$D \subseteq \mathbb{R}^{n_o} \times \mathbb{R}^{n}$, $\mathcal{X} = \{x \mid (x,y) \in D\}$, $\mathcal{Y} = \{y \mid (x,y) \in D\}$.

feed-forward network:
$$\begin{cases} h^{\ell+1} = x^\ell W^{\ell+1} + b^{\ell+1} \\ x^{\ell+1} = \phi(h^{\ell+1}) \end{cases} , \quad \begin{cases} W_{ij}^{\ell+1} = \frac{\sigma_w}{\sqrt{n_\ell}} \omega_{ij}^\ell \\ b_j^\ell = \sigma_b \beta_j^\ell \end{cases} \quad \text{(NTK parameterization)}$$

at init, $w_{ij}^\ell, \beta_j^\ell \sim N(0,1)$

$\theta^\ell \equiv \text{vec}(\{W^\ell, b^\ell\}) \in \mathbb{R}^{(n_\ell(n_{\ell-1}+1))}$, $\quad \theta = \text{vec}\left(\bigcup_{\ell=1}^{L+1} \theta^\ell\right)$ (vector of all network parameters)

$f_t(x) \equiv h^{L+1}(x) \in \mathbb{R}^k$ (logits at time $t$)

$\ell(\hat{y}, y) : \mathbb{R}^k \times \mathbb{R}^k \to \mathbb{R}$, $\quad \ell(\hat{y}, y) = \frac{1}{2} \|\hat{y} - y\|_2^2$ (MSE loss)

We want to solve $\quad \min_\theta \mathcal{L} = \min_\theta \sum_{(x,y) \in D} \ell(f_t(x,\theta), y)$

For continuous time GD,
$$\theta_{t+1} \leftarrow \theta_t - \eta \nabla_\theta \mathcal{L} \quad \text{where} \quad \nabla_\theta \mathcal{L} = \nabla_\theta f_t(x)^T \nabla_{f_t(x)} \mathcal{L}$$
$$\Rightarrow \Delta \theta_t = -\eta \nabla_\theta \mathcal{L}$$

$$\frac{d\theta}{dt} = -\eta \nabla_\theta f_t(x)^T \nabla_{f_t(x)} \mathcal{L}$$

$$\frac{df_t(x)}{dt} = \nabla_\theta f_t(x) \cdot \frac{d\theta}{dt} = -\eta \nabla_\theta f_t(x) \nabla_\theta f_t(x)^T \nabla_{f_t(x)} \mathcal{L}$$

$$= -\eta \hat{\Theta}_t(x,x) \nabla_{f_t(x)} \mathcal{L} \quad \text{where} \quad f_t(x) = \text{vec}\left([f_t(x)]_{x \in \mathcal{X}}\right) \in \mathbb{R}^{k|D|+1} \text{ (concatenated logits for all examples)}$$

$$\frac{df_t(x)}{dt} = \sum_{\ell=1}^{L+1} \frac{\partial f_t(x)}{\partial \theta^\ell} \frac{d\theta^\ell}{dt}$$

$$= \nabla_\theta f_t(x) \frac{d\theta}{dt}$$

and $\hat{\Theta}(x,x) \in \mathbb{R}^{k|D| \times k|D|}$ s.t. $\quad \hat{\Theta}(x,x) = \nabla_\theta f_t(x) \nabla_\theta f_t(x)^T = \sum_{\ell=1}^{L+1} \nabla_{\theta^\ell} f_t(x) \nabla_{\theta^\ell} f_t(x)^T$

(tangent kernel at time $t$)

$$f_t^{lin}(x) \equiv f_0(x) + \underbrace{\nabla_\theta f_0(x)\big|_{\theta=\theta_0} \omega_t}_{\text{change to init. value during training}} \quad (\omega_t \equiv \theta_t - \theta_0)$$
$\underbrace{\phantom{f_0(x)}}_{\text{const.}}$

$$\xrightarrow{\quad} \frac{d\omega_t}{dt} = -\eta \nabla_\theta f_0(x)^T \nabla_{f_t^{lin}(x)} \mathcal{L} \quad \text{since} \quad \Delta\omega_t = \theta_t - \theta_0 = \Delta t \cdot \frac{d\theta}{dt}(0) = -\Delta t \eta \nabla_\theta f_0(x)^T \nabla_{f_t^{lin}(x)} \mathcal{L}$$

$$\frac{df_t^{lin}(x)}{dt} = -\eta \hat{\Theta}_0(x,x) \nabla_{f_t^{lin}(x)} \mathcal{L}$$

Thm 2.1 (informal)

Assume: $n_1 = \cdots = n_L = n$, $\lambda_{min}(\Theta) > 0$

For GD with learning rate $\eta < \eta_{critical} = \dfrac{2}{\lambda_{min}(\Theta) + \lambda_{max}(\Theta)}$,

$\forall x \in \mathbb{R}^{n_0}$ s.t. $\|x\|_2 \leq 1$,

with probability arbitrarily close to 1 over random init.,

$$\sup_{t \geq 0} \|f_t(x) - f_t^{lin}(x)\|_2, \quad \sup_{t \geq 0} \frac{\|\Theta_t - \Theta_0\|_2}{\sqrt{n}}, \quad \sup_{t \geq 0} \|\hat{\Theta}_t - \hat{\Theta}_0\|_F = O\left(\frac{1}{\sqrt{n}}\right) \quad \text{as} \quad n \to \infty$$

↳ intuition.

- in the 'lazy' regime, individual weights barely change, but          (big picture)
  they collectively conspire to provide a finite change in the final output
  ⟶ the network is perfectly described by a first-order approximation
     because the individual weight changes are so small.

  - this is the intuition for why the empirical NTK ($\hat{\Theta}_t$) should stay constant throughout training
  - stability of NTK is easier to show tabula rasa since there is no cumulation

- real network: $\dot{f}_t = -\eta \hat{\Theta}_t \nabla \mathcal{L}$ (empirical NTK)

  linear approx.: $\dot{f}_t^{lin} = -\eta \hat{\Theta}_0 \nabla \mathcal{L}$ (const. NTK from init)

  ⟶ difference b/w the two is entirely determined by $\|\hat{\Theta}_t - \hat{\Theta}_0\|_F$ (error for derivatives)

  ⟹ $\|f_t - f_t^{lin}\| \leq$ some function of $\|\hat{\Theta}_t - \hat{\Theta}_0\|$

  ↳ using Gronwall's Inequality-type arguments (bounding functions using pre-existing
                                                              bounds for the derivatives)

  $$U'(t) \leq \beta(t) U(t) \implies U(t) \leq U(0) \exp\left(\int_0^t \beta(s)\,ds\right)$$

  $$U(t) \leq \alpha(t) + \int_0^t \beta(s)\, U(s)\,ds \implies U(t) \leq \alpha(t) \exp\left(\int_0^t \beta(s)\,ds\right) = \text{some function of } \|\hat{\Theta}_t - \hat{\Theta}_0\|$$

  $\propto \text{dist}(f_t, f_t^{lin})$          $\propto \|\hat{\Theta}_t - \hat{\Theta}_0\|_{op}$

                                                  the integral is compounding error

  $\propto \int_0^t \|\hat{\Theta}_s - \hat{\Theta}_0\|_{op} \cdot \left[\text{training error of } f^{lin} \text{ at time } s\right]ds$