

# DENOISING AUTOENCODERS (DAE)

Jérémie Cabessa

Laboratoire DAVID, UVSQ

# INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un **espace latent (latent space)** et être capable de les reconstruire ces data.
- ▶ Applications en compression de data (data compression).
- ▶ Autre application des autoencdeurs: le **débruitage des data (denoising)**.
- ▶ Pour cela, on utilise des **denoising autoencoders**.

# INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un **espace latent (latent space)** et être capable de les reconstruire ces data.
- ▶ Applications en compression de data (data compression).
- ▶ Autre application des autoencdeurs: le **débruitage des data (denoising)**.
- ▶ Pour cela, on utilise des **denoising autoencoders**.

# INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un **espace latent (latent space)** et être capable de les reconstruire ces data.
- ▶ Applications en compression de data (data compression).
- ▶ Autre application des autoencdeurs: le **débruitage des data (denoising)**.
- ▶ Pour cela, on utilise des **denoising autoencoders**.

# INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un **espace latent (latent space)** et être capable de les reconstruire ces data.
- ▶ Applications en compression de data (data compression).
- ▶ Autre application des autoencdeurs: le **débruitage des data (denoising)**.
- ▶ Pour cela, on utilise des **denoising autoencoders**.

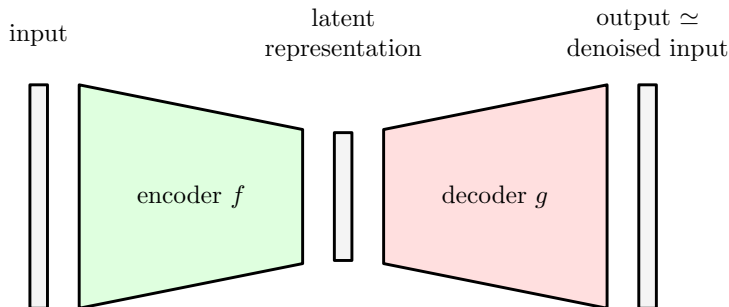
# INTRODUCTION

- ▶ Principe similaire à celui d'un autoencodeur.
- ▶ Au lieu de viser la reconstruction des data, on vise leur débruitage.

# INTRODUCTION

- ▶ Principe similaire à celui d'un autoencodeur.
- ▶ Au lieu de viser la reconstruction des data, on vise leur débruitage.

## ARCHITECTURE ENCODEUR-DÉCODEUR





# ARCHITECTURE ENCODEUR-DÉCODEUR

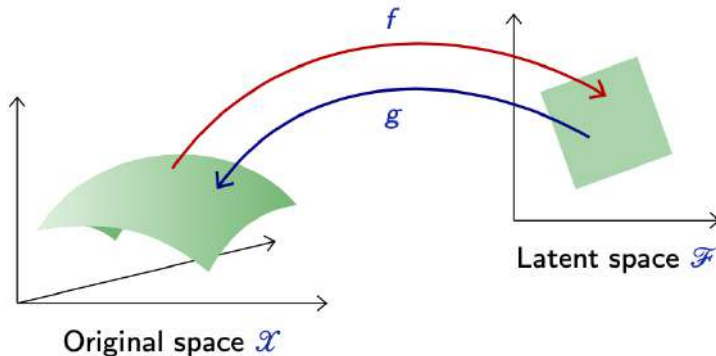


Figure taken from [Fleuret, 2022].

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Pour cela, on dispose de *data propres* comme *targets*.
- ▶ Si  $\mathbf{x}$  et  $\tilde{\mathbf{x}}$  représentent une data propre (clean) et bruitée (noisy), respectivement, où  $\tilde{\mathbf{x}}$  est une perturbation de  $\mathbf{x}$ , alors on veut:

$$\mathcal{N}(\tilde{\mathbf{x}}; \Theta_f, \Theta_g) = g(f(\tilde{\mathbf{x}}; \Theta_f); \Theta_g) \simeq \mathbf{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Pour cela, on dispose de *data propres* comme *targets*.
- ▶ Si  $\mathbf{x}$  et  $\tilde{\mathbf{x}}$  représentent une data propre (clean) et bruitée (noisy), respectivement, où  $\tilde{\mathbf{x}}$  est une perturbation de  $\mathbf{x}$ , alors on veut:

