

Class Notes on GANs Models

Jack Li

October 24, 2024

Contents

1	Introduction	1
2	Basic Math	2
2.1	Probability Theory	2
2.2	Sampling from a Probability Distribution: Why $P(x) dx$?	2
2.2.1	Mathematical Expression	2
2.2.2	Why $P(x) dx$?	2
2.2.3	Example of Sampling	2
2.3	Standard Sampling (Non-Differentiable)	2
2.4	Reparameterization Trick (Differentiable)	3
2.5	Gumbel-Softmax	3
2.5.1	Definition	3
2.5.2	Proof of Differentiability	3
2.6	Concrete Distribution	4
2.6.1	Definition	4
2.6.2	Differentiability	4
2.7	Pros and Cons	4
2.8	Use Cases	4
2.8.1	Gumbel-Softmax	4
2.8.2	Concrete Distribution	4
2.9	Equations Summary	4
2.9.1	Gumbel-Softmax Equation	4
2.9.2	Concrete Distribution Equation	5
2.10	Linear Algebra	5
2.11	Optimization	5
3	PyTorch Basics	5
3.1	Tensors	5
4	PyTorch training gradients	5
4.1	Autograd	6
4.2	Building Neural Networks	7
4.3	Loss Functions	8
5	GANs Models	9
5.1	Basic GAN	9
5.2	DCGAN	9
5.3	WGAN	9
5.4	CycleGAN	9
6	Conclusion	9

1 Introduction

Provide an introduction to GANs and their importance in machine learning.

2 Basic Math

2.1 Probability Theory

- Definitions of probability, random variables, expectation, etc.

2.2 Sampling from a Probability Distribution: Why $P(x) dx$?

In continuous probability theory, the **probability density function** (PDF) $P(x)$ represents the density of the probability at a particular point x . However, the actual probability of the random variable X falling within a small interval around x , say $[x, x + dx]$, is given by the product of the PDF at x and the small interval dx .

2.2.1 Mathematical Expression

The probability of X falling within the interval $[x, x + dx]$ is approximately:

$$P(X \in [x, x + dx]) \approx P(x) dx$$

Where: - $P(x)$ is the probability density function evaluated at x , - dx is an infinitesimally small interval around x .

2.2.2 Why $P(x) dx$?

- The PDF $P(x)$ by itself does not give the actual probability for any specific value of x , because for continuous random variables, the probability of any single point is zero:

$$P(X = x) = 0 \quad \text{for continuous variables.}$$

- Instead, the **probability** is found by integrating the PDF over an interval:

$$P(X \in [a, b]) = \int_a^b P(x) dx$$

For a very small interval dx , the integral simplifies to:

$$P(X \in [x, x + dx]) \approx P(x) dx$$

This shows that the product $P(x) dx$ gives the probability mass within the small region $[x, x + dx]$.

2.2.3 Example of Sampling

When sampling a value x from a distribution, the probability that a sample lies within a small range $[x, x + dx]$ is proportional to $P(x) dx$. For example, in the case of a Gaussian (Normal) distribution, the probability density is given by:

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

To find the probability of the random variable X falling within the range $[x, x + dx]$, we compute:

$$P(X \in [x, x + dx]) \approx P(x) dx$$

2.3 Standard Sampling (Non-Differentiable)

In the standard sampling method, we sample directly from a normal distribution, which does not allow gradients to propagate:

```
1 # Non-differentiable direct sampling
2 mu = torch.tensor(0.0, requires_grad=True)
3 sigma = torch.tensor(1.0, requires_grad=True)
4
5 # Direct sampling from a normal distribution
6 z = torch.normal(mu, sigma)
7
```

```

8 # Loss function
9 loss = (z - 5) ** 2
10 loss.backward() # This breaks the gradient flow

```

This will not compute gradients correctly because the sampling is a discrete, non-differentiable operation.

