

# Machine Learning: Homework 2

## Overfitting

Ana Mestre

### Problem

We set up an experimental framework to study various aspects of overfitting. The input space is  $X = [-1,1]$  with uniform input probability density,  $P(x) = \frac{1}{2}$ . We consider the two models  $H_2$  and  $H_{10}$ . The target function is a polynomial of degree  $Q_f$ , which we write as  $f(x) = \sum a_q * L_q(x)$ , where  $L_q(x)$  are the Legendre polynomials.

We use the Legendre polynomials because they are a convenient orthogonal basis for the polynomials on  $[-1,1]$ .

The data set is  $D = (x_1, y_1), \dots, (x_N, y_N)$ , where  $y_n = f(x) + \sigma \epsilon_n$  and  $\epsilon_n$  are i.i.d. Standard Normal random variables.

For a single experiment, with specified values for  $Q_f$ ,  $N$ ,  $\sigma$ , generate a random degree- $Q_f$  target function by selecting coefficients  $a_q$  independently from a standard Normal distribution, rescaling them so that  $\text{Ex}[f(x)^2] = 1$ . Generate a data set, selecting  $x_1, \dots, x_N$  independently from  $P(x)$  and  $y_n = f(x) + \sigma \epsilon_n$ . Let  $g_2$  and  $g_{10}$  be the best fit hypotheses to the data from  $H_2$  and  $H_{10}$ , respectively, with respective out-of-sample errors  $E_{\text{out}}(g_2)$  and  $E_{\text{out}}(g_{10})$ .

**(a) Why do we normalize f ? [Hint: how would you interpret  $\sigma$  ?]**

Knowing the following statements:

- Our function is:  $f(x) = \sum a_q * L_q(x)$  and our output is:  $y_n = f(x) + \sigma \epsilon_n$ , from this we can interpret  $\sigma$  as the noise value of the function. The noise increases as the  $\sigma$  increases.
- $\text{Var}[f] = E[f^2] - E[f]^2$
- All Legendre polynomials of degree = 0 have mean 0 on  $[-1,1]$ , which means that the mean of the function  $f(x)$  on  $[-1,1]$  is  $E[f(x)] = a_0$ .
- The variance of  $f(x)$  on  $[-1, 1]$  is given by:

$$\text{Var}(f(x)) = \sum \frac{a_q^2}{2q+1}$$

Taking these statements into consideration and knowing that we want to get a function really close to a normal distribution, we have to normalize the coefficients such that  $E(f^2) = 1$ . The way of doing so is by dividing a by:

$$2\sqrt{\sum \frac{a_q^2}{2q+1}}$$

```
# Get rescaled aq
def norm_aqs(Qf):
    aqs = np.random.standard_normal(Qf + 1)
    summ = sum([aqs[q]**2/(2 * q + 1) for q in range(Qf + 1)])
    return [aqs / math.sqrt(2 * summ)]
```

**(b) How can we obtain g2, g10? [Hint: pose the problem as linear regression]**

We build polynomial models of order 2 and 10. We use them along with the dataset D to fit a linear regression model (H2 and H10, respectively).

```
def polRegression(x, y, Qf, color, label, title="", plot = True):
    x_reshape = x.reshape(-1, 1)

    poly_feat = PolynomialFeatures(degree = Qf, include_bias = True)
    lin_reg = LinearRegression()
    x_ = poly_feat.fit_transform(x_reshape)
    lin_reg.fit(x_, y)
    y_pred = lin_reg.predict(x_)

    if plot:
        plotResults(x, y, y_pred, Qf, label, title, color)

    return x, y_pred, lin_reg, poly_feat
```

**(c) How can we compute E\_out analytically for a given g10?**

We have to calculate:  $E_{\text{out}}(g_{10}) = E[(g_{10}(x) - f(x))^2]$ , which can be analyzed by the bias-variance decomposition.

Finding an  $f$  that generalizes to points outside of the training set can be done with any of the countless algorithms used for supervised learning. It turns out that whichever function  $f$  we select, we can decompose its expected error on an unseen sample  $x$  as follows:

$$E[(y - \hat{f}(x))^2] = (\text{Bias}[\hat{f}(x)])^2 + \text{Var}[\hat{f}(x)] + \sigma^2$$

where

$$\text{Bias}[\hat{f}(x)] = E[\hat{f}(x)] - E[f(x)]$$

and

$$\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2.$$

(d) Vary  $Q_f$ ,  $N$ ,  $\sigma$  and for each combination of parameters, run a large number of experiments, each time computing  $E_{\text{out}}(g_2)$  and  $E_{\text{out}}(g_{10})$ . Averaging these out-of-sample errors estimates of the expected out-of-sample error for the given learning scenario  $(Q_f, N, \sigma)$  using  $H_2$  and  $H_{10}$ . Let:

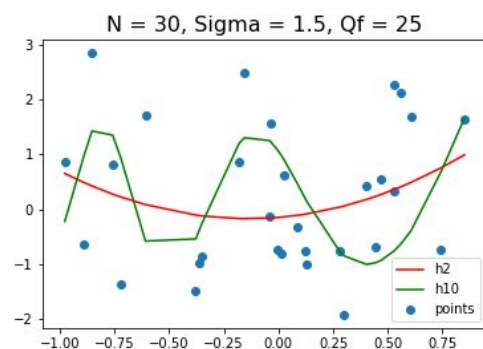
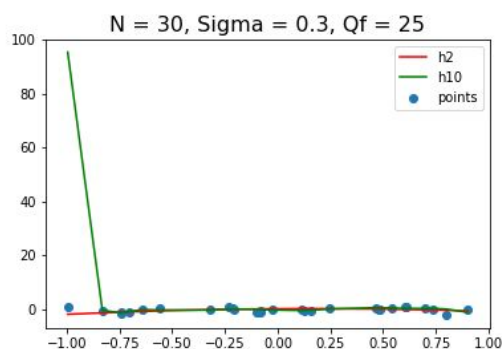
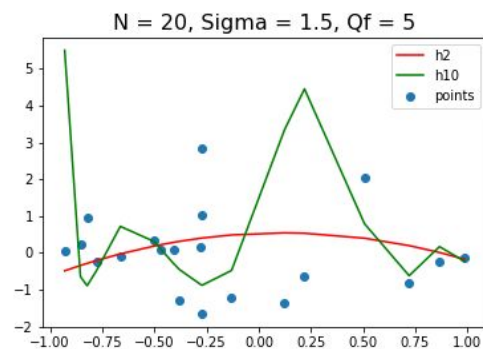
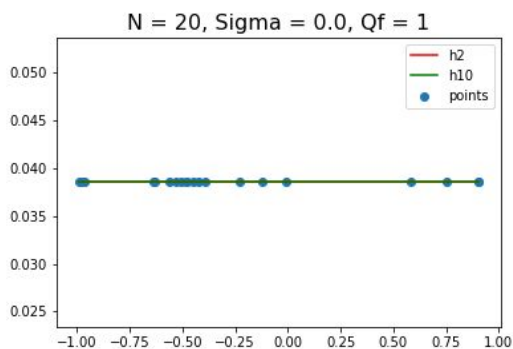
$$E_{\text{out}}(H_2) = \text{average over experiments}(E_{\text{out}}(g_2)),$$

$$E_{\text{out}}(H_{10}) = \text{average over experiments}(E_{\text{out}}(g_{10})).$$

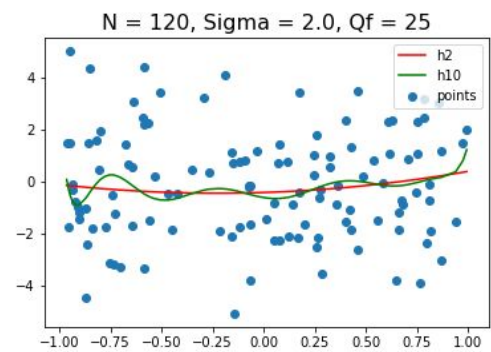
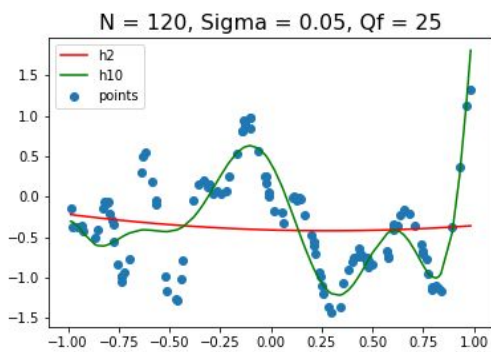
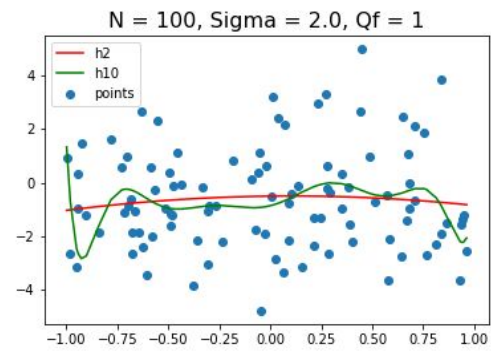
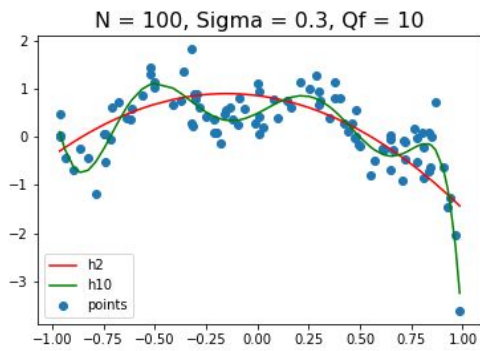
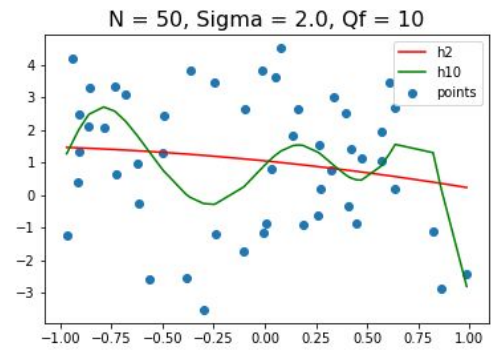
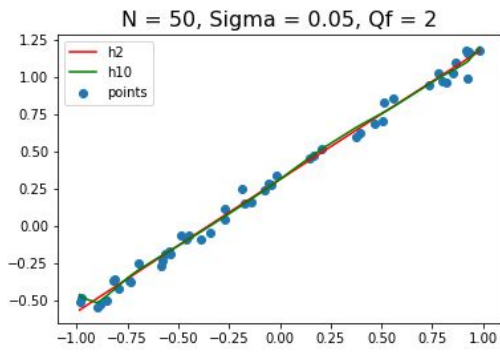
Define the overfit measure  $E_{\text{out}}(H_{10}) - E_{\text{out}}(H_2)$ . When is the overfit measure significantly positive (i.e. overfitting is serious) as opposed to significantly negative?

