

# Kernel methods and Deep Learning II

Johannes Pitz

Technische Universität München

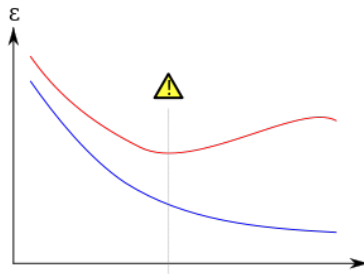
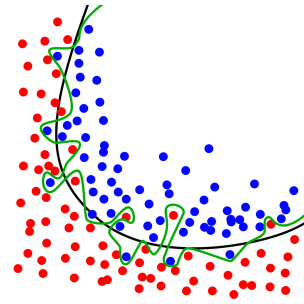
*johannes.pitz@tum.de*

April 24, 2020

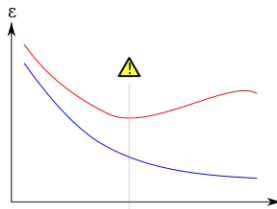
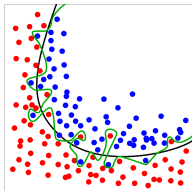
# Overview

- 1 Motivation
- 2 Experiment
  - Setup
  - Results
- 3 Theoretical Bound
- 4 Adding Noise
  - Setup
  - Results on the Norm
  - Results with more Noise
- 5 More Results
- 6 Conclusion

# Overfitting



# Overfitting



- Expect drop in test performance
- However, highly over parameterized setting
  - Overfit on training data, still get good generalization
- Recently a lot of attention in deep neural networks

# Recent Work

- Understanding deep learning requires rethinking generalization. Zhang et al.(2016).
  - Trained state-of-the-art CNN's on random labels
  - Model family/regularization techniques don't explain the small difference between training and test performance.
- Theory of Deep Learning III: explaining the non-overfitting puzzle. Poggio et al. (2017).
  - Extend the properties of SGD for linear networks to deep non linear networks  
(1.  $\exists$  optimal early stopping, 2. classification: min norm  $\implies$  maximum margin)

# This paper

"To Understand Deep Learning We Need to Understand Kernel Learning"  
Mikhail Belkin, Siyuan Ma, Soumik Mandal.

- Kernel Learning  $\implies$  a two-layer neural network
- Check if depth is needed

# Reproducing Kernel Hilbert Space $\mathcal{H}$

Let  $K(\mathbf{x}, \mathbf{z}) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be positive definite.

Given data  $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}, \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Then there exists a RKHS  $\mathcal{H}$  of functions on  $\mathbb{R}^d$ .

$$f(\cdot) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \cdot)$$

$$f^* = \arg \min_{f \in \mathcal{H}, f(\mathbf{x}_i) = y_i} \|f\|_{\mathcal{H}}$$

Where  $\|f\|_{\mathcal{H}}^2 = \langle \boldsymbol{\alpha}, \mathbf{K} \boldsymbol{\alpha} \rangle = \sum_{ij} \alpha_i \mathbf{K}_{ij} \alpha_j$ ,  
and  $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

# Interpolation

By the classical representer theorem:

$$f^*(\cdot) = \sum_{i=1}^N \alpha_i^* K(\mathbf{x}_i, \cdot)$$

$$\text{where } \boldsymbol{\alpha}^* = \mathbf{K}^{-1} \mathbf{y}$$

Note:

$$f^*(\mathbf{X}) = \mathbf{K} \boldsymbol{\alpha}^*$$

$$= \mathbf{K} \mathbf{K}^{-1} \mathbf{y}$$

$$= \mathbf{y}$$

$$f^*(\mathbf{x}_i) = y_i$$

---

Solving linear System is close to  $O(N^3)$



# Stochastic Gradient Descent

$$f^*(\cdot) = \sum_{i=1}^N \alpha_i^* K(x_i, \cdot)$$

Minimization:

$$\begin{aligned} \alpha^* &= \arg \min_{\alpha \in \mathbb{R}} \sum_{j=1}^N l(f(x_j), y_j) \\ &= \arg \min_{\alpha \in \mathbb{R}} \sum_{j=1}^N l\left(\sum_{i=1}^N \alpha_i K(x_i, x_j), y_j\right) \end{aligned}$$

