

• HW 2 Grades out

Generative Models

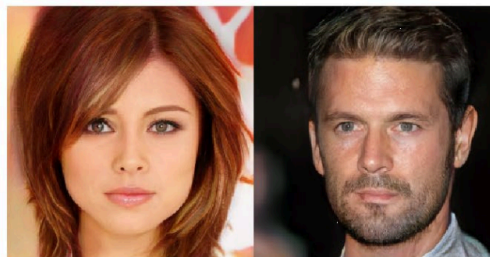


Distribution learning

- learn a distribution P over X ;
- evaluate $P(x)$, $x \in X$
- Sample from distribution



Training
Data(CelebA)

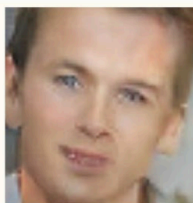


Model Samples (Karras et.al.,
2018)

4 years of progression on Faces



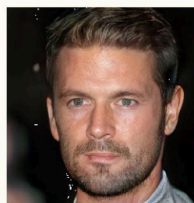
2014



2015



2016



2017

Brundage et al.,
2017

Image credits to Andrej Risteski

Distribution learning



BigGAN, Brock et al '18

Distribution learning

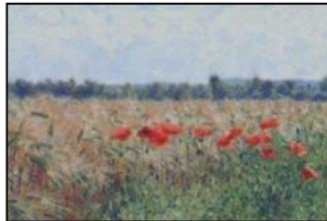
Conditional generative model $P(\text{zebra images} | \text{horse images})$



Style Transfer



Input Image



Monet



Van Gogh

Image credits to Andrej Risteski

Distribution learning

Source
actor



Target
actor



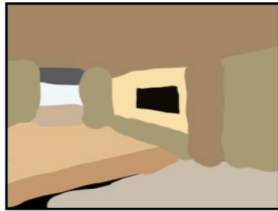
Real-time Reenactment



Reenactment Result

Real-time
reenactment

Generative model



Generative model
of realistic images



Stroke paintings to realistic images

[Meng, He, Song, et al., ICLR 2022]

“Ace of Pentacles”



Generative model
of paintings



Language-guided artwork creation

<https://chainbreakers.kath.io> @RiversHaveWings


Generative model

$$p(x)$$

$$p(y)$$




High
probability



Generative model
of traffic signs

Low
probability



Outlier detection

[Song et al., ICLR 2018]

Desiderata for generative models

$P_{\theta}(\cdot)$, θ is parameter

- **Probability evaluation:** given a sample x , it is computationally efficient to evaluate the probability of this sample.

