



Lecture 19 and last 1





- · project results ?
- Evaluations 🖌
- DP Mixtures Models <del><</del>
- Manifold Learning <del><</del>
- 55 -

Project < Report [Experimental 10%.]





# **Double Descent**

Beyond the Bias-Variance trade-off

## STAT 535+LPL2019 + STAT 527

Marina Meila University of Washington



- Classical regime p < N</p>
- Modern/Deep Learning/High dimensional regime N > n
  - Think N fixed, p increases, gamma=p/N
  - Training error = 0 (interpolation)
  - Test error decreases with p (or gamma)



2. 
$$\|f\| \otimes \|a\| = \hat{f} \otimes \|a\|$$
  
+ Theorem If  $f^{true} \otimes \|a\| = \hat{f} + f^{true} +$ 

**Theorem 1.** Fix any  $h^* \in \mathcal{H}_{\infty}$ . Let  $(x_1, y_1), \ldots, (x_n, y_n)$  be independent and identically distributed random variables, where  $x_i$  is drawn uniformly at random from a compact cube  $\Omega \subset \mathbb{R}^d$ , and  $y_i = h^*(x_i)$  for all *i*. There exists absolute constants A, B > 0 such that, for any interpolating  $h \in \mathcal{H}_{\infty}$  (i.e.,  $h(x_i) = y_i$  for all *i*), so that with high probability



# + Main intuition [Belkin et al.]



- The target function h\* is (mostly) smooth
   i.e. ||h\*||<sub>RKHS</sub> is small
- p > N, no noise, hence h<sub>p</sub> interpolates data
- Train to minimize | |h<sub>p</sub>| | subject to 0 training error
- Then ||h<sub>p</sub>|| will decrease with p!

### + Theorem



**Theorem 1.** Fix any  $h^* \in \mathcal{H}_{\infty}$ . Let  $(x_1, y_1), \ldots, (x_n, y_n)$  be independent and identically distributed random variables, where  $x_i$  is drawn uniformly at random from a compact cube  $\Omega \subset \mathbb{R}^d$ , and  $y_i = h^*(x_i)$  for all *i*. There exists absolute constants A, B > 0 such that, for any interpolating  $h \in \mathcal{H}_{\infty}$  (i.e.,  $h(x_i) = y_i$  for all *i*), so that with high probability

$$\sup_{x \in \Omega} |h(x) - h^*(x)| < A e^{-B(n/\log n)^{1/d}} \left( \|h^*\|_{\mathcal{H}_{\infty}} + \|h\|_{\mathcal{H}_{\infty}} \right).$$

# + Linear regression<sub>∞</sub> [Hastie, Montanari, Rosset, Tibshirani 2019] ∞

- Linear, nonlinear features behave the same way
- Model correct, misspecified
- Noise level sigma affects asymptotic error
- and optimal N/n



Double descent is not regularization

Figure 1: Asymptotic risk curves for the linear feature model, as a function of the limiting aspect ratio  $\gamma$ . The risks for min-norm least squares, when SNR = 1 and SNR = 5, are plotted in black and red, respectively. These two match for  $\gamma < 1$  but differ for  $\gamma > 1$ . The null risks for SNR = 1 and SNR = 5 are marked by the dotted black and red lines, respectively. The risk for the case of a misspecified model (with significant approximation bias, a = 1.5 in (13)), when SNR = 5, is plotted in green. Optimally-tuned (equivalently, CV-tuned) ridge regression, in the same misspecified setup, has risk plotted in blue. The points denote finite-sample risks, with n = 200,  $p = [\gamma n]$ , across various values of  $\gamma$ , computed from features X having i.i.d. N(0, 1) entries. Meanwhile, the "x" points mark finite-sample risks for a nonlinear feature model, with n = 200,  $p = [\gamma n]$ , d = 100, and  $X = \varphi(ZW^T)$ , where Z has i.i.d. N(0, 1) entries, W has i.i.d. N(0, 1/d) entries, and  $\varphi(t) = a(|t| - b)$  is a "purely nonlinear" activation function, for constants a, b. The theory predicts that this nonlinear risk should converge to the linear risk with p features (regardless of d). The empirical agreement between these two—and the agreement in finite-sample and asymptotic risks—is striking.



