

# MACHINE LEARNING IN BIOINFORMATICS

## MODEL SELECTION AND REGULARIZATION

Philipp Benner

*philipp.benner@bam.de*

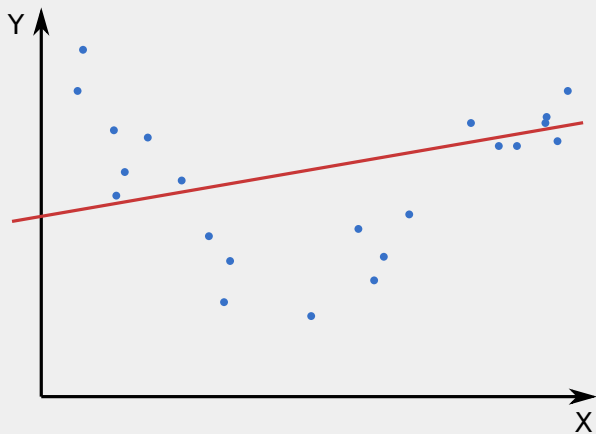
VP.1 - eScience

Federal Institute of Materials Research and Testing (BAM)

April 25, 2024

# MODEL SELECTION PROBLEM

# MODEL SELECTION



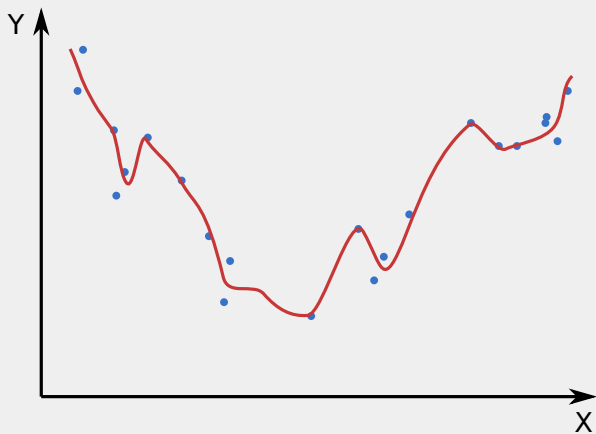
Linear model class

# MODEL SELECTION



Quadratic model class

# MODEL SELECTION



Polynomial model class

# **BIAS-VARIANCE DECOMPOSITION AND TRADEOFF**

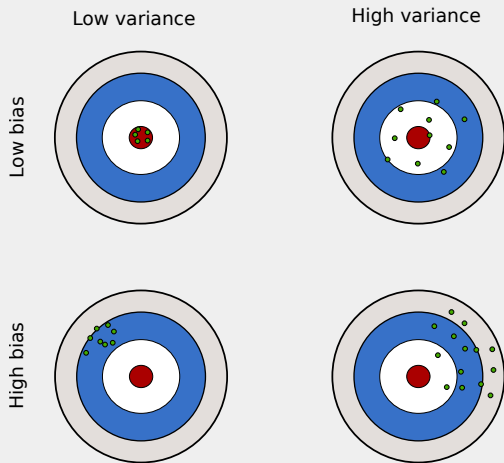
# BIAS-VARIANCE DECOMPOSITION

- Let  $\mathbf{Y}$ ,  $\mathbf{X}$  and  $\epsilon$  be random variables such that  $\mathbf{Y} = f(\mathbf{X}) + \epsilon$ , with  $\mathbb{E}[\epsilon] = 0$  and  $\text{var}[\epsilon] = \sigma^2$
- Assume that  $\hat{f}_D$  has been estimated on some training data  $D = (X, y)$ , where  $X$  is a matrix of  $n$  observations from  $\mathbf{X}$  and  $y$  a vector of  $n$  observations from  $\mathbf{Y}$
- At a query point  $x$  we have

$$\mathbb{E}_{\mathbf{Y}, D}[(\mathbf{Y} - \hat{f}_D(x))^2] = \underbrace{[\mathbb{E}_D \hat{f}_D(x) - f(x)]^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}_D [\hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x)]^2}_{\text{Variance}} + \sigma^2$$

