

Modeling tasks & losses

a) Choosing physically meaningful metrics

- Task:
- dataset
 - functional space (architecture)
 - optimization criterion \leftarrow modeling the errors

regression
 $\| \hat{y} - y \|^2$
 prediction error
 ground truth $\in \mathbb{R}^d$

classification
 $\hat{y} = \begin{pmatrix} \hat{p}_1 \\ \hat{p}_2 \\ \vdots \end{pmatrix}$
 $y = \begin{pmatrix} p_1 \\ p_2 \\ \vdots \end{pmatrix}$
 $\sum \hat{p}_i = 1$
 $\hat{y} =$ regression on probabilities

$\rightarrow L^2$ norm? meaningful

ex: $p_1 = 0,75$ 73%
 $\pm 0,01$ 72-74%

$p'_1 = 0,001$ 0% - 1,1% not possible \rightarrow 1 chance over 1000
 $\pm 0,01$ \rightarrow 1 chance over 100

Criteria to compare probabilities

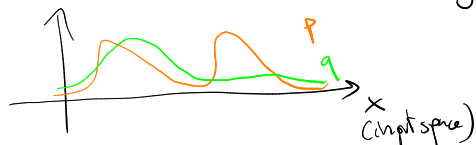
$$|\hat{p}_2 - p_1| = 0,01$$

Information theory $\Rightarrow -\log p \geq 0$ \leftarrow output of the network
 $-\log \hat{p}_2 - (-\log p_2) = \left| \log \frac{p_1}{\hat{p}_1} \right|$ softmax

Kullback-Leibler: $KL(p||q) = \int p \log \frac{p}{q} \geq 0$
 $= \mathbb{E} \left[\log \frac{p(x)}{q(x)} \right]$
 $= 0$ only if $p=q$

not symmetric: $KL(p||q) \neq KL(q||p)$

$$KL(q||p) = \int q \log \frac{q}{p} = \mathbb{E} \left[-\log \frac{p}{q} \right]$$



\hookrightarrow make drop: $KL(q||p)$

$$KL(p||q) = \int p \log \frac{p}{q} = \int p \log p - \int p \log q$$

$-\underbrace{H(p)}_{\text{entropy of } p}$ $\underbrace{\int p \log q}_{\text{cross-entropy}}$

p : target distribution
 q_θ : output "

$$\inf_{\theta} KL(p||q_\theta)$$

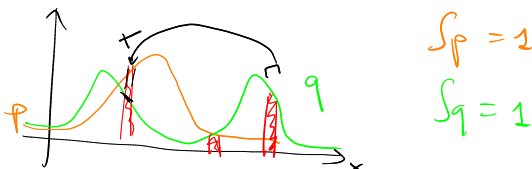
$$\theta_{t+1} = \theta_t - \epsilon \nabla_{\theta} KL(p||q_\theta)$$

easy to implement!

$\sum_{i \in \text{minibatch}}$ - network output before softmax (x_i)

other issue!
 if at same point x ,
 $p(x) \rightarrow$ not $q(x)$ to
 $\Rightarrow KL(q||p)$ is too
 \hookrightarrow careful about bin discretization

Optimal transport aka Earth Mover Distance



$S_p = 1$
 $S_q = 1$

$\inf_M \text{Total Labor}(M, c) \rightarrow$ optimal transport

Transportation Matrix

quantity of earth moved from $q(x_2)$ to $p(x_2)$

M

$\text{proj } M / \text{rows} = q$
 $\text{proj } M / \text{cols} = p$
 $M \geq 0$

Cost: $c(x_1, x_2)$

$$\text{ex: } = \|x_1 - x_2\|$$

$$\text{Total Labor}(M, c) = \int_{x_1, x_2} M_{x_1, x_2} c(x_1, x_2) dx_1 dx_2$$

how much earth I move from location x_1 to location x_2
 cost of this move (length of transport)

- in dim 1: exact solution in linear complexity (with cumulative distrib.)
- in dim > 2: hard problem
 - ⇒ approximate solution: Sinkhorn algorithm [Marco Cuturi, 2013]
 - + Gabriel Peyré 2019 ...

