

Regression

Motivation

$$H^L = \{h_{w,b} : \mathbf{w} \in \mathbb{R}^d \text{ and } b \in \mathbb{R}\} \quad |H^L| = \infty$$

$$h_{w,b}(x_i) = \langle \mathbf{w}, x_i \rangle + b$$

H^L is one of the most useful families of hypothesis classes.

Many models that are used in practice rely on linear predictors.

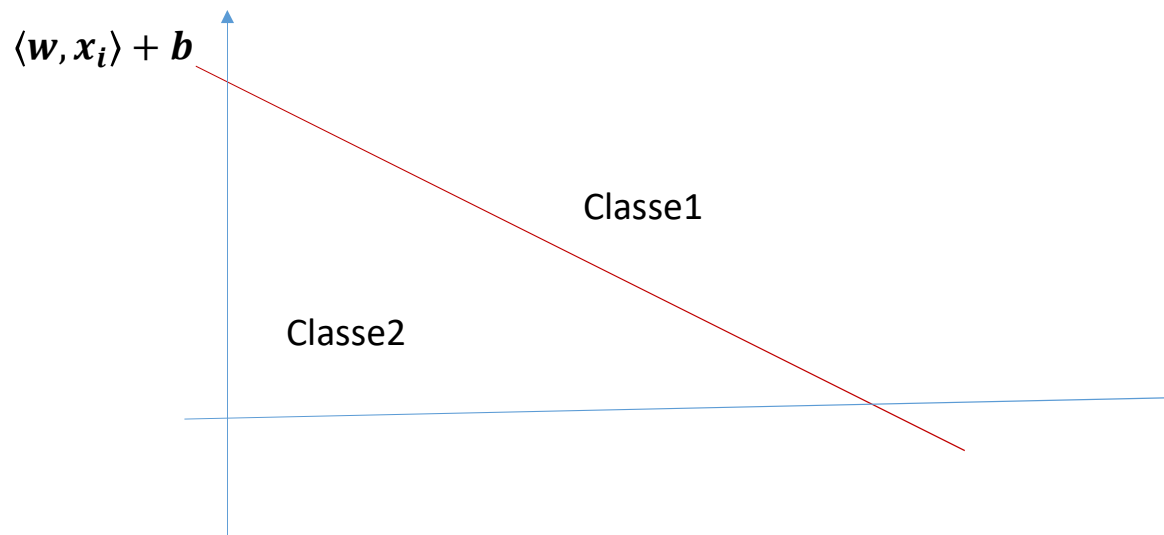
There exist many types of linear models, including:

- Perceptron(classification). $\mathbf{y} \in \{0, 1\}$
- Linear regression. $\mathbf{y} \in \mathbb{R}$
- Logistic regression. output *probability* $p \in (0, 1)$ and input $\mathbf{y} \in \{-1, +1\}$

Motivation

Perceptron(classification). $h_{w,b}(x_i) = \text{sign}(\langle w, x_i \rangle + b)$ $y \in \{-1, +1\}$

Linear regression. $y_i \in \mathbb{R}$ $h_{w,b}(x_i) = \langle w, x_i \rangle + b = y_i$



Motivation

Let's define the class of affine functions $L=L_d$: input = x

$$H^{L_d} = \{h_{w,b} : (\mathbf{w}, \mathbf{b}) \in \mathbb{R}^{d+1}\}, \quad |H^{L_d}| = \infty$$

Where:

$$h_{w,b}(x_i) = \langle \mathbf{w}, \mathbf{x}_i \rangle + b = \sum_{j=1}^d w_j x_i^j + b = y_i$$

To simplify the notation, we will integrate the bias as an extra coordinate into w :

$$x_i = (x_i^1, \dots, x_i^d) \in \mathbb{R}^d \rightarrow x_i = (1, x_i) \in \mathbb{R}^{d+1}, w \rightarrow (b, w) \in \mathbb{R}^{d+1}$$

$$h_w(x_i) = \sum_{j=0}^d w_j x_i^j \text{ avec } x_i^0 = 1 \text{ and } w_0 = b$$

