

How does the GAN make $\text{\textbackslash pmodel} \approx \text{\textbackslash pdata}$?

The Generator Loss function we constructed is a proxy to achieve the goal

$$\text{\textbackslash pmodel} \approx \text{\textbackslash pdata}$$

That is: the distribution of samples produced by the Generator is (approximately) the same as the "true" distribution

- we note that we don't know the "true" $\text{\textbackslash pdata}$
 - we only have available a sample and those the training set defines an *empirical* distribution

There are several ways to quantify

$$\text{\textbackslash pmodel} \approx \text{\textbackslash pdata}$$

One choice would be the minimization of KL Divergence

- $\text{\textbackslash KL}(\text{\textbackslash pdata} || \text{\textbackslash pmodel})$

An alternative, still using KL Divergence

- $\text{\textbackslash KL}(\text{\textbackslash pmodel} || \text{\textbackslash pdata})$

Which is a better choice ?

In order to answer the question, we begin with a few preliminaries

- Definition of KL Divergence
- Proving
 - minimizing KL Divergence increases log-likelihood

Definition of KL Divergence

As a reminder of the definition of KL Divergence

$$\begin{aligned}\text{KL}(p||q) &= -\sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) \\ &= \sum_x p(x) * (\log p(x) - \log q(x)) \\ &= \textcolor{red}{\mathbb{E}_{\mathbf{x} \sim p}(\log p(x) - \log q(x))}\end{aligned}$$

You can see that it is

- the point-wise difference between the (log) probability of \mathbf{x} in distributions p and q
- averaged over the distribution of $\mathbf{x} \sim p$

and thus is a point-wise measure of the dis-similarity of the two distributions.

We note that the KL Divergence is *not symmetric*

$$\text{\textbackslash KL}(\text{\textbackslash pdata}||\text{\textbackslash pmodel}) \neq \text{\textbackslash KL}(\text{\textbackslash pmodel}||\text{\textbackslash pdata})$$

so the two choices are different.

- both are expectations
- but over *different distributions*

KL Divergence leads to Maximum Likelihood Estimation

We now show that using $\text{KL}(\text{\textbackslash pdata} || \text{\textbackslash pmodel})$ as a loss function

- results in a estimation of the model distribution $\text{\textbackslash pmodel}$
- that is the Maximum Likelihood estimator of the training examples (represented by $\text{\textbackslash pdata}$)

That is

- $\text{\textbackslash pmodel}$ is the best explanation of the training dataset $\text{\textbackslash pdata}$

Choosing pmodel to Minimize gives

$$\begin{aligned}\text{\textbackslash KL}(\text{\textbackslash pdata} || \text{\textbackslash pmodel}) &= \int_{\text{\textbackslash x}} \text{\textbackslash pdata}(\text{\textbackslash x}) \left(\log \frac{\text{\textbackslash pdata}(\text{\textbackslash x})}{\text{\textbackslash pmodel}(\text{\textbackslash x})} \right) d\text{\textbackslash x} \\ &= \text{\textbackslash E}_{\text{\textbackslash x} \in \text{\textbackslash pdata}} \log(\text{\textbackslash pdata}(\text{\textbackslash x})) - \log(\text{\textbackslash pmodel}(\text{\textbackslash x}))\end{aligned}$$

minimizing KL

$$\approx \text{\textbackslash E}_{\text{\textbackslash x} \in \text{\textbackslash pdata}} - \log(\text{\textbackslash pmodel}(\text{\textbackslash x}))$$

So minimizing \textbackslash KL is equivalent to

- minimizing the Negative Log Likelihood
 - in other words: *maximizing* the Log Likelihood
-

Choosing the KL Divergence

The first choice

$$\text{KL}(\text{pdata} \parallel \text{pmodel}) = \mathbb{E}_{x \sim \text{pdata}} (\log \text{pdata}(x) - \log \text{pmodel}(x))$$

maximizes $\log(\text{pmodel}(x))$ for $x \in \text{pdata}$

- pmodel assigns high probability to Real examples
- model creates Real examples with high probability

By way of analogy with measures for Classification

- the expectation over $\text{\textbackslash p}data$ emphasizes Recall over Precision

We can achieve high Recall

- by reducing chance of False Negatives (FN)
- even if it increases chance of False Positives (FP)

In the GAN context this means

reducing FN $\rightsquigarrow \text{\textbackslash pmodel}$ assigns high probability to each training example in
increasing FP $\rightsquigarrow \text{\textbackslash pmodel}$ assigns high probability to $\text{\textbackslash x} \notin \text{\textbackslash pdata}$

The second choice $\text{KL}(\text{pmodel} \parallel \text{pdata})$

$$\text{KL}(\text{pmodel} \parallel \text{pdata}) = \mathbb{E}_{x \sim \text{pmodel}} (\log \text{pmodel}(x) - \log \text{pdata}(x))$$

maximizes $\text{pdata}(x)$ for $x \in \text{pmodel}$

- emphasizes that synthetic examples are "realistic"
 - highly probable, as defined by the empirical distribution (training data)
 pdata

This ("realistic examples") might be the more desirable property than "high fidelity" to the training data.

