

10907 Pattern Recognition

Lecturers

Prof. Dr. Ivan Dokmanić (ivan.dokmanic@unibas.ch)

Tutors

Felicitas Haag (felicitas.haag@unibas.ch)

Alexandra Spitzer (alexandra.spitzer@unibas.ch)

Cheng Shi (cheng.shi@unibas.ch)

Vinith Kishore (vinith.kishore@unibas.ch)

Problem set 4

Math

Exercise 1 (Conditional expectations and a new kind of orthogonality - ★). Let (X, Y) be two real-valued jointly distributed random variables and define their “ \mathbb{P} inner product” as

$$\langle X, Y \rangle_{\mathbb{P}} = \mathbb{E}(XY).$$

Check quickly that this is indeed an inner product. The corresponding norm is given as

$$\|X\|_{\mathbb{P}}^2 = \langle X, X \rangle_{\mathbb{P}} = \mathbb{E}(X^2).$$

Let now

$$\hat{X} = \arg \min \mathbb{E}[X | Y]$$

be the MMSE estimator of X given Y . Show that the estimator error $\hat{X} - X$ is orthogonal to Y in the geometry induced by $\langle \cdot, \cdot \rangle_{\mathbb{P}}$, that is to say

$$\langle \hat{X} - X, Y \rangle_{\mathbb{P}} = 0$$

Try to interpret this in words. How can we think of the conditional expectation?

Hint: First show that

$$\mathbb{E}[\mathbb{E}[XY | Y]] = \mathbb{E}[Y\mathbb{E}[X | Y]] \quad \text{and} \quad \mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X].$$

You can assume that all the needed joint and conditional probability density functions exist.

Exercise 2 (Cascading convolutions - ★★). Let \mathbf{A} and \mathbf{B} be two circulant matrices of size $N \times N$. Show in two different ways that \mathbf{AB} and \mathbf{BA} are also circulant:

- Using direct calculations for matrix multiplication;
- Using the convolution–multiplication rule for discrete Fourier transforms (\equiv the fact that circular convolutions are diagonalized by the DFT).

Exercise 3 (Filling in the steps from lecture notes - ★★). Let \mathbf{w} be a random image of size $N \times N$ with pixels $w[m, n]$. Let also

$$a_{\mathbf{ww}}[k, \ell] = \mathbb{E}[w[m, n]w[(m+k)_N, (n+\ell)_N]]$$

be its (periodic) autocorrelation function at shift (k, ℓ) (NB: if \mathbf{w} were complex-valued we'd have to conjugate the shifted image). We implicitly assumed that the distribution of \mathbf{w} is (wide-sense) stationary. Show that the DFT of $a_{\mathbf{ww}}$ is its power-spectral density (PSD), defined as

$$S_{\mathbf{ww}}[u, v] = \mathbb{E} |W[u, v]|^2$$

where \mathbf{W} is the discrete Fourier transform of \mathbf{w} (itself a random vector).

Show furthermore that under stationarity the different frequencies are uncorrelated; that is, that

$$\mathbb{E} W[u, v] \overline{W[u', v']} = 0$$

when $(u, v) \neq (u', v')$, with overline denoting complex conjugation.



University
of Basel

Exercise 4 (Fun with CNNs - ★★). The following listing defines a simple U-Net mapping images to images (adapted from <https://github.com/clemkao/u-net>):

```

1  import torch
2  import torch.nn as nn
3  from typing import List
4
5  class DoubleConv(nn.Module):
6      """(Conv -> BN -> ReLU) x 2, 'same' spatial size for k=3"""
7      def __init__(self, in_ch, out_ch, dilation=1):
8          super().__init__()
9          pad = dilation
10         self.net = nn.Sequential(
11             nn.Conv2d(in_ch, out_ch, kernel_size=3, padding=pad, dilation=dilation,
12                 ↪ bias=False),
13             nn.BatchNorm2d(out_ch),
14             nn.ReLU(inplace=True),
15             nn.Conv2d(out_ch, out_ch, kernel_size=3, padding=pad, dilation=dilation,
16                 ↪ bias=False),
17             nn.BatchNorm2d(out_ch),
18             nn.ReLU(inplace=True),
19         )
20
21     def forward(self, x):
22         return self.net(x)
23
24 class Down(nn.Module):
25     """Downsample by 2 then DoubleConv"""
26     def __init__(self, in_ch, out_ch, dilation=1):
27         super().__init__()
28         self.pool = nn.MaxPool2d(2)
29         self.conv = DoubleConv(in_ch, out_ch, dilation=dilation)
30
31     def forward(self, x):
32         return self.conv(self.pool(x))
33
34 class Up(nn.Module):
35     """Upsample by 2 then concat skip and DoubleConv"""
36     def __init__(self, in_ch, out_ch):
37         super().__init__()
38         self.up = nn.ConvTranspose2d(in_ch, in_ch // 2, kernel_size=2, stride=2)
39         self.conv = DoubleConv(in_ch, out_ch)
40
41     def forward(self, x, skip):
42         x = self.up(x)
43         # pad if needed (odd sizes)
44         diffY = skip.size(2) - x.size(2)
45         diffX = skip.size(3) - x.size(3)
46         if diffY or diffX:
47             x = nn.functional.pad(x, [diffX // 2, diffX - diffX // 2,
48                 ↪ diffY // 2, diffY - diffY // 2])
49         x = torch.cat([skip, x], dim=1)
50         return self.conv(x)
51
52 class OutConv(nn.Module):
53     def __init__(self, in_ch, out_ch):
54         super().__init__()
55         self.conv = nn.Conv2d(in_ch, out_ch, kernel_size=1)
56
57     def forward(self, x):

