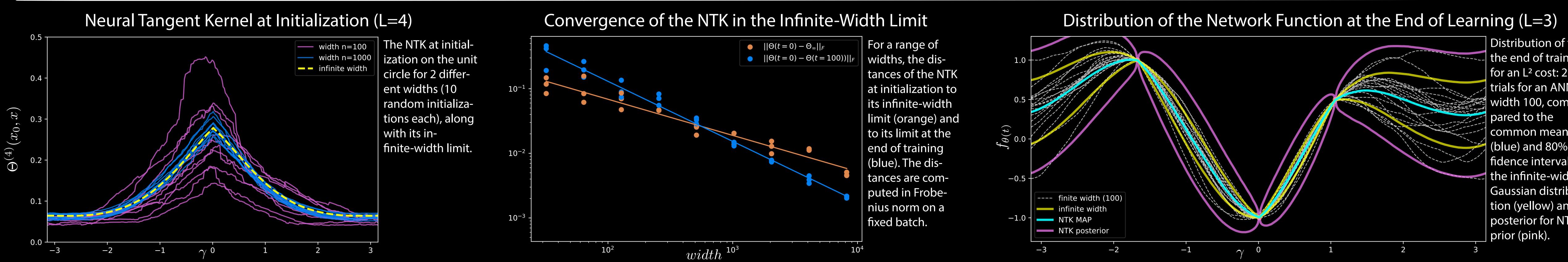
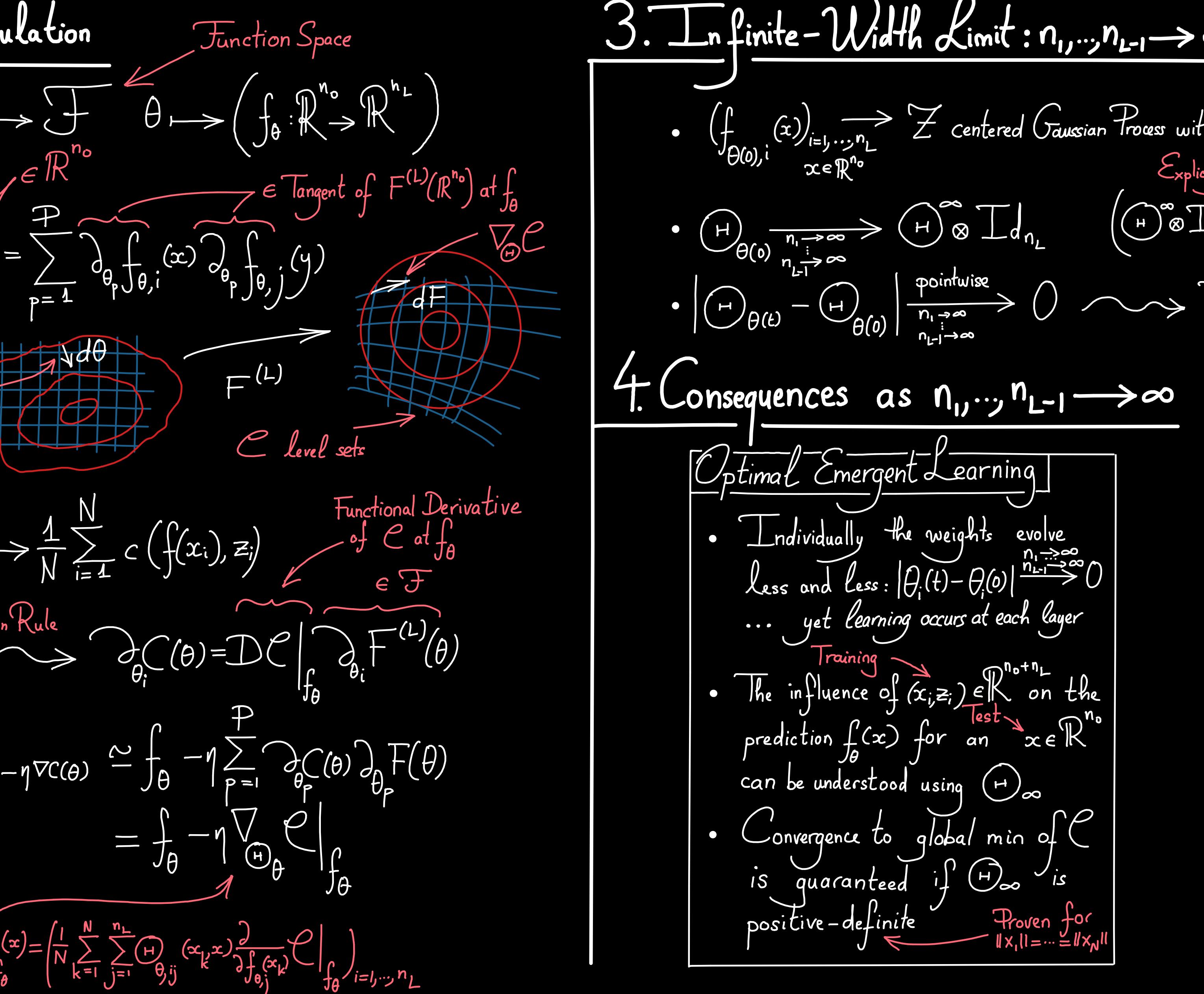
The training of Artificial Neural Networks (ANNK) involves
a prior highly non-convex optimization, yet in practice one
can observe impressive convergence and generalization group tes.
Question How to writhwarhigh understand the dynamics
of the training of an ANN as it becomes lage?
1. ANN Training: Architecture and Optimization
• Fully connected neural net with 1+1 lagrees of widths
$$n_0, ..., n_L$$
 and nonlinearity or: $\mathbb{R} \rightarrow \mathbb{R}^n$
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$



