# Adversarial Attack and Learning

**Anosh Billimoria, [billimor@rhrk.uni-kl.de](mailto:billimor@rhrk.uni-kl.de), @anbilli (Mattermost)**

**Setup:** The constraints on $\delta$ can be summarized where, $|\delta|$ is the norm of the noise, $\epsilon$ and $\alpha$ are integers.

(i) …. $|\delta| = |\delta|_\infty + \alpha|\delta|_2$

such that, (ii) …. $|\delta| \leq \epsilon$

**Procedure:** To keep matters simple LeNet-5 CNN is selected.

**Constraint Simplification:** To bound $|\delta|$ within desired $\epsilon$ and we estimate the $|\delta|_\infty$ using gradient method and then clipping the result. To that end here we have selected some methods to modify the constraints over the norm using some simplifications. Namely,

1. Since it is not specified the noise is randomly drawn from a uniform distribution between -$\epsilon$ and $\epsilon$ of size equal to the image size (28,28)
2. The $|\delta|_2$ is calculated and $\delta$ is normalized such that $|\delta|_2 = 1$. Making (i).
   (iii) …. $|\delta'| = |\delta'|_\infty + \alpha.1$
   where $\delta'$ is the new normalised version. We can modify (ii) similarly.
   (iv) …. $|\delta'|_\infty \leq \epsilon - \alpha$
   dividing by $\alpha$ in (ii) and take $\alpha^* = \frac{1}{\alpha}$, we get,
   (v) …. $\alpha^*|\delta'|_\infty \leq \alpha^*\epsilon - 1$
3. By reducing the problem to simply depend on the $|\delta'|_\infty$ of new normalised $\delta'$. We can apply gradient estimation method and clip results within ($\alpha^*\epsilon - 1$)
4. Here, we can note that to replicate $\alpha^*$ is multiplied with norm-infinity. We can also use it as step-size during the estimation.

**Effect of varying $\epsilon$ and $\alpha$:** Now we perturb the images.

1. Constant $\epsilon$ and varying $\alpha^*$: We observe that $\alpha$ value can affect the transparency of the digit. A lower value almost blacks out any grey pixels, while a higher value causes extreme alterations making the digit incomprehensible even for humans.
2. Varying $\epsilon$: A larger $\epsilon$ distorts images beyond recognition and leaves lesser room for selecting appropriate $\alpha^*$ to balance out the effect. On the other hand, a smaller $\epsilon$ can allow images to have more contrast for a wider range of $\alpha^*$.
3. We can also say that choosing a just right $\epsilon$ can offer better tuning opportunities on $\alpha^*$, especially since our modification to $\alpha^*$ makes it act like step-size for the infinity-norm estimate.

**Training:** We have selected $\alpha^*$= 0.4, $\epsilon$ = 0.7, as the adversarial training parameters.

1. **Choice of $\alpha^*$ and $\epsilon$:** For this combination the images are still recognizable by the human eye, yet the accuracy dropped to two-thirds of the original. Since the values are well-balanced with respect to each other this combination can potentially train a more universally robust model.
2. **First training:** We decided to retrain the model and add an adversarial image to the dataset every tenth batch using the above values. Validation is performed on the adversarial images from the test set. The results are interesting since testing on other combinations of $\epsilon$ and $\alpha^*$ tend to respond with better accuracy and lower loss than before.
3. **Second training:** Train a newly initialized model with every tenth batch is replaced with an adversarial one. The results of which are not impressive probably due to using a simple model. It may be beneficial to modify the model slightly to incorporate more robustness.