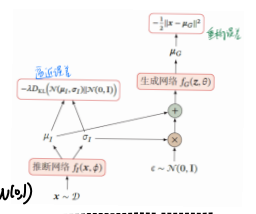


**AE 自编码器**  
 $\phi: X \rightarrow F, \text{e.g. } h = \sigma(W_1 x + b_1)$   
 $\psi: F \rightarrow X, \text{e.g. } x' = \sigma(W_2 h + b_2)$   
 $\phi, \psi = \arg \min_{\phi, \psi} \|x - (\psi \circ \phi)(x)\|^2$   
 $= \|x - x'\|^2 = \|x - \sigma(W_2 \sigma(W_1 x + b_1) + b_2)\|^2$

**VAE 变分自编码器**  
 推理神经网络:  $q(z|x; \phi) = N(z; \mu_1, \sigma_1^2 I)$   
 $h = \sigma(W^{(1)} x + b^{(1)})$   
 $\mu_1 = W^{(2)} h + b^{(2)}$   
 $\sigma_1^2 = \text{softplus}(W^{(3)} h + b^{(3)})$   
 生成网络:  $p(x|z; \theta) = N(x; \mu_2, \sigma_2^2 I)$

总目标函数:  $\max_{\theta, \phi} E_{L, D} \log \frac{p(x|z; \theta) p(z; \theta)}{q(z; \phi)}$   
 $= \max_{\theta, \phi} E_{z \sim q(z; \phi)} \left[ \log \frac{p(x|z; \theta) p(z; \theta)}{q(z; \phi)} \right]$   
 $= \max_{\theta, \phi} E_{z \sim q(z; \phi)} \left[ \log p(x|z; \theta) - \text{KL}(q(z|x; \phi) \| p(z; \theta)) \right]$   
 ①  $\frac{1}{M} \sum_{m=1}^M \log p(x^{(m)}; \theta), z^{(m)} = \mu_1 + \sigma_1 \epsilon, \epsilon \sim N(0, I)$



重参数化技巧, 使得  $z^{(m)}$  对  $\mu_1$  和  $\sigma_1$  可求导  
 $J(\phi, \theta) = \sum_{m=1}^M \left( \frac{1}{M} \sum_{m=1}^M \log p(x^{(m)}; z^{(m)}, \theta) - \text{KL}(q(z|x^{(m)}; \phi) \| p(z; \theta)) \right)$   
 采样  $\epsilon \sim N(0, I)$ , 计算  $z = \mu_1 + \sigma_1 \epsilon, \mu_2 = f_\theta(z; \theta)$   
 $J(\phi, \theta) = -\frac{1}{2} \|x - \mu_2\|^2 - \lambda \text{KL}(q(z|x; \phi) \| p(z; \theta))$

GAN 流程图:  $x \sim D \rightarrow$  判别网络  $D(x, \theta) \rightarrow 1/0$  生成  
 $z \sim N(0, I) \rightarrow$  生成网络  $G(z, \theta) \rightarrow$  判别网络  $D(x, \theta) \rightarrow 1/0$  生成  
 判别网络的目标函数:  $\min_{\theta} -E_x [y \log p(y=1|x) + (1-y) \log p(y=0|x)]$   
 生成网络的目标函数:  $\min_{\theta} E_{z \sim p(z)} [\log(1 - D(G(z; \theta); \theta))] + \max_{\theta} E_{z \sim p(z)} [\log(D(G(z; \theta); \theta))]$

在最优判别器下, 生成网络的目标函数简化为  $L(G|D^*) = E_{x \sim p(x)} [\log D^*(x)] + E_{x \sim p_G(x)} [\log(1 - D^*(x))]$   
 $= E_{x \sim p(x)} \left[ \log \frac{p(x)}{p(x) + p_G(x)} \right] + E_{x \sim p_G(x)} \left[ \log \frac{p_G(x)}{p(x) + p_G(x)} \right] = 2 \text{JS}(p, p_G) - 2 \log 2$   
 其中  $\text{JS}(p, p_G) = \frac{1}{2} \text{KL}(p \| \frac{p+p_G}{2}) + \frac{1}{2} \text{KL}(p_G \| \frac{p+p_G}{2})$   
 还可以计算  $L(G|D^*) = E_{x \sim p_G(x)} [\log D^*(x)] = E_{x \sim p_G(x)} \left[ \log \frac{p(x)}{p(x) + p_G(x)} \right]$   
 $= E_{x \sim p_G(x)} \left[ \log \frac{p(x)}{p_G(x)} \right] + E_{x \sim p_G(x)} \left[ \log \frac{p_G(x)}{p(x) + p_G(x)} \right] = -\text{KL}(p_G \| p) + E_{x \sim p_G(x)} [\log(1 - D^*(x))]$   
 $= -\text{KL}(p_G \| p) + 2 \text{JS}(p, p_G) - 2 \log 2 - E_{x \sim p_G(x)} [\log D^*(x)]$

