

STAT 391

5/23/23

Lecture

17

it's not too late!

Participation

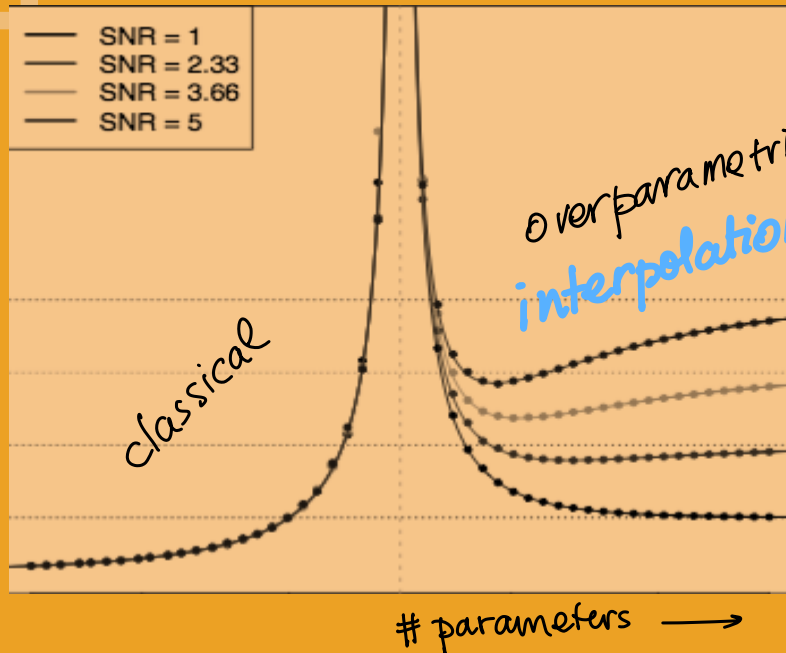
Regression - double descent

Logistic regression

HW 6 + b. posted
due Friday 6/1
optional

All HW + b. graded
on Friday 6/1

Exam week
Exam reviews / SP
+ b. scheduled / MHP



Double Descent

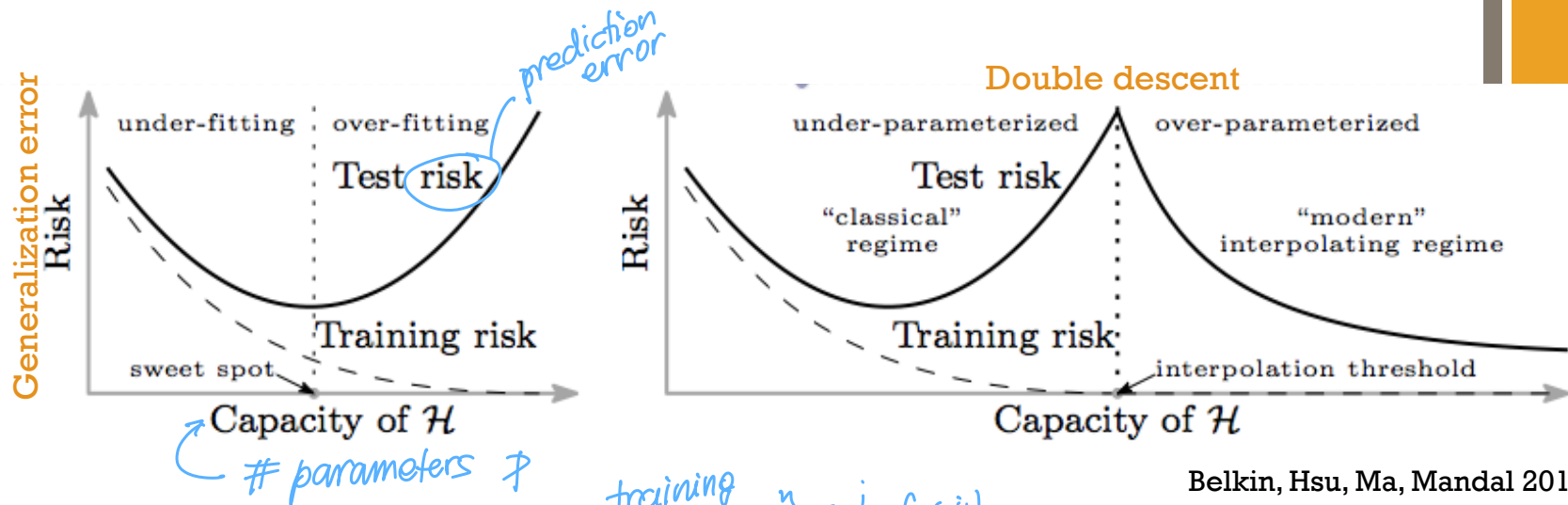
Beyond the Bias-
Variance trade-off

STAT 535+LPL2019

Marina Meila
University of Washington



+ What is observed



Belkin, Hsu, Ma, Mandal 2018

- Classical regime $p < N$

- 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 γ)