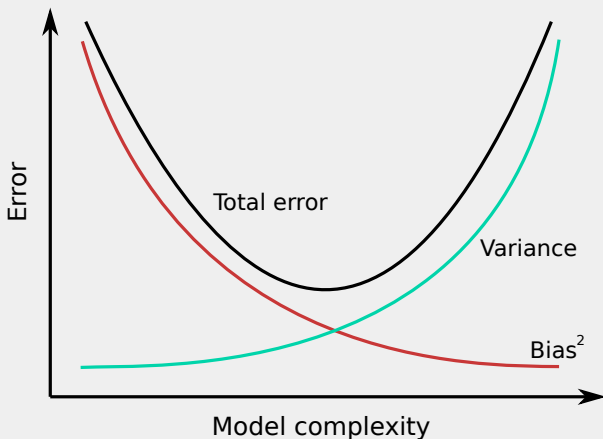
- bias: Is there a bias towards a particular kind of solution (e.g. linear model)? (inductive bias)
- variance: How much does the estimated model change if you train on a different data set? (overfitting)

# BIAS-VARIANCE DECOMPOSITION



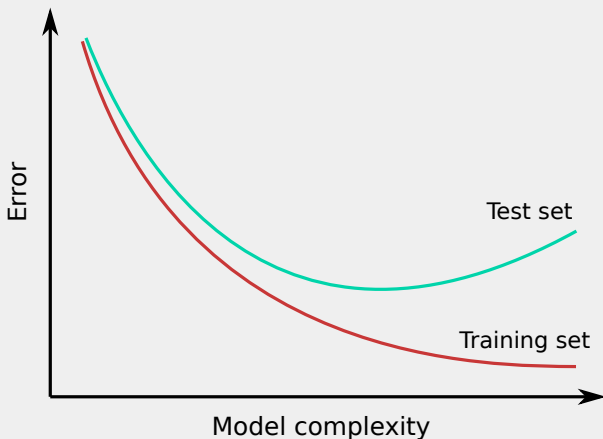


# BIAS-VARIANCE DECOMPOSITION



<sup>o</sup>Note that here we average over multiple data sets. On a single data set we might observe bumps when increasing model complexity

# BIAS-VARIANCE DECOMPOSITION



<sup>o</sup>Note that here we average over multiple data sets. On a single data set we might observe bumps when increasing model complexity

# BIAS-VARIANCE DECOMPOSITION - LESSONS LEARNED

- Every model comes with a bias
- More complex models have a smaller bias but larger variance
- A bias is required to reduce the variance, but introducing a good bias requires domain knowledge
- Classical statistics often uses unbiased estimators, which is nowadays often questioned
- Keep in mind: There is no free lunch!<sup>1</sup>

---

<sup>1</sup>The *no free lunch theorem* [Wolpert and Macready, 1997] tells us that there exists no generic model that works well on all domains, but we need to tailor our models to the data at hand in order to introduce a model bias, which reduces variance.

# COMPLEXITY MEASURES

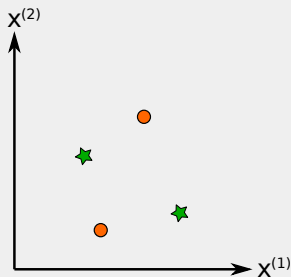
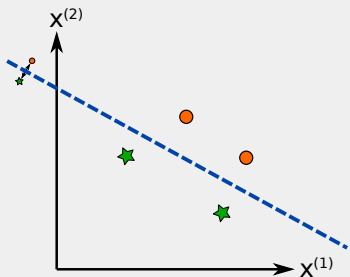
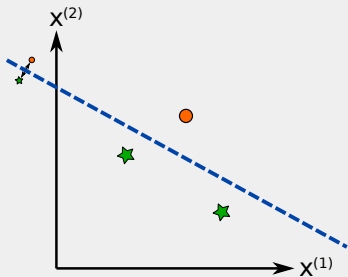
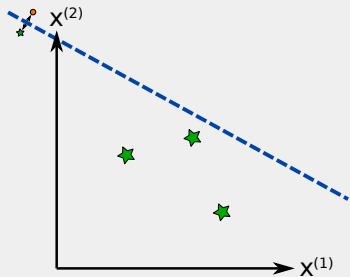
## VC-Dimension (Vapnik Chervonenkis)

Let  $\mathbb{F}_p$  be a set of classifiers on an  $n$ -dimensional input space. The VC-dimension  $VC(\mathbb{F}_p)$  is defined as the maximum number of points that can be correctly classified by at least one member of  $\mathbb{F}_p$ .

