

## Convolutional Neural Networks

Deep Learning (DL) systems are usually have 4 main components - 1) **Feature Extractor** (responsible for capturing hierarchical data patterns; I<sub>e</sub> edges to help the network 'understand' the data - CNNs), 2) **Task Specific Head** (component tailored to the problem), 3) **Parameter Optimization**. For a non-constant, bounded, monotonically increasing function, for any  $\epsilon > 0$ , and any conv,  $f$  defined on a compact of  $\mathbb{R}^n$ , there exists such that:  $F(x) = \sum_i w_i \phi(x + b_i)$  with  $|f(x) - f(x)| < \epsilon$ ; is high;  $C$  means for more accuracy, need more layers = exp. more compute. DNNs help to generalise accurately grow exponentially to optimise a Lipschitz cont. func.  $F$ :  $R^d \rightarrow R$  with  $1 \times 1$  accuracy, we need  $O(\epsilon^{-d})$  samples. e.g if in 2D space you want to approx a func; and you need 10 samples for  $\epsilon$  in 1D, you'll need 100 in 2D. For 12 Pixels<sup>3</sup>: 3 channels = 36M elements, for  $O(1)$  you need  $10^6$  36M.

**Mass Concentration in High Dimensions:** in higher D, most mass of objects resides near the boundaries; the most data points lie near the edges of the data space. Consider in 1D, the  $n$ th DNN needs to learn the original data distribution which is local or shrinked. The volume of the rind relative to original vol is  $= 1/n$ . As a function of the rate of growth is  $(1+n)^{-n}$ .

With no shrinking ( $\epsilon=1$ ) and noting that  $a$  is decreasing, the initial growth rate of the rind is  $n$ . The volume of the rind initially grows much faster in times faster than the rate at which the object is being shrunk. In higher D, tiny changes in distance cause massive changes in volume. Minuscule changes in datasets can massively change the volume in D1.

**Length-Half:** How much do you need to shrink object in all D to have 1/vol?

Solve  $\epsilon^n = 1/2$ : In 1D: 0.35; so half the data lies within 35/2 = 18% of its diameter from boundary. In 3D: 1.02...; half the data is in 1% of diam from b.

Distance measures lose effectiveness in measuring dissimilarities in high D.

**Invariance:** no matter how the input is transformed, the output is invariant

**Equivariance:** the output is transformed in the same way as the input

**Shift-Invariance:** response is the same when input (like in pixels) are shifted.

Important in CNNs, which have SI, which is why they're good at image tasks.

**Shift-Equivariance:** the response shifts by the same amount as the input shifts.

I.e. for a model that segments an image, marking cat pixels "1" and not mark "0".

ie shift and then mark is the same as mark and then shift.  $s'(x) = s(x)$

**Locality:** You don't have to look far away from a location ( $i,j$ ) to find info about it

The assumptions of locality and translation invariance have guided model dev.

**Correlation:** involves sliding a subimage across a base image and comparing pixel values. While shift-invariant and local, it doesn't always find the best

Correlation: it's based on the mean of the local receptive field. If we had FCNs each unit is connected to each node in the previous layer. So every single input feature influences every neuron. Very flexible, can capture any relation, but very expensive and easily overfits. Consider 36M inputs (image); we end with 36M<sup>2</sup> params which is astronomical. More sparsely connected nets are better; these paved the way for CNNs.

**Weight Sharing:** A subset of weights are identical and used across multiple connections. **Parity:** sharing makes CNNs very powerful!

$$(f * g)(t) = \int_0^t f(t-u)g(u)du = \int_0^t f(t-u)du \cdot g(t) = f(t) \cdot g(t)$$

Convolutions:  $f * g$  is the  $\text{corr}(f, g)$  of  $f$  and  $g$   $\in \mathbb{R}^{\infty \times \infty}$

of  $f$ , and doesn't involve flipping the kernel).

Both Conv and Corr are similar for ML models. Both aim for local pattern detection thru shared weights. We use discrete conv on arrs.

