

Machine Learning & Pattern Recognition

SONG Xuemeng

sxmustc@gmail.com

<http://xuemeng.bitcron.com/>

Linear Regression

age	23 years
annual salary	NTD 1,000,000
year in job	0.5 year
current debt	200,000

Training dataset: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$;

Features of the i -th customer: $\mathbf{x}_i = (x_{i1} \ x_{i2} \ \dots \ x_{id})^T$; (Column vector)

The **ground truth** of the credit limit for the i -th customer: $y_i \in \mathbb{R}$.

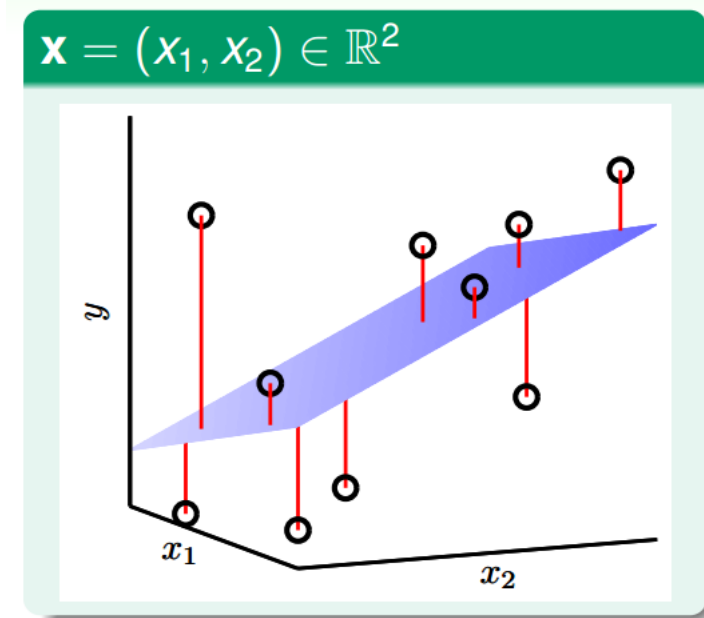
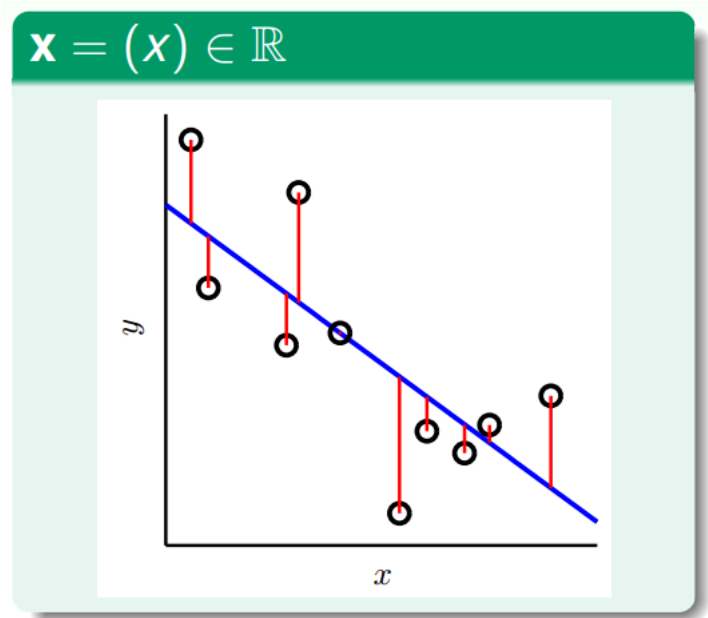
Linear regression: $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b = \sum_{j=1}^d w_j x_{ij} + b$, where $\mathbf{w} = (w_1 \ w_2 \ \dots \ w_d)^T \in \mathbb{R}^d$

For simplicity, the bias b can be merged into the weight \mathbf{w} :

$$h(\mathbf{x}_i) = \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i \quad \begin{aligned} \hat{\mathbf{w}} &= (b; \mathbf{w}) = (b \ w_1 \ w_2 \ \dots \ w_d) \in \mathbb{R}^{d+1} \\ \hat{\mathbf{x}}_i &= (1; x_{i1}; x_{i2}; \dots; x_{id}) \in \mathbb{R}^{d+1} \end{aligned}$$

Linear Regression

Linear regression hypothesis: $h(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i = \sum_{j=0}^d w_j x_{ij}$, $x_{i0} = 1$



Linear regression: find lines/hyperplanes with small **residuals**

Empirical Error

We usually prefer to minimize the objective function where the expectation is taken across the **data generating distribution** p_{data} rather than just over the finite training set:

$$J^*(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim p_{data}} L(h(\mathbf{x}, \boldsymbol{\theta}), y)$$

However, in most cases, we do not know p_{data} but only have a training set of samples. One simplest way to convert the machine learning problem back into an optimization problem is to minimize the expected loss on the training set.

$$J(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{P}_{data}} L(h(\mathbf{x}, \boldsymbol{\theta}), y)$$

Replacing the **true** distribution $p_{data}(\mathbf{x}, y)$ with the **empirical** distribution $\hat{P}_{data}(\mathbf{x}, y)$ defined by the training set.