Hence, the class of affine functions is called « homogenous affine functions »

$$H^{L_d} = \mathbf{L}_d = \{\mathbf{h}_w : \mathbf{w} \in \mathbb{R}^{d+1}\} \rightarrow |\mathbf{L}_d| = \infty$$

$$\mathbf{h}_w(x_i) = y_i$$

Motivation

Therefore, we can generate different hypothesis classes H^L , defining different models, by using the composition of φ over L_d such that: ($\mathbf{h}_w \in L_d$)

$$\varphi : \mathbb{R} \rightarrow Y$$

- Perceptron(classification):

$$\varphi_p(x) = \text{sign}(x) \text{ and } Y = \{-1, +1\}$$

$$H_p = \text{sing}(\varphi_p \circ L_d = \{\varphi_p \circ \mathbf{h}_w(\mathbf{x}) : \mathbf{h}_w \in L_d\})$$

- Linear regression:

$$\varphi_{reg}(x) = \text{Id}(x) = x \text{ and } Y = \mathbb{R}$$

$$H_{reg} = \varphi_{reg} \circ L_d = \{\varphi_p \circ \mathbf{h}_w(\mathbf{x}) : \mathbf{h}_w \in L_d\} = \{\mathbf{h}_w(\mathbf{x}) : \mathbf{h}_w \in L_d\}$$

Logistic regression:

$$\varphi_{sig}(x) = \frac{1}{1+e^{-x}} \text{ and } Y = \{-1, +1\}, \varphi_{sig}(\mathbf{h}_w(\mathbf{x}_j^i) = \sum_{i=0}^d w_i x_j^i) = \frac{1}{1+e^{-\sum_{i=0}^d w_i x_i}}$$

$$H_{sig} = \varphi_{sig} \circ L_d(\mathbf{x}) = \frac{1}{1 + e^{-\sum_{i=0}^d w_i x_i}}$$

Linear Regression

Definition:

Linear regression is a type of model used for regression tasks by studying the relationship between some explanatory variables and some real valued outcome.

Here we have: $x = (1, x) \in X$

$$X \subset \mathbb{R}^{d+1} \text{ for some } d$$

And

$$Y = \mathbb{R}$$

Objective:

Learn a linear predictor $h_w \in L_d$ that best approximate the relationship between our variables:

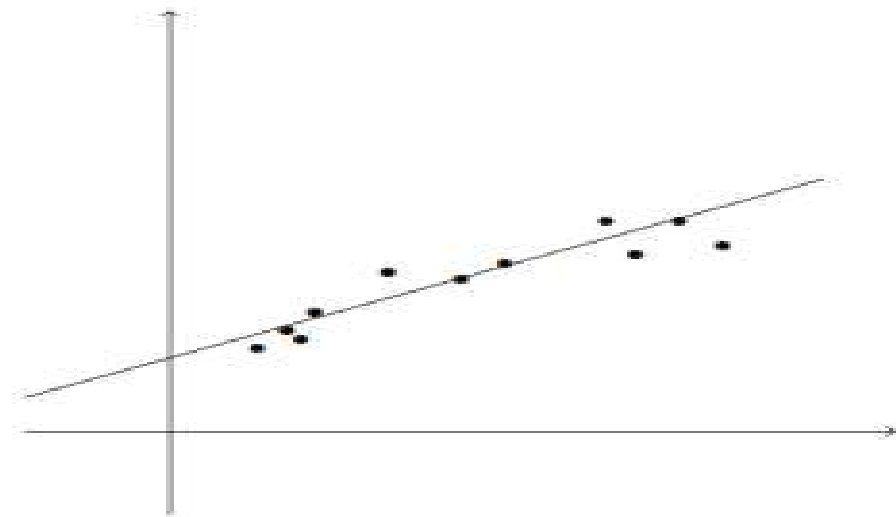
$$\begin{aligned} h_w: \mathbb{R}^{d+1} &\longrightarrow \mathbb{R} \\ x &\longrightarrow h_w(x) = w^t x = y \end{aligned}$$

