Peter Bartlett CS, Statistics, Simons Institute for the Theory of Computing UC Berkeley

November 2, 2019





1/26

Tension:

computation versus statistics optimization versus estimation

• A central feature of machine learning (past, present, and future).

Tension:

computation versus statistics optimization versus estimation

• A central feature of machine learning (past, present, and future).

Outline

Deep learning

Tension:

computation versus statistics optimization versus estimation

• A central feature of machine learning (past, present, and future).

Outline

- Deep learning
 - Nonparametric statistical methodology

Tension:

computation versus statistics optimization versus estimation

• A central feature of machine learning (past, present, and future).

Outline

- Deep learning
 - Nonparametric statistical methodology
 - Benign overfitting
 - Implicit regularization

Tension:

computation versus statistics optimization versus estimation

• A central feature of machine learning (past, present, and future).

Outline

- Deep learning
 - Nonparametric statistical methodology
 - Benign overfitting
 - Implicit regularization

Occupational complexity of estimation problems

Computer Vision



Easture Name	Value
readure realine	Table .
Description	{ "type": 0, "captions": [{ "text": "a man swimming in a pool of water", "confidence": 0.7850108124440484)] }
Tags	[['name': 'water', 'confidence': 0.999644279479805.); 'name': 'sport', 'confidence': 0.95049235583707.); 'name': 'smining', 'confidence': 0.90628182883010, 'hint': 'sport', !('name': 'poot', 'confidence': 0.8787588477134705.) ('name': 'mater sport', 'confidence': 0.631849467754364, 'hint': 'sport'.)]
Image Format	Jpeg
Image Dimensions	1500 x 1155
Clip Art Type	0 Non-clipart
Line Drawing Type	0 Non-LineDrawing
Black & White Image	False



a green jacket. a white horse. a man on a horse. two people riding horses. man wearing a green jacket the helmet is black. Horwn horse with white mane. white van parked on the street. a paved sidewalk. green and yellow jacket. a helmet on the head, white horse with white face.



bus parked on the street, a city street scene. front windshield of a bus, man vaking on sidewalk, a silver car parked on the street, a city scene, a green traffic light, a building in the background, the bus has a number, a large building, a brick building, red brick building with windows, a blue sign with a with earrow, while lines on the read.



a man on a skateboard. man riding a bicycle, orange cone on the ground. man riding a bicycle, two people riding a skateboard. red helmet on the man. skateboard on the ground, white shirt with red and white stripes, orange and white cone, trees are behind the people.

(microsoft.com) (Johnson et al, 2016

Speech Recognition



(bgr.com) (~ 4 / 26

・ロン ・四マ ・ヨマ ・ヨマ

Natural language processing

N



(google.com, skype.com)

Deep learning to date: an era of craftsmanship

• Engineering solutions to practical problems, striving for good performance on benchmark datasets.

Deep learning to date: an era of craftsmanship

- Engineering solutions to practical problems, striving for good performance on benchmark datasets.
- Ingredients: big data + fast computation (GPUs/TPUs) + ... ?

Deep learning to date: an era of craftsmanship

- Engineering solutions to practical problems, striving for good performance on benchmark datasets.
- Ingredients: big data + fast computation (GPUs/TPUs) + ... ?

Deep learning theory

• There are huge gaps in our understanding of deep learning.

Deep learning to date: an era of craftsmanship

- Engineering solutions to practical problems, striving for good performance on benchmark datasets.
- Ingredients: big data + fast computation (GPUs/TPUs) + ... ?

Deep learning theory

- There are huge gaps in our understanding of deep learning.
- Practitioners stumbled on methodology that contradicts established statistical wisdom.

Deep learning to date: an era of craftsmanship

- Engineering solutions to practical problems, striving for good performance on benchmark datasets.
- Ingredients: big data + fast computation (GPUs/TPUs) + ... ?

Deep learning theory

- There are huge gaps in our understanding of deep learning.
- Practitioners stumbled on methodology that contradicts established statistical wisdom.
- Key challenge: develop the theoretical foundations that will allow us to understand this new methodology and to improve it.

Dept of Systems Engineering, ANU



Brian Anderson

John Moore

A Digression: Model Reference Adaptive Control

United States Patent Office

A321.229 MODEL REFERENCE ADAPTIVE CONTROL SYSTEM Ress Kone, Neth BASSIE Henry Philip Whielow, Prostatplans, Mass. sing and Yornso tobers, Osco Hill, Ma. subjects to Massuchastin Institute of Technology, Cambridge, Mess, a corporation of Missochetty, and the second sec

Filed Jun, 22, 1962, Ser, No. 168,584 7 Claims, (Cl. 318-18)

the performance of the control vestors automatically is

The October 1960 issue of "Electro-Technology" monterpretation. This invention relates in a more restricted matic adjustments to certain parameters in the control loop or loops to compensate for changes in plant or signal characteristics or in both plant and signal characteristics At rage 119 of this article, Mathias and Van Nice dis close a model reference adaptive fight control system deseloped and tested by H. P. Whitsher of the Masuchasette situte of Technology. The present invention relates to ingerovaments in adaptive control systems of this so-called "MTT" type.

