

PURDUE

BME 64600 – 001 and ECE 60146 – 001

Midterm #2, Spring 2024

NAME _____

PUID _____

Exam instructions:

- You have 75 minutes to work the exam.
- This is a closed-book and closed-note exam. You may not use or have access to your book, notes, any supplementary reference, a calculator, or any communication device including a cell-phone or computer.
- You may not communicate with any person other than the official proctor during the exam.

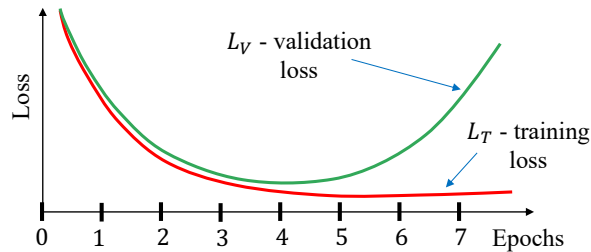
To ensure Gradescope can read your exam:

- Write your full name and PUID above and on the top of every page.
- Answer all questions in the area designated for each problem.
- Write only on the front of the exam pages.
- DO NOT run over to the next question.

Name/PUID: _____ **Key**

Problem 1. (30pt) Training, generalization, and regularization

You are training a deep neural network, and you get the result pictured below.



a) What number of epochs is likely to result in the best parameter estimates? Why?

Solution: 4 epochs because this is the value at which the validation loss is approximately minimized.

b) What is likely to happen if you choose the parameter estimates from 7 epochs? Why?

Solution: In this case, the parameter estimates will likely be overfit to the training data, so they will not generalize well.

c) What does this curve suggest about the capacity of your model?

Solution: The curve suggests that the model's capacity may be a bit too high given the amount of available training data.

d) Given your answer to c) above and assuming that you can not change the model, then what other modification to the training might be helpful? Why?

Solution: It might be helpful to introduce regularization of the parameters into the loss function. For example, l_1 regularization can enforce sparsity in the model, thereby removing the unnecessary parameters, and effectively reducing the model order.

e) If you increase the training data by a factor of 10, then what would you expect to happen to the **training** loss curve?

Solution: Since it is more difficult to fit a larger amount of data, I would expect the training loss to go up.

f) If you increase the training data by a factor of 10, then what would you expect to happen to the **validation** loss curve?

Solution: Since the larger amount of training data will reduce overfitting, I would expect the validation loss to go down.

Name/PUID: _____

Problem 2. (25pt) An image is a point in a space and can have a distribution.

Let $X \in [0, 1]^{N^2}$ be an $N \times N$ image that takes values in the range $[0, 1]$. Furthermore, we will assume it is a random vector with distribution $X \sim p(x)$. Also assume that $Y \in \mathfrak{R}^{N^2}$ is an image and that $Y = AX$ for some matrix A .

a) What is the shape of A ?

Solution: A is a $N^2 \times N^2$ matrix.

b) If X is a 1024×1024 image (i.e., 1 Mega pixel), and A is stored with single precision floats, then how much memory does it take to directly store the matrix A on a standard computer?

Solution: Since A is then a $(2^{10})^2 \times (2^{10})^2$ matrix, and since a single precision float is 4 bytes on a standard computer, we know that A will take 2^{42} bytes or equivalently 4 Terra Bytes of storage.

c) Let X_i be a single pixel in X . Then write an expression for the marginal distribution of $X_i \sim p_i(x_i)$ in terms of $p(x)$.

Solution:

$$p_i(x_i) = \int_{\mathfrak{R}} \cdots \int_{\mathfrak{R}} p(x) dx_0 \cdots dx_{i-1} dx_{i+1} \cdots dx_{N^2-1} .$$

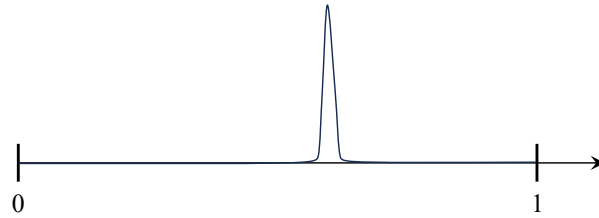
d) Write an expression for $p_i(x_i|x_0, \cdots, x_{i-1}, x_{i+1}, \cdots, x_{N^2-1})$, the conditional probability distribution of X_i given $(X_0 \cdots, X_{i-1}, X_{i+1}, \cdots, X_{N^2-1})$.

