Machine Learning Theory 2025 Lecture 7

Tim van Erven

- Complexity of classification vs regression
- Neural networks
- Bias-variance trade-off and double descent
- Towards an explanation

Binary Classification

Sample complexity of agnostic PAC-learnability determined by VC-dimension:

$$m_{\mathcal{H}}(\epsilon, \delta) pprox rac{\mathsf{VCdim}(\mathcal{H}) + \ln(1/\delta)}{\epsilon^2}$$

For some (not all!) hypothesis classes, VCdim(H) = nr. of parameters:

- Linear predictors: $\mathcal{H} = \{h_w(X) = \operatorname{sign}(\langle w, X \rangle) : w \in \mathbb{R}^d\}$
- Axis-aligned rectangles
- ► ...

Regression

$$\mathcal{H}_1^B = \{h_w(X) = \langle w, X \rangle : w \in \mathbb{R}^d, \|w\|_1 \leq B\}.$$

Theorem (Lasso Estimator)

Consider linear regression with $\ell(h, X, Y) = \frac{1}{2}(Y - \langle w, X \rangle)^2$ for $X \in [-1, +1]^d$, $Y \in [-1, +1]$. Then \mathcal{H}_1^B is agnostically PAC-learnable by ERM with sample complexity

$$m(\epsilon, \delta) \leq c_B \frac{\ln(2d) + \ln(2/\delta)}{\epsilon^2}$$

for some constant $c_B > 0$ that depends only on B.

General pattern for regression tasks:

- Complexity of hypothesis class depends on bound B on norm ||w|| of parameters
- (and sometimes weakly on number of parameters d)

Difference between Linear Regression and Linear Classification

Linear Classification:

- Not Lipschitz in w: tiny change in w can flip prediction $h_w(X)$
- Measure of complexity: number of parameters d

Linear Regression:

- Lipschitz in w: tiny change in w implies tiny change in $h_w(X)$
- Main measure of complexity: norm constraint B

Deep Learning / Neural Networks



Geogle			8
Translate		substantiant tent of the	з
English Novich Event Intercept -	5	France English (back - Normal	
Deep neural networks are making it possible to finally KABI the big promises that antificial intelligence has been making since the 160ks.	×	Diepe neurale networken maken het mogelijk om eindelijk de grote belaften na te komen die kunstmatige intelligentie sinds de janen zweig heeft gemaakt.	
6 = *		±0.e< ≠uero	

Machine translation







Self-driving cars



Class of **non-convex** functions parametrized by matrices $w = (A_1, \ldots, A_m)$:

Fully connected network: $\mathcal{H} = \{h_w(X) = A_m \sigma A_{m-1} \cdots \sigma A_1 X : w \in \mathcal{W}\},\$

with activation function $\sigma(z)$ applied component-wise to vectors. E.g.

• Rectified linear unit (ReLU): $\sigma(z) = \max\{0, z\}$

• Sigmoid:
$$\sigma(z) = 1/(1 + e^{-z})$$



Class of **non-convex** functions parametrized by matrices $w = (A_1, \ldots, A_m) \in \mathbb{R}^d$:

Fully connected network: $\mathcal{H} = \{h_w(X) = A_m \sigma A_{m-1} \cdots \sigma A_1 X : w \in \mathbb{R}^d\},\$

with activation function $\sigma(z)$ applied component-wise to vectors. E.g.

Rectified linear unit (ReLU): σ(z) = max{0, z}

• Sigmoid:
$$\sigma(z) = 1/(1 + e^{-z})$$



VC-dimension dependence on nr. of parameters *d*:

Conclusion: need sample size $m \gg nr$. of parameters to learn

Class of **non-convex** functions parametrized by matrices $w = (A_1, \ldots, A_m) \in \mathbb{R}^d$:

Fully connected network: $\mathcal{H} = \{h_w(X) = A_m \sigma A_{m-1} \cdots \sigma A_1 X : w \in \mathbb{R}^d\},\$

with activation function $\sigma(z)$ applied component-wise to vectors. E.g.