Linear Regression

Popular/historical error measure:

squared error $err(\hat{y} - y) = (\hat{y} - y)^2$

$$E(\mathbf{w}) = \sum_{i=1}^m \frac{(h(\mathbf{x}_i) - y_i)^2}{\mathbf{w}^T \mathbf{x}_i}$$

Next: How to minimize $E(\mathbf{w})$?

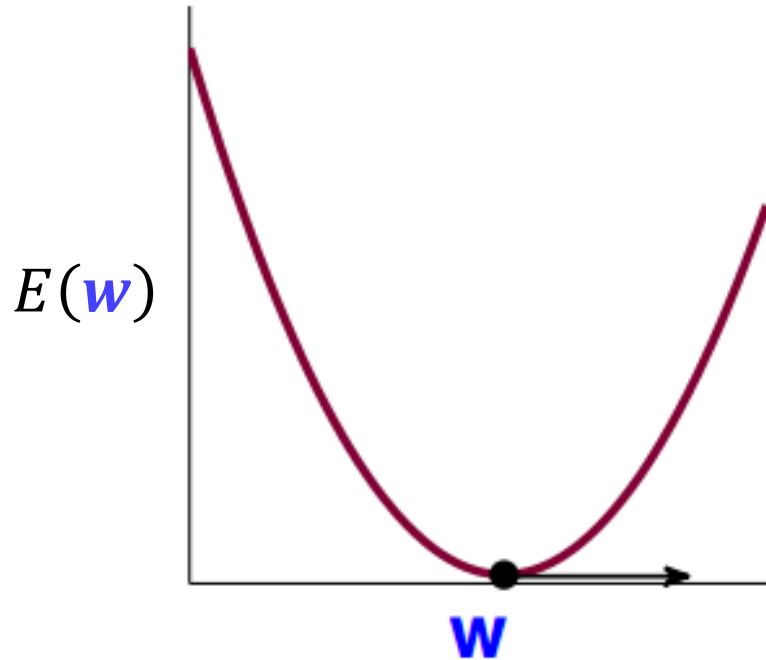
Matrix Form of $E(\mathbf{w})$

$$\begin{aligned} E(\mathbf{w}) &= \sum_{i=1}^m (h(\mathbf{x}_i) - y_i)^2 = \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i)^2 = \sum_{i=1}^m (\mathbf{x}_i^T \mathbf{w} - y_i)^2 \\ &= \left\| \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} - y_1 \\ \mathbf{x}_2^T \mathbf{w} - y_2 \\ \vdots \\ \mathbf{x}_m^T \mathbf{w} - y_m \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} - & - & \mathbf{x}_1^T & - & - \\ - & - & \mathbf{x}_2^T & - & - \\ & & \vdots & & \\ - & - & \mathbf{x}_m^T & - & - \end{bmatrix} \mathbf{w} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \right\|^2 \\ &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{md} \end{pmatrix} \in \mathbb{R}^{m \times (d+1)}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

Matrix Form of $E(\mathbf{w})$

$$\min E(\mathbf{w}) = \min \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$



- $E(\mathbf{w})$: continuous, differentiable, convex
- Necessary condition of 'best' \mathbf{w} .

$$\nabla E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_0}(\mathbf{w}) \\ \frac{\partial E}{\partial w_1}(\mathbf{w}) \\ \vdots \\ \frac{\partial E}{\partial w_d}(\mathbf{w}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \checkmark \text{ Not possible to 'roll down'}$$

Task: find the \mathbf{w}^* such that $\nabla E(\mathbf{w}^*) = 0$

The Gradient $\nabla E(\mathbf{w})$

$$\min_{\mathbf{w}} E(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = \underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}_A - 2 \underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{y}}_b + \underbrace{\mathbf{y}^T \mathbf{y}}_c$$

One w only

$$E(w) = (aw^2 - 2bw + c)$$

$$\nabla E(w) = 2aw - 2b$$

Vector w

$$E(\mathbf{w}) = (\mathbf{w}^T \mathbf{A} \mathbf{w} - 2\mathbf{w}^T \mathbf{b} + c)$$

$$\nabla E(\mathbf{w}) = 2\mathbf{A}\mathbf{w} - 2\mathbf{b}$$

$$\nabla E(\mathbf{w}) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y})$$

Optimal Linear Regression Weights

Task: find \mathbf{w}^* such that $\nabla E(\mathbf{w}^*) = 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$