Given arr  $u_{t,a}$  and  $w_{t,a}$ , the conv is the fcn s.t.  $\sum_a u_{t,a}w_{t,a}$

Mathematically, conv = cross-correlation between  $f(x)$  and  $g(x)$  and  $g(x)$

**Properties:** 1) **Commutativity:**  $f * g = g * f$ . 2) **Associativity:**  $f * (g * h) = (f * g) * h$

3) **Distributivity:**  $f * (g+h) = (f * g) + (f * h)$ . 4) Assoc w/ scalar mult:  $a(f * g) = (af) * g$

CNN Building Blocks: 1) **Convolutional Layers:** core of the CNN, where kernels and filters slide across input producing feature maps. Catches local patterns and details like edges and textures, making them very suited for image tasks

2) **Activation Functions:** introduces non-linearity, allowing the model to learn complex mappings from input to output. ReLU is commonly used as its good

3) **Pooling Layers:** Takes the input of the previous layers and summarises/simplifies it. This reduces the spatial dimensions of the feature maps.

4) **Batch Normalization:** standardises the data. Mean-Pooling etc.

5) **Fully Connected Layers:** used towards the end of the model. Each neuron connects to each in prev layer. Integrates learned features for final prediction.

6) **Flatten Layers:** Bridge b/w conv layers and fully connected layers by converting 2D feature maps into 1D vectors.

7) **Dropout:** Regularization technique: randomly zero some neuron weights. Reduces overfitting and makes the network more robust.

8) **Batch Normalization:** This keeps the distribution of each layer's outputs stable, aiding in faster and more reliable training convergence.

9) **Softmax Layer:** commonly found at the end. It converts the final layer's outputs into a probability distribution, useful for classification tasks.

10) **Loss Function:** measures training perf for backpropagation. Cross-entropy loss is commonly used for classification tasks in CNNs. Losses quantify how well the model's predictions align with the true labels.

11) **Optimizer:** Adjusts the network weights based on gradients found during backprop. SGD (Stochastic Gradient Descent) and Adam are popular choices.

**Kernels:** For a typical NN you'd have took an image and flattened it into an  $X \times Y \times Z$  sized tensor. But these do not take advantage of our assumptions:

Locality and Translation Invariance! Kernels do, and thus CNNs do by doing a  $X, Y, Z \times [I, J, K]$  conv to make feature maps; I = kernel height, J width, K depth. The resulting feature map is of size  $(X-I+1) \times (Y-J+1) \times (Z-W+1)$  padding = 0 and if we're using 1 filter - conv layers - will use single filters / kernels at the same time. I.e. all the same dimension. But the feature's shape can then learn to on. To present the (X, Y) size we can use padding, we can zero pad the edges of the image with  $[I/2, J/2]$  pixels. This is helpful if our input is quite small as we want lots and lots of conv. layers (which otherwise decrease the size of the image).

At its core a CNN is just a seq. of conv. layers with activation functions.

Computational Complexity: assume a  $5 \times 5$  convolution. It involves 25 multiplications + add ops per input element. But if we use two consecutive  $3 \times 3$  convolutions instead which can be constructed to produce the an approximation of the same output, they only take 18 total ops per element. But the second conv has to work on even fewer elements because the shrinking thanks to the first conv as well - so the efficiency gains are even better! We also go from 25 to 18 param

**Separable Convolutions:** You can approximate a  $5 \times 5$  conv with a  $1 \times 5$  conv followed by a  $5 \times 5$  conv - if the  $5 \times 5$  conv is separable. It works well only if the original conv is approximable. Example: Consider an image of dim (D, M, N). If  $M=3 \times 3$  and  $D=3$ , the  $1 \times 5$  conv then  $5 \times 1$  conv is  $3 \times 2 \times 3 \times 3 \times 5 \times 1 + 2 \times 2 \times 3 \times 5 \times 3 \times 1 = 75602$  ops. But:  $28 \times 28 \times 3 \times 5 \times 3 \times 1 = 176401$ . Pooling: Permutation invariance - involves aggregating/downsampling. Reduces the amount of a layer hierarchy, aids shift invariance. Max-Pooling/Avg-Pooling is applied to each channel separately, e.g. RGB on images.

