

Generative Models

- Desired properties :
1. Easy to sample $x_{\text{new}} \sim p(x)$
 2. Easy to evaluate likelihood $p(x)$
 3. Optionally extract latent features

Classical distribution (eg. Gaussian) doesn't capture the complexity of real-world data.
Mixtures (GMM) can but the number of components grows exponentially as the dimension increases.

1. Normalizing Flows

NF can model complex distribution by applying a transformation on a simple distribution $f(z) = x$

Probability in output space should be same as input space by normalizing $\left| \frac{\partial z}{\partial x} \right|$

Change of variable

f valid and $x = f(z)$

$$p_z(x) = p_z(z) \underbrace{\left| \det \left(\frac{\partial z}{\partial x} \right) \right|}_{\substack{\text{volume:} \\ \text{normalizes} \\ \text{probs}}} = p_z(f^{-1}(x)) \underbrace{\left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|}_{\text{Jacobian of } f^{-1}}$$

valid $\begin{cases} \text{invertible: input and output same dimension, monotonic} \\ \text{differentiable: (continuously differentiable) Jacobian exists at any point} \end{cases}$

stacking:

$$p_x(x) = p_0(z_0) \prod_{i=1}^n \left| \det \left(\frac{\partial f_i^{-1}(z_i)}{\partial z_i} \right) \right|$$

log prob (numerical stability):

$$\log p_x(x) = \log p_0(z_0) + \sum_{i=1}^n \log \left| \det \left(\frac{\partial f_i^{-1}(z_i)}{\partial z_i} \right) \right|$$

Reverse Parametrization

Goal: Evaluate $p_x(x)$ at any point x . x domain

$$g_\phi(x) = f^{-1}(x) \Rightarrow p_x(x) = p_0(g_\phi(x)) \left| \det \left(\frac{\partial g_\phi(x)}{\partial x} \right) \right| \text{ parametrize } f^{-1}.$$

Once learned g_ϕ 's, given a x we can go back to z_0 . Then compute $p_0(z_0)$ and the det's to have $p_x(x)$

φ learned by maximizing likelihood with given dataset.

Forward Parameterization

Goal: sample from $p_2(x)$. Given z , sample from $p_1(z)$ and also give $p_2(x)$
 $f_\theta(z) = x \Rightarrow p_2(x) = p_1(z) \left| \det \left(\frac{\partial f_\theta(z)}{\partial z} \right) \right|^{-1}$ z domain

Once learned f_θ , sample from $z_0 \sim p_0(z_0)$ and compute $z_i = f_\theta(z_{i-1})$ until k to have z_k . We can also have the density for this particular sample by computing intermediate $\left| \det \left(\frac{\partial f_\theta(z_{i-1})}{\partial z_{i-1}} \right) \right|$

Det Jacobian

Diagonal: g is applied element-wise, $\det(J) = \prod \frac{\partial g_i(x_i)}{\partial x_i}$ $O(D)$

Triangular: $g(x) = \begin{bmatrix} g_1(x_1) \\ g_2(x_1, x_2) \\ \vdots \end{bmatrix}$, $\det(J) = \prod \frac{\partial g_i(x_i)}{\partial x_i}$ $O(D)$

full: $g(x) = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \end{bmatrix}$, decompose $J = LU$, $\det(J) = \det(L) \det(U)$ $O(D^3)$

2. Variational Inference

Latent Variable Model

$p_\theta(x)$ where x high dim but it can be described with latent z .
That is $p(x|z)$ simple but $p(x)$ complex.

$$p_\theta(x) = \int p_\theta(x|z) dz = \int p_\theta(x|z) p_\theta(z) = E_{z \sim p_\theta(z)} [p_\theta(x|z)].$$

- Like in GMM:
1. $z \sim p_\theta(z)$
 2. $x \sim p_\theta(x|z)$

Tasks:

Inference

Given a sample x , find posterior distribution over z :

$$p_\theta(z|x) = \frac{p_\theta(x|z) p_\theta(z)}{p_\theta(x)}$$

Learning

Given dataset X , find θ that maximizes log likelihood

$$\max_{\theta} \log P_{\theta}(X) = \max_{\theta} \frac{1}{N} \sum_i^n \log p_{\theta}(x_i)$$

↳ usually intractable (integral), can't compute gradients

special case: N^f . Tractable with reverse parameterization. But somewhat constrained.

ELBO

We can define a lower bound g of f (or a collection of lower bounds) and try to maximize it.

$$\max_{\theta} f(\theta) \geq \max_{\theta} g(\theta).$$

for $p_{\theta}(x)$ we can also find a lower bound:

$$\log p_{\theta}(x) = \underbrace{E_{z \sim q(z)} [\log \frac{p_{\theta}(x, z)}{q(z)}]}_{\text{ELBO} = L(\theta, q)} + \text{KL}(q(z) || p_{\theta}(z|x))$$

which is tight if $\text{KL}(\cdot) = 0 \Leftrightarrow q(z) = p_{\theta}(z|x)$

EM is an alternative optimization based on this property:

- E-step: set $q(z) = p_{\theta}(z|x)$, with fixed θ from last iteration
- M-step: $\theta^{\text{new}} = \arg \max_{\theta} L(\theta, q)$, with q fixed

Optimizing ELBO

parameterize $q(z) \rightarrow q_{\phi}(z)$ to be able to optimize it.

$$\max_{\theta, \phi} L(\theta, \phi) : \text{need to compute } \nabla_{\theta} L(\theta, \phi) \text{ and } \nabla_{\phi} L(\theta, \phi)$$

$\nabla_{\theta} L(\theta, \phi)$

Assuming ϕ known and fixed: $E_{z \sim q_{\phi}(z)} [\log p_{\theta}(x, z)] = \frac{1}{S} \sum_i^S \log p_{\theta}(x_i)$

Expectation is approximated with Monte Carlo

$$\nabla_{\theta} E_{\mathbb{E}_{q_{\phi}(z)}[f_{\theta}(z)]} = \nabla_{\theta} \int q_{\phi}(z) f_{\theta}(z) dz = \int q_{\phi}(z) \nabla_{\theta} f_{\theta}(z) dz = E_{\mathbb{E}_{q_{\phi}(z)}} [\nabla_{\theta} f_{\theta}(z)] =$$

$$\frac{1}{N} \sum_i^N \nabla_{\theta} f_{\theta}(z_i)$$

$$\nabla_{\phi} L(\theta, \phi)$$

Assuming θ known and fixed: $E_{\mathbb{E}_{q_{\phi}(z)}[h_{\phi}(z)]} = \int q_{\phi}(z) h_{\phi}(z) dz$

$$h_{\phi}(z) = \log p_{\theta}(x, z) - \log q_{\phi}(z)$$

here we can't push gradient inside integral as in ∇_{θ} because the expectation depends on ϕ .

Reparameterization Trick: Sampling from some base distribution that doesn't depend on ϕ and then apply transformation

$$\Rightarrow E_{\mathbb{E}_{q_{\phi}(z)}[h_{\phi}(z)]} = \int q_{\phi}(z) h_{\phi}(z) dz = \int b(\epsilon) h_{\phi}(T(\epsilon, \phi)) d\epsilon =$$

$$E_{\mathbb{E}_{b(\epsilon)}[h_{\phi}(T(\epsilon, \phi))]}$$

Now as in ∇_{θ} we can push gradients inside and approximate expectation with MC.

Mean field assumption

When optimizing over the entire dataset X , to lower bound $p_{\theta}(X)$ we need to consider $q(z) \quad z \in \mathbb{R}^{N \times L}$

In practice we often simplify $q(z) = \prod q(z_i)$ which simplifies the ELBO and allows to approximate gradient with mini-batch for efficiency.

Implications:

- we can't capture dependences between dimensions (not important if data iid)
- less expressive but easier to optimize and infer

3. Variational Autoencoder

$$\text{ELBO} = E_{\mathbb{E}_{q(z)}[\log p_{\theta}(x|z) + \log p_{\theta}(z) - \log q_{\phi}(z)]} \quad \text{many choices for } p_{\theta}(x|z), p_{\theta}(z), q_{\phi}(z)$$

Prior $p_{\theta}(z)$

$z \in \mathbb{R}^L$, pick simple distributions for prior.
 $p(z) = N(z|0, I)$

Conditional likelihood $p_\theta(x|z)$

$z \in \mathbb{R}^L$, for every z we have to produce a different distribution.
Then pick a parametric distribution whose parameters are produced by some function f_ϕ that takes z as input.

$$\text{e.g. } f_\phi: \mathbb{R}^L \rightarrow \mathbb{R}^D \\ p(x|z) = N(x | \mu = f_\phi(z), \Sigma) \quad x \in \mathbb{R}^D$$

f_ϕ is the decoder in neural networks that converts latent z into params of $p_\theta(x|z)$

Variational distribution $q_\phi(z)$

Dataset X consist of N $x^{(i)} \in \mathbb{R}^D$ each corresponding to different latent $z^{(i)}$ with mean field assumption:

$$q_\phi(z^{(1)}, \dots, z^{(N)}) = \prod_i q_\phi(z^{(i)})$$

For simplicity $q_\phi(z^{(i)}) = N(z^{(i)} | \mu^{(i)}, \Sigma^{(i)}) \quad \phi^{(i)} = \{\mu^{(i)}, \Sigma^{(i)}\}$.

The choice of $q_\phi(z)$ should be similar to $p_\theta(z)$ as it minimize ELBO (KL) when $q_\phi(z)$ is close to $p_\theta(z)$. Additionally, KL between normals have closed form solution.

Every $x^{(i)}$ should map to the optimal $\phi^{(i)}$. We train an encoder g_λ that finds the optimal mapping.

Summary

$$L(\psi, \lambda) = \mathbb{E}_{z \sim q_\phi(z)} [\log p_\theta(x|z)] - KL(q_\phi(z) || p(z)) \quad \theta = f_\phi(z) \quad \phi = g_\lambda(x).$$

1. Compute $\phi = g_\lambda(x)$
2. Compute ML estimate (often done with a single sample):
 - draw $z' \sim q_\phi(z)$ reparametrization
 - Compute $\theta = f_\phi(z')$
 - Compute ELBO $L(\psi, \lambda) = (\log p_\theta(x|z') - KL(q_\phi(z') || p(z))$
3. Backprop gradients ($\nabla_\psi, \nabla_\lambda$)
4. Update weights

Once trained (of optimal ψ of decoder) we can generate samples:

- $z' \sim p(z) = N(z | 0, I)$
- $x \sim p_\theta(x|z) = N(x | f_\phi(z'), \Sigma)$

4. Generative Adversarial Network

Two competing network models:

Generator: takes noise sample $z \sim p(z)$ and transform it to $x = f_\theta(z)$.

more flexible than NF as $f_\theta(z)$ can be any function

Discriminator: distinguish generated data and real data

Problem: Sampling is easy but density (likelihood) evaluation is hard.

