

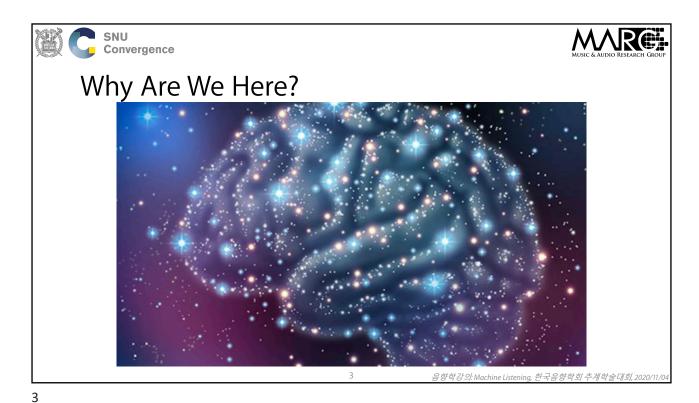


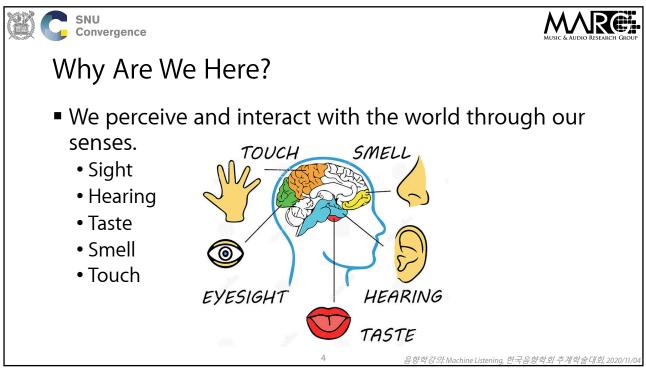
# 인공청각지능: 소리에서 의미로

Machine Listening: from sound to meaning

서울대학교 지능정보융합학과 음악오디오연구실 이교구











# Why Are We Here?

- Unfortunately, research in hearing (audition) is outnumbered by research in sight (vision) in every measure:
  - Number of researchers/scholars/students
  - Number of books/papers/articles
  - Number of academic communities
  - And more...

5

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

5





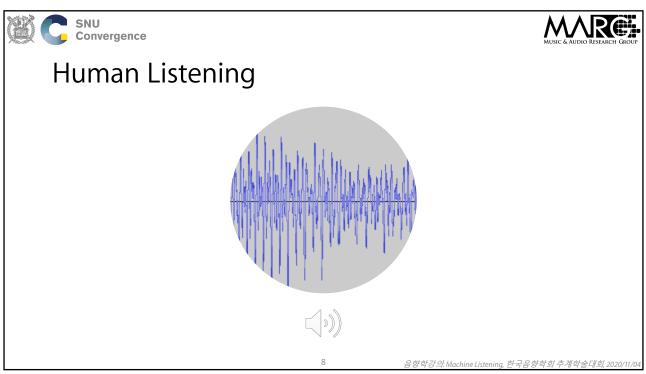
# Why Are We Here?

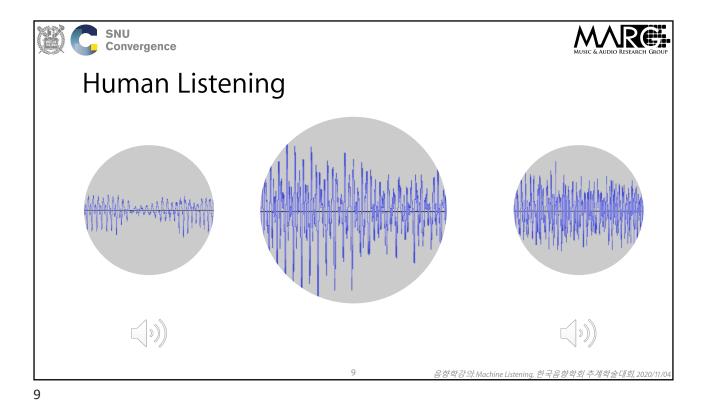


6

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0







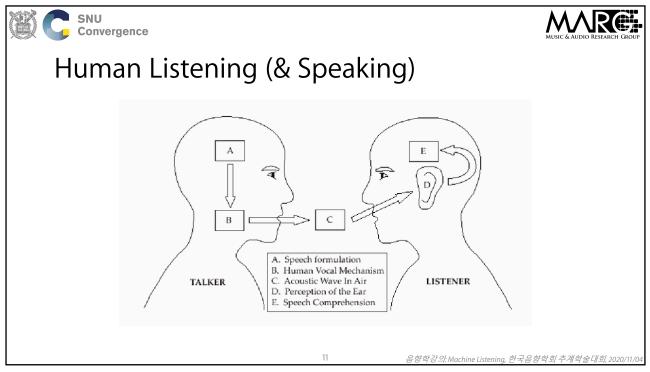
### Convergence

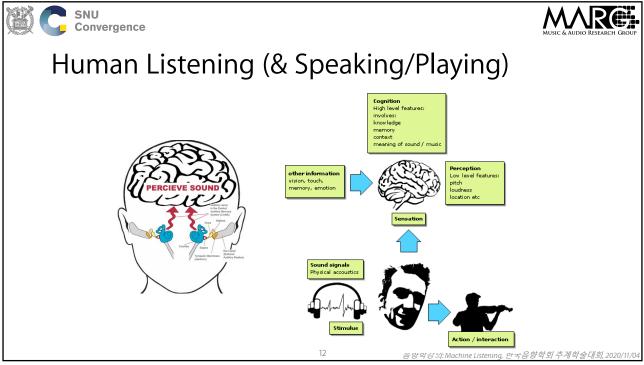
Human Listening

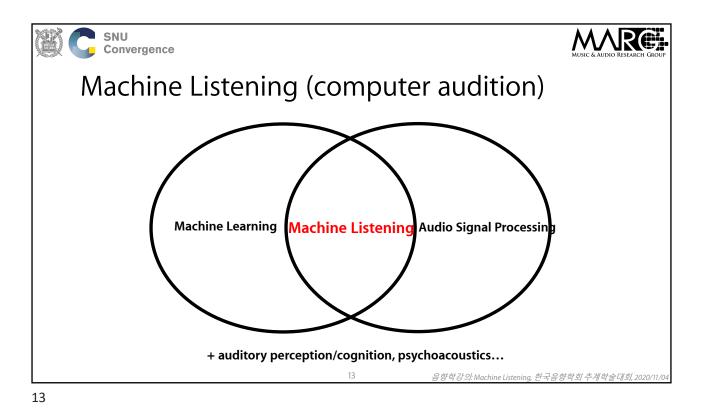
Arabic speech English speech (Adult directed) (Infant directed)

Wherearethesilences between words?

The segmentation problem Where are the silences between words?











# Fundamentals of Digital Audio Signal Processing

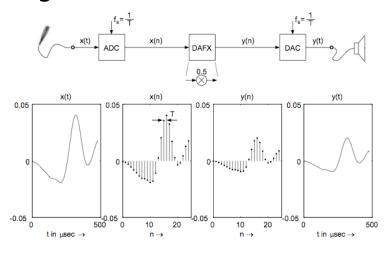
14

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0-





# **Digital Signals**



·향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

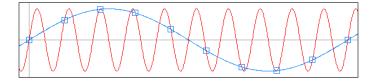
15





# Sampling

- According to the sampling theorem:  $f_s > 2f_{max}$  (Nyquist limit)
- Otherwise there is another, lower-frequency, signal that share samples with the original signal (aliasing)

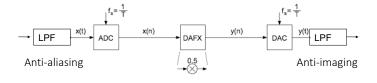


Related to the wagon-wheel effect: www.michaelbach.de/ot/mot wagonWheel/index.html



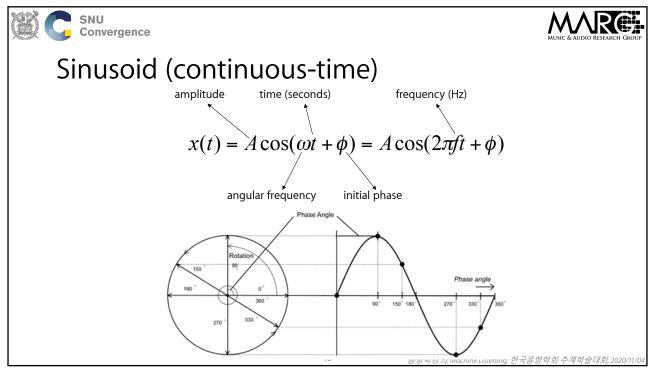


# Anti-aliasing[imaging]



17

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

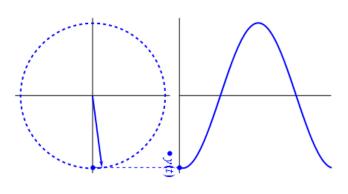






# Sinusoid (continuous-time) in action

$$x(t) = A\cos(\omega t + \phi) = A\cos(2\pi f t + \phi)$$



19

『향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

19





# Sinusoid (discrete-time)

• Since 
$$t = \frac{n}{R}$$
,  $x(n) = A\cos(\frac{2\pi fn}{R} + \phi)$ ,

where R =sampling rate







#### **Fourier Theorem**

 Fourier theorem says: any periodic signal can be decomposed into a sum of sinusoids

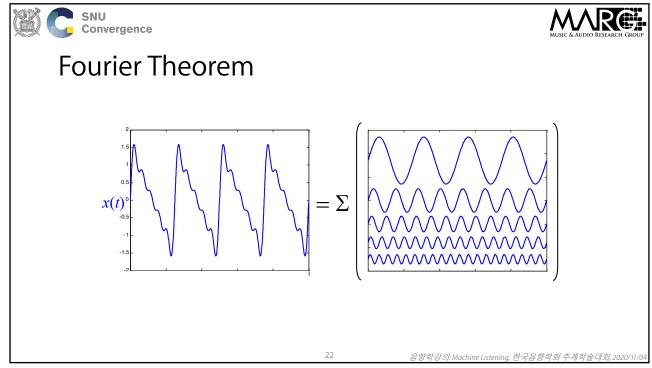
$$x(t) = a_0 + \sum_{k=1}^{\infty} a_k \cos(2\pi k f t + \phi_k)$$

- f is called fundamental and 2f, 3f, ... are harmonics of fundamental
- Sequence of sinusoids with harmonic frequency is harmonic series



21

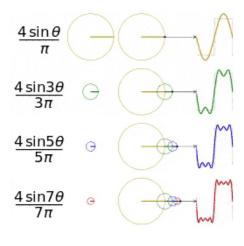
음향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/04







#### **Fourier Theorem**



23

『향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

23





# Fourier Transform (FT)

■ The Fourier transform of a continuous-time signal x(t) may be defined as

$$X(\omega) = FT[x(t)] = \int_{-\infty}^{\infty} x(t)e^{-j\omega t}dt, \quad \omega \in (-\infty, \infty)$$





#### Discrete Fourier Transform (DFT)

• The Discrete Fourier Transform of a signal x(n) may be defined as

$$X(k) = DFT[x(n)] = \sum_{n=0}^{N-1} x(n)e^{-j2\pi nk/N}, \quad k = 0,1,...,N-1$$

• The resulting N samples X(k) are complex-valued:

$$\begin{split} X(k) &= X_R(k) + j X_I(k) \\ &\left| X(k) \right| = \sqrt{X_R^2(k) + X_I^2(k)} \quad : \text{magnitude} \\ \varphi(k) &= \arctan \frac{X_I(k)}{X_R(k)} \quad : \text{phase} \end{split}$$

$$k = 0,1,...,N-1$$

25

용향학강의: Machine Listenina. 한국음향학회 추계학술대회 2020/11/04

25

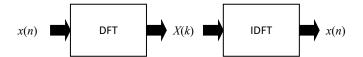




#### Inverse DFT (IDFT)

■ The DFT allows perfect reconstruction of a signal x(n) from its DFT X(k) via inverse DFT defined as:

$$x(n) = IDFT[X(k)] = \frac{1}{N} \sum_{k=0}^{N-1} X(k)e^{j2\pi nk/N}, \quad n = 0, 1, ..., N-1$$



26

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





#### Fast Fourier Transform (FFT)

- An algorithm to efficiently compute the DFT is known as the Fast Fourier Transform (FFT) and its inverse as the IFFT
- Computational complexity
  - N-point DFT:  $O(N^2)$
  - N-point FFT:  $O(N \log N)$
- The FFT is so fast that even time-domain operations, like convolution, can be performed faster using FFT and IFFT instead:

$$(x*h)(n) = \sum_{m=0}^{N-1} x(m)h(n-m)$$

$$x(n) \longrightarrow X(k) \longrightarrow X(k) \longrightarrow X(k)H(k) \longrightarrow IFFT \longrightarrow (x*h)(n)$$

$$h(n) \longrightarrow H(k) \longrightarrow X(k)H(k) \longrightarrow IFFT \longrightarrow (x*h)(n)$$

27

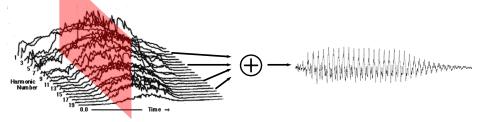
·향학강의: Machine Listenina. 한국음향학회 추계학술대회, 2020/11/04