bistochastic
(set of samples with same mass)

$$\text{Total cost} + \epsilon H(M)$$

↑ entropy
small value

⇒ optimal solution of the form:

$$M = \begin{pmatrix} \sqrt{\kappa} \\ \sqrt{\kappa} \end{pmatrix} \begin{matrix} \left[\frac{-c_{ij}}{\epsilon} \right] \\ \left[\frac{-c_{ij}}{\epsilon} \right] \end{matrix} \begin{pmatrix} \sqrt{\kappa} \\ \sqrt{\kappa} \end{pmatrix}$$

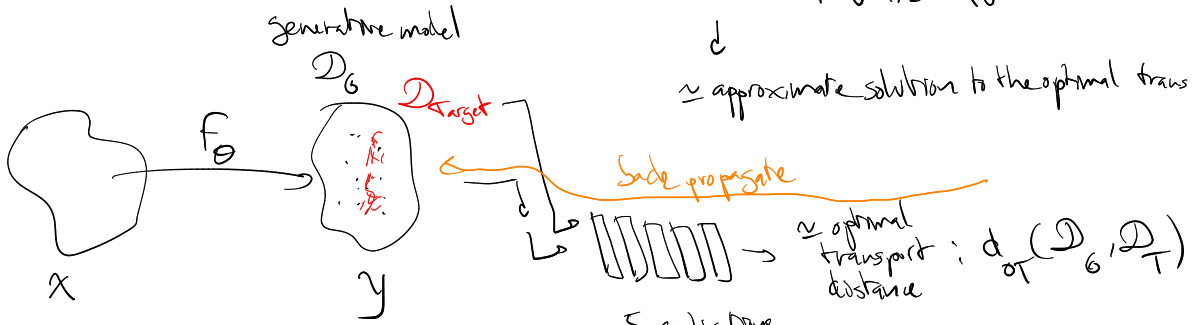
↑ fixed

initialize u, v with $(1, 1, \dots, 1)$

iterative scheme $\left\{ \begin{array}{l} u \cdot / = \kappa v \\ v \cdot / = \kappa u \end{array} \right.$

↓

≈ approximate solution to the optimal transport pb

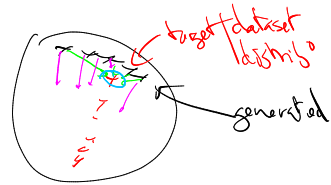


5 applications of the process

neural network with fixed weights differentiable (simple)

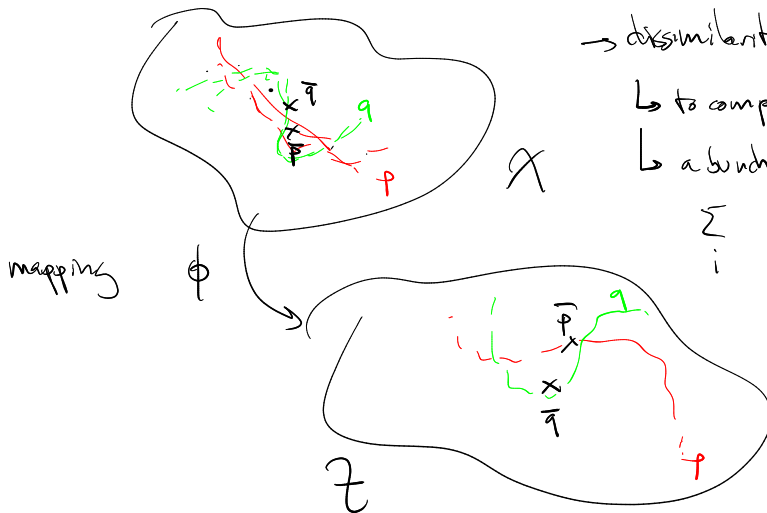
use case:

- if classification: cross-entropy is fine
- general distribution to be matched to a target one (continuous space) & metric in that space → O.T.



Minimum Mean Discrepancy (MMD)

to compare 2 distributions p, q



→ dissimilarity between p & q ?

