

Revisiting Model Complexity in the Wake of Overparameterized Machine Learning

Pratik Patil

University of California, Berkeley

MDS 2024

Based on joint work with Jin-Hong Du and Ryan Tibshirani

<https://pratikpatil.io/papers/model-complexity.pdf>

Overparametrization in machine learning

Modern machine learning models typically fit 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 problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

Overparametrization in machine learning

Modern machine learning models typically fit 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 problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

Overparametrization in machine learning

Modern machine learning models typically fit 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 problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

Overparametrization in machine learning

Modern machine learning models typically fit 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 problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

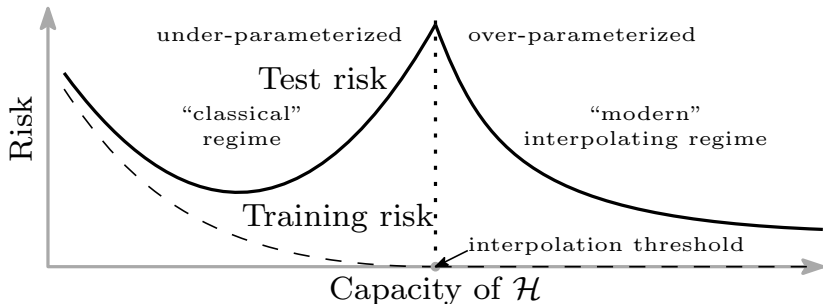
Overparametrization in machine learning

Modern machine learning models typically fit 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 problem
- **Generalization**: despite overfitting, models generalize well in practice

This talk is about generalization aspect in overparameterized learning.

The double/multiple descent phenomenon



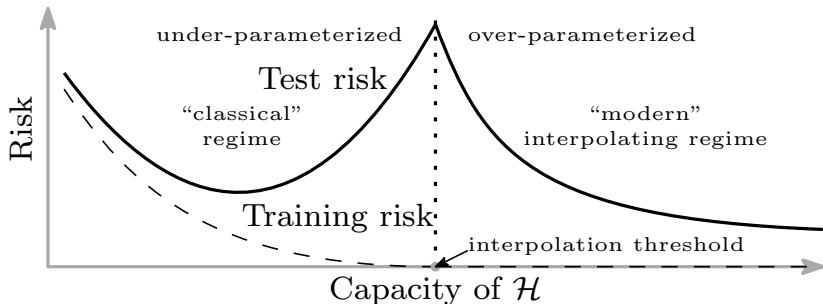
"Reconciling modern machine learning practice and the bias variance tradeoff"

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"?

The double/multiple descent phenomenon



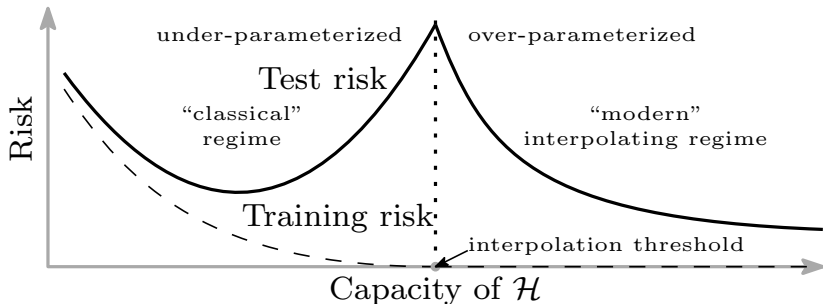
"Reconciling modern machine learning practice and the bias variance tradeoff"

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"?

The double/multiple descent phenomenon



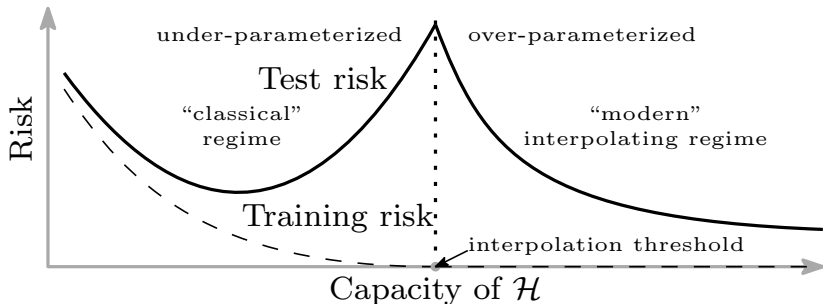
"Reconciling modern machine learning practice and the bias variance tradeoff"

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"?