27

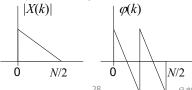




#### Revisiting the Fourier Theorem



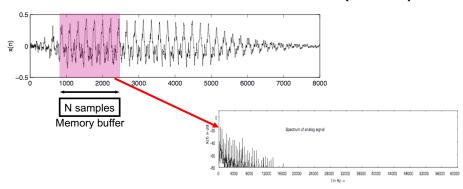
- Any periodic signal can be described by a sum of a series of sinusoids with timevarying amplitudes and phases
- Thus a complex spectrum is just a snapshot of those sinusoids' parameters







# Short-time Fourier Transform (STFT)



- For block processing, a short segment is sent to a buffer and processed as a block
- The DFT done in this way is called the short-time Fourier Transform or STFT

29

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

29





# Frequency Resolution

- Determined by how many sinusoids are used to describe a spectrum
- In N-point DFT or FFT, the frequency resolution is given by

$$\Delta f = f_s/N$$

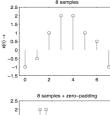
- $\ ^{\blacksquare}$  Intuitively, we can have finer frequency resolution if we increase N
- However, this results in poorer temporal resolution => tradeoff between frequency[spectral] and time[temporal] resolution

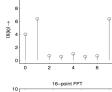


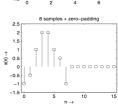


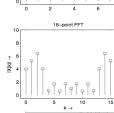
# Zero-Padding

- lacktriangledown A possible solution to increase frequency resolution without increasing N *i.e.*, without losing time resolution
- Add zero-valued samples thus doesn't change spectrum itself to yield better spectral resolution









8 to 15 k → 10 15 k → 10 15 m + 1

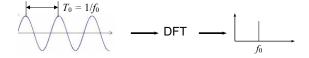
31





#### DFT of a Sinusoid

 In theory the DFT of a pure sinusoid results in a single sharp line at the frequency of the sinusoid

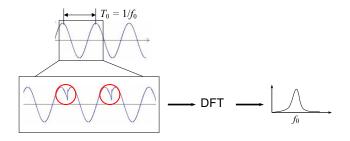






# Spectral Leaking

- In practice, unless we perform  $f_0$ -synchronous analysis, there are discontinuities (sharp changes) at the segment boundaries that introduce some noise. Thus the spectral line around  $f_0$  is smeared.
- This is known as spectral leaking



3

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

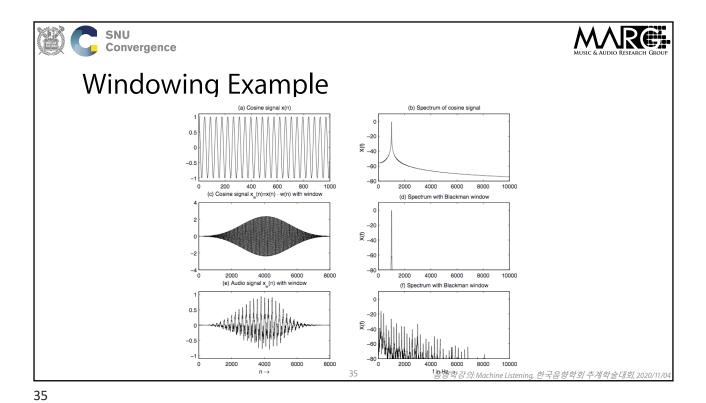
33





# Windowing

- In order to reduce spectral leaking, we need to avoid abrupt changes between the segment boundaries
- This is done by multiplying a window whose amplitude gradually reaches zero at both ends, thus guaranteeing the continuity of a segmented signal when repeated
- Possible windows are: rectangular, hamming, hann, Blackman, Gaussian, and so on.

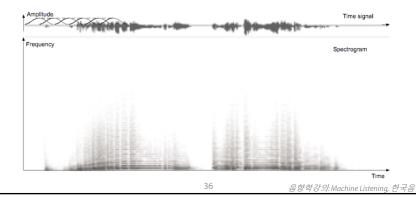






# 2-D Time-Frequency Representation

- Using the STFT, independent DFTs are calculated on windowed segments
- The segments usually overlap to compensate for the loss of temporal resolution
- Produces a 2-D spectrogram







#### **Acoustic Features**

37

음향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/0

37





#### **Audio Feature Extraction**

- Raw audio samples are:
  - Noisy
  - Redundant
  - Computationally inefficient
  - Not a good input to audio applications
- Need to convert them to more robust, compact yet meaningful representations called *audio features* or acoustic features

38

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





#### Acoustic (audio) Features

- Spectral features
  - Spectral low-level features: spectral centroid, spectral flatness measure, spectral flux, etc.
  - Spectral envelope: LPCs, MFCCs, etc.
- Temporal features
  - ZCR (zero crossing rate), tempo histogram, novelty function, etc.
- Tonal features
  - PCP (pitch class profile) or chroma, chromagram, tonal centroid, etc.

39

음향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/0

39





# **Spectral Low-level Features**





# Spectral Centroid (SC)

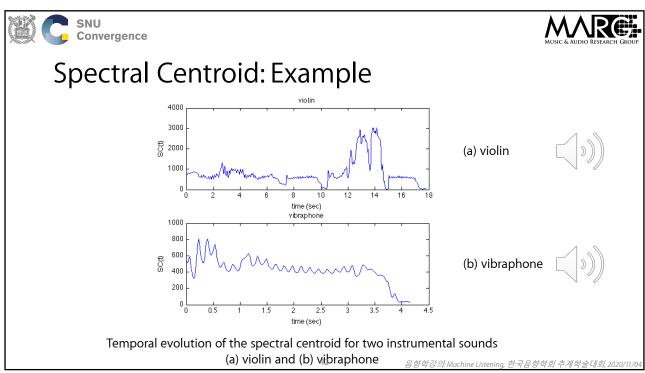
$$SC = \frac{\sum_{k=0}^{N/2} f_k |X(k)|^2}{\sum_{k=0}^{N/2} |X(k)|^2},$$

where  $f_k$  is the center frequency of kth bin and |X(k)| is the DFT

- Characterizes the center of gravity of the (power) spectrum
- Usually associated with the "sharpness / dullness (or brightness / darkness)" of a sound

41

음향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/04







## Spectral Spread (SS)

$$SS = \sqrt{\frac{\sum_{k=0}^{N/2} (f_k - SC)^2 |X(k)|^2}{\sum_{k=0}^{N/2} |X(k)|^2}}$$

- Measure of the average spread of the spectrum in relation to its centroid
- Noisy, broadband signals have high SS while tonal sounds show lower SS

43

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

43





#### Spectral Flatness Measure (SFM)

$$SFM_{b} = \frac{\sqrt[N_{b}]{\prod_{k_{b}}|X(k_{b})|^{2}}}{\frac{1}{N_{b}}\sum_{k_{b}}|X(k_{b})|^{2}}, \quad k_{b} = k_{l}, k_{l} + 1, ..., k_{u}$$

where  $N_b$  is the number of spectral bins in a subband or  $N_b = k_u - k_l + I$ 

- Reflects how "flat" a signal's power spectrum is
- Calculated as the ratio of the geometric mean and the harmonic mean
- Usually computed per spectral band (critical bands or bark bands, etc.), thus SFM<sub>b</sub>
- Flatness for the whole spectrum is the average of the subband measures





# Harmonic Spectral Centroid (HSC)

$$HSC = \frac{\sum_{h=1}^{N_h} f_h A_h}{\sum_{h=1}^{N_h} A_h}$$

where  $f_h$  and  $A_h$  are the frequency and the amplitude of the  $h_{th}$  harmonic, respectively

- Measure of the amplitude-weighted mean of the harmonic (spectral) peaks of the spectrum
- Compared to SC, HSC focuses only on harmonic (spectral) peaks, which are more musically meaningful in general
- Harmonic spectral spread (HSS) is similarly defined

45

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

45





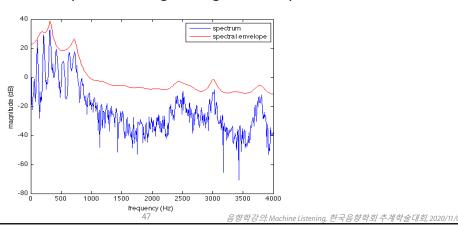
# **Spectral Envelope**





# What is the Spectral Envelope?

 Spectral envelope is a smoothed version of the spectrum that preserves its overall shape while neglecting its fine spectral structure



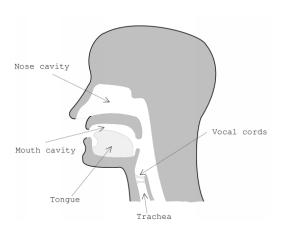
47





# Human Speech System

- Vocal cords act as an oscillator, which generates a spectrally rich source signal
- Everything else is filter: vocal tract, mouth/nose cavity, tongue
- Thus called "source-filter" model
- These filters define the shape of the spectral envelope



18

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0-





# Spectral Envelope Estimation

- A few popular techniques to estimate the spectral envelope
  - Channel vocoder: estimates the amplitude of the signal within several frequency bands
  - Linear prediction: estimates the parameters (or filter coefficients) of a filter that approximates the spectrum
  - Cepstrum analysis: inverse-FFT the log-spectrum and low-pass filters it to obtain the envelope

49

우향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

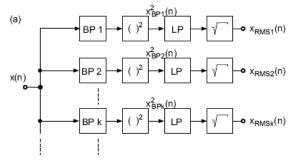
49





# Channel Vocoder (1)

- Filters a signal with a bank of bandpass filters
- Calculates RMS of each bandpassed signal
- The more filters used, the finer spectral envelope estimated



50

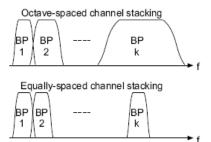
음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





# Channel Vocoder (2)

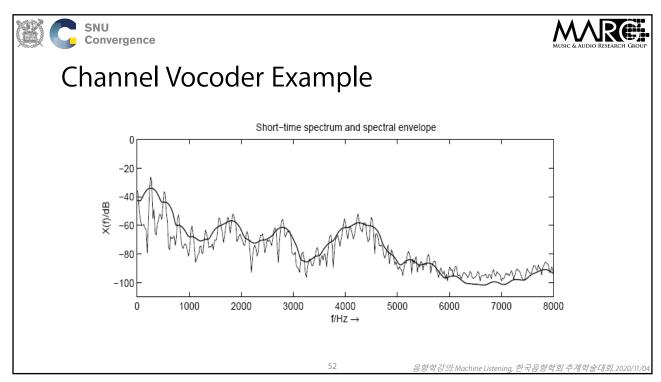
 In the frequency domain, multiply the spectrum with the filters' frequency response and square-root the sum of each filter's output



 The filterbank can be either linearly or logarithmically spaced (e.g., constant-Q or Mel-scale filter bank)

51

유향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04







# Linear Predictive Coding (1)

• Linear predictive coding is a source-filter model that approximates the way a sound is generated as an excitation (a pulse train or noise) passing through an all-pole resonant filter

Resonant filter
(spectral envelope model)

Synthesized sound

- Widely used in speech and music applications
- Reduces large amount of data (e.g., N samples) to a few filter coefficients while preserving the overall shape

53

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

53





# Linear Predictive Coding (2)

■ The nth sample x(n) is extrapolated, i.e., predicted by a linear combination of p past samples:

$$x(n) \approx \hat{x}(n) = \sum_{k=1}^{p} a_k x(n-k)$$

The residual error is given by

$$e(n) = x(n) - \hat{x}(n) = x(n) - \sum_{k=1}^{p} a_k x(n-k)$$

and we want to minimize this error

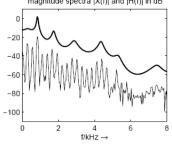




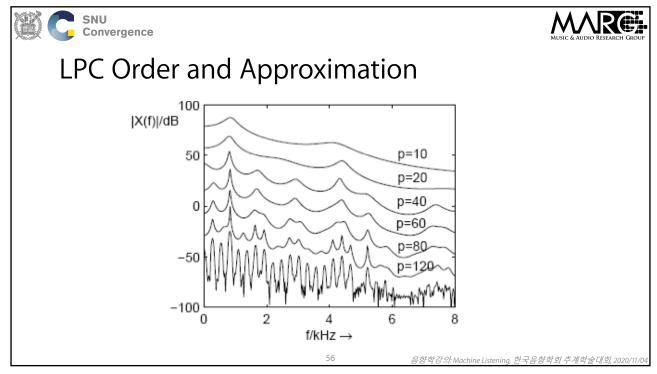
# Linear Predictive Coding (3)

- The parameters  $a_k$ 's are called linear predictive coefficients (LPCs)
- The filter represented by these coefficients is a resonant filter and its frequency response represents the spectral envelope

■ The higher the filter order p, the closer the approximation is to the signal's spectrum  $_{\text{magnitude spectra} \mid X(f) \mid \text{and} \mid H(f) \mid \text{in dB}}$ 