Invertible/positive definite $\mathbf{X}^T \mathbf{X}$

- Unique solution

$$\mathbf{w}^* = (\underbrace{\mathbf{X}^T \mathbf{X}}_{\text{pseudo-inverse } \mathbf{X}^\dagger})^{-1} \mathbf{X}^T \mathbf{y}$$

pseudo-inverse \mathbf{X}^\dagger

- Often the case because
 $N \gg d + 1$

Singular $\mathbf{X}^T \mathbf{X}$

- Many optimal solutions
- One of the solution
 - Define \mathbf{X}^\dagger in other ways

Linear Regression Algorithm

1. From \mathcal{D} , construct input matrix \mathbf{X} and output vector \mathbf{y} by

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{md} \end{pmatrix} \in \mathbb{R}^{m \times (d+1)}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

2. Calculate pseudo-inverse

$$\mathbf{X}^\dagger \in \mathbb{R}^{(d+1) \times m}$$

3. Return $\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y} \in \mathbb{R}^{(d+1)}$

Simple and efficient with **good \mathbf{X}^\dagger routine**

Exercise

After getting \mathbf{w}^* , we can calculate the predictions $\hat{y}_n = (\mathbf{w}^*)^T \mathbf{x}_n$. If all \hat{y}_n are collected in a vector $\hat{\mathbf{y}}$ similar to how we form \mathbf{y} , what is the matrix formula of $\hat{\mathbf{y}}$?

- 1 \mathbf{y}
- 2 $\mathbf{X}\mathbf{X}^T \mathbf{y}$
- 3 $\mathbf{X}\mathbf{X}^\dagger \mathbf{y}$
- 4 $\mathbf{X}\mathbf{X}^\dagger \mathbf{X}\mathbf{X}^T \mathbf{y}$

Exercise

After getting \mathbf{w}^* , we can calculate the predictions $\hat{y}_n = (\mathbf{w}^*)^T \mathbf{x}_n$. If all \hat{y}_n are collected in a vector $\hat{\mathbf{y}}$ similar to how we form \mathbf{y} , what is the matrix formula of $\hat{\mathbf{y}}$?

- 1 \mathbf{y}
- 2 $\mathbf{X}\mathbf{X}^T\mathbf{y}$
- 3 $\mathbf{X}\mathbf{X}^\dagger\mathbf{y}$
- 4 $\mathbf{X}\mathbf{X}^\dagger\mathbf{X}\mathbf{X}^T\mathbf{y}$

Reference Answer: 3

Note that $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}^*$. Then, a simple substitution of \mathbf{w}^* reveals the answer.

Heart Attack Prediction Problem

age	40 years
gender	male
blood pressure	130/85
cholesterol level	240
weight	70

heart disease? **yes**

Binary classification:

Ideal $f(\mathbf{x}) = \text{sign}(p(+1|\mathbf{x}) - 0.5) \in \{-1, +1\}$

Heart Attack Prediction Problem

age	40 years
gender	male
blood pressure	130/85
cholesterol level	240
weight	70

heart attack? **80% risk**

'Soft' Binary classification:

$$f(x) = p(+1|x) \in [0,1]$$

Soft Binary classification:

Target function $f(x) = p(+1|x) \in [0,1]$

Ideal data

$$\begin{pmatrix} \mathbf{x}_1, y'_1 = 0.9 = P(+1|\mathbf{x}_1) \\ \mathbf{x}_2, y'_2 = 0.2 = P(+1|\mathbf{x}_2) \\ \vdots \\ \mathbf{x}_N, y'_N = 0.6 = P(+1|\mathbf{x}_N) \end{pmatrix}$$

Actual data

$$\begin{pmatrix} \mathbf{x}_1, y_1 = \circ \sim P(y|\mathbf{x}_1) \\ \mathbf{x}_2, y_2 = \times \sim P(y|\mathbf{x}_2) \\ \vdots \\ \mathbf{x}_N, y_N = \times \sim P(y|\mathbf{x}_N) \end{pmatrix}$$