The double/multiple descent phenomenon



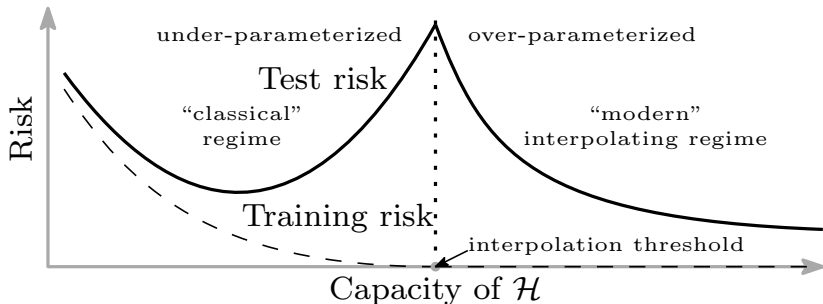
"Reconciling modern machine learning practice and the bias variance tradeoff"

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"?

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”?

Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
 - Hastie, Montanari, Rosset, Tibshirani, 2019
 - Belkin, Hsu, Xu, 2019
 - Muthukumar, Vodrahalli, Sahai, 2019
 - Bartlett, Long, Lugosi, Tsigler, 2019
 - Mei, Montanari, 2019
- Kernel regression
 - Liang, Rakhlin, 2018
 - Liang, Rakhlin, Zhai, 2019
- Local methods
 - Belkin, Hsu, Mitra, 2018
 - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
 - Hastie, Montanari, Rosset, Tibshirani, 2019
 - Belkin, Hsu, Xu, 2019
 - Muthukumar, Vodrahalli, Sahai, 2019
 - Bartlett, Long, Lugosi, Tsigler, 2019
 - Mei, Montanari, 2019
- Kernel regression
 - Liang, Rakhlin, 2018
 - Liang, Rakhlin, Zhai, 2019
- Local methods
 - Belkin, Hsu, Mitra, 2018
 - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
 - Hastie, Montanari, Rosset, Tibshirani, 2019
 - Belkin, Hsu, Xu, 2019
 - Muthukumar, Vodrahalli, Sahai, 2019
 - Bartlett, Long, Lugosi, Tsigler, 2019
 - Mei, Montanari, 2019
- Kernel regression
 - Liang, Rakhlin, 2018
 - Liang, Rakhlin, Zhai, 2019
- Local methods
 - Belkin, Hsu, Mitra, 2018
 - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
 - Hastie, Montanari, Rosset, Tibshirani, 2019
 - Belkin, Hsu, Xu, 2019
 - Muthukumar, Vodrahalli, Sahai, 2019
 - Bartlett, Long, Lugosi, Tsigler, 2019
 - Mei, Montanari, 2019
- Kernel regression
 - Liang, Rakhlin, 2018
 - Liang, Rakhlin, Zhai, 2019
- Local methods
 - Belkin, Hsu, Mitra, 2018
 - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

Recent theoretical developments

Understanding generalization of interpolators in simpler settings:

- Linear regression
 - Hastie, Montanari, Rosset, Tibshirani, 2019
 - Belkin, Hsu, Xu, 2019
 - Muthukumar, Vodrahalli, Sahai, 2019
 - Bartlett, Long, Lugosi, Tsigler, 2019
 - Mei, Montanari, 2019
- Kernel regression
 - Liang, Rakhlin, 2018
 - Liang, Rakhlin, Zhai, 2019
- Local methods
 - Belkin, Hsu, Mitra, 2018
 - Belkin, Rakhlin, Tsybakov, 2018
- and many more ...

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

What do we currently understand?

- In nearly all applications, current practice suggests we should design models to be massively **overparametrized**
- Once trained (typically by SGD), these models **interpolate** the training data (achieve zero training error)
- Still they are capable of having (often do have) **good test error**

Current understanding of this? In full theoretical rigor, not great.

However, the story is fairly well-understood for linear models, kernel models, and random feature models. See, e.g., nice monographs:

- Bartlett, Montanari, and Rakhlin (2021), “Deep learning: a statistical viewpoint”
- Belkin (2021), “Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation”

Goals of this work

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 ℓ_2/ℓ_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**.

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 ℓ_2/ℓ_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**.

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**.

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**.

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 ℓ_2/ℓ_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**.

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**.

