

Week 1: Intro to GANs

Generative Models:

- * What are generative models?
- * Types of generative models
 - GANs included!

Generative Models vs. Discriminative Models

Discriminative Models

Used for classification in ML. They learn how to distinguish between classes such as dogs and cats. Often called classifiers. Takes a set of features X , such as having a wet nose, puffs.

$$\begin{array}{c} \text{Features} \\ X \end{array} \rightarrow \begin{array}{c} \text{CLASS} \\ Y \end{array}$$

They try to model the probability,

$$P(Y|X)$$

Class Y , given a set of features.

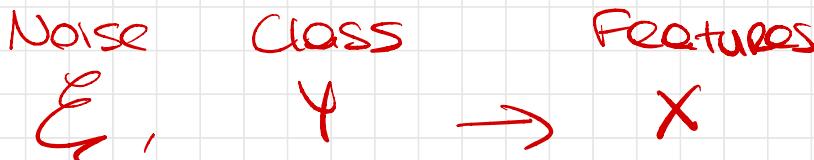
Generative Models

Tries to learn how to make a realistic representation of some class. They take some random input represented by the noise here, ϵ

ϵ can take a random number 3, -5, 2.6... the noise represents a random set of values going into the generative model. The generative model sometimes takes in a class Y such as a dog.

↳ Generative Models:

It's goal is to generate a set of features x that look like a realistic dog.



Why we need this noise?

↳ Generative Models try to capture the probability distribution of x .
The noise is to ensure that what's generated isn't actually the same dog each time.

$$\rightarrow P(x | y)$$

features

Realistic and diverse representation.

dog

Types of Generative Models

* Variational Autoencoders:

WORKS with two model "encoder" and "decoder". These are typically a neural network.

Encoder → learns by feeding realistic images into the encoder.

Encoder's job is to find a good way of representing that image in latent space. And the suitable point in the latent space can be represented by this vector of numbers

6.2, -3, 21.

Decoder → We put this vector to decoder and its goal is to reconstruct the realistic image that the encoder saw before.

After training, we lop off the encoder and we can pick random points in the latent space and the decoder will have learned to produce a realistic image of a dog.

⚠ The part we covered is the autoencoder part. Variational part adds some noise into this whole model and training process.

Instead of having the encoder encode the image into a single point in that latent space, the encoder encodes the image into a whole distribution and then samples a point on that distribution to feed into the decoder to then produce a realistic image. This adds a little bit of noise since different points can be sampled on this distribution.

* GANs



produce realistic image

Generator
+ Decoder

[1.2]
[3]
[-5]

Discriminator



Real or fake?



Generative models try to mimic the distribution of the data you're using to train them on.

GANS

Generators

* Generates fake
that look real

Discriminators

* learns to distin-
guish real from
fake

both
real and
fake images
in a pile

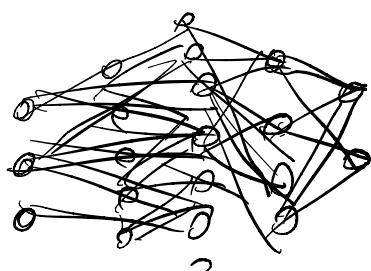
each of those
↑ nodes

(θ)

Parameters

After prediction
tell it whether
it was real or
fake

update parameters



Output \hat{y} → Cost func.

Discriminator

X →

Y

* Discriminator is a classifier

* learn the probability of class Y given

features X ($P(Y|X)$)

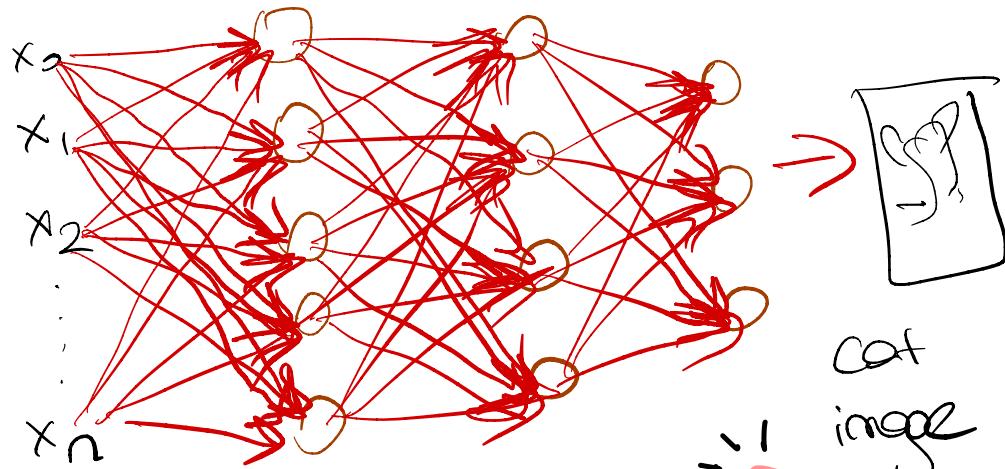
* The probabilities are the feedback of generator

Generator

If you trained it from the class of a cat, then the generator will do some computations and output a representation of a cat that looks real.



$$\begin{bmatrix} 1 \\ 2.5 \\ -7 \\ 10 \\ 0.3 \\ -8 \end{bmatrix}$$



Noise vector + class y

(E)

↳ sometimes if we have more than one class

* It learns the mimic that distribution of features x from the class of your data. In order to produce different outputs each time it takes random features as input.

feature → class

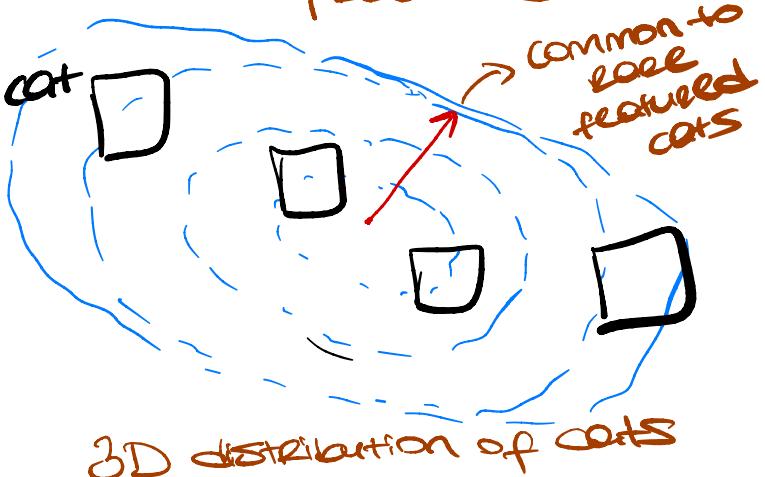
$$P(X|Y)$$

↳ conditional probability

Probability distributions:

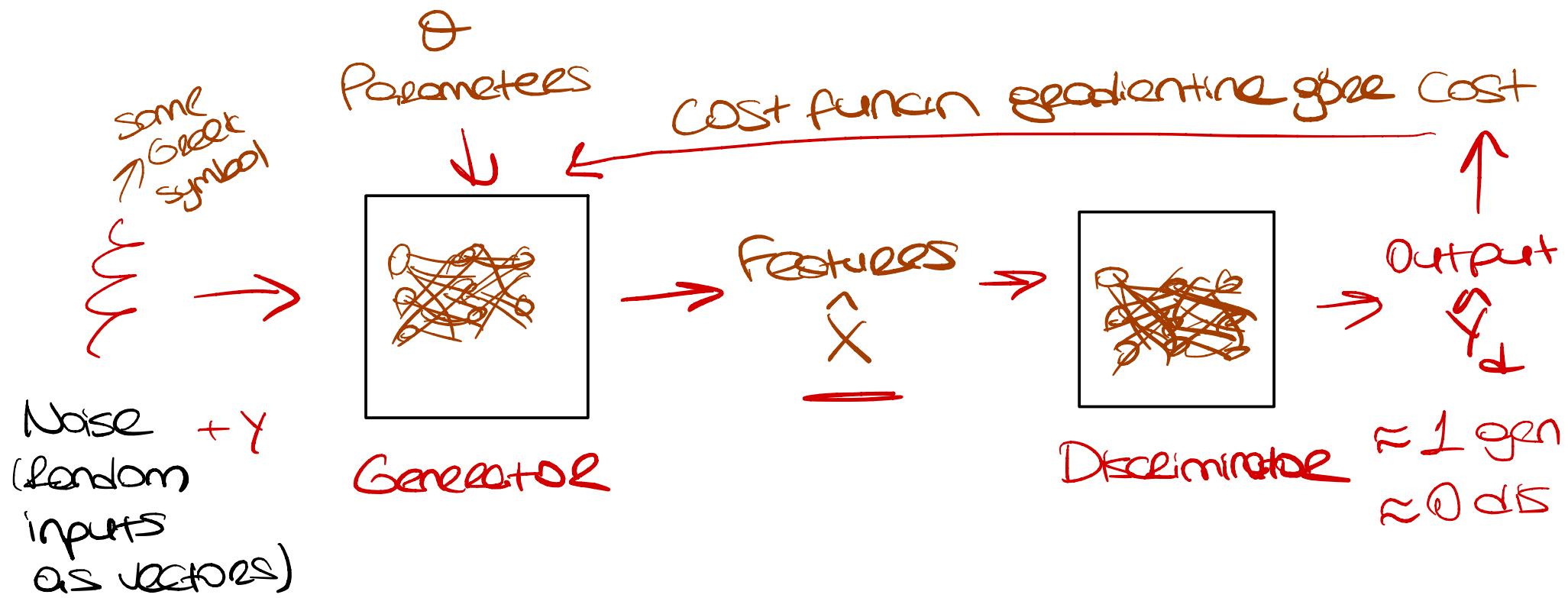
$$P(x)$$

we don't need label bc we now it's cat



All the features of a cat in the world

Generator: Learning



→ If it's good enough, you freeze the parameters of generator and you can generate images from new random noises.

BCE Cost Function (Binary Cross Entropy)

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h(x^{(i)}, \theta) + (1-y^{(i)}) \log(1-h(x^{(i)}, \theta))]$$

features

① $y^{(i)}$ $\log h(x^{(i)}, \theta)$
 Prediction $(0, 1)$

② $(1-y^{(i)})$ $\log(1-h(x^{(i)}, \theta))$
 true label $(0, 1)$

↓
 Average loss of the whole batch

① $y^{(i)} \log h(x^{(i)}, \theta) \rightarrow$ the product of true label y , times the log of the prediction.

| $y^{(i)}$ | $h(x^{(i)}, \theta)$ | $y^{(i)} \log h(x^{(i)}, \theta)$ |
|-----------|----------------------|-----------------------------------|
| 0 | any | 0 |
| 1 | 0.99 | ~ 0 (close to 0) |
| 1 | ~ 0 | $-\infty$ |

↓
Relevant when the label is 1.

bc of the log

* It makes it 0 if your prediction is good. And it makes it negative infinite when it's bad.

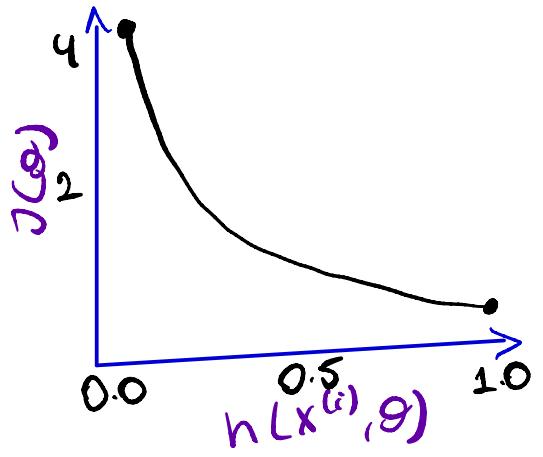
| $y^{(i)}$ | $h(x^{(i)}, \theta)$ | $(1-y^{(i)}) \log(1-h(x^{(i)}, \theta))$ |
|-----------|----------------------|--|
| 0 | any | 0 |
| 0 | 0.01 | ~ 0 |
| 0 | ~ 1 | $-\infty$ |

→ Relevant when the label is 0.

