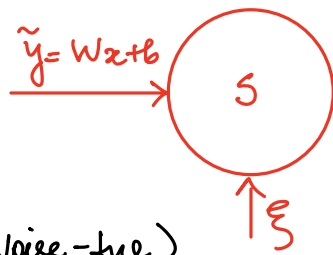


"TriDeNT"



Consider input: (Noise-free)

$$\tilde{y} = wx + b$$

Subject to (gaussian) noise

e.g. $\xi(t) = N(0; \sigma)$
 $y = \tilde{y} + \xi$

Neuron is in state S given by:

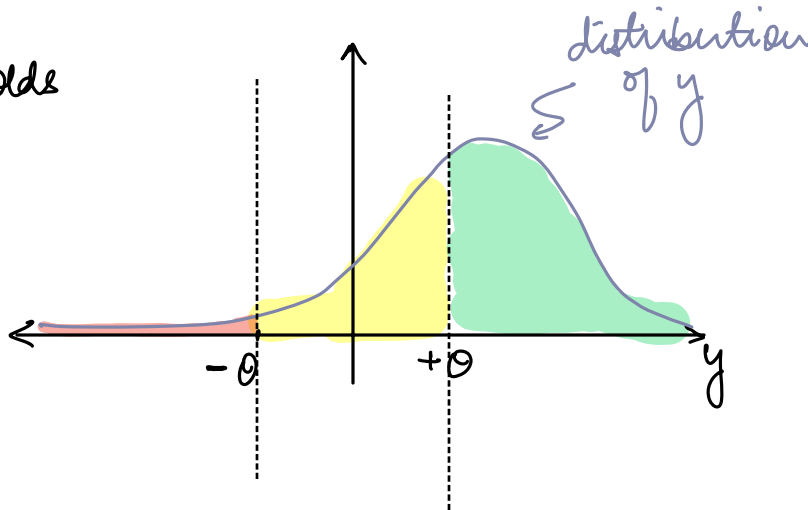
$$S = h(y)$$

The goal is to find "Expected activation" of the neuron.

Let the neuron take states $S \in \{s_1, s_2, \dots, s_k\}$

The neuron assumes a state $S = s_i$ as a function of the noisy input arriving at it y . Let that function be $h(y)$

θ 's: thresholds



for a $k=3$:

$$S = \{-1, 0, +1\}$$

we can define thresholds as $\pm\theta$.

as an example: say h is the activation:

$$h(y) = \begin{cases} -1 & ; -\infty < y < -\theta \\ 0 & ; -\theta < y < \theta \\ 1 & ; +\theta < y < +\infty \end{cases}$$

$\mathbb{E}[h(y)] \doteq$ Expected value of the state of the neuron.

In general, if the **controllable** inputs are $\tilde{y} = Wx + b$;

$P(S=s_i | \tilde{y}) \rightarrow$ probability of finding a neuron in state s_i
given **NOISE FREE** input \tilde{y} .

Let $\theta \in \{\theta_0, \theta_1, \dots, \theta_{k-1}, \theta_k\}$ be the thresholds. $\theta_0 = -\infty$; $\theta_k = +\infty$

$$P(S=s_i | \tilde{y}) = P(y \in [\theta_{i-1}, \theta_i])$$

$$= \int_{\theta_{i-1}}^{\theta_i} P(y | \tilde{y}) dy \quad (\text{due to noise, } y \text{ can overshoot or undershoot a threshold})$$

Define $\int_{-\infty}^{\theta} P(y | \tilde{y}) dy = f(\theta)$ (CDF of $P(y | \tilde{y})$)

$$\therefore P(S=s_i | \tilde{y}) = f(\theta_i) - f(\theta_{i-1})$$

$$\therefore E[(s | \tilde{y})] = \sum_{i=1}^k s_i P(S=s_i | \tilde{y})$$

$$E[(s | \tilde{y})] = \sum_{i=1}^k s_i (f(\theta_i) - f(\theta_{i-1}))$$

if $s \in \{-1, 0, 1\}$; $\theta_1 = -\theta$, $\theta_2 = +\theta$

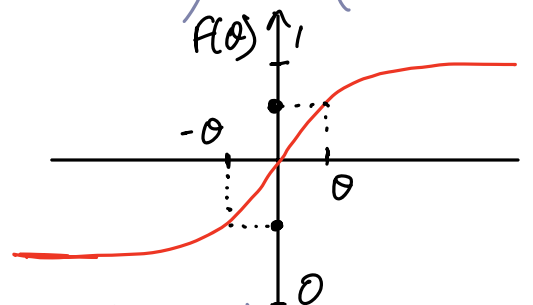
$$E[(s | \tilde{y})] = -1 \cdot (f(-\theta) - f(-\infty)) + 0 \cdot (f(\theta) - f(-\theta)) + 1 \cdot (f(\infty) - f(\theta))$$

$$= -f(-\theta) + \underbrace{f(\infty)}_1 - f(\theta)$$

$$= 1 - (f(\theta) + f(-\theta))$$

if $f(\cdot)$ is sigmoidal function (erf, tanh, ...)

