# Machine Learning & Pattern Recognition

**Xin Xin(辛鑫)**

**xinxin@sdu.edu.cn**

**https://xinxin-me.github.io/**

# Linear Regression

# Linear Regression

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

**Training dataset:** $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;

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

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

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

For simplicity, the bias $b$ can be merged into the weight $w$:

$$h(x_i) = \widehat{w}^T \widehat{x}_i$$

$\widehat{w} = (b; w) = (b \ w_1 \ w_2 \ \dots \ w_d) \in \mathbb{R}^{d+1}$
$\widehat{x}_i = (1; \ x_{i1}; x_{i2}; \dots; x_{id}) \in \mathbb{R}^{d+1}$

# Linear Regression

**To-be-learned parameter**
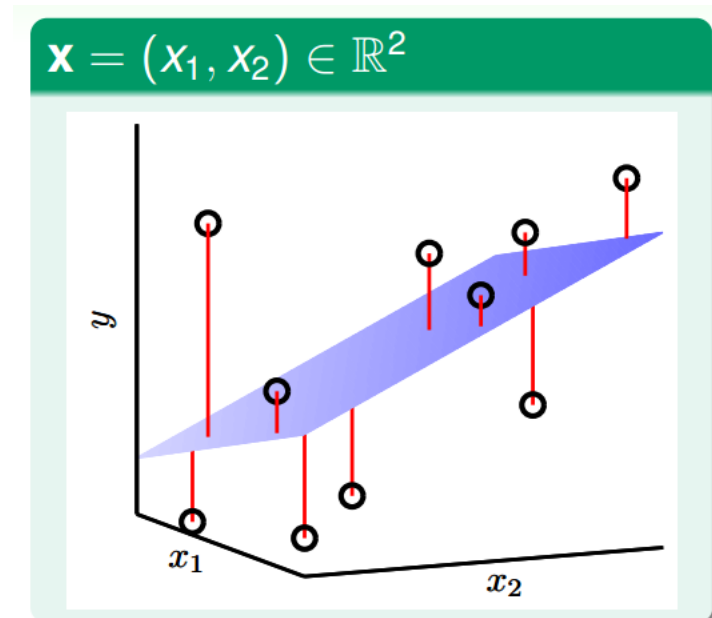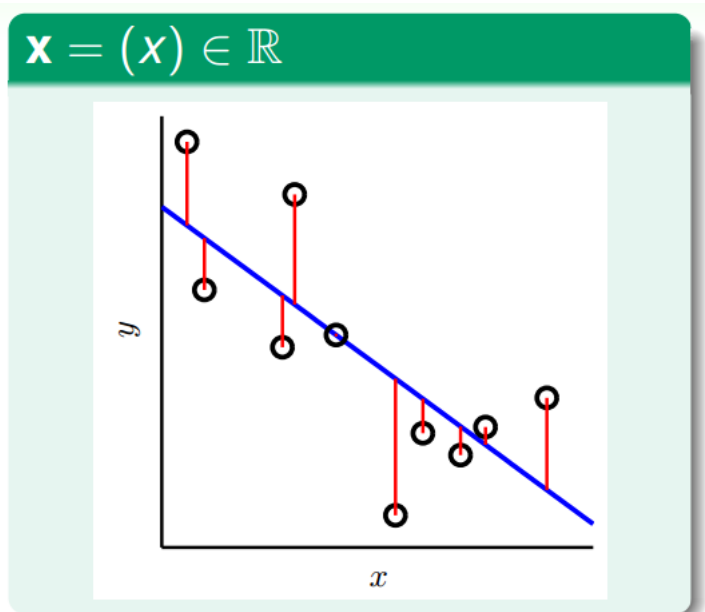
**Linear regression hypothesis:** $h(\boldsymbol{x_i}) = \boldsymbol{w}^T \boldsymbol{x_i} = \sum_{j=0}^{d} w_j x_{ij}, x_{i0} = 1$

Linear regression: find lines/hyperplanes with small residuals.

Popular/historical squared error measure:
$$L(h(\boldsymbol{x}), y) = (\hat{y} - y)^2$$



$\mathbf{x} = (x) \in \mathbb{R}$



$\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$

# Empirical Error

We 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:

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

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.

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

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

# Linear Regression

Popular/historical error measure:
squared error $L(h(\boldsymbol{x}), y) = (\hat{y} - y)^2$

$$E(\boldsymbol{w}) = \sum_{i=1}^{m} (\underbrace{h(\boldsymbol{x_i})}_{\boldsymbol{w}^T \boldsymbol{x_i}} - y_i)^2$$

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

# Matrix Form of $E(\boldsymbol{w})$

$loss = \sum_{i=1}^{m}(\boldsymbol{w^T}\boldsymbol{x_i} - y_i)^2 + \lambda\|\boldsymbol{w}\|^2, \quad \boldsymbol{w} = (w_0, w_1, \ldots, w_d)^T$

$$E(\boldsymbol{w}) = \sum_{i=1}^{m}(h(\boldsymbol{x_i}) - y_i)^2 = \sum_{i=1}^{m}(\boldsymbol{w^T}\boldsymbol{x_i} - y_i)^2 = \sum_{i=1}^{m}(\boldsymbol{x_i^T}\boldsymbol{w} - y_i)^2$$

$$= \left\| \begin{matrix} \boldsymbol{x_1^T}\boldsymbol{w} - y_1 \\ \boldsymbol{x_2^T}\boldsymbol{w} - y_2 \\ \vdots \\ \boldsymbol{x_m^T}\boldsymbol{w} - y_m \end{matrix} \right\|^2 = \left\| \begin{bmatrix} -\,-\boldsymbol{x_1^T}-\,- \\ -\,-\boldsymbol{x_2^T}-\,- \\ \vdots \\ -\,-\boldsymbol{x_m^T}-\,- \end{bmatrix} \boldsymbol{w} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \right\|^2$$

$l_2\text{-}norm\ \|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_d^2}$

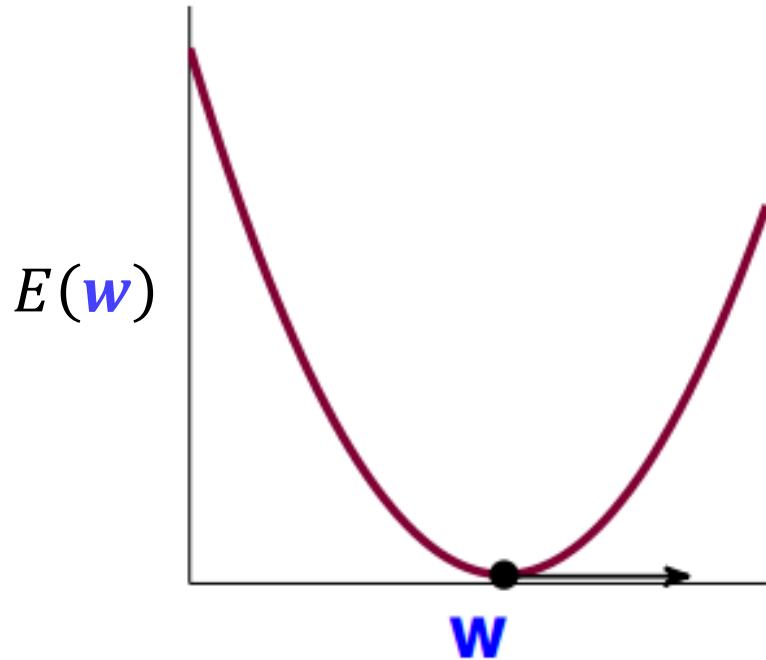$$= \|\boldsymbol{Xw} - \boldsymbol{y}\|^2$$

The subscript '2' is usually omitted.

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

# Matrix Form of $E(\boldsymbol{w})$

A continuous, twice differentiable function of several variables is convex on a convex set if and only if its Hessian matrix is positive semidefinite on the interior of the convex set.

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

$E(\boldsymbol{w})$



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

$$\nabla E(\boldsymbol{w}) = \begin{bmatrix} \dfrac{\partial E}{\partial w_0}(\boldsymbol{w}) \\ \dfrac{\partial E}{\partial w_1}(\boldsymbol{w}) \\ \vdots \\ \dfrac{\partial E}{\partial w_d}(\boldsymbol{w}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

✓ Not possible to 'roll down'

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

# The Gradient $\nabla E(w)$

$$\min_{w} E(w) = \|Xw - y\|^2 = (Xw - y)^T (Xw - y) = w^T \underbrace{X^T X}_{A} w - 2w^T \underbrace{X^T y}_{b} + \underbrace{y^T y}_{c}$$

## One w only

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

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

## Vector w

$$E(w) = (w^T A w - 2w^T b + c)$$

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

# Derivatives

| | Differentiate | | |
| --- | --- | --- | --- |
| w.r.t | scalar | vector | matrix |
| scalar | scalar | vector | matrix |
| vector | vector | matrix | |
| matrix | matrix | | |

- **scalar –scalar:** e.g., $\dfrac{d}{dx}x^2 = 2x$

- **scalar-vector:** e.g., $f(\boldsymbol{x})$ is a scalar function of vector $\boldsymbol{x}$

$$\mathrm{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \qquad \frac{df}{d\mathbf{x}} = \begin{bmatrix} \dfrac{\sigma f}{\sigma x_1} \\ \vdots \\ \dfrac{\sigma f}{\sigma x_d} \end{bmatrix}$$

- **scalar-matrix:** $f(\boldsymbol{A})$ is a scalar function and $m \times n$ matrix $A$

$$\frac{df}{d\mathbf{A}} = \begin{bmatrix} \dfrac{\sigma f}{\sigma a_{11}} & \cdots & \dfrac{\sigma f}{\sigma a_{1d}} \\ \vdots & \ddots & \vdots \\ \dfrac{\sigma f}{\sigma a_{m1}} & \cdots & \dfrac{\sigma f}{\sigma a_{mn}} \end{bmatrix}$$

https://en.wikipedia.org/wiki/Matrix_calculus

# Matrix Calculus

- Numerator layout: lay out according to $y$ and $x^T$. (Jacobian formulation)
- Denominator layout: lay out according to $y^T$ and $x$. (Hessian formulation)

**Numerator layout:**
分子布局

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \cdots & \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{bmatrix}$$

**Denominator layout:**
分母布局

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \cdots & \frac{\partial y_n}{\partial x} \end{bmatrix}$$

# Commonly Used Derivatives

- $$\frac{d}{dx}(Ax) = A^T$$

- $$\frac{dx}{dx} = I$$

- $$\frac{dy^T x}{dx} = \frac{dx^T y}{dx} = y$$

- $$\frac{d}{dx}(x^T A x) = \begin{cases} (A + A^T)x & \text{If } \mathbf{A} \text{ square} \\ 2Ax & \text{If } \mathbf{A} \text{ symmetric} \end{cases}$$

# The Gradient $\nabla E(w)$

$$\min_{w} E(w) = \|Xw - y\|^2 = (Xw - y)^T (Xw - y) = \underbrace{w^T X^T X w}_{A} - \underbrace{2 w^T X^T y}_{b} + \underbrace{y^T y}_{c}$$

## One w only

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

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

## Vector **w**

$$E(w) = (w^T A w - 2 w^T b + c)$$

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

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

# Optimal Linear Regression Weights

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

## Invertible/positive definite $\boldsymbol{X}^T \boldsymbol{X}$

- Unique solution

$$\boldsymbol{w}^* = \underbrace{(\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T}\, \boldsymbol{y}$$

pseudo-inverse $\boldsymbol{X}^\dagger$

**Note:** $\boldsymbol{X}^\dagger \boldsymbol{X} = \boldsymbol{I}$ , but $\boldsymbol{X} \boldsymbol{X}^\dagger \neq \boldsymbol{I}$

If $\boldsymbol{X}$ is square and invertible, $\boldsymbol{X}^\dagger = \boldsymbol{X}^{-1}$.

