

DEEP LEARNING FOR SPEECH & LANGUAGE

Winter Seminar UPC TelecomBCN, 24 - 31 January 2017



Instructors



Antonio Bonafonte J. Adrián Rodríguez Fonollosa Marta R. Costa-jussà Javier Hernando Santiago Pascual Elisa Sayrol Xavier Giró

Organizers



Image Processing Group
Signal Theory and Communications Department



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

+ info: [TelecomBCN.DeepLearning.Barcelona](https://www.telecombcn.com/deeplearning-barcelona)

[\[course site\]](#)

Day 2 Lecture 5

Generative Adversarial Networks



Santiago Pascual



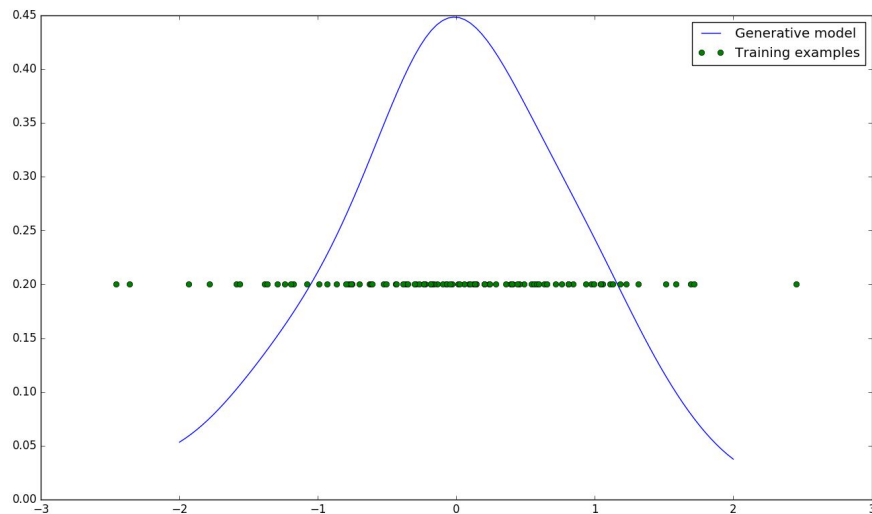
Outline

1. What are Generative models?
2. Why Generative Models?
3. Which Generative Models?
4. Generative Adversarial Networks
5. Applications to Images
6. Applications to Speech

What are Generative Models?

We want our model with parameters $\theta = \{\text{weights, biases}\}$ and outputs distributed like P_{model} to estimate the distribution of our training data P_{data} .

Example) $y = f(\mathbf{x})$, where y is scalar, make P_{model} similar to P_{data} by training the parameters θ to maximize their similarity.



What are Generative Models?

Key Idea: our model cares about what distribution generated the input data points, and we want to mimic it with our probabilistic model. **Our learned model should be able to make up new samples from the distribution, not just copy and paste existing samples!**

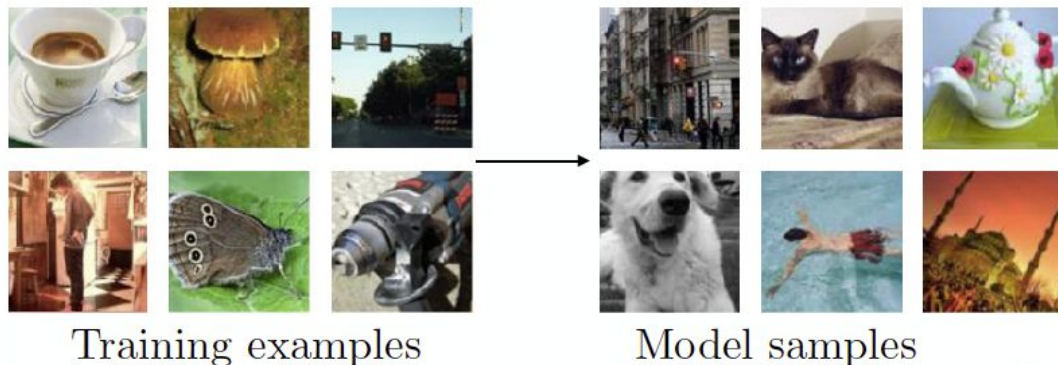


Figure from [NIPS 2016 Tutorial: Generative Adversarial Networks \(I. Goodfellow\)](#)

Why Generative Models?