Solution:

$$p_i(x_i|x_0, \cdots, x_{i-1}, x_{i+1}, \cdots, x_{N^2-1}) = \frac{p(x_0, \cdots, x_{N^2-1})}{\int_{\mathfrak{R}} p(x_0, \cdots, x_{i-1}, x_i, x_{i+1}, \cdots, x_{N^2-1}) dx_i}$$

e) Notice that both $p_i(x_i)$ and $p_i(x_i|x_0, \cdots, x_{i-1}, x_{i+1}, \cdots, x_{N^2-1})$ are density functions for X_i . If someone gives you the following plot of a density function, would you guess that it

was the marginal distribution or the conditional distribution of X_i ? Why?



Solution: I would guess that it is the conditional distribution because it is relatively narrow. If it was the marginal distribution, then I would expect it to be broad since a pixel might take any value between 0 and 1. However, the conditional distribution of a pixel given its neighbors usually has very little uncertainty since the pixel should be similar to its neighbors.

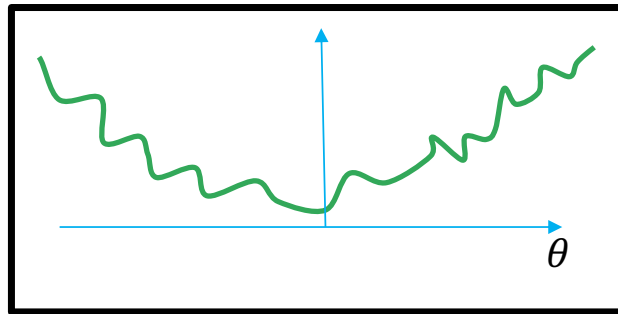
Name/PUID: _____

Problem 3. (25pt) Stochastic Gradient Descent

Consider the loss function given by,

$$L(\theta) = \sum_{k=0}^{K-1} \|x_k - f_{\theta}(y_k)\|^2$$

that we would like to minimize as a function of θ using stochastic gradient descent.



Furthermore, let the full set of training data be partitioned into subsets, S_b , for $b = 0, \dots, K/N_b$, where N_b is the batch size, and define the partial loss function,

$$L(\theta; S_b) = \sum_{k \in S_b} \|x_k - f_{\theta}(y_k)\|^2.$$

Also, assume you use an algorithm such as the one below.

$$v \leftarrow 0$$

Repeat until converged:

For $b = 1$ to N_b :

$$d \leftarrow -\nabla_{\theta} L(\theta; S_b)$$

$$v \leftarrow \gamma v + \alpha(1 - \gamma)d$$

$$\theta \leftarrow \theta + v^t$$

a) Name two advantages of using a smaller batch size for this problem?

Solution: Two advantages of using a smaller batch size are:

1. Each iteration will be much faster (but the amount of time for an epoch will remain approximately the same).
2. The SGD algorithm will be more robust to local minimum since the “noise” in the gradient update will help to “knock it out” of local minimum.

b) Name a disadvantage of using a smaller batch size for this problem?

Solution: A disadvantage of using a smaller batch size is that SGD will tend to “hunt” around the exact solution. So the final estimate of θ will be a bit noisy.

c) What happens when $\gamma = 0$?

Solution: In this case, there is no momentum, and the step is exactly proportional to the negative gradient.

d) What is it called when we set $\gamma > 0$, and what is the advantage of doing this?

Solution: This is called “momentum”. The advantage of using momentum is that SGD tends to do a better job of jumping out of local minimum.

e) If at every iteration, $d = \mathbf{1}$, then what is $\lim_{iterations \rightarrow \infty} v$?

Solution: We can find the solution to this by solving for the equilibrium given by

$$\begin{aligned}v &= \gamma v + \alpha(1 - \gamma)d \\v &= \gamma v + \alpha(1 - \gamma)\mathbf{1} \\(1 - \gamma)v &= \alpha(1 - \gamma)\mathbf{1} \\v &= \alpha\mathbf{1}\end{aligned}$$

So asymptotically, the step size is α times the negative gradient.

Name/PUID: _____

Problem 4. (25pt) Residual training of denoisers

Our goal is to train a denoiser to remove additive Gaussian white noise from images of size $N \times N$. To do this, we generate training data with the form

$$Y_k = X_k + W_k ,$$

where $W_k \sim N(0, \sigma^2 I)$ are i.i.d. Gaussian random vectors of dimension N^2 for $k = 0, \dots, K - 1$. Furthermore, let $W_{k,i}$ denote the i^{th} component of W_k .

a) Calculate a simple expression for $R(i, j) = E[W_{k,i}W_{k,j}]$. What does this tell you about the components of W_k ?

