

Random Thoughts on Machine Learning

A Layman's Perspective

YU LIU¹

May 7, 2016

¹https://github.com/ohliumliu/thoughts_learning.git

Contents

1	Gradient descent: a first-order approximation	3
2	Regularization: the “hidden” facts?	5
2.1	Motivation of regularization	6
2.2	Linear model	6
2.3	Polynomial model	8
3	Principle component analysis	11
3.1	Moments of inertia	12
4	Entropy cost: Boltzmann or Shannon?	13
4.1	Boltzmann	14
5	Feature scaling: the Π theorem	15
5.1	Buckingham’s contribution	16
6	k-mean: the Harmonic mean	17
6.1	Interpretation in terms of harmonius average	18

List of Figures

List of Tables

Preface

This is a collection of ongoing thoughts I gathered on my journey to better understand machine learning. I have been trying to relate with my past training in math and physics and these thoughts, naive as it may appear to professionals, might be helpful to get my feet into the door of this field.

1

Gradient descent: a first-order approximation

to be done.¹

¹footnote

2

Regularization: the “hidden” facts?

2.1 Motivation of regularization

It is well known that introducing too many features in a model will almost always lead to perfect score with the training and validation set, but miserable failure with the test set. In another word, a feature-rich model has a very small bias (toward a particular model), but a very high variance (when used for prediction). The high variance probably comes from the observation that the confidence region of the fitted parameters is big enough to overwhelm the average value.

If there is no way or rationale to reduce the number of features, regularization is often applied to solve the problem of over-fitting. I have a feeling that this amounts to the introduction of extra data points to the modeling problem, and this chapter is a summary of my attempt to formulate this feeling.

2.2 Linear model

I will start with the simple case of a linear model with a feature space $\vec{x} \equiv (x_1, \dots, x_n)^T$ and the corresponding set of coefficient vector $\vec{\theta} \equiv (\theta_1, \dots, \theta_n)^T$. The prediction of the model is simply linear with a constant offset θ_0 :

$$h(\vec{x}; \vec{\theta}) = \theta_0 + \vec{\theta} \cdot \vec{x} \quad (2.1)$$

Given m data points $(\vec{x}^{(m)}, y^{(m)})$ with $\vec{x}^{(m)} = (\vec{x}_1^{(m)}, \dots, \vec{x}_n^{(m)})^T$, the cost function without regularization reads

$$J(\vec{\theta}) = \sum_{i=1}^m \left[y^{(i)} - h(\vec{x}^{(i)}; \vec{\theta}) \right]^2. \quad (2.2)$$

With regularization, an extra term is added to penalize non-zero values of θ_i ,

$$J(\vec{\theta}) = \sum_{i=1}^m \left[y^{(i)} - h(\vec{x}^{(i)}; \vec{\theta}) \right]^2 + \lambda \sum_{i=0}^n \theta_i^2, \quad (2.3)$$

As an aside, we can vectorize the cost function by first defining an error vector $(m \times 1)$,

$$\vec{\epsilon} = \vec{y} - \left[\vec{\theta}_0 + \vec{\theta}^T \cdot \mathbf{X} \right], \quad (2.4)$$

where $\vec{\theta}_0 \equiv \theta_0 \times (1, \dots, 1)$ and \mathbf{X} is an $n \times m$ matrix which each column as $\vec{x}^{(m)}$. In this way, the cost function is

$$J(\vec{\theta}) = \vec{\epsilon} \cdot \vec{\epsilon} + \lambda(\theta_0^2 + \vec{\theta} \cdot \vec{\theta}). \quad (2.5)$$

It's also worth noting that the penalty term could take other forms. For example, in Lasso regression, the L_1 norm (absolute values) of θ_i is used.

Going back to the regularization process, the factor λ is a tunable parameter to control the biasness of the regression. I will try to show that, in this context of simple linear regression, the introduction of regularization term seems to be equivalent to the “artificial” addition of extra data points.

In order to demonstrate this point, we need to find n data points that can lead to the extra penalty term $\lambda(\theta_0^2 + \vec{\theta} \cdot \vec{\theta})$.

For example, we want to find $(\vec{u}^{(1)}, v^{(1)})$ that satisfies

$$\left[v^{(1)} - (\theta_0 + \vec{\theta} \cdot \vec{v}^{(1)}) \right]^2 = \lambda \theta_1^2. \quad (2.6)$$

An obvious choice that is independent θ_i is $\vec{u}^{(1)} = (\sqrt{\lambda}, 0, \dots)$ and $v^{(1)} = \theta_0$. In general, we would like to have $\vec{u}^{(i)} = \{\sqrt{\lambda} \delta_{ij}\}_j (j = 1, \dots, n)$ and $v^{(i)} = \theta_0$ ¹. We need to have n points.

With this choice of points, the regularization cost can be compared with adding some extra points to the training set. In the space spanned by the features, these points share the same predicted value or θ_0 and distributed on each of the independent axis. It also includes the origin. The adoption of regularization essentially claims that there is no dependence of the prediction on any of the features. The hypothesis, of course, is just the opposite, hence the hypothesis is associated with some “bias”. The combination of agnostic regularization and the biased hypothesis reminds me of the central idea of Bayes statistics. In Bayes, the prior knowledge has to be “unbiased”, and the training set makes the posterior knowledge more “biased”. Here, the two parts are mixed together in the hope to achieve a balance.