Niquist Sampling Theorem: a continuous signal can be completely represented by its samples and fully reconstructed if it is sampled at least twice as fast as its highest frequency component. By aliasing via pooling/conv we lose info - we

no longer adhere to the theorem. It's found that Convs/Pooling only improves results when used properly and carefully.

**Image Frequency:** An image has high frequency if it "oscillates" a lot in colour. High freq images have lots of edges, textures and fine details (high gradients).

Max-Pooling breaks shift equivariance. A solution is to use a Gaussian Blur to downsample. Sadly, CNNs are not rotation invariant; but this can be achieved with spherical harmonics (or remember Hog7).

CNNs are largely deformation invariant -  $f(x) = f(D_t x)$  and this makes them amaze-

ing for image classification.

**Flattening Layers:** connect conv and fully conv layers, as conv layers output a

3D tensor of learned features, and FC layers expect a 1D tensor. Flattening layers use 1x1 convs to do this; i.e. they take a  $4 \times 4 \times 6$  tensor into a  $1024 \times 1$

**Activation Functions and Loss**

Activation Functions: they determine how neurons in the network activate with a set of inputs. Introduces non-linearity, allows learning complex behaviour.

1) **Naive Activation:** if the input is above a specific value, declare the neuron activated (output = 1). Not differentiable, so cannot backprop.

2) **Linear Activation:**  $c(x)$ : Constant grad, no relation to x during backprop.

3) **Sigmoid Activation:** output =  $1 / (1 + e^{-x})$ , nonlinear - allows us to stack layers. Continuous output, differentiable. Curve is steep between 2 and -2 giving a lot of feedback during training, allowing for quick learning.

4) **Tanh** ( $tanh(x) = e^{2x} - 1 / e^{2x} + 1$ ): It's chosen when we want a central sigmoid with output from -1 to 1 rather than 0 to 1, and steeper gradients than sigmoid.

5) **ReLU** ( $max(0, x)$ ): Nonlinear, especially thanks to the sharp corner at origin.

Combos of ReLU are non linear, Unbounded, useful for some types of data but causes activation explosion in relus. ReLU often outputs 0, because its response to negatives which causes some neurons to be ignored making the net more computationally efficient. But, this sparsity causes the dying ReLU problem: if a neuron output is always 0 then the gradient is too: preventing that portion of the NN from being updated through backprop. ReLU is used a lot due to efficiency and effectiveness in many applications.

6) **Leaky ReLU:** output =  $x$  for  $x \geq 0$ ,  $0.01x$  for  $x < 0$ . Has a small slope for negative values (nonzero gradient) avoiding the dying ReLU problem.

7) **PReLU:** output =  $x$  for  $x < 0$ ,  $0.01x$  for  $x > 0$ , is learnable, and either one is used for all channels or separate one for each input channel (more nuance behaviour during training at the cost of extra computational complexity).

All the ReLUs are scale invariant which makes them suited for some tasks.

8) **SoftPlus**  $1 / \theta * \ln(1 + e^x)$ : This is a smooth (easier to differentiate)

approximation of ReLU with only positive output. The higher  $\theta$  is, the closer to ReLU it is; and we default to linear if  $x > \theta$ . Smoother, differentiable alternative to ReLU that is desirable when we want positive activations and numerical stability.

9) **ELU** ( $max(0, x) - min(0, e^{-\alpha} * (x - \theta))$ ): It is element wise, operating one each input element independently. ELU can go below 0 and is parameterized by  $\theta$ . It's scale variant and messes with gradients (thanks to being able to output negative values) - this helps the network converge quicker. It's good when we want a balance of smoothness, differentiability, and the ability to have a mean around 0. This is a unique blend of features.

10) **CELU** ( $max(0, x) - min(0, e^{-|x|} * (c - |x|))$ ): ELU variant: element wise.