---

Use EigenPro Optimizer to speed up training

# Kernel Layer

Layer (type)	Output Shape	Param #
indexed-feat- (InputLayer)	(None, 785)	0
lambda_1 (Lambda)	(None, 784)	0
kernel_embedding_2 (KernelEm	(None, 10000)	7840000
trainable (Dense)	(None, 10)	100000
Total params: 7,940,000		
Trainable params: 100,000		
Non-trainable params: 7,840,000		

# Kernels

Gauss (smooth):

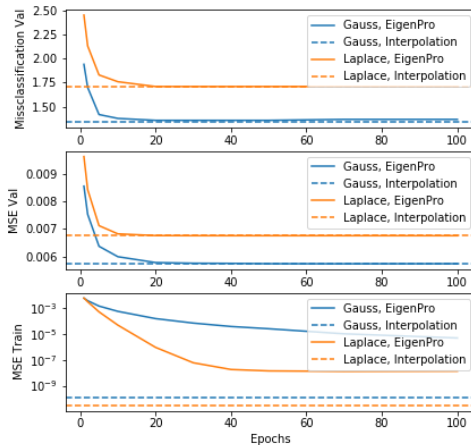
$$K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$$

Laplace (non-smooth):

$$K(x, y) = \exp\left(-\frac{\|x - y\|}{\sigma}\right)$$

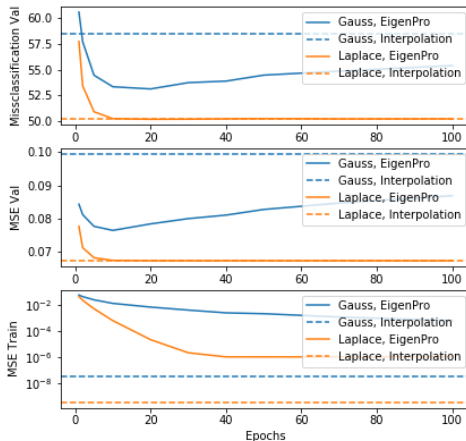
# Reproduced MNIST

MNIST, 40000 Training Samples



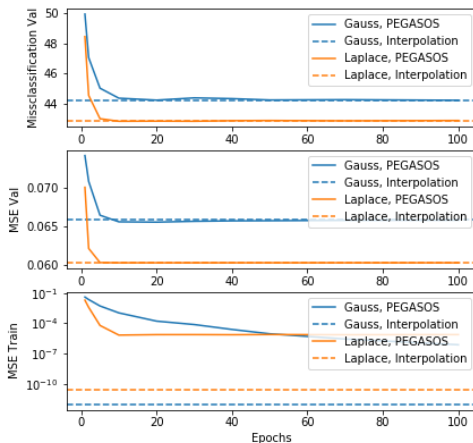
# Reproduced CIFAR-10 grey-scale

CIFAR, 40000 Training Samples

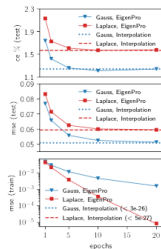


# Reproduced CIFAR-10 color

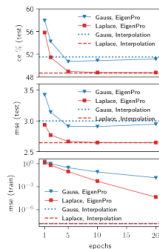
color, 40000 Training Samples



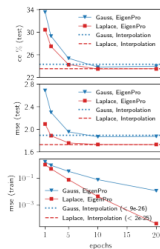
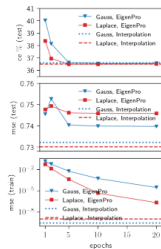
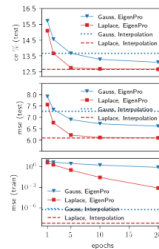
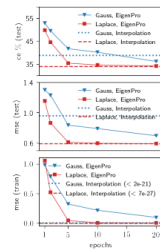
## Belkin et al.



