

9/1/15

①

DEEP LEARNING : BACKPROP

① Basic Setup / Notation

→ (Row) Input

$\vec{x} \in \mathbb{R}^d$ (continuous input)

$\in \{0, 1\}$ (discrete input)

} sometimes both

→ Scalars: a, b, c, i, j, k, x, y

(small choice)

Vectors: $\vec{x}, \vec{y}, \vec{w}, \vec{h}$

(vec notation)

Matrices: A, B, W, X, H

(capital choice)

(or sets or Random Variables)

d : #dimension

N : #inputs

→ Output / Target / Labels

y or $\hat{y} \in \mathbb{R}^k$ (k -dim regression)

$\in \{0, 1\}$ (Binary classification)

$\in \{1, 2, \dots, k\}$ (Multi-class classification)

$\alpha \in \{0, 1\}^k$

s.t. $\sum_{c=1}^k y_c = 1$

} 1-hot encoding for multiclass classification

→ Unknown

Target Function

ft: $X \rightarrow Y$

↑
input space

↑
Output Space

→ Given: Dataset
or "Training Data"

$$D = \{ (\vec{x}_i, y_i), \dots, (\vec{x}_N, y_N) \}$$

sample or data-point

→ Goal: Given D (input outputs) from unknown f ,
predict $f^{\text{gt}}(\vec{x})$ on NEW \vec{x}

→ Approach: [All of Supervised Learning]

① Pick a model class / Hypothesis Space

$$H = \{ g : X \rightarrow Y \}$$

a set of mappings from $X \rightarrow Y$

[A lot of ML is about what kind of H "makes sense"
for one task]

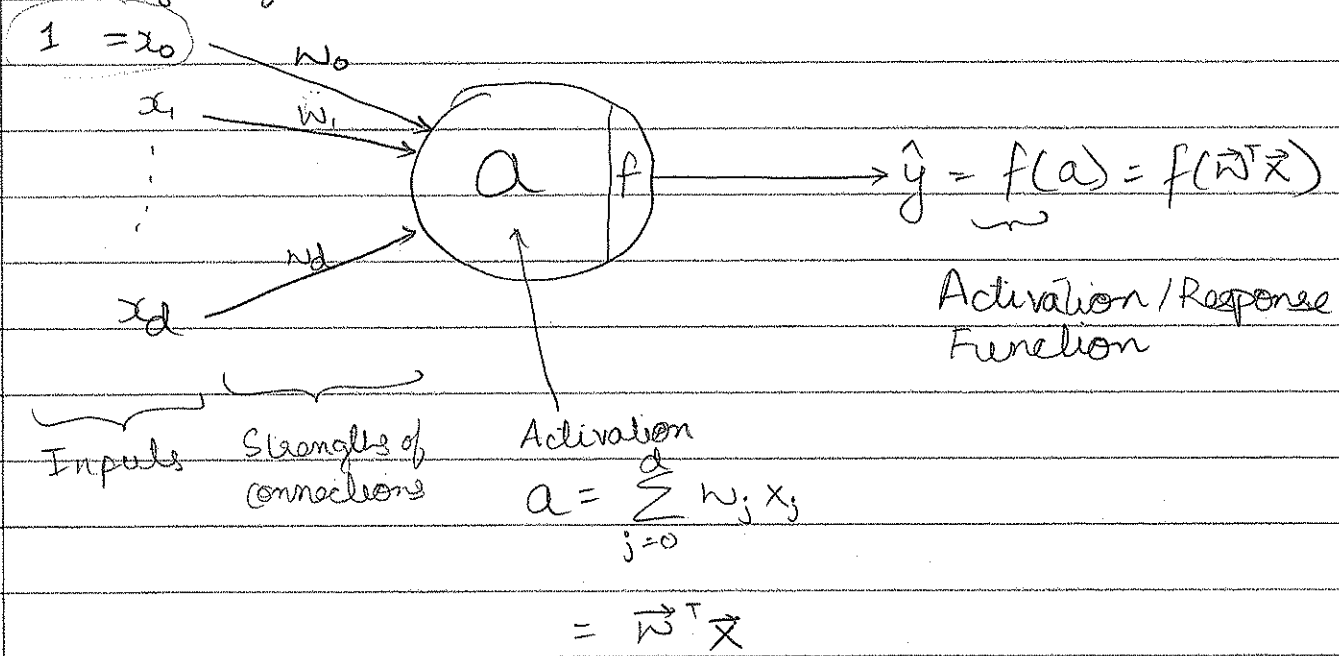
② Search / Optimize a loss function
over H to find "best" $g \in H$

$$\Rightarrow \text{Learning} = \underset{g \in H}{\text{argmin}} \text{Loss}(g; D)$$

usually $\sum_{i=1}^N L_i(y_i^{\text{gt}}, g(\vec{x}_i))$

② Artificial "Neuron"

usually [bias feature]