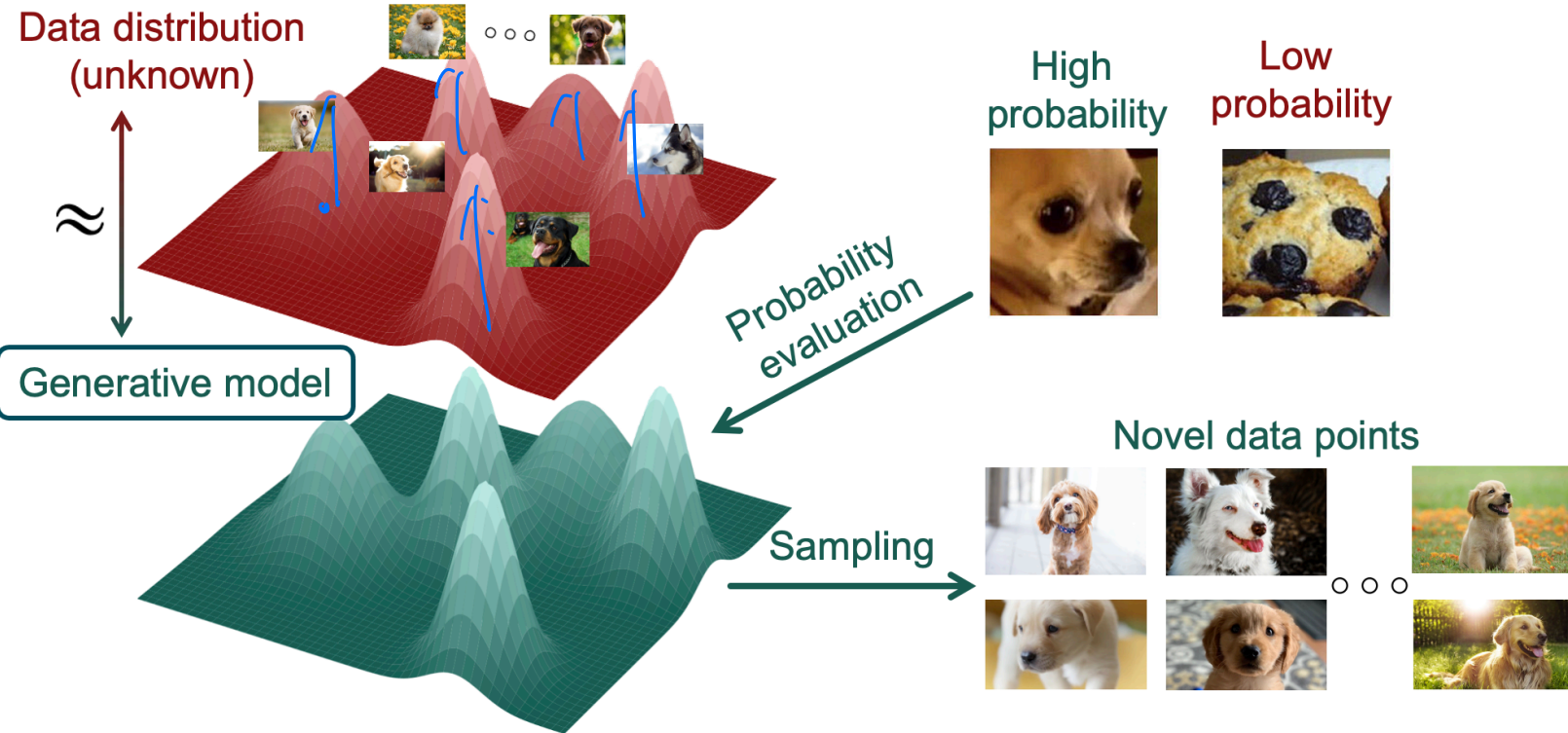
- **Flexible model family:** it is easy to incorporate any neural network models.

$P_{\theta}(\cdot)$: CNN, transformer, supervised. $x \mapsto f(x)$, $P_{\theta}(x) \geq 0$, $\sum_x P_{\theta}(x) = 1$

- **Easy sampling:** it is computationally efficient to sample a data from the probabilistic model.