- Model very complex and high-dimensional distributions.
- Be able to generate realistic synthetic samples
 - possibly perform data augmentation
 - simulate possible futures for learning algorithms
- Fill the blanks in the data
- Manipulate real samples with the assistance of the generative model
 - Edit pictures (exemplified later)

Current most known options:

1. Generative Adversarial Networks
2. Autoregressive methods (WaveNet) → stay tuned, next chapters
3. Variational Autoencoder

Generative Adversarial Networks (GAN)

We have a pair of networks, Generator (G) and Discriminator (D):

- They “fight” against each other during training → **Adversarial Training**
- G mission: make its pdf, P_{model} , as much similar as possible to our training set distribution P_{data} → Try to make predictions so realistic that D can't distinguish
- D mission: distinguish between G samples and real samples



Adversarial Training ([Goodfellow et al. 2014](#))

We have networks **G** and **D**, and **training set with pdf Pdata**. Notation:

- $\theta(\mathbf{G}), \theta(\mathbf{D})$ (Parameters of model **G** and **D** respectively)
- $\mathbf{x} \sim P_{\text{data}}$ (M-dim sample from training data pdf)
- $\mathbf{z} \sim N(0, I)$ (sample from prior pdf, e.g. N-dim normal)
- $\mathbf{G}(\mathbf{z}) = \tilde{\mathbf{x}} \sim P_g$ (M-dim sample from G network)

D network receives \mathbf{x} or $\tilde{\mathbf{x}}$ inputs \rightarrow decides whether input is real or fake. It is **optimized to learn: \mathbf{x} is real (1), $\tilde{\mathbf{x}}$ is fake (0) (binary classifier)**.

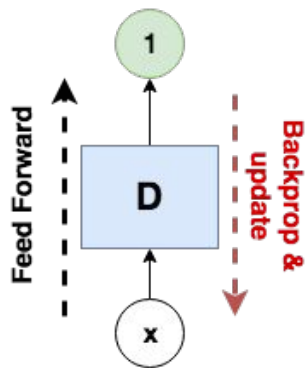
$$J^{(D)}(\theta^{(D)}, \theta^{(G)}) = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(\mathbf{G}(\mathbf{z}))) .$$

G network maps sample \mathbf{z} to $\mathbf{G}(\mathbf{z}) = \tilde{\mathbf{x}} \rightarrow$ it is **optimized to maximize D mistakes**.

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(\mathbf{G}(\mathbf{z}))$$

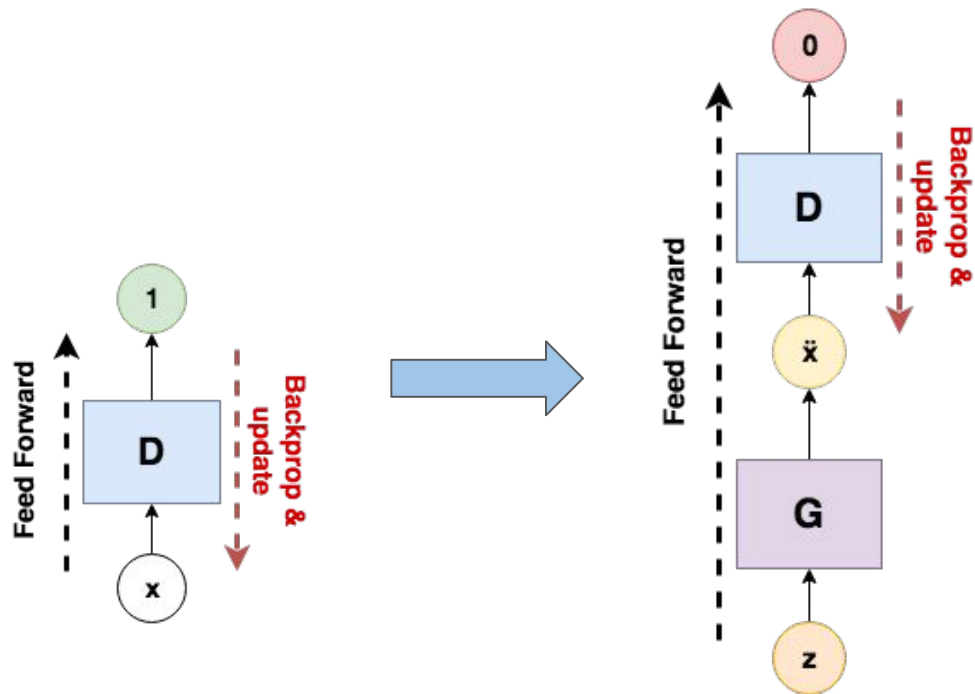
Adversarial Training (batch update)

- Pick a sample \mathbf{x} from training set
- Show \mathbf{x} to **D** and update weights to output 1 (real)