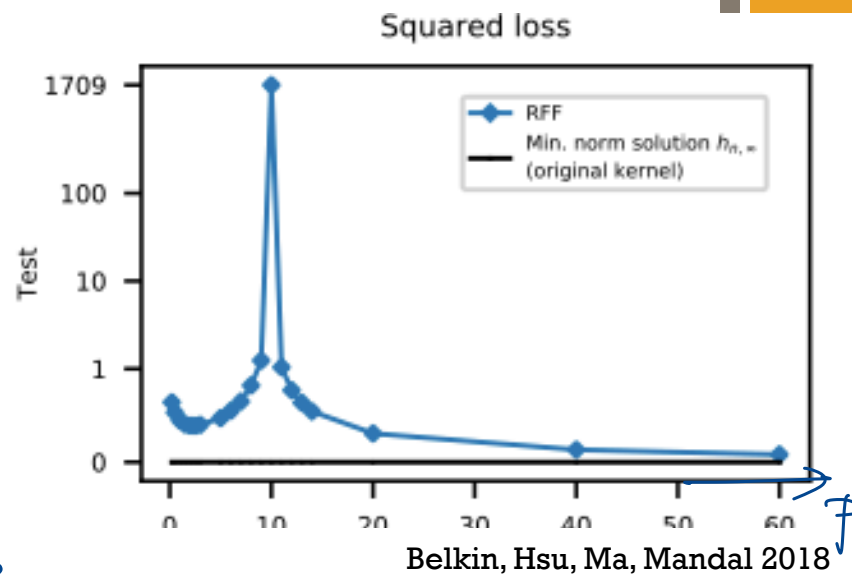
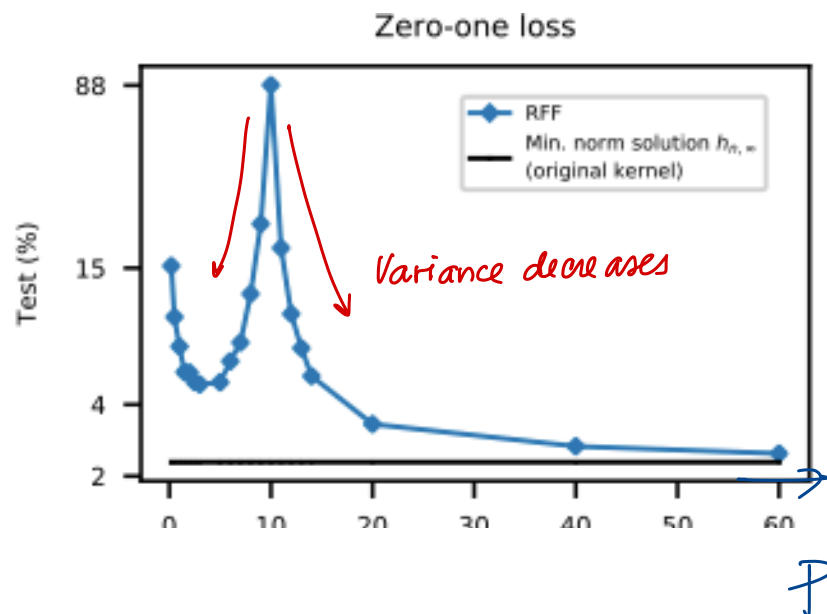
$$\text{training err} = \frac{1}{n} \sum_{i=1}^n (y^i - f_{\theta}(x^i))^2$$

$\uparrow \mathcal{F}_p$ $\# \theta = p$

$$\text{test err} = \frac{1}{n_v} \sum_{j=1}^{n_v} (y^j - f_{\theta}(x^j))^2$$

$\uparrow x^j, y^j \in \mathcal{D}^{\text{validation}}$ “Squared Loss”

+ What is observed



■ Double descent curves for the generalization error

- Random Fourier Features (RFF) *SRM*
- ReLU 2 layer networks (with random first layer weights)
- Random Forests, l2-Adaboost
- Linear regression

