

# Deep Learning Midterm

April 1, 2020

Total number of points is 105. Maximum score is achieved with any score between 100 and 105 points.

**Tejaishwarya Gagadam, tg1779**

# 1 LSTMs don't solve vanishing gradient [10 pts]

## Background: RNNs

Consider initial state  $\mathbf{h}_0$  and input sequence  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ . The RNN state  $\mathbf{h}_t$  is a function of the previous hidden state and current input, usually

$$\mathbf{h}_t = \sigma(\mathbf{W}_1 \mathbf{h}_{t-1} + \mathbf{W}_2 \mathbf{x}_t)$$

Usually, loss  $\mathcal{L}_t$  at time  $t$  is a function of the state  $\mathbf{h}_t$ . We minimise total loss  $\mathcal{L} = \sum_t \mathcal{L}_t$ . The gradient  $\partial \mathcal{L} / \partial \mathbf{W} = \sum_t \partial \mathcal{L}_t / \partial \mathbf{W}$  and  $\partial \mathcal{L}_\tau / \partial \mathbf{W}$  depends on  $\partial \mathbf{h}_\tau / \partial \mathbf{h}_t$  for all  $\tau > t$ , where  $\partial \mathbf{h}_\tau / \partial \mathbf{h}_t$  is what vanishes or explodes due to repeated application of matrix  $\mathbf{W}_1$ . We say that RNNs can't learn long dependencies when gradients vanish or explode from  $\mathcal{L}_\tau$  back to  $\mathbf{h}_t$  (and therefore to  $\mathbf{W}_1$ ) for  $\tau \gg t$ , due to vanishing or exploding in  $\partial \mathbf{h}_\tau / \partial \mathbf{h}_t$ .

## Background: LSTMs

The LSTM maintains an auxiliary cell state  $\mathbf{c}$  that helps copy the hidden state  $\mathbf{h}$  forward in time as  $\mathbf{x}$  is processed. Letting  $\mathbf{f}_t, \mathbf{i}_t, \tilde{\mathbf{c}}_t$  and  $\mathbf{o}_t$  be nonlinear functions of  $\mathbf{h}_{t-1}$  and  $\mathbf{x}_t$  and letting  $\odot$  denote element-wise product, the LSTM hidden state computation for  $\mathbf{c}$  and  $\mathbf{h}$  is defined by

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t; \quad (1)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \quad (2)$$

In the update of the cell, the *don't forget* gate  $\mathbf{f}_t$  determines the influence of previous cell states on the next. The *input* gate  $\mathbf{i}_t$  determines the influence of new information on the cell.  $\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t$  and  $\tilde{\mathbf{c}}_t$  are defined via usual RNN recurrences:

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_{gh} \mathbf{h}_{t-1} + \mathbf{W}_{gx} \mathbf{x}_t + \mathbf{b}_g);$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fh} \mathbf{h}_{t-1} + \mathbf{W}_{fx} \mathbf{x}_t + \mathbf{b}_f);$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_{ih} \mathbf{h}_{t-1} + \mathbf{W}_{ix} \mathbf{x}_t + \mathbf{b}_i);$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{oh} \mathbf{h}_{t-1} + \mathbf{W}_{ox} \mathbf{x}_t + \mathbf{b}_o).$$

## Do LSTMs Still Have Vanishing Gradients?

Assume that losses  $\mathcal{L}_t$  at each time step are a function of hidden state  $\mathbf{h}_t$  (eq. 2). Consider the gradient  $\partial \mathcal{L}_t / \partial \mathbf{W}$  for all eight matrices  $\mathbf{W}$ . Because there is a path  $\mathbf{W} \rightarrow \mathbf{c} \rightarrow \mathbf{h} \rightarrow \mathcal{L}_t$ , the cell update (eq. 1) affects the gradients  $\partial \mathcal{L}_t / \partial \mathbf{W}$ . The gradient for  $\mathbf{W}$  is accumulated through both terms in the sum of (eq. 1).

