

Neural Tangent Kernel

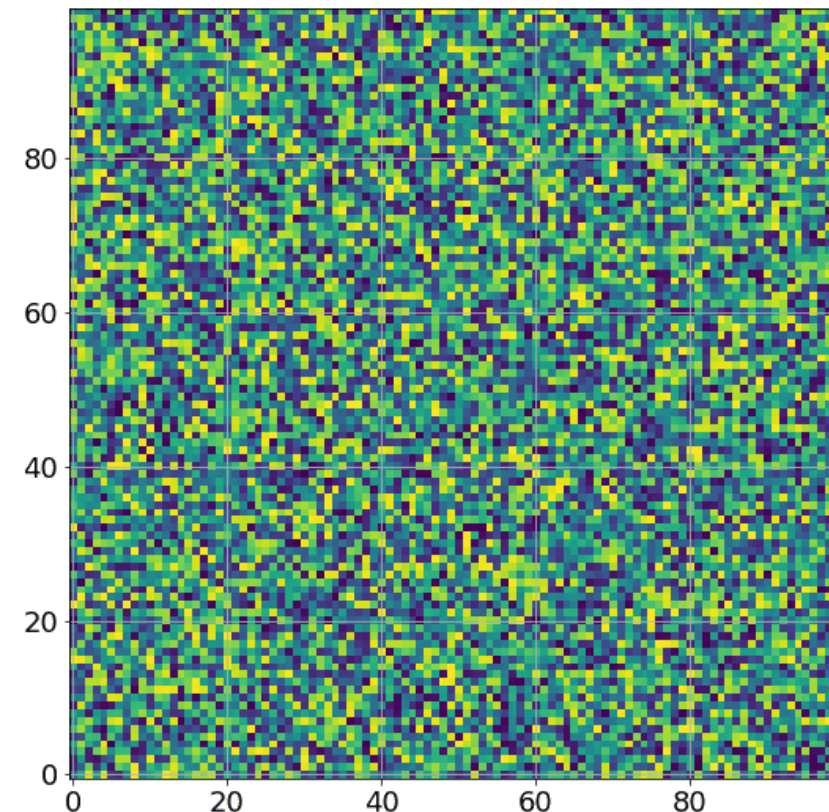