Continuing with our Recall versus Precision analogy, this measure

- increases Precision by reducing False Positives
 - examples generated by \pmodel are likely according to \pdata

So it seems as if the second choice $\text{KL}(\text{\textit{pmodel}}||\text{\textit{pdata}})$ may be more desirable.

But we don't know the true $\text{\textit{pdata}}$!

- we only have an empirical sample: the training dataset
- so, in practical terms: we can't maximize it

Thus, practical considerations lead us to the first choice.

Jensen-Shannon Divergence

We have observed that the KL divergence is *not* symmetric

$$\text{\textcolor{red}{KL}}(P||Q) \neq \text{\textcolor{red}{KL}}(Q||P)$$

because the expectations are taken over different distributions.

An alternative measure of similarity of two distributions is the Jensen-Shannon Divergence (JSD)

$$\begin{aligned}\text{JSD}(P||Q) &= \text{JSD}(Q||P) \\ &= \frac{1}{2} \text{KL} \left(P \parallel \frac{P+Q}{2} \right) + \\ &\quad \frac{1}{2} \text{KL} \left(Q \parallel \frac{P+Q}{2} \right)\end{aligned}$$

This measure is

- symmetric
- is a kind of mixture of $\text{KL}(P||Q)$ and $\text{KL}(Q||P)$.

Huszár (<https://arxiv.org/pdf/1511.05101.pdf>) has a Generalized JSD which interpolates between the two terms

$$\begin{aligned}\text{JSD}_\pi(P||Q) &= \text{JSD}(Q||P) \\ &= \pi \text{KL}(P || \pi P + (1 - \pi)Q) + \\ &\quad (1 - \pi) \text{KL}(Q || \pi P + (1 - \pi)Q)\end{aligned}$$

The Generalized JSD

- **Not** symmetric although

$$\text{JSD}_\pi(P||Q) = \text{JSD}_{1-\pi}(Q||P)$$

In []:

Huszar shows that, for small values of π

$$\frac{\text{JSD}_\pi(P||Q)}{\pi} \approx \text{KL}(P || Q)$$

and

$$\frac{\text{JSD}_{1-\pi}(P||Q)}{1 - \pi} \approx \text{KL}(Q || P)$$

In the first case

- $\text{JSD}_\pi(P||Q)$ is proportional to Maximum Likelihood

In the second case

- $\text{JSD}_{1-\pi}(P||Q)$ is proportional to $\text{KL}(Q || P)$

In implementing Generalized JSD

- The Discriminator is trained (as usual) on a mix of real and fake examples
 - But *not* in equal numbers
 - π is fraction of samples from Q
 - $(1 - \pi)$ is fraction of samples from P
 - $\pi < \frac{1}{2}$: real samples over represented
 - $\pi > \frac{1}{2}$: biased toward Q
- Explains why we often see training with Generator updated twice for each update of Discriminator ?

Adversarial Training and the Jensen-Shannon Divergence

The Discriminator Loss loss_D

- summed over all examples
 - (ignoring the $\frac{1}{2}$ from the previous presentation where we assumed equal number of Real and Fake)

is

$$\text{loss}_D = - \left(\mathbb{E}_{\mathbf{x}^{\text{ip}} \in \text{pdata}} \log D(\mathbf{x}^{\text{ip}}) + \mathbb{E}_{\mathbf{x}^{\text{ip}} \in \text{pmodel}} \log (1 - D(\mathbf{x}^{\text{ip}})) \right)$$

We also showed that the optimal Discriminator results in

$$D^*(\mathbf{x}) = \frac{\text{pdata}(\mathbf{x})}{\text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})}$$

In []:

Plugging $D^*(\mathbf{x})$ into loss_D (Goodfellow Equation):

$$\begin{aligned}
 \text{loss}_D &= - \left(\mathbb{E}_{\mathbf{x} \in \text{pdata}} \log \frac{\text{pdata}(\mathbf{x})}{\text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})} + \mathbb{E}_{\mathbf{x} \in \text{pmodel}} \log \frac{\text{pmodel}(\mathbf{x})}{\text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})} \right) \\
 &= - (\text{KL}(\text{pdata} || \text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})) \\
 &\quad + \text{KL}(\text{pmodel} || \text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x}))) \\
 &= - \left(\log 4 + \text{KL}(\text{pdata} || \frac{\text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})}{2}) \right. \\
 &\quad \left. + \text{KL}(\text{pmodel} || \frac{\text{pmodel}(\mathbf{x}) + \text{pdata}(\mathbf{x})}{2}) \right) \\
 &= - (\log 4 + 2 * \text{JSD}(\text{pdata} || \text{pmodel}))
 \end{aligned}$$

The above equations shows that

- minimizing KL Divergence (second line above)
- under the assumption that the Discriminator can train to be the **optimal** adversary

results in $\backslash \text{loss}_D$ becoming equivalent to Jensen-Shannon Distance (last line above)

So solving the minimax optimally minimizes the JSD divergence between $\backslash \text{pdata}$ and $\backslash \text{pmodel}$.

In [2]: `print("Done")`

Done