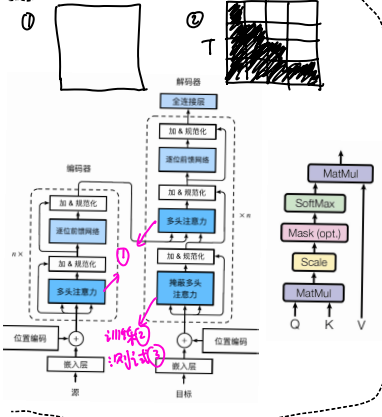
WGAN 用 Wasserstein 距离衡量分布, 解决 ① ② 模式崩塌, 即生成图片都很像  
 ① Total Variation:  $\delta(p, p_G) = \max_A |p(A) - p_G(A)|$   
 ② KL distance:  $\text{KL}(p, p_G) = \int \log \left( \frac{p(x)}{p_G(x)} \right) p(x) dx$   
 ③ JS-divergence:  $\text{JS}(p, p_G) = \text{KL}(p \| p_m) + \text{KL}(p_G \| p_m)$  其中  $p_m = \frac{p+p_G}{2}$   
 Wasserstein-1, 或 Earth-Mover distance:  $W(p, p_G) = \inf_{\gamma \in \Pi(p, p_G)} \int \|x - y\| \gamma(x, y)$   
 $= \int \max_x |F(x) - F_G(x)| dx$

Kantorovich-Rubinstein 对偶  $W(p, p_G) = \sup_{\|f\|_{Lip} \leq 1} E_{x \sim p}[f(x)] - E_{x \sim p_G}[f(x)]$ , 其中  $\|f\|_{Lip}$  是 f 的 Lipschitz 常数. 根据对偶的优化目标:  $\max_{f \in \mathcal{F}} \{ E_{x \sim p}[f(x)] - E_{x \sim p_G}[f(x)] \}$   
 加入正则化惩罚项:  $\min_{f, w} E_{z \sim p(z)} [f_w(G_\theta(z))] - E_{x \sim p(x)} [f_w(x)] + \lambda E_z [\| \nabla_x f_w(x) \|^2]$ , 其中  $\hat{x}$  和  $t x_1 + (1-t)x_2$  同分布,  $x_1 \sim p, x_2 \sim p_G$  with  $z \sim p(z), t \sim U(0,1)$

**注意力机制**  
 高斯核函数  $K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$   
 $f(x) = \sum_{i=1}^n \alpha(x, x_i) y_i$   
 $\alpha(x, x_i) = \frac{\exp(-\frac{1}{2}(x-x_i)^2)}{\sum_{j=1}^n \exp(-\frac{1}{2}(x-x_j)^2)}$   
 $= \sum_{j=1}^n \text{Softmax}(-\frac{1}{2}(x-x_j)^2) y_j$

加性注意力:  $q \in R^d, k \in R^d$   
 $A(q, k) = W_V \tanh(W_q q + W_k k) \in R$   
 点积注意力:  $q, k \in R^d$   
 $a(q, k) = \frac{q \cdot k}{\sqrt{d}}$  除 d 是为了保持 scale 不变  
 批量化的公式:  $\text{Softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V \in R^{n \times n}$   
 其中  $Q \in R^{n \times d}, K \in R^{m \times d}, V \in R^{m \times n}$

