香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

# Logistic Regression, Exponential Family

Junxian He

Sep 12, 2024

# Classification



CAT

Labels are discrete

# Logistic Regression

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)})$ for $i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$.
Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x)$$

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)})$ for $i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$.
Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x)$$   Link Function

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)}) \text{ for } i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$. Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x) \quad \text{Link Function}$$

There are many options of $g$....

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)})$ for $i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$. Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x)$$   Link Function

There are many options of $g$....

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

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)})$ for $i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$. Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x) \quad \text{Link Function}$$
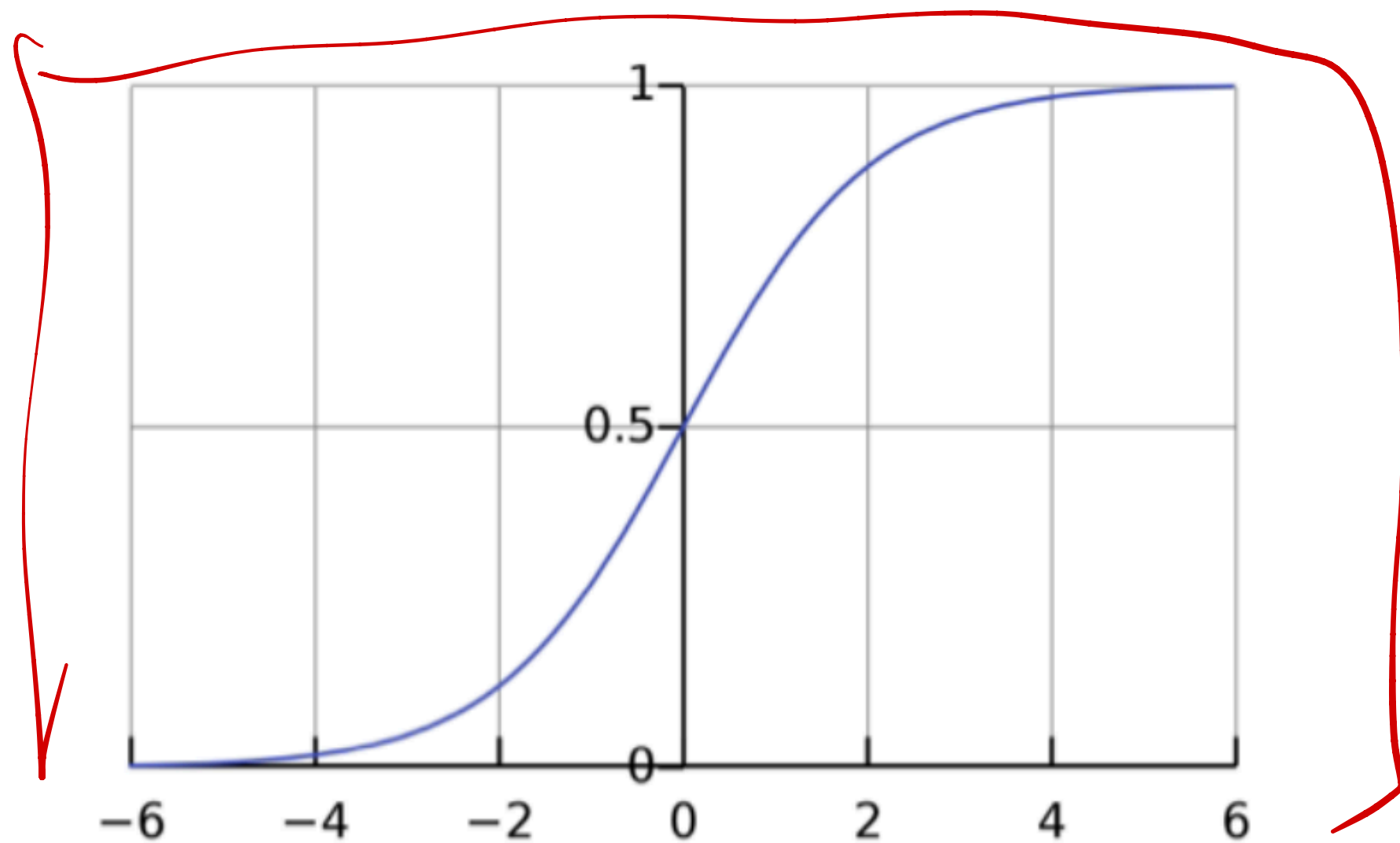
There are many options of $g$....



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

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)})$ for $i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$.
Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x) \quad \text{Link Function}$$

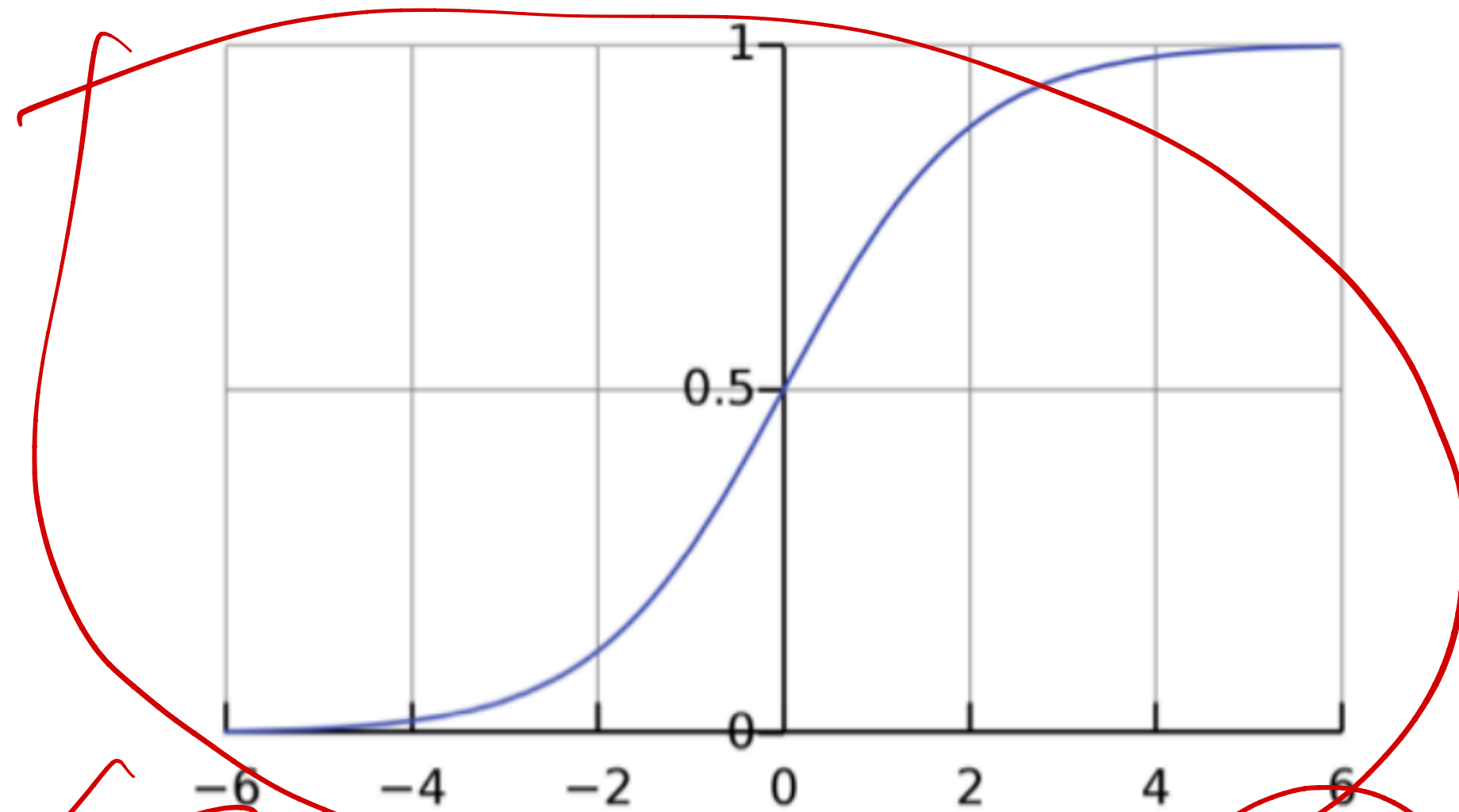There are many options of $g$….

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

How do we interpret $h_\theta(x)$?

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

3

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)}) \text{ for } i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$.
Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x) \quad \text{Link Function}$$

There are many options of $g$….

Logistic Function



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

How do we interpret $h_\theta(x)$?