(a) MNIST



(b) CIFAR-10

(c) SVHN ( $2 \cdot 10^4$  subsamples)(d) TIMIT ( $5 \cdot 10^4$  subsamples)(e) HINT-S ( $2 \cdot 10^4$  subsamples)

(f) 20 Newsgroups

# Recap

- Overfit on Training Data, but still reasonable generalization
  - Apparently works also in shallow networks
  - Not inherent to depth of deep neural nets
- Even more surprising:
  - Authors show that any known generalization bound goes to infinity assuming some label noise



# Theoretical Bound

$\Omega \subset \mathbb{R}^d$ ,  $\mathcal{H}$  is RKHS

Probability measure  $P$  on  $\Omega \times \{-1, 1\}$   
 $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$  sampled from  $P$   
 is a labeled dataset, with non zero label noise.

$$B_R = \{f \in \mathcal{H}, \|f\|_{\mathcal{H}} < R\} \subset \mathcal{H}$$

$h \in \mathcal{H}$   $t$ -overfits the data  $\iff$  it achieves zero classification error, and  
 $\forall_i y_i h(\mathbf{x}_i) > t > 0$

---

Interpolation requires  $\forall_i y_i = h(\mathbf{x}_i)$

# Theoretical Bound

## Theorem

Any  $h \in \mathcal{H}$  corresponding to a Gaussian kernel that  $t$ -overfits the data satisfies with high probability

$$\|h\|_{\mathcal{H}} > Ae^{Bn^{\frac{1}{d}}}$$

for some constants  $A, B > 0$  depending on  $t$

---

The results also apply to other classes of smooth kernels

# Theoretical Bound - Proof

$$l(f(\mathbf{x}), y) = \max(t - yf(\mathbf{x}), 0)$$

hinge loss with margin  $t$

$V_\gamma(B_R)$  be fat shattering dimension of the function space  $B_R$  with the parameter  $\gamma$ .

By Anthony, Bartlett in *Neural network learning: Theoretical foundation*  
 $\exists C_1, C_2 > 0$  such that with high probability  $\forall f \in B_R$ :

$$\left| \frac{1}{n} \sum_i l(f(\mathbf{x}_i), y_i) - \mathbb{E}_P[l(f(\mathbf{x}), y)] \right| \leq C_1 \gamma + C_2 \sqrt{\frac{V_\gamma(B_R)}{n}}$$

# Theoretical Bound - Proof

$$\left| \frac{1}{n} \sum_i l(f(\mathbf{x}_i), y_i) - \mathbb{E}_P[l(f(\mathbf{x}), y)] \right| \leq C_1 \gamma + C_2 \sqrt{\frac{V_\gamma(B_R)}{n}}$$

Since  $y$  is not a deterministic function of  $x$  (label noise)

$$\mathbb{E}_P[l(f(\mathbf{x}), y)] > 0$$

Now fix  $\gamma > 0$  such that  $C_1 \gamma < \mathbb{E}_P[l(f(\mathbf{x}), y)]$

Suppose  $h \in B_R$   $t$ -overfits the data. Then by construction

$$\frac{1}{n} \sum_i l(h(\mathbf{x}_i), y_i) = 0$$

# Theoretical Bound - Proof

$$0 < \mathbb{E}_P[l(f(\mathbf{x}), y)] - C_1\gamma < C_2\sqrt{\frac{V_\gamma(B_R)}{n}}$$

$$\frac{n}{C_2^2}(\mathbb{E}_P[l(f(\mathbf{x}), y)] - C_1\gamma)^2 < V_\gamma(B_R)$$