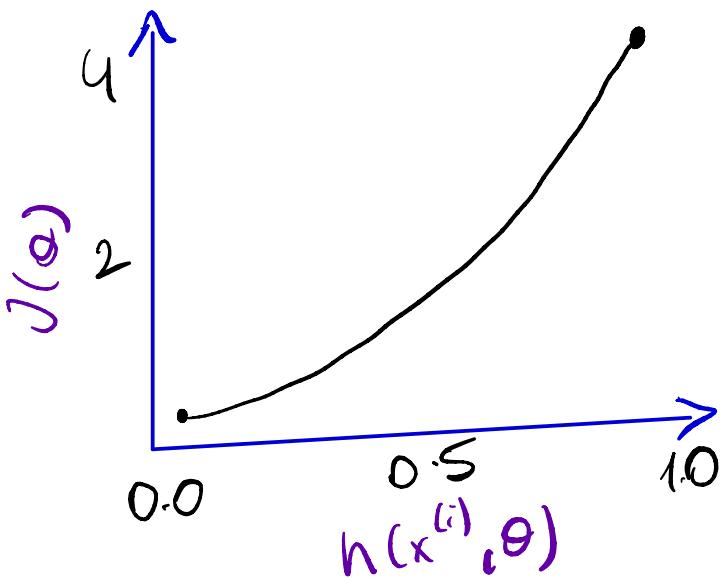


* logarithm of a value between 1-0 is calculated, which returns that negative result. That's why we want that negative term.

$$\frac{y^{(i)} \log h(x^{(i)}, \theta)}{\rightarrow \text{relevant when the label } 1}$$



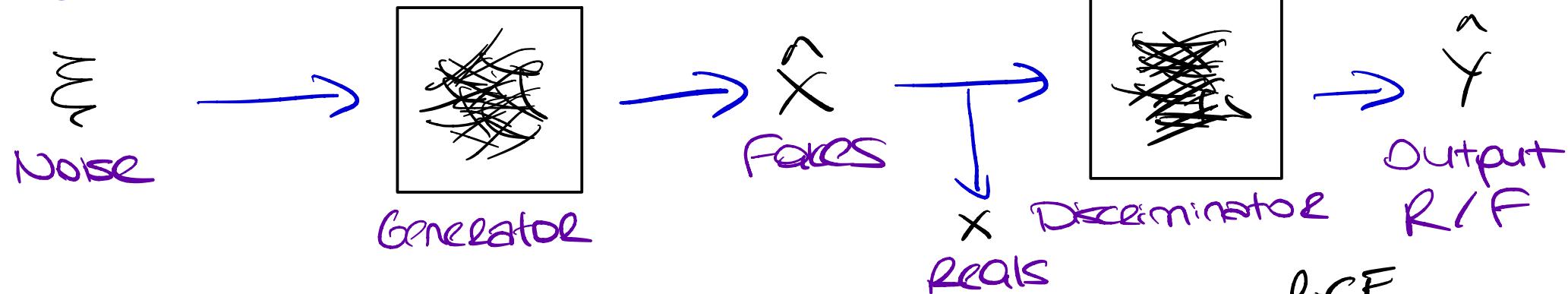
$$\frac{(1-y^{(i)}) \log (1-h(x^{(i)}, \theta))}{\rightarrow \text{relevant label } 0.}$$



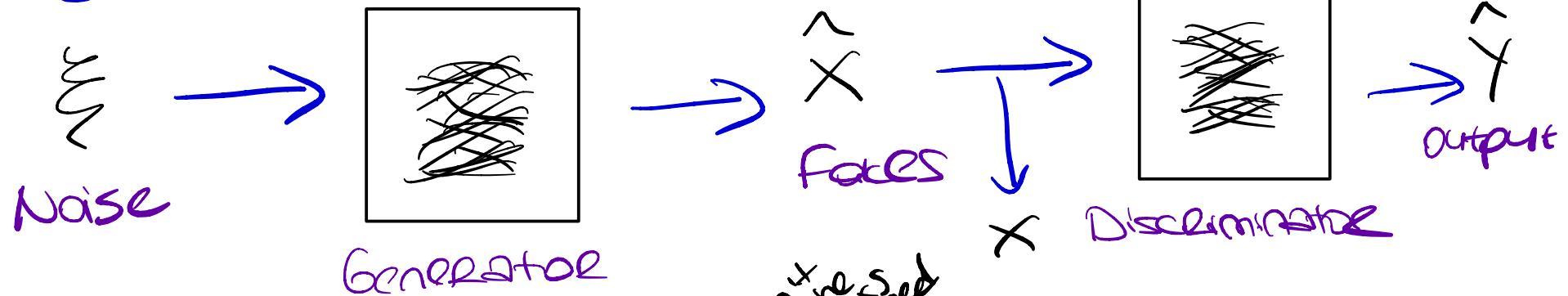
SUMMARY

- * The BCE cost function has two parts that relevant for each class.
 - * Close to zero when the label and the prediction are similar.
 - * Approaches infinity when the label and the prediction are different.
- ⇒ The BCE loss is performed across a mini batch of several examples and it takes the average of all those n examples.

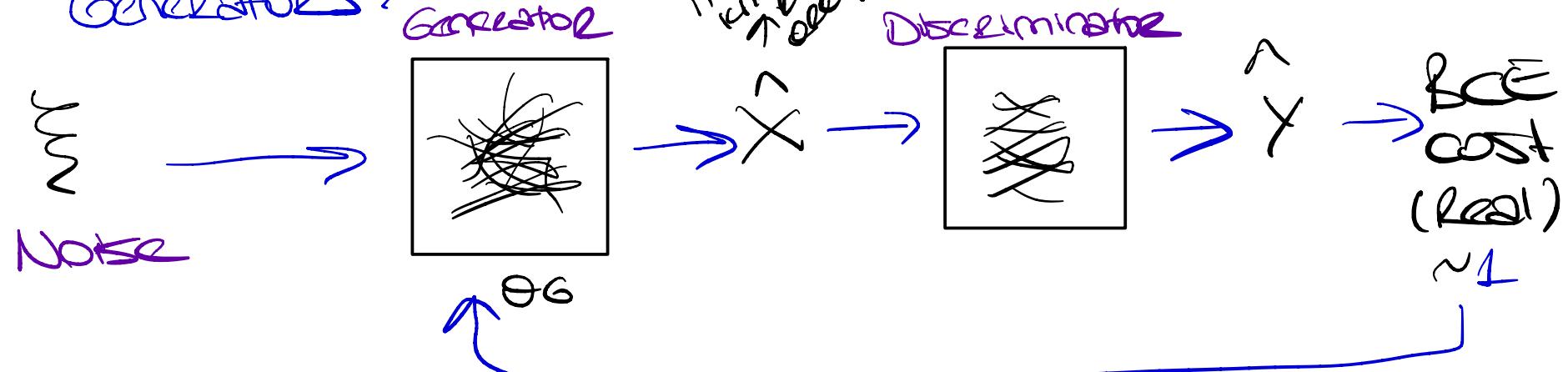
GANs Model:



Discriminator:



Generators:



$W^{(j)}$ → column vectors

↳ when the input of a tanh is far from zero, its gradient is very close to zero, because of the flat edges of tanh.)

