

# An Interpolation Perspective on Modern Machine Learning

Sept 2018

Purdue Workshop on Approximation and ML

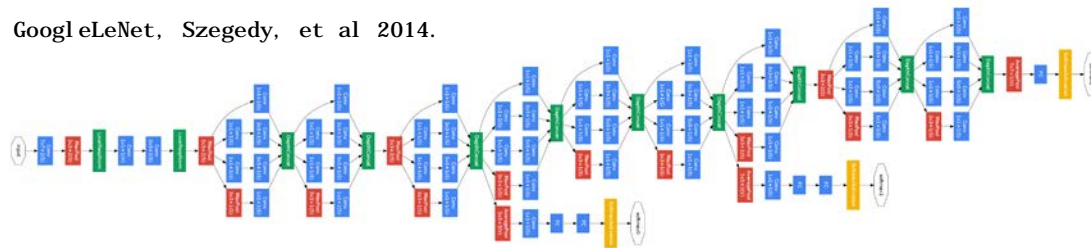
Mikhail Belkin, Ohio State University,  
Department of Computer Science and Engineering,  
Department of Statistics

Collaborators: Siyuan Ma, Soumik Mandal, Raef Bassily,  
Daniel Hsu, Partha Mitra, Alexander Rakhlin, Alexandre  
Tsybakov

---

Machine Learning/AI is becoming a backbone of commerce and society.

GoogleLeNet, Szegedy, et al 2014.



**The fog of war:**

What is new and what is important?  
Isolate and analyze components.

---



# Supervised ML

---

Data  $(x_i, y_i), i = 1..n, x_i \in \mathbb{R}^d, y_i \in \{-1,1\}$

**Goal**: construct  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , that “**generalizes**” to unseen data.

Empirical risk minimization (basis for most algorithms):

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum L(f(x_i), y_i)$$

$(f(x_i) - y_i)^2$ , e.g.

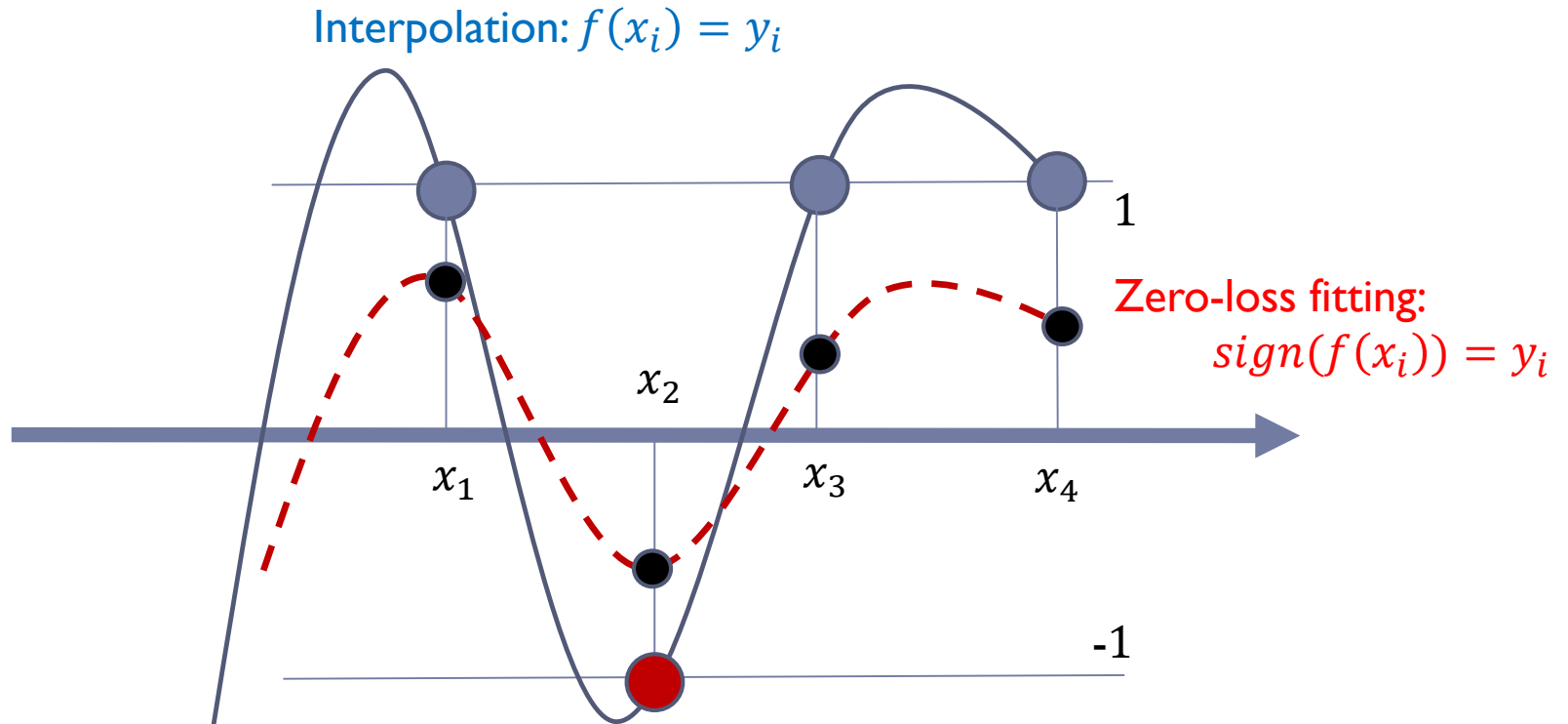
Optimized using **SGD**.



...

# Interpolation

Data  $(x_i, y_i)$ ,  $i = 1..n$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$



# Accepted wisdom

---

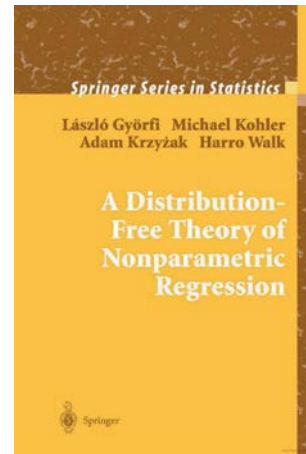
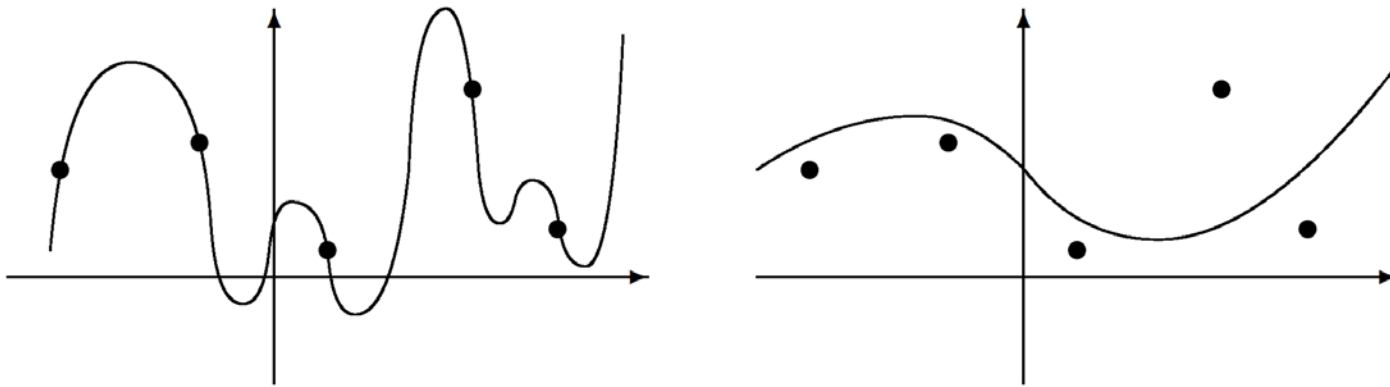


Figure 2.3. The estimate on the right seems to be more reasonable than the estimate on the left, which interpolates the data.