```



```

56         return self.conv(x)
57
58 class SimpleUNet(nn.Module):
59     def __init__(self, in_ch=1, out_ch=1, base=32, depth=4, dilations: List[int] =
    ↪ None):
60         super().__init__()
61         if dilations is None:
62             dilations = [1] * depth
63         assert len(dilations) == depth
64
65         ch = [base * (2 ** i) for i in range(depth)] # e.g., 32, 64, 128, 256
66
67         # encoder
68         self.inc = DoubleConv(in_ch, ch[0], dilation=dilations[0])
69         self.downs = nn.ModuleList([
70             Down(ch[i], ch[i+1], dilation=dilations[i+1]) for i in range(depth-1)
71         ])
72
73         # decoder
74         self.ups = nn.ModuleList([
75             Up(ch[i], ch[i-1]) for i in range(depth-1, 0, -1)
76         ])
77
78         self.outc = OutConv(ch[0], out_ch)
79
80     def forward(self, x):
81         xs = [self.inc(x)]
82         for d in self.downs:
83             xs.append(d(xs[-1]))
84         y = xs[-1]
85         for i, up in enumerate(self.ups, start=1):
86             y = up(y, xs[-i-1])
87         return self.outc(y)

```

- Give a formula for the total number of network parameters in terms of the four hyperparameters (`in_ch`, `out_ch`, `base`, `depth`). How much do you get for `in_ch = 3`, `out_ch = 1`, `base = 32`, `depth = 4`? Which parts dominate the parameter count?
- Compute the output receptive field R ; that is, given a pixel in the output image, which pixels (how many) in the input image can it depend on?
- How would you increase the receptive field without increasing the number of parameters too much? (*Hint*: check out dilated convolutions.)
- Is `SimpleUNet` linear?
- Explain what will happen if we replace the ReLU activations with linear activations $r(x) = x$.
- We like convolutions because they are translation invariant. `SimpleUNet` is a convolutional neural net. Is it strictly translation invariant (ignoring boundary effects)?

Coding

Exercise 5 (Audio Denoising with Wiener and CNNs - ★★). In this coding exercise we will use two strategies for audio denoising: a simple Wiener filter and a convolutional neural network. We'll hear that the CNN does better, but a Wiener filter will also be okay given that we will show it much less data: we will only use a single noisy audio clip and some magic.

In both cases we will use the discrete Fourier transform (DFT) to pass to the frequency domain, but instead of applying the DFT to the entire length of the signal we will split the audio into (overlapping) frames and apply the transform to the frames. This is called the (discrete) short-time Fourier transform (STFT). Using a single huge DFT of the entire length of the audio signal has a number of drawbacks. The main one is that the DFT tells us what frequencies are present in the signal, but not where they are. For audio that matters. This is why the DFT is computed on short overlapping frames: then we at least know that the frequencies belong to a given frame. Another benefit is that we can use the same processing for arbitrarily long audio signals since the frame size is fixed.

We will use the dataset available at <https://datashare.ed.ac.uk/handle/10283/2791>. This dataset contains clean speech audio files and versions corrupted with different types of noise at different signal-to-noise ratios (SNRs)—not only additive white Gaussian noise. Make sure to listen to a couple of samples to get a feel for the data.