- **Answer:** From the LSTM's auxiliary cell state definition: (eq. 1), we can see that the term  $\mathbf{f}_t \odot \mathbf{c}_{t-1}$  contains the *don't forget gate*. This gate helps the system decide at every time step, if the previous state's information should be remembered or forgotten. This term is credited to solving the vanishing gradient problem in Recurrent Neural Networks.

Back-propagation through time shows:

$$\begin{aligned}\frac{\partial \mathcal{L}_t}{\partial \mathbf{W}} &= \frac{\partial \mathcal{L}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} \cdots \frac{\partial \mathbf{c}_2}{\partial \mathbf{c}_1} \frac{\partial \mathbf{c}_1}{\partial \mathbf{W}} \\ &= \frac{\partial \mathcal{L}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{c}_t} \left[ \prod_{t=2} \frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} \right] \frac{\partial \mathbf{c}_1}{\partial \mathbf{W}}\end{aligned}$$

The term in the square brackets, as shown in the equation above, causes the gradients to either explode or vanish. In LSTM, due to the auxiliary cell state and the hidden state, as seen in equations (eq. 1) and (eq. 2) this term in the square bracket will yield an additive gradient. Since they are additive in nature, the gradients will not explode/vanish.

$$\begin{aligned}\frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} &= \sigma'(\mathbf{W}[\mathbf{h}_{t-1}, \mathbf{x}_t]) \mathbf{W} \mathbf{o}_{t-1} \odot \tanh'(\mathbf{c}_{t-1}) \mathbf{c}_{t-1} \\ &\quad + \mathbf{f}_t \\ &\quad + \sigma'(\mathbf{W}[\mathbf{h}_{t-1}, \mathbf{x}_t]) \mathbf{W} \mathbf{o}_{t-1} \odot \tanh'(\mathbf{c}_{t-1}) \tilde{\mathbf{c}}_t \\ &\quad + \sigma'(\mathbf{W}[\mathbf{h}_{t-1}, \mathbf{x}_t]) \mathbf{W} \mathbf{o}_{t-1} \odot \tanh'(\mathbf{c}_{t-1}) \mathbf{i}_t\end{aligned}$$

However, if one of the four additive terms, as seen in the equation above, is more dominating than the others, the larger terms magnitude will outweigh the others. So the resultant will resemble the method that vanilla RNN takes to backpropagate through time, which means that there is still a possibility that this term might vanish/blow up!

## 2 Optimisation [25 pts]

Consider minimising the following loss function over  $\theta$ :

$$L(\hat{\mathbf{Y}}_\theta(\mathbf{X}), \mathbf{y}) = \frac{1}{N} \sum_i \ell(\hat{\mathbf{y}}_\theta(x_i), y_i)$$

For example, it could be multi-class classification where  $\hat{\mathbf{y}}_\theta(x_i)$  is a vector of probabilities for the output classes.

(a) **Answer:** The Gradient Descent (GD) Equation is

$$\begin{aligned}\theta_{k+1} &\leftarrow \theta_k - \gamma_k \frac{\partial L(\hat{\mathbf{Y}}_\theta(\mathbf{X}), \mathbf{y})}{\partial \theta_k} \\ \theta_{k+1} &\leftarrow \theta_k - \gamma_k \frac{1}{N} \frac{\partial (\sum_i \ell(\hat{\mathbf{y}}_\theta(x_i), y_i))}{\partial \theta_k}\end{aligned}$$

where,

- $\gamma_k$  is the step size
  - $\theta_{k+1}$  is the updated value after  $k$ -th iteration
  - $\theta_k$  is the updated value before  $k$ -th iteration
- (b) The Stochastic Gradient Descent (SGD) Equation, assuming a mini batch of size  $M$  is:

$$\theta_{k+1} \leftarrow \theta_k - \gamma_k \frac{1}{|M_i|} \sum_{j \in M_i} \frac{\partial \ell_j(\hat{\mathbf{y}}_\theta(x_j), y_j)}{\partial \theta_k}$$