↳ to compare their means: $\|\bar{p} - \bar{q}\|$

↳ a bunch of simple descriptors of distrib^o: F^i

$$\sum_i \|F^i_p - F^i_q\|^2$$

much higher-dim

$F^i = \phi^i$
means matching
⇒ close match

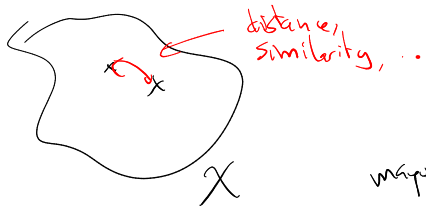
$\phi^i =$ density around point x_i

do not compute all ϕ^i explicitly

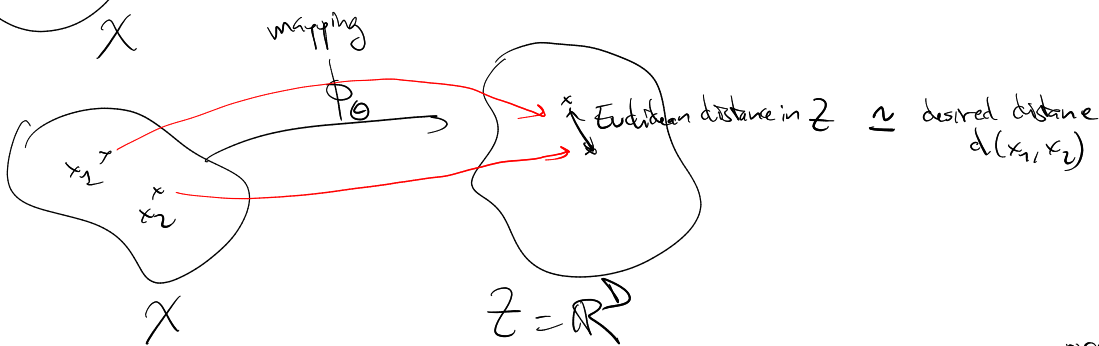
↳ use a kernel $k(p, q)$

[Gretton]

2) Modeling the task: metric learning



[Siamese networks, W+R 1554]



Dataset: triplets (x, p, n)

example of interest \rightarrow "positive": sample "close to" x

"negative": sample "far away" from x

desired property:

$$d(x, p) \leq d(x, n) - m$$

predefined margin $(\in \mathbb{R}^+)$

Criterion: $\sum_{(x, p, n) \in \text{dataset}} \max(d(x, p) - d(x, n) + m, 0)$

"triplet loss"

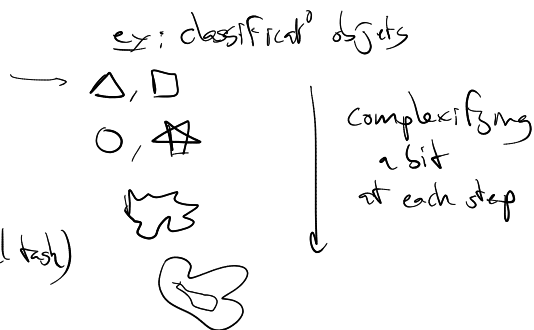
\rightarrow if negative: do nothing

\rightarrow if ≥ 0 : do smthg

Curriculum learning: from easy task to harder task [Bengio et al, 2005]

- hard ML task

\Rightarrow series of tasks, $T_1 = \text{easy}$
 $T_2 = \text{harder}$
 \vdots
 $T_N = \text{really hard (real task)}$



- change: the distribution of samples

real: $p(x)$

easy \rightarrow less easy \rightarrow difficult

at time t :

$$q_t(x) = w_t(x) p(x)$$

at $t \rightarrow 0$: $w_0(x) = 0$ if difficult x
 1 if easy

with time: $w_t(x) \uparrow$

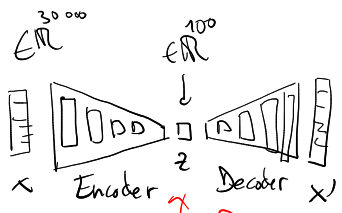
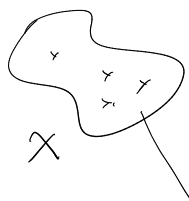
at final step $t=N$: $w_f(x) = 1$