## Singular $\boldsymbol{X}^T \boldsymbol{X}$

- Define $\boldsymbol{X}^\dagger$ in other ways (e.g., SVD).
- Add regularization
  - E.g., $l_2$ norm      $\lambda > 0$

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

$$\nabla E(\boldsymbol{w}^*) = 2(\boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w} + \lambda \boldsymbol{w} - \boldsymbol{X}^T \boldsymbol{y}) = \boldsymbol{0}$$

$$\underline{(\boldsymbol{X}^T \boldsymbol{X} + \lambda \boldsymbol{I})} \boldsymbol{w} = \boldsymbol{X}^T \boldsymbol{y}$$

Invertible?

# Linear Regression Algorithm

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

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

**2. Calculate pseudo-inverse**

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

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

Simple and efficient（？） with **good $X^{\dagger}$**

# Logistic Regression

# Heart Attack Prediction Problem

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

heart disease? **yes**

Binary classification:
Ideal $f(x) = sign(p(+1|x) - 0.5) \in \{-1, +1\}$

# 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

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

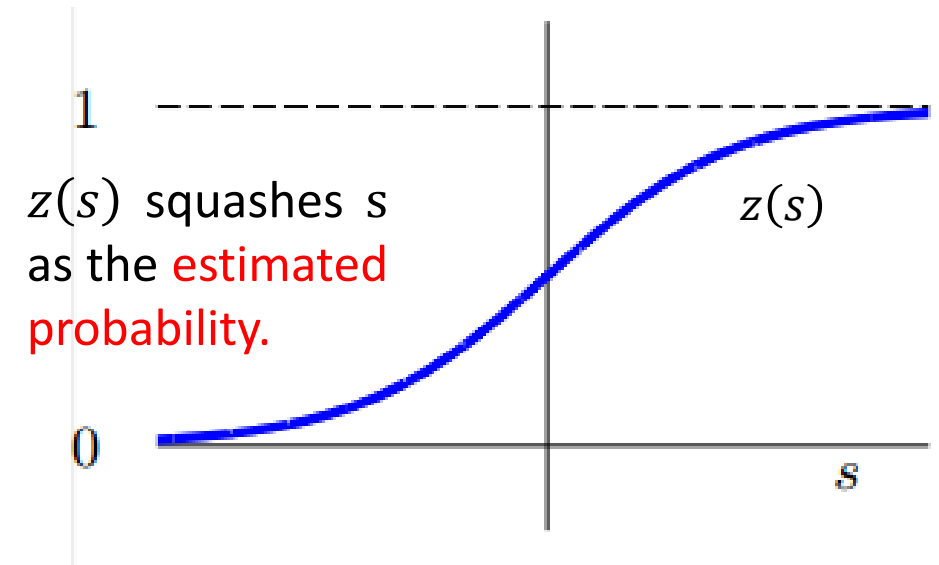Same data as hard binary classification, different **target function**

# Logistic Hypothesis

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

Let $\boldsymbol{x_i} = (x_{i0}, x_{i1}, x_{i2}, ..., x_{id})$ be the features of the patient, calculate a weighted 'risk score':
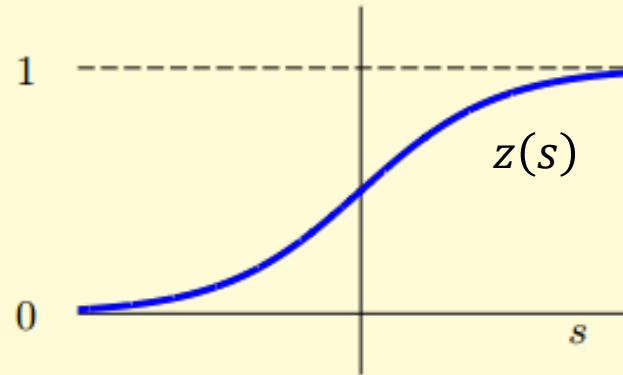
$$s = \sum_{j=0}^{d} w_j x_{ij} = \boldsymbol{w}^T \boldsymbol{x_i},$$

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

$z(s)$ squashes s as the estimated probability.

Logistic hypothesis: $h(\boldsymbol{x_i}) = z(\boldsymbol{w}^T \boldsymbol{x_i})$

# Logistic Function

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



smooth, monotonic, sigmoid function of $s$

Bound $\quad z(s) \in [0,1] \qquad z(-\infty) = 0 \qquad z(0) = 0.5 \qquad z(\infty) = 1$

Symmetric $\qquad 1 - z(s) = z(-s)$

Gradient $\qquad z'(s) = z(s)(1 - z(s))$

Logistic regression use $h(x) = z(w^T x)$ to approximate the target $f(x) = p(+1|x)$

# 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 $\text{sign}\left(h(\mathbf{x}) - \frac{1}{2}\right)$. What is the equivalent formula for the binary classification prediction?

1. $\text{sign}\left(\mathbf{w}^T\mathbf{x} - \frac{1}{2}\right)$
2. $\text{sign}\left(\mathbf{w}^T\mathbf{x}\right)$
3. $\text{sign}\left(\mathbf{w}^T\mathbf{x} + \frac{1}{2}\right)$
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 $\text{sign}\left(h(\mathbf{x}) - \frac{1}{2}\right)$. What is the equivalent formula for the binary classification prediction?

1. $\text{sign}\left(\mathbf{w}^T\mathbf{x} - \frac{1}{2}\right)$

2. $\text{sign}\left(\mathbf{w}^T\mathbf{x}\right)$

3. $\text{sign}\left(\mathbf{w}^T\mathbf{x} + \frac{1}{2}\right)$

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.
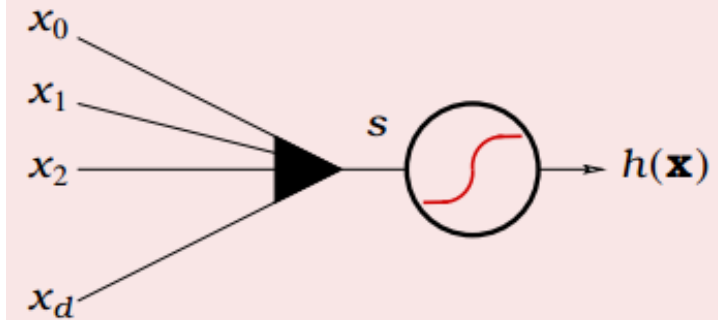
# Linear Models



How to define the cost (error) function for logistic regression?
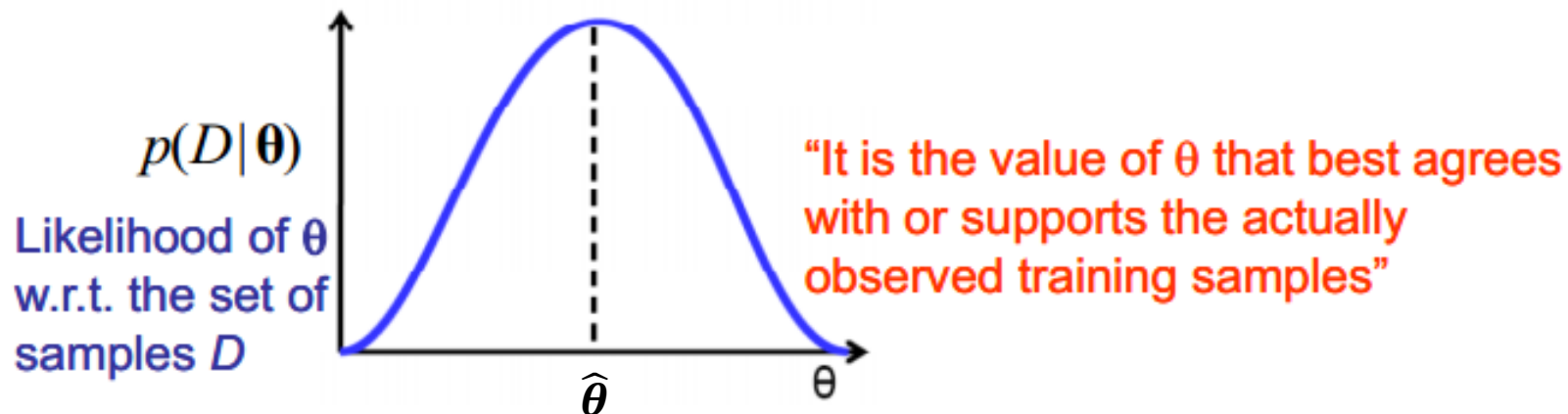
# Maximum-Likelihood Estimation

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

ML estimate parameters $\theta$ maximizes $p(\mathcal{D}|\theta)$    $\mathcal{D}$ is an i.i.d set

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

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



$p(D|\theta)$

Likelihood of θ w.r.t. the set of samples $D$

$\widehat{\theta}$    θ

"It is the value of θ that best agrees with or supports the actually observed training samples"

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

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

**Likelihood that $h$ generates $\mathcal{D}$**

$$p(x_1)h(x_1)$$
$$p(x_2)(1 - h(x_2))$$
$$\vdots$$
$$p(x_m)(1 - h(x_m))$$

- Target function:
  $f(x) = p(+1|x)$

- If $h \approx f$, then likelihood $(h) \approx$ that using $(f)$

# Likelihood of Logistic Regression

Goal: $\arg\max\limits_{h} likelihood(h)$         $likelihood(h) = \prod\limits_{i=1}^{m} p(\boldsymbol{x_i})p(y|\boldsymbol{x_i})$

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

$$likelihood(h) = \prod_{i=1}^{m} p(\boldsymbol{x_i})p(y_i|\boldsymbol{x_i})$$

$$= p(\boldsymbol{x_1})h(\boldsymbol{x_1})p(\boldsymbol{x_2})(1 - h(\boldsymbol{x_2})) \cdots p(\boldsymbol{x_m})(1 - h(\boldsymbol{x_m}))$$

# Likelihood of Logistic Regression

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

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

$$likelihood(h) = \prod_{i=1}^{m} p(\boldsymbol{x_i})p(y_i|\boldsymbol{x_i})$$

$$= p(x_1)h(\boldsymbol{x_1})p(x_2)(1 - h(\boldsymbol{x_2})) \cdots p(x_m)(1 - h(\boldsymbol{x_m}))$$

We remove all the $p(\boldsymbol{x_i})$ which remains the same for all the hypothesis $h$ .

