② Artificial "Neuron"

usually {bias feature}



$$\hat{y} = f(a) = f(\vec{w}^T\vec{x})$$

Activation / Response Function
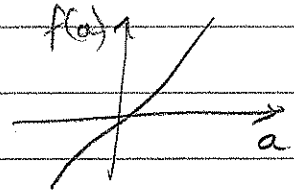
Inputs — Strengths of connections — Activation

$$a = \sum_{j=0}^{d} w_j x_j$$

$$= \vec{w}^T \vec{x}$$

Many different activation functions

→ Linear:  $f(a) = a$

  ⟹ $\hat{y} = \vec{w}^T\vec{x}$  [Linear Regression]



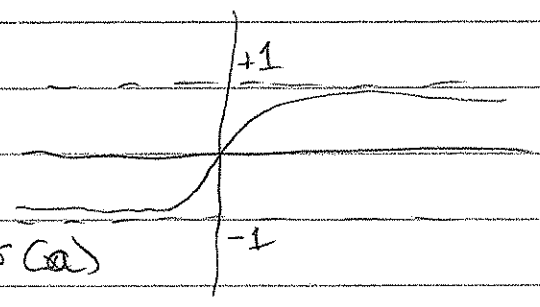→ Logistic  $f(a) = \dfrac{1}{1+e^{-a}} = \sigma(a)$   sigmoid



$$\hat{y} = \frac{1}{1+e^{w^Tx}}$$   | Logistic Regression
$P(Y=1 \mid \vec{x}, \vec{w})$
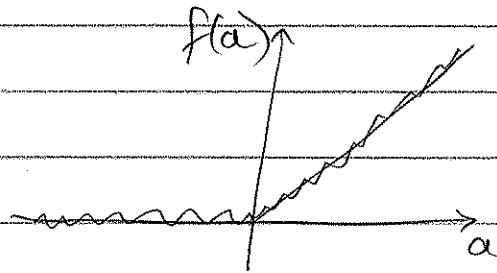
→ Tanh

$$f(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$



For hidden units in NN
we always prefer tanh over $\sigma(a)$
why?

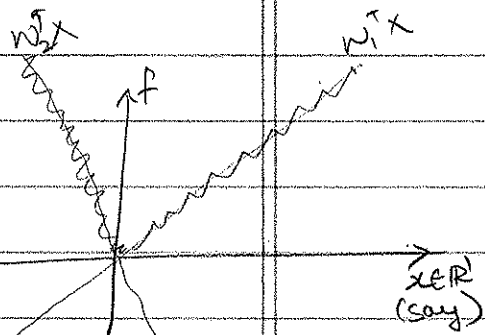$\rightarrow$ ReLU [Rectified Linear Unit]

$$f(a) = \max \{0, a\}$$

In hidden layer of deep NN, this is always preferred over $\sigma(a)$ or $\tanh(a)$. Why??

$\rightarrow$ Maxout

$$f(a) = \max \left\{ \overbrace{\vec{W_1}^T x}^{a_1}, \overbrace{\vec{W_2}^T x}^{a_2} \right\}$$

Each neuron has (say) 2 weights $\vec{W_1}, \vec{W_2}$

Take max activation.

ReLU is a special case of this. How?

③ Loss Functions

→ functions of both parameters $\vec{w}$ & training data

① Log-Loss / Cross-Entropy / Maximum-Likelihood/
KL-Divergence

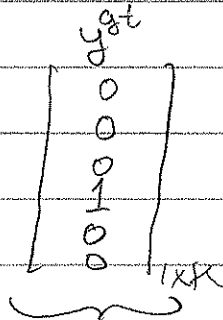$$L(\vec{w}; D) = \sum_{i=1}^{N} L_i(\vec{w}) \quad \} \text{Decomposable Loss}$$

where $L_i(\vec{w}) = -\log p(y_i^{gt} \mid \vec{x}_i, \vec{w})$

How much prob does your
model assign to GT labels?

$\equiv$ negative log-likelihood for
this sample

→ Why is this called Cross-Entropy? And where is
the KL divergence coming in?