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 ℓ_2/ℓ_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)

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 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 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 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 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 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 (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Degrees of freedom in statistics

Degrees of freedom (df) is a classical topic in statistics, dating back to Mallows (1973), Stein (1981), Efron (1986)

Given an estimator \hat{f} of the regression function (i.e., $\hat{f}(x)$ estimates $\mathbb{E}[y|x]$), trained on (x_i, y_i) , $i = 1, \dots, n$, we define

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, \hat{f}(x_i))$$

where $\sigma^2 = \text{Var}(y|x)$, and each x_i is treated as fixed

Key fact. For \hat{f} the **least squares** estimator of response $Y \in \mathbb{R}^n$ on feature matrix $X \in \mathbb{R}^{n \times p}$, with p linearly independent columns,

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \text{tr}(\text{Cov}(Y, P_X Y)) = \text{tr}(P_X) = p$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Connection to fixed-X error

This definition is intrinsically connected to **fixed-X** test error:

$$\underbrace{\mathbb{E}_{Y, Y^*} \left[\frac{1}{n} \sum_{i=1}^n (y_i^* - \hat{f}(x_i))^2 \right]}_{\text{ErrF}(\hat{f})} = \underbrace{\mathbb{E}_Y \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_F(\hat{f})} + \frac{2\sigma^2}{n} \text{df}(\hat{f})$$

The difference between test and training error is called **optimism**:

$$\text{OptF}(\hat{f}) = \text{ErrF}(\hat{f}) - \text{train}_F(\hat{f})$$

This nicely motivates degrees of freedom and connects it to practice

Also points to a big **shortcoming** ... fixed-X error is not as relevant to modern stat/ML practice which is driven by **random-X** error

$$\text{ErrR}(\hat{f}) = \mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]$$

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be *very different*.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Fixed-X \neq random-X, especially for interpolators

Classically—for smooth functions, in low dimensions—we can tie fixed-X and random-X errors together (e.g., by empirical process theory)

But beyond this—nonsmooth functions, or high dimensions—they can be **very different**.

Epitomized by generalizing interpolators:

- $\text{ErrR}(\hat{f}) \rightarrow 0$, but
- $\text{ErrF}(\hat{f}) = \mathbb{E}[n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2] = 2\sigma^2$

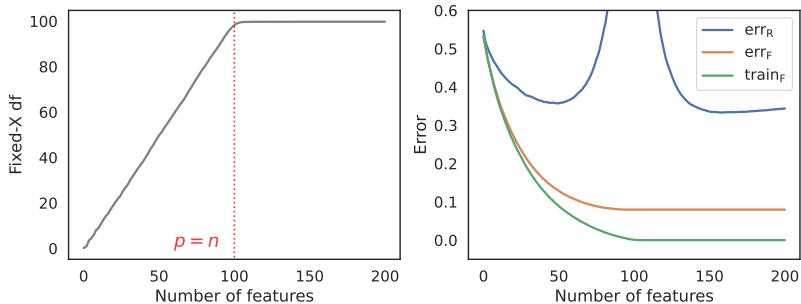
Similarly, for any interpolator \hat{f} :

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(y_i, y_i) = n$$

which is useless!

Example: ridgeless least squares df

Classical, fixed-X df for ridgeless least squares, with $n = 100$:



Can we generalize df to random-X, and beyond?

Goals:

- Generalize df so that it connects to **random-X** error, and gives meaningful answers for **any estimator** (even interpolators)
- Allow for **decomposition** of df according to some user-specified components

Disclaimer: *what follows is embarassingly simple ...*

Can we generalize df to random-X, and beyond?

Goals:

- Generalize df so that it connects to **random-X** error, and gives meaningful answers for **any estimator** (even interpolators)
- Allow for **decomposition** of df according to some user-specified components

Disclaimer: *what follows is embarassingly simple ...*

Can we generalize df to random-X, and beyond?

Goals:

- Generalize df so that it connects to **random-X** error, and gives meaningful answers for **any estimator** (even interpolators)
- Allow for **decomposition** of df according to some user-specified components

Disclaimer: *what follows is embarassingly simple ...*

Can we generalize df to random-X, and beyond?

Goals:

- Generalize df so that it connects to **random-X** error, and gives meaningful answers for **any estimator** (even interpolators)
- Allow for **decomposition** of df according to some user-specified components

Disclaimer: *what follows is embarassingly simple ...*

Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

Discussion

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

What is complexity?