$P_{\theta}(x)$, $x \sim P_{\theta}(\cdot)$

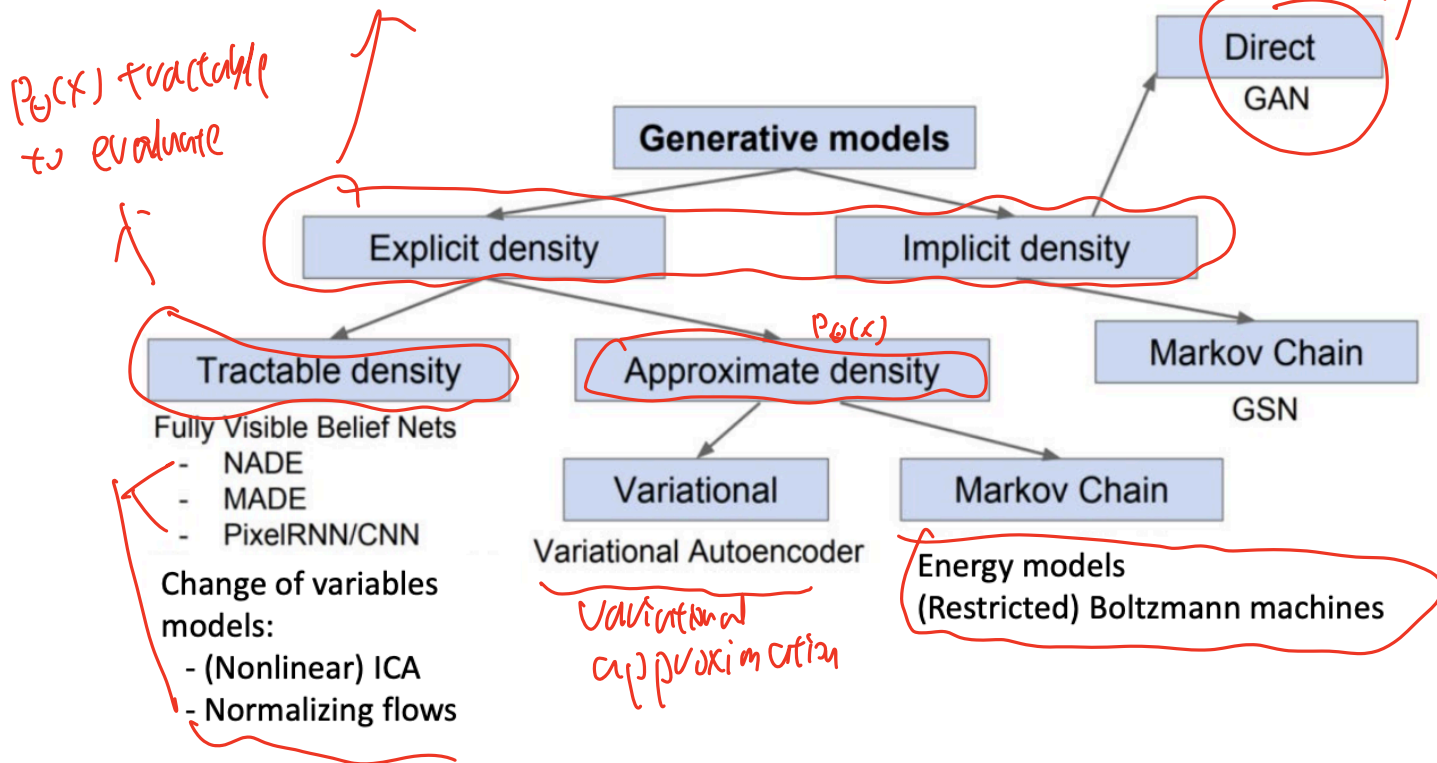
Desiderata for generative models



Taxonomy of generative models

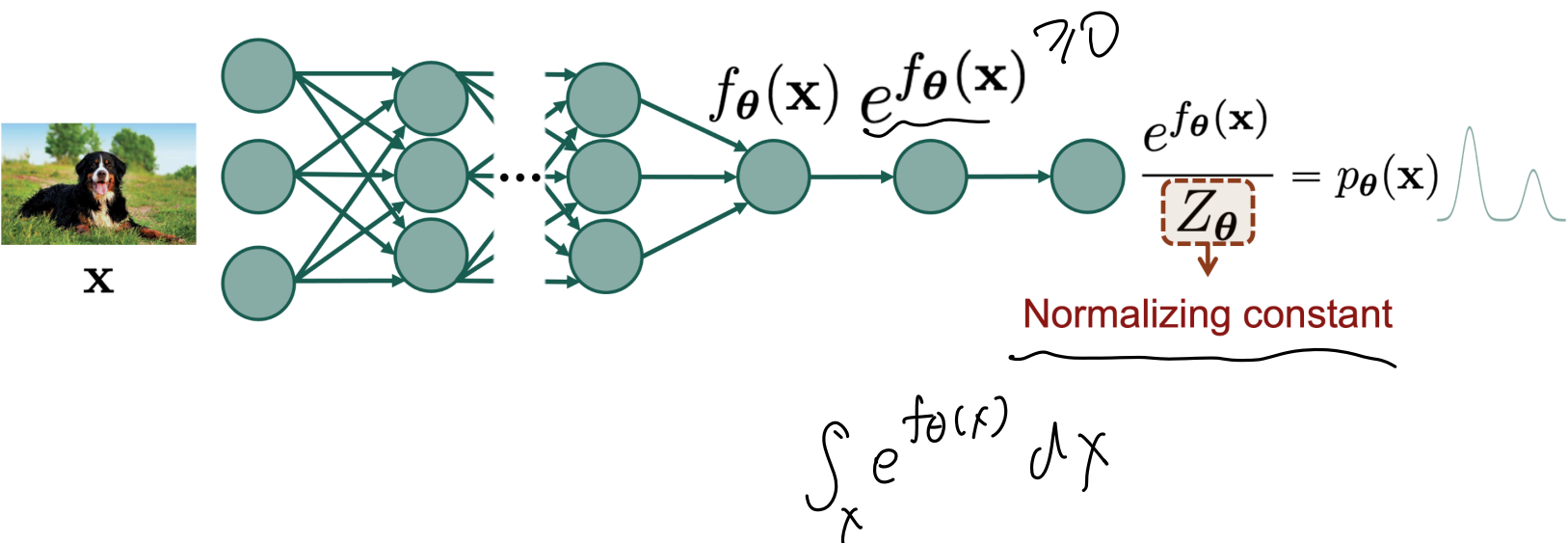
Given x , explicit expression for $p_\theta(x)$

