

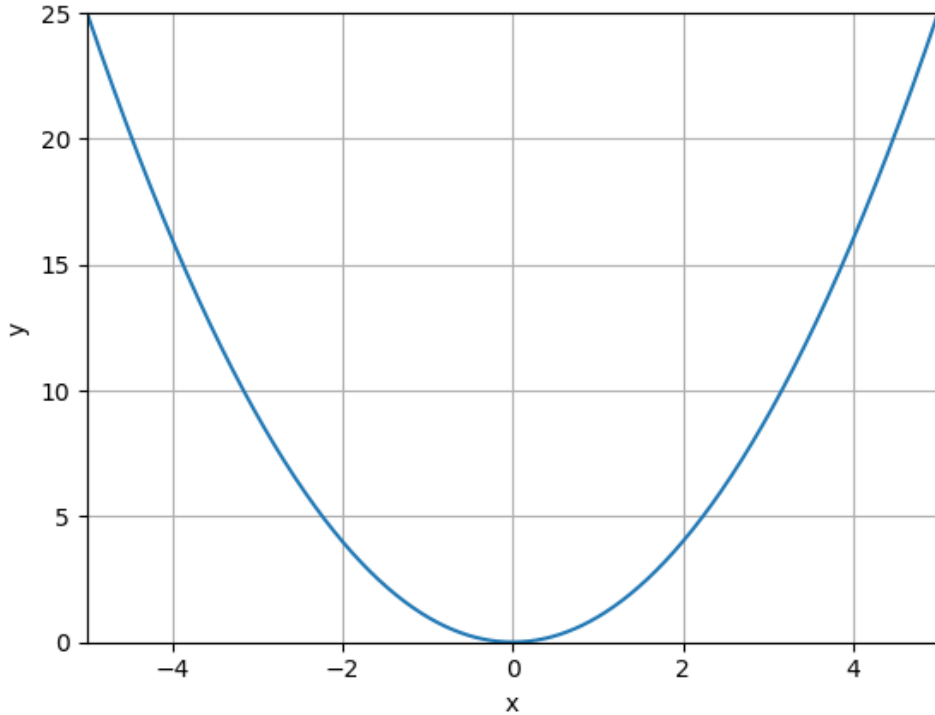
Blixtkurs i flervariabelsanalys

- 1 dimension:
 - Givet: en funktion $y = f(x)$, dvs $f: R^1 \rightarrow R^1$ (från skalär till skalär).
 - Derivataoperatorn definieras som $\frac{\delta}{\delta x} f(x) = \frac{\delta f(x)}{\delta x} = f'_x(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$.

En liten, men viktig, detalj:

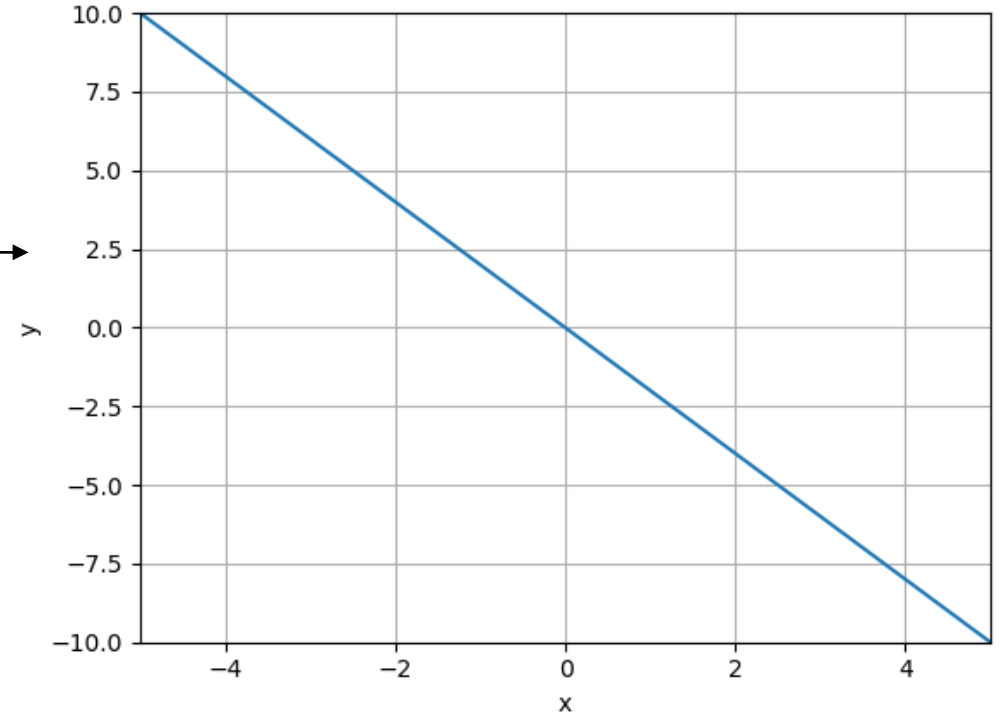
- $\left. \frac{\delta f(x)}{\delta x} \right|_{x=a} = f'_x(a)$ betecknar värdet för derivatafunktionen i punkten a . Vi kan inte skriva detta som $\frac{\delta f(a)}{\delta x}$, eftersom $\frac{\delta f(a)}{\delta x} = \frac{\delta}{\delta x} f(a) = \frac{\delta}{\delta x} (\text{constant}) = 0$.