- More refined analysis includes noise, non-linearity, data dimension n, ridge regularization lambda [Mei, Montanari 2019]
- When is global minimum in overparametrized regime?
- Enough data N/n > 1
- lambda  $\rightarrow$  0 ( or min-norm LS)
- p >> N
- SNR || beta ||/noise > 1
- Bias, Variance strictly decreasing with p/N to > 0 limit

#### CSE 547/STAT 548

#### Non-linear dimension reduction: an introduction

Marina Meilă

Department of Statistics University of Washington

January 2022

Marina Meilă (Statistics)

Manifold Learning Intro

January 2022 1 / 71

< □ > < □ > < □ > < □ > < □ >

#### Outline

Manifold Learning / Non-lin dimension reduction

イロト イヨト イヨト

- What is manifold learning good for?
- 2 Manifolds, Coordinate Charts and Smooth Embeddings
- Son-linear dimension reduction algorithms
  - Local PCA
  - PCA, Kernel PCA, MDS recap
  - Principal Curves and Surfaces (PCS)
  - Embedding algorithms
- Metric preserving manifold learning Riemannian manifolds basics
  - Metric Manifold Learning Intuition
  - Mathematical defihitons
  - Estimating the Riemannian metric
- Choice of neighborhood radius
  - What graph? Radius-neighbors vs. k nearest-neighbors
  - What neighborhood radius/kernel bandwidth?

#### Who needs manifold learning?



イロト イ団ト イヨト イヨト

#### Spectra of galaxies measured by the Sloan Digital Sky Survey (SDSS)





• Preprocessed by Jacob VanderPlas and Grace Telford • n = 675,000 spectra  $\times D = 3750$  dimensions



embedding by James McQueen

< □ > < □ > < □ > < □ > < □ >

#### Molecular configurations





#### When to do (non-linear) dimension reduction

- n = 698 gray images of faces in D = 64 × 64 dimensions
- head moves up/down and right/left
- With only two degrees of freedom, the faces define a 2D manifold in the space of all 64 × 64 gray images



 $x \in \mathbb{R}^{b} \xrightarrow{F} y \in \mathbb{R}^{d}$  $e \mathbb{R}^{m}$ D >>m≥d F embedding = map to lowerdim so that 7 smooth and F smooth d=1 D Not continuous d=2

distortion

#### G for Sculpture Faces

- n = 698 gray images of faces in  $D = 64 \times 64$  dimensions
- head moves up/down and right/left

stretch

Compress

#### Corrections for 3 embeddings of the same data



Isomap





Laplacian Eigenmaps 🛌 🖘 🚊 🛛 외 ( 이

#### Isomap vs. Diffusion Maps



#### Isomap

- Preserves geodesic distances
  - $\bullet\,$  but only when  ${\cal M}$  is flat and "data" convex
- Computes all-pairs shortest paths  $\mathcal{O}(n^3)$
- Stores/processes dense matrix





• t-SNE, UMAP visualization algorithms

January 2022 35 / 71

#### Metric Manifold Learning

#### Wanted

- eliminate distortions for any "well-behaved"  $\mathcal{M}$
- and any any "well-behaved" embedding  $\phi(\mathcal{M})$
- in a tractable and statistically grounded way

#### Idea

```
Given data \mathcal{D} \subset \mathcal{M}, some embedding \phi(\mathcal{D}) that preserves topology
(true in many cases)
```

- Estimate distortion of  $\phi$  and correct it! (see pred/next slide)
- The correction is called the pushforward Riemannian Metric g

< □ > < □ > < □ > < □ > < □ >



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

э.