### ■ Examples:

- ▶ Linear classifier on  $\mathbb{R}^p$ :  $VC = p + 1$
- ▶ SVM with RBF kernel:  $VC = \infty$
- ▶ Neural network with  $n_e$  edges,  $n_v$  nodes and sigmoid activation function:  $\Omega(n_e^2) < VC < \mathcal{O}(n_e^2 n_v^2)$   
[Shalev-Shwartz and Ben-David, 2014, Section 20.4]

# COMPLEXITY OF CLASSIFIERS - VC DIMENSION



## Degrees of Freedom (DF) [Efron, 1986]

The **degrees of freedom** of an estimate  $\hat{y} = \hat{f}(X)$  is defined as

$$\text{df}(\hat{y}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{cov}(\hat{y}_i, y_i) = \frac{1}{\sigma^2} \text{tr cov}(\hat{y}, y),$$

where

- $X$  denotes a fixed set of  $n$  covariates of dimension  $p$
- $y = (y_1, \dots, y_n)$  is a vector of  $n$  observations from

$$\mathbf{Y} = f(X) + \epsilon$$

for some function  $f$ , assuming  $\mathbb{E}[\epsilon] = 0$  and  $\text{var}[\epsilon] = \sigma^2$

---

<sup>1</sup>df is normalized by the magnitude of the aleatory uncertainty ( $\sigma^2$ )

- Degrees of freedom for the OLS estimate:

$$\begin{aligned} \text{df}(\hat{y}) &= \frac{1}{\sigma^2} \text{tr} \text{cov}(\hat{y}, y) \\ &= \frac{1}{\sigma^2} \text{tr} \text{cov} \left( X(X^\top X)^{-1} X^\top y, y \right) \\ &= \frac{1}{\sigma^2} \text{tr} \left( X(X^\top X)^{-1} X^\top \right) \text{cov}(y, y) \\ &= \text{tr} \left( X(X^\top X)^{-1} X^\top \right) \\ &= p \end{aligned}$$

- $\text{df}(\hat{y}) = p$ , i.e. the number of parameters, assuming independent feature vectors (i.e. columns of  $X$ )
- This result holds for  $p < n$

---

$X(X^\top X)^{-1} X^\top$  is the hat matrix  $H \in \mathbb{R}^{n \times n}$ , hence  $\text{df}(\hat{y}) = \text{rank}(H)$



- Ridge regression is defined as

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2 + \lambda \|\theta\|_2^2$$

for some regularization strength  $\lambda \geq 0$

- The ridge estimator has

$$\text{df}(\hat{y}) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

degrees of freedom, where  $(d_j)_j$  are the singular values of  $X$

- Increasing  $\lambda$  decreases model complexity

- There is some criticism about used DF as measure of model complexity [Janson et al., 2015]
- In some cases, we also need  $X$  to be random [Luan et al., 2021]
- We will see other measures when turning to model selection

# MODEL SELECTION

# MODEL SELECTION APPROACHES

- A measure of accuracy or fit, such as the mean squared error (MSE), is not enough: Increasing model complexity will always lead to a better fit
- Estimating a model requires to minimize both
  - ▶ **in-sample-error** (loss on training data), and
  - ▶ **out-of-sample-error** (generalization error)
- Cross-validation (CV) estimates generalization error on left-out samples<sup>2</sup>
- Traditional statistics: Combine measure of accuracy (in-sample-error) with a penalty for complexity

---

<sup>2</sup>Heavy hyperparameter tuning using CV can lead to overfitting and requires to select a final holdout set

# MODEL SELECTION APPROACHES - LOO-CV

- Leave-one-out Cross-Validation (LOO-CV) at iteration  $i = 1, 2, \dots, n$ :
  - ▶ Compute estimate on data set without the  $i$ -th sample
  - ▶ Compute prediction error on the  $i$ -th sample
- Report the average prediction error over all  $n$  samples
- PRESS statistic (predicted residual error sum of squares):

$$\text{PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{-i})^2$$

where  $\hat{y}_{-i}$  is the prediction for the  $i$ -th sample where the model has been estimated on all but the  $i$ -th sample

# MODEL SELECTION APPROACHES - PRESS

- LOO-CV is very costly for large data sets and complex models
- $k$ -fold CV with  $k = 5$  or  $k = 10$  is often used in practice
- For (ridge) linear regression with mean squared error we can efficiently compute LOO-CV [Cook, 1977]

$$\begin{aligned}\text{PRESS} &= \sum_{i=1}^n (y_i - \hat{y}_{-i})^2 \\ &= \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1 - H_{ii})^2}\end{aligned}$$

- The matrix

$$H = X(X^T X + \lambda I)^{-1} X^T$$

is called the hat matrix, because it puts a hat on  $y$ , i.e.  $\hat{y} = Hy$

# MODEL SELECTION APPROACHES

- LOO-CV is computationally very expensive
- $k$ -fold CV is cheaper, but uses a large fraction of the data for testing
- Model performance could be better if this data was used for training
- Overfitting if we use CV for testing too many models (requires final hold out data)
- **Can we do model selection by using all data for training?**

