

# Revisiting Model Complexity in the Wake of Overparameterized Learning

Pratik Patil

Carnegie Mellon University

TOPML 2022

*Based on joint work with Ryan Tibshirani*

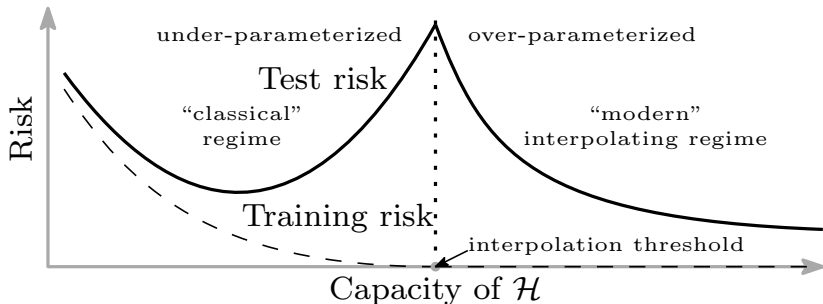
# Overparameterization in machine learning

Modern machine learning involves fitting a huge number of parameters. Such **overparameterization** seems to be useful for:

- **Representation**: allows rich, expressive models for diverse real data
- **Optimization**: simple, local optimization methods often find near-optimal solutions to empirical risk minimization on training data
- **Generalization**: despite overfitting, models generalize well in practice

This talk is broadly about the generalization aspect.

## The double/multiple descent phenomenon



"Reconciling modern machine learning practice and the bias variance tradeoff"

Belkin, Hsu, Ma, Mandal, 2018

- The phenomenon is dubbed "double descent" in the risk curve.
- This trend holds for many model classes including linear regression, kernel regression, random forest, boosting, neural networks, etc.

Is there an equivalent "reparameterization" of "overparameterization"?

## Conceptual and theoretical challenges

Current machine learning practice suggests:

- we should design models to be highly **overparametrized**
- train until nearly zero training error (i.e., **interpolate** training data)
- the trained models still can **generalize** well and have good test error

Main statistical challenges for us:

- benign overfitting conceptually breaks classical bias-variance tradeoff
- trained models fall outside the realm of uniform convergence
- one of core questions is how to correctly measure **model complexity**

This work attempts to answer the following questions:

- Is there a good measure of model complexity for predictive models?
- How to compare model complexity of different (near) interpolators?

## Takaway points

- Propose measures of model complexity that are:
  - **algorithm-specific** and applies for any prediction algorithm
  - produce a number between **0 and  $n$**  (the number of observations)
- Two variants of model complexities are:
  - **emergent** model complexity that depends on the prediction algorithm as well as underlying the regression function
  - **intrinsic** model complexity that depends on the prediction algorithm only and its adaptability to pure noise
- Results when applied to some illustrative examples:
  - **min  $\ell_2/\ell_1$ -norm interpolators**: the complexity measures maximized when  $n \approx p$  and typically decreases as  $p$  increases beyond  $n$
  - we can **reparameterize** every overparameterized model into an equivalent underparameterized model in terms of risk behavior

Based on extension of ideas from **optimism theory** and **degrees of freedom**.

# Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

Examples

Discussion

## Degrees of freedom at large

**Degrees of freedom** means different things in different fields, but they're more or less similar. There is a core concept behind it.

It essentially has to do with the dimension or **effective number of parameters** of “something”.

- In mechanics, that something = mechanical system.
- In physics and chemistry, something = physical system.
- In statistics, something = prediction procedure.

A nice commonality is that we can usually guess the degrees of freedom based on intuition, at least qualitatively.

## Degrees of freedom in statistics

Degrees of freedom in statistics is defined, intuitively, as the **effective number of parameters** used by a prediction procedure.

While this seems vague, it has a precise definition for a broad class of estimation problems. We will define this shortly.

Why is this an important concept? Why you would ever go to the trouble of describing degrees of freedom?

Essentially, it provides a way to put two different prediction procedures on **equal footing**.