**位置编码**  
 $P_i, z_j = \sin\left(\frac{i}{10000 \frac{2j}{d}}\right), P_i, z_{j+1} = \cos\left(\frac{i}{10000 \frac{2j}{d}}\right)$   
 $\begin{pmatrix} \cos(\delta w_j) & \sin(\delta w_j) \\ -\sin(\delta w_j) & \cos(\delta w_j) \end{pmatrix} \begin{pmatrix} P_i, z_j \\ P_i, z_{j+1} \end{pmatrix} = \begin{pmatrix} P_i + \delta, z_j \\ P_i + \delta, z_{j+1} \end{pmatrix}$   
 最优 Arm 与当前 Arm 的 gap  
 记  $\Delta_n(v) = U^n(v) - \mu_n(v)$   
 记  $T_n(t) = \sum_{s=1}^t I(A_s = n)$ , 即使使用 Arm n 的总次数  
 期望  $E[T_n(t)] = \sum_{s=1}^t \Delta_n E[I(A_s = n)]$  是超参数  
 ETC: Explore then Commit  
 Explore: 每个 Arm 都尝试 m 次  
 Commit: 选则 Explore 阶段从最大的  
 $A_t = \begin{cases} (t \bmod k) + 1, & \text{if } t \leq mk \\ \arg \max_i \hat{\mu}_i(mk), & \text{if } t > mk \end{cases}$   
 其中  $\hat{\mu}_i(t) = \frac{1}{T_i(t)} \sum_{s=1}^t I(A_s = i) X_s$   
 $T_i(t)$  就等于  $m$  (进使得  $R_n \leq \Delta + C\sqrt{n}$ )



**Bandit**  
 遗憾值  $R_n = n \max_a \mu_a - E \left[ \sum_{t=1}^n X_t \right]$   
 越小越好 最优回报 当前策略实际回报  
 Stochastic bandit  
 每个臂的期望:  $\mu_a(v) = \int_{-\infty}^{\infty} x dP_a(x)$   
 记  $\mu^*(v) = \max_a \mu_a(v)$   
 对任意策略  $\pi$ , 目标函数为遗憾值最小  
 $R_n(\pi, v) = n \mu^*(v) - E \left[ \sum_{t=1}^n X_t \right]$

当  $k=2$  时,  $\Delta_1 = 0, \Delta_2 = 0$ , 则  $R_n \leq m \Delta + (n-2m) \Delta \exp(-m \Delta / 4) \approx C \sqrt{n}$   
 取  $m = \max \left\{ 1, \left\lceil \frac{4}{\Delta} \log \left( \frac{n \Delta}{4} \right) \right\rceil \right\}$ , 则可得到  $R_n \leq \min \{ n \Delta, \Delta + \frac{4}{\Delta} (1 + \max \{ 0, \log \left( \frac{n \Delta}{4} \right) \}) \}$   
 若  $\Delta$  未知, 则最优的界为  $R_n = O(\sqrt{n})$   
 若 ETC 对  $\Delta$  未知, 则 ETC 通常比 UCB 更差

**Upper Confidence Bound, UCB**  
 若  $X_t$  服从  $\mu=1$  的 1-Subgaussian 分布  
 令  $\hat{\mu} = \frac{1}{n} \sum_{t=1}^n X_t$ , 则  $P(\mu > \hat{\mu} + \sqrt{\frac{2 \log(1/\delta)}{n}}) \leq \delta$ , for all  $\delta \in (0,1)$   
 定义 UCB 指标: 其中  $T_i(t-1)$  是 t 时刻前用 Arm i 的次数  
 $UCB_i(t-1, \delta) = \begin{cases} \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} & \text{if } T_i(t-1) > 0 \\ \hat{\mu}_i(t-1) & \text{o.w.} \end{cases}$   
 算法: Input:  $k, \delta$ , for  $t=1, \dots, n$ ,  
 Choose  $A_t = \arg \max_i UCB_i(t-1, \delta)$   
 Observe reward  $X_t$  and update  $UCB_i$   
 比 ETC 的优越:  
 a. 不需要事先知道 Suboptimality gap  
 b. 当超过两个臂时, 表现更好  
 c. 算法设计可以依赖于 horizon  $n$ .

**贝尔曼方程 Bellman equation**  
 $V_\pi(s; \epsilon) = E_\pi [R_t + \gamma V_\pi(S_{t+1}) | S_t = s; \epsilon]$   
 $V_\pi(s) = \sum_a \pi(a|s) \left\{ \sum_{s'} P(s'|s, a) r(s, a, s') + \gamma \sum_{s'} P(s'|s, a) V_\pi(s') \right\}$   
 $V_\pi(s) = E_{a \sim \pi(\cdot|s)} [Q_\pi(s, a)]$   
 $Q_\pi(s, a) = \sum_{s'} P(s'|s, a) r(s, a, s') + \gamma \sum_{s'} P(s'|s, a) V_\pi(s')$