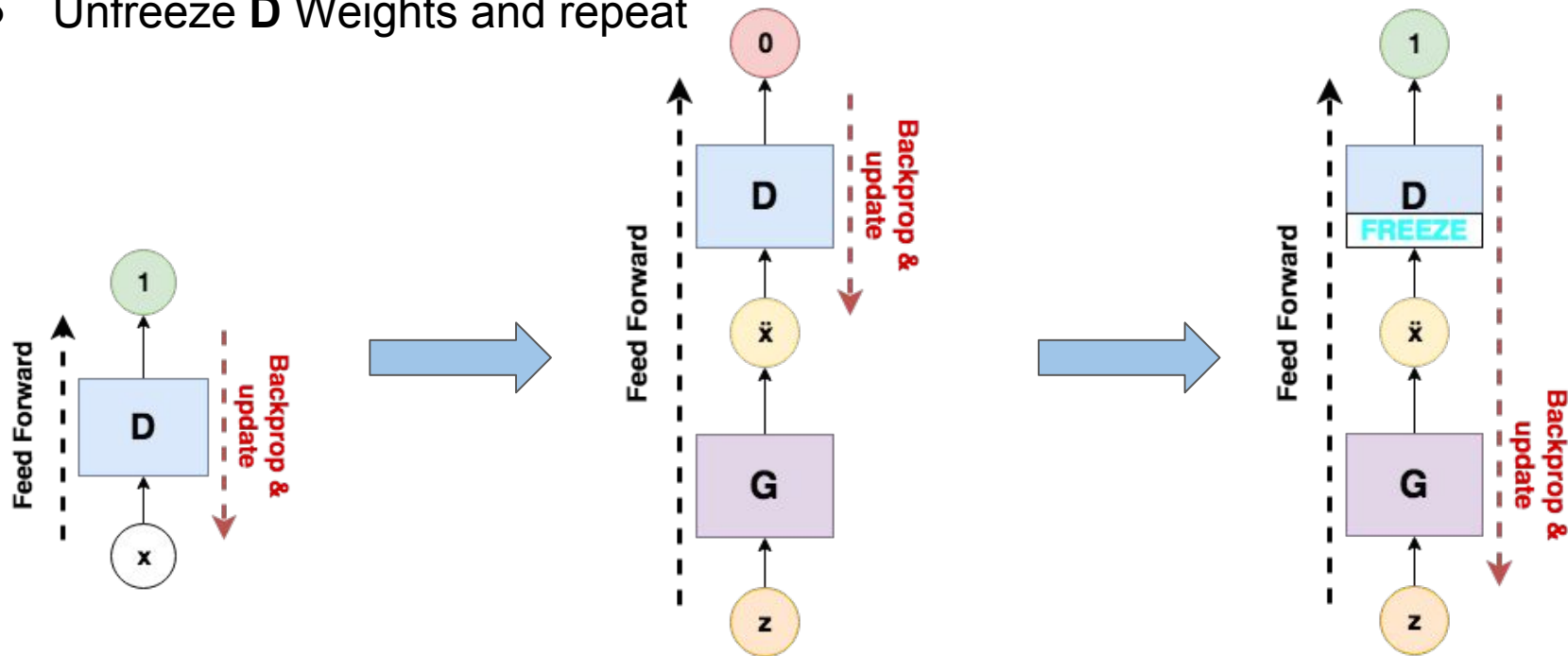
Adversarial Training (batch update)

- **G** maps sample \mathbf{z} to $\tilde{\mathbf{x}}$
- show $\tilde{\mathbf{x}}$ and update weights to output 0 (fake)



Adversarial Training (batch update)

- Freeze **D** weights
- Update **G** weights to make **D** output 1 (just **G** weights!)
- Unfreeze **D** Weights and repeat



Adversarial Training analogy

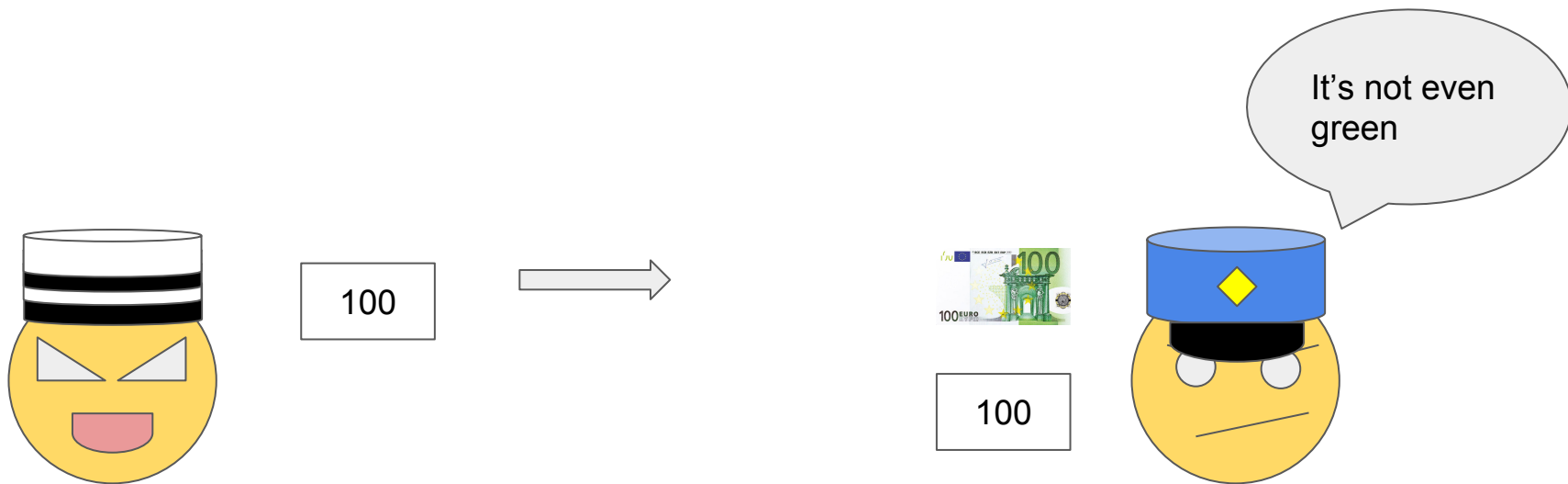
Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.