$z \sim \mathcal{N}(0, I)$
 $x = G_\theta(z)$



Key challenge for building generative models

distribution $\left\{ \begin{array}{l} p_{\theta}(x) \geq 0 \\ \int_x p_{\theta}(x) dx = 1 \end{array} \right.$



Key challenge for building generative models

$p_\theta(x)$

Approximating the normalizing constant

- Variational auto-encoders [Kingma & Welling 2014, Rezende et al. 2014]
- Energy-based models [Ackley et al. 1985, LeCun et al. 2006]



Inaccurate probability evaluation

Using restricted neural network models

- Autoregressive models [Bengio & Bengio 2000, van den Oord et al. 2016]
- Normalizing flow models [Dinh et al. 2014, Rezende & Mohamed 2015]



Restricted model family

Generative adversarial networks (GANs)

- Model the generation process, not the probability distribution [Goodfellow et al. 2014]



Cannot evaluate probabilities

Training generative models

$$\{x_1, \dots, x_n\} \sim p$$
$$p_\theta(\cdot)$$

- **Likelihood-based:** maximize the likelihood of the data under the model (possibly using advanced techniques such as variational method or MCMC):

$$\max_{\theta} \sum_{i=1}^n \log p_{\theta}(x_i)$$

- Pros:

- **Easy training:** can just maximize via SGD.
- **Evaluation:** evaluating the fit of the model can be done by evaluating the likelihood (on test data).

- Cons:

- **Large models needed:** likelihood objective is hard, to fit well need very big model.
- **Likelihood encourages averaging:** produced samples tend to be blurrier, as likelihood encourages “coverage” of training data.

Training generative models

- **Likelihood-free:** use a **surrogate loss** (e.g., GAN) to train a discriminator to differentiate real and generated samples.
- Pros:
 - **Better objective, smaller models needed:** objective itself is learned - can result in visually better images with smaller models.
- Cons:
 - **Unstable training:** typically min-max (saddle point) problems.
 - **Evaluation:** no way to evaluate the quality of fit.

Generative Adversarial Nets



Implicit Generative Model

$$p_{\theta}(x)$$

- **Goal:** a sampler $g(\cdot)$ to generate images
- A simple generator $g(z; \theta)$:
 - $z \sim N(0, I)$
 - $x = g(z; \theta)$ deterministic transformation

Likelihood-free training:

- Given a dataset from some distribution p_{data}
- Goal: $g(z; \theta)$ defines a distribution, we want this distribution $\approx p_{data}$
- Training: minimize $D(g(z; \theta), p_{data})$
 - D is some distance metric (not likelihood)

- Key idea: Learn a differentiable D : neural net

metrics

- (1) Kullback-Leibler Divergence
KL-Div
- (2) Total Variation
- (3) Wasserstein distance
- (4) Jensen-Shannon Divergence
- (5) Integral probability metric (IPM)

GAN (Goodfellow et al., '14)

binary classification

- Parameterize the discriminator $D(\cdot; \phi)$ with parameter ϕ
- **Goal:** learn ϕ such that $D(x; \phi)$ measures how likely x is from p_{data}
 - $D(x, \phi) = 1$ if $x \sim p_{data}$
 - $D(x, \phi) = 0$ if $x \not\sim p_{data}$
 - a.k.a., a binary classifier
- GAN: use a neural network for $D(\cdot; \phi)$
- **Training:** need both negative and positive samples
 - Positive samples: just the training data *$\{x_1, \dots, x_N\} \sim p_{data}$*
 - Negative samples: use our sampler $g(z; \theta)$ (can provide infinite samples).
 $D(g(z; \theta), \phi) \approx 0$
- **Overall objectives:**
 - Generator: $\theta^* = \max_{\theta} D(g(z; \theta); \phi)$ *Given ϕ*
 - Discriminator uses MLE Training:
 $\phi^* = \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim g(\cdot)} [\log(1 - D(\hat{x}; \phi))]$ *$g \sim p_{data}$*