$$\mathcal{N}(\tilde{\mathbf{x}}; \Theta_f, \Theta_g) = g(f(\tilde{\mathbf{x}}; \Theta_f); \Theta_g) \simeq \mathbf{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Pour cela, on dispose de *data propres* comme *targets*.
- ▶ Si  $\mathbf{x}$  et  $\tilde{\mathbf{x}}$  représentent une data propre (clean) et bruitée (noisy), respectivement, où  $\tilde{\mathbf{x}}$  est une perturbation de  $\mathbf{x}$ , alors on veut:

$$\mathcal{N}(\tilde{\mathbf{x}}; \Theta_f, \Theta_g) = g(f(\tilde{\mathbf{x}}; \Theta_f); \Theta_g) \simeq \mathbf{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Pour cela, on dispose de *data propres* comme *targets*.
- ▶ Si  $\mathbf{x}$  et  $\tilde{\mathbf{x}}$  représentent une data propre (clean) et bruitée (noisy), respectivement, où  $\tilde{\mathbf{x}}$  est une perturbation de  $\mathbf{x}$ , alors on veut:

$$\mathcal{N}(\tilde{\mathbf{x}}; \Theta_f, \Theta_g) = g(f(\tilde{\mathbf{x}}; \Theta_f); \Theta_g) \simeq \mathbf{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Pour cela, on dispose de *data propres* comme *targets*.
- ▶ Si  $\mathbf{x}$  et  $\tilde{\mathbf{x}}$  représentent une data propre (clean) et bruitée (noisy), respectivement, où  $\tilde{\mathbf{x}}$  est une perturbation de  $\mathbf{x}$ , alors on veut:

$$\mathcal{N}(\tilde{\mathbf{x}}; \Theta_f, \Theta_g) = g(f(\tilde{\mathbf{x}}; \Theta_f); \Theta_g) \simeq \mathbf{x}$$

# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\mathbf{X}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des data propre et bruitées, respectivement.

- Un exemple de bruitage serait:

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$

# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\mathbf{X}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des data propre et bruitées, respectivement.

- Un exemple de bruitage serait:

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$



# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\mathbf{X}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des data propre et bruitées, respectivement.

- Un exemple de bruitage serait:

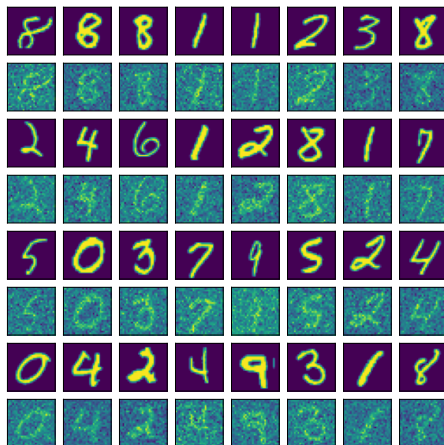
$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

- Data bruitées passées au DAE.



# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
class DenoisingAutoencoder_CNN(nn.Module):
    """Implements a CNN denoising autoencoder"""

    def __init__(self):
        """constructor"""
        super().__init__()

        # N, 1, 28, 28
        self.encoder = nn.Sequential(
            # -> N, 16, 14, 14
            nn.Conv2d(1, 16, 3, stride=2, padding=1),
            nn.ReLU(),
            # -> N, 32, 7, 7
            nn.Conv2d(16, 32, 3, stride=2, padding=1),
            nn.ReLU(),
            # -> N, 64, 1, 1
            nn.Conv2d(32, 64, 7)
        )