Exempel



$$y = f(x) = x^2$$

$$\boxed{-\frac{\delta}{\delta x}}$$



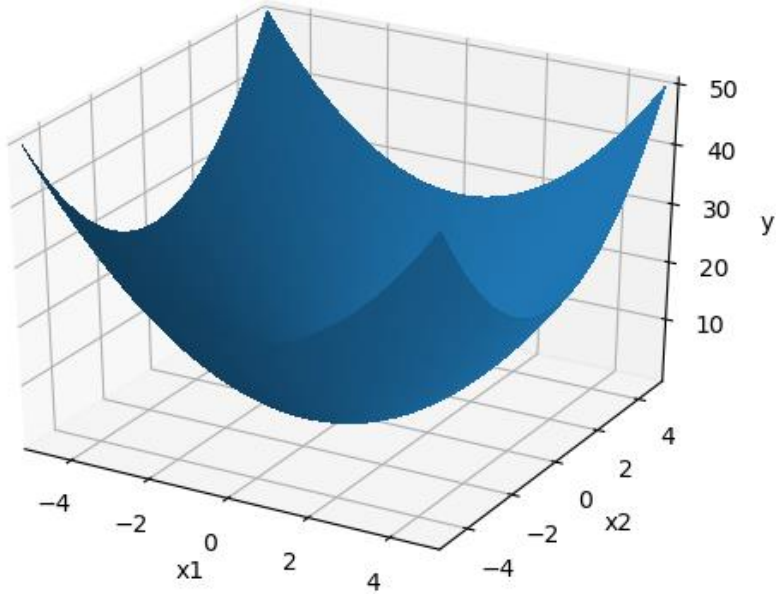
$$y = -f'_x(x) = -2x$$

- $\frac{\delta}{\delta x} f(x) = \frac{\delta f(x)}{\delta x} = f'_x(x)$ är en funktion $f'_x: R^1 \rightarrow R^1$ (från skalär till skalär).
- Notera: den negativa derivatans tecken för en punkt a anger i vilken riktning längs x-axeln från punkten $x = a$ man bör förflytta sig för att nå ett (lokalt) minimum.

