

Random Musings on a Sunday Morning

(extraordinarily basic stuff that's worth reflecting upon now and then)

David Meyer

dmm@{1-4-5.net,uoregon.edu}

13 Oct 2015

1 Introduction

One of the interesting things about machine learning is that it affords many different ways to look at the same data and hence underlying phenomena (aside: why exactly should this be?). For example, the Logistic Function is really a special case of a Conditional Random Field and PCA is a special case of a Linear AutoEncoder (with no regularization). This document looks at a different correspondence. Here we'll look at why minimizing error (a sum), maximizing probability (a product), and minimizing energy (in an energy based model, for example, a Restricted Boltzmann Machine) are all really the same thing. Note that this document is likely to have many errors.

2 Minimizing Cost is a Sum

2.1 Linear Regression Cost Function

For linear regression, our *hypothesis* $h_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x})$, where $g(z) = z$ and \mathbf{x} and $\boldsymbol{\theta}$ are $N \times 1$ column vectors (hence $\boldsymbol{\theta}^T$ is a $1 \times N$ row vector)¹:

¹While $g(z)$ is linear here, in other models g can be a non-linearity such as the *sigmoid* function.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\boldsymbol{\theta}^T = [\theta_1, \theta_2, \dots, \theta_n]$$

which implies that

$$\boldsymbol{\theta}^T \mathbf{x} = [\theta_1, \theta_2, \dots, \theta_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n \theta_i x_i = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad (1)$$

Note that Equation 1 is the *algebraic interpretation* of the dot product of $\boldsymbol{\theta}^T$ and \mathbf{x} (also called the inner or scalar product). An alternate, *geometric interpretation* of the dot product in Euclidian space, where a geometrical object possesses both a magnitude and a direction is

$$\boldsymbol{\theta}^T \mathbf{x} = \|\boldsymbol{\theta}^T\| \|\mathbf{x}\| \cos(\alpha) \quad (2)$$

where α is the angle between the vectors $\boldsymbol{\theta}^T$ and \mathbf{x} , and $\|\boldsymbol{\theta}^T\|$ and $\|\mathbf{x}\|$ are the magnitudes of the vectors $\boldsymbol{\theta}^T$ and \mathbf{x} respectively (that is, the L^2 -norm²). Note that Equation 2 implies that

$$\cos(\alpha) = \frac{\boldsymbol{\theta}^T \mathbf{x}}{\|\boldsymbol{\theta}^T\| \|\mathbf{x}\|} = \frac{\sum_{i=1}^n \theta_i^T x_i}{\sqrt{\sum_{i=1}^n (\theta_i^T)^2} \times \sqrt{\sum_{i=1}^n (x_i)^2}} \quad (3)$$

²The general form of the p -norm is $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$

which is also known as the *cosine similarity* between $\boldsymbol{\theta}^T$ and $\mathbf{x}[1]$.

Getting back to our hypothesis, we see that $h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$. Note that by convention θ_0 is treated specially; in particular, it is not part of the vector $\boldsymbol{\theta}$. Now, given this hypothesis, our cost function can be written as a *sum*:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})^2 \quad (4)$$

The goal of machine learning then is to find the parameters $\boldsymbol{\theta}$ such that loss/error function $J(\boldsymbol{\theta})$ is minimized (note that which error is minimized, and when, is a topic unto itself). For linear regression, Equation 4 is a *convex* optimization objective.

2.2 Logistic Regression Cost Function

For logistic regression, our hypothesis $h_{\boldsymbol{\theta}}(\mathbf{x})$ is slightly different, as shown in Equations 5, 6 and 7.

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) \quad (5)$$

Here $g(z)$ is the *logistic* or *sigmoid* function³, and is defined as follows

$$g(z) = \frac{1}{1 + e^{-z}} \quad (6)$$

Putting Equations 5 and 6 together we get

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}} \quad (7)$$

While it seems like we could perhaps use the same loss function as we did in linear regression (Equation 4), it turns out that this loss function is *non-convex* when applied to logistic regression, i.e., when $g(z) = \frac{1}{1+e^{-z}}$. As a result we typically use some version of *cross-entropy* as the loss function:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \left[\sum_{i=1}^n y^{(i)} \log h_{\boldsymbol{\theta}}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(x^{(i)})) \right] \quad (8)$$

The regularized version of Equation 8 is

$$J(\boldsymbol{\theta}) = \frac{1}{n} \left[\sum_{i=1}^n y^{(i)} \log h_{\boldsymbol{\theta}}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(x^{(i)})) \right] + \frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2 \quad (9)$$

³BTW, the logistic function is a special case of a *Conditional Random Field*

2.3 Deriving the Loss/Error function for Logistic Regression

Recall that this loss function linear regression was

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (10)$$

and that this loss function was a *non-convex* optimization objective for logistic regression. To find a convex loss function for logistic regression we first define a *Cost* function:

$$Cost(h_{\theta}(x^{(i)}), y^{(i)}) = (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (11)$$

so that

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n Cost(h_{\theta}(x^{(i)}), y^{(i)}) \quad (12)$$

Now we can ask what our *Cost* function look like. Well, if we predict 1 and $y = 1$, the the cost should be close to 0 (because we predicted the right value). Alternatively, if we predict 0 and $y = 1$, then the cost should converge on ∞ (that is, we penalize the prediction). On the other hand, if we predict 0 and $y = 0$, the the cost should again be close to 0 as we predicted the right value. Similarly, if we predict 1 and $y = 0$, the cost should again converge on ∞ . These two cases are captured in Equation 13.

$$Cost(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \quad (13)$$

Note that if you combine the two cases in Equation 13 together, you get Equation 8, the *cross entropy*.