Key Idea: as **D** is trained to detect fraud, its parameters learn discriminative features of “what is real/fake”. As backprop goes through **D** to **G** there happens to be information leaking about the requirements for bank notes to look real. This makes **G** perform small corrections by little steps to get closer and closer to what a real sample would be.

- **Caveat:** this means GANs are not suitable for discrete tokens predictions, like words, because in that discrete space there is no “small change” criteria to get to a neighbour word (but can work in a word embedding space for example)

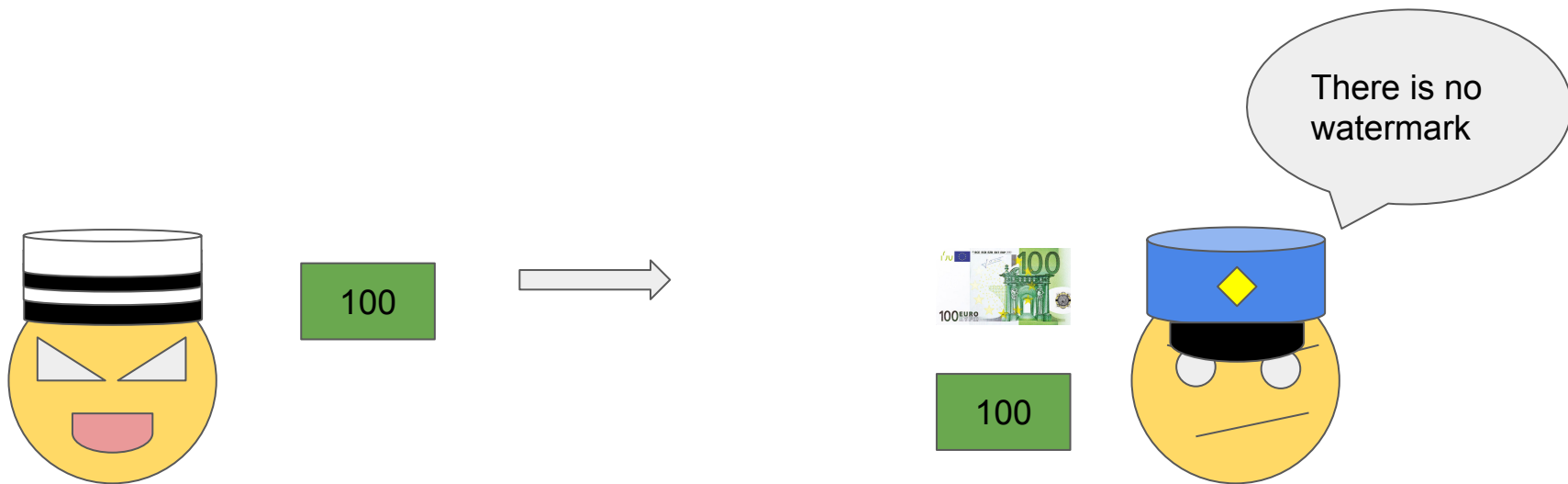
