

Generative models I

Nikita Kazeev

Generative models

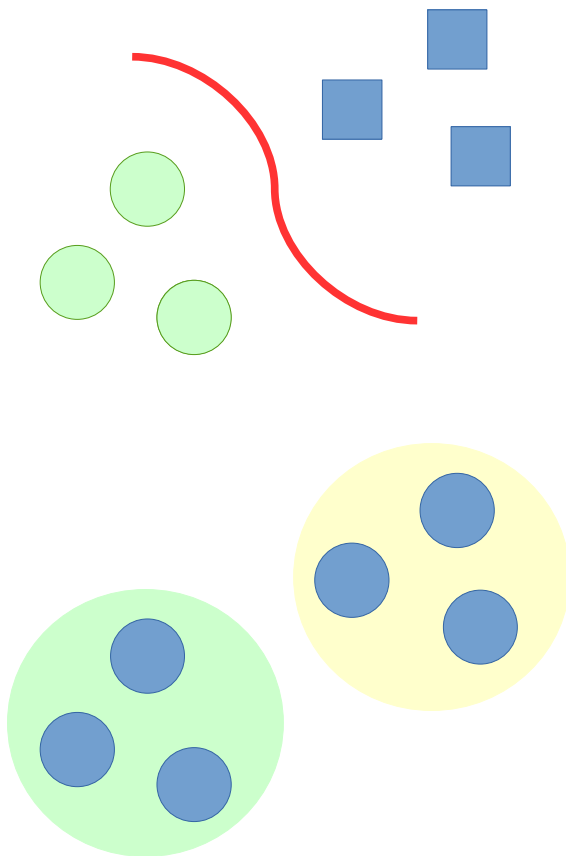
Dataset: (x, y) pairs

Regression & classification

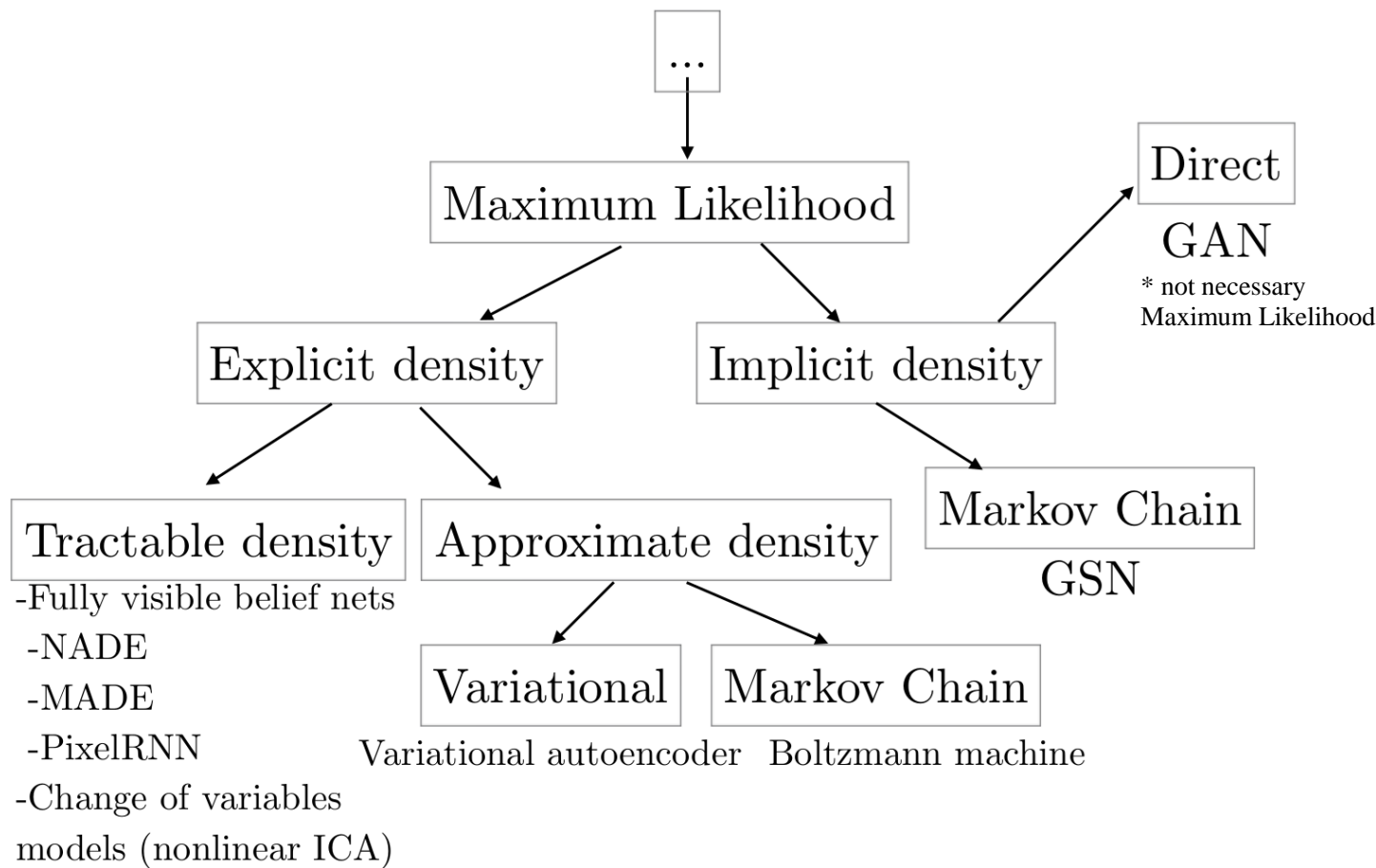
- › Learn mapping $x \rightarrow y$

Generative models

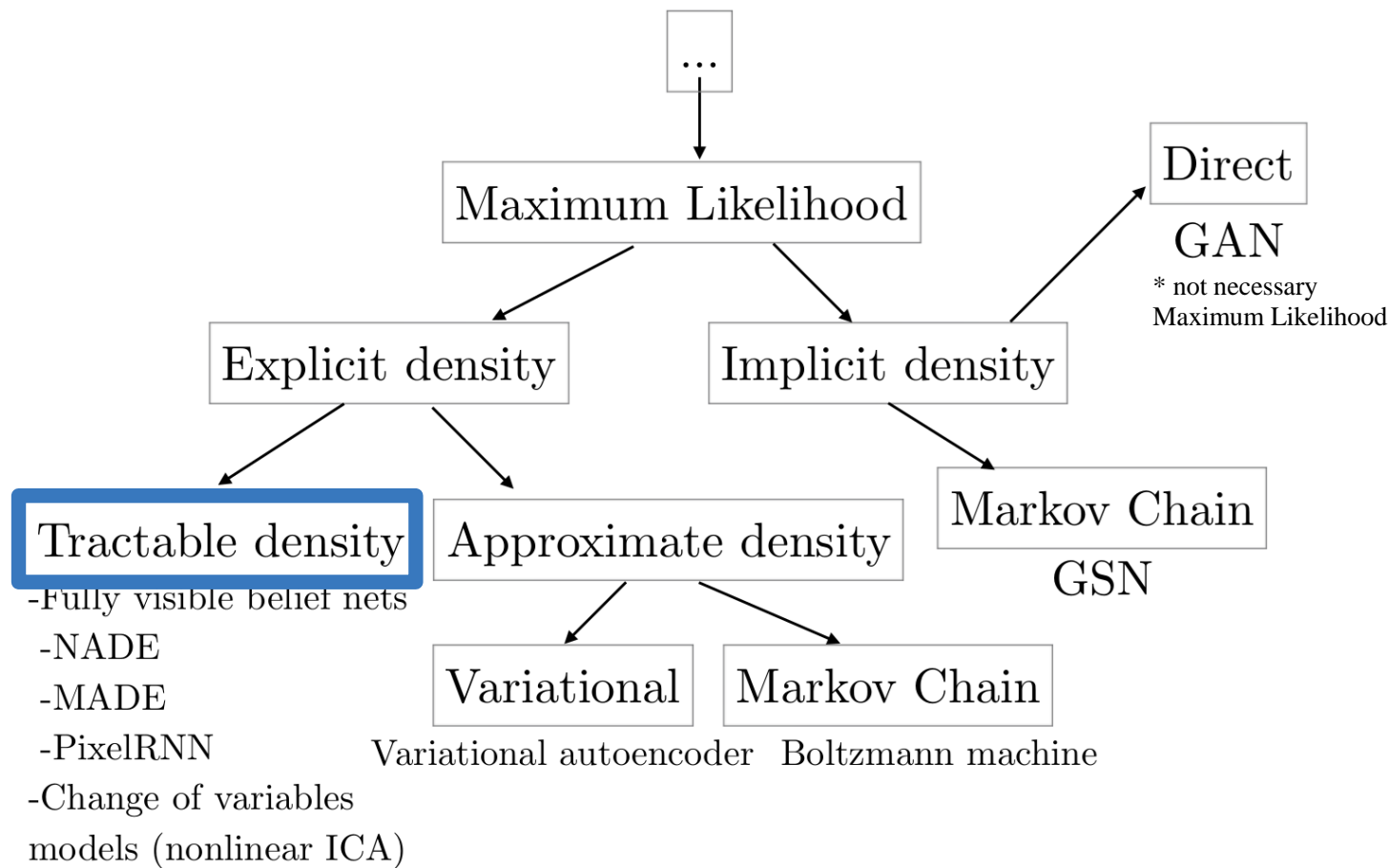
- › Learn to sample $P(y|x)$
- › X can be null, then it's unsupervised learning