*observed
for many classes of
predictors*

■ With and without noise

Ex: 1) $x = \text{image}$ $y = 1$ iff contains car $\{(x^{(i)}, y^{(i)})\}$ training set
 $f_{\theta} \in \mathcal{F}_p$ neural net
 $p = \# \text{weights}$ or $\# \text{layers}$ in neural net

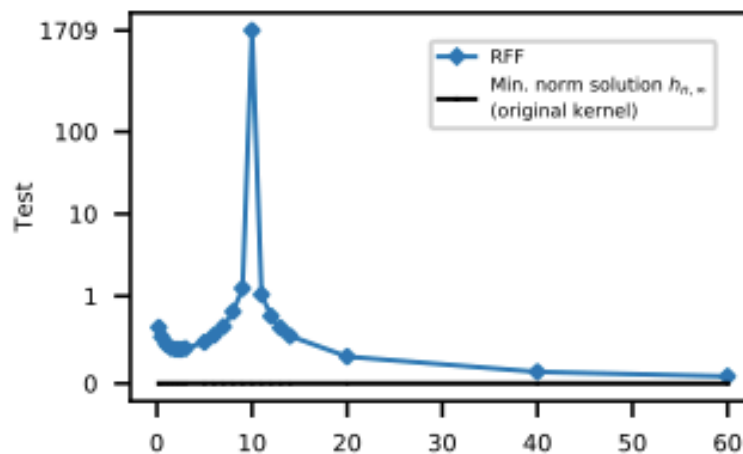
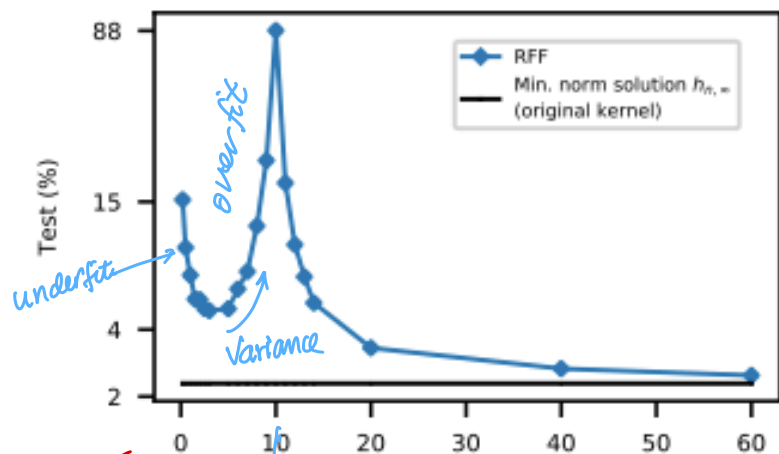
Two cases of increasing p

+ "NN"

RFF = $\mathcal{F} = \{f_{\theta}\}$ 2) $f_{\theta} = \Phi^T x$
 $x, \theta \in \mathbb{R}^p$ $p \uparrow \Leftrightarrow \text{resolution of } x \text{ increases}$