2.4 Reparameterization Trick (Differentiable)

In the reparameterization trick, we express the random variable z as a function of μ , σ , and a noise term ϵ , which allows backpropagation to compute gradients:

$$z = \mu + \sigma \cdot \epsilon$$

```

1 # Differentiable sampling using reparameterization trick
2 mu = torch.tensor(0.0, requires_grad=True)
3 log_sigma = torch.tensor(0.0, requires_grad=True)
4 sigma = torch.exp(log_sigma)
5
6 # Sample epsilon from N(0, 1)
7 epsilon = torch.randn_like(sigma)
8
9 # Reparameterization: z = mu + sigma * epsilon
10 z = mu + sigma * epsilon
11
12 # Loss function
13 loss = (z - 5) ** 2
14 loss.backward() # Gradient flows through mu and sigma

```

Here, ϵ is sampled from a standard normal distribution $\mathcal{N}(0, 1)$, and the parameters μ and σ can be updated by gradient descent because the whole process is now differentiable.

2.5 Gumbel-Softmax

2.5.1 Definition

The Gumbel-Softmax trick is a method for sampling from a categorical distribution in a differentiable way. It provides a continuous approximation to a categorical distribution, making it suitable for use in backpropagation.

Let $\pi_1, \pi_2, \dots, \pi_k$ be the probabilities of each category in a categorical distribution. The Gumbel-Softmax trick works by introducing Gumbel noise g_i for each category i , sampled from a Gumbel distribution:

$$g_i = -\log(-\log(u_i)), \quad u_i \sim \text{Uniform}(0, 1)$$

The continuous approximation to the categorical distribution is given by the softmax function:

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)}$$

where τ is the temperature parameter controlling the sharpness of the approximation. As $\tau \rightarrow 0$, the distribution becomes more discrete (closer to one-hot), and as $\tau \rightarrow \infty$, the distribution becomes more uniform.

2.5.2 Proof of Differentiability

Let $\pi = (\pi_1, \pi_2, \dots, \pi_k)$ be the parameter of the categorical distribution. The reparameterization with Gumbel noise is differentiable because:

$$\frac{\partial y_i}{\partial \pi_j} = \frac{\partial}{\partial \pi_j} \left(\frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \right)$$

Since g_i is independent of π_j , the function remains differentiable.

2.6 Concrete Distribution

2.6.1 Definition

The Concrete distribution is similar to the Gumbel-Softmax but is defined for both binary and multinomial distributions. The Concrete distribution for binary random variables introduces a continuous relaxation by using the logistic sigmoid function for binary variables.

For a binary random variable with probability p , the Concrete distribution samples:

$$z = \frac{\log(p) + g}{\tau}$$

where $g \sim \text{Gumbel}(0, 1)$ and τ is the temperature parameter. The output z is passed through the sigmoid function to produce a value between 0 and 1:

$$y = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

2.6.2 Differentiability

As in the Gumbel-Softmax case, the Concrete distribution is differentiable because the sampling is reparameterized using Gumbel noise, and the final step uses differentiable functions like the softmax or sigmoid.

2.7 Pros and Cons

Method	Pros	Cons
Gumbel-Softmax	<ul style="list-style-type: none">- Differentiable approximation of categorical variables.- Control over discreteness via temperature τ.- Easy to implement for multiclass problems.	<ul style="list-style-type: none">- Approximation error due to continuous relaxation.- Gradients vanish as $\tau \rightarrow 0$.- Limited to categorical (one-hot) outputs.
Concrete Distribution	<ul style="list-style-type: none">- Differentiable for both binary and categorical variables.- More flexibility than Gumbel-Softmax for binary decisions.	<ul style="list-style-type: none">- Same vanishing gradient issue as Gumbel-Softmax for low τ.- Still an approximation of a discrete distribution, not exact.

Table 1: Comparison of Gumbel-Softmax and Concrete distribution

2.8 Use Cases

2.8.1 Gumbel-Softmax

- **Generative Models**: Used in **Variational Autoencoders (VAEs)** for categorical latent variables, allowing for discrete sampling while maintaining differentiability. - **Reinforcement Learning**: In policy gradient methods, Gumbel-Softmax can be used for discrete action selection.

