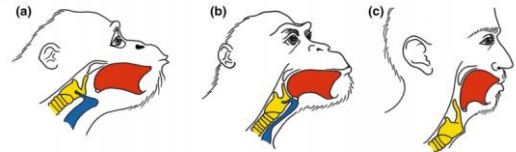
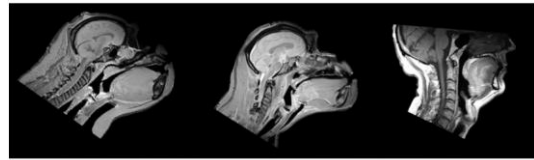


Technologia Mowy

Wykład

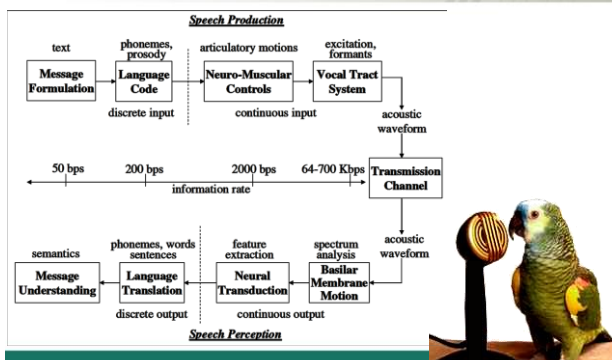
Jakub Gałka, Katedra Elektroniki © 2011-2016
jgałka@agh.edu.pl

Pochodzenie mowy - ewolucja



Trends in Cognitive Sciences

Mowa a informacja



Komunikacja - Informacja

XXI wiek – czy to wiek danych, czy wiek informacji?



Mowa a technologia



Mowa a technologia

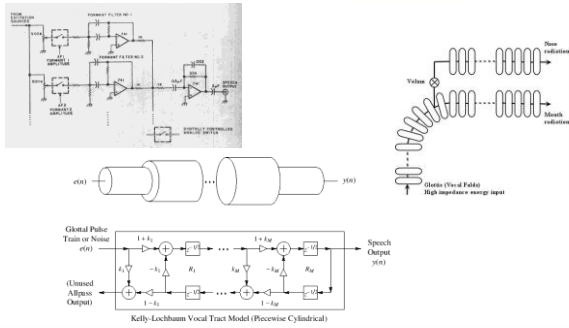
- Tranzystor i Telekomunikacja
- Rozwój techniki i technologii
- Maszyny cyfrowe
- Przetwarzanie sygnałów i w konsekwencji mowy (Shannon, Viterbi, Markow)



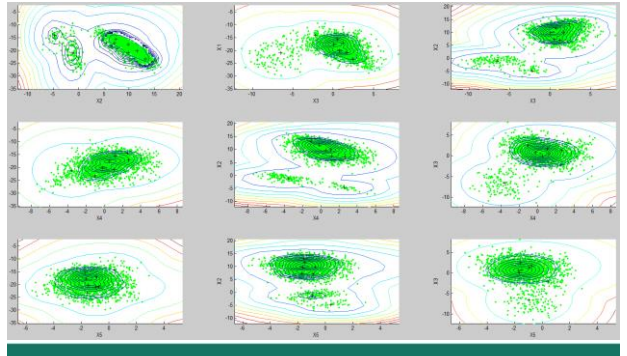
- Telefonia cyfrowa – komórkowa



Modele traktu głosowego (Voder, 1939, <http://www.youtube.com/watch?v=e5gQBei-z-c>)

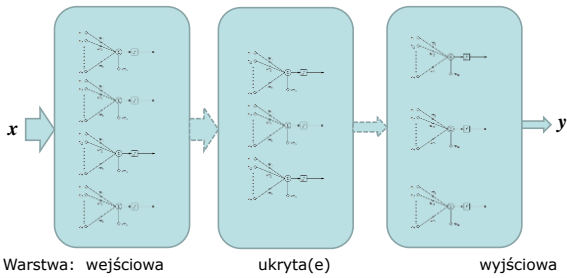


Do niedawna - status-quo DSP i modelowanie statystyczne



Sieci neuronowe

Frank Rosenblatt (1928-1971)