## Definition under squared error loss

Consider data  $(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$ ,  $i = 1, \dots, n$  such that  $y_i = f(x_i) + \varepsilon_i$  where  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  is regression function,  $\varepsilon_i$  has mean 0 and variance  $\sigma^2$ .

Let  $\mathcal{A}$  be any fitting algorithm that maps  $\{(x_i, y_i)\}_{i=1}^n \xrightarrow{\mathcal{A}} \hat{f}$ .

The **degrees of freedom** of predictor  $\hat{f}$  is defined as

$$\text{dfF}(\hat{f}) = \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i)) / \sigma^2 = \text{tr} [\text{Cov}(y, \hat{f}(X))] / \sigma^2,$$

where  $y$ : response vector,  $X$ : feature matrix,  $\hat{f}(X)$ : predicted response

Where does squared error loss come into play?

$$\underbrace{\mathbb{E} \left[ \sum_{i=1}^n (\tilde{y}_i - \hat{f}(x_i))^2 \right]}_{\text{fixed-X prediction error} =: \text{ErrF}(\hat{f})} - \underbrace{\mathbb{E} \left[ \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{expected training error} =: \text{ErrT}(\hat{f})} = 2\sigma^2 \text{dfF}(\hat{f})$$

## Degrees of freedom in linear regression

- Suppose  $p \leq n$  and  $X$  has full (column) rank, and we take  $\hat{f}$  to be **ordinarily least squares** predictor  $\hat{f}(X) = X\hat{\beta}$ , where

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 = (X^T X)^{-1} X^T y.$$

$$\text{dfF}(\hat{f}) = \text{tr}[\text{Cov}(y, X\hat{\beta})]/\sigma^2 = \text{tr}[\sigma^2 X(X^T X)^{-1} X^T]/\sigma^2 = p.$$

- Suppose  $p \geq n$  and  $X$  has full (row) rank, and we take  $\hat{f}$  to be **min  $\ell_2$ -norm least squares** predictor  $\hat{f}(X) = X\hat{\beta}$ , where

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{\|\beta\|_2 : X\beta = y\} = (X^T X)^\dagger X^T y.$$

$$\text{dfF}(\hat{f}) = \text{tr}[\text{Cov}(y, X\hat{\beta})]/\sigma^2 = \text{tr}[\sigma^2 X^T (X X^T)^{-1} X]/\sigma^2 = n.$$

Thus,  $\text{dfF}(\hat{f})$  is  $p$  for  $p \leq n$ , but is always  $n$  for  $p \geq n$  (not meaningful).

This is **fixed-X** degrees of freedom. How to extend for **random-X setting**?

# Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

Examples

Discussion

## Re-interpreting fixed-X degrees of freedom

Fixed-X degrees of freedom is a standard algorithm specific measure of complexity, but no notion of random-X degrees of freedom we know of.

We cast fixed-X degrees of freedom from a **different perspective**.

- Define fixed-X optimism of  $\hat{f}$  by  $\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{ErrT}(\hat{f})$ .
- Consider the following family of “reference” models:
  - $\mathcal{A}^{\text{ref}}$  is the **least squares** reference algorithm,
  - $(U_k, v)$  is random design with  $k$  features, and noise with level  $\sigma^2$ .
- Recall that  $\text{dfF}(\mathcal{A}^{\text{ref}}(U_k, v)) = k$  so long as  $\text{rank}(U_k) = k$ .
- Thus, for a fitting procedure  $\hat{f} = \mathcal{A}(X, y)$ ,  $\text{dfF}(\hat{f})$  is also equal to the value of  $k$  that satisfy the following relation:

$$\text{OptF}(\mathcal{A}(X, y)) = \text{OptF}(\mathcal{A}^{\text{ref}}(U_k, v)) \quad (\text{dfF})$$

## Emergent random-X degrees of freedom

“Matching optimism” interpretation can be extended to random-X setting and leads to the definition of random-X degrees of freedom.