↳ neuron number

* tanh always works better than sigmoid when used in hidden layers.
↳ exception is for the output layer in classification.

* Logistic regression doesn't have a hidden layer. If you initialize weights to zero, the first example x fed into the f will output zero but the derivatives of the logistic regression depends on the input x , which is not zero. So, at the second iteration, the weights' values follows x 's distribution and are different from each other if x is not a constant vector.

$W^{(k)}$ → row is the number of neurons of k^{th} layer
→ column is the " " " input of k^{th} .

* Output of the tanh is between -1 and 1. It thus centers the data which makes the learning simpler for the next layer.

hidden layer
P with tanh

weight
initialization

* The use of random numbers helps to "break the symmetry" between all the neurons allowing them to compute different functions. When using small random numbers the values $z^{[k]}$ will be close to zero thus the activation values will have a large gradient speeding up the training process.

WEEK 2

ACTIVATIONS : are functions that take any real number as input also known as its domain and outputs a number in a certain range using a non-linear differentiable function

Why non-linear and differentiable?

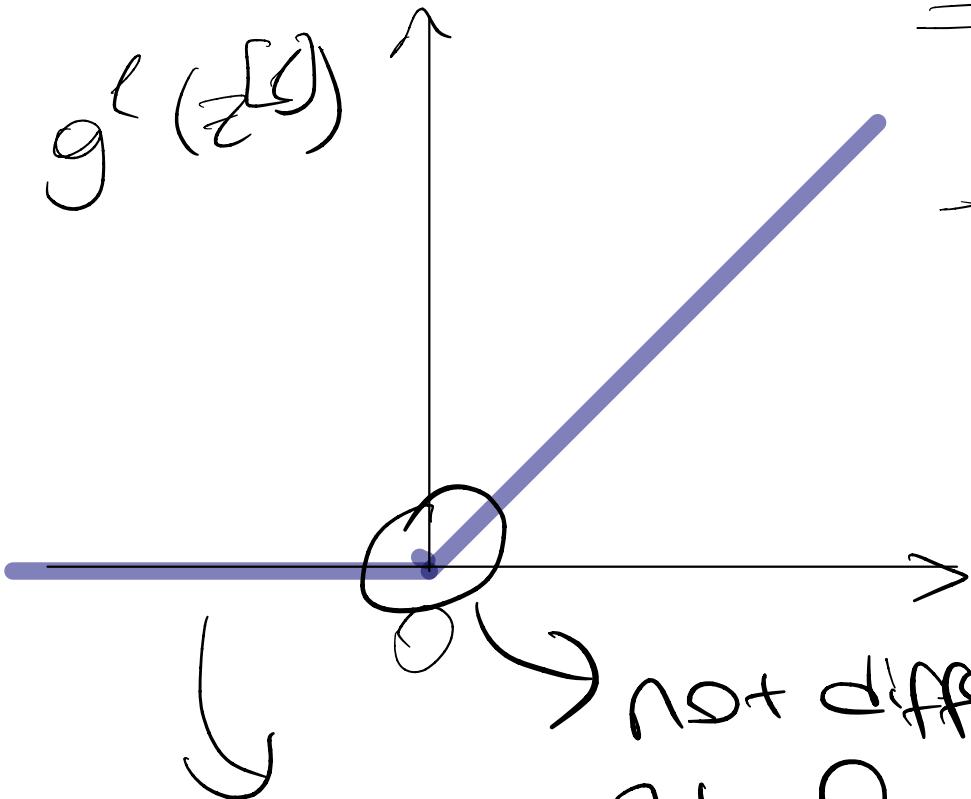
- 1) Differentiable for backpropagation because you use backpropagation for training and updating its parameters, so it needs to be able to differentiate and provide gradient to previous layer.
- 2) It needs to be non-linear, the features you compute within the neural networks can be complex.

↳ So you need non-linear differentiable activation functions to take advantage of deep learning models and to have the layers stack onto each other to construct a complex neural network to model complex nonlinear functions.

* Non-linearity ensures that your linear layers don't reduce to a single layer

Common Activation Functions

Relu → Rectified Linear Unit



$$= \max(0, z^l)$$

→ no-negative values are allowed.

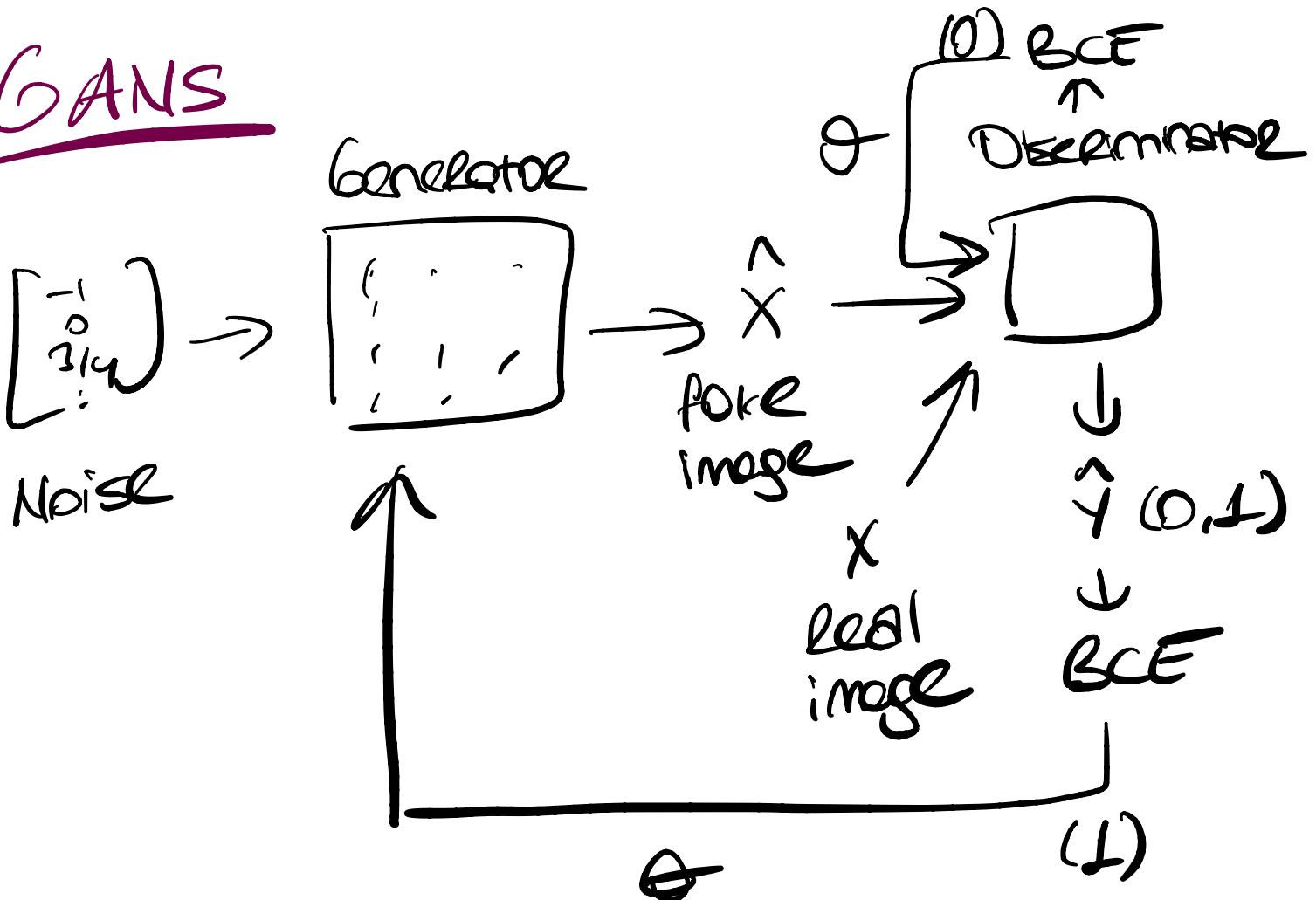
this part
z is negative
always has a derivative
equal to zero