# MODEL SELECTION APPROACHES - DF

- Assume again the following model

$$\mathbf{Y} = f(X) + \epsilon$$

where  $X \in \mathbb{R}^{n \times p}$  is a fixed set of  $n$  predictors and  $\mathbf{Y} \in \mathbb{R}^n$

- Setup is very similar to the bias-variance decomposition, but  $X$  is now fixed
- Let  $\mathbf{Y}_t \in \mathbb{R}^n$  a vector of  $n$  independent observations and  $\hat{f}_{\mathbf{Y}_t}$  an estimate on the training set  $(X, \mathbf{Y}_t)$ , then [Efron, 1986]

$$\underbrace{\mathbb{E}_{\mathbf{Y}, \mathbf{Y}_t} \left\| \mathbf{Y} - \hat{f}_{\mathbf{Y}_t}(X) \right\|_2^2}_{\text{expected prediction error}} = \underbrace{\mathbb{E}_{\mathbf{Y}_t} \left\| \mathbf{Y}_t - \hat{f}_{\mathbf{Y}_t}(X) \right\|_2^2}_{\text{expected training error}} + 2\sigma^2 \text{df}(\hat{f})$$



- This motivates the following model selection criterium [Mallows, 2000]

$$\underbrace{\|y_t - \hat{f}_{y_t}(X)\|_2^2}_{\text{training error}} + \underbrace{2\sigma^2 df(\hat{f})}_{\text{complexity penalty}}$$

- The more complex a model, the larger the penalty
- If two models fit the data equally well, we select the simpler one (Occam's razor)

# MODEL SELECTION APPROACHES - BAYES APPROACH

- Assume we have a set of models  $(m_i)_i$
- In a probabilistic setting we evaluate the probability of a model  $m_i$  given data  $x$ , i.e. using Bayes theorem

$$\text{pr}(m_i | x) = \frac{\text{pr}(x | m_i)\text{pr}(m_i)}{\sum_j \text{pr}(x | m_j)\text{pr}(m_j)} = \frac{\text{pr}(x | m_i)\text{pr}(m_i)}{\text{pr}(x)}$$

- We compare two models  $m_i$  and  $m_j$  using

$$\frac{\text{pr}(m_i | x)}{\text{pr}(m_j | x)} = \frac{\frac{\text{pr}(x | m_i)\text{pr}(m_i)}{\text{pr}(x)}}{\frac{\text{pr}(x | m_j)\text{pr}(m_j)}{\text{pr}(x)}} = \frac{\text{pr}(x | m_i)\text{pr}(m_i)}{\text{pr}(x | m_j)\text{pr}(m_j)}$$

because  $\text{pr}(x)$  drops

# MODEL SELECTION APPROACHES - BAYES FACTOR

- With a uniform prior over models we arrive at the Bayes factor [Kass and Raftery, 1995]

$$\frac{\text{pr}(x | m_i)}{\text{pr}(x | m_j)}$$

- Hence, in Bayesian model selection, we evaluate a model  $m$  based on its *marginal likelihood*

$$\text{pr}(x | m) = \int_{\theta} \text{pr}(x | \theta, m) \text{pr}(\theta | m) d\theta$$

where  $\theta$  are the model parameters

- The marginal likelihood is often difficult to evaluate, even numerically!

## MODEL SELECTION APPROACHES - BIC

- The marginal likelihood is tractable only for very simple models
- As an alternative, we use approximations of the marginal likelihood
- The **Bayes information criterion (BIC)** is such an approximation. Let  $x$  contain  $n$  samples and assume that  $n \gg p$ , then

$$\text{pr}(x | m) \approx \exp \left\{ -\frac{1}{2} \text{BIC}(x; m) \right\}$$
$$\text{BIC}(x; m) = -2 \log \text{pr}(x | \hat{\theta}, m) + p \log(n)$$

where  $\hat{\theta}$  refers to the maximum likelihood estimate and  $p$  to the number of parameters