# Who is afraid of over-fitting?

Zero loss classifiers produce near optimal results.

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes	yes	100.0	89.05
		yes	no	100.0	89.31
		no	yes	100.0	86.03
		no	no	100.0	85.75

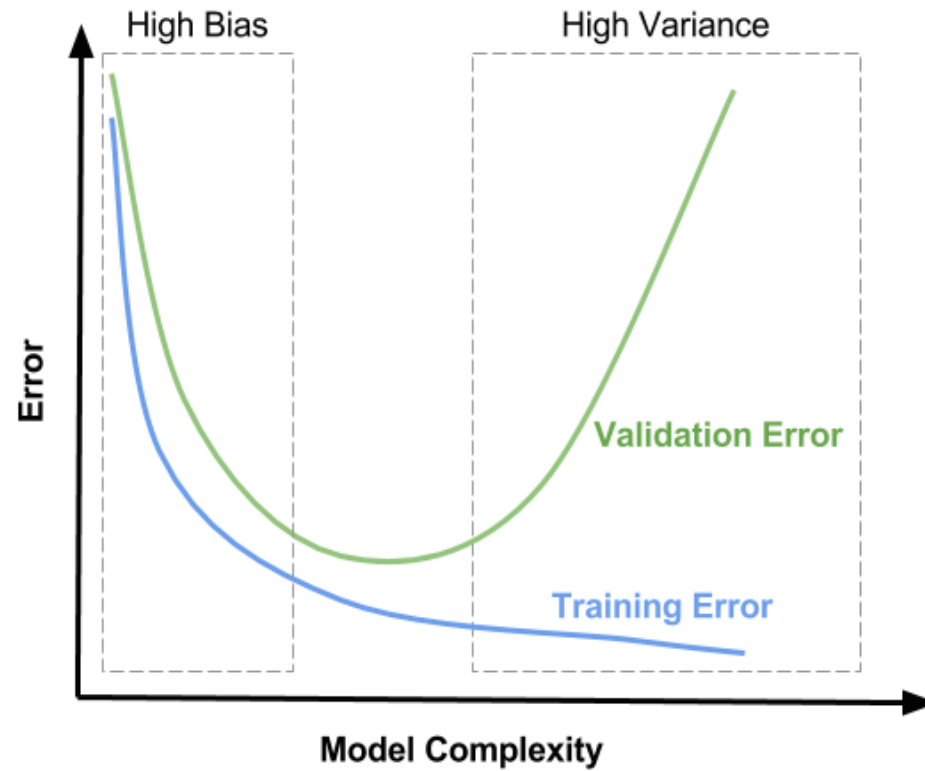
[CIFAR 10, from *Understanding deep learning requires rethinking generalization*, Zhang, et al, 2017]

Ruslan Salakhutdinov's tutorial on deep learning (Simons Institute, Berkeley, 2017):

The best way to solve the problem from practical standpoint is you build a very big system. If you remove any of these regularizations like dropout or L2, **basically you want to make sure you hit the zero training error**. Because if you don't, you somehow waste the capacity of the model.



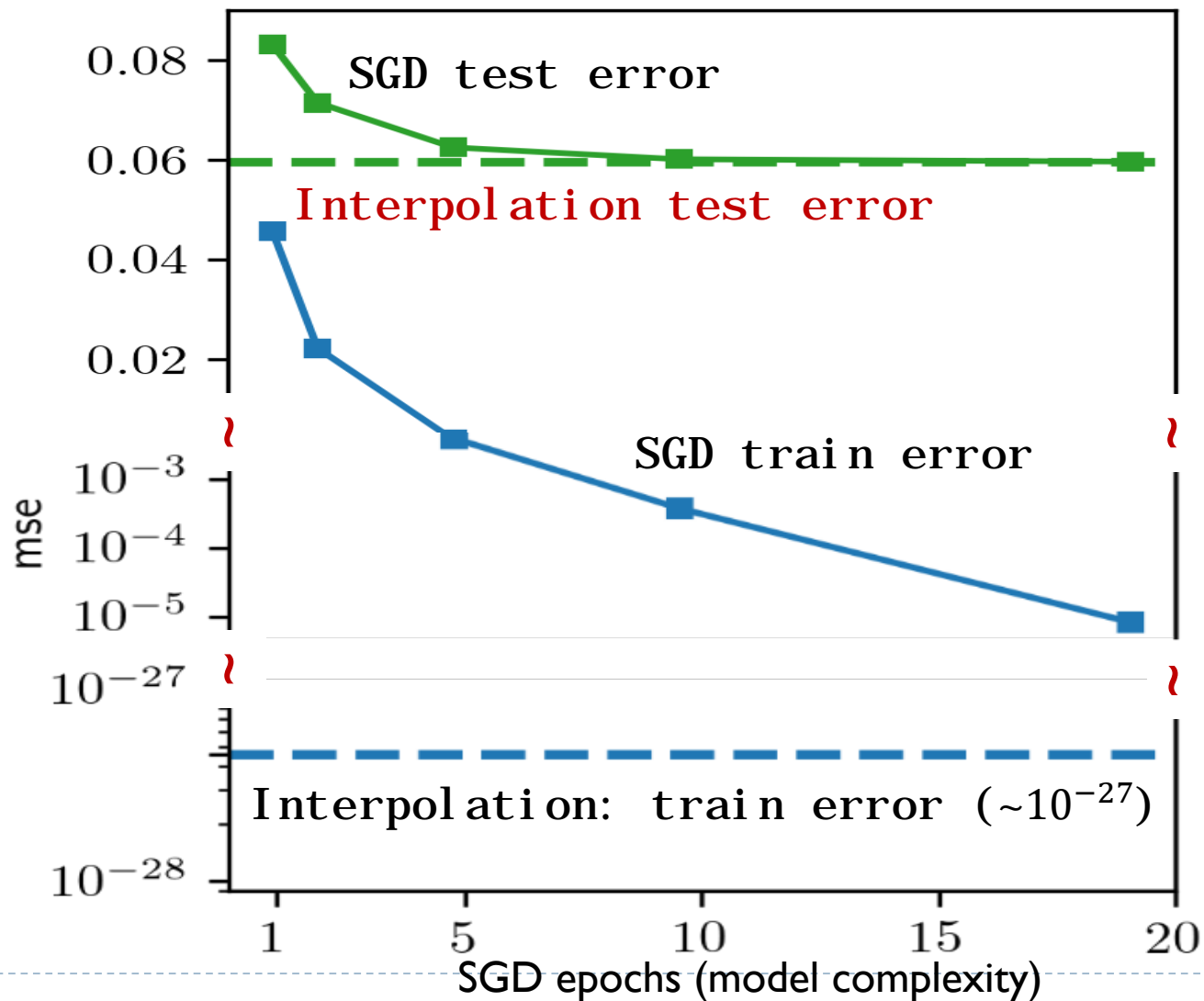
# Theory



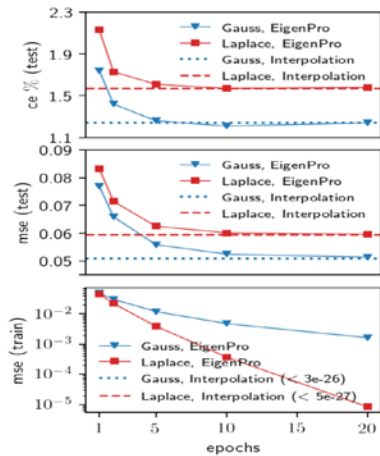
From <https://www.learnopencv.com/bias-variance-tradeoff-in-machine-learning/>



