

Machine Learning for Applied Physics and High Energy Physics ⁴

Lecture 3 The bias-variance trade-off and error representation

3.1 Bias-variance trade-off.

In this lecture we slip into the central principle that underlies much of machine learning: the bias-variance trade-off. We will discuss it in the context of continuous predictions, such as those involved in a regression problem.

Let us consider a data set

$$\mathcal{D} = \left\{ \bar{x}_n, \bar{y}_n \right\} \quad n = 1, \dots, N$$

consisting of N pairs of features and labels. Let us assume that the true data is generated from a noisy model

$$\bar{y} = f(\bar{x}) + \bar{\epsilon}$$

where ϵ is normally distributed with mean zero and standard deviation σ_ϵ . Let us suppose that we have a statistical procedure (e.g. least square regression) for forming a predictor $f(\bar{x}, \hat{\bar{\theta}})$ that gives the prediction of our model for a new data point x . This estimator is chosen by minimizing a cost function

$$\mathcal{L}_2(\bar{y}, f(\bar{x}, \bar{\theta})) = \frac{1}{N} \sum_{n=1}^N (\bar{y}_n - f(\bar{x}_n, \bar{\theta}))^2 \equiv \text{MSE}(\bar{\theta})$$

in such a way that

$$\hat{\bar{\theta}} = \underset{\bar{\theta}}{\text{argmin}} \text{MSE}(\bar{\theta})$$

The best-fit parameters $\hat{\bar{\theta}}$ are therefore a function of the data set $\mathcal{D} = \{ \bar{x}_n, \bar{y}_n \}$. We would obtain a different $\text{MSE}(\bar{\theta})$

if we had a different data set in a Universe of possible data sets obtained by drawing N samples from the true data distribution. We denote an expectation value over all of these datasets as $\mathbb{E}_{\mathcal{D}}$. We would also like to average over different instances of the "noise" ϵ and we denote the expectation value over the noise by \mathbb{E}_{ϵ} . We can therefore decompose the expected generalisation error as

$$\begin{aligned}\mathbb{E}_{\mathcal{D}, \epsilon} [l_2(\bar{y}, f(\bar{x}, \hat{\theta}_{\mathcal{D}}))] &= \mathbb{E}_{\mathcal{D}, \epsilon} \left[\frac{1}{N} \sum_{n=1}^N (\bar{y}_n - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))^2 \right] \\ &= \mathbb{E}_{\mathcal{D}, \epsilon} \left[\frac{1}{N} \sum_{n=1}^N (\bar{y}_n - f(\bar{x}_n) + f(\bar{x}_n) - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))^2 \right] \\ &= \frac{1}{N} \sum_{n=1}^N \underbrace{\mathbb{E}_{\epsilon} [(\bar{y}_n - f(\bar{x}_n))^2]}_{\sigma_{\epsilon}^2} + \frac{1}{N} \mathbb{E}_{\mathcal{D}, \epsilon} [(f(\bar{x}_n) - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))^2] \\ &\quad + \frac{2}{N} \underbrace{\mathbb{E}_{\epsilon} [(\bar{y}_n - f(\bar{x}_n))]}_0 \mathbb{E}_{\mathcal{D}} [(f(\bar{x}_n) - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))] \end{aligned}$$

0 because we summed a Gaussian noise with mean zero

$$= \frac{1}{N} \sum_{n=1}^N \sigma_{\epsilon}^2 + \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathcal{D}} [(f(\bar{x}_n) - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))^2]$$

We can further decompose the second term as

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} [(f(\bar{x}_n) - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}))^2] &= \mathbb{E}_{\mathcal{D}} [\{f(\bar{x}_n) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\theta}_{\mathcal{D}})] + \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\theta}_{\mathcal{D}})] - f(\bar{x}_n, \hat{\theta}_{\mathcal{D}})\}^2] \\ &= \mathbb{E}_{\mathcal{D}} [\{f(\bar{x}_n) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\theta}_{\mathcal{D}})]\}^2] + \mathbb{E}_{\mathcal{D}} [\{f(\bar{x}_n, \hat{\theta}_{\mathcal{D}}) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\theta}_{\mathcal{D}})]\}^2] \end{aligned}$$

$$+ 2 \mathbb{E}_{\mathcal{D}} \left[\left\{ f(\bar{x}_n) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\} \underbrace{\left\{ f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}}) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\}}_0 \right]^2 \\ = \left\{ f(\bar{x}_n) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\}^2 + \mathbb{E}_{\mathcal{D}} \left[\left\{ f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}}) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\}^2 \right]$$