GAN (Goodfellow et al., '14)

- Generator $G(z; \theta)$ where $z \sim N(0, I)$

- Generate realistic data

- Discriminator $D(x; \phi)$

- Classify whether the data is real (from p_{data}) or fake (from G)

- Objective function:

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \left[\mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))] \right]$$

- Training procedure:

- Collect dataset $\{(x, 1) \mid x \sim p_{data}\} \cup \{(\hat{x}, 0) \mid \hat{x} \sim g(z; \theta)\}$

- Train discriminator

$$D : L(\phi) = \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))]$$

- Train generator $G : L(\theta) = \mathbb{E}_{z \sim N(0, I)} [\log D(G(z; \theta), \phi)]$

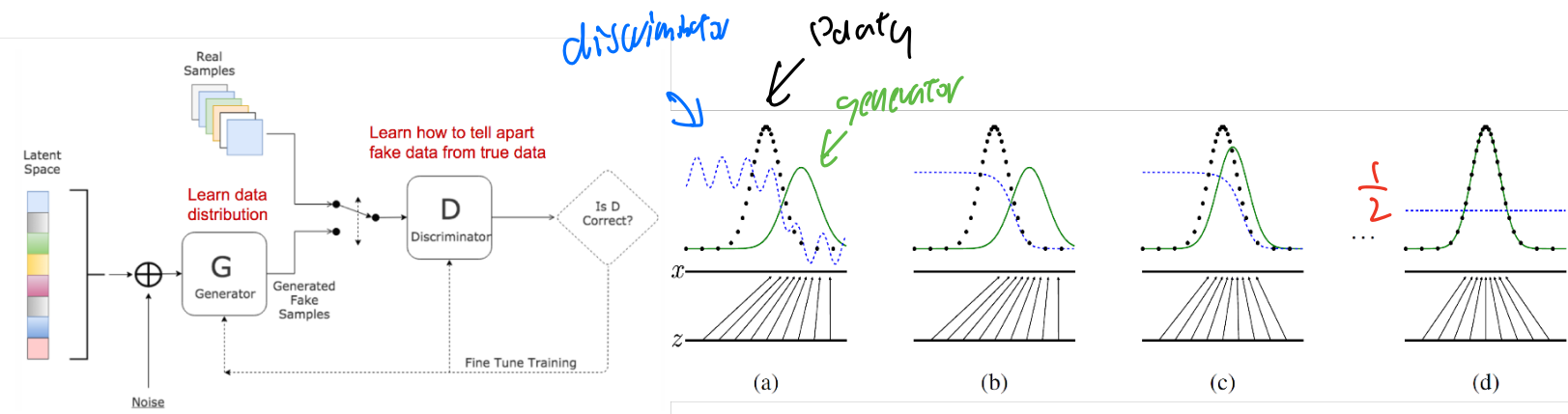
- Repeat

$1 \approx 0_{data}$
 $0 \approx 1$

GAN (Goodfellow et al., '14)

- Objective function:

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{data}} [\log D(x; \phi)] + \mathbb{E}_{\hat{x} \sim G} [\log(1 - D(\hat{x}; \phi))]$$



Math Behind GAN

$$L(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x; \phi)] + \mathbb{E}_{x \sim G(\theta)} [\log(1 - D(x; \phi))]$$