# MODEL SELECTION APPROACHES - BIC

- Let  $\mathbf{Y}$  and  $\epsilon$  be two random variables such that  $\mathbf{Y} = f(X) + \epsilon$
- Let  $f_{\hat{\theta}}$  denote a maximum likelihood estimate on some training data
- For  $\epsilon \sim \text{Normal}(0, \sigma^2)$  the BIC is related to the mean squared error with complexity penalty

$$\begin{aligned} \text{BIC}(\mathbf{x}; \mathbf{m}) &= \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - f_{\hat{\theta}}(x_i))^2 + p \log(n) + C_n \\ &\propto \frac{1}{\sigma^2} \|\mathbf{y} - f_{\hat{\theta}}(\mathbf{x})\|_2^2 + p \log(n) \end{aligned}$$

where  $C_n$  is a constant depending on  $n$ , which can be dropped for model comparison

# MODEL SELECTION APPROACHES - FIC

- BIC assumes  $n \gg p$  and therefore depends only on the number of parameters
- Fisher Information Approximation (FIA) [Ly et al., 2017]:

$$\begin{aligned} \text{pr}(\mathbf{x} | m) &\approx \exp \{-\text{FIA}(\mathbf{x}; m)\} \\ \text{FIA}(\mathbf{x}; m) &= \underbrace{-\log \text{pr}(\mathbf{x} | \hat{\theta}, m) + \frac{p}{2} \log \left( \frac{n}{2\pi} \right)}_{\text{BIC like term}} + \log C_m \\ C_m &= \underbrace{\int_{\theta} \sqrt{\det \mathcal{I}_m(\theta)} d\theta}_{\text{Geometric complexity}} \end{aligned}$$

where  $\mathcal{I}_m$  denotes the *Fisher information matrix*

- $C_m$  is essential if  $n \gg p$  is not given [Cheema and Sugiyama, 2020]

# HOW DO WE CONTROL MODEL COMPLEXITY?

- Regularization (e.g. ridge regression):
  - ▶ Constrain the feasible set of parameter values
  - ▶ Keep the number of parameters in the model constant, but allow them to become zero
- Number of parameters:
  - ▶ A good approximation of model complexity if  $n < p$
  - ▶ For  $n > p$  we saw that the optimization problem has many solutions
    - In deep neural networks, the gradient descent method can act similar to a regularizer
    - Model complexity can decrease when adding more parameters (double descent)

# REGULARIZATION



# $l_k$ -PENALIZED REGRESSION

Objective function

$$\omega(\theta) = -\log \text{pr}_\theta(\mathbf{y}) \quad (\text{maximum likelihood}), \text{ or}$$

$$\omega(\theta) = \|\mathbf{y} - \mathbf{X}\theta\|_2^2 \quad (\text{linear regression})$$

Regularized estimate with  $l_k$ -norm penalty

$$\hat{\theta} = \begin{cases} \arg \min_{\theta} & \omega(\theta) \\ \text{subject to} & \|\theta\|_k^k = \Lambda \end{cases}$$

where

$$\|\theta\|_k = \left( \sum_{j=2}^p |\theta_j|^k \right)^{1/k}$$

---

<sup>2</sup>Remember that we do not regularize the bias or y-intercept  $\theta_0$

# $l_k$ -PENALIZED REGRESSION

Identify saddle points of Lagrangian

$$\mathcal{L}(\theta, \lambda) = \omega(\theta) + \lambda(\|\theta\|_k^k - \Lambda)$$

In practice, we do not work with  $\Lambda$ , but set  $\lambda$  such that the classification performance is optimal, i.e. we work with the Lagrangian

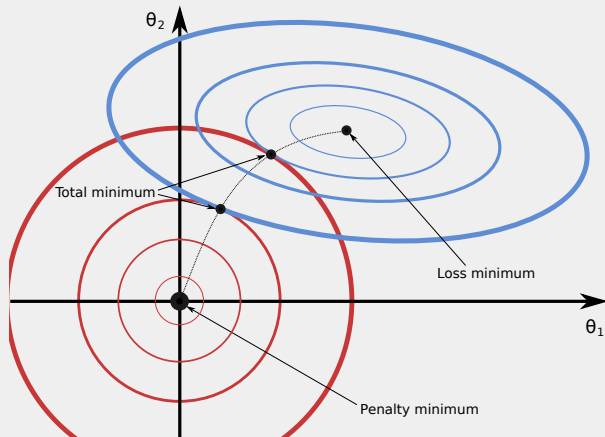
$$\hat{\theta}(\lambda) = \arg \min_{\theta} \omega(\theta) + \lambda \|\theta\|_k^k$$