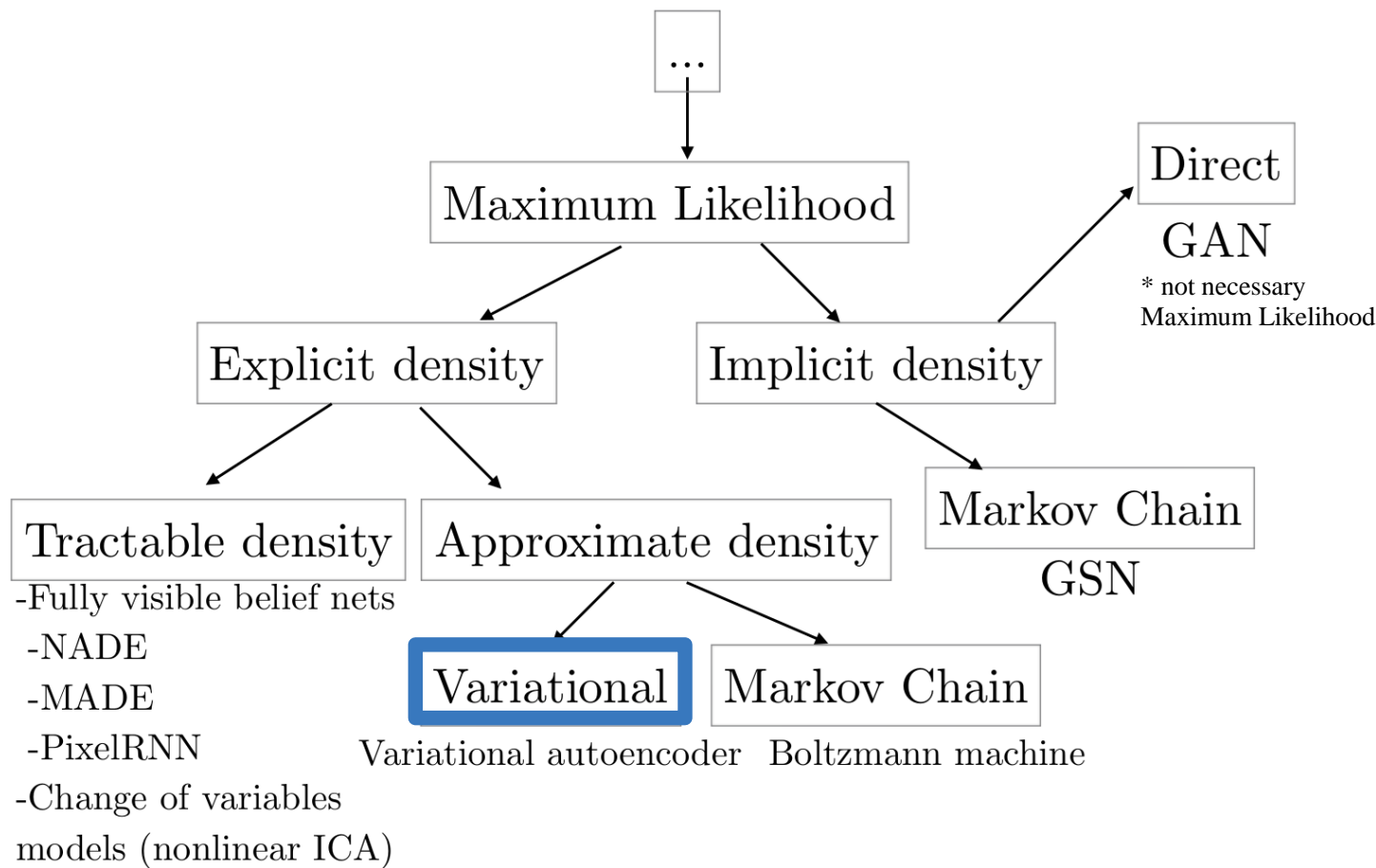
Generative Model Taxonomy



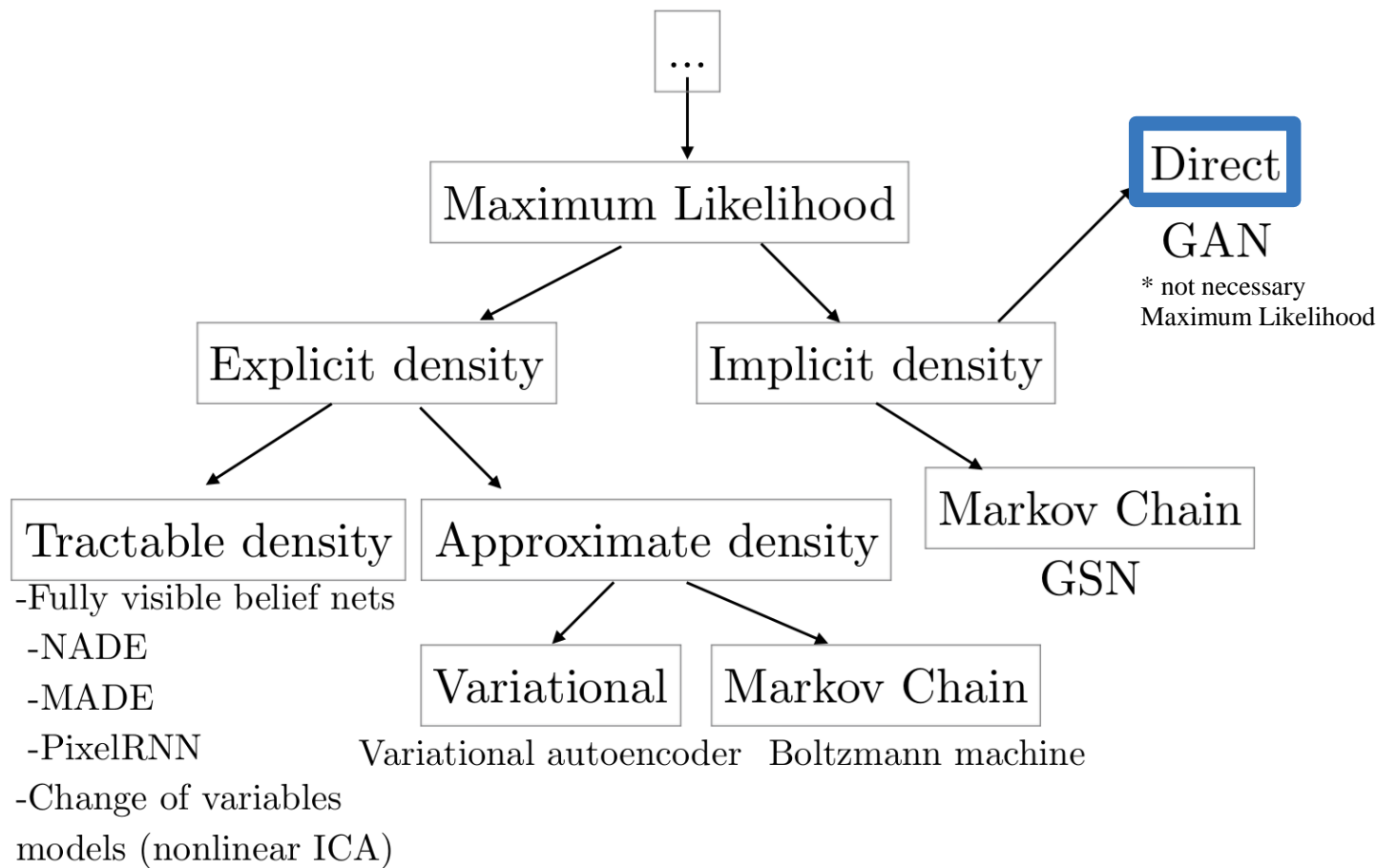
Generative Model Taxonomy



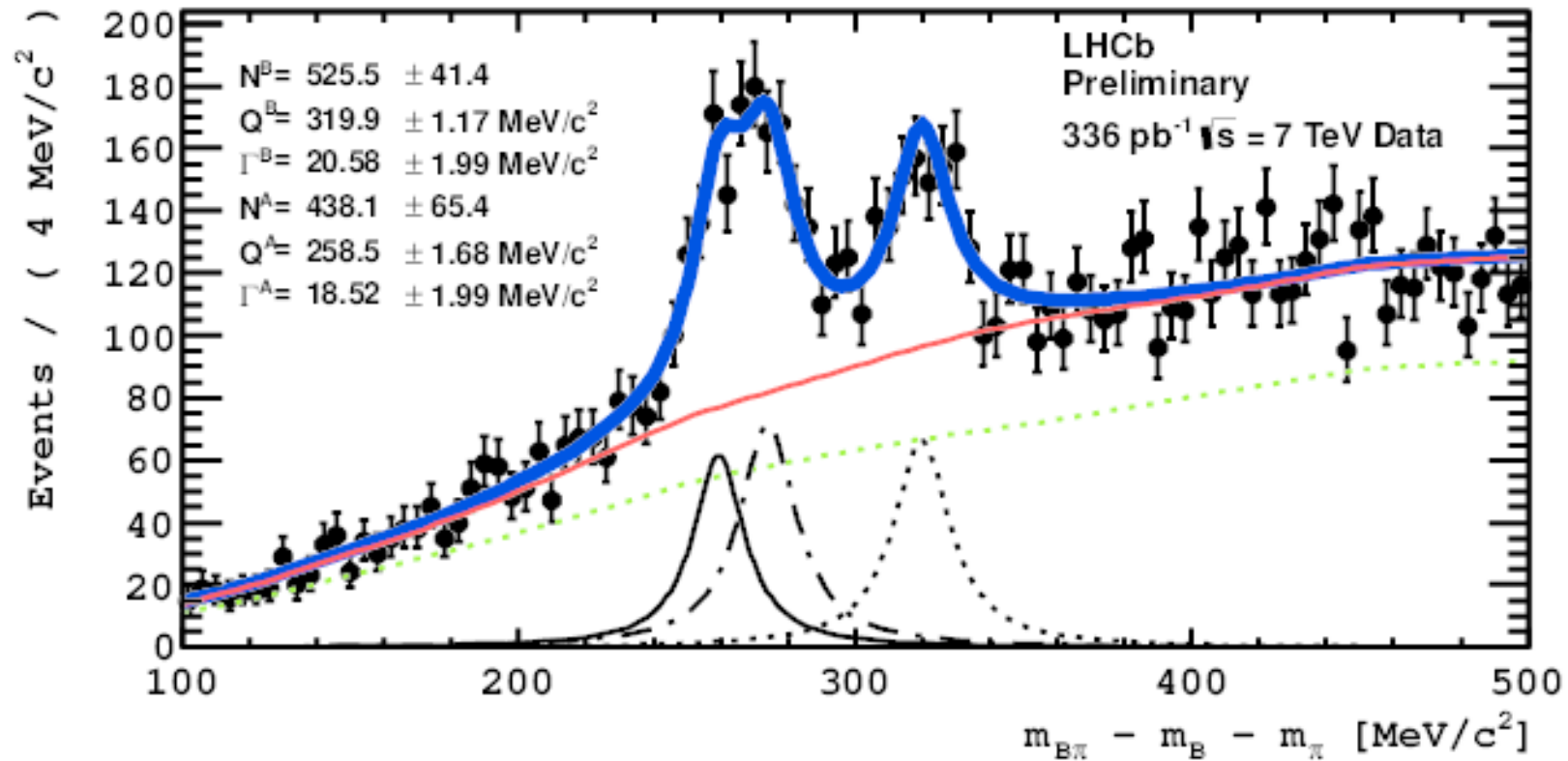
Generative Model Taxonomy



Generative Model Taxonomy



Good ol' distribution fitting



