# Računalni vid

#### uvod u generativno modeliranje

Siniša Šegvić UniZg-FER D307

# Agenda

Part 1: recent advances in generative recognition

- □ algorithms for generating complex data
- applications of generative models

Part 2: energy-based models and normalizing flows

- energy-based models
- generating complex data with bijective models
- affine coupling and other bijective operations

Part 3: generative recognition for dense anomaly detection

- generating synthetic training data
- finding anomalies according to local density

#### INTRO: RECOGNITION

#### **Discriminative recognition** recovers P(Y = y | X = x)

□ default flavour of machine learning

established technology, exciting applications

Generative recognition recovers p(X = x, Y = y)

- □ or, equally interesting, p(X = x | Y = y) or p(X = x)
- called *generative* since sampling from them generates **synthetic data** in the input space
- somewhat eclipsed by the success of discriminative approaches
- rapidly developing, exciting applications



[unizg-fer-dl1]



[delic21sem]

#### INTRO: GENERATIVE VS DISCRIMINATIVE

Both approaches produce **probabilistic** output in each input datum, however:

- □ discriminative models produce distributions over targets
- generative models produce distributions over inputs
  - □ how to normalize the distribution, i.e. ensure that it sums (integrates) to 1?

generative recognition is tough!



# INTRO: GENERATIVE VS DISCRIMINATIVE (2)

Intuitively, it is easier to tell the painter than to actually do the painting:



[public domain]

It is easier to distinguish composers than to develop musical ideas.

In words of the infamous food critic Anton Ego:

... in the grand scheme of things, the average piece of (bad food) is probably more meaningful than our criticism designating it so.

#### INTRO: GENERATIVE VS DISCRIMINATIVE (3)

Things get especially tough when the data is complex:

- □ language: 15-20 words per sentence, 170000 total words (English)
- □ vision: at least 3×64×64 components, 256 values per component
- □ generative recognition has to consider integrals over thousands of dimensions in order to normalize  $p(\mathbf{x})$
- Close encounters with the curse of dimensionality!



#### INTRO: TASKS

Generative approaches aim at density of the training data

Three main tasks of a generative model:

- $\Box$  generate synthetic data points by sampling  $p(\mathbf{x})$ 
  - trade off: quality vs coverage
  - useful for content creation and improving discriminative models
- □ transform data-points to a factorized latent representation
  - useful for content editing
- evaluate density  $p(\mathbf{x})$ 
  - only generative models with explicit density can do that
  - □ useful for anomaly detection, compression

#### INTRO: APPROACHES

Generative algorithms come in many flavours and many ways to appreciate them:

Algorithm	latent	bias	sampling	density
Energy	unable	coverage	slow	unnormalized
VAE	easy	coverage	fast	ELBO
Autoreg.	unable	coverage	slow	exact
GAN	unable	quality	fast	unable
NFlow	easy	coverage	fast	exact
Diffusion	easy	(coverage)	slow	ELBO
Score	unable	coverage	slow	(unable)



#### INTRO: DENSITY ESTIMATION

Comparison of generative algorithms with respect to density estimation:

algorithm / density formulation	tractability	efficiency
$p_{EBM}(\mathbf{x}) = e^{-E(\mathbf{x})} / \int_{\mathbf{x}} e^{-E(\mathbf{x})}$	only inference	fast
$m{ ho}_{V\!AE}(\mathbf{x}) = \int_{\mathbf{z}} m{ ho}(\mathbf{x} \mathbf{z}) m{ ho}(\mathbf{z}) d\mathbf{z}$	intractable (ELBO)	fast
$oldsymbol{p}_{\mathrm{ar}}(\mathbf{x}) = oldsymbol{p}(x_1) \prod_{i=2}^{HW} oldsymbol{p}(x_i   \mathbf{x}_{< i})$	tractable	slow
$p_{GAN}(\mathbf{x}) = ?$ (implicit density)	unavailable	
$m{ ho}_{\mathit{flow}}(\mathbf{x}) = m{ ho}_{\mathit{z}}(\mathbf{f}(\mathbf{x})) \cdot  \det(rac{\partial \mathbf{f}}{\partial \mathbf{x}}) $	tractable	fast
$m{ ho}_{diff}(\mathbf{x}) = \int m{ ho}(\mathbf{x} \mathbf{x}^{(1)})$		
$\prod_{t=1}^{T-1} p(\mathbf{x}^{(t)} \mathbf{x}^{(t+1)}) \pi(\mathbf{x}^T) d\mathbf{x}^{(1T)}$	intractable (ELBO)	slow



#### **APPLICATIONS: CONTENT CREATION**

Most generative algorithms are able to generate data.