Continuously differentiable when  $a$  not equal to 1. Builds on ELU w/ more mathematical rigour to suit certain applications and optim scenarios.

11) **SELU** ( $scale * max(0, x) + min(0, e^{-\alpha} * (x - \theta))$ ): It comes with predefined a and scale which have been meticulously optimized. SELU allows for internal normalization,  $a$  and scale are solutions to a fixed  $p$  eq that maintains mean 0 and var 1 across layers. Normalizing activations in an NN can happen at 3 levels. The first is input norm, where we scale input features,  $e.g.$  grayscale pixel vals, into a specific range, such as 0 to 1. The second level is batch normalization, a technique specifically designed for NNs to stabilize the learning process. SELU shines at the third level, which is internal normalization. SELU ensures that mean & var of activations are preserved one layer to the next. SELU has pos and neg outputs to have a mean of 0. It has gradients close to 0 which helps normalization. SELU avoids the "dying unit" problem that plagues. SELU offers a way to build DNNs without worrying too much about manual normalization. Keeps internal NN stats stable.

12) **GELU**:  $\text{GELU}(x)$  is the Gaussian CDF. This adds probability to the activation process. This can cause regularization; good for training & perf.

13) **ReLU6**:  $min(0, x, 6)$  - this caps the ReLU to a value. This saturates activation preventing exploding growth and numerical instability. The max value of 6 must be configured for different saturation levels.

14) **LogSigmoid**  $log(1/(1 + e^{-x}))$ : Operates elsewhere. Typically used in cost functions rather than activations.

15) **SoftMin**:  $e^{-x} / \sum_i e^{-x_i}$  - this applies softmax on an n-dim tensor. Rescales the sum of the output lie between [0,1] and sum to 1. Acts as a method to convert a set of num into a representation that resembles a prob dist; valuable in scenarios where you want relative weights or preferences among alternatives. SoftMin's role lies in rescaling and transforming input tensors into prob dists, contributing to the multi-dim non-linearities of NNs and enables the modelling of various factors in the data.

16) **SoftMax**:  $e^{-x_i} / \sum_j e^{-x_j}$  - Common and important. Like softmax, rescales elements  $b \in [0,1]$  and sum to 1. Key use case: convert raw scores (logits) generated by a NN into meaningful class probability scores.

17) **LogSoftmax**: log( $\text{SoftMax}$ ). Beneficial as it gives numerical stability. Good for loss functions for NNs, rather than in activations.

18) **SIREN** (Sinusoidal Representation Network) activation:

$$\phi(x) = W_{\phi} \cdot (\sin(\omega_1 x_1 + \dots + \omega_n x_n) + b_{\phi})$$

Periodic activation - useful but can be hard to use. Useful for dealing with implicit representations - i.e. finding a cont. func representing sparse input, such as images, waves. Great for spatial/temporal derivatives of signals.

Error/Loss Functions: quantify how well the network performed, allowing for backprop where we can update model params. Loss functions do this by measuring the distance between the produced output and ground truth.

1) **L1 Norm** (MSE): Take samples from same classes close and different

and compare. Objective: Distance for the good pair has to be smaller than distance to the bad pair. Actual distance does not need to be small, just smaller. Used for metric learning and Siamese networks

(minibatch),  $l = \lVert x_i - x_j \rVert_1$  - (x, y) - penalisles small outliers more; good for noisy data

Not differentiable at 0, so can use smooth versions.