# Likelihood of Logistic Regression

$$likelihood(h) = \prod_{i=1}^{m} p(\boldsymbol{x_i})p(y_i|\boldsymbol{x_i}) \propto \prod_{i=1}^{m} p(y_i|\boldsymbol{x_i})$$

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

Bernoulli distribution

$$likelihood(h) \propto \prod_{i=1}^{m} p(y_i|\boldsymbol{x_i}) = \prod_{i=1}^{m} h(\boldsymbol{x_i})^{y_i}(1 - h(\boldsymbol{x_i}))^{(1-y_i)}$$

# Log-Likelihood of Logistic Regression

**Negative Log-likelihood**

$$\min_h E(h) = \sum_{i=1}^{m} -(y_i \ln h(\boldsymbol{x_i}) + (1 - y_i) \ln(1 - h(\boldsymbol{x_i})))$$

Cross-entropy loss

**Cross-entropy**

$$H(\textcolor{red}{p}, \textcolor{blue}{q}) = -\sum_{x} \textcolor{red}{p(x)} \log(\textcolor{blue}{q(x)})$$

$$\textcolor{red}{p} \in \{y, 1 - y\}$$
$$\textcolor{blue}{q} \in \{h(\boldsymbol{x}), 1 - h(\boldsymbol{x})\}$$

**Negative Log-likelihood**

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

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

# Minimize $E(w)$

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

Cross-entropy loss



$E(\boldsymbol{w})$

$\boldsymbol{W}$

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

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

# **Matrix Calculus**

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

Identities: scalar-by-vector $\dfrac{\partial y}{\partial \mathbf{x}} = \nabla_{\mathbf{x}} y$

| Condition | Expression | Numerator layout, i.e. by $\mathbf{x}^T$; result is row vector | Denominator layout, i.e. by $\mathbf{x}$; result is column vector |
|---|---|---|---|
| $a$ is not a function of $\mathbf{x}$ | $\dfrac{\partial a}{\partial \mathbf{x}} =$ | $\mathbf{0}^{\top}$ [4] | $\mathbf{0}$ [4] |
| $a$ is not a function of $\mathbf{x}$, $u = u(\mathbf{x})$ | $\dfrac{\partial au}{\partial \mathbf{x}} =$ | | $a\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial(u+v)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial u}{\partial \mathbf{x}} + \dfrac{\partial v}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial uv}{\partial \mathbf{x}} =$ | | $u\dfrac{\partial v}{\partial \mathbf{x}} + v\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial g(u)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial f(g(u))}{\partial \mathbf{x}} =$ | | $\dfrac{\partial f(g)}{\partial g}\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $\mathbf{u} = \mathbf{u}(\mathbf{x})$, $\mathbf{v} = \mathbf{v}(\mathbf{x})$ | $\dfrac{\partial(\mathbf{u}\cdot\mathbf{v})}{\partial \mathbf{x}} = \dfrac{\partial \mathbf{u}^{\top}\mathbf{v}}{\partial \mathbf{x}} =$ | $\mathbf{u}^{\top}\dfrac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^{\top}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$  • assumes numerator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\mathbf{v} + \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}\mathbf{u}$  • assumes denominator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ |

# **Gradient** $\nabla E(w)$

$$\nabla E(\boldsymbol{w}) = \sum_{i=1}^{m}\left[-y_i\boldsymbol{x_i} + \frac{e^{\boldsymbol{w^T x_i}}}{1 + e^{\boldsymbol{w^T x_i}}}\boldsymbol{x_i}\right] = \sum_{i=1}^{m}\left[z(\boldsymbol{w^T x_i}) - y_i\right]\boldsymbol{x_i} = 0$$

- $\nabla E(\boldsymbol{w})$ is a non-linear equation of $\boldsymbol{w}$
  - ➢ It is hard to derive the closed form solution. :-(

# Gradient $\nabla E(w)$

$$\nabla E(\boldsymbol{w}) = \sum_{i=1}^{m} \left[ -y_i \boldsymbol{x_i} + \frac{e^{\boldsymbol{w}^T x_i}}{1 + e^{\boldsymbol{w}^T x_i}} \boldsymbol{x_i} \right] = \sum_{i=1}^{m} \left[ z(\boldsymbol{w}^T \boldsymbol{x_i}) - y_i \right] \boldsymbol{x_i} = 0$$

- Apply the iterative optimization to the logistic regression.

# Iterative Optimization

# Optimization Methods

- Optimization: either minimize or maximize some function $f(x)$ by altering $x$.
- In most cases, optimization refers to the minimization of $f(x)$.

**Maximization** $f(x)$ $\Longleftrightarrow$ **Minimization** $-f(x)$

- $f(x)$: objective function, cost function, loss function, error function.
- The value that minimize $f(x)$: $x^* = \arg\min f(x)$.

# Optimization Methods

- **Deterministic Optimization**
  - The data for the given problem are known accurately.

- **Stochastic Optimization**
  - Refers to a collection of methods for minimizing or maximizing an objective function when randomness is present.

# Deterministic Optimization

- First-order methods: methods that use only the gradient.

- Second-order methods: methods that also use the Hessian matrix.

$$H(f)_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\boldsymbol{x})$$

$\boldsymbol{x}$ : multiple input dimensions.

# Taylor Approximation

**Expansion at $x_0$**

$$f(x) = \frac{f(x_0)}{0!} + \frac{f'(x_0)}{1!}(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \ldots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n + R_n(x)$$

**Examples**

$$e^x = 1 + \frac{1}{1!}x + \frac{1}{2!}x^2 + \frac{1}{3!}x^3 + o\left(x^3\right)$$

$$\ln(1 + x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 + o\left(x^3\right)$$

$$\sin x = x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 + o\left(x^5\right)$$

# Gradient Descent [Cauchy 1847]

- Motivation: to <span style="color:red">minimize</span> the local <span style="color:red">first-order Taylor approximation</span> 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).

# 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/stationary points: Points where $f'(x) = 0$    驻点



An illustration of gradient descent.

# Interpretation

- At each iteration, consider the expansion

$$f(x) \approx f(x_t) + \nabla f(x_t)^T (x - x_t) + \frac{1}{2\eta_t} \|x - x_t\|^2$$

Linear approximation of $f$    Proximity term with weight $\frac{1}{2\eta_t}$

- Quadratic approximation, replacing usual $\nabla^2 f(x)$ by $\frac{1}{\eta_t} I$:

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

# Interpretation

$$f(x) \approx f(x_t) + \nabla f(x_t)^T (x - x_t) + \frac{1}{2\eta_t} \|x - x_t\|^2$$



Blue point is $x_t$, red point is $x_{t+1}$.

# Global VS Local Minimum

- Global minimum: a point that obtains the absolute lowest value of $f(x)$.

- Local minimum: a point where $f(x)$ is lower than at all neighboring points.