At the optimum we must have

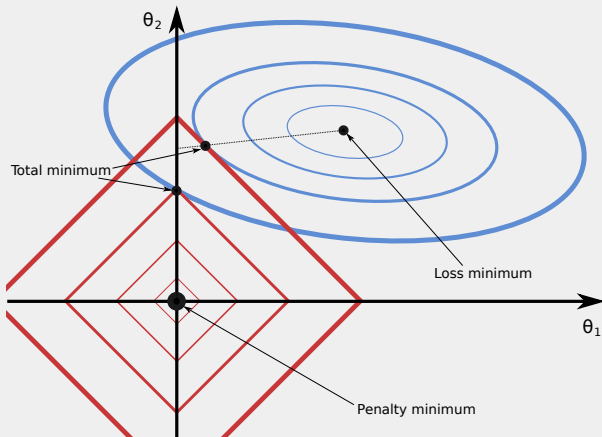
$$\nabla_{\theta} \omega(\theta) + \lambda \nabla_{\theta} \|\theta\|_k^k = \mathbf{0}$$

i.e. the gradients of  $\omega(\theta)$  and  $\lambda \|\theta\|_k^k$  must point to opposite directions

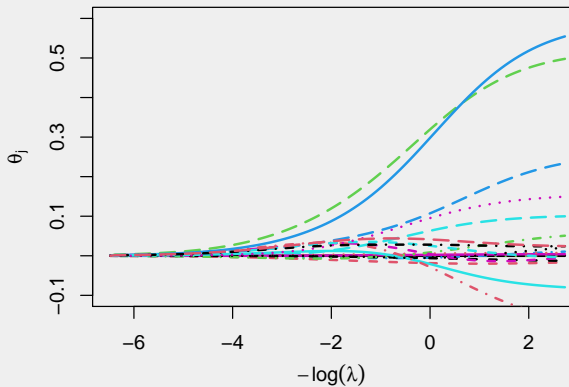
# REGULARIZATION - K=2



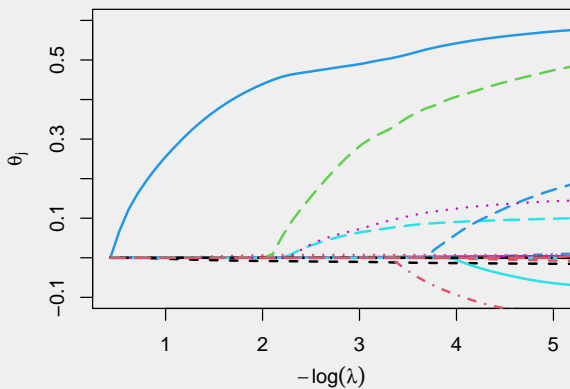
# REGULARIZATION - K=1



# REGULARIZATION PATHS - $K=2$

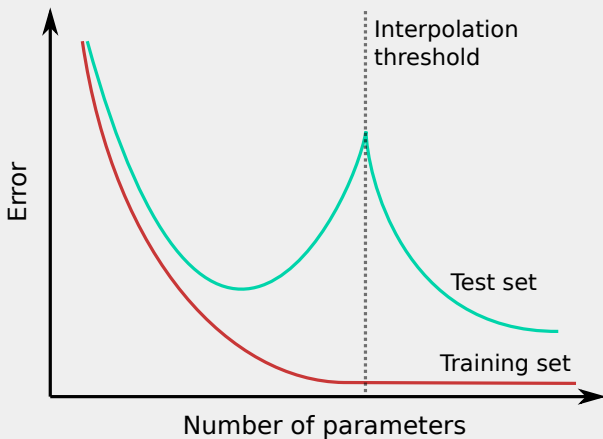


# REGULARIZATION PATHS - $K=1$



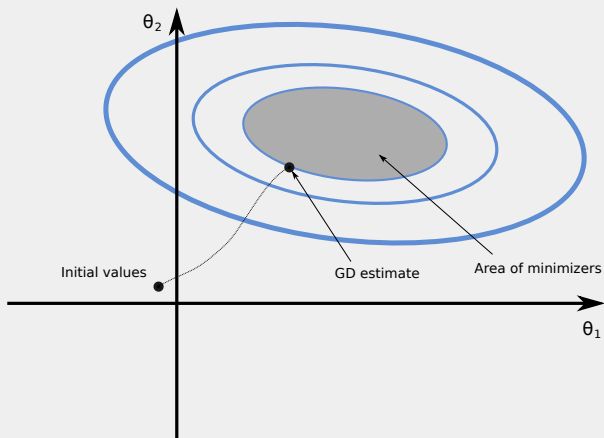
# **IMPLICIT REGULARIZATION AND DOUBLE DESCENT**