Lecture 14: Dirichlet Process Mixtures in a nutshell

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

March 9, 2018

#### • The Chinese Restaurant Process $\rightarrow$ generates cluster labels $\varepsilon$ cluster parameters $\mu_{k_{1}} \geq \varepsilon$

#### The Chinese Restaurant Process

- Given parameters  $\alpha > 0$ ,  $G_0$ , with  $G_0$  a continuous measure on measurable space  $(\Theta, \mathcal{B}).$ PkH ~ Go L
- Assume we already have samples  $\theta_{1:n} \in \Theta$ .
- The probability of  $\theta_{n+1}$  is then

$$\theta_{n+1} | \theta_{1:n} \sim \sum_{k=1}^{K} \frac{n}{n+\alpha} \delta_{\theta_k} + \frac{\alpha}{n+\alpha} G_0.$$
 (1)

In the above, K represents the number of distinct values among the n samples Note: all distinct O's are sampled from Go!  $\theta_{1:n}$ .

- This defines a Chinese Restaurant Process (CRP). It is easy to see that the process is exchangeable. I likelihood invariant to ordering 1:w
- One can also prove that for  $n \to \infty$ ,  $\theta_{1:n} \to G$  where  $G \sim DP(\alpha, G_0)$ .

$$\Theta_{k} = (\mu_{k}, \Sigma_{k})$$
 cluster parame

 $\Pr[\text{new table}] = \frac{\alpha}{n+\alpha} \approx \frac{1}{n}$  $\in [\# \text{tables}] \leq \alpha \left( 1 + \frac{1}{2} + \dots + \frac{1}{n} \right) \neq \infty \text{ fmn}$ 



#### **Dirichlet Process**

- A Dirichlet Process (DP) is distribution over measures.
- Let  $(\Theta, \mathcal{B})$ ,  $\alpha, G_0$  be as above.
- We say that the random function G is drawn from  $DP(\alpha, G_0)$  iff

for any partition  $B_{1:K} \subset \mathcal{B}$  of  $\Theta$ ,  $G(B_{1:K}) \sim Dirichlet(\alpha G_0(B_{1:K}))$ . (2)

#### Dirichlet Process Mixture

- Given:  $DP(\alpha, G_0)$ , family of distributions  $\{f_{\theta}\}$  on  $\mathcal{X}$ .
- ▶ For *i* = 1, 2, . . . *n*

$$\theta_i \sim CRP(\alpha, G_0, \theta_{1:i-1})$$
(3)
  
 $x_i \sim f_{\theta_i}$ 
(4)

<□ > < @ > < E > < E > E のQ @

#### Estimation of DP Mixture by Gibbs sampling

Input  $\alpha$ ,  $G_0$ ,  $\{f\}$ ,  $\mathcal{D} = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$ State cluster assignments  $c_i$ , i = 1 : n, parameters  $\theta_k$  for all distinct kIterate 1. for i = 1 : n(reassign data to clusters) 1.1 if  $n_{c_i} = 1$  delete this cluster and its  $\theta_{c_i}$ 1.2 resample  $c_i$  by  $c_i = \begin{cases} existing k & w.p \propto \frac{n_k}{n-1+\alpha} f(x_i, \theta_k) \\ new cluster & w.p \frac{\alpha}{n-1+\alpha} \int f(x_i, \theta) G_0(\theta) d\theta \end{cases}$ (5) 1.3 if  $c_i$  is new label, sample a new  $\theta_{c_i}$  from  $f_{\theta} G_0(\theta)$ 2. (resample cluster parameters) for  $k \in \{c_{1:n}\}$ 2.1 sample  $\theta_k$  from posterior  $f_{\theta_k} \propto G_0(\theta) \prod_{i \in C_k f(x_i, \theta)}$ this can be computed in closed form if  $G_0$  is conjugate prior

#### Output a state with high posterior