

Machine Learning Notes

Sebastian Gaume

June 11, 2020

Contents

1	Algorithms	2
1.1	Gradient Descent	2
2	Linear Regression	3
2.1	Hypothesis and Cost Function	3
2.2	The Normal Equation	3
2.3	Gradient Descent	3
2.4	Regularised Linear Regression	3
2.4.1	The Normal Equation for Regularised Linear Regression	4
2.4.2	Gradient Descent for Regularised Linear Regression	4
3	Logistic Regression	5
3.1	Hypothesis and Cost Function	5
3.2	Decision Boundaries	5
3.3	Multi-Class Problems	5
3.4	Gradient Descent	6
3.5	Regularised Logistic Regression	6
3.5.1	Gradient Descent for Regularised Logistic Regression	6

Algorithms

Algorithm 1: The gradient descent algorithm

Output: The optimised parameter vector.

$$n \leftarrow 0;$$
$$\theta \leftarrow thetaInit;$$
while $n < nIters$ **do**
$$\theta \leftarrow \theta - \alpha \nabla J(\theta) ;$$

```
// J(theta) is the cost function we are optimising against
```

$$n \leftarrow n + 1;$$

end

return θ

Chapter 2

Linear Regression

Linear regression models are used in regression problems of an arbitrary number of features and with features which can be functions of other features. In these problems, $\mathbf{y}^{(i)} \in \mathbb{R}$ and can represent any quantity. Our goal is to predict, based on features $\mathbf{x}^{(i)}$, what output value a datum corresponds to. For this model, the output of the hypothesis function is the predicted value of $\mathbf{y}^{(i)}$ given $\mathbf{x}^{(i)}$, parameterised by $\boldsymbol{\theta}$.

2.1 Hypothesis and Cost Function

In a linear regression model, our hypothesis is of the form

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \quad h_{\boldsymbol{\theta}}(\mathbf{x}) \in \mathbb{R} \quad (2.1)$$

or alternatively,

$$\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{X}) = \mathbf{X}\boldsymbol{\theta}, \quad \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{X}) \in \mathbb{R}^m \quad (2.2)$$

and we use a mean square difference cost function of the form

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m \left[(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2 \right] = \frac{1}{2m} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}). \quad (2.3)$$

2.2 The Normal Equation

One way to minimise the cost function in linear regression is to do so analytically using the normal equation. This says that

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (2.4)$$

gives us the $\boldsymbol{\theta}$ that minimises the cost function $J(\boldsymbol{\theta})$. This technique is useful up until $n \geq 10^4$, as for these relatively low numbers of features it is more efficient than gradient descent and doesn't require us to choose a learning rate α .

2.3 Gradient Descent

We can also minimise the cost function using gradient descent. To do this, we need the gradient of $J(\boldsymbol{\theta})$ which is

$$\nabla_j J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \cdot \mathbf{x}_j^{(i)} \right] \quad (2.5)$$

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{m} \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad (2.6)$$

and so we can just use this in our gradient descent algorithm and it will optimise $\boldsymbol{\theta}$. Note that all linear regression hypotheses are convex and so gradient descent will always converge with a correct learning rate and enough iterations to the global optimum.

2.4 Regularised Linear Regression

To avoid overfitting our training data by adding higher order terms to our hypothesis, we can use regularisation, which penalises larger weights for each term causing the algorithm to prefer lower valued weights and thus reducing

the problem of overfitting. To do this, we need to change our cost function to

$$J(\boldsymbol{\theta}) = \frac{1}{2m}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2. \quad (2.7)$$

The extra term penalises higher values of θ_j for all of $\boldsymbol{\theta}$ except θ_0 by convention. It also introduces a new hyperparameter λ which controls how much the algorithm penalises the higher values. In general, we want $\lambda \geq 1$ but not too large as if it's too large it will cause the model to underfit the training data.

2.4.1 The Normal Equation for Regularised Linear Regression

To minimise our new cost function for regularised linear regression, we have the same two options as before, but we need to modify them slightly. For the normal equation, this means we instead find that

$$\boldsymbol{\theta} = \left(\mathbf{X}^\top \mathbf{X} + \lambda \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix} \right)^{-1} \mathbf{X}^\top \mathbf{y} \quad (2.8)$$

finds the $\boldsymbol{\theta}$ that minimises the cost function $J(\boldsymbol{\theta})$. Note that the matrix that is being multiplied by λ is the \mathbb{R}^{n+1} identity matrix except the $I_{1,1}$ entry is 0 instead of 1.

2.4.2 Gradient Descent for Regularised Linear Regression

We also need to change the $\nabla J(\boldsymbol{\theta})$ we use in our gradient descent algorithm if we choose to use regularised linear regression. The new gradient is given by

$$\nabla_0 J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \right] \quad (2.9)$$

$$\nabla_j J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \cdot \mathbf{x}_j^{(i)} \right] + \frac{\lambda}{m} \theta_j \quad \text{for } j \neq 0 \quad (2.10)$$

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{m} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) + \frac{\lambda}{m} \begin{bmatrix} 0 \\ \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \\ \vdots \\ \boldsymbol{\theta}_n \end{bmatrix} \quad (2.11)$$

and we simply use these in our gradient descent algorithm instead of the previous versions.

