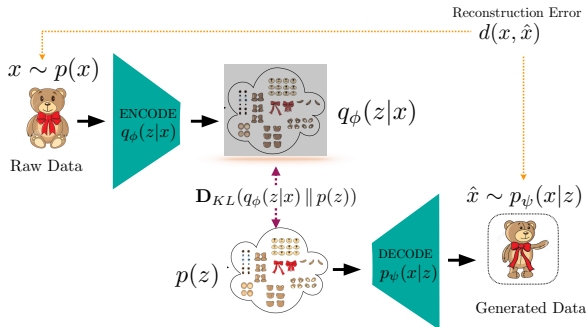


# 11-695: AI Engineering Variational AutoEncoder II

LTI/SCS

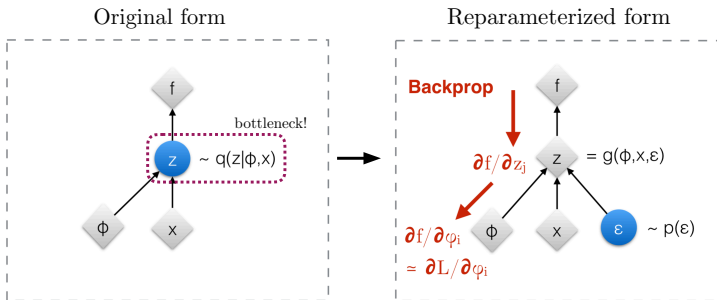
Spring 2020

- 1 Variational AutoEncoder: Review
- 2 Coding VAE
- 3 Some Results
- 4 Tuning VAE (Active Research Topic)



- Still AutoEncoder, with KL term for variational  $\rightarrow$  VAE
- Note: Decoding phase is from  $p(z)$ , and not the auto-encoded probability  $q(z|x)$
- ELBO =  $\mathbb{E}_{q_\phi(z|x)}[\log p(x|z)] - \mathbf{D}_{KL}(q_\phi(z|x) \parallel p(z))$
- $p(z)$  is random, so now we can actually *generate* data

# But, not yet fully differentiable



- Sampling  $z \sim q_\phi(z|x) = q(z|x, \phi)$  is not continuous
- Trick: push the sampling to the input level of the graph
  - $\epsilon \sim p(\epsilon) = \mathbf{N}(\mathbf{0}, \mathbf{I})$
  - $z = g_\phi(z, \epsilon)$ , a differentiable neural network now!

- 1 Variational AutoEncoder: Review
- 2 Coding VAE
- 3 Some Results
- 4 Tuning VAE (Active Research Topic)

## encoder.py

```
1 encoder_net = tf.keras.Sequential([
2     conv(base_depth, 5, 1),
3     conv(base_depth, 5, 2),
4     conv(2 * base_depth, 5, 1),
5     conv(2 * base_depth, 5, 2),
6     conv(4 * latent_size, 7, padding="VALID"),
7     tf.keras.layers.Flatten(),
8     tf.keras.layers.Dense(2 * latent_size, activation=None),
9 ])
10
11 def encoder(images):
12     images = 2 * tf.cast(images, dtype=tf.float32) - 1
13     net = encoder_net(images)
14     return tfd.MultivariateNormalDiag(
15         loc=net[..., :latent_size], # mean
16         scale_diag=tf.nn.softplus(net[..., latent_size:] + # variance
17                                   _softplus_inverse(1.0)),
18         name="code")
```

- New Module: [Tensorflow Probability](#) [tfd.distributions.MultivariateNormalDiag](#)

- For simple cases, can use Feedforward NNs instead of CNN
- Tying  $\mu$  and  $\sigma$  weights, each is half of the output.
- Other option not using TFD:
  - $\epsilon \sim p(\epsilon) = \mathbf{N}(\mathbf{0}, \mathbf{I})$
  - Then  $f(\phi, \epsilon) = \mu + \sigma * \epsilon$