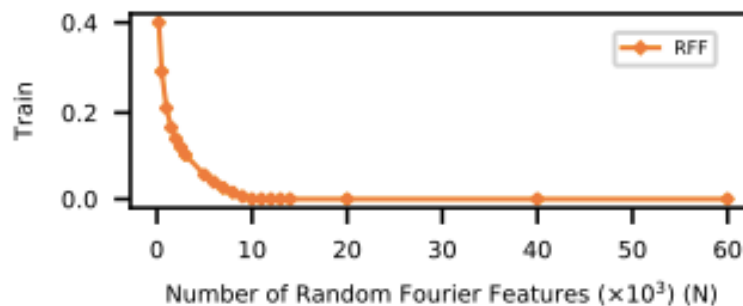
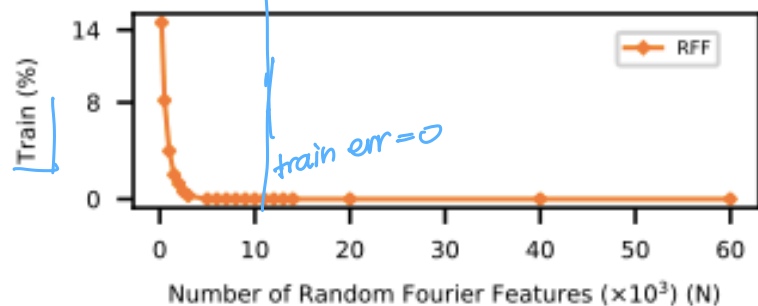
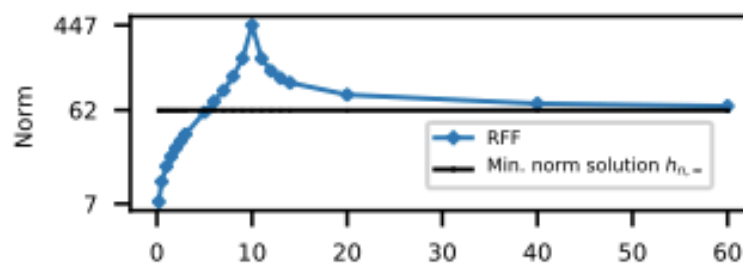
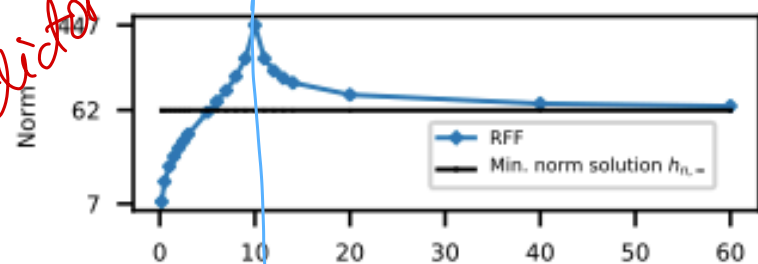
Zero-one loss

Squared loss



$\|v\|^2 = \sum_j v_j^2$
 $\| \text{neural network} \|^2 = \sum_{j=1}^p w_j^2$
 $w_{1:p} = \text{weights of neural network}$

norm of predictor



1. Double descent $\exists \Rightarrow$ Interpolation $\Leftrightarrow f_{\hat{\theta}}(x^i) = y^i$ for $i=1:n$

2. in Interpolation regime

- $\|\hat{\theta}\| \rightarrow \text{to min} > 0$
hypothesis $\|\hat{\theta}\| \text{ small} \Rightarrow \text{Var } \hat{\theta} \text{ small}$
- \exists infinite number θ so that $f_{\theta}(x^{1:n}) = y^{1:n}$

Intuition: $\|\theta\| \text{ small} \Leftrightarrow f_{\theta} \text{ smooth}$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$$

$$\sum_{j=1}^p |w_j| = 1$$

$$w^0 = \left[\frac{1}{p} \quad \dots \quad \frac{1}{p} \right]^T \text{ constant vector} = \text{smoothest!!}$$

$$\|w^0\|^2 = p \cdot \frac{1}{p^2} = \frac{1}{p} \leq \|w\|^2 \quad w \in \mathbb{R}^p$$

$$\sum |w_j| = 1$$

$$\min_{\sum |w_j|=1} \|w\|^2 \Rightarrow \frac{1}{p}$$

$\frac{1}{p} \downarrow$ for $p \uparrow \Rightarrow \text{smoother for } p \uparrow$

3. $\|\theta\| \text{ small} \Rightarrow f_{\theta} \text{ smooth} \Rightarrow \text{low variance!!}$

