$$p(y_1, y_2 \cdots y_m) x_1 \cdots x_m) = \prod_{i=1}^{m} P(y_i/x_i)$$

$$\max \log \left( \prod_{i=1}^{m} P(y_i/x_i) \right)$$

$$\cancel{\min -\log (} \quad \min \sum_{i=1}^{m} -p \log P(y_i/x_i)$$

$$y(\theta) = -\log p(y/x)$$
$$\downarrow$$
$$w_1 w_2 \cdots \quad = -\log(\sigma(2y-1)z)$$

$$= \zeta(-(2y-1)z)$$

not activation

softplus

cost func $\longleftarrow$

$$= \zeta((1-2y)z)$$

@ $y = 1$

$$\zeta((1-2y)z) = \zeta(-z)$$

when $z \gg 0$

$$\zeta(-z) \to 0$$

When $z \ll 0$

$$\zeta(-z) \to \text{large}$$

@ $y = 0 \quad z \gg 0$

$$\zeta(z) \Rightarrow \text{large} \quad (w^{(k+1)} \text{ changes quite alot)} \quad \nabla \text{also large}$$

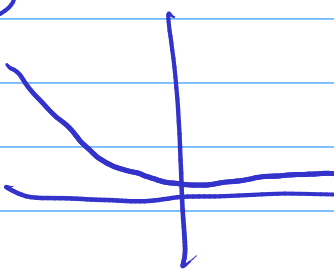$$z \ll 0$$
$$\zeta(z) \to 0 \qquad w\text{'s wont change}$$

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

$$\frac{d}{dz}\sigma(z) = \sigma(1-\sigma)$$

$$\frac{\partial y(\theta)}{\partial z} =$$

$$\sigma((1-2y)z)(1-2y)$$

(Take 1 example, treat as SGD)
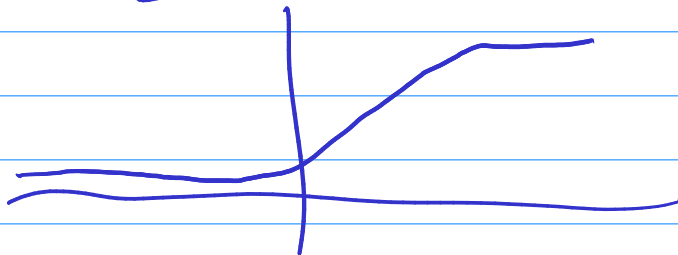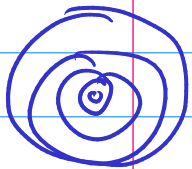
why is this cost functions any better than MSE

Pt

$$\max \log P(y_1, y_2 \cdots y_m \mid x_1, x_2 \cdots x_m)$$

$$= \max \log \prod_{i=1}^{m} P(y_i ? x_i)$$

$$= \max \log \left[ \sigma(z_i)^{y_i \cdot p} (1 - \sigma(z_i))^{1 - \log y_i} \right]$$

$$= \max \sum \left( y_i \log \sigma(z_i) + (1-y) \log (1 - \sigma(z_i)) \right)$$

$$= \min \left[ \sum - \left( y_i \log \sigma(z_i) + (1-y) \log (1 - \sigma(z_i)) \right) \right]$$



$$= \min - \left[ y \log \sigma(z) + (1-y) \log (1 - \sigma(z)) \right]$$