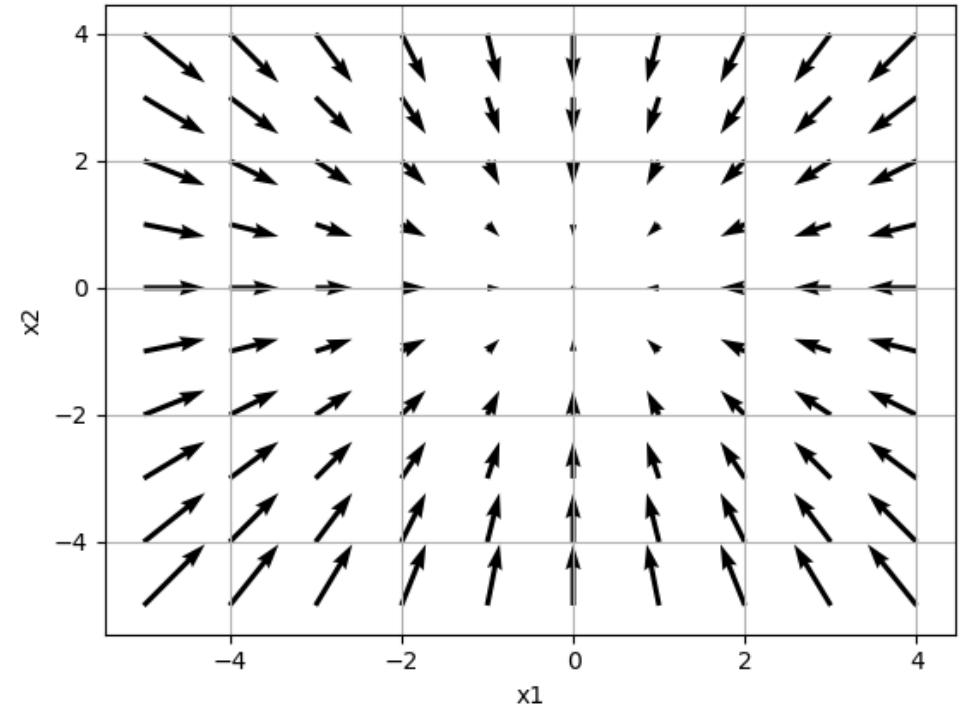
Blixtkurs i flervariabelsanalys

- $n \geq 2$ dimensioner:
 - Givet: en funktion $y = f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$, dvs $f: R^n \rightarrow R^1$ (från n -vektor till skalär)
 - Gradientoperatören $\nabla_{\mathbf{x}}$ ("nabla") ges av $\nabla_{\mathbf{x}} = \nabla_{x_1, x_2, \dots} = \left(\frac{\delta}{\delta x_1}, \frac{\delta}{\delta x_2}, \dots, \frac{\delta}{\delta x_n} \right)^T$
 - Funktionen f 's gradient är $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left(\frac{\delta f(\mathbf{x})}{\delta x_1}, \frac{\delta f(\mathbf{x})}{\delta x_2}, \dots \right)^T = (f'_{x_1}(\mathbf{x}), f'_{x_2}(\mathbf{x}), \dots)^T$
 - $\nabla_{\mathbf{x}} f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_{val}} = \nabla_{\mathbf{x}} f(x_1, x_2, \dots)|_{(x_1, x_2, \dots)=(a, b, \dots)}$ betecknar gradientvektorns faktiska värde i punkten $\mathbf{x}_{val} = (a, b, \dots)$.

Exempel



$$\longrightarrow \boxed{-\nabla_x} \longrightarrow$$



$$y = f(x_1, x_2) = x_1^2 + x_2^2$$

$$-\nabla_x f(x_1, x_2)|_{x_1=a, x_2=b} = (-2a, -2b)$$

- Gradienten är en funktion $\nabla_x f: R^n \rightarrow R^n$ (från n-vektor till n-vektor)
- Gradienten i en punkt pekar i den riktning vartåt funktionen ökar snabbast. Des storlek representerar hur stor ökningen är. Den negativa gradienten anger därför hur man bör förflytta sig i rymden av möjliga input $\mathbf{x} = (x_1, x_2, \dots)$ för att nå ett (lokalt) minimum.

Grundläggande optimering

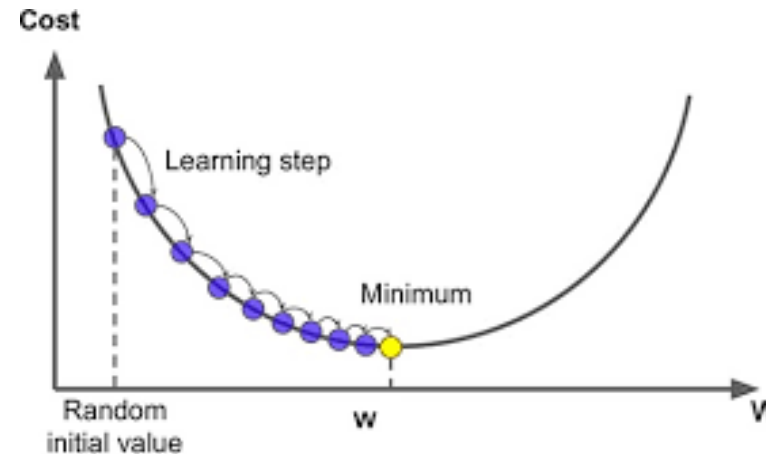
- Givet: En cost/loss-funktion $C = L(\mathbf{x}, \mathbf{w}) = L(x_1, x_2, \dots, w_1, w_2, \dots)$.
- Problem: Välj element i \mathbf{w} för att göra C så litet som möjligt.
 - Dvs \mathbf{x} är här att betrakta som konstant(er)
- Den faktiska lösningen $\mathbf{w}^{opt} = (w_1^{opt}, w_2^{opt}, \dots)$ finns bland lösningarna till ekvationssystemet:

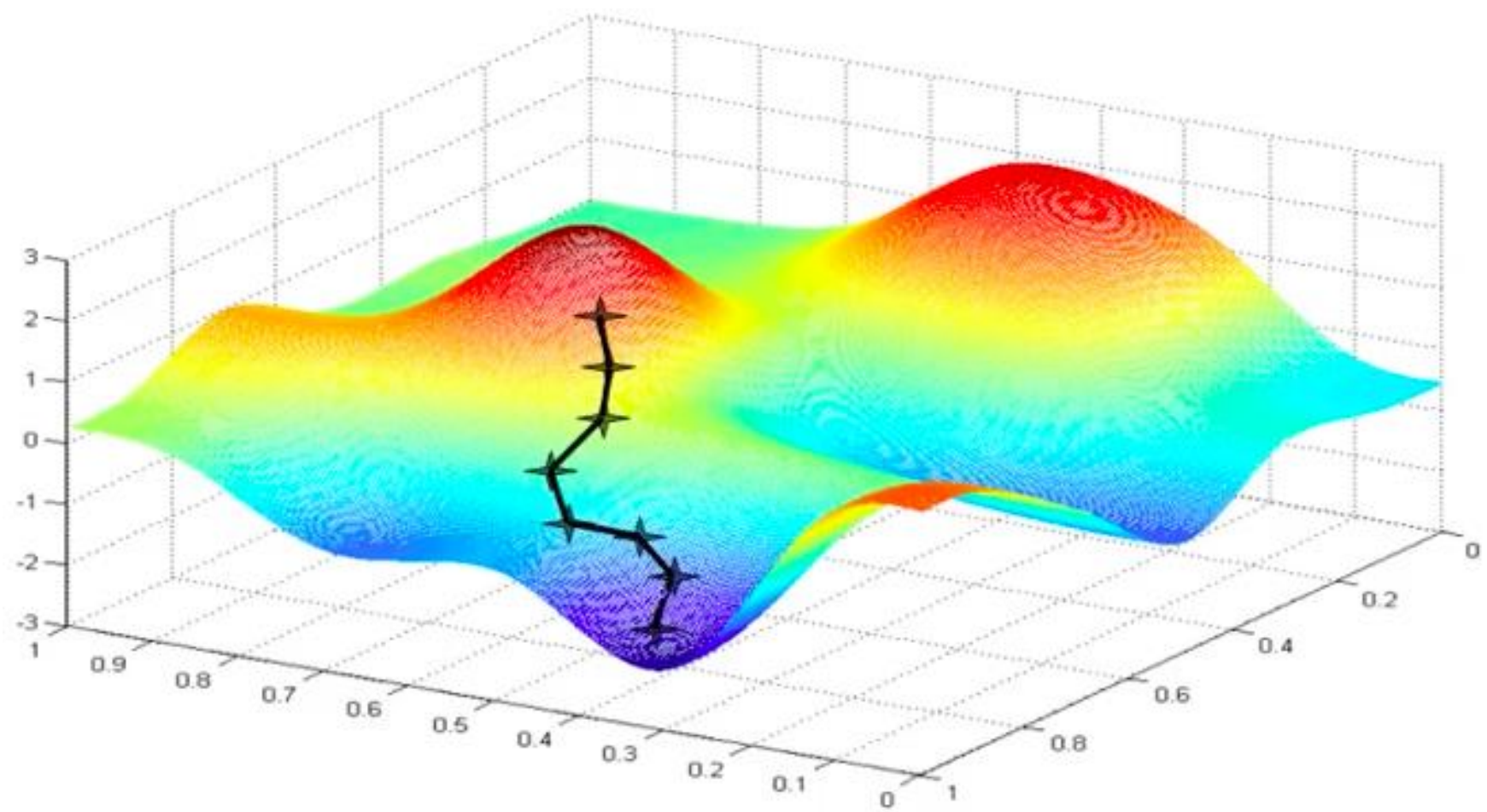
$$\nabla_{\mathbf{w}} L|_{\mathbf{w}=\mathbf{w}^{opt}} = \mathbf{0} \leftrightarrow \begin{cases} \left. \frac{\delta L}{\delta w_1} \right|_{w_1=w_1^{opt}} = 0 \\ \left. \frac{\delta L}{\delta w_2} \right|_{w_2=w_2^{opt}} = 0 \\ \vdots \end{cases}, \text{ ty gradienten är noll i alla minima/maxima.}$$

