00}} + \left[\frac{\delta \delta K_{01}}{2K_{00}} \langle z_0(z_+ + z_-) F \rangle_{K_{00}} + \frac{1}{2} \left(\delta \delta K_{02} + \frac{(\delta K_{02})^2}{2K_{00}} \right) (z_0^2 - K_{00}) (z_+ - z_-)^2 F \right]_{K_{00}} \right]$$



$$\begin{cases}
 K_{00}^{(l+1)} = C_b + C_w \langle \sigma \sigma \rangle_{K_{00}^{(l)}} = K_{00}^* + \chi_{||} (K_{00}^*) \Delta K_{00}^{(l)} + O(\Delta^2) = C_b + C_w g(K_{00}^{(l)}) \\
 \delta K_{[1]}^{(l+1)} = \left(\frac{\delta K_{[1]}^{(l)}}{K_{00}^{(l)}} \langle z \sigma' \sigma \rangle_{K_{00}^{(l)}} \right) C_w = \frac{1}{2} R^{(l+1)} = \chi_{||} (K_{00}^{(l)}) \delta K_{[1]}^{(l)} \\
 \delta \delta K_{[2]}^{(l+1)} = \left(\delta \delta K_{[2]}^{(l)} \langle \sigma' \sigma' \rangle_{K_{00}^{(l)}} + \left(\frac{\delta K_{[1]}^{(l)}}{2 K_{00}^{(l)}} \right)^2 \langle (z^2 - K_{00}^{(l)}) \sigma' \sigma' \rangle_{K_{00}^{(l)}} \right) C_w^2 = \chi_{\perp} (K_{00}^{(l)}) \delta \delta K_{[2]}^{(l)} + h(K_{00}^{(l)}) (\delta K_{[1]}^{(l)})^2
 \end{cases}$$

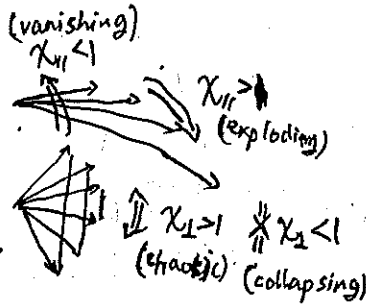
$$\begin{cases}
 g(K) = \langle \sigma \sigma \rangle_{2\omega R(0,K)} \\
 \chi_{||}(K) = C_w g'(K) = \frac{C_w}{2K^2} \langle \sigma \sigma (z^2 - K) \rangle_K = \frac{C_w}{R} \langle z \sigma' \sigma \rangle_K = C_w \langle \sigma' \sigma + \sigma' \sigma' \rangle_K \\
 \chi_{\perp}(K) = C_w \langle \sigma' \sigma' \rangle_K \\
 h(K) = \frac{C_w}{4K^2} \langle \sigma' \sigma' (z^2 - K) \rangle_K = \frac{1}{2} \chi_{\perp}'(K) = \frac{C_w}{2K} \langle z \sigma' \sigma' \rangle_K
 \end{cases}$$

$\left\{ \begin{array}{l} K_{00}^{(l)} \doteq \text{midpoint norm}^2 \\ \delta K_{[1]}^{(l)} = \frac{1}{2} R^{(l)} : \Delta(\text{norm}^2) \\ \delta \delta K_{[2]}^{(l)} = \frac{1}{4} D^{(l)} : \text{norm}^2(\Delta) \end{array} \right.$



criticality equations

$$\begin{aligned}
 K_{00}^* &= C_b + C_w g(K_{00}^*) \\
 \chi_{||}(K_{00}^*) &= 1 \\
 \chi_{\perp}(K_{00}^*) &= 1
 \end{aligned}$$



$$\begin{cases}
 R^{(l+1)} = \chi_{||} R^{(l)} \\
 D^{(l+1)} = \chi_{\perp} D^{(l)} + h(R^{(l)})^2 \\
 K_{00}^{(l)} = \langle \| z_{00}^{(l)} / \sqrt{m} \|^2 \rangle + O(\Delta^2)
 \end{cases}$$

$$\sigma = \sum \frac{\sigma_k}{k!} z^k \quad \text{with } \sigma(0)=0, \sigma'(1) \neq 0$$

$$\begin{cases}
 g = \langle \sigma \sigma \rangle_K = \sigma_1^2 [K + a_1 K^2 + a_2 K^3 + O(K^4)] \\
 \chi_{||} = C_w \sigma_1^2 [1 + 2a_1 K + 3a_2 K^2 + O(K^3)] \\
 \chi_{\perp} = C_w \sigma_1^2 [1 + b_1 K + O(K^2)] \\
 h = \frac{1}{2} C_w \sigma_1^2 (b_1 + O(K))
 \end{cases}$$

$$K_{00}^{(l)} \sim -\frac{t}{a_1 l}, \quad \delta K_{[1]}^{(l)} \sim \frac{\delta}{l^2}, \quad \delta \delta K_{[2]}^{(l)} \sim \frac{\delta^2}{l^{b_1/2a_1}}$$

If $a_1 = b_1$, then we have angle preservation on a ball.

Criticality at

$$C_b = 0, \quad C_w = V \sigma_1^2$$

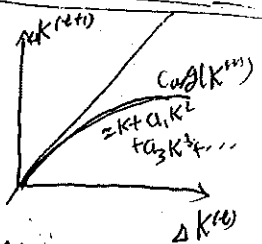
$$h = \frac{1}{2} b_1 + O(K)$$

$$\delta K^{(l)} = \frac{1}{a_1} \frac{1}{l} + O\left(\frac{1}{l^2}\right)$$

$$V^{(l)} = \frac{2}{3a_1^2} \frac{1}{l} + O\left(\frac{1}{l^2}\right)$$

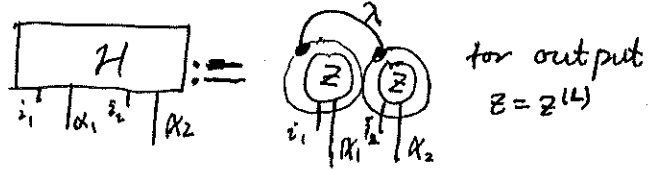
$$\frac{V^2}{n(K)^2} = \frac{2}{3} \frac{1}{n} + O\left(\frac{1}{n}\right)$$