We call

$$\text{Bias}^2 = \frac{1}{N} \sum_{n=1}^N \left\{ f(\bar{x}_n) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\}^2$$

the measure of the deviation of the expectation value of an estimator from the true value

We call (as usual)

$$\text{Var} = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathcal{D}} \left[\left\{ f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}}) - \mathbb{E}_{\mathcal{D}} [f(\bar{x}_n, \hat{\vec{y}}_{\mathcal{D}})] \right\}^2 \right]$$

the measure of the fluctuations of the estimator due to finite-sample effects.

We call noise

$$\text{Noise} = \frac{1}{N} \sum_{n=1}^N \sigma_{\mathcal{E}}^2$$

We therefore arrive at the following decomposition

$$\mathbb{E}_{\mathcal{D}, \mathcal{E}} [\ell_2(\bar{\vec{y}}, f(\bar{\vec{x}}, \hat{\vec{y}}_{\mathcal{D}}))] = \text{Bias}^2 + \text{Var} + \text{Noise}$$

The minimum of $\mathbb{E}_{\mathcal{D}, \mathcal{E}}$, which is sought by good ML algorithms, does not generally correspond to the minimum of each of its components.

3.2 Overview of Bayesian inference

4

To solve a problem using Bayesian methods, we have to specify two functions

- the likelihood $p(\mathcal{D} / \vec{\theta})$, which describes the probability of observing a data set \mathcal{D} for given values of the model parameters $\vec{\theta}$;
- the prior $p(\vec{\theta})$, which describes any (prior) knowledge about the parameters before we collect the data.

The posterior distribution follows from Bayes' theorem

$$p(\vec{\theta} / \mathcal{D}) = \frac{p(\vec{\theta}) p(\mathcal{D} / \vec{\theta})}{\int d\vec{\theta}' p(\vec{\theta}') p(\mathcal{D} / \vec{\theta}')}$$

Many common statistical procedures can be cast as MAXIMUM LIKELIHOOD ESTIMATION (MLE). In MLE one chooses the best-fit parameters $\hat{\vec{\theta}}$ that maximise the likelihood (or equivalently, the minimise the negative log-likelihood)

$$\vec{\theta}_{MLE} = \underset{\vec{\theta}}{\operatorname{argmax}} \ln p(\mathcal{D} / \vec{\theta}) = \underset{\vec{\theta}}{\operatorname{argmin}} \left\{ \underbrace{-\ln p(\mathcal{D} / \vec{\theta})}_{NLL(\vec{\theta})} \right\}$$

In other words, MLE consists in maximising the probability of seeing the observed data, given a generative model.

The prior distribution is genuinely Bayesian.

- If we do not have any specialised knowledge of $\vec{\theta}$ before seeing the data, we should select an "uninformative" prior
- If we do have knowledge of $\vec{\theta}$ before seeing the data,

we should choose an informative prior.

Using informative priors tends to decrease the variance of the posterior distribution while, potentially, increasing its bias. This is beneficial if the decrease in variance is larger than the increase in bias. A commonly used prior is the Gaussian prior:

$$p(\vec{\theta}/\mathcal{D}) = \prod_j \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} \theta_j^2}$$

which is used to express the assumption that many of the parameters will be small.

We often summarise the posterior probability distribution with a single estimator:

- the mean $\langle \vec{\theta} \rangle = \int d\vec{\theta} \vec{\theta} p(\vec{\theta}/\mathcal{D})$ (Bayes estimate)
- the mode $\hat{\vec{\theta}}_{\text{MAP}} = \arg \max_{\vec{\theta}} p(\vec{\theta}/\mathcal{D})$ (MAXIMUM A POSTERIOR)

The reason being that a probability distribution (which is often multi-dimensional) is an object difficult to manipulate.

The mean or MAP estimators come with a credibility interval which quantifies the uncertainty associated to them. The credibility interval $100(1-\alpha)\%$ is the interval that contains a fraction $1-\alpha$ of the posterior probability.

CREDIBILITY INTERVAL

CONFIDENCE LEVEL

(Bayesian)

(frequentist)

In the expression of the Gaussian prior distribution, λ is a hyperparameter (or nuisance variable). One could

define another prior distribution for λ , usually using an ϵ uninformative prior, and to average the posterior distribution over all choices of λ . This is called a hierarchical prior.