

Applied Machine Learning

Introduction

- Professor is a physicist in high-energy physics
- We will not go so much into theory
- The exam will be a ML project
- ML is the capacity of a computer to do a task without being explicitly programmed
- AI contains ML, which contains DL (deep learning)
 - ML started in 1980, DL in 2010
- Strong AI is really far
- ML can learn faster and with lower latency than humans
- It is useful for tasks that humans cannot or don't want to do
- Why today? Data available and Cloud computing
- ML can be supervised, unsupervised and reinforcement learning
- Supervised: I know some real solutions
 - It is a regression or classification problem
 - Regression: continuous
 - Classification: discrete
- Unsupervised: no label on the data
 - I use clustering algos
 - I want to find some structure in the data
 - I can get groups, but I don't know the meaning of these groups

Univariate linear regression

- I can define a cost function that measures the average distance of the real outcomes from my regression
- I want to choose the parameters θ s that minimize the cost function $J(\theta_1, \theta_2, \dots, \theta_n)$
 - In a linear regression the cost function has 2 parameters (!)
 - * Intercept and angular coefficient
- To minimise a function I can use a gradient descent algo
 - It is an iterative process
 - For now, only local minima, no global
 - It uses an aggressiveness factor α , which is how big every step is
 - * If too small it is too slow
 - * If it is too large I can miss a minimum
 - * α is referred to as a hyperparameter
 - It refers to the learning, not to the problem
 - When updating θ s, all of them must be updated simultaneously
- The minimization algo can be analytical or iterative
 - An analytical solution to univariate linear regression exists
 - In ML the analytical version does not scale well
 - GD is the iterative approach
- The iterative update of θ is done by subtracting to its previous value α times the partial derivative of the cost function with respect to θ
 - If the derivative is positive θ decreases, if negative increases, if 0 doesn't change
 - The magnitude of the change is proportional to the derivative at that point (!)
- In a linear regression the cost function is always a convex quadratic: the only minimum is the global minimum (!)
- Batch GD: start from any point and apply GD until I get to a minimum
 - It is batch since at every iteration I evaluate the cost function for the whole batch of datapoints

Multivariate linear regression

- The real world is multivariate (!)

- Nonetheless, univariate is useful for understanding concepts
- I have one θ for each x , plus θ_0
 - θ_0 is a bit uncomfortable, since it is different from the others (no x associated!)
 - To make things easier, I introduce $x_0 = 1$ that multiplies θ_0
 - This means that I have $n+1$ dimensional vectors if n is the number of independent variables
 - In this way, I have a vector of x s and a vector of θ s
 - I can represent the whole multivariate function as a product of the x vector with the traspose of the θ vector
 - $h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \boldsymbol{\theta}^T \mathbf{x}$
- The different variables can have different magnitudes, and I want to account for this
 - To correct, I will do feature scaling
 - I divide the data for the highest value for that variable
 - My data becomes all in the range 0-1
 - Outliers can skew my features: I remove them
 - More generally I want to be in the -1/+1 range since x_0 is already 1
 - I need to rescale also features which are really small
- A different way can be to do mean normalisation
 - I subtract the mean and divide for the range (max-min) or stdev

Learning rate

- The selection of α is important for determining if the GD converges, and if it does how much does it take
- How do I determine if the GD has converged?
 - I can decide a threshold decrease, i.e. if J decreases of less than 10^{-3} in one iteration I stop
- If I see a strange behaviour (divergence, bouncing around) the first thing to try is to decrease α
- But what values for α ?
 - First try in factor 10 steps: 0.0001, 0.001, 0.01, 1, 10, ...
 - Then go to a factor 3

Polynomial regression

- It is the simplest non-linear model but it can fit really complicated behaviours
- I can create features: instead of using x , why not e^x ?
 - I can make linear dependencies which are not linear
- I can reduce any polynomial regression to a linear by adding new features (!)
 - I can use x and x^2 instead of only x

Classification

- Classification problems can be binary or multiclass
- Linear regression is not good for pure classification problems
 - My problem is in nature not linear
 - I want an output in the range 0-1, not a continuous one
- Logistic regression: a classification algorithm
 - It is a sigmoid or logistic function that outputs in the 0-1 range
 - It is a function of the regression function itself $\theta^T x$
 - * $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$
 - I can interpret it as a probability of belonging to class $y=1$, given the measurement x and the parametrization θ
 - * $h_{\theta}(x) = p(y = 1|x, \theta)$
 - In general the logistic function takes any range of values, e.g. outputs of a function, and reports it in the range 0-1

- The output is the probability of the input belonging to class 1, and the probability of belonging to 0 is its complementary
- I am defining with the logistic a decision boundary that discriminates 1 and 0 outputs
- The decision boundary is not decided by the data, but by our hypothesis
 - It is a product of the model we use
- The decision boundary is not necessarily linear
 - By using higher order polynomials I can have circles and more complex boundaries
- The cost function for the logistic regression cannot be the argument of the logistic
 - If we apply GD on the initial function that is plugged into the logistic, there is no guarantee of convergence
 - This cost function is not convex (!)
- We can define this cost function for a single element y

$$cost(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & y = 1 \\ -\log(1 - h_{\theta}(x)) & y = 0 \end{cases}$$
- This can be rewritten as
 - $cost(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$
- The total cost function J is then
 - $J(\theta) = \frac{1}{n} \sum_{i=1}^m cost(h_{\theta}(x_i), y_i)$
 - $J(\theta) = \frac{1}{m} \sum_{i=1}^m -y_i \log(h_{\theta}(x_i)) - (1 - y_i) \log(1 - h_{\theta}(x_i))$
- The GD algorithm for classification is identical to that for linear regression
 - The only difference is the h itself, so our hypothesis
 - The process for optimizing the descent is the same

Alternatives to GD

- GD is not the only possibility, there is also conjugate gradients and other approaches
- Other approaches are more opaque, there are libraries that provide them but they are difficult to understand

Multiclass classification

- One-vs-all approach: I decompose the problem in several binary classifications
 - I assign a class to 1 and all the other datapoints to 0
 - I determine the decision boundary
 - I repeat with the second class and so on
 - Now we know the probability that a datapoint belongs to each of the classes
 - Our prediction is the class that gives me the highest probability

Overfitting

- I have overfitting when my model does not generalize
- How to reduce overfitting
 - Reduce the number of features
 - * This is risky since I can lose useful information
 - Tune down the weight of features
- I don't need to specify how small a feature should be (!)
- I can modify my cost function so that the cost for a feature is really high, and thus gets tuned down by the GD
 - I can add the square of the parameters that I want to tune down to the cost function
 - I add the square because the parameters can be negative (!)
 - In this way I penalise when they get too big
 - This term is added to the sum of squared distances of the previous cost function
- I get a cost function that is a tradeoff between fitting and avoidance of overfitting
- I can do this by introducing the regularization hyperparameter λ

- It is a multiplier to the sum of squares of all the parameters
 - * I am actually excluding θ_0 from this
 - It penalizes the cost function when parameters get too large
- By tuning λ I can modify the behaviour of my model
 - When λ is really large I go towards underfitting
 - When it is too small I have overfitting

Improving performance

- Do not over-optimize the model: if needed try to increase the amount of data
 - Not always easy!
- Another possibility: tune down or remove features!
- Maybe your dataset is not descriptive enough: more features!