# IMPLICIT REGULARIZATION - DOUBLE DESCENT

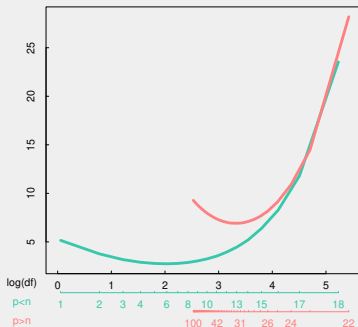
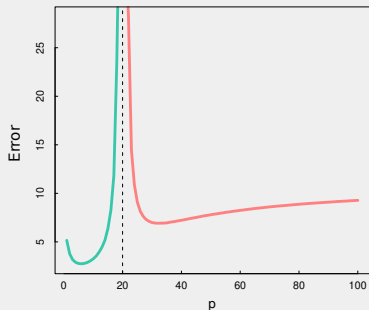




# IMPLICIT REGULARIZATION - DOUBLE DESCENT

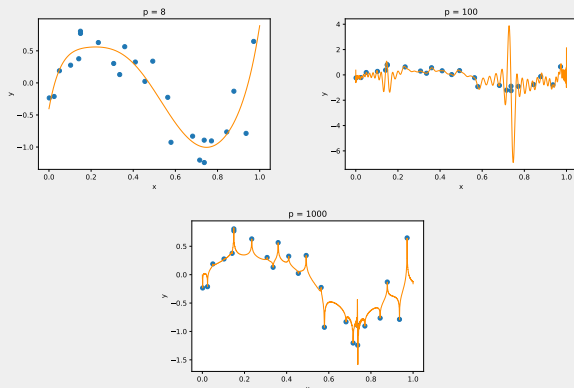


# MINIMUM $\ell_2$ -NORM ESTIMATE - DF



<sup>2</sup>Requires a more advanced definition of DF that treats  $X$  as random variable [Luan et al., 2021]

# IMPLICIT REGULARIZATION



**Figure:** Fitting degree  $d = p - 1$  Legendre polynomials. For  $p > n$  the solution with the smallest  $\ell_2$ -norm is used.

---

<sup>2</sup>Legendre polynomials are quite useful, since their absolute value is bounded by one.

## TAKE HOME MESSAGES

- Expected performance is the sum of training performance and model complexity
- Complex models require regularization to prevent overfitting
- The number of parameters does not correspond to the complexity of a model
- Increasing the number of features can reduce model complexity if a min- $\ell_2$ -norm estimator is used
- If we have complex data and cannot make any assumptions on the generating process, we might be better off with an overparametrized model using regularization (success behind deep learning)

## MORE REFERENCES





- Akaike information criterion (AIC)  
[Akaike, 1974, Cavanaugh and Neath, 2019]
- Bayesian information criterion (BIC) [Schwarz, 1978]
- Deviance information criterion (DIC)  
[Spiegelhalter et al., 2002]
- Fisher Information Approximation (FIA) [Rissanen, 1996, Grünwald, 2007, Cheema and Sugiyama, 2020]
- Degrees of freedom (DF)  
[Tibshirani, 2015, Gao and Jojic, 2016, Luan et al., 2021]
- Implicit regularization and double descent  
[Hastie et al., 2022, Luan et al., 2021, Derezhinski et al., 2020, Kobak et al., 2020]

- Sections 3.4, 7.3, 7.6, 7.7 and 7.9 [Hastie et al., 2009]

*"All models are wrong, but some are useful."*





*[Moody, 1991]*

# REFERENCES I





-  AKAIKE, H. (1974).  
**A NEW LOOK AT THE STATISTICAL MODEL IDENTIFICATION.**  
*IEEE transactions on automatic control*, 19(6):716–723.
-  CAVANAUGH, J. E. AND NEATH, A. A. (2019).  
**THE AKAIKE INFORMATION CRITERION: BACKGROUND, DERIVATION, PROPERTIES, APPLICATION, INTERPRETATION, AND REFINEMENTS.**  
*Wiley Interdisciplinary Reviews: Computational Statistics*, 11(3):e1460.
-  CHEEMA, P. AND SUGIYAMA, M. (2020).  
**DOUBLE DESCENT RISK AND VOLUME SATURATION EFFECTS: A GEOMETRIC PERSPECTIVE.**  
*arXiv preprint arXiv:2006.04366*.
-  COOK, R. D. (1977).  
**DETECTION OF INFLUENTIAL OBSERVATION IN LINEAR REGRESSION.**  
*Technometrics*, 19(1):15–18.