Try the choices  $Q_f \in \{1, 2, \dots, 100\}$ ,  $N \in \{20, 25, \dots, 120\}$ ,  $\sigma^2 \in \{0, 0.05, 0.1, \dots, 2\}$ . Explain your observations.

Here are some results (find more attached to the assignment's submission):



## Machine learning: overfitting homework



Example of results for N = 30 (more to be found attached on this assignment's submission)

N	Qf	Sigma	Mean E2	Mean E10	Overfit
30	1	0.0	0.0	0.0	0.0
30	1	0.05	0.0022	0.0026	0.0004
30	1	0.3	0.0718	0.6541	0.5823
30	1	1.0	1.1326	148.4168	147.2842
30	1	1.5	1968	56.4844	54.5164
30	1	2.0	4.3613	283.6016	279.2403
30	2	0.0	0.0	0.0	0.0
30	2	0.05	2	0.0087	0.0067
30	2	0.3	0.0863	0.5751	0.4888
30	2	1.0	1.0646	21.2056	20141
30	2	1.5	2.0266	190137	188.1104
30	2	2.0	2.1119	8.4165	6.3046
30	5	0.0	158	0.0	-158
30	5	0.05	0.0999	0.0058	-0.0941
30	5	0.3	229	0.7638	0.5348
30	5	1.0	1.0857	10.9593	9.8736
30	5	1.5	2327	10.9999	8.6729
30	5	2.0	3884	5.2269	1.3429
30	10	0.0	0.8404	0.0	-0.8404
30	10	0.05	0.2276	0.0323	-0.1953
30	10	0.3	0.2623	0.0991	-0.1632
30	10	1.0	1.0807	180.9269	179.8462
30	10	1.5	2293	2.3501	0.0571
30	10	2.0	3.2291	13.0701	9841
30	25	0.0	0.5303	6.9601	6.4298

30	25	0.05	0.2686	8.4549	8.1863
30	25	0.3	0.1809	0.2357	0.0548
30	25	1.0	1.3466	11.9003	10.5537
30	25	1.5	2.1861	5.2774	3.0913
30	25	2.0	4.5349	15.8969	11362

By taking a look at the results obtained, but especially to the table of output errors, we can observe that usually the 2nd-degree model is the one that obtains better results. This can be seen notably when the noise level is increased, that is when our sigma value is higher. Therefore, we can verify what we have studied in class, the simpler hypothesis is the one that obtains better results.

In addition, we can observe that also for higher noise values (1.5, 2 in this case), we obtain a overfit measure significantly positive and it tends to be significantly negative when the noise values are low (0, 0.05). However, there are more cases of positive overfitting than negative overfitting.

**(e) Why do we take the average over many experiments? Use the variance to select an acceptable number of experiments to average over.**

As we are using random variables in order to create our datasets for training, we need to perform different experiments for the sake of approximating better to the real value. We can see that for some tests using the same parameters, the output error is different. This specially happens when the amount of data is small. Therefore, after executing a different number of times every experiment, we compute the mean of these errors for getting a more realistic error.