By Belkin in *Approximation beats concentration?*

$$V_\gamma(B_R) < O(\log^d(\frac{R}{\gamma}))$$

$$\frac{n}{C_2^2}(\mathbb{E}_P[l(f(\mathbf{x}), y)] - C_1\gamma)^2 < \log^d(\frac{R}{\gamma})$$

$$Ae^{Bn^{\frac{1}{d}}} < R \text{ where } A, B > 0$$



# Theoretical Bound - Consequences

Most available bounds are of the form:

$$\underbrace{\left| \frac{1}{n} \sum_i l(f(\mathbf{x}_i), y_i) \right|}_{\text{Train Loss}} - \underbrace{\mathbb{E}_P[l(f(\mathbf{x}), y)]}_{\text{Expected Test Loss}} \leq C_1 + C_2 \frac{\|f\|_{\mathcal{H}}^\alpha}{n^\beta}$$

where  $C_1, C_2, \alpha, \beta \geq 0$

$\|h\|_{\mathcal{H}} > Ae^{Bn^{\frac{1}{d}}}$  yields trivial bound

# Probability Distribution

## Proposition

Let  $P$  be a multiclass probability distribution on  $\Omega \times \{1, \dots, k\}$ . Let  $P_\epsilon$  be the same distribution, with a  $\epsilon$  fraction of the labels flipped at random with equal probability for all labels.

- 1 The Bayes Optimal Classifier  $c^*$  for  $P_\epsilon$  is the same the Bayes Optimal Classifier for  $P$
- 2 The error rate is

$$P_\epsilon(c^*(x) \neq y) = \epsilon \frac{k-1}{k} + (1-\epsilon)P(c^*(x) \neq y)$$

# Synthetic datasets

$$y \in \{0, 1\}$$

$$x_2 \sim \mathcal{U}(-1, 1), \dots, x_{50} \sim \mathcal{U}(-1, 1)$$

Synthetic 1, Separable classes:

$$x_1 \in \begin{cases} \mathcal{N}(0, 1), & \text{if } y = 1 \\ \mathcal{N}(10, 1), & \text{otherwise} \end{cases}$$

Synthetic 2, Non-Separable classes:

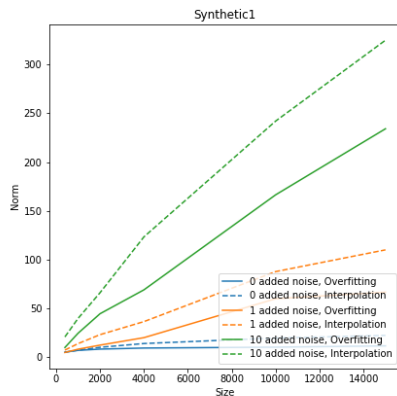
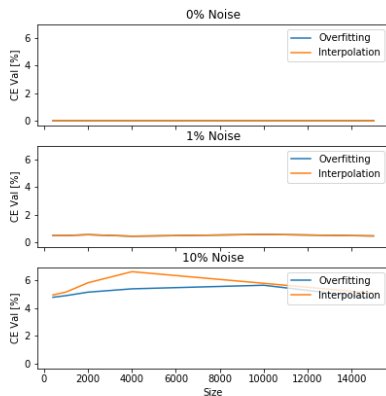
$$x_1 \in \begin{cases} \mathcal{N}(0, 1), & \text{if } y = 1 \\ \mathcal{N}(2, 1), & \text{otherwise} \end{cases}$$