$$f(-\theta) = 1 - f(\theta)$$



$$\therefore \mathbb{E}[(s|\tilde{y})] = 1 - (f(0) + 1 - f(0))$$

$$\boxed{\mathbb{E}[(s|\tilde{y})] = 0}$$

The expected state is always zero (?)

- Bringing in the effect of input:

$$\mathbb{E}[(s|\tilde{y})] = \mathbb{E}[(h(y)|\tilde{y})]$$

and if $h(\cdot)$ is simply thresholding function; under additive noise θ_i

$$h(y) \rightarrow \int_{\theta_{i-1}}^{\theta_i} P(y|\tilde{y}) dy; \quad \forall i=1, \dots, k$$

$$\therefore \frac{d}{dy} h(y) = \frac{d}{dy} \int_{\theta_{i-1}}^{\theta_i} P(y|\tilde{y}) dy$$

Fundamental theorem of calculus:

$$\boxed{\frac{d}{dy} h(y) = P(y=\theta_i|\tilde{y}) - P(y=\theta_{i-1}|\tilde{y})}$$

Now consider change in the expected state w.r.t. input

$$\begin{aligned} \frac{d}{dy} \mathbb{E}[(s|\tilde{y})] &= \frac{d}{dy} \mathbb{E}[(\underbrace{h(y)}_{s_i}|\tilde{y})] \\ &= \frac{d}{dy} \sum_{i=1}^k s_i P(y \in [\theta_{i-1}, \theta_i]|\tilde{y}) \end{aligned}$$

$$= \sum_{i=1}^k s_i \frac{d}{dy} \int_{\theta_{i-1}}^{\theta_i} p(y|\tilde{y}) dy$$

$$\frac{d}{dy} E[s|\tilde{y}] = \sum_{i=1}^k s_i [p(y=\theta_i|\tilde{y}) - p(y=\theta_{i-1}|\tilde{y})]$$

$$= h(y) \cdot \underbrace{[P(y=\theta_i|\tilde{y}) - P(y=\theta_{i-1}|\tilde{y})]}_{\frac{d}{d\tilde{y}} h(y)} \cdot \int_{\theta_{i-1}}^{\theta_i} P(y|\tilde{y}) dy$$

$$\frac{d}{dy} E[(s|\tilde{y})] = h(y) \cdot [P(y=\theta_i|\tilde{y}) - P(y=\theta_{i-1}|\tilde{y})] \cdot [f(\theta_i) - f(\theta_{i-1})]$$

- New credit assignment can be done with chain rule:

$$\frac{d}{dw} E[(s|\tilde{y})] = \frac{d}{dy} E[(s|\tilde{y})] \cdot \frac{dy}{dw}$$