$$p_1 < p_2 < p_3$$

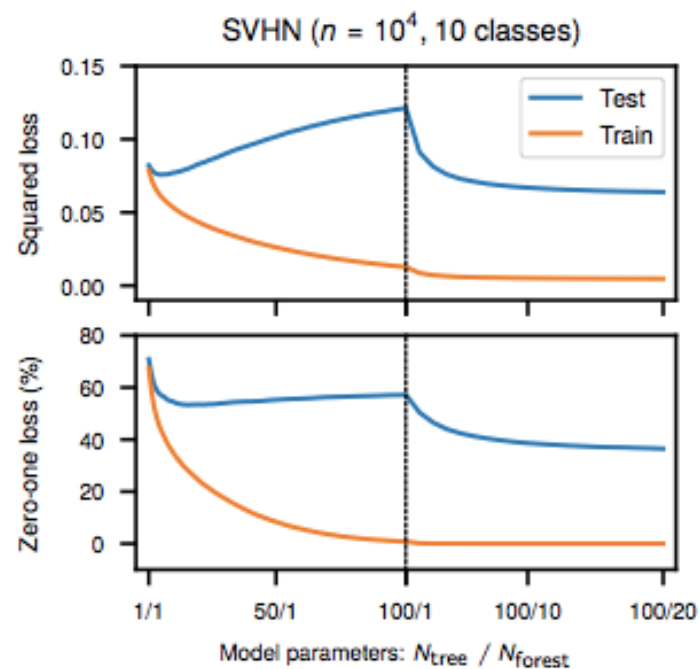
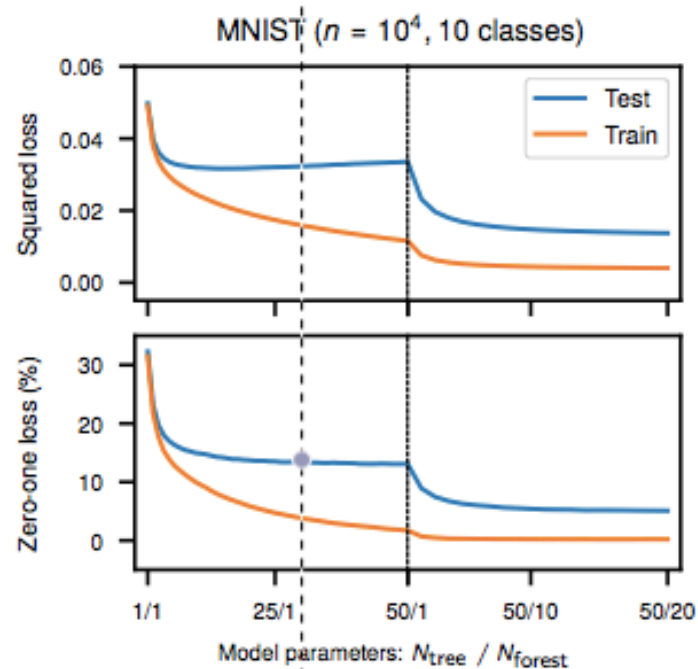
$$\mathcal{F}_{p_1} \subset \mathcal{F}_{p_2} \subset \mathcal{F}_{p_3} \dots$$

more models to choose from!!

4. If training algorithm can choose smoothest $f_{\theta} \in \mathcal{F}_p \Rightarrow f_{\theta_{p_2}}$ smoother

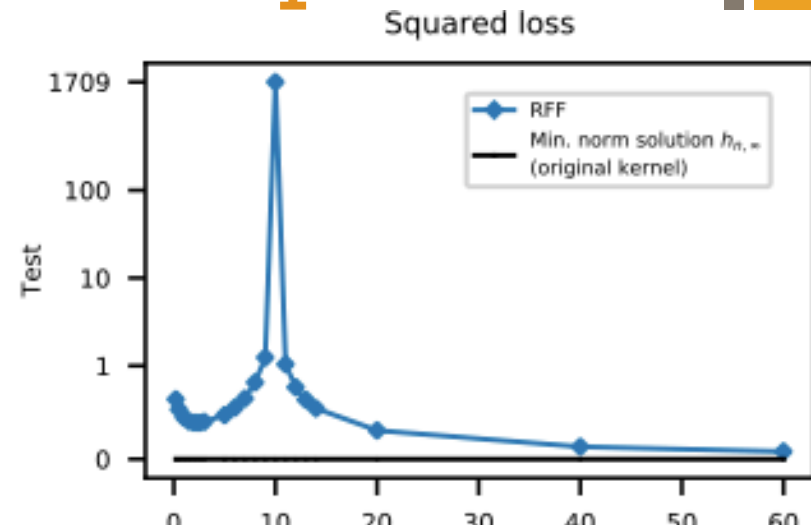
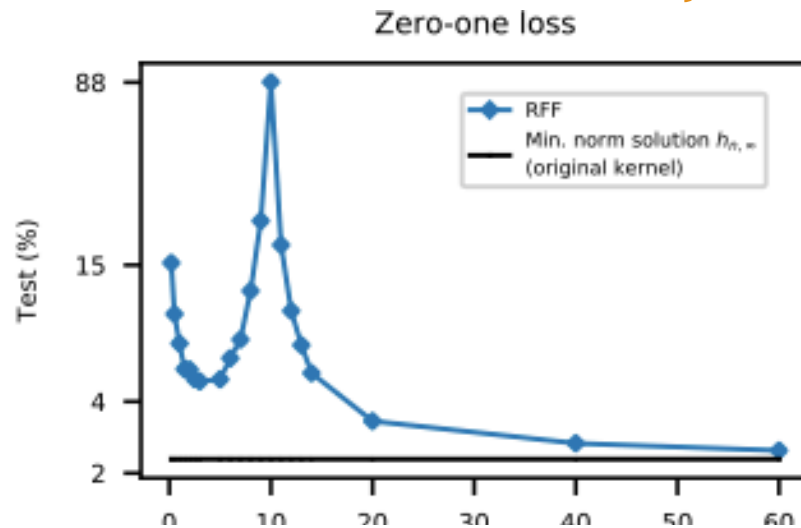
than $f_{\theta_{p_1}}$

