

4F13: Probabilistic Machine Learning

Summarized from C. Rasmussen & D. Krueger lectures, Michaelmas 2021

Oussama Chaib

October 2021

Contents

1	Introduction	2
1.1	Modelling data	2
1.1.1	Purpose of models	2
1.1.2	Origin of models	2
1.1.3	Priors	2
1.1.4	Components of a model	2
1.1.5	Practical modelling	3
1.2	Linear in the parameters regression	3
1.2.1	Least squares approach	3
1.2.2	Probabilistic approach: Likelihood and concept of noise	4
1.3	Comments from QnA	5

1 Introduction

1.1 Modelling data

1.1.1 Purpose of models

The purpose of models is:

- Making predictions
- Generalizing: interpolation, extrapolation
- Generating more data from a similar distribution as the training set
- Compressing and summarizing data
- Interpreting statistical relationships in data
- Evaluating the relative probability of a hypothesis on data

1.1.2 Origin of models

The origin of models can be:

- **First principles:** (i.e: Newtonian mechanics model, high level of accuracy)
- **Observations and data:** (i.e: annual production of timber depending on climate and geographical factors)

Definition – Machine learning is a broad term that covers theory and practice of mathematical models which to a significant degree rely on data.

1.1.3 Priors

Every model relies on priors:

- **Knowledge**
- **Assumptions** (could be true or false)
- **Simplifying assumptions** (not necessarily true, but good enough – i.e: the mistake associated with the assumption is fairly small even though it might not be necessarily true)

1.1.4 Components of a model

Time series have:

- Unobserved/hidden/latent variables ($x(t)$, $x(t - 1)$)
- Observations (shaded $y(t)$, $y(t - 1)$)
- Parameters to link everything
 - Transitions (between latent variables)
 - Emissions (from a latent variable to an observation)

Note: The number of latent variables increases with the number of observations, but the number of parameters doesn't!

– **Learning/training models:** "What to do with all this data?"

Depending on the data, some models include: inference, estimation, sampling, and marginalization.

1.1.5 Practical modelling

1. Treat (training) the unobserved quantities (latent variables, observations, parameters)
2. Make predictions based on test cases, interpret the trained model (can we figure out what the model is trying to tell us about the data?)
3. Evaluate the accuracy of the data
4. Model selection and criticism (choose the right model or variant of the model, identify limitations)

There is not "true" or "correct" model – "*All models are wrong, but some are useful*" - George E.T. Box

1.2 Linear in the parameters regression

Let's start off with a dataset $D = x_i, y_{i=1}^N$. From a dataset of N points, we would like to infer the coordinate y_* at x_* . A simple model to do that would be polynomial **linear in the parameters** regression:

$$f_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M$$

where w_i are the corresponding weights of the polynomial, and the **parameters** of the model.

– **Relevant questions:**

Model structure: Should we choose a polynomial? What degree M should we choose?

Parameters: What values of w_i do we choose?

1.2.1 Least squares approach

Let's find the "best" polynomial (degree M and weights w_i) according to the least squares approach (minimizing the variance or sum of squared error e_i^2):

$$e_i(x)^2 = (y_i(x_i) - f_w(x_i))^2$$

$$E(x) = \sum_{i=1}^N e_i^2$$

– **Some notations:**

- **Training target** (input data that we would like to fit):

$$y = [y_1, \dots, y_N]$$

- **Prediction** (the functions that will attempt to fit the data y):

$$f = [f_w(x_1), \dots, f_w(x_N)]$$

- **Errors:**

$$e = f - y = [e_1, \dots, e_N]$$

- **Parameters (what we're solving for!):** The **weights** of f :

$$w = [w_0, \dots, w_M]$$

The sum of errors we would like to minimize is (defined as a variance to minimize):

$$E(w) = ||e||^2 = e^T e = (f - y)^T (f - y)$$

We define the **basis function** of the model (linear in the parameters model) as a matrix that contains all the x^j (the matrix contains a list of all polynomials of unknown x with maximum order from 0 to M) :

$$\Phi = \Phi_{i,j} = \begin{bmatrix} \phi_0(x_1) & \dots & \phi_M(x_1) \\ \dots & \dots & \dots \\ \phi_0(x_M) & \dots & \phi_M(x_M) \end{bmatrix}$$

$$\text{where } \phi_j(x) = x^j$$

Ultimately, we can define the prediction function as a **linear in parameters** model that is the product of the *basis function* and the *weights*:

$$\boxed{f = \Phi \cdot w}$$

To minimize the sum of squared errors $E(w)$, the gradient of the function should be equal to the zero vector:

$$\frac{\partial E(w)}{\partial w} = 2\Phi^T \Phi w - 2\Phi^T y = "0"$$

The weight vector \hat{w} that minimizes $E(w)$:

$$\boxed{\hat{w} = (\Phi^T \Phi)^{-1} \Phi^T y}$$

The error vector is minimal when it is orthogonal with all columns of Φ :

$$\Phi^T \cdot e = 0$$

When we go to higher values of M, we are able to fit a function that minimizes the error term but we may run into the problem of **overfitting**. In that case, **additional assumptions** are needed.

1.2.2 Probabilistic approach: Likelihood and concept of noise

(Probabilistic view of what was done in the previous section. We're doing the exact same thing but using probabilities, only to end up getting equivalent weights)

Definition: The **likelihood** is the probability of the data given the parameters.

Example:

- $p(y|w, \sigma^2)$ is the probability of the observed data given the weights and noise.
- $L(w) \propto p(y|w, \sigma^2)$ is the **likelihood** of the weights.

The **maximum likelihood** of the weights:

$$\hat{w} = \operatorname{argmax}(L(w)) = \operatorname{argmax}(e^{-\frac{E(w)}{2\sigma_{noise}^2}}) = \operatorname{argmin}(E(w))$$

which is equivalent to the result that we got using least squares. But we still haven't solved the problem: we still overfit!

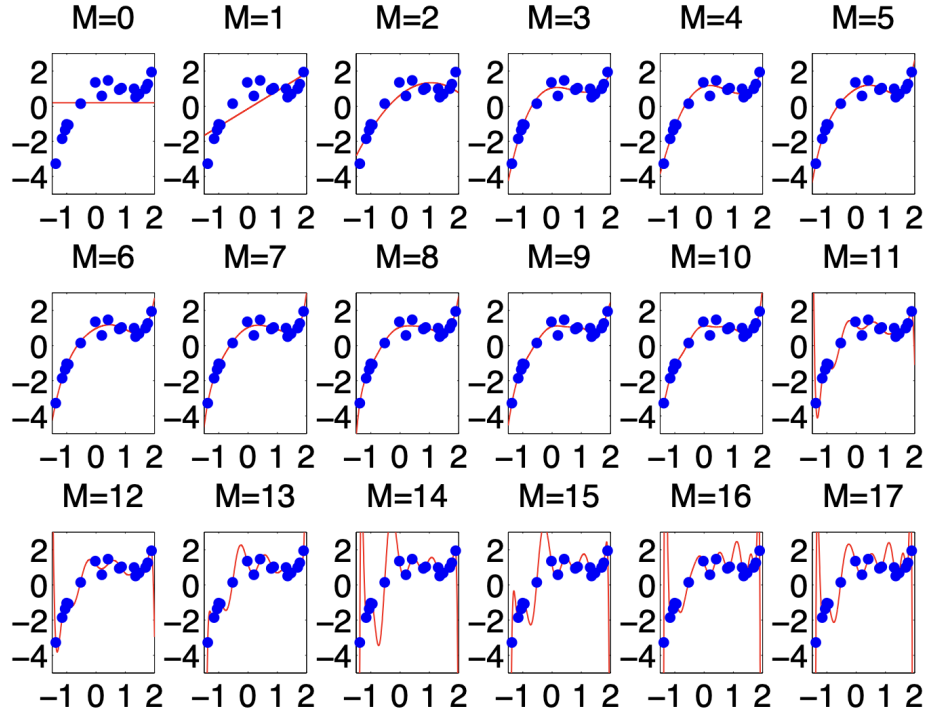


Figure 1: Overfitting illustrated. In this example, we are solving for the weights for each case of M

1.3 Comments from QnA

– Why is the error vector e minimal if it's orthogonal to all columns of ϕ ?

ϕ is a fixed function, it doesn't have any parameters in it. Still linear in parameters because the product of w and ϕ matrices is linear (even at high polynomial orders).

– What is Euclidian geometry?

Simple geometries, straight lines, basic shapes (including circle).

– Bayesian methods, why are they not as popular?

Not so sure about that. Maximum likelihood (used often, single value that best explains the data, quite popular/successful).