Density Estimation by Comparison

Test hypothesis that true data distribution $p^*(x)$ and generated data distribution $p_\theta(x)$ are equal

$$H_0: p^* = p_\theta \text{ vs } p^* \neq p_\theta \quad L(\theta, \phi)$$

Density Difference Approaches

Moment Matching

$p^* = p_\theta$ if their expectation of any statistical tests $s(x)$ equals . Approximated by MC.
e.g. $s(x) = \begin{bmatrix} x^1 \\ x^2 \\ e^x \\ \vdots \end{bmatrix}$ $E_{p^*}(x) = E_{p_\theta}(x)$? $E_{p^*}(x^2) = E_{p_\theta}(x^2)$? ..

Max Mean Discrepancy

Generalization of moment matching

Ratio-based Approaches

$$r^*(x) = \frac{p^*(x)}{p_\theta(x)} \quad \text{but true } r^*(x) \text{ hard to compute, instead approximate by } r_\phi(x).$$

Class Probability Estimation

Estimate density ratio via classifier which can distinguish between real and fake data.

- $Y=1 \Rightarrow \text{true}, Y=0 \Rightarrow \text{fake} \Rightarrow p(x|Y=1) = p^*(x), p(x|Y=0) = p_\theta(x)$.

- $p(Y=1) = \pi$

$$r^*(x) = \frac{p^*(x)}{p_\theta(x)} = \frac{p(x|Y=1)}{p(x|Y=0)} = \frac{\frac{p(Y=1|x)p(x)}{p(Y=1)}}{\frac{p(Y=0|x)p(x)}{p(Y=0)}} = \frac{p(Y=1|x)p(Y=0)}{p(Y=0|x)p(Y=1)} = \frac{p(Y=1|x)}{p(Y=0|x)} \frac{(1-\pi)}{\pi}$$

Discriminator $D_\phi(x) = p(Y=1|x)$. CE loss for learning:

$$L_{\theta, \phi} = \mathbb{E}_{(x, y) \sim p(x, y)} [-y \log [D_\phi(x)] - (1-y) \log [1 - D_\phi(x)]]$$

1. Optimize discriminator:

$$\phi^*(\theta) = \arg \min_{\phi} L_{\theta, \phi} . \text{ wlog } t_1 = 0.5 \Rightarrow r^*(x) \approx \frac{D_\phi(\theta)(x)}{1 - D_\phi(\theta)(x)}$$

$$\text{Ratio loss} = \frac{L_{\theta, \phi}(x, 1) + L_{\theta, \phi}(f_\theta(z), 0)}{2}$$

2. Optimize generator:

$$\text{drive } p(Y=1|x) = p(Y=0|x) \quad \text{ie } r^*(x) = 1 \quad y=0$$

$$\theta^* = \arg \max_{\theta} L_{\theta, \phi^*}(\theta) . \text{ Generative loss} = L_{\theta, \phi^*}(f_\theta(z), 0)$$

Minimax game (bilevel optimization) $\min_{\phi} \max_{\theta} \Pi(-L_{\theta, \phi})$
alternating optimization

Robustness

ML models are optimized following simple metrics (CE, L₁...) \Rightarrow prone to specifically handcrafted inputs (attacks)

It's a counterexample to the hypothesis that neural networks can learn meaningful representations capturing the semantics

Adversarial Examples

Perturbation set $P(x)$ which doesn't change the semantic of x (small perturbation) but changes the classification result.

$$\tilde{x} \in P(x) \Rightarrow f(x) = y \text{ and } f(\tilde{x}) \neq y$$

$$P_{\epsilon,p}(x) = \{\tilde{x} : \|\tilde{x} - x\|_p < \epsilon\}$$

Attack variants

Poisoning Attack

Modify training set so classifier has some desired property by the attacker after training

Targeted Attack

Attacker wants some sample classified as some specific class

Untargeted Attack

Attacker wants a sample misclassified as any class different than the correct one.

Construction of Adversarial Examples

$$\tilde{x}_x^* = \underset{\tilde{x} \in P(x)}{\operatorname{argmax}} \ell_{0/1}(f(\tilde{x}), y)$$

$\ell_{0/1}$, zero/one loss which could have undefined gradient.

$$\tilde{x}_x^* = \underset{\tilde{x} \in P(x)}{\operatorname{argmax}} L(f(\tilde{x}), y)$$

Use CE loss as alternative

Projected Gradient Descent

After each gradient step update the sample and project it back to the valid domain defined by $P(x)$.

$$x_{t+1} = \Pi(x_t + \eta \in \nabla_x L(f(x_t), y))$$

if $f(\tilde{x})$ not convex, cannot find global minimum

Fast Gradient-Sign Method

$$\hat{x} = \Pi(x + \eta \cdot \text{sign}(\nabla_x L(f(x), y)))$$

when norm defined on the radius ϵ , setting $\eta = \epsilon$ yields valid perturbation with only one step

Alternative formulation

$$\min_{\tilde{x}} D(x, \tilde{x}) \text{ st } l_{0/1}(f(\tilde{x}), y) \geq 0. \quad D = L_p \text{ dist between } x, \tilde{x}$$

\downarrow Lagrange

$$\min_{\tilde{x}} D(x, \tilde{x}) + \lambda \cdot L(\tilde{x}, y)$$

find min perturbed sample which yields misclassification

Robust Training

Training process that aims to produce models that are robust to adversarial attack.

Optimize robust loss (worst-case loss)

$$R_{\text{rob}} = \mathbb{E}_{(x,y) \in P_{\text{data}}} [\sup_{\tilde{x} \in \mathcal{X}} l(f(\tilde{x}), y)]$$

Adversarial training

use adversarial examples as proxy for worst-case perturbations to train robust classifiers.

$$\text{Perform SGD on } R_{\text{rob}}. \Rightarrow \nabla_{\theta} R_{\text{rob}} = \underbrace{\mathbb{E}_{(x,y) \in P_{\text{data}}} [\nabla_{\theta} (\sup_{\tilde{x} \in \mathcal{X}} l(f(\tilde{x}), y))]}_{\text{gradient of the worst case loss}}$$

Dunkin's Theorem

$\Delta(\theta)$ set of \tilde{x} where sup is attained. If $|\Delta(\theta)|=1$ then sup is differentiable:

$$\nabla_{\theta} (\sup_{\tilde{x} \in \Delta(\theta)} l(f(\tilde{x}), y)) = \nabla_{\theta} l(f(\tilde{x}_0^*), y) \quad \text{i.e. use adv. example as proxy}$$

Training outline:

1. Sample $(\tilde{x}_i, y_i) \sim P_{\text{data}}$

2. Find \tilde{x}_i with FGSM which has high loss $l(f(\tilde{x}_i), y_i)$

3. Update weights $\theta \leftarrow \theta - \eta \nabla_{\theta} l(f(\tilde{x}_i), y_i)$

- Cons:
- 10x slow down if we use powerful attack on step 2.
 - Lower accuracy on resulting model
 - No theoretical guarantee (certificate)

Certifiable Robustness

Prove that the classifier's prediction doesn't change within a radius (adversarial-fade)

Exact Verification

Algorithm that return YES \Leftrightarrow no adversarial examples.

Designed for NN with ReLU activations.

However exact verification for ReLU activation is NP-complete. But still tractable if network size is small-medium

As MILP

Exact certification of ReLU network is possible with mixed integer linear programming: LP polynomial but MILP NP-complete.

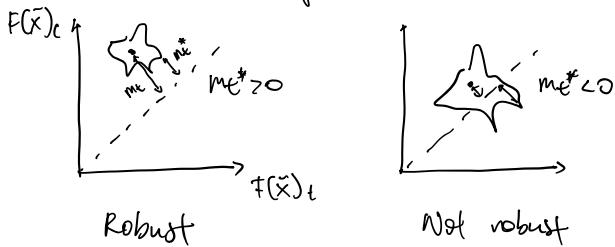
$m_t = f(x)_c - f(x)_t$ $f(x)_c$ class with highest prob (logit).

Worst case margin: $m_t^* = \min_{\tilde{x}} f(\tilde{x})_c - f(\tilde{x})_t$ s.t. $\|\tilde{x} - x\|_p \leq \epsilon$

- $m_t^* > 0$: classification can't be changed from c^* to t

- if $\forall t \neq c$ we have $m_t^* \geq 0 \rightarrow$ certified robustness

- if $\exists t \neq c$ with $m_t^* < 0 \rightarrow$ exist adversarial example $\tilde{x} \in P_G$)



$$M_{\ell}^* = \max_{\tilde{x}} [\hat{x}^{(u)}]_{c^*} - [\hat{x}^{(l)}]_t \quad \hat{x}^{(l)} \text{ pre-relu at layer } \ell$$

s.t. $\|\tilde{x} - x\|_p \leq \epsilon$

$$\begin{aligned} y^{(0)} &= \tilde{x} \\ \hat{x}^{(\ell)} &= w_\ell y^{(\ell-1)} + b_\ell \quad \forall \ell \\ y^{(\ell)} &= \text{ReLU}(\hat{x}^{(\ell)}) \quad \forall \ell = 1 \dots L-1 \end{aligned}$$

Encoding LP

$$\begin{aligned} p=0 &\rightarrow x_i - \hat{x}_i \leq \epsilon \quad \forall i & p=2 & \text{straight forward} \\ \hat{x}_i - x_i \leq \epsilon \quad \forall i & & p=1 & \text{quadratic programming} \end{aligned}$$

Encoding ReLU

Add binary variable a : $y_i \leq a_i x_i$ and $y_i \geq 0$ and $y_i \geq x_i$

But this adds an extra variable, making it a quadratic programming ($a_i x_i$)

If we have an upper and lower bound $[l_i, u_i]$ on the input x_i :

$$y_i = \text{ReLU}(x_i) \Leftrightarrow \begin{aligned} 1) \quad y_i &\leq x_i - l_i(1-a_i) \quad \forall i && \text{still linear and integer} \\ 2) \quad y_i &\leq a_i \cdot u_i \\ 3) \quad y_i &\geq x_i \\ y_i &\geq 0 \\ a_i &\in \{0, 1\} \end{aligned}$$

- 1) $a_i=1 \Rightarrow y_i=x_i$ $a_i=0 \Rightarrow y_i \leq x_i - l_i$
- 2) $a_i=1 \Rightarrow y_i \leq u_i$ $a_i=0 \Rightarrow y_i=0$
- 3) ~~|||||~~

To get $[l^{(\ell)}, u^{(\ell)}]$ on every layer with interval arithmetic:

$$\begin{aligned} u^{(\ell)} &= [w_\ell]_+ u^{(\ell-1)} - [w_\ell]_- l^{(\ell-1)} + b_\ell & [w]_+ &= \max(w, 0) & u^{(0)} &= x_i + \epsilon \\ l^{(\ell)} &= [w_\ell]_+ l^{(\ell-1)} - [w_\ell]_- u^{(\ell-1)} + b_\ell & [w]_- &= \min(-w, 0) & l^{(0)} &= x_i - \epsilon \end{aligned}$$

$$w = \begin{bmatrix} -1 & u \\ 3 & -2 \end{bmatrix} \quad w_+ = \begin{bmatrix} 0 & u \\ 3 & 0 \end{bmatrix} \quad w_- = \begin{bmatrix} -1 & 0 \\ 0 & -2 \end{bmatrix}$$

We only need the case where $l^{(\ell)} \leq 0 \leq u^{(\ell)}$ as we don't know a_i beforehand.
Otherwise it's $a_i=1$ for $0 \leq l^{(\ell)} \leq u^{(\ell)}$ or $a_i=0$ for $l^{(\ell)} \leq u^{(\ell)} \leq 0$.