Stabilizacja (stagnacja?) technologii

- Przetwarzanie
- Modelowanie statystyczne
- Konwencjonalny Machine Learning
- Zastosowania domenowe (specjalizowane)



Nowe nadzieje, nowe środowisko społeczno-techniczne

- IoT, BIG DATA
- Crowdsourcing
- WEB 2.0
- Deep Learning

DNN for Speech
10k hours of training data
100 training samples
1000s on a GPU cluster

• Kolejny przełom ?



Kolejny przełom ?

Deep Speech: Scaling up end-to-end speech recognition

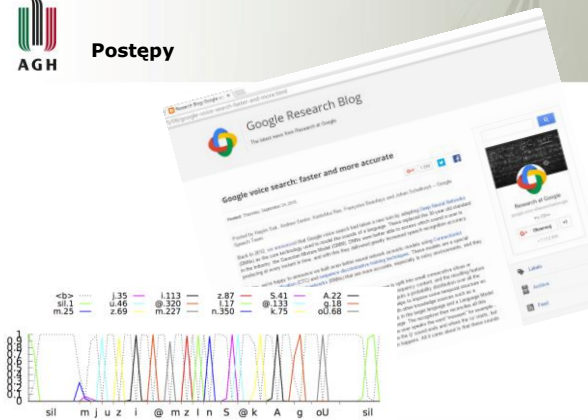
Avri Hamon, Carl Case, Jared Casper, Bryan Catanzano, Greg Diamos, Erich Elsen, Ryan Prager, Sanjeev Satishbh, Shobho Sengupta, Adam Coates, Andrew Y. Ng
Baidu Research - Silicon Valley AI Lab

Abstract
We present a state-of-the-art speech recognition system based on end-to-end learning. Our architecture is significant

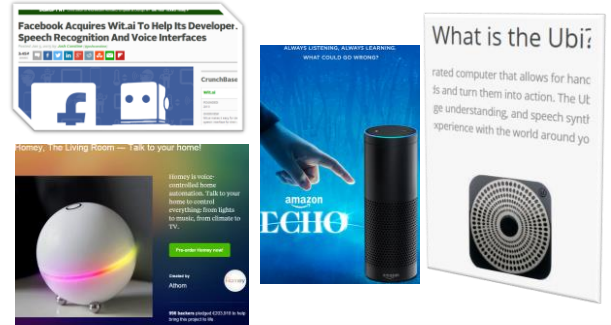
System	Clean (94)	Noisy (82)	Combined (176)
Apple Dictation	14.24	43.76	26.73
Bing Speech	11.73	36.12	22.05
Google API	6.64	30.47	16.72
wit.ai	7.94	35.06	19.41
Deep Speech	6.56	19.06	11.85



Postępy



Sukcesy, klapy, Startup'y, korporacje



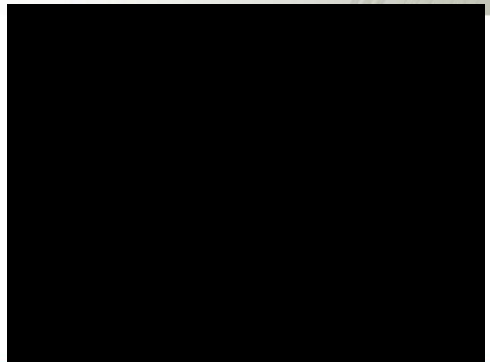
Podstawowe oznaczenia

- x – zmienna skalarna
- \mathbf{x} – zmienna wektorowa
- X – stała, macierz
- \mathbf{X} – zbiór
- $x_i, x(i), X(i,j)$ – element wektora, macierzy
- α, β, \dots – parametr skalarny, wektorowy
- $y=f(x), \mathbf{z}=\mathbf{g}(\mathbf{x})$ – funkcje
- \mathbf{x}^T, X^T – transpozycja
- itd.