Same data as hard binary classification, different **target function**

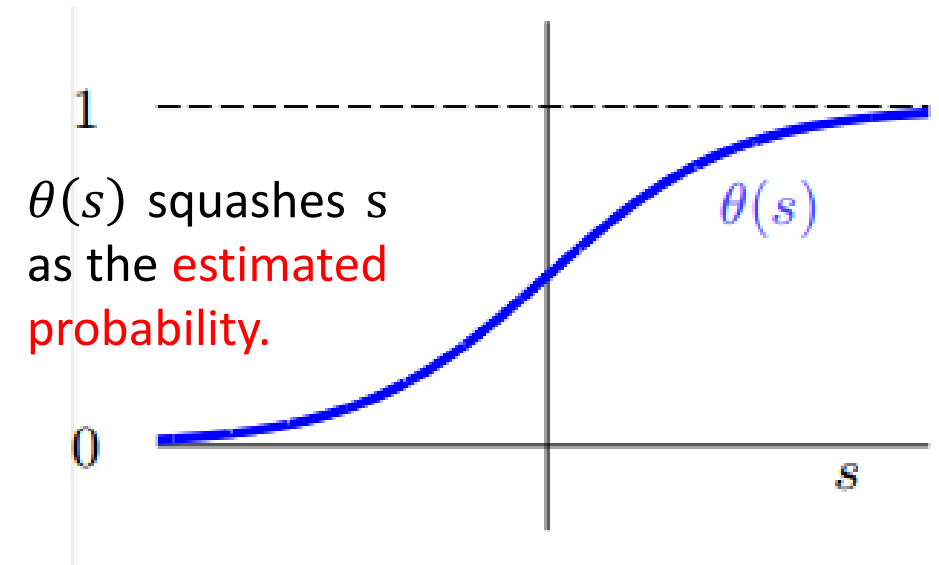
Logistic Hypothesis

age	40 years
gender	male
blood pressure	130/85
cholesterol level	240

Let $\mathbf{x}_i = (\mathbf{x}_{i0}, x_{i1}, x_{i2}, \dots, x_{id})$ be the features of the patient, calculate a weighted 'risk score':

$$s = \sum_{j=0}^d w_j x_{ij} = \mathbf{w}^T \mathbf{x}_i,$$

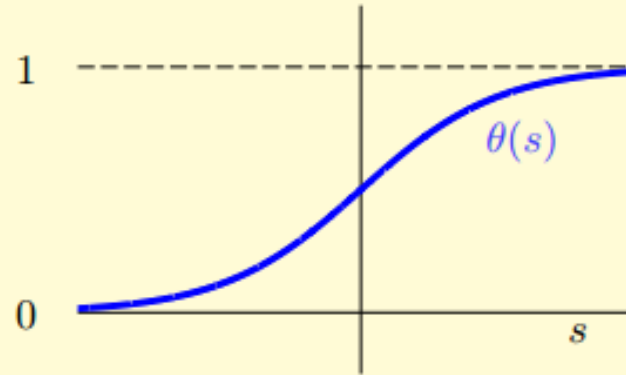
Convert the score to estimated probability by logistic function $\theta(s)$.



$$\text{Logistic hypothesis: } h(\mathbf{x}_i) = \theta(\mathbf{w}^T \mathbf{x}_i)$$

Logistic Function

$$\theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$



smooth, monotonic,
sigmoid function of s

Bound	$\theta(s) \in [0,1]$	$\theta(-\infty) = 0$	$\theta(0) = 0.5$	$\theta(\infty) = 1$
Symmetric	$1 - \theta(s) = \theta(-s)$			
Gradient	$\theta'(s) = \theta(s)(1 - \theta(s))$			

Logistic regression use $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$ to approximate the target $f(\mathbf{x}) = p(+1|\mathbf{x})$

Exercise

Logistic Regression and Binary Classification

Consider any logistic hypothesis $h(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$ that approximates $P(y|\mathbf{x})$. 'Convert' $h(\mathbf{x})$ to a binary classification prediction by taking $\text{sign}(h(\mathbf{x}) - \frac{1}{2})$. What is the equivalent formula for the binary classification prediction?

- 1 $\text{sign}(\mathbf{w}^T \mathbf{x} - \frac{1}{2})$
- 2 $\text{sign}(\mathbf{w}^T \mathbf{x})$
- 3 $\text{sign}(\mathbf{w}^T \mathbf{x} + \frac{1}{2})$
- 4 none of the above

Exercise

Logistic Regression and Binary Classification

Consider any logistic hypothesis $h(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$ that approximates $P(y|\mathbf{x})$. 'Convert' $h(\mathbf{x})$ to a binary classification prediction by taking $\text{sign}(h(\mathbf{x}) - \frac{1}{2})$. What is the equivalent formula for the binary classification prediction?