        # N, 64, 1, 1
        self.decoder = nn.Sequential(
            # -> N, 32, 7, 7
            nn.ConvTranspose2d(64, 32, 7),
            nn.ReLU(),
            # N, 16, 14, 14 (N, 16, 13, 13 without output_padding)
            nn.ConvTranspose2d(32, 16, 3,
                               stride=2, padding=1, output_padding=1),
            nn.ReLU(),
            # N, 1, 28, 28 (N, 1, 27, 27)
            nn.ConvTranspose2d(16, 1, 3,
                               stride=2, padding=1, output_padding=1),
            nn.Sigmoid()
        )

    def forward(self, x):
        """forward function"""

        encoded_data = self.encoder(x)
        decoded_data = self.decoder(encoded_data)

        return decoded_data
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
self.encoder = nn.Sequential(  
    nn.Linear(28 * 28, 128),  
    nn.ReLU(),  
    nn.Linear(128, 64),  
    nn.ReLU(),  
    nn.Linear(64, 12),  
    nn.ReLU(),  
    nn.Linear(12, 2) # -> N, 2 only!  
)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
self.decoder = nn.Sequential(  
    nn.Linear(2, 12),  
    nn.ReLU(),  
    nn.Linear(12, 64),  
    nn.ReLU(),  
    nn.Linear(64, 128),  
    nn.ReLU(),  
    nn.Linear(128, 28 * 28),  
    nn.Sigmoid()  
)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
model = DenoisingAutoencoder_CNN()  
criterion = nn.MSELoss()  
optimizer = torch.optim.Adam(model.parameters(),  
                               lr=1e-3,  
                               weight_decay=1e-5)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

- **Training:** on calcule la MSE loss entre les images des data bruitées (images of noisy data) et les data propres (clean data, targets!).