What makes up a **measure of complexity**? Two things:

- a metric
- a reference model

That is, the metric **assigns a number** to the given model \hat{f} , and the reference model \hat{f}^{ref} **provides units**, so we can interpret the metric

A bit more detail:

- Metric should be “negatively oriented” for complexity—smaller values mean less complex
- Reference model should be something whose parameters we are “happy to count”

Re-interpreting degrees of freedom

We can actually **re-interpret** the classical definition of df in this light, with: $\text{metric} = \text{OptF}$, and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, *we will define* $df(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $df(\hat{f}) = 5.3$ means \hat{f} has the same fixed-X optimism as least squares with 5.3 parameters

Re-interpreting degrees of freedom

We can actually **re-interpret** the classical definition of df in this light, with: $\text{metric} = \text{OptF}$, and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, *we will define* $df(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $df(\hat{f}) = 5.3$ means \hat{f} has the same fixed- X optimism as least squares with 5.3 parameters

Re-interpreting degrees of freedom

We can actually **re-interpret** the classical definition of df in this light, with: $\text{metric} = \text{OptF}$, and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, *we will define* $df(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $df(\hat{f}) = 5.3$ means \hat{f} has the same fixed- X optimism as least squares with 5.3 parameters

Re-interpreting degrees of freedom

We can actually **re-interpret** the classical definition of df in this light, with: $\text{metric} = \text{OptF}$, and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, *we will define* $df(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $df(\hat{f}) = 5.3$ means \hat{f} has the same fixed-X optimism as least squares with 5.3 parameters

Re-interpreting degrees of freedom

We can actually **re-interpret** the classical definition of df in this light, with: metric = OptF, and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, *we will define* $\text{df}(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $\text{df}(\hat{f}) = 5.3$ means \hat{f} has the same fixed-X optimism as least squares with 5.3 parameters



Re-interpreting degrees of freedom

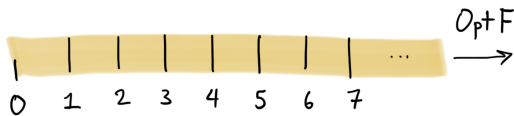
We can actually **re-interpret** the classical definition of df in this light, with: metric = OptF , and $\hat{f}^{\text{ref}} = \text{least squares}$

That is, we will define $\text{df}(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptF}(\hat{f}) = \text{OptF}(\hat{f}_k^{\text{ls}})$$

where \hat{f}_k^{ls} is the least squares estimator of the given $Y \in \mathbb{R}^n$ on a feature matrix $X \in \mathbb{R}^{n \times k}$ with k linearly independent columns

Think about it this way: $\text{df}(\hat{f}) = 5.3$ means \hat{f} has the same fixed- X optimism as least squares with 5.3 parameters



Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Lifting this idea to random-X

Naturally we can lift this to the **random-X** setting. We define

$$\text{OptR}(\hat{f}) = \underbrace{\mathbb{E}_{X, Y, x_0, y_0} [(y_0 - \hat{f}(x_0))^2]}_{\text{ErrR}(\hat{f})} - \underbrace{\mathbb{E}_{X, Y} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]}_{\text{train}_R(\hat{f})}$$

We will again use \hat{f}_k^{ls} for reference class: least squares on k features. However, in the random-X setting, it matters:

- **what** these features are (law of x_i)
- **how** they relate to the response (law of $y_i|x_i$)

(Recall, none of this mattered in the fixed-X setting ...)

So what data generating distribution should we use for our reference model in the random-X setting?

Universal asymptotics

Consider i.i.d. samples (x_i, y_i) , $i = 0, \dots, n + 1$, from the following standard random matrix theory (RMT) model:

- $x_i = \Sigma^{1/2} z_i$, where $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic with eigenvalues bounded away from 0 and ∞ ; and $z_i \in \mathbb{R}^p$ has i.i.d. coordinates with zero mean, unit variance, and finite 4th moment
- $y_i = x_i^\top \beta + \epsilon_i$, where ϵ_i has zero mean, unit variance, and is independent of x_i

Theorem. Under the standard RMT model above, as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma < 1$, the random- X optimism least squares satisfies:¹

$$\text{OptR}(\hat{f}_p^{\text{ls}}) - \sigma^2 \left(\frac{n}{n-p} - \frac{n-p}{n} \right) \rightarrow 0$$

¹Convergence is actually almost sure wrt training features $X \in \mathbb{R}^{n \times p}$; need uniform integrability to get convergence in expectation.