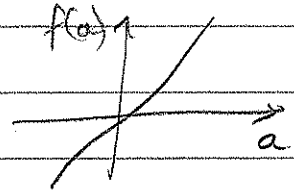


Many different activation functions

→ Linear:

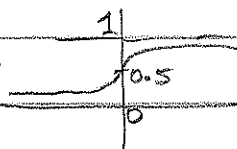
$$f(a) = a$$

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



→ Logistic

$$f(a) = \frac{1}{1 + e^{-a}} = \overset{\text{sigmoid}}{\sigma(a)}$$

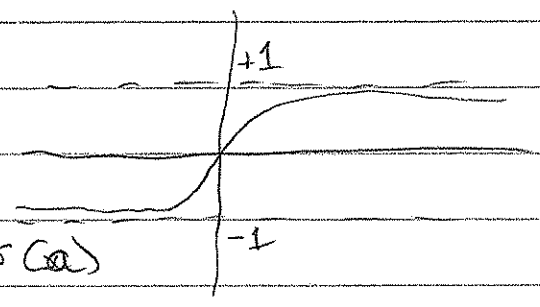


$$\hat{y} = \frac{1}{1 + e^{-\vec{w}^T \vec{x}}}$$

Logistic Regression
 $P(Y=1 | \vec{x}, \vec{w})$

→ Tanh

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

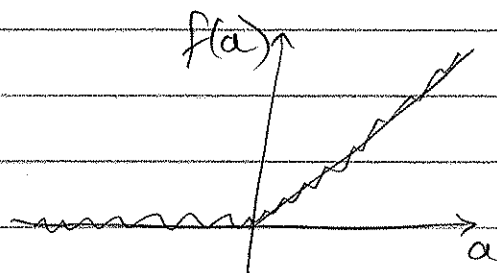


For hidden units in NN

we always prefer tanh over $\sigma(a)$
why?

→ ReLU [Rectified Linear Unit]

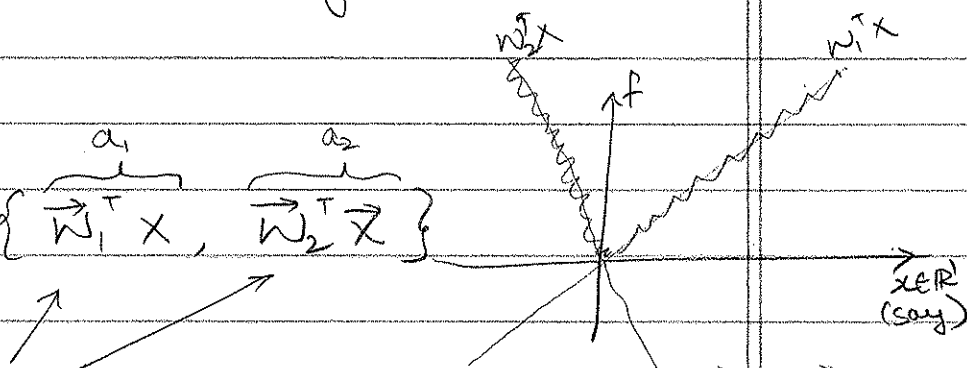
$$f(a) = \max\{0, a\}$$



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

→ 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} to training data

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

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

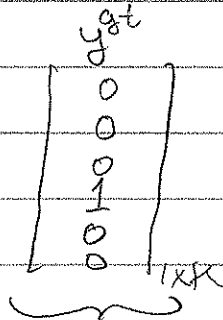
where $L_i(\vec{w}) = -\log P(y_i^{gt} | \vec{x}_i, \vec{w})$

How much prob does your model assign to GT labels?