$$g^{(l)} = \frac{\sigma_1^2}{a_1} \frac{1}{l} + O\left(\frac{1}{l^2}\right)$$



NTK

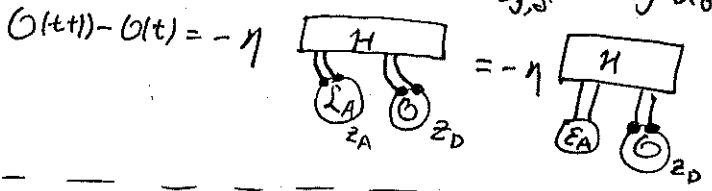
$$H_{ij\alpha\beta} = \lambda_{\mu\nu} \partial_{\theta_{i\mu}} z_{i\alpha} \partial_{\theta_{j\nu}} z_{j\beta}$$



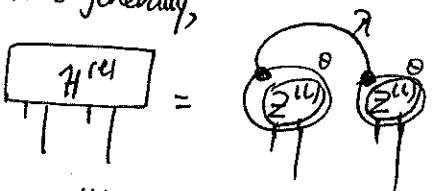
$$E_{z_i, \alpha} := \partial_{z_i, \alpha} L_A$$

$$\theta_{\mu}^{(t+1)} = \theta_{\mu}^{(t)} - \eta \lambda_{\mu\nu} \partial_{\theta_{\nu}} L_A \quad +O(\eta^2)$$

$$O^{(t+1)} = O^{(t)} - \eta \partial_{z_{i,\alpha}} L_A \partial_{z_{j,\beta}} H_{ij\alpha\beta} \quad +O(\eta^2)$$



more generally,



$$\lambda = \begin{cases} \delta \lambda_b^{(L)} \\ \delta \delta \frac{\lambda_w^{(L)}}{n_{L+1}} \end{cases}$$

$$H_{i_1 i_2 \alpha_1 \alpha_2}^{(L)} = \lambda_{\mu\nu} \partial_{\theta_{i_1 \mu}} z_{i_1 \alpha_1}^{(L)} \partial_{\theta_{i_2 \nu}} z_{i_2 \alpha_2}^{(L)}$$

H means "at initialization"

$$H_{i_1 i_2 \alpha_1 \alpha_2}^{(L+1)} = \delta z_{i_2}^{(L)} \left[\lambda_b^{(L+1)} + \frac{\lambda_w^{(L+1)}}{n_L} (\delta_{\alpha_1}^{(L)} \cdot \delta_{\alpha_2}^{(L)}) \right] \text{ new weights}$$

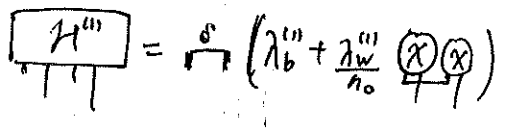
$$+ \sum_{j_1 j_2} H_{j_1 j_2 \alpha_1 \alpha_2}^{(L)} W_{j_1 i_1}^{(L+1)} (\sigma_{j_1}^{(L)}) W_{j_2 i_2}^{(L+1)} W_{j_1 j_2}^{(L+1)} (\sigma_{j_1}^{(L)}) (\sigma_{j_2}^{(L)})$$

chain rule forward prop.

RN flow $\hat{G}_{\alpha_1 \alpha_2}^{(1)} = C_b + C_w (\hat{x}_{\alpha_1} \hat{x}_{\alpha_2}) / N_0$

layer 1 $P(z^{(1)}, \hat{H}^{(1)} | D) = P_{\hat{G}^{(1)}}(z^{(1)}) \prod_{i, \alpha} S(\hat{H}_{i, \alpha}^{(1)} \delta H_{i, \alpha}^{(1)})$

$$H_{\alpha_1 \alpha_2}^{(1)} = \lambda_b^{(1)} + \lambda_w^{(1)} (\hat{x}_{\alpha_1} \hat{x}_{\alpha_2}) / N_0$$



$$\langle \hat{H}^{(2)} \rangle = \delta H^{(2)} \quad H^{(2)} = \lambda_b^{(2)} + \lambda_w^{(2)} \langle \sigma \sigma \rangle_{G^{(1)}}$$

$$H_{\alpha_1 \alpha_2}^{(2)} = \lambda_b^{(2)} + \lambda_w^{(2)} \langle \sigma_{\alpha_1}^{(1)} \sigma_{\alpha_2}^{(1)} \rangle_{G^{(1)}} + C_w^{(2)} \langle H_{\alpha_1 \alpha_2}^{(1)} \sigma_{\alpha_1}^{(1)} \sigma_{\alpha_2}^{(1)} \rangle_{G^{(1)}}$$

$$\text{Cov}[\hat{H}^{(2)}, \hat{H}^{(2)}] = \frac{1}{N_1} \left(\Pi \Pi A^{(2)} + \Pi \Pi B^{(2)} + \Pi \Pi C^{(2)} \right)$$

$$B^{(2)} = (C_w^{(2)})^2 \langle H_{\alpha_1 \alpha_2}^{(1)} \sigma_{\alpha_1}^{(1)} \sigma_{\alpha_2}^{(1)} H_{\alpha_3 \alpha_4}^{(1)} \sigma_{\alpha_3}^{(1)} \sigma_{\alpha_4}^{(1)} \rangle_{G^{(1)}}$$

$$B_{\alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(2)} = (C_w^{(2)})^2 H_{\alpha_1 \alpha_2}^{(1)} \otimes H_{\alpha_3 \alpha_4}^{(1)} \langle \sigma_{\alpha_1}^{(1)} \sigma_{\alpha_2}^{(1)} \sigma_{\alpha_3}^{(1)} \sigma_{\alpha_4}^{(1)} \rangle_{G^{(1)}}$$