However, some optimize for coverage while others optimize for quality [lucas19neurips].





Normalizing flows - high coverage (left); GANs - high quality (right).



[grcic21neurips]



[sehwag22cvpr] GMCV → Applications 10/59

#### **APPLICATIONS: BOOSTING DIVERSITY**

Recent work leverages computational power to boost diversity without sacrificing quality:



[sehwag22cvpr]

Other recent work allows to favour either quality (left) or diversity (right):



[humayun22cvpr]

#### **APPLICATIONS: CONDITIONAL GENERATION**

For practical pupposes, we are mostly interested in conditional generation

A popular recent approach connects language embeddings with generative vision.

This is what I got by feeding "a photo of a white cat on a unicycle" to DALL-e:





# APPLICATIONS: CONDITIONAL GENERATION (2)

It also works the other way round (from images to text):



GT: A young boy in the park throwing a frisbee.

L-Verse: A young boy throwing a green frisbee in a lush green park.



GT: A laptop and a cell phone on a table.

L-Verse: A collection of electronic devices and cords sitting on top of a shower curtain over the bathtub table.



GT: A small bathroom is shown from a door.

L-Verse: A bathroom with a next to a toilet.

[kim22cvpr]

# Applications: conditional generation (3)

One can also perform arithmetic operations on visual semantics and display results in text:





A cow's milk.

[tewel22cvpr]

#### **APPLICATIONS: CONTENT EDITING**

#### Conditional generation, visual replace, extrapolation by editing the VQ VAE latent:

(a) Class-conditional Image Generation





– Flamingo –

(b) Image Manipulation



(c) Image Extrapolation





<sup>[</sup>chang22cvpr]

# APPLICATIONS: CONTENT EDITING (2)

Another instance of visual replace:



# APPLICATIONS: CONTENT EDITING (3)

#### Super resolution by leveraging Style GAN latent:



(f) GT [zhong22cvpr]

# Applications: content editing (4)

Colorization and other inverse problems (inpainting, medical image reconstruction):



[song22blog]

#### **APPLICATIONS: ANOMALY DETECTION**

Detect images (or pixels) that are unrelated to the training data

Popular benchmark: Segment Me If You Can [chan21neurips]

The task is to detect pixels that do not belong to any of the 19 road-driving classes:



[https://segmentmeifyoucan.com/]

# EBMs: LIKELIHOOD

Consider  $E_{\theta}(\mathbf{x})$  as a differentiable scalar function of the input:

- $\Box$  we define energy-based density as  $p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z(\theta)}$
- $\square$  the normalizing constant is:  $\textit{Z}(\theta) = \int_{x} \exp(-\mathrm{E}_{\theta}(x)) \textit{d}x$

□ this distribution often appears in physics where it is known as Boltzmann distribution.

We can train such model by optimizing negative log likelihood on training data:

$$\square \mathcal{L}(\theta) = -\frac{1}{N} \sum_{i} \log p_{\theta}(\mathbf{x}_{i})$$

□ this loss is equivalent to KL divergence between the data distribution and  $p_{\theta}(\mathbf{x})$ □ many generative models use this loss (those who do can **recover the density**!)

Conceptually, this is the simplest formulation of a generative model!

## EBMs: LEARNING

We usually optimize  $-\log p_{\theta}(\mathbf{x})$  due to being additive wrt iid data:

$$\begin{split} -\log p_{\theta}(\mathbf{x}) &= -\log \frac{\exp(-\mathrm{E}_{\theta}(\mathbf{x}))}{Z_{\theta}} = \mathrm{E}_{\theta}(\mathbf{x}) + \log Z_{\theta} \\ &= \mathrm{E}_{\theta}(\mathbf{x}) + \log \int_{\mathbf{x}'} \exp(-\mathrm{E}_{\theta}(\mathbf{x}')) d\mathbf{x}' \end{split}$$

The above criterion is minimized when we decrease the energy of the current datum and increase the energy in **all other data**.

The **all other data** part is problematic:

□ we are supposed to tweak the model in all possible images in each iteration

 $\Box$  there are  $256^{32 \cdot 32 \cdot 3}$  images in CIFAR 10 (a number with thousands of zeros...)