Dźwięki i ich postrzeganie

http://www.ted.com/talks/julian_treasure_the_4_ways_sound_affects_us.html

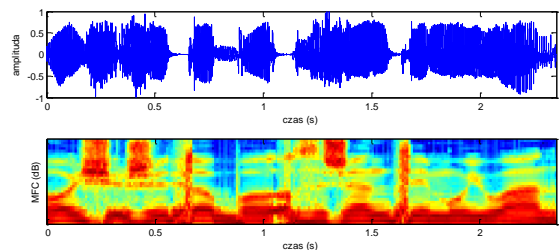


Mowa

- Szybkie i precyzyjne przekazywanie informacji – rozwój gatunku
- Proces złożony / wielopoziomowy
- Ewolucja mowy – optymalizacja <http://uvafon.hum.uva.nl/bart/papers/deBoerEncyclopedia2006.pdf>
- Pragmatyzm (środowisko)
- Anatomia i fizjologia
- Komunikaty, informacje, pojęcia, abstrakty, myślenie
- Wspólna wiedza

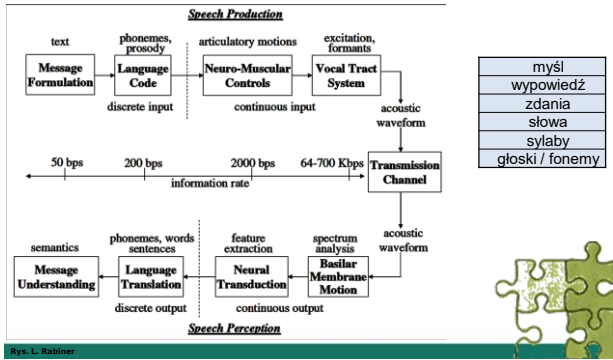


Czym jest mowa

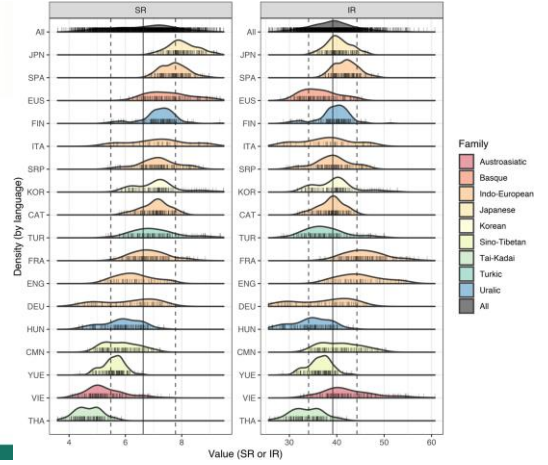




Zjawisko mowy – proces komunikacji



Rys. L. Rabiner

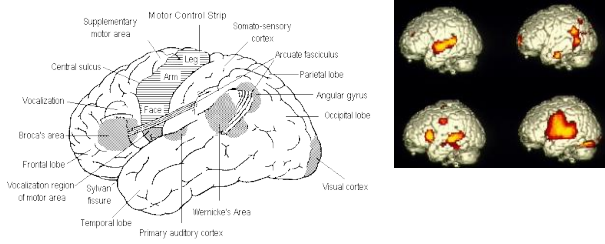


Rys. 1



Ośrodki mowy w mózgu

- Obszary Broca i Wernicke'go

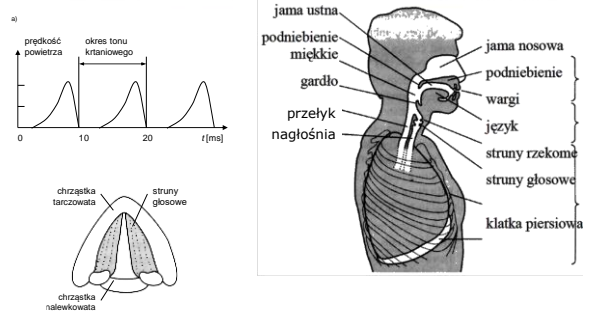


- Więcej ciekawostek: www.talkingbrains.org

Rys.