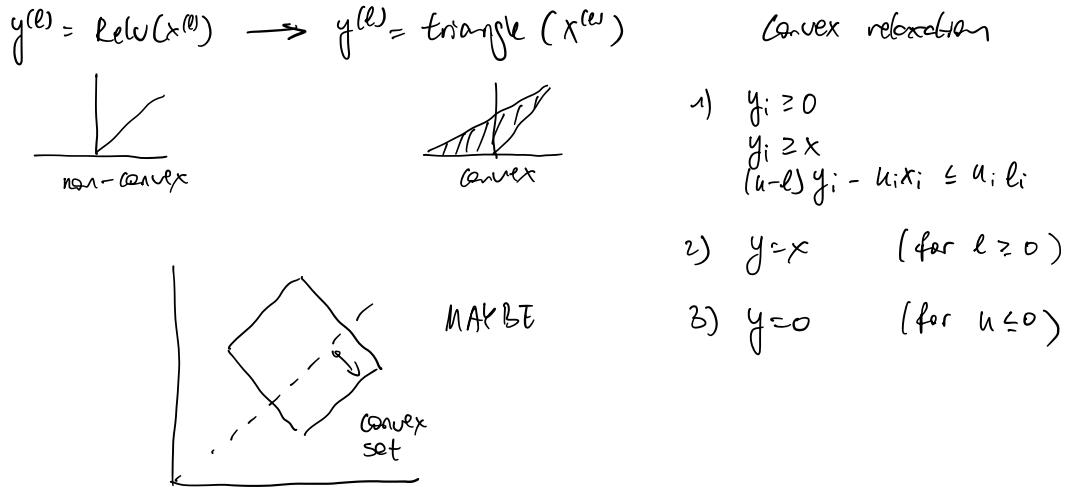
Transfer the interval, faster the convergence

Convex Relaxation

Relax the question from exact certification to allow the answer **MAYBE**. Therefore instead of solving for true M_t^* , we solve for \underline{M}_t^* which is a lower bound.

$$0 \leq \underline{M}_t^* \leq M_t^* \quad \begin{array}{ll} \text{if } \underline{M}_t^* > 0 & \rightarrow \text{Adv free} \\ \text{else} & \rightarrow \text{Maybe} \end{array}$$

ReLU has the cause of NP-complete, if we relax it, the algorithm becomes polynomial.



For the case of **MAYBE**, we actually have $\hat{x} = \arg\min$ which we can feed it into the model and see output to prove non-robustness

The tightness of $[l, u]$ has an impact on the frequency of **MAYBE**

Robust Training

M_t actually depends on the weight θ . $\rightarrow M_t(\theta)$
 We can use this for robust training

$$\ell(f_\theta(\hat{x}), y) = \max_{t \neq y} \underbrace{\max (f_\theta(\hat{x})_t - f_\theta(\hat{x})_y, 0)}_{-\underline{M}_t^*(\theta)} \quad \begin{array}{l} \text{Correctly classified} \rightarrow t=0 \\ \text{else} \rightarrow t = \underline{M}_t^*(\theta) \end{array}$$

However even with $-\underline{M}_t^*(\theta)$ we have to solve an LP and computing θ is expensive. But we can take advantage of the duality, each \hat{x} in the dual is a lower bound for the primal problem and the gradient is easily computable. $\underline{M}_t^*(\theta) = g(d')$
 (as it's just some analytical func) $h_\theta(x') \geq h(x^*) = g(d^*) \geq g(d')$.

Lipschitz Continuity

A function is said to be k -Lipschitz continuous if:

$$|f(x_1) - f(x_2)| \leq k \cdot \|x_1 - x_2\| \quad \text{for any } x_1, x_2 \in X. \quad k = \text{Lipschitz constant}$$

The change of output space (logits) is less than k times change in input space.

Lipschitz in NN

For a stacking of functions (NN), we can have an upper bound for the Lipschitz constant.

$$L = L(F) \leq \prod_{l=1}^L L(f_l) \quad \text{for functions in the stack.}$$

For a fully connected layer $f(x) = Wx + b$:

$$L(f) = \sup_{a \neq 0} \frac{\|Wa\|_p}{\|a\|_p} \quad a = x_1 - x_2 \quad (\text{easily deduced from Lipschitz formula and FC})$$

if $p=1 \rightarrow L(f) = \max_i |w_i|$ between columns of W
 $p=\infty \rightarrow \max_i \sum_j |w_{ij}|$ rows
 $p=2 \rightarrow \max_i \sqrt{\sum_j w_{ij}^2}$ singular value of W

(for conv layers it can be approximated efficiently)

Most activation functions (ReLU, sigmoid ...) have Lipschitz const of almost 1. ReLU has exactly 1.

Therefore, for the entire network:

$$L(F) = \prod_{l=1}^L \|W_l\|_p$$

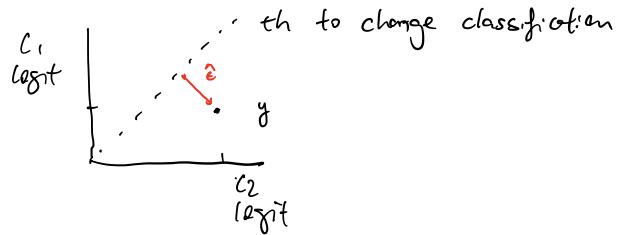
so by regularizing the norm of the weights we can enforce any desired constant

Certifying Robustness

With known Lipschitz constant of the NN.

for an input x and logit $y = f(x)$. We can compute maximum norm ϵ of perturbation δ : $\|\delta\|_p \leq \epsilon$ for which we can guarantee robustness.

i. Compute min perturbation norm $\hat{\epsilon} = \|\hat{f}\|_p$ in logit space to change the classification



2. Compute perturbation norm $\epsilon = \|\delta\|_p$ in input space based on Lipschitz constant.

$$\epsilon = \|x - \bar{x}\|_p \leq \frac{\hat{\epsilon}}{k} \Rightarrow k \cdot \|x - \bar{x}\|_p \leq \hat{\epsilon} \Rightarrow \|f(x) - f(\bar{x})\|_p \leq \hat{\epsilon}$$

Certifying that for any $\bar{x} \in P(x) = \{\bar{x} : \|\bar{x} - x\|_p \leq \frac{\hat{\epsilon}}{k}\}$ the robustness is guaranteed.

For smaller k , the network is more robust. However, it's a trade off with expressiveness.

Randomized Smoothing

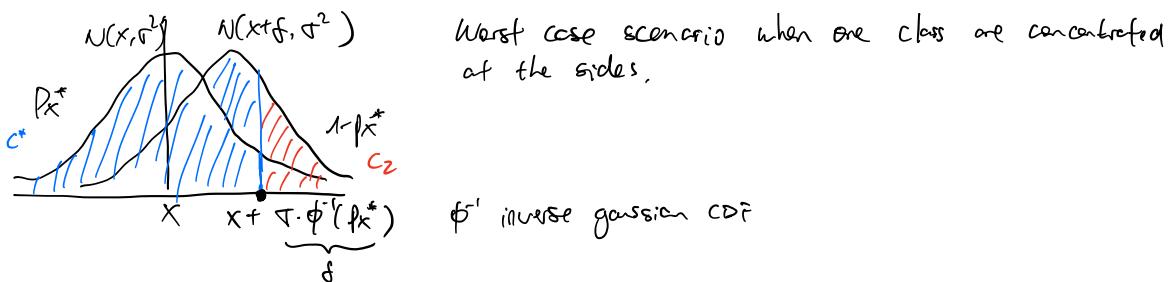
The base classifier is transformed into a smoothed one via adding gaussian noise $\delta \sim N(0, \sigma^2 I)$.

$f(x) = \arg_{c^*} \mathbb{P}(x)_c$ the class prediction of base classifier.

$g(x)_c = \mathbb{E}_\delta [\mathbb{I}(f(x+\delta)=c)] = \mathbb{P}_\delta(f(x+\delta)=c)$ logit of class c for smoothed classifier.

With $c^* = \arg_{c^*} g(x)_c$ and $p_x^* = g(x)_{c^*}$, we want to certify:

$$\arg_{c^*} g(x+\delta)_c = c^* \quad \forall \|\delta\|_2 \leq r$$



For any $\|\delta\|_2 < T \cdot \phi^{-1}(p_x^*)$ the classifier is robust (prediction still the same). Larger T leads to larger radius, but reduces p_x^* by introducing too much noise.

\hat{P}_x^*

$A = \mathbb{I}[f(x+\epsilon) = c^*] \in \{0, 1\}$ Bernoulli random variable of observing c^* with prob \hat{P}_x^*

We can compute an one-sided confidence interval $(1-\alpha)$ for \hat{P}_x^* .