□ we refer to such situations as intractable.

#### **EBMs: GRADIENTS**

A closer look at the gradients (note the Leibniz rule) confirms the fears:

$$\begin{split} \frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} &= -\frac{\partial \mathbf{E}_{\theta}(\mathbf{x})}{\partial \theta} - \frac{1}{Z_{\theta}} \int_{\mathbf{x}'} \exp(-\mathbf{E}_{\theta}(\mathbf{x}')) \frac{\partial(-\mathbf{E}_{\theta}(\mathbf{x}'))}{\partial \theta} d\mathbf{x}' \\ &= -\frac{\partial \mathbf{E}_{\theta}(\mathbf{x})}{\partial \theta} + \int_{\mathbf{x}'} p_{\theta}(\mathbf{x}') \frac{\partial \mathbf{E}_{\theta}(\mathbf{x}')}{\partial \theta} \\ &= \mathbb{E}_{p_{\theta}(\mathbf{x}')} \left[ \frac{\partial \mathbf{E}_{\theta}(\mathbf{x}')}{\partial \theta} \right] - \frac{\partial \mathbf{E}_{\theta}(\mathbf{x})}{\partial \theta} \end{split}$$

Note that the above equations show the gradients of the positive log-likelihood.

Applying them will decrease the energy of the current datum and increase the energy in **all other data**.

# EBMs: practice

Practical implementations approximate intractable expectation  $\mathbb{E}$  with MCMC (n=1!):

$$-\frac{\partial \log p_{\theta}(\mathbf{x})}{\partial \theta} = \frac{\partial \mathcal{E}_{\theta}(\mathbf{x})}{\partial \theta} - \frac{\partial \mathcal{E}_{\theta}(\mathbf{x}^{s})}{\partial \theta}$$

We have seen that conceptual simplicity does not imply convenient implementation:

- □ hard learning: the gradients of the loss involve differentiation of an intractable integral
  - you can't iterate over all possible images!
  - even if you could, you couldn't keep all activations that are required for efficient backprop!

□ MCMC approximation requires slow sampling through random walks:

- per-dimension (Gibbs)
- □ hill-climbing (Langevin)
- this appears feasible only for small images

# EBMs: JOINT ENERGY [GRATWOHL20ICLR]

Idea: consider discriminative logits  $s = f_{\theta}(x)$  as logarithm of unnormalized joint density:

$$s_y = \log \hat{p}_{\theta}(\mathbf{x}, y) = \log p_{\theta}(\mathbf{x}, y) + \log Z(\theta)$$

This is an energy-based model since we can consider  $s_y$  as  $-E_{\theta}(\mathbf{x}, y)$ 

$$p_{ heta}(\mathbf{x}, y) = rac{\exp(s_y)}{Z( heta)} = rac{\exp(-\mathrm{E}_{ heta}(\mathbf{x}, y))}{Z( heta)}$$

We can easily recover unnormalized density of the data by marginalizing out y:

$$p_{\theta}(\mathbf{x}) = \sum_{y} p_{\theta}(\mathbf{x}, y) = \frac{\sum_{y} \exp(s_{y})}{Z(\theta)}$$

This is again energy-based density since we can say  $-E(\mathbf{x}) = \log \sum \exp s_y$ 

 $GMCV \rightarrow EBMs$  (5) 24/59

# EBMs: JOINT ENERGY (2)

This view delivers the same class posteriors as in the standard discriminative formulation:

$$\begin{aligned} \mathsf{P}_{\theta}^{\text{disc}}(y \mid \mathbf{x}) &:= \mathsf{softmax}(\mathbf{s}) = \frac{\exp(s_y)}{\sum_k \exp(s_k)} \\ \mathsf{P}_{\theta}^{\text{JEM}}(y \mid \mathbf{x}) &= \frac{p_{\theta}(\mathbf{x}, y)}{p_{\theta}(\mathbf{x})} = \frac{\exp(s_y)/Z(\theta)}{\sum_k \exp(s_k)/Z(\theta)} = \frac{\exp(s_y)}{\sum_k \exp(s_k)} \end{aligned}$$

Both formulations introduce assumptions to be confirmed through training

softmax vs log-unnormalized-joint

The hybrid algorithm must enforce consistent  $\hat{p}(x)$  across data through two loss terms:

- $\Box$  generative:  $-\log p_{\theta}(\mathbf{x_i}) = \log Z(\theta) \log \sum_{y} \exp(s_y)$
- $\Box$  discriminative:  $-\log P_{\theta}(y_i \mid \mathbf{x_i})$
- □ this can succeed since softmax has a spare degree of freedom

### EBMs: SUMMARY

A probability density function must integrate to 1

□ this innocent fact causes much pain!