- Saddle points: some critical points are neither maxima or minima. 鞍点
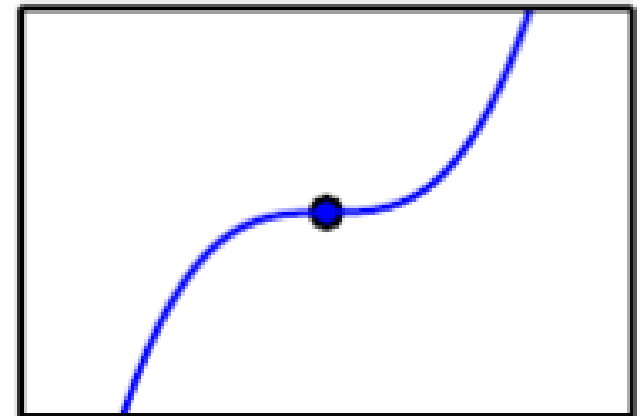
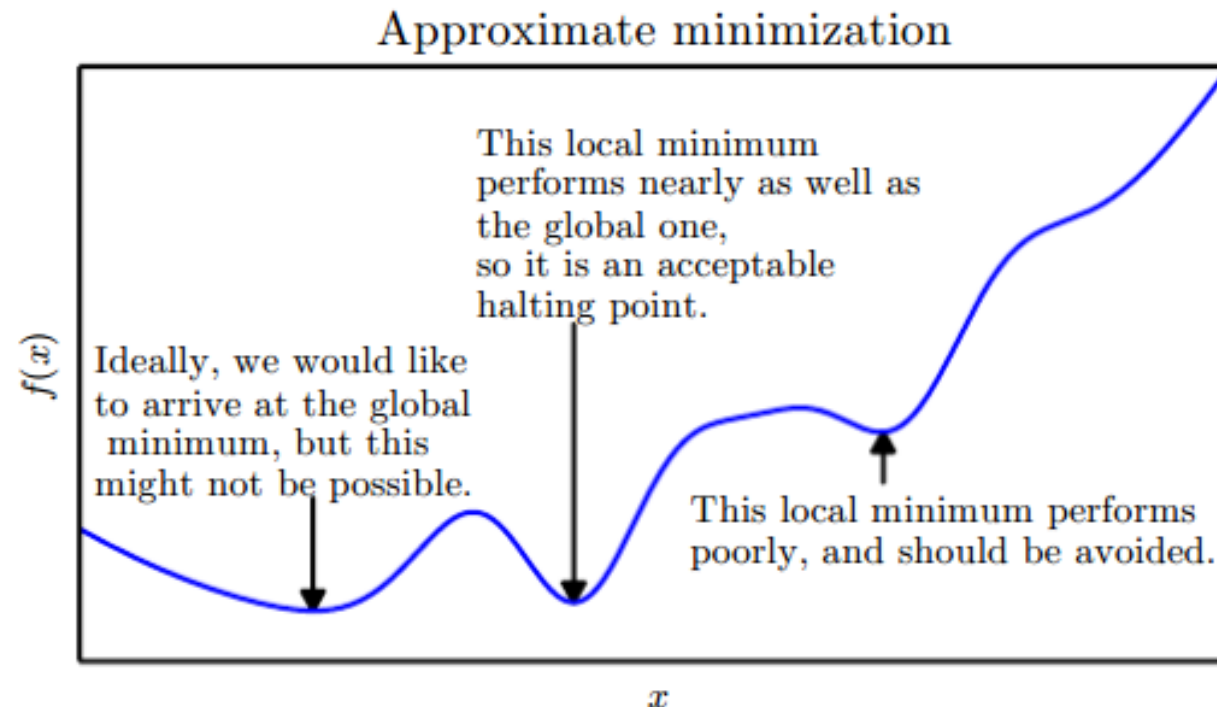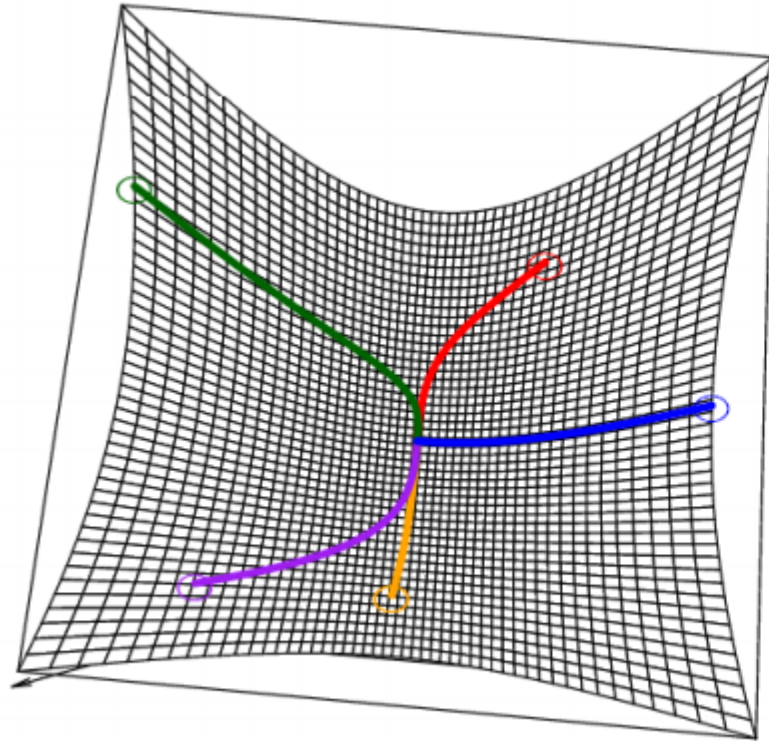Types of critical points

| Minimum | Maximum | Saddle points |

# Global VS Local Minimum

- Global minimum: a point that obtains the absolute lowest value of $f(x)$.

- Local minimum: a point where $f(x)$ is higher than at all neighboring points.

- Saddle points: some critical points are neither maxima or minima. 鞍点



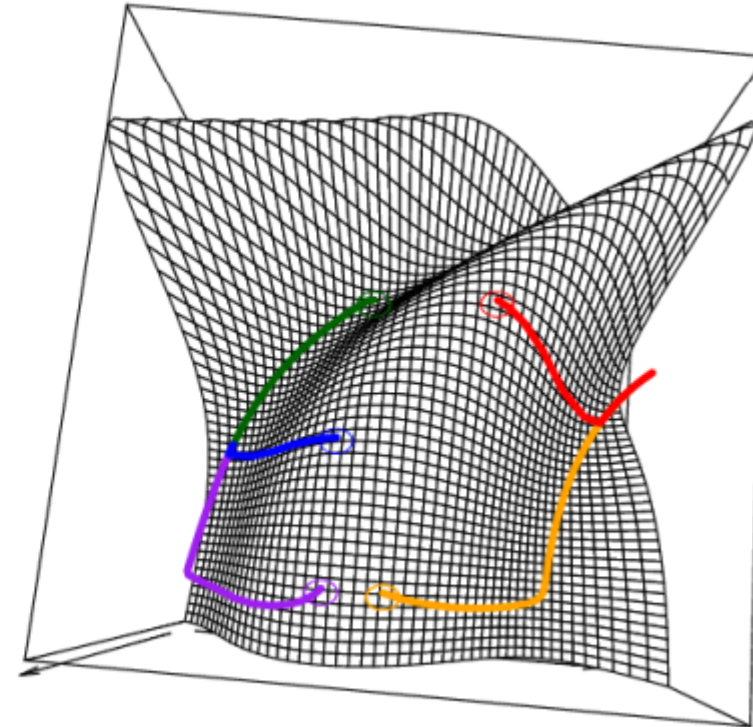Approximate minimization

# Different Starting Points

- Gradient Descent with different starting points are illustrated in different colors.



(a) Convex function          (b) Non-convex function

- (a): Strictly convex function:  Converge to the global optimum.
- (b): Non-convex function: Different paths may end up at different local optima.
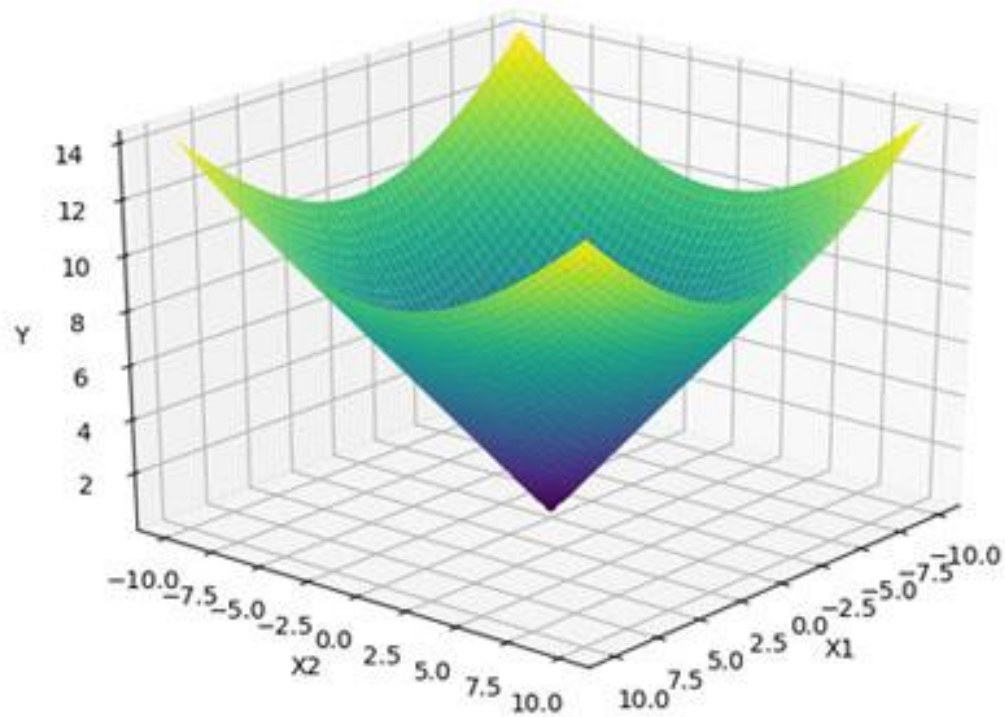
# Gradient Descent [Cauchy 1847]

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

- Gradient Descent requires a step size $\eta$ controlling the amount of gradient updated to the current point at each iteration.

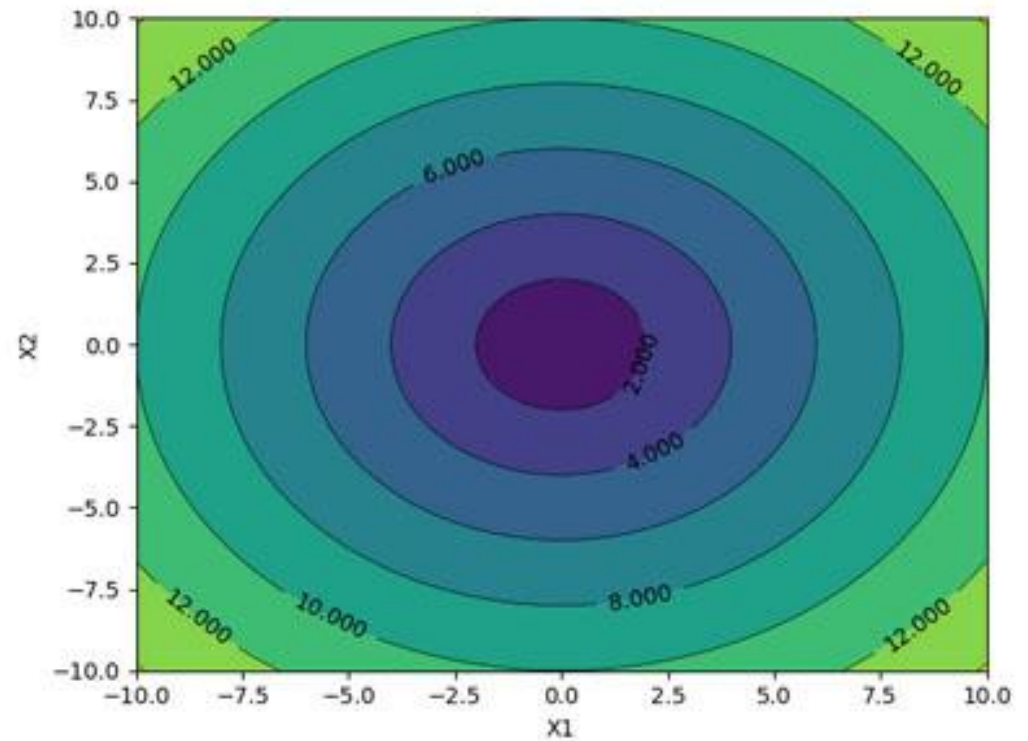- It is naïve to set $\eta_t = \eta$ for all iterations.

How to choose step sizes?