$$\therefore y = \underbrace{Wx + b}_{\tilde{y}} + \xi$$

$$\therefore \frac{dy}{dw} = x \left(+ \frac{d\xi}{dw} \right)$$

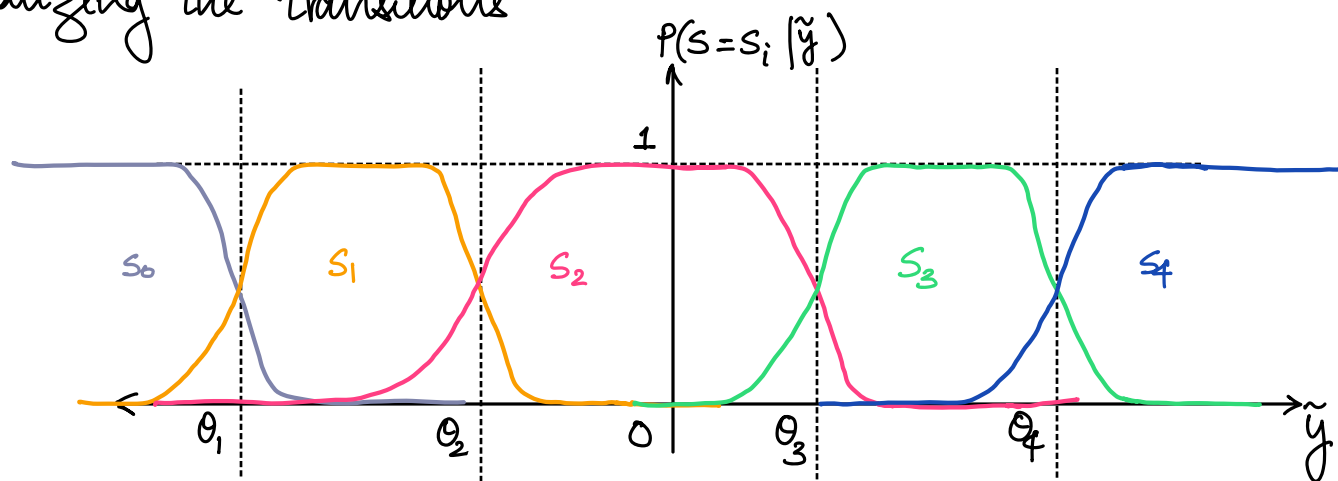
$\therefore \xi$ does not depend on w

$$\therefore \frac{d}{dw} E[(s|\tilde{y})] = \frac{d}{dy} E[(s|\tilde{y})] \cdot x \quad (\text{single weight case})$$

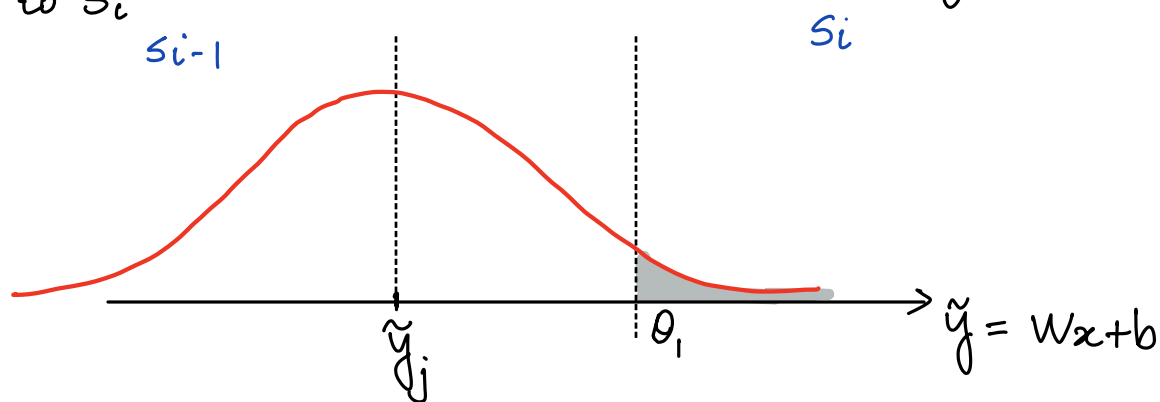
→ This would turn into outer product update if W is matrix & x & y are vectors

$$\Delta W = -\underset{\substack{\uparrow \\ \text{learning} \\ \text{rate}}}{\eta} \frac{d}{dw} E[(s|\tilde{y})] \otimes \vec{x}$$