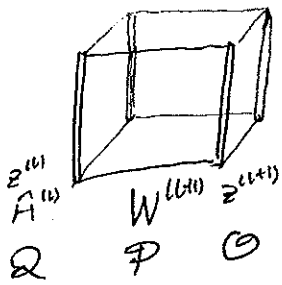
$$A^{(2)} = \text{Cov}[\hat{\Sigma}^{(2)}, \hat{\Sigma}^{(2)}]_{G^{(1)}}$$

$$\hat{\Sigma}^{(2)} = \lambda_w^{(2)} (\sigma^{(1)} \sigma^{(1)}) + C_w^{(2)} H^{(1)} \sigma^{(1)} \sigma^{(1)}$$

$$\text{Cov}[\hat{z}^{(2)}, \hat{z}^{(2)}] = \frac{1}{N_1} \Pi \Pi \left(\lambda_w^{(2)} C_w^{(2)} \text{Cov}[\sigma \sigma, \sigma \sigma] + (C_w^{(2)})^2 \text{Cov}[\sigma \sigma, H^{(1)} \sigma \sigma] \right) + \frac{1}{N_1} (\Pi + \Pi) (C_w^{(2)})^2 \langle \sigma \sigma \sigma \sigma \rangle H^{(1)}$$

$$= \frac{1}{N_1} \left(\Pi \Pi D^{(2)} + \Pi \Pi F_{\alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(2)} + \Pi \Pi F_{\alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(2)} \right)$$

$$D^{(2)} = C_w^{(2)} \text{Cov}[\sigma \sigma, \hat{\Sigma}^{(2)}] \quad F_{\alpha_1 \alpha_2 \alpha_3 \alpha_4}^{(2)} = (C_w^{(2)})^2 H^{(1)} \langle \sigma \sigma \sigma \sigma \rangle$$



$$\langle Q(z^{(l)}, \hat{H}^{(l)}) \cdot P(W^{(l)}) \cdot O(z^{(l+1)}) \rangle = e^{\frac{C_w^{(l+1)}}{2n_c} \|J\|_2^2} \langle Q(z^{(l)}, \hat{H}^{(l)}) \cdot \langle O(z^{(l+1)} + C_w \frac{1}{n_c} \sum_j J_{ij} \sigma_j^{(l)}) \rangle \rangle$$

$$P(W^{(l)}) = e^{\sum_j J_{ij} W_{ij}^{(l+1)}}$$

$$\langle Q O \rangle = \langle Q(z^{(l)}, \hat{H}^{(l)}) \cdot \langle O(z^{(l+1)}) \rangle \rangle_{G^{(l+1)}}$$

$$\langle Q W_{ij} O \rangle = \sum_{\alpha} \langle Q \langle \partial_{z_{i\alpha}^{(l+1)}} O \rangle \rangle_{G^{(l+1)}} \cdot \frac{C_w^{(l+1)}}{n_c}$$

$$\langle Q W_{i_2 j_2} W_{i_4 j_4} O \rangle = \langle Q \boxed{W} \boxed{W} O \rangle$$

$$= \frac{C_w^{(l+1)}}{n_c} \langle \langle Q \langle O \rangle \rangle \rangle + \left(\frac{C_w^{(l+1)}}{n_c}\right)^2 \langle Q \langle \langle \sigma_{i_2}^{(l+1)} \sigma_{j_2}^{(l+1)} - \sigma_{i_4}^{(l+1)} \sigma_{j_4}^{(l+1)} \rangle \rangle O \rangle_{G^{(l+1)}}$$

$$\text{Cov}[Q, \boxed{W} \boxed{W} O] = \frac{C_w^{(l+1)}}{n_c} \langle \langle \langle Q, \langle O \rangle \rangle \rangle \rangle + \left(\frac{C_w^{(l+1)}}{n_c}\right)^2 \text{Cov}[Q, \langle \langle \sigma_{i_2}^{(l+1)} \sigma_{j_2}^{(l+1)} - \sigma_{i_4}^{(l+1)} \sigma_{j_4}^{(l+1)} \rangle \rangle O \rangle_{G^{(l+1)}}$$

$$\frac{1}{n_c} \boxed{D}^{(l+1)} = \frac{\lambda_w^{(l+1)}}{C_w^{(l+1)}} \langle \langle \Delta G \rangle \rangle_{G^{(l+1)}} + C_w^{(l+1)} \langle \langle \Delta G \rangle \rangle_{G^{(l+1)}} \langle \langle \hat{H}^{(l)} \rangle \rangle_{G^{(l+1)}}$$

$$\langle \langle \hat{G} \rangle \rangle = C_b^{(l+1)} + \frac{C_w^{(l+1)}}{n_c} \langle \langle \sigma \sigma \rangle \rangle$$

$$\frac{1}{n_c} \boxed{F}^{(l+1)} = \left(\frac{C_w^{(l+1)}}{n_c}\right)^2 \langle \langle \sigma \sigma \sigma \sigma \rangle \rangle_{G^{(l+1)}} \langle \langle \hat{A}^{(l)} \rangle \rangle_{G^{(l+1)}}$$

$$\text{Cov}[O(z^{(l)}), \hat{H}^{(l)}] = \frac{1}{2n_{l-1}} \langle \langle \langle \langle \sigma \sigma \rangle \rangle O(z^{(l)}) \rangle \rangle_{G^{(l)}}$$

$$+ \frac{1}{n_{l-1}} \langle \langle \langle \langle \sigma \sigma \rangle \rangle O(z^{(l)}) \rangle \rangle_{G^{(l)}}$$

$$\langle \langle \hat{\Sigma} \rangle \rangle = \lambda_w^{(l)} \langle \langle \sigma \sigma \rangle \rangle + C_w^{(l)} \langle \langle \hat{H}^{(l)} \rangle \rangle \langle \langle \sigma \sigma \rangle \rangle$$

$$\text{Cov}[F(z_{i_1, A_1}), G(z_{i_2, A_2})] \quad (i_1 \neq i_2)$$

$$\frac{1}{4n_{l-1}} \langle \langle \langle \langle \sigma \sigma \rangle \rangle F(z_{A_1}) \rangle \rangle_{G^{(l)}} \langle \langle \langle \langle \sigma \sigma \rangle \rangle G(z_{A_2}) \rangle \rangle_{G^{(l)}} + O\left(\frac{1}{n^2}\right)$$

$$\langle \langle \langle \langle \sigma \sigma \rangle \rangle \rangle_{G^{(l)}}$$