When the model sees a new datum, the training should raise its density

- □ however, distribution learning is a zero-sum game
- $\hfill\square$  if you give some probability to  $x_{127},$  you must take the same amount from the others!

Energy-based models address this by minimizing  $E_{\theta}(\mathbf{x}) + \log Z_{\theta}$ 

 $\Box$  in practice this involves approximation  $Z_{\theta} = \exp(-E_{\theta}(\mathbf{x}^s))$  (MCMC with n=1 )

This can also be addressed by embedding the unit-integral constraint into the model itself

this is the point where we meet normalizing flows!

# NFs: FROM EBMs TO FLOWS

An EBM model  $f_{\theta_{EBM}}$  maps high-dimensional inputs x into scalar energy z  $\Box$  training EBMs is hard since the predicted energy gives unnormalized density.

A normalizing flow  $f_{\theta_{NF}}^{-1}$  bijectively maps high-dimensional inputs  $\mathbf{x}$  into decorelated latent vectors  $\mathbf{z}$  such that  $p_{\mathbf{z}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ .

 $\Box$  as a reward for being clever,  $p_x(\mathbf{x})$  is easily recovered through change of variables

The learned distribution can be easily sampled by evaluating random noise  ${\bf z}$  in the forward direction:  ${\bf \hat x}=f_{\theta_{\rm NF}}({\bf z}),$ 

Requirements (no free lunch):

 $\ \ \Box \ \ f_{\theta} \text{ must be bijective} \Rightarrow \text{dim}(z) \text{ must equal dim}(x),$ 

 $\Box \det \nabla \mathbf{f}_{\theta}^{-1}$  must be easy to compute.

# NFs: DEFINITION

A normalizing flow f<sub>θ</sub> bijectively maps high-dimensional z into high-dimensional x
 □ consequently, it can recover p(x) through change of variables:

$$\boldsymbol{\rho}_{\boldsymbol{x}}(\mathbf{x}) = |\det \frac{\partial \mathbf{z}}{\partial \mathbf{x}}| \cdot \boldsymbol{\rho}_{\boldsymbol{z}}(\mathbf{z}) = |\det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x})| \cdot \boldsymbol{\rho}_{\boldsymbol{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x}))$$

- $\square$  distribution of the latents is hardwired:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
  - □ fully factorized latents represent independent factors of variation
- $\square$   $p(\mathbf{x})$  integrates to 1 **by design** (even with random weights!)

We usually use the absolute determinant so we can neglect the direction of integration.

Non-bijective flows are also feasible (however, the learning gets more involved).

#### NFs: CHANGE OF VARIABLES VS INTEGRATION BY SUBSTITUTION (1D)

Let us apply the substitution rule for integration of  $p_z(z)$  while assuming x = f(z):

$$P_{z}(z \in S) = \int_{z \in S} p_{z}(z) dz = \int_{x \in f(S)} p_{z}(f^{-1}(x)) \cdot |\frac{dz}{dx}| \cdot dx = \int_{x \in f(S)} p_{z}(f^{-1}(x)) \cdot |\nabla f^{-1}(x)| \cdot dx$$

However, we can also express  $P(z \in S)$  directly in terms of  $p_x(x)$ :

$$P_z(z \in S) = P_x(x \in f(S)) = \int_{x \in f(S)} p_x(x) dx$$

The two equations hold for any S. Hence, the comparison of the two right-hand sides shows that the change of variables is closely related to **integration by substitution**:

$$p_x(x) = p_z(f^{-1}(x)) \cdot |\nabla f^{-1}(x)|$$

 $GMCV \rightarrow NFs$  (3) 29/59

#### NFs: example

Assume we have a random variable  $Z \sim U(0, 1)$ .

Assume we have a dependent random variable  $X = f(Z) = 2 \cdot Z$ .

Obviously,  $p_Z(z) = 1$  for all  $z \in [0, 1]$ . But how much is  $p_X(0.4)$ ?

Let us find out by changing the variable:  $p_X(x) = p_Z(f^{-1}(x)) \cdot |\nabla f^{-1}(x)|$ 

We recover  $p_x(0.4)$  as  $p_z(f^{-1}(0.4)) \cdot |\nabla f^{-1}(0.4)| = 0.5$ .

