# * Memory Networks
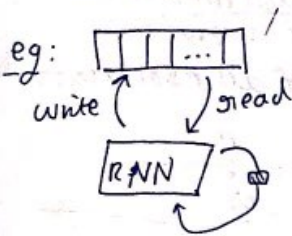
- NNs capture implicit knowledge but not explicit knowledge.
- Hence, these were invented to work on explicit

eg:



write ↑   ↓ read

Weston 2014

```
RNN
```

///. read chap 11   ⎫  ( GoodFellow et. al)
• read chap 12   ⎬  Practical methodologie,
                  ⎭  Application,

# * Auto encoder

Encoder   Decoder
   f          g
x→input

$f(x) = h$

$g(h) = \tilde{x}$

min $L(x, \tilde{x})$

min $L(x, g(f(x)))$

when $L = 0$

• $g(f(x)) = x$   (identity)

i.e $x \longrightarrow h \longrightarrow x$
              |
            code

" Uses:- ① Dimensionality redn    ($h$ is low dim)
                              dim $(h) <<$ dim $(x)$

② Information retrieval
   → easier to search in lower dimension
   → efficient to retrieve similar info.

③ Anomaly detection

④ Image denoising

/// 1. Under complete vs over complete
   autoencoder ↘ not useful

2. purposefully alter the model so that we get $\tilde{x}$ which is similar to $x$ but not same.

/// $L(x, g(f(x))) + \Omega(h)$          Regularized autoencoder
   $\Omega(h) = \lambda \sum_i |h_i|$        Sparse autoencoders

/// Denoising auto encoder (DAE)

$x + z \longrightarrow A E \longrightarrow x$          $L(x, g(f(\tilde{x})))$

$\underset{\tilde{x}}{}$

$c(x) = \tilde{x}$ (adding noise)

$\in N(x, \sigma^2 I)$                    If force AE to learn salient feature of data.

$x \underset{c}{\longrightarrow} \tilde{x} \underset{f}{\longrightarrow} h \underset{g}{\longrightarrow} x$

$\tilde{x} = x + 3$
    where $3 \sim N(0, \xi)$  (Gaussian).

/// Regularization based on restricting derivatives

$\cdot L(x, g(f(x))) \ast + \Omega(x, h)$          $h = f(x)$

$\cdot \Omega(x, h) = \lambda \sum_i (\nabla_x h_i)^2$          $\nabla_x h = \dfrac{\partial f(x)}{\partial x}$

Contractive
   Auto Encoder (CAE)          $\left|\dfrac{\partial f(x)}{\partial x}\right|_F^2$ : Frobenius norm

/// Chap 13 - Good Fellow

/// Deep Generative Models
- Variational autoencoder (VAE)
- Generative adversarial networks (GAN)
                    ① Need to assume some structure
- MCMC!
- Variational inference          ② Inference and estimation technique are not efficient

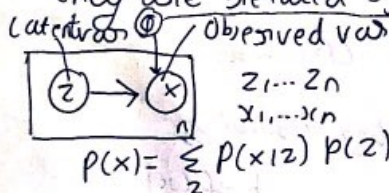$(\lambda-1)^2 - \alpha^2$      $\lambda-1=\pm\alpha$        $\lambda-1=\pm\alpha$        $\frac{1+2}$

③ Approximations are required, which lead to suboptimal solutions.

---

• NN can approximate complex functions
• Backpropagation can learn those functions efficiently over large scale data.

---

- Deep feature consistent variational autoencoder      2016
- Hierarchical variational autoencoder      2017
- Auto-encoding variational Bayes,
                    kingma, welling      ICLR 2014

- Tutorial on VAE    Doersch 2016

// Hierarchical models
 - multiple random vars
 - they are related by a dependency gph
Latent var ⓩ / observed var      global var

[diagram: box with (z) → (x)]

$z_1 \cdots z_n$
$x_1, \cdots x_n$

$P(x) = \sum\limits_{z} P(x|z)\, P(z)$

Algo : ① get $\theta$ // all params
        ② for $i = 1$ to $n$
            a. sample $z_i$ from $P(z)$
            b. sample $x_i$ from $P(x|z=z_i)$

// - Gaussian mixture model
   - MCMC
   - VI

// Inference: finding value of latent variable
Estimation: Finding value of parameter

• $P(x) = \int P(x|z) \cdot P(z)\, dz$   (we want to learn this)

★ so we use params $\theta$

$P(x|\theta) = \int P(x|z,\theta)\, P(z|\theta)\, dz$

① How to find a suitable model for $z$?
② How to compute the integration

• $z \sim N(0, I_1)$      (example)
• $x \sim N(f(z,\theta), \sigma^2 I)$
$f$ can be any complicated function

Core idea:   1. $z \sim N(0, I_1)$
   2. implement $f$ through a NN
   3. $x \sim N(f(z), \sigma^2 I)$
   Apply backprop to learn $f$

// Goal • maximize $P(x|z)$
    • But $P(z|x)$ is intractable
• Need $z$ so that when sampled from $P(x|z)$,
   $x$ is similar to real data.
• $D(\varphi(z) \| P(z|x))$ is small
   KL divergence
• $D(\varphi(z) \| P(z|x)) = E_{z \sim \varphi}\left[\log \varphi(z) - \log P(z|x)\right]$

   $= E_{z \sim \varphi}\left[\log \varphi(z) - \log P(x|z) - \log P(z) + \log P(x)\right]$