Linear Regression

Example:

$$h_w(x) = w^t x_i = y_i \Rightarrow L_S(w) = \frac{1}{m} \sum_{i=1}^m d(w^t x_i, y_i) = \frac{1}{m} \sum_{i=1}^m (w^t x_i - y_i)^2 \approx \|w^t x - y\|^2$$
$$\min_{w \in \mathbb{R}^{d+1}} L_S(w) \approx 0$$

Predicting the weight
of a baby as a function of his age
and weight at birth.
Here, $d = 1$.



Linear Regression

The hypothesis class for linear regression model:

In linear regression model, we have:

$$\varphi_{regL}(x) = Id(x) = x \Rightarrow \varphi_{regL} \circ \mathbf{h}_w(x) = Id(\mathbf{h}_w(x)) = \mathbf{h}_w(x)$$

The hypothesis class of linear regression predictors is simply the set of linear functions:

$$H_{regL} = \varphi \circ L_d = L_d = \{\varphi_{regL} \circ \mathbf{h}_w : \mathbf{h}_w \in L_d\}$$

$$\mathbf{H}_{reg} = \{\mathbf{h}_w : x \mapsto \langle \mathbf{w}, x \rangle : \mathbf{w} \in \mathbb{R}^{d+1}\} \rightarrow |\mathbf{H}_{reg}| \approx \infty$$

Best Metric

- $\|w^t x - y\|_2 = \sqrt{\sum_{i=1}^{\infty} (w^t x_i - y_i)^2} \leftrightarrow \sum_{i=1}^{\infty} (w^t x_i - y_i)^2$
- $\|w^t x - y\|_1 = \sum_{i=1}^{\infty} |w^t x_i - y_i|$
- $\|w^t x - y\|_{\infty} = \max(|w^t x_i - y_i|)$
- $L_S(w) = \frac{1}{m} \sum_{i=1}^m d(w^t x_i, y_i) \rightarrow \min_{w \in \mathbb{R}^{d+1}} L_S(w) \approx 0$
- $d_2(w^t x_i, y_i) = (w^t x_i - y_i)^2$
- $d_1(w^t x_i, y_i) = |w^t x_i - y_i| \rightarrow ? \text{ diff: } 1, 2$
- $\varepsilon = 0,02$
- $\text{Min } f(x) = |x| \rightarrow |x| = \varepsilon = 0,02 \rightarrow x = \pm 0,02$ (subgradient algorithm by yourself)
- $\text{Min } g(x) = x^2 \rightarrow x^2 = \varepsilon = 0,02 \rightarrow x = \sqrt{0,02} = 0,14$
- $|w^t x_i - y_i| = 0,0002 \rightarrow w^t x_i - y_i = \pm 0,0002$
- $(w^t x_i - y_i)^2 = 0,0002 \rightarrow w^t x_i - y_i = \sqrt{0,0002} = 0,014$

Linear Regression

The loss function for linear regression model:

It measures how much the model should be penalized for the discrepancy between $h_w(x)$ and y . One common way is to use the squared-loss function:

$$0 \approx d(h_w(x), y) = l(h_w, (x, y)) = (h_w(x) - y)^2$$

For this loss function, the empirical risk is called the Mean Squared Error:

$$L_S(h_w) = E_{\text{empi}} \left(l(h_w, (x, y)) \right) = \frac{1}{m} \sum_{i=1}^m (h_w(x_i) - y_i)^2$$

Notice:

There are a variety of other loss functions that one can use, for example, the absolute value loss function:

$$l(h_w, (x, y)) = |h_w(x) - y| \Rightarrow \partial l(h_w, (x, y)) \text{ is a set if } h_w(x) = y$$

Linear Regression

The learning algorithm for linear regression model:

The learning algorithm follows ERM_H learning rule.