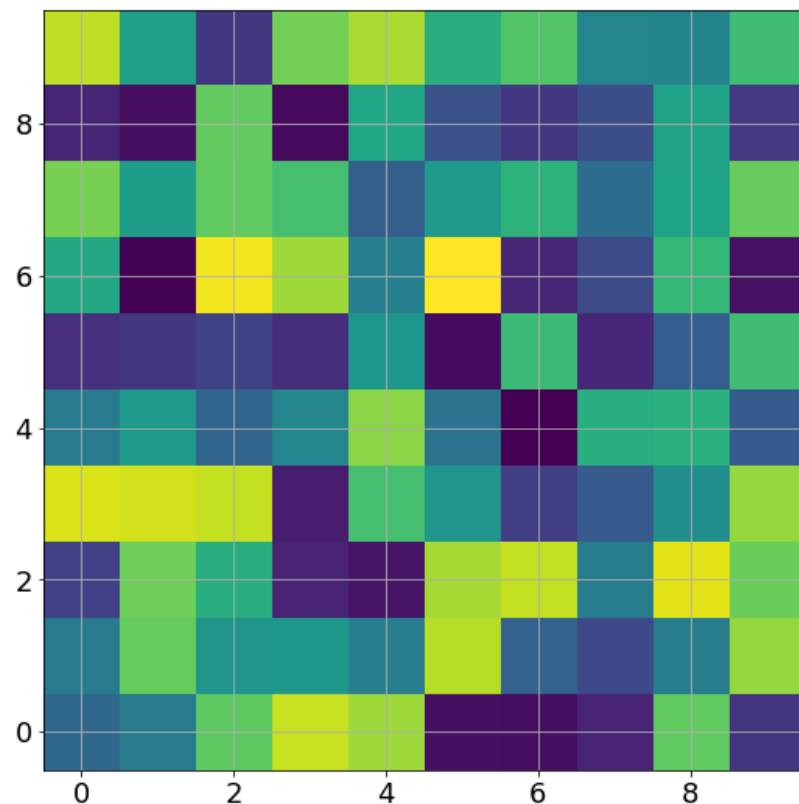
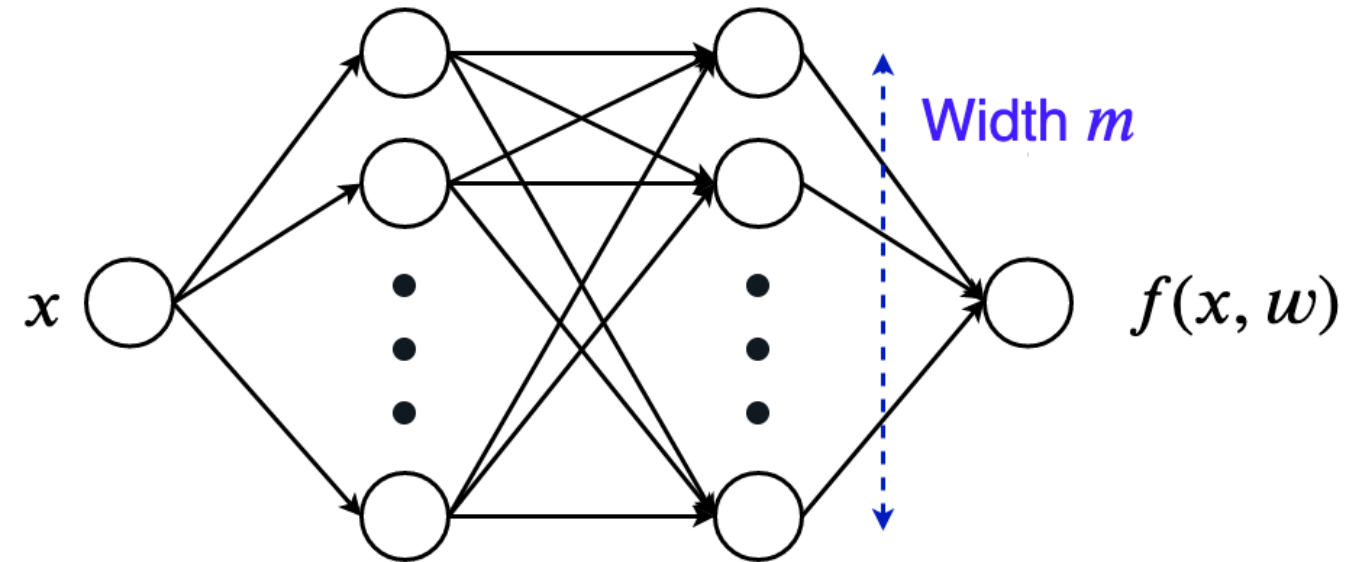
Джаин Никита БПМИ172