   $\left(P(z|x) = \dfrac{P(x|z)\, P(z)}{P(x)}\right)$

• $\log P(x) - D\left(\varphi(z) \| P(z|x)\right)$

$= E_{z \sim \varphi}\left[\log P(x|z)\right] - D\left(\varphi(z) \| P(z)\right)$

       $=$ core eqn of VAE

// $\varphi(z|x)$ : Recognition model
$\log f(x) - D(\varphi(z|x) \| P(z|x)) = E_{z \sim \varphi}\left[\log P(x|z)\right]$
        $- D\left(\varphi(z|x) \| P(z)\right)$

$$\varphi(z/x) = N(\mu(x), \Sigma(x))$$

$$x \rightarrow \boxed{E} \begin{array}{c} \nearrow \textcircled{\mu} \searrow \\ \searrow \textcircled{\Sigma} \nearrow \end{array} \boxed{D} \rightarrow \tilde{x}$$

## * GAN

$x \sim P$

P is unknown,
· so we use
some $\varphi$ and sample from it.

$\hookrightarrow \boxed{D} \longrightarrow$ real or fake
Discriminator

- If D is accurate, then we would know when $\varphi$ is close to P when D gives 'real' for generated samples.
- G is the generator, $G \rightarrow \tilde{x}$
- G and D would be NNs

$$\textcircled{z} \rightarrow \boxed{G} \rightarrow \textcircled{\tilde{x}} \xrightarrow{\textcircled{x}} \boxed{D} \rightarrow score \in \{0,1\}$$

If score is high if D say real o/w score is low
· z is sampled from a simple dist.

$z \sim N(0, I)$, or $U(a,b)$

· $G(z, \theta_g)$, $D(x, \theta_d)$
Need to find $\theta_g$ and $\theta_d$
* - Find $\theta_d$ s.t D is acc
 - Find $\theta_g$ s.t G is acc

// D always penalizes output of G, so that
G should always try to be better.
- Goal of D: $D(x) \uparrow$ and $D(G(z)) \downarrow$
- Goal of G: $D(G(z)) \uparrow$

/// $V(D,G) = E_{x \sim p_{data}}[\log D(x)]$
$\qquad\qquad + E_{\tilde{z} \sim p_z}[\log(1 - D(G(\tilde{z})))]$

2 players — D and G

$$\min_{G} \max_{D} V(D, G) \quad ///$$

* Algo:-

For a number of iteration do

    for h step do
      - sample $\{z^1, \ldots, z^m\}$ from $P_z$
      - sample $\{x^1, \ldots, x^m\}$ from $P_{data}$
      - update discriminator $D(x; \theta_d)$

    end for

    - sample $\{z^1, \ldots, z^m\}$ from $P_z$
    - update generator $G(z; \theta_g)$

end for

$$/// \quad E_{x \sim P_{data}}[\log D(x)] = \int_{-\infty}^{\infty} \log D(x) \, P_{data}^{(x)} \, dx \quad -①$$

$$E_{z \sim P_z}[\log(1 - D(G(z)))] = \int_{-\infty}^{\infty} \log(1 - D(G(z))) \, P_z(z) \, dz \quad -②$$

///

$z \sim P_z$
$\tilde{x} \sim P_g$ — represent $G(z)$

hence, $E_{\tilde{x} \sim P_g}[\log(1 - D(\tilde{x}))] = E_{z \sim P_z}[\log(1 - D(G(z)))]$

So,
$$V(D, G) = \int \log D(x) \, P_{data}^{(x)} \, dx + \int \log(1 - D(\tilde{x})) \, P_g(\tilde{x}) \, d\tilde{x}$$

/// max $a \log y + b \log(1-y)$, here $p = P_{data} \cup P_g$

D will attain max at $\dfrac{a}{a+b}$

$$D(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)}$$

$D \geq 1/2$ when $P_{data} = P_g$

$$C(G) = \max_{D} V(G, D)$$

$$= E_{x \sim p_d}[\log D^*_G(x)] + \cancel{E_{z \sim p_z}[\log D^*_G(x)]}$$

$$+ E_{z \sim p_z}[\log(1 - D^*_G(x))]$$

$$= E_{x \sim p_d}[\log D^*_G(x)] + E_{\tilde{x} \sim p_g}[\log(1 - D^*_G(\tilde{x}))]$$

$$\bullet \quad D^*_{(x)} = \frac{p_d(x)}{p_d(x) + p_g(x)} \qquad \text{where } D_* = \max(\dots)$$

$$\Rightarrow C(G) = E_{x \sim p_d}\left[\log \frac{p_d(x)}{p_d(x) + p_g(x)}\right]$$

$$+ E_{\tilde{x} \sim p_g}\left[\log \frac{p_g(\tilde{x})}{p_d(\tilde{x}) + p_g(\tilde{x})}\right]$$

$$= KL\left[p_d \| \frac{(p_d + p_g)}{2}\right] + KL\left[p_g \| \frac{(p_d + p_g)}{2}\right]$$

we want $\min C(G)$

/// later: $D^*(x) = 1/2 \; \forall x$

$$C(G) = \log\frac{1}{2} + \log\frac{1}{2} = 2\log(1/2)$$

///


///

/// imp :. Transpose convolution
$\equiv$ inv of ~~down~~ sampling $\ell$
down

· Radford et al. — architecture
· mode collapse · unrolled GAN