These systems comparise a network of elements each having on input and an output and a number of samming performance function relating its input quantity to its an corpet quartity. The Whitaker system makes it possible to design a con-

trel system which adjusts its own controllable parameters featings in the presence of changing operating character- 45 stics. The novel feature is a reference model which stores the system's specifications and permits closed-loop error function measured during the normal operating ro-10

Optimum or fully adapted performances are achieved spending to a specified performance inlet. Use of the model permits design flexibility, since the model can be model permits occups involves, more on an the vehicle go and can eabilit nan-linear characteristics if the systems ary relatively loose the system model can be crude and simple, but, on the other hand, granter control of per-formance, and hence grader floatility, can be achieved at by designing a model of increased sochistication. self-insproving process in these price-ort systems has emproyed test puries, simpling of its over quantum, cross of services for a minimum point of an error function. reachieve adaptation within two or three time constants of the dominant response mades of the protect.

It is an object of the instant lavantien to make a significant subsction in the convergence time and to stimi-Features of the invention are that adaptation is contin-

3.221.229 Patented Nov. 30, 1965

abad with simple equipment. Another feature of the investion is the extension of

the model reference principle. In addition to the model system, which is codinatily in the form of an electrica network, this investion features in the adaptation circuit of a filter which is the reciprocal of certain of the forward to path elements, one such filter being required for each of

the adaptable parameters of the control system features and admetance of the impaction will be aware hended from the following specifications and annexed

FIG. 1 is a block diagram of the model reference adap-

tive centrel system; FIG. 2 is a mathematical block diagram of the system WEIG I

FIG. 3 is a graph showing typical variation of integral squared enterior within the adjustable marameter:

FIGS. 4e and 4b are generalized control system block

te invention; PEC, 6 is a mathematical block diagram of an alreraft coll control system to which the invention is specied by

wing of specifies encouple: PROS. 76, 76, 76, 26, and 76 are graphs of error functions for the adaptive reli system of PIO. 6; FROS. 84, 84, 96, and 84 are graphs of the system re-spense of the adaptive system of FRO. 6 for supersonic

ght conditions; FIGS, 9a, 9b, 9c, and 9d are muths of the system represe for setsonic flight conditions;

FIG. 11 is a block diagram illustrating another alterna-

FIG. 1 is a simplified functional disprom of an adaptive centere. The dotted box 19 incloses the stabilization and the system 11 operator, in output quantity indicating unit 13, and a feedback path 14. The stabilization and directional control system II comprises various gyroscopes accelifion, and acumters which are necessary to conver the vilot's directional commands 15 to motions of the sir craft control surfaces. Equipment, perhaps involving gyroscoper, is included in the corput quantity indicating verted into suitable electrical signals for feedback 14 to the control restors 11. A dotted box 24 encloses the elements of the adaptation unit 29 which comprise an error signal comparator 28, a performance analyser 22, and an adjusting serve 23, the output 24 of which adjusts the variable persenties of the stabilization and directional control writers if. The performance analyzer 22 receives as inputs the Response Error Quantity from the comcommon symmers 12, as hereconstret expansed. Associated with the adaptation multi 20 is the performance reflerence model 34 designed to respond to the directional control commands 15 as the alternit 12 is desired to respond. Typically the model 30 will itself be subject to reference artises M whereis take-off weight, charges is balance and other yeth variables are accountedated. A second

seedback pairs 34 carries she electrical institutions of opwarater 21 where they may be compared to the electrical entputs 35 of the performance reference model. In these posterny there are bee different kinds of orner

quantities. One is the true system error quantity, [e]

Whitaker's "MIT Rule" circa 1961

Dept of Systems Engineering, ANU





John Moore

A Digression: Model Reference Adaptive Control

United States Patent Office

A321.229 MODEL REFERENCE ADAPTIVE CONTROL SYSTEM Ress Kone, Neth BASSIE Heavy Philip Whielow, Prostanglans, Mass. sing and Yornso tobers, Osco Hill, Ma. subjects in Massichastin Initions of Technology, Cambridge, Mess, a corporation of Missochetty, and and the second message and the second secon

Filed Jun, 22, 1962, Ser, No. 168,584 7 Claims, (Cl. 318-18)

the performance of the control vestors automatically is

The October 1960 issue of "Electro-Technology" monterpretation. This invention relates in a more restricted matic adjustments to certain parameters in the control loop or loops to compensate for changes in plant or signal characteristics or in both plant and signal characteristics At rage 119 of this article, Mathias and Van Nice dis close a model reference adaptive fight control system deseloped and tested by H. P. Whitsher of the Masuchasette situte of Technology. The present invention relates to ingerovaments in adaptive control systems of this so-called "MTT" type.