2.8.2 Concrete Distribution

- **Binary Decision Making**: Useful in models where decisions are binary, like in the **Binary Variational Autoencoder (Binary VAE)**. - **Binary Latent Variables**: Concrete distributions work well for models with binary latent variables, where a differentiable approximation is required.

2.9 Equations Summary

2.9.1 Gumbel-Softmax Equation

For categorical variables, the Gumbel-Softmax approximation is given by:

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)}$$

where $g_i \sim \text{Gumbel}(0, 1)$.

2.9.2 Concrete Distribution Equation

For binary variables, the Concrete distribution is given by:

$$z = \frac{\log(p) + g}{\tau}, \quad g \sim \text{Gumbel}(0, 1)$$

and the relaxed binary sample is:

$$y = \frac{1}{1 + \exp(-z)}$$

2.10 Linear Algebra

- Vectors, matrices, eigenvalues, eigenvectors, etc.

2.11 Optimization

- Gradient descent, stochastic gradient descent, etc.

3 PyTorch Basics

3.1 Tensors

- Definition and operations on tensors.

4 PyTorch training gradients

Code 1: PyTorch training gradients

```
1  #SECTION: Gradient computation
2
3  # Step 1: Define a simple model
4  model = nn.Linear(1, 1)
5  optimizer = optim.SGD(model.parameters(), lr=0.01)
6
7  # Dummy input and target
8  input = torch.tensor([[1.0]], requires_grad=True)
9  target = torch.tensor([[2.0]])
10
11 # Step 2: Print the initial parameters
12 print("Initial parameters:")
13 for param in model.parameters():
14     print(param.data)
15
16 # Step 3: Forward pass
17 output = model(input)
18 loss = (output - target).pow(2).mean()
19
20 # Step 4: Zero the gradients
21 optimizer.zero_grad()
22
23 # Step 5: Backward pass
24 loss.backward()
25
26 # Step 6: Update the parameters
27 optimizer.step()
28
29 # Step 7: Print the parameters after the update
30 print("\nParameters after one training step:")
31 for param in model.parameters():
32     print(param.data)
```

1. Initialize Parameters:

- Assume initial weights w and bias b are both 0.

- Model: $y = wx + b$

2. Forward Pass:

- Compute the output: $\hat{y} = wx + b$
- Given input $x = 1.0$ and target $y = 2.0$:

$$\hat{y} = 0 \cdot 1.0 + 0 = 0$$

3. Compute Loss:

- Loss function: Mean Squared Error (MSE)

$$\text{Loss} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

- For our single data point:

$$\text{Loss} = (0 - 2.0)^2 = 4.0$$

4. Backward Pass (Gradient Calculation):

- Compute gradients of the loss with respect to w and b :

$$\frac{\partial \text{Loss}}{\partial w} = 2(\hat{y} - y)x = 2(0 - 2.0) \cdot 1.0 = -4.0$$

$$\frac{\partial \text{Loss}}{\partial b} = 2(\hat{y} - y) = 2(0 - 2.0) = -4.0$$

5. Parameter Update:

- Using Stochastic Gradient Descent (SGD) with learning rate $\eta = 0.01$:

$$w_{\text{new}} = w - \eta \frac{\partial \text{Loss}}{\partial w} = 0 - 0.01 \cdot (-4.0) = 0.04$$

$$b_{\text{new}} = b - \eta \frac{\partial \text{Loss}}{\partial b} = 0 - 0.01 \cdot (-4.0) = 0.04$$

6. Updated Parameters:

- After one training step, the new parameters are:

$$w = 0.04, \quad b = 0.04$$

Summary - Initial parameters: $w = 0, b = 0$ - After one training step: $w = 0.04, b = 0.04$

4.1 Autograd

- Automatic differentiation in PyTorch.