- Rectified linear unit (ReLU): $\sigma(z) = \max\{0, z\}$
- Sigmoid: $\sigma(z) = 1/(1 + e^{-z})$



VC-dimension dependence on nr. of parameters d:

Conclusion: need sample size $m \gg nr$. of parameters to learn

Class of non-convex functions parametrized by matrices $w = (A_1, ..., A_{-}) \in \mathbb{D}^d$. Fully connect With activati Rectified Sigmoid: $\sigma(z) = 1/(1 + e^{-z})$ Rectified With activati

Bias-Variance Trade-off and the Double Descent Phenomenon

Classical Bias-Variance Trade-off



Double Descent Phenomenon



Varying the number of hidden units in a two-layer neural network
Classification: MNIST hand-written digits data with 10 classes

Double Descent Phenomenon



Varying the number of hidden units in a two-layer neural network
Classification: MNIST hand-written digits data with 10 classes

Towards an Explanation

- 1. Large margins turn classification into regression
- 2. Explaining double descent

Classifiers as Real-valued Functions



- Margin = Yh(X), where $Y \in \{-1, +1\}$
- Larger margin > 0: more confident correct classification

Classifiers as Real-valued Functions



NB Real-valued classifiers. E.g. $h_w(X) = \langle w, X \rangle$. Prediction is sign(h(X))

- Margin = Yh(X), where $Y \in \{-1, +1\}$
- Larger margin > 0: more confident correct classification
- Common loss functions encourage finding large margin solutions:

logistic loss: $\ln(1 + e^{-Yh(X)})$ squared loss for classification: $(Y - h(X))^2 = (1 - Yh(X))^2$



0/1-loss $\leq \gamma$ -Lipschitz loss $\leq \gamma$ -large margin loss



 $0/1\text{-}{\rm loss} \leq \gamma\text{-}{\rm Lipschitz}\ {\rm loss} \leq \gamma\text{-}{\rm large}\ {\rm margin}\ {\rm loss}$

$$\begin{split} \mathcal{L}_{\mathcal{D}}^{0/1}(h_{\mathcal{S}}) &\leq \mathcal{L}_{\mathcal{D}}^{\mathsf{Lipschitz}}(h_{\mathcal{S}}) \\ &\leq \mathcal{L}_{\mathcal{S}}^{\mathsf{Lipschitz}}(h_{\mathcal{S}}) + 2 \,\mathbb{E}[\mathcal{R}(\ell^{\mathsf{Lipschitz}},\mathcal{H},\mathcal{S})] + \sqrt{\frac{\mathsf{ln}(4/\delta)}{2m}} \quad \text{w.p.} \geq 1 - \delta \\ &\leq \mathcal{L}_{\mathcal{S}}^{\mathsf{large margin}}(h_{\mathcal{S}}) + 2 \,\mathbb{E}[\mathcal{R}(\ell^{\mathsf{Lipschitz}},\mathcal{H},\mathcal{S})] + \sqrt{\frac{\mathsf{ln}(4/\delta)}{2m}} \end{split}$$

Theorem

Let $h_S \in \mathcal{H}$ be the output of a learning algorithm. Then, with probability at least $1 - \delta$,

$$L_{\mathcal{D}}^{0/1}(h_S) \leq L_S^{\gamma\text{-large margin}}(h_S) + 2 \mathbb{E}[\mathcal{R}(\ell^{\gamma\text{-Lipschitz}}, \mathcal{H}, S)] + \sqrt{rac{\ln(4/\delta)}{2m}}$$

Theorem

Let $h_S \in \mathcal{H}$ be the output of a learning algorithm. Then, with probability at least $1 - \delta$,

$$L_{\mathcal{D}}^{0/1}(h_S) \leq L_S^{\gamma\text{-large margin}}(h_S) + 2 \mathbb{E}[\mathcal{R}(\ell^{\gamma\text{-Lipschitz}}, \mathcal{H}, S)] + \sqrt{rac{\ln(4/\delta)}{2m}}$$

1. If h_S has margin $\geq \gamma$ on (most of) S, then $L_S^{\gamma-\text{large margin}}(h_S)$ is small

Theorem

Let $h_S \in \mathcal{H}$ be the output of a learning algorithm. Then, with probability at least $1 - \delta$,

$$L_{\mathcal{D}}^{0/1}(h_{\mathcal{S}}) \leq L_{\mathcal{S}}^{\gamma\text{-large margin}}(h_{\mathcal{S}}) + 2 \operatorname{\mathbb{E}}[\mathcal{R}(\ell^{\gamma\text{-Lipschitz}},\mathcal{H},\mathcal{S})] + \sqrt{rac{\ln(4/\delta)}{2m}}$$

If h_S has margin ≥ γ on (most of) S, then L_S^{γ-large margin}(h_S) is small
Lipschitz loss is ¹/_γ-Lipschitz, so can apply contraction lemma:

$$\mathcal{R}(\ell^{\mathsf{Lipschitz}}, \mathcal{H}, \mathcal{S}) \leq \frac{1}{\gamma} \mathcal{R}(\{(h(X_1), \dots, h(X_m)) : h \in \mathcal{H}\})$$

- So small changes in h imply small changes in loss
- We have turned the classification problem into a regression task!
- Complexity of H can be controlled by some norm on h.