- 1 $\text{sign}(\mathbf{w}^T \mathbf{x} - \frac{1}{2})$
- 2 $\text{sign}(\mathbf{w}^T \mathbf{x})$
- 3 $\text{sign}(\mathbf{w}^T \mathbf{x} + \frac{1}{2})$
- 4 none of the above

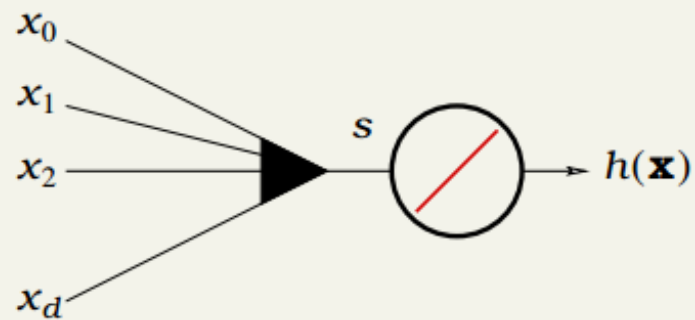
Reference Answer: 2

When $\mathbf{w}^T \mathbf{x} = 0$, $h(\mathbf{x})$ is exactly $\frac{1}{2}$. So thresholding $h(\mathbf{x})$ at $\frac{1}{2}$ is the same as thresholding $(\mathbf{w}^T \mathbf{x})$ at 0.

Linear Models

linear regression

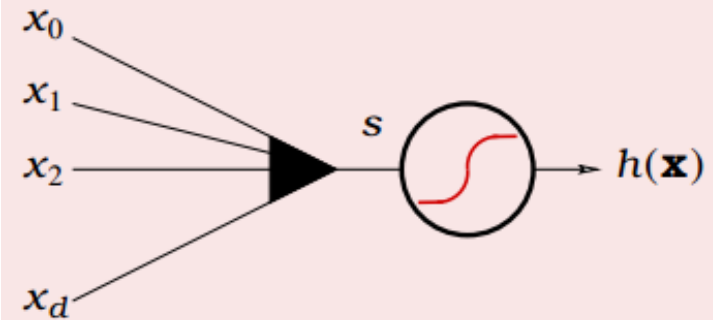
$$h(\mathbf{x}) = s$$



friendly err = squared
(easy to minimize)

logistic regression

$$h(\mathbf{x}) = \theta(s)$$



err = ?

How to define the cost (error) function for logistic regression?

Logistic Regression-- $y \in \{0,1\}$

Target function:

$$f(x) = p(+1|x) \quad \Leftrightarrow \quad p(y|x) = \begin{cases} f(x) & \text{for } y = 1 \\ 1 - f(x) & \text{for } y = 0 \end{cases}$$

Consider $\mathcal{D} = \{(x_1, +), (x_2, -), \dots, (x_m, -)\}$

Likelihood that h generates \mathcal{D}

Maximum-Likelihood Estimation

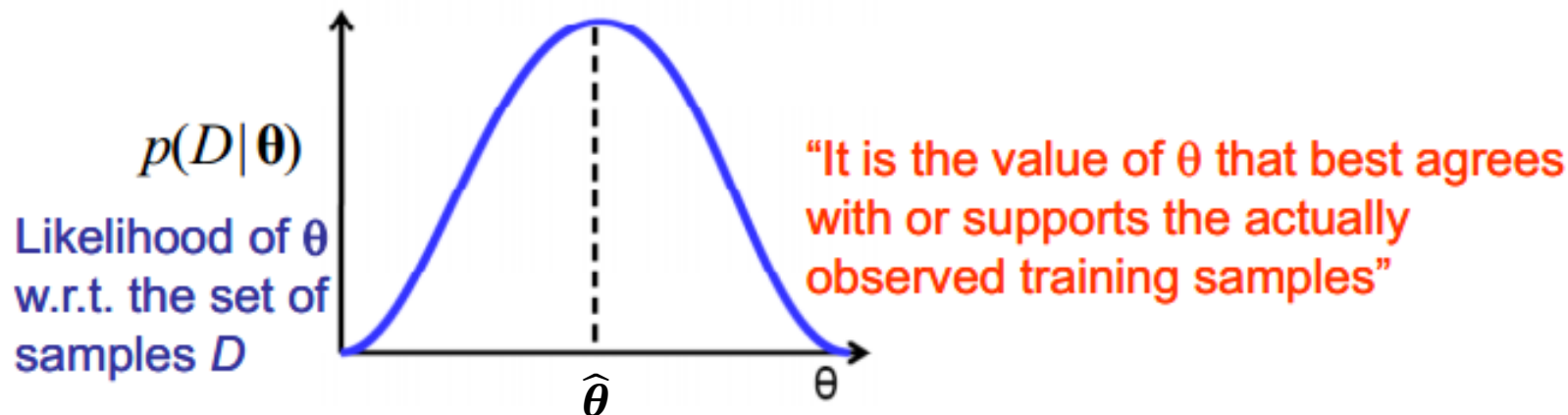
Given a dataset $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$, where the n samples are drawn **independently** from **identical** distribution $p(x|\theta)$, estimate parameters θ .