The active selection `gradient.norm(2, dim = 1)` is a PyTorch operation that computes the L2 norm (Euclidean norm) of the `gradient` tensor along a specified dimension. In this case, the dimension specified is `dim = 1`.

$$\Theta = \underset{\Theta}{\operatorname{argmin}} \frac{1}{B} \sum_{i=1}^B \left[D(z_i, \Theta) - D(y_i, \Theta) \right] + \lambda \left(\left\| \frac{\partial D(y, \Theta)}{\partial y} \right\| - 1 \right)^2$$

Detailed Explanation:

1. L2 Norm (Euclidean Norm):

- - The L2 norm of a vector is a measure of its magnitude and is calculated as the square root of the sum of the squares of its components. Mathematically, for a vector v , the L2 norm is given by $\|v\|_2 = \sqrt{\sum v_i^2}$.
- - In PyTorch, the `norm` function can compute various types of norms, with the L2 norm being specified by the argument 2.

2. Dimension Specification ($dim = 1$):

- - The *dim* argument specifies the dimension along which the norm is computed. In a multi-dimensional tensor, this allows you to compute norms along specific axes.
- - For example, if *gradient* is a 2D tensor (matrix) with shape $[batchsize, numfeatures]$, setting $dim = 1$ means that the norm is computed for each row independently. This results in a tensor of shape $[batchsize]$, where each element is the L2 norm of the corresponding row in the original tensor.

Code 2: PyTorch gradient sampling example

```
1  # Define the sampling function
2  def sample_function(x):
3      return torch.sin(x)
4
5
6  # NOTE: Generate sample points with requires_grad=True, and need requires_grad=True
7  # Generate sample points with requires_grad=True
8  x = torch.tensor(
9      np.linspace(0, 2 * np.pi, 100), dtype=torch.float32, requires_grad=True
10 )
11
12 # Define f by sampling from the sample_function
13 f = sample_function(x)
14
15 # Compute the gradient of f with respect to x
16 grad = torch.autograd.grad(outputs=f, inputs=x, grad_outputs=torch.ones_like(f))
17
18 print(f"The gradient of f(x) = sin(x) at x = {x} is {grad[0]}")
```

4.2 Building Neural Networks

- Layers, activation functions, loss functions, etc.