+ Boosted decision trees





Double descent, the case $p > N$

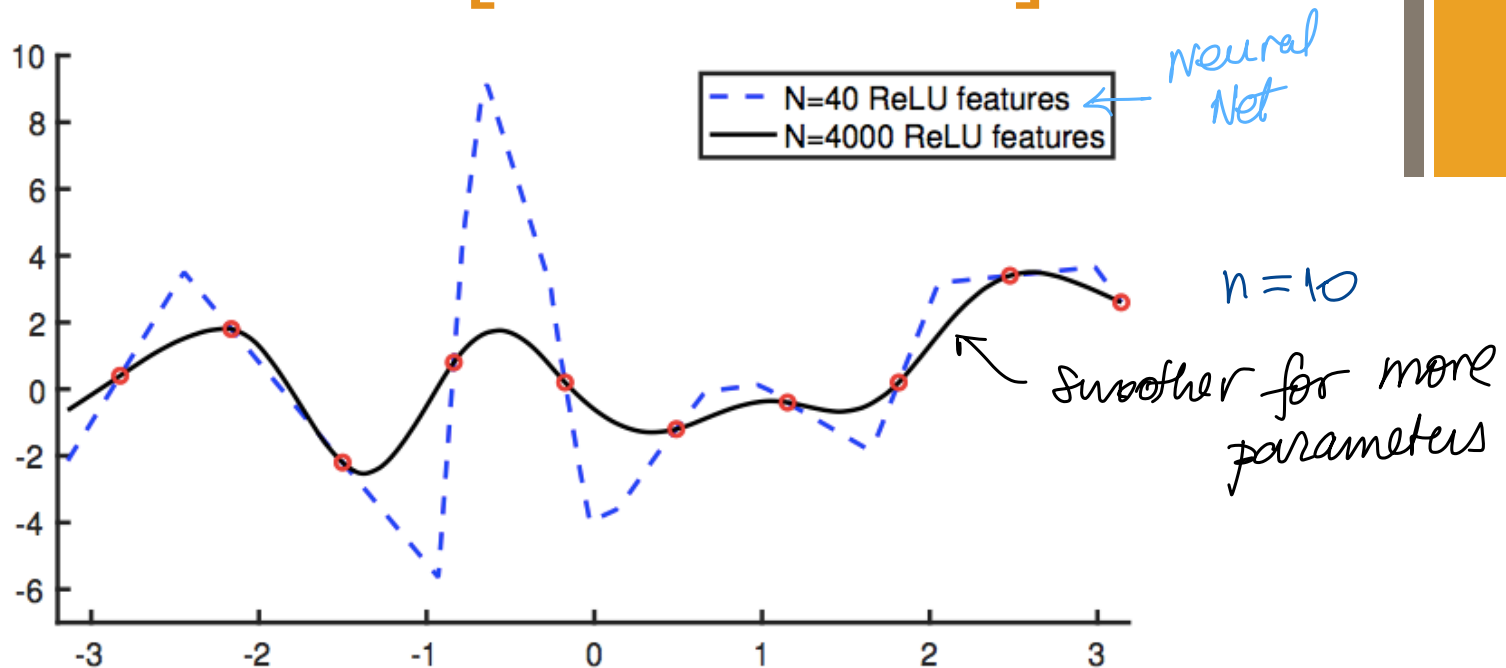


$p \rightarrow$
parameters

Belkin, Hsu, Ma, Mandal 2018

- Model $y = \langle \phi(x), \beta \rangle$
- Large N (cover a compact data domain)
- Features random
- Min-norm solution β^*