ML estimate parameters θ maximizes $p(\mathcal{D}|\theta)$

\mathcal{D} is an i.i.d set

$$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D}|\theta)$$

$$p(\mathcal{D}|\theta) = \prod_{k=1}^n p(x_k|\theta)$$



Logistic Regression-- $y \in \{0,1\}$

Target function:

$$f(x) = p(+1|x) \quad \Leftrightarrow \quad p(y|x) = \begin{cases} f(x) & \text{for } y = 1 \\ 1 - f(x) & \text{for } y = 0 \end{cases}$$

Consider $\mathcal{D} = \{(x_1, +), (x_2, -), \dots, (x_m, -)\}$

Likelihood that h generates \mathcal{D}

$$\begin{aligned} & p(x_1)h(x_1) \\ & p(x_2)(1 - h(x_2)) \\ & \vdots \\ & p(x_m)(1 - h(x_m)) \end{aligned}$$

- If $h \approx f$, then likelihood $(h) \approx$ that using (f)
- Probability using (f) is usually large

Likelihood of Logistic Regression

Goal: $\arg \max_h \text{likelihood}(h)$ $\text{likelihood}(h) = \prod_{i=1}^m p(\mathbf{x}_i) p(y|\mathbf{x}_i)$

Consider $\mathcal{D} = \{(\mathbf{x}_1, +), (\mathbf{x}_2, -), \dots, (\mathbf{x}_m, -)\}$

$$\begin{aligned} \text{likelihood}(h) &= \prod_{i=1}^m p(\mathbf{x}_i) p(y_i|\mathbf{x}_i) \\ &= p(\mathbf{x}_1) h(\mathbf{x}_1) p(\mathbf{x}_2) (1 - h(\mathbf{x}_2)) p(\mathbf{x}_m) (1 - h(\mathbf{x}_m)) \end{aligned}$$

We remove all the $p(\mathbf{x}_i)$ which remains the same for all the hypothesis h .

Likelihood of Logistic Regression

Goal: $\arg \max_h \text{likelihood}(h)$ $\text{likelihood}(h) = \prod_{i=1}^m p(\mathbf{x}_i) p(y|\mathbf{x}_i)$

Consider $\mathcal{D} = \{(\mathbf{x}_1, +), (\mathbf{x}_2, -), \dots, (\mathbf{x}_m, -)\}$

$$p(y_i|\mathbf{x}_i) = \begin{cases} h(\mathbf{x}_i) & \text{for } y_i = 1 \\ 1 - h(\mathbf{x}_i) & \text{for } y_i = 0 \end{cases} \iff p(y_i|\mathbf{x}_i) = h(\mathbf{x}_i)^{y_i} (1 - h(\mathbf{x}_i))^{(1-y_i)}$$

Bernoulli distribution

$$\text{likelihood}(h) \propto \prod_{i=1}^m p(y_i|\mathbf{x}_i) = \prod_{i=1}^m h(\mathbf{x}_i)^{y_i} (1 - h(\mathbf{x}_i))^{(1-y_i)}$$

Log-Likelihood of Logistic Regression

Negative Log-likelihood

$$\min_{\mathbf{h}} E(\mathbf{h}) = \sum_{i=1}^m -(y_i \ln \mathbf{h}(\mathbf{x}_i) + (1 - y_i) \ln(1 - \mathbf{h}(\mathbf{x}_i)))$$

Cross-entropy loss

Cross-entropy

$$H(\mathbf{p}, \mathbf{q}) = - \sum_x \mathbf{p}(\mathbf{x}) \log(\mathbf{q}(\mathbf{x}))$$

$\mathbf{p} \in \{y, 1 - y\}$
 $\mathbf{q} \in \{\mathbf{h}(\mathbf{x}), 1 - \mathbf{h}(\mathbf{x})\}$