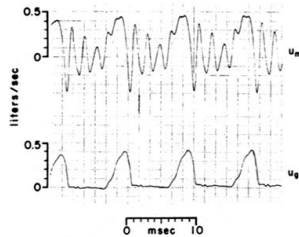
Aparat mowy



Rys. 1 – Wikimedia Commons



Aparat mowy – struny głosowe



Głos	Rodzaj	[Hz]	Od	Do
bas	męski	80	320	
baryton	męski	100	400	
tenor	męski	120	480	
alt	żeński	160	640	
mezzosopran	żeński	200	800	
sopran	żeński	240	960	

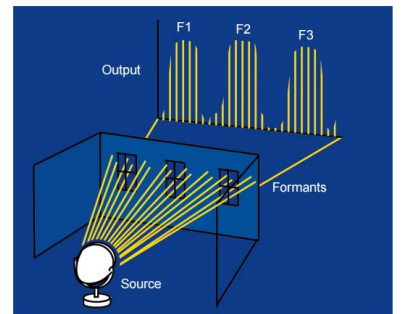
Struny głosowe (endoskop): http://www.youtube.com/watch?v=Z_ZGq1tZn8

Rys. 1 – Wikimedia Commons



Dziwięk mowy

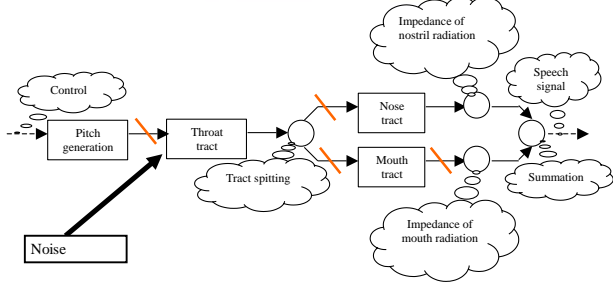
- Analiza widma



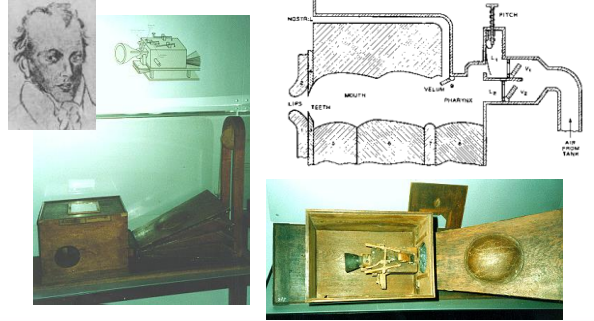
trends in Cognitive Sciences



Trakt głosowy



Modele traktu głosowego Von Kempelen (1791), Riesz (1937)



Modele traktu głosowego

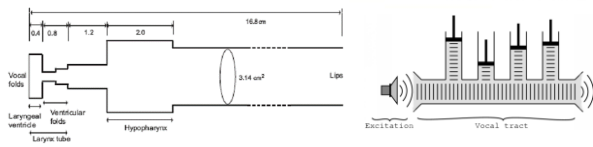
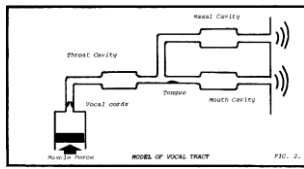
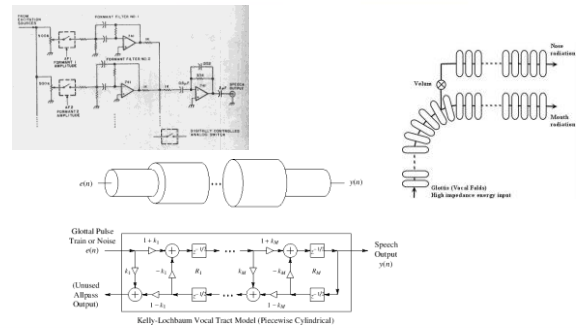