Universal asymptotics

Consider i.i.d. samples (x_i, y_i) , $i = 0, \dots, n + 1$, from the following standard random matrix theory (RMT) model:

- $x_i = \Sigma^{1/2} z_i$, where $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic with eigenvalues bounded away from 0 and ∞ ; and $z_i \in \mathbb{R}^p$ has i.i.d. coordinates with zero mean, unit variance, and finite 4th moment
- $y_i = x_i^\top \beta + \epsilon_i$, where ϵ_i has zero mean, unit variance, and is independent of x_i

Theorem. Under the standard RMT model above, as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma < 1$, the random- X optimism least squares satisfies:¹

$$\text{OptR}(\hat{f}_p^{\text{ls}}) - \sigma^2 \left(\frac{n}{n-p} - \frac{n-p}{n} \right) \rightarrow 0$$

¹Convergence is actually almost sure wrt training features $X \in \mathbb{R}^{n \times p}$; need uniform integrability to get convergence in expectation.

Universal asymptotics

Consider i.i.d. samples (x_i, y_i) , $i = 0, \dots, n + 1$, from the following standard random matrix theory (RMT) model:

- $x_i = \Sigma^{1/2} z_i$, where $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic with eigenvalues bounded away from 0 and ∞ ; and $z_i \in \mathbb{R}^p$ has i.i.d. coordinates with zero mean, unit variance, and finite 4th moment
- $y_i = x_i^\top \beta + \epsilon_i$, where ϵ_i has zero mean, unit variance, and is independent of x_i

Theorem. Under the standard RMT model above, as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma < 1$, the random- X optimism least squares satisfies:¹

$$\text{OptR}(\hat{f}_p^{\text{ls}}) - \sigma^2 \left(\frac{n}{n-p} - \frac{n-p}{n} \right) \rightarrow 0$$

¹Convergence is actually almost sure wrt training features $X \in \mathbb{R}^{n \times p}$; need uniform integrability to get convergence in expectation.

Universal asymptotics

Consider i.i.d. samples (x_i, y_i) , $i = 0, \dots, n + 1$, from the following standard random matrix theory (RMT) model:

- $x_i = \Sigma^{1/2} z_i$, where $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic with eigenvalues bounded away from 0 and ∞ ; and $z_i \in \mathbb{R}^p$ has i.i.d. coordinates with zero mean, unit variance, and finite 4th moment
- $y_i = x_i^\top \beta + \epsilon_i$, where ϵ_i has zero mean, unit variance, and is independent of x_i

Theorem. Under the standard RMT model above, as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma < 1$, the random- X optimism least squares satisfies:¹

$$\text{OptR}(\hat{f}_p^{\text{ls}}) - \sigma^2 \left(\frac{n}{n-p} - \frac{n-p}{n} \right) \rightarrow 0$$

¹Convergence is actually almost sure wrt training features $X \in \mathbb{R}^{n \times p}$; need uniform integrability to get convergence in expectation.

Universal asymptotics

Consider i.i.d. samples (x_i, y_i) , $i = 0, \dots, n + 1$, from the following standard random matrix theory (RMT) model:

- $x_i = \Sigma^{1/2} z_i$, where $\Sigma \in \mathbb{R}^{p \times p}$ is deterministic with eigenvalues bounded away from 0 and ∞ ; and $z_i \in \mathbb{R}^p$ has i.i.d. coordinates with zero mean, unit variance, and finite 4th moment
- $y_i = x_i^\top \beta + \epsilon_i$, where ϵ_i has zero mean, unit variance, and is independent of x_i

Theorem. Under the standard RMT model above, as $n, p \rightarrow \infty$ with $p/n \rightarrow \gamma < 1$, the random- X optimism least squares satisfies:¹

$$\text{OptR}(\hat{f}_p^{\text{ls}}) - \sigma^2 \left(\frac{n}{n-p} - \frac{n-p}{n} \right) \rightarrow 0$$

¹Convergence is actually almost sure wrt training features $X \in \mathbb{R}^{n \times p}$; need uniform integrability to get convergence in expectation.

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**:
 $df^e(f) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e.,
 $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**: $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**: $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**:
 $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e.,
 $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**:
 $\text{df}^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e.,
 $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**: $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**: $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**:
 $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e.,
 $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Emergent random-X df

Using this we can define a quantity we call **emergent random-X df**: $df^e(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original data}\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on the original data
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “emergent” random-X df as it refers to the complexity that “emerges” from the optimism of \hat{f} on data at hand. The LHS generally has both **bias and variance** components, whereas the RHS is **pure variance**

Illustration: least squares as reference

Random-X optimism for least squares, with $n = 50$:

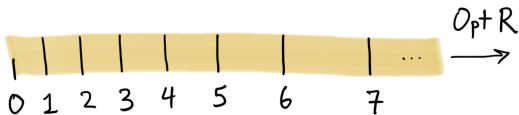
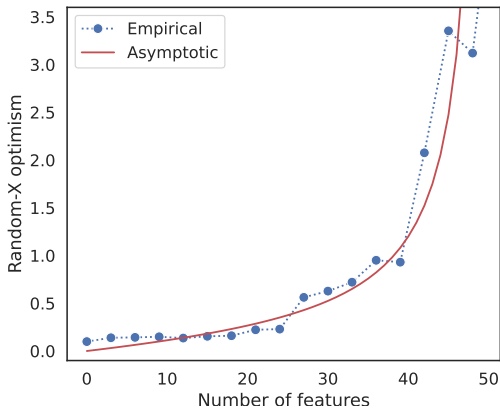
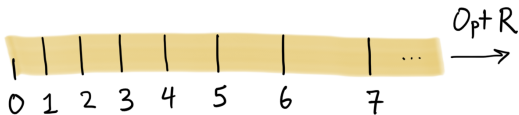
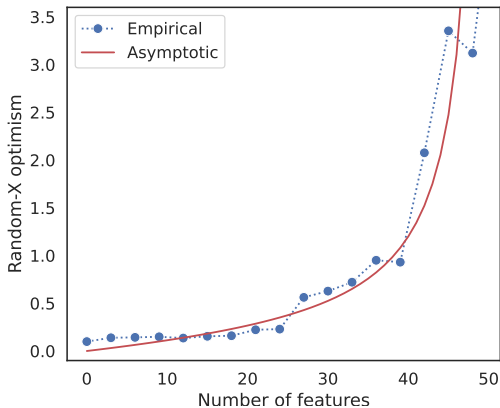


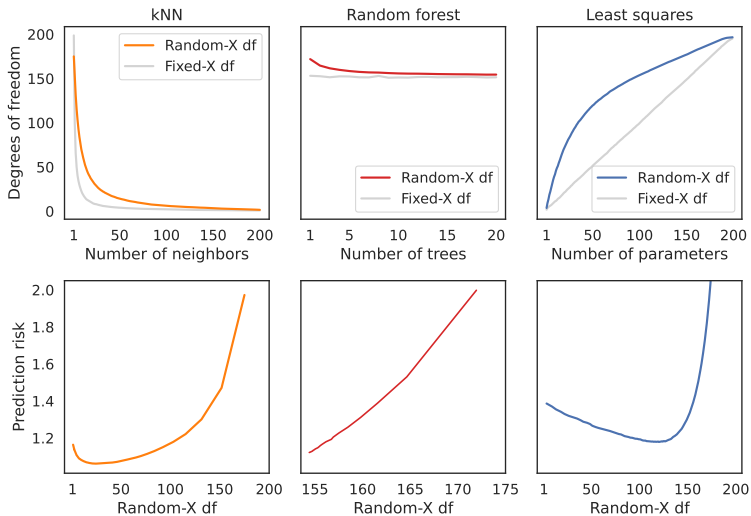
Illustration: least squares as reference

Random-X optimism for least squares, with $n = 50$:



Example: emergent random-X df

Ex with $n = 200$ samples, $d = 300$ features, $\mathbb{E}[y|x]$ nonlinear in x :



Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**:

$df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**:
 $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Intrinsic random-X df

We can also define a different quantity called **intrinsic random-X df**: $df^i(\hat{f}) = k$, for the number $k \geq 0$ such that

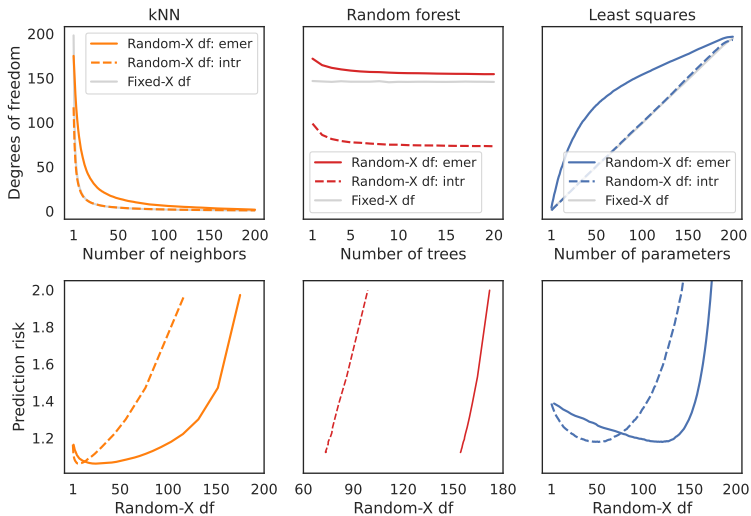
$$\text{OptR}(\hat{f}; \{\text{original } x, \text{ noisy } y\}) = \text{OptR}(\hat{f}_k^{\text{ls}}; \{\text{RMT data}\})$$