Entropy $H(q_t) \rightarrow$ increasing with time

Generative Models



1) Auto-encoders



@ optimal criterion:

$$\|x' - x\|$$

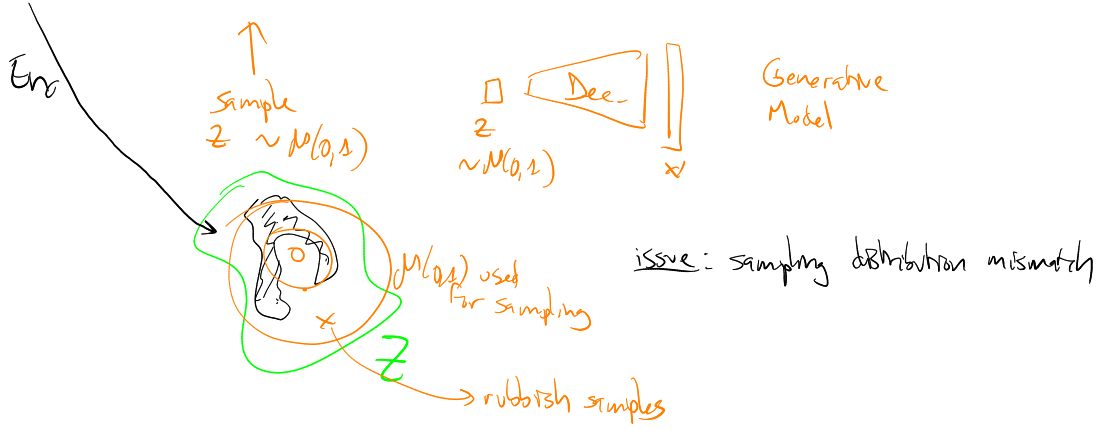
\hookrightarrow best: identity

reconstruction error

\rightarrow issue: which metric?

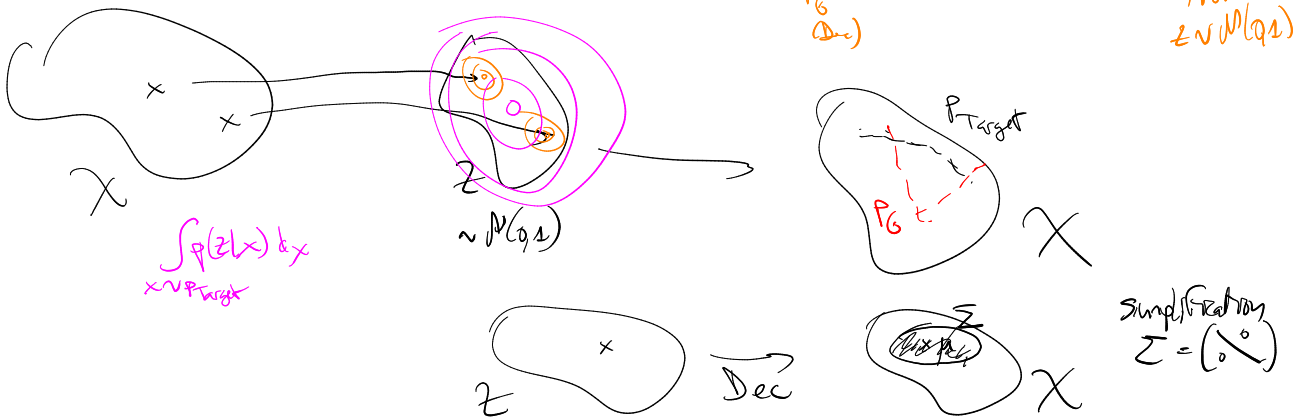
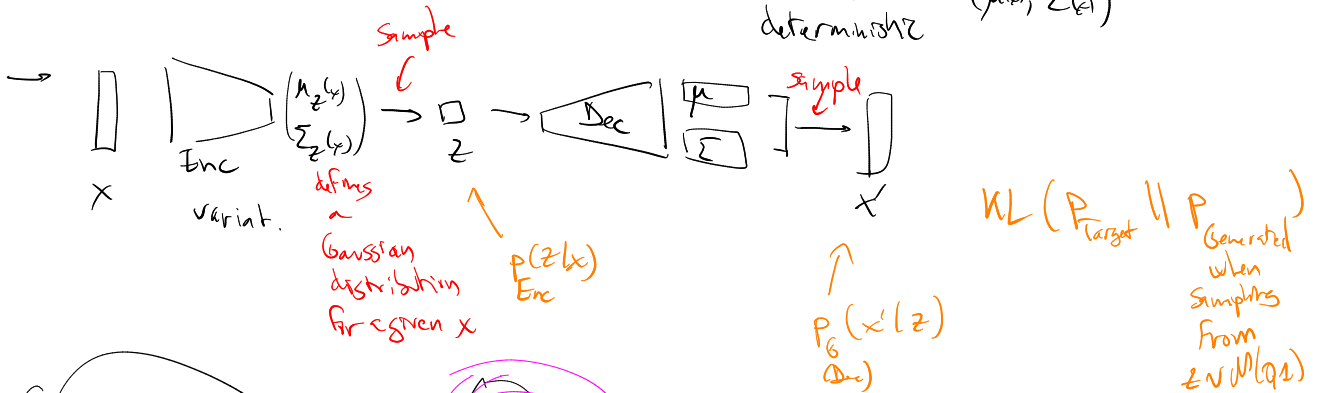
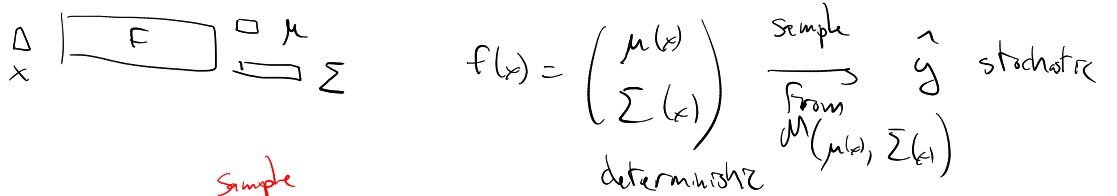
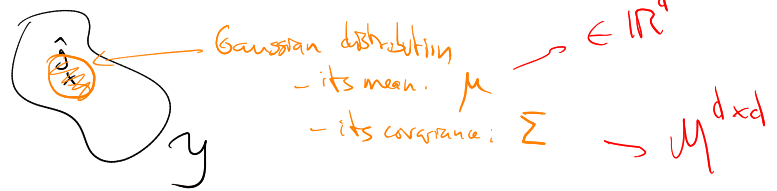
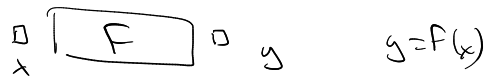
(or VGG activities: perceptual loss)

