
Contents

1	Introduction	3
1.1	Supervised learning	3
1.2	Unsupervised learning	3
1.3	Reinforcement learning	3
2	Supervised learning	5
2.1	Linear classifiers	5
2.2	Features	7
3	Logistic Regression	9
3.1	Linear logistic classifier	9
3.2	Gradient descent	10

Chapter 1

Introduction

1.1 Supervised learning

Given a dataset of pair

$$\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\} \quad (1.1)$$

we wish to establish a relationship between $x^{(i)}$ and $y^{(i)}$. Typically, $x^{(i)} \in \mathbb{R}^d$ is a representation of input, **feature representation**. Based on the format of the output we can have different types of supervised learning:

Classification when the set of possible values of $y^{(i)}$ is discrete (small finite set). If there two possible values then the classification problem is *binary* otherwise, it is called *multi-class*.

Regression when the set of possible values of $y^{(i)}$ is continuous (or a large finite set). That is, $y^{(i)} \in \mathbb{R}^k$.

1.2 Unsupervised learning

Given a dataset we wish to find some patterns or structures in it. There are several types of unsupervised learning

Density estimation The data is i.i.d from some distribution $p_X(x)$. The goal is to predict the probability $p_X(x^{(n+1)})$.

Clustering the goal is to find a partitioning of the sample data that groups together samples that are similar. Clustering is sometimes used in density estimation.

Dimensionality reduction the goal is to re-represent the same data in \mathbb{R}^l where $l < d$.

1.3 Reinforcement learning

The goal is to learn a mapping from input values to output values without a direct supervision signal. There is no training set specified *a priori*. Instead, the learning problem is framed as an agent interacting with an environment. Looking at input as our states and output as a transition between states, we can assign a reward value $r^{(i,j)}$ to each such transition. We aim to find a policy π that maximizes the long-term sum or average of rewards.

Chapter 2

Supervised learning

To predict, we come up with a **hypothesis**. A hypothesis is a parametrized function that maps input to output

$$y = h(x; \theta), \quad h \in \mathcal{H}, \theta \in \Theta$$

we then wish to find the parameters θ that matches our data well. One way to evaluate how well our hypothesis predicts is to introduce a **loss function** (or *cost function*), $L(a, a_h)$ where a, a_h are in the members output set and function assigns a value to how close our prediction a_h when the actual value is a . We wish that our hypothesis to have the least loss on new data.

$$\mathcal{E}_n(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}; \theta), y^{(i)})$$

One way to do this is minimize the loss function on the training data

$$\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}; \theta), y^{(i)})$$

There are several types of loss function

0-1 Loss

$$L(a, a_h) = \begin{cases} 0 & \text{if } a = a_h \\ 1 & \text{otherwise} \end{cases}$$

Squared loss

$$L(a, a_h) = (a - a_h)^2$$

Linear loss

$$L(a, a_h) = |a - a_h|$$

Asymmetric loss

The model we use, typically, select h and we need to minimize the loss (or any other optimization) on the θ so that our prediction *fits* data. To determine a good θ we need algorithms, *learning algorithms*.

2.1 Linear classifiers

A linear classifier has the following form

$$h(x; \theta, \theta_0) = \text{sign}(\theta^T x + \theta_0) \quad \theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$$

2.1.1 Random linear classifier

Algorithm 1: rand_lin_classifier (\mathcal{D}, k)

```

for  $j = 1 \rightarrow k$  do
     $\theta^{(j)} = \text{Random}(\mathbb{R}^d)$ 
     $\theta_0^{(j)} = \text{Random}(\mathbb{R})$ 
 $j^* = \operatorname{argmin}_{1 \leq j \leq k} \mathcal{E}\left(h\left(x, \theta^{(j)}, \theta_0^{(j)}\right)\right)$ 
return  $(\theta^{(j^*)}, \theta_0^{(j^*)})$ 

```

2.1.2 Perceptron

Algorithm 2: perceptron (\mathcal{D}, T)

```

 $\theta = 0$ 
 $\theta_0 = 0$ 
for  $t = 1 \rightarrow T$  do
    for  $i = 1 \rightarrow n$  do
        if  $y^{(i)}(\theta^T x^{(i)}) + \theta_0 \leq 0$  then
             $\theta = \theta + y^{(i)} x^{(i)}$ 
             $\theta_0 = \theta_0 + y^{(i)}$ 
return  $(\theta, \theta_0)$ 

```

By adding another dimension to our data set we can simplify our prediction to pass through the origin

$$x' = [x_1 \ \dots \ x_n \ 1], \quad \theta' = [\theta \ \theta_0]$$

$$\implies \theta'^T x' = \theta^T x + \theta_0$$

Therefore, one can also simplify the Line 2 to the following

Algorithm 3: perceptron (\mathcal{D}, T)

```

 $\theta = 0$ 
for  $t = 1 \rightarrow T$  do
    for  $i = 1 \rightarrow n$  do
        if  $y^{(i)}(\theta^T x^{(i)}) + \theta_0 \leq 0$  then
             $\theta = \theta + y^{(i)} x^{(i)}$ 
return  $\theta$ 