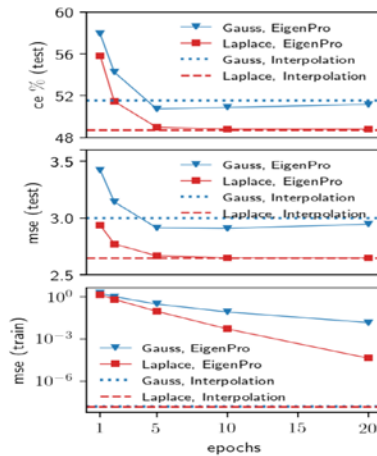
# Practice (SGD, kernel machine)



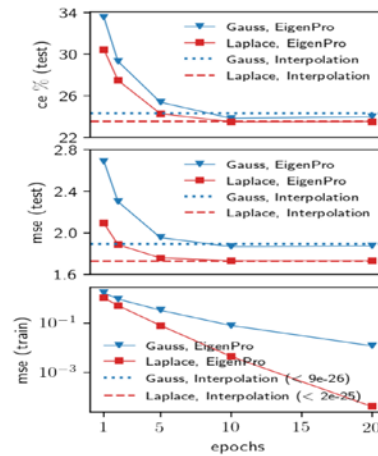




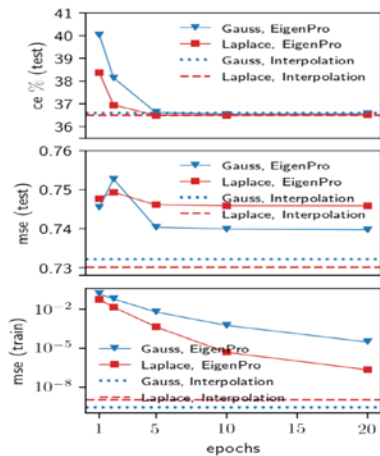
(a) MNIST



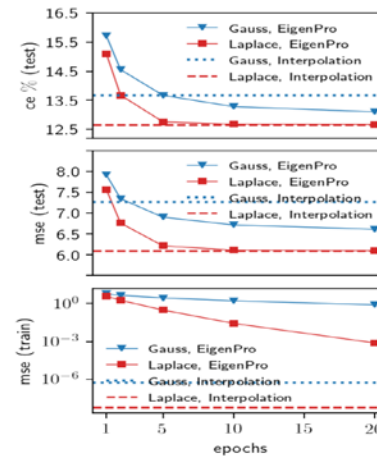
(b) CIFAR-10



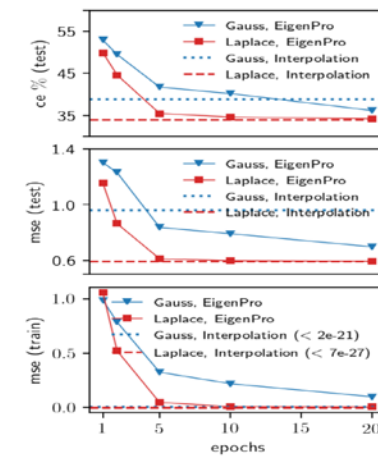
(c) SVHN ( $2 \cdot 10^4$  subsamples)



(d) TIMIT ( $5 \cdot 10^4$  subsamples)



(e) HINT-S ( $2 \cdot 10^4$  subsamples)



(f) 20 Newsgroups

[B., Ma, Mandal, ICML 18]

# A new phenomenon?

---

Interpolated classifiers produce near optimal results.

Deep Neural Networks (Zhang, et al, 17, but also much earlier)

Kernel Machines (our work)

Random Forests (PERT, Cutler, Zhao, 2001)

Adaboost (Schapire, et, al. 1998)

Not always recognized as such

(e. g. , regularization with very small  $\lambda$ ).

---



# A new phenomenon?

---

Leo Breiman, 1995

From **Reflections after refereeing papers for NIPS:**

- ▶ *Why don't heavily parametrized networks overfit the data?*
- ▶ *What is the effective number of parameters?*
- ▶ *When doesn't backpropagation head for poor local minima?*
- ▶ *When should one stop backpropagation and use the current parameters?*

**We are finally closing on some answers.**

---



# Two key questions

---

1. Why do interpolated classifiers generalize?

2. Why interpolate?

---



# The challenge of interpolation (approximation/statistics)

---

- Do we have theory?
  - ❑ Not much help from existing theory.
  - ❑ Interpolated classifiers are robust to label noise.
  - ❑ Unlikely to understand neural networks until (convex) kernel machines are understood.

[B., Ma, Mandal ICML 18]

- Moving forward: provable (near-optimal) generalization for methods that interpolate.

[B., Hsu, Mitra, 18] [B., Rakhlin, Tsybakov, 18]

---



# This talk: the power of interpolation (optimization)

---

## ➤ Optimization under interpolation

- ❑ Why is SGD so efficient in modern learning?
- ❑ Exponential convergence of mini-batch SGD.

[Ma, Bassily, B., ICML 18]

## ➤ Application: Fast and simple kernel machines for large data.

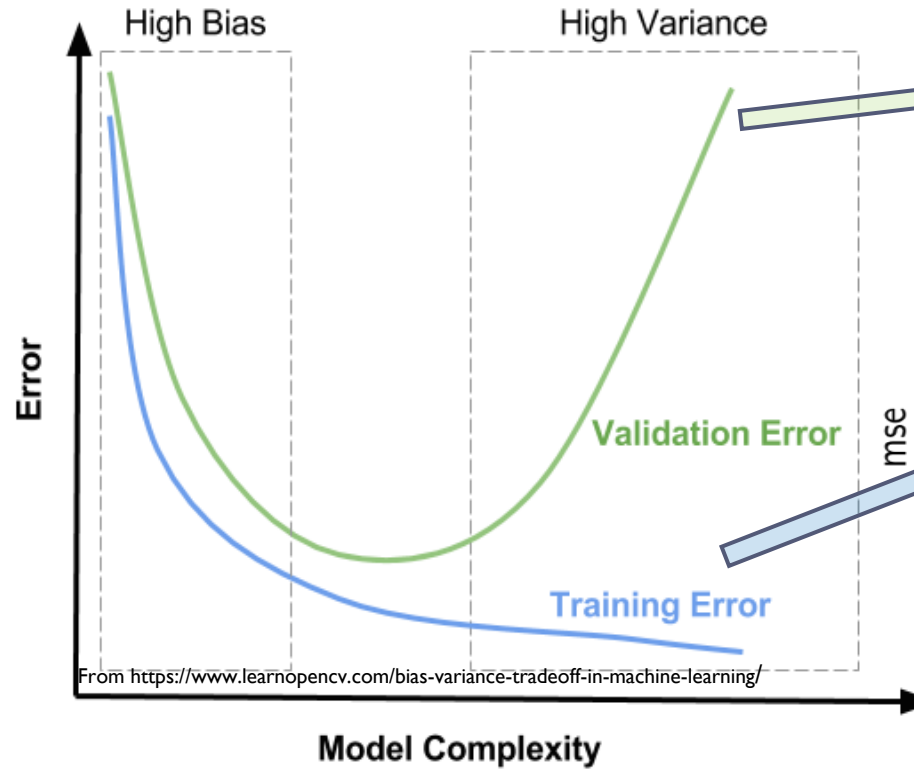
- ❑ **Simplicity:** no regularization or loss function
- ❑ Learning kernels that adapt to GPU.
- ❑ Automatic mini-batch/step size selection.