$x \rightarrow$ VGG
 $x' \rightarrow$ VGG



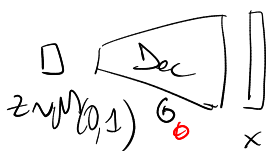
Variational Auto-Encoders (VAE)

- variational setup:
- networks not deterministic, stochastic
 - output: distribution



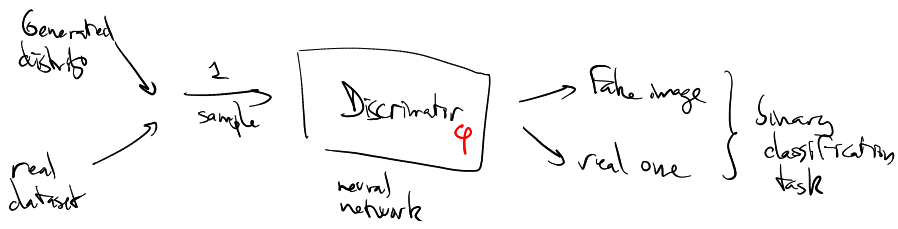
2) Adversarial settings: GAN: Generative Adversarial Networks

↳ DANN: Domain-adversarial NN



→ quantify how good outputs are





standard task
 ↓
 train ϕ

Goal: generate samples (with G_θ) that no-one can distinguish from real samples



same loss for G_θ & D_ϕ
 set with opposite sign

→ no convergence guarantee

get if converges, Nash equilibrium
 ⇒ $P_G = P_T$

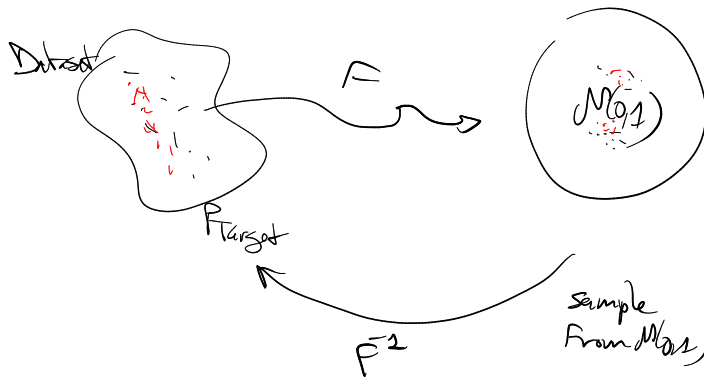
↳ Kolmogorov test to distinguish 2 distributions