Adversarial Training analogy

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



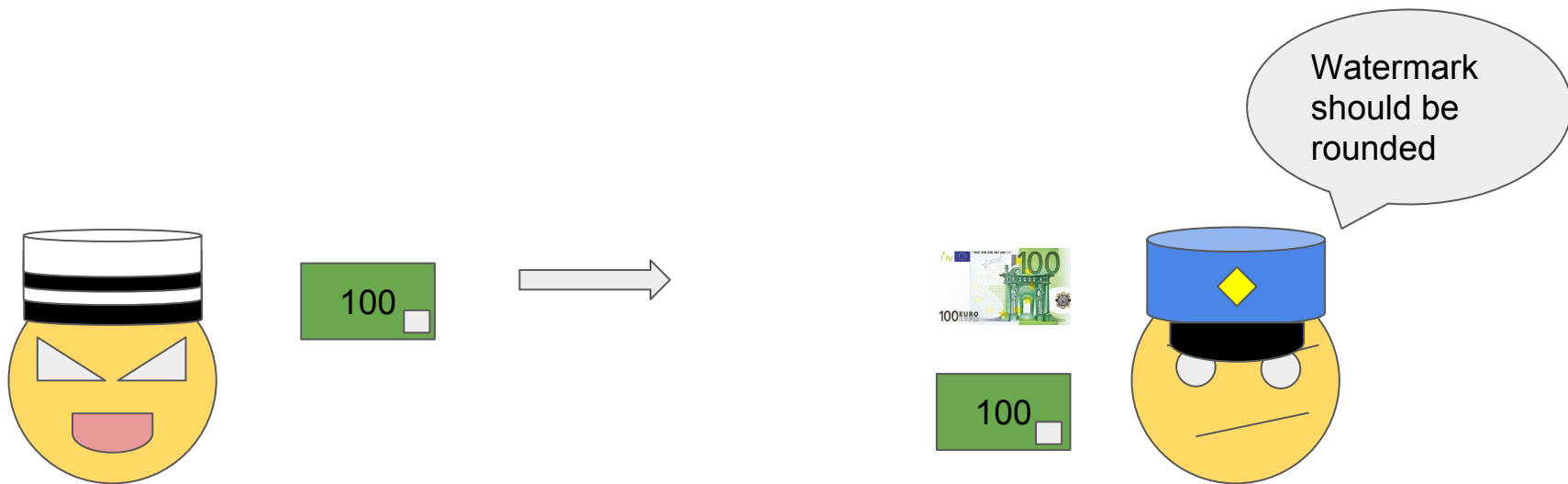
Adversarial Training analogy

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



Adversarial Training analogy

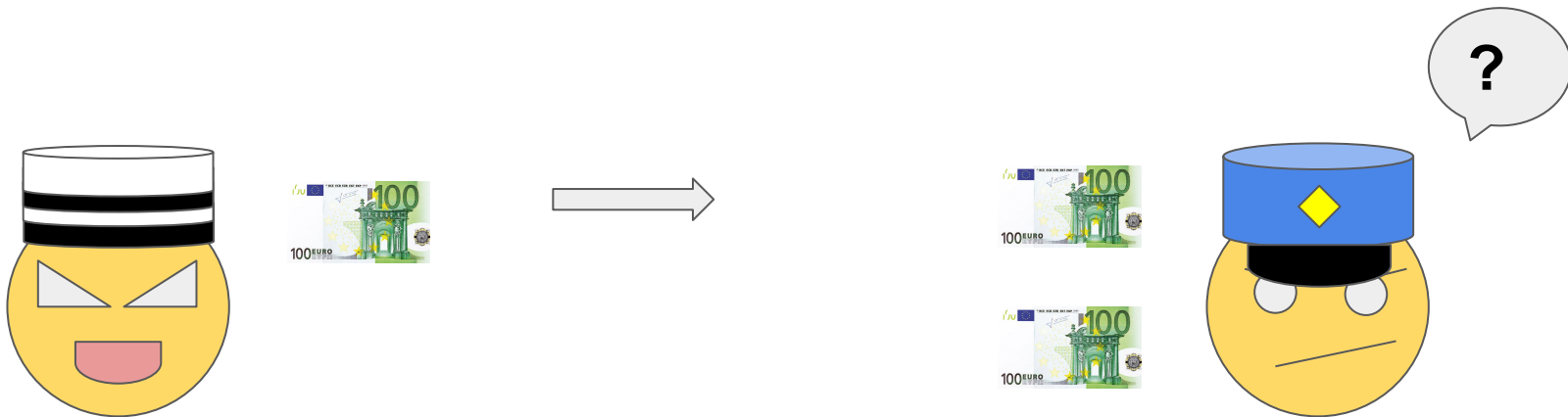
Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.



Adversarial Training analogy

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.

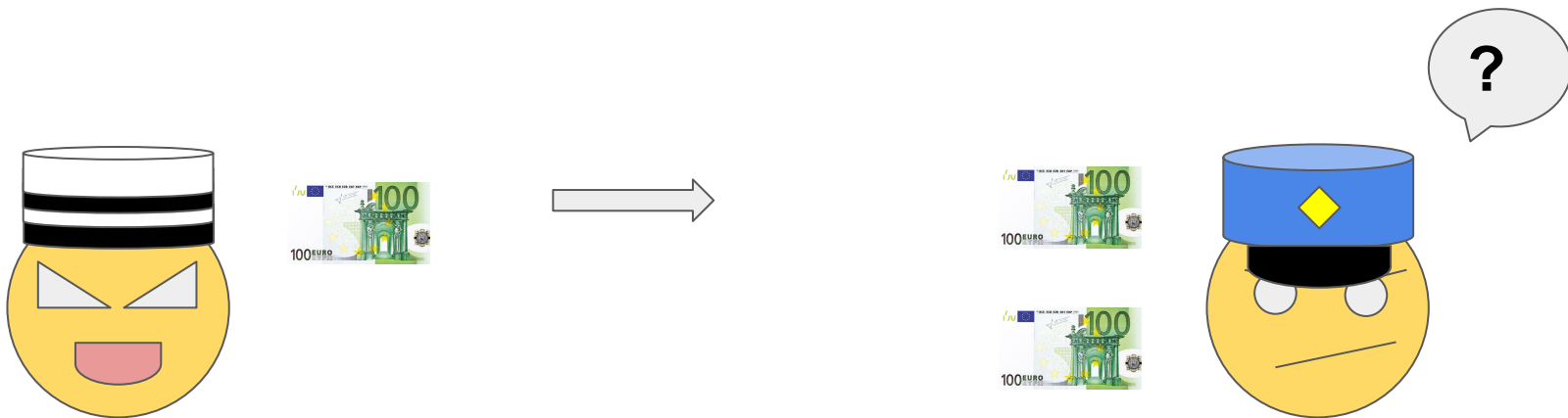
After enough iterations, and if the counterfeiter is good enough (in terms of **G** network it means “has enough parameters”), the police should be confused.