Chapter 3

Logistic Regression

Logistic regression models are used in binary classification problems, although they can be extended to multi-class classification problems. In these problems, $\mathbf{y}^{(i)} \in \{0, 1\}$, where 0 represents the negative case and 1 represents the positive case, though these can then represent any given classes. Our goal is to predict, based on features $\mathbf{x}^{(i)}$, which class a datum is part of. For this model, the output of the hypothesis function is $P(\mathbf{y}^{(i)} = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta})$.

3.1 Hypothesis and Cost Function

In a logistic regression model, our hypothesis is of the form

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}), \quad h_{\boldsymbol{\theta}}(\mathbf{x}) \in \mathbb{R} \quad (3.1)$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad (3.2)$$

or equivalently

$$\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{X}) = g(\mathbf{X}\boldsymbol{\theta}), \quad \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{X}) \in \mathbb{R}^m \quad (3.3)$$

as this ensures our hypothesis is between 0 and 1 for all possible inputs. We then use a cost logarithmic cost function which rewards $h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$ for being closer to $\mathbf{y}^{(i)}$ than the other possible value of $\mathbf{y}^{(i)}$ as it can only hold the values of 0 and 1. This cost function is

$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[\mathbf{y}^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - \mathbf{y}^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] \quad (3.4)$$

$$= \frac{1}{m} (-\mathbf{y}^T \log(g(\mathbf{X}\boldsymbol{\theta})) - (\mathbf{1} - \mathbf{y})^T \log(\mathbf{1} - g(\mathbf{X}\boldsymbol{\theta}))) \quad (3.5)$$

and it is, like that of linear regression, convex. This means we can be sure gradient descent will always converge to the global optimum when used to minimise the function.

3.2 Decision Boundaries

Because the output of the hypothesis function gives us $P(\mathbf{y}^{(i)} = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta})$, we should predict that $\mathbf{y}^{(i)} = 1 \iff h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq 0.5$. This means that if we want to visualise what our model will predict the class of a given datum is to plot the decision boundary of the hypothesis function given the optimised $\boldsymbol{\theta}$. This decision boundary is given by

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = 0.5 \quad (3.6)$$

or equivalently

$$\boldsymbol{\theta}^T \mathbf{x} = 0 \quad (3.7)$$

and it is the boundary that separates the regions where the model will predict \mathbf{y} is 1 and where it will predict \mathbf{y} is 0. For a small enough number of features this can be plotted and interpreted by a human to help give insight into what the model is doing.

3.3 Multi-Class Problems

We can also use logistic regression in problems where there are more than 2 classes. To do this, we use a ‘one-vs-all’ approach, whereby if we have n different classes, we train n different classifiers, one for each class. We then have n different hypotheses, $h_{\boldsymbol{\theta}}^{(i)}(\mathbf{x})$ which correspond to the classes $i \in \{1, 2, 3, \dots, n\}$. We train them such that $h_{\boldsymbol{\theta}}^{(i)}(\mathbf{x})$

is a binary classifier using a logistic regression modifier with the positive case being that \mathbf{x} is a member of class i and the negative case being that \mathbf{x} is a member of any other class. As such, to find the class \mathbf{x} is most likely to be a part of, we just compare the values of all the different $h_{\theta}^{(i)}(\mathbf{x})$, and see which is highest. If $h_{\theta}^{(j)}(\mathbf{x})$ outputs the highest value, we predict that \mathbf{x} is a member of class j .

3.4 Gradient Descent

To fit our logistic regression model, we need to use an optimisation algorithm such as gradient descent. It is worth noting that there are other algorithms which do the same thing faster, however their implementations are best left to those who specialise in numerical methods, however they only require the cost function and its gradient to be written by the user, so we only need to know the same things as we do for gradient descent. For logistic regression, the gradient of the cost function turns out to be the same as that of linear regression's cost function, except the hypothesis function is switched out, ie

$$\nabla_j J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \cdot \mathbf{x}_j^{(i)} \right] \quad (3.8)$$

$$\nabla J(\theta) = \frac{1}{m} \mathbf{X}^T (g(\mathbf{X}\theta) - \mathbf{y}). \quad (3.9)$$

3.5 Regularised Logistic Regression

Like with linear regression, we can include higher order features in our model to get a better fit model, however this can cause our model to overfit the training data. To solve this, we can once again use regularisation to penalise higher values for the weights of our parameters, to discourage over-fitting. For logistic regression, this makes the cost function

$$J(\theta) = \frac{1}{m} (-\mathbf{y}^T \log(g(\mathbf{X}\theta)) - (\mathbf{1} - \mathbf{y})^T \log(\mathbf{1} - g(\mathbf{X}\theta))) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2. \quad (3.10)$$

3.5.1 Gradient Descent for Regularised Logistic Regression

To use gradient descent, or other similar optimisation algorithms, for regularised logistic regression, we need to find the gradient of our new cost function, which turns out to be

$$\nabla_0 J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \right] \quad (3.11)$$

$$\nabla_j J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[(h_{\theta}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) \cdot \mathbf{x}_j^{(i)} \right] + \frac{\lambda}{m} \theta_j \quad \text{for } j \neq 0 \quad (3.12)$$

$$\nabla J(\theta) = \frac{1}{m} \mathbf{X}^T (g(\mathbf{X}\theta) - \mathbf{y}) + \frac{\lambda}{m} \begin{bmatrix} 0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}. \quad (3.13)$$