- Visualizing the transitions



Consider a threshold θ_i around which the model goes from state s_{i-1} to s_i

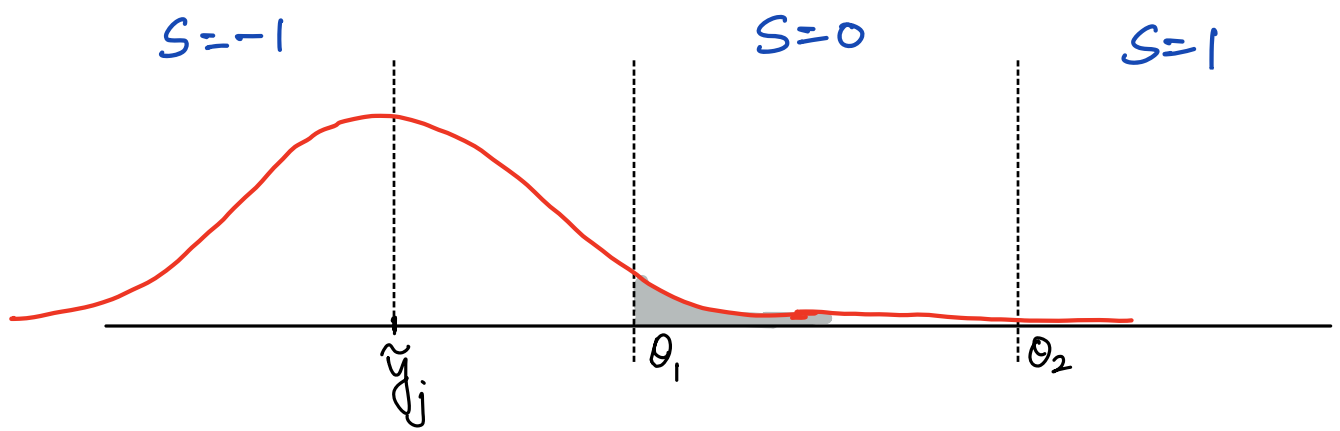


Let \tilde{y}_j be the input for j^{th} sample datum. As the decision variable $y_j = \tilde{y}_j + \xi_j$, the probability of the model switching to state s_i is:

$$P(y > \theta_1 | \tilde{y}) = \frac{1}{\sqrt{2\pi} \sigma} \int_{\theta_1}^{\infty} e^{-\frac{(u - y_j)^2}{2\sigma^2}} du$$

$$\text{and } P(y < \theta_1 | \tilde{y}) = \frac{1}{\sqrt{2\pi} \sigma} \int_{-\infty}^{\theta_1} e^{-\frac{(u - y_j)^2}{2\sigma^2}} du.$$

for $s \in \{-1, 0, +1\}$



$$\mathbb{E}[s|\tilde{y}_j] = (-1) \cdot P(y < \theta_1 | \tilde{y}_j) + 0 \cdot (P(\theta_1 < y < \theta_2 | \tilde{y}_j)) + (+1) \cdot P(y > \theta_2 | \tilde{y}_j)$$

$$= (-1) \cdot \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\theta_1} e^{-\frac{(u-\tilde{y}_j)^2}{2\sigma^2}} du + \frac{1}{\sqrt{2\pi}\sigma} \int_{\theta_2}^{\infty} e^{-\frac{(v-\tilde{y}_j)^2}{2\sigma^2}} dv$$

$$= \underbrace{\frac{1}{\sqrt{2\pi}\sigma}}_c \left[\int_{\theta_2}^{\infty} e^{-\frac{(v-\tilde{y}_j)^2}{2\sigma^2}} dv - \int_{-\infty}^{\theta_1} e^{-\frac{(u-\tilde{y}_j)^2}{2\sigma^2}} du \right]$$

$$\frac{d}{d\tilde{y}_j} \mathbb{E}[s|\tilde{y}_j] = c \cdot \left[- \left\{ e^{-\infty} - e^{-\frac{(\theta_2-\tilde{y}_j)^2}{2\sigma^2}} \right\} - (-1) \left\{ e^{-\frac{(\theta_1-\tilde{y}_j)^2}{2\sigma^2}} - e^{\infty} \right\} \right]$$

$$\frac{d}{d\tilde{y}_j} \mathbb{E}[s|\tilde{y}_j] = c \cdot \left[\mathcal{N}(\theta_2; \tilde{y}_j, \sigma) + \mathcal{N}(\theta_1; \tilde{y}_j, \sigma) \right]$$

- Discretizing the gradients

→ Sigmoid everywhere, computing der. → \mathbb{R}

→ Experiments: 4 figures

1. Net architecture, overview

↓
- 3 panel idea

- Net topology

- Backprop.

→ 3 result figures.

→ How does this scale?

↳ Not sure how to show that.

label $\begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline \end{array}$

