

CMPT 419/726 — Machine Learning

Salehen Shovon Rahman

September 17, 2015

Personal notes regarding machine learning.

Polynomial Curve Fitting

HERE'S AN EXAMPLE MACHINE LEARNING PROBLEM: try to find the best polynomial that can potentially fit a set of data points, and have it be fit as best as possible. This is known as the polynomial curve fitting problem, and it's a supervised regression learning problem.

The Problem

SUPPOSE we are given a training set of N observations, (x_1, x_2, \dots, x_N) and (t_1, t_2, \dots, t_N) , $x_i, t_i \in \mathbb{R}$. We want to find a polynomial $y(x)$ that fits these data the best.

Let's start out by defining a $y(x, \mathbf{w})$.

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{i=1}^M (w_i x^i)$$

How do we measure success? Or, a better question, for what values of the coefficients \mathbf{w} will yield the best results? To answer that, we define an error function E .

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2$$

We then use the $\arg \min_x f(x)$ function to find the value for the parameter that yields the lowest value in a given function.

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$$

$\min_x f(x)$ finds the lowest possible value for the expression $f(x)$, while $\arg \min_x f(x)$ finds the value for x where $f(x)$ would be the lowest.

So, in other words, we want to find a \mathbf{w} such that $E(\mathbf{w})$ is the lowest among the set of all possible values of \mathbf{w} .

EXCEPT, the attempt at finding a value \mathbf{w}^* such that $E(\mathbf{w}^*) = 0$ can become problematic.

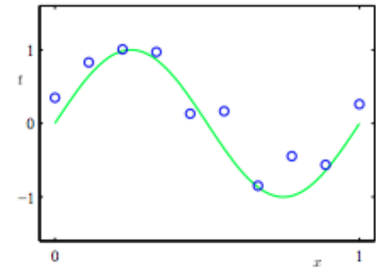


Figure 1: An example data set where we want to fit a polynomial curve into.

Earlier, we mentioned that we had an initial set of training data. However, for most cases, when trying to fit the polynomial such that $E(\mathbf{w}^*) = 0$, for the training set, we risk having it so that when a test data set is introduced, the error function yields a high value. This is known as *overfitting*.

In the end of the day, although we want the curve to fit the data as best as possible, we also want a *generalization* derived from the given training.

BUT FIRST, before we go ahead with finding a good generalization, for convenience, instead of just using the error function E , we use the root-mean-square (RMS) error function, defined by

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

According to the text book, *Pattern Recognition and Machine Learning*, the reason why we are using RMS, is the following:

[The] division by N allows us to compare different sizes of data sets on an equal footing, and the square root ensures that E_{RMS} is measured on the same scale (and in the same units) as the target variable t . (p. 7)

Now, to actually tune our \mathbf{w} for better generalization, we can split our training data into two sets: training set and validation set. In the case of finding the polynomial, the training set can be used to find each $w_i \in \mathbf{w}$, and the validation set is used to optimize the complexity, which can be represented by M (the size of \mathbf{w}), or a λ , which will be discussed later.

There are several techniques used to control overfitting.

THE FIRST TECHNIQUE to avoid overfitting is cross-validation.

Here, we group the data into separate sets. We first “train” our parameters to a union of all the separated set, while leaving one out. Then we optimize by including the one we initially excluded. Afterwards, we “train” again with a new union of our sets, while leaving yet another one out, but including the one that we initially left out, all the way until no sets are left to “leave out”.

AND THEN, there’s regularization for controlling over-fitting.

Notice how the oscillation increases as M increases? This is because the magnitudes of the coefficients in \mathbf{w} increases as M increases.

In order to avoid high coefficient magnitudes, we can “penalize” them using a modified error function, by adding a $\frac{\lambda}{2} \|\mathbf{w}\|^2$ term to the original error function.

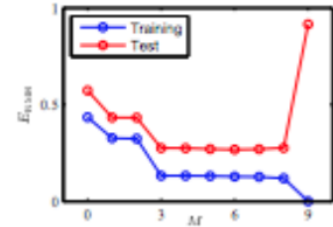


Figure 2: As we can see, the first few polynomials of degree $N < 9$ fit the data fine, even when test data is introduced to the training set, but misses the mark entirely when $N = 9$. This is the result of overfitting.

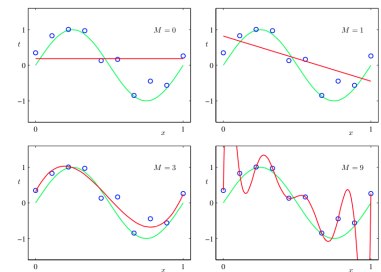


Figure 3: Visually, we see that the oscillation increases as M increases.