Least squares:

Least squares is the algorithm that solves the ERM_H problem for the hypothesis class of linear regression predictors with respect to squared loss.

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} L_S(h_w) = \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} \left(\frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \right)$$

To solve this problem, we calculate the gradient of the objective function and compare it to zero. That is, we need to solve:

$$\nabla L_S(h_w) = \frac{2}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i) x_i = 0 \rightarrow \nabla^2 L_S(h_w) = \frac{2}{m} \sum_{i=1}^m x_i (x_i)^T = \text{cte}$$

Linear Regression

We can rewrite the problem as the problem :

$$\nabla L_S(h_w) = \frac{2}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i) x_i = 0 \Rightarrow \sum_{i=1}^m (\langle w, x_i \rangle x_i) - \sum_{i=1}^m y_i x_i = 0$$

- $\Rightarrow Aw - b = 0 \Rightarrow Aw = b$

- $\Rightarrow \nabla^2 L_S(h_w) = A$

Where:

$$A = \left(\sum_{i=1}^m x_i \cdot x_i^T \right), \quad b = \sum_{i=1}^m y_i x_i$$

Linear Regression

$$Aw = b$$

Or in matrix $A(d \times d)$ (symetric) form: $A = (\sum_{i=1}^m x_i \cdot x_i^T)$, $b = \sum_{i=1}^m y_i x_i$

$$\bullet A^T = A = XX^T = \begin{pmatrix} \vdots & & \vdots \\ x_1 & \dots & x_m \\ \vdots & & \vdots \end{pmatrix} \begin{pmatrix} \vdots & & \vdots \\ x_1 & \dots & x_m \\ \vdots & & \vdots \end{pmatrix}^T : X = \begin{pmatrix} \vdots & & \vdots \\ x_1 & \dots & x_m \\ \vdots & & \vdots \end{pmatrix}$$

And

$$\bullet b = \begin{pmatrix} \vdots & & \vdots \\ x_1 & \dots & x_m \\ \vdots & & \vdots \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$$

If A is invertible, then the solution to the ERM problem is:

$$w = A^{-1}b$$

Example 1

- $x_1^T = (1,0,0), x_2^T = (1,1,0), x_3^T = (0,1,0)$
- $A = \left(\sum_{i=1}^3 x_i \cdot x_i^T\right) = x_1 \cdot x_1^T + x_2 \cdot x_2^T + x_3 \cdot x_3^T$
- $= \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (1,0,0) + \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} (1,1,0) + \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} (0,1,0)$
- $= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix}$
- A isn't inve

Example2

- $x_1^T = (1,0,0), x_2^T = (1,1,0), x_3^T = (0,0,1)$
- $A = \left(\sum_{i=1}^3 x_i \cdot x_i^T\right) = x_1 \cdot x_1^T + x_2 \cdot x_2^T + x_3 \cdot x_3^T$
- $= \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (1,0,0) + \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} (1,1,0) + \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} (0,0,1)$
- $= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$
- A is inv

Linear Regression

If A is not invertible, we require a few standard tools from linear algebra.

A is not invertible when the training data do not cover the entire space of \mathbb{R}^d .

Even if A is not invertible, we can always find a solution to the system:

$$Aw = b$$

because b is in the range of A .

Indeed, since A is symmetric, then we can write it using its eigenvalue decomposition as:

$$A = VDV^T$$

Where:

D is a diagonal matrix.

V is an orthonormal matrix (because $V^T V = I$ which is a $d \times d$ matrix).

