

# Bias-variance decomposition

## Linear models for regression

Francesco Corona

# **Bias-variance decomposition**

## **Linear models for regression**

## Bias-variance decomposition

The use of maximum likelihood, or least squares, can lead to severe over-fitting

- ▶ if complex models are trained using data sets of limited size

Limiting the number of basis functions to avoid over-fitting has the side effect of limiting the flexibility of the model to capture interesting trends in the data

Regularisation terms can control over-fitting for models with many parameters

- ▶ How to determine a suitable value for the regularisation coefficient  $\lambda$ ?

Seeking the solution that minimises the regularised error function with respect to both the weight vector  $\mathbf{w}$  and the regularisation coefficient  $\lambda$  is clearly not the right approach since this leads to the unregularised solution with  $\lambda = 0$

## Bias-variance decomposition (cont.)

The over-fitting phenomenon is an unfortunate property of maximum likelihood

- ▶ It does not arise when we marginalise over parameters in a Bayesian setting

It is instructive to first consider a frequentist viewpoint of model complexity

- ▶ **bias-variance trade-off**

We introduce the concept only in the context of linear basis function models

## Bias-variance decomposition (cont.)

When we discussed decision theory for regression problems, the decision stage consists of choosing a specific estimate  $y(\mathbf{x})$  of the target  $t$  for each input  $\mathbf{x}$

We can do this using a loss  $L(t, y(\mathbf{x}))$ , so that the average/expected loss is

$$\mathbb{E}[L] = \int \int L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$$

Various loss functions for regression lead to a corresponding optimal prediction

- ▶ once we are given the conditional density  $p(t|\mathbf{x})$

## Bias-variance decomposition (cont.)

A common loss function in regression problems is the **squared loss function**

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2 \quad \implies \quad \mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dx dt$$

Squared loss function (decision theory)  $\neq$  sum-of-squares error function (ML)

Squared loss function

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2$$

- ▶ Optimal prediction  $h(\mathbf{x})$  is given by the conditional expectation  $\mathbb{E}[t|\mathbf{x}]$

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x})dt \quad (1)$$

## Bias-variance decomposition (cont.)

We also obtained: 
$$\mathbb{E}[L] = \int (y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int (\mathbb{E}[t|\mathbf{x}] - t)^2 p(\mathbf{x}) d\mathbf{x}$$

It is minimised when  $y(\mathbf{x})$ , in the first term, equals  $\mathbb{E}[t|\mathbf{x}]$

The second term is independent of  $y(\mathbf{x})$ , arises from the noise  $\varepsilon$

- ▶ The variance of the distribution of  $t$ , averaged over  $\mathbf{x}$
- ▶ It is the intrinsic variability of the target variable
- ▶ The minimum achievable value of the expected loss

## Bias-variance decomposition (cont.)

The expected squared loss function can be written also in another form

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dx dt \quad (2)$$

With an infinite supply of data and unlimited computational resources

- ▶ we could find the regression function  $h(\mathbf{x})$  to any accuracy

In practice, we only have a data set  $\mathcal{D}$  with a finite number  $N$  of points

- ▶  $h(\mathbf{x})$  is not known exactly



## Bias-variance decomposition (cont.)

If we model  $h(\mathbf{x})$  using a parametric function  $y(\mathbf{x}, \mathbf{w})$  with parameter vector  $\mathbf{w}$

- ▶ the uncertainty in our model is expressed through a posterior distribution over  $\mathbf{w}$  (Bayesian perspective)

A frequentist treatment makes a point estimate of  $\mathbf{w}$  based on the data set  $\mathcal{D}$

- ▶ the uncertainty of this estimate is expressed through a large number of data sets each of size  $N$  and each drawn independently from distribution  $p(t, \mathbf{x})$

For any set  $\mathcal{D}$ , we learn our algorithm and get a prediction function  $y(\mathbf{x}; \mathcal{D})$

- ▶ Different data sets, different functions
- ▶ Different functions, different values of the squared loss

The performance of a learning algorithm is assessed by averaging over sets

## Bias-variance decomposition (cont.)

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dx dt$$

Consider the integrand of the first term of the expected squared loss, it becomes

$$(y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}))^2 \quad (3)$$

for a particular data set  $\mathcal{D}$  and it has to be averaged over the ensemble of sets

Before taking its expectation wrt  $\mathcal{D}$ , add and subtract the quantity  $\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]$

$$(y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}))^2$$

## Bias-variance decomposition (cont.)

Expanding, we obtain

$$\begin{aligned} & \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 \\ &= \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] \right)^2 + \left( \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 \\ & \quad + 2 \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] \right) \left( \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right) \quad (4) \end{aligned}$$

And taking the expectation with respect to  $\mathcal{D}$ , it gives

$$\mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right)^2 \right] = \left( \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 + \mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] \right)^2 \right]$$

## Bias-variance decomposition (cont.)

$$\mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right)^2 \right] = \underbrace{\left( \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] \right)^2 \right]}_{\text{variance}} \quad (5)$$

The expected squared difference between  $y(\mathbf{x}; \mathcal{D})$  and the regression function  $h(\mathbf{x})$  can be expressed as the sum of two terms

- ▶ The first term, squared **bias**, represents the extent to which the average prediction over all data sets differs from the desired regression function
- ▶ The second term, **variance**, measures the extent to which the solutions for individual data sets vary around their average, and hence measures the extent to which function  $y(\mathbf{x}; \mathcal{D})$  is sensitive to the particular data set