Since the learning process depends on the derivative to provide important information on how to update the weight with some values and their weights

Stop learning. All the components will be affected as well.

This is dying relu problem.

GANS

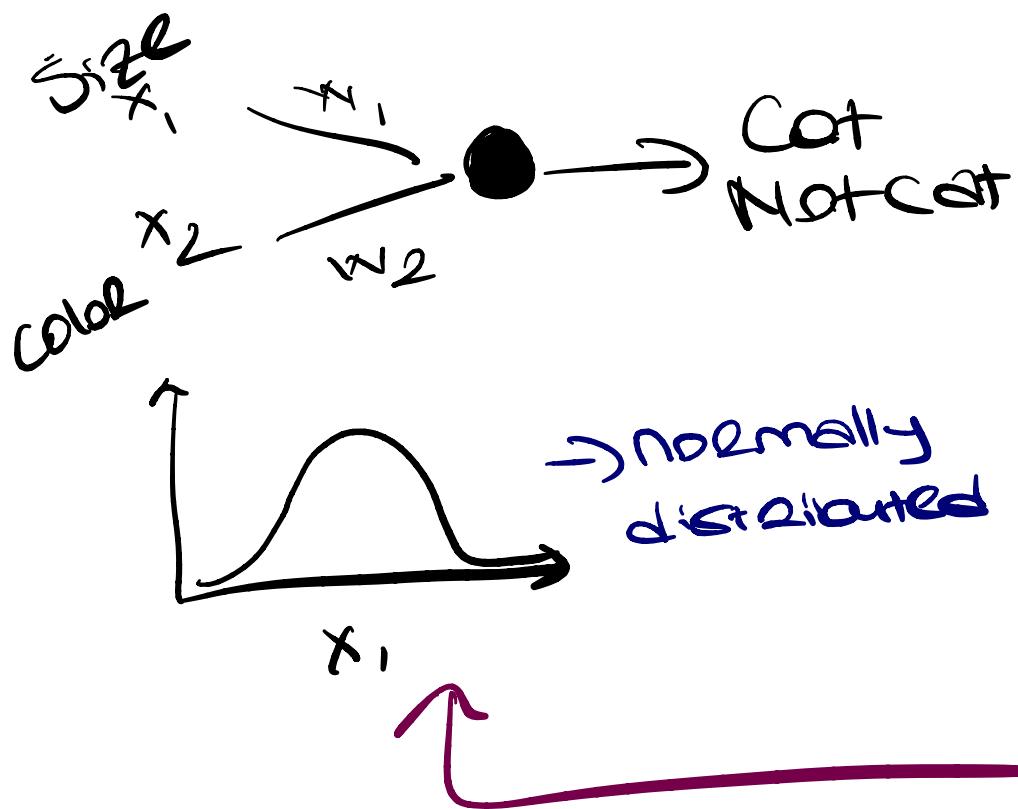


Batch Normalization

* Normalization

* Internal Covariate Shift

* Batch Normalization

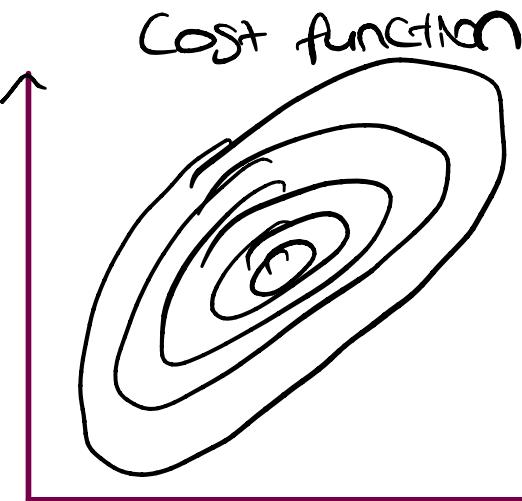


Skew Σ
little bit
towards higher
values \uparrow



(↳ higher mean
(and SD))

Comparing the
same value on
two features is
not possible

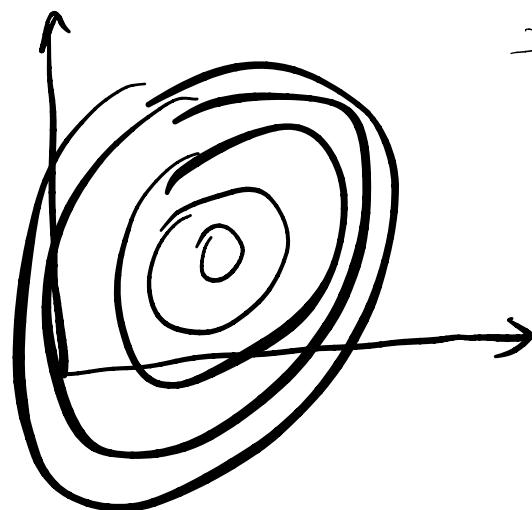


* Different distributions across inputs, cost function will be **elongated**.