□ 
$$z(X = 0.4) = f^{-1}(0.4) = 0.2$$
  
□  $p_z(Z = 0.2) = 1$ 

 $\Box \nabla f^{-1}(0.4) = 0.5$ 

 $\hfill\square$   $\textit{p}_{\textit{X}}(0.4) = 0.5$  fits since P(Z < 1) = P(X < 2) = 1 .

### NFs: change of variables (N-D)

The relations remain very similar even in the multivariate case:

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x})) \mid \det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x}) \mid$$

Instead of the derivative, now we have a determinant of the Jacobian:

 determinant is a product of eigenvalues: it reflects the local volumetric stretching of the linear transformation implied by the matrix

If the transformation is composite ( $\mathbf{f}_{\theta} = \mathbf{f}_1 \circ \mathbf{f}_2 \circ ... \circ \mathbf{f}_L$ ), then the composite determinant is a product of determinants of the individual Jacobians:

$$|\det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x})| = \prod_{\ell} |\det \nabla \mathbf{f}_{\ell}^{-1}(\mathbf{z}_{\ell})|, \quad \text{where } \mathbf{z}_1 = \mathbf{x}.$$

# NFs: MODEL

A normalizing flow  $\mathbf{f}_{\theta}$  bijectively maps the latent  $\mathbf{z}$  to the input  $\mathbf{x}$ 

 $\Box$  "normalizing": a complex distribution  $p_x(x)$  is transformed into a normalized one  $p_z(z)$ :

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x})) \mid \det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x}) \mid$$

□ "flow": the transformation is performed through a number of differentiable steps:

$$\mathbf{f}_{\theta} = \mathbf{f}_1 \circ \mathbf{f}_2 \circ \dots \circ \mathbf{f}_L$$

It is a generative model since it can:

- $\Box$  evaluate the density  $p_{\mathbf{x}}(\mathbf{x})$
- $\hfill\square$  generate new data by sampling from z:  $\mathbf{x}$  =  $\mathbf{f}(\mathbf{z})$
- $\square$  perform mapping to the latent space:  $\mathbf{z} = \mathbf{f}^{-1}(\mathbf{x})$

### NFs: Loss

Normalizing flows are trained to maximize the likelihood given the training data:

$$\begin{aligned} \mathcal{L}(\theta | \mathbf{x}_i) &= -\log p(\mathbf{x}_i) = -\log p_{\mathbf{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x}_i)) - \log \mid \det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x}_i) \mid \\ &= -\log p_{\mathbf{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x}_i)) - \sum_{\ell} \log \mid \det \nabla \mathbf{f}_{\ell}^{-1}(\mathbf{z}_{i\ell}) \mid \end{aligned}$$

□ the term with prior density pushes the latents to 0

 $\square$  we usually have  $\mathbf{z} \sim \mathcal{N}(0,1)$ 

□ the term with the determinants opposes and prevents the collapse

This objective can be optimized with the usual variants of SGD

### NFs: AFFINE COUPLING

The most popular computational layer for normalizing flows:

- $\square$  the input representation  ${\bf z}$  is split (e.g. mapwise) into two subsets  ${\bf z}_1$  and  ${\bf z}_2$
- □ the two subsets are separately processed as follows [dinh17iclr]:

$$\begin{aligned} \mathbf{z}_1' &= \mathbf{z}_1 \\ \mathbf{z}_2' &= \mathbf{z}_2 \odot \exp s_{\theta_s}(\mathbf{z}_1) + t_{\theta_t}(\mathbf{z}_1) \end{aligned}$$

This formulation supports the inverse (generative) pass:

$$\begin{split} \mathbf{z}_1 &= \mathbf{z}_1' \\ \mathbf{z}_2 &= [\mathbf{z}_2' - t_{\theta_t}(\mathbf{z}_1')] \oslash \exp \textit{s}_{\theta_s}(\mathbf{z}_1') \end{split}$$

The coupling network  $(s_{\theta_s}, t_{\theta_t})$  can have arbitrary structure and complexity.