Adversarial Training analogy

Imagine we have a counterfeiter (**G**) trying to make fake money, and the police (**D**) has to detect whether money is real or fake.

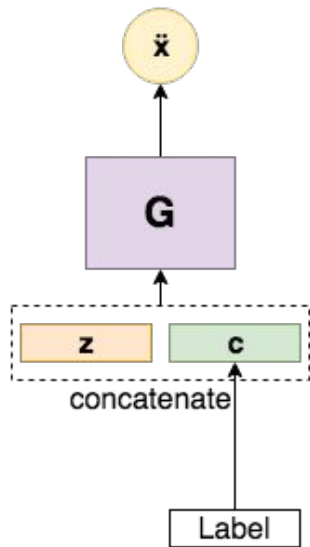
After enough iterations, and if the counterfeiter is good enough (in terms of **G** network it means “has enough parameters”), the police should be confused.



Conditioned GANs

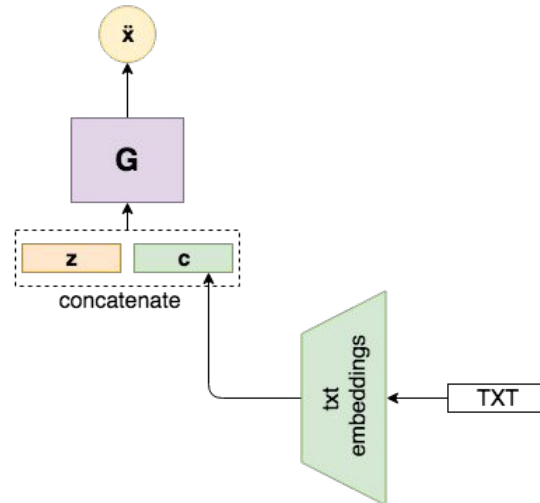
GANs can be conditioned on other info extra to \mathbf{z} : text, labels, etc..

\mathbf{z} might capture random characteristics of the data, variabilities of possible futures, whilst \mathbf{c} would condition the deterministic parts



For details on ways to condition GANs:

[Ways of Conditioning Generative Adversarial Networks \(Wack et al.\)](#)



GAN Applications

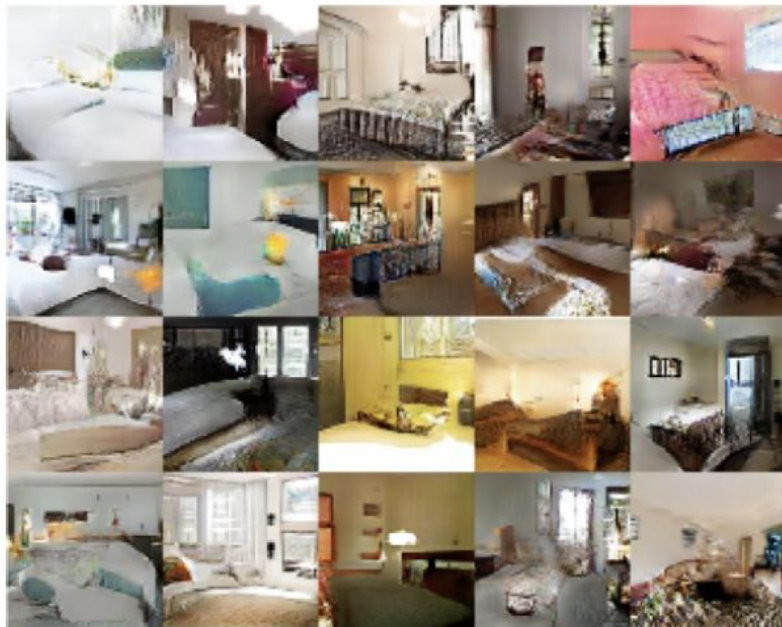
So far GANs have been extensively used in computer vision tasks:

- Generating images/generating video frames
- Unsupervised feature extraction/learning representations
- Manipulating images (in a photoshop advanced level)
- Image coding/Super Resolution
- Transferring image styles

However we have been working on advances for speech generation!