So the changes to the weights relating to each of the inputs will have kind of a different effect of varying impact on this cost functions.

* This makes training difficult and slow. And highly dependent on how your weights are initialized.

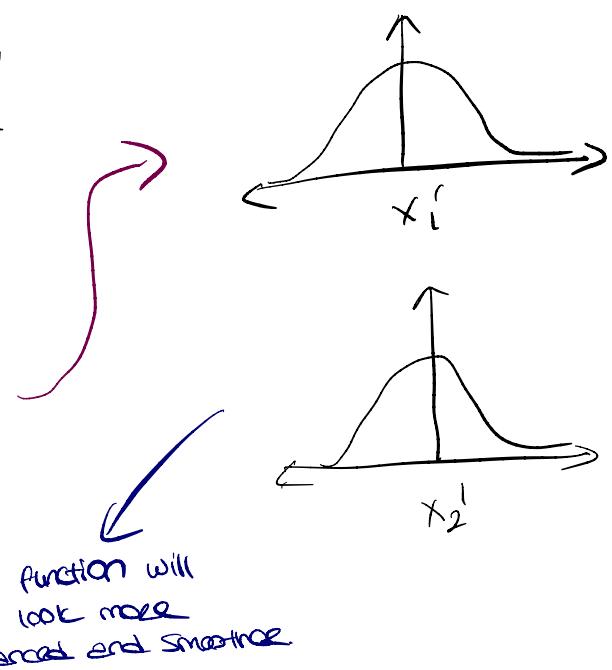
- If train or test data has different distributions, kinds of shifts or changes in some way then the form of cost function could change too.



→ Little bit more round and location of the minimum could also move. This known as "**covariance shift**".



What if the input variables x_1 and x_2 are normalized?



cost function will also look more balanced and smoother

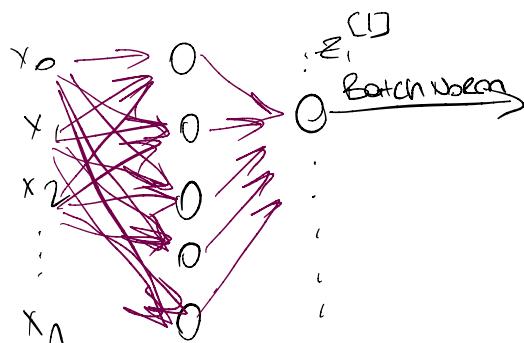
Normalized in same way, around of 0

and a standard deviation of 1.

- For training data, you train 2 batch. You take the mean and standard deviation and shift it to be around mean 0 and sd of 1.

- For the test data, look at the statistics that we're gathered over time as you went through the training set and use those to center the test data to be closer to the training data.

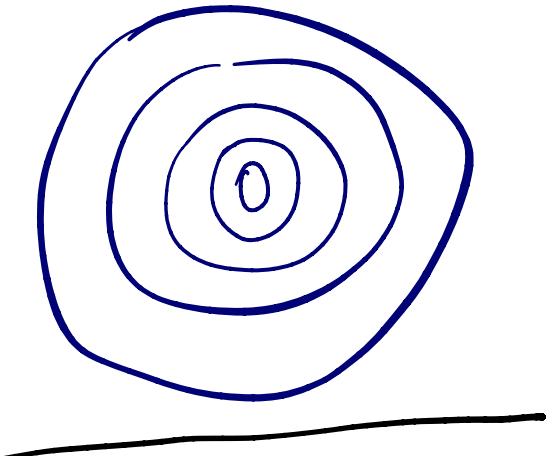
* Covariant Shift Shouldn't be a problem if you just make sure that the distribution of your data set is similar to the task you're modelling.



All the weights that effect this activation value are updated. So all of the weights are updated. The distribution of values contained in that activation changes.

→ This makes training process difficult due to the shifts.

→ B.N normalizes all these internal nodes based on statistics calculated for each input batch.

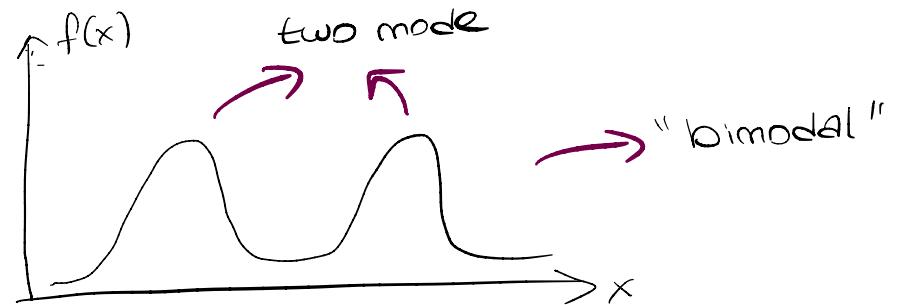
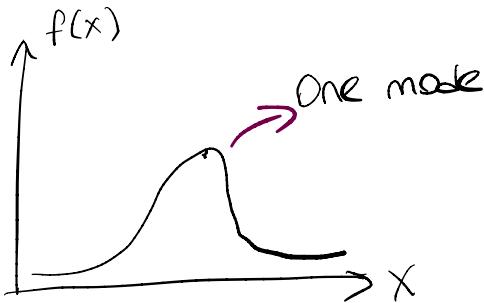


Smoothing the cost function and reducing the covariance shift.

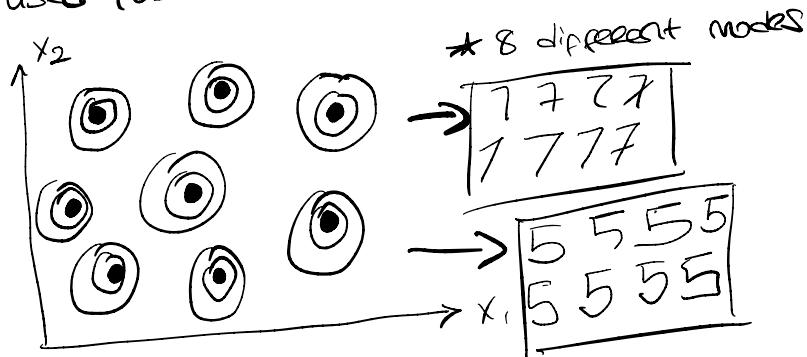
→ Internal Covariance Shift

Mode Collapse:

- * Any peak on the density function is a mode.