Equivalence principle: All parameters ~~types~~ ^{types} should contribute equally to learning.

$$\hat{H}_{\mu\nu}^{(l)} = \sum_{k \in \mathcal{L}} \sum_{\mu, \nu} \frac{\partial z_{i\alpha}^{(l)}}{\partial \theta_{\mu}^{(l)}} \frac{\partial z_{j\beta}^{(l)}}{\partial \theta_{\nu}^{(l)}} \lambda_{\mu\nu}^{(l)}$$

Each of $l' \in \mathcal{L}$, $\sum_{k, k'} \frac{\partial z_{i\alpha}^{(l')}}{\partial b_k^{(l')}} \frac{\partial z_{j\beta}^{(l')}}{\partial b_{k'}^{(l')}} \lambda_{b, k, k'}^{(l')}$

$$\sum_{k, k'} \frac{\partial z_{i\alpha}^{(l')}}{\partial W_{k, k'}^{(l')}} \frac{\partial z_{j\beta}^{(l')}}{\partial W_{k, k'}^{(l')}} \lambda_{W, k, k'}^{(l')}$$

should have the same order of magnitude for each of $l' = 1, 2, \dots, \ell$

$$\Delta L = -\eta \hat{H}^{(L)} + O(\eta^2)$$

so trainability means $\Delta L \sim O(\eta)$, or, $\hat{H}^{(L)} \sim O(1)$.

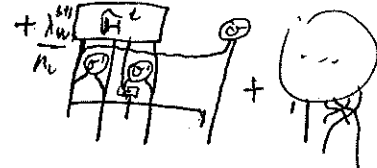
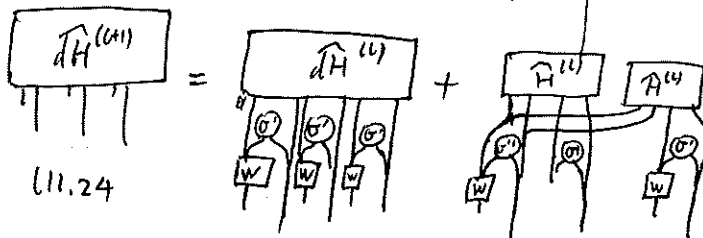
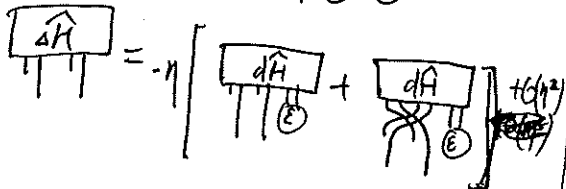
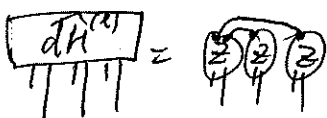
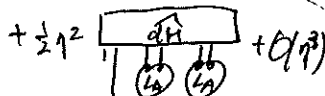
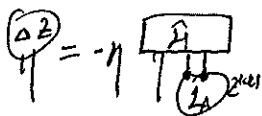
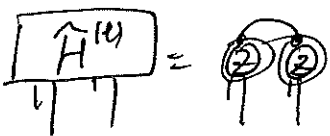
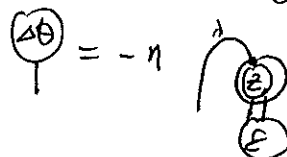
NTK learning.

~~$$\Delta L \sim O(\eta^2)$$~~

$$\Delta z_{i\alpha}^{(l)} = -\eta \sum_{j\beta} \hat{H}_{i\alpha, j\beta}^{(l)} \frac{\partial z_{j\beta}^{(l)}}{\partial \theta_{i\alpha}^{(l)}} \Delta \theta_{i\alpha}^{(l)}$$

$$+ \frac{1}{2} \eta^2 \sum_{j_1 \beta_1, j_2 \beta_2} d\hat{H}_{i\alpha, j_1 \beta_1, j_2 \beta_2} \frac{\partial z_{j_1 \beta_1}^{(l)}}{\partial \theta_{i\alpha}^{(l)}} \Delta \theta_{j_1 \beta_1}^{(l)} \frac{\partial z_{j_2 \beta_2}^{(l)}}{\partial \theta_{i\alpha}^{(l)}} \Delta \theta_{j_2 \beta_2}^{(l)} + O(\eta^3)$$

$$\Delta \theta_{i\alpha}^{(l)} = \sum_{j\beta} \lambda_{i\alpha, j\beta} \frac{\partial L}{\partial z_{j\beta}^{(l)}} \frac{\partial z_{j\beta}^{(l)}}{\partial \theta_{i\alpha}^{(l)}} = -\eta \sum_{j\beta} \lambda_{i\alpha, j\beta} \frac{\partial z_{j\beta}^{(l)}}{\partial \theta_{i\alpha}^{(l)}} \mathcal{E}_{j\beta}^{(l)}$$



Cov [z, ...]	z	H	dH	ddzH	dddH
z	1	0	1/n	1/n^2	1/n^3
H	1/n	1/n	1/n	1/n	1/n
dH	1/n	1/n	1/n	1/n	1/n
ddzH	1/n^2	1/n^2	1/n^2	1/n^2	1/n^2
dddH	1/n^3	1/n^3	1/n^3	1/n^3	1/n^3