Good ol' distribution fitting

Given

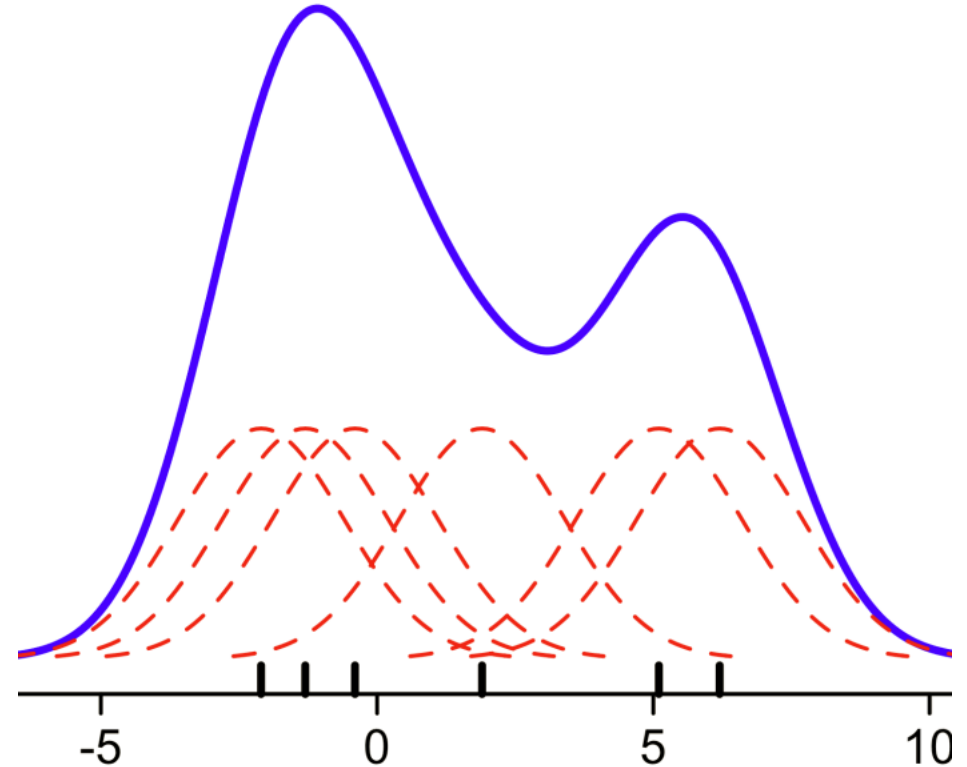
- › data points $x_1, \dots, x_n \in \mathbb{R}_m$
- › a parametrized distribution it's supposed to come from $P(x|\theta)$

Find a set of parameters to maximize the empirical likelihood:

- ›