Least squares loss

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (f(\bar{x}_i, \mathbf{w}) - \bar{y}_i)^2$$

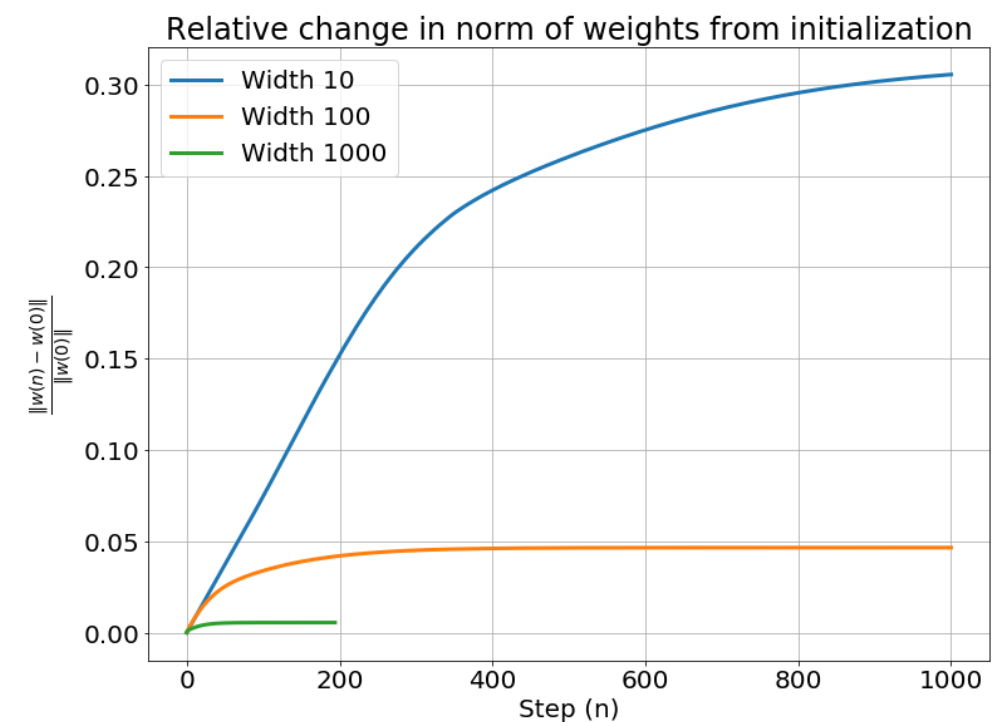
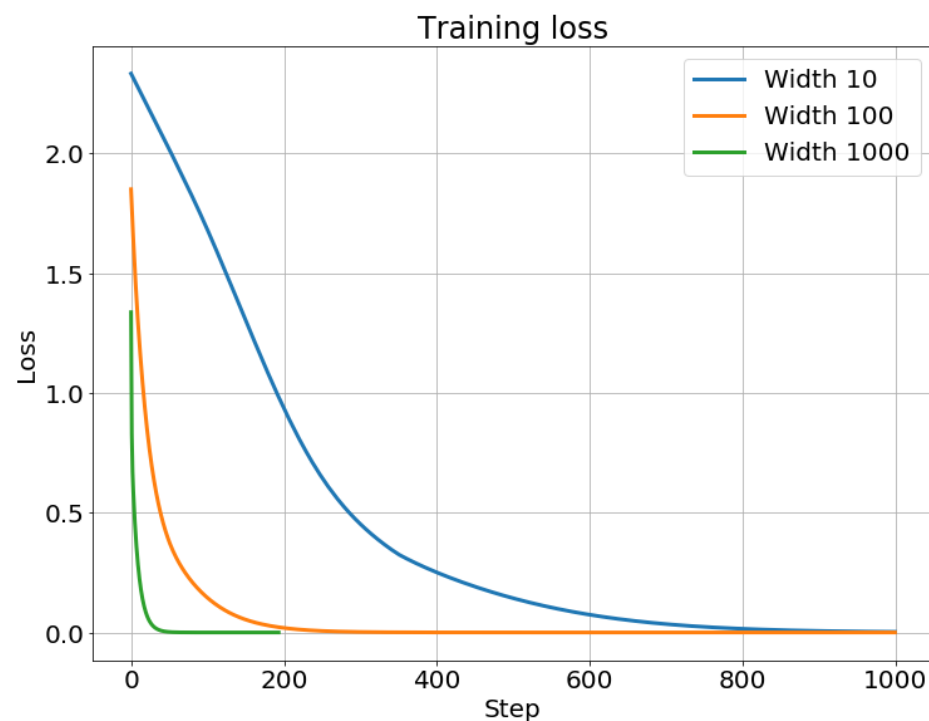
Пусть $\mathbf{y}(\mathbf{w})_i = f(\bar{x}_i, \mathbf{w})$

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}\|_2^2$$



Least squares loss

Рассмотрим простую модель с 1-D входом и выходом, с 2 полносвязными слоями:



Можно заметить, что с ростом ширины слоев относительное изменение нормы весов почти не меняется

$$\frac{\|\mathbf{w}(n) - \mathbf{w}_0\|_2}{\|\mathbf{w}_0\|_2}$$

Linear approximation

Мы можем линейно приблизить нашу модель разложив ее по Тейлору:

$$f(x, \mathbf{w}) \approx f(x, \mathbf{w}_0) + \nabla_{\mathbf{w}} f(x, \mathbf{w}_0)^T (\mathbf{w} - \mathbf{w}_0)$$

Мы линейно аппроксимировали нашу модель относительно весов.

Данное приближение можно представить, как обычную линейную модель с отображением $\phi(x)$

$$\phi(x) = \nabla_{\mathbf{w}} f(x, \mathbf{w}_0)$$

Насколько справедливо использовать такое приближение?