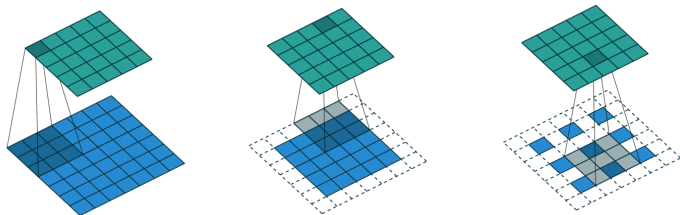
- For simple cases, can use Feedforward NNs instead of CNN
- Tying  $\mu$  and  $\sigma$  weights, each is half of the output.
- Other option not using TFD:
  - $\epsilon \sim p(\epsilon) = \mathbf{N}(\mathbf{0}, \mathbf{I})$
  - Then  $f(\phi, \epsilon) = \mu + \sigma * \epsilon$
- Usually we model  $\log(\sigma)$  instead of  $\sigma$ 
  - Then  $f(\phi, \epsilon) = \mu + e^{\log(\sigma)} * \epsilon$



## decoder.py

```
1 decoder_net = tf.keras.Sequential([
2     deconv(2 * base_depth, 7, padding="VALID"),
3     deconv(2 * base_depth, 5),
4     deconv(2 * base_depth, 5, 2),
5     deconv(base_depth, 5),
6     deconv(base_depth, 5, 2),
7     deconv(base_depth, 5),
8     conv(output_shape[-1], 5, activation=None),
9 ])
10
11 def decoder(codes):
12     original_shape = tf.shape(input=codes)
13     # Collapse the sample and batch dimension and convert to rank-4 tensor for
14     # use with a deconv network.
15     codes = tf.reshape(codes, (-1, 1, 1, latent_size))
16     logits = decoder_net(codes)
17     logits = tf.reshape(
18         logits, shape=tf.concat([original_shape[:-1], output_shape], axis=0))
19     return tfd.Independent(tfd.Bernoulli(logits=logits),
20                           reinterpreted_batch_ndims=len(output_shape),
21                           name="image")
```

- New Module: `tfd.distributions.Independent`



- We need to use deconv op (`tf.keras.layers.Conv2DTranspose`)
- Some demos of ConvTranspose: [▶ 1](#) [▶ 2](#) [▶ 3](#)
- Apply Bernoulli distribution to decoder's output

## prior.py

```
1 if mixture_components == 1:
2     # See the module docstring for why we don't learn the parameters here.
3     return tfd.MultivariateNormalDiag(
4         loc=tf.zeros([latent_size]),
5         scale_identity_multiplier=1.0)
6
7 loc = tf.compat.v1.get_variable(
8     name="loc", shape=[mixture_components, latent_size])
9 raw_scale_diag = tf.compat.v1.get_variable(
10    name="raw_scale_diag", shape=[mixture_components, latent_size])
11 mixture_logits = tf.compat.v1.get_variable(
12    name="mixture_logits", shape=[mixture_components])
13
14 return tfd.MixtureSameFamily(
15     components_distribution=tfd.MultivariateNormalDiag(
16         loc=loc,
17         scale_diag=tf.nn.softplus(raw_scale_diag)),
18     mixture_distribution=tfd.Categorical(logits=mixture_logits),
19     name="prior")
```

- New Module: `tfd.distributions.MixtureSameFamily`



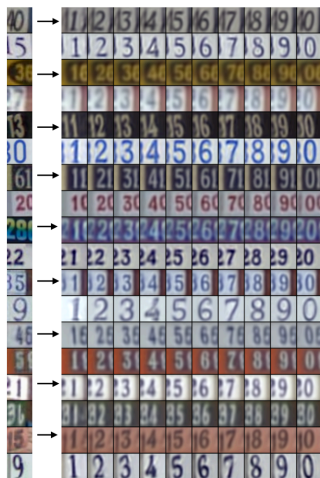
- Remember in Bayesian Statistics what is the role of prior?
- Here: mixture of Gaussian

## model.py

```
1 approx_posterior = encoder(features)
2 approx_posterior_sample = approx_posterior.sample(params["n_samples"])
3 decoder_likelihood = decoder(approx_posterior_sample)
4
5 # 'distortion' is just the negative log likelihood.
6 distortion = -decoder_likelihood.log_prob(features)
7 avg_distortion = tf.reduce_mean(input_tensor=distortion)
8
9 rate = tfd.kl_divergence(approx_posterior, latent_prior)
10 avg_rate = tf.reduce_mean(input_tensor=rate) # KL term
11
12 # total loss: NLL + KL
13 elbo_local = -(rate + distortion)
14 elbo = tf.reduce_mean(input_tensor=elbo_local)
15 loss = -elbo
16
17 # train_op
18 optimizer = tf.compat.v1.train.AdamOptimizer(learning_rate)
19 train_op = optimizer.minimize(loss, global_step=global_step)
```