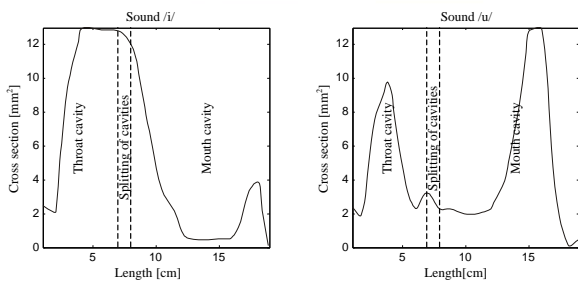
Figure 3: Sketch of the vocal tract model



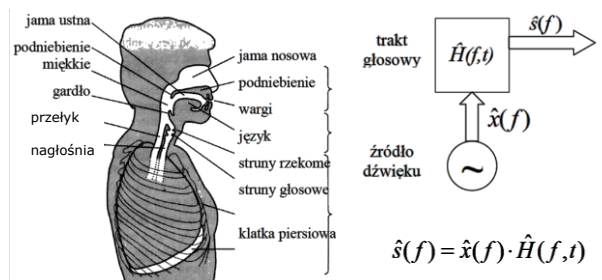
Modele traktu głosowego (Voder, 1939, <http://www.youtube.com/watch?v=e5gQBel-z-c>)



Profil traktu głosowego

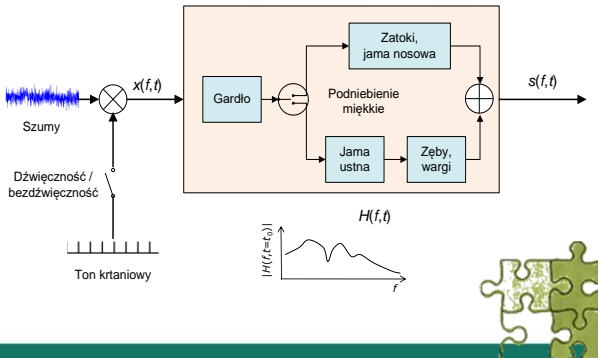


Model sygnałowy źródło-filtr



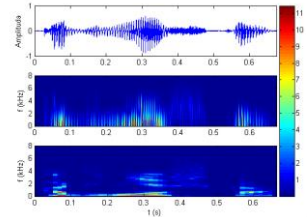
Rys. 8. Tadeuszewicz, Sygnal Nowy

Model źródło-filtr, c. d.

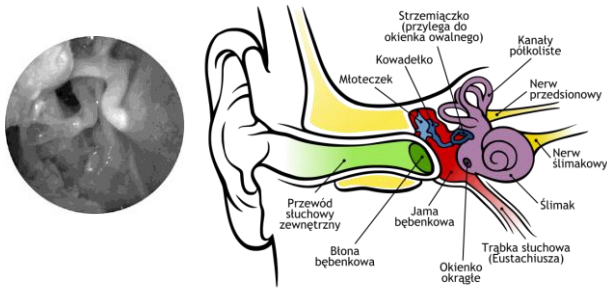


Akustyka mowy

- Dopasowany system komunikacji
- Fizyka, teoria informacji,...
- Zakres częstotliwości:
 - PCM 4000Hz
 - Std. 8000Hz
 - CD 22050Hz
- Dynamika
 - min. 3, 4 bity
 - $6.5 \cdot 3 \sim 20$ dB