- (c) We prefer smaller batch sizes not only because they converge sooner, but also because it results in regularization. Smaller batch sizes result in noisy steps, which can help us skip over unwanted local optimas. Since deep neural networks are easily prone to over-fitting, small batches can offer regularization and reduce the test error.
- (d) A good batch size for  $K$  class classification would be slightly more than or equal to the number of classes - that is slightly more than or equal to  $K$ . This is because we want the mini batches to represent the entire dataset in each batch, so this is a way to enforce that that condition.
- (e)
- Stochastic Gradient Descent (SGD) is noisy compared to Gradient Descent (GD). In the initial stages of training, the noise experienced due to SGD is small when compared to the information in the gradient, since have not accumulated gradients yet - so the steps taken in SGD are virtually comparable to GD's.
  - Due to this noise produced by SGD, it can help stepping over unwanted shallow local optimas - a phenomenon known as annealing, whereas there are high chances to get stuck in a local optimas when we use GD. However, if the loss function has only one global optima, both SGD and GD will converge to this point - SGD being more computationally efficient.
- (f)
- In standard GD, we update the weights based on the accumulated errors over all the samples in the training set. However, in SGD we update the weights with every training sample incrementally (starting with a randomly picked one). So SGD results in faster convergence, since it updates the weights more frequently.
  - Since the path to find the global minima/maxima is quite noisy for SGD, it can gloss over local optimas. However, GD could be stuck in a local optima, and assume that it reached a global optima.
- (g) Here is the “stochastic heavy ball interpretation” of the SGD with momentum formula. In the following, let  $\mathbf{w}$  be weights,  $\gamma$  be learning rate, and  $f$  be loss function:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \gamma_t \nabla f_i(\mathbf{w}_t) + \beta_t (\mathbf{w}_t - \mathbf{w}_{t-1})$$

- The  $\beta_t$  term refers to the *dampening factor*. It is called a heavy ball, because the use of these additional terms  $\beta_t(\mathbf{w}_t - \mathbf{w}_{t-1})$  will help stabilize the oscillatory movement seen in SGD. Similar to how a heavy ball rolls down a hill, and since this ball will have momentum it will not experience immediate changes in its path due to external forces/changes in the space.
- The appropriate range of  $\beta_t$  is  $[0,1)$ . If  $\beta_t$  was  $< 0$ , the algorithm is simply standard Gradient Descent. If  $\beta_t$  was  $> 1$ , the equation will blow up eventually.
- A smaller  $\beta_t$  value would mean a higher impact due to external changes. A large  $\beta_t$  value refers to lesser reaction due to these changes. In practice, it is advisable to use  $\beta_t = 0.9$  or  $0.99$ .

(h) Here are the equations for RMS Prop and Adam:

RMS Prop

$$\begin{aligned}\mathbf{v}_{t+1} &= \alpha \mathbf{v}_t + (1 - \alpha) \nabla f_i(\mathbf{w}_t)^2 \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \gamma \frac{\nabla f_i(\mathbf{w}_t)}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}\end{aligned}$$

- $\mathbf{v}_{t+1}$  is the  $2^{nd}$  moment estimate
- $\gamma$  is the learning rate
- $\epsilon$  helps avoid a 0 in the denominator (it is very small quantity, of the order  $10^{-7}$ ), which will make minimal impact.
- $\nabla f_i(\mathbf{w}_t)^2$  ensures that the each vector element is squared individually
- In the first equation we are calculating the exponential moving average to estimate the noisy quantity that changes over time. This adaptive method aims to up weight newer values, and assign lower weights to the values in the past - since new values give us more information. This method does this by exponentially down-weighting the older values at every step by a constant  $\alpha$ .
- $\alpha$  is a constant between 0 and 1. This works as a dampening agent, and lowers the older values till they are no longer important in the exponential moving average.
- The second equation is a modified update equation. It normalizes the gradient by Root Mean Square of the second moment ( $+ \epsilon$ ). If the gradient is a small value, it is equivalent to dividing the gradient by the standard deviation.
- We don't subtract the mean because the method maintains an exponential moving average for a non-central second moment.