Many real life distributions used for GANs are multimodal.



↳ Multimodal

- * Sounds like things collapsing on one mode or fewer modes and the modes disappearing.



Discriminator

Good at identifying fake images but generates images that look like 1 and 7.

This means the discriminator is at of local minima of cost function.

Then this information is passed on to the generator.

It generates 1 and 7's.

| | | |
|---|---|---|
| 1 | 1 | 7 |
| 1 | 7 | 1 |

Passes to the discriminator

Missclassifies every picture except for one maybe.

Passes to the generator

Sees that discriminator's weakness is with 1.

→ Produce images

| | | | |
|---|---|---|---|
| 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 |

↓
Collapsing to a single mode.



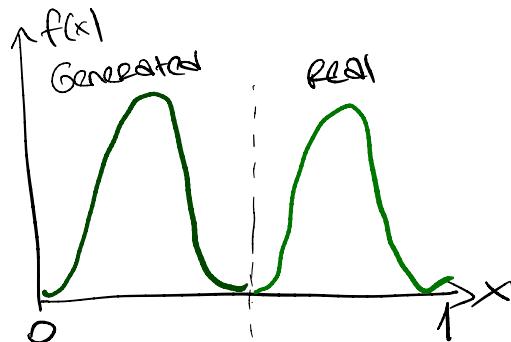
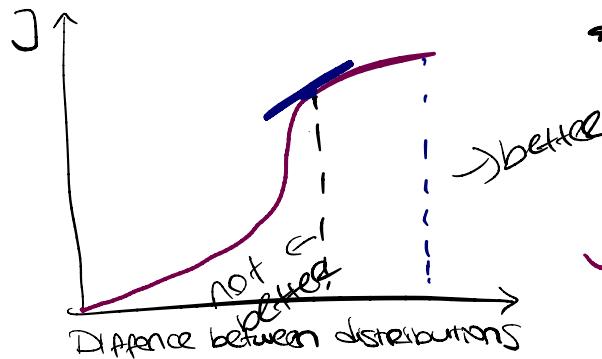
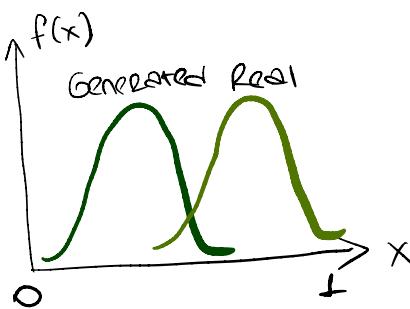
The generator sees this and looks at the feedback from the discriminator and gets a good idea of how to fool the discriminator in the next round.

Problem with BCE Loss

- Discriminator is easier to train.
- Generator is difficult to train.

↳ Complex output

But beginning in training this isn't such a problem. Because discriminator is not that good.



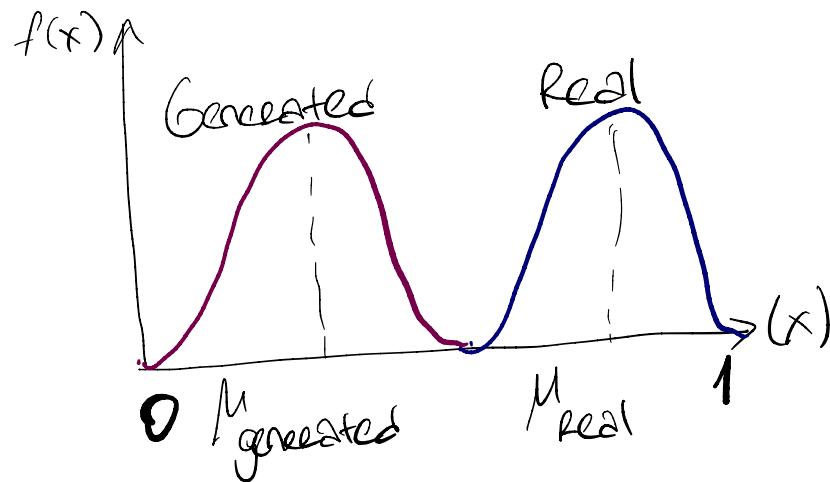
- * Discriminator improves during training and sometimes improves more easily than the generator, that underlying cost function will have those flat regions when the distributions are very different from one another

* It has problem distinguishing generated from real. As a result, it's able to give useful feedback in the form of a non-zero gradient back to the generator.

→ As discriminator getting better it'll start giving less informative feedback. It might give gradients close to zero and becomes unhelpful to the generator. Because the generator doesn't know how to improve.

How the vanishing gradient problem will arise.

Earth Mover's Distance (EMD)



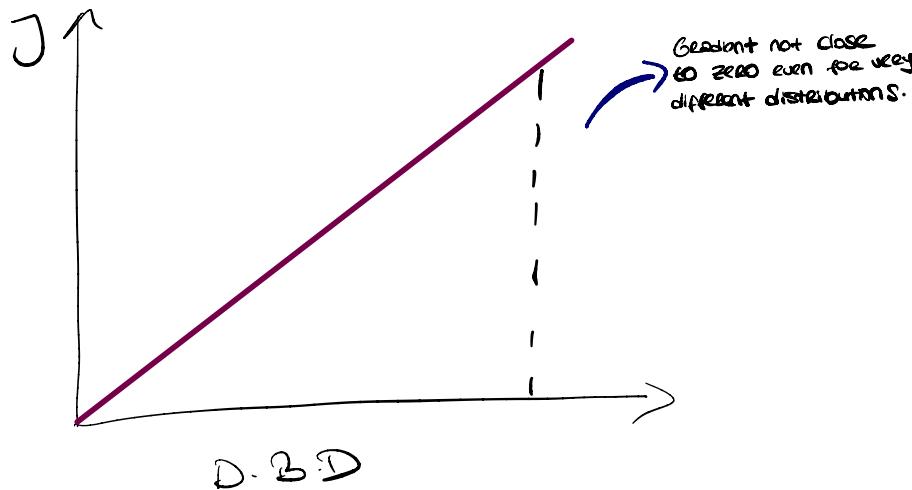
The function depends on the distance and amount moved.

* What EMD does is it measures how different these two distributions are, by estimating the amount of effort it takes to make the generated distributions equal to the real.