These systems comparise a network of elements each having on input and an output and a number of samming performance function relating its input quantity to its an corpet quartity. The Whitaker system makes it possible to design a con-

trel system which adjusts its own controllable parameters featings in the presence of changing operating character- 45 stics. The novel feature is a reference model which stores the system's specifications and permits closed-loop error function measured during the normal operating ro-10

Optimum or fully adapted performances are achieved spending to a specified performance inlet. Use of the model permits design flexibility, since the model can be model permits occups involves, more on an the vehicle go and can eabilit nan-linear characteristics if the systems ary relatively loose the system model can be crude and simple, but, on the other hand, granter control of per-formance, and hence gradet fields in the achieved at by designing a model of increased sochistication. self-insproving process in these price-ort systems has emproyed test puries, simpling of its over quantum, cross of services for a minimum point of an error function. reachieve adaptation within two or three time constants of the dominant response mades of the protect.

It is an object of the instant lavantien to make a significant subsction in the convergence time and to stimi-Features of the invention are that adaptation is contin-

3.221.229 Patented Nov. 30, 1965

- abad with simple equipment. Another feature of the investion is the extension of the model reference principle. In addition to the model system, which is codinatily in the form of an electrica network, this investion features in the adaptation circuit of a filter which is the reciprocal of certain of the forward to path elements, one such filter being required for each of
- the adaptable parameters of the control system 06features and admetance of the impaction will be aware hended from the following specifications and annexed
- FIG. 1 is a block diagram of the model reference adaptive centrel system;
- FIG. 2 is a mathematical block diagram of the system WEIG I
- FIG. 3 is a graph showing typical variation of integral squared enterior within the adjustable marameter:
- FIGS. 4e and 4b are generalized control system block
- te invention; PEC, 6 is a mathematical block diagram of an alreraft coll control system to which the invention is specied by
- wing of specifies encouple: PROS. 76, 76, 76, 26, and 76 are graphs of error functions for the adaptive reli system of PIO. 6; FROS. 84, 84, 96, and 84 are graphs of the system re-spense of the adaptive system of FRO. 6 for supersonic
- ght conditions; FIGS, 9a, 9b, 9c, and 9d are muths of the system represe for setsonic flight conditions;
- FIG. 11 is a block diagram illustrating another alterna-
- FIG. 1 is a simplified functional disprom of an adaptive centere. The dotted box 19 incloses the stabilization and the system 11 operator, in output quantity indicating unit 13, and a feedback path 14. The stabilization and directional control system II comprises various gyroscopes accelifion, and acumters which are necessary to conver the vilot's directional commands 15 to motions of the sir craft control surfaces. Equipment, perhaps involving gyroscoper, is included in the corput quantity indicating verted into suitable electrical signals for feedback 14 to the control restors 11. A dotted box 24 encloses the elements of the adaptation unit 29 which comprise an error signal comparator 28, a performance analyser 22, and an adjusting serve 23, the output 24 of which adjusts the variable persenties of the stabilization and directional control writers if. The performance analyzer 22 receives as inputs the Response Error Quantity from the comcommon symmers 12, as hereconstret expansed. Associated with the adaptation multi 20 is the performance reflerence model 34 designed to respond to the directional control commands 15 as the alternit 12 is desired to respond. Typically the model 30 will itself be subject to reference
- artises M whereis take-off weight, charges is balance and other yeth variables are accountedated. A second seedback pairs 34 carries she electrical institutions of opwarater 21 where they may be compared to the electrical entputs 35 of the performance reference model.
- In these posterny there are bee different kinds of orner quantities. One is the true system error quantity, [e]

Whitaker's "MIT Rule" circa 1961

$$heta_{k+1} = heta_k - \gamma rac{d\hat{J}(heta_k)}{d heta}; \ J = \int (y_r(t) - y_p(t))^2 \ dt.$$

Dept of Systems Engineering, ANU



Brian Anderson

John Moore

Deep compositions of nonlinear functions, trained on data with SGD

$$h = h_L \circ h_{L-1} \circ \cdots \circ h_1,$$

$$h_i : x \mapsto \sigma(W_i x)$$

Deep compositions of nonlinear functions, trained on data with SGD

$$h = h_L \circ h_{L-1} \circ \cdots \circ h_1,$$

 $h_i : x \mapsto \sigma(W_i x)$
with, e.g., $\sigma(v)_i = rac{1}{1 + \exp(-v_i)},$



Deep compositions of nonlinear functions, trained on data with SGD

$$h = h_L \circ h_{L-1} \circ \cdots \circ h_1,$$

$$h_i : x \mapsto \sigma(W_i x)$$

with, e.g., $\sigma(v)_i = \frac{1}{1 + \exp(-v_i)},$
 $\sigma(v)_i = \max\{0, v_i\},$





Deep compositions of nonlinear functions, trained on data with SGD

$$h = h_L \circ h_{L-1} \circ \cdots \circ h_1,$$

$$h_i : x \mapsto \sigma(W_i x)$$

with, e.g.,
$$\sigma(v)_i = \frac{1}{1 + \exp(-v_i)}, \qquad \sigma(v)_i = \max\{0, v_i\}$$

or "resnets," "max pooling," "attention," ...