Adam

$$\begin{aligned}\mathbf{m}_{t+1} &= \beta \mathbf{m}_t + (1 - \beta) \nabla f_i(\mathbf{w}_t) \\ \mathbf{v}_{t+1} &= \alpha \mathbf{v}_t + (1 - \alpha) \nabla f_i(\mathbf{w}_t)^2 \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \gamma \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}\end{aligned}$$

(replaces numerator of RMS prop with  $\mathbf{m}_t$ )

- $\mathbf{m}_{t+1}$  is the momentum's exponential moving average
- $\mathbf{v}_{t+1}$  is the  $2^{nd}$  moment estimate
- Adaptive Moment Estimation (Adam) is essentially RMSprop + momentum.
- In the first equation, the momentum is updated to an exponential moving average
- We choose a step-size and we don't have to change it, since it is being taken care of by  $\beta$ .
- The momentum is utilized to update weights instead of mini-batch gradients as seen in RMS Prop. Since the momentum values are scaled down, it has a smoother descent trajectory compared to RMS Prop.
- Adam accommodates bias correction as well (not shown in the equation above) unlike RMS Prop, hence Adam is preferred over RMS Prop for problems that have an oscillatory trajectory.
- RMS Prop takes up more computational efficiency compared to Adam. Adaptive methods (like, RMS Prop and Adam) tend to generalize poorly, compared to SGD. RMS Prop is more suitable for sparse problems

(i) `if val_loss > min(val_losses[: - NONMONO]) AND LR_DECAY > 0:`  
`optimiser.param_groups[0]['lr'] /= LR_DECAY`

- Depending on hyperparameter `NONMONO`, when the desired condition is in the vicinity, the algorithm wants to know if we over-shot the optima, or we are on the right track. It updates the learning rate accordingly, to ensure that we reach our optimal point, without missing it.
- The conditional statement checks if the current loss is greater than the minimum loss encountered until now. If it is not, the new loss would be the lowest loss yet and we maintain the same learning rate, since we are heading in the right direction.
- If the condition is met, that is, the current epoch loss is greater than the lowest loss recorded until now, the algorithm updates the learning

rate to a smaller value by dividing it with LR\_DECAY. This change in the learning rate ensures that smaller steps are taken towards the optima and we don't over-shoot it. The direction is adjusted accordingly.

### 3 Energy based models [25 pts]

An unconditional EBM (suitable for self-supervised learning) is a scalar-valued function  $F(\mathbf{y})$  that outputs low values if  $\mathbf{y}$  looks like samples from the training set (i.e. is in the region of high training data density), and higher values if  $\mathbf{y}$  is outside the region of high training data density. A conditional EBM (suitable for supervised learning and prediction-based self-supervised learning) is a function  $F(\mathbf{x}, \mathbf{y})$  that takes low values if  $\mathbf{x}$  and  $\mathbf{y}$  are “compatible” (as defined by the training set), and higher values if they are not.

The energy function of an EBM is parameterised by a trainable parameter vector  $\mathbf{w}$ . Given a training set, a good loss functional  $\mathcal{L}(\mathbf{w})$ , when minimised, must shape the energy function so as to give lower energy to regions of high training data density, and higher energies outside.

- (a) **Answer:** Given a  $\mathbf{x}$ , we can infer a compatible  $\mathbf{y}$  using the energy function:  $F(\mathbf{x}, \mathbf{y})$ .  $F(\mathbf{x}, \mathbf{y})$  determines the level of dependency between (x,y) pairs. That is an input  $\mathbf{x}$  produces the most compatible  $\hat{\mathbf{y}}$ , selected from a set of points  $\mathbf{y}$ , where the value of  $F(\mathbf{x}, \mathbf{y})$  is lowest.

$$\hat{\mathbf{y}} = \operatorname{argmin}_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}) \quad (3)$$

- (b) Give the expression for  $\mathbb{P}(\mathbf{y} \mid \mathbf{x})$  assuming a Gibbs-Boltzmann distribution. [3 pts]

$$\mathbb{P}(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(-\beta F(\mathbf{x}, \mathbf{y}))}{\int_{\mathbf{y}} \exp(-\beta F(\mathbf{x}, \mathbf{y}))} \quad (4)$$

where  $\beta$  is a positive constant.