$$\max_{\theta} L(\theta|x) = \max_{\theta} \prod_{i=1}^n P(x_i|\theta)$$

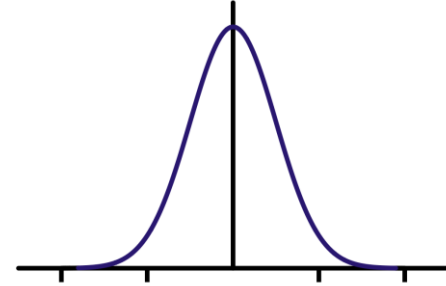
Kernel density

What if we place many Gaussians on the data points and call their sum a PDF?



Kernel density

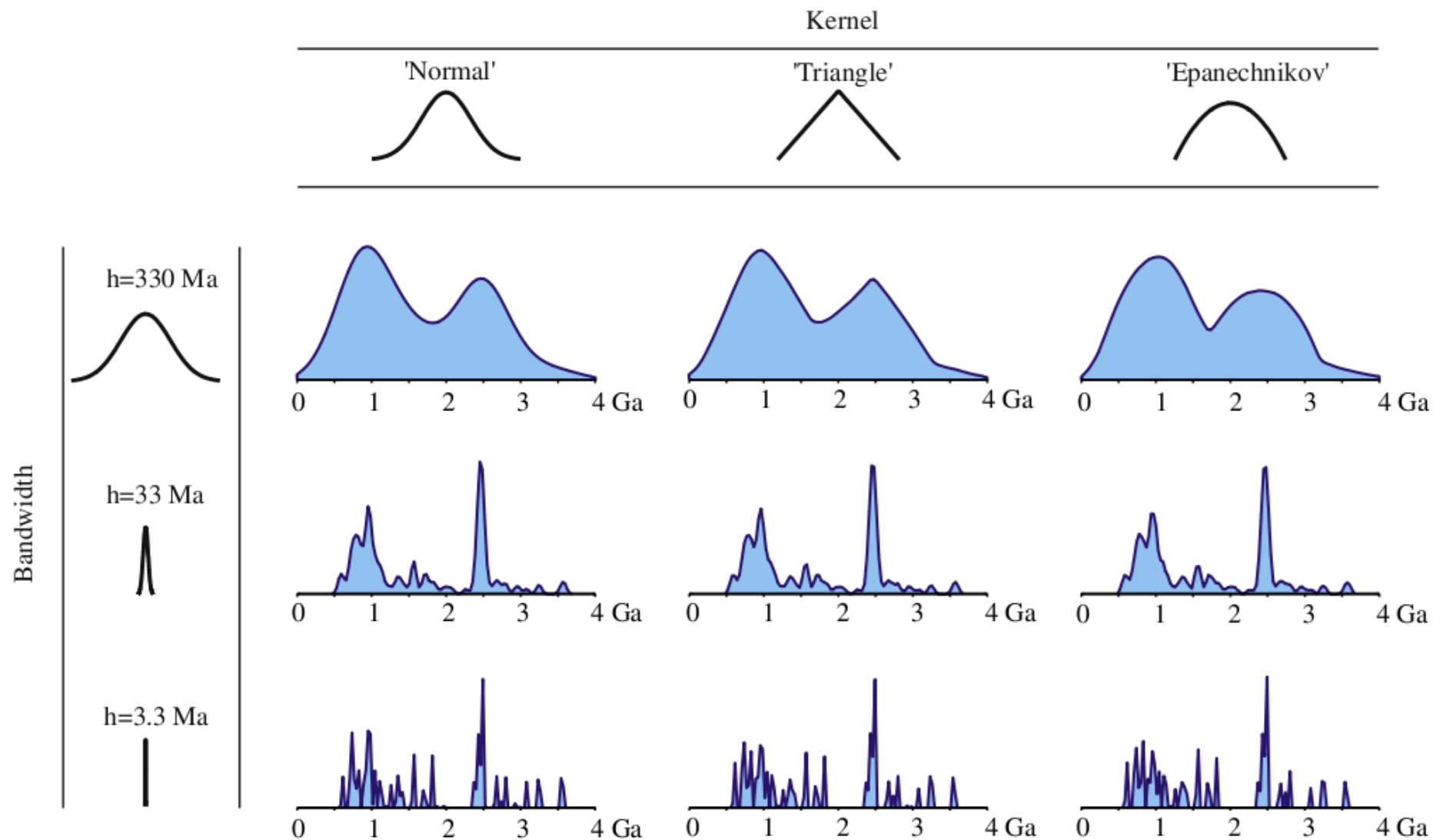
Kernel function, usually looks like this:



$$P_{\text{KDE}}(x) = \sum_i K(x - x_i) / N$$

Number of points in the training dataset

Sum over all points in the training dataset



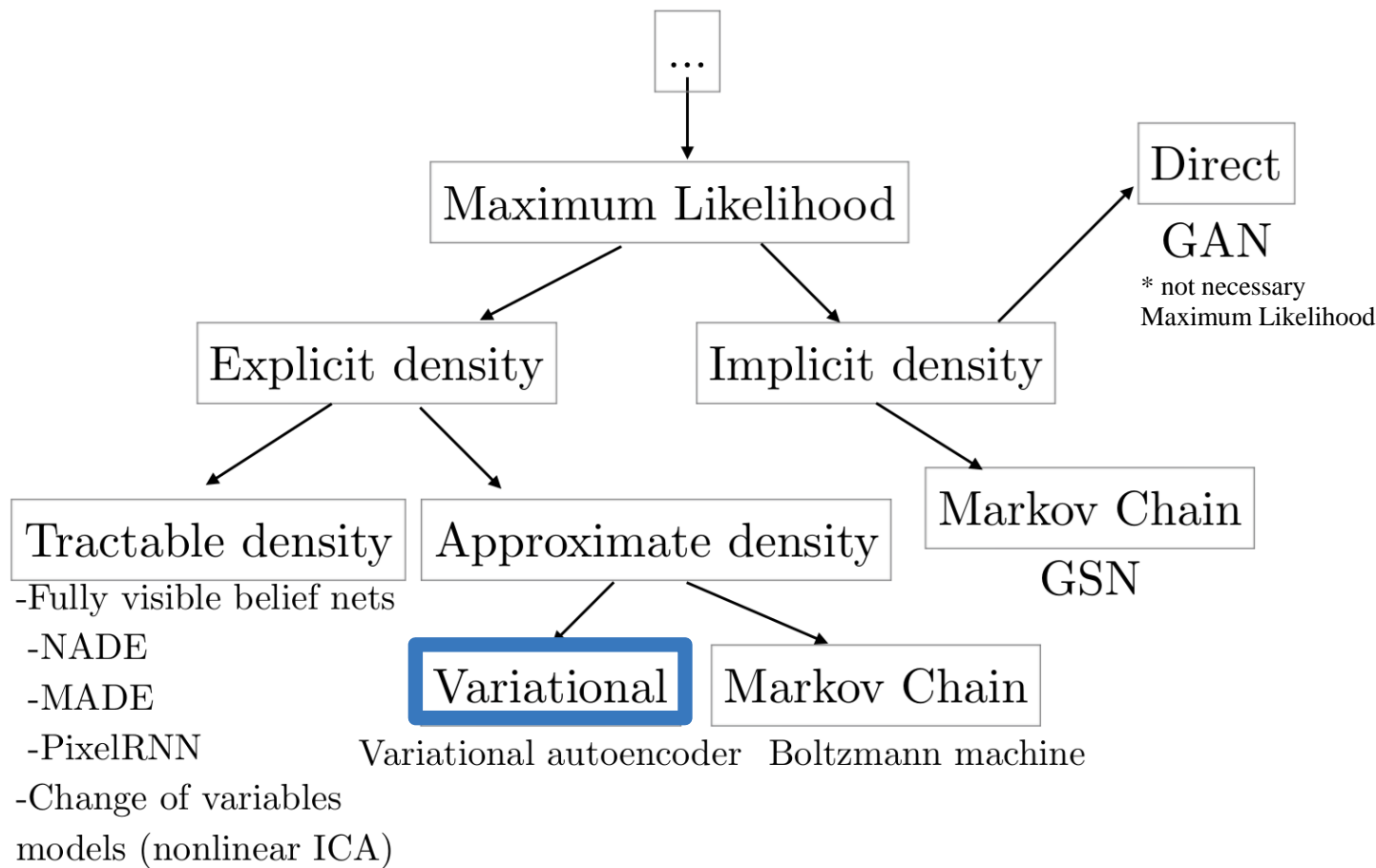
Kernel Density vs Histogram

Kernel Density	Histogram
Smooth PDF	Discrete binned PDF
With number of data points approaching infinity, the value in a point approaches the convolution of the PDF with the kernel function	With number of data points approaching infinity, the value in a bin approaches the unbiased mean PDF in that bin
No easy way to estimate the uncertainty	Straightforward uncertainty estimation of bin values
User-defined parameter: kernel shape and width	User-defined parameter: bins
Requires storing the full training dataset Finite support kernels allow for KDTree optimization	Fast, memory efficient

Kernel Density: summary

- Go-to way for an easy 1-2D PDF approximation
- Applicable in the same cases as histograms
- Has heuristic parameters: kernel shape
- Memory expensive
- Doesn't scale for high dimensions
- Nice demo: <https://mathisonian.github.io/kde/>

Generative Model Taxonomy

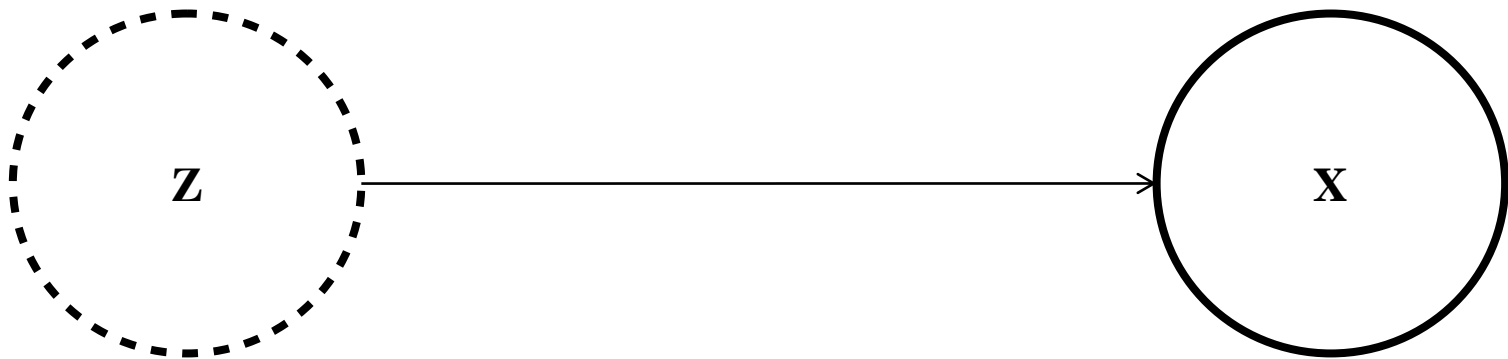


Variational autoencoders (VAE)