## REFERENCES II

-  DEREZINSKI, M., LIANG, F. T., AND MAHONEY, M. W. (2020).  
**EXACT EXPRESSIONS FOR DOUBLE DESCENT AND IMPLICIT REGULARIZATION VIA SURROGATE RANDOM DESIGN.**  
*Advances in neural information processing systems*, 33:5152–5164.
-  EFRON, B. (1986).  
**HOW BIASED IS THE APPARENT ERROR RATE OF A PREDICTION RULE?**  
*Journal of the American statistical Association*, 81(394):461–470.
-  GAO, T. AND JOJIC, V. (2016).  
**DEGREES OF FREEDOM IN DEEP NEURAL NETWORKS.**  
*arXiv preprint arXiv:1603.09260*.
-  GRÜNWARD, P. D. (2007).  
**THE MINIMUM DESCRIPTION LENGTH PRINCIPLE.**  
MIT press.


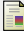
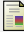
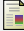
## REFERENCES III

-  HASTIE, T., MONTANARI, A., ROSSET, S., AND TIBSHIRANI, R. J. (2022).  
**SURPRISES IN HIGH-DIMENSIONAL RIDGELESS LEAST SQUARES INTERPOLATION.**  
*The Annals of Statistics*, 50(2):949–986.
-  HASTIE, T., TIBSHIRANI, R., AND FRIEDMAN, J. (2009).  
**THE ELEMENTS OF STATISTICAL LEARNING: DATA MINING, INFERENCE, AND PREDICTION.**  
Springer Science & Business Media.
-  JANSON, L., FITHIAN, W., AND HASTIE, T. J. (2015).  
**EFFECTIVE DEGREES OF FREEDOM: A FLAWED METAPHOR.**  
*Biometrika*, 102(2):479–485.
-  KASS, R. E. AND RAFTERY, A. E. (1995).  
**BAYES FACTORS.**  
*Journal of the american statistical association*, 90(430):773–795.




## REFERENCES IV

-  KOBAK, D., LOMOND, J., AND SANCHEZ, B. (2020).  
**THE OPTIMAL RIDGE PENALTY FOR REAL-WORLD HIGH-DIMENSIONAL DATA CAN BE ZERO OR NEGATIVE DUE TO THE IMPLICIT RIDGE REGULARIZATION.**  
*J. Mach. Learn. Res.*, 21:169–1.
-  LUAN, B., LEE, Y., AND ZHU, Y. (2021).  
**PREDICTIVE MODEL DEGREES OF FREEDOM IN LINEAR REGRESSION.**  
*arXiv preprint arXiv:2106.15682*.
-  LY, A., MARSMAN, M., VERHAGEN, J., GRASMAN, R. P., AND WAGENMAKERS, E.-J. (2017).  
**A TUTORIAL ON FISHER INFORMATION.**  
*Journal of Mathematical Psychology*, 80:40–55.
-  MALLOW, C. L. (2000).  
**SOME COMMENTS ON CP.**  
*Technometrics*, 42(1):87–94.

# REFERENCES V

-  MOODY, J. (1991).  
**THE EFFECTIVE NUMBER OF PARAMETERS: AN ANALYSIS OF GENERALIZATION AND REGULARIZATION IN NONLINEAR LEARNING SYSTEMS.**  
*Advances in neural information processing systems*, 4.
-  RISSANEN, J. J. (1996).  
**FISHER INFORMATION AND STOCHASTIC COMPLEXITY.**  
*IEEE transactions on information theory*, 42(1):40-47.
-  SCHWARZ, G. (1978).  
**ESTIMATING THE DIMENSION OF A MODEL.**  
*The annals of statistics*, pages 461-464.
-  SHALEV-SHWARTZ, S. AND BEN-DAVID, S. (2014).  
**UNDERSTANDING MACHINE LEARNING: FROM THEORY TO ALGORITHMS.**  
Cambridge university press.

## REFERENCES VI

-  SPIEGELHALTER, D. J., BEST, N. G., CARLIN, B. P., AND VAN DER LINDE, A. (2002).  
**BAYESIAN MEASURES OF MODEL COMPLEXITY AND FIT.**  
*Journal of the royal statistical society: Series b (statistical methodology)*, 64(4):583–639.
-  TIBSHIRANI, R. J. (2015).  
**DEGREES OF FREEDOM AND MODEL SEARCH.**  
*Statistica Sinica*, pages 1265–1296.
-  WOLPERT, D. H. AND MACREADY, W. G. (1997).  
**NO FREE LUNCH THEOREMS FOR OPTIMIZATION.**  
*IEEE transactions on evolutionary computation*, 1(1):67–82.