# NFs: AFFINE COUPLING (2)

Recall the forward pass through the affine coupling module:

$$\begin{aligned} \mathbf{z}_1' &= \mathbf{z}_1 \\ \mathbf{z}_2' &= \mathbf{z}_2 \odot \exp \mathbf{s}_{\theta_s}(\mathbf{z}_1) + t_{\theta_t}(\mathbf{z}_1) \end{aligned}$$

The Jacobian is triangular; its determinant is the product along the diagonal:

$$\mathsf{det}\frac{\partial[\mathbf{z}_1',\mathbf{z}_2']}{\partial[\mathbf{z}_1,\mathbf{z}_2]} = \prod_j \exp s_{\theta_s}(\mathbf{z}_{1j}) = \exp \sum_j s_{\theta_s}(\mathbf{z}_{1j}) \; .$$

The subsequent coupling module will revert the splits:

 $\hfill\square$  this time  $\mathbf{z}_2$  will pass through unchanged

# NFs: Affine Coupling (3)

Affine coupling transforms each input coordinate with a scalar affine transform:

□ the transform depends on the coordinate and the example!

The operation as a whole is non-linear:

□ it can transform a banana-shaped distribution into a normal distribution

This is related to B-cos networks where the non-linearity is introduced by per-example scaling

# NFs: GLOW [KINGMA18NEURIPS]

The Glow unit consists of affine coupling, 1x1 convolution and ActNorm.

1x1 convolution corresponds to a matrix multiplication at the feature level:

- □ inverse is easy after LU decomposition
- □ Jacobian determinant is a product of diagonal elements of U ( $L_{ii} = 1$ )

ActNorm is a kind of batchnorm that allows very small batches:

- initialize population statistics on the first mini-batch
- □ afterwards, treat them as regular parameters!
- can learn with batch size 1
- practical flows tend to require a lot of GPU memory :-(

# NFs: MORE DETAILS

The latent of the normalizing flow never collapses to z = 0 since that would be penalized by the determinant of the flow Jacobian:

 $p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{f}_{\theta}^{-1}(\mathbf{x})) \mid \det \nabla \mathbf{f}_{\theta}^{-1}(\mathbf{x}) \mid$ 

Normalizing flows require careful implementation in order to avoid numerical instability.

Normalizing flows can also be built on top of standard residual architectures (non-bijective flows) [behrmann19icml]:

- □ iterative inverse with fixed-point algorithm (price: 5-10 forward passes)
  - $\square$  key condition: Lipschitz condition of residual blocks Lip(g) < 1
- Jacobian determinant also iteratively approximated
  - Skilling-Hutchinson trace estimator

# NFs: SUMMARY

Properties:

- $\square$  exact inference:  $p(\mathbf{x}) = p(\mathbf{z}) \prod_{\ell} |\det \nabla \mathbf{f}_{\ell}^{-1}(\mathbf{z}_{\ell})|$
- $\Box$  fast generation  $\mathbf{x} = \mathbf{f}(\mathbf{z})$
- $\square$  fast mapping to the latent space  $\mathbf{z} = \mathbf{f}^{-1}(\mathbf{x})$

Requirements on computational units ("layers"):

- □ bijective, support for forward and inverse transformations
- efficient evaluation of the Jacobian determinant

Rationale:

- when the model sees news data, it increases their density
- □ this rearranges the entire surface of the model (everybody gets affected!)
- the density always integrates to 1 by construction

# NFs: pros and cons

Advantages of normalizing flows with respect to other generative models:

- □ training without MCMC sampling (vs EBM, diffusion)
- no mode collapse (vs GAN)
- □ exact density (vs VAE, diffusion)
- □ fast generation (vs autoregressive, diffusion)
- moderate computational complexity (vs diffusion)

Weaknesses:

- □ inefficient use of capacity (vs EBM)
- □ large training footprint (vs EBM)
- □ no translational equivariance (vs EBM)

#### ANOMALY DETECTION: ABOUT ANOMALIES

**Definition**: an observation which arouses suspicions of being unrelated to the process that generates training data [hawkins80book].



[yang21arxiv]

Anomalous data points are related (or also known as) outliers, out-of-distribution samples or novelties [ruff21pieee].