- (Vi antar här att funktionen L är kontinuerlig och deriverbar överallt, dvs inga “hopp” eller “kanter” i funktionsytan).

Grundläggande optimering

- Gradient Descent:
 - Hur hittar man ner från ett berg om man är blind och har minnesförlust? Följ lutningen.
 - Ett iterativt alternativ för hitta en lösning utan att behöva lösa ekvationssystem.
 - För komplicerade kostnadsfunktioner som beror olinjärt på miljontals parametrar kan det helt saknas metoder för att lösa $\nabla_w L = \mathbf{0}$ analytiskt.
 - Att hitta gradienten i en given punkt är dock alltid en (jämförelsevis) enkel operation.
- Algoritmen:
 1. Initialisera w_0 (exempelvis slumpmässigt).
 2. Uppdatera $w_{t+1} \leftarrow w_t - \alpha \cdot \nabla_w L(x, w)|_{w=w_t}$, (där α är det vi kallar learning rate eller step size).
 3. Repetera steg 2 ovan n gånger (n stort). w_n är (approximativt) lösningen.
- Metoden är matematiskt garanterad hamna godtyckligt nära ett lokalt minimum efter ett ändligt antal n steg givet tillräckligt litet α !





Grundläggande optimering

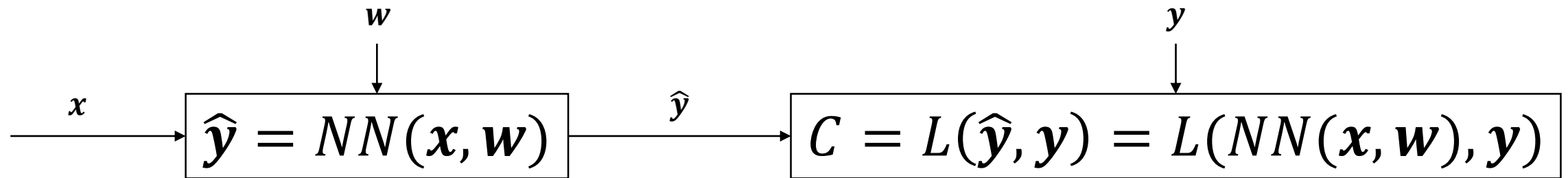
- Stochastic Gradient Descent, SGD:
 - Problem: Kostnadsfunktionen $L(\mathbf{x}, \mathbf{w})$ kanske inte är "känd" i bemärkelsen att den beror på input \mathbf{x} som inte är konstant, utan dragen från en (okänd) slumpfördelning.
 - Lösning: För $I \geq 1$ observationer av \mathbf{x} (skrivet \mathbf{x}^{obs}), beräkna medelvärdet av kostnaden:

$$\hat{L}(\mathbf{x}^{obs}, \mathbf{w}) = \frac{\sum_{i=1}^I L(\mathbf{x}_i^{obs}, \mathbf{w})}{I} \approx L(\mathbf{x}, \mathbf{w})$$

- Nu kan vi använda ungefär samma metod som innan för att minimera \hat{L} :
 1. Initialisera \mathbf{w}_0
 2. Generera observationer \mathbf{x}^{obs}
 3. Uppdatera $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \alpha \cdot \nabla_{\mathbf{w}} \hat{L}(\mathbf{x}^{obs}, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_t}$
 4. Repetera steg 2-3 ovan.
- Detta minimerar även L över stora mängder observerade \mathbf{x} .