The strength of the agnostic regularization is adjusted by λ . In the interpretation above, it translates to the distance of the extra points from the origin. In the hyper-plane spanned by the features and the prediction, these extra points, $(\vec{u}^{(i)}, v^{(i)}) (i = 0, \dots, n)$, define a hyperplane which is “flat” (agnostic). The size of the plane is $\sqrt{\lambda}$. If λ is bigger, it will be more important in the regression because it covers a larger range of feature values. When it is smaller, it will not be very significant compared with the original training points. In the extreme case that $\lambda \sim 0$, these extra points gather around $(0, \theta_0)$, no longer serving as extra constraints.

An immediate application of this interpretation is the choice of λ . It should be based on the distribution of feature values. Because there is only one lambda available, it is obviously beneficial to perform feature scaling. Or else, there is no way to choose a universal λ that is appropriate for all features. Of course, we could use a different λ for different θ , but this is adding to the complexity of the model furthermore.

¹You must have realized that the choice of $v^{(i)} = \theta_0$ violates our original requirement that it should be independent of θ_i . To simplify the discussion, let's assume θ_0 is not to be optimized. Another way to put it is to shift y by $\langle y \rangle$ before regression and scale the features by the mean as well, then the model shouldn't involve the constant term.

In practice, it may be helpful to study the response of θ_i with respect to the change in λ . As λ increase, some coefficients may be more resilient than others, and these must be more important features. For those θ_i that approaches zero quickly, we might remove them just as well.

Another potential development is to dig deeper into the Bayes interpretation. Is there a way to adjust λ based on training set in iterations: starting from a large λ , as more training data are fed to the program, λ is reduced. I am not sure if this makes sense.

2.3 Polynomial model

Without increasing feature set, one can increase the complexity of the model, for example, by increasing the polynomial order. To simplify the discussion, let's assume that the features have been properly scaled so that it is reasonable to use the following hypothesis,

$$h(\vec{\theta}, x) = \sum_{i=1}^n \theta_i x^i, \quad (2.7)$$

or, in matrix format,

$$h(\vec{\theta}; x) = \vec{\theta} \cdot \begin{pmatrix} x \\ x^2 \\ \vdots \\ x^n \end{pmatrix}. \quad (2.8)$$

Following the same idea used in the linear case, in order to interpret the regularization term as the cost function associated with some extra points, we need to find one or more points (x_i, y_i) independent of θ_i that satisfies the following:

$$\sum_{i=1}^m \left[h(\vec{\theta}, x_m) - y_m \right]^2 = \sum_{i=1}^n \theta_i^2. \quad (2.9)$$

Since these points should reflect our unbiased knowledge, it seems only reasonable to choose a vanishing y_m . In fact, we should choose

$$y_m = \sum_{i=1}^n \theta_i \langle x^n \rangle, \quad (2.10)$$

which reduces to 0 only for some cases even though $\langle x \rangle = 0$ due to feature scaling.

For now, with $y_m = 0$, Eq. 2.9 is reduced to the following matrix format,

$$\|\vec{\theta} \cdot \mathbf{X}\|^2 = \|\vec{\theta}\|^2, \quad (2.11)$$

where

$$\mathbf{X} = \begin{pmatrix} x_1 & x_2 & \cdots & x_m \\ x_1^2 & x_2^2 & \cdots & x_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & \cdots & x_m^n \end{pmatrix}. \quad (2.12)$$

It is not clear to me at this point how to find x_i that is independent of θ_i . In a simpler case where $m = n$, \mathbf{X} needs to be a rotation matrix to satisfy Eq. 2.11. This is because Eq. 2.11 only requires that the norm of the l.h.s and r.h.s is the same, and rotation keeps the norm ⁽²⁾ If we recall the choice of extra point in the linear case, $\mathbf{X} = \sqrt{\lambda}I$, and its null space is the whole feature space. This is stronger than Eq. 2.11. Of course, we could have chosen a more general rotation matrix at that time.

However, it is probably impossible to choose x_i such that \mathbf{X} is a rotation matrix, simply because of the strong relation between different rows of \mathbf{X} . In an almost trivial case of $n = 2$, it we can choose $x_1 = 1$ and $x_2 = -1$ to satisfy Eq. 2.11. In particular, the cost function associated with these two points is

$$J = (\theta_1 + \theta_2)^2 + (\theta_1 - \theta_2)^2 \quad (2.13)$$

$$= 2(\theta_1^2 + \theta_2^2). \quad (2.14)$$

If we attempt $m > n$, at least for each choice of $\vec{\theta}$, we should be able to find a set of x_i that can satisfy Eq. 2.11. But can we find such a set that is independent of $\vec{\theta}$?

Playing with another trivial choice of $(1, 0)$ means that its cost is $(\theta_1 + \theta_2 + \cdots + \theta_n)^2$. We cannot use this term for regularization, because it doesn't acutally penalize big θ_i . In another word, adding $(1, 0)$ as an extra point for regularization doesn't disfavor higher order term, because higher order polynomial can still go through $(1, 0)$.

2.4 Support Vector Machine

²It's been implicit that $\lambda = 1$. If not, we have to associate a uniform weight of λ with each of these extra points.

3

Principle component
analysis: Moment of inertia
and rotational spectroscopy

3.1 Moments of inertia

4

Entropy cost: Boltzmann or
Shannon?

4.1 Boltzmann

5

Feature scaling: the Π theorem

5.1 Buckingham's contribution

6

k-mean: the Harmonic mean

6.1 Interpretation in terms of harmonic average