음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0







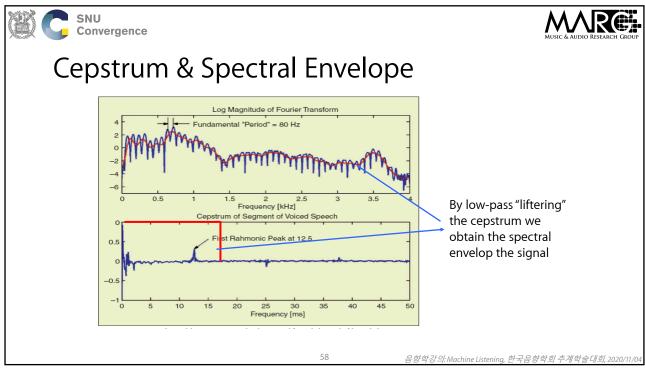
# Cepstrum Analysis

- Cepstrum is the result of taking the FFT of the log-spectrum as if it were a signal
- Measures the rate of change in different spectral bands
- The name cepstrum was coined by Bogert *et al.* (1963) by reversing the first four letters of the spectrum (similarly for quefrency alanysis and liftering, etc.)
- For a real signal x(n), the real cepstrum is calculated as follows:

$$c_R(n) = IFFT(\log(|X(k)|))$$

57

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04







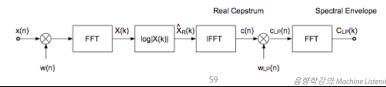
# Spectral Envelope Estimation

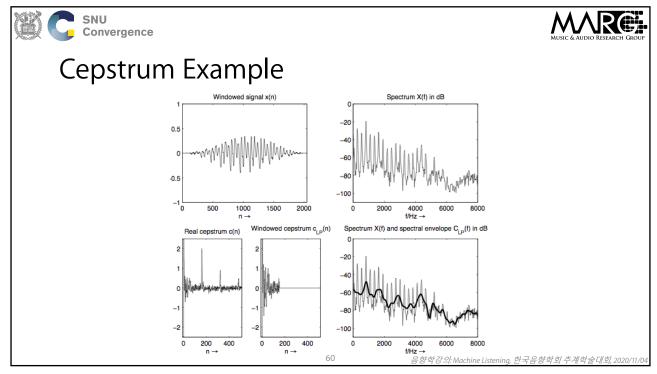
Using a low-pass window of the form:

$$\omega_{LP}(n) = \begin{cases} 1 & n = 0, N_1 \\ 2 & 1 \leq n \leq N_1 \\ 0 & N_1 < n \leq N-1 \end{cases}$$

we can low-pass the cepstrum and obtain the spectral envelope by:

$$\begin{split} c_{LP}(n) &= c_R(n) \cdot \omega_{LP}(n) \\ C_{LP}(k) &= FFT[c_{LP}(n)] \end{split}$$

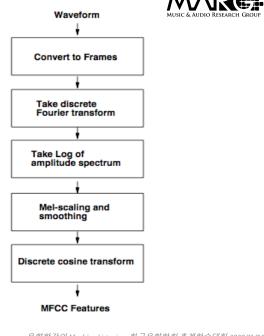






#### **MFCC**

- Mel-frequency Cepstral Coefficients (MFCCs) are a variation of cepstrum, motivated by human perception (Logan, 2000)
- Most extensively used in speech and music applications (e.g., speech recognition, genre classification, instrument recognition, etc.), due to its ability to compactly represent the spectral characteristics (just ~13 coefficients)



음향학강의: <u>M</u>achine Listening, 한국음향학회 추계학술대회, 2020/11/04

61





#### Mel Scale

- The Mel scale is a non-linear perceptual scale of pitches judged to be equidistant
- Approximately linear below 1 kHz and logarithmic above
- 1 kHz corresponds to 1000 Mel (reference point)
- With the Mel scale, a 1000-Mel tone should sound as twice as high as a 500-Mel tone (this is not true with linear frequency Hz)

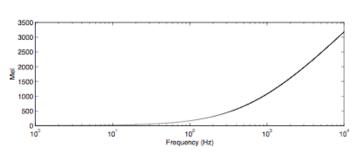




# Mel vs. Linear Frequency

■ The relation between Mel and Hz is given by

$$m = 1127.01048\log(1 + f/700)$$
$$f = 700(e^{m/1127.01048} - 1)$$



음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

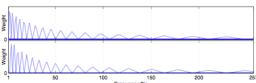
63



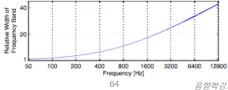


# Mel-frequency Spectrum

To convert a linear spectrum to Mel we can use a filterbank of overlapping triangular windows:



• Such that the width d of each window increases according to the Mel scale, and the height of each triangle is 2/d







## Decorrelation of Mel-scale Spectrum

- The resulting Mel-scale spectral vectors are highly correlated with each other; i.e. highly redundant
- Thus a more efficient representation of the log-spectrum can be obtained by applying a transform that decorrelates those vectors (Rabiner and Juang, 93)
- This decorrelation is commonly approximated by means of the Discrete Cosine Transform (DCT)
- The DCT is similar to a DFT but only for real numbers. It has the property that most of its energy is concentrated on a few initial coefficients (thus effectively compressing the spectral info)

$$X_{DCT}(k) = \sqrt{\frac{2}{N}} \sum_{N=0}^{N-1} x(n) \cos[\frac{\pi}{N}(n+\frac{1}{2})k]$$

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

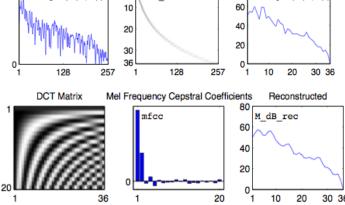
65





#### Fast Computation of MFCCs Power Spectrum [dB]

10\*log10(P(:,i))



Triangular Filter Matrix

mel filter

Mel Power Spectrum [dB]

10\*log10(M(:,i))

MFCCs roughly model certain characteristics of human auditory perception: the nonlinear perception of loudness and frequency and spectral masking (Pampalk, 2006)





# **Tonal Features**

67

음향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/0

67





# **Tonality**

- Very important attribute in (Western) tonal music
- Explain the relationship among the tones
- Several musical attributes are closely related with tonality
  - Scale
  - Key
  - Pitch
  - Interval
  - Chord

68

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





# Tonal Features vs. Spectral Features

- Spectral features
  - Good for describing certain spectral characteristics (e.g., sharpness, noisiness, etc.)
  - Good for representing sonic texture or timbre by capturing overall frequency magnitude response (e.g., LPCs, cepstral coefficients, MFCCs)
  - Not good for tonal analysis: pitch- or tone-relevant information gets lost
- Tonal features retain tonal structure in musical audio
  - Tonal relations
  - Interval relations

69

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

69





#### **Constant-Q Transform**





#### Constant-Q Transform

• In DFT, the center frequency  $f_k$  of the frequency bin is given by

$$f_k = \frac{f_s}{N}k, \quad k = 0,1,...,N-1$$

where  $f_s$  is the sampling rate and N is the DFT size

- Therefore, all the frequency bins are linearly spaced
- However, musical scale as well as human hearing mechanism are logarithmic
- Brown proposed the constant-Q transform whose frequency resolution conforms to equal-tempered scale (1990)
- Well suited for pitch-related analysis

71

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

71





#### Constant-Q Transform (cont'd)

■ In constant-Q transform the kth spectral frequency is defined as

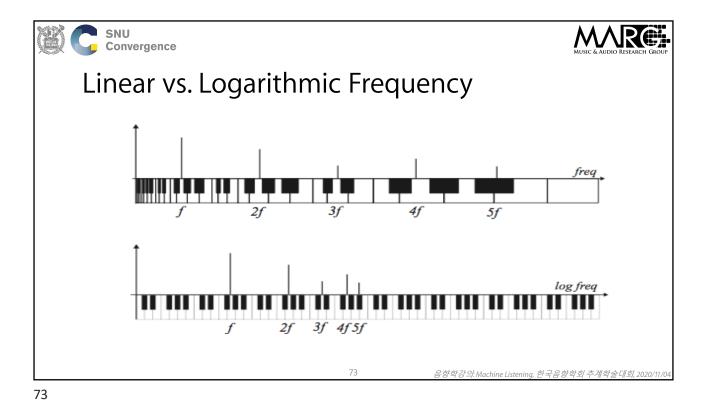
$$f_k = (2^{1/B})^k f_{\min}, \quad k = 0,1,...,N-1,$$

where B is the number of bins in an octave and  $f_{\min}$  is the minimum frequency set by user

 It is called "constant-Q" because Q or "quality factor" is constant along the frequency axis, which is defined as

$$Q = \frac{f_k}{f_w}, \quad k = 0,1,...,N-1,$$

where  $f_w$  is the filter width





Convergence



#### Computation of CQ Transform

 CQ transform can be obtained from the DFT using logarithmically-spaced filterbank

$$X_{cq}(k) = \frac{1}{N(k)} \sum_{n=0}^{N(k)-1} x(n) w(n,k) e^{-j2\pi Qn/N(k)}$$

$$N(k) = f_s Q / f_k$$

■ That uses a variable window length to obtain more resolution at lower frequencies and less at higher (logarithmic distribution of bins in frequency)

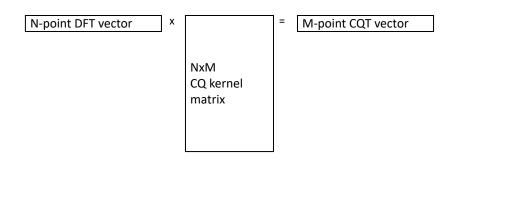


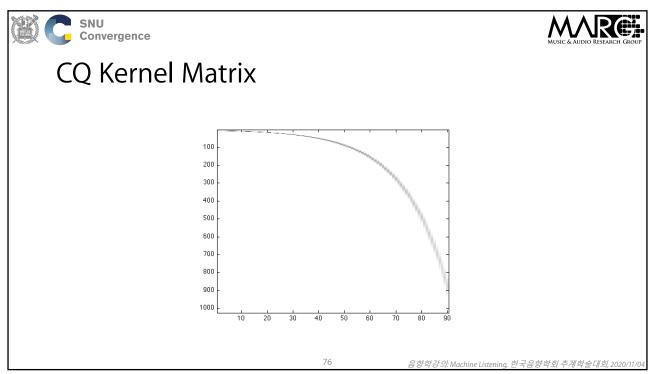


음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

#### Computation of CQT (cont'd)

 CQ transform can be efficiently computed using a CQ kernel which is a 2-d matrix that maps the DFT to the CQT









#### **Chroma**

77

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

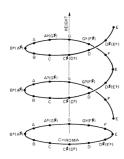
77

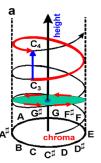




#### Pitch Helix

• The pitch helix is a pitch space where linear pitch is wrapped around a cylinder, thus modeling the special relationship that exists between octave intervals





- Two dimensions
  - Height: naturally organizes absolute pitches from low to high
  - Chroma: represents the inherent circularity of pitch (*relative* relationship between pitch classes)

78

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





#### Chroma (aka Pitch Class Profile)

- Good for describing relative pitch relationship, disregarding absolute pitch height
- Very useful for harmony analysis, key and chord, in particular
- A key and/or a chord can be described as a function of its pitch classes
- Almost universal feature for key/chord estimation applications
- Chroma as audio feature first introduced by Fujishima (1999)

79

유향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

79



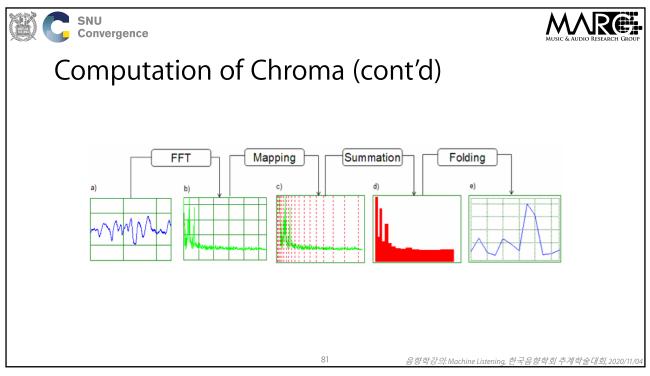


#### Computation of Chroma

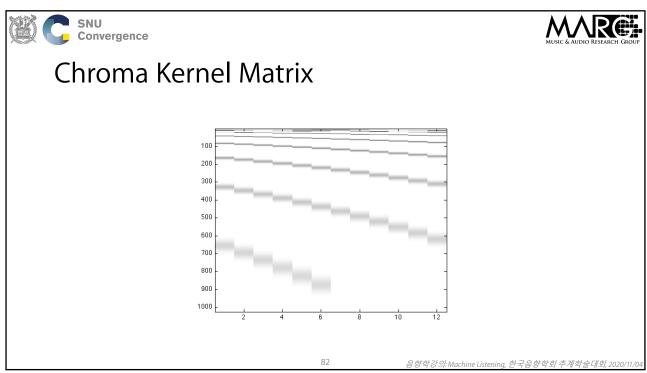
Easily computed from the constant-Q transform by collapsing it to an octave, or

$$Chroma(b) = \sum_{m=0}^{M-1} \left| X_{CQ}(b+mB) \right|,$$

where  $X_{CQ}(k)$  is the CQ transform, M is the total number of octaves of interest, B is the number of chroma bins in an octave, and b=1,2,...,B is the chroma bin index











#### **Machine Learning**

83

우향학강의: Machine Listening, 한국음향학회 주계학술대회, 2020/11/04

83





#### Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998) Well-posed Learning Problem: A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

84

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0-





#### Machine learning algorithms:

- Supervised learning
- Unsupervised learning

Others: Reinforcement learning, recommender systems.