If we sample n times and $A=1$ in m/n then $P(A=m) = \text{Binomial}(n, \hat{P}_x^*)$. To be conservative we take the lower bound $\underline{\hat{P}}_x^*$ from the confidence interval for the Binomial

$$\underline{\hat{P}}_x^* \quad \overbrace{\qquad}^{5\%} \quad \hat{P}_x^*$$

In practice we first use small set of samples to estimate c^* and then use large set to compute \hat{P}_x^* with confidence interval

Robust Training

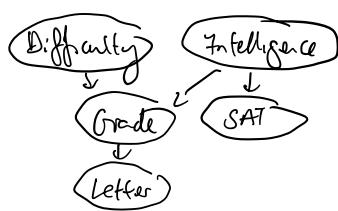
Randomized smoothing lends itself to robust training by adding data augmentation at training. The training set is augmented with random gaussian noise added to the original samples. This way, the model is more robust to noise.

Sequential Data

Bayesian Network

Natively to model sequential data $\Pr(X_1, \dots, X_T) \quad X_t \in \{1, \dots, k\}$ we would need k^{T-1} parameters.

However, with Bayesian networks we can make the model more compact with assumptions about the data. The model can be represented in a directed graph where each node is a conditional distribution.



$$\Pr(X_1, \dots, X_M) = \prod_i^M \Pr(X_i | \text{Parents}(X_i)) \quad \text{Parents} = \text{parents}$$

Then each $\Pr(X_i | \text{Parents}(X_i))$ is parameterized by θ_i .

And the overall # variables is reduced from k^{T-1} to $\sum_i \theta_i$

$$\Pr(D, I, G, S, L) = P(D) \cdot P(I) \cdot P(G|D, I) \cdot P(S|I), P(L|G)$$

This defines a generative process if we start sampling from the parent and following the graph get the joint probability

Autoregressive

$$\text{AR}(p) : X_t = c + \sum_{i=1}^p q_i X_{t-i} + \epsilon_t \quad q_i \text{ params, } c \text{ constant, } \epsilon_t \sim N(0, \sigma^2). \\ \text{The variable } X_t \text{ are non-iid}$$

$$\Pr(X_t | X_{t-1}, \dots, X_{t-p}) \sim N(c + \sum q_i X_{t-i}, \sigma^2) \quad q, c \text{ shared through time}$$

$$\mu(t) = E[X_t] \quad \text{mean} \quad \gamma(t, i) = \text{Cov}(X_t, X_{t-i}) \rightarrow p(\epsilon, i) = \frac{\text{Cov}(X_t, X_{t-i})}{\sqrt{\text{Var}(X_t)} \sqrt{\text{Var}(X_{t-i})}} \quad t \in [-1, 1] \quad \text{Pearson}$$

Stationary

If's said to be stationary if:

1. $E[X_t] = E[X_{t-i}] = \mu \quad \forall t, i$ mean is constant
2. $\text{Cov}(X_t, X_{t-i}) = \gamma_i \quad \forall t, i$ cov only depends on the lagged value
3. $E[|X_t|^2] < \infty \quad \forall t$ finite

mean and autocov can be estimated by averaging measures over time

$$E(X_t) = \mu = \frac{c}{1 - \sum_{i=1}^p \varphi_i} \quad \text{Var}(X_t) = \gamma_0 = \sum_{i=1}^p \varphi_i \gamma_i + \sigma^2 \quad \text{Cov}(X_t, X_{t-i}) = \gamma_i = \sum_{j=1}^p \varphi_j \gamma_{i-j} \quad \varphi_i = \frac{\gamma_i}{\gamma_0}$$

AR(p) is stationary iff roots of characteristic polynomial $\Phi(L) = 1 - \sum_{i=1}^p \varphi_i L^i$ lie outside the unit circle.

$$\begin{aligned} \text{for } p=1 &\rightarrow |\varphi_1| < 1 \\ p=2 &\rightarrow \begin{aligned} \varphi_1 + \varphi_2 &< 1 \\ \varphi_1 - \varphi_2 &< 1 \\ |\varphi_2| &< 1 \end{aligned} \\ &\dots \end{aligned}$$

The value doesn't explode ($\rightarrow \infty$) as the sequence goes due to the restriction on φ_i

Lecming

Opt 1. Parameters can be learned with least squares regression (normal equation):

$$\begin{bmatrix} \varphi_1 \\ \vdots \\ \varphi_p \end{bmatrix} = (X^T X)^{-1} X^T Y \quad X = \begin{bmatrix} 1 & \cdots & x_0 \\ 1 & \cdots & x_1 \\ \vdots & & \vdots \\ 1 & \cdots & x_{p-1} \end{bmatrix} \quad Y = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_p \end{bmatrix}$$

Opt 2. Use Yule-Walker equations:

1. Estimate moments $\gamma_0, \gamma_1, \dots, \gamma_p$ (Covs)
2. Inverse Yule-Walker matrix to estimate $\varphi_1, \dots, \varphi_p$
3. Use γ_0 equation to estimate σ^2

$$\begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_p \end{bmatrix} = \begin{bmatrix} \gamma_0 & \cdots & \gamma_{p-1} \\ \vdots & & \vdots \\ \gamma_{p-1} & \cdots & \gamma_0 \end{bmatrix} \begin{bmatrix} \varphi_1 \\ \vdots \\ \varphi_p \end{bmatrix}$$

Yule-Walker
matrix

Note: data splitting (test, train, val) should preserve time order

Markov Chain

Markov assumption:

$$P(X_t | X_{t-1}, \dots, X_1) = P(X_t | X_{t-1}) \quad X_t \text{ discrete } \in \{1, \dots, K\}$$

$$P(X_1 = i_1, X_2 = i_2, \dots, X_T = i_T) = P(X_1 = i_1) \prod_{t=1}^{T-1} P(X_{t+1} = i_{t+1} | X_t = i_t)$$

$$= \pi_{i_1} \cdot A_{i_1, i_2} \cdots A_{i_{T-1}, i_T} \quad (A^{(t)} = A \text{ time homogeneous assumption})$$

$T \in \mathbb{R}^{K \times K}$ prior probs $A \in \mathbb{R}^{K \times K}$ transition matrices . #params: $(K-1) + (T-1)K(K-1)$

Therefore it can be interpreted as a state machine

$$A_{ij}(n) = P(X_{t+n} = j | X_t = i) \quad \text{prob of getting from } i \text{ to } j \text{ in } n \text{ steps}$$
$$= A^n$$

$$\pi_j(t) = \Pr(X_t = j) = \pi A^{(t-1)} \quad \text{prob of getting to } j \text{ in step } t$$

Learning

We can use MLE to learn π and A (by counting).

$$A_{ij} = \frac{N_{C_{ij}}}{\sum_j N_{C_{ij}}} \quad \pi_k = \frac{L(k)}{\sum_k L(k)}$$

Hidden Markov Chain

Markov model not expressive enough:

1. Markov assumption doesn't fit real world
2. Latent variable

HMM is a sequence composed of latent variables:

$$P(X_{t+1} | Z_1, Z_2, \dots, Z_T, X_1, \dots, X_T) = P(X_{t+1} | Z_{t+1})$$
$$P(Z_{t+1} | Z_1, \dots, Z_T) = P(Z_{t+1} | Z_t)$$

$$P(Z_1 = z_1, Z_2 = z_2, \dots, Z_T = z_T, X_1 = x_1, \dots) = P(Z_1 = z_1; \pi) \prod_{t=1}^{T-1} P(Z_{t+1} = z_{t+1} | Z_t, \theta^{(t+1)}) \prod_{t=1}^T P(X_t = x_t | Z_t = z_t; \phi^{(t)})$$

$$P(Z_t = i) = \pi_i$$
$$P(Z_{t+1} = j | Z_t = i) = A_{ij}^{(t+1)}$$
$$P(X_{t+1} | Z_{t+1} = i) = B_{ij}^{(t+1)}$$
$$\# \text{params: } \pi \quad A \quad B \quad \text{As variables can share params}$$

Inference

Seek to find information about posterior $P(Z_{1:T} | X_{1:T})$. 3 types:

- Filtering: computes $P(Z_t | X_{1:t})$ incrementally as data streams (online). \rightarrow Forward
- Smoothing: " based on past and future data (offline). \rightarrow Backward
- MAP: computes $\underset{Z_{1:T}}{\operatorname{argmax}} P(Z_{1:T} | X_{1:T})$, mode of posterior. \rightarrow Viterbi

Forward Algorithm

Incrementally compute $P(Z_t | X_{1:t})$, online. With Bayes' Rule

$$P(Z_t=k | X_{1:t}) = \frac{P(Z_t=k, X_{1:t})}{\sum_j P(Z_t=j, X_{1:t})} = \frac{\alpha_t(k)}{\sum_j \alpha_t(j)} \quad \alpha_t(k) = P(Z_t=k, X_{1:t})$$

1. Compute α_0 (initialization)

$$\alpha_0(k) = P(Z_1=k) P(X_1 | Z_1=k) = B_{k,1} \text{B}_{k,1}$$

2. Given α_t , compute α_{t+1} (recursion)

$$\alpha_{t+1}(k) = P(X_{t+1} | Z_{t+1}=k) \sum_j P(Z_{t+1}=k | Z_t=j) P(Z_t=j, X_{1:t}) = B_{k,(x_{t+1})} \sum_j A_{jk} \alpha_t(j)$$

$$\alpha_{t+1} = B_{:(x_{t+1})} \odot (A^T \alpha_t)$$

Backward Algorithm

Incrementally compute $P(Z_t | X_{1:T})$, offline. With Bayes' Rule:

$$P(Z_t=k | X_{1:T}) = \frac{P(Z_t=k, X_{1:T}) P(X_{t+1:T} | Z_t=k)}{\sum_j P(Z_t=j, X_{1:T})} \propto \alpha_t(k) \beta_t(k) \quad \beta_t(k) = P(X_{t+1:T} | Z_t=k)$$

1. Compute β_T (initialization)

$$\beta_T(k) = 1$$

2. Given β_{t+1} , compute β_t (recursion)

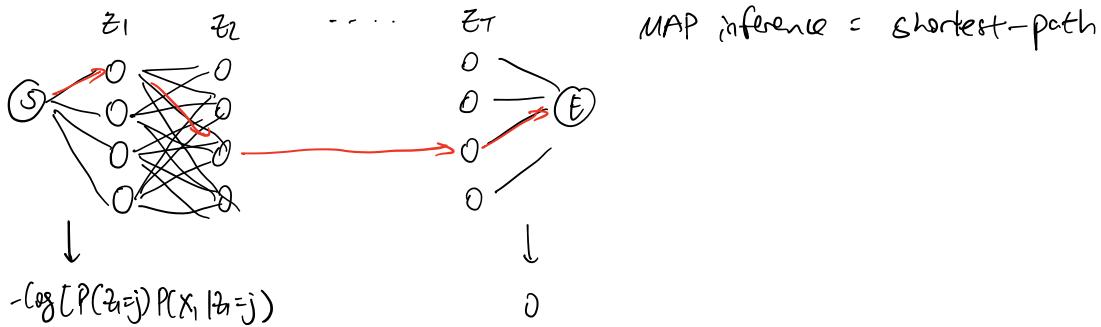
$$\beta_t(j) = \sum_k^n A_{jk} B_{k,T} \beta_{t+1}(k) \quad \beta_t = A(B_{:(x_{T+1})} \odot \beta_{T+1})$$

Viterbi

MAP inference in HMM. Given $X_{1:T}$, find Z_1, \dots, Z_T .

$$\arg \max_z P(Z_{1:T} | X_{1:T}) = \arg \max_z \log(P(z_1) P(x_1 | z_1)) + \sum \log(P(z_t | z_{t-1}) P(x_t | z_t)).$$

Each $\log[P(z_t | z_{t-1}) P(x_t | z_t)]$ depends on z_{t-1} and z_t . This can be viewed as a bipartite graph with edge weight $(i \rightarrow j) = -\log [P(z_t=j | z_{t-1}=i) P(x_t | z_t=j)]$



Learning

We want to learn $\Theta = \{\pi, A, B\}$ from observed sequence X . By solving $\max_{\Theta} \log P_{\Theta}(X)$ via variational inference (ELBO optimization).