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0	0.19	0.82	0.31	0.35
w_1		-1.27	7.99	232.37
w_2			-25.43	-5321.83
w_3			17.37	48568.31
w_4				-231639.30
w_5				640042.26
w_6				-1061800.52
w_7				1042400.18
w_8				-557682.99
w_9				125201.43

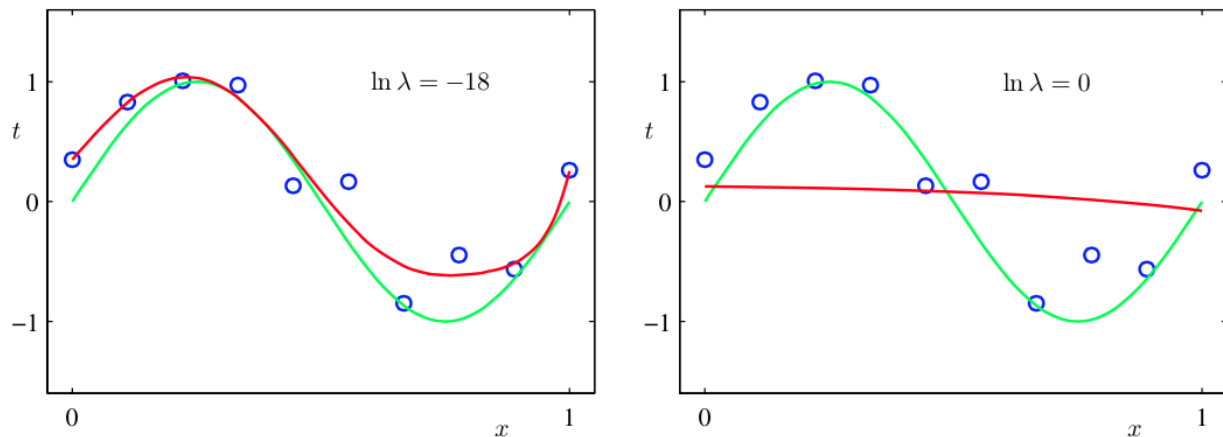
Table 1: For higher values of M , we see the magnitude of w_i increasing

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = E(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

If we are to now apply the error function to our trials, we see that the coefficient are no longer large.

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 9$
w_0	0.35	0.35	0.13
w_1	232.37	4.74	-0.05
w_2	-5321.83	-0.77	-0.06
w_3	48568.31	-31.93	-0.05
w_4	-231639.30	-3.89	-0.03
w_5	640042.26	55.28	-0.02
w_6	-1061800.52	41.32	-0.01
w_7	1042400.18	-45.95	-0.00
w_8	-557682.99	-91.53	0.00
w_9	125201.43	72.68	0.01

Table 2: For higher values of λ , the coefficient magnitudes are much lower, possibly even near 0



As you can see here, even for $M = 9$, previously, the polynomial curve deviated wildly, which could have potentially yielded high error values relative to new potential data

FINALLY, there's the third option: just get more data! A rule of thumb is that the number of datapoints should not be any less than

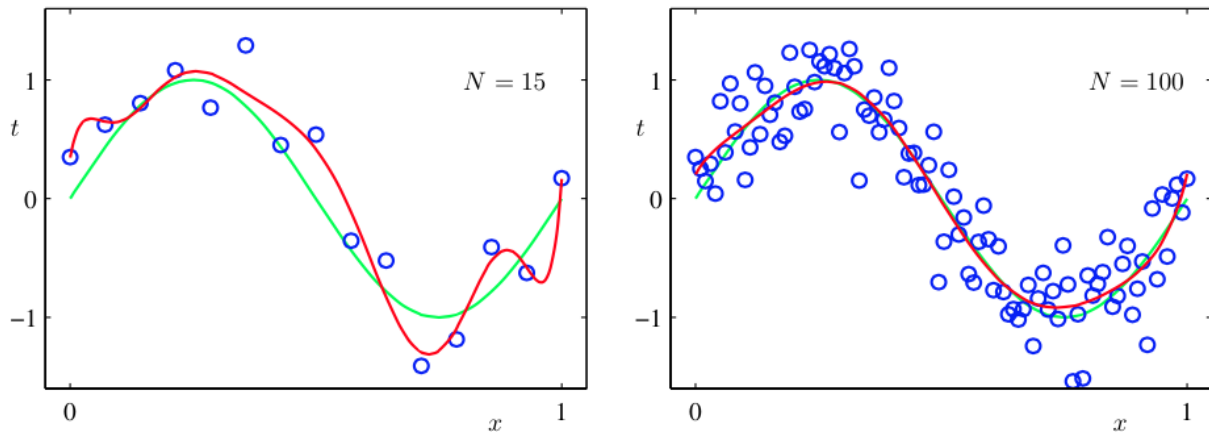


Figure 4: On the left, we see an *attempt* at having the curve fitted to the data we have so far; on the right, we see more data added, and the curve appears show a much better approximation.

five times the number of adaptive parameters.