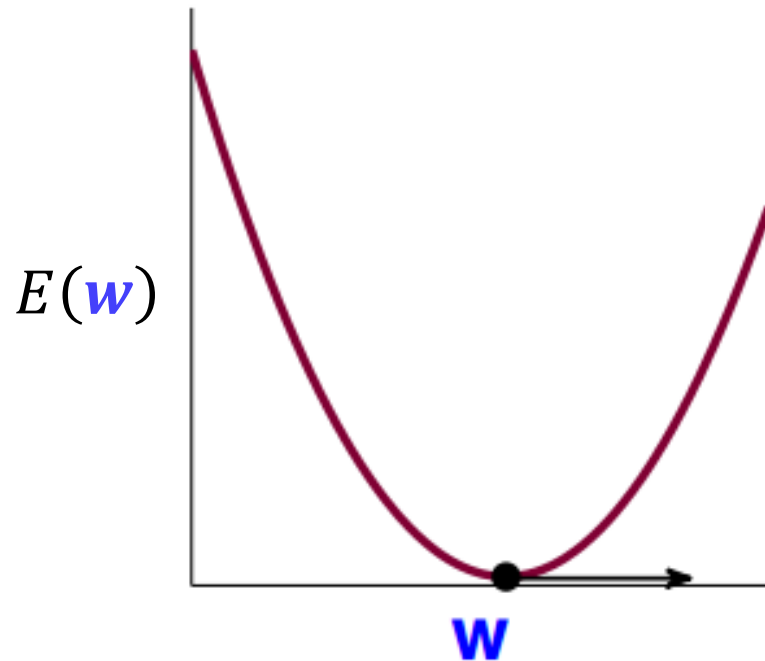
Negative Log-likelihood

$$\min_{\mathbf{w}} \sum_{i=1}^m \left[-y_i \ln \left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} \right) - (1 - y_i) \ln \left(\frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}_i}} \right) \right]$$
$$\min_{\mathbf{w}} \sum_{i=1}^m \left[-y_i \mathbf{w}^T \mathbf{x}_i + \ln(1 + e^{\mathbf{w}^T \mathbf{x}_i}) \right]$$

Minimize $E(w)$

$$\min_w E(w) = \sum_{i=1}^m \left[-y_i w^T x_i + \ln(1 + e^{w^T x_i}) \right]$$

Cross-entropy loss



$E(w)$: continuous, differentiable, twice-differentiable, **convex**
We want to find the valley

$$\nabla E(w) = 0$$

Gradient $\nabla E(\mathbf{w})$

$$\nabla E(\mathbf{w}) = \sum_{i=1}^m \left[-y_i \mathbf{x}_i + \frac{e^{\mathbf{w}^T \mathbf{x}_i}}{1 + e^{\mathbf{w}^T \mathbf{x}_i}} \mathbf{x}_i \right] = \sum_{i=1}^m [\theta(\mathbf{w}^T \mathbf{x}_i) - y_i] \mathbf{x}_i = 0$$

- $\nabla E(\mathbf{w}) = 0 \iff \begin{cases} \theta(\mathbf{w}^T \mathbf{x}_i) = 1, \text{ if } y_i = 1 \\ \theta(\mathbf{w}^T \mathbf{x}_i) = 0, \text{ if } y_i = 0 \end{cases} \iff \begin{cases} \mathbf{w}^T \mathbf{x}_i \rightarrow \infty, \text{ if } y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i \rightarrow -\infty, \text{ if } y_i = 0 \end{cases}$
 - The data must be **linearly separable**. :- (
- $\nabla E(\mathbf{w})$ is a non-linear equation of \mathbf{w}
 - It is hard to derive the **closed form** solution. :- (

Gradient Descent [Cauchy 1847]

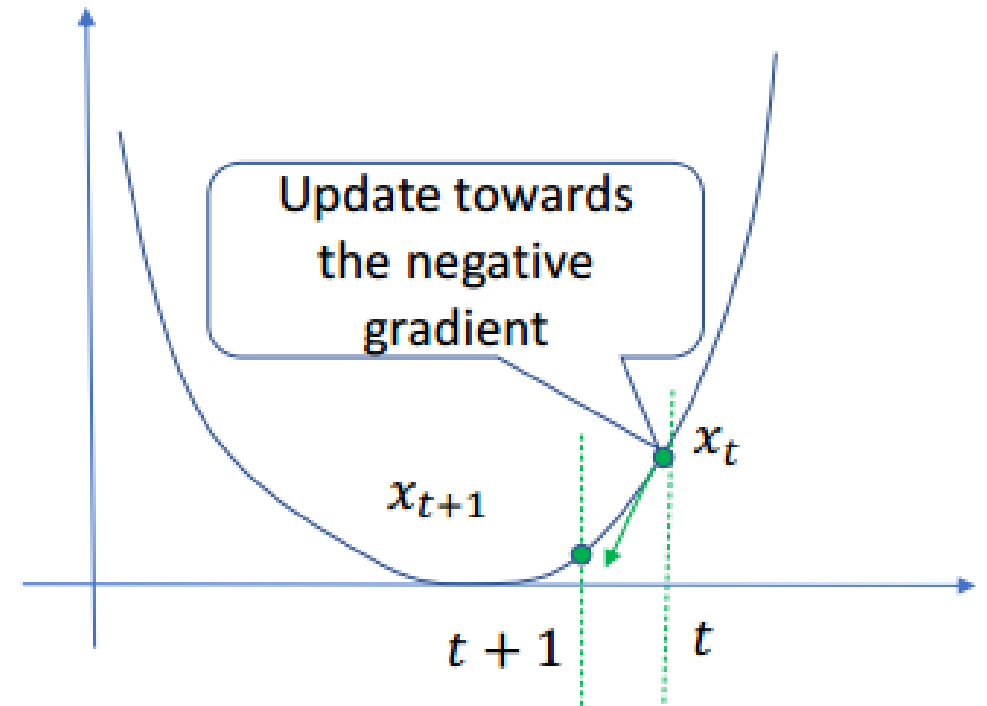
- Motivation: to **minimize** the local **first-order Taylor approximation** of f

$$\min_x f(x) \approx \min_x f(x_t) + \nabla f(x_t)^T (x - x_t)$$

- Update rule:

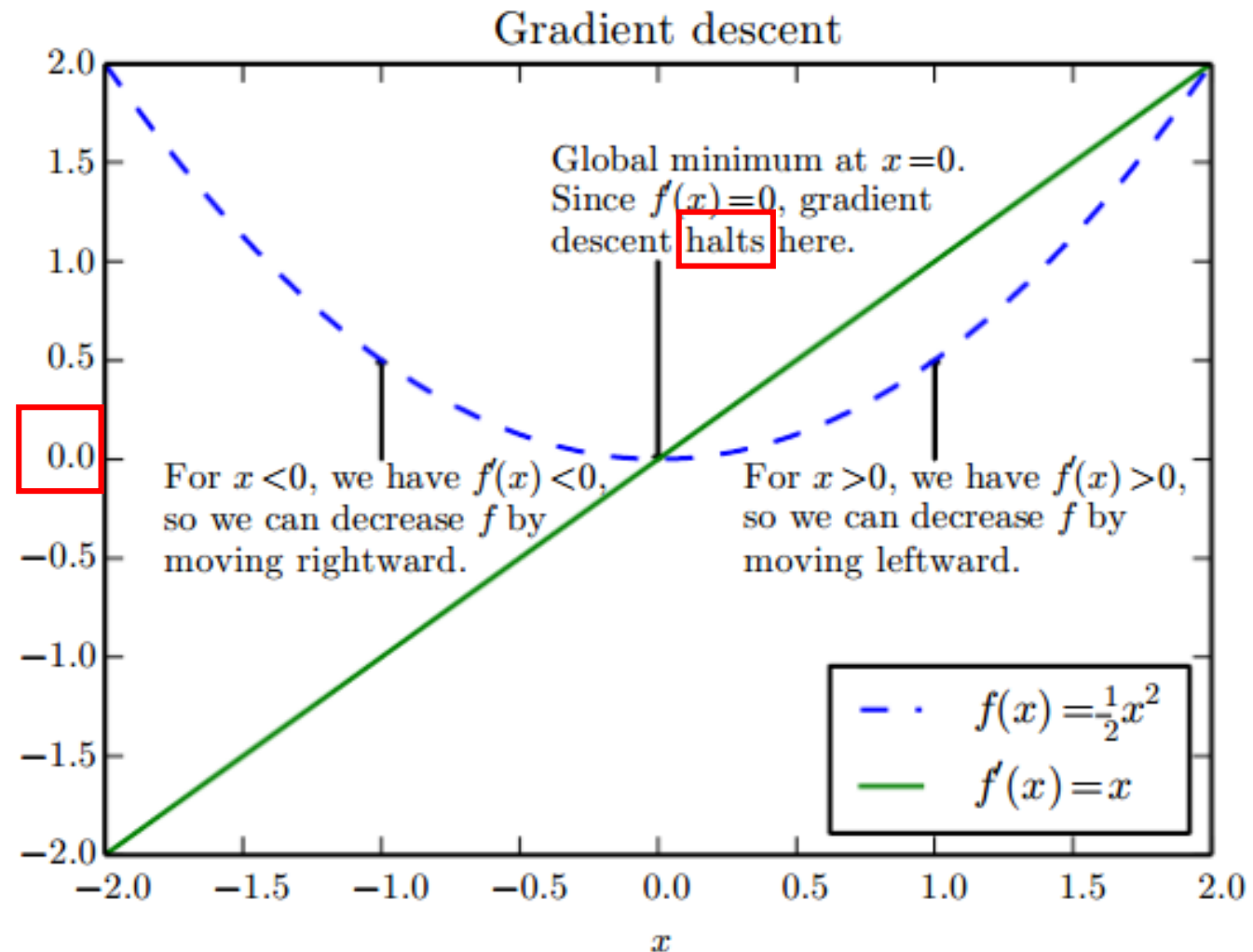
$$x_{t+1} = x_t - \eta_t \nabla f(x_t)$$

Where $\eta_t > 0$ is the step-size (learning rate).



Gradient Descent--Interpretation

- Reduce $f(x)$ by moving x in small steps with opposite sign of the derivative.
- Update rule:
$$x_{t+1} = x_t - \eta_t \nabla f(x_t)$$
- Critical points/stationary points:
Points where $f'(x) = 0$



An illustration of gradient descent.