STFT processing Let $x[n]$ be a single-channel (mono) audio clip of length N (and thus duration N/f_s with f_s the sampling frequency). As mentioned in the exercise introduction, we'll split x into overlapping frames of length M_w with hop size H (the number of samples between the starts of consecutive frames). This is illustrated in Figure 1.

The number of frames is

$$K = \left\lceil \frac{N - M_w}{H} \right\rceil + 1.$$

Each frame is multiplied with a “window” $w[m]$, $m = 0, \dots, M_w - 1$ which rises and falls gradually; see Fig. .2. You do not have to worry about this since we use a ready-made implementation of STFT (the windowing prevents abruptly cutting off at the frame boundaries, which

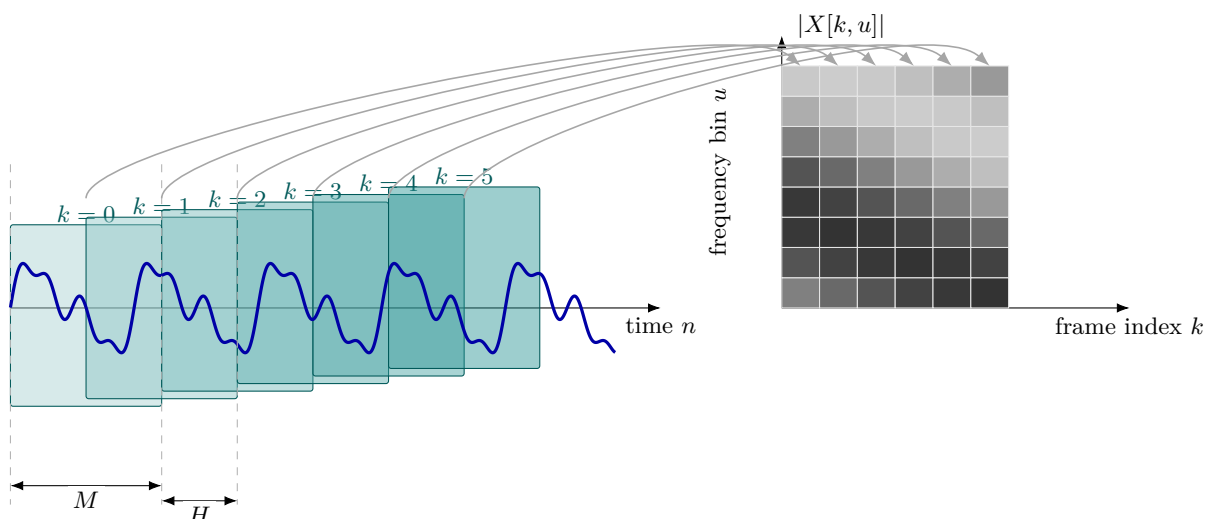


Figure 1: Short-time Fourier transform (STFT): overlapping windows of length M_w slide with hop size H across a signal. Each frame k yields a spectrum across frequency bins u , forming the time–frequency representation $X[k, u]$.

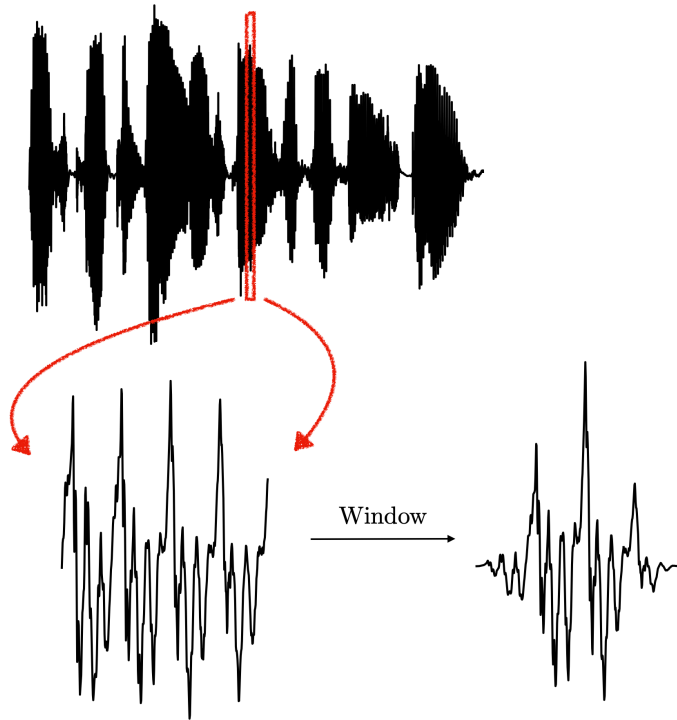


Figure 2: Effect of windowing a speech sample. We extract a 512 length sample segment from a clean speech signal and apply a Hann window, which tapers the edges to reduce spectral leakage.