Generating images/frames

Deep Conv. GAN (DCGAN) effectively generated 64x64 RGB images in a single shot. For example bedrooms from LSUN dataset.



([Radford et al. 2015](#))

Generating images/frames conditioned on captions

this small bird has a pink breast and crown, and black primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma



this magnificent fellow is almost all black with a red crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen



This small blue bird has a short pointy beak and brown on its wings



This bird is completely red with black wings and pointy beak



A small sized bird that has a cream belly and a short pointed bill



A small bird with a black head and wings and features grey wings



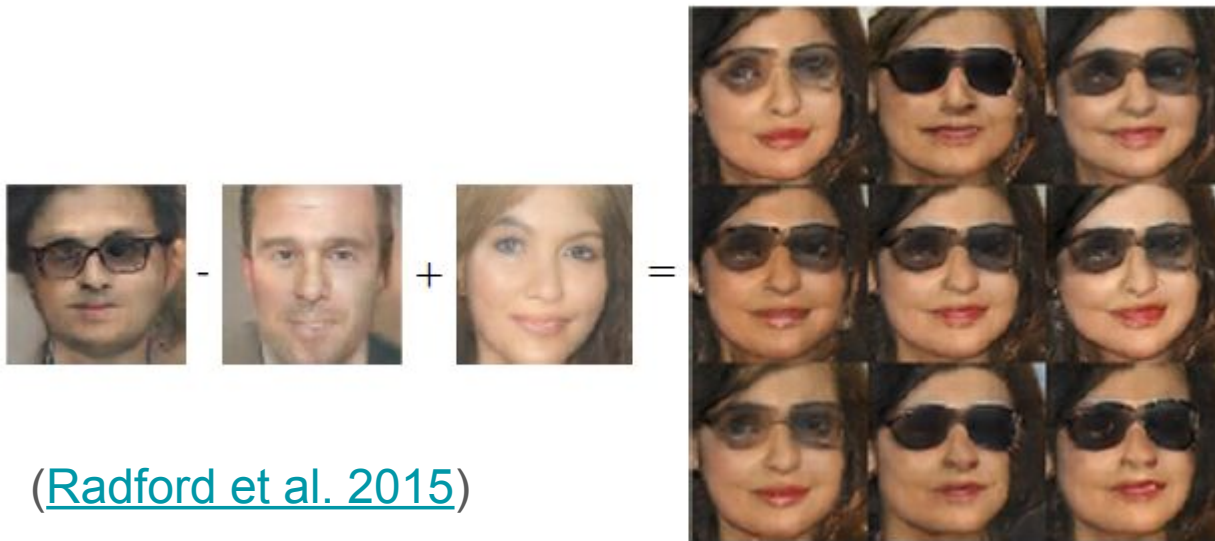
([Reed et al. 2016b](#))

([Zhang et al. 2016](#))

Unsupervised feature extraction/learning representations

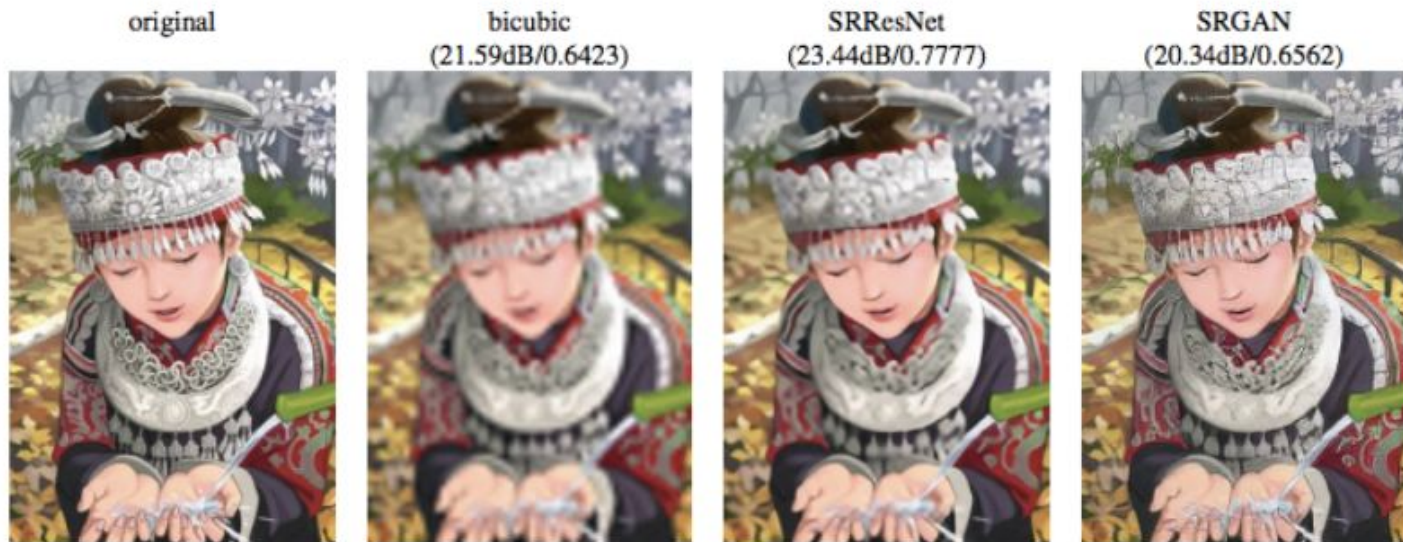
Similarly to word2vec, GANs learn a distributed representation that disentangles concepts such that we can perform operations on the data manifold:

$$v(\text{Man with glasses}) - v(\text{man}) + v(\text{woman}) = v(\text{woman with glasses})$$



([Radford et al. 2015](#))

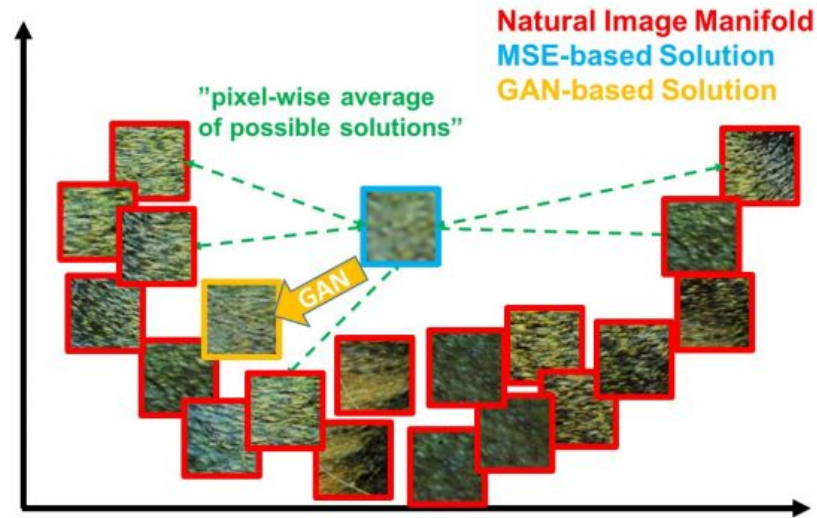
Image super-resolution



([Ledig et al. 2016](#))

Bicubic: not using data statistics. SRResNet: trained with MSE. SRGAN is able to understand that there are multiple correct answers, rather than averaging.

Image super-resolution



([Ledig et al. 2016](#))

Averaging is a serious problem we face when dealing with complex distributions.

Manipulating images and assisted content creation



<https://youtu.be/9c4z6YsBGQ0?t=126>

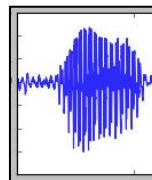
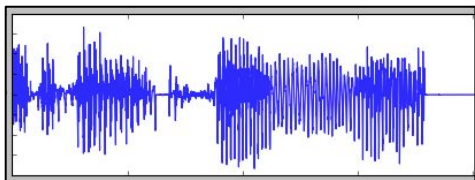
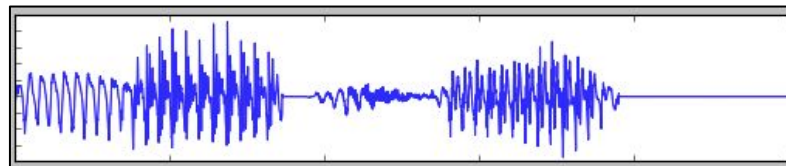
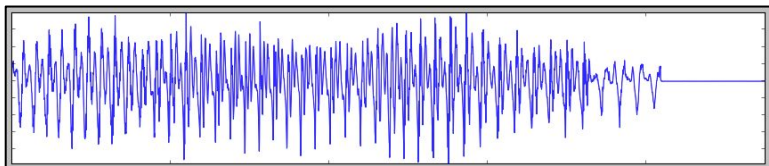


<https://youtu.be/9c4z6YsBGQ0?t=161>

([Zhu et al. 2016](#))

Waveforms generation

We have done recent advances in generating speech signals with GAN models. There are none existent systems doing so until now. Current line of research.



(S.Pascual, A.Bonafonte, J.Serrà)

Caveats

Where is the downside...?

Well GANs are tricky and hard to train! We do not want to minimize a cost function. Instead we want both networks to reach a Nash equilibria (saddle point).

Because of extensive experience within the GAN community (with some does-not-work-frustration from time to time), you can find some tricks and tips on how to train a GAN here: <https://github.com/soumith/ganhacks>

An open problem of GANs is the evaluation of their generation quality: there are no established objective metrics to do so → We look (or hear) the generated samples to know when to stop training or how are we doing.

Thanks ! Q&A ?

Strongly recommended reading: [NIPS 2016 Tutorial: Generative Adversarial Networks \(I. Goodfellow\)](#)



[@Santty128](#)