We shall provide some intuition to support these definitions

## Bias-variance decomposition (cont.)

$$\mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right)^2 \right] = \left( \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 + \mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] \right)^2 \right]$$

Expected squared difference between  $y(\mathbf{x}; \mathcal{D})$  and the regression function  $h(\mathbf{x})$

- ▶ when considering only a single input value  $\mathbf{x}$

Substituting in  $\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$

$$\text{expected loss} = (\text{BIAS})^2 + \text{VARIANCE} + \text{noise} \quad (6)$$

$$(\text{BIAS})^2 = \int \left( \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \quad (7)$$

$$\text{VARIANCE} = \int \mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}; \mathcal{D})] \right)^2 \right] p(\mathbf{x}) d\mathbf{x} \quad (8)$$

$$\text{noise} = \int \int (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (9)$$

## Bias-variance decomposition (cont.)

We decomposed the expected loss into (integrated) bias, (integrated) variance and a constant noise term, but our goal is the same: We want to minimise it

There is a trade-off between bias and variance:

- ▶ flexible models will have low bias and high variance
- ▶ rigid models will have high bias and low variance

$$(\text{BIAS})^2 = \int \left( \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{VARIANCE} = \int \mathbb{E}_{\mathcal{D}} \left[ \left( y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] \right)^2 \right] p(\mathbf{x}) d\mathbf{x}$$

The model with optimal predictive capability is the one with the best balance

## Bias-variance decomposition (cont.)

As an example, we consider the usual data from a sinusoidal function

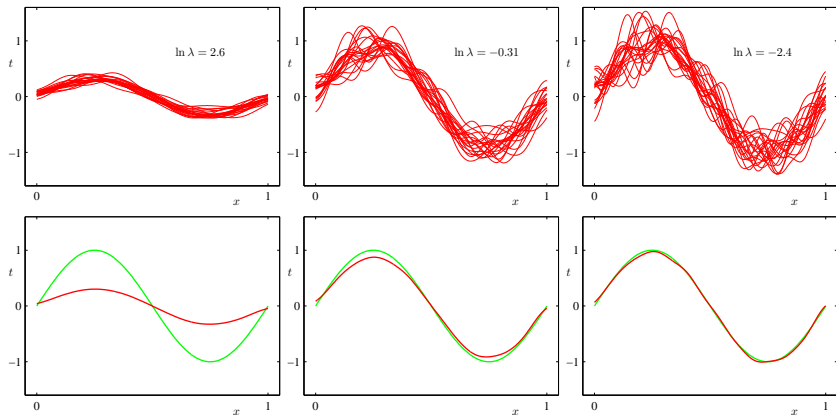
- ▶  $l = 1, \dots, L$  datasets  $\mathcal{D}^{(l)}$ , each with  $N = 25$  points,  $L = 100$
- ▶ The points of each  $\mathcal{D}^{(l)}$  are iid from  $h(x) = \sin(2\pi x)$

For each  $\mathcal{D}^{(l)}$ , we fit a model with 24 Gaussian basis ( $M = 25$  parameters)

- ▶ We minimised the regularised error  $\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$
- ▶ The resulting parameter vector is  $\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$
- ▶ We use  $\mathbf{w}^{(l)}$  to get a predictive function  $y^{(l)}$

All this, for different values of the regularisation parameter  $\lambda$

## Bias-variance decomposition (cont.)



- ▶ Large  $\lambda$  (left), low variance but high bias
- ▶ Small  $\lambda$  (right), low bias but high variance



## Bias-variance decomposition (cont.)

In this case, averaging many solutions turned out to be a beneficial procedure

$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x) \quad \rightsquigarrow \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] \quad (10)$$

The integrated<sup>1</sup> squared bias and the integrated variance are given by

$$\begin{aligned} (\text{BIAS})^2 &= \frac{1}{N} \sum_{n=1}^N \left( \bar{y}(x_n) - h(x_n) \right)^2 \\ &\rightsquigarrow \int \left( \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x) \right)^2 p(x) dx \quad (11) \end{aligned}$$

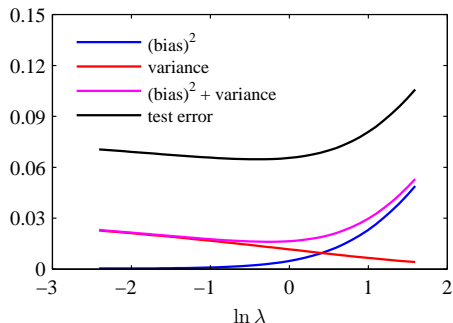
$$\begin{aligned} \text{VARIANCE} &= \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \left( y^{(l)}(x_n) - \bar{y}(x_n) \right)^2 \\ &\rightsquigarrow \int \mathbb{E}_{\mathcal{D}} \left[ \left( y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] \right)^2 \right] p(x) dx \quad (12) \end{aligned}$$

<sup>1</sup>Integration over  $x$  weighted by the distribution  $p(x)$  is approximated by a finite sum over points draw from that distribution

## Bias-variance decomposition (cont.)

Plot of squared bias and variance, together with their sum

- ▶ Also shown is the average test set error for a test set size of 1000 points



The minimum of  $(\text{BIAS})^2 + \text{VARIANCE}$  occurs around a value  $\ln \lambda = -0.31$

It is close to the value that gives the minimum error on the test data