$P(x)$ is the distribution of the data (e. g. images)

$P(z)$ is some easily samplable distribution (e. g. Gaussian)

What if we could find some mapping $F(z)$ so that $P(F(z)) = P(x)$?



VAE

Given

- › data points $x_1, \dots, x_n \in \mathbb{R}_m$
- › a distribution $P(z)$
- › a parametrized mapping $P(x|z, \theta)$

Find a set of parameters θ to maximize the empirical likelihood:

- ›
$$L(\theta) = \prod_i P(x_i|\theta) = \prod_i \int P_\theta(x_i|z)P(z)dz$$

VAE

$$L(\theta) = \prod_i \mathbb{E}_z [P_\theta(x_i|z)]$$

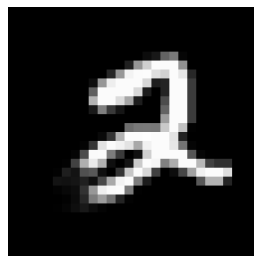
How to optimize?

Let x be any training example
Problem: for early stages

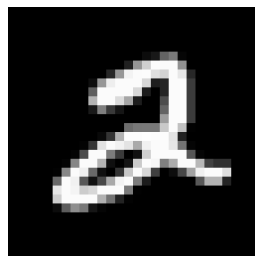
$$P_\theta(x|z)$$

is very sparse, most pixel
combinations are
not like training images

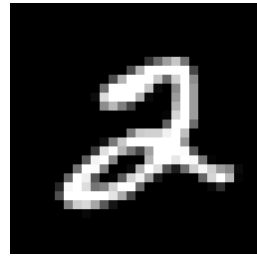
Semantic difference is not like pixel
difference. To make them aligned, we need to
be strict.



MSE=0.0387



original



MSE=0.2693

VAE: training

Let's add $Q_{\theta}(z|x)$ – distribution of z values likely produce x

If we were able to sample $Q_{\theta}(z|x)$, computation of $\mathbb{E}_{z \sim Q} P_{\theta}(x|z)$ is easy

At the same time, we are interested in the $P_{\theta}(x)$

Kullback–Leibler divergence

two functions are - the smaller the divergence, t

effectively minimising this when doing a Max Likely

$$D_{KL}(P \parallel Q) = \int \log \left(\frac{P(x)}{Q(x)} \right) P(x) dx$$

Asymmetric

Possibly infinite

Has roots in information theory

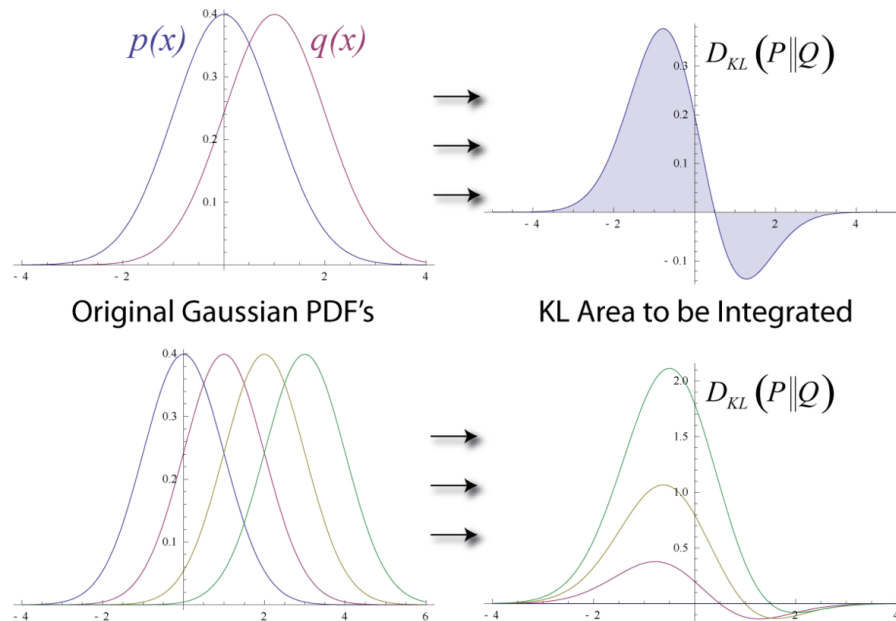


Image: Wikipedia

VAE training: KL divergence

$$D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z|x)] = \mathbb{E}_{z \sim Q}[\log Q_{\theta}(z|x) - \log P_{\theta}(z|x)]$$

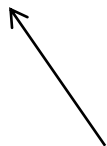
applying Bayes rule to $P(z|x)$:

$$D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z|x)] = \mathbb{E}_{z \sim Q}[\log Q_{\theta}(z|x) - \log P_{\theta}(x|z) - \log P(z)] + \log P_{\theta}(x)$$

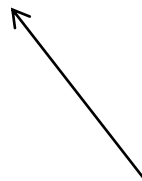
$$\log P_{\theta}(x) - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z|x)] = \mathbb{E}_{z \sim Q}[\log P_{\theta}(x|z)] - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z)]$$



We want to maximize the likelihood of the model



Minimize the error
 Q approximation



This can be optimized via SGD!

VAE: Gaussian parametrization

Let's use Gaussian Q:

$$Q_{\theta}(z|x) = N(z|\mu_{\theta}(x), \sigma_{\theta}(x))$$

This will be an NN



This too



VAE: Gaussian parametrization

Let's use Gaussian Q:

$$Q_{\theta}(z|x) = N(z|\mu_{\theta}(x), \sigma_{\theta}(x))$$

This will be an NN



This too

And Gaussian P(z):

$$P(z) = N(z|0,1)$$

This way D is analytically computable:

$$-D_{\text{KL}}(Q_{\theta}(z|x) \parallel P_{\theta}(z)) = \frac{1}{2} \sum (1 + \log(\sigma_{\theta}^2(x)) - \mu_{\theta}^2(x) - \sigma_{\theta}^2(x))$$

VAE training

$$\log P_{\theta}(x) - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z|x)] = \mathbb{E}_{z \sim Q}[\log P_{\theta}(x|z)] - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z)]$$

$$-D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z)] = \frac{1}{2} \sum (1 + \log(\sigma_{\theta}^2(x)) - \mu_{\theta}^2(x) - \sigma_{\theta}^2(x))$$

$$\mathbb{E}_{z \sim Q}[\log P_{\theta}(x|z)] = ?$$

VAE training

$$\log P_{\theta}(x) - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z|x)] = \mathbb{E}_{z \sim Q}[\log P_{\theta}(x|z)] - D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z)]$$

$$-D_{\text{KL}}[Q_{\theta}(z|x) \parallel P_{\theta}(z)] = \frac{1}{2} \sum (1 + \log(\sigma_{\theta}^2(x)) - \mu_{\theta}^2(x) - \sigma_{\theta}^2(x))$$

Gaussian?