In any event, in both cases we fit the parameters θ to the model by minimizing a **sum** such as Equation 4 or 8.

3 Maximizing Probability is a Product

Basic assumption: all of this analysis depends on the assumption, as you will see, that the error $\epsilon^{(i)}$ is Gaussian, so Danger Will Robinson!⁴ That said, assume that the target variables and the inputs are related via the equation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \quad (14)$$

⁴In reality it could be any parametric distribution. In this case I know what the simplifying assumptions are/mean so I can write the density down.

and thus

$$e^{(i)} = y^{(i)} - \theta^T x^{(i)} \quad (15)$$

where $\epsilon^{(i)}$ is an error term that captures either unmodeled effects (such as if there are some important features that we left out of the model, or just random noise). Interestingly this is intuitive as $\epsilon^{(i)} = y^{(i)} - \theta^T x^{(i)}$. Assume also that the $\epsilon^{(i)}$'s are IID (Independent and Identically Distributed) according to some Gaussian distribution with mean $\mu = 0$ and variance σ^2 , i.e., $\mathcal{N}(0, \sigma^2)$ (note that this argument also relies on the *Central Limit Theorem* and the *Law of Large Numbers* [2]). We typically use the notation $X \sim \mathcal{N}(\mu, \sigma^2)$ to indicate that the random variable X is normally distributed with mean μ and variance σ^2 (note that if $X \sim \mathcal{N}(0, 1)$ we say that X follows the *standard normal* distribution).

The general form of the probability density function (pdf) of $\mathcal{N}(\mu, \sigma^2)$ is⁵

$$P(x) = f(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right], \quad -\infty < x < \infty \quad (16)$$

Here $\mu = \mathbb{E}[X]$, i.e., the mean (and mode), $\sigma^2 = \text{var}[X]$, and $\sqrt{2\pi\sigma^2}$ is a normalization constant that ensures that the density f integrates to 1.

Now, assuming $e^{(i)} \sim \mathcal{N}(0, \sigma^2)$, we can write the density of $\epsilon^{(i)}$ as⁶:

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right] \quad (17)$$

so that

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right] \quad (18)$$

Rewriting this in vector/sum notation

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \sigma) = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2\right] \quad (19)$$

⁵ $\exp(x)$ is defined to be e^x

⁶We can always "normalize" our data so that $\mu = 0$ by computing the z score, $z = \frac{x - \mu}{\sigma}$

and taking advantage of the fact that the sum of exponential powers is product of exponentials⁷ we get

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left[-\frac{1}{2\sigma^2}(y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2\right] \quad (20)$$

Note that the *likelihood* of a set of parameter values $\boldsymbol{\theta}$ given outcomes x is equal to the probability of those observed outcomes given those parameter values, that is

$$\mathcal{L}(\boldsymbol{\theta} | x) = P(x | \boldsymbol{\theta}) \quad (21)$$

In the discrete case

$$\mathcal{L}(\boldsymbol{\theta} | x) = p_{\boldsymbol{\theta}}(x) = P_{\boldsymbol{\theta}}(X = x) \quad (22)$$

Thus we can see that Equation 20 is essentially the likelihood of $\boldsymbol{\theta}$ given X , and can be rewritten as

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \boldsymbol{\theta}) \quad (23)$$

which familiarly implies

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2}{2\sigma^2}\right] \quad (24)$$

That is, maximizing the likelihood is a **product**. Going the other way, maximizing the log likelihood $\ell(\boldsymbol{\theta})$ gives

$$\ell(\boldsymbol{\theta}) = \log \mathcal{L}(\boldsymbol{\theta}) \quad (25)$$

$$= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2}{2\sigma^2}\right] \quad (26)$$

$$= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2}{2\sigma^2}\right] \quad (27)$$

$$= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2 \quad (28)$$

Hence maximizing $\ell(\boldsymbol{\theta})$ gives the same answer as minimizing

$$\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \boldsymbol{\theta}^T x^{(i)})^2 \quad (29)$$

⁷ $a^{(b+c)} = a^b \times a^c$

which turns out to be the original least squares cost function $J(\boldsymbol{\theta})$ (Equation 4).

4 Minimizing Energy

So we've seen that minimizing error (a sum) is roughly equivalent to maximizing the probability (or likelihood), a product. It turns out that minimizing energy in a physical system such as a system of springs is the same thing! In this section we'll look at minimizing energy functions such as used by Restricted Boltzmann Machines (RBMs). Energy-based probabilistic models (e.g., RBMs) define a probability distribution through an energy function, as follows:

$$p(x) = \frac{e^{-E(x)}}{Z} \quad (30)$$

where the normalizing factor Z , called a partition function and is typically computationally intractable. Z is defined as follows

$$Z = \sum_{i=1} e^{-E(x)} \quad (31)$$

Using these definitions we can write the expression for the *likelihood* \mathcal{L} of θ and \mathcal{D} as

$$\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{N} \sum_{x^{(i)} \in \mathcal{D}} \log p(x^{(i)}) \quad (32)$$

$$\ell(\theta, \mathcal{D}) = -\mathcal{L}(\theta, \mathcal{D}) \quad (33)$$

which we can minimize using stochastic gradient descent (SGD) where the gradient is $-\frac{\partial \log p(x^{(i)})}{\partial \theta}$ and $\boldsymbol{\theta}$ are the model parameters.

References

- [1] The Wikipedia. Cosine similarity.
- [2] The Wikipedia. Law of large numbers.