Nonparametric Statistical Learning Methodology

Approximation

What kinds of functions can these compositions approximate well?

Nonparametric Statistical Learning Methodology

Approximation

What kinds of functions can these compositions approximate well?

Estimation

How can we effectively trade off complexity with sample size requirements?

Approximation

What kinds of functions can these compositions approximate well?

Estimation

How can we effectively trade off complexity with sample size requirements?

Optimization

How can we *efficiently* find a prediction rule that fits the data well?

(More precisely, how can we efficiently find a prediction rule that gives a good balance of complexity and fit to the data?)

Approximation

What kinds of functions can these compositions approximate well?

Estimation

How can we effectively trade off complexity with sample size requirements?

Optimization

How can we *efficiently* find a prediction rule that fits the data well?

(More precisely, how can we efficiently find a prediction rule that gives a good balance of complexity and fit to the data?)

Deep learning appears to give favorable trade-offs between these competing issues.

Approximation

What kinds of functions can these compositions approximate well?

Estimation

How can we effectively trade off complexity with sample size requirements?

Optimization

How can we *efficiently* find a prediction rule that fits the data well? (More precisely, how can we efficiently find a prediction r

(More precisely, how can we efficiently find a prediction rule that gives a good balance of complexity and fit to the data?)

Deep learning appears to give favorable trade-offs between these competing issues. We don't understand why.

Estimation

Typically, we aim for a trade-off between

• Fit to the training data,

• Complexity $\Omega(f)$ of a prediction rule

Estimation

Typically, we aim for a trade-off between

• Fit to the training data,

$$\hat{\mathcal{R}}(f) := \frac{1}{n} \sum_{i=1}^n \ell_i(f)$$

• Complexity $\Omega(f)$ of a prediction rule

Estimation

Typically, we aim for a trade-off between

• Fit to the training data, e.g.,

$$\hat{R}(f) := \frac{1}{n} \sum_{i=1}^{n} \ell_i(f) = \frac{1}{n} \sum_{i=1}^{n} (f(x_i) - y_i)^2$$

• Complexity $\Omega(f)$ of a prediction rule

Estimation

Typically, we aim for a trade-off between

• Fit to the training data, e.g.,

$$\hat{R}(f) := \frac{1}{n} \sum_{i=1}^{n} \ell_i(f) = \frac{1}{n} \sum_{i=1}^{n} (f(x_i) - y_i)^2$$

- Complexity $\Omega(f)$ of a prediction rule, e.g.,
 - Number of parameters
 - Norm of parameter vector
 - Norm of function in a reproducing kernel Hilbert space,
 - Bandwidth of smoothing kernel,
 - ...

Estimation

Typically, we aim for a trade-off between

• Fit to the training data, e.g.,

$$\hat{R}(f) := \frac{1}{n} \sum_{i=1}^{n} \ell_i(f) = \frac{1}{n} \sum_{i=1}^{n} (f(x_i) - y_i)^2,$$

- Complexity $\Omega(f)$ of a prediction rule, e.g.,
 - Number of parameters
 - Norm of parameter vector
 - Norm of function in a reproducing kernel Hilbert space,
 - Bandwidth of smoothing kernel,
 - ...

This is especially important for nonparametric methods, that is, those for which the number of parameters grows with the sample size.

Estimators for Inverse Problems: Convex Regularization

Found Comput Math (2012) 12:805-849 DOI 10.1007/s10208-012-9135-7



The Convex Geometry of Linear Inverse Problems

Venkat Chandrasekaran - Benjamin Recht -Pablo A. Parrilo - Alan S. Willsky

Received: 2 December 2010 / Revised: 25 February 2012 / Accepted: 3 July 2012 / Published online: 16 October 2012 (0 SFoCM 2012

Abstract In applications throughout science and engineering one is often faced with the challenge of ovlary and in-posed inverge poblem, where the number of available measurements is smaller than the dimension of the model to be estimated. However in many practical situations of interest, models are constrained structurally so that they only have a few degrees of freedom relative to their ambient dimension. This paper provides a general framework to counter thatoms of simplicity in courses penally interpreting a general framework to counter that one of the simplicity of the simple models considered in the work problems. The class of simple models considered includes those formed as the sum of a few atoms from some (possibly infinite) elementary stores jet; example

Communicated by Emmanuel Candès.