Solution:

$$R(i, j) = E[W_{k,i}W_{k,j}] = \sigma^2 \delta(i - j)$$

b) How should we select the training images X_k ? Why?

Solution: We should select K images that are typical of the images we plan on denoising. This will allow the neural network to learn the distribution of both the images (i.e., the prior distribution) and the noise.

c) If we are going to train a deep neural network, $f_\theta(y)$, to perform this task, then what is the best way to train it?

Solution: It is best to train the neural network, $f_\theta(y)$, to estimate the noise rather than to estimate the image.

d) Write out the loss function for training the deep neural network in terms of X_k, Y_k, W_k and f_θ .

Solution:

$$L(\theta) = \frac{1}{K} \|W_k - f_\theta(Y_k)\|^2$$

e) What is the advantage of estimating the noise over estimating the signal?

Solution: The advantage of estimating the noise is that it has less dynamic range, so it is easier to train the neural network to accurately estimate it. Moreover, estimating the noise can be viewed as a skipped connection, so this has the advantage of mitigating the vanishing gradient problem.

Name/PUID: _____

Problem 5. (20pt) Convolution blocks and their adjoint

Consider the function $f : \mathfrak{R}^N \rightarrow \mathfrak{R}^N$ that implements 1D (true) convolution with the kernel $[w_0, \dots, w_{2p}]$ for some p and uses a “same” boundary condition.

Also, let $g = f^t$ be the adjoint function of f .

a) Is f a linear function? Provide a proof.

Hint: A function f is linear if $\forall x, y \in \mathfrak{R}^N$ and $\forall \alpha, \beta \in \mathfrak{R}$ it is always the case that $f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$.

Solution: Yes, f is linear because

$$f(\alpha x + \beta y) = w * (\alpha x + \beta y) \tag{1}$$

$$= \alpha w * x + \beta w * y \tag{2}$$

$$= \alpha f(x) + \beta f(y) . \tag{3}$$

b) For the case that $N = 5$ and $P = 1$, write out a matrix A so that $f(x) = Ax$.

Solution:

$$A = \begin{bmatrix} w_1 & w_0 & 0 & 0 & 0 \\ w_2 & w_1 & w_0 & 0 & 0 \\ 0 & w_2 & w_1 & w_0 & 0 \\ 0 & 0 & w_2 & w_1 & w_0 \\ 0 & 0 & 0 & w_2 & w_1 \end{bmatrix}$$

c) For the case that $N = 5$ and $P = 1$, write out a matrix B so that $g(y) = By$.

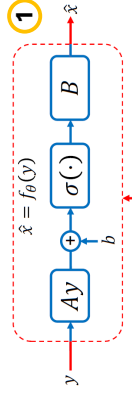
Solution:

$$B = \begin{bmatrix} w_1 & w_2 & 0 & 0 & 0 \\ w_0 & w_1 & w_2 & 0 & 0 \\ 0 & w_0 & w_1 & w_2 & 0 \\ 0 & 0 & w_0 & w_1 & w_2 \\ 0 & 0 & 0 & w_0 & w_1 \end{bmatrix}$$

d) What is the interpretation of g ?

Solution: g is convolution with the time-reversed kernel, $[w_{2p-1}, \dots, w_0]$.

Single Layer NN Graphically:



Where: $f_{\theta}(y) = B\sigma(Ay + b)$

M-dimensional simplex:

$$S^M = \left\{ x \in \mathbb{R}^M, \forall i, x_i \geq 0 \text{ and } 1 = \sum_{i=0}^{M-1} x_i \right\}$$

Softmax: $\sigma_i(z) = \frac{e^{z_i}}{\sum_j e^{z_j}}$

Gradient of the Softmax:

$$[\nabla \sigma(z)]_{i,j} = \frac{1}{\sum_k e^{z_k}} \left(e^{z_i} \delta_{i,j} - \frac{e^{z_i} e^{z_j}}{\sum_k e^{z_k}} \right)$$

MSE Loss:

$$L_{MSE}(\theta) = \frac{1}{K} \sum_{k=0}^{K-1} \|x_k - f_{\theta}(y_k)\|^2$$

Gradient Descent and line search:

Repeat until converged {
 $d \leftarrow -\nabla L(\theta)$
 $\alpha^* \leftarrow \arg \min_t \{L(\theta + \alpha d^t)\}$
 $\theta \leftarrow \theta + \alpha^* d^t$
}

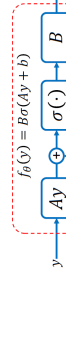
Loss gradient (MSE):

$$-\nabla_{\theta} L_{MSE}(\theta) = \frac{2}{K} \sum_{k=0}^{K-1} (x_k - f_{\theta}(y_k)) \nabla_{\theta} f_{\theta}(y_k)$$

Einstein Notation: Delta Functions:

$$\delta^i_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

$$G^{k_j} = G^k_i \delta^i_j$$



Gradients:

w.r.t. y $[\nabla_y f]^t = B^t_{i_1} [\nabla \sigma]_{i_1} A^t_{i_2} y^t$
w.r.t. A $[\nabla_A f]^{i_1}_{j_1, j_2} = B^{i_1}_{i_2} [\nabla \sigma]^{i_2}_{j_1} y^t_{j_2}$
w.r.t. B $[\nabla_B f]^{i_1}_{j_1, j_2} = B^{i_1}_{i_2} [\nabla \sigma]^{i_2}_{j_1}$
w.r.t. B $[\nabla_B f]^{i_1}_{j_1, j_2} = \delta^{i_1}_{i_2} \sigma_{j_1}$

Adjoint of the Loss Gradient:

$$[\nabla L_{MSE}(\theta)]^t = -\frac{2}{K} \sum_{k=0}^{K-1} [\nabla f_{\theta}(y_k)]^t (x_k - f_{\theta}(y_k))$$



Back Propagation:

$$g_{m,k} = \nabla_{\theta_m} \|x_k - f_{\theta}(y_k)\|^2 = -2(x_k - f_{\theta}(y_k)) \nabla_{\theta_m} f_{\theta}(y_k)$$

MAE:

$$\|a - b\|_1 = \sum_i |a_i - b_i|$$

MC Cross Entropy:

$$-\sum_i a_i \log b_i$$

Back Prop for Loss Functions:

MSE: $\epsilon^t_k = -2(x_k - f_{\theta}(y_k))$
MAE: $\epsilon^t_k = -\text{sign}(x_k - f_{\theta}(y_k))$
CE: $\epsilon^t_k = \frac{x_k}{f_{\theta}(y_k)}$

Convolution "Same" Boundary:

$$x_{(j_1, j_2)} = \sum_{k_1=-1}^{k_1=1} \sum_{k_2=-1}^{k_2=1} z_{(j_1 - k_1, j_2 - k_2)} w_{(k_1, k_2)} * z_{(j_1, j_2)}$$

Convolution "Valid" Boundary:

$$x_{(j_1, j_2)} = \sum_{k_1=-1}^{k_1=1} \sum_{k_2=-1}^{k_2=1} z_{(j_1 + 1 - k_1, j_2 + 1 - k_2)} w_{(k_1, k_2)} * z_{(j_1, j_2)}$$

2D Conv w/ vector input:

$$x_{(j_1, j_2)} = \sum_{i=0}^{p-1} \sum_{k_1=-p-k_2}^{-p-k_2} z_{(j_1 - i, j_2 - k_2)} w_{(i, j_2)} * z_{(j_1, j_2)}$$

2D Conv w/ vector input/output:

$$[g]_{j_1, j_2, j_3} = w_{(j_1, j_2), i} j_3 * z_{(j_1, j_2), i}$$

CNN Average Pooling:

$$x_{(j_1, j_2)} = \max_{0 \leq k_1 < p} \max_{0 \leq k_2 < p} z_{(d_1 j_1 + k_1, d_2 j_2 + k_2)}$$

CNN Max Pooling:

$$x_{(j_1, j_2)} = \max_{0 \leq k_1 < p} \max_{0 \leq k_2 < p} z_{(d_1 j_1 + k_1, d_2 j_2 + k_2)}$$

Update Directions:

Gradient step: \rightarrow Replace A
 $A \leftarrow A + \alpha d^t$ with b or B

Overfitting Fixes:

Regularization - adds penalties to weights, encouraging them to be small or collapse to 0
 $S(\theta) = -\log p(\theta)$
 $S(\theta) = \|\theta\|^2$
 $S(\theta) = \|\theta\|_1$