EigenPro 2.0 [Ma, B., NIPS 17, 18]

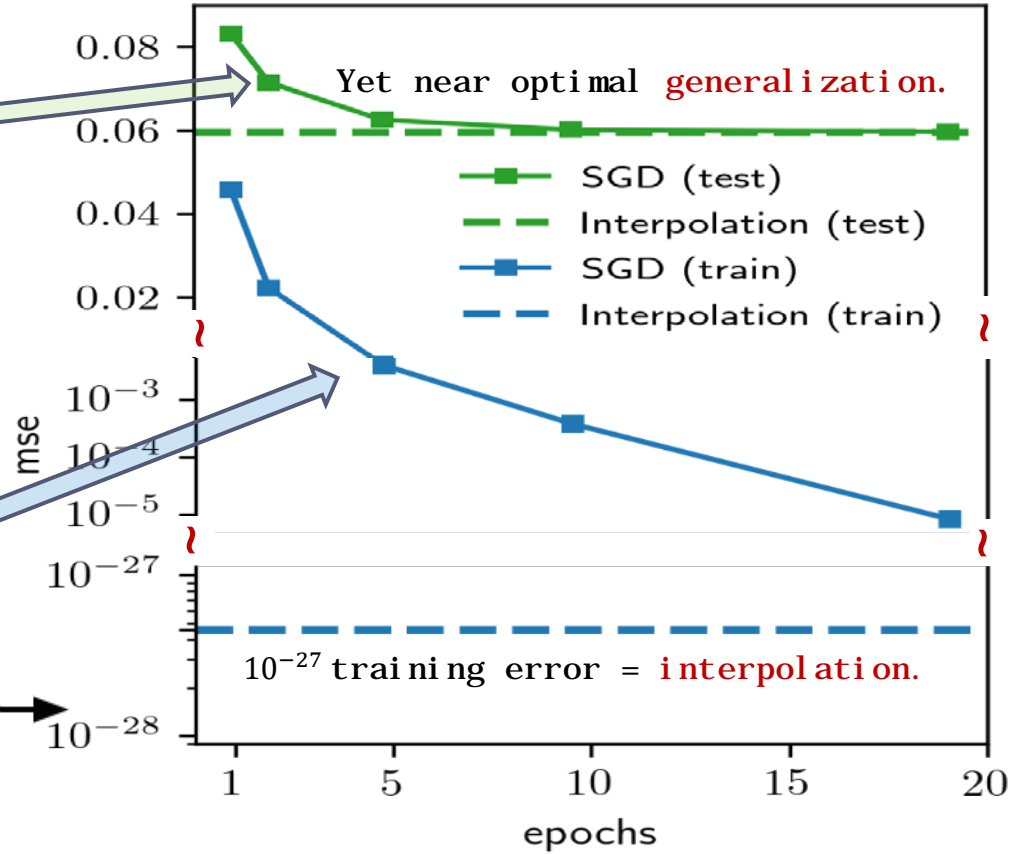
---



## Theory



## Practice (kernel machine, MNIST)



# Generalization bounds (interpolation)

---

## Basic bound:

VC-dim, fat shattering, Rademacher, Covering numbers, etc...

Expected loss

Empirical loss

Model or function complexity, e.g., VC or  $\|f\|_{\mathcal{H}}$

$$E(L(f^*, y)) \leq \frac{1}{n} \sum L(f^*(x_i), y_i) + O^* \left( \sqrt{\frac{c}{n}} \right)$$

Interpolation = 0

Can such bounds be useful?

---

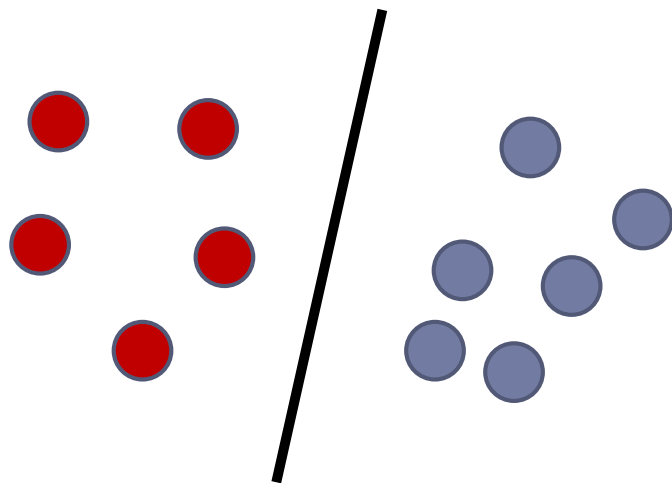




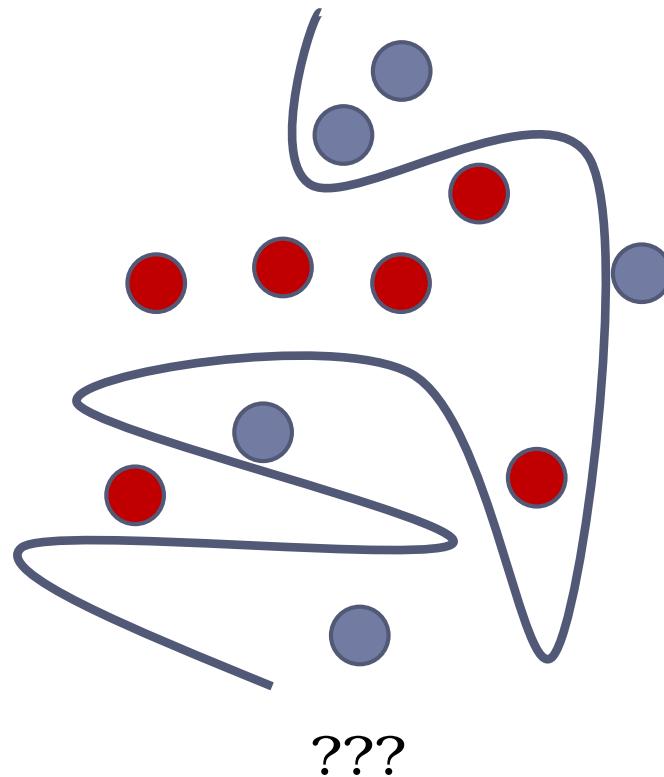
## Model complexity of interpolation?

---

Low complexity, does not grow with data size.



High complexity, grows linearly with data size.



Many classical results available.

Margin bounds,  
[Schapire, et al 98], etc.

---

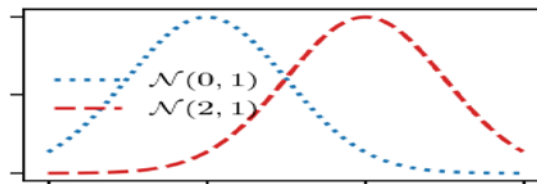


## How to test model complexity?

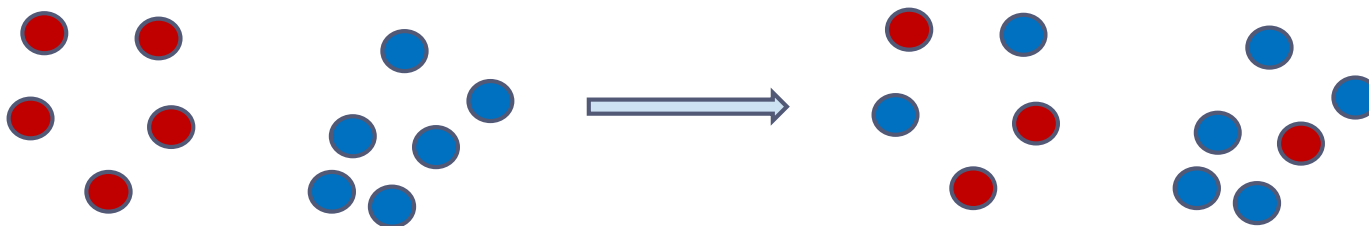
---

Two ways:

1. Synthetic data.



2. Add **label noise**.



Model complexity grows but Bayes opt. does not change!  
Expect **overfitting** to become severe as model complexity grows.