Cov [z, ...]	z	H	dH	ddH
z	1			
H	1/n			
dH		0(1/n)		
ddH				

$$\Delta z = -\eta \left(\text{Diagram 1} \right) + \frac{1}{2} \eta^2 \left(\text{Diagram 2} \right) - \frac{1}{6} \eta^3 \left(\text{Diagram 3} \right)$$

$$= -\eta \left(\text{Diagram 4} \right) + \frac{1}{2} \eta^2 \left(\text{Diagram 5} \right) - \frac{1}{6} \eta^3 \left(\text{Diagram 6} \right)$$

generalized Newton's method

$$\Delta \theta = -\eta \hat{K}^{-1} \Delta L$$

$$\Delta \theta_{i\alpha} = -\eta \sum_{j\beta, k\gamma} \lambda_{i\alpha, j\beta, k\gamma} \frac{\partial z_{j\beta}}{\partial \theta_{i\alpha}} \frac{\partial z_{k\gamma}}{\partial \theta_{i\alpha}} \mathcal{E}_{j\beta, k\gamma}$$

ReLU universality NTK

$$A_2 = \frac{a_1^2 + a_2^2}{2}, \quad A_4 = \frac{a_1^4 + a_2^4}{2}$$

$$\begin{cases} \langle \sigma^4 \rangle_{K^*} = 3A_4 K^{(4)} \\ \langle \sigma\sigma\sigma\sigma \rangle = A_4 K^{(4)} \\ \langle \sigma^4 \rangle = A_4 \end{cases}$$

$$\begin{cases} \langle \sigma^2 \rangle = A_2 K^{(2)} \\ \langle \sigma^4 \rangle = A_4 \\ g = A_2 K^{(2)} \\ \chi_1 = C_w A_2 \\ \chi_2 = C_w A_2 \\ h = 0 \end{cases}$$

point-correlations at criticality

$$\begin{cases} \chi_1 = 1 \\ \chi_2 = 1 \\ K^{(4)} \text{ not explode} \end{cases} \Rightarrow \begin{cases} C_b = 0 \\ C_w A_2 = 1 \\ K^{(4)} = K^{(2)} = \frac{1}{A_2} \left(\frac{\sum \chi_i^2}{N_0} \right) = \frac{1}{A_2} \|\bar{\chi}\|_2^2 \\ V^{(4)} = (4-1) \left(\frac{3A_4}{A_2^2} - 1 \right) (K^*)^2 \\ g = K^* A_2 \\ h = 0 \end{cases}$$

NTK at criticality

$$\begin{cases} \Theta^{(4)} = \sum_{s=1}^{\ell} (\lambda_b^{(s)} + \lambda_w^{(s)} K^* A_2) \\ D^{(4)} = \frac{\ell(\ell-1)}{2} \times \left[\lambda_b \left(\frac{A_4}{A_2^2} - 1 \right) + \lambda_w A_2 \left(\frac{4A_4}{A_2^3} - 2 \right) \right] K^* \\ F^{(4)} = \frac{\ell(\ell-1)}{2} \times \frac{A_4}{A_2^2} (\lambda_b + \lambda_w A_2 K^*) K^* \\ B^{(4)} = \frac{\ell(\ell+1)(2\ell+1)}{6} \times \frac{A_4}{A_2^2} (\lambda_b + \lambda_w A_2 K^*)^2 \\ A^{(4)} = \frac{1}{3} \ell^3 \times \left(\frac{A_4}{A_2^2} - 1 \right) \lambda_b^2 + 3 \left(\frac{A_4}{A_2^2} - 1 \right) \lambda_b \lambda_w A_2 K^* + (5A_4 - 3A_2^2) \lambda_w^2 (K^*)^2 + O(\ell) \end{cases}$$

scaling laws.

$$\begin{aligned} \frac{\frac{1}{\ell} A^{(4)}}{\Theta^{(4)2}} &\sim \frac{\ell}{\ell} > \frac{\frac{1}{\ell} B^{(4)}}{\Theta^{(4)2}} &\sim \frac{\ell}{\ell} \\ \frac{\frac{1}{\ell} D^{(4)}}{K^{(4)} \Theta^{(4)}} &\sim \frac{\ell}{\ell} > \frac{\frac{1}{\ell} F^{(4)}}{K^{(4)} \Theta^{(4)}} &\sim \frac{\ell}{\ell} \end{aligned}$$

$K^* = 0$ universality (tanh)

$$\begin{cases} \sigma_i = \sum \frac{\sigma_n}{n!} z^n, \quad \sigma_0 = 0 \\ \langle \sigma^4 \rangle = 3\sigma_1^4 + K^2 + O(K^3) \\ \langle \sigma\sigma\sigma\sigma \rangle = \sigma_1^4 K + O(K^2) \\ \langle \sigma^4 \rangle = \sigma_1^4 + O(K) \end{cases}$$

$$C_b = 0, \quad C_w = \sqrt{\sigma_1^2}$$

$$\begin{aligned} \langle K^{(4)} \rangle &= \frac{1}{a_1} \frac{1}{\ell} + O\left(\frac{\ell}{\ell}\right) \\ V^{(4)} &= \frac{2}{3a_1^2} \frac{1}{\ell} + O\left(\frac{\ell}{\ell^2}\right) \\ \frac{V^{(4)}}{\ell \langle K^{(4)} \rangle^2} &= \frac{2\ell}{3\ell} + O\left(\frac{\ell}{\ell}\right) \end{aligned}$$

$$\frac{1}{\ell} \chi_1^{(s)} \sim \ell^{-\frac{b}{a_1}}$$

$$\Theta^{(4)} = \sum_{s=1}^{\ell} \left(\lambda_b^{(s)} + \lambda_w^{(s)} \frac{\sigma_1^2}{-a_1} \frac{1}{\ell} \right) \left(\frac{1}{\ell} \right)^{p_2} + \dots$$

$$\begin{aligned} b_1/a_1 &:= p_1 \\ &= \frac{\sigma_3 + \sigma_2^2}{\sigma_3 + \frac{3}{2}\sigma_2^2} \end{aligned}$$