Linear approximation

Насколько справедливо использовать такое приближение?

$$\mathbf{w}_1 = \mathbf{w}_0 - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_0)$$

При достаточно маленьком значении η , ошибка всегда будет уменьшаться. Тогда:

$$\text{net change in } y(\mathbf{w}) \lesssim \|y(\mathbf{w}_0) - \bar{y}\|$$

$$\text{distance } d \text{ moved in } \mathbf{w} \text{ space} \approx \frac{\text{net change in } y(\mathbf{w})}{\text{rate of change of } y \text{ w.r.t } \mathbf{w}} = \frac{\|y(\mathbf{w}_0) - \bar{y}\|}{\|\nabla_{\mathbf{w}} y(\mathbf{w}_0)\|}$$

$$\text{relative change in model Jacobian} \approx \frac{d \cdot \text{rate of change of Jacobian}}{\text{norm of Jacobian}} = \frac{d \cdot \|\nabla_{\mathbf{w}}^2 y(\mathbf{w}_0)\|}{\|\nabla_{\mathbf{w}} y(\mathbf{w}_0)\|} = \|y(\mathbf{w}_0) - \bar{y}\| \frac{\|\nabla_{\mathbf{w}}^2 y(\mathbf{w}_0)\|}{\|\nabla_{\mathbf{w}} y(\mathbf{w}_0)\|^2} = \kappa(\mathbf{w}_0) \ll 1$$

Если выполняется последнее условие, то мы можем гарантировать, что наше линейное приближение близко к самой модели.

Gradient flow

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_k)$$

Перепишем как:

$$\frac{\mathbf{w}_{k+1} - \mathbf{w}_k}{\eta} = - \nabla_{\mathbf{w}} L(\mathbf{w}_k)$$

Разницу весов можно проинтерпретировать, как изменение весов в момент времени:

$$\frac{d\mathbf{w}(t)}{dt} = - \nabla_{\mathbf{w}} L(\mathbf{w}(t))$$

Подставляя ls loss, получим:

$$\dot{\mathbf{w}} = - \nabla_{\mathbf{y}(\mathbf{w})} (\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}})$$

Теперь мы можем вывести изменение модели за итерацию:

$$\dot{\mathbf{y}}(\mathbf{w}) = \nabla_{\mathbf{y}(\mathbf{w})}^T \dot{\mathbf{w}} = - \nabla_{\mathbf{y}(\mathbf{w})}^T \nabla_{\mathbf{y}(\mathbf{w})} (\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}})$$

Часть выражения, выделенная красным - NTK

Gradient flow

$$H(w) = \nabla y(w)^T \nabla y(w)$$

Если наша модель близка к линейной аппроксимации ($\kappa(w_0) \ll 1$), то Якобиан не сильно меняется:

$$\nabla y(w(t)) \approx \nabla y(w_0) \implies H(w(t)) \approx H(w_0)$$

То есть, мы можем свести динамику тренировки к решению простого ДУ

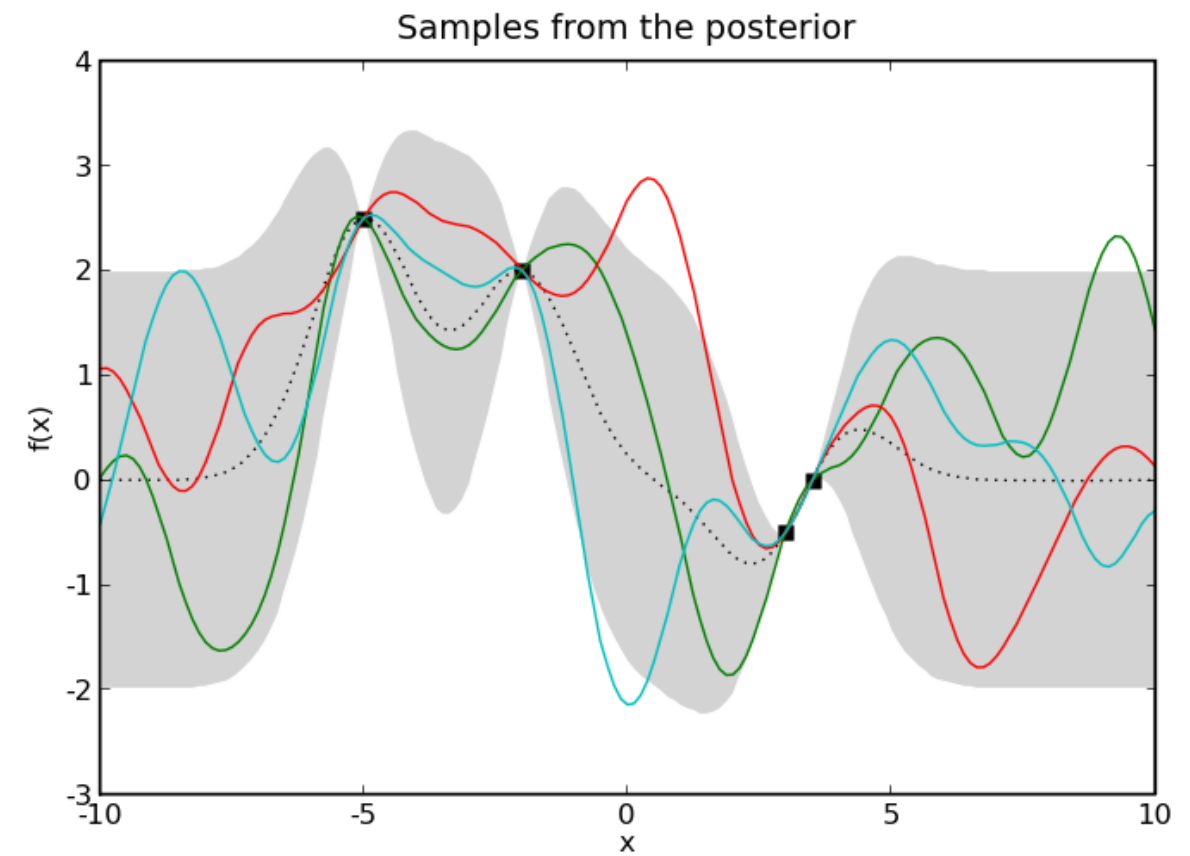
$$\dot{y}(w) = -H(w_0)(y(w) - \bar{y})$$