V. Chandrasekaran (⊠) Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125, USA e-mail: venkint 0%altech.eda

B. Recht Computer Sciences Department, University of Wisconsin, Madison, WI 53706, USA e-mail: brecht@cs.wisc.edu

P.A. Parillo - A.S. Willsky Laboratory for Information and Decision Systems, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

P.A. Parrilo e-mail: parrilo@mit.edu

A.S. Willsky e-mail: willsky@mit.edu

Estimators for Inverse Problems: Convex Regularization

Found Comput Math (2012) 12:805-849 DOI 10.1007/s10208-012-9135-7



The Convex Geometry of Linear Inverse Problems

Venkat Chandrasekaran - Benjamin Recht -Pablo A. Parrilo - Alan S. Willsky

Received: 2 December 2010 / Revised: 25 February 2012 / Accepted: 3 July 2012 / Published online: 16 October 2012 (0 SFoCM 2012

Abstract In applications throughout science and engineering one is often faced with the challenge of ovlary and in-posed inverge poblem, where the number of available measurements is smaller than the dimension of the model to be estimated. However in many practical situations of interest, models are constrained structurally so that they only have a few degress of freedom relative to their ambient dimension. This paper functions, resulting in convex espiratizations studies to historical underlearning of the structure of the structure of the structure of the structure of the vecue problems. The class of simple models considered includes those formed as the sum of a few atoms from some (possible) infinite) clementary atomic set; exam-

Communicated by Emmanuel Candès.

V. Chandrasekaran (⊠) Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125, USA e-mail: venkint 0%altech.eda

B. Recht Computer Sciences Department, University of Wisconsin, Madison, WI 53706, USA e-mail: becht@cs.wisc.edu

P.A. Parillo - A.S. Willsky Laboratory for Information and Decision Systems, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

P.A. Parrilo e-mail: parrilo@mit.edu

A.S. Willsky e-mail: willsky@mit.edu

Estimating simple structure

e.g., sparse linear regression, low rank tensors, rankings from surveys, ...

イロト 不得下 イヨト イヨト 二日

11/26

Estimators for Inverse Problems: Convex Regularization

Found Comput Math (2012) 12:805-849 DOI 10.1007/s10208-012-9135-7



The Convex Geometry of Linear Inverse Problems

Venkat Chandrasekaran - Benjamin Recht -Pablo A. Parrilo - Alan S. Willsky

Received: 2 December 2010 / Revised: 25 February 2012 / Accepted: 3 July 2012 / Published online: 16 October 2012 0 SFoCM 2012

Abstract In applications throughout science and engineering one is often fued with the challenge of obvious mill posed inverge problem, where the number of available measurements is smaller than the dimension of the model to be estimated. However in many practical situations of interest, media are constrained structurally so that they only have a few degrees of freedom relative to their ambient dimension. This paper functions, resulting in correct optimization solutions to linear, underletermined in verse problems. The class of simple models considered indicates those formed as the sum of a few atoms from some (possibility linitial) elementary atomic set; examtions from some from some (possibility linitial) elementary atomic set; exam-

Communicated by Emmanuel Candès.

V. Chandrasekaran (⊠) Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125, USA e-mail: venkint 0%altech.eda

B. Recht Computer Sciences Department, University of Wisconsin, Madison, WI 53706, USA e-mail: brecht@cs.wisc.edu

P.A. Parillo - A.S. Willsky Laboratory for Information and Decision Systems, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

P.A. Parrilo e-mail: parrilo@mit.edu

A.S. Willsky e-mail: willsky@mit.edu

Estimating simple structure

e.g., sparse linear regression, low rank tensors, rankings from surveys, ...

Minimum norm fit: $\Omega(f) = \|f\|_{\mathcal{A}}$

 $\begin{array}{ll} \min & \Omega(f) \\ \text{s.t.} & \hat{R}(f) \leq \epsilon. \end{array}$

- 32

11/26

Regularization: balancing \hat{R} and Ω

```
\min_f a\Omega(f) + \hat{R}(f),
```


Regularization: balancing \hat{R} and Ω

 $\min_f a\Omega(f) + \hat{R}(f), \quad ext{or} \quad \min \ \hat{R}(f)$ s.t. $\Omega(f) \leq b,$

Regularization: balancing \hat{R} and Ω

$$\begin{array}{ll} \min_{f} a\Omega(f) + \hat{R}(f), \quad \text{or} \quad \min \ \hat{R}(f) & \text{or} \quad \min \ \Omega(f) \\ & \text{s.t.} \ \Omega(f) \leq b, & \text{s.t.} \ \hat{R}(f) \leq c. \end{array}$$

・ロト ・回ト ・ヨト ・ヨト

3

12/26

Regularization: balancing \hat{R} and Ω

$$\min_{f} a\Omega(f) + \hat{R}(f), \quad \text{or} \quad \min \ \hat{R}(f) \quad \text{or} \quad \min \ \Omega(f) \\ \text{s.t.} \ \Omega(f) \le b, \quad \text{s.t.} \ \hat{R}(f) \le c.$$

Classical analysis: uniform convergence

$$\sup \left\{ R(f) - \hat{R}(f) : \Omega(f) \le b \right\} \le \epsilon,$$

$$R(f) := \mathbb{E}\ell(f) \quad \frac{1}{n} \sum_{i=1}^{n} \ell_i(f) =: \hat{R}(f).$$