Depth-scaling of learning rates

$$\lambda_b^{(4)} = \ell^{-p_2} \tilde{\lambda}_b, \quad \lambda_w^{(4)} = \ell^{-p_2+1} \tilde{\lambda}_w \Rightarrow \Theta^{(4)} = \ell^{-p_2+1} \left(\tilde{\lambda}_b + \frac{\sigma_1^2}{a_1} \tilde{\lambda}_w \right) \ell^{-p_2+1}$$

$$\begin{aligned} F^{(4)} &= \frac{1}{5-p_2} (\dots) \left(\frac{1}{\ell} \right)^{p_2-1} \\ D^{(4)} &= \frac{-2}{4-(a_1)} (\dots) \left(\frac{1}{\ell} \right)^{p_2+1} \\ B^{(4)} &= \frac{1}{3} (\dots)^2 \left(\frac{1}{\ell} \right)^{2p_2-3} \\ A^{(4)} &= (\dots) \left(\frac{1}{\ell} \right)^{2p_2-3} \end{aligned}$$

(assumes $p_1 < 5$) universality (9.81) ~ (9.86)

$$\begin{aligned} \chi_1 &= 1 - 2/\ell + \dots \\ \chi_2 &= 1 - p_2/\ell + \dots \\ C_w g &= \frac{1}{-a_1} \frac{1}{\ell} + \dots \\ h &= \frac{1}{2} b_1 + \dots \end{aligned}$$

$$\begin{aligned} C_w^2 \langle \sigma\sigma\sigma\sigma \rangle &= \frac{1}{a_1} \frac{1}{\ell} + \dots \\ C_w^2 \langle \sigma^4 \rangle &= 1 + \dots \end{aligned}$$

Scaling is still $\frac{1}{\ell} A^{(4)} / \Theta^{(4)2} \sim \frac{1}{\ell}, \dots$

$$\begin{aligned}
 \hat{Z}(w) &= \hat{Z} - \underbrace{\left[\begin{array}{c} \hat{Z}-y \\ \hat{A} \end{array} \right]}_{\text{NTK kernel regression}} + \left(\begin{array}{c} d\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} d\hat{A} \\ \hat{A} \end{array} \right) \left[\begin{array}{c} Z_A \\ \hat{A} \end{array} \right] (Z_A (Z-y) \otimes 2) \\
 &+ \left[\begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} \right] (Z_{IA} (Z-y) \otimes 2) + \left(\begin{array}{c} d\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} d\hat{A} \\ \hat{A} \end{array} \right) \left[\begin{array}{c} Z_B \\ \hat{A} \end{array} \right] (Z_B (Z-y) \otimes 2) \\
 &+ \left[\begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} \right] (Z_{IB} (Z-y) \otimes 2) + \left[\begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} \right] (Z_{IA} (Z-y) \otimes 3) \\
 &+ \left[\begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} - \begin{array}{c} dd\hat{A} \\ \hat{A} \end{array} \right] (Z_{IB} (Z-y) \otimes 3)
 \end{aligned}$$

$=$ (∞ width NTK kernel regression) + ($\mathcal{O}(\frac{1}{n})$ correction to NTK fluctuation kernel regression)

$+ [d\hat{A}] (Z_A, Z_B \text{ projections of initial error of prediction})$

$+ [dd\hat{A}, dd\hat{A}] (Z_{IA}, Z_{IB}, Z_{IA}, Z_{IB} \text{ projections of initial errors of prediction})$

$$Z^I(\omega) = \eta \sum_{s=0}^{\infty} F(s) - \eta \sum_{s=0}^{\infty} F(s) \begin{array}{|c|} \hline H \\ \hline A \\ \hline \end{array}$$

$$\eta \sum_{s=0}^{\infty} F(s) = -\frac{1}{6} \eta^2 \begin{array}{|c|} \hline dd_{II}H \\ \hline C(\omega) \\ \hline \end{array} + (7) + (6)$$

$$C(\omega) = X_{III} (z-y)^{\otimes 3}$$

$$(2) = \frac{1}{2} \eta \left(\begin{array}{|c|} \hline dH \\ \hline (X_{II}(z-y)^{\otimes 1}) \\ \hline \end{array} - \begin{array}{|c|} \hline dd_{II}H + 2dd_{III}H \\ \hline (X_{II} \otimes I - X_{III})(z-y)^{\otimes 3} \\ \hline \end{array} \right)$$

$$H^2(t) = - \left(\begin{array}{|c|} \hline dH \\ \hline \end{array} + \begin{array}{|c|} \hline H \\ \hline \end{array} \right) \quad (1) \quad (6)$$

$$+ \left(\begin{array}{|c|} \hline dd_{II}H \\ \hline \end{array} + \begin{array}{|c|} \hline dd_{III}H \\ \hline \end{array} + \begin{array}{|c|} \hline dd_{II}H \\ \hline \end{array} + \begin{array}{|c|} \hline cd_{III}H \\ \hline \end{array} + \begin{array}{|c|} \hline dd_{II}H \\ \hline \end{array} + \begin{array}{|c|} \hline dd_{III}H \\ \hline \end{array} \right)$$

$$\left(\begin{array}{|c|} \hline Y \\ \hline \end{array} \right) = \dots \quad (a(t) \otimes (z-y) - b(t))$$

$$-\frac{1}{2} \left(\begin{array}{|c|} \hline dd_{II}H \\ \hline \end{array} + \begin{array}{|c|} \hline dd_{III}H \\ \hline \end{array} + 2dd_{III}H \right)$$

$$\left(\begin{array}{|c|} \hline Y \\ \hline \end{array} \right) b(t)$$

~~Y_1 = X_{II} \otimes H~~
~~Y_2 = H \otimes H~~
~~Y_3 = H \otimes H \otimes H~~
~~Y_4 = H^2 \otimes X_{II} + (H \otimes X_{II}) X_{III}~~

$$Y_1 = X_{II} \begin{array}{|c|} \hline H \\ \hline \end{array} - \begin{array}{|c|} \hline X_{III} \\ \hline \end{array}$$

$$Y_2 = \begin{array}{|c|} \hline H \\ \hline \end{array} \begin{array}{|c|} \hline H \\ \hline \end{array} - \begin{array}{|c|} \hline X_{III} \\ \hline \end{array}$$