In an EM-fashion, as we can set $q(z) = p_{\Theta}(z|X)$ for a fixed Θ .

$$\mathbb{E}_{P(z_{1:T}|x_{1:T}, \Theta^{old})} [\log P(x_{1:T}, z_{1:T}, \Theta)] = \sum_k P(z_1=k|x_{1:T}, \Theta^{old}) \log \pi_k + \\ \sum_{i,j} \sum_t P(z_t=i, z_{t+1}=j|x_{1:T}, \Theta^{old}) \log A_{ij} + \\ \sum_i \sum_t P(z_t=i|x_{1:T}, \Theta^{old}) I(x_t=j) \log B_{ij}$$

Thanks to forward-backward algorithm, we have a closed-form solution for the blue terms. Thus, Θ can easily be computed as the argmax

We don't pose mean field assumption as we want to keep dependencies introduced by the sequence

Continuous Data

X can also be continuous instead of discrete. The only difference is that the previous categorical distribution $P(x_{t+1}=x | z_{t+1}=i) = B_{ij}$ becomes gaussian $N(x | \mu_i, \Sigma_i)$

The algorithms for inference stays the same but with this change to gaussian.

For learning, we want μ_i, Σ_i instead of B_{ij} . The setting is very similar to GMM

Relation to other models:

		Z		
		NO	DISC	CONT
DISC		Markov	HMM	/
CONT	NO	AR	HMM	kalman Filter

Neural Network Approaches

Use NN to solve word sequence tasks. But need meaningful representation for words.

Word2Vec

Naïve one-hot encoding is highly dimensional, sparse and assumed independence between words.

However, we want similar words to have similar embedding vectors. To be aware of the context, a solution is to have a sliding window around the word $\{x_{i-3}, \dots, x_{i-1}, x_i, \dots, x_{i+1}, \dots, x_{i+3}\}$ and this will be the context. We can count the times 2 words appear in the same context and this defines the co-occurrence matrix (symmetric) which has the property of similar words to have similar vectors, but still high-dim. \rightarrow SVD dim-reduction.

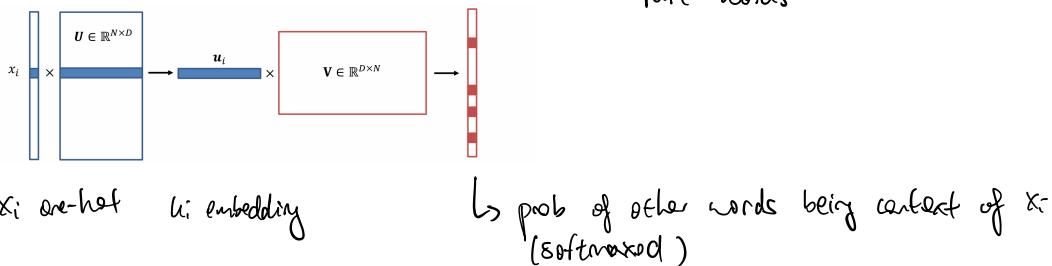
A different way to get word vector is with NN. Two approaches:

(CBOW)

predict words from its surrounding. Not good for rare words

Skip-gram

predict context from current word. Slower to train but better for rare words



x_i one-hot u_i embedding

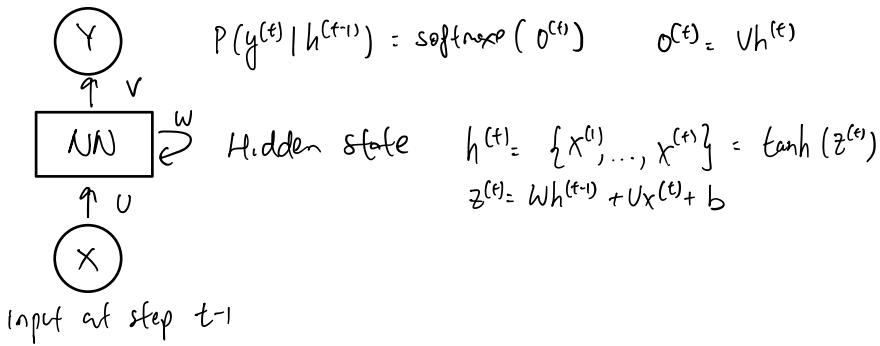
\hookrightarrow prob of other words being context of x_i (softmaxed)

Train to find optimal params U and V . But can choose $U=V$ for less params (also less expressiveness)

NN Architectures

For sequences we need layers that support input of dynamic length.

RNN



parameters $\theta = \{U, W, V\}$ train with negative log likelihood loss.

Backprop:

$$\frac{\partial L}{\partial V} = \sum (\hat{y}^{(t)} - y^{(t)}) (h^{(t)})^T \quad \frac{\partial L}{\partial W} = \sum \text{diag}(1 - h^{(t)})^2 \frac{\partial L}{\partial h^{(t)}} (h^{(t-1)})^T$$

$$\frac{\partial L}{\partial U} = \sum \text{diag}(1 - h^{(t)})^2 \frac{\partial L}{\partial h^{(t)}} X^{(t)T}$$

Can't retain large dependencies \rightarrow vanishing/exploding gradient

GRU

to alleviate vanishing gradient problem. Not all information should be fully taken into account \rightarrow gating mechanism. Adds gate (cell) that allows gradient to flow

$$\begin{aligned}
 z^{(t)} &= \sigma(W_z [h^{(t-1)}, x^{(t)}]) \\
 r^{(t)} &= \sigma(W_r [h^{(t-1)}, x^{(t)}]) \quad \text{Gates} \\
 \tilde{h}^{(t)} &= \tanh(W [r^{(t)} \odot h^{(t-1)}, x^{(t)}]) \quad \text{Simple RNN update - gives candidate state} \\
 h^{(t)} &= (1 - z^{(t)}) \odot h^{(t-1)} + z^{(t)} \odot \tilde{h}^{(t)}
 \end{aligned}$$

How much to take from previous state vs. candidate state

LSTM

Improved version of GRU.

$$\begin{aligned}
 f^{(t)} &= \sigma(W_f [h^{(t-1)}, x^{(t)}]) \quad \text{Forget gate} \\
 i^{(t)} &= \sigma(W_i [h^{(t-1)}, x^{(t)}]) \quad \text{Input gate} \\
 o^{(t)} &= \sigma(W_o [h^{(t-1)}, x^{(t)}]) \quad \text{Output gate} \\
 c^{(t)} &= f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot \tanh(W [h^{(t-1)}, x^{(t)}]) \quad \text{Simple RNN update - LSTM treats it as an input} \\
 h^{(t)} &= o^{(t)} \odot \tanh(c^{(t)})
 \end{aligned}$$

Update hidden state (now the output) using a cell state

Convolutions

Sometimes we only need the information around the current word, ie a sliding window or context. In this case, conv comes handy.

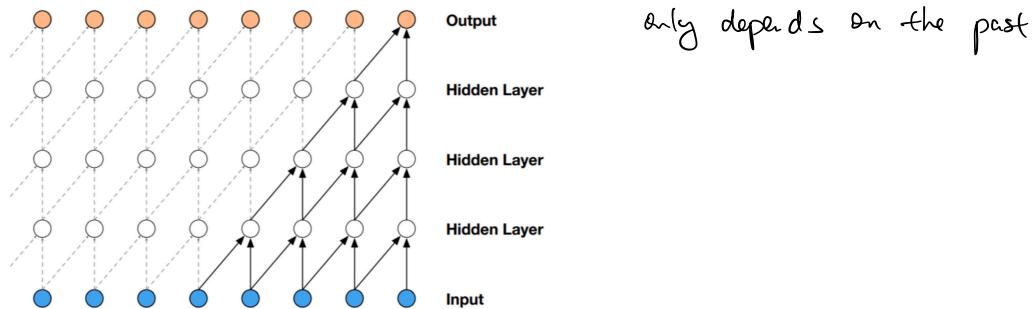
$$(f * g)(t) = \int_{-\infty}^{+\infty} f(\tau)g(t - \tau) d\tau$$

Different design of kernel (filter) have different effects: smoothing, edge detection... often this is learned by training.

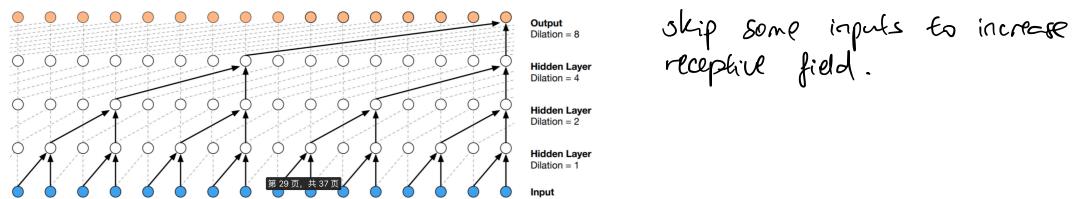
WaveNet

Architecture that uses 1-D conv to model speech. It uses spatial conv to ensure causality and increase receptive field.

Causal conv



Dilated conv



Transformer

use attention mechanism to learn where to give attention. Faster than recurrent architectures as it's parallelizable. It's a stack of autoencoder using self-attention.