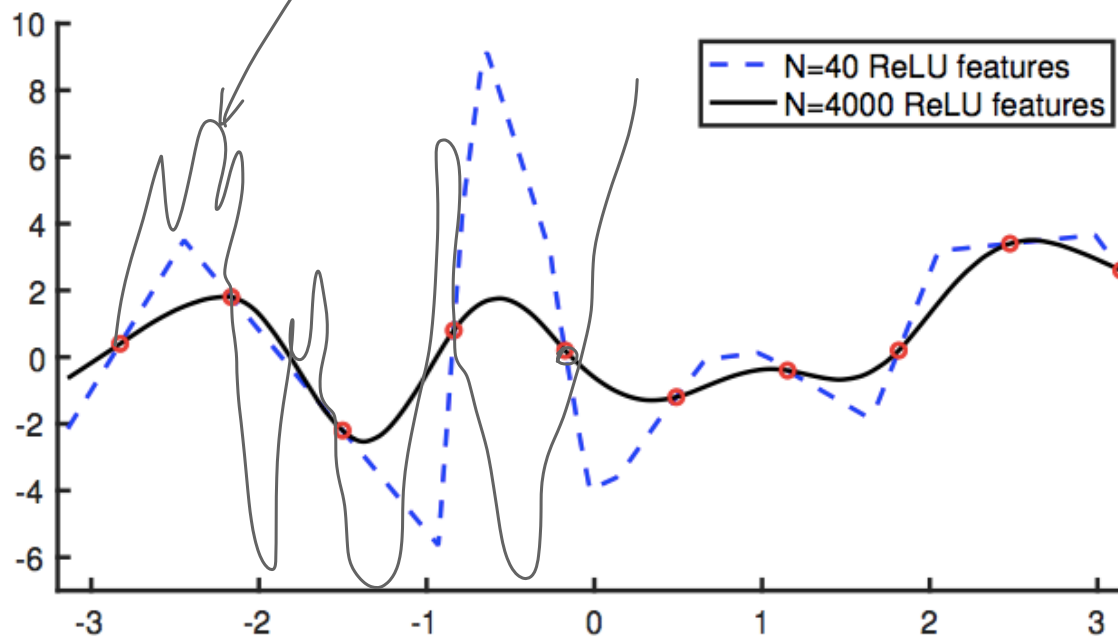
+ Main intuition [Belkin et al.]



- The target function h^* is (mostly) smooth
 - i.e. $||h^*||_{RKHS}$ is small
- $p > N$, no noise, hence h_p interpolates data
- Train to minimize $||h_p||$ subject to 0 training error
- Then $||h_p||$ will decrease with p !

\exists many f 's that interpolate!

+ Main intuition [Belkin et al.]



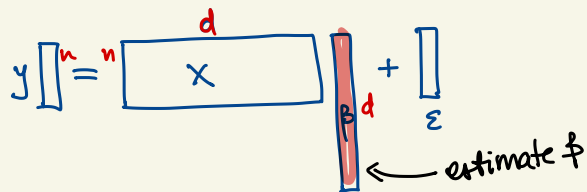
Surprising fact:
gradient descent
finds smooth (est)
 f_θ !!

- The target function h^* is (mostly) smooth
 - i.e. $||h^*||_{\text{RKHS}}$ is small
- $p > N$, no noise, hence h_p interpolates data
- Train to minimize $||h_p||$ subject to 0 training error
- Then $||h_p||$ will decrease with p !

Linear regression

• True model $y = (\beta^*)^T x + \varepsilon \iff y^i = (\beta^*)^T x^i + \varepsilon^i \quad i=1:n$
 $\mathcal{D} = \{(x^i, y^i)\}_{i=1:n}$

• $n < d$ $x, \beta, \beta^* \in \mathbb{R}^d$

$y \begin{bmatrix} n \\ \end{bmatrix} = \begin{bmatrix} d \\ x \end{bmatrix} \beta_d + \begin{bmatrix} \varepsilon \end{bmatrix}$


• $y = X\beta$ has ∞ solutions

• $X^T X \in \mathbb{R}^{d \times d}$ with rank $n < d$

can't solve this way!

if $n > d \Rightarrow \hat{\beta}^{ML} = (X^T X)^{-1} X^T y$

• So gradient descent instead

$\min_{\beta} -\ell(\beta) = \min_{\beta} \sum_{i=1}^n (\bar{y}^i - \bar{\beta}^T x^i)^2$

Exercise

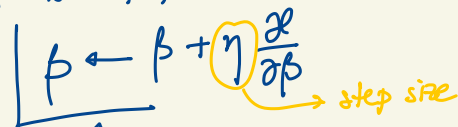
$-\frac{\partial \ell}{\partial \beta} = 2(X^T X \beta - X^T y)$