---

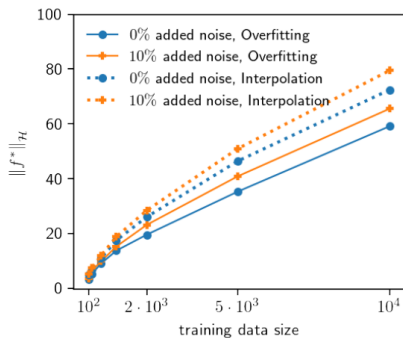
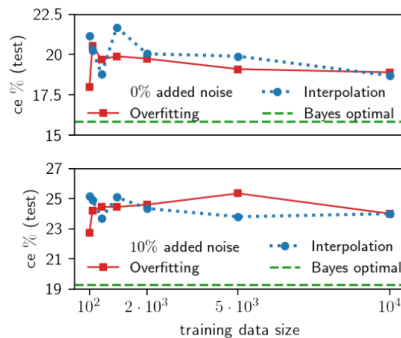
Bayes Optimal Error close to 0%, around 15.9% respectively  
 $y \in \{-1, 1\}$  in the original paper



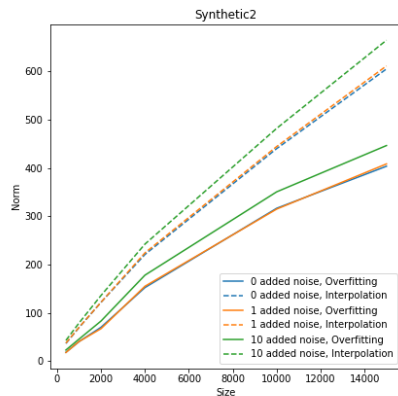
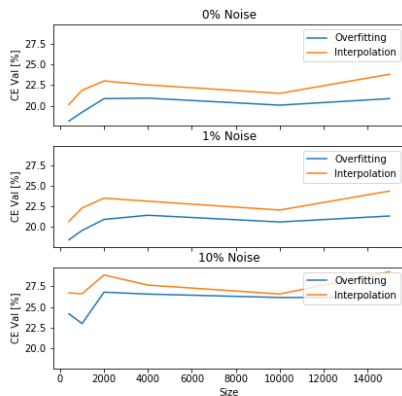
# Reproduced Synthetic 1



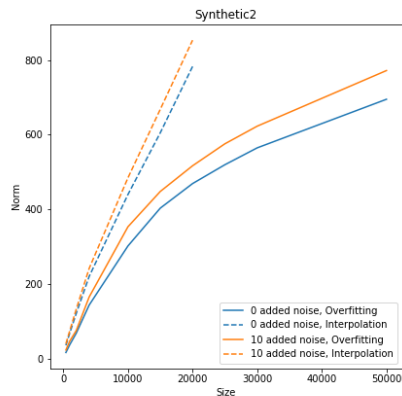
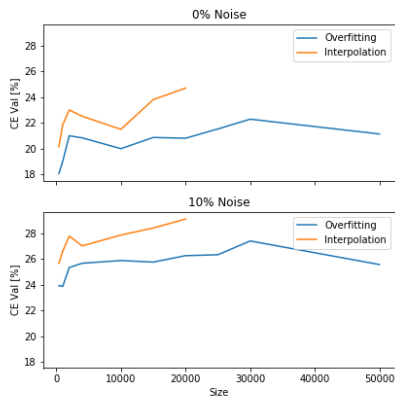
# Synthetic 2



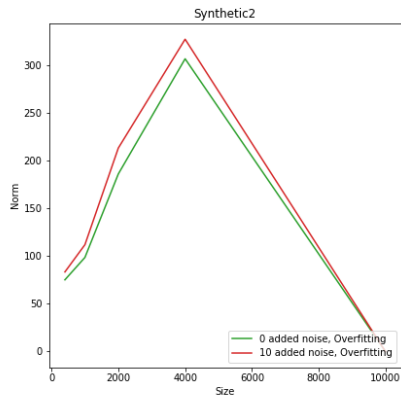
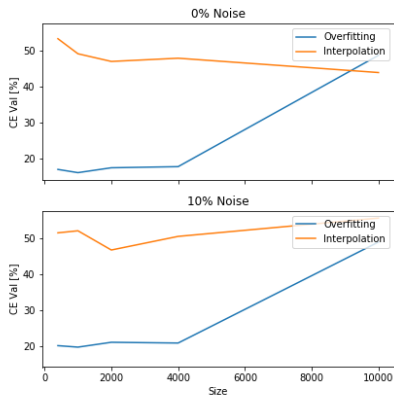
# Reproduced Synthetic 2