# Fixed Step Sizes

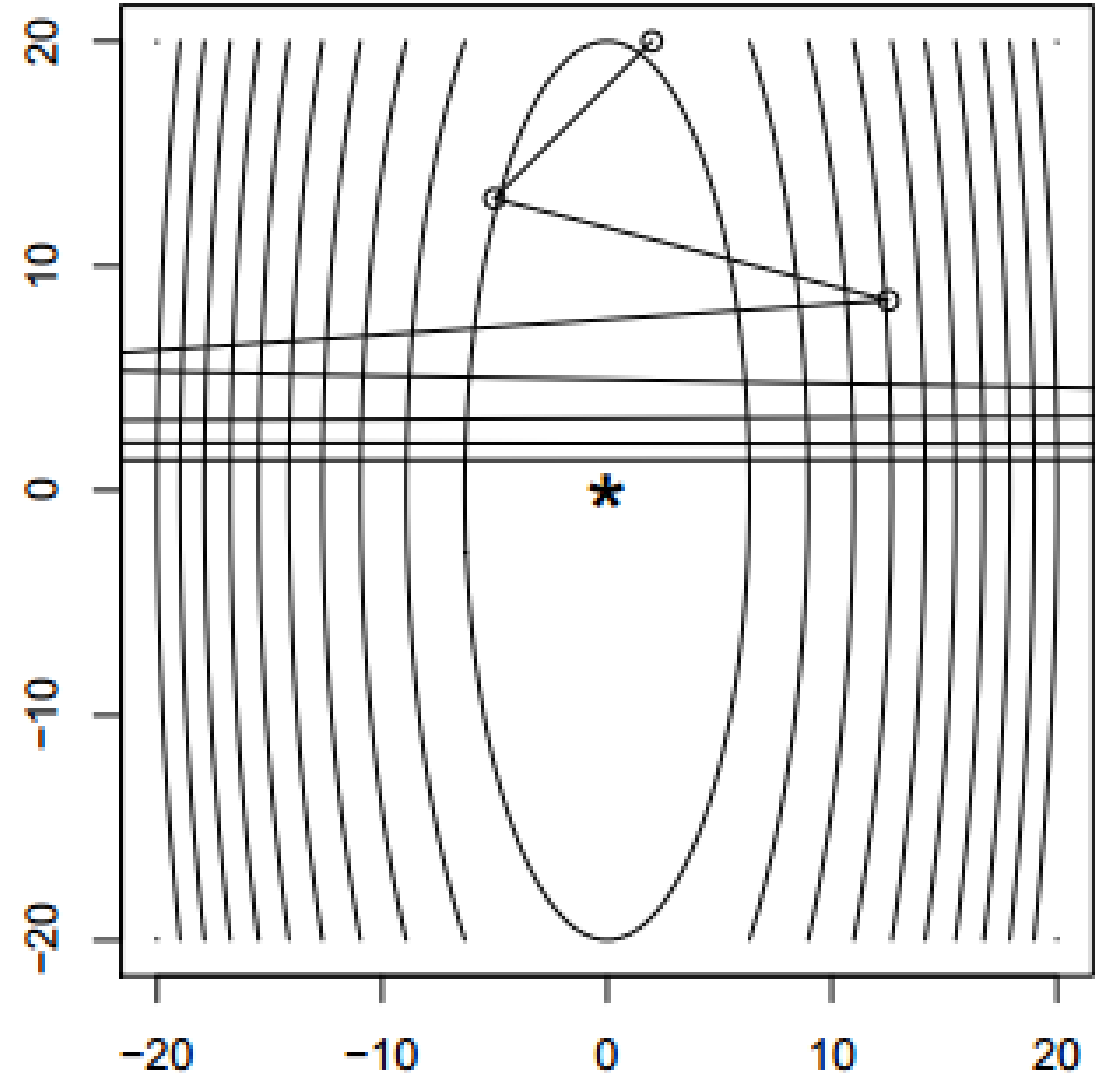Considering $f(x) = (10x_1^2 + x_2^2)/2$



3D Plot



Contour Plot

# Fixed Step Sizes

Considering $f(x) = (10x_1^2 + x_2^2)/2$

If $\eta$ is too big, can lead to divergence.

- The learning function oscillates away from the optimal point.

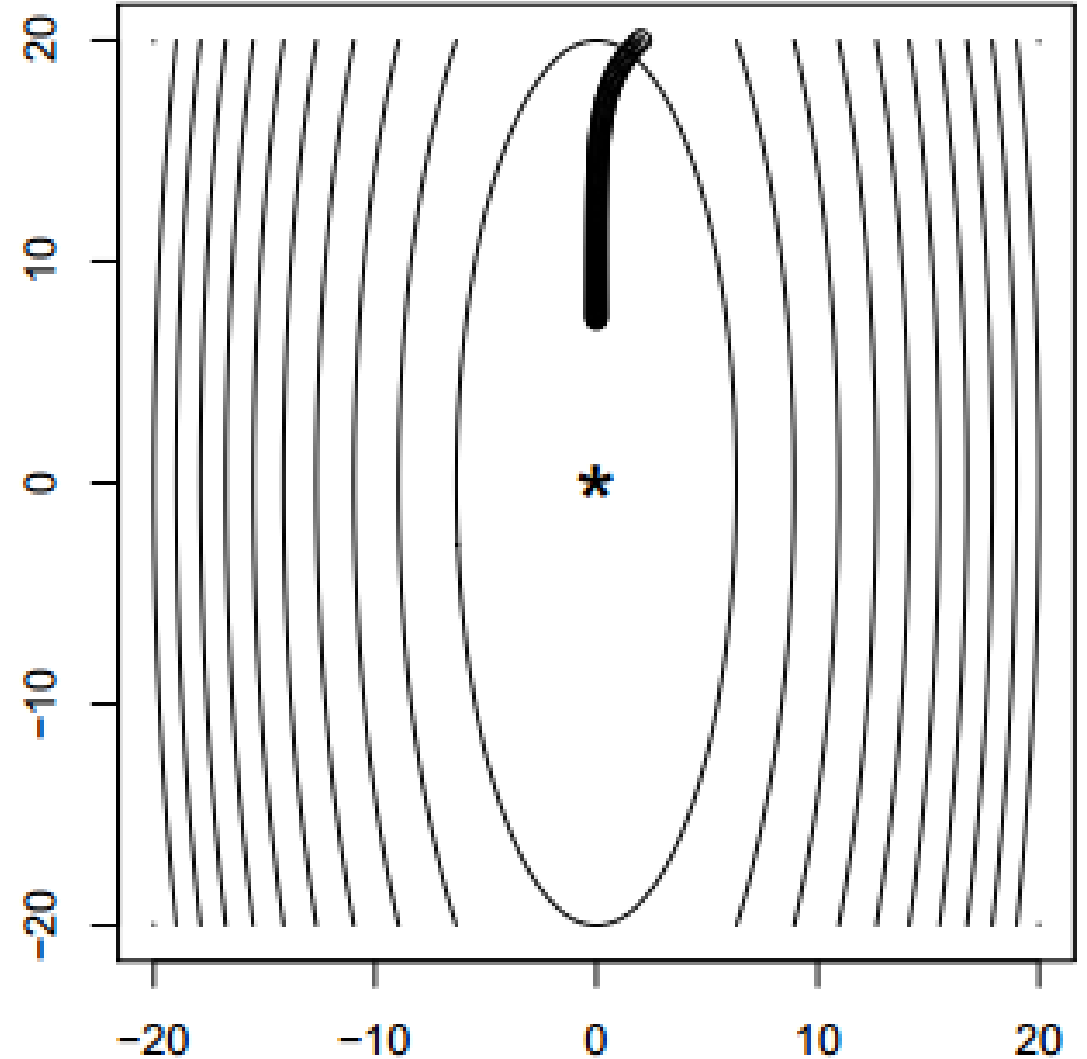- As shown, it oscillates after 8 steps.

# Fixed Step Sizes

Considering $f(x) = (10x_1^2 + x_2^2)/2$

If $\eta$ is too small, takes longer time for the function to converge.
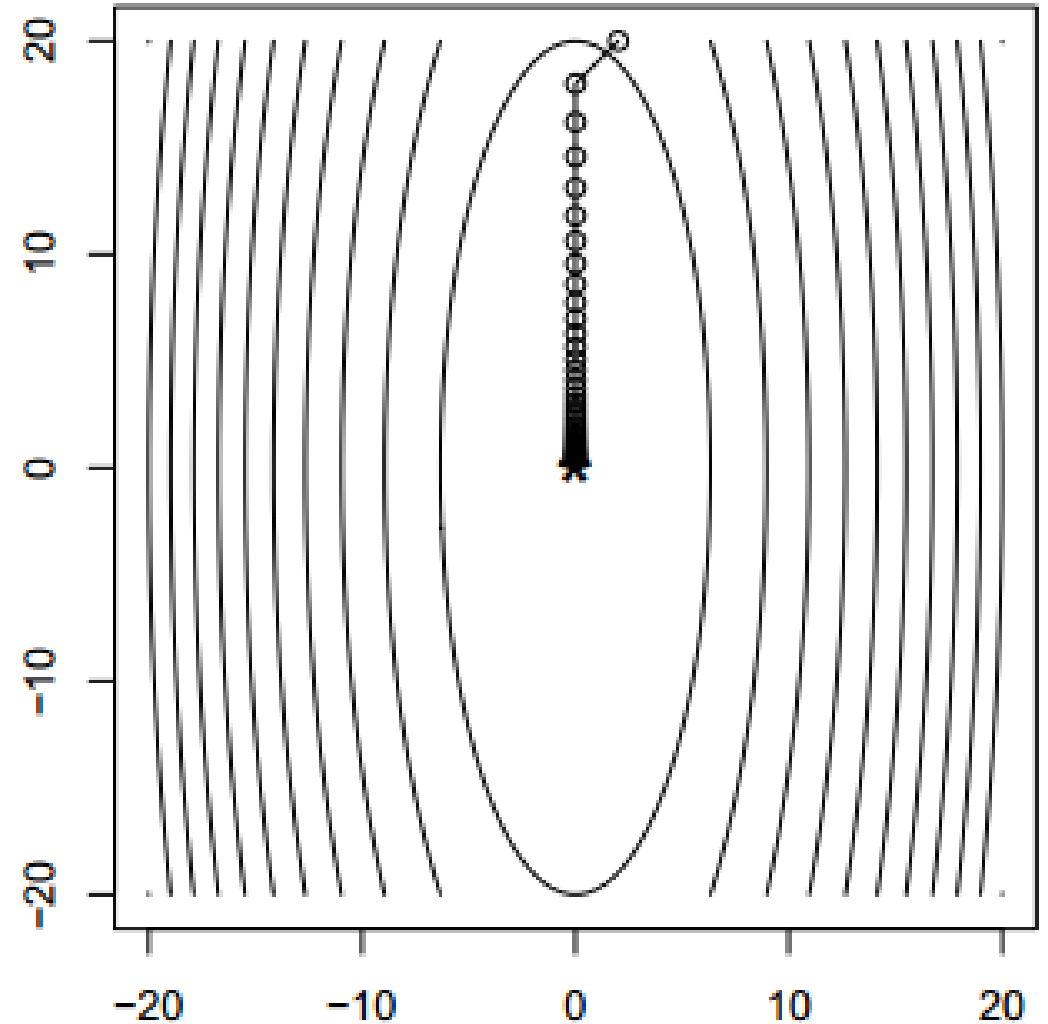- As shown, GD after 100 steps.