$$u(t) = u(0)e^{-H(w_0)t}, \text{ где } u = y(w) - \bar{y}$$

Gaussian process

$$f(\mathbf{x}) \sim \mathcal{GP}(\cdot | \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}'))$$

$$\text{Cov}(f(x), f(x')) = K(x, x')$$



Neural tangent kernel

Proposition 1. *For a network of depth L at initialization, with a Lipschitz nonlinearity σ , and in the limit as $n_1, \dots, n_{L-1} \rightarrow \infty$, the output functions $f_{\theta,k}$, for $k = 1, \dots, n_L$, tend (in law) to iid centered Gaussian processes of covariance $\Sigma^{(L)}$, where $\Sigma^{(L)}$ is defined recursively by:*

$$\begin{aligned}\Sigma^{(1)}(x, x') &= \frac{1}{n_0} x^T x' + \beta^2 \\ \Sigma^{(L+1)}(x, x') &= \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma^{(L)})} [\sigma(f(x)) \sigma(f(x'))] + \beta^2,\end{aligned}$$

taking the expectation with respect to a centered Gaussian process f of covariance $\Sigma^{(L)}$.

Neural tangent kernel

Theorem 1. *For a network of depth L at initialization, with a Lipschitz nonlinearity σ , and in the limit as the layers width $n_1, \dots, n_{L-1} \rightarrow \infty$, the NTK $\Theta^{(L)}$ converges in probability to a deterministic limiting kernel:*

$$\Theta^{(L)} \rightarrow \Theta_{\infty}^{(L)} \otimes Id_{n_L}.$$

The scalar kernel $\Theta_{\infty}^{(L)} : \mathbb{R}^{n_0} \times \mathbb{R}^{n_0} \rightarrow \mathbb{R}$ is defined recursively by

$$\begin{aligned}\Theta_{\infty}^{(1)}(x, x') &= \Sigma^{(1)}(x, x') \\ \Theta_{\infty}^{(L+1)}(x, x') &= \Theta_{\infty}^{(L)}(x, x') \dot{\Sigma}^{(L+1)}(x, x') + \Sigma^{(L+1)}(x, x'),\end{aligned}$$

where

$$\dot{\Sigma}^{(L+1)}(x, x') = \mathbb{E}_{f \sim \mathcal{N}(0, \Sigma^{(L)})} [\dot{\sigma}(f(x)) \dot{\sigma}(f(x'))],$$

taking the expectation with respect to a centered Gaussian process f of covariance $\Sigma^{(L)}$, and where $\dot{\sigma}$ denotes the derivative of σ .