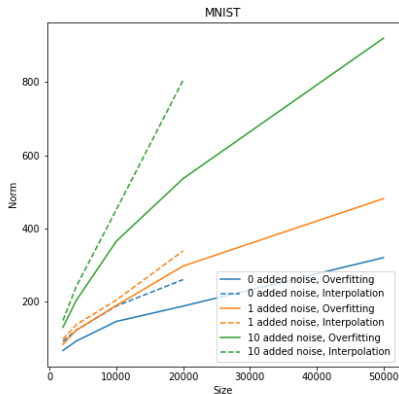
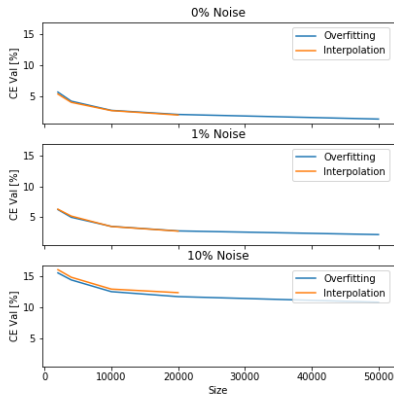
# Extended Synthetic 2



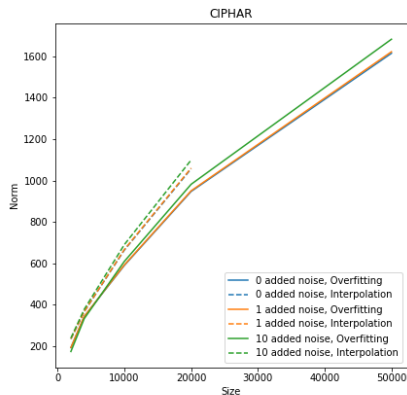
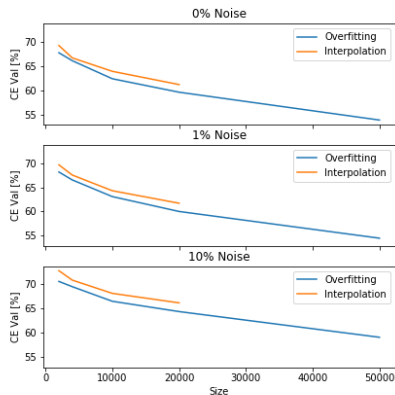
# $d = 4$ ; Synthetic 2