Extensions: Wasserstein GAN (optimal transport)

⇒ same setup, with one constraint on DISCR.

$$W(p, q) = \sup_{\text{Lip } F \leq 1} \mathbb{E}[F(x)] - \mathbb{E}[F(y)]$$

3) Normalizing Flows



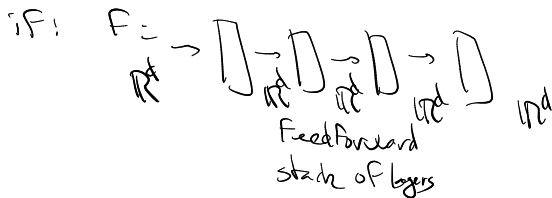
transforming P_{target} into $\mathcal{N}(0,1)$ normal distrib

Constraints

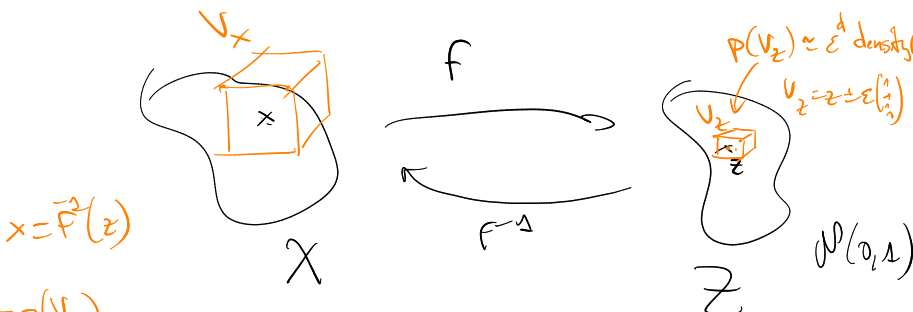
- F needs to be invertible (and computationally easy)
- need to compute fast det J J: Jacobian of
- some dimensions in input/output space

$$J_{(r)} = \frac{\partial F}{\partial x}$$

← huge matrix



⇒ same number of neurons in every layer = dim input space



$$x = F^{-1}(z)$$

$$p(x) = p(z)$$

density $p(x) = p(z) / \text{volume change}$

compute exactly $p_G(x) \neq x$

$$KL(P_{\text{target}} \parallel P_{\text{Generated}})$$

$$KL(p \parallel q) = \int p \log \frac{p}{q}$$

$$p(z) dz = p(z)$$

$P_{\theta}(x) = P_{\theta_0}(z) \times \text{volume ratio}$ $\rightarrow \det \left| \frac{\partial F^{-1}}{\partial z} \right| = \det \left(\frac{\partial F}{\partial x} \right)^{-1} = \frac{1}{\det \frac{\partial F}{\partial x}}$

$\log P_{\theta}(x) = \log P_{\theta_0}(z) - \log \det J_{(F)}$

Example of architectures:



! not expressive at all

$y^i = \lambda^i x^i$

$\lambda^i \neq 0$

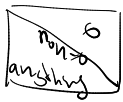
$y = \Lambda x$

$\Lambda^{-1} = \begin{pmatrix} 1/\lambda_1 & & 0 \\ & \ddots & \\ 0 & & 1/\lambda_n \end{pmatrix}$

$\frac{dy}{dx} = \Lambda \quad \det J = \det \Lambda = \prod \lambda_i$

! activation function: invertible! Leaky ReLU

Triangular matrices:



"auto-regressive"