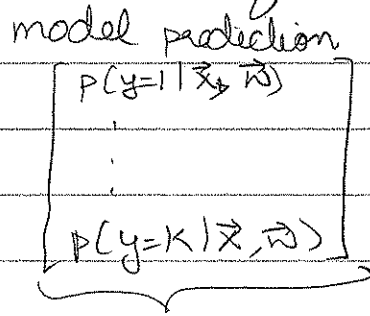
≡ negative log-likelihood for this sample

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

Consider Multiclass-classification w/ 1-HOT encoding



$p^{gt}(y)$
[delta distribution]



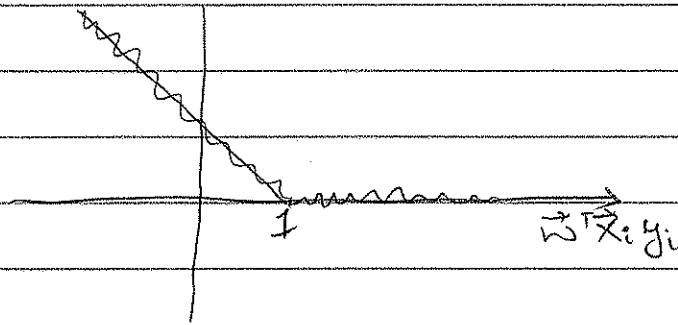
$\hat{P}(y)$
[Model distribution]

$$KL(p^{gt} \parallel \hat{p}) = -\sum_{y=1}^K p^{gt}(y) \log \hat{p}(y)$$

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

② Hinge-Loss [for binary-classification]

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



④ Detour: Matrix/Vector differentiation

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

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

Tens^{0,1,2}

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

[Gradient] $\frac{\partial y}{\partial \vec{x}} = \left[\frac{\partial y}{\partial x_1} \dots \frac{\partial y}{\partial x_d} \right]$ denominator = dim 2
= row-vector

[Jacobian Matrix] $\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \vdots \\ \frac{\partial y_i}{\partial x_j} \\ \vdots \end{bmatrix}_{k \times d}$

Easy to prove: $\rightarrow \frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{x}} = \left[\frac{\partial (\vec{w}^T \vec{x})}{\partial w_1} \dots \frac{\partial (\vec{w}^T \vec{x})}{\partial w_d} \right] = \vec{x}^T$

$\rightarrow \frac{\partial (\vec{w}^T A \vec{x})}{\partial \vec{x}} = \vec{w}^T A$

$\rightarrow \vec{y} = A \vec{x} \quad \frac{\partial \vec{y}}{\partial \vec{x}} = A$

⑤ Chain Rule

some books/people use opposite order

$$\rightarrow \text{Function Composition: } L(x) = (f \circ g)(x) \\ = f(g(x))$$

Chain Rule:

\rightarrow [Most General Notation]

$$D_x (f \circ g) = \underbrace{D_{g(x)} f}_{\text{total derivative}} \circ D_x g$$

\rightarrow [More concrete notation for scalars]

$$L'(x) = f'(g(x)) g'(x)$$

\rightarrow [With intermediate variables]

$$y = g(x)$$

$$z = f(y)$$

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

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

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

$$\frac{\partial L_i}{\partial w} = \begin{bmatrix} 1 \\ p \end{bmatrix} \cdot \underbrace{\begin{bmatrix} -1 & \dots & -e^{-a} \\ (1+e^a)^2 & \dots & \end{bmatrix}}_{p \cdot (1-p)} \cdot x^T = (1-p)x^T$$

→ Multivariate Chain Rule

$$g: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

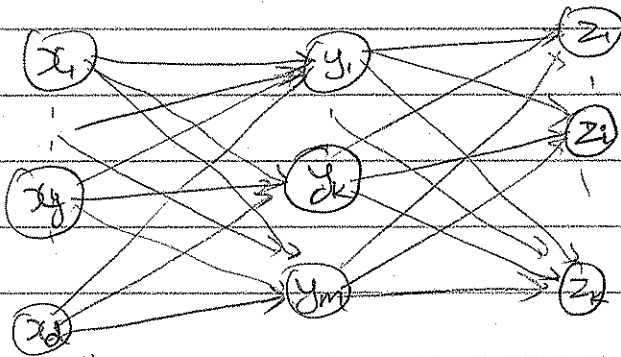
$$f: \mathbb{R}^m \rightarrow \mathbb{R}^k$$

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

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

→ Chain Rule: $D_{\vec{x}}(f \circ g) = D_{\vec{y}} f \circ D_{\vec{x}} g$
 [Abstract form holds]
 But what does this mean??

Visualize:



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

"outcome" $\frac{\partial z_i}{\partial x_j}$
 "knob" $\frac{\partial y_k}{\partial x_j}$
 all intermediate variables $\frac{\partial z_i}{\partial y_k}$
 how these intermediate vars affect outcome!
 how my "knob" affects intermediate variable