Linear Regression

Let's define D^+ to be the diagonal matrix such that:

$$\begin{cases} D_{i,i}^+ = 0 & \text{if } D_{i,i} = 0 \\ D_{i,i}^+ = \frac{1}{D_{i,i}} & \text{if } D_{i,i} \neq 0 \end{cases} \Rightarrow DD^+ = I \setminus D_{i,i} = 0?$$

Now, define:

$$A^+ = VD^+V^T \text{ and } \hat{\mathbf{w}} = A^+ \mathbf{b}$$

Let v_i denote the i th column of V . Then we have:

$$A\hat{\mathbf{w}} = AA^+ \mathbf{b} = VDV^TVD^+V^T \mathbf{b} = VDD^+V^T \mathbf{b} = VV^T \mathbf{b} = \sum_{i: D_{i,i} \neq 0} v_i v_i^T \mathbf{b}$$

This means that $A\hat{\mathbf{w}}$ is the projection of \mathbf{b} on the space of vectors v_i for which $D_{i,i} \neq 0$.

Linear Regression

Since the linear space of $(x_1, \dots, x_m) \in \mathbb{R}^m$ is the same as the linear space of those $\{v_i\}$.

And, since b is in the linear space of x_i .

We obtain that:

$$A\hat{w} = b$$

Then, \hat{w} is a solution of $Aw = b$.

Linear Regression for polynomial regression tasks $x \in \mathbb{R}$

Some learning tasks call for nonlinear predictors, such as polynomial predictors.

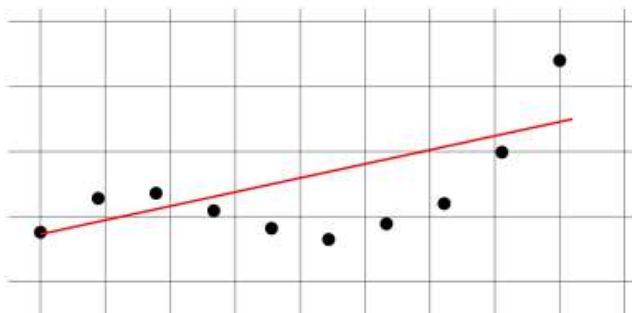
Let's consider a one dimensional(one feature) polynomial function of degree n : A_α $n = \alpha$

- $P_1(x) = w_0 + w_1x$
- $P_n(x) = w_0 + w_1x + w_2x^2 + \dots + w_nx^n \Rightarrow P_w(Z) = w_0 * 1 + w_1z_1 + w_2z_2 + \dots + w_nz_n$

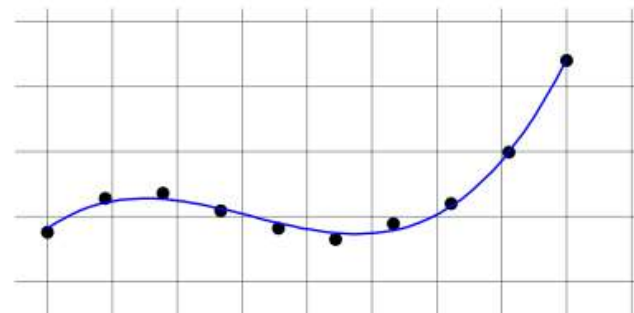
With $z = (1, z_1, \dots, z_n) \leftarrow z = (1, x, \dots, x^n) \in \{1\} \times \mathbb{R}^n \Leftrightarrow \text{space}(x) \simeq \text{space}(z)$

Where $w = (w_0, \dots, w_n)$ is a vector of coefficients of size $n + 1$.

$$P_1(x) = w_0 + w_1x$$



$$P_3(x) = w_0 + w_1x + w_2x^2 + w_3x^3$$



Linear Regression for polynomial regression tasks

We will focus on the class of one dimensional, n-degree, polynomial regression hypotheses. Therefore, the class of polynomial hypotheses is:

$$H_{poly}^n = \{P_{w,n=\alpha}: X \subseteq \mathbb{R} \mapsto \mathbb{R}: n \in \mathbb{N}^*, w \in \mathbb{R}^{n+1}\} \Rightarrow |H_{poly}^n| = \infty, A_\alpha = P_{w,n=\alpha}$$

Where $P_{w,n=\alpha}$ is a one dimensional polynomial of degree n , parameterized by a vector of coefficients (w_0, \dots, w_n) .