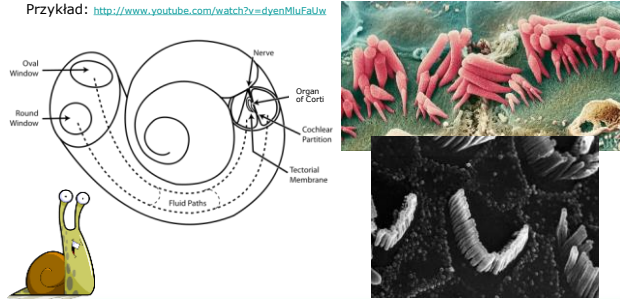
Narząd słuchu



Rys. Wikimedia Commons, Chitta L. Brockmann

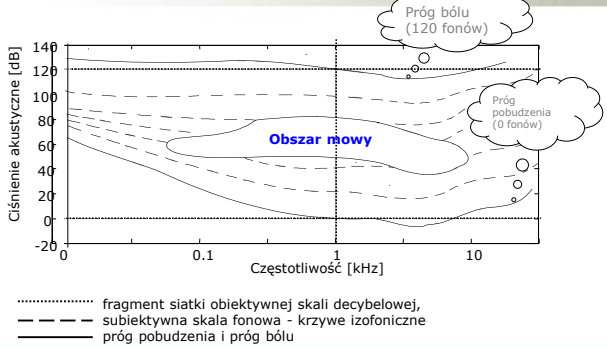
Ślimak, narząd Cortiego

Aparat słuchowy człowieka: <http://www.youtube.com/watch?v=7Q-adw-HyrQ>
 Narząd Cortiego - zasada działania: <http://www.youtube.com/watch?v=xMUJ5CCoW6Y> ()
 Przykład: <http://www.youtube.com/watch?v=dyenMfuFalw>



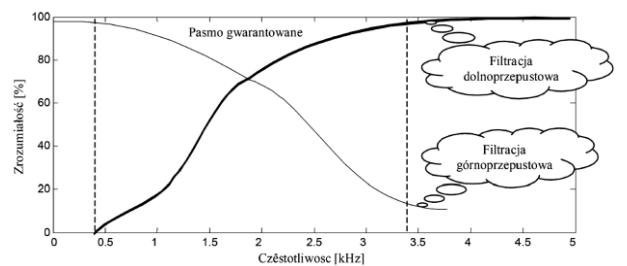
Rys. Wikimedia Commons, www.sciencephoto.com

Krzywe izofoniczne Fletchera-Munsona obszar mowy



Rys. M. Kepiński

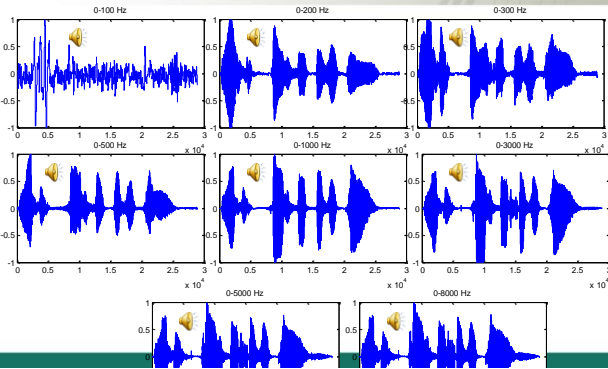
Pasmo a zrozumiałość



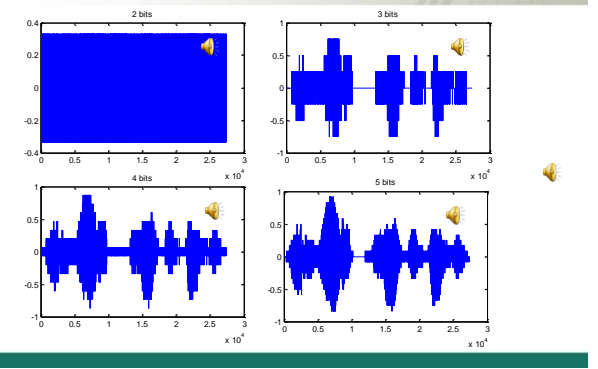
Rys. M. Kepiński