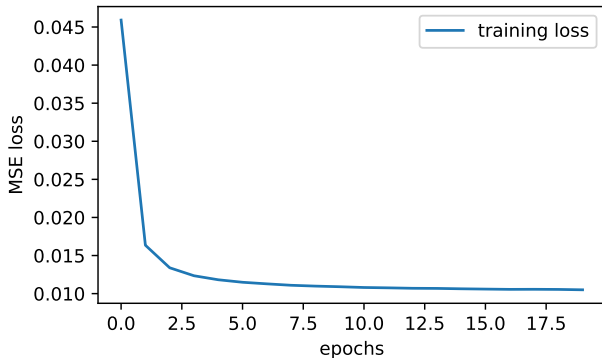
```
for (inputs, _) in train_loader:           # clean inputs           (1)

    inputs_noisy = add_noise(inputs)        # noisy inputs
    outputs = model(inputs_noisy)           # images of noisy inputs (2)

    # loss between images of noisy inputs and clean inputs
    loss = criterion(outputs, inputs) # loss between (2) and (1)
    train_loss.append(loss.item())

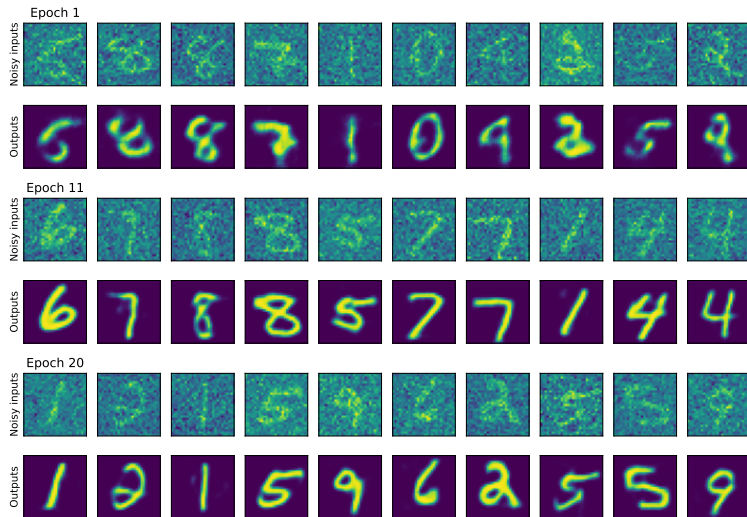
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL





## IMPLÉMENTATION: DAE CONVOLUTIONNEL



# REMARQUES

- ▶ Au fur et à mesure de l'entraînement, le DAE débruite les data de mieux en mieux.
- ▶ Le DAE est capable d'apprendre à débruiter les data.

# REMARQUES

- ▶ Au fur et à mesure de l'entraînement, le DAE débruite les data de mieux en mieux.
- ▶ Le DAE est capable d'apprendre à débruiter les data.

# PARADIGME NOISE2NOISE

- ▶ Jusqu'ici, on a appris à débruiter des data à partir de data propres (loss between images of noisy data and clean data).
- ▶ Mais il se peut que l'on ne dispose pas de data propres (clean data) comme targets pour entraîner notre DAE.
- ▶ De manière surprenante, on peut apprendre à débruiter des data à partir de data elles-mêmes bruitées!
- ▶ C'est le paradigme "noise to noise".

# PARADIGME NOISE2NOISE

- ▶ Jusqu'ici, on a appris à débruiter des data à partir de data propres (loss between images of noisy data and clean data).
- ▶ Mais il se peut que l'on ne dispose pas de data propres (clean data) comme targets pour entraîner notre DAE.
- ▶ De manière surprenante, on peut apprendre à débruiter des data à partir de data elles-mêmes bruitées!
- ▶ C'est le paradigme "noise to noise".

# PARADIGME NOISE2NOISE

- ▶ Jusqu'ici, on a appris à débruiter des data à partir de data propres (loss between images of noisy data and clean data).
- ▶ Mais il se peut que l'on ne dispose pas de data propres (clean data) comme targets pour entraîner notre DAE.
- ▶ De manière surprenante, on peut apprendre à débruiter des data à partir de data elles-mêmes bruitées!
- ▶ C'est le paradigme "noise to noise".

# PARADIGME NOISE2NOISE

- ▶ Jusqu'ici, on a appris à débruiter des data à partir de data propres (loss between images of noisy data and clean data).
- ▶ Mais il se peut que l'on ne dispose pas de data propres (clean data) comme targets pour entraîner notre DAE.
- ▶ De manière surprenante, on peut apprendre à débruiter des data à partir de data elles-mêmes bruitées!
- ▶ C'est le paradigme “**noise to noise**”.

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶ Une fois encore,  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Mais dans ce cas, on ne dispose pas de data propres, mais d'*autres data bruitées* comme *targets*.
- ▶ Si  $\bar{x}$  et  $\tilde{x}$  représentent deux versions bruitées indépendantes d'une même data propre inconnue  $x$ , alors on a:

$$\mathcal{N}(\tilde{x}; \Theta_f, \Theta_g) = g(f(\tilde{x}; \Theta_f); \Theta_g) \simeq \bar{x}$$



# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶ Une fois encore,  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Mais dans ce cas, on ne dispose pas de data propres, mais d'*autres data bruitées* comme *targets*.
- ▶ Si  $\bar{x}$  et  $\tilde{x}$  représentent deux versions bruitées indépendantes d'une même data propre inconnue  $x$ , alors on a:

$$\mathcal{N}(\tilde{x}; \Theta_f, \Theta_g) = g(f(\tilde{x}; \Theta_f); \Theta_g) \simeq \bar{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶ Une fois encore,  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Mais dans ce cas, on ne dispose pas de data propres, mais d'*autres data bruitées* comme *targets*.
- ▶ Si  $\bar{x}$  et  $\tilde{x}$  représentent deux versions bruitées indépendantes d'une même data propre inconnue  $x$ , alors on a:

$$\mathcal{N}(\tilde{x}; \Theta_f, \Theta_g) = g(f(\tilde{x}; \Theta_f); \Theta_g) \simeq \bar{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶ Une fois encore,  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Mais dans ce cas, on ne dispose pas de data propres, mais d'*autres data bruitées* comme *targets*.
- ▶ Si  $\bar{x}$  et  $\tilde{x}$  représentent deux versions bruitées indépendantes d'une même data propre inconnue  $x$ , alors on a:

$$\mathcal{N}(\tilde{x}; \Theta_f, \Theta_g) = g(f(\tilde{x}; \Theta_f); \Theta_g) \simeq \bar{x}$$

# DENOISING AUTOENCODER (DAE)

- ▶ L'**encodeur** est un deep neural network  $f(\cdot; \Theta_f)$ .
- ▶ Le **décodeur** est un deep neural network  $g(\cdot; \Theta_g)$ .
- ▶ La composition de  $f(\cdot; \Theta_f)$  et  $g(\cdot; \Theta_g)$  forme un réseau de neurones  $\mathcal{N}(\cdot; \Theta_f, \Theta_g)$ .
- ▶ Une fois encore,  $\mathcal{N}$  reçoit des inputs bruitées et son but est de les débruiter. Mais dans ce cas, on ne dispose pas de data propres, mais d'*autres data bruitées* comme *targets*.
- ▶ Si  $\bar{x}$  et  $\tilde{x}$  représentent deux versions bruitées indépendantes d'une même data propre inconnue  $x$ , alors on a:

$$\mathcal{N}(\tilde{x}; \Theta_f, \Theta_g) = g(f(\tilde{x}; \Theta_f); \Theta_g) \simeq \bar{x}$$

# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\bar{\mathbf{x}}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des deux types de data bruitées, respectivement.

- Un exemple de bruitage serait:

$$\bar{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon'_i \quad \text{où} \quad \epsilon'_i \sim \mathcal{N}(\mu', \Sigma')$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$

# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\bar{\mathbf{x}}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des deux types de data bruitées, respectivement.

- Un exemple de bruitage serait:

$$\bar{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon'_i \quad \text{où} \quad \epsilon'_i \sim \mathcal{N}(\mu', \Sigma')$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$

# DENOISING AUTOENCODER (DAE)

- Pour cela, on minimise la mean squared error (MSE):

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g) = \frac{1}{N} \sum_{i=1}^N \|\bar{\mathbf{x}}_i - g(f(\tilde{\mathbf{x}}_i; \Theta_f); \Theta_g)\|^2$$

où  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  représente la concaténation des deux types de data bruitées, respectivement.

- Un exemple de bruitage serait:

$$\bar{\mathbf{x}}_i = \mathbf{x}_i + \epsilon_i \quad \text{où} \quad \epsilon_i \sim \mathcal{N}(\mu, \Sigma)$$

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon'_i \quad \text{où} \quad \epsilon'_i \sim \mathcal{N}(\mu', \Sigma')$$

- L'entraînement du réseau consiste alors à trouver les poids  $\hat{\Theta}_f$  et  $\hat{\Theta}_g$  qui satisfont

$$\hat{\Theta}_f; \hat{\Theta}_g = \arg \min_{\Theta_f; \Theta_g} \mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta_f, \Theta_g)$$

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
class DenoisingAutoencoder_CNN(nn.Module):
    """Implements a CNN denoising autoencoder"""

    def __init__(self):
        """constructor"""

        super().__init__()

        # N, 1, 28, 28
        self.encoder = nn.Sequential(
            # -> N, 16, 14, 14
            nn.Conv2d(1, 16, 3, stride=2, padding=1),
            nn.ReLU(),
            # -> N, 32, 7, 7
            nn.Conv2d(16, 32, 3, stride=2, padding=1),
            nn.ReLU(),
            # -> N, 64, 1, 1
            nn.Conv2d(32, 64, 7)
        )

        # N, 64, 1, 1
        self.decoder = nn.Sequential(
            # -> N, 32, 7, 7
            nn.ConvTranspose2d(64, 32, 7),
            nn.ReLU(),
            # N, 16, 14, 14 (N, 16, 13, 13 without output_padding)
            nn.ConvTranspose2d(32, 16, 3,
                               stride=2, padding=1, output_padding=1),
            nn.ReLU(),
            # N, 1, 28, 28 (N, 1, 27, 27)
            nn.ConvTranspose2d(16, 1, 3,
                               stride=2, padding=1, output_padding=1),
            nn.Sigmoid()
        )

    def forward(self, x):
        """forward function"""

        encoded_data = self.encoder(x)
        decoded_data = self.decoder(encoded_data)

        return decoded_data
```



# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
self.encoder = nn.Sequential(  
    nn.Linear(28 * 28, 128),  
    nn.ReLU(),  
    nn.Linear(128, 64),  
    nn.ReLU(),  
    nn.Linear(64, 12),  
    nn.ReLU(),  
    nn.Linear(12, 2) # -> N, 2 only!  
)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
self.decoder = nn.Sequential(  
    nn.Linear(2, 12),  
    nn.ReLU(),  
    nn.Linear(12, 64),  
    nn.ReLU(),  
    nn.Linear(64, 128),  
    nn.ReLU(),  
    nn.Linear(128, 28 * 28),  
    nn.Sigmoid()  
)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

```
model = DenoisingAutoencoder_CNN()  
criterion = nn.MSELoss()  
optimizer = torch.optim.Adam(model.parameters(),  
                               lr=1e-3,  
                               weight_decay=1e-5)
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL

- **Training:** on calcule la MSE loss les images des data bruitées (images of noisy data) et les data bruitées initiales (noisy data, targets!).

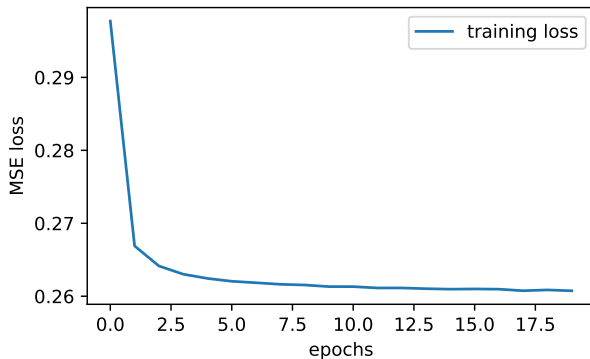
```
for (inputs, _) in train_loader:

    inputs_noisy1 = add_noise(inputs) # noisy inputs! (1)
    inputs_noisy2 = add_noise(inputs) # other noisy inputs
    outputs = model(inputs_noisy2)    # images of noisy inputs (2)

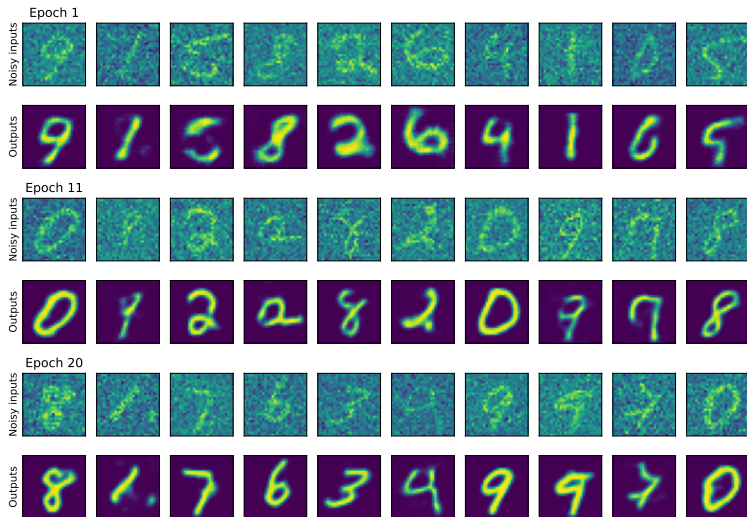
    # loss between noisy inputs and images of other noisy inputs
    loss = criterion(outputs, inputs_noisy1) # loss between (2) and (1)
    train_loss.append(loss.item())

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

# IMPLÉMENTATION: DAE CONVOLUTIONNEL



## IMPLÉMENTATION: DAE CONVOLUTIONNEL



# REMARQUES

- ▶ Bien que le DAE soit entraîné sur des data bruitées uniquement, au fur et à mesure de son entraînement, il apprend à débruiter ces data de mieux en mieux.
- ▶ Le DAE est capable de débruiter les data, même sans avoir été entraîné sur des data propres!
- ▶ Ceci n'a rien de magique, et nous allons le montrer dans les slides suivants.
- ▶ L'hypothèse d'indépendance des deux bruits est cruciale.