In that case, we have:

$$X \subseteq \mathbb{R} \quad \text{and} \quad Y \subseteq \mathbb{R}$$

One way to learn the class H_{poly}^n is by reduction to the problem of linear regression.

Linear Regression for polynomial regression tasks

To translate a polynomial regression problem to a linear regression problem, we define the mapping:

$$\psi_n: \mathbb{R} \rightarrow \mathbb{R}^{n+1}$$

Such that:

$$\psi_n(x) = (1, x, x^2, \dots, x^n)^T$$

Then, we have that: $P_{w,n} \circ \psi_n(x) = P_{w,n}(\psi_n(x))$

$$P_{w,n}(\psi_n(x)) = w_0 + w_1x + w_2x^2 + \dots + w_nx^n = \langle w, \psi_n(x) \rangle \Rightarrow \nabla_w P_{w,n} = \psi_n(x)$$

Finally, we can find the optimal vector of coefficients w by using the **Least Squares Algorithm**.

Logistic Regression

Definition:

Logistic regression is a type of model used for classification tasks by studying the relationship between some explanatory variables and some binary outcome.

Here we have:

$$X \subset \mathbb{R}^d \text{ for some } d \text{ and } Y = \{-1, +1\}$$

Objective:

Learn a linear predictor that best approximate the relationship between our variables:

$$h_w: \mathbb{R}^d \rightarrow [0,1]$$

We can interpret $h_w(x)$ as the probability that the label of x is 1 or -1 :

$$\mathbf{P(y|x) = P(y = 1 \vee y = -1|x) = P(y = 1|x) + P(y = -1|x)=1}$$

$$\mathbf{h_w(x) = P(y = 1|x)=1-P(y = -1|x)}$$

Reminder

we have:

1. $X \subset \mathbb{R}^d$ for some d and $Y = \{-1, +1\}$
2. $X \subset \mathbb{R}^d$ for some d and $Y \in \mathbb{R}$
3. $X \subset \mathbb{R}^d$ for some d and $Y \in [0,1]$

- If $A \cap B = \emptyset \Leftrightarrow A$ and B are disjoint $\Rightarrow P(A \cup B) = P(A) + P(B)$
- A and B are independent $\Leftrightarrow P(A \cap B) = P(A).P(B) \Leftrightarrow P(A|B) = P(A)$

$X \subset \mathbb{R}^d$ for some d and $y \in Y = \{-1, +1\} \rightarrow y = -1$ or $y = 1$

• **Output:**

$$\begin{array}{ccc} \varphi_{sig} \circ h_w & \mathbb{R}^d & \rightarrow [0, 1] \\ \mathbf{x} & \rightarrow & \varphi_{sig}(h_w(\mathbf{x})) = \mathbf{P}(\mathbf{y}|\mathbf{x}) \end{array}$$

- $\varphi_{sig}(h_w(\mathbf{x})) = \mathbf{P}(\mathbf{y} = 1|\mathbf{x}) \Rightarrow \mathbf{P}(\mathbf{y} = -1|\mathbf{x})$
- $\varphi_{sig}(h_w(\mathbf{x})) = \mathbf{P}(\mathbf{y} = -1|\mathbf{x}) \Rightarrow \mathbf{P}(\mathbf{y} = 1|\mathbf{x})$
- $\mathbf{h}_w(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = y(y\mathbf{w}^T \mathbf{x} > 0) \rightarrow \mathbf{w} \in \mathbb{R}^{d+1}$
- $\varphi_{sig}(\mathbf{x}) = \frac{1}{1+e^{-x}} \in [0, 1]$

Logistic Regression

The hypothesis class for logistic regression model:

In logistic regression model, we have:

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

The hypothesis class of logistic regression predictors is the composition of a sigmoid function over the set of linear functions:

$$H_{sig} = \varphi \circ L_d$$

$$H_{sig} = \{\varphi_{sig}(\mathbf{h}_w): x \mapsto \varphi_{sig}(\langle w, x \rangle) = \frac{1}{1 + e^{-\langle w, x \rangle}} : w \in \mathbb{R}^{d+1}\}$$

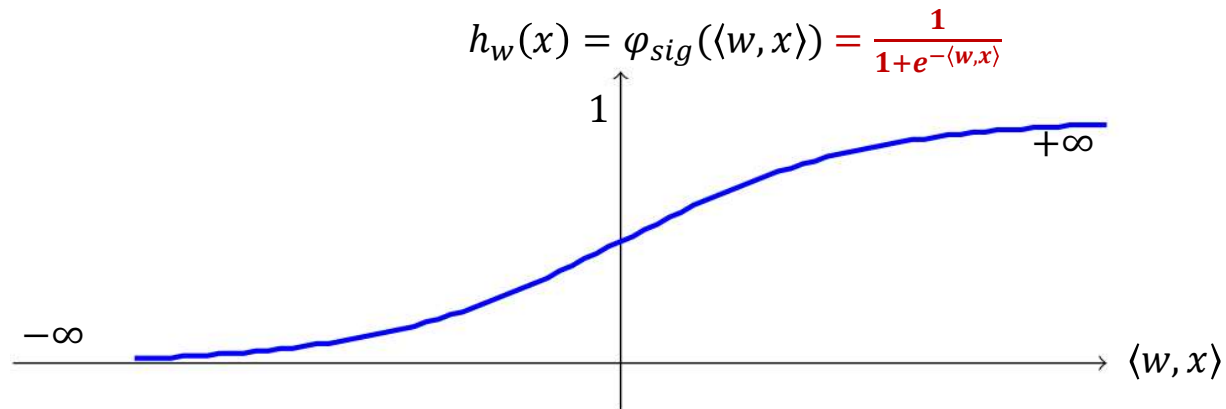
$$|H_{sig}| = \infty$$

Logistic Regression

The name « sigmoid » means «S-shaped », referring to the plot of this function shown in the figure:

• $e^{+\infty} = +\infty, e^{-\infty} = 0 \Rightarrow$

- if $\langle w, x \rangle = +\infty \rightarrow -\langle w, x \rangle = -\infty \rightarrow \frac{1}{1+e^{-\infty}} = 1$
- if $\langle w, x \rangle = -\infty \rightarrow -\langle w, x \rangle = +\infty \rightarrow \frac{1}{1+e^{+\infty}} = 0$



Logistic Regression

Logistic regression Vs Perceptron:

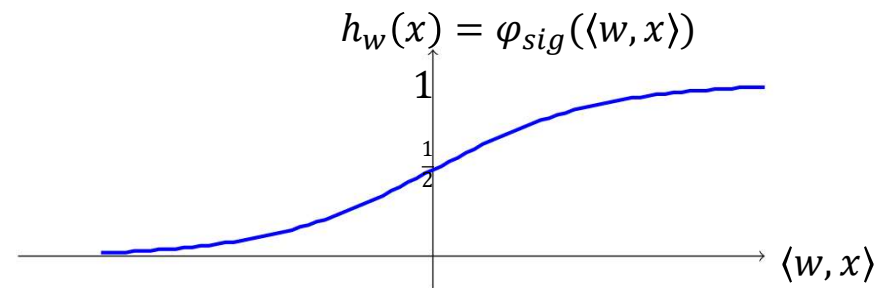
Whenever, $|\langle w, x \rangle|$ is large, the predictions of logistic regression hypothesis and perceptron hypothesis are similar.

However, whenever $|\langle w, x \rangle|$ is close to zero, we have that:

$$\varphi_{sig}(\langle w, x \rangle) \approx \frac{1}{2} \text{ and } \varphi_p(\langle w, x \rangle) = \text{sign}(\langle w, x \rangle)$$

The logistic regression hypothesis is not sure about the value of the label.