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
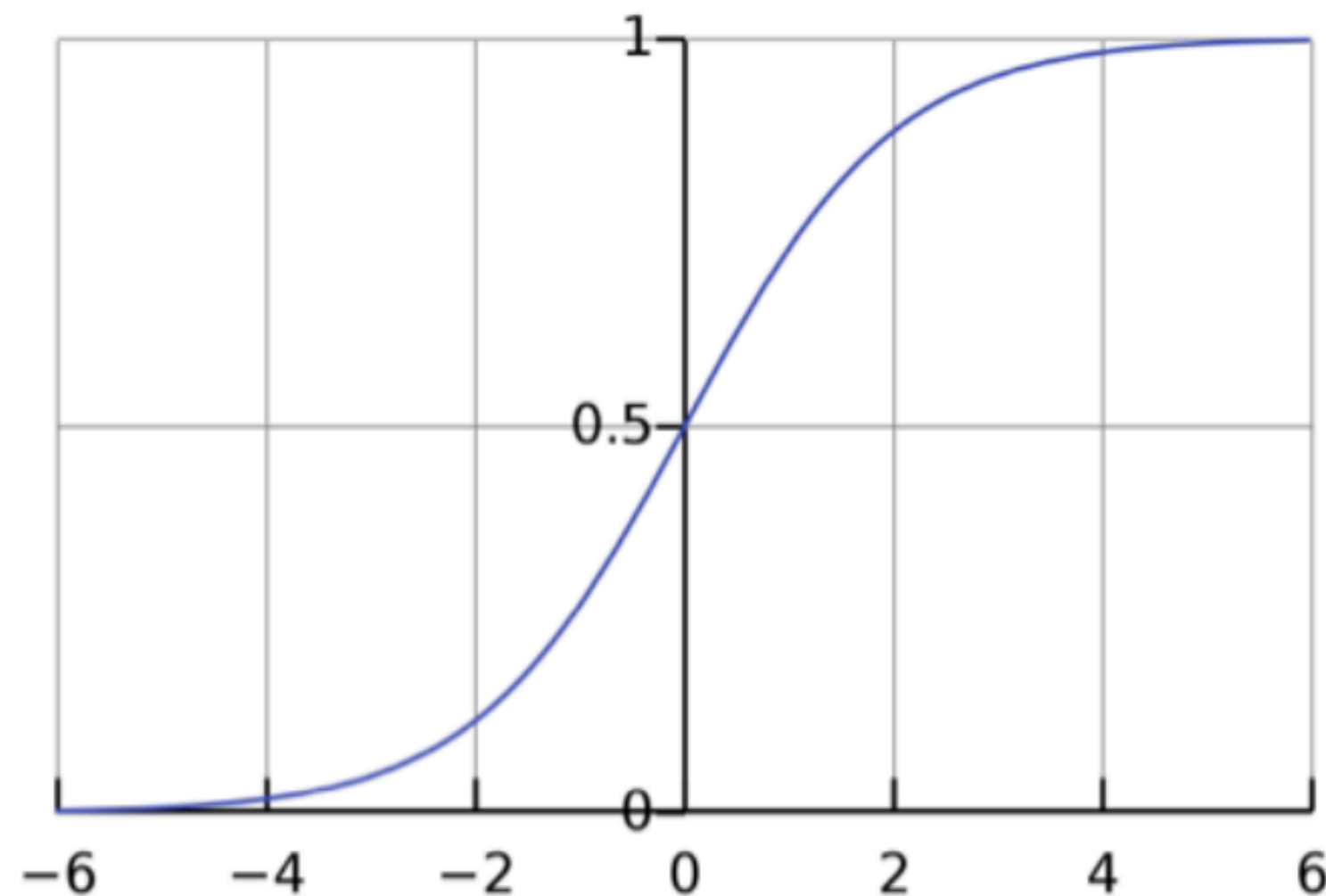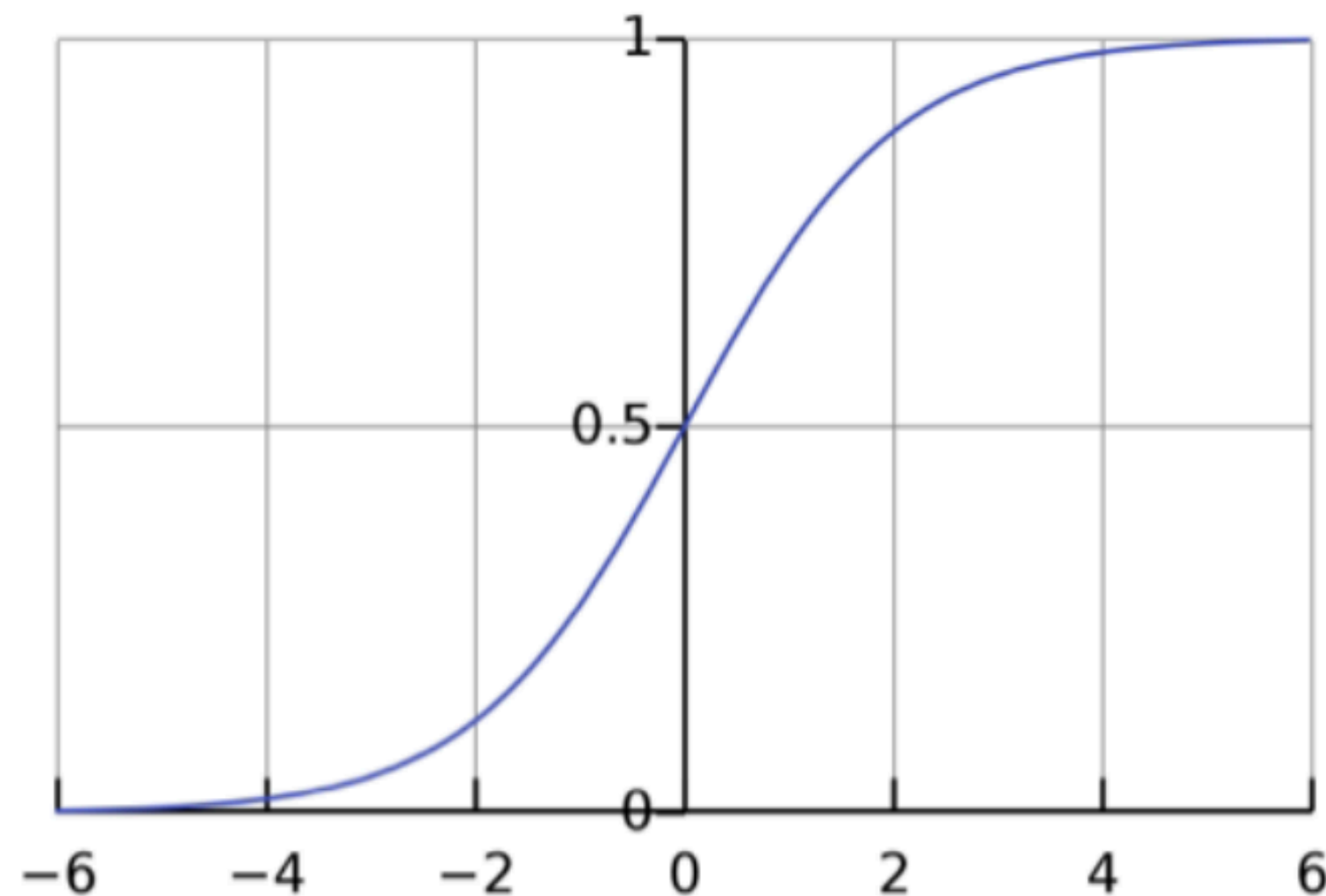$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

# Logistic Regression

Given a training set $\{(x^{(i)}, y^{(i)}) \text{ for } i = 1, \ldots, n\}$ let $y^{(i)} \in \{0, 1\}$.
Want $h_\theta(x) \in [0, 1]$. Let's pick a smooth function:

$$h_\theta(x) = g(\theta^T x) \quad \text{Link Function}$$

There are many options of $g$....

Logistic Function

$$g(z) = \frac{1}{1 + e^{-z}}. \quad \text{Sigmoid Function}$$

How do we interpret $h_\theta(x)$?

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

# Logistic Regression

Let's write the Likelihood function.  Recall:

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

Least Mean Square

Maximum likelihood estimation

4

# Logistic Regression

Let's write the Likelihood function. Recall:

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$

$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

$$h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

Then,

IID

$$L(\theta) = P(y \mid X; \theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$

Maximum likelihood estimation

4

# Logistic Regression

Let's write the Likelihood function. Recall:

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

Then,

$$L(\theta) = P(y \mid X; \theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$

$$P(y \mid x) = \begin{cases} h_\theta(x) & y = 1 \\ 1 - h_\theta(x) & y = 0 \end{cases}$$

We want to express "if-then" logics, how?

Maximum likelihood estimation

4

# Logistic Regression

Let's write the Likelihood function. Recall:

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

Then,

$$L(\theta) = P(y \mid X; \theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$

$$= \prod_{i=1}^{n} h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}}$$

We want to express "if-then" logics, how?

$y=1$

$y=0$

Maximum likelihood estimation

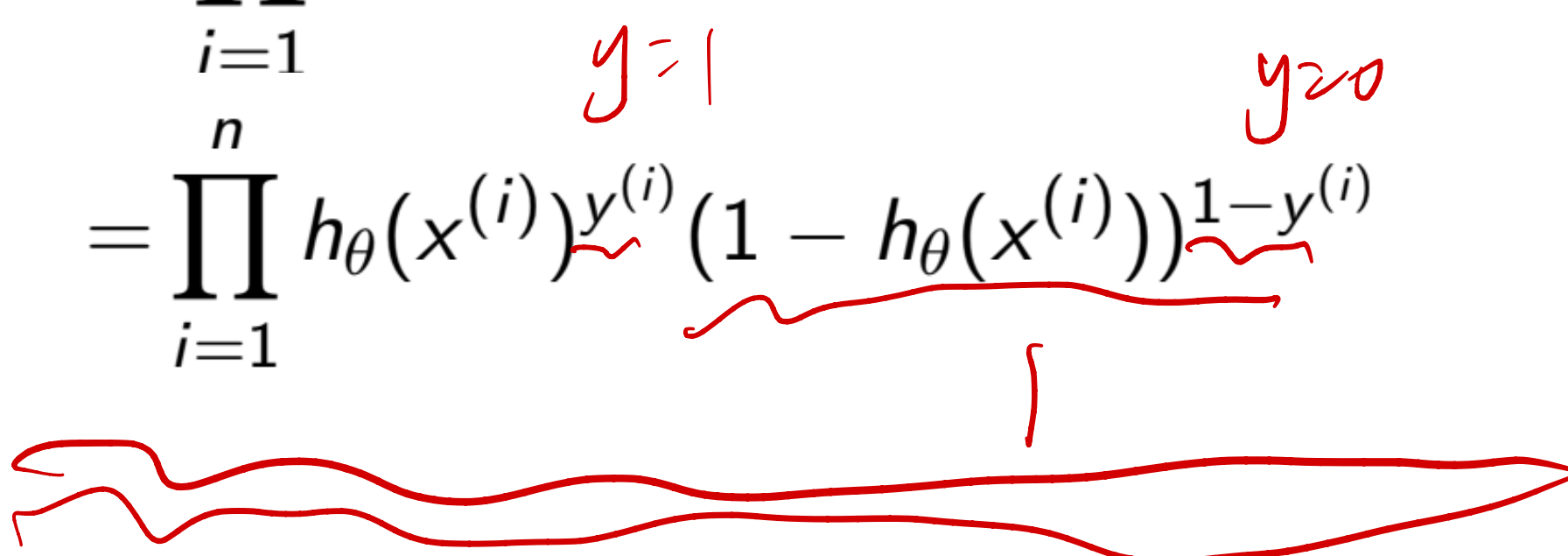# Logistic Regression

Let's write the Likelihood function. Recall:

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

Then,

$$L(\theta) = P(y \mid X; \theta) = \prod_{i=1}^{n} p(y^{(i)} \mid x^{(i)}; \theta)$$

$$= \prod_{i=1}^{n} h_\theta(x^{(i)})^{y^{(i)}}(1 - h_\theta(x^{(i)}))^{1-y^{(i)}}$$

We want to express "if-then" logics, how?

$$\theta = \arg\max_\theta \log L(\theta)$$

Taking logs to compute the log likelihood $\ell(\theta)$ we have:

$$y=1$$

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^{n} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$$

Maximum likelihood estimation

4

# Derivative of Logistic Function

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

$$
\begin{aligned}
g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\
&= \frac{1}{(1 + e^{-z})^2} \left( e^{-z} \right) \\
&= \frac{1}{(1 + e^{-z})} \cdot \left( 1 - \frac{1}{(1 + e^{-z})} \right) \\
&= g(z)(1 - g(z)).
\end{aligned}
$$

$$g(z)(1 - g(z)] \quad \frac{1}{g(z)}$$

$$l(\theta) = y \log h_\theta(x) + (1-y) \log(1 - h_\theta(x))$$

$$\frac{\partial}{\partial \theta_j} l(\theta) = y \cdot \frac{1}{g(\theta^T x)} \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$h_\theta(x)$  logistic function

$= g(\theta^T x)$

$$+ (1-y) \frac{1}{1 - g(\theta^T x)} \left( - \frac{\partial g(\theta^T x)}{\partial \theta_j} \right)$$

$$g'(z) = g(z)(1 - g(z))$$

$$= \left( y \cdot \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$$\frac{\partial \theta^T x}{\partial \theta_j} = x_j$$

$$= \left( y \cdot \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) x_j$$

$$= \left( y(1 - g(\theta^T x)) - (1-y) g(\theta^T x) \right) x_j = \left( y - g(\theta^T x) \right) x_j$$

# Gradient Descent

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$$= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x$$

$$= \left( y(1 - g(\theta^T x)) - (1-y)g(\theta^T x) \right) x_j$$

$$= (y - h_\theta(x)) x_j$$

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

# Gradient Descent

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$$= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x$$

$$= \left( y(1 - g(\theta^T x)) - (1-y)g(\theta^T x) \right) x_j$$

$$= (y - h_\theta(x)) x_j$$

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

Looks identical to LMS update rule in linear regression

# Gradient Descent

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$$= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x$$

$$= \left( y(1 - g(\theta^T x)) - (1-y)g(\theta^T x) \right) x_j$$

$$= (y - h_\theta(x)) x_j$$

$$G \, L \, M_S$$

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

Looks identical to LMS update rule in linear regression

Is this coincidence?