<注><注><注><注><注><注><注><注><注><</p>

< 🗇 >

Regularization: balancing \hat{R} and Ω

$$\min_{f} a\Omega(f) + \hat{R}(f), \quad \text{or} \quad \min \ \hat{R}(f) \quad \text{or} \quad \min \ \Omega(f) \\ \text{s.t.} \ \Omega(f) \le b, \quad \text{s.t.} \ \hat{R}(f) \le c.$$

Classical analysis: uniform convergence

$$\sup\left\{R(f)-\hat{R}(f):\Omega(f)\leq b
ight\}\leq\epsilon,$$

so that $R(\hat{f}):=\mathbb{E}\ell(\hat{f})pproxrac{1}{n}\sum_{i=1}^n\ell_i(\hat{f})=:\hat{R}(\hat{f}).$

<日、<日、<日、<日、<日、<日、<日、<日、<日、<日、<日、<12/26

Regularization: balancing \hat{R} and Ω

$$\min_{f} a\Omega(f) + \hat{R}(f), \quad \text{or} \quad \min \ \hat{R}(f) \quad \text{or} \quad \min \ \Omega(f) \\ \text{s.t.} \ \Omega(f) \le b, \quad \text{s.t.} \ \hat{R}(f) \le c.$$

Classical analysis: uniform convergence

$$\sup\left\{R(f) - \hat{R}(f) : \Omega(f) \le b
ight\} \le \epsilon,$$

so that $R(\hat{f}) := \mathbb{E}\ell(\hat{f}) \approx rac{1}{n}\sum_{i=1}^{n}\ell_i(\hat{f}) =: \hat{R}(\hat{f}).$

To exploit uniform convergence, we should consider $c \ge R^* - \epsilon$, where $R^* = \min_f R(f)$.

12 / 26

<回> <き> <き> <き> = =

Regularization: balancing \hat{R} and Ω

$$\begin{split} \min_{f} \ a\Omega(f) + \hat{R}(f), \quad \text{or} \quad \min_{f} \ \hat{R}(f) & \text{or} \quad \min_{f} \ \Omega(f) \\ & \text{s.t.} \ \Omega(f) \leq b, & \text{s.t.} \ \hat{R}(f) \leq c. \end{split}$$

Classical analysis: uniform convergence

$$\sup \left\{ R(f) - \hat{R}(f) : \Omega(f) \le b \right\} \le \epsilon,$$

so that $R(\hat{f}) := \mathbb{E}\ell(\hat{f}) \approx \frac{1}{n} \sum_{i=1}^{n} \ell_i(\hat{f}) =: \hat{R}(\hat{f}).$

To exploit uniform convergence, we should consider $c \ge R^* - \epsilon$, where $R^* = \min_f R(f)$.



12/26

◆□→ ◆注→ ◆注→ □注□

• Deep networks can achieve zero training error (for regression loss)

- Deep networks can achieve zero training error (for regression loss)
- ... with near state-of-the-art performance

- Deep networks can achieve zero training error (for regression loss)
- ... with near state-of-the-art performance
- ... even for noisy problems $(R^* \gg 0)$.

- Deep networks can achieve zero training error (for regression loss)
- ... with near state-of-the-art performance
- ... even for noisy problems $(R^* \gg 0)$.
- No tradeoff between fit to training data and complexity!

- Deep networks can achieve zero training error (for regression loss)
- ... with near state-of-the-art performance
- ... even for noisy problems $(R^* \gg 0)$.
- No tradeoff between fit to training data and complexity!
- Deep networks seem to operate in the overfitting regime $(\hat{R}(f) \ll R^*)$ but still predict accurately.

- Deep networks can achieve zero training error (for regression loss)
- ... with near state-of-the-art performance
- ... even for noisy problems $(R^* \gg 0)$.
- No tradeoff between fit to training data and complexity!
- Deep networks seem to operate in the overfitting regime $(\hat{R}(f) \ll R^*)$ but still predict accurately.
- A new statistical phenomenon: benign overfitting.

Statistical Wisdom and Overfitting

"... interpolating fits... [are] unlikely to predict future data well at all."



"... a function which interpolates the data ... is not a reasonable estimate."



see also (B. and Rakhlin, Simons Institute, May 2019) ・ロン ・四 と ・ ヨ と ・

A new statistical phenomenon: good prediction with very small training error for regression loss

- Statistical wisdom says a prediction rule should not fit too well.
- But deep networks are trained to fit noisy data perfectly, and they predict well.





m. 35 – May 12, 2017 In program sims to entend the reach and impact of CS theory within machine were, by the address based questions in developing areas at plastice, vancing the alignment to teach of machine searcing, and putting widely-used unities on a first theoretical foundation.



Foundations of Deep Learning Mer 22 - Avr. 8, 2019

This program will being together researchers from academia and industry to develop empirically relevant theoretical foundations of deep teaming, with the sim of guiding the real-world use of deep teaming.