85

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

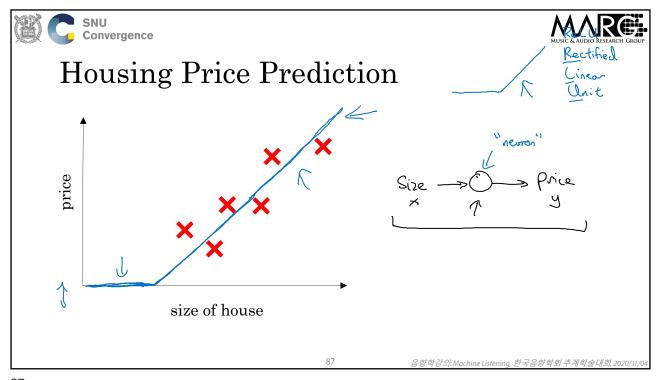
85

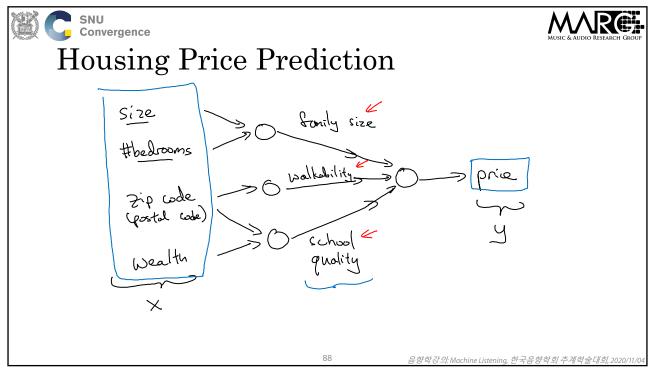


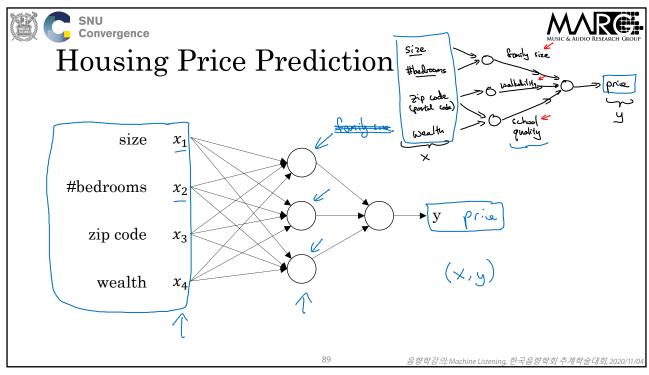
# Introduction to Deep Learning

# What is a Neural Network?

Lecture slides are from: https://www.coursera.org/learn/deep-neural-network



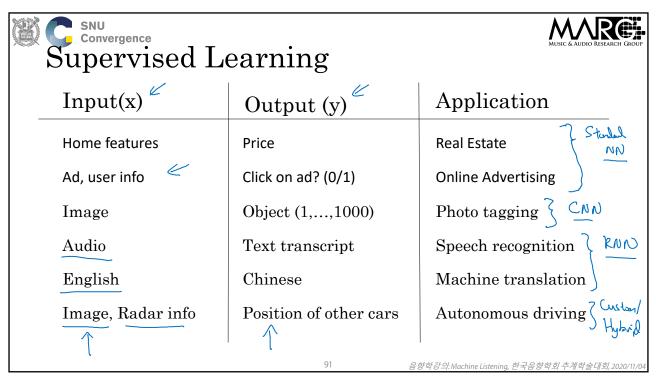


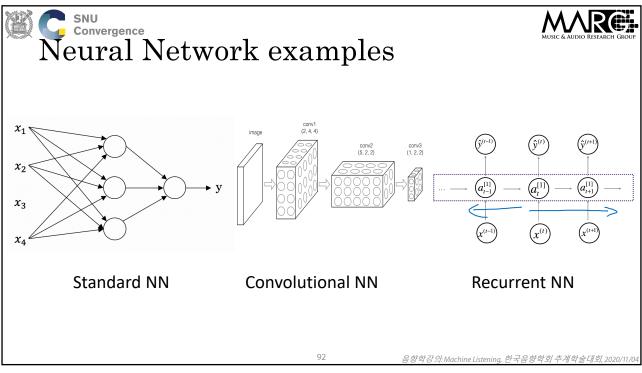




# Introduction to Deep Learning

Supervised Learning with Neural Networks









#### Supervised Learning

#### Structured Data

Size	#bedrooms	 Price (1000\$s)
2104 1600	3	400 330
2400 : 3000	3 : 4	369 : 540

User Age	Ad Id	 Click
41	93242	1
80	93287	0
18	87312	1
:	:	<b>:</b>
27	71244	1

#### **Unstructured Data**





Audio

**Image** 

Four scores and seven years ago...

Text

음향학강의:Machine Listening, 한국음향학회 추계학술대회, 2020/11/0-

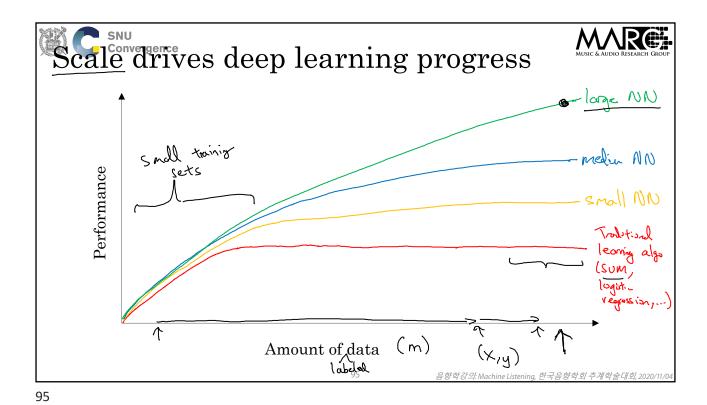
93

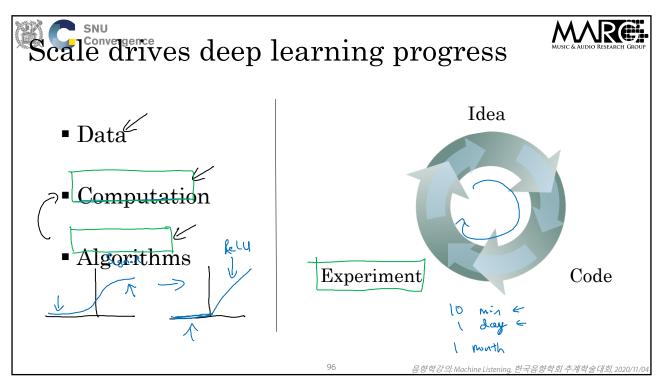
93



## Introduction to Neural Networks

Why is Deep Learning taking off?



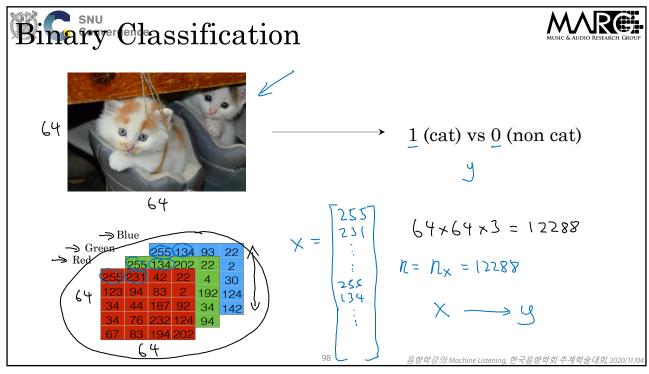




#### Basics of Neural Network Programming

#### **Binary Classification**

deeplearning.ai



Notation:

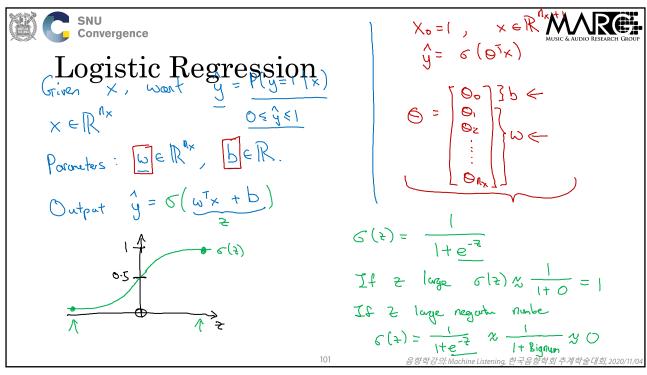
(x,y) 
$$x \in \mathbb{R}^{n_x}$$
,  $y \in \{0,1\}$ 
 $m \in \mathbb{R}^{n_x}$ ,  $y \in \{0,1\}$ 
 $m \in \mathbb{R}^{n_$ 



deeplearning.ai

#### Basics of Neural Network Programming

**Logistic Regression** 

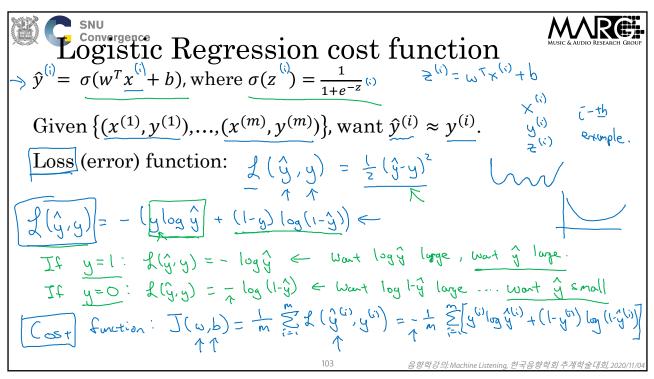




deeplearning.ai

#### Basics of Neural Network Programming

Logistic Regression cost function

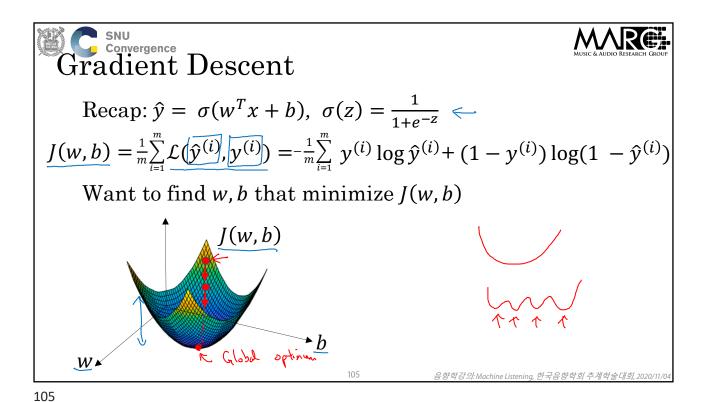


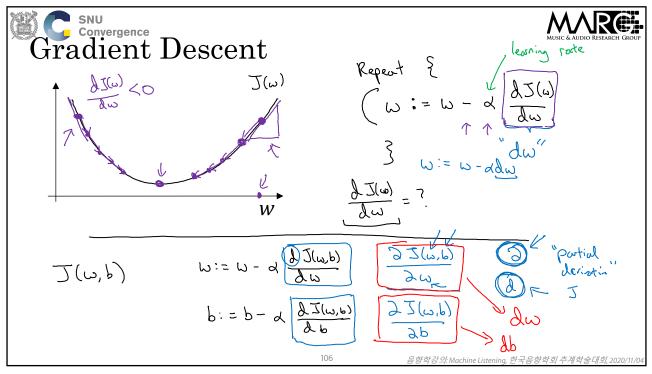


deeplearning.ai

#### Basics of Neural Network Programming

**Gradient Descent** 







deeplearning.ai

#### Basics of Neural Network Programming

#### Logistic Regression Gradient descent

107





#### Logistic regression recap

$$\Rightarrow z = w^T x + b$$

$$\Rightarrow \hat{y} = a = \sigma(z)$$

$$\Rightarrow \mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

$$\frac{\chi_1}{\omega_1}$$

$$\frac{\omega_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

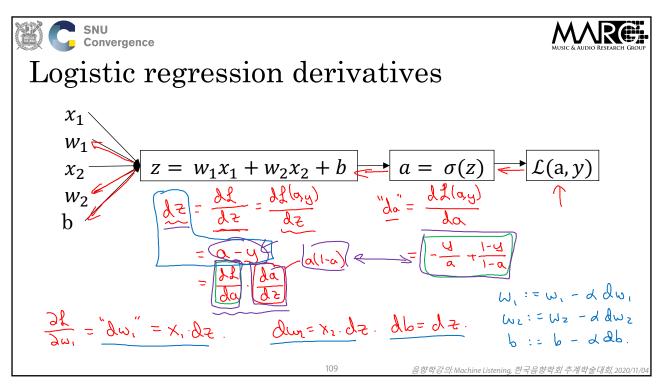
$$\frac{\chi_2}{\chi_2}$$

$$\frac{\chi_1}{\chi_2}$$

$$\frac{\chi_2}{\chi_2}$$

108

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0-



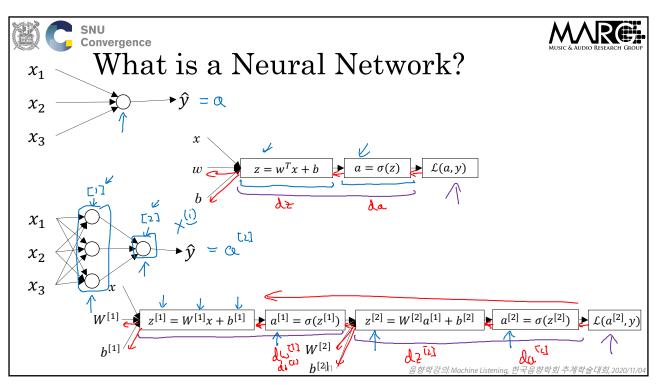


#### deeplearning.ai

#### One hidden layer Neural Network

# Neural Networks Overview

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04



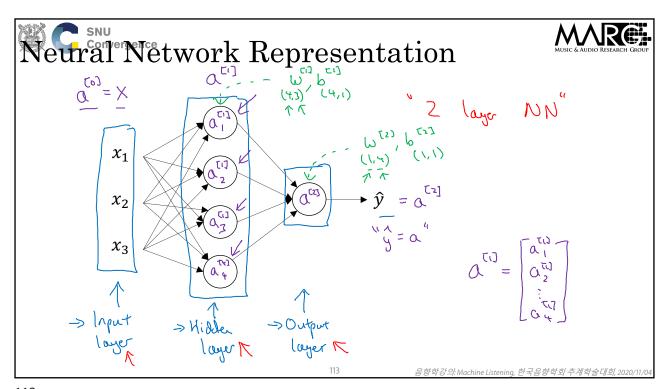


deeplearning.ai

#### One hidden layer Neural Network

#### Neural Network Representation

음향학강의: Machine Listening, 한국음향학회 추계학술대회 2020/11/04





#### deeplearning.ai

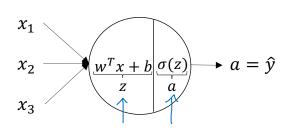
#### One hidden layer Neural Network

#### Computing a Neural Network's Output

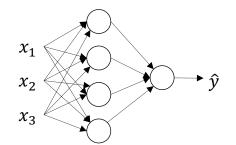
음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

