

# Overfitting, regularization, and complexity

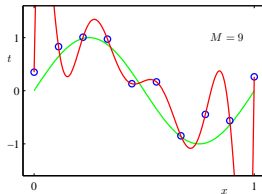
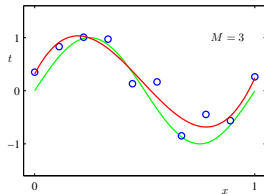
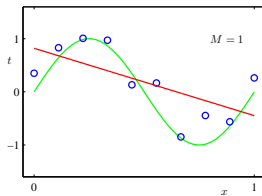
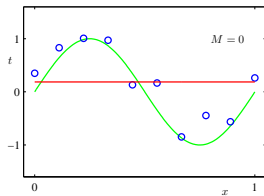
MATH-412 - Statistical Machine Learning

# Polynomial regression and overfitting

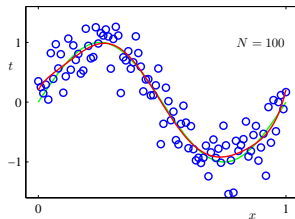
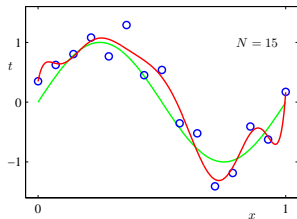
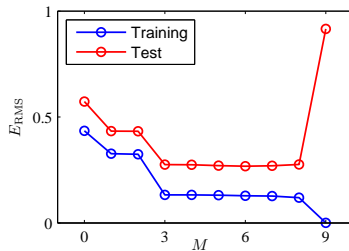
# Polynomial regression : an instance of linear regression

Model of the form  $Y = w_0 + w_1X + w_2X^2 + \dots + w_pX^p + \varepsilon$

$$\min_{\mathbf{w}} \frac{1}{2n} \sum_{i=1}^n (y_i - (w_0 + w_1x_i + w_2x_i^2 + \dots + w_px_i^p))^2 \quad \text{with } p = M$$



# Overfitting : symptoms and characteristics



# Overfitting and generalization

In ML/stats we care about the **generalization ability** of the predictor

*Fitting perfectly the data* does not always entail *lack of generalization*...

- e.g., deep neural networks in computer vision.

But fitting perfectly the data is a problem if

- the data is **noisy** and the model “fits the noise” or
- to be able to fit the data training the model learned is too “**complex**”.

How do we measure **complexity**?...

# Regularization

# Tikhonov regularization



Andrey N. Tikhonov  
(1906 - 1993)

$$\min_{f \in S} \hat{\mathcal{R}}_n(f) + \lambda \|f\|^2$$

- $\lambda$  is the regularization coefficient or hyperparameter

Is the problem now well-posed ?

If  $\hat{\mathcal{R}}_n$  is convex

- ⇒ The objective is strongly convex and coercive for any  $\lambda > 0$
- ⇒ The solution exists and is unique.
- ⇒  $\lambda \mapsto \hat{f}_\lambda$  is a continuous function

If  $\hat{\mathcal{R}}_n$  is bounded below

- ⇒ The objective is coercive for any  $\lambda > 0$
- ⇒ At least a solution exists

If  $\hat{\mathcal{R}}_n$  is  $\mathcal{C}^2$  with bounded curvature

- ⇒ Regularization eliminates small local minima.

## Ridge regression

Is obtained by applying Tikhonov regularization to OLS regression.