Progress in Benign Overfitting

• Simplicial interpolation (\approx nearest neighbor)

(Belkin, Hsu, Mitra, 2018)

• Nadaraya-Watson estimator with singular kernels

(Belkin, Hsu, Mitra, 2018; Belkin, Rakhlin, Tsybakov, 2018)

• Random matrix theory asymptotics $(d \asymp n)$ for linear regression, random nonlinear features

(Hastie, Montanari, Rosset, Tibshirani, 2019; Mei, Montanari, 2019; Belkin, Hsu and Xu, 2019)

• Certain reproducing kernel Hilbert spaces

(Liang and Rakhlin, 2018; Rakhlin and Zhai, 2018; Liang, Rakhlin, Zhai, 2019)

 Minimum norm linear regression: tight upper and lower bounds for finite sample, arbitrary dimension (B., Long, Lugosi, Tsigler, 2019)







(B., Long, Lugosi, Tsigler, 2019)

Characterizing benign overfitting in linear regression

For $\ell(f) = (f(x) - y)^2$,







(B., Long, Lugosi, Tsigler, 2019)

Characterizing benign overfitting in linear regression

For $\ell(f) = (f(x) - y)^2$, $\Omega(x \mapsto \langle x, \theta \rangle) = \|\theta\|_2$,







Characterizing benign overfitting in linear regression

For $\ell(f) = (f(x) - y)^2$, $\Omega(x \mapsto \langle x, \theta \rangle) = \|\theta\|_2$, and $\binom{x}{y} = \Phi z$ where Φ

is a bounded linear operator and z has subgaussian, independent entries,







Characterizing benign overfitting in linear regression

For $\ell(f) = (f(x) - y)^2$, $\Omega(x \mapsto \langle x, \theta \rangle) = \|\theta\|_2$, and $\begin{pmatrix} x \\ y \end{pmatrix} = \Phi z$ where Φ is a bounded linear operator and z has subgaussian, independent entries,

$$c_1\left(\frac{d^*}{n}+\frac{n}{R_{d^*}}+\phi\left(\frac{1}{n}\right)\right) \leq \mathbb{E}R(\hat{f})-R^* \leq c_2\left(\frac{d^*}{n}+\frac{n}{R_{d^*}}+\frac{1}{\sqrt{n}}\right).$$

where $d^* = \min\{d : r_d \ge c_3n\}$, r_d and R_d are effective ranks of the covariance of x in the subspace orthogonal to the d highest variance directions, and ϕ is increasing.







Characterizing benign overfitting in linear regression

For $\ell(f) = (f(x) - y)^2$, $\Omega(x \mapsto \langle x, \theta \rangle) = \|\theta\|_2$, and $\begin{pmatrix} x \\ y \end{pmatrix} = \Phi z$ where Φ is a bounded linear operator and z has subgaussian, independent entries,

$$c_1\left(rac{d^*}{n}+rac{n}{R_{d^*}}+\phi\left(rac{1}{n}
ight)
ight) \leq \mathbb{E}R(\hat{f})-R^*\leq c_2\left(rac{d^*}{n}+rac{n}{R_{d^*}}+rac{1}{\sqrt{n}}
ight),$$

where $d^* = \min\{d : r_d \ge c_3n\}$, r_d and R_d are effective ranks of the covariance of x in the subspace orthogonal to the d highest variance directions, and ϕ is increasing.

That is, benign overfitting occurs iff there is a subspace where the covariance has small magnitude, high dimension, and low eccentricity.

イロン 不得 とうほう イヨン

Progress in Benign Overfitting

• Simplicial interpolation (\approx nearest neighbor)

- (Belkin, Hsu, Mitra, 2018)
- Nadaraya-Watson estimator with singular kernels

(Belkin, Hsu, Mitra, 2018; Belkin, Rakhlin, Tsybakov, 2018)

 Random matrix theory asymptotics (*d* ≍ *n*) for linear regression, random nonlinear features

(Hastie, Montanari, Rosset, Tibshirani, 2019; Mei, Montanari, 2019; Belkin, Hsu and Xu, 2019)

Certain reproducing kernel Hilbert spaces

(Liang and Rakhlin, 2018; Rakhlin and Zhai, 2018; Liang, Rakhlin, Zhai, 2019)

 Minimum norm linear regression: tight upper and lower bounds for finite sample, arbitrary dimension (B., Long, Lugosi, Tsigler, 2019)

Benign Overfitting in Deep Networks?

Deep learning mysteries

- Benign overfitting
- Implicit regularization
- Omputational complexity of estimation





Implicit Regularization

 Stochastic gradient descent finds deep networks satisfying the (overfitting) constraint, and these deep networks predict accurately.



Implicit Regularization

- Stochastic gradient descent finds deep networks satisfying the (overfitting) constraint, and these deep networks predict accurately.
- What is the regularizer Ω?



Implicit Regularization

- Stochastic gradient descent finds deep networks satisfying the (overfitting) constraint, and these deep networks predict accurately.
- What is the regularizer Ω?
- The boundaries between the key issues of *optimization, estimation, and approximation* are blurred.

Progress in Implicit Regularization

- Linear. $f: x \mapsto \langle \theta, x \rangle$: $\Omega(f) = \|\theta \theta_0\|$.
- Polynomial. θ_i replaced by θ_i^{α} : $\Omega(f)$ like a Huber norm.

(Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro, 2017)

Logistic regression

(Soudry, Hoffer, Srebro, 2017)

Linear convolutional: Ω(f) penalizes norm of Fourier transform.

(Gunasekar, Lee, Soudry, Srebro, 2018)

Progress in Implicit Regularization

- Linear. $f: x \mapsto \langle \theta, x \rangle$: $\Omega(f) = \|\theta \theta_0\|$.
- Polynomial. θ_i replaced by θ_i^{α} : $\Omega(f)$ like a Huber norm.

(Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro, 2017)

Logistic regression

(Soudry, Hoffer, Srebro, 2017)

• Linear convolutional: $\Omega(f)$ penalizes norm of Fourier transform.

(Gunasekar, Lee, Soudry, Srebro, 2018)

Implicit Regularization in Deep Networks?

Deep learning mysteries

- Benign overfitting
- Implicit regularization

2 Computational complexity of estimation

Computational complexity of estimation

Mean estimation

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample

Computational complexity of estimation

Mean estimation

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

• What is the computational complexity?

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{rac{d}{n}} + \sqrt{rac{\log(1/\delta)}{n}}
ight)$$
 (optimal).

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{rac{d}{n}} + \sqrt{rac{\log(1/\delta)}{n}}
ight)$$
 (optimal).

• With weaker assumptions? e.g., just a second moment?

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right)$$
 (optimal).

• With weaker assumptions? e.g., just a second moment? The empirical mean is sensitive to large outliers.
Mean estimation

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right)$$
 (optimal).

- With weaker assumptions? e.g., just a second moment? The empirical mean is sensitive to large outliers.
- In one dimension, we can compute a median of means.

Mean estimation

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right)$$
 (optimal).

- With weaker assumptions? e.g., just a second moment? The empirical mean is sensitive to large outliers.
- In one dimension, we can compute a median of means.

$$\epsilon = O\left(\sqrt{rac{\log(1/\delta)}{n}}
ight)$$
 (optimal).

Mean estimation

Consider estimating the mean μ of a distribution in \mathbb{R}^d from an *n*-sample: Pr $(\|\hat{\mu} - \mu\| \ge \epsilon) \le \delta$.

- What is the computational complexity?
- With subgaussian data, the empirical mean suffices.

$$\epsilon = O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right)$$
 (optimal).

- With weaker assumptions? e.g., just a second moment? The empirical mean is sensitive to large outliers.
- In one dimension, we can compute a median of means.

$$x = O\left(\sqrt{rac{\log(1/\delta)}{n}}
ight)$$
 (optimal).

● In ℝ^d?

Mean estimation with bounded second moments

Mean estimation with bounded second moments

 The optimal sample complexity for estimating the mean with heavy-tailed (i.e., just second moments) data in R^d is the same as with subgaussian data. (Lugosi and Mendelson, 2017)

Mean estimation with bounded second moments

- The optimal sample complexity for estimating the mean with heavy-tailed (i.e., just second moments) data in R^d is the same as with subgaussian data. (Lugosi and Mendelson, 2017)
- Sum-of-squares machinery provides an efficient estimator.

(Hopkins, 2018)

Mean estimation with bounded second moments

- The optimal sample complexity for estimating the mean with heavy-tailed (i.e., just second moments) data in \mathbb{R}^d is the same as with subgaussian data. (Lugosi and Mendelson, 2017)
- Sum-of-squares machinery provides an efficient estimator. But $O(n^{24})$.

(Hopkins, 2018)

Mean estimation with bounded second moments

- The optimal sample complexity for estimating the mean with heavy-tailed (i.e., just second moments) data in R^d is the same as with subgaussian data. (Lugosi and Mendelson, 2017)
- Sum-of-squares machinery provides an efficient estimator. But $O(n^{24})$.
- A simple descent-based method has run-time $O(n^4 + n^2 d)$.

Fast Mean Estimation with Sub-Gaussian Rates. (Cherapanamjeri, Flammarion, B., 2019.



(Hopkins, 2018)

Many open problems

◆□ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ → < □ →

• Hardness of estimation?

- Hardness of estimation?
- Information theoretic boundary versus computational boundary?

イロト イポト イヨト イヨト

25 / 26

- Hardness of estimation?
- Information theoretic boundary versus computational boundary?
- (Fine-grained) reductions between estimation problems?

- Hardness of estimation?
- Information theoretic boundary versus computational boundary?
- (Fine-grained) reductions between estimation problems?
- Canonical hard estimation problems? (What is the SAT of estimation?)

c.f., e.g., (Berthet and Rigollet, 2013), (Brennan and Bresler, 2019)

Machine Learning: Computation versus Statistics

Deep learning mysteries

- Benign overfitting
- Implicit regularization

Opposition Computational complexity of estimation

Machine Learning: Computation versus Statistics

Deep learning mysteries

- Benign overfitting
- Implicit regularization

Omputational complexity of estimation