- KL Divergence distance: [tfp.distributions.kl\\_divergence](#)

- 1 Variational AutoEncoder: Review
- 2 Coding VAE
- 3 Some Results
- 4 Tuning VAE (Active Research Topic)



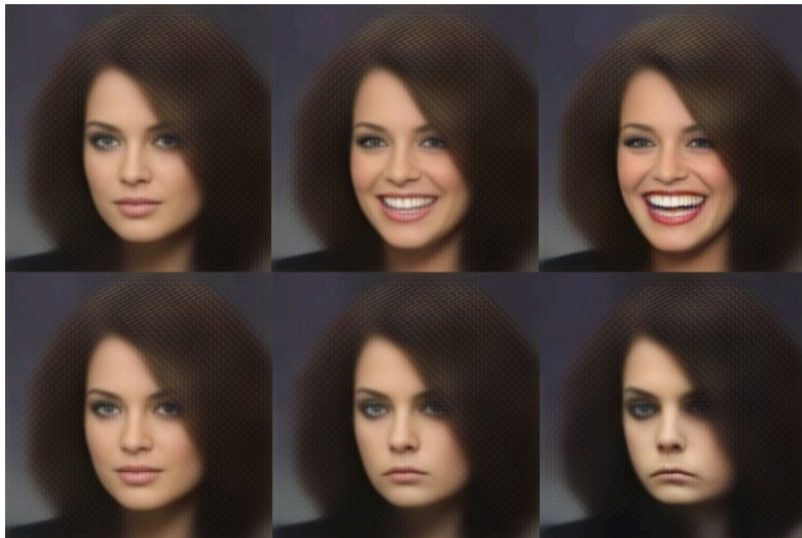
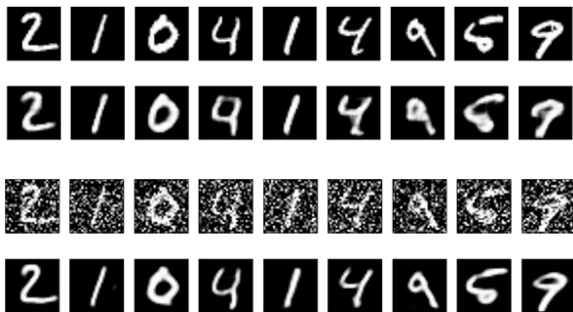


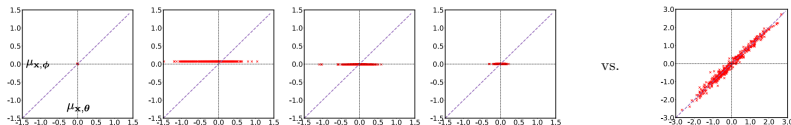
Image credit: Tom White



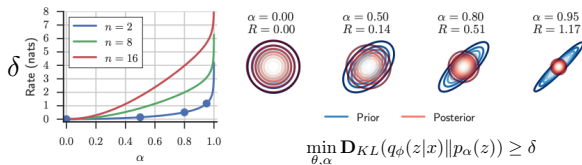


- ▶ Demo (D. P. Kingma)
- ▶ VAE with keras
- ▶ VAE with tf
- ▶ Convolutional VAE with tf 2
- ▶ VAE with tensorflow\_probability

- 1 Variational AutoEncoder: Review
- 2 Coding VAE
- 3 Some Results
- 4 Tuning VAE (Active Research Topic)



- Variational posterior  $q(z|x)$  collapse to prior  $p(z)$  for  $\forall x$ :  $\mathbf{D}_{KL} \approx 0$
- Can also happen:  $p(z|x) \approx p(z) \forall x$
- Also means the model ignores the latent features, producing outputs not resembling inputs.



- Anneal weight for  $\mathbf{D}_{KL}$  (Sam Bowman et al. 2015)
- Learn (dynamic) prior and discretize latent space (Aaron van den Oord et al. 2017)
- Use autoregressive flow to learn prior and decoding (Xi Chen et al 2017)
- Carefully init variational params and refine (Yoon Kim et al 2018)
- Aggressively update encoder more (Junxian He et al. 2019)
- $\delta$ -VAE: Lower bound  $\mathbf{D}_{KL}$  (Ali Razavi et al. 2019)

Image credit: Ali Razavi et al. 2019