**Markov Decision Tree, MDP**: 状态空间, 动作空间, 状态转移函数, 奖励函数, 折扣因子  
 动作价值函数: 依赖于  $S_t, a_t$  和  $\pi$   
 $Q_\pi(S_t, a_t) = E_{S_{t+1}, A_{t+1}, \dots, S_n, A_n} [U_t | S_t = s_t, A_t = a_t]$   
 最优动作价值函数: 消除策略  $\pi$  的影响, 只评价  $S_t$  和  $a_t$  的好坏  
 $Q^*(s_t, a_t) = \max_{\pi} Q_\pi(s_t, a_t) \Rightarrow \pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} Q^*(s, a) \\ 0 & \text{o.w.} \end{cases}$   
 状态价值函数: 消除动作  $a_t$  的影响, 只依赖于  $\pi$  和  $s_t$   
 $V_\pi(s_t) = E_{A_t \sim \pi(\cdot|s_t)} [Q_\pi(s_t, A_t)] = \sum_a \pi(a|s_t) Q_\pi(s_t, a)$   
 最优状态价值函数: 消除  $\pi$  和  $a_t$  的影响, 只依赖于  $s_t$   
 $V^*(s_t) = E_{a \sim \pi^*(\cdot|s_t)} Q^*(s_t, a) = \max_a Q^*(s_t, a)$   
 定义  $\pi^* = \arg \max_{\pi} V_\pi(s)$ ,  $\forall s \in S$   
 $V_\pi(s_t)$  也可写成  $E_{A_t, S_{t+1}, A_{t+1}, \dots, S_n, A_n} [U_t | S_t = s_t]$

最优贝尔曼方程:

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a) \Rightarrow \pi^*(a|s) = \begin{cases} 1 & a = \operatorname{argmax}_a Q^*(s, a) \\ 0 & \text{o.w.} \end{cases}$$

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) r(s, a, s') + \sum_{s'} P(s'|s, a) \gamma V^*(s')$$

$$= \sum_{s'} P(s'|s, a) r(s, a, s') + \sum_{s'} P(s'|s, a) \gamma \max_{a'} Q^*(s', a')$$

计算状态价值函数  $V_\pi$

$$\vec{V}_\pi = R_\pi + \gamma \tilde{P}_\pi \vec{V}_\pi \Rightarrow \vec{V}_\pi = (I - \gamma \tilde{P}_\pi)^{-1} R_\pi$$

(例):  $P(s'|s, a=1) = \begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.7 & 0.2 & 0.1 \\ 0.6 & 0.2 & 0.2 \end{pmatrix}$ ,  $P(s'|s, a=2) = \begin{pmatrix} 0.1 & 0.6 & 0.3 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.4 & 0.5 \end{pmatrix}$

$T(s, a, s') \equiv r(s) = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}$ , 则  $\tilde{P}_\pi(s'|s) = \sum_{a \in A} \pi(a|s) \cdot P(s'|s, a)$

若  $\pi(a|s) = \begin{cases} \frac{1}{2} & a=1 \\ \frac{1}{2} & a=2 \end{cases}$ , 则  $\tilde{P}_\pi(s'|s) = 0.5 \times \begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.7 & 0.2 & 0.1 \\ 0.6 & 0.2 & 0.2 \end{pmatrix} + 0.5 \times \begin{pmatrix} 0.1 & 0.6 & 0.3 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.4 & 0.5 \end{pmatrix}$

若  $\pi(a|s) = \begin{cases} \frac{1}{2} & a=1 \\ \frac{1}{2} & a=2 \end{cases}$ , 则  $\tilde{R}_\pi(s) = \sum_{a \in A, s' \in S} \pi(a|s) \cdot P(s'|s, a) \cdot T(s, a, s') = \begin{cases} \sum_{a \in A} \pi(a|s) \cdot r(s, a) & \text{若依赖于 } s' \text{ 但依赖于 } a \\ r(s) & \text{若依赖于 } s' \text{ 和 } a \end{cases}$

基础算法: ① 随机生成策略  $\pi$ , ② 计算价值函数  $V_\pi$

③ 更新策略  $\pi(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [r(s, a, s') + \gamma V_\pi(s')]$

④ 重复②和③至收敛, 返回确定性策略  $\pi$  (状态转移已知)

若状态转移未知: off-policy 策略 (DQN (连续), SARSA, on-policy 策略)