3) **Smooth L1 Loss**:  $l = 0.5 \cdot |x - y| + k \cdot |x - y| - 0.5 \cdot \text{bal}_1$  when we transition between 0 and 1

and -add ops per input element. But if we use two consecutive  $3 \times 3$  convolutions instead which can be constructed to produce the an approximation of the same output, they only take 18 total ops per element. But the second conv has to work on even fewer elements because the shrinking thanks to the first conv as well - so the efficiency gains are even better! We also go from 25 to 18 param

4) **Negative Log Likelihood Loss:** Assum net output is log likelihoods. We want to make the desired class output big, others as small as possible  $l(x,y) = -\sum_i l_i(y_i)$ ,  $l_i = -w_i \ln y_i$ ,  $w_i = \text{weight}[i]$ ,  $1/c - i$ : ignore index

Reductions are used to aggregate the NLL scores per class:

- if using mean reduction:  $l(x,y) = \sum_i l_i / N$  in 1. (x, y) in N

- if using sum reduction:  $l(x,y) = \sum_i l_i$  in 1. N

We use weights to diff importance to classes e.g. imbalanced dataset

5) **Cross Entropy Loss:** combines LogSoftmax and NLLloss

• Useful for classification problems with C classes:

$$loss(x, c) = -\log(\sum_i x[i] \cdot \delta(c)[i]) + \log(\sum_i x[i])$$

• Classes can be weighted. Loses avg across observations for each minibatch.  $\delta(c)$  is the correct class.  $\delta$  esries = all logits.

6) **Binary Cross Entropy Loss:** As above but designed for 2 classes:

$$l(x,y) = -\lvert x - y \rvert \cdot \ln \lvert x - y \rvert$$
,  $l(x,y) = -\lvert x - y \rvert \cdot \ln \lvert x - y \rvert$

• if you don't have probabilities in range [0,1] you need to use "WithLogs"

$y = \text{true label } (0,1)$ ,  $x = \text{predicted prob of class } 1$

no longer adhere to the theorem. It's found that Convs/Pooling only improves results when used properly and carefully.

**Image Frequency:** An image has high frequency if it "oscillates" a lot in colour. High freq images have lots of edges, textures and fine details (high gradients).

Max-Pooling breaks shift equivariance. A solution is to use a Gaussian Blur to downsample. Sadly, CNNs are not rotation invariant; but this can be achieved with spherical harmonics (or remember Hog7).

CNNs are largely deformation invariant -  $f(x) = f(D_t x)$  and this makes them amaze-

ing for image classification.

**Flattening Layers:** connect conv and fully conv layers, as conv layers output a

3D tensor of learned features, and FC layers expect a 1D tensor.

Pooling: Permutation invariance - involves aggregating/downsampling. Reduces the amount of a layer hierarchy, aids shift invariance. Max-Pooling/Avg-Pooling is applied to each channel separately, e.g. RGB on images.

Niquist Sampling Theorem: a continuous signal can be completely represented by its samples and fully reconstructed if it is sampled at least twice as fast as its highest frequency component. By aliasing via pooling/conv we lose info - we

Batch too small the noise level is counterprod, affecting convergence. At inference time we have fixed learned Gamma and Beta and we use running avg of mu and sigma from training to normalise the features.

Dense Layer: One norm for all. Conv layer: One norm per channel.

ResNet: Solves vanishing gradients in deep nets, ensuring effective grad prop.

Deep nets can model more complex funcs but can diverge from optim func instead of improving. Ideally func classes are nested meaning deeper layers should result instead of changing layer sizes.

• Primary role of strided downsampling

• Convolutional layer ( $p$ ):  $\frac{\partial}{\partial x} f(x)$  =  $\frac{\partial}{\partial x} f(x - \Delta x)$

Invariance vs equivariance: variance means

Invariance: the good pair must be smaller than dist to the bad pair. Actual dist doesn't need to be small. Used in metric learning and Siamese networks

10) **Cosine Embedding Loss:** if  $y = 1 - \cos(x - y)$ ,  $y = \max(0, 1 - \cos(x - y))$

Measures whether two inputs are similar or not. Similar to a normalised Euclidean distance but focuses on angles b/w the inputs not magnitude.

**Network Architectures**

LeNet: First designed for low-res, b&w image recog, specifically for digits.