* with EMD there's no such a ceiling to the zero and one. So cost function continues to grow regardless how far apart these distributions are.

* Doesn't have flat regions when the distributions are very different.

* EMD solves the problem associated with BCE.



Wasserstein Loss

- W-Loss approximates the Earth Moon's distance.
- * Generator's trying to maximize the cost, discriminator's trying to minimize the cost.

→ This is often called
minimax game!

$$\min_d \max_g - \left[E(\log(d(x))) + E(1 - \log(d(g(z)))) \right] \rightarrow \begin{array}{l} \text{Simplified} \\ \text{BCE loss} \end{array}$$

\downarrow \downarrow

How bad disc
classifies real
observations.

How bad disc
classifies fake
observations.

W1-Loss: the function calculates the difference between the expected values of the predictions of the discriminator. Called "critic".

$$\min_g \max_c \mathbb{E}(c(x)) - \mathbb{E}(c(g(z)))$$

Discriminator wants to minimize the distance between its thoughts on the real vs fake.

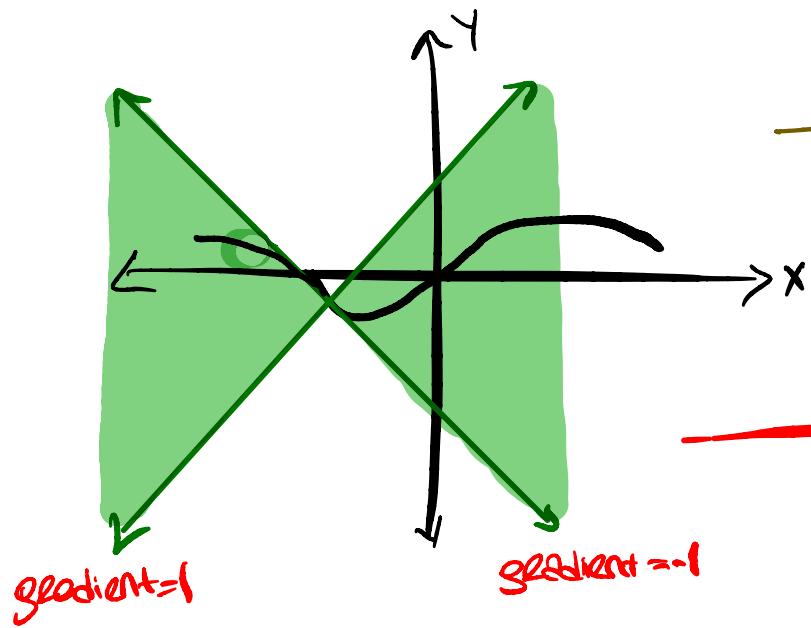
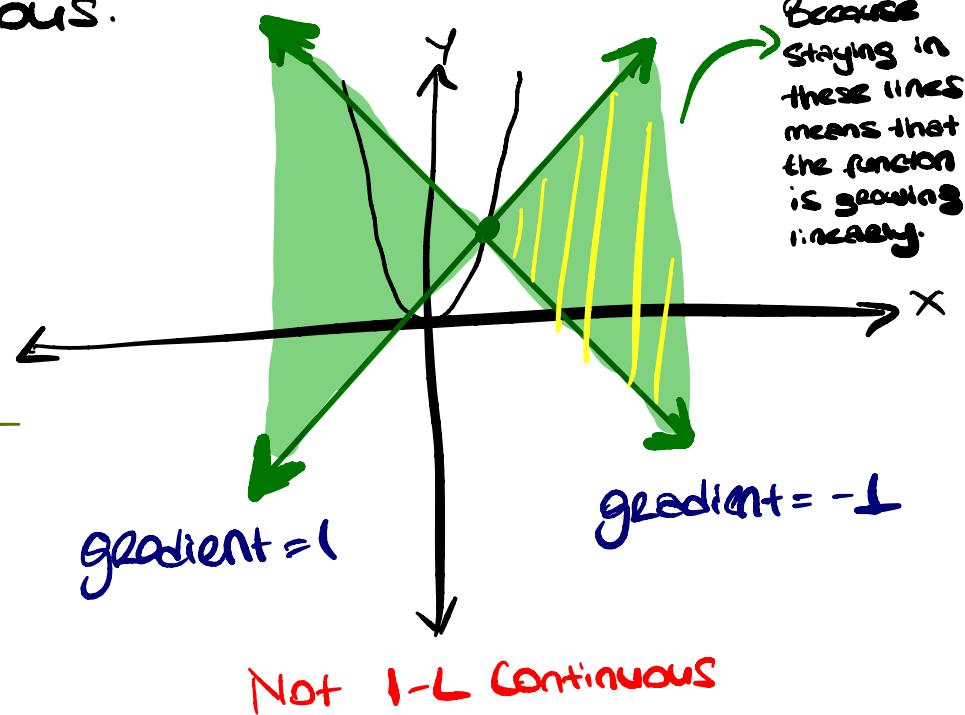
Generator wants to minimize this bc it want the discriminator to think that its fake images are as close as possible to the reals.

1-Lipschitz Continuous (1-L continuous)

Critic needs to be 1-L Continuous.

The norm of it's gradient
needs to be at most one
for every point.

You want to
make sure that
growth of this
function never
goes out of bounds
from these lines.



* This condition on the critics is important for VI-LOSS bc it assures that the VI-LOSS is not only continuous and differentiable but also that it doesn't grow too much.

1-L Enforcement on Critic

1-L Continuous $\rightarrow \|\nabla f(x)\|_2 \leq 1 \rightarrow$ Norm of gradient at most L_2 norm

Gradient critic image

Weight Clipping

Forces the weights of the critic to a fixed interval. After you update the weights during gradient descent, you'll clip any weight outside of the desired interval. The weights over that interval either too high or too low, will be set to the maximum or the minimum amount allowed.



It has downside!

Forcing the weights of the critic to a limited range of values could limit the critics ability to learn and ultimately for the gradient to perform because if the critic can't take on many different parameter values, it's weights can't take on many different values.

Gradient Penalty

$$\min_{\mathcal{G}} \max_{\mathcal{C}} \mathbb{E}(c(x)) - \mathbb{E}(c(g(z))) + \lambda \text{reg}$$

regularization

HS penalizes the critic when its gradient norm is higher than one.

Checking the gradients at every possible point of the feature space, not practical. Sample some points by interpolating between real and fake examples.

random interpolation