$z = \begin{array}{|c|c|c|} \hline .01 & .9 & .03 \\ \hline \end{array}$

$$z[p] = .9$$

$$\sum_0 \max(0, .01 - .9) + \max(\quad)$$

$$\frac{d}{dy} \underbrace{E[h(y) | \tilde{y}]}_{\downarrow} = \underbrace{\frac{dE}{dh}}_{\underbrace{\quad}_{\checkmark}} \cdot \underbrace{\frac{dh}{dy}}_{\checkmark}$$

$$= \underbrace{\sum_{i=1}^k h_i(y) P(h_i(y) | \tilde{y})}_{\checkmark}$$

$$\frac{dE}{dh_j} \cdot \frac{dh}{dy} = \frac{d}{dh_j} \sum_{i=1}^k h_i(y) P(h(y)=s_i | \tilde{y}) \cdot \frac{dh_i}{dy}$$

$$= \underbrace{P(h_j(y) | \tilde{y})}_{s_j} \cdot \frac{dh}{dy}$$

$$y \in \{\theta_{j-1}, \theta_j\}$$

C.R.V

$$E[h(y) | \tilde{y}] = \int_{-\infty}^{\infty} h(y) P(h(y) | \tilde{y}) dh$$

$$\begin{aligned}
 \frac{\partial E}{\partial h} &= \int \frac{\partial}{\partial h} (h(y) P(h(y) | \tilde{y})) dh \\
 &= \int [P(h(y) | \tilde{y}) + h(y) \frac{\partial}{\partial h} P(h(y) | \tilde{y})] dh \\
 &= 1 + h(y) \frac{\partial}{\partial h} P(h(y) | \tilde{y})
 \end{aligned}$$

D.R.V

$$E[P(h(y) | \tilde{y})] = \sum_{i=1}^K h_i(y) P(h = h_i(y) | \tilde{y})$$

$$\frac{\partial E}{\partial h_j} = \frac{\partial}{\partial h_j} \sum_i h_i(y) P(h = h_i | \tilde{y}) = P(h = h_j | \tilde{y}) + \underbrace{h_j(y) \frac{\partial}{\partial h_j} P(h = h_j | \tilde{y})}_{\text{what's happening to this term?}}$$

what's happening to this term?

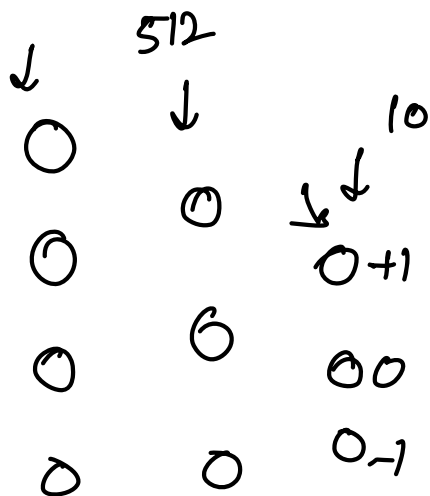
$$E[(s | \tilde{y})] = \sum_i s_i P(y | \tilde{y})$$

$$\frac{d}{d\tilde{y}} \int e^{-\frac{(u - \tilde{y})^2}{2\sigma^2}} du$$

$$\begin{aligned}
&= \int \frac{d}{d\tilde{y}} e^{-\frac{(u-\tilde{y})^2}{2\sigma^2}} du \\
&= \int e^{-\frac{(u-\tilde{y})^2}{2\sigma^2}} \cdot \frac{1}{2\sigma^2} \cdot (-2)(u-\tilde{y}_j) du \\
&= (-1) \int \frac{(u-\tilde{y}_j)}{\sigma^2} e^{-\frac{(u-\tilde{y})^2}{2\sigma^2}} du \\
&= -1 \cdot \int x \cdot e^{-x^2} dx
\end{aligned}$$

$\xrightarrow{-x^2}$
 $\xrightarrow{-x^2}$

784



1. Expected state
2. N/w topology
3. Loss: Cross entropy

$l \rightarrow \text{label}$

$\hat{y} \rightarrow \text{model} \rightarrow \text{softmax}(w_f x_f + b_f)$

$$\hat{\mathcal{L}} = - \sum_{c=1}^{10} l_c \log p_c$$

$$\boxed{0 \mid 0 \mid 1}$$

$$\boxed{.01 \mid .01 \mid 0.9}$$

$$\boxed{.9 \mid .01 \mid .01}$$

$$= 0 \cdot \log 0.01 +$$

$$0 \log 0.01 +$$

$$1 \log \underbrace{(0.9)}_{\approx 1} \approx 0$$

$$\frac{0 \cdot \log_{10} .9}{10} + \frac{0 \cdot \log_{10} .01}{10}$$

$$+ 1 \cdot \log_{10} .01$$

↖ -ve number

$$(-1)(-2) = 2$$

$$\sum_{c=1}^{10} p_c \log \frac{p_c}{l_c + \epsilon}$$

Potential Paper title(s)

1. Deep learning with ternary stochastic neurons
2. Deep stochastic Ternary nns