Towards an Explanation

- 1. Large margins turn classification into regression
- 2. Explaining double descent

A Potential Explanation



A Potential Explanation



Proposed explanation: suppose learning alg roughly behaves as

among **ERM solutions** $h_S \in \underset{h \in \mathcal{H}}{\operatorname{arg min}} L_S(h)$

choose solution with smallest norm $||h_S||_{??}$

Below int. threshold: ERM unique \rightarrow classical bias-variance trade-off Above int. threshold: larger $\mathcal{H} \rightarrow$ more ERM solutions \rightarrow smaller $||h_S||_{??}$

A Potential Explanation



Proposed explanation: suppose learning alg roughly behaves as

among ERM solutions $h_S \in \underset{h \in \mathcal{H}}{\operatorname{arg min}} L_S(h)$

choose solution with smallest norm $||h_S||_{??}$

Below int. threshold: ERM unique \rightarrow classical bias-variance trade-off Above int. threshold: larger $\mathcal{H} \rightarrow$ more ERM solutions \rightarrow smaller $||h_S||_{??}$

L_S for e.g. logistic or squared loss (encouraging large margin)

Different norm depending on manifestation of double descent

Double Descent for Neural Networks Again



[Belkin, Hsu, Ma, Mandal, 2019]

 Classification: CIFAR-10 32x32 images from 10 classes, e.g. airplanes, cats, dogs

Double Descent for Neural Networks Again



[Belkin, Hsu, Ma, Mandal, 2019]

 Classification: CIFAR-10 32x32 images from 10 classes, e.g. airplanes, cats, dogs

Which norm $||h_S||_{??}$?

Implicitly induced by optimization algorithm!

Exist proposals in the literature to characterize norm.
E.g. using neural tangent kernel [Jacot, Gabriel, Hongler, 2018]

Double Descent: Not Just for Neural Networks

[Belkin et al., 2019] reproduce double descent phenomenon on e.g. MNIST:



Random Fourier features: linear model over N randomly generated basis functions that approximate a certain (reproducing kernel) Hilbert space as $N \rightarrow \infty$

Double Descent: Not Just for Neural Networks

[Belkin et al., 2019] reproduce double descent phenomenon on e.g. MNIST:



Random forests: ensembles of decision trees

 Complexity controlled by number of leaves per tree and by number of trees

Recent Alternative Explanation [Curth, Jeffares, v.d. Schaar, 2023]: Need More Careful Parameter Counting

In all non-deep learning experiments by [Belkin et al., 2019]:

- Below interpolation threshold m: increase model complexity along dimension 1
- Above interpolation threshold *m*: increase model complexity along dimension 2

Examples:

Random forests: [Belkin et al., 2019] increase depth of single tree up to *m* (complexity dimension 1). Then average additional trees above *m* (complexity dimension 2).



Recent Alternative Explanation [Curth, Jeffares, v.d. Schaar, 2023]: Need More Careful Parameter Counting

In all non-deep learning experiments by [Belkin et al., 2019]:

- Below interpolation threshold m: increase model complexity along dimension 1
- Above interpolation threshold *m*: increase model complexity along dimension 2

Examples:

- Random forests: [Belkin et al., 2019] increase depth of single tree up to m (complexity dimension 1). Then average additional trees above m (complexity dimension 2).
- N Random Fourier features: equivalent to least squares on basis with dimension min(m, N), obtained by unsupervised dimensionality reduction. [Curth, Jeffares, v.d. Schaar, 2023]
 - Nr. of least squares parameters is min(m, N) (complexity dimension 1).
 - Quality of dimensionality reduction improves with *m* (complexity dimension 2).

Recent Alternative Explanation [Curth, Jeffares, v.d. Schaar, 2023]: Need More Careful Parameter Counting

In all non-deep learning experiments by [Belkin et al., 2019]:

- Below interpolation threshold m: increase model complexity along dimension 1
- Above interpolation threshold *m*: increase model complexity along dimension 2



Double descent happens because experiments stitch together two independent U-curves!

Conclusion

- Exciting new attempts to understand the double descent phenomenon observed in deep learning, random Fourier features, random forests, etc.
- Crucial to understand true model complexity rather than counting parameters.
- Analysis involves tools like Rademacher complexity that you have learned in this course.
- Whether proposed explanations can be fully formalized for deep learning remains to be seen...
- In any case, the role of optimization algorithms in determining effective model complexity provides a fascinating new frontier for understanding the classical bias-variance trade-off!