 $GMCV \rightarrow$  Anomaly detection 41/59

# ANOMALY DETECTION: ABOUT ANOMALIES (2)

There are several flavours of anomaly detection:

- pointwise vs groupwise
- contextual pointwise vs contextual groupwise
- □ low level (texture) vs high level (semantics)







[ruff21pieee]

#### ANOMALY DETECTION: OVERVIEW

Three main approaches to express anomaly score *s* of the sample x

- i) discriminative: arbitrary scalar  $s_{cla} = f(\mathbf{x})$
- ii) generative: probability density  $s_{pdf} = p(\mathbf{x})$
- iii) reconstructive:  $\textit{s}_{\mathrm{rec}} = \| \mathbf{x} \textit{f}_{\mathrm{dec}}(\textit{f}_{\mathrm{enc}}(\mathbf{x}) \|$



[ruff21pieee] GMCV  $\rightarrow$  Anomaly detection (3) 43/59

# ANOMALY DETECTION: OVERVIEW (2)

Three main approaches to express anomaly score s of the sample x

- i) discriminative: arbitrary scalar  $s_{cla} = f(\mathbf{x})$ 
  - □ problems: feature collapse, negative training data
  - $\Box$  related to dataset posterior  $P(\mathcal{D}_{in}|\mathbf{x}) = \sigma(s_{cla})$
- ii) generative: probability density  $s_{pdf} = p(\mathbf{x})$ 
  - □ problem: semantic anomalies (if applied to the data)
  - □ problems: feature collapse (if applied to semantic features)
- iii) reconstructive:  $s_{\rm rec} = \|\mathbf{x} f_{\rm dec}(f_{\rm enc}(\mathbf{x})\|$ 
  - problems: generalization in inliers, "identity" shortcut
  - $\square$  related to dataset posterior:  $\textit{P}(\textit{s}_{\rm rec} | \mathcal{D}_{\rm in}, x) \sim e^{-\textit{s}_{\rm rec}^2}$

#### ANOMALY DETECTION: REALITY CHECK

Direct application of estimated density (either flows or pixel-cnn) to detection of semantic outliers fails miserably:



[serra20iclr]

 $GMCV \rightarrow Anomaly detection (5) 45/59$ 

### ANOMALY DETECTION: ROLE OF COMPLEXITY

Recovered densities wildly depend on image complexity:

- $\hfill\square$  consider images with lower compressed lengths  $L({\rm x})$
- □ e.g. MNIST, poliglot, constant (the simple ones)
- these images score higher in spite of being outliers





[krizhevsky09tr]



- CIFAR10 (Train)
- CIFAR10 (Test)
- Constant (Test)
- Omniglot (Test)
- MNIST (Test)
- FashionMNIST (Test)
- SVHN (Test)
- CIFAR100 (Test)
- CelebA (Test)
- FaceScrub (Test)
- TinyImageNet (Test)
- TrafficSign (Test)
- Noise (Test)

[serra20iclr]

#### ANOMALY DETECTION: INAPPROPRIATE BIAS

Possible explanation: maximum likelihood training is unable to recover semantic anomalies since generative models know nothing about semantics.

- visualization of internal activations suggest a similar reaction for inliers and outliers
- □ however, they train only with generative loss  $L = -\log p(\mathbf{x})$
- □ chances improve when sharing features with a discriminative task [zhang00eccv]





(b) ImageNet input, in-distribution



(c) CelebA input, OOD [kirichenko20neurips]

 $GMCV \rightarrow$  Anomaly detection (7) 47/59

#### ANOMALY DETECTION: HYBRID OPEN-SET RECOGNITION

Open-set performance can be improved by evaluating the density of semantic features

- □ discriminative and generative predictions share the latent features
- □ the anomaly score corresponds to inverse density
- both losses affect the shared features (hybrid recognition)



<sup>[</sup>zhang20eccv]

#### DENSE ANOMALIES: GENERATIVE APROACHES

Detect regions with low pixel-level density:

- □ apply any generative model to 1x1 feature windows [blum21ijcv]
- □ find a way to train EBM without sampling DenseHybrid [grcic22eccv]

Leverage generative modeling in non-density based approaches:

- □ detect reconstruction errors in the resynthesized image [lis19iccv]
- discriminative training with jointly generated synthetic negatives NFlowJS [grcic21visapp, grcic21arxiv]









[https://segmentmeifyoucan.com/] GMCV  $\rightarrow$  Dense anomalies 49/59

#### DENSE ANOMALIES: PIXEL-LEVEL DENSITY

Dense density estimation: recover probability density as if in a sliding window

Desireable properties: efficient inference, equivariance to translation

- □ VAE fast, can be equivariant
- EBM fast, equivariant (but very hard training)
- pixel-cnn slow, not equivariant ("linear" factorization)
- □ flow fast, not equivariant
- □ diffusion, score-based slow, not equivariant
- GAN fast, can be equivariant (but no explicit density)