# REMARQUES

- ▶ Bien que le DAE soit entraîné sur des data bruitées uniquement, au fur et à mesure de son entraînement, il apprend à débruiter ces data de mieux en mieux.
- ▶ Le DAE est capable de débruiter les data, même sans avoir été entraîné sur des data propres!
- ▶ Ceci n'a rien de magique, et nous allons le montrer dans les slides suivants.
- ▶ L'hypothèse d'indépendance des deux bruits est cruciale.



# REMARQUES

- ▶ Bien que le DAE soit entraîné sur des data bruitées uniquement, au fur et à mesure de son entraînement, il apprend à débruiter ces data de mieux en mieux.
- ▶ Le DAE est capable de débruiter les data, même sans avoir été entraîné sur des data propres!
- ▶ Ceci n'a rien de magique, et nous allons le montrer dans les slides suivants.
- ▶ L'hypothèse d'indépendance des deux bruits est cruciale.

# REMARQUES

- ▶ Bien que le DAE soit entraîné sur des data bruitées uniquement, au fur et à mesure de son entraînement, il apprend à débruiter ces data de mieux en mieux.
- ▶ Le DAE est capable de débruiter les data, même sans avoir été entraîné sur des data propres!
- ▶ Ceci n'a rien de magique, et nous allons le montrer dans les slides suivants.
- ▶ L'hypothèse d'indépendance des deux bruits est cruciale.

## NOISE 2 NOISE: PREUVE

- Soient  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  des data bruitées provenant de data  $\mathbf{X}$  par adjonction de bruits *indépendants, additifs* et *non-biaisés* (i.e. *moyennes nulles*)  $\epsilon$  et  $\delta$ , e.g.:

$$\begin{aligned}\bar{\mathbf{X}} &= \mathbf{X} + \epsilon \quad \text{où} \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma) \\ \tilde{\mathbf{X}} &= \mathbf{X} + \delta \quad \text{où} \quad \delta \sim \mathcal{N}(\mathbf{0}, \Sigma')\end{aligned}$$

- Alors la minimisation de la mean squared error (MSE)

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta) = \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \bar{\mathbf{x}}_i\|^2$$

induit un réseau de neurones  $\mathcal{N}$  capable de débruiter les data.

- Comment est-ce possible?

## NOISE 2 NOISE: PREUVE

- Soient  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  des data bruitées provenant de data  $\mathbf{X}$  par adjonction de bruits *indépendants, additifs et non-biaisés* (i.e. *moyennes nulles*)  $\epsilon$  et  $\delta$ , e.g.:

$$\begin{aligned}\bar{\mathbf{X}} &= \mathbf{X} + \epsilon \quad \text{où} \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma) \\ \tilde{\mathbf{X}} &= \mathbf{X} + \delta \quad \text{où} \quad \delta \sim \mathcal{N}(\mathbf{0}, \Sigma')\end{aligned}$$

- Alors la minimisation de la mean squared error (MSE)

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta) = \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \bar{\mathbf{x}}_i\|^2$$

induit un réseau de neurones  $\mathcal{N}$  capable de débruiter les data.

- Comment est-ce possible?

## NOISE 2 NOISE: PREUVE

- Soient  $\bar{\mathbf{X}}$  et  $\tilde{\mathbf{X}}$  des data bruitées provenant de data  $\mathbf{X}$  par adjonction de bruits *indépendants, additifs* et *non-biaisés* (i.e. *moyennes nulles*)  $\epsilon$  et  $\delta$ , e.g.:

$$\begin{aligned}\bar{\mathbf{X}} &= \mathbf{X} + \epsilon \quad \text{où} \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma) \\ \tilde{\mathbf{X}} &= \mathbf{X} + \delta \quad \text{où} \quad \delta \sim \mathcal{N}(\mathbf{0}, \Sigma')\end{aligned}$$

- Alors la minimisation de la mean squared error (MSE)

$$\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta) = \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \bar{\mathbf{x}}_i\|^2$$

induit un réseau de neurones  $\mathcal{N}$  capable de débruiter les data.

- Comment est-ce possible?