Input is  $32 \times 32$  pixel grayscale, then conv1 layer makes  $6 \times 28 \times 28$  feature maps, then avg pool layer reduces dimensionality to  $14 \times 14$ , 6 channels (pooling applied to each channel, not across them so num channels unchanged). Then conv2 layer  $\rightarrow 10 \times 16$ , channels. Avg pooling  $\rightarrow 5 \times 16$ , channels. Fully connected layers, first 120 units, second 84 channels.  $\text{softmax} \rightarrow 10$  classes.

• Skip connections:  $f(x) = x + f(x)$  (residual connection)

• Layer Normalization:  $\frac{x - \mu}{\sigma}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

• Layer Norm:  $\frac{x - \mu}{\sqrt{\epsilon + \sigma^2}}$  (mean and std deviation)

&lt;p

Divergences: Compares two probability distribution. Given a set of probability distributions  $P$  on a random variable  $X$ , a divergence is defined as a function  $D(P||Q)$ :  $P \rightarrow R$  such that  $D(P||Q) \geq 0$  and  $D(P||P) = 0$ .

- Kullback-Leibler (KL) Divergence:  $KL[p||q] = \int p(x) \log \frac{p(x)}{q(x)} dx$ ,  $p, q \in \mathcal{P}$
- Jensen-Shannon (JS) Divergence:  $JS(p||q) = \frac{1}{2} KL(p||m) + \frac{1}{2} KL(q||m)$  where  $m(x) = \frac{1}{2}(p(x) + q(x))$

Probabilistic Graphical Models (PGMs): Directed acyclic graphs that describe the factorisation structure of a joint probability distribution.

$$p(x_1, \dots, x_D) = \prod_i p(x_i | pa(x_i))$$

(a)  (b)  (c) 

$p(x, z) = p(x|z)p(z)$ ,  $p(x, y, z) = p(x|y, z)p(y|z)p(z)$ ,  $p(x, y, z, w) = p(x|y, z, w)p(y|z, w)p(z|w)p(w)$

Jensen's Inequality: For a convex function  $f$ , then for any probability distribution  $p(x)$ :

$$\mathbb{E}_p[f(x)] \leq f(\mathbb{E}_p[x])$$

with equality iff  $f$  is linear on  $p(x)$ .  $\mathbb{E}_p[x]$  is a data measure.

LO(GAN): Makes binary class' task to help learning of gen've model  $p_\theta(x)$

to fit data dist  $p_{data}(x)$ . Done by labelling datapoints from data dist "real" and from model as "fake". So joint dist  $p(x, y)$  as follows for binary class' n task:

$$p(x, y) = p(x|y)p(y), \quad p(y) = \text{Bern}(0.5), \quad p(x|y) = p_{data}(x), \quad y = 0 \quad (3)$$

Fitting this to  $p_\theta(x, y)$  with  $p_\theta(x) = D_\phi(x)$  to  $p(y|p_\theta(x))$  by MLE obj

$$\phi'(\theta) = \arg \max \mathcal{L}(\theta, \phi), \quad \mathcal{L}(\theta, \phi) := \mathbb{E}_{p_{data}(x)}[\log D_\phi(x)] + \mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))] \quad (4)$$

There is dependence on generative model parameter  $\theta$ , since "data dist":  $\mathbb{P}(x, y)$  of bin class' n task depends on  $p_\theta(x)$ . Training of this gen model tries to fool discriminator by minimising log probability of making the right decisions:

$$\theta'(\phi) = \arg \min \mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))] \quad (5)$$

2-player GAN gen&discr train obj:  $\min_{\theta, \phi} \mathcal{L}(\theta, \phi)$   $\quad (6)$

Term in obj for  $p_\theta(x)$  is  $\mathbb{E}_{p_\theta(x)}[\log(1 - D_\phi(x))]$ , in practice approx by Monte Carlo:

$$\mathbb{E}_{p_\theta(x)}[1 - D_\phi(x)] \approx \log(1 - D_\phi(x)), \quad x \sim p_\theta(x) \quad (7)$$