# Reproduced MNIST

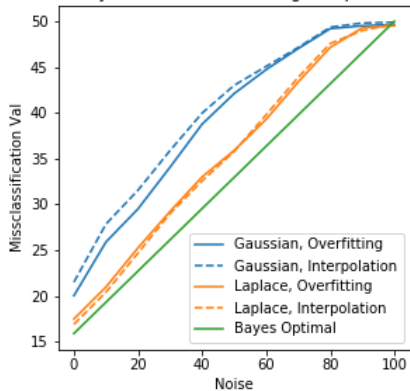


# Tried also CIFAR-10

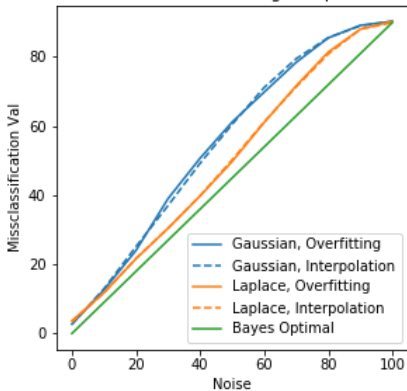


# Reproduced

Synthetic2, 10000 Traing Samples



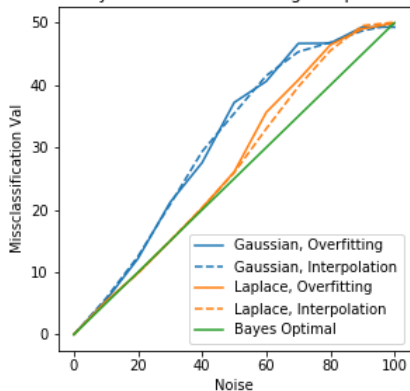
MNIST, 10000 Traing Samples



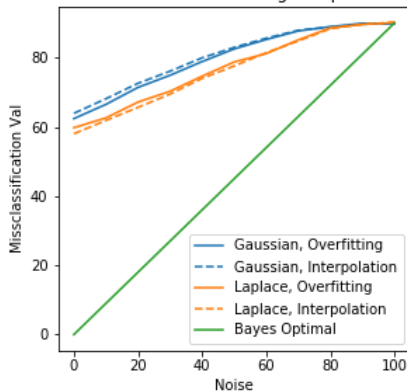


# Tried also

Synthetic1, 10000 Traing Samples



CIPHAR, 10000 Traing Samples



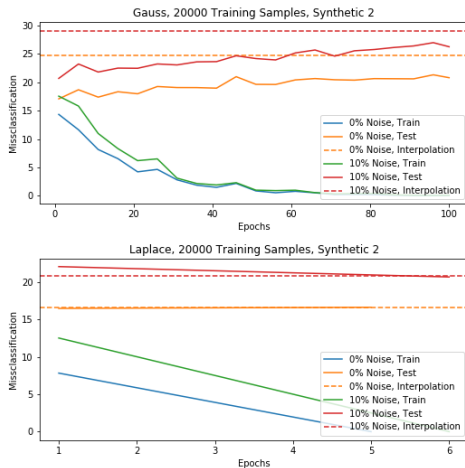
# Recap

- Norm increases with training data size
- Still yields reasonable generalization on test data
- Even when noise is close to 1

(Consistently) better than any known bound can guarantee

⇒ Need to find new bounds!

# Early Stopping would have helped



# Computational Reach

Gaussian Kernel:

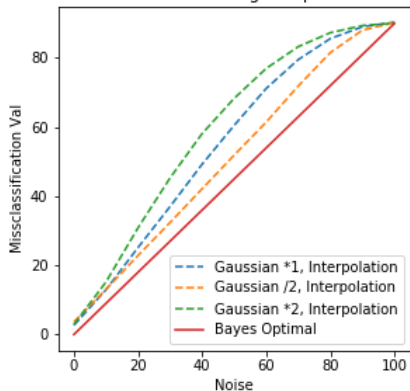
Label	MNIST	CIFAR-10	Synthetic 1	Synthetic 2
Original	8	>300	1	205
Random	>300	>300	106	269

Laplaceian Kernel:

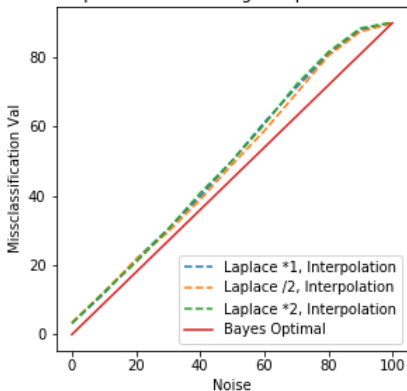
Label	MNIST	CIFAR-10	Synthetic 1	Synthetic 2
Original	3	5	1	5
Random	7	5	4	4

# Bandwidth

Gaussian, 10000 Training Samples (MNIST)



Laplace, 10000 Training Samples (MNIST)



# Conclusion

## Practical Considerations:

- Laplace Kernel ( $\sim \text{relu}$ ) often outperforms Gaussian
- Model scales with Training Data size

## We have seen:

- Kernel Learning == shallow Neural Networks
  - Interesting theoretical results
  - Surprisings experimental results
- Hope for new bounds