# Fixed Step Sizes

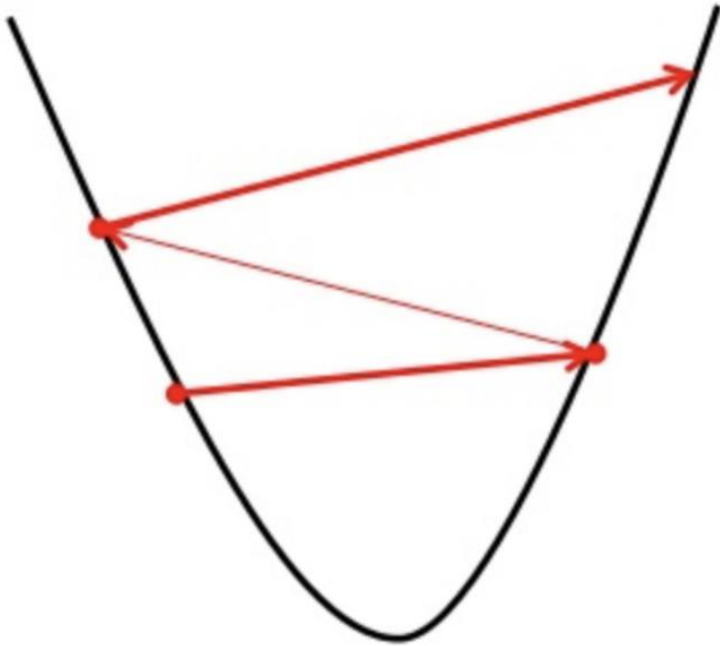Considering $f(x) = (10x_1^2 + x_2^2)/2$

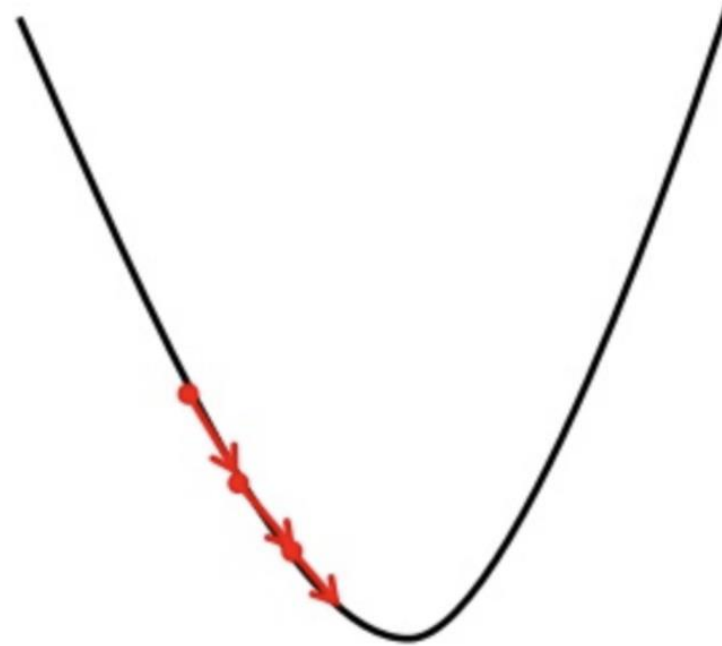Same example, gradient descent after 40 appropriately sized steps.

# Fixed Step Sizes

Considering $f(x) = x^2/2$

Big learning rate

Small learning rate

# Deterministic Optimization

- First-order methods: methods that use only the gradient.

- Second-order methods: methods that also use the Hessian matrix.

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a function taking as input a vector $x \in \mathbb{R}^n$ and outputting a scalar $f(\boldsymbol{x}) \in \mathbb{R}$; if all second partial derivatives of $f$ exist and are continuous over the domain of the function, then the Hessian matrix $\boldsymbol{H}$ of $f$ is a square $n \times n$ matrix, usually defined as follows.

$$H = \nabla^2 f = \begin{bmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1 x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_1 x_n} \\ \dfrac{\partial^2 f}{\partial x_2 x_1} & \dfrac{\partial^2 f}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial^2 f}{\partial x_n x_1} & \dfrac{\partial^2 f}{\partial x_n x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \text{, or} \qquad H_{ij} = \dfrac{\partial^2 f}{\partial x_i x_j}$$

# Newton's Methods

- Motivation: to minimize the local <span style="color:red">second-order Taylor</span> approximation of $f$.

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) \approx \min_{\boldsymbol{x}} f(\boldsymbol{x}_t) + \nabla f(\boldsymbol{x}_t)^T (\boldsymbol{x} - \boldsymbol{x}_t) + \frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}_t)^T \nabla^2 f(\boldsymbol{x}_t)(\boldsymbol{x} - \boldsymbol{x}_t)$$

- Take the derivative of $\boldsymbol{x}$ on both side, we have,

$$\frac{df(\boldsymbol{x})}{d\boldsymbol{x}} = \nabla f(x_t) + \nabla^2 f(x_t)(x - x_t) = \mathbf{0}$$

- Update rule: suppose $\nabla^2 f(x_t)$ is positive definite,

$$\boldsymbol{x} = \boldsymbol{x}_t - [\nabla^2 f(x_t)]^{-1} \nabla f(x_t)$$

# Newton's Methods

- Motivation: to minimize the local <span style="color:red">second-order Taylor</span> approximation of $f$.

$$\min_x f(x) \approx \min_x f(x_t) + f'(x_t)(x - x_t) + \frac{1}{2}f''(x_t)(x - x_t)^2$$

- Take the derivative of $x$ on both side, we have,

$$f'(x) = f'(x_t) + f''(x_t)(x - x_t) = 0$$

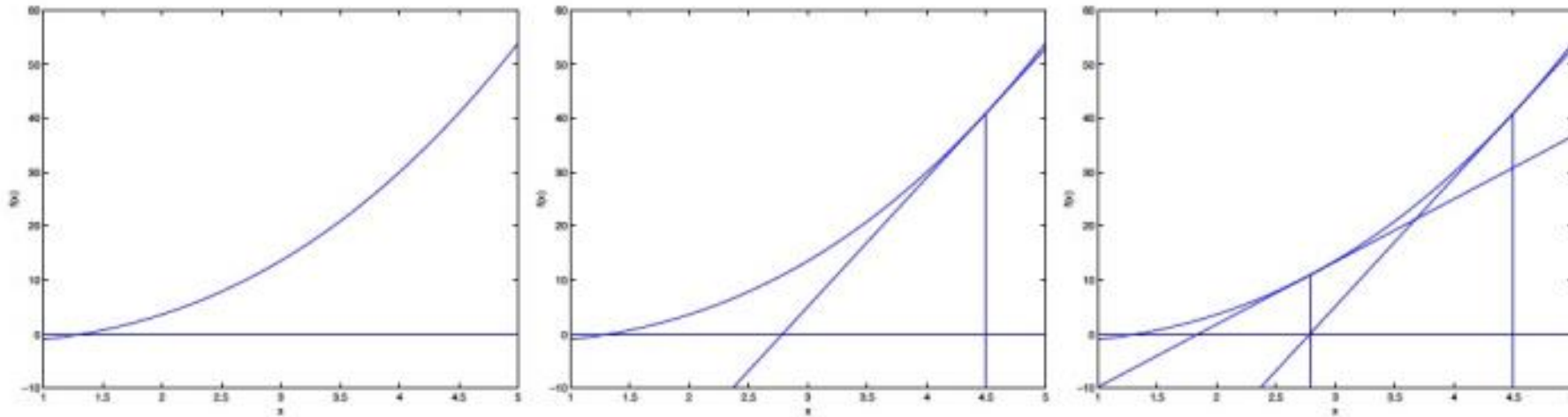- Update rule: suppose $f''(x_t) \neq 0$,

$$x = x_t - \frac{f'(x_t)}{f''(x_t)}$$

# Newton's Methods

- In numerical analysis, Newton's Methods is to find successively better approximations to the roots of a real-valued function, $(i.e, f(z) = 0)$.

$$z = z_t - \frac{f(z_t)}{f'(z_t)}$$



- In optimization, we want to find the stationary point $f'(x_t) = 0$, i.e.,

$$x = x_t - \frac{f'(x_t)}{f''(x_t)}$$

# Newton's Methods

- **Advantage:**

  ➢ More accurate local approximation of the objective,

  ➢ The convergence is much faster.

- **Disadvantage:**

  ➢ Need to compute the second derivatives

  ➢ Need to compute the inverse of Hessian (time/storage consuming)

# Go back to logistic regression

# **Gradient $\nabla E(w)$**

$$\nabla E(w) = \sum_{i=1}^{m}\left[-y_i x_i + \frac{e^{w^T x_i}}{1+e^{w^T x_i}} x_i\right] = \sum_{i=1}^{m}\left[z(w^T x_i) - y_i\right] x_i = X^T(\hat{y} - y)$$

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{md} \end{pmatrix} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{pmatrix} \in \mathbb{R}^{m \times (d+1)}, \hat{y} = \begin{pmatrix} z(w^T x_1) \\ z(w^T x_2) \\ \vdots \\ z(w^T x_m) \end{pmatrix} \in \mathbb{R}^m, y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

- Apply the Newton's method to the logistic regression,

$$x = x_t - [\nabla^2 f(x_t)]^{-1} \nabla f(x_t) \implies w = w_t - H(w_t)^{-1} \nabla E(w_t)$$