Discrete Probability Density:

$$P\{X \in A\} = \sum_{i \in A} p_i$$

Expected Mean, Variance:

$$\mu = E[X] = \sum_{i=0}^{M-1} i p_i$$

$$E[f(X)] = \sum_{i=0}^{M-1} f(i) p_i$$

Conditional Density:

$$P(X = i | Y = j) = \frac{P(X=i, Y=j)}{P(Y=j)}$$

Conditional Expectation:

$$E[f(X)|Y] = \sum_{i=0}^{M-1} f(i) P_{X|Y}(i|Y)$$

Bayesian View:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\sum_{i=0}^{M-1} \frac{N_i}{N} \log \theta_i \right\} = \frac{N_i}{N}$$

MAP Estimate:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log P_{Y|X}(y|X) - \log P_X(x) \right\}$$

Bias and Variance:

$$\text{bias}^2 = E\{[\hat{X}] - X\|^2\}$$

$$\text{Var} = E\{[\hat{X}] - E[\hat{X}]|^2\}$$

MAP Estimate:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log P_{Y|X}(y|X) - \log P_X(x) \right\}$$

Bayesian View:

$$P_{Y|X}(x|Y) = \frac{P_Y(y) \prod_{k=0}^{M-1} P_{X_k}(x_k)}{\sum_{k=0}^{M-1} P_{Y|X}(y|X_k) P_X(x_k)}$$

Bayesian View:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x|y) - \log p(y) \right\}$$

Regression:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x, y) \right\}$$

MLE:

$$\hat{\theta} = \arg \max_{\theta} \left\{ \log p_{\theta}(x, y) \right\}$$

Bayesian View:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x|y) - \log p(y) \right\}$$

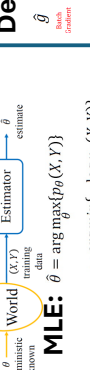
Bayesian View:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x, y) \right\}$$

Bayesian View:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x, y) \right\}$$

Frequentist View:



MLE:

$$\hat{\theta} = \arg \max_{\theta} \left\{ \log p_{\theta}(x, y) \right\}$$

Regression:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x, y) \right\}$$

Classification (w/ 1-hot encoding):

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log p_{\theta}(x|y) - \log p(y) \right\}$$

Multinomial Distro:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\sum_{i=0}^{M-1} \frac{N_i}{N} \log \theta_i \right\} = \frac{N_i}{N}$$

Bayesian View:

$$P_{Y|X}(x|Y) = \frac{P_Y(y) \prod_{k=0}^{M-1} P_{X_k}(x_k)}{\sum_{k=0}^{M-1} P_{Y|X}(y|X_k) P_X(x_k)}$$

Bayes Law:

$$P_{Y|X}(x|Y) = \frac{P_Y(y) \prod_{k=0}^{M-1} P_{X_k}(x_k)}{\sum_{k=0}^{M-1} P_{Y|X}(y|X_k) P_X(x_k)}$$

Min Mean Square Error:

$$\hat{X} = E[X|Y] = \sum_{x=0}^{M-1} x P_{X|Y}(x|Y)$$

MAP Estimate:

$$\hat{\theta} = \arg \min_{\theta} \left\{ -\log P_{Y|X}(y|X) - \log P_X(x) \right\}$$

Bias and Variance:

$$\text{bias}^2 = E\{[\hat{X}] - X\|^2\}$$

$$\text{Var} = E\{[\hat{X}] - E[\hat{X}]|^2\}$$

Bayesian View:

$$\text{MSE}_{\theta} = \text{Var}_{\theta} + (\text{bias}_{\theta})^2$$

Bayesian View:

$$\text{bias}^2 = E\{[\hat{X}] - X\|^2\}$$

Bayesian View:

$$\text{Var} = E\{[\hat{X}] - E[\hat{X}]|^2\}$$

Bayesian View:

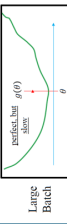
$$\text{MSE} = \text{Var} + \text{bias}^2$$

Bayesian View:

$$\text{MSE}_{\theta} = \text{Var}_{\theta} + (\text{bias}_{\theta})^2$$

Stochastic Gradient Descent:

init $v \leftarrow 0$; $r \leftarrow 0$;
init $t \leftarrow 0$
Repeat until converged {
Repeat $b = 1$ to B {
 $t \leftarrow t + 1$
 $d \leftarrow -\nabla L(\theta; S_b)$
 $v \leftarrow \beta_1 v + (1 - \beta_1) d$
 $r \leftarrow \beta_2 r + (1 - \beta_2) d^2$
 $\hat{v} \leftarrow v / (1 - \beta_1^t)$
 $\hat{r} \leftarrow r / (1 - \beta_2^t)$
 $\theta \leftarrow \theta + \alpha(\hat{v} + \epsilon) - \hat{r}$
}