Let D^* , G^* be the solution to

optimal D^* , for a given x

$$L_x(D) = p_{\text{data}}(x) \cdot \log D(x) + P_G(x) \cdot \log(1 - D(x))$$

$$\Rightarrow D^*: \frac{\partial L}{\partial D} = 0 \quad \text{first-order condition}$$

$$\Rightarrow \frac{p_{\text{data}}(x)}{D^*(x)} - \frac{P_G(x)}{1 - D^*(x)} = 0, \quad \forall x$$

$$\Rightarrow \underline{D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + P_G(x)}}$$

$$\Rightarrow \text{if } P_G = p_{\text{data}} \Rightarrow D^*(x) = 0.5$$

Math Behind GAN

Consider optimal generator h^* given optimal D

$$\mathcal{L}(\Theta, \Phi) = \mathbb{E}_{x \sim p_{\text{data}}} \left[\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right] + \mathbb{E}_{x \sim h(\phi)} \left[\log \frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right]$$

$$= \mathbb{E}_{x \sim p_{\text{data}}} \left[\log \frac{p_{\text{data}}(x)}{\frac{p_{\text{data}}(x) + p_G(x)}{2}} \right] + \mathbb{E}_{x \sim h(\phi)} \left[\log \frac{p_G(x)}{\frac{p_{\text{data}}(x) + p_G(x)}{2}} \right] - \log 2$$

$$= \text{KL} \left(p_{\text{data}} \parallel \frac{1}{2} (p_{\text{data}} + p_G) \right) + \text{KL} \left(p_G \parallel \frac{1}{2} (p_{\text{data}} + p_G) \right)$$

2 · Jensen-Shannon Divergence (JS())

KL-Divergence and JS-Divergence

$$\bullet \text{KL}(p||q) = \mathbb{E}_{x \sim p} \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$

asymmetric

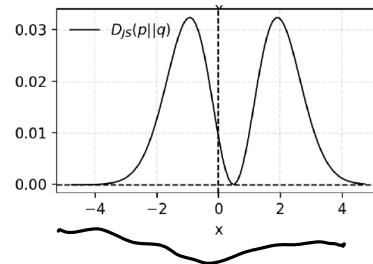
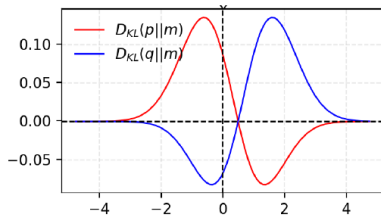
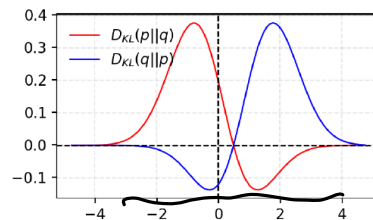
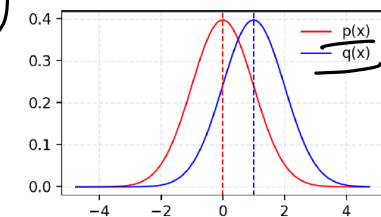
$$\text{KL}(p||q) \neq \text{KL}(q||p)$$

$$\bullet \text{JSD}(p||q)$$

$$= \frac{1}{2} \left(\text{KL} \left(p || \frac{1}{2}(p+q) \right) + \text{KL} \left(q || \frac{1}{2}(p+q) \right) \right)$$

symmetric, $\text{JSD} \geq 0$, $\text{JSD}(p||q) = 0 \Leftrightarrow p = q$

$\sqrt{\text{JSD}}$ satisfies triangle inequality



Math Behind GAN

\Rightarrow Given optimal D^*

$$\mathcal{L}(\theta) = 2 \cdot \text{JSD}(P_G \parallel \frac{1}{2}(P_G + P_{data})) - \log 4$$

goal: find P_G minimize $\mathcal{L}(\theta)$

$$\text{JSD} \geq 0$$

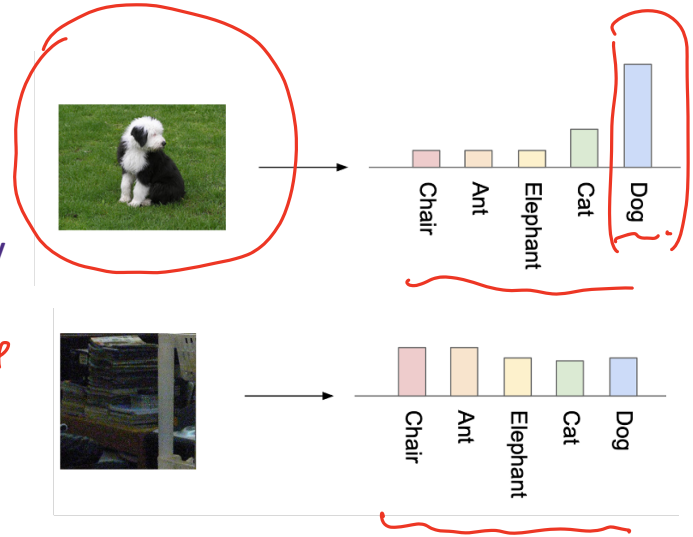
$$\Rightarrow \text{minimizer } \text{JSD}(P_G \parallel \frac{1}{2}(P_G + P_{data})) = 0$$