The perceptron hypothesis always outputs a deterministic prediction $\{-1, +1\}$, even if $|\langle w, x \rangle|$ is very close to zero.



Logistic Regression

The loss function for logistic regression model:

It measures how bad it is to predict some $h_w(x) \in [0,1]$ given that the true label is $y = \{\pm 1\}$.

Clearly, we want that:

$$P(y|x) = \begin{cases} h_w(x) & \text{if } y = +1 \\ 1 - h_w(x) & \text{if } y = -1 \end{cases} \Rightarrow P(y|x) = P(y = 1|x) + P(y = -1|x) = 1$$

to be large.

We have:

$$P(y = 1|x) = h_w(x) = \frac{1}{1+e^{-\langle w, x \rangle}} \quad \text{and} \quad P(y = -1|x) = 1 - h_w(x) = \frac{1}{1+e^{\langle w, x \rangle}}$$

Generally:

$$P(y|x) = \frac{1}{1 + e^{-y\langle w, x \rangle}}$$

Logistic Regression

It is clear that the loss function will increase monotonically if the probability $P(y|x)$ decreases.

This implies that, it will increase monotonically if $1 + e^{-y\langle w, x \rangle}$ increases.

Therefore, the loss function used in logistic regression penalizes h_w based on the log of $1 + e^{-y\langle w, x \rangle}$, that is:

$$l(h_w, (x, y)) = \log(1 + e^{-y\langle w, x \rangle})$$

(recall that the log is a monotonic function).

Therefore, given a training set $S = (x_1, y_1), \dots, (x_m, y_m)$, the **ERM** problem associated with logistic regression is:

$$\operatorname{argmin}_w L_S(h_w) = \operatorname{argmin}_{w \in \mathbb{R}^d} \left(\frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i \langle w, x_i \rangle}) \right)$$

$$\mathbf{w} \in \mathbb{R}^{d+1}, \quad \mathbf{x} = (1, x_1, \dots, x_d)$$

- $\ln(x)' = \frac{1}{x}$ if $x \neq 0$

- $l(\mathbf{h}_{\mathbf{w}}, (\mathbf{x}, y)) = \log(1 + e^{-y\langle \mathbf{w}, \mathbf{x} \rangle})$

- $\nabla_{\mathbf{w}} l(\mathbf{h}_{\mathbf{w}}, (\mathbf{x}, y)) = \begin{pmatrix} \frac{\partial \log(1 + e^{-y\langle \mathbf{w}, \mathbf{x} \rangle})}{\partial w_0} \\ \vdots \\ \frac{\partial \log(1 + e^{-y\langle \mathbf{w}, \mathbf{x} \rangle})}{\partial w_d} \end{pmatrix} = \frac{1}{1 + e^{-y\langle \mathbf{w}, \mathbf{x} \rangle}} \begin{pmatrix} -ye^{-y\langle \mathbf{w}, \mathbf{x} \rangle} \\ \vdots \\ -y\mathbf{x}_d e^{-y\langle \mathbf{w}, \mathbf{x} \rangle} \end{pmatrix}$

- $\frac{1}{m} \sum_{i=1}^m \nabla_{\mathbf{w}} \log(1 + e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}) = \frac{1}{m} \sum_{i=1}^m \frac{1}{1 + e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle}} \begin{pmatrix} -y_i e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle} \\ \vdots \\ -y_i \mathbf{x}_{i,d} e^{-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle} \end{pmatrix}$

- $\nabla_{\mathbf{w}}^2 l(\mathbf{h}_{\mathbf{w}}, (\mathbf{x}, y))$

Logistic Regression

Notice:

It is clear that the loss function of the logistic regression is a convex function with respect to w .

So, the ERM_H problem for logistic regression model can be solved using a gradient descent algorithm.