# Multi-Label Classification



{Cat, dog, dragon, fish, pig}

*Language    Models*

# Multi-Label Classification

Given a training set $\{(x^{(1)}, y^{(1)}), \cdots, (x^{(n)}, y^{(n)})\}, y^{(i)} \in \{1, 2, \cdots, k\}$,
we aim to model the distribution $p(y \mid x; \theta)$

# Multi-Label Classification

Given a training set $\{(x^{(1)}, y^{(1)}), \cdots, (x^{(n)}, y^{(n)})\}$, $y^{(i)} \in \{1, 2, \cdots, k\}$ , we aim to model the distribution $p(y|x;\theta)$

Categorical distribution, $p(y = k|x;\theta) = \phi_k$

$\phi_k$

$$\text{s.t. } \sum_{i=1}^{k} \phi_i = 1$$

$\phi_k \in [0,1]$

# Multi-Label Classification

Given a training set $\{(x^{(1)}, y^{(1)}), \cdots, (x^{(n)}, y^{(n)})\}$, $y^{(i)} \in \{1, 2, \cdots, k\}$, we aim to model the distribution $p(y \,|\, x; \theta)$

Categorical distribution, $p(y = k \,|\, x; \theta) = \phi_k$

$$\text{s.t. } \sum_{i=1}^{k} \phi_i = 1$$

$\phi_i = \theta_i^T x$ ?   $\notin [0, 1)$

# Softmax Function

# Softmax Function

Softmax: $\mathbb{R}^k \rightarrow \mathbb{R}^k$

# Softmax Function

Softmax: $\mathbb{R}^k \to \mathbb{R}^k$

$$\mathrm{softmax}(t_1, \ldots, t_k) = \begin{bmatrix} \dfrac{\exp(t_1)}{\sum_{j=1}^{k} \exp(t_j)} \\ \vdots \\ \dfrac{\exp(t_k)}{\sum_{j=1}^{k} \exp(t_j)} \end{bmatrix}$$

# Softmax Function

$K^{\#} = \arg\max \{t_1, t_2 - \sim t_k)$

$= \sim \qquad \nsim$

Softmax: $\mathbb{R}^k \to \mathbb{R}^k$

$logit$

$$\text{softmax}(t_1, \dots, t_k) = \begin{bmatrix} \dfrac{\exp(t_1)}{\sum_{j=1}^k \exp(t_j)} \\ \vdots \\ \dfrac{\exp(t_k)}{\sum_{j=1}^k \exp(t_j)} \end{bmatrix}$$

$\sum \phi_i = 1$

$exp()$ is monotonic

$t_i > t_j, \quad \phi_i > \phi_j$
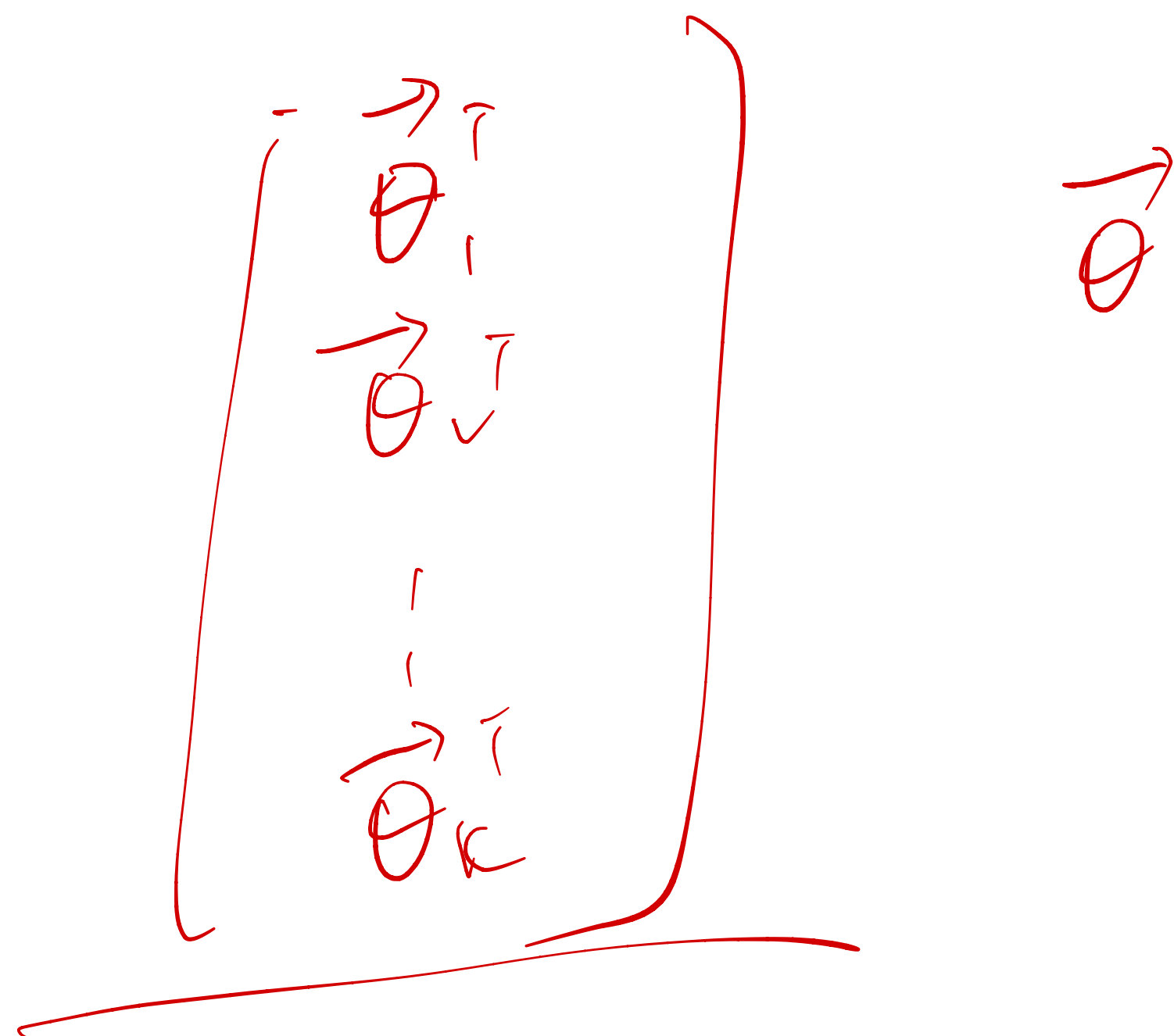
The denominator is a normalization constant

$exp(t_i) > 0 \qquad 0 < \phi_i \le 1$

9

# Multi-Label Classification

# Multi-Label Classification

Let $(t_1, \ldots, t_k) = (\theta_1^\top x, \cdots, \theta_k^\top x)$

$$NN(x) = (t_1 \: -- \: - \: t_k)$$

$$\begin{bmatrix} \vec{\theta_1}^\top \\ \vec{\theta}^\top \\ \vdots \\ \vec{\theta_k}^\top \end{bmatrix} \qquad \vec{\theta}$$

# Multi-Label Classification

Let $(t_1, \ldots, t_k) = (\theta_1^\top x, \cdots, \theta_k^\top x)$

$$\begin{bmatrix} P(y = 1 \mid x; \theta) \\ \vdots \\ P(y = k \mid x; \theta) \end{bmatrix} = \mathrm{softmax}(t_1, \cdots, t_k) = \begin{bmatrix} \dfrac{\exp(\theta_1^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \\ \vdots \\ \dfrac{\exp(\theta_k^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \end{bmatrix}$$

*CNN    LSTM*

*transformer*

10

# Multi-Label Classification

Let $(t_1, \ldots, t_k) = (\theta_1^\top x, \cdots, \theta_k^\top x)$

$$
\begin{bmatrix} P(y = 1 \mid x; \theta) \\ \vdots \\ P(y = k \mid x; \theta) \end{bmatrix} = \mathrm{softmax}(t_1, \cdots, t_k) = \begin{bmatrix} \dfrac{\exp(\theta_1^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \\ \vdots \\ \dfrac{\exp(\theta_k^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \end{bmatrix}
$$

$$
P(y = i \mid x; \theta) = \phi_i = \frac{\exp(t_i)}{\sum_{j=1}^{k} \exp(t_j)} = \frac{\exp(\theta_i^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)}
$$

# **Multi-Label Classification**

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log\left(\frac{\exp(t_y)}{\sum_{j=1}^{k}\exp(t_j)}\right) = -\log\left(\frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k}\exp(\theta_j^\top x)}\right)$$

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log \left( \frac{\exp(t_y)}{\sum_{j=1}^{k} \exp(t_j)} \right) = -\log \left( \frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \right)$$

$$\ell(\theta) = \sum_{i=1}^{n} -\log \left( \frac{\exp(\theta_{y^{(i)}}^\top x^{(i)})}{\sum_{j=1}^{k} \exp(\theta_j^\top x^{(i)})} \right)$$

$n$   # data samples

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log\left(\frac{\exp(t_y)}{\sum_{j=1}^{k} \exp(t_j)}\right) = -\log\left(\frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)}\right)$$

$$\ell(\theta) = \sum_{i=1}^{n} -\log\left(\frac{\exp(\theta_{y^{(i)}}^\top x^{(i)})}{\sum_{j=1}^{k} \exp(\theta_j^\top x^{(i)})}\right)$$ Negative log likelihood