最优贝尔曼方程:  $Q^*(s_t, a_t) = E_{s_{t+1} \sim P(\cdot|s_t, a_t)} [R_t + \gamma \max_a Q^*(s_{t+1}, a) | s_t, a_t, a_t]$

定义损失函数  $L(w) = \frac{1}{2} [Q(s_t, a_t; w) - \hat{y}_t]^2$ , 其中  $\hat{y}_t = r_t + \gamma \max_a Q(s_{t+1}, a; w_{now})$

梯度为  $\nabla_w L(w) = (\hat{q}_t - \hat{y}_t) \nabla_w Q(s_t, a_t; w)$ , 其中  $\hat{q}_t = Q(s_t, a_t; w_{now})$

SARSA on policy, 不能使用经验回放

$\hat{y}_t = r_t + \gamma Q(s_{t+1}, a_{t+1})$ ,  $\tilde{a}_{t+1} \sim \pi_{now}(\cdot | s^{(t+1)})$  (Q-Learning 有 max, 而 SARSA 没有)

梯度:  $\nabla A x = A^T$ ,  $\nabla x^T A = A$ ,  $\nabla x^T A x = (A + A^T)x$ ,  $\nabla \|x\|^2 = 2x$  清除梯度:  $x \cdot \text{grad\_zero}()$

线性回归解析解:  $\hat{w} = (X^T X)^{-1} X^T Y$ , 无交叉项损失  $-y_i \cdot \log(P(y_i)) - (1 - y_i) \cdot \log(1 - P(y_i))$

知名函数性质: ① 连续可导 (允许少数点不可导) ② 非线性函数 ③ 导数计算简单 ④ 导数数值域合适

Q-Learning 行为策略和目标策略不同

① 采样:  $a_t = \begin{cases} \operatorname{argmax}_a Q^{(t-1)}(s_t, a) & \text{with Pr } 1-\epsilon \\ \text{uniformly sampling} & \text{with Pr } \epsilon \end{cases}$

② 更新:  $Q(s_t, a_t) = (1 - \alpha) Q(s_t, a_t) + \alpha \hat{y}_t$   
其中  $\hat{y}_t = r_t + \gamma \max_a Q^{(t-1)}(s_{t+1}, a)$

③ 返回最后的动作价值函数  $Q^{(T)}$

前馈神经网络

$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$   
 $a^{(l)} = f_l(z^{(l)})$

AdaGrad:  $G_t = \sum_{s=1}^t g_s \odot g_s$   
 $\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{G_t}} \odot g_t$

RMSprop:  $G_t = \rho G_t + (1 - \rho) g_t \odot g_t$   
 $\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{G_t}} \odot g_t$

动量法:  $\Delta \theta_t = \rho \Delta \theta_{t-1} - \alpha g_t = -\alpha \sum_{s=1}^t \rho^{t-s} g_s$

Adam:  $M_t = \beta M_{t-1} + (1 - \beta) g_t$ ,  $G_t = \beta_2 G_{t-1} + (1 - \beta_2) g_t \odot g_t$

随机初始化:  $N(0, \sigma^2)$  或  $[r, r]$ ,  $r = \sqrt{\sigma^2}$

反向传播

$\delta^{(l)} = \frac{\partial L(y, \hat{y})}{\partial z^{(l)}} = \frac{\partial L}{\partial z^{(l+1)}} \cdot \frac{\partial z^{(l+1)}}{\partial a^{(l)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}}$

$\frac{\partial L(y, \hat{y})}{\partial w^{(l)}} = \delta^{(l)} (a^{(l-1)})^T$ ,  $\frac{\partial L(y, \hat{y})}{\partial b^{(l)}} = \delta^{(l)}$

MSE:  $L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$  Logit:  $L(y, \hat{y}) = -y \hat{y} + \log(1 + e^{\hat{y}})$

$\Delta \theta_t = -\frac{\alpha}{\sqrt{G_t} + \epsilon} \hat{M}_t$

Logistic:  $\sigma(x) = \frac{1}{1 + e^{-x}}$ ,  $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$

Tanh:  $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ ,  $\tanh'(x) = 1 - \tanh^2(x)$

$\tanh(x) = 2\sigma(2x) - 1$

ReLU(x) =  $\begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$  ① 计算高效, ② 单侧抑制, 避免饱和, ③ 稀疏性, ④ 缓解梯度消失