Probability Theory

Random Variables

THE TERM “random variable” is can be misleading. A random variable is not a variable at all. Instead, it’s a mapping from a *random event* to a possible outcome.

So here’s an example random variable:

$$X = \begin{cases} 1 & \text{if heads} \\ 0 & \text{if tails} \end{cases}$$

As we can see here, X is a mapping from the set {heads, tails}, to the set {1, 0}.

More informally, a random variable is a “label” to a given element of a set of possible events.

To add to the confusion, when we say that we have an event from our random variable (for example, an event from X defined above), we express it as $X = 0$ or $X = 1$, etc.

The end purpose of random variables is to derive a probability distribution given a value $x \in \{\text{Range of } X\}$. So, going back to the above head/tail example, if we wanted to express the probability of getting heads, we write $p(X = 1)$, likewise, for expressing the probability of getting tails, we write $p(X = 0)$, and so and so forth.

The purpose of using random variables is for us to easily be able to write inequalities inside our probability expression. So, for instance, let’s define a random variable Y that contains a mapping of

a list of combinations of two dices to the sum of their faces. We can easily express the probability of a roll, for example, a value greater than 10, like so: $p(Y > 10)$.

Figures and Tables

Figure 5 is an example of the `figure*` environment and figure 6 is an example of the normal `figure` environment.

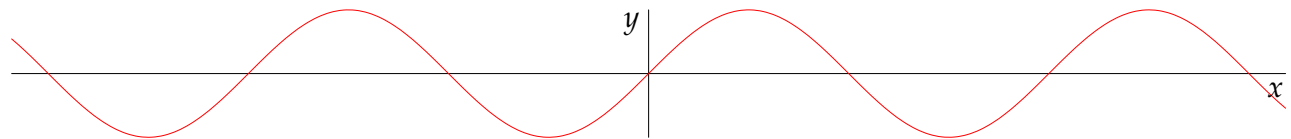


Figure 5: This graph shows $y = \sin x$ from about $x = [-10, 10]$. Notice that this figure takes up the full page width.

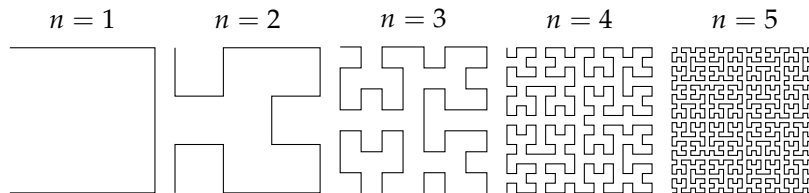


Figure 6: Hilbert curves of various degrees n . Notice that this figure only takes up the main textblock width.

Table 3 shows table created with the `booktabs` package. Notice the lack of vertical rules—they serve only to clutter the table’s data.

Margin	Length
Paper width	8 $\frac{1}{2}$ inches
Paper height	11 inches
Textblock width	6 $\frac{1}{2}$ inches
Textblock/sidenote gutter	$\frac{3}{8}$ inches
Sidenote width	2 inches

Table 3: Here are the dimensions of the various margins used in the Tufte-handout class.

Full-width text blocks

In addition to the new float types, there is a `fullwidth` environment that stretches across the main text block and the sidenotes area.

```
\begin{fullwidth}
Lorem ipsum dolor sit amet...
\end{fullwidth}
```

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus

rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Typography

Typefaces

If the Palatino, Helvetica, and Bera Mono typefaces are installed, this style will use them automatically. Otherwise, we'll fall back on the Computer Modern typefaces.

Letterspacing

This document class includes two new commands and some improvements on existing commands for letterspacing.

When setting strings of ALL CAPS or SMALL CAPS, the letter-spacing—that is, the spacing between the letters—should be increased slightly.¹ The `\allcaps` command has proper letterspacing for strings of FULL CAPITAL LETTERS, and the `\smallcaps` command has letterspacing for SMALL CAPITAL LETTERS. These commands will also automatically convert the case of the text to upper- or lowercase, respectively.

The `\textsc` command has also been redefined to include letterspacing. The case of the `\textsc` argument is left as is, however. This allows one to use both uppercase and lowercase letters: THE INITIAL LETTERS OF THE WORDS IN THIS SENTENCE ARE CAPITALIZED.

Installation

To install the Tufte- \LaTeX classes, simply drop the following files into the same directory as your `.tex` file:

```
tufte-book.cls
tufte-common.def
tufte-handout.cls
tufte.bst
```

More Documentation

For more documentation on the Tufte- \LaTeX document classes (including commands not mentioned in this handout), please see the sample book.

Support

The website for the Tufte- \LaTeX packages is located at <https://github.com/Tufte-LaTeX/tufte-latex>. On our website, you'll find links to our svn repository, mailing lists, bug tracker, and documentation.