$$Y_3 = \begin{array}{|c|} \hline H \\ \hline \end{array} \begin{array}{|c|} \hline H \\ \hline \end{array} \begin{array}{|c|} \hline H \\ \hline \end{array} - \begin{array}{|c|} \hline X_{III} \\ \hline \end{array} \begin{array}{|c|} \hline H \\ \hline \end{array} - \begin{array}{|c|} \hline H \\ \hline \end{array} \begin{array}{|c|} \hline X_{III} \\ \hline \end{array} + \begin{array}{|c|} \hline X_{III} \\ \hline \end{array} = \begin{array}{|c|} \hline Y_4 \\ \hline \end{array}$$

$$Y_4 = \begin{array}{|c|} \hline H \\ \hline \end{array} \begin{array}{|c|} \hline X_{II} \\ \hline \end{array} - \begin{array}{|c|} \hline X_{III} \\ \hline \end{array} = \begin{array}{|c|} \hline Y_2 \\ \hline \end{array} \begin{array}{|c|} \hline H \\ \hline \end{array} - \begin{array}{|c|} \hline X_{III} \\ \hline \end{array}$$

$$Z_A = Y_2$$

$$Z_B = Y_2 + \frac{1}{2} X_{II}$$

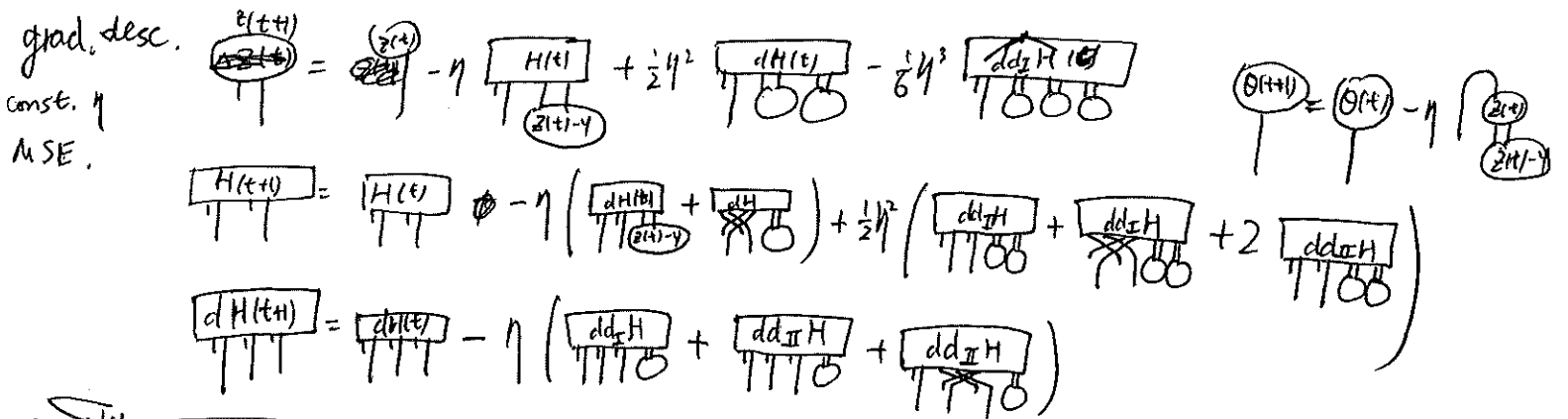
$$Z_{IA} = -Y_3 - \frac{1}{2} Y_4$$

$$Z_{IB} = -Y_3 - \frac{1}{2} (Y_3 + Y_4) - \frac{1}{6} X_{III}$$

$$Z_{IIA} = -Y_3$$

~~$$Z_{III} = Y_3$$~~

$$Z_{IIB} = - \left(\begin{array}{|c|} \hline Y_2 \\ \hline \end{array} + \begin{array}{|c|} \hline X_{II} \\ \hline \end{array} + \begin{array}{|c|} \hline X_{III} \\ \hline \end{array} \right) - \eta (Y_3 + Y_4)$$



~~width~~

$(z^F(t) - y)_{i\alpha} = \sum_{j\alpha_2} U_{ij\alpha_2}(t) (z(0) - y)_{j\alpha_2}$

$U(t) = [I - \eta \hat{H}]^t \hat{H}_{i\alpha_2, j\alpha_2}$

$z^F(t)_{i\alpha} = z(0)_{i\alpha} - \sum_{j\alpha_2, k\alpha_2} \hat{H}_{ij\alpha_2} (\hat{H}^{-1})_{jk\alpha_2} (z - z^F(t))_{k\alpha_2}$

$\hat{H}_{i\alpha_2, j\alpha_2} = \frac{\partial z^F}{\partial z(0)}$

$U(t) = [I - \eta \hat{H}] [I - \eta \hat{H}] \dots$

$z^F(t) = z^F(t) + z^I(t)$ free dynamic + interacting dynamic

$H(t) = \hat{H} + H^I(t)$ free dynamic is due to \hat{H} , interacting dynamic is due to $d\hat{H}, dd_{\hat{H}}, dd_{\hat{H}}$

dynamic free + interaction

$\eta \sum_{s=0}^{t-1} (z^F(s) - y) = \hat{H}^{-1} (z - z^F(t)) = a(t)$

$\epsilon^F(t) := z^F(t) - y = [I - \eta \hat{H}]^t (z - y)$

$\epsilon^F(\infty) = 0$

$dH(t) = d\hat{H} - (dd_{\hat{H}} + dd_{\hat{H}} + dd_{\hat{H}})$

$a(t) := \eta \sum_{s=0}^{t-1} \epsilon^F(s)$

$b(t) := \eta \sum_{s=0}^{t-1} \epsilon^F(s) \otimes \epsilon^F(s) = \chi_{II} (z - y) \otimes (z - y) - \epsilon^F(t) \otimes \epsilon^F(t)$

$\chi_{II} = \eta (I \otimes I - (I - \eta \hat{H}) \otimes (I - \eta \hat{H}))^{-1}$

$c(t) := \eta \sum_{s=0}^{t-1} \epsilon^F(s) \otimes z = \chi_{III} (z - y) \otimes z - \epsilon^F(t) \otimes z$