Neural tangent kernel

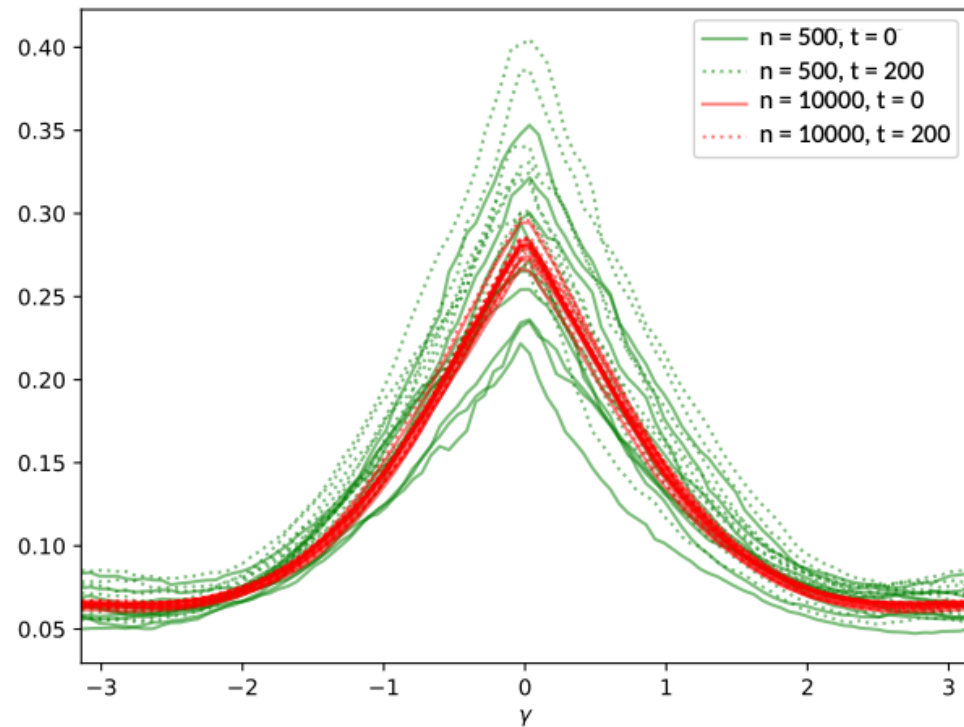


Figure 1: Convergence of the NTK to a fixed limit for two widths n and two times t .

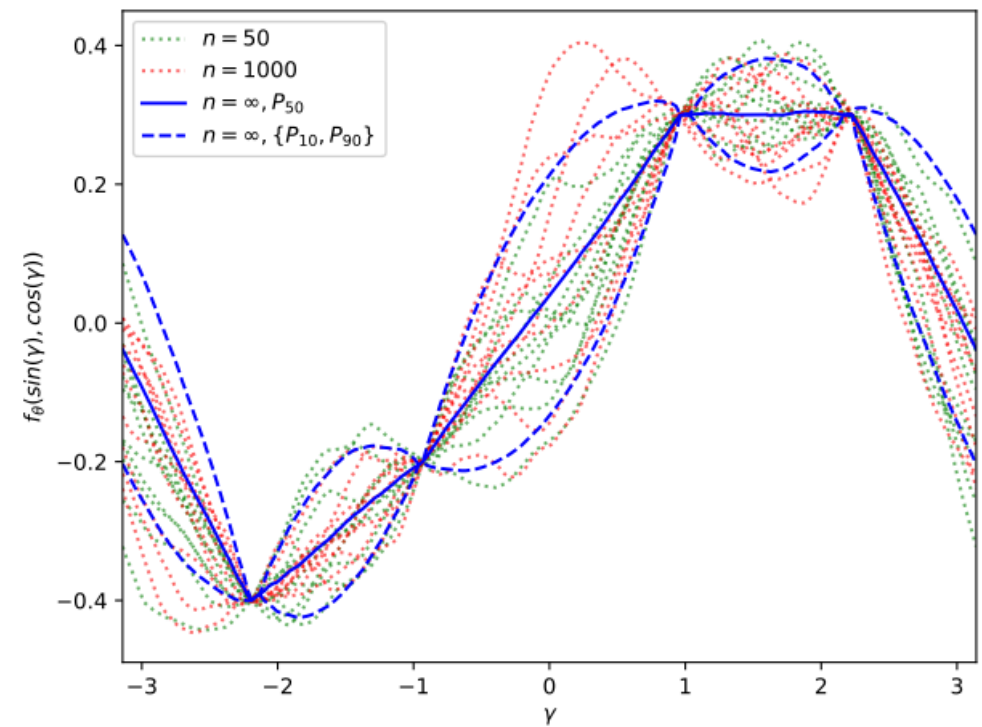


Figure 2: Networks function f_θ near convergence for two widths n and 10th, 50th and 90th percentiles of the asymptotic Gaussian distribution.

Источники

- <https://arxiv.org/pdf/1806.07572.pdf>
- <https://rajatvd.github.io/NTK/>