- Need to solve,

$$H = \nabla^2 E(w) = \frac{\nabla E(w)}{\nabla w} = ?$$

# Gradient $\nabla E(w)$

$$\nabla E(w) = \sum_{i=1}^{m}\left[z(w^T x_i) - y_i\right]x_i$$

$$H = \nabla^2 E(w) = \frac{\nabla E(w)}{\nabla w}$$

$$H = \sum_{i=1}^{m}\frac{\nabla\{z(w^T x_i)x_i\}}{\nabla w}$$

Identities: vector-by-vector $\dfrac{\partial y}{\partial x}$

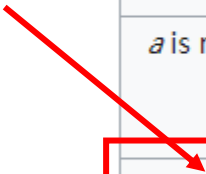| Condition | Expression | Numerator layout, i.e. by y and $x^T$ | Denominator layout, i.e. by $y^T$ and x |
|---|---|---|---|
| **a** is not a function of **x** | $\dfrac{\partial \mathbf{a}}{\partial \mathbf{x}} =$ | **0** | |
| | $\dfrac{\partial \mathbf{x}}{\partial \mathbf{x}} =$ | **I** | |
| **A** is not a function of **x** | $\dfrac{\partial \mathbf{Ax}}{\partial \mathbf{x}} =$ | $\mathbf{A}$ | $\mathbf{A}^\top$ |
| **A** is not a function of **x** | $\dfrac{\partial \mathbf{x}^\top \mathbf{A}}{\partial \mathbf{x}} =$ | $\mathbf{A}^\top$ | $\mathbf{A}$ |
| $a$ is not a function of **x**, $\mathbf{u} = \mathbf{u(x)}$ | $\dfrac{\partial a\mathbf{u}}{\partial \mathbf{x}} =$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | |
| $a = a(\mathbf{x})$, $\mathbf{u} = \mathbf{u(x)}$ | $\dfrac{\partial a\mathbf{u}}{\partial \mathbf{x}} =$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \mathbf{u}\dfrac{\partial a}{\partial \mathbf{x}}$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \dfrac{\partial a}{\partial \mathbf{x}}\mathbf{u}^\top$ |
| **A** is not a function of **x**, $\mathbf{u} = \mathbf{u(x)}$ | $\dfrac{\partial \mathbf{Au}}{\partial \mathbf{x}} =$ | $\mathbf{A}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\mathbf{A}^\top$ |
| $\mathbf{u} = \mathbf{u(x)}$, $\mathbf{v} = \mathbf{v(x)}$ | $\dfrac{\partial (\mathbf{u} + \mathbf{v})}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ | |
| $\mathbf{u} = \mathbf{u(x)}$ | $\dfrac{\partial \mathbf{g(u)}}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{g(u)}}{\partial \mathbf{u}}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\dfrac{\partial \mathbf{g(u)}}{\partial \mathbf{u}}$ |
| $\mathbf{u} = \mathbf{u(x)}$ | $\dfrac{\partial \mathbf{f(g(u))}}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{f(g)}}{\partial \mathbf{g}}\dfrac{\partial \mathbf{g(u)}}{\partial \mathbf{u}}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\dfrac{\partial \mathbf{g(u)}}{\partial \mathbf{u}}\dfrac{\partial \mathbf{f(g)}}{\partial \mathbf{g}}$ |

# Gradient $\nabla E(w)$

$$\nabla E(w) = \sum_{i=1}^{m} [z(w^T x_i) - y_i] x_i$$

$$H = \nabla^2 E(w) = \frac{\nabla E(w)}{\nabla w}$$

$$H = \sum_{i=1}^{m} \frac{\nabla\{z(w^T x_i)x_i\}}{\nabla w}$$

$a: z(w^T x_i)$

$u(w): x_i$

$\dfrac{\nabla z(w^T x_i)}{\nabla w}$ is a scalar –by-vector problem.

Identities: vector-by-vector $\dfrac{\partial y}{\partial x}$

| Condition | Expression | Numerator layout, i.e. by y and $x^T$ | Denominator layout, i.e. by $y^T$ and x |
|---|---|---|---|
| **a** is not a function of **x** | $\dfrac{\partial a}{\partial x} =$ | 0 | |
| | $\dfrac{\partial x}{\partial x} =$ | I | |
| **A** is not a function of **x** | $\dfrac{\partial Ax}{\partial x} =$ | $A$ | $A^\top$ |
| **A** is not a function of **x** | $\dfrac{\partial x^\top A}{\partial x} =$ | $A^\top$ | $A$ |
| $a$ is not a function of **x**, **u** = **u**(**x**) | $\dfrac{\partial au}{\partial x} =$ | $a\dfrac{\partial u}{\partial x}$ | |
| $a = a(x)$, **u** = **u**(**x**) | $\dfrac{\partial au}{\partial x} =$ | $a\dfrac{\partial u}{\partial x} + u\dfrac{\partial a}{\partial x}$ | $a\dfrac{\partial u}{\partial x} + \dfrac{\partial a}{\partial x}u^\top$ |
| **A** is not a function of **x**, **u** = **u**(**x**) | $\dfrac{\partial Au}{\partial x} =$ | $A\dfrac{\partial u}{\partial x}$ | $\dfrac{\partial u}{\partial x}A^\top$ |
| **u** = **u**(**x**), **v** = **v**(**x**) | $\dfrac{\partial(u+v)}{\partial x} =$ | $\dfrac{\partial u}{\partial x} + \dfrac{\partial v}{\partial x}$ | |
| **u** = **u**(**x**) | $\dfrac{\partial g(u)}{\partial x} =$ | $\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial x}$ | $\dfrac{\partial u}{\partial x}\dfrac{\partial g(u)}{\partial u}$ |
| **u** = **u**(**x**) | $\dfrac{\partial f(g(u))}{\partial x} =$ | $\dfrac{\partial f(g)}{\partial g}\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial x}$ | $\dfrac{\partial u}{\partial x}\dfrac{\partial g(u)}{\partial u}\dfrac{\partial f(g)}{\partial g}$ |

# Gradient $\nabla E(w)$

$\dfrac{\nabla z(w^T x_i)}{\nabla w}$ is a scalar –by-vector problem

Identities: scalar-by-vector $\dfrac{\partial y}{\partial x} = \nabla_x y$

| Condition | Expression | Numerator layout, i.e. by $x^T$; result is row vector | Denominator layout, i.e. by x; result is column vector |
|---|---|---|---|
| $a$ is not a function of $\mathbf{x}$ | $\dfrac{\partial a}{\partial \mathbf{x}} =$ | $\mathbf{0}^\top$ [4] | $\mathbf{0}$ [4] |
| $a$ is not a function of $\mathbf{x}$, $u = u(\mathbf{x})$ | $\dfrac{\partial au}{\partial \mathbf{x}} =$ | | $a\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial(u+v)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial u}{\partial \mathbf{x}} + \dfrac{\partial v}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial uv}{\partial \mathbf{x}} =$ | | $u\dfrac{\partial v}{\partial \mathbf{x}} + v\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial g(u)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial f(g(u))}{\partial \mathbf{x}} =$ | | $\dfrac{\partial f(g)}{\partial g}\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $\mathbf{u} = \mathbf{u}(\mathbf{x})$, $\mathbf{v} = \mathbf{v}(\mathbf{x})$ | $\dfrac{\partial(\mathbf{u}\cdot\mathbf{v})}{\partial \mathbf{x}} = \dfrac{\partial \mathbf{u}^\top \mathbf{v}}{\partial \mathbf{x}} =$ | $\mathbf{u}^\top\dfrac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^\top\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ <br> • assumes numerator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\mathbf{v} + \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}\mathbf{u}$ <br> • assumes denominator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ |

# Gradient $\nabla E(w)$

$\dfrac{\nabla z(w^T x_i)}{\nabla w}$ is a scalar –by-vector problem     $u: w^T x_i$     $z: g$     $\dfrac{\nabla z(w^T x_i)}{\nabla w} = z(w^T x_i)z(-w^T x_i)x_i$

Identities: scalar-by-vector $\dfrac{\partial y}{\partial \mathbf{x}} = \nabla_{\mathbf{x}} y$

| Condition | Expression | Numerator layout, i.e. by x<sup>T</sup>; result is row vector | Denominator layout, i.e. by x; result is column vector |
|---|---|---|---|
| $a$ is not a function of $\mathbf{x}$ | $\dfrac{\partial a}{\partial \mathbf{x}} =$ | $\mathbf{0}^{\top}$ [4] | $\mathbf{0}$ [4] |
| $a$ is not a function of $\mathbf{x}$, $u = u(\mathbf{x})$ | $\dfrac{\partial au}{\partial \mathbf{x}} =$ | | $a\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial(u+v)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial u}{\partial \mathbf{x}} + \dfrac{\partial v}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$, $v = v(\mathbf{x})$ | $\dfrac{\partial uv}{\partial \mathbf{x}} =$ | | $u\dfrac{\partial v}{\partial \mathbf{x}} + v\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial g(u)}{\partial \mathbf{x}} =$ | | $\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $u = u(\mathbf{x})$ | $\dfrac{\partial f(g(u))}{\partial \mathbf{x}} =$ | | $\dfrac{\partial f(g)}{\partial g}\dfrac{\partial g(u)}{\partial u}\dfrac{\partial u}{\partial \mathbf{x}}$ |
| $\mathbf{u} = \mathbf{u}(\mathbf{x})$, $\mathbf{v} = \mathbf{v}(\mathbf{x})$ | $\dfrac{\partial(\mathbf{u} \cdot \mathbf{v})}{\partial \mathbf{x}} = \dfrac{\partial \mathbf{u}^{\top}\mathbf{v}}{\partial \mathbf{x}} =$ | $\mathbf{u}^{\top}\dfrac{\partial \mathbf{v}}{\partial \mathbf{x}} + \mathbf{v}^{\top}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$  • assumes numerator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\mathbf{v} + \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}\mathbf{u}$  • assumes denominator layout of $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}, \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ |

# Gradient $\nabla E(w)$

$$\nabla E(w) = \sum_{i=1}^{m} \left[ z(w^T x_i) - y_i \right] x_i$$

$$H = \nabla^2 E(w) = \frac{\nabla E(w)}{\nabla w}$$

$$H = \sum_{i=1}^{m} \frac{\nabla \{ z(w^T x_i) x_i \}}{\nabla w}$$

$a: z(w^T x_i)$

$u(w): x_i$

$$\frac{\nabla z(w^T x_i)}{\nabla w} = z(w^T x_i) z(-w^T x_i) x_i$$

$$H = \sum_{i=1}^{m} x_i z(w^T x_i) z(-w^T x_i) x_i^T$$

Identities: vector-by-vector $\frac{\partial y}{\partial x}$

| Condition | Expression | Numerator layout, i.e. by y and $x^T$ | Denominator layout, i.e. by $y^T$ and x |
|---|---|---|---|
| **a** is not a function of **x** | $\dfrac{\partial \mathbf{a}}{\partial \mathbf{x}} =$ | **0** | |
| | $\dfrac{\partial \mathbf{x}}{\partial \mathbf{x}} =$ | **I** | |
| **A** is not a function of **x** | $\dfrac{\partial \mathbf{A}\mathbf{x}}{\partial \mathbf{x}} =$ | $\mathbf{A}$ | $\mathbf{A}^\top$ |
| **A** is not a function of **x** | $\dfrac{\partial \mathbf{x}^\top \mathbf{A}}{\partial \mathbf{x}} =$ | $\mathbf{A}^\top$ | $\mathbf{A}$ |
| $a$ is not a function of **x**, $\mathbf{u} = \mathbf{u}(\mathbf{x})$ | $\dfrac{\partial a\mathbf{u}}{\partial \mathbf{x}} =$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | |
| $a = a(\mathbf{x}), \mathbf{u} = \mathbf{u}(\mathbf{x})$ | $\dfrac{\partial a\mathbf{u}}{\partial \mathbf{x}} =$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \mathbf{u}\dfrac{\partial a}{\partial \mathbf{x}}$ | $a\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \dfrac{\partial a}{\partial \mathbf{x}}\mathbf{u}^\top$ |
| **A** is not a function of **x**, $\mathbf{u} = \mathbf{u}(\mathbf{x})$ | $\dfrac{\partial \mathbf{A}\mathbf{u}}{\partial \mathbf{x}} =$ | $\mathbf{A}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\mathbf{A}^\top$ |
| $\mathbf{u} = \mathbf{u}(\mathbf{x}), \mathbf{v} = \mathbf{v}(\mathbf{x})$ | $\dfrac{\partial (\mathbf{u} + \mathbf{v})}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}} + \dfrac{\partial \mathbf{v}}{\partial \mathbf{x}}$ | |
| $\mathbf{u} = \mathbf{u}(\mathbf{x})$ | $\dfrac{\partial \mathbf{g}(\mathbf{u})}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{g}(\mathbf{u})}{\partial \mathbf{u}}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\dfrac{\partial \mathbf{g}(\mathbf{u})}{\partial \mathbf{u}}$ |
| $\mathbf{u} = \mathbf{u}(\mathbf{x})$ | $\dfrac{\partial \mathbf{f}(\mathbf{g}(\mathbf{u}))}{\partial \mathbf{x}} =$ | $\dfrac{\partial \mathbf{f}(\mathbf{g})}{\partial \mathbf{g}}\dfrac{\partial \mathbf{g}(\mathbf{u})}{\partial \mathbf{u}}\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}$ | $\dfrac{\partial \mathbf{u}}{\partial \mathbf{x}}\dfrac{\partial \mathbf{g}(\mathbf{u})}{\partial \mathbf{u}}\dfrac{\partial \mathbf{f}(\mathbf{g})}{\partial \mathbf{g}}$ |

# Gradient $\nabla E(w)$

$$\nabla E(\boldsymbol{w}) = \sum_{i=1}^{m}\left[-y_i \boldsymbol{x_i} + \frac{e^{\boldsymbol{w}^T \boldsymbol{x_i}}}{1+e^{\boldsymbol{w}^T \boldsymbol{x_i}}}\boldsymbol{x_i}\right] = \sum_{i=1}^{m}\left[z(\boldsymbol{w}^T \boldsymbol{x_i}) - y_i\right]\boldsymbol{x_i} = \boldsymbol{X}^T(\widehat{\boldsymbol{y}} - \boldsymbol{y})$$

$$\boldsymbol{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{md} \end{pmatrix} = \begin{pmatrix} \boldsymbol{x_1^T} \\ \boldsymbol{x_2^T} \\ \vdots \\ \boldsymbol{x_m^T} \end{pmatrix} \in \mathbb{R}^{m \times (d+1)}, \widehat{\boldsymbol{y}} = \begin{pmatrix} z(\boldsymbol{w}^T \boldsymbol{x_1}) \\ z(\boldsymbol{w}^T \boldsymbol{x_2}) \\ \vdots \\ z(\boldsymbol{w}^T \boldsymbol{x_m}) \end{pmatrix} \in \mathbb{R}^m, \boldsymbol{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

$$\boldsymbol{H} = \nabla^2 E(\boldsymbol{w}) = \frac{\nabla E(\boldsymbol{w})}{\nabla \boldsymbol{w}} = \sum_{i=1}^{m} \boldsymbol{x_i} z(\boldsymbol{w}^T \boldsymbol{x_i}) z(-\boldsymbol{w}^T \boldsymbol{x_i}) \boldsymbol{x_i^T} = \boldsymbol{X}^T \boldsymbol{R} \boldsymbol{X}$$

$\boldsymbol{R} \in \mathbb{R}^{m \times m}$ is a diagonal matrix with elements $R_{ii} = z(\boldsymbol{w}^T \boldsymbol{x_i}) z(-\boldsymbol{w}^T \boldsymbol{x_i})$

- Apply the Newton's method to the logistic regression,

$$\boldsymbol{w} = \boldsymbol{w}_t - \boldsymbol{H}(\boldsymbol{w}_t)^{-1} \nabla E(\boldsymbol{w}_t)$$

# Compare with Linear Regression

For the linear regression with the sum-of-squares error function, we have,

$$E(\boldsymbol{w}) = \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|^2 = (\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})^T (\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y})$$

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

$$\boldsymbol{H} = \nabla^2 E(\boldsymbol{w}) = \frac{\nabla E(\boldsymbol{w})}{\nabla \boldsymbol{w}} = \boldsymbol{X}^T \boldsymbol{X}$$

$\boldsymbol{H}$ is a constant: the error function is quadratic.

Apply the Newton's method to the linear regression,

$$\boldsymbol{w} = \boldsymbol{w}_t - \boldsymbol{H}(\boldsymbol{w}_t)^{-1} \nabla E(\boldsymbol{w}_t)$$

$$\boldsymbol{w} = \boldsymbol{w}_t - (\boldsymbol{X}^T \boldsymbol{X})^{-1}(\boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}_t - \boldsymbol{X}^T \boldsymbol{y}) = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y} \quad \text{Closed-form}$$

The Newton method gives the exact solution in one step.

# Summary

## Linear Regression

- ➤ **Problem**
  - Use hyperplanes to approximate real values
- ➤ **Error (Cost) function**
  - Least square
  - $E(\boldsymbol{w})$: continuous, differentiable, **convex**
- ➤ **Algorithm**
  - ➤ Analytic solution with pseudo-inverse

# Summary

## Logistic Regression

➤ **Problem**
- $P(+1|x)$ as target and $z\left(\boldsymbol{w^T x_i}\right)$ as hypotheses

➤ **Error (Cost) Function**
- Negative log-likelihood (cross-entropy)
- $E(\boldsymbol{w})$: continuous, differentiable, twice-differentiable, **convex**

➤ **Optimization**
- Iterative methods, e.g., Gradient descent, Newton's method

# Exercise

$y \in \{0,1\}$

Target function:
$f(x) = p(+1|\boldsymbol{x})$ $\iff$ $p(y|\boldsymbol{x}) = \begin{cases} h(\boldsymbol{x}) & \text{for } y = 1 \\ 1 - h(\boldsymbol{x}) & \text{for } y = 0 \end{cases}$

$y \in \{-1,1\}$

Target function:
$f(x) = p(+1|\boldsymbol{x})$ $\iff$ $p(y|\boldsymbol{x}) = \begin{cases} h(\boldsymbol{x}) & \text{for } y = 1 \\ 1 - h(\boldsymbol{x}) & \text{for } y = -1 \end{cases}$

Can you derive the objective function?

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

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

$$h(\boldsymbol{x_i}) = P(+1|\boldsymbol{x_i}) \qquad \Leftrightarrow \qquad p(y|\boldsymbol{x_i}) = \begin{cases} \textcolor{red}{h(\boldsymbol{x_i})} & \text{for } y = +1 \\ \textcolor{blue}{1 - h(\boldsymbol{x_i})} & \text{for } y = -1 \end{cases}$$

$$\Leftrightarrow \quad p(y|\boldsymbol{x_i}) = \begin{cases} \textcolor{red}{h(\boldsymbol{x_i})} & \text{for } y = +1 \\ \textcolor{blue}{h(-\boldsymbol{x_i})} & \text{for } y = -1 \end{cases} \qquad \Leftrightarrow \qquad p(y|\boldsymbol{x_i}) = h(y\boldsymbol{x_i})$$

$$\boxed{1 - z(s) = z(-s)}$$

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

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

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

$$likelihood(h) = \prod_{i=1}^{m} p(\boldsymbol{x_i})p(y_i|\boldsymbol{x_i}) = p(\boldsymbol{x_1})\textcolor{red}{h(\boldsymbol{x_1})}p(\boldsymbol{x_2})\textcolor{blue}{h(-\boldsymbol{x_2})} \dots p(\boldsymbol{x_m})\textcolor{blue}{h(-\boldsymbol{x_m})}$$

$$\max_{h} likelihood(h) \propto \prod_{i=1}^{m} p(y_i|\boldsymbol{x_i}) = \prod_{i=1}^{m} h(y_i\boldsymbol{x_i}) = \prod_{i=1}^{m} \theta(y_i\boldsymbol{w^T}\boldsymbol{x_i})$$

$$\min_{\boldsymbol{w}} -\sum_{i=1}^{m} \ln \theta(y_i\boldsymbol{w^T}\boldsymbol{x_i}) \qquad \Leftrightarrow \qquad \min_{\boldsymbol{w}} -\sum_{i=1}^{m} \ln 1/(1 + e^{-y_i\boldsymbol{w^T}\boldsymbol{x_i}})$$

Cross-entropy loss for $y \in \{-1,1\}$

$$\boxed{\min_{\boldsymbol{w}} -\frac{1+y_i}{2}\sum_{i=1}^{m} \ln \frac{1}{1 + e^{-\boldsymbol{w^T}\boldsymbol{x_i}}} - \frac{1-y_i}{2}\sum_{i=1}^{m} \ln \frac{1}{1 + e^{\boldsymbol{w^T}\boldsymbol{x_i}}}}$$

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

**Cross-entropy**
$$H(p, q) = -\sum_{x} p(x)\log(q(x)) \qquad p \in \left\{\frac{1+y_i}{2}, \frac{1-y_i}{2}\right\}$$
$$q \in \{h(x), 1 - h(x)\}$$

$$\min_{w} -\frac{1+y_i}{2} \sum_{i=1}^{m} \ln \frac{1}{1 + e^{-w^T x_i}} - \frac{1-y_i}{2} \sum_{i=1}^{m} \ln \frac{1}{1 + e^{w^T x_i}}$$

**Simplified function**
$$\min_{w} -\sum_{i=1}^{m} \ln 1/(1 + e^{-y_i w^T x_i})$$