$$\Rightarrow P_G = P_{data}$$

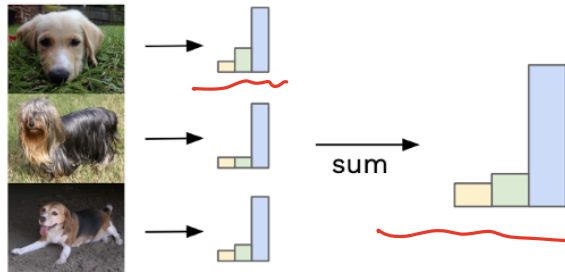
$$\cdot \mathcal{L}^* = -\log 4$$

Evaluation of GAN

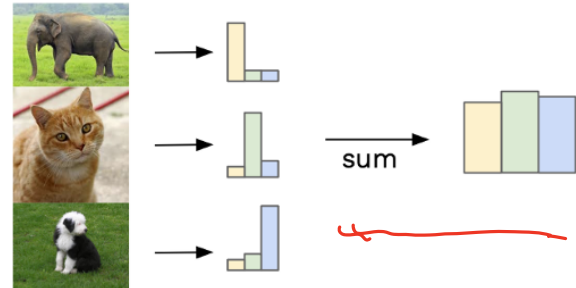
- No $p(x)$ in GAN.
- Idea: use a trained classifier $f(y | x)$:
- If $x \sim p_{data}$, $f(y | x)$ should have low entropy
 - Otherwise, $f(y | x)$ close to uniform.
- Samples from G should be diverse: *on average marginally*
 - $p_f(y) = \mathbb{E}_{x \sim G}[f(y | x)]$ close to uniform.



Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



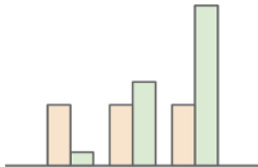
Evaluation of GAN

- **Inception Score (IS, Salimans et al. '16)**

- Use Inception V3 trained on ImageNet as $f(y|x)$

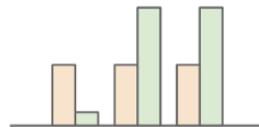
- $$IS = \exp \left(\mathbb{E}_{x \sim G} \left[\underbrace{KL(f(y|x) || p_f(y))}_{\text{low-entropy}} \right] \right)$$
- Higher the better $\underbrace{\hspace{10em}}_{\text{marginal high entropy}}$

High KL divergence



Ideal situation

Medium KL divergence



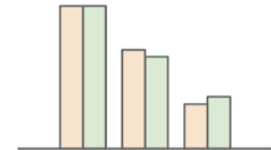
Generated images are not distinctly one label

Low KL divergence



Generated images are not distinctly one label

Low KL divergence



Generator lacks diversity

Label distribution

Marginal distribution

Comments on GAN

- Other evaluation metrics:
 - Fréchet Inception Distance (FID): Wasserstein distance between Gaussians
- Mode collapse:
 - The generator only generate a few type of samples.
 - Or keep oscillating over a few modes.
- Training instability:
 - Discriminator and generator may keep oscillating
 - Example: $-xy$, generator x , discriminatory. NE: $x = y = 0$ but GD oscillates.
 - No stopping criteria.
 - Use Wsserstein GAN (Arjovsky et al. '17):
$$\min_G \max_{f: \text{Lip}(f) \leq 1} \mathbb{E}_{x \sim p_{data}} [f(x)] - \mathbb{E}_{\hat{x} \sim p_G} [f(\hat{x})]$$
 - And need many other tricks...