Comp of dist  $p_\theta(x)$  not required for eval of obj, instead can define directly the sampling process of  $p_\theta(x)$  which also defines dist  $p_\theta(x)$  in an implicit way:

$$x \sim p_\theta(x) \iff z \sim p(z), \quad x = G_\theta(z) \quad (8)$$

... often we set  $p(z) = \mathcal{N}(z, 0)$ ;  $I$

Jenson-Shannon divergence minimisation: To justify 2-player obj (6) we show that with infinite capacity for both gen and discr, the global optimum of gen is  $g(x) = p_{data}(x)$ . For fixed gen  $p_\theta(x)$ , we compute gradient of Gen obj wrt  $\phi$ :

$$\nabla_\phi \mathcal{L}(\theta, \phi) = \int \left( \frac{\partial p_{data}(x)}{\partial D_\phi(x)} - \frac{p_\theta(x)}{1 - D_\phi(x)} \right) \nabla_\phi D_\phi(x) dx \quad (9)$$

Given infinite capacity of the discr, setting  $\nabla_\phi \mathcal{L}(\theta, \phi) = 0$  results in

$$\frac{\partial p_\theta(x)}{\partial D_\phi(x)} = \frac{\partial p_{data}(x)}{\partial D_\phi(x)} \implies D_\phi'(x) = \frac{p_\theta(x)}{p_\theta(x) + p_{data}(x)} \quad (10)$$

Plugging in the optimal discriminator to the GAN objective:

$$\mathcal{L}(\theta, \phi^*) = \mathbb{E}_{p_{data}(x)} \left[ \log \frac{p_{data}(x)}{p_\theta(x) + p_{data}(x)} \right] + \mathbb{E}_{p_\theta(x)} \left[ \log \frac{p_\theta(x)}{p_\theta(x) + p_{data}(x)} \right] \quad (11)$$

$$= \mathbb{E}_{p_{data}(x)} \left[ \log \frac{p_{data}(x)}{\frac{1}{2}(p_\theta(x) + p_{data}(x))} \right] + \mathbb{E}_{p_\theta(x)} \left[ \log \frac{p_\theta(x)}{\frac{1}{2}(p_\theta(x) + p_{data}(x))} \right] - 2 \log 2$$

$$= 2 \left[ \frac{1}{2} \mathbb{E}_{p_{data}(x)} \left[ \frac{1}{2}(p_\theta(x) + p_{data}(x)) \right] + \frac{1}{2} \mathbb{E}_{p_\theta(x)} \left[ \frac{1}{2}(p_\theta(x) + p_{data}(x)) \right] \right] - 2 \log 2$$

$$:= JS[p_{data}(x)||p_\theta(x)] \quad (11)$$

We quantify the error between the lower bound and  $\log p_\theta(x)$  as follows

$$\log p_\theta(x) = \log p_\theta(x|z)p(z) = \log p_\theta(x) - \mathbb{E}_{p_\theta(x)}[\log \frac{p_\theta(x)}{p_\theta(x|z)p(z)}]$$

$$= \log p_\theta(x) + \mathbb{E}_{p_\theta(x)}[\log \frac{p_\theta(x|z)p(z)}{p(z)}] \quad (\text{Bayes' rule})$$

$$= \mathbb{E}_{p_\theta(x)}[\log p_\theta(x|z) - KL(p_\theta(x||p(z))] = \mathcal{L}(x, q, \theta)$$

By improving lower bound  $q(z)$  goes to the exact posterior  $p_\theta(x|z)$ ,  $q, q(z)$  is a tractable posterior which we use the KL divergence to fit, since  $p_\theta(x|z)$  is intractable. Optimal var lower bound wrt model params doesn't *ALWAYS* optimise log likelihood, as KL divergence could grow, so we optimise wrt  $\phi$ , too.

