

Logistic $\sigma(x) = \frac{1}{1+e^{-x}} \in (0, 1)$

$$\sigma'(x) = \sigma(x)(1-\sigma(x))$$

Hypobolic tangent

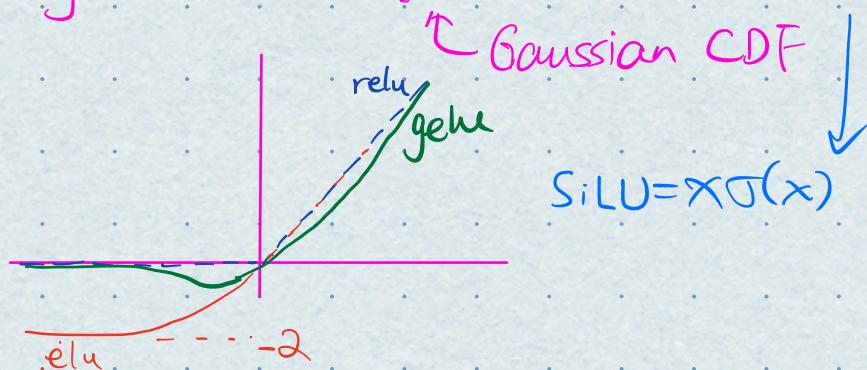
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \in (-1, 1)$$

$$\tanh'(x) = 1 - \tanh^2(x)$$

ReLU $\text{relu}(x) = \max(0, x)$

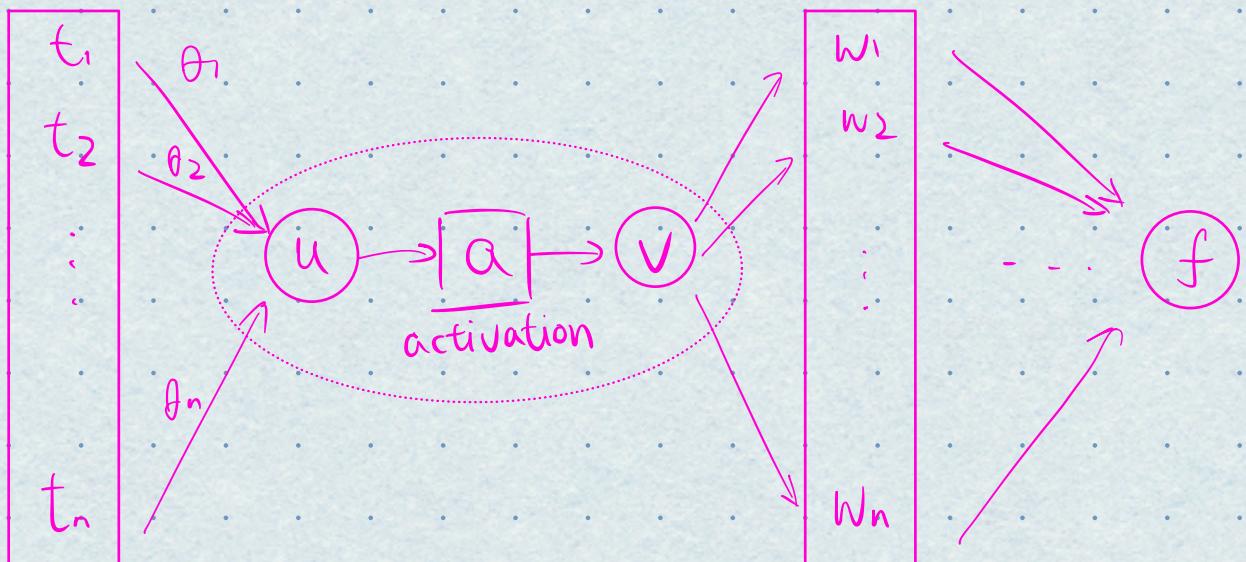
Leaky ReLU $f(x) = \begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$

GEU $\text{gelu}(x) = x\Phi(x) \approx x\sigma(1.702x)$



ELU $\text{elu}(x) = \begin{cases} 2(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$

Back propagation



$$\frac{\partial f}{\partial \theta_1} = \underbrace{\frac{\partial f}{\partial v}}_{?} \underbrace{\frac{\partial v}{\partial u}}_{a'(\cdot)} \underbrace{\frac{\partial u}{\partial \theta_1}}_{t_1}$$

$$\frac{\partial f}{\partial v} = \sum_j \underbrace{\frac{\partial f}{\partial w_j}}_{\text{compute on } w_j} \frac{\partial w_j}{\partial v}$$

$$\frac{\partial f}{\partial \theta_i} = \left(\sum_j \frac{\partial f}{\partial w_j} \frac{\partial w_j}{\partial v} \right) \frac{\partial v}{\partial \theta_i} = \underbrace{\left(\sum_j \left(\sum_k \frac{\partial f}{\partial z_k} \frac{\partial z_k}{\partial w_j} \frac{\partial w_j}{\partial v} \right) \frac{\partial v}{\partial \theta_i} \right)}_{\text{repeated terms}}$$

$$\frac{\partial f}{\partial \theta_k} = \left(\sum_j \frac{\partial f}{\partial w_j} \frac{\partial w_j}{\partial v} \right) \frac{\partial v}{\partial \theta_k} = \left(\sum_j \left(\sum_k \frac{\partial f}{\partial z_k} \frac{\partial z_k}{\partial w_j} \frac{\partial w_j}{\partial v} \right) \frac{\partial v}{\partial \theta_k} \right)$$

(Shannon) Entropy:

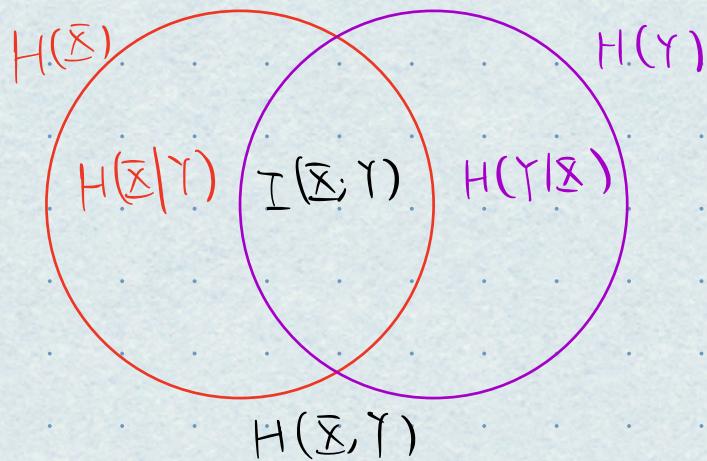
$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i)$$

Conditional entropy:

$$H(Y|X) = \sum_x P_X(x) H(Y|X=x)$$

$$H(Y|X=x) = - \sum_y P_{Y|X}(y|x) \log P_{Y|X}(y|x)$$

$$\Rightarrow H(Y|X) = \sum_{x,y} P(x,y) \log \frac{P(x)}{P(x,y)}$$



Chain rule: $H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i | X_1, \dots, \hat{X}_i)$

Mutual information

$$I(X;Y) = \sum_{x,y} P_{X,Y}(x,y) \log \frac{P_{X,Y}(x,y)}{P_X(x)P_Y(y)}$$

$$I(X;Y) = D_{KL}(P_{X,Y} || P_X P_Y)$$

Relative entropy (KL divergence)

$$D_{KL}(P || Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

Cross entropy

$$\begin{aligned} CE(P, Q) &= -\sum_x P(x) \log Q(x) \\ &= H(P) + D_{KL}(P || Q) \end{aligned}$$

Diversity : Hill number

$$D^q = \left(\sum_{i=1}^R p_i^q \right)^{\frac{1}{q-1}}$$

$$D^1 = \exp \left\{ - \sum_{i=1}^R p_i \ln(p_i) \right\} \text{ Shannon index}$$

Gradient descent methods in NN

Cost function: $J(\theta)$, $\theta \in \mathbb{R}^d$ NN params

Gradient: $\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}$

Generally speaking, $\theta^{(t+1)} = \theta^{(t)} - \Delta^{(t)}$

① Gradient descent

$$\Delta^{(t)} = \eta \nabla_{\theta} J(\theta^{(t)}) \quad \leftarrow \eta: \text{learning rate}$$

② Momentum

$$\Delta^{(t)} = \gamma \Delta^{(t-1)} + \eta \nabla_{\theta} J(\theta^{(t)})$$

γ : decay factor; 0.9

③ Adagrad

$$\Delta^{(t)} = \frac{\eta}{\sqrt{G^{(t)} + \epsilon I}} \nabla_{\theta} J(\theta^{(t)})$$

$$G^{(t)} = \sum_{i=1}^t \left[\frac{\partial J}{\partial \theta_i} \right]^2$$

$$i=1 \quad \left[\frac{\partial J}{\partial \theta_d} \right]_{\theta=\theta^{(i)}}$$

④ RMS prop

$$\Delta^{(t)} = \frac{n}{\sqrt{G^{(t)} + \epsilon I}} \nabla_{\theta} J(\theta^{(t)})$$

$$G^{(t)} = \gamma G^{(t-1)} + (1-\gamma) \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}_{\theta=\theta^{(t)}}^T$$

0.9 ↗

⑤ Adam

$$m^{(t)} = \beta_1 m^{(t-1)} + (1-\beta_1) \nabla_{\theta} J(\theta^{(t)})$$

$$v^{(t)} = \beta_2 v^{(t-1)} + (1-\beta_2) (\nabla_{\theta} J(\theta^{(t)}))^2$$

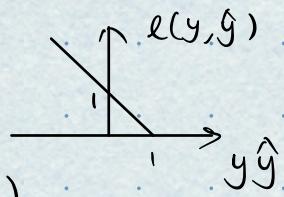
$$\hat{m}^{(t)} = \frac{m^{(t)}}{1-\beta_1^t}, \quad \hat{v}^{(t)} = \frac{v^{(t)}}{1-\beta_2^t}$$

element wise

$$\Delta^{(t)} = \frac{n}{\sqrt{G^{(t)} + \epsilon}} \hat{m}^{(t)}$$

Typical values: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

Hinge loss



$$\text{Binary: } l(y, \hat{y}) = \max_{\{\pm 1\}} (0, 1 - y\hat{y})$$

N-class: $\cup_{i \in \{-\infty, +\infty\}}$

$$\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N)$$

$$l(y, \hat{y}) = \sum_{i \neq y} \max(0, 1 + \hat{y}_i - \hat{y}_y)$$

Hinge embedding loss

$$l(x, y) = \begin{cases} x & \text{if } y=1 \\ \max(0, \Delta - x) & \text{if } y=-1 \end{cases}$$

x measures the distance between a pair instances
 $y = 1/-1$ means the pair is similar/dissimilar

Cross entropy loss

$$\text{Binary: } l(y, \hat{y}) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\text{N-class: } \hat{y} = (\hat{y}_1, \dots, \hat{y}_N)$$

$$l(y, \hat{y}) = - \sum_{i=1}^n \mathbb{1}_{y=i} \log \hat{y}_i$$

Logistic loss

$$\text{Binary: } \ln(1 + e^{-y\hat{y}}) \xleftarrow[\text{on NN output with } \sigma(\cdot)]{} \text{equal to CE loss applied}$$