Tangent Kernel: Convergence and Generalization in Neural Network withur Jacot-Guillarmod * Franck Gabriel * Clément Hongler * EPF

2. ANN Training: Functional Farmer
• Neural Realization Function:
$$F^{(1)} R^{\mp}$$

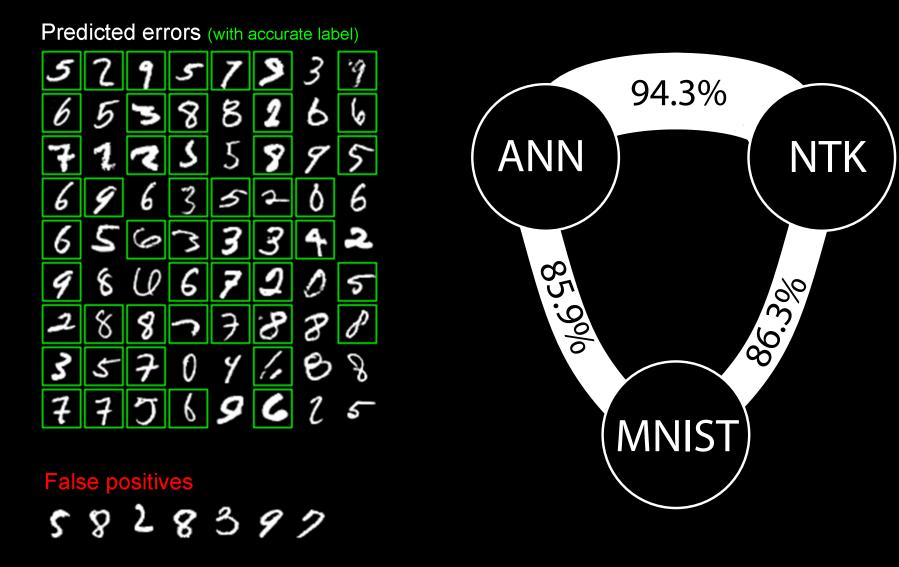
• Neural Tangent Kernel: $\bigoplus_{B,ij} (x,y) =$
• Initialization: $F^{(1)}(B) \in F$
Non-Gaussian
• Junctional Cost: $C: F \rightarrow R$ for
• Functional Cost: $C: F \rightarrow R$ for
• Training Loss: $(= C \circ F^{(1)})$
• Gradient Descent Step: $f_0 \rightarrow f_0$
 \longrightarrow Flow: $\partial_t f_{\theta(t)} = -\sum_{\theta(t)} \int_{f_{\theta(t)}} \int_{f_{\theta(t$



Distribution of f_{θ} at the end of training or an L² cost: 20 trials for an ANN of width 100, comcommon mean and 80% confidence intervals of infinite-width Gaussian distribution (yellow) and the posterior for NTK prior (pink).

Error Prediction on MNIST using the NTK

Predicted errors (w



35652808 Prediction Coincidences

The NTK can be used to predict the misclassifications that an ANN of width 500 with L=4 will make on MNIST. t yields a very good prediction on which wrong labels will be chosen by the ANN.

Thanks to a new object, the Neural Tangent Kernel,
we can precisely describe the evolution, convergence,
and generalization of ANN's of large width
Theorem In the infinite-width limit, the ANN
function follows a kernel gradient descent
with respect to the limiting NTK
A limit:
$$n_{1,1}...,n_{L-1} \rightarrow \infty$$

Veal, 1996 [2]
Z centered Gaussian Process with $Co(Z_i(x), Z_i(y)) = \sum Cx_i y/\delta_{ij}$
At Initialization
 $H \otimes Id_{n_L}$ $(H) \otimes Id_{n_L}$ $(x, y) = (H) \otimes (x, y) \delta_{ij}$
pointwise
 $n_{1 \to \infty}$ $0 \rightarrow \partial_t \int_{\theta(t)}^{\infty} = -\nabla_{H} \otimes_{Id_{n_L}}^{\infty} \int_{\theta(t)}^{\infty} During Training$

for any t>0 Bayesian Max A Posteriori $\mathcal{N}(0, \Theta \otimes \mathbb{I}_{d_n})$ prior Kernel Ridge Regression $\lambda \rightarrow 0_{+}$ regularization ned by the (H) PCA: along $V_i \propto \lambda_i$ ection PCA eigenvalue

Kernel PCA of MNIST with respect to the NTK The 2nd and 3rd kernel principal comonents of MNIST with respect to the NTK (the first component is essentially constant) with L=4. The three first eigen-00 0 values are 41.03, 1.88, and 1.46. 4444 9 94777

 \mathcal{D}_{a}

A. Jacot, F. Gabriel, C. Hongler al Tangent Kernel: Convergenc Generalization in Neural Networks, NeurIPS 2018, arXiv: 1806.07572 [2] R.M. Neal, Bayesian Learning for Neural Networks, 1996 Acknowledgments ERC CG CRITICAL (Gabriel) ERC SG CONSTAMIS) Blavatnik Family Foundation { (Hongler)

Latsis Family Foundation