among other things helps eliminate some audible artifacts). The short-time Fourier transform (STFT) of $x[n]$ is then defined as

$$X[k, u] = \sum_{m=0}^{M_w-1} x[m + kH] w[m] e^{-j2\pi um/N_{\text{fft}}}, \quad k = 0, \dots, K-1, \quad u = 0, \dots, N_{\text{fft}}-1$$

with $M_w \leq N_{\text{fft}}$. Here, $X[k, u]$ is the STFT coefficient for frame k and frequency bin u . It is really just a DFT of the samples of x with indices from kH to $kH + M_w - 1$, multiplied by the window $w[m]$.

A simple two-step inverse STFT (iSTFT) is:

$$\begin{aligned} \tilde{x}_k[m] &= \frac{1}{N_{\text{fft}}} \sum_{u=0}^{N_{\text{fft}}-1} X[k, u] e^{j2\pi um/N_{\text{fft}}}, \quad m = 0, \dots, M_w - 1, \\ \hat{x}[n] &= \frac{1}{C} \sum_{k=0}^{K-1} w[n - kH] \tilde{x}_k[n - kH], \quad \text{with } 0 \leq n - kH < M_w, \end{aligned}$$

where the window w satisfies the so-called constant-overlap-add (COLA) condition (i.e., $\sum_k w^2[n - kH] = C$ is constant in n). The first line simply inverts the DFT on each frame; the second line overlaps and adds the frames back together to form the reconstructed signal $\hat{x}[n]$. You don't need to know these details about the STFT on the exam. We are giving them here for completeness. You do need to understand what is going on here conceptually.

STFT turns 1D audio into a 2D time-frequency representation which can be processed as an image. Since STFT is complex valued, we often only work on the magnitude spectrogram $|X[k, u]|$. An example is in Figure 3. This is exactly what we will do with the CNN below—we will treat $|X[k, u]|$ as an image and process it by 2D convolutions.

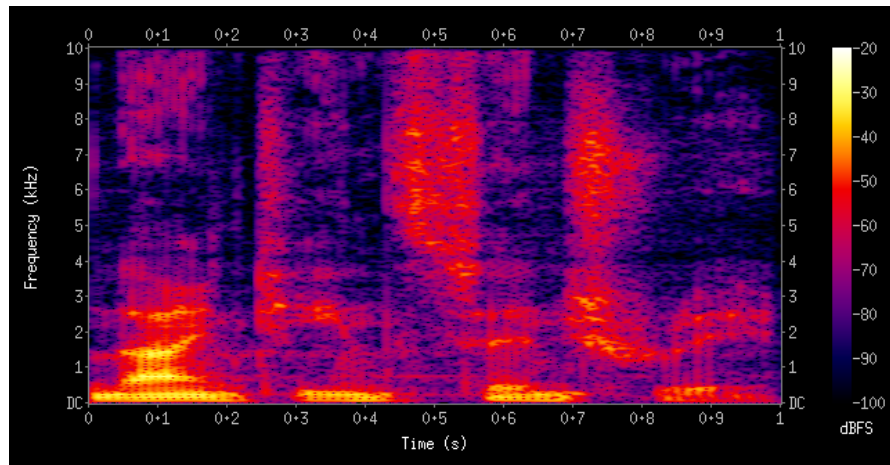


Figure 3: Spectrogram of the spoken words “nineteenth century”. Frequencies are shown increasing up the vertical axis, and time on the horizontal axis. The legend to the right shows that the color intensity increases with the magnitude (power). Work by Aquegg from <https://en.wikipedia.org/wiki/Spectrogram#/media/File:Spectrogram-19thC.png>.

Wiener filter The first task is to denoise audio with a Wiener filter in the STFT domain using one gain per frequency, applied to every frame. We assume additive noise $y[n] = x[n] + w[n]$ with $X = \text{STFT}(x)$, $Y = \text{STFT}(y)$, $W = \text{STFT}(w)$. Just for fun we will assume that we only have access to the noisy signal y and its STFT Y . This means that we should estimate the per-frequency power spectral densities (PSDs) from the noisy clip alone by aggregating across frames k , then form the Wiener gain $G[u]$.