$\frac{\partial \ell}{\partial \beta} = \left[\frac{\partial \ell}{\partial \beta_j} \right]_{j=1:d}$

Gradient descent Alg

Init $\beta \leftarrow 0$

For $t = 1, 2, \dots$ until convergence

$\beta \leftarrow \beta + \eta \frac{\partial \ell}{\partial \beta}$


Out $\hat{\beta} = \beta$ at convergence

Properties of $\hat{\beta}$

Ex (Linear algebra, not easy)
Prove Prop 1, 2, 3

XX^T , $X^T X$ same
singular value

Prop 1

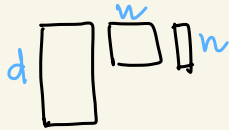
$$y = X \hat{\beta}$$

interpolation

Prop 2

$$\hat{\beta} = X^T K^{-1} y$$

$K = XX^T$ Gram matrix
 $n \times n$, full rank



Prop 3

$$\|\hat{\beta}\| = \arg \min_{y=X\beta} \|\beta\|$$

min norm interpolating
solution

logistic Regression

used for: Classification

$$x \in \mathbb{R}^d$$

$$y \in \{0, 1\} \text{ "Label"}$$

$$\mathcal{D} = \{(x^i, y^i)_{i=1:n}\}$$

Model

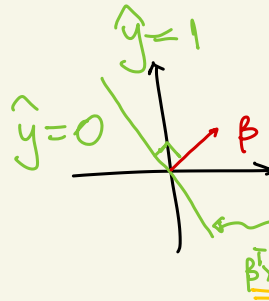
$$f(x) = \ln \underbrace{\frac{P_{y|x}(y=1|x)}{P_{y|x}(y=0|x)}}_{\substack{\text{meaning} \\ \text{odds} \in (0, \infty) \\ \text{log odds} \in (-\infty, \infty)}} = \underline{\beta^T x} \quad \beta \in \mathbb{R}^d$$

parametrization

Classification with $f = \beta^T x$ (β known)

$$\hat{y} = \frac{\text{sign } f(x) + 1}{2} \iff f(x) > 0 \iff \underbrace{P(y=1|x)}_{y|x} > \underbrace{P(y=0|x)}_{y|x}$$

Probabilistic classifier



linear classifier \iff
decision boundary
linear

- $P_{y|x}[y=1|x] = p$

$$\boxed{f = \ln \frac{p}{1-p} = \ln \frac{P_{y|x}(y=1|x)}{P_{y|x}(y=0|x)}} \Rightarrow e^f = \frac{p}{1-p} \Rightarrow p = \frac{e^f}{1+e^f} \quad e^{(0/1)}$$

$f = \beta^T x$

$$\boxed{p = \frac{e^f}{1+e^f} \quad e^{(0/1)}$$

$$1-p = \frac{1}{1+e^f}$$

$$p + (1-p) = 1$$

- Estimating β by Max likelihood

likelihood

$$L(\beta) = P[y^{1:n} | x^{1:n}, \beta] = \prod_{i=1}^n P_{y|x}[y^i | x^i, \beta] = \prod_{i=1}^n \frac{e^{y^i \beta^T x^i}}{1 + e^{\beta^T x^i}}$$

\downarrow \downarrow
 p $1-p$
 for $y^i=1$ for $y^i=0$

confidence
of
classification

$$P_{y|x}[y|x] = \frac{e^{y f}}{1 + e^f}$$

log likelihood

$$l(\beta) = \sum_{i=1}^n [y^i \beta^T x^i - \ln(1 + e^{\beta^T x^i})]$$

maximize
by Gradient
Ascent

$$\ell(\beta) = \sum_{i=1}^n \left[y_i \overset{f}{\beta^T x_i} - \ln(1 + e^{\beta^T x_i}) \right]$$

Gradient $\frac{\partial \ell}{\partial \beta} \in \mathbb{R}^d$

$$\frac{\partial \ell}{\partial \beta_j} = \frac{\partial \ell}{\partial f} \cdot \frac{\partial f}{\partial \beta_j}$$

$$\frac{\partial \ell}{\partial \beta} = \left[\frac{\partial \ell}{\partial \beta_j} \right]_{j=1:d} = \frac{\partial \ell}{\partial f} \left[\frac{\partial f}{\partial \beta_j} \right]_{j=1:d}$$

$$\frac{\partial f}{\partial \beta} (\beta^T x) = x \leftarrow \text{Exercise prove this}$$