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log\left(\frac{\exp(t_y)}{\sum_{j=1}^{k}\exp(t_j)}\right) = -\log\left(\frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k}\exp(\theta_j^\top x)}\right)$$

$$\ell(\theta) = \sum_{i=1}^{n} -\log\left(\frac{\exp(\theta_{y^{(i)}}^\top x^{(i)})}{\sum_{j=1}^{k}\exp(\theta_j^\top x^{(i)})}\right)$$  Negative log likelihood

Cross-entropy loss    $\ell_{\mathrm{ce}} : \mathbb{R}^k \times \{1, \ldots, k\} \to \mathbb{R}_{\geq 0}$

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log \left( \frac{\exp(t_y)}{\sum_{j=1}^{k} \exp(t_j)} \right) = -\log \left( \frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k} \exp(\theta_j^\top x)} \right)$$

$$\ell(\theta) = \sum_{i=1}^{n} -\log \left( \frac{\exp(\theta_{y^{(i)}}^\top x^{(i)})}{\sum_{j=1}^{k} \exp(\theta_j^\top x^{(i)})} \right) \qquad \text{\textcolor{red}{Negative log likelihood}}$$

Cross-entropy loss $\qquad \ell_{\text{ce}} : \mathbb{R}^k \times \{1, \ldots, k\} \to \mathbb{R}_{\geq 0}$

$$\ell_{\text{ce}}((t_1, \ldots, t_k), y) = -\log \left( \frac{\exp(t_y)}{\sum_{j=1}^{k} \exp(t_j)} \right)$$

11

# Multi-Label Classification

$$-\log p(y \mid x, \theta) = -\log\left(\frac{\exp(t_y)}{\sum_{j=1}^{k}\exp(t_j)}\right) = -\log\left(\frac{\exp(\theta_y^\top x)}{\sum_{j=1}^{k}\exp(\theta_j^\top x)}\right)$$

$$\ell(\theta) = \sum_{i=1}^{n} -\log\left(\frac{\exp(\theta_{y^{(i)}}^\top x^{(i)})}{\sum_{j=1}^{k}\exp(\theta_j^\top x^{(i)})}\right)$$

Negative log likelihood

Cross-entropy loss $\qquad \ell_{\mathrm{ce}} : \mathbb{R}^k \times \{1,\ldots,k\} \to \mathbb{R}_{\geq 0}$

$$t_1 - - - t_k \, , \quad y$$

$$\ell_{\mathrm{ce}}((t_1,\ldots,t_k),y) = -\log\left(\frac{\exp(t_y)}{\sum_{j=1}^{k}\exp(t_j)}\right) \qquad \ell(\theta) = \sum_{i=1}^{n} \ell_{\mathrm{ce}}((\theta_1^\top x^{(i)},\ldots,\theta_k^\top x^{(i)}),y^{(i)})$$

# The Derivative

# The Derivative

$$\frac{\partial \ell_{\text{ce}}(t, y)}{\partial t_i} = \phi_i - 1\{y = i\}$$

$$g(z)$$

$$1(e) = \begin{cases} 1 & e \text{ is true} \\ 0 & e \text{ is false} \end{cases}$$

$$\ell_{ce}\, \big( [t_1 \cdots - t_k], \, y \big) = -\log \underbrace{\frac{\exp(t_y)}{\sum\limits_{j=1}^{k} \exp(t_j)}}$$

$$\frac{\partial \ell_{ce}}{\partial t_i} = \frac{1}{-\phi_y} \cdot \frac{\partial}{\partial t_i} \frac{\exp(t_y)}{\sum\limits_{j=1}^{k} \exp(t_j)}$$

$$\left[ \frac{\exp(t_y)}{\sum\limits_{}^{k} \exp(t_j)} \right]$$

$$\phi_i = \frac{\exp(t_i)}{\sum \cdots}$$

$$= \frac{1}{-\phi_y} \left[ \boxed{\frac{\partial}{\partial t_i} \exp(t_y)} \times \frac{1}{\sum\limits_{j=1}^{k} \exp(t_j)} + \exp(t_y) \cdot \frac{\partial}{\partial t_i} \left[ \frac{1}{\sum\limits_{j=1}^{k} \exp(t_j)} \right] \right]$$

<span style="color:red">if - then $\qquad y = i$</span>

$$= -\frac{1}{\phi_y} \left[ \exp(t_y) \cdot \frac{-\exp(t_i)}{\left( \sum\limits_{j=1}^{k} \exp(t_j) \right)^2} + \right.$$

$$\left[ \frac{\exp(t_i)}{\sum\limits_{j=1}^{k} \exp(t_j)} \quad i = y \right.$$

$$\qquad \qquad \qquad \qquad 0 \qquad \text{if } y$$

$$= \begin{cases} \phi_i - 1 & i = y \\ \phi_i & i \neq y \end{cases} \qquad \color{red}{- \phi_y \cdot \phi_i}$$

$$\color{red}{+ \begin{cases} \phi_i & i = y \\ 0 & i \neq y \end{cases}}$$

# The Derivative

$$\frac{\partial \ell_{\mathrm{ce}}(t, y)}{\partial t_i} = \phi_i - 1\{y = i\} \qquad \phi_i = \frac{\exp(t_i)}{\sum_{j=1}^{k} \exp(t_j)}$$

# The Derivative

$$\frac{\partial \ell_{\mathrm{ce}}(t, y)}{\partial t_i} = \phi_i - 1\{y = i\} \qquad\qquad \phi_i = \frac{\exp(t_i)}{\sum_{j=1}^{k} \exp(t_j)}$$

Chain rule

$$\frac{\partial \ell_{\mathrm{ce}}((\theta_1^\top x, \dots, \theta_k^\top x), y)}{\partial \theta_i} = \frac{\partial \ell(t, y)}{\partial t_i} \cdot \frac{\partial t_i}{\partial \theta_i} = (\phi_i - 1\{y = i\}) \cdot x$$

$$t_i = \theta_i^\top x$$

*back propagation*

$\theta_i$ is only related to $t_i$

$$t_i = \frac{\theta_i^\top x}{\sum \theta_i^\top x}$$

12

$$\frac{\partial l_{ce}(t_1 \cdots - t_K)}{\partial \theta_i} = \frac{\partial l_{ce}}{\partial t_i} \cdot \frac{\partial t_i}{\theta_i} \qquad t_i = f(\theta_i)$$

For every $j$, $t_j \doteq g_j(\theta_i)$

$$\frac{\partial l_{ce}}{\partial \theta_i} = \sum_j \frac{\partial l_{ce}}{\partial t_j} \cdot \frac{\partial t_j}{\partial \theta_i}$$

# The Derivative

$$\frac{\partial \ell_{\mathrm{ce}}(t, y)}{\partial t_i} = \phi_i - 1\{y = i\} \qquad\qquad \phi_i = \frac{\exp(t_i)}{\sum_{j=1}^{k} \exp(t_j)}$$

## Chain rule

$$\frac{\partial \ell_{\mathrm{ce}}((\theta_1^\top x, \ldots, \theta_k^\top x), y)}{\partial \theta_i} = \frac{\partial \ell(t, y)}{\partial t_i} \cdot \frac{\partial t_i}{\partial \theta_i} = (\phi_i - 1\{y = i\}) \cdot x$$

$$\frac{\partial \ell(\theta)}{\partial \theta_i} = \sum_{j=1}^{n} (\phi_i^{(j)} - 1\{y^{(j)} = i\}) \cdot x^{(j)}$$

$$a\, LM$$

$$\theta_i \leftarrow \theta_i + \partial \sum_{j=1}^{n} (1\{y^{(i)} = i\} - \phi_i^{(j)}) \cdot x^{(j)}$$

$$\theta \leftarrow \theta_i + \partial \sum_{j=1}^{n} (y^{(i)} - h_\theta(x^{(i)})) x_i$$

# The Derivative

$$\frac{\partial \ell_{\text{ce}}(t, y)}{\partial t_i} = \phi_i - 1\{y = i\} \qquad \phi_i = \frac{\exp(t_i)}{\sum_{j=1}^{k} \exp(t_j)}$$

Chain rule

$$\frac{\partial \ell_{\text{ce}}((\theta_1^\top x, \ldots, \theta_k^\top x), y)}{\partial \theta_i} = \frac{\partial \ell(t, y)}{\partial t_i} \cdot \frac{\partial t_i}{\partial \theta_i} = (\phi_i - 1\{y = i\}) \cdot x$$

$$\frac{\partial \ell(\theta)}{\partial \theta_i} = \sum_{j=1}^{n} (\phi_i^{(j)} - 1\{y^{(j)} = i\}) \cdot x^{(j)}$$

Intuitive explanation of the rule?

$> 0$

$< 0$

$\forall \downarrow$

$\theta^{new} = \theta^{old} - \alpha \cdot \text{coeffient} \cdot \vec{x}$

$\vec{\theta_i}$ to label $i$

$y^{(j)} = i \quad (\phi_i^{(j)} - 1(y^{(j)} = i)) < 0$

$y^{(j)} \neq i \quad \text{coeffi---} \geq 0$

12

$$\theta_i^{new} = \theta_i^{old} - [\phi_i^{old} - \mathbb{1}_{(y=i)}] x$$

$$t_i = \theta_i^{\top} \boxed{x}$$

$$\begin{cases} \underline{t_i^{new}} = \theta_i^{old^{\top}} x - [\phi_i^{old} - \mathbb{1}_{(y=i)}] \underline{\boxed{x^{\top} x}} \\ t_i^{old} = \theta_i^{old^{\top}} x \end{cases}$$

$> 0$

$\downarrow$

$$t_i^{new} > t_i^{old}$$

$$\begin{cases} y=i \quad < 0 \\ y \neq i \quad > 0 \end{cases}$$

$$\begin{cases} t_i^{new} > t_i^{old} \quad y=i \\ t_i^{new} < t_i^{old} \quad y \neq i \end{cases}$$

# Another Optimization Method — Newton's Method

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.

$$\nabla_\theta l(\theta) = 0$$

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$. $\qquad \nabla_\theta l(\theta) = 0$

▶ This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$

# Another Optimization Method —
# Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.       $\nabla_\theta l(\theta) = 0$