- Define **random-X optimism** of  $\hat{f}$  by  $\text{OptR}(\hat{f}) = \text{ErrR}(\hat{f}) - \text{ErrT}(\hat{f})$ , where  $\text{ErrR}(\hat{f}) = \mathbb{E}[(y_0 - \hat{f}(x_0))^2]$  is the random-X prediction error.
- We thus define the random-X degrees of freedom,  $\text{dfR}(\hat{f})$ , of any predictor  $\hat{f} = \mathcal{A}(X, y)$ , as the value of  $k$  for which the following relation holds:

$$\text{OptR}(\mathcal{A}(X, y)) = \text{OptR}(\mathcal{A}^{\text{ref}}(U_k, v)) \quad (\text{dfR, emergent})$$

Recall here:

- $\mathcal{A}^{\text{ref}}$  is the **least squares** reference algorithm,
- $(U_k, v)$  is random design with  $k$  features, and noise with level  $\sigma^2$ .

We call the measure **emergent** random-X degrees of freedom.

## Intrinsic random-X degrees of freedom

- The emergent random-X degrees of freedom,  $\text{dfR}(\hat{f})$ , depends of both the the predictor  $\hat{f}$  and the underlying regression function  $f$ .
- When matching optimisms, the observed random-X optimism of  $\hat{f}$  consists of bias, which may inflate the degrees of freedom.
- We thus also define **intrinsic** random-X degrees of freedom, denoted by  $\text{dfR}^i$ , as the  $k$  for which the following relation holds:

$$\text{OptR}(\mathcal{A}(X, \nu)) = \text{OptR}(\mathcal{A}^{\text{ref}}(U_k, \nu)) \quad (\text{dfR}, \text{intrinsic})$$

The intrinsic random-X degrees of freedom measures the **inherent** complexity of the predictor  $\hat{f}$  in terms of overfitting to “pure noise”.

## Computing random-X degrees of freedom

We have defining equations for random-X degrees of freedom, but how do we actually compute the solution in practice?

We need three numbers:

- $\text{OptR}(\mathcal{A}(X, y))$ : emergent random-X optimism  
Either supplied by user or via **cross or held-out validation**
- $\text{OptR}(\mathcal{A}(X, v))$ : intrinsic random-X optimism  
**Simulate** at noise level if known, otherwise average
- $\text{OptR}(\mathcal{A}^{\text{ref}}(X, v))$ : reference optimism  
Either simulate or use **invariant asymptotic limit** that holds under quite generic conditions on the random design matrix and noise:

$$\frac{\text{OptR}(\mathcal{A}^{\text{ref}}(U_k, v))}{\sigma^2} \rightarrow \frac{1 - (1 - \xi)^2}{1 - \xi} \text{ as } n, p \rightarrow \infty \text{ and } p/n \rightarrow \xi \in (0, 1),$$

## Computing random-X degrees of freedom

Under the asymptotic limit, solving for the random-X degrees of freedom by matching optimisms leads:

$$\text{dfR}(\hat{f})/n \rightarrow 1 + \frac{\psi}{2} - \sqrt{1 + \frac{\psi^2}{4}}.$$

where  $\psi$  is the normalized random-X optimism of the given predictor.

Remarks:

- There is a unique number in  $[0, n]$  satisfying the desired relations.
- This is an interpretable range for random-X degrees of freedom:
  - The least complex predictor has dfR of 0,
  - The most complex predictor has dfR of  $n$ , as if saturated model.



# Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

Examples

Discussion

## Ordinary, min $\ell_2$ -norm and min $\ell_1$ -norm least squares

Given  $(X, y)$ , consider predictors of the form  $\hat{f}(x) = x^\top \hat{\beta}$ , where