$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Normal equation

$$\left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right) \mathbf{w} = \frac{1}{n} \mathbf{X}^\top \mathbf{y}$$

- Thus with unique solution :

$$\hat{\mathbf{w}}^{(\text{ridge})} = \frac{1}{n} \left( \frac{1}{n} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Shrinkage effect



## Linear vs affine regression and regularization

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \quad \text{vs} \quad f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b = \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}$$

With

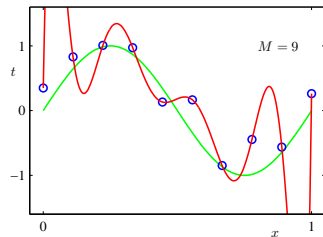
$$\tilde{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

- ... an affine model in dimension  $p$  is a linear model in dimension  $p + 1$
- Working with  $(\mathbf{w}, b)$  vs  $\tilde{\mathbf{w}}$  is equivalent **when we don't regularize** and otherwise not, because usually  $b$  is not regularized :

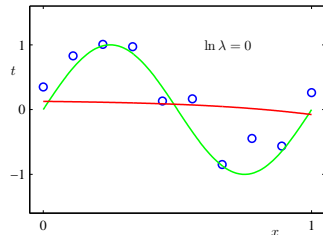
$$\min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w} + b\mathbf{1}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

# Polynomial regression with ridge

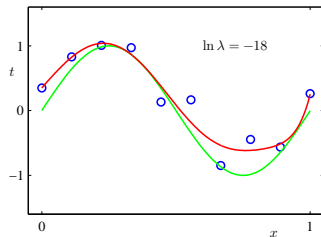
No regularization



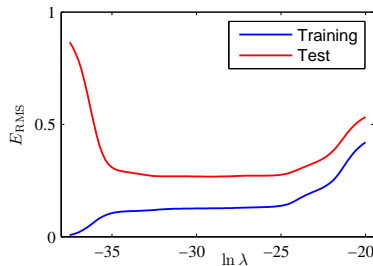
$\lambda = 1$



$10^{-8} \leq \lambda \leq 10^{-7}$



Regression with :  
 $n = 10$  examples and  
a polynomial of degree  $M = 9$ ,  
*but with ridge regularization.*



# Complexity

# Controlling the complexity of the hypothesis space

## Explicit control

- number of variables
- max degree for polynomial functions
- degree and # of knots for spline functions
- max resolution in wavelet approximations.
- bandwidth in RKHS

## Implicit control

- with regularization,
- using Bayesian formulations
- via the learning/optimization algorithm
- randomization
- ...

The complexity of the predictor often results from a compromise between fitting and increasing complexity.

**Problem of model selection** : How to choose the right level of complexity ?

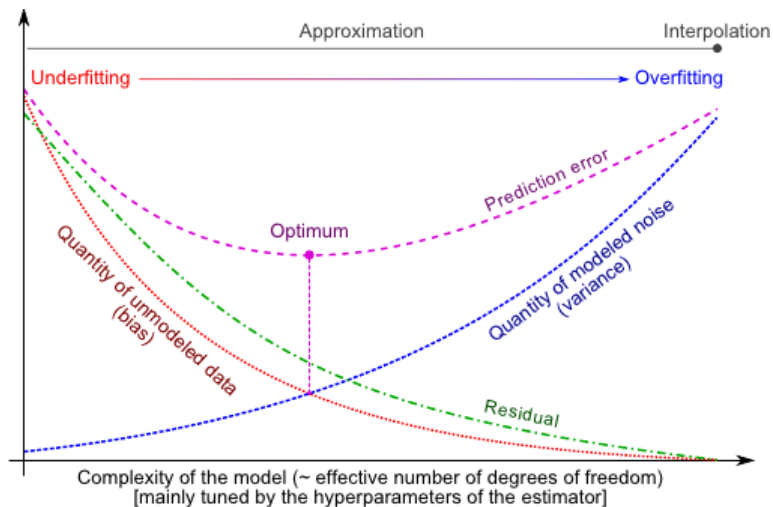
## Risk decomposition : approximation-estimation trade-off

- $f^* = \text{target function}$
- $f_S^* = \operatorname{argmin}_{f \in S} \mathcal{R}(f)$
- $\hat{f}_S = \text{predictor/estimator in } S$

$$\underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*)}_{\text{excess risk}} = \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{estimation error}} + \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{approximation error}}$$

- Sometimes also called “bias-variance tradeoff”

# Approximation-estimation tradeoff



The view that there is a necessarily a compromise between fitting well the training data and learning a too complex model to generalize has been challenged by neural networks...