▶ This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$

Solution to a linear equation

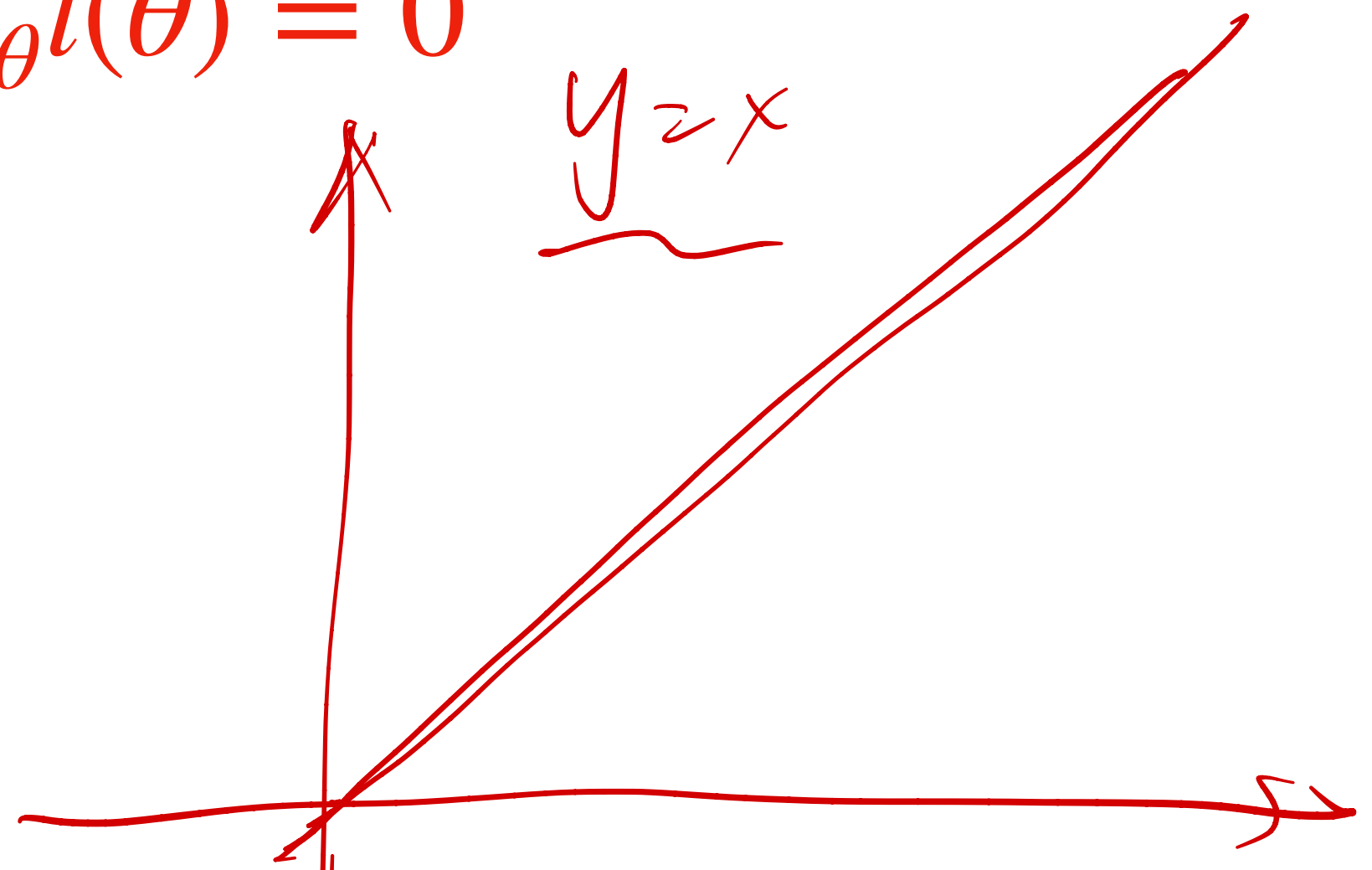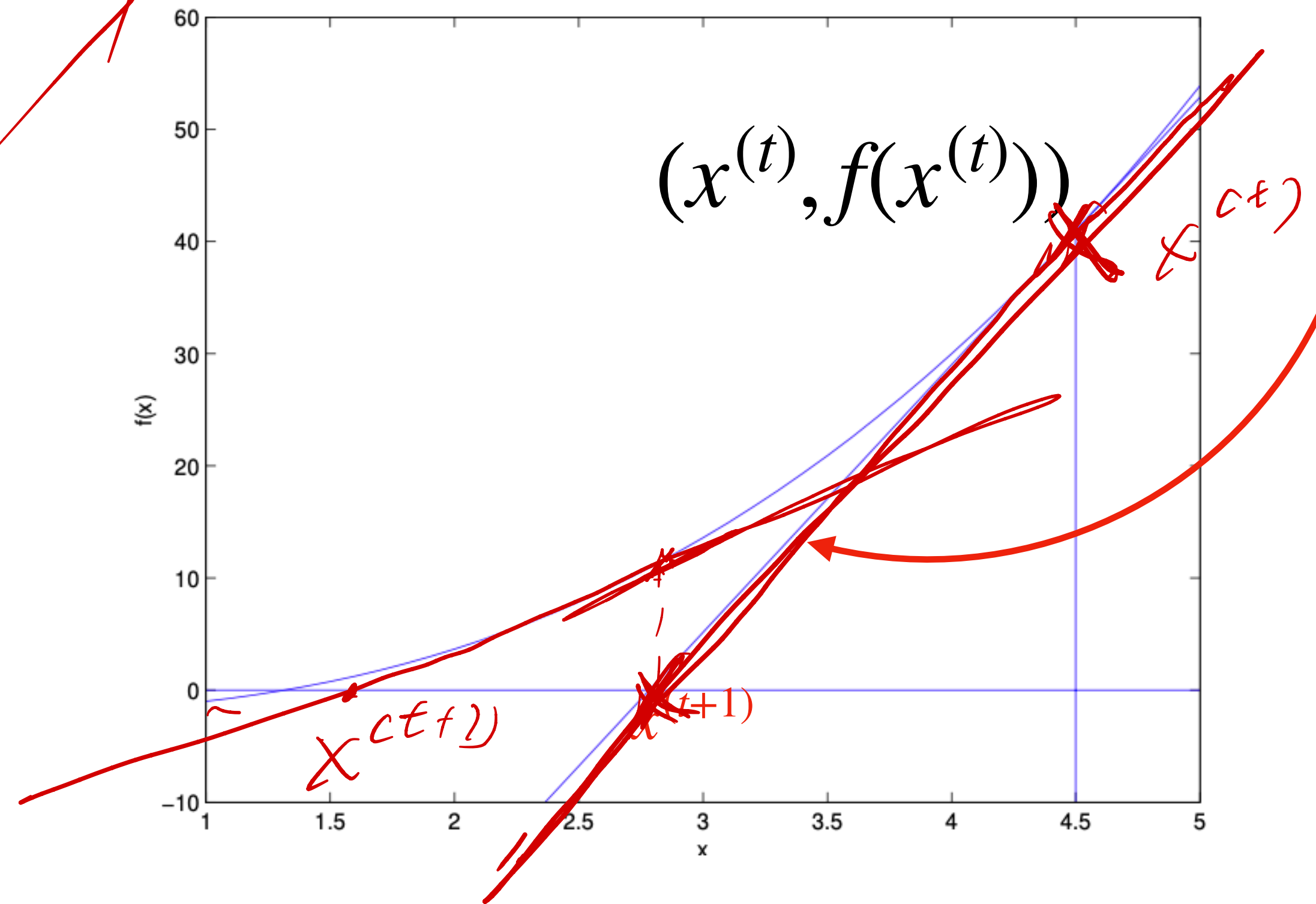$$f'(x^{(t)})x^{(t+1)} + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = 0$$

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \rightarrow \mathbb{R}$ find $x$ s.t. $f(x) = 0$.

$\nabla_\theta l(\theta) = 0$

► This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$

Solution to a linear equation

$$f'(x^{(t)})x^{(t+1)} + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = 0$$