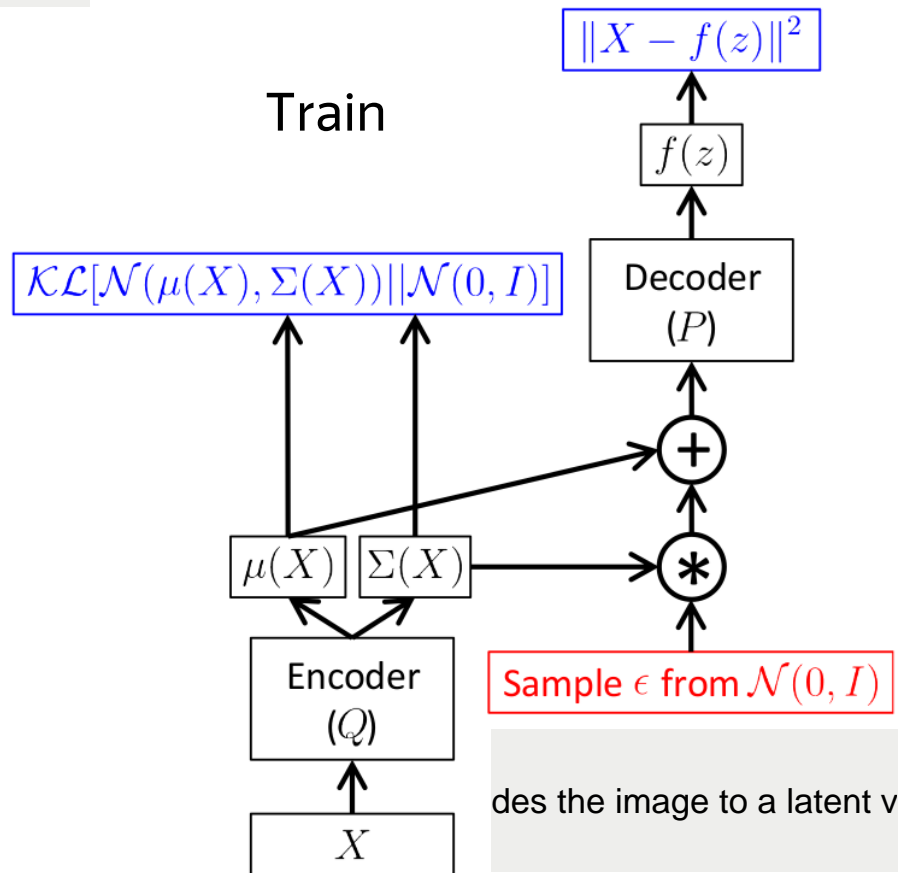
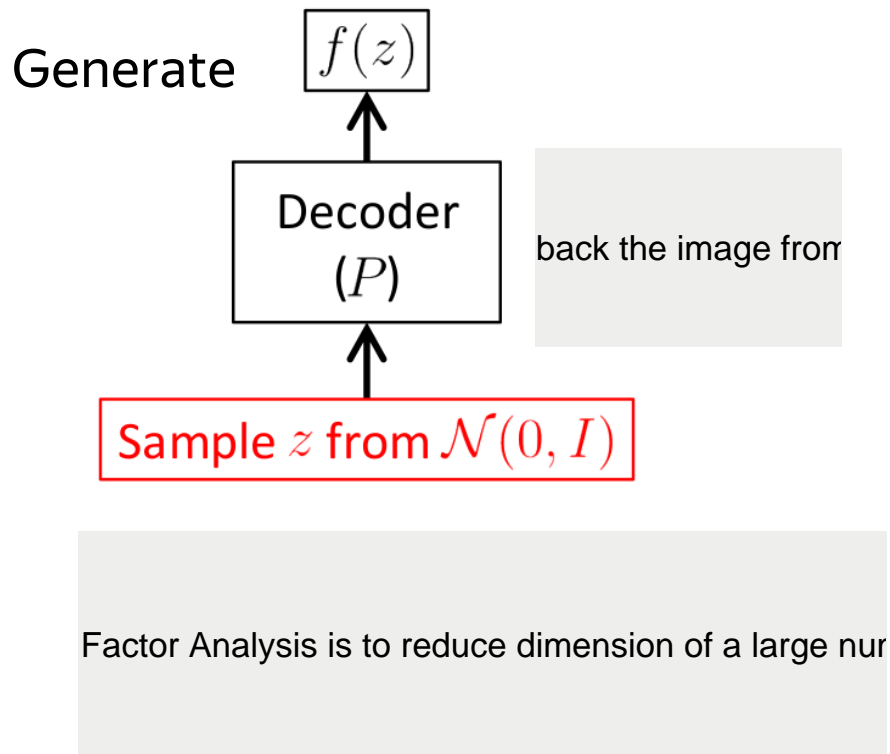
$$\mathbb{E}_{z \sim Q}[\log P_{\theta}(x|z)] = \mathbb{E}_{z \sim Q}[\log N(x|f_{\theta}(z), h^2)] = C - \frac{1}{2} \|x - f_{\theta}(z)\|^2 / h^2$$