- (c) The deterministic and probabilistic expressions for the free energy  $F(\mathbf{x}, \mathbf{y})$  as a function of  $E(\mathbf{x}, \mathbf{y}, \mathbf{z})$  are given by:

When the model depends on a latent variable  $\mathbf{z}$ :

$$\mathbb{P}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{\exp(-\beta F(\mathbf{x}, \mathbf{y}, \mathbf{z}))}{\int_{\mathbf{y}} \exp(-\beta F(\mathbf{x}, \mathbf{y}, \mathbf{z}))}$$

When we marginalize:

$$\begin{aligned} \mathbb{P}(\mathbf{y} \mid \mathbf{x}) &= \int_{\mathbf{z}} \mathbb{P}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) \\ \mathbb{P}(\mathbf{y} \mid \mathbf{x}) &= \frac{\int_{\mathbf{y}} \exp(-\beta E(\mathbf{x}, \mathbf{y}, \mathbf{z}))}{\int_{\mathbf{y}} \int_{\mathbf{z}} \exp(-\beta E(\mathbf{x}, \mathbf{y}, \mathbf{z}))} \\ &= \frac{\exp(-\beta(-\frac{1}{\beta} \log \int_{\mathbf{z}} \exp(-\beta E(\mathbf{x}, \mathbf{y}, \mathbf{z}))))}{\int_{\mathbf{y}} \exp(-\beta E(\mathbf{x}, \mathbf{y}, \mathbf{z}))} \end{aligned}$$



$$= \frac{\exp(-\beta F_\beta(\mathbf{x}, \mathbf{y}))}{\int_{\mathbf{y}} \exp(-\beta F_\beta(\mathbf{x}, \mathbf{y}))}$$

$F_\beta(\mathbf{x}, \mathbf{y})$  is the Free Energy

$$F_\beta(\mathbf{x}, \mathbf{y}) = -\frac{1}{\beta} \log \int_{\mathbf{z}} \exp(-\beta E(\mathbf{x}, \mathbf{y}, \mathbf{z}))$$

- (d) The expression of the negative log-likelihood loss, which is designed to maximise the likelihood that our model assigns to the  $\mathbf{y}^{(i)}$  given the  $\mathbf{x}^{(i)}$  in the training set:  $\prod_i \mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$  can be derived as follows:

$$\mathbb{P}(\mathbf{y}^1, \dots, \mathbf{y}^N | \mathbf{x}^1, \dots, \mathbf{x}^N, \mathbf{w}) = \prod_i \mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})$$

using maximum likelihood we try to find the values of the trainable parameter vector  $\mathbf{w}$  which will maximize this product (or which minimizes the negative log of this product)

$$-\log(\mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})) = \sum_i^N -\log(\mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}))$$

$$-\log \mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) = \sum_i^N \beta F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) + \log \int_{\mathbf{y}} \exp(-\beta F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y}))$$

$$\mathcal{L}(\mathbf{w}) = \frac{1}{N} \sum_i^N \left[ F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) + \frac{1}{\beta} \log \int_{\mathbf{y}} \exp(-\beta F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y})) \right]$$

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})}{\partial \mathbf{w}} - \int_{\mathbf{y}} \frac{\partial F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y})}{\partial \mathbf{w}} \mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$$

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = F_{\mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) + F_{\beta, \mathbf{w}}(\mathbf{x}^{(i)}, \mathbf{y})$$

where  $\mathbb{P}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$  is the expression from Gibbs-Boltzmann distribution (shown above).

- (e) **Noise Contrastive Estimator's** objective function is a contrastive method, which tries to maximize the cost function without any constraints on normalization which allows us to operate using un-normalized functions.