$y = 0$

View it as a equation of $x^{(t+1)}$, and $x^{(t)}$ is a constant

$y = f'(x^{(t)})x + f(x^{(t)}) - x^{(t)}f'(x^{(t)})$

$y = x$

# Another Optimization Method — Newton's Method

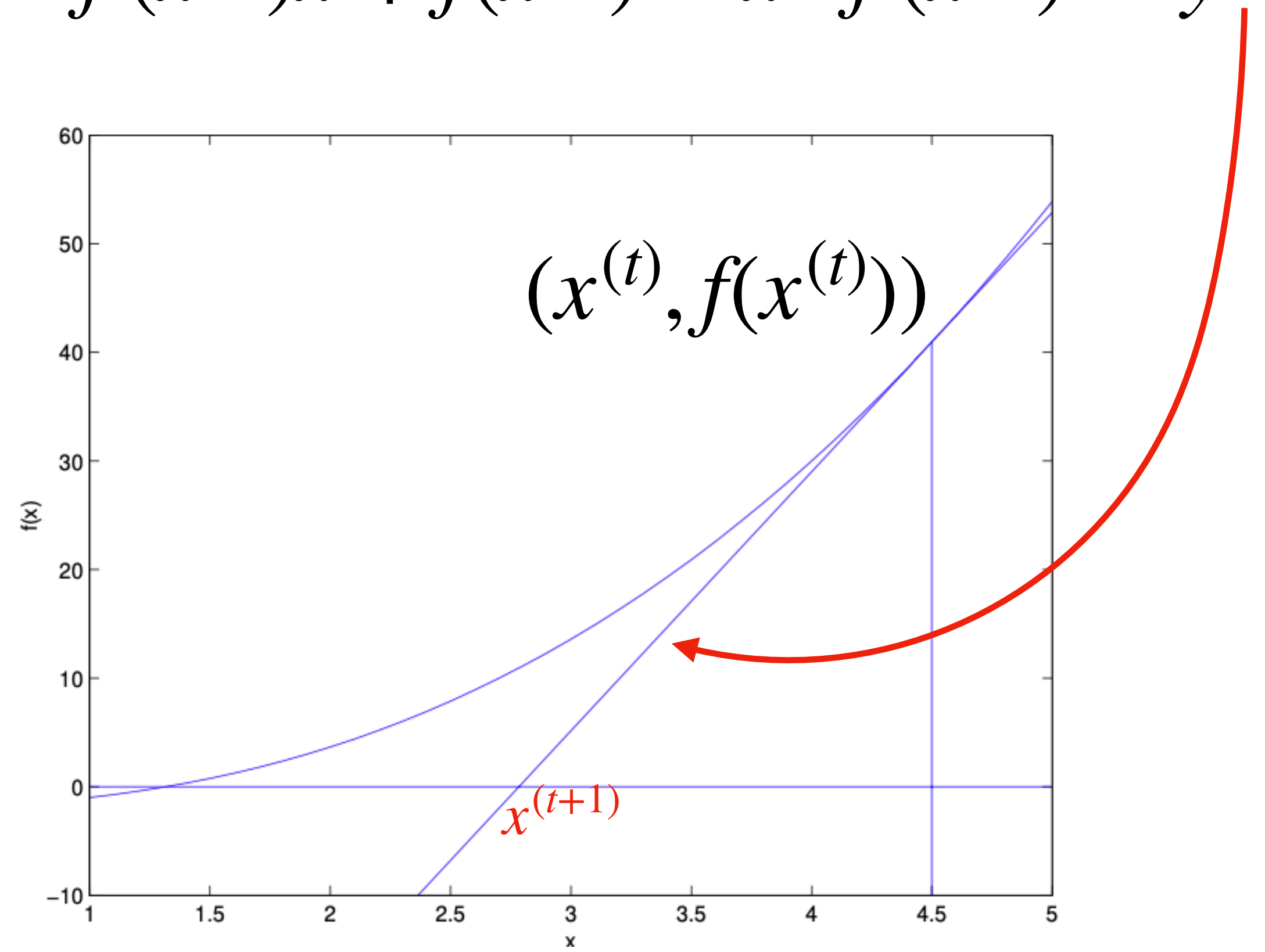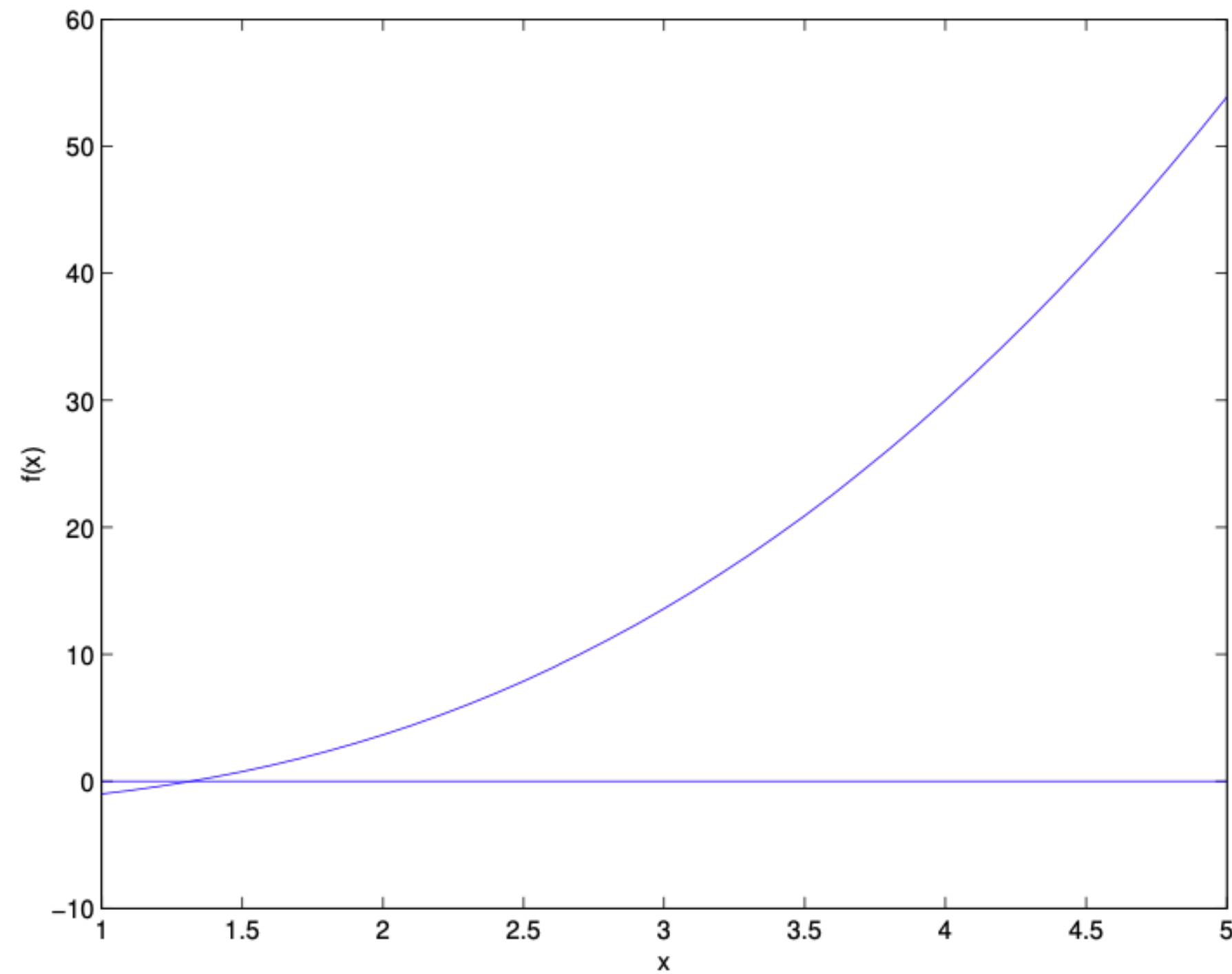$$f'(x^{(t)})x + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = y$$

$(x^{(t)}, f(x^{(t)}))$

$(x^{(t)}, f(x^{(t)}))$

$x^{(t)}$

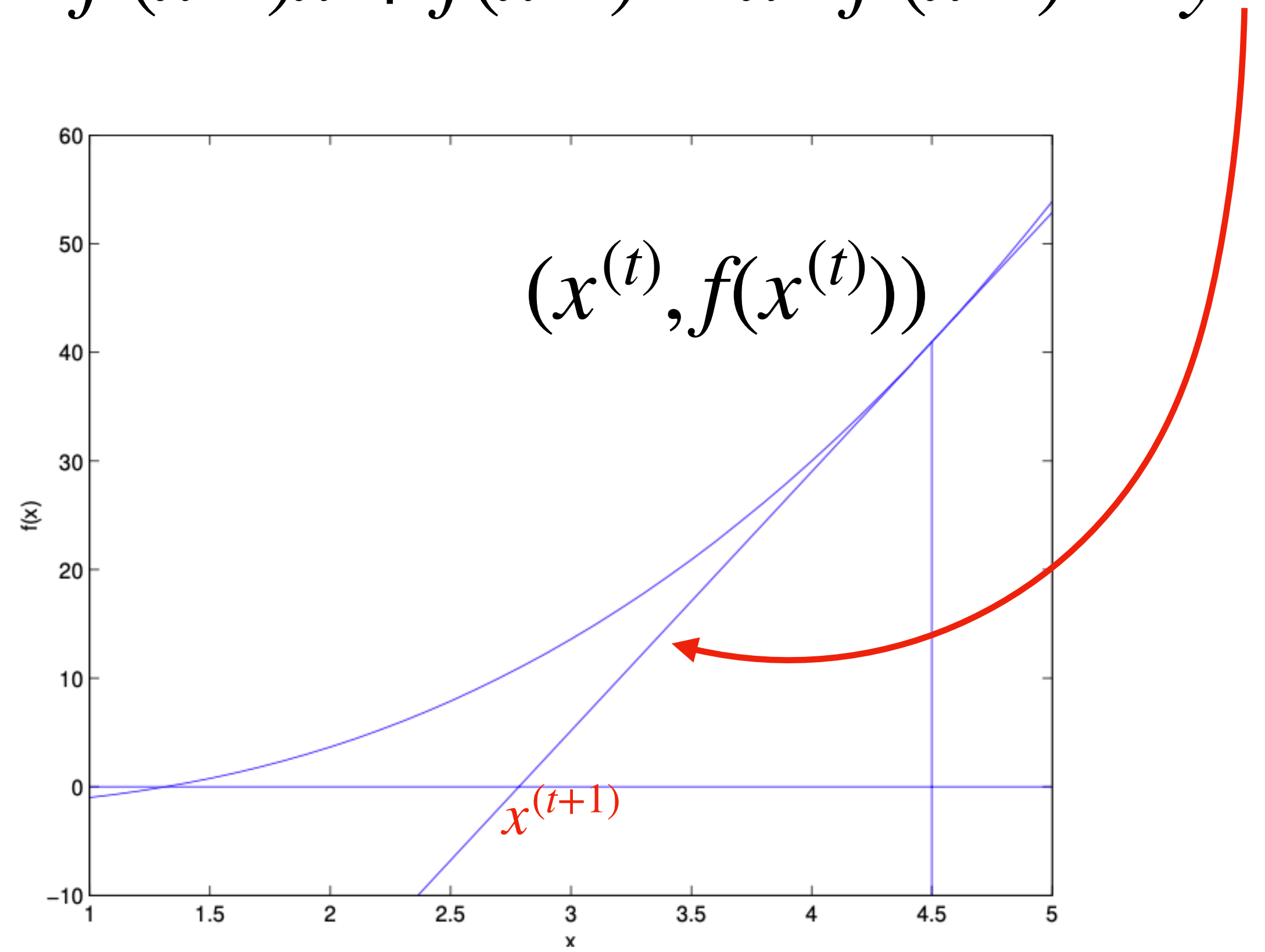$x^{(t+1)}$

$x^{(t+1)}$

# Another Optimization Method — Newton's Method

$$f'(x^{(t)})x + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = y$$

$(x^{(t)}, f(x^{(t)}))$

$x^{(t+1)}$

# Another Optimization Method — Newton's Method

$$f'(x^{(t)})x + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = y$$



f(x)