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E} [\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned}\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E} [\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] &\end{aligned}$$

et donc

$$\begin{aligned}\hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2]\end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned}\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E} [\epsilon]^T \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}]}_{=0} + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] &\end{aligned}$$

et donc

$$\begin{aligned}\hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2]\end{aligned}$$



## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E}[\epsilon]^T \mathbb{E}[\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}]}_{=0} + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E} [\epsilon]^T \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}]}_{=0} + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] &= \\ \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - \underbrace{2\mathbb{E} [\epsilon]^T \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}]}_{=0} + \mathbb{E} [\|\epsilon\|^2] &= \\ \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] &= \\ \mathbb{E} \left[ \|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2\mathbb{E} \left[ \epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2 \underbrace{\mathbb{E} [\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2 \right] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] &= \\ \mathbb{E} \left[ \|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2\mathbb{E} \left[ \epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2 \underbrace{\mathbb{E}[\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2 \right] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] &= \\ \mathbb{E} \left[ \|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2\mathbb{E} \left[ \epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2 \underbrace{\mathbb{E}[\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2 \right] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2]$ .

De plus, on a

$$\begin{aligned} & \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] = \\ & \mathbb{E} [\|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2] = \\ & \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2] = \\ & \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\mathbb{E} [\epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})] + \mathbb{E} [\|\epsilon\|^2] = \\ & \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] - 2\underbrace{\mathbb{E} [\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} [\|\epsilon\|^2] = \\ & \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} [\|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2] + \mathbb{E} [\|\epsilon\|^2] \\ &= \arg \min_{\Theta} \mathbb{E} [\|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2] \end{aligned}$$

## NOISE 2 NOISE: PREUVE

- En fait,  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  est un estimateur de  $\mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right]$ .

De plus, on a

$$\begin{aligned} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] &= \\ \mathbb{E} \left[ \|\mathcal{N}(\mathbf{X} + \delta; \Theta) - (\mathbf{X} + \epsilon)\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) - \epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2\mathbb{E} \left[ \epsilon^T (\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}) \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] - 2 \underbrace{\mathbb{E}[\epsilon]^T}_{=0} \mathbb{E} [\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X}] + \mathbb{E} \left[ \|\epsilon\|^2 \right] &= \\ \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \end{aligned}$$

et donc

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta} \mathbb{E} \left[ \|\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \bar{\mathbf{X}}\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\mathbf{X} + \delta; \Theta) - \mathbf{X})\|^2 \right] + \mathbb{E} \left[ \|\epsilon\|^2 \right] \\ &= \arg \min_{\Theta} \mathbb{E} \left[ \|(\mathcal{N}(\tilde{\mathbf{X}}; \Theta) - \mathbf{X})\|^2 \right] \end{aligned}$$



## NOISE 2 NOISE: PREUVE

- ▶ Ainsi, minimiser la MSE loss  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  avec des data bruitées comme targets

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \bar{\mathbf{x}}_i\|^2$$

revient à minimiser la MSE loss  $\mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta)$  (originale) avec des data propres comme targets

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \mathbf{x}_i\|^2$$

- ▶ En d'autres termes, l'entraînement du DAE sur des targets bruitées équivaut à l'entraînement du DAE sur des targets propres, ce qui engendre donc un "débruiteur".

## NOISE 2 NOISE: PREUVE

- ▶ Ainsi, minimiser la MSE loss  $\mathcal{L}(\bar{\mathbf{X}}, \tilde{\mathbf{X}}; \Theta)$  avec des data bruitées comme targets

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \bar{\mathbf{x}}_i\|^2$$

revient à minimiser la MSE loss  $\mathcal{L}(\mathbf{X}, \tilde{\mathbf{X}}; \Theta)$  (originale) avec des data propres comme targets

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \|\mathcal{N}(\tilde{\mathbf{x}}_i; \Theta) - \mathbf{x}_i\|^2$$

- ▶ En d'autres termes, l'entraînement du DAE sur des targets bruitées équivaut à l'entraînement du DAE sur des targets propres, ce qui engendre donc un "débruiteur".

## BIBLIOGRAPHIE



Fleuret, F. (2022).  
Deep Learning Course.