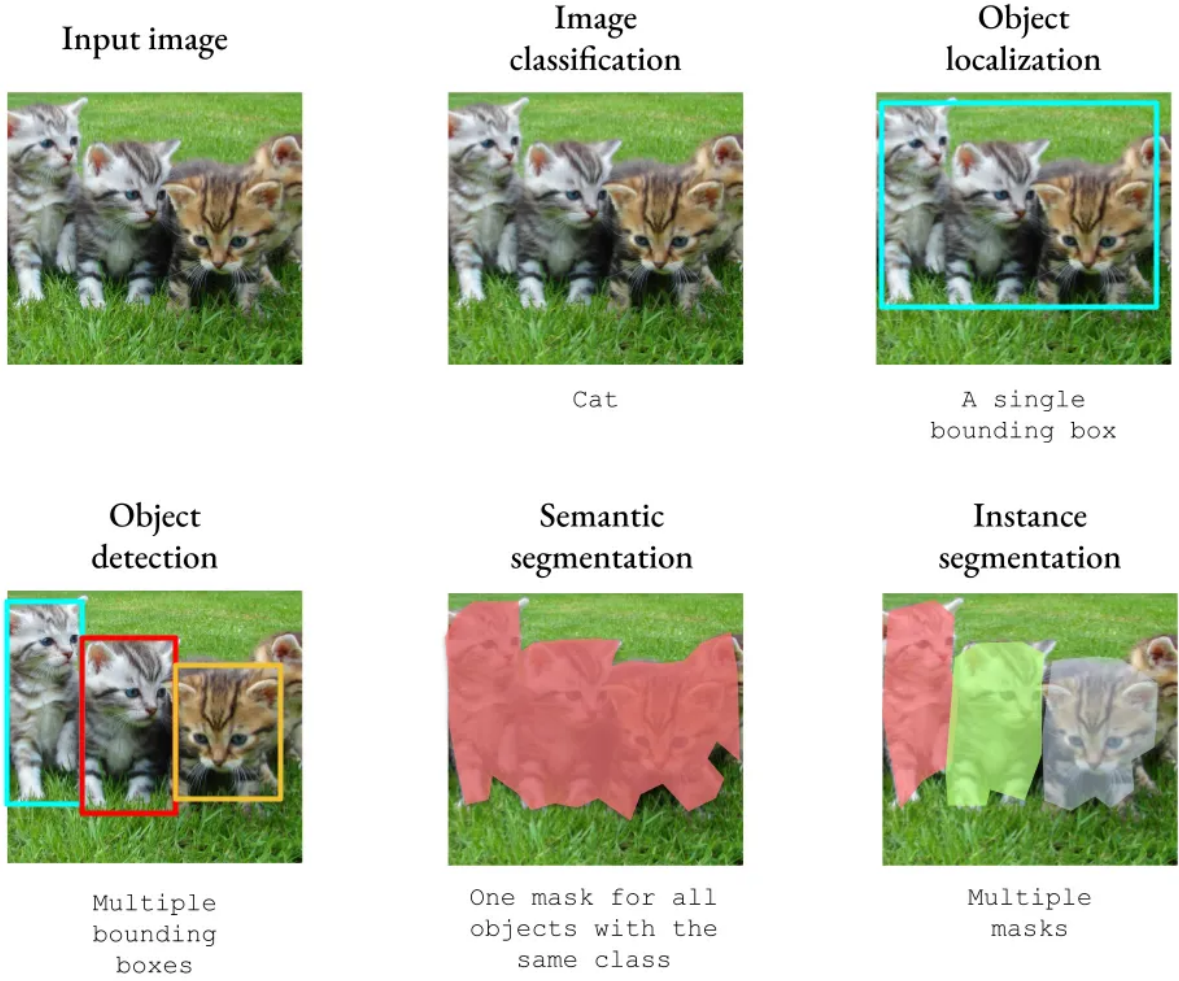


Figure 1: An example image illustrating segmentation.

4.3 Loss Functions

— PyTorch Conv2d Equation

The output size of a Conv2d layer can be calculated using the following equation:

$$\text{Output Size} = \left\lfloor \frac{\text{Input Size} + 2 \times \text{Padding} - \text{Kernel Size}}{\text{Stride}} \right\rfloor + 1$$

Where: - Input Size is the size of the input feature map (height or width). - Padding is the number of zero-padding added to both sides of the input. - Kernel Size is the size of the convolution kernel (height or width). - Stride is the stride of the convolution.

PyTorch ConvTranspose2d Equation

The output size of a ConvTranspose2d (transposed convolution) layer can be calculated using the following equation:

$$\text{Output Size} = (\text{Input Size} - 1) \times \text{Stride} - 2 \times \text{Padding} + \text{Kernel Size} + \text{Output Padding}$$

Where: - Input Size is the size of the input feature map (height or width). - Stride is the stride of the convolution. - Padding is the number of zero-padding added to both sides of the input. - Kernel Size is the size of the convolution kernel (height or width). - Output Padding is the additional size added to the output (usually used to ensure the output size matches a specific value).

5 GANs Models

5.1 Basic GAN

- Architecture: Generator and Discriminator.
- Loss functions: Minimax game.
- Training process.

5.2 DCGAN

- Architecture: Convolutional layers.
- Improvements over basic GAN.
- Training tips.

5.3 WGAN

- Wasserstein distance.
- Critic network.
- Gradient penalty.

5.4 CycleGAN

- Architecture: Cycle consistency loss.
- Applications: Image-to-image translation.

6 Conclusion

Summarize the key points and discuss future directions.