- LHS: random-X opt of the given model \hat{f} on original features x_i , but with **pure noise** for the response $y_i \sim N(0, \sigma^2)$
- RHS: random-X opt of least squares \hat{f}_k^{ls} on “RMT data”, i.e., $x_i = \Sigma^{1/2} z_i$, $y_i = x_i^\top \beta + \epsilon_i$, as before
- RHS admits **simple asymptotic approximation**: $\sigma^2 \left(\frac{n}{n-k} - \frac{n-k}{n} \right)$

We call this “intrinsic” random-X df as it refers to the complexity internal to the model \hat{f} , regardless of the data. Note that both the LHS and RHS reflect **pure variance**

Example: intrinsic random-X df

Ex with $n = 200$ samples, $d = 300$ features, $\mathbb{E}[y|x]$ nonlinear in x :

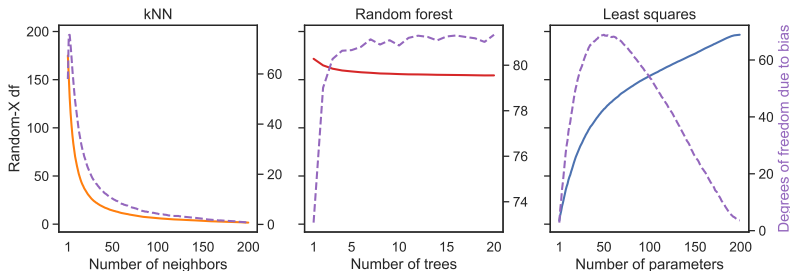


Bias attribution

By subtracting emergent and intrinsic degrees of freedom, we are left with the **df due to bias**:

$$df^{\text{bias}}(\hat{f}) = df^e(\hat{f}) - df^i(\hat{f})$$

Back to our example:

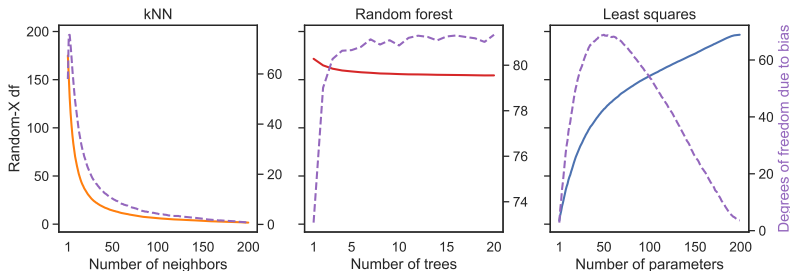


Bias attribution

By subtracting emergent and intrinsic degrees of freedom, we are left with the **df due to bias**:

$$df^{\text{bias}}(\hat{f}) = df^e(\hat{f}) - df^i(\hat{f})$$

Back to our example:



General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $df^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (df^{S \cup \{e_i\}}(\hat{f}) - df^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = df^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $df^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (df^{S \cup \{e_i\}}(\hat{f}) - df^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = df^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

General decomposition

Generic decompositions are possible. Let e_1, \dots, e_m be any list of user-chosen **error components**. E.g., $e_1 = \text{bias}$, $e_2 = \text{covariate shift}$, and so on

For any subset $S \subseteq \{e_1, \dots, e_m\}$, let $\text{df}^S(\hat{f}) = k$, for the number $k \geq 0$ such that

$$\text{OptR}(\hat{f}; \{\text{data subject to } S\}) = \text{OptR}(\hat{f}_k^S; \{\text{RMT data}\})$$

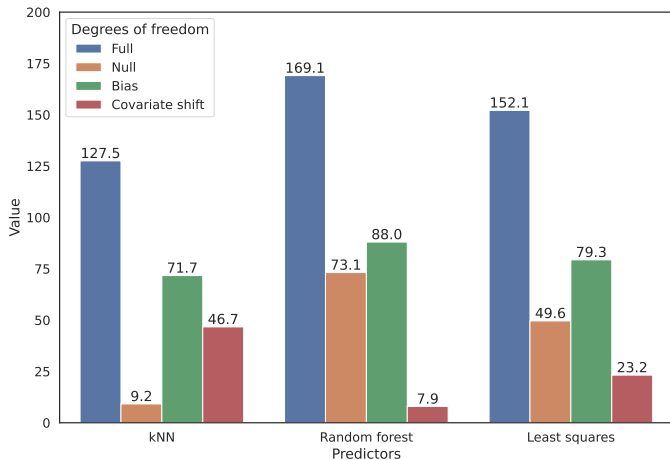
We can then **attribute** d_i df to each error component e_i , as follows:

$$d_i = \sum_{S \subseteq \{e_1, \dots, e_m\} \setminus \{e_i\}} \frac{|S|!(m - |S| - 1)!}{m!} (\text{df}^{S \cup \{e_i\}}(\hat{f}) - \text{df}^S(\hat{f}))$$