$y = \Lambda x$

$\det J = \prod \lambda_i$

inverse by Gauss method

$\begin{cases} y_1 = \lambda_1 x_1 \Rightarrow x_1 = \frac{1}{\lambda_1} y_1 \\ y_2 = \lambda_2 x_2 + m_{12} x_1 \Rightarrow x_2 = \frac{1}{\lambda_2} (y_2 - m_{12} x_1) \end{cases}$

iterative process $\rightarrow O(N^2)$

Real-NVP: variation on this



triangular Jacobian

split x into z : $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1 \\ x_2 + f(x_1) + g(x_2) \end{pmatrix}$

easy to inverse + triangular Jacobian



diag: F
stoch: g

Orthogonal matrices

$\Rightarrow \det = \pm 1 \rightarrow$ how to train?

becomes $O = S_1 S_2 S_3 \dots S_N$

\downarrow
 $Id - \frac{v_i v_i^T}{\|v_i\|^2}$

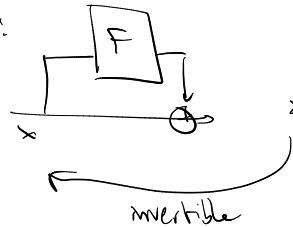
(Householder decomposition)

any matrix can be decomposed

as Orthogonal \times $\in \mathbb{R}$

\Rightarrow full expressivity

ResNet: each block:



$x + F(x)$

constraint $\|F\| \ll 1$

$Id + \text{small}$ is invertible

KL loss: $KL(P_{\text{Target}} \| P_{\theta})$ \leftarrow generated distrib^o (over X)

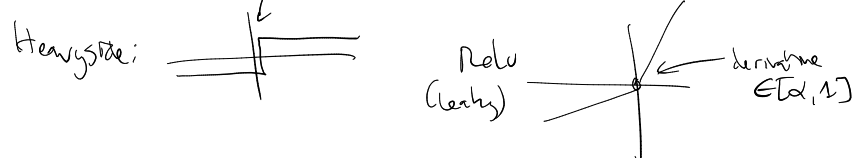
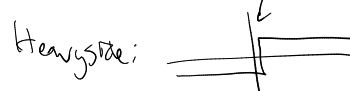
$= \int_{x \in X} P_{\text{Target}}(x) \log \frac{P_{\text{Target}}(x)}{P_{\theta}(x)} dx$

$= \text{const} - \int_{x \in X} P_{\text{Target}}(x) \log P_{\theta}(x) dx$

$\mathbb{E}[\log P_{\theta}(x)]$ $x \sim P_{\text{Target}}$ \rightarrow simple minibatch of real examples

approximate $\mathbb{E}_{P_{\text{Target}}}$ by $\frac{1}{N} \sum_i \text{ minibatch}$

\downarrow
 $-\frac{1}{N} \sum_i \log P_{\theta}(x_i)$



Stacking layers: if $F = F_3 \circ F_2 \circ F_1$:

and associated Jacobians $J_3 = \frac{dF_3(z)}{dz}$ $J_2 = \frac{dF_2(z)}{dz} \dots$

$$J_F = J_3 \times J_2 \times J_1$$

$$\det J_F = \det J_3 \times \det J_2 \times \det J_1$$

$$\log \det J_F = \sum_{\text{layers } l} \log \det J_l$$

⚠ do not forget activation functions

⚠ Check the expressivity of your architecture

ex: stacking only lower-triangular matrices: $\begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix} \times \begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix} \times \begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix}$ is a lower-triangular matrix \Rightarrow no full expressivity

↳ stacking lower-upper-triangular matrices alternately:

$\begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix} \times \begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix} \times \begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix} \times \begin{bmatrix} \times & & \\ & \times & \\ & & \times \end{bmatrix}$ mixes better the variables \rightarrow Full matrix

↳ Theorem: full expressivity of PLU decomposition: can represent any matrix
 permutation \rightarrow

but: P not optimizable by gradient descent

\Rightarrow better use QR decomposition (and stack), and add non-linear layers (Leaky ReLU, e.g.)

* Full expressivity proven for some non-linear blocks (real-NVP, iResNet...)
 upon large hidden layers / many blocks enough assumption

$\begin{pmatrix} \times & & \\ & \times & \\ & & \times \end{pmatrix}$
 diagonal (= coefficient-wise) function

det Jacobian (Leaky ReLU) = $\prod x_i$
 $= \underbrace{\times}_{\# \text{activations} < 0} \times \underbrace{1}_{\# \text{activations} > 0}$

$\log \det J_{\text{Leaky ReLU Layer}} = (\# \text{activations}) \times \log x$