#### Dense anomalies: sliding $1 \times 1$ window

Apply per-layer flows to frozen features of a standard semantic segmentation model

embedding density [blum21ijcv]

The inference considers normalized feature density with respect to the layer statistics:  $\overline{N}(z_{\ell}^{(i)}) = \log p(z_{\ell}^{(i)}) - \frac{1}{N} \sum_{k} \log p(z_{\ell}^{(k)})$ 

normalized contributions are suitable for ensembling

 $\Box$   $\ell$  denotes the layer, *i* denotes the pixel:

Strength: combines principled density estimation with feature semantics

Weakness 1: vulnerability to feature collapse due to frozen features

Weakness 2: does not exploit negative training data

#### DENSE ANOMALIES: IMAGE RESYNTHESIS

Approach [lis19iccv, vojir21iccv, dibiase21cvpr]:

- 1. perform standard semantic segmentation
- 2. resynthesize input by generative image-to-image translation
- 3. detect anomalous pixels as reconstructions errors



<sup>[</sup>lis19iccv]

Strength: rather principled, can detect all kinds of anomalies

Weakness 1: diverse inliers and wrong predictions lead to false positive anomalies

Weakness 2: works only on the road, fails in non-standard road pixels

Weakness 3: rather slow, unsuitable for real-time

# DENSE ANOMALIES: TWO HEADS (BOTH DISCRIMINATIVE)

Idea: detect outliers with a discriminative prediction head

- □ craft mixed-content images by pasting noisy negatives into regular training images
- leverage negative images from the ImageNet dataset

Input image x

train the dense closed-set classifier only on positive pixels





#### DENSE ANOMALIES: DENSE HYBRID: IDEA

Idea: improve the two-head approach with per-pixel density estimates:

- □ pro: generative and discriminative anomaly detection exhibit different failure modes
- □ con: requires negative data (use case for sythetic negatives)
- con: most generative models are not translation equivariant
  - GANs and VAEs can not deliver exact density
  - hybrid EBM is a method of choice for this task



#### DENSE ANOMALIES: SYNTHETIC NEGATIVES

Training with pasted noisy negative samples produces great outlier detection performance

- □ hard to evaluate the performance
- □ some test anomalies may have been seen during training...

We wish to address this issue by replacing real negative samples with synthetic ones

Question: how to co-train a generative model in order to produce negative examples which could teach the discriminative model to better recognize anomalies?

# DENSE ANOMALIES: SYNTHETIC NEGATIVES (2)

We pose the following requirements on model parameters  $\theta$ :

- $\Box$  high density in inliers  $p_{\theta}(x)$
- $\Box$  high discriminative entropy in generated data  $P(y|\mathbf{f}_{\theta}(z))$

Such learning generates samples at the border of the training distribution [lee18iclr]



[lee18iclr]

The dicriminative model can be trained to predict high uncertainty in these samples!

 $GMCV \rightarrow Dense$  anomalies (8) 56/59

# DENSE ANOMALIES: SYNTHETIC NEGATIVES (3)

We adapt the joint learning scheme for dense prediction:

- □ we use a normalizing flow instead of GAN (arbitrary resolution, better coverage)
- we contribute a robust loss that accounts for generative noise
- we propose a suitable optimization procedure for joint learning



# CONCLUSION

- □ Generative recognition has experienced a lot of exciting recent progress
  - we are proud of our systems although they are not intelligent in the strong sense.
- □ Image is a collection of easily counterfeited pixels
  - □ digital photographs should not be treated as a reflection of the reality
  - verification of integrity possible only in presence of cryptographic signatures
  - important implications for our society
- Semantic anomaly detection can not be properly addressed in absence of semantic supervision.
- Open-set recognition appears easier than four years ago
  - □ it is not unlikely that soon it will be considered as solved.

# Thank you for your attention!

Questions?

This presentation would not have been possible without insightful ideas and hard work of Matej Grcić, Jakob Verbeek, Ivan Krešo, Marin Oršić, Petra Bevandić, Josip Šarić, Ivan Grubišić, Marin Kačan, Iva Sović, Nenad Markuš and Jelena Bratulić.

This research has been supported by Croatian Science Foundation (MULTICLOD, ADEPT), ERDF (DATACROSS, A-UNIT, SAFETRAM), NVidia Academic Hardware Grant Program, Rimac automobili, Microblink, Gideon brothers, Romb technologies, Promet i prostor, Končar, UniZg-FPZ, and VSITE.