$(x^{(t)}, f(x^{(t)}))$

$x^{(t+1)}$

# Another Optimization Method — Newton's Method

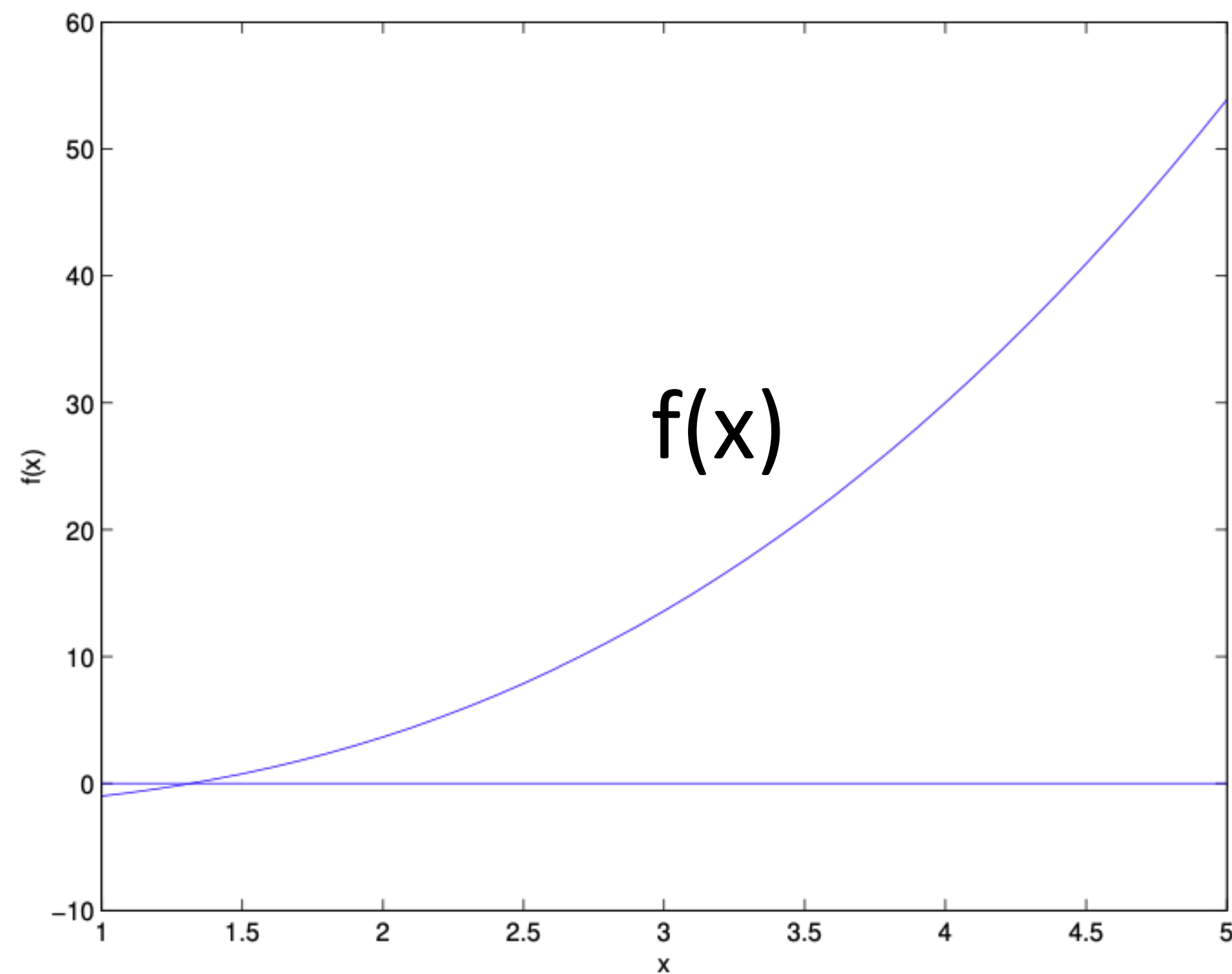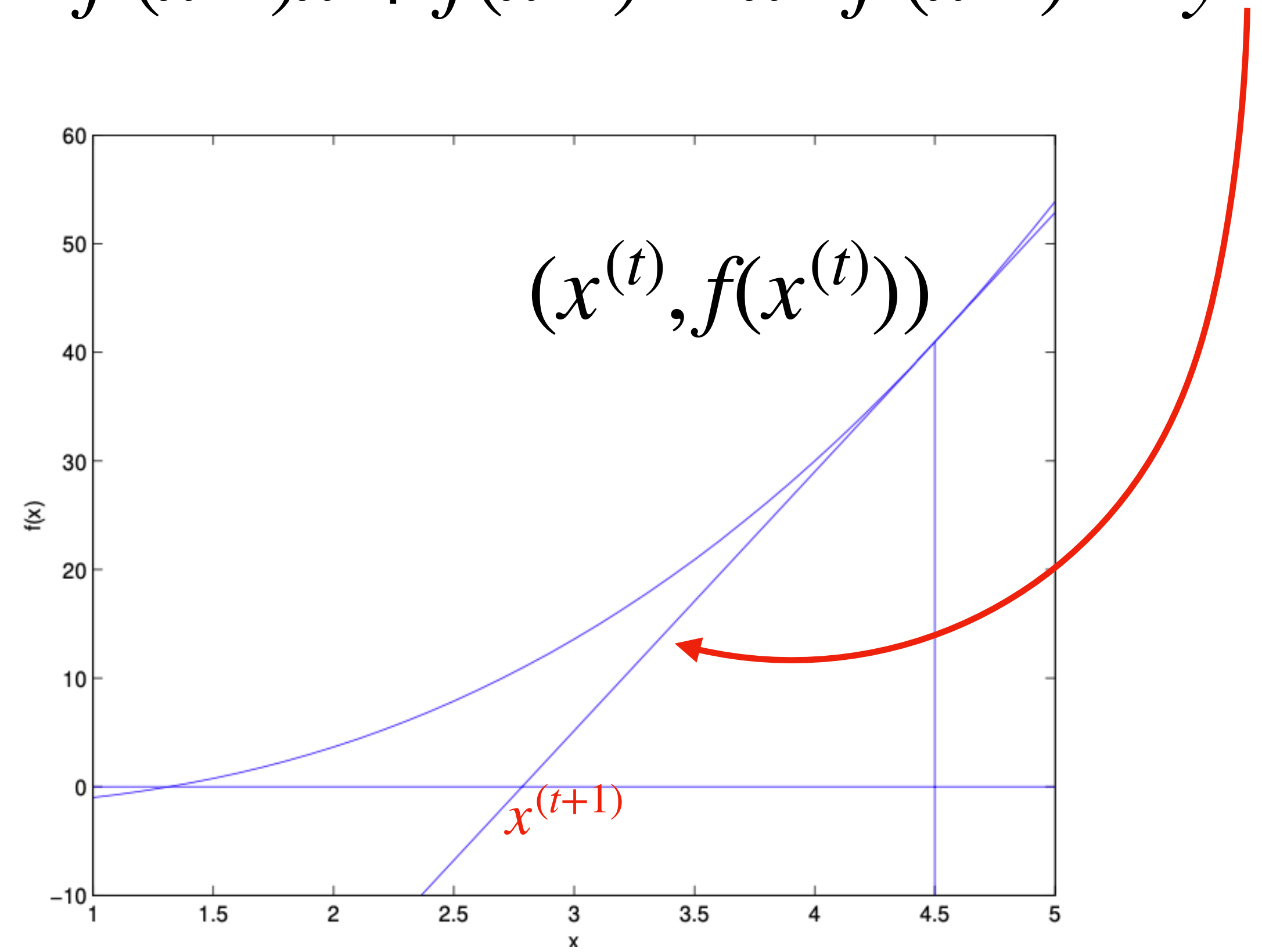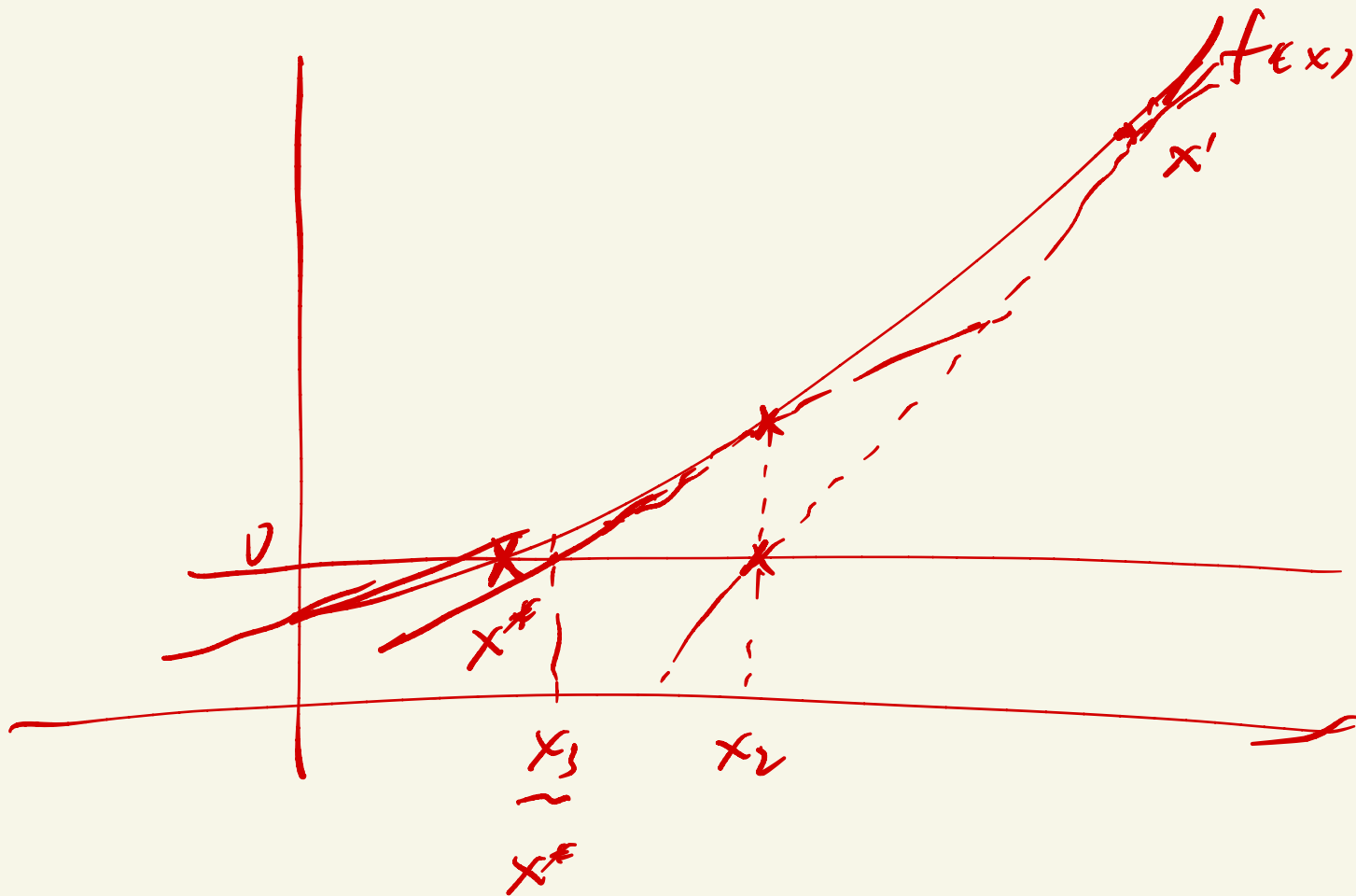$$f'(x^{(t)})x + f(x^{(t)}) - x^{(t)}f'(x^{(t)}) = y$$

$(x^{(t)}, f(x^{(t)}))$

f(x)

x*

$x^{(t+1)}$

$f(x)$

$x'$

$v$

$x^{\#}$

$x_3$

$\sim$

$x^{\#}$

$x_2$

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.

$$\nabla_\theta l(\theta) = 0$$

▶ This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$

$$\theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

second-order derivative

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.     $\nabla_\theta l(\theta) = 0$

▶ This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})} \qquad \theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}.$$

▶ It may converge *very* fast (quadratic local convergence!)   Requires fewer iterations

# Another Optimization Method — Newton's Method

Given $f : \mathbb{R}^d \to \mathbb{R}$ find $x$ s.t. $f(x) = 0$.  $\nabla_\theta l(\theta) = 0$

▶ This is the update rule in 1d

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})} \qquad \theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}.$$

▶ It may converge *very* fast (quadratic local convergence!)  Requires fewer iterations

▶ For the likelihood, i.e., $f(\theta) = \nabla_\theta \ell(\theta)$ we need to generalize to a vector-valued function which has:

$$\theta^{(t+1)} = \theta^{(t)} - \left(H(\theta^{(t)})\right)^{-1} \nabla_\theta \ell(\theta^{(t)}).$$

in which $H_{i,j}(\theta) = \frac{\partial}{\partial \theta_i \partial \theta_j} \ell(\theta).$  second-order

15

local minimum

# Exponential Family

# Exponential  Family

- Exponential family unifies inference and learning for many important models

# Exponential Family

# Exponential Family

**Rough Idea** *"If P has a a special form, then inference and learning come for free"*

$$P(y; \eta) = b(y) \exp\left\{\eta^T T(y) - a(\eta)\right\}.$$

Here $y$, $a(\eta)$, and $b(y)$ are scalars. $T(y)$ same dimension as $\eta$.