$S(u, v)$ : is a similarity metric between two feature maps/vectors/Convolution Net outputs as a measure of their cosine similarity.  $f(u), g(u)$  can be

thought of as independent layers on top of convolutional feature extractors/layers.

Using mini-batches, each batch will consist of one similar/positive pair, and many negative/dissimilar pairs. Using these we can measure the similarity between the transformed feature vector  $\mathbf{I}_t$  and the other feature vectors in the mini-batch, one positive and the rest negative pairs. Then we compute the score of a softmax-like function on the positive pair. So, maximizing a softmax score means that we minimize the rest of the scores. Which aligns with what we want from our Energy-Based Model. The final loss function, as shown below, pushes the energy down on similar/positive pair locations, and pushes up on dissimilar/negative pair locations.

$$h(\mathbf{v}_I, \mathbf{v}_{I^t}) = \frac{\exp(\frac{s(\mathbf{v}_I, \mathbf{v}_{I^t})}{\tau})}{\exp(\frac{s(\mathbf{v}_I, \mathbf{v}_{I^t})}{\tau}) + \sum_{I' \in D_N} \exp(\frac{s(\mathbf{v}_I, \mathbf{v}_{I'})}{\tau})}$$

$$L_{NCE}(\mathbf{I}, \mathbf{I}^t) = -\log[h(f(\mathbf{v}_I), g(\mathbf{v}_{I^t}))] - \sum_{I' \in D_N} \log[1 - h(g(\mathbf{v}_{I^t}), f(\mathbf{v}_{I'}))]$$

- (f) Another class of Energy-Based Model training is called Architectural methods. These methods are different ways of limiting the information capacity. For example,

- **K-Means Clustering:**

- This is an unsupervised learning algorithm where we cluster the data in various groups (k clusters) based on their nearest means (of the cluster centroids).
- K-means uses a energy-based model with the energy function  $F(\mathbf{x}, \mathbf{y}) = ||\mathbf{y} - \mathbf{W}\mathbf{z}||^2$ , where  $\mathbf{z}$  is a 1-hot vector. With the values of  $\mathbf{y}$  and  $\mathbf{k}$ , we can infer the  $\mathbf{k}$  possible columns of  $\mathbf{W}$  which minimizes the reconstruction error/energy function.
- This method ensures that the volume of the low energy regions are bounded by the number of prototypes  $\mathbf{k}$ .
- To train the algorithm, we use an architectural method where we use latent variable  $\mathbf{z}$ , to choose the column of  $\mathbf{W}$  closest to  $\mathbf{y}$  and then try to move even closer by taking a gradient step and repeat the process.
- We can observe how we don't intend to push the energy levels up anywhere (like we do in contrastive methods), but we try to push down the energy values in certain regions.

- **Variational Auto Encoders:**

- VAEs are generative models which have well defined latent spaces.

- This Architectural method builds an energy function  $F(\mathbf{x}, \mathbf{y})$  which have low energy regions by the means of regularization.
- For VAE, the loss function is defined as:  $l(\mathbf{x}, \hat{\mathbf{x}}) = l_{reconstructionerror} + \beta l_{KL}(\mathbf{z}, N(0, I_d))$ .
- The regularization term is on the latent layer, which enforces a bottleneck of information using some specific Gaussian structure on the latent space. We use a penalty term  $l_{KL}(\mathbf{z}, N(0, I_d))$ . Without this regularization term, VAE will resemble an autoencoder, which usually leads to over-fitting and won't exhibit the generative properties we need.

- **Sparse Coding:**

- Sparse coding can be defined as an unconditional regularized latent-variable energy based model. It approximates the data with a piecewise linear function.
- The energy function can be defined as:  $F(\mathbf{x}, \mathbf{y}) = \|\mathbf{y} - W\mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_{L1}$ .
- This Architectural model also tries to lower the energy values in certain regions (determined by the latent variable  $\mathbf{z}$ ), by using a regularization term.