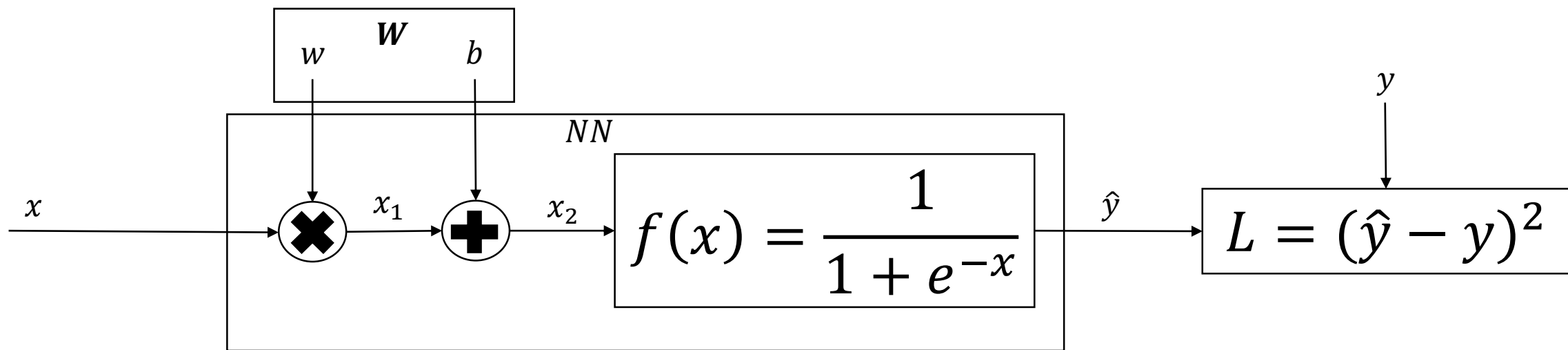
Neuronnät

- (nästan) alla modeller har följande struktur:



- Ett neuronnät $NN(\cdot)$ är en *universal function approximator* som kan lära sig efterlikna någon funktion $f: \mathbf{x} \rightarrow \mathbf{y}$. Exempelvis kan \mathbf{x} vara en bild, $f(\mathbf{x}) = \mathbf{y}$ är 1 om bilden föreställer en katt, 0 annars. f antas *a priori* existera och vara beräkningsbar.
- Vi vill minimera kostnaden C , som mäter hur lika $\hat{\mathbf{y}}$ och \mathbf{y} är, med avseende på nätverkets justerbara parametrar \mathbf{w} . När C är minimerad lär $NN(\cdot, \mathbf{w})$ "likna" f .
- Vi kan tillämpa SGD, där "observationerna" utgörs av exempel på par (\mathbf{x}, \mathbf{y}) . Förhoppningen är att för någon konfiguration av \mathbf{w} kommer nätverket att beräkna den funktion som (\mathbf{x}, \mathbf{y}) är exempel på input resp. output från.

Exempel:



- Vi använder SGD för att uppdatera w och b :

1. Initialisera $\mathbf{w}_0 = (w_0, b_0)$
2. Generera ett exempel x med motsvarande y
3. Uppdatera $(w_{t+1}, b_{t+1}) \leftarrow (w_t, b_t) - \alpha \cdot \nabla_{w,b} (\hat{y} - y)^2 \big|_{w=w_t, b=b_t}$
4. Repetera steg 2-3 ovan.

- För att kunna beräkna $\nabla_{w,b} L = \left(\frac{\delta L}{\delta w}, \frac{\delta L}{\delta b} \right)$ behöver vi använda kedjeregeln:

$$\bullet \frac{\delta L}{\delta w} = \frac{\delta L}{\delta \hat{y}} \frac{\delta \hat{y}}{\delta x_2} \frac{\delta x_2}{\delta x_1} \frac{\delta x_1}{\delta w} \text{ och } \frac{\delta L}{\delta b} = \frac{\delta L}{\delta \hat{y}} \frac{\delta \hat{y}}{\delta x_2} \frac{\delta x_2}{\delta b}$$

- Detta kallas *backpropagation*