- Ordinary least squares

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2$$

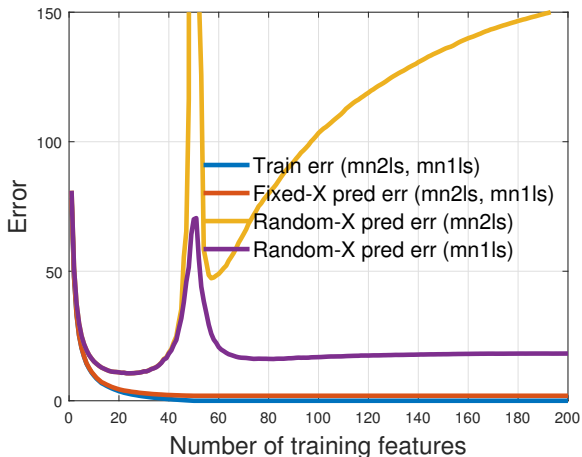
- Min  $\ell_2$ -norm least squares (mn2ls)

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{\|\beta\|_2 : y = X\beta\}$$

- Min  $\ell_1$ -norm least squares (mn1ls)

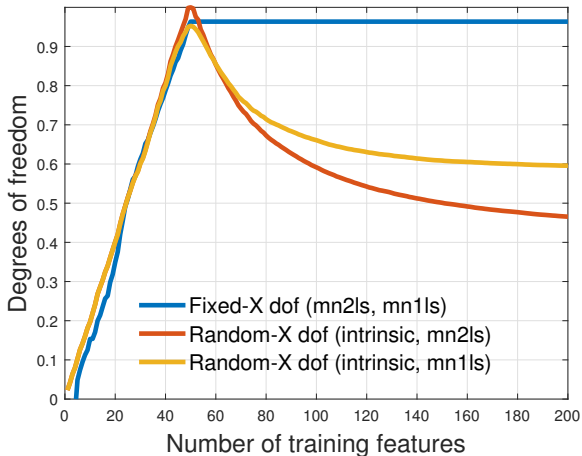
$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{\|\beta\|_1 : y = X\beta\}$$

## Errors (train, fixed-X, random-X) versus model size



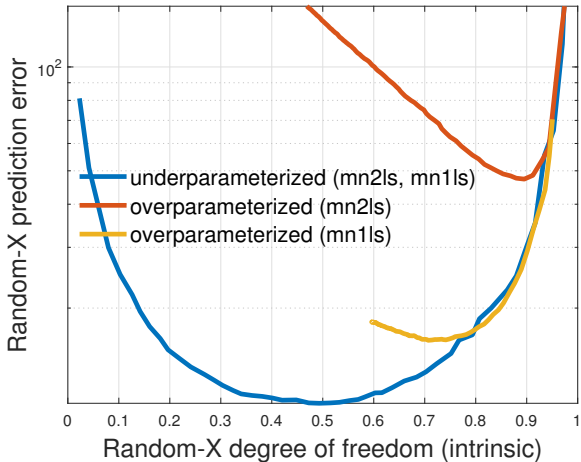
- Fixed data with  $n = 50$  and response non-linear in  $p = 200$  features
- Model class: estimators fitted on nested subsets of 1 to 200 features
- Train, fixed-X: descent then constant; random-X: double descent

## Degrees of freedom (fixed-X, random-X) versus model size



- Fixed data with  $n = 50$  and response non-linear in  $p = 200$  features
- Model class: estimators fitted on nested subsets of 1 to 200 features
- Fixed-X: increase then constant; random-X: increase then decrease

## Prediction error (random-X) versus model size



- Same setup as before
- Underparameterized:  $U$ -curve; overparameterized: also  $U$ -curve!
- Punchline: map from overparameterized to underparameterized

# Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

Examples

Discussion

## Discussion and future directions

A high-level view of the work:

- Suppose we are given a family of models for which we want a complexity measure under a **specific error metric**.
- Construct a family of **“reference” models** spanning same optimisms.
- Find the model in the reference family that is closest to the **observed optimism**. Declare complexity as complexity of that reference model.

Key relation:

$$\text{OptR}(\hat{f}) = \text{OptR}(\hat{f}^{\text{ref}})$$

Going beyond ...

- Attribute total complexity to various components: bias, variance, covariate shift, etc.
- Other error metrics beyond squared error
- Unsupervised setting?

Thanks for listening!

Questions/comments/thoughts?