### Neural Network Representation



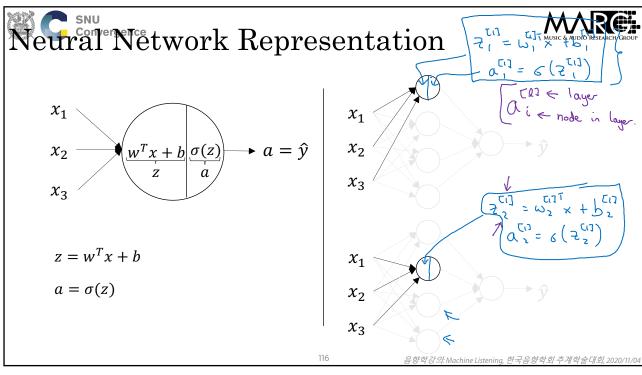


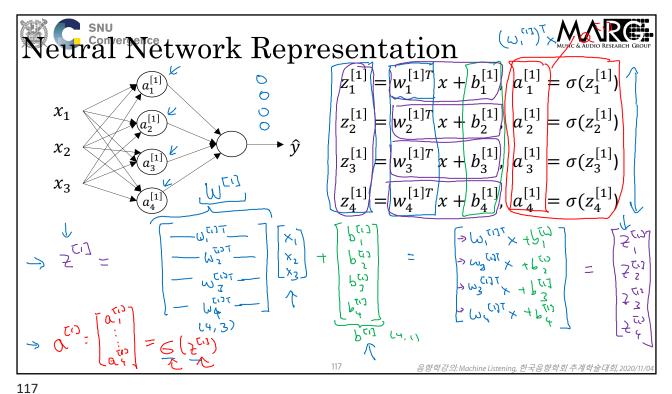
$$z = w^T x + b$$
$$a = \sigma(z)$$



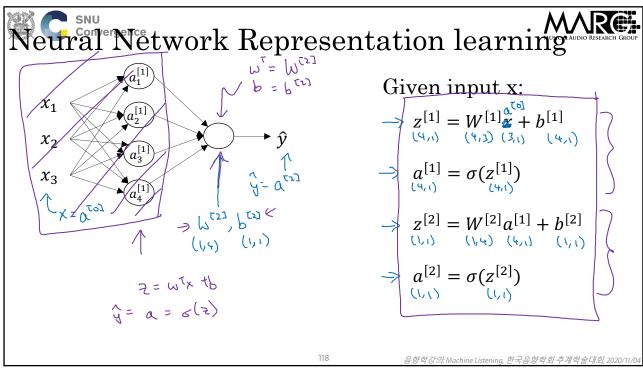
115

우향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04





11/





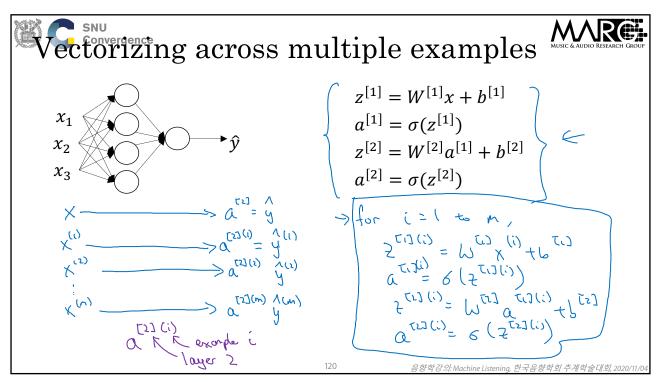
deeplearning.ai

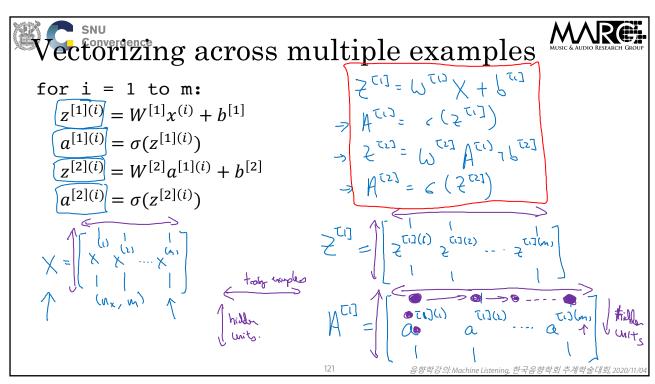
#### One hidden layer Neural Network

# Vectorizing across multiple examples

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

119





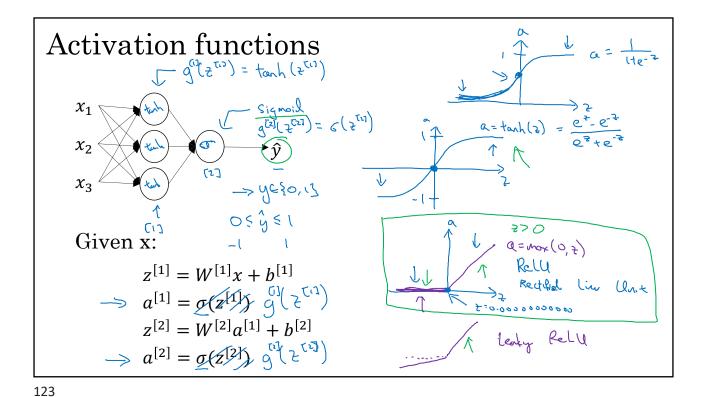


#### One hidden layer Neural Network

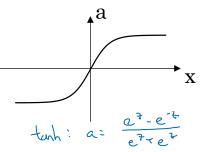
#### **Activation functions**

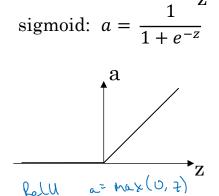
deeplearning.ai

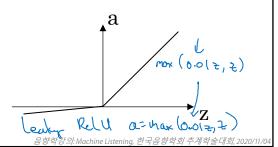
음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04



Pros and cons of activation functions









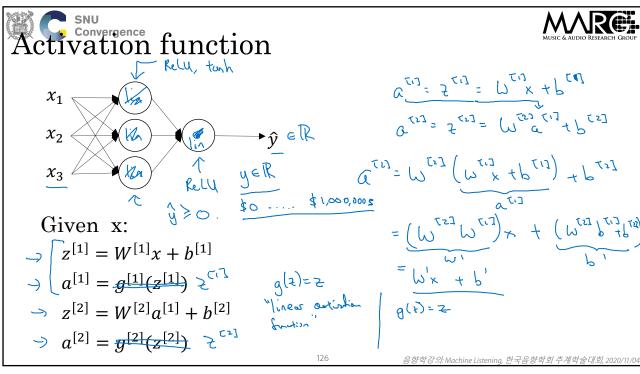
deeplearning.ai

#### One hidden layer Neural Network

# Why do you need non-linear activation functions?

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

125





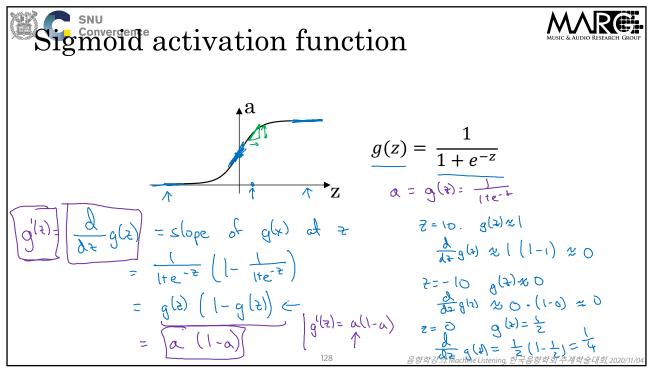
deeplearning.ai

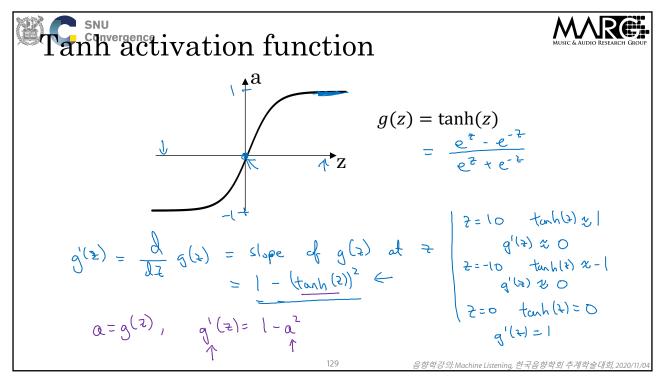
#### One hidden layer Neural Network

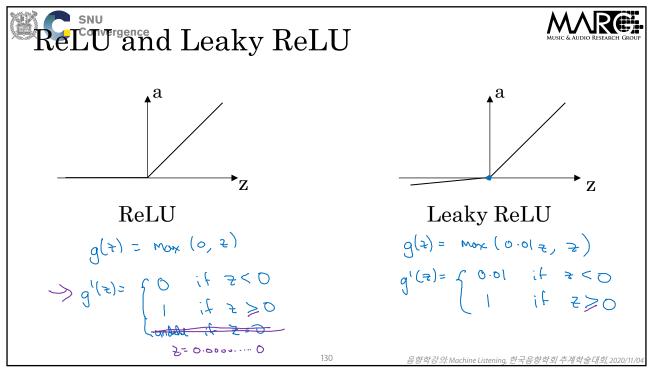
# Derivatives of activation functions

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

127









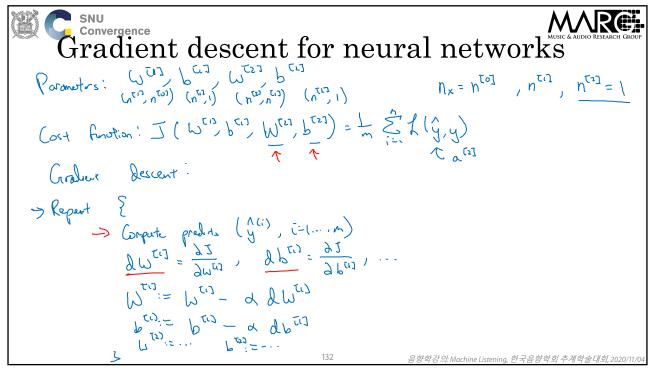
deeplearning.ai

#### One hidden layer Neural Network

# Gradient descent for neural networks

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

13:







$$V_{LSJ} = \partial_{LSJ} (S_{LSJ}) = O(S_{LSJ})$$

$$V_{LSJ} = P_{LSJ} V_{LSJ} + P_{LSJ}$$

$$S_{LSJ} = P_{LSJ} (S_{LSJ}) \leftarrow$$

$$S_{LSJ} = P_{LSJ} (S_{LSJ}) \leftarrow$$

$$S_{LSJ} = P_{LSJ} (S_{LSJ}) \leftarrow$$

Formulas for computing derivatives

Formulas for computing derivatives

$$\begin{cases}
\sum_{i=1}^{C1} = \sum_{i=1}^{C1} X_i + \sum_{i=1}^{C1} X_i
\end{cases}$$

$$\begin{cases}
\sum_{i=1}^{C1} = \sum_{i=1}^{C1} X_i + \sum_{i=1}^{C1} X_i
\end{cases}$$