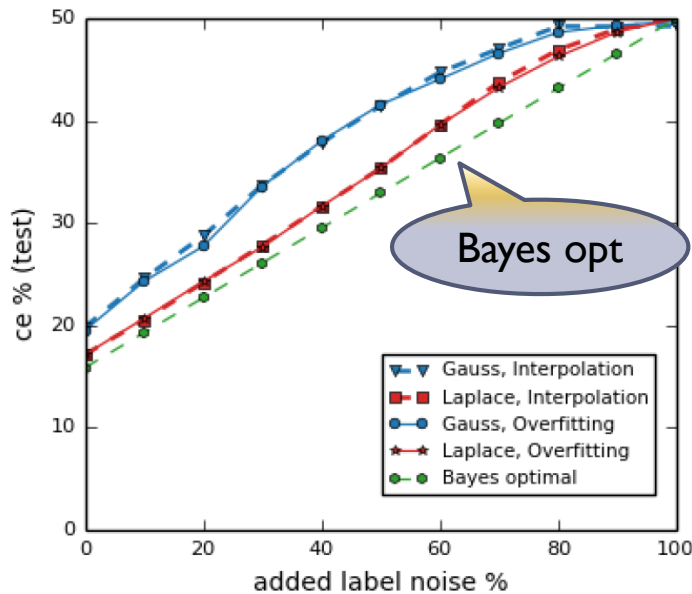
---



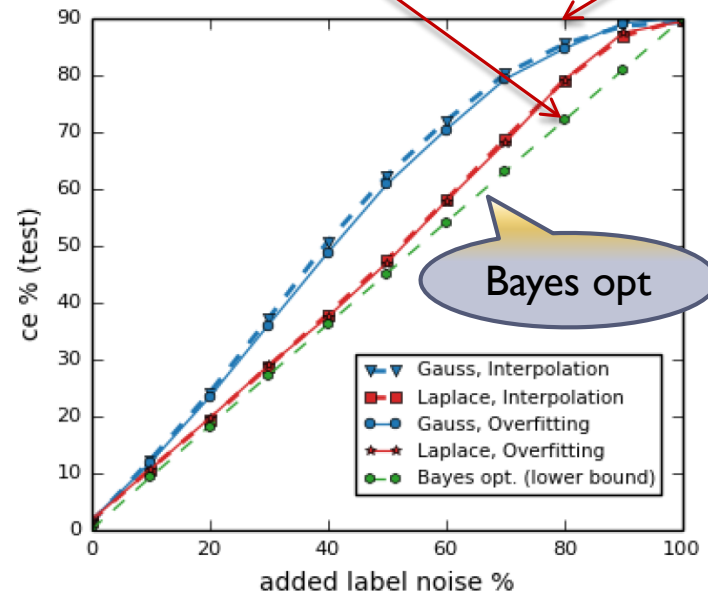
# Robustness to noise

What kind of **generalization bound** could work here?  
(hopefully **correct** but **nontrivial**)

$$0.7 < O^* \left( \sqrt{\frac{c(n)}{n}} \right) < 0.9$$



(a) Synthetic-2



(b) MNIST

---

$$0.7 < O^* \left( \sqrt{\frac{c(n)}{n}} \right) < 0.9 \quad n \gg 1$$

There are no bounds like this!  
Not clear whether they can exist mathematically.

e.g., can it be true if  $c(n) = \phi(\|f\|_{\mathcal{H}})$ ?



# The challenge of interpolation

---

- Do we have theory?
  - ❑ Not much help from existing theory.
  - ❑ Interpolated classifiers are robust to label noise.
  - ❑ Unlikely to understand neural networks until (convex) kernel machines are understood.

[B., Ma, Mandal ICML 18]

- Moving forward: provable (near-optimal) generalization for methods that interpolate.

[B., Hsu, Mitra, 18] [B., Rakhlin, Tsybakov, 18]

---



# Theoretical analyses fall short

---

- ▶ **VC-dimension/Rademacher complexity/covering bounds.**
  - ▶ Cannot deal with interpolated classifiers when Bayes risk is non-zero.
  - ▶ Empirical risk is zero – hard to bound the gap.
- ▶ **Regularization-type analyses (Tikhonov, early stopping, etc.)**
  - ▶ Diverge as  $\lambda \rightarrow 0$  for fixed  $n$ .
- ▶ **Algorithmic stability.**
  - ▶ Does not apply when empirical risk is zero, expected risk non-zero.
- ▶ **Classical smoothing methods (i.e., Nadaraya–Watson).**
  - ▶ Most classical analyses do not support interpolation.  
(But Hilbert Regression Scheme [Devroye, et al, 1998], 1-NN, our recent results)



# A way forward?

---

1-nearest neighbor classifier is very suggestive.

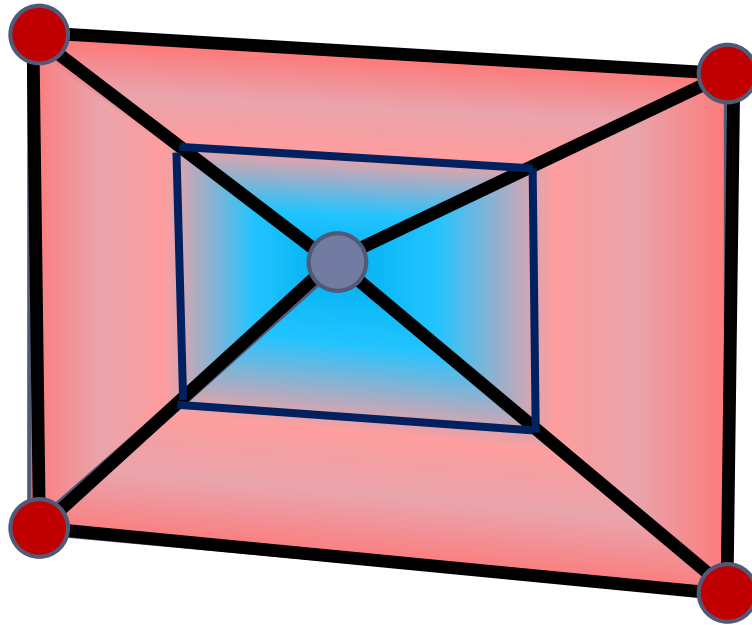
Interpolating classifier with a non-trivial (sharp!) performance guarantee (twice the Bayes risk).

- ▶ No margin assumptions.
- ▶ Analysis not based on uniform bounds.
- ▶ Estimate generalization, not the generalization gap.



# Simpli c i a l i n t e r p o l a t i o n

---



1. Tri angul ate.
2. Li nearl y i n t e r p o l a t e
3. Threshol d

[B., Hsu, Mi tra, 18]