As in lecture, we'll use the formula

$$G[u] = \frac{\text{SNR}[u]}{1 + \text{SNR}[u]} = \frac{S_{\mathbf{xx}}[u]}{S_{\mathbf{xx}}[u] + S_{\mathbf{ww}}[u]}.$$

To estimate $S_{\mathbf{ww}}[u]$ we will use the fact that speech is relatively “sparse” in the STFT representation. This means that many $[k, u]$ bins $X[k, u]$ of speech are very small and there are only a handful of large coefficients. We have seen something similar when summing up the frequency components of an image starting with the largest ones (which were mostly low frequencies). This means that for many frequency bins u , the lower percentiles of the magnitudes $|Y[k, u]|$ across frames k will correspond to noise-only coefficients. We can thus estimate the noise PSD using a quantile estimator and then use the average power to estimate the noisy signal PSD. Your task is to implement this Wiener filter in the provided Python files.

Here is how it works: In the provided python file `wiener.py`, implement the following steps in the function `wiener_enhance_freqavg()`: For each frequency u :

- Collect the magnitudes $|Y[k, u]|$ across all frames k .
- Estimate the noise PSD as $\hat{S}_{\mathbf{ww}}[u] = (\text{quantile}_k\{|Y[k, u]|\})^2$; (e.g., $q = 10\%$).
- Estimate the noisy signal PSD $\hat{S}_{\mathbf{yy}}[u]$ as the average power across frames:

$$\hat{S}_{\mathbf{yy}}[u] = \frac{1}{K} \sum_{k=0}^{K-1} |Y[k, u]|^2.$$

- Estimate the clean signal PSD as $\hat{S}_{\mathbf{xx}}[u] = \max\{\hat{S}_{\mathbf{yy}}[u] - \hat{S}_{\mathbf{ww}}[u], 0\}$. This is motivated by the formula $S_{\mathbf{yy}} = S_{\mathbf{xx}} + S_{\mathbf{ww}}$ which holds for additive noise uncorrelated with the signal.

- Compute the Wiener filter $G[u] = \hat{S}_{\mathbf{xx}}[u] / (\hat{S}_{\mathbf{xx}}[u] + \hat{S}_{\mathbf{ww}}[u])$ and apply this same $G[u]$ to all frames: $\hat{S}[k, u] = G[u]Y[k, u]$.

Finally, invert with iSTFT. Use the provided STFT wrappers `stft(x, ...)` and `istft(X, ...)` from `audio_helper.py`. Please do not change its signature or return type. Return a 1-D `torch.Tensor` at 16 kHz, same dtype and device as the input. For numerical stability, use a small `eps` in divisions. An optional `gain_floor` may be applied by clamping G from below.

When you are done with implementing your `wiener_enhance_freqavg()` function, you can apply your Wiener filter to a few noisy clips and listen to the results. You can listen and sanity-check your result by running `eval_wiener_and_export(...)` in `audio_denoising.py`, which calls your `wiener_enhance_freqavg` and prints the SNR improvement while saving a few WAVs to a folder "out". You will hear that sometimes it works quite well and sometimes it leaves a lot to be desired. For what kind of noise does it work well? How is this related to the assumptions used to estimate the noise PSD? Try playing with the quantile q used for noise PSD estimation.

Convolutional neural network We will now compare this Wiener filter with a very simple convolutional neural network (CNN). The CNN is provided as `TinyDenoiser` in `models.py`, with training and evaluation entry points `train_cnn_denoiser(...)` and `eval_cnn_and_export(...)` in `audio_denoising.py`. You do not need to modify the CNN; it is there to let you compare to the Wiener result and to listen to outputs.

This is what the code does: The CNN will take as input the magnitude spectrogram $|Y[k, u]|$ of the noisy audio and output an estimate of the clean magnitude spectrogram $|\hat{X}[k, u]|$. The network treats the magnitude spectrogram as a single-channel image, where the time frames correspond to one spatial dimension and the frequency bins to the other spatial dimension. The phase of the noisy audio $Y[k, u]$ will be reused for reconstruction; the CNN does not use it at all. The CNN architecture is very simple: it consists of a few convolutional layers with ReLU activations, followed by a final convolutional layer with a sigmoid activation to ensure that the output magnitudes are non-negative.

Again, use the provided code to apply the learned non-linear filter to a couple of noisy clips and listen to the results. How do they compare with the Wiener filter? This is a simple version of the kind of audio processing implemented in the devices you're using every day!