$$\begin{cases}
\sum_{i=1}^{C1} = \sum_{i=1}^{C1} X_i
\end{cases}$$

$$\begin{cases}
\sum_{i=1}^{C1$$

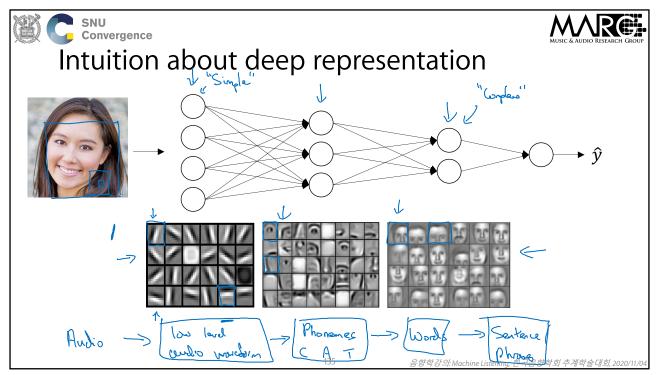


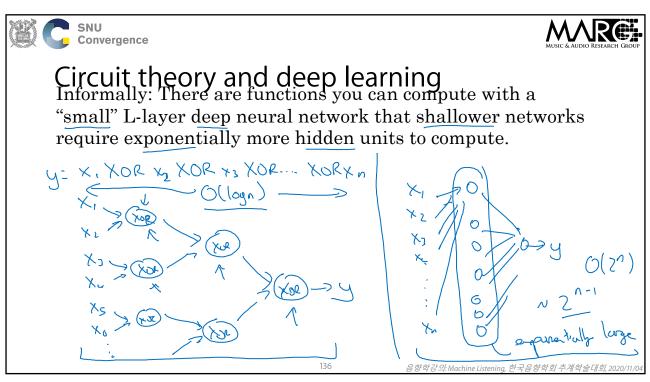
#### deeplearning.ai

#### Deep Neural **Networks**

#### Why deep representations?

음향학강의: Machine Listening, 한국음향학회 추계학술대회,







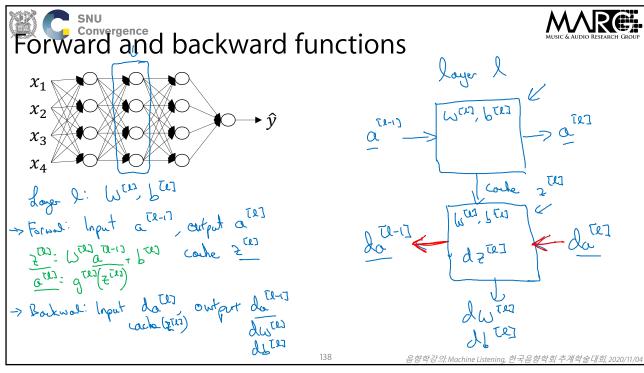
deeplearning.ai

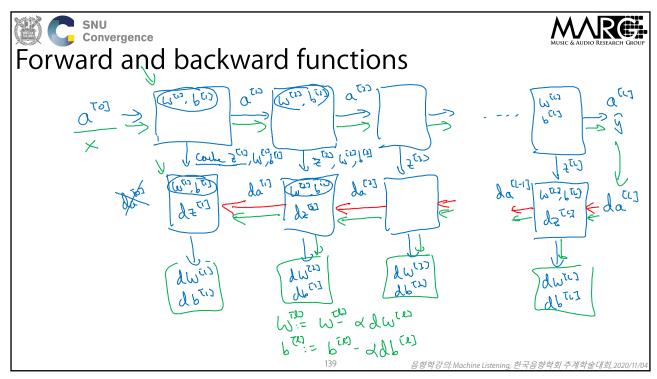
#### Deep Neural Networks

# Building blocks of deep neural networks

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

137



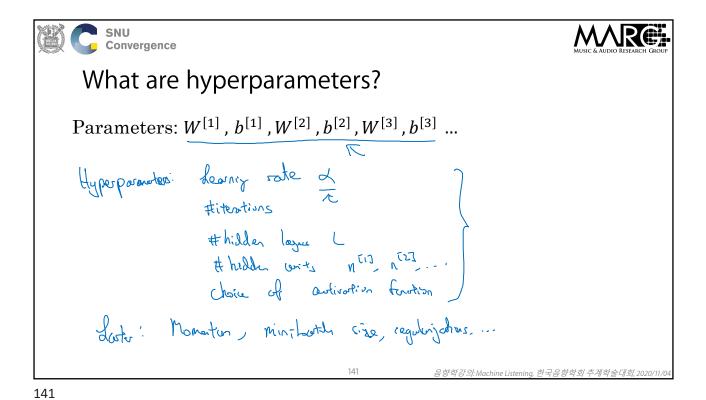


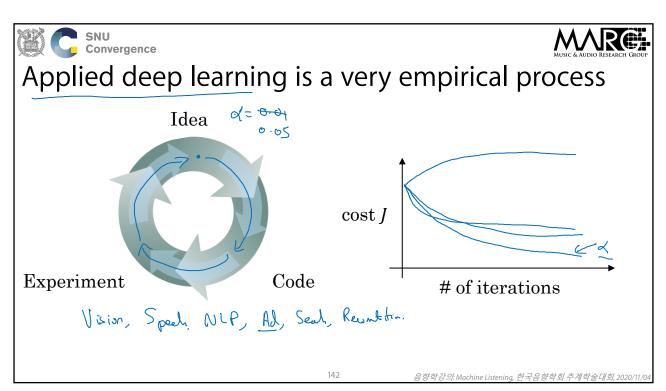


#### Deep Neural Networks

#### Parameters vs Hyperparameters

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04

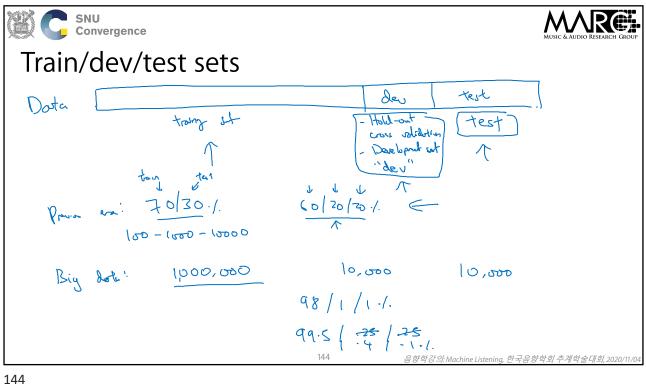






#### Setting up your ML application

#### Train/dev/test sets







#### Mismatched train/test distribution

Training set: Cat pictures from webpages Cat pictures from users using your app

Dev/test sets:

tran / der

tran / der

tran / der

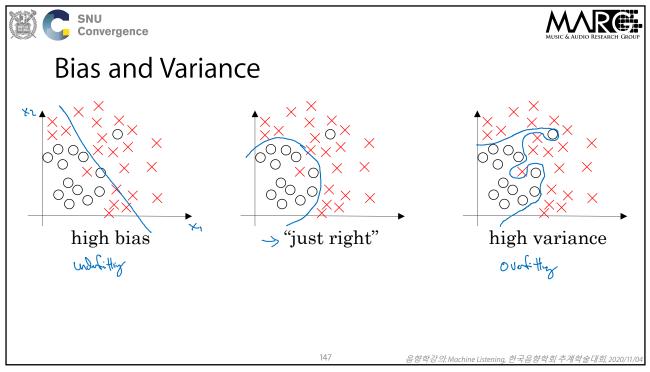
Not having a test set might be okay. (Only dev set.)

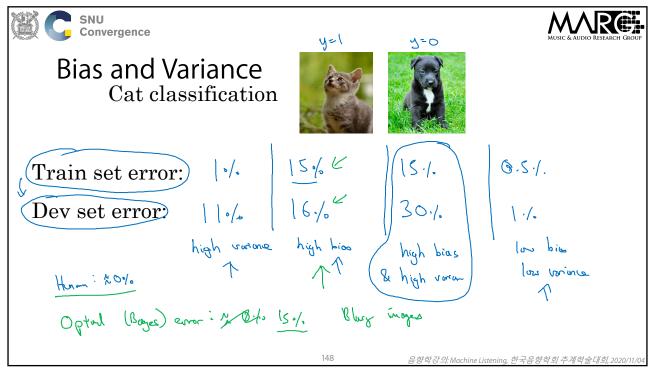
145

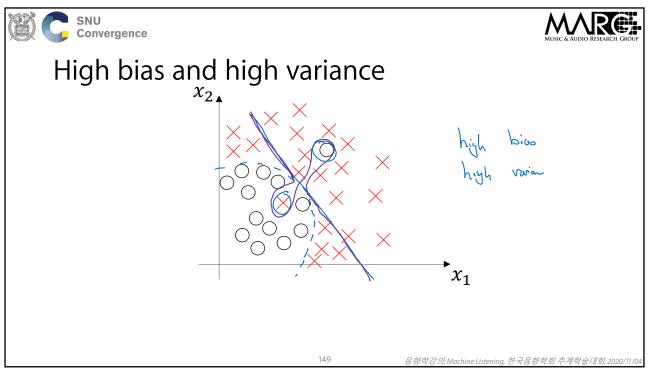


Setting up your ML application

Bias/Variance



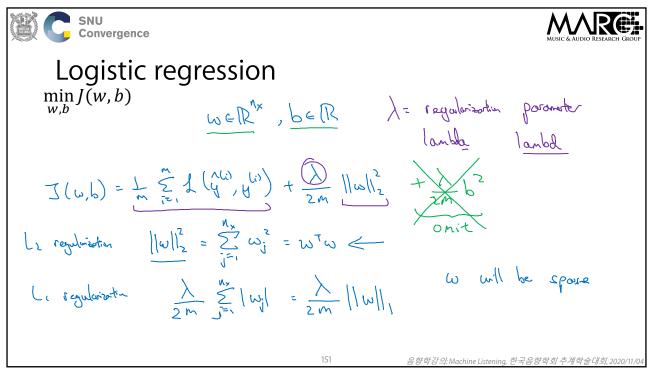






#### Regularizing your neural network

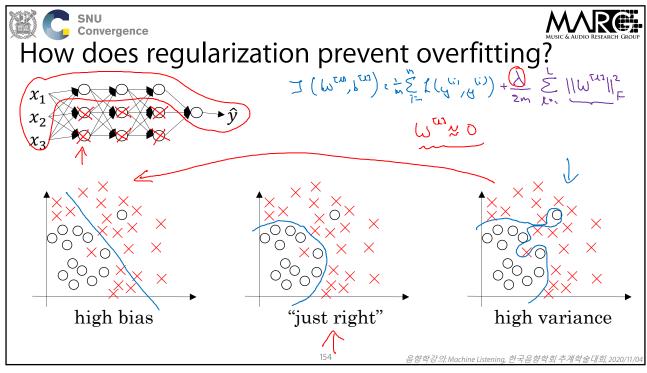
#### Regularization





#### Regularizing your neural network

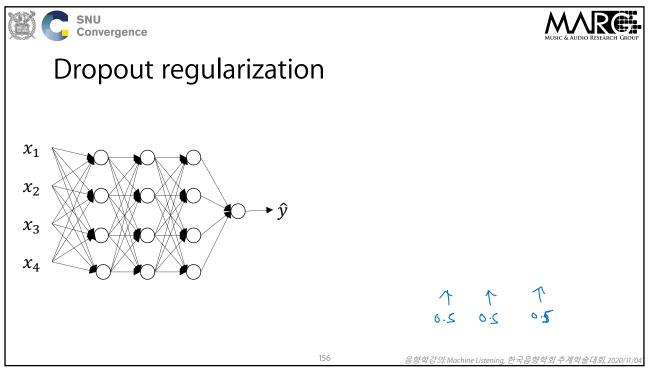
# Why regularization reduces overfitting