The weighting is computed by softmaxed dot prod of the embeddings.

$$Z = \text{softmax} \left(\frac{Qk^T}{\sqrt{dk}} \right) V$$

weights

Q : query (current embedding)
 k : key (all embeddings)
 V : values (embedding)

↔
same dims

Self-attention: current embedding also taken into account

Positional encoding: describe position/index of a word by a vector. The result is a positional encoding matrix. It is a way that transformers use to represent order.

Temporal Point Process

We are interested in event data which are discrete events in continuous time and irregular intervals

TTP models the distribution of these events $t = \{t_1, \dots, t_N\}$ on the interval $[0, T]$ and gives the likelihood $p(t|t_1, \dots, t_N)$

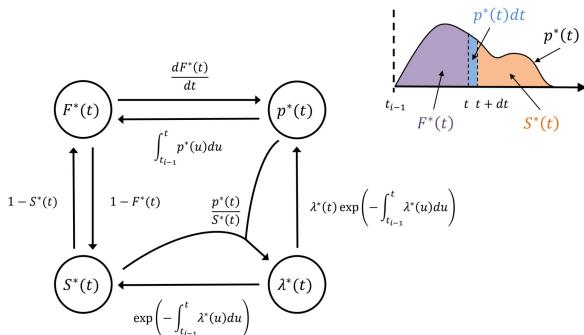
$$p^*(t) = p(t|H(t)) \quad \text{conditional density} \quad H(t) = \{t_j < t\}$$

\hookrightarrow prob of next event happening in $[t, t+dt]$ and no event $\{t_{i-1}, t\}$

$$F^*(t) = \int_{t_{i-1}}^t p^*(u) du \quad \text{CDF. Prob that next event happens in } [t_{i-1}, t)$$

$$S^*(t) = 1 - F^*(t) \quad \text{Survival func. Prob next event doesn't happen before } t$$

$$\lambda^*(t) dt = \frac{p^*(t) dt}{S^*(t)} \quad \text{conditional intensity. Prob event in } [t, t+dt] \text{ given no event in } [t_{i-1}, t). \# \text{event/time unit}$$



$$P(\{t_1, t_2, t_3\}) = p^*(t_1) p^*(t_2) p^*(t_3) S^*(T) = \lambda^*(t_1) \lambda^*(t_2) \lambda^*(t_3) \exp\left(-\int_0^T \lambda^*(u) du\right)$$

Models

Define TTP in terms of $\lambda^*(t)$:

1. Easy pre-defined behavior
2. Easy to combine different TTP
3. Efficient sampling

Homogeneous Poisson Process

Simplest model with const. intensity $\lambda^*(t) = \mu$.

$p^*(t) = \mu e^{-\mu(t-t_{i-1})}$ inter-event time follows exponential distribution

Easy to simulate HPP by sampling inter-event time $T \sim \text{Exp}(\mu)$
 $T = F^{-1}(\mu) = \frac{1}{\mu} \log(1-u)$ $u \sim \text{Unif}(0,1)$

Inhomogeneous Poisson Process

Intensity changes over time $\lambda^*(t) = g(t) \geq 0$ depend. of the history.

Number of event in $[a, b] \subseteq [0, T]$ follows Poisson distribution:

$$N([a, b]) \sim \text{Poisson}\left(\underbrace{\int_a^b g(t) dt}_{\text{expected events}}\right)$$

IIPP with different intensities can be combined and it's still IIPP

Simulating ZPP is hard, but we can:

1. Find upper bound $\mu_0 \geq g(t)$ $\forall t$
2. Simulate by HPP with μ_0
3. keep each t_i with prob $g(t_i)/\mu_0$

Hawkes Process

$$\lambda^*(t) = \mu + \alpha \sum_{t_j < t} k \omega(t - t_j) \quad \text{self-exciting (burst) event occurrence}$$

Learning

We can learn the (parametric) intensity function $\lambda^*(t)$ of the models by maximizing the log-likelihood of observed sequences.

$$\log P(\{t_1, \dots, t_N\}) = \sum_{i=1}^N \log \lambda^*(t_i) - \int_0^T \lambda^*(u) du$$

HPP has closed form solution. Hawkes / IIPP can be optimized with convex optimization or grad. descend

RNN

We could also model $p^*(t)$ directly instead of $\lambda^*(t)$ with RNN, as it gives more flexibility.

We feed the inter-event time τ_i as inputs and the hidden state h_i will keep track of the history. Then the output θ_i is a parameter for the distribution $p^*(t)$. $p^*(t) = p(t | h(t)) = p(t | h_i)$

Instead of t_i , p^* can be defined over T_i . This way we only need an non-negative distribution (Exp, Gamma...), a mixture or NF

Graphs

Graph $G = (V, E)$ V nodes $E \subseteq V \times V$ edges, can be represented by adjacency matrix $A \in \{0, 1\}^{M \times M}$ (instead of binary could also be weights)

Properties

Real-world graph are not random (Erdős-Renyi model):

N nodes, each edge has same prob p of being added.

Instead it's better explained with a power law distribution $g(x) = Ax^{-r}$ over the degree of nodes. In $\log(\text{degree})$ vs $\log(\text{rank})$ it forms a line.

Clustering coefficient

Social network has a lot of triangles.

$$\text{clustering coeff } C = \frac{3 \cdot \# \text{triangles}}{\# \text{connected triplets}}$$

Efficient triangle counting

$$\# \text{triangles} = \frac{1}{6} \text{trace}(A^3) = \frac{1}{6} \sum_i \lambda_i^3 \quad \lambda_i \text{ eigenvalues of } A$$

λ follows power law, so we only need the top few. These can be computed with power iteration:

$$\text{until convergence: } v = \frac{xv}{\|xv\|} \quad v = \text{eigenvector, initialized randomly}$$

$$\text{for the next eigenvect: } x = x - \lambda vv^T \text{ and repeat}$$

Small world phenomenon

$E \ll N^2$ $E = O(N^k)$ $k \approx 1.4$ real-world graph very sparse.
The average length of path to reach other person is 6

Characteristic path length

Starting from a node, compute min path to reach other nodes. Take avg consider path length for all starting nodes. Take median of the averages

$$\text{median} \left\{ \frac{1}{|V|} \sum_{v \in V} \text{len}(\text{Path}(v, v_j)) \right\}$$

Average diameter

$$\frac{1}{|V|} \sum_{v \in V} \frac{1}{|V|} \sum_{v \in V} \text{len}(p_{\min}(v, v_j)) \quad \text{mean instead of median}$$

Hop-plot

$N_h(u)$ the number of nodes reachable from u via h hops

$$N_h = \sum_{u \in V} N_h(u)$$

plot N_h vs h

Effective diameter (Eccentricity)

$$\min \{k \in \mathbb{N} \mid N_k \geq 0.9 |V|^2\} \quad \begin{matrix} \text{min number of hops in which some fraction (90\%)} \\ \text{of nodes (connected) can reach each other} \end{matrix}$$

Generative Models

Erdős - Renyi

$$p \in [0, 1] \quad A_{ij} \sim \text{Ber}(p)$$

Degree distribution

prob of vertex having degree $k \Rightarrow p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k} \approx \frac{z^k e^{-z}}{k!} \quad z = p(N-1)$
 follow a Poisson distribution, but real-world power law

Diameter

$\log(N) / \log(z) \quad z = \text{avg degree of nodes.}$
 Grows slowly with N , but real-world almost constant

Clustering Coefficient

$p = z/(N-1)$ connection probability.
 But real-world has community structures.

Planted Partition Model

Generalization of ER, nodes from same community have higher connectivity than outer nodes.

$$\Pr(A_{ij}=1 \mid z_i, z_j) = \begin{cases} p & z_i=z_j \\ q & z_i \neq z_j \end{cases} \quad p > q$$

Stochastic Block Model

Generalization of PPM with arbitrary num and size of communities, and varying edge densities

$$\pi = [\pi_1, \dots, \pi_k] \text{ community proportions } p(z_i=k) = \pi_k$$

$$\eta = [] \text{ edge prob between communities } P(A_{ij} | z_i, z_j) = \text{Ber}(\eta_{z_i z_j})$$

Preferential Attachment Model

Generate the network continuously. Start with a subset of nodes and grow by adding new nodes and edges. Prob connecting nodes is proportional to the current degree of nodes (rich gets richer)

Initial Attractiveness

Implementation of PAM. When a new node is added, m new edges are added each (u, v) with prob $A_{uv} = A + \text{in-deg}(v)$.

A constant for all nodes (the attractiveness)

Degree Distribution

power law with $\gamma = 2 + \frac{A}{m}$

Diameter

for $m \geq 2 \rightarrow O\left(\frac{\log N}{\log \log N}\right)$ matches small world effect but still increasing

Average degree

constant over time. But real-world increasing.

Configuration Model

pre define degree for each node (as stubs), randomly select stub pair and connect them.

Deep Models

Learn real-world graph properties from data.

NetGAN

Use GAN to learn distribution over random walks with LSTMs, initialize

states with $\sim N(0, 1\sigma)$ and generate sequence of nodes. Each step's outputs are softmaxed and fed into categorical dist
 Raden walk : $P(w_{j+1} = v_i | w_j) = \frac{1}{D} \sum_{v_i \in N(w_j)} \text{softmax}(w_i)$ to pick the node.

Reparameterization trick to make the sampling differentiable for the categorical distr.