$\chi_{III} = \eta (I \otimes z - (I - \eta \hat{H}) \otimes z)^{-1}$

$\chi_I = \eta (I - (I - \eta \hat{H}))^{-1} = \hat{H}^{-1}$

$z^I(t+1) = z^I(t) - \eta \hat{H} z^I(t) + \eta F(t)$

$F(t) := -H^I(t) + \frac{1}{2} \eta dH(t) - \frac{1}{6} \eta^2 dd_{\hat{H}}$

$z^I(t+1) = (I - \eta \hat{H}) z^I(t) + \eta F(t)$

$z^I(0) = 0$

$z^I(t) = \eta \sum_{s=0}^{t-1} (F(s) - \hat{H} z^I(s)) = \eta (F(t-1) + U^I(t-2) + \dots + U^I(t) F(0))$

$\eta (z^I(t-1) + \dots + z^I(0)) = \hat{H}^{-1} (\eta (F(t-1) + \dots + F(0)) - z^I(t))$

$z^I(t) = \eta \sum_{s=0}^{t-1} \left(F(s) - \hat{H} z^I(s) \right) + z^I(t) + O(\frac{1}{t})$

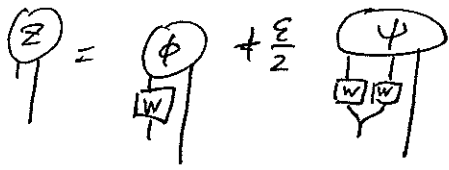
$a(\infty) = \chi_I (z - y)$

$b(\infty) = \chi_{II} (z - y) \otimes z$

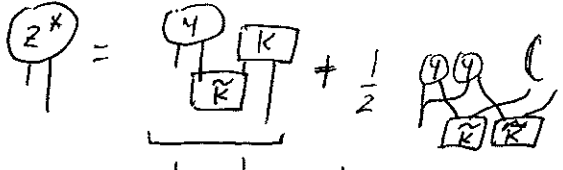
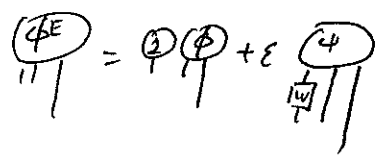
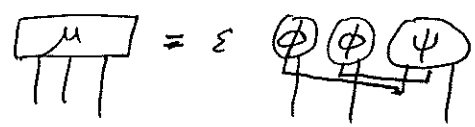
$c(\infty) = \chi_{III} (z - y) \otimes z$

$z^I(t) \sim \eta e^{-t\eta \hat{H}}$ for training set.

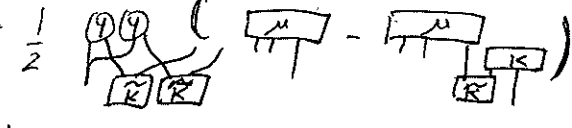
meta kernel, feature learning..



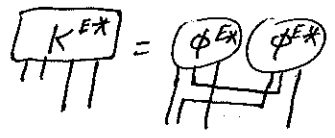
$$\phi_{ij}^E(X_j) = \phi_j(X_j) + \epsilon \sum_k W_{ik} \psi_{kj}(X_j)$$



kernel learning

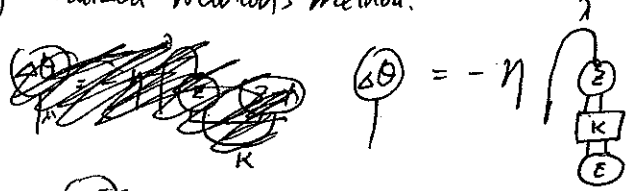


metakernel feature learning

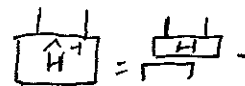
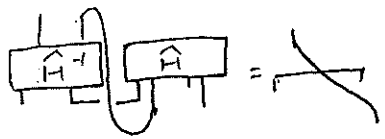
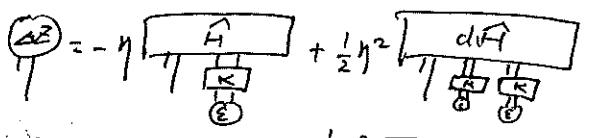
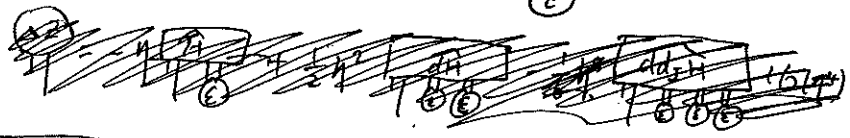


$$W^* = W^E + W^I$$

generalized Newton's method.



$$\Delta \theta_{\mu} = -\eta \sum \lambda_{\mu} \frac{\partial z_{id_1}}{\partial \theta_{\mu}} \epsilon^{j d_1} K_{ij}^{d_1 d_2}$$



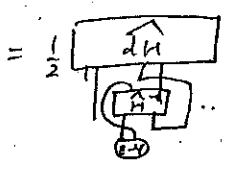
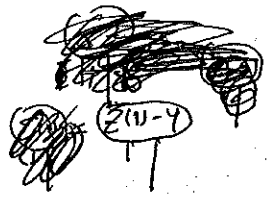
$$-\frac{1}{6} \eta^3 \frac{d^3 H}{d z^3} + O(\eta^4)$$

step 1. $K \approx \frac{1}{\eta} \hat{H}^{-1}$. This removes the $-\eta \hat{H}^{-1} K E$ term.

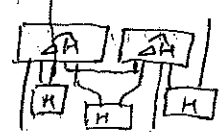
step 2: order $O(\frac{1}{\eta})$ update to Z , $K \approx \frac{1}{\eta} \hat{H}$

To order $O(\frac{1}{\eta})$, we have

$$\eta K^{(2)} = \hat{H}$$



$$= \frac{1}{2} \frac{d^2 H}{d z^2} = \frac{1}{6} \frac{d^3 H}{d z^3} + O(\frac{1}{\eta}) = \frac{1}{2} \frac{d^2 H}{d z^2} - \frac{1}{6} \frac{d^3 H}{d z^3}$$



I quit. see (ao. 79)