Consider Multiclass-classification w/ 1-HOT encoding

$y^{gt}$      model prediction

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}_{1 \times K} \qquad \begin{bmatrix} p(y=1 \mid \vec{x}, \vec{w}) \\ \vdots \\ \vdots \\ p(y=K \mid \vec{x}, \vec{w}) \end{bmatrix}$$

$\underbrace{\quad\quad}_{p^{gt}(y)}$        $\underbrace{\quad\quad}_{\hat{p}(y)}$

[delta distribution]    [Model distribution]

$$KL(p^{gt} \| \hat{p}) = -\sum_{y=1}^{K} p^{gt}(y) \log \hat{p}(y)$$
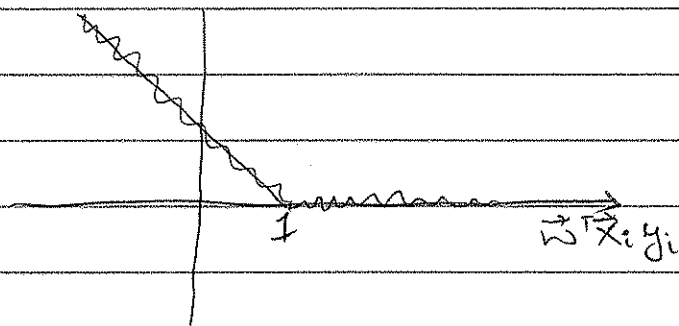
$$= -\log p(y = y_i^{gt} \mid \vec{x}_i, \vec{w})$$

② Hinge-Loss  [for binary-classification]

$$L_i(\vec{w}) = \max\{0, 1 - \vec{w}^T\vec{x}_i y_i\} \qquad \text{where } y_i \in \{+1, -1\}$$

④ Detour: Matrix / Vector differentiation

|   | S | V | M |
|---|---|---|---|
| S | $\frac{\partial y}{\partial x}$ | $\frac{\partial y}{\partial \vec{x}}$ | $\frac{\partial y}{\partial X}$ |
| V | $\frac{\partial \vec{y}}{\partial x}$ | $\frac{\partial \vec{y}}{\partial \vec{x}}$ | |
| M | $\frac{\partial Y}{\partial x}$ | Tensor | |

$$x, y \in \mathbb{R}^1$$
$$\vec{x} \in \mathbb{R}^d$$
$$\vec{y} \in \mathbb{R}^k$$

Convention: $\quad \dfrac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_k}{\partial x} \end{bmatrix}$ $\quad \downarrow$ numerator = dim 1
= col-vector

[Gradient] $\quad \dfrac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_d} \end{bmatrix}$ denominator = dim 2
= row-vector

[Jacobian Matrix] $\quad \dfrac{\partial \vec{y}}{\partial \vec{x}} = i \begin{bmatrix} & & j & \\ & & \downarrow & \\ & \frac{\partial y_i}{\partial x_j} & & \\ & & & \end{bmatrix}$ 

$k \times d$

Easy to prove: $\to \dfrac{\partial(\vec{w}^T\vec{x})}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial(\vec{w}^T\vec{x})}{\partial w_1} & \cdots & \frac{\partial(\vec{w}^T\vec{x})}{\partial w_d} \end{bmatrix}$

$$= x^T$$

$\to \dfrac{\partial(\vec{w}^T A \vec{w})}{\partial \vec{w}} = 2w^T A$

$\to \quad \vec{y} = A\vec{x} \qquad \dfrac{\partial \vec{y}}{\partial \vec{x}} = A$

## ⑤ Chain Rule

→ Function Composition: $L(x) = (f \circ g)(x)$
$$= f(g(x))$$

Chain Rule:

→ [Most General Notation]
$$D_x (f \circ g) = D_{g(x)} f \circ D_x g$$
$\underbrace{\qquad}_{\text{total derivative}}$

→ [More concrete notation for scalars]
$$L'(x) = f'(g(x)) \, g'(x)$$

→ [With intermediate variables]
$$y = g(x)$$
$$z = f(y)$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Example: $\quad L_i(w) = -\log\left(\frac{1}{1+e^{-w^T x_i}}\right) \quad [\text{For } y_i = +1]$

$$= \left( \underbrace{-\log(\cdot)}_{\frac{\partial L}{\partial p}} \circ \underbrace{\frac{1}{1+e^{-(\cdot)}}}_{\frac{\partial p}{\partial a}} \circ \underbrace{x^T(\cdot)}_{\frac{\partial a}{\partial w}} \right)(w)$$

$$\frac{\partial L_i}{\partial w} = \left[\frac{1}{P}\right] \cdot \underbrace{\left[\frac{-1}{(1+e^a)^2} \cdot -e^{-a}\right]}_{P \cdot (1-P)} \cdot x^T \qquad = (1-P)x^T$$