1. Draw扁滑 noise $g_k = -\log(-\log(u))$ $u \sim U(0, 1)$
2. Sum it to $\log \pi_k$ (softmaxed LSTM output)
3. Take k with highest val. $v = \arg\max_k [g_k + \log \pi_k]$

Argmax from 3 still not differentiable. \rightarrow Replace by softmax + temperature τ

$$V \approx V^* = \tau \left(\frac{\log \pi}{\tau} \right) \quad \pi, g \text{ vectors}$$

Clustering

Find cluster by graph data instead of attributes. Vertices in same cluster are connected by many edges and few for different cluster.
 Can be categorized in overlapping or non-overlapping

But NP-hard problem if formulated as optimization over some objective function $Q(c)$, s.t. $C_1 \cup \dots \cup C_n = V$ and $C_i \cap C_j = \emptyset$ (discrete)

Spectral Clustering

Cuts

Minimize edge weights between clusters \rightarrow min cut (Ford Fulkerson (1 poly)).
 But min cut tends to cut small vertex sets \rightarrow only consider inter-cluster edges.

$$\text{Cut} : \min_{C_1, \dots, C_n} \sum_i \text{Cut}(C_i, V \setminus C_i)$$

$$\text{Ratio cut} : \min_{C_1, \dots, C_n} \sum_i \frac{\text{Cut}(C_i, V \setminus C_i)}{|C_i|}$$

$$\text{Normalized ratio cut} : \min_{C_1, \dots, C_n} \sum_i \frac{\text{Cut}(C_i, V \setminus C_i)}{\text{Vol}(C_i)} \quad \text{Vol}(C_i) = \sum_{v \in C_i} \deg(v)$$

$$\text{Cut}(C_i, C_j) = \sum_{\substack{v_i \in C_i \\ v_j \in C_j}} w(v_i, v_j)$$

Finding optimal for various cuts is NP-hard

Graph Laplacian

$L = D - W$ D = degree matrix W = adjacency (weighted) matrix

$$L_{uv} = \begin{cases} -w_{uv} & \text{if } (u, v) \in E \\ \deg(v) & \text{if } u = v \\ 0 & \text{else} \end{cases} \quad \text{for any vector } f : f^T L f = \frac{1}{2} \sum_{(u, v) \in E} w_{uv} (f_u - f_v)^2$$

- L is symmetric and PSD. With n nonnegative eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

- Laplacian is additive $L_{G \cup H} = L_G + L_H$ if G, H graphs on same vertex set
- Laplacian of graph is sum of laplacian of edges. $L_G = \sum_{\text{edges } e} L(e)$
- The number of eigenvalues = 0 correspond to the number of connected components.
- Algebraic connectivity of the graph is $\lambda_2(L)$. λ_2 big \rightarrow graph well connected

The spectrum of L encodes useful information. However, different graphs can have same spectrum (cospectral graphs)

Spectral Clustering

With $f_{C_1}[i] = \begin{cases} 1 & \text{if } v_i \in C_1 \\ -1 & \text{else} \end{cases}$ we can get min cut by minimizing $f_{C_1}^T L f_{C_1}$.
But min cut is not ideal. \rightarrow ratio cut

For $k=2$

$$f_{C_1}[i] = \begin{cases} \sqrt{\frac{|C_1|}{|V|}} & \text{if } v_i \in C_1 \\ -\sqrt{\frac{|C_1|}{|V|}} & \text{else} \end{cases}$$

1. $\sum f_{C_1}[i] = 0 \rightarrow f_{C_1} \perp \vec{1}$
 2. $f_{C_1}^T f_{C_1} = \|f_{C_1}\|_2^2 = |V| \rightarrow$ constant norm
 3. $f_{C_1}^T L f_{C_1} = |V| \cdot \text{ratio cut}$

Need to drop discreteness of f_{C_1} (relaxation) to make the optimization work.

$$\min_{f_{C_1} \in \mathbb{R}^{|V|}} \{f_{C_1}^T L f_{C_1}\} \text{ st. } f_{C_1} \perp \vec{1} \text{ and } \|f_{C_1}\| = \sqrt{|V|}$$

$\lambda_2(L)$ is always 0 with eigenvect $\vec{1}$. \rightarrow The second smallest eigenvect is f_{C_1} . And the sign of the elements in f_{C_1} tells the cluster. Alternatively, k-means could be performed on these values.

For $k \geq 2$

$$h_C[i] = \begin{cases} \frac{1}{\sqrt{|C_i|}} & \text{if } v_i \in C_i \\ 0 & \text{else} \end{cases} \quad H = [h_{C_1}; h_{C_2}; \dots] \quad h_{C_i} \text{ as columns}$$

$$1. H^T H = I \quad (\text{orthonormal}) \\ 2. h_{C_i}^T L h_{C_i} = \frac{\text{cut}(C_i, V \setminus C_i)}{|C_i|}, \quad h_{C_i}^T L h_{C_i} = (H^T L H)_{ii}$$

$$\Rightarrow \text{RatioCut}(C_1, \dots, C_k) = \text{trace}(H^T L H)$$

$\min_{H \in \mathbb{R}^{V \times k}} \text{trace}(H^T L H) \text{ st } H^T H = I$. Optimal $H =$ first K smallest eigenvect of L .

$$H = \begin{bmatrix} \uparrow & \uparrow \\ X_1 & \cdots & X_k \\ \downarrow & & \downarrow \end{bmatrix} \begin{matrix} v_1 \\ \vdots \\ v_n \end{matrix} \quad v_i \text{ doesn't count, as it's constant anyways}$$

The rows gives the embedding vector for each node. performing k-means on these gives the clustering result.

Remarks

Could change $L \rightarrow L_{\text{sym}} = D^{-1/2} L D^{-1/2}$ for normalized Laplacian.

In general spectral clustering gives good result although no guarantee to optimal due to relaxation.

For numerical data, can build similarity matrix (based on some similarity function) instead of L .

Probabilistic Approaches

Perform inference on generative models to find cluster assignment. Advantages: noise handling, capture uncertainty and generation of new data.

Planted Partition Model

$$P(A|z) = \prod_{i \in j} \underbrace{[p^{A_{ij}}(1-p)^{1-A_{ij}}]}_{\text{same cluster}}^{\mathbb{I}(z_i=z_j)} \underbrace{[q^{A_{ij}}(1-q)^{1-A_{ij}}]}_{\text{different cluster}}^{\mathbb{I}(z_i \neq z_j)}$$

likelihood funct. n

$$\text{Inference: } z^* = \underset{z \in \{0,1\}^N}{\operatorname{argmax}} P(A|z) \quad \text{MLE.}$$

No closed form solution. But if we assume $|C_1| = |C_2|$ then:

$$P(A|z) \propto \left(\frac{(1-p)^q}{(1-q)p}\right)^{\text{Ent}(z)} \quad \text{Ent}(z) = \# \text{inter-cluster edges}$$

Maximizing $P(A|z)$ \Leftrightarrow minimizing $\text{Ent}(z)$ \Leftrightarrow spectral clustering

Stochastic Block Model

Inference:

$$P(z|A, \eta, \pi) = \frac{P(A|z, \eta, \pi) P(z|\pi)}{P(A|\eta, \pi)} \quad \text{But } P(A|\eta, \pi) \text{ intractable} \rightarrow \text{variational inference.}$$

$$\approx q(z|\psi)$$

Learning (η, π) :

$$\text{If } z, A \text{ observed} \Rightarrow \text{simple counting} \quad \pi_{ik} = \frac{n_{ik}}{N} \quad \eta_{uv} = \frac{\# \text{observed edges}}{\# \text{possible edges}} = \frac{\sum_i \mathbb{I}(z_i=u) \mathbb{I}(z_i=v)}{P_{uv}}$$

$$P_{uv} = \begin{cases} \binom{n_u}{2} & u=v \\ N_u N_v & u \neq v \end{cases}$$

If only A observed: optimize ELBO w.r.t $q(z)$
and with fixed z optimize η, π

Node Embedding

Raw graph data not suitable for standard ML models \rightarrow learn function $q: V \rightarrow \mathbb{R}^d$ map discrete nodes to vectors

Applications: clustering, semi-supervised classification, link prediction.

Methods

Spectral clustering

As seen before, the matrix built with the top K eigenvectors of L gives the embedding of the nodes preserving cluster structure. Spectral embedding

Deepwalk

Train word2vec (graph $\approx NLP$) by sampling many random walks from the graph. Each random walk can be treated as a sentence in NLP to train the network. Nodes close to each other in the graph have similar embedding

Graph2Gauss

Learn mapping from nodes to a gaussian distribution $f_\theta(v) \rightarrow N(\mu_{\theta(v)}, \Sigma_{\theta(v)})$.
Loss: $\forall u \sum_{v \sim u} E_{uv} + e^{-E_{uv}} E_{uv} = KL(f_\theta(u) || f_\theta(v)) \rightarrow$ ranking loss: v is closer to u than w

Distance in embedding space relates to distance in graph, and variance tells the certainty

Graph Embedding

Aggregate all node embeddings from the graph via avg, sum... and that's the graph embedding which transforms a discrete graph to embedding.

Ranking

The web is a graph, a page is important if many important pages point to it.

Page Rank

Each node's importance is the sum of its in-votes. Each node's vote is proportional to its importance and divided by the number of votes made $r_j = \sum_{i \sim j} \frac{r_i}{d_i}$ d_i = out-degree

To get the importance of each node we can set the system of equations and solve it. To make the system unique we can force uniqueness by adding $\sum r_i = 1$.

To solve it efficiently:

- form the transition matrix M (similar to adjacency matrix but divided by out-degree).
- " rank vector r (the unknown importances)
- $r = Mr \rightarrow r$ is eigenvector of M with eigenvalue = 1 (the largest one, since M column stochastic)

Markov Chain

A random walk in the web can be considered a the sequence generated in the Markov chain

$$\text{Score of page } i = r_i = \lim_{t \rightarrow \infty} \pi_t(x_t=i) \quad \nearrow \text{transition matrix}$$

$$r = \lim_{t \rightarrow \infty} \pi_t(t) \quad \text{limiting distribution}$$

$$\pi_t(\omega) = \pi_t(\omega)B \quad \text{stationary distribution}$$

if B row stochastic $\Rightarrow \lim_{t \rightarrow \infty} \pi_t(t) = \lim_{t \rightarrow \infty} \pi_t B^{(t-1)} \Rightarrow$
at $t \rightarrow \infty$ it doesn't matter from where we start
if exists it's equal to rank vector r

If has to meet some properties:

- Irreducible: possible to get any state from any state
- Aperiodic: $\Pr(X_n=i | X_0=i) > 0 \quad \forall n \geq 0$ the chain is aperiodic if all nodes aperiodic
(able get back to the node at any time point)

real-world webs might not satisfy these. \rightarrow Allow random teleport.
With prob β follow link, with prob $1-\beta$ jump random page

$$r_j = \sum_{i: s_j} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{N} \quad // \quad A = \beta M + (1-\beta) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{N \times N} \quad \text{in matrix notation}$$

$$\Rightarrow \text{solution } r = Ar \quad (B = A^T, \quad r = \pi B^\infty)$$

Topic Specific Page Rank

Generic popularity biased against specific topics. Therefore instead of the generic popularity we can measure the popularity within a topic.

$$r = \beta Mr + (1-\beta)\pi \quad \sum \pi_i = 1 \quad . \quad \text{Now teleport } \pi_i \text{ is only restricted in a subset } S.$$

$$\text{solution: } r = (1-\beta)(I - \beta M)^{-1}\pi \quad O(N^3)$$

\hookrightarrow only change π to match a different query

We are now measuring the rank against the subset S (could be the user query for ex.)