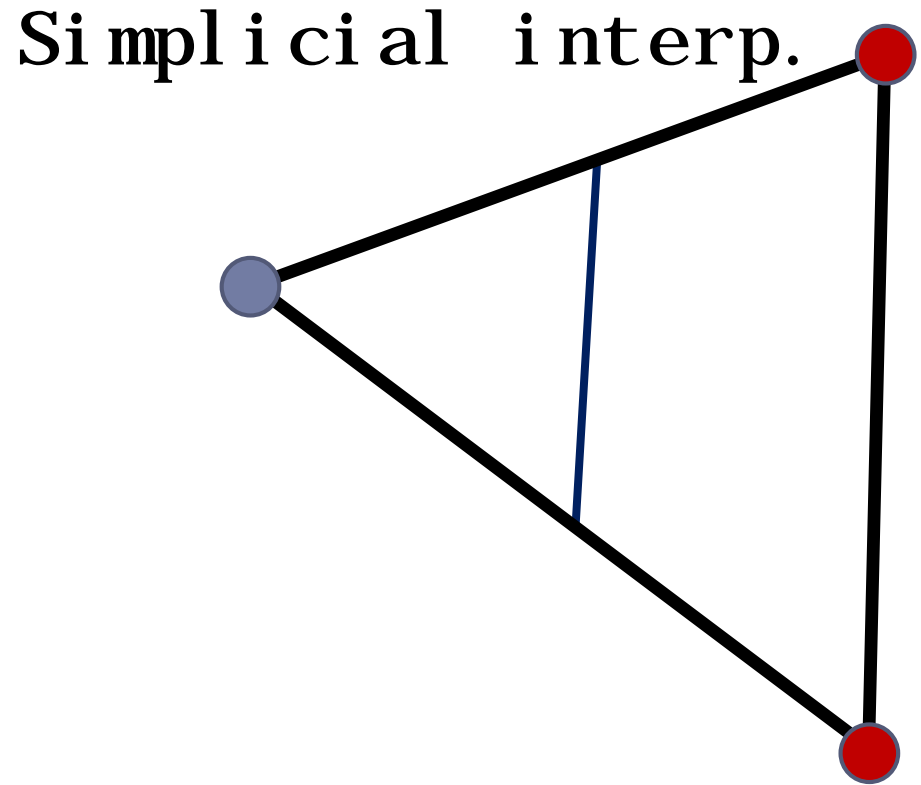
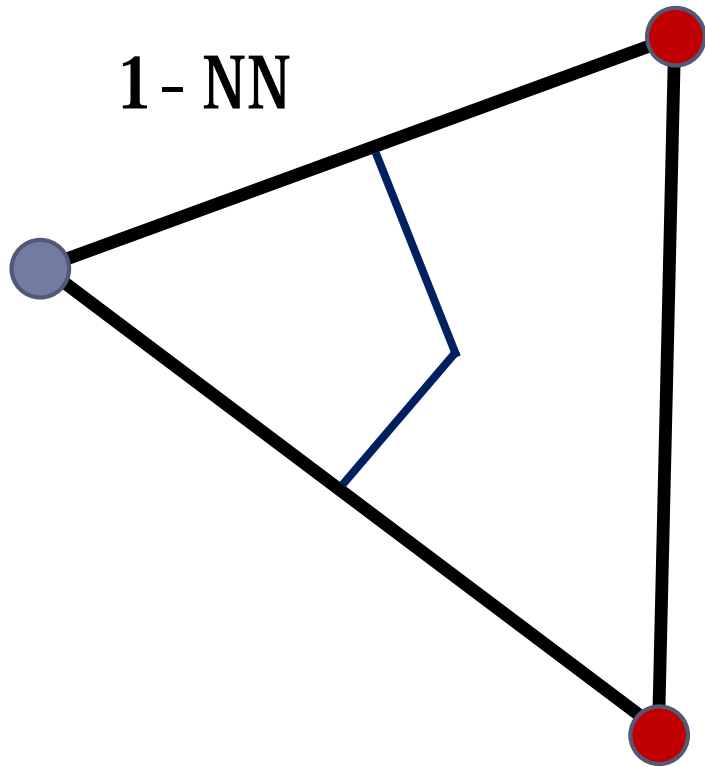
---





# 1-NN vs simplicial interpolation

---



# Nearly Bayes optimal

---

**Theorem:** (dimension  $d$ )

$$E(L(SI)) < \left(1 + \frac{1}{2^d}\right) \times \text{Bayes Risk}$$

Classical bound:

$$E\left(L(1_{NN})\right) < 2 \times \text{Bayes Risk}$$

The blessing of dimensionality!

# Interpolated k-NN schemes

---

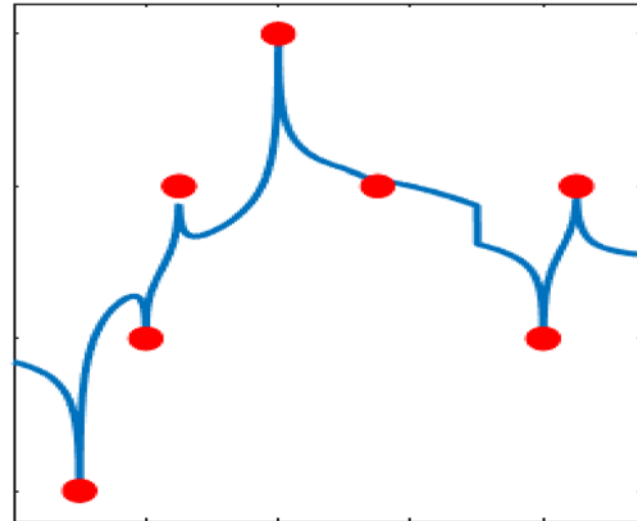
$$y(x) = \text{sign} \left( \sum w(x, x_i) y_i \right)$$

Singular kernel, e. g.  $w(x, z) = -\ln ||x - z||$

**Theorem:**

Weighted (interpolated) k-NN schemes with certain singular kernels are consistent.

Moreover, statistically optimal!



---

[B., Hsu, Mitra, 18] [B., Rakhlin, Tsybakov, 18]



# This talk: the power of interpolation (optimization)

---

## ➤ Optimization under interpolation

- ❑ Why is SGD so efficient in modern learning?
- ❑ Exponential convergence of mini-batch SGD.

[Ma, Bassily, B., ICML 18]

## ➤ Application: Fast and simple kernel machines for large data.

- ❑ **Simplicity:** no regularization or loss function
- ❑ Learning kernels that adapt to GPU.
- ❑ Automatic mini-batch/step size selection.

EigenPro 2.0 [Ma, B., NIPS 17, 18]

---



# Stochastic Gradient Descent

---

$$w^* = \operatorname{argmin}_w L(w) = \operatorname{argmin}_w \frac{1}{n} \sum L_i(w)$$

$= (f_w(x_i) - y_i)^2$ , e.g.

**SGD Idea:** optimize  $\sum L_i(w)$ ,  $m$  at a time.

Error after  $t$  steps

SGD:  $1/t$

GD:  $e^{-t}$

Why use SGD?

**All** major neural network architectures use SGD.

---



# The Power of Interpolation

---

**Key observation:**

Interpolation  $f_{w^*}(x_i) = y_i \Rightarrow \forall_i L_i(w^*) = 0$

Implies exponential convergence.

$(w \cdot 1 - 1)^2 + (w \cdot 1 + 1)^2$  [non-interpolation]

[SGD oscillates, adaptive step size/slow convergence]

$(w_1 \cdot 1 - 1)^2 + (w_2 \cdot 1 + 1)^2$  [interpolation]

[fixed step size, exponential convergence]

---



# Exponential convergence of $m$ -SGD

---

Convex loss function  $L$  ( $\lambda$ -smooth,  $\alpha$ -strongly convex),  
 $L_i$  ( $\beta$ -smooth).