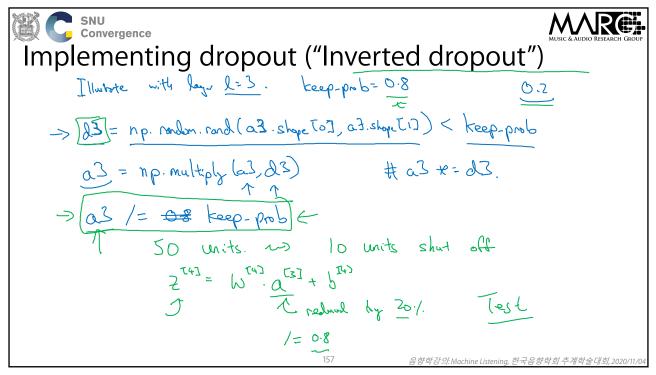


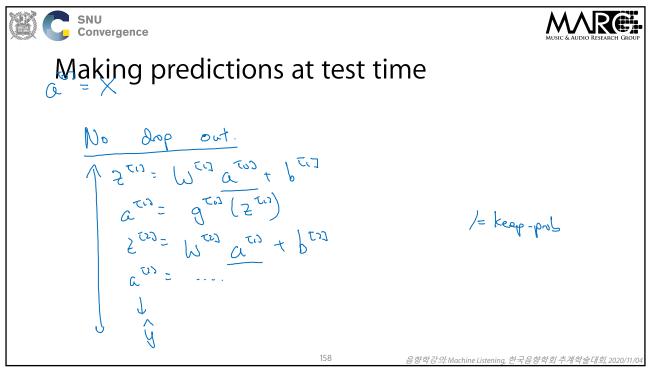


#### Regularizing your neural network

# Dropout regularization



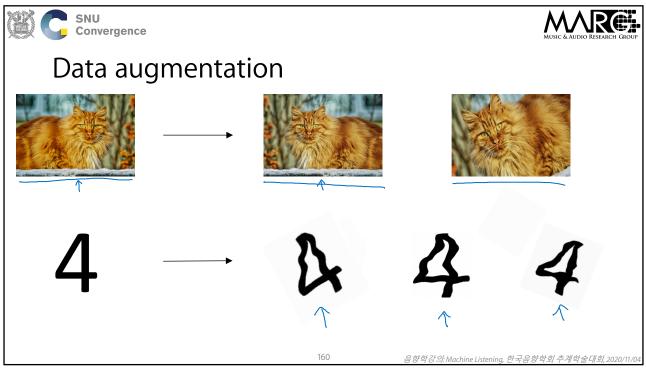


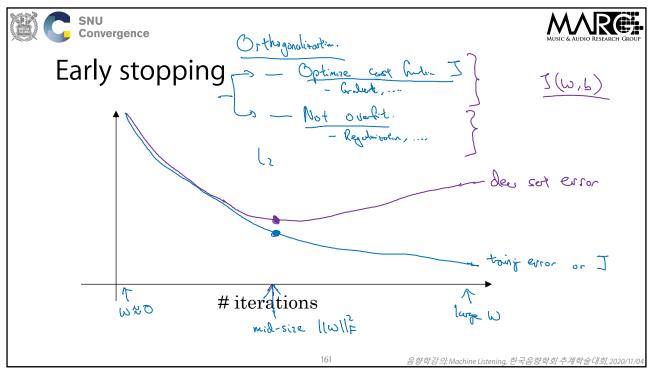




#### Regularizing your neural network

## Other regularization methods

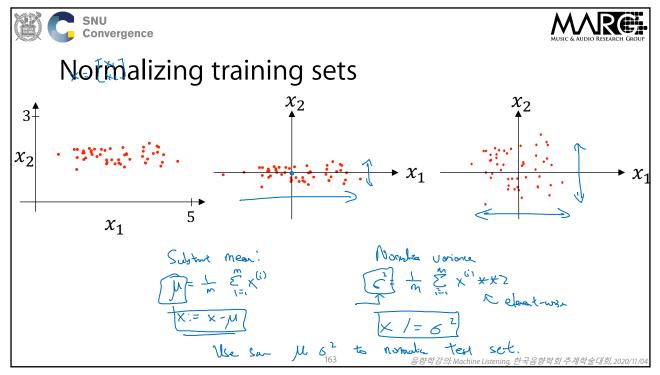


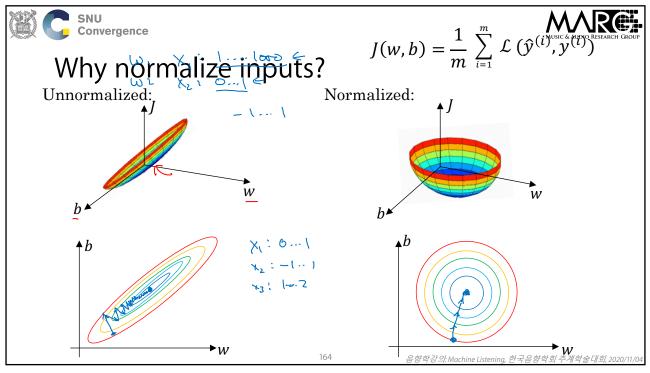




Setting up your optimization problem

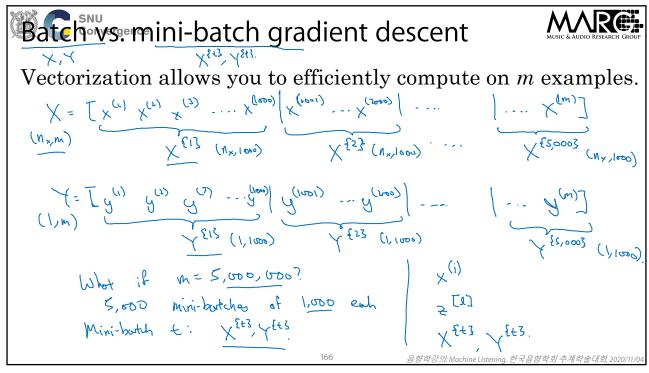
Normalizing inputs

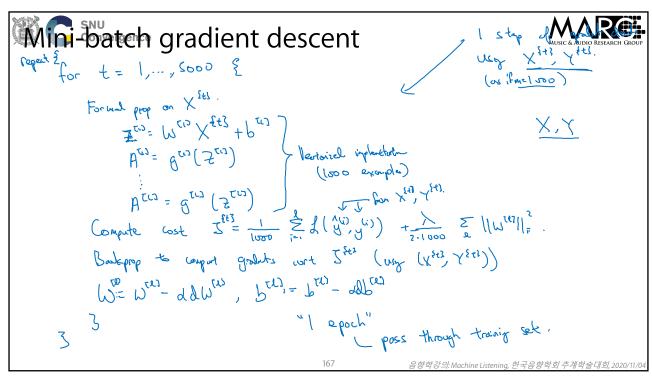






# Mini-batch gradient descent

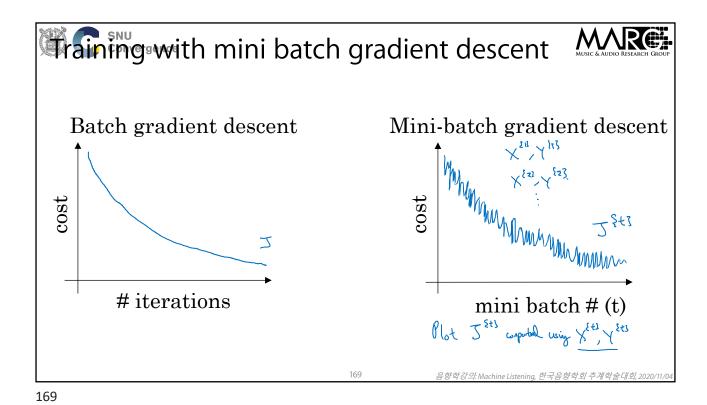


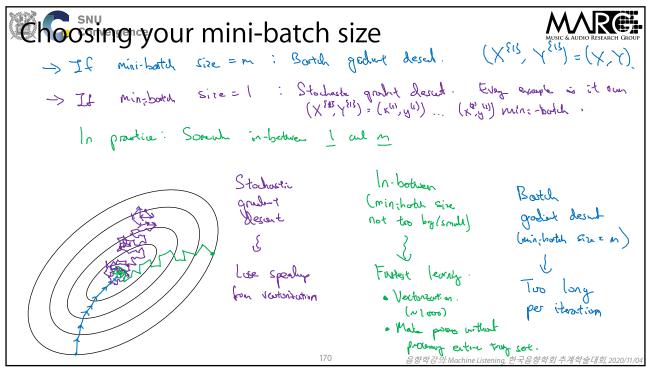


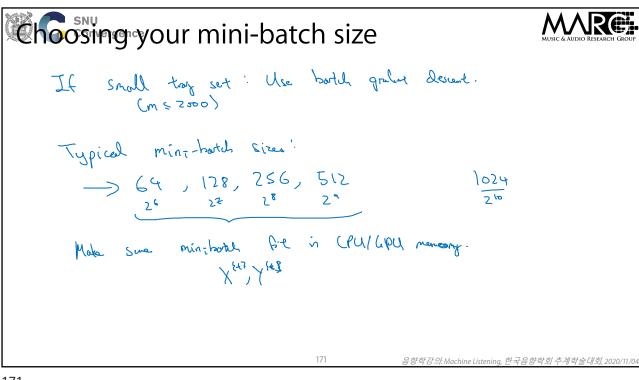


#### Optimization Algorithms

Understanding mini-batch gradient descent



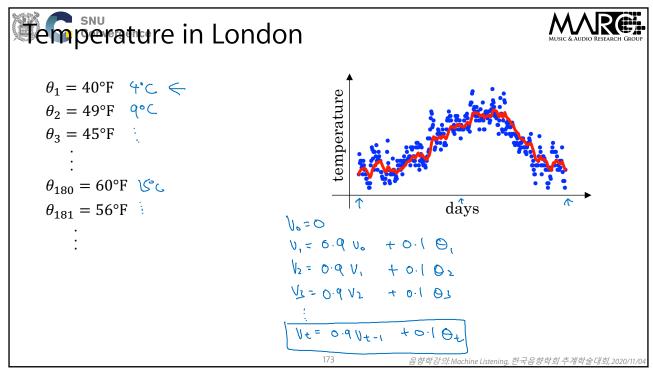


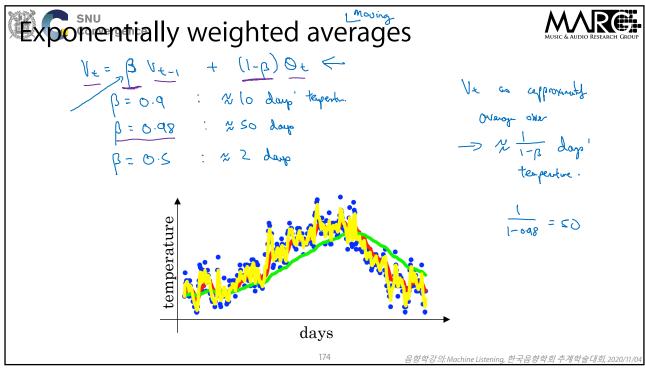




### Optimization Algorithms

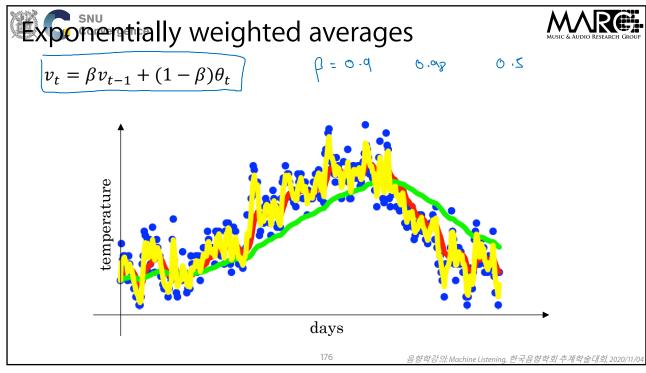
# Exponentially weighted averages

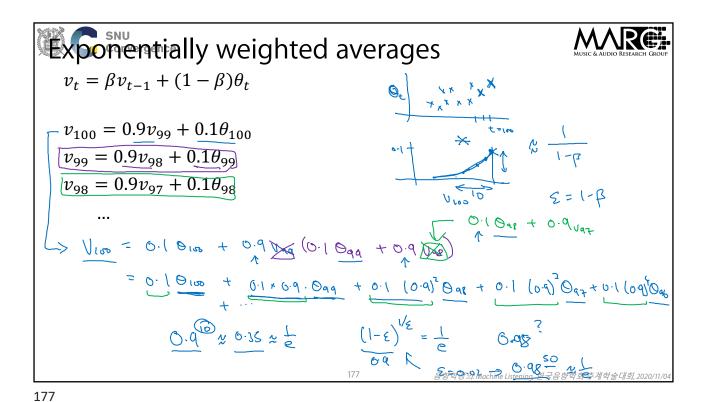


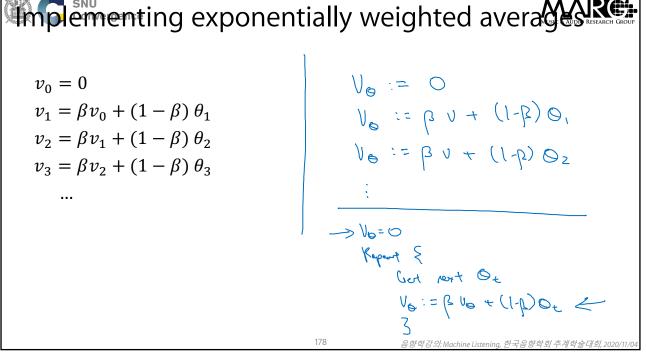




# Understanding exponentially weighted averages

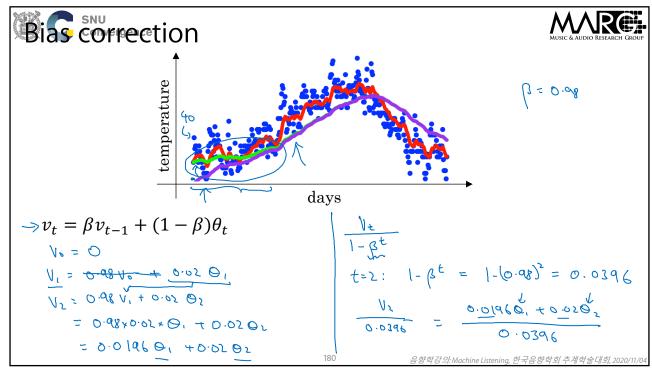






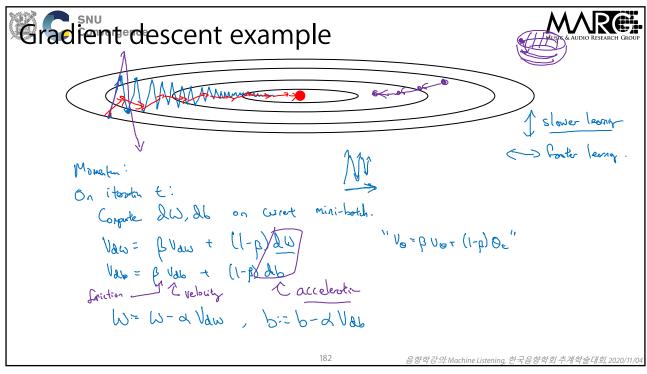


# Bias correction in exponentially weighted average





## Gradient descent with momentum



#### Implementation details



Van= 0, Vab=0

On iteration *t*:

Compute dW, db on the current mini-batch

$$\Rightarrow v_{dW} = \beta v_{dW} + M \beta dW$$

Vaw=BVaw+ dW <

$$\Rightarrow v_{db} = \beta v_{db} + (1 - \beta) \underline{db}$$

$$W = W - \alpha v_{dW}, \ b = b - \alpha v_{db}$$

law / pt

Hyperparameters:  $\alpha, \beta$ 

 $\beta = 0.9$ 

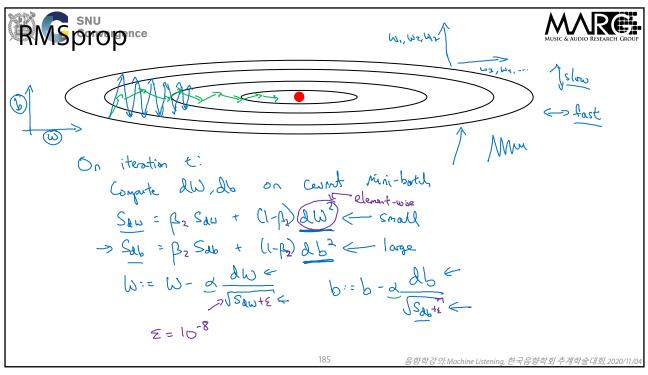
ONE OF OST X O graduty
183 으라라 간 아 Marking listening 하고요하하

183



Optimization Algorithms

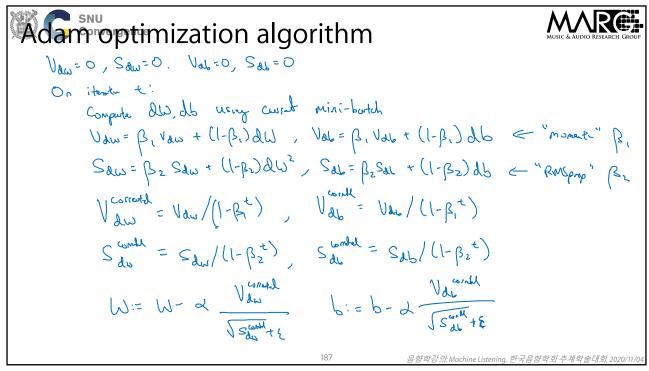
**RMSprop** 

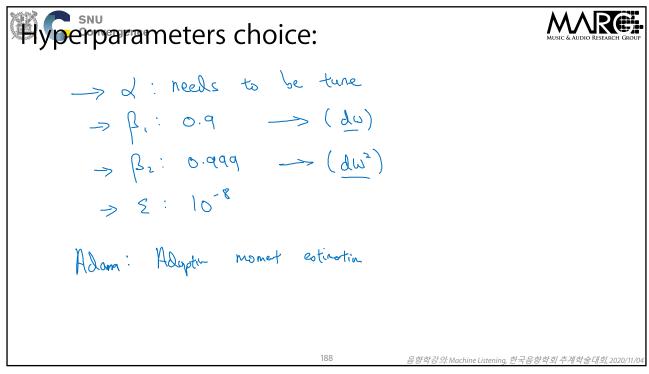




### Optimization Algorithms

# Adam optimization algorithm









#### Applications of Machine Listening

189

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0

189



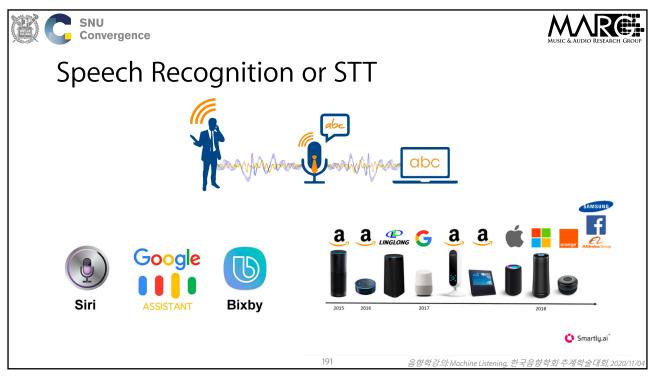


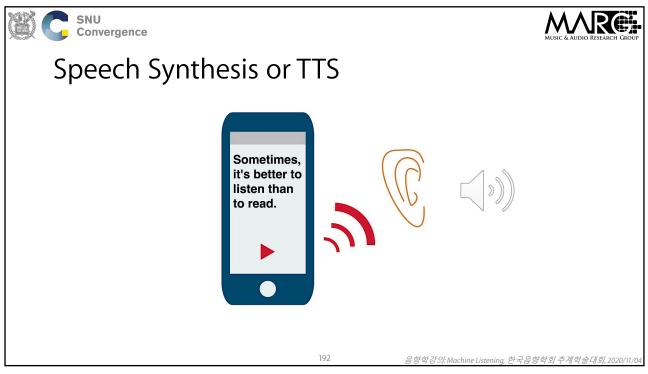
#### Machine Listening Applications

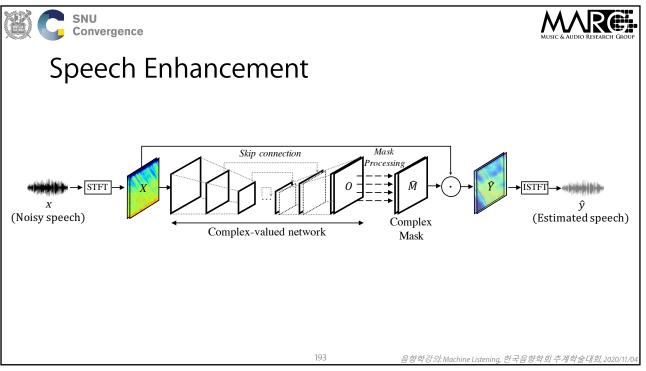
- Speech recognition (speech-to-text or STT)
- Speaker identification/verification
- Emotion recognition
- Speech enhancement
- Source separation
- Automatic music transcription
- Genre/Artist identification
- Music identification
- Lyric-audio alignment
- Speech synthesis (text-to-speech or TTS)
- Singing voice synthesis
- And many more...

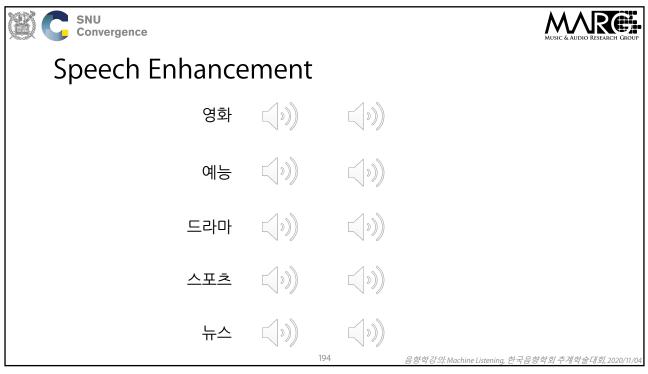
190

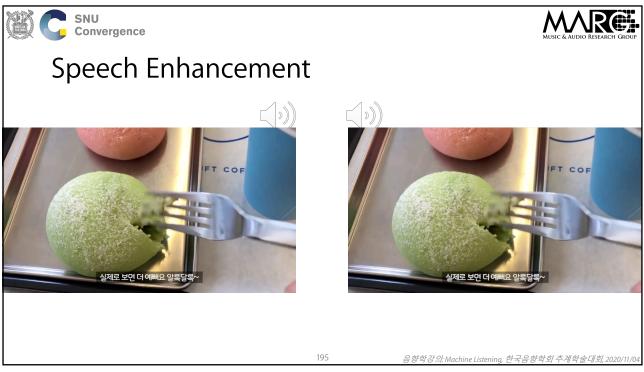
음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/0





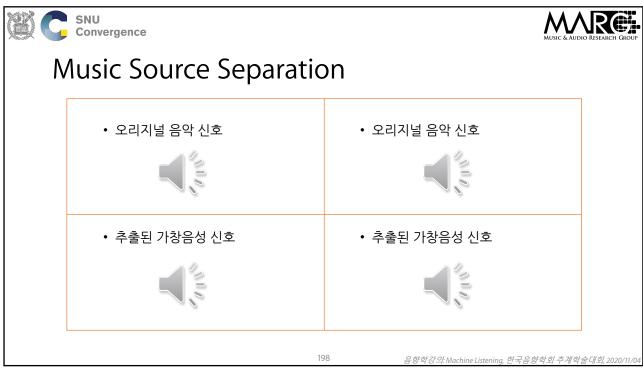


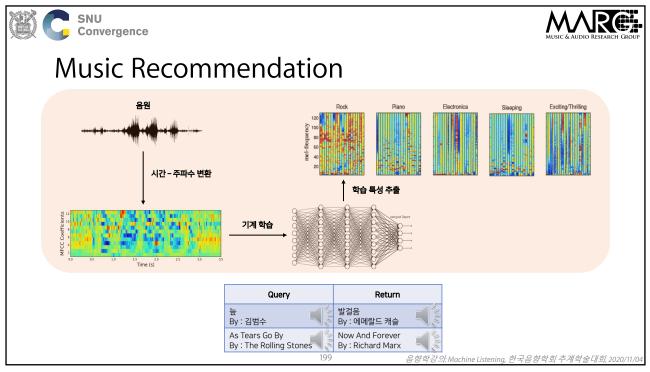


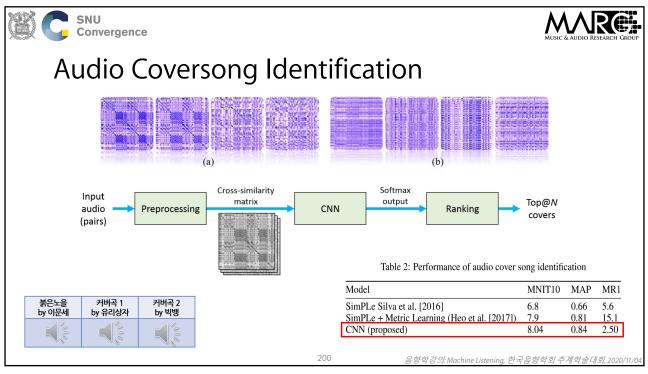


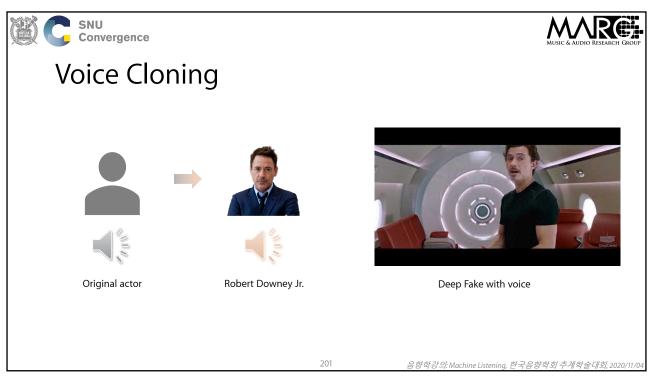




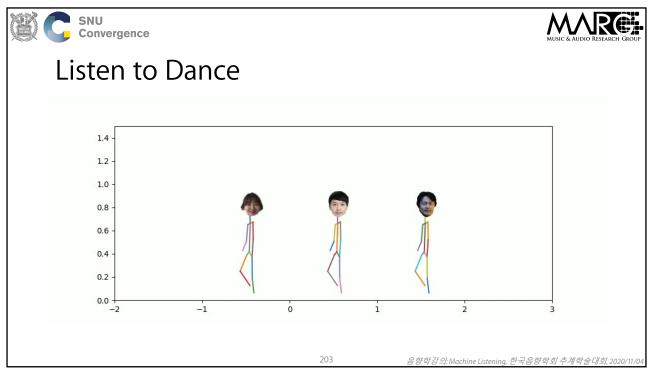


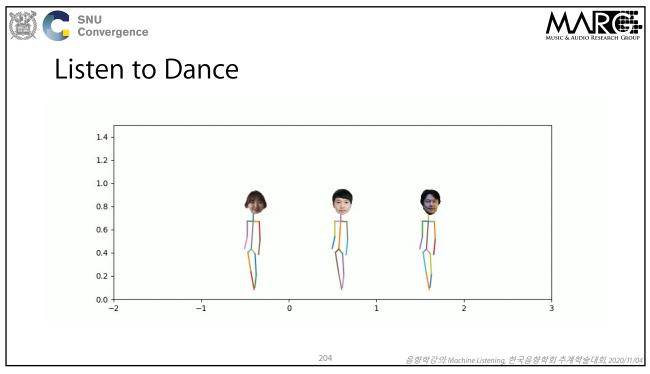
















#### Questions?

서울대학교 지능정보융합학과 음악오디오연구실 이교구 kglee@snu.ac.kr

205

음향학강의: Machine Listening, 한국음향학회 추계학술대회, 2020/11/04