Constant, doesn't depend on f_{θ}

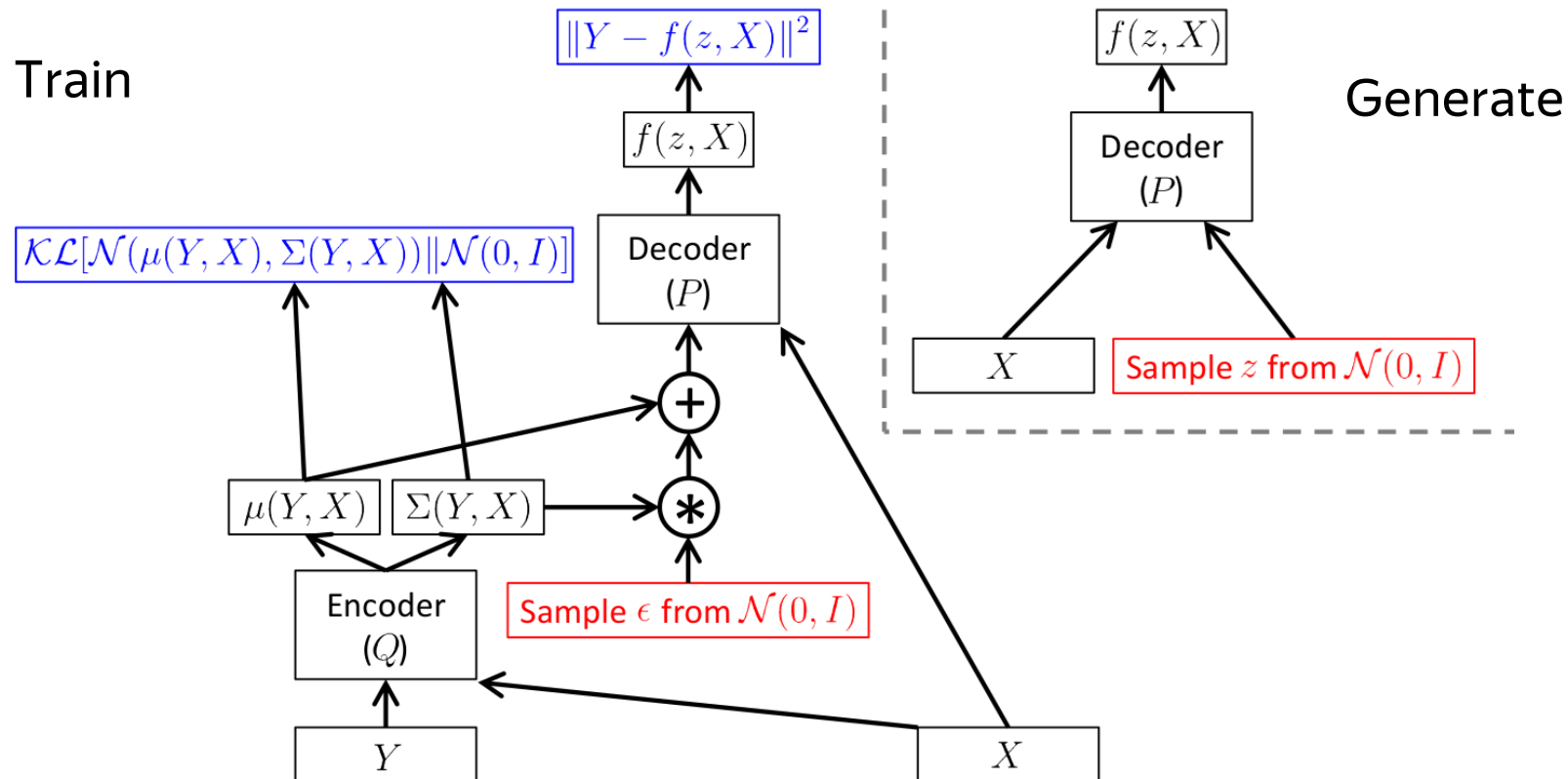
Parameter

VAE as feedforward NN

** see: <http://kvfrans.com/variational-autoencoders-explained/>



Conditional VAE



VAE: summary

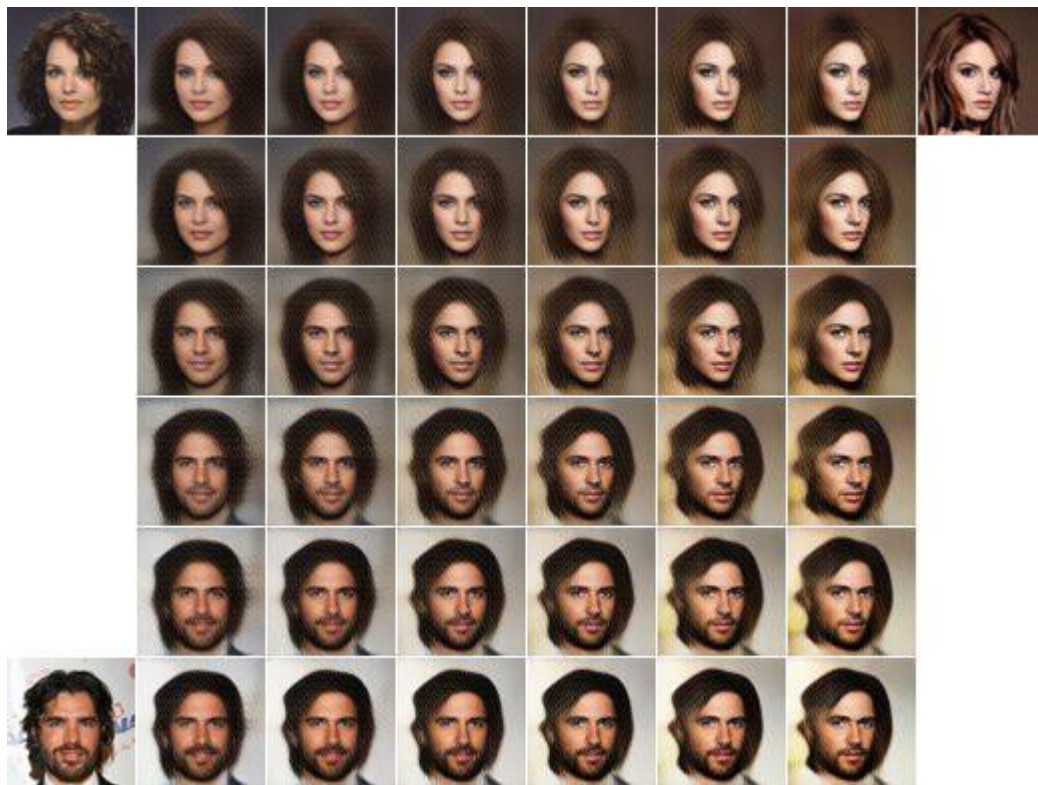
General-purpose generative model

Lots of image generation fun

Latent space fun

The Gaussian approximations may be too rigid

Subjectively worse sample image quality than GAN



Feel free to drop a line

Nikita Kazeev

HSE, Rome Sapienza, YSDA, trace amounts of
Yandex proper



kazeevn@yandex-team.ru



[telegram.me/kazeevn](https://t.me/kazeevn)

Backup



What is probability?

DID THE SUN JUST EXPLODE?
(IT'S NIGHT, SO WE'RE NOT SURE.)

THIS NEUTRINO DETECTOR MEASURES
WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY
BOTH COME UP SIX, IT LIES TO US.
OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.

DETECTOR! HAS THE
SUN GONE NOVA?

ROLL
YES.



FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT
HAPPENING BY CHANCE IS $\frac{1}{36} = 0.027$.

SINCE $p < 0.05$, I CONCLUDE
THAT THE SUN HAS EXPLODED.



BAYESIAN STATISTICIAN:

BET YOU \$50
IT HASN'T.