```

A dataset \mathcal{D} is **linearly separable**-through the origin if there is some θ such that

$$y^{(i)} \theta^T x^{(i)} > 0 \quad \forall i$$

The **margin** of a labeled data point (x, y) with respect to a separator (hyperplane) θ, θ_0 is

$$y \cdot \frac{\theta^T x + \theta_0}{\|\theta\|}$$

which basically quantifies how well θ approximates the data point (x, y) in a data set \mathcal{D} . Also the margin of \mathcal{D} w.r.t θ, θ_0 is the minimum of all margins:

$$\min_i y^{(i)} \cdot \frac{\theta^T x^{(i)} + \theta_0}{\|\theta\|}$$

Theorem 2.1 (Perceptron convergence theorem). *If there exists a vector θ^* such that the margin of database with respect to θ^* is greater than $\gamma > 0$ and then norm $\|x^{(i)}\| \leq R$ for some R then perceptron will make at most $\left(\frac{R}{\gamma}\right)^2$ updates/mistakes.*

Proof. put a increasing lower bound on the cosine of the angle. ■

2.2 Features

2.2.1 Transformation

As we saw we can transform linearly separable dataset to another linearly separable dataset but without an offset. What happens if the original dataset is not linearly separable? For example, *xor dataset*:

$$\mathcal{D} = \{((-1, -1), -1), ((-1, 1), 1), ((1, -1), -1), ((1, 1), 1)\}$$

is not linearly separable in 2 dimensions. A transformation that might be applicable here is **polynomial basis**. A polynomial basis transformation of order k , transforms a feature $x \in \mathbb{R}^d$ to

$$\phi(x) = (x_1^{\alpha_1} \dots x_d^{\alpha_d}), \quad \sum_{i=1}^d \alpha_i = k, \forall i, \alpha_i \geq 0$$

which has $\binom{k+d-1}{d-1}$ dimension.

2.2.2 Representation

we can represent a discrete feature as

1. numeric
2. thermometer code (a vector of m booleans where $1 \dots j$ bits are on and the rest are off)
3. one-hot (a vector of m booleans where j_{th} bit is on and the rest are off)
4. factoring (group information of a feature based its structure maybe)

For numeric feature we would like to standardize as follow

$$\tilde{x}_j = \frac{x_j - \bar{x}_j}{\sigma}$$

Chapter 3

Logistic Regression

In machine learning we wish to optimize a function like $J(\Theta)$. Usually a function in form

$$\begin{aligned} J(\Theta) &= \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}(h(x^{(i)}; \theta), y^{(i)}) \right) + \lambda R(\theta) \\ &= \mathcal{E}_n + \lambda R(\theta) \end{aligned}$$

where R is the **regularization function** and λ is a hyperparameter.

3.1 Linear logistic classifier

The problem of minimizing 0-1 Loss problem is NP-hard. A problem with sign is that incremental change are hard find because of the discrete nature of the function hence, to smooth out the sign function we use sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Equivalently, we want to make classifier that predict +1 when $\sigma(\theta^T x + \theta_0) > 0.5$ and -1 otherwise.

Loss on all data is inversely related to the probability that θ, θ_0 assign to the data. Assuming that points in data set are independent.

$$\begin{aligned} g^{(i)} &= \sigma(\theta^T x + \theta_0) \\ p^{(i)} &= \begin{cases} g^{(i)} & \text{if } y^{(i)} = 1 \\ 1 - g^{(i)} & \text{if } y^{(i)} = 0 \end{cases} \end{aligned}$$

and we wish to maximize the probability

$$\prod_{i=1}^n p^{(i)} = \prod_{i=1}^n (g^{(i)})^{y^{(i)}} (1 - g^{(i)})^{(1-y^{(i)})}$$

Using the log-likelihood

$$\Rightarrow \mathcal{L}_{LL}(p) = \sum_{i=1}^n y^{(i)} \log(g^{(i)}) + (1 - y^{(i)}) \log(1 - g^{(i)})$$

3.2 Gradient descent

Algorithm 4: gradient descent $(f, \nabla f, \theta_{\text{init}}, \eta, \epsilon)$

 $\theta^{(0)} = \theta_{\text{init}}$ $t = 0$ **repeat** $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla f(\theta^{(t-1)})$ **until** $|f(\theta^{(t)}) - f(\theta^{(t-1)})| < \epsilon$ **return** θ

Theorem 3.1. *If f is convex, for any desired accuracy ϵ there is some η such that gradient descent will converge to θ within ϵ of the optimum.*