Variational Autoencoder: We define the  $q$  distribution as  $q(z) = q_\phi(x|z)$  generally using a Neural Network:

$$q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$$

$\mu_\phi(x)$ ,  $\sigma_\phi(x)$  = NN( $x$ ,  $\phi$ )

The optimisation object is then given by:  $\phi': \theta \mapsto \arg \max \mathcal{L}(\theta, \phi)$

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{data}(x)} \left[ \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p(z)] \right] \quad (12)$$

Analytical KL term: If  $q_\phi(z|x)$  and  $p(z)$  are assumed to be factorized gaussians we can write the KL divergence term analytically (d = latent dimension) and  $\mathcal{L}$  is simply the sum of all values in  $x$ :

$$KL[q_\phi(z|x) || p(z)] = \frac{1}{2} (\|\mu_\phi(x)\|_2^2 + \|\sigma_\phi(x)\|_2^2 - 2 \log \sigma_\phi(x, 1) - d)$$

Monte Carlo Estimation: the optm obj above relies on evaluating NN transforms for all possible  $z$  to calculate expectation - can be replaced by MC approximations. 1) Draw samples  $z$  from  $q_\phi(z|x)$  (paramised Gaussian) 2) Feed these to the decoder  $p_\theta(x|z)$  which gives the parameters of the likelihood dist of  $x$ . i.e. mean var for Gaussian, the 3) plug those params and sample from  $p(x|z)$  4) Compute log prob of  $p(x|z)$  (bad abt) so  $D_\phi(x) = 0$  for  $x \sim p_\theta(x|z)$ . Also assume gen model implicitly defines  $D_\phi(x) = G_\theta(z)$ .  $D_\phi(x)$  often defined using sigmoid activation (sigmoid( $d_\phi(x)$ ) =  $(1 + \exp(-x))^{-1}$ ) - the last layer, i.e.  $D_\phi(x) = \text{sigmoid}(d_\phi(x))$  with  $d_\phi(x)$  param of a neural net. So  $D_\phi(x) \approx 0$  when  $d_\phi(x) \approx 0$  at the start of GAN training  $d_\phi(x) \rightarrow -\infty$  for  $x \sim p_\theta(x|z)$ . So gradients of the 2 objectives wrt Theta are

$$\nabla_{\theta, \phi} \mathcal{L}(\theta, \phi) = \mathbb{E}_{p_{data}(x)}[\log(1 - D_\phi(x))] - \mathbb{E}_{p_\theta(x)}[\log(1 + \exp(-d_\phi(G_\theta(z)))] \quad (13)$$

$$= \nabla_{\theta, \phi} \mathcal{L}(\theta, \phi) = -\nabla_{\theta, \phi} \mathcal{L}(\theta, \phi) = \mathbb{E}_{p_\theta(x)}[\underbrace{\log(1 + \exp(-d_\phi(G_\theta(z)))]}_{\approx 1}] \quad (14)$$

The alt obj adds the vanishing gradient problem of the OG one (5) at the start of training, hence the name "non-saturated objective".

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

$$= -\mathbb{E}_{p_\theta(x)}[\log D_\phi(x)] - \log 2$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

$$= -\mathbb{E}_{p_\theta(x)}[\log D_\phi(x)] - \log 2$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$

Shows that max alt "non-saturated obj" is equiv to min'ing an f-divergence

$b$  w/ model  $p_\theta(x|z)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

This allows us to sample from Gaussian dist in a way that is differentiable - used for step 1 above. Conditional VAE is to gen images conditioned on other other discrete/continuous info. Corresponds to learning den model  $p_\theta(x|y)$  to approx  $p(x|y)$ ,  $y$  is random var for  $i.e.$   $y \sim \mathcal{N}(0, I)$

Another justification of the alt obj is given by deriving optimal solution of the generator, given optimal discriminator. Define  $f(t) = \log(1 + t^{-1}) - \log 2$ , in which  $f(t)$  is convex and  $f(1) = 0$ . Then can define an f-divergence as

$$D_f[p_\theta(x)||p_{data}(x)] := \int p_\theta(x) \left[ \frac{p_{data}(x)}{p_\theta(x)} \right] dx$$

$$= \int p_\theta(x) \log \left( 1 + \frac{p_\theta(x)}{p_{data}(x)} \right) dx - \log 2 \quad (15)$$