# Exponential Family

**Rough Idea** *"If P has a a special form, then inference and learning come for free"*

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

$\eta$: natural parameter or canonical parameter

Here $y$, $a(\eta)$, and $b(y)$ are scalars. $T(y)$ same dimension as $\eta$.

# Exponential Family

**Rough Idea** *"If P has a a special form, then inference and learning come for free"*

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

$\eta$: natural parameter or canonical parameter

Here $y$, $a(\eta)$, and $b(y)$ are scalars. $T(y)$ same dimension as $\eta$.

$T(y)$ is called the **sufficient statistic**.

$b(y)$ is called the **base measure** – does *not* depend on $\eta$.

$a(\eta)$ is called the **log partition function** – does *not* depend on $y$.

# Exponential Family

**Rough Idea** *"If P has a a special form, then inference and learning come for free"*

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

*not related*

$\eta$: natural parameter or canonical parameter

Here $y$, $a(\eta)$, and $b(y)$ are scalars. $T(y)$ same dimension as $\eta$.

holds all information the data provides with regard to the unknown parameter values

$T(y)$ is called the **sufficient statistic**.

$b(y)$ is called the **base measure** – does *not* depend on $\eta$.

$a(\eta)$ is called the **log partition function** – does *not* depend on $y$.

*normalization*

# Exponential Family

**Rough Idea** *"If P has a a special form, then inference and learning come for free"*

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

$\eta$: natural parameter or canonical parameter

Here $y$, $a(\eta)$, and $b(y)$ are scalars. $T(y)$ same dimension as $\eta$.

holds all information the data provides with regard to the unknown parameter values

$T(y)$ is called the **sufficient statistic**.

$b(y)$ is called the **base measure** – does *not* depend on $\eta$.

$a(\eta)$ is called the **log partition function** – does *not* depend on $y$.

$$1 = \sum_y P(y; \eta) = e^{-a(\eta)} \sum_y b(y) \exp \left\{ \eta^T T(y) \right\}$$

$$\implies a(\eta) = \log \sum_y b(y) \exp \left\{ \eta^T T(y) \right\}$$

18

# Example: Bernoulli

Bernoulli random variable is an event (say flipping a coin) then:

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

Logistic binary

$$P(y|x) = h_\theta(x)^y (1 - h_\theta(x))^{1-y}$$

$$\phi = h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Example: Bernoulli

Bernoulli random variable is an event (say flipping a coin) then:

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

How do we put it in the required form?

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

# Example: Bernoulli

Bernoulli random variable is an event (say flipping a coin) then:

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

How do we put it in the required form?

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

$$
\begin{aligned}
p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\
&= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\
&= \exp \left( \left( \log \left( \frac{\phi}{1 - \phi} \right) \right) y + \log(1 - \phi) \right)
\end{aligned}
$$

# Example: Bernoulli

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}$$

$$
\begin{aligned}
p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\
&= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\
&= \exp \left( \left( \log \left( \frac{\phi}{1 - \phi} \right) \right) y + \log(1 - \phi) \right)
\end{aligned}
$$

# Example: Bernoulli

$$\eta = \theta^T x$$

GLM

$$
\begin{aligned}
p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\
&= \exp(y \log \phi + (1-y)\log(1-\phi)) \\
&= \exp\left(\left(\log\left(\frac{\phi}{1-\phi}\right)\right) y + \log(1-\phi)\right)
\end{aligned}
$$

$$P(y; \eta) = b(y)\exp\left\{\eta^T T(y) - a(\eta)\right\}$$

So then:

$$\eta = \log\frac{\phi}{1-\phi}, \quad T(y) = y, \quad a(\eta) = -\log(1-\phi).$$

$$b(y) = 1$$

$$\eta = \log\frac{\phi}{1-\phi}$$

$$\phi = \frac{1}{1+e^{-\theta^T x}}$$

$$= \eta = \theta^T x$$

Generalized linear models

20

# Example: Bernoulli

$$
\begin{aligned}
p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\
&= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\
&= \exp\left(\left(\log\left(\frac{\phi}{1 - \phi}\right)\right) y + \log(1 - \phi)\right)
\end{aligned}
$$

$$
P(y; \eta) = b(y) \exp\left\{\eta^T T(y) - a(\eta)\right\}
$$

So then:

$$
\eta = \log \frac{\phi}{1 - \phi}, \; T(y) = y, \; a(\eta) = -\log(1 - \phi).
$$

$$
b(y) = 1
$$

We need to show $a(\eta)$ is a function of $\log \dfrac{\phi}{1 - \phi}$

# Example: Bernoulli

# Example: Bernoulli

We first observe that:

$$\eta = \log \frac{\phi}{1 - \phi} \implies e^{\eta}(1 - \phi) = \phi$$

$$e^{\eta} = (e^{\eta} + 1)\phi \implies \phi = \frac{1}{1 + e^{-\eta}}$$

# Example: Bernoulli

We first observe that:

$$\eta = \log \frac{\phi}{1 - \phi} \implies e^\eta (1 - \phi) = \phi$$

$$e^\eta = (e^\eta + 1)\phi \implies \phi = \frac{1}{1 + e^{-\eta}}$$

Now, we plug into $\log(1 - \phi)$ and we verify:

$$a(\eta) = \log(1 - \phi) = \log \frac{e^{-\eta}}{1 + e^{-\eta}} = -\log(1 + e^\eta).$$

# Example: Bernoulli

We first observe that:

$$\eta = \log \frac{\phi}{1 - \phi} \implies e^{\eta}(1 - \phi) = \phi$$

$$e^{\eta} = (e^{\eta} + 1)\phi \implies \phi = \frac{1}{1 + e^{-\eta}}$$

Now, we plug into $\log(1 - \phi)$ and we verify:

$$a(\eta) = \log(1 - \phi) = \log \frac{e^{-\eta}}{1 + e^{-\eta}} = -\log(1 + e^{\eta}).$$

We have verified Bernoulli distribution is in the exponential family

# Example: Gaussian with Fixed Variance $\sigma^2 = 1$

# Example: Gaussian with Fixed Variance $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(y - \mu)^2 \right\}.$$

Can we put it in the exponential family form?

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

# Example: Gaussian with Fixed Variance $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(y - \mu)^2 \right\}.$$

Can we put it in the exponential family form?

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

Multiply out the square and group terms:

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -y^2/2 \right\} \exp \left\{ \mu y - \frac{1}{2}\mu^2 \right\}.$$

# Example: Gaussian with Fixed Variance $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(y - \mu)^2\right\}.$$

Can we put it in the exponential family form?

$$P(y; \eta) = b(y) \exp\left\{\eta^T T(y) - a(\eta)\right\}.$$

Multiply out the square and group terms:

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left\{-y^2/2\right\} \exp\left\{\mu y - \frac{1}{2}\mu^2\right\}.$$

$$\eta = \mu, \; T(y) = y, \; a(\eta) = \frac{1}{2}\eta^2.$$

# Example: Gaussian with Fixed Variance $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2}(y - \mu)^2 \right\}.$$

Can we put it in the exponential family form?

$$P(y; \eta) = b(y) \exp \left\{ \eta^T T(y) - a(\eta) \right\}.$$

Multiply out the square and group terms:

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -y^2/2 \right\} \exp \left\{ \mu y - \frac{1}{2}\mu^2 \right\}.$$

In all the exponential family distribution we work with in the course, T(y) = y

$$\eta = \mu, \; T(y) = y, \; a(\eta) = \frac{1}{2}\eta^2.$$

22

# An Observation

# An Observation

Notice that for a Gaussian with mean $\mu$ we had

$$\eta = \mu, \ T(y) = y, \ a(\eta) = \frac{1}{2}\eta^2.$$

# An Observation

Notice that for a Gaussian with mean $\mu$ we had

$$\eta = \mu, \ T(y) = y, \ a(\eta) = \frac{1}{2}\eta^2.$$

$$\partial_\eta a(\eta) = \eta = \mu = \mathbb{E}[y] \text{ and } \partial_\eta^2 a(\eta) = 1 = \sigma^2 = \text{var}(y)$$

# An Observation

Notice that for a Gaussian with mean $\mu$ we had

$$\eta = \mu,\ T(y) = y,\ a(\eta) = \frac{1}{2}\eta^2.$$

$$\partial_\eta a(\eta) = \eta = \mu = \mathbb{E}[y] \text{ and } \partial_\eta^2 a(\eta) = 1 = \sigma^2 = \mathrm{var}(y)$$

Is this true for general?

# Log Partition Function

Yes! Recall that

$$a(\eta) = \log \sum_y b(y) \exp\left\{\eta^T T(y)\right\}$$

# Log Partition Function

Yes! Recall that

$$a(\eta) = \log \sum_y b(y) \exp\left\{\eta^T T(y)\right\}$$

Then, taking derivatives

$$\nabla_\eta a(\eta) = \frac{\sum_y T(y) b(y) \exp\left\{\eta^T T(y)\right\}}{\sum_y b(y) \exp\left\{\eta^T T(y)\right\}} = \mathbb{E}[T(y); \eta]$$

# Many Other Exponential Models

▶ There are many canonical exponential family models:

- ▶ Binary $\mapsto$ Bernoulli
- ▶ Multiple Classses $\mapsto$ Multinomial
- ▶ Real $\mapsto$ Gaussian
- ▶ Counts $\mapsto$ Poisson
- ▶ $\mathbb{R}_+ \mapsto$ Gamma, Exponential
- ▶ Distributions $\mapsto$ Dirichlet

# Thank You!
## Q & A