Trust Rank

topic-sensitive page rank with S as a set of trusted pages

Semi-Supervised Learning

Label Propagation

If's the scenario where only a set of nodes have known label.
We need to assume homophily (same labels nodes are likely connected)

Binary case

$$V = S \cup U \quad S = \text{labeled set} \quad U = \text{unlabeled set}$$
$$W \in \mathbb{R}^{(M \times N)} \quad \text{weighted adjacency matrix}$$

$$\min_y E(y) = \min_y \frac{1}{2} \sum_{ij} w_{ij} (y_i - y_j)^2 \quad \hat{y} = \text{MAP inference} \quad p(y_u | W, y_L) = \frac{1}{Z} \exp(-E(y))$$

$$\Rightarrow \min_{y \in \{0,1\}^{|U|}} \underbrace{\frac{1}{2} \sum_{ij} w_{ij} (y_i - y_j)^2}_{\text{smoothness}} \quad \text{st. } y_i = \hat{y}_i \quad \forall i \quad \hat{y} = \text{gt labels}$$
$$Z = \sum_y \exp(-E(y)) \quad \text{normalizing const.}$$

$$\Rightarrow \min_{y \in \mathbb{R}^{|U|}} y^T L y \quad \text{st. } y_i = \hat{y}_i \quad \forall i \quad L = D - W$$

$$\begin{bmatrix} \hat{y}_s \\ \hat{y}_u \end{bmatrix}^T \begin{bmatrix} L_{ss} & L_{su} \\ L_{us} & L_{uu} \end{bmatrix} \begin{bmatrix} \hat{y}_s \\ \hat{y}_u \end{bmatrix} \quad \nabla_{y_u} = 2 L_{us} \hat{y}_s + 2 L_{uu} y_u \Leftrightarrow y_u = -L_{uu}^{-1} L_{us} \hat{y}_s$$

Generalization

$$k \text{ labels} \rightarrow y_u \in \mathbb{R}^k \text{ i.e. class } \kappa \quad E(y) = \sum_{i,j} w_{ij} (y_i - y_j)^T (y_i - y_j)$$

$$\text{Neural effect} \rightarrow \text{Compatibility matrix } H \quad E(y) = \sum_{i,j} w_{ij} (y_i - y_j)^T H (y_i - y_j)$$

Non-graph data → construct graph with feature vector

LP is an instance of transductive learning, where we only try to predict the label of the set U , in contrast with inductive learning where we learn function mapping.

Semi-supervised learning is a more generic principle, the idea is that we use both labeled and unlabeled data for learning (exploit connectivity as few data is labeled)

Graph Neural Network

Exploit graph structure as well as attributes by differentiable message passing

for each node:

1. Gather message from neighbors $m_v^{(u)} = \sum_{u \in N(v)} M(h_v^{(u-1)}, h_u^{(u-1)}, e_{uv})$
2. Update $h_v^{(u)} = U(h_v^{(u-1)}, m_v^{(u)})$

M, U MLPs $h_v^{(u)}$ hidden state e_{uv} edge weight

Network weights are shared across the whole graph
 Message passing steps (propagation) ≈ layers of standard network

for classification task, we perform softmax on the last layer (step) of the node $h_v^{(u)}$ (it has to be k-dim). Then the loss is just cross entropy:

$$\{w^u, b^u, p^u\}_{u \in \mathcal{U}} - \sum_{v \in \text{ccl}} \sum_{c \in \mathcal{C}} y_{vc} \log p_{vc} \quad y_{vc} = \text{one-hot gt}, \quad p_{vc} = \text{softmaxed output}$$

\hookrightarrow labeled nodes

Relation with CNN

It's similar to performing convolution in images where each pixel (node) is connected to the neighbors.

$$X^{(t+1)} = X^t + A X^t = (A + I) X^t \quad \text{Conv in matrix form}$$

$$H^{(l+1)} = \sigma \left[\underbrace{\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}}_{\text{message passing}} \underbrace{H^l W^l}_{\text{feature transform}} \right] \quad \text{graph conv}$$

Multi-graph

Setup where we have multiple graphs as input and we have one output per graph, i.e. molecule classification.

The idea is very similar to graph embedding. We get the hidden state of all nodes $H_{G_i} = [h_1^{(u)}, \dots, h_M^{(u)}]$ and aggregate via nlp $R(H_{G_i})$, the simplest case could be avg/sum.

Limitations of GNN

Expressiveness, Robustness (adversarial manipulation), scalability (mpn expensive)

Overview

Isomorphism Problem

Best algorithm to determine if structure is the same for two graphs is NP.

Weisfeiler-Lehman Test: 1. Label all nodes \downarrow
 (WL)
 2. Collect neighboring labels (concat)
 3. Relabel each node by hashing collected labels
 4. Compare result of each node.

If not 4 → not isomorphic
 Else → might be or not

GNN is known to not solve the isomorphism problem. In fact, it's less powerful than WL (case of mean, max aggregation func). With more complex func it could get same power as WL.

Oversmoothing

Influence between the nodes is proportional to prob of visiting from a random walk of length k ($k = \text{max length}$). The more central \rightarrow more influence \rightarrow smoother

If $k \rightarrow \infty$: the influence becomes prop. to stationary distribution in MC \rightarrow Deep GNN unable to preserve local neighborhood info.

$$r^{(t+1)} = A D^{-1} r^{(t)} \quad H^{(t+1)} = \nabla(\cdot) \quad \text{if we push the non-linearity one layer up, the similarity becomes higher.}$$

Page-rank GNN

To preserve neighborhood \rightarrow as in pr, add teleport mechanism. \rightarrow PPND

Personalized Propagation of Neural Network

1. separate transformation and propagation.

$$H_{i,:}^{(0)} = f_0(x_{i,:}) \in \mathbb{R}^k \quad H^{(0)} \in \mathbb{R}^{N \times k}$$

2. Introduce personalized teleportation

$$R = \alpha \underbrace{(I - (1-\alpha) \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}})^{-1}}_{\in \mathbb{R}^{N \times N} \text{ weightings}} \underbrace{H^{(0)}}_{\in \mathbb{R}^{N \times k} \text{ embeddings}} \quad \text{exact sol.}$$

$$3. z = \text{softmax}(R)$$

Approximate exact sol $z \approx H^{(k)}$. Otherwise it's not scalable.

Now $\alpha H^{(0)}$ is fixed, we only need to iteratively compute $(1-\alpha)(\cdot)$ term, as it's sparse this is doable, propagation is done one $\mathbb{R}^{N \times k}$ matrix instead of high dim hidden state.

$$H^{(t+1)} = (1-\alpha) \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^t + \alpha H^{(0)}$$

Robustness

Adversarial attack now can happen both on node attributes or on connectivity (structural attacks)

$$X \quad A$$

Adversarial Attack

Perturbations measured via non-convex L_0 -norm.

$(A', x') \approx (A, x)$ hard to find measure of perturbation that do not change semantics

Unlabeled data also used for training \rightarrow poisoning attack.

Nefattack

$$\text{Linearize } z = \text{softmax}(\hat{A}x(\hat{A}xw^{(1)})w^{(2)}) \rightarrow \hat{A}^2xw'$$

for:
 1. Structure perturbation $\nabla_{\hat{A}} L(\hat{A}^2C)$ pick optimal
 2. Feature perturbation $\nabla_{\hat{x}} L(C, xC_2)$ perturbing at each step

Poisoning attack 1% correctness from 90% with clean data.

Certificates

Convex relaxation as before \rightarrow answer YES / MAYBE

Certify attribute perturbations via LP

Certify structure perturbations via jointly constraint bilinear program

Certify PPNP model (structure perturbation) via quadratically constrained linear program.

Polynomial if local budget ($\max x$ perturbations per node) but NP if global budget ($\max x$ perturbations overall)

Randomized Smoothing

We want to Certify the smooth classifier: $\|x - x'\|_0 \leq r : c^* = \arg \max_c g(x)_c$
 \rightarrow L_0 -norm radius certification.

But how to add noise to adjacency matrix? $\rightarrow n^2$ Bernoulli random variables

$$P(\tilde{x}_i | x) = \begin{cases} p & \tilde{x}_i = 1-x \text{ flip} \\ 1-p & \tilde{x}_i = x \text{ not flip} \end{cases}$$

But still problems: real-world graphs are sparse ($m \ll n^2$) therefore picking right p almost impossible.

\rightarrow Sparsity-aware random sampling: each elem of A modelled by Bernoulli with data dependent probability.

$$P(\tilde{x}_i | x) = \begin{cases} p_i x_i p_a^{1-x_i} & \tilde{x}_i = 1-x \\ (1-p_i)x_i(1-p_a)^{1-x_i} & \tilde{x}_i = x \end{cases} \quad \text{If } x_i = 0 \text{ flip with prob } p_a \text{ else with prob } p_b.$$

$$\text{Now } g(x)_c = (\dots) \stackrel{\text{defn}}{=} \sum_{\substack{\tilde{x} \\ \text{s.t. } f(\tilde{x})=c}} P(\tilde{x}|x) \stackrel{\text{iid}}{=} \sum_{\substack{\tilde{x} \\ \dots}} \prod_{i=1}^{n^2} P(\tilde{x}_i|x) \quad \text{sum of probs. of } \tilde{x} \text{ that leads to class } c.$$

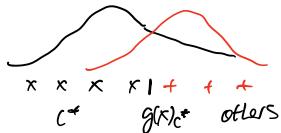
To find the radius we need to calculate the worst-case base classifier.
 This can be done with linear programming:

$$\min_{\mathbf{h}} \sum_i h_i P(\bar{x}^{(i)} | x) \quad \text{s.t. } \sum_i h_i P(\bar{x}^{(i)} | x) = g(x)_c^* \quad \text{and} \quad 0 \leq h_i \leq 1$$

\mathbf{h} is the worst-case classifier and h_i indicates whether $h(\bar{x}^{(i)})$ outputs c^* , $g(x)_c^*$ the budget. \mathbf{h} is in the set H which produce the lowest prob.

The LP can be solved optimally by greedy approach:

1. Initialize all $h_i = 0$
2. Compute likelihood ratios $\eta_i = P(x^i | x) / P(\bar{x}^i | x)$ and sort \rightarrow indices j_1, j_2, \dots s.t. $\eta_{j_1} \geq \eta_{j_2} \geq \dots$
3. Set $h_{j_1} = 1$ while $g(x)_c^*$ not used up.



Now we are in the same situation as gaussian variables.

Technically we don't need to compute all h (as there are infinite), instead we can group them by some η_i .