This is an instance of a **Shapley value**. Therefore it obeys all of the Shapley axioms; in particular, efficiency: $\sum_{i=1}^m d_i = \text{df}^e(\hat{f})$

Example: bias and covariate shift

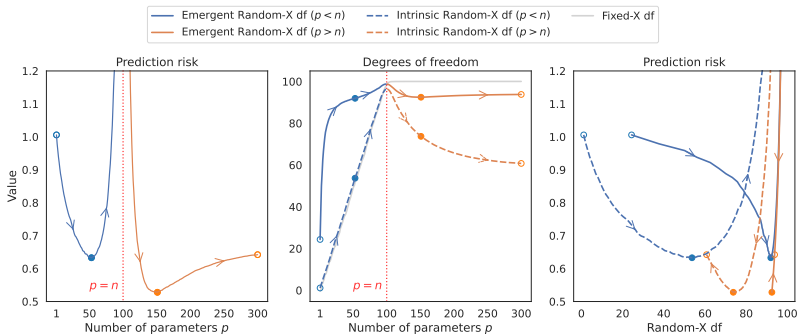
Ex with $n = 200$ samples, $d = 300$ features, $\mathbb{E}[y|x]$ nonlinear in x , and covariate shift ($\Sigma \rightarrow \tilde{\Sigma}$):



Double descent, revisited

We can use random-X df (any flavor) to **reparametrize** error curve for models with **double descent**. The df map is not monotone, but it shows that in the overparametrized regime, the effective number of parameters can actually be small

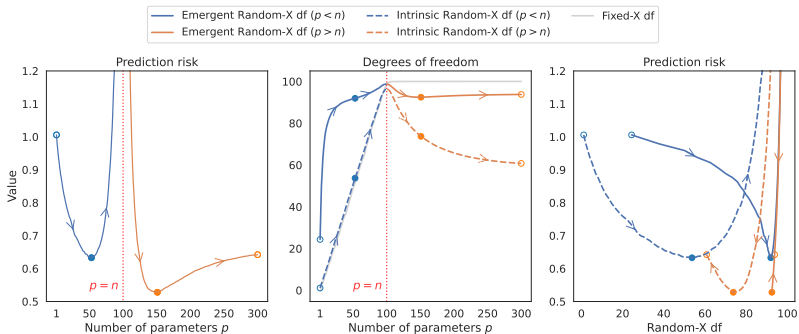
Back to our running example:



Double descent, revisited

We can use random-X df (any flavor) to **reparametrize** error curve for models with **double descent**. The df map is not monotone, but it shows that in the overparametrized regime, the effective number of parameters can actually be small

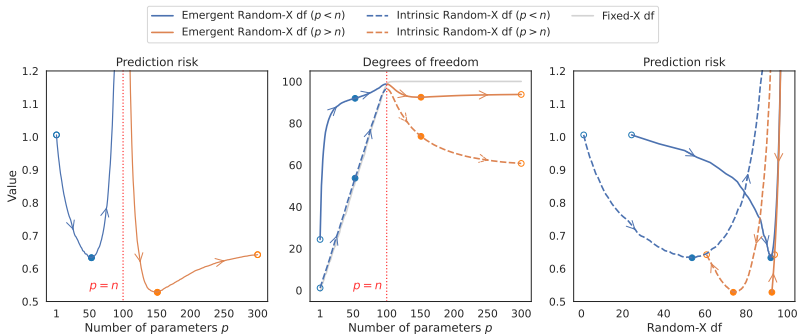
Back to our running example:



Double descent, revisited

We can use random-X df (any flavor) to **reparametrize** error curve for models with **double descent**. The df map is not monotone, but it shows that in the overparametrized regime, the effective number of parameters can actually be small

Back to our running example:



Outline

Degree of freedom (fixed-X setting)

Degrees of freedom (random-X setting)

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 . . .

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

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 ...

- Other error metrics beyond squared error
- Unsupervised setting?

Thanks for listening!

Questions/comments/thoughts?