SGD w/ Momentum:

Input parameters:
 ϵ - small number to regularize division
 α - momentum parameter
Learned parameters:
 γ, β - scaling and offset parameters
Estimated parameters:
 μ_1, σ_1^2 - learned mean and variance parameters
Parameters saved for inference:
 μ_1, σ_1^2 - learned mean and variance parameters
 γ, β, ϵ - same as above

Batch Inference:

For each batch (batch) {
 $\mu_B \leftarrow \frac{1}{K_B} \sum_{k=0}^{K_B-1} z_k$
 $\sigma_B^2 \leftarrow \frac{1}{K_B-1} \sum_{k=0}^{K_B-1} (z_k - \mu_B)^2$
For each input x {
 $\hat{z} \leftarrow \frac{z_k - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$
 $y_k \leftarrow \gamma \hat{z} + \beta$
}

Training:

Initialize $\mu_1 = 0, \sigma_1^2 = 0$
For each $k \in S_b$ {
 $\hat{z}_k \leftarrow \frac{z_k - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$
 $y_k \leftarrow \gamma \hat{z}_k + \beta$
}

Positional Encoding:

$$P(t, 2i) = \sin\left(\frac{2\pi i}{\lambda^t}\right)$$

$$P(t, 2i+1) = \cos\left(\frac{2\pi i}{\lambda^t}\right)$$

GRU:

$$\text{Loss}(\theta) = \sum_{k=0}^{K-1} \|y_k - g_{\theta}(f_{\theta}(y_k))\|^2$$

Autoencoder:



Generative Models:

Generate random vectors that "look" like ground truth
 $f_{\theta_d}(y) \approx P(\text{Class} = R|y)$
where
 $R_{\theta_d}(y) = \frac{P_R(y)}{P_R(y) + P_G(y)}$
 $R_{\theta_d}(y) = \frac{P_R(y)}{P_R(y) + P_G(y)}$

Discriminator Loss Function:

$$d(\theta_d, \theta_g) = \frac{1}{K} \sum_{k=0}^{K-1} [-\log f_{\theta_d}(y_k) - \log(1 - f_{\theta_d}(y_k))]$$

G/D Loss Function:

$$d(\theta_d, \theta_g, \theta_d) = E[-\log f_{\theta_d}(y)] + E[-\log(1 - f_{\theta_d}(y))]$$

2 Agents:

$$d(\theta_d, \theta_g) = \min_{\theta_d} d(\theta_d, \theta_g)$$

$$d(\theta_d, \theta_g) = \max_{\theta_d} d(\theta_d, \theta_g)$$

Nash Equilibrium:

$$d(\theta_d, \theta_g) = E[-\log f_{\theta_d}(y)] + E[-\log(1 - f_{\theta_d}(y))]$$

Zero Sum:

$$d(\theta_d, \theta_g) = \min_{\theta_d} d(\theta_d, \theta_g)$$

$$d(\theta_d, \theta_g) = \max_{\theta_d} d(\theta_d, \theta_g)$$

Dropout:

For each layer l {
- or each layer l {
 $\tilde{y}^{(l)} \leftarrow \tau^{(l)} * y^{(l)}$
 $z^{(l+1)} \leftarrow W^{(l+1)} * \tilde{y}^{(l)}$
 $y^{(l+1)} \leftarrow f(z^{(l+1)})$
}

Dropout:

For each batch {
Dropout - For each layer l {
randomly drop nodes }
- or each layer l {
 $\tilde{y}^{(l)} \leftarrow \tau^{(l)} * y^{(l)}$
 $z^{(l+1)} \leftarrow W^{(l+1)} * \tilde{y}^{(l)}$
 $y^{(l+1)} \leftarrow f(z^{(l+1)})$
}

Dropout:

For each layer l {
- or each layer l {
 $\tilde{y}^{(l)} \leftarrow \tau^{(l)} * y^{(l)}$
 $z^{(l+1)} \leftarrow W^{(l+1)} * \tilde{y}^{(l)}$
 $y^{(l+1)} \leftarrow f(z^{(l+1)})$
}