**Theorem 1** [exponential convergence of  $m$ -SGD]

$$E L(w_{t+1}) \leq \frac{\lambda}{2} (1 - \eta^*(m)\alpha)^t \|w_1 - w^*\|$$

$$\eta^*(m) = \frac{m}{\beta + \lambda(m-1)}$$

[Ma, Bassily, B., ICML 18]

Related work ( $m = 1$ ): [Strohmer, Vershynin 09] [Moulines, Bach, 11]  
[Needell, Srebro, Ward, 14]

---

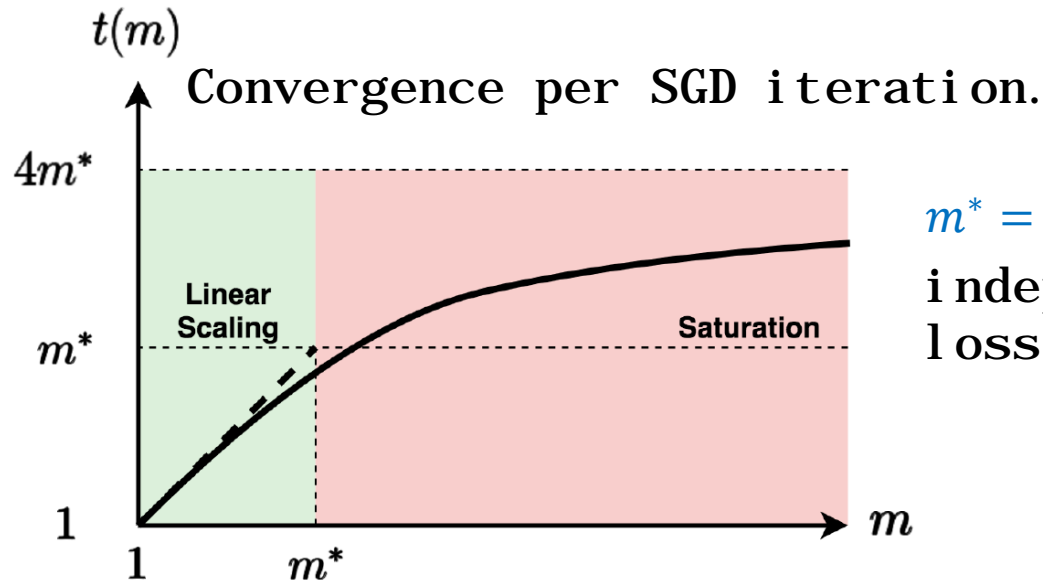


# Mini batch size?

**Theorem 2:** Critical size  $m^* = \frac{\beta}{\lambda}$  [optimal **fixed** step size]

[1. **linear scaling:**  $m \leq m^*$ ] One step  $m$ -SGD  $\approx m$  steps of 1-SGD

[2. **saturation:**  $m \geq m^*$ ] One step  $m$ -SGD  $\sim$  one step of **full GD**



$m^* = \frac{\beta}{\lambda}$  (nearly) data independent. Quadratic loss function:

$$m^* \sim \frac{\text{tr } H}{\lambda_1(H)}$$

$O(n)$  computational gain over GD



# Why mini batch?

---

**GPU:** fast highly parallel **matrix x matrix** products

→ limits algorithms available

Algorithmic requirements:

**matrix x matrix** products +

limited amount of other computation (CPU)

Full parallel minibatch computation  $1 \ll m \ll n$

(much larger than 1-SGD, much smaller than full GD)



Nvidia Titan X GPU (from nvidia.com)



# Interpolation in modern ML

How do you interpolate?

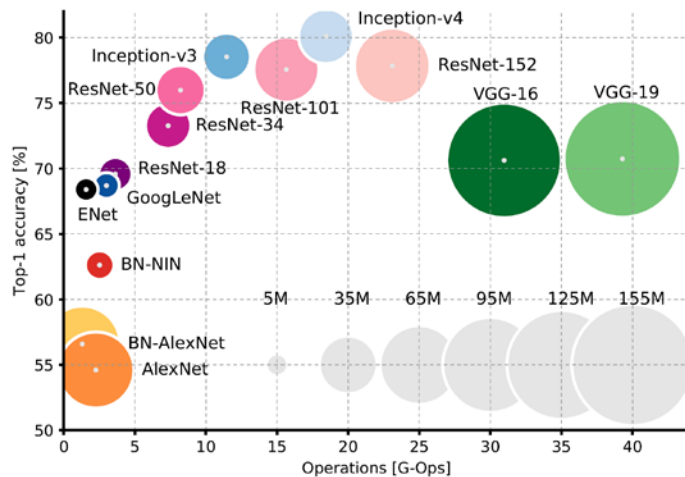
Think about a  $m \times n$  linear system:  $Ax = b$ .

Need rank  $\text{rank} A \geq n$  (at least as many parameters as equations).

Over-parametrization # parameters  $\geq$  # training data.

More parameters  $\rightarrow$  easier to interpolate (even non-convex).

From Canziani, et al., 2017.



Ruslan Salakhutdinov's tutorial on deep learning (Simons Institute, Berkeley, 2017):

*The best way to solve the problem from practical standpoint is you **build a very big system**. If you remove any of these regularizations like dropout or L2, **basically you want to make sure you hit the zero training error**. Because if you don't, you somehow waste the capacity of the model.*

# SGD in modern ML


---

Systematic over-parametrization.

# parameters  $\gg$  # training data

Over-parametrization  $\rightarrow$  interpolation  $\rightarrow$  fast SGD:

SGD  $O(n)$  computational gain ( $n \sim 10^5$ ) over GD  
+ GPU implementation  $\sim 100$  over CPU.



We used to do that..

Combined: SGD on GPU  $\sim 10^7$  faster than GD on CPU!

---



# This talk: the power of interpolation

---

## ➤ Optimization under interpolation

- ❑ Why is SGD so efficient in modern learning?
- ❑ Exponential convergence of mini-batch SGD.

[Ma, Bassily, B., ICML 18]

## ➤ Application: Fast and simple kernel machines for large data.

- ❑ Simplicity: no regularization or loss function
- ❑ Learning kernels that adapt to GPU.
- ❑ Automatic mini-batch/step size selection.

EigenPro 2.0 [Ma, B., NIPS 17, 18]

---



# Kernel machines

---

$\mathcal{H}$  (Reproducing Kernel Hilbert space),  
p. d. kernel  $K(x, z)$ , (e. g.,  $K(x, y) = e^{-\frac{\|x-y\|^2}{\sigma^2}}$ )

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}, f(x_i) = y_i} \|f\|_{\mathcal{H}}$$

Representer Theorem ( $K(x, y) = e^{-\frac{\|x-y\|^2}{\sigma^2}}$ , for example)

$$f^*(x) = \sum_i \alpha_i K(x_i, x), \quad \alpha = K^{-1}y$$

Minimum norm interpolation.

---



# Kernel learning

---

Beautiful classical  
statistical/mathematical theory

RKHS Theory [Aronszajn, ..., 50s]

Potential functions method [Izerman, ..., 60s]

Splines [Parzen, Wahba, ..., 1970-80s]

Kernel machines [Vapnik, ..., 90s]

Very attractive setting:

- Convex
- Analytically tractable

Can be viewed as a 2-layer neural net.

---



# Kernel Interpolation

---

No loss functions -- no regularization:

$$K \alpha^* = y$$

Direct inversion: cost  $n^3$  (does not map to GPU).

Gradient descent:  $\alpha^{(t)} = \alpha^{(t-1)} - \eta(K\alpha^{(t-1)} - y)$

[Richardson, Landweber, ...]

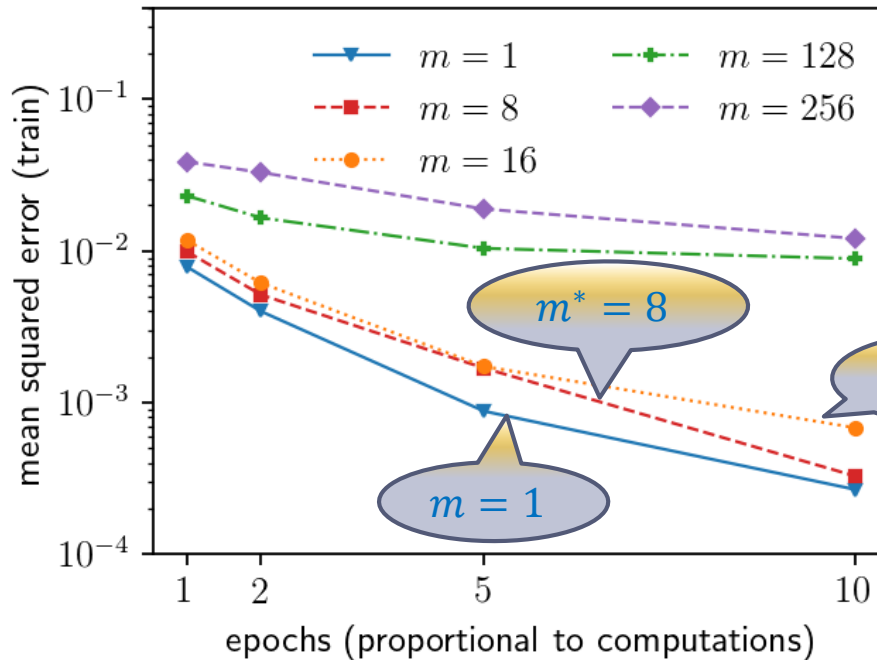
Cost  $n^2$  per iteration. GPU compatible.

How much gain from SGD?

---



# Real data example (Gaussian kernel)



Critical size  $m^* \approx \frac{\beta(K)}{\lambda_1(K)} + 1 \approx 8$ .

SGD acceleration factor over GD  $\sim 10^4$  (sequential computation)

But: parallelization beyond  $m^* = 8$  has little effect.

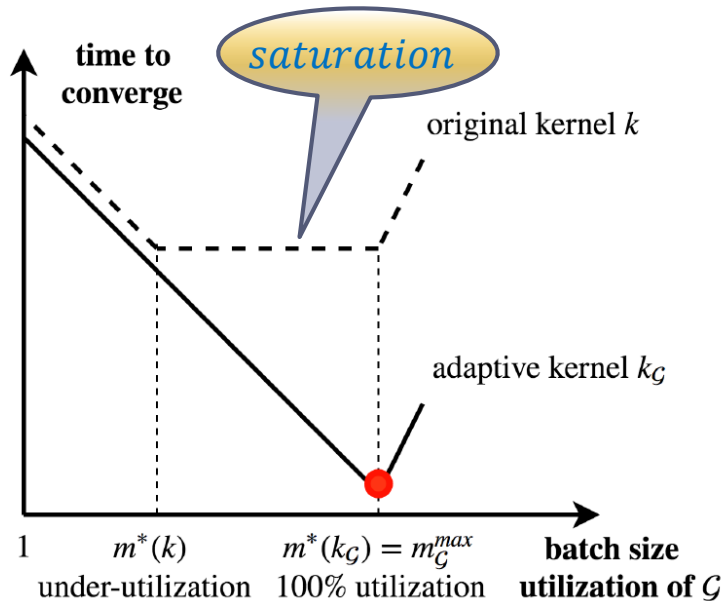
(a) MNIST (Gaussian,  $\sigma = 5$ )

$\beta = 1, \lambda_1 = 0.15, m^* \approx 8$



# Controlling parallelism

**Problem:** Parallelization is controlled by  $\frac{1}{\lambda_1}$ . (cf. convergence of gradient descent)



**Idea:** Make  $\lambda_1$  smaller to fully utilize parallel resource (GPU) without changing the original solution.

# Eigenvalue control

**Idea:** Change  $\lambda_1$  to fully utilize parallel resource (GPU).

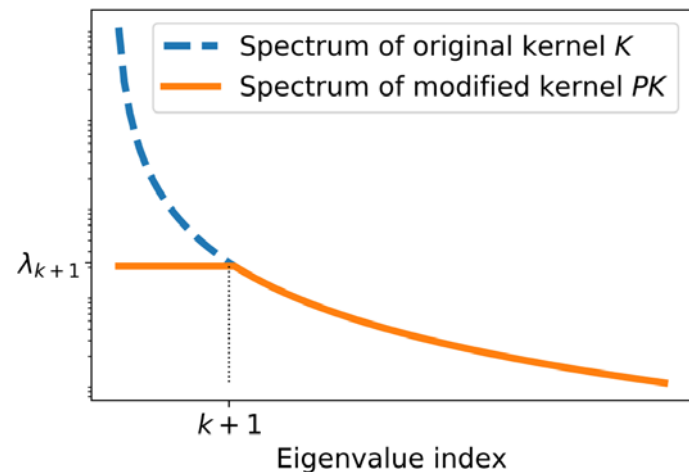
We have a tool: EigenPro kernel.

Original kernel:

$$K(x, z) = \sum_{i=1}^{\infty} \lambda_i e_i(x) e_i(z)$$

EigenPro kernel:

$$K_{EiP}(x, z) = \sum_{i=1}^k \lambda_{k+1} e_i(x) e_i(z) + \sum_{i=k+1}^{\infty} \lambda_i e_i(x) e_i(z)$$



[Ma, B. NIPS 2017]

# Comparisons to state-of-the-art

Dataset	Size	Our Method (use 1 GTX Titan Xp)		Results of Other Methods		
		error	GPU time	resource time	error	reference
MNIST	$1 \cdot 10^6$ / $6.7 \times 10^6$	0.67%	<b>21 m</b>	4.8 h on 1 GTX Titan X	0.70%	EigenPro [MB17]
				1.1 h on 1344 AWS vCPUs	0.72%	PCG [ACW16]
				less than 37.5 hours on 1 Tesla K20m	0.85%	[LML+14]
ImageNet†	$1.3 \times 10^6$	20.6%	<b>40 m</b>	4 h on 1 Tesla K40c	20.7%	FALKON [RCR17]
				-	19.9%	Inception-ResNet-v2 [SIVA17]
TIMIT	$1.1 \cdot 10^6$ / $2 \cdot 10^6$	31.6%	34 m (4 epochs)	3.2 h on 1 GTX Titan X	31.7%	EigenPro [MB17]
				1.5 h on 1 Tesla K40c	32.3%	FALKON [RCR17]
				512 IBM Blue Gene/Q cores	33.5%	Ensemble [HAS+14]
		32.0%	17 m (2 epochs)	7.5 h on 1024 AWS vCPUs	33.5%	BCD [TRVR16]
				multiple AWS g2.2xlarge instances	32.4%	DNN [MGL+17]
				multiple AWS g2.2xlarge instances	30.9%	SparseKernel [MGL+17] (use learned features)
SUSY	$4 \cdot 10^6$	19.7%	<b>48 s</b>	6 m on 1 GTX Titan X	19.8%	EigenPro [MB17]
				4 m on 1 Tesla K40c	19.6%	FALKON [RCR17]
				36 m on IBM POWER8	$\approx 20\%$	Hierarchical [CAS16]



## Interactive ML

---

Dataset	Size	Our method	LibSVM
TIMIT	$1 \cdot 10^5$	<b>15 s</b>	1.6 h
SVHN	$7 \cdot 10^4$	<b>13 s</b>	3.8 h
MNIST	$6 \cdot 10^4$	<b>6 s</b>	9 m
CIFAR-10	$5 \cdot 10^4$	<b>8 s</b>	3.4 h

Smaller datasets take seconds.

No optimization parameters to select.

---



# Take away messages

---

## The Power of Interpolation:

1. Explains fast SGD in modern ML.
2. Allows for much simpler analysis.
3. Leads to very efficient kernel machines, adaptive to modern hardware.



---

## Deep learning:

 overparametrization  
interpolation  generalizes well!  
fast SGD  
GPU

+ convolutional structures

---



# The challenge of interpolation

---

## Why do interpolated classifiers generalize?

Ubiquitous: Deep Neural Networks, Kernel Machines, Random Forests, Adaboost

Generalization is probably not (primarily) determined by:

- ▶ Non-convexity
- ▶ Regularization
- ▶ Loss functions
- ▶ Deep architectures
- ▶ Specific properties of optimization algorithms

Inductive bias is clearly important.

Time to revisit high-dimensional statistics!

---