Leaky ReLU(x) =  $\begin{cases} x & x \geq 0 \\ \gamma x & x < 0 \end{cases}$ ,  $\gamma$  可取 0.01

Softplus(x) =  $\log(1 + e^x)$ , 求导后为  $\sigma(x)$

Maxout(x) =  $\max_k z_k$ , 其中  $z_k = W_k^T X + b_k$

$a^{(l)} = f(\sum_{i=1}^{M_{l-1}} w_i^{(l)} a_i^{(l-1)})$ , 若  $w$  和  $a$  均值为 0, 且  $i, id$ . 则  $\text{Var}(a^{(l)}) = M_{l-1} \cdot \text{Var}(w_i^{(l)}) \cdot \text{Var}(a_i^{(l-1)})$ , 因此  $\text{Var}(w_i^{(l)}) = \frac{1}{M_{l-1}}$

考虑反向传播, 则  $\text{Var}(w_i^{(l)}) = \frac{1}{M_{l-1}}$ , 综合为  $\text{Var}(w_i^{(l)}) = \frac{2}{M_{l-1} + M_l}$ . 因此 Xavier 对 Logistic:  $r = \sqrt{\frac{6}{M_{l-1} + M_l}}$ ,  $\sigma^2 = 16 \times \frac{2}{M_{l-1} + M_l}$ . 对 Tanh:  $r = \sqrt{\frac{6}{M_{l-1} + M_l}}$ ,  $\sigma^2 = \frac{2}{M_{l-1} + M_l}$

He 初始化: ReLU 激活时,  $r = \sqrt{\frac{2}{M_{l-1}}}$ ,  $\sigma^2 = \frac{2}{M_{l-1}}$

BN:  $\hat{z}^{(l)} = \frac{z^{(l)} - E[z^{(l)}]}{\sqrt{\text{Var}(z^{(l)}) + \epsilon}}$ ,  $\gamma + \beta$

$\hat{x}_n = \frac{x_n - \min_i x_i}{\max_i x_i - \min_i x_i}$  标准化  $\hat{x}_n = \frac{x_n - \mu}{\sigma}$

$\mu = \frac{1}{M_l} \sum_{i=1}^{M_l} z_i^{(l)}$ ,  $\sigma^2 = \frac{1}{M_l} \sum_{i=1}^{M_l} (z_i^{(l)} - \mu)^2$

卷积:  $y_t = \sum_{k=1}^K w_k x_{t-k+1}$ ,  $y_{ij} = \sum_{u=1}^U \sum_{v=1}^V w_{uv} x_{i-u, j-v+1}$

交换性: 对  $X$  两端各补 0-1 和 -1-0 个零, 得到  $\tilde{X}_{M+2U-2, N+2V-2}$ , 则  $W \otimes \tilde{X} = \tilde{X} \otimes W$

$\frac{\partial f(y)}{\partial w_{uv}} = \sum_{i=1}^{M-U+1} \sum_{j=1}^{N-V+1} \frac{\partial f(y)}{\partial y_{ij}} x_{u+i-1, v+j-1}$  即  $\frac{\partial f(y)}{\partial W} = \frac{\partial f(y)}{\partial Y} \otimes X$

$Z^p = \sum_{d=1}^D W^p \otimes X^d + b^p$ ,  $Y^p = f(Z^p)$ ,  $f$  用 ReLU, 参差向量:  $P \times D \times U \times V + P$

反卷积

$CX = Y, X = C^T Y$

stride=1 时,  $4^3 \rightarrow 2$

stride > 1 时,  $5^3 \rightarrow 2$

$5 = o' = s(i'-1) + k = 2(2-1) + 3$

LSTM

$z_t = \tanh(W_z x_t + U_z h_{t-1} + b_z)$   
 $i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$   
 $c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$   
 $h_t = o_t \odot \tanh(c_t)$

GRU:  $h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$   
update:  $z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$   
reset:  $r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$   
 $\tilde{h}_t = \tanh(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$

梯度消失:  $\delta_{t,k} = \prod_{\tau=k}^t (\text{diag}(f'(z_\tau)) U^T) \delta_{t,k}$

$\delta_{t,k} = \sum_{\tau=k}^t \delta_{t,\tau} \frac{\partial L_{total}}{\partial U} = \sum_{\tau=k}^t \delta_{t,\tau} h_{\tau-1}^T$

loss.backward()