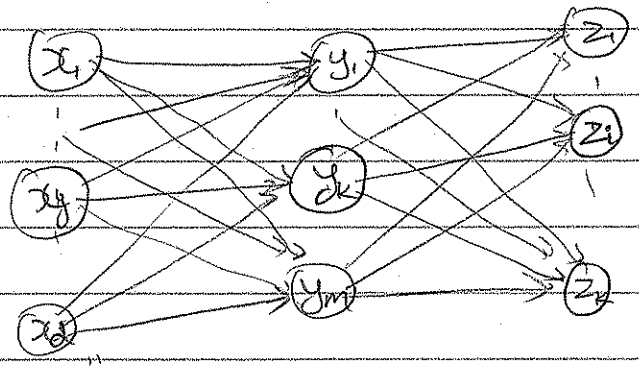→ Multivariate Chain Rule

$$g: \mathbb{R}^d \to \mathbb{R}^m$$
$$f: \mathbb{R}^m \to \mathbb{R}^k$$

$$L(\vec{x}) = (f \circ g)(\vec{x})$$

$$\vec{y} = g(\vec{x}) \qquad \vec{z} = f(\vec{y})$$

→ Chain Rule: $\quad D_{\vec{x}}(f \circ g) = D_{g(x)} f \circ D_{\vec{x}}(g)$

[Abstract form holds]

But what does this mean??

Visualize:



$$\underbrace{\frac{\partial z_i}{\partial x_j}}_{} = \boxed{\sum_k} \cdot \boxed{\frac{\partial z_i}{\partial y_k}} \cdot \boxed{\frac{\partial y_k}{\partial x_j}}$$

"outcome"

"knob"

all intermediate variables

how those intermediate vars affect outcome!

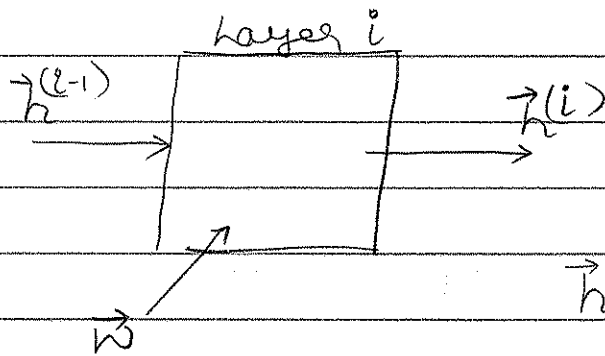how my "knob" affects intermediate variable

Formally,

Jacobian relationship

$$J_{f \circ g} = (J_f \circ g) \; J_g$$

$$i \left[ -- \frac{\partial z_i}{\partial x_j} \quad \right]_{k \times d} = i \left[ -- \frac{\partial z_i}{\partial y_k} \quad \right]_{k \times m} k \left[ -- \frac{\partial y_k}{\partial x_j} \quad \right]_{m \times d}$$

→ What if my $\vec{x}, \vec{y}, \vec{z}$ are tensors?
   → string up into vectors & proceed
   → Matlab notation $x\text{-vec} = X(:);$
   → Trust me, this is the cleanest way

→ In Neural Nets $\vec{z} \in \mathbb{R}^1$ (Loss) $L(\vec{w})$

Layer i

$\vec{h}^{(i-1)} \longrightarrow \boxed{\phantom{xxxx}} \longrightarrow \vec{h}^{(i)}$

$\vec{w}$

$$\vec{h}^{(i)} = g(\vec{h}^{(i-1)}, \vec{w})$$

$$L(\vec{w}) = f(h^{(i)})$$

$$\frac{\partial L}{\partial \vec{w}} = \frac{\partial L}{\partial \vec{h}^{(i)}} \frac{\partial \vec{h}^{(i)}}{\partial \vec{w}} \qquad \frac{\partial L}{\partial \vec{w}} = \left\langle \frac{\partial L}{\partial out}, \frac{\partial out}{\partial \vec{w}} \right\rangle$$

$$\left[ \phantom{xxxx} \right]_{1 \times d} \left[ \phantom{xxxx} \right]_{|h^i| \times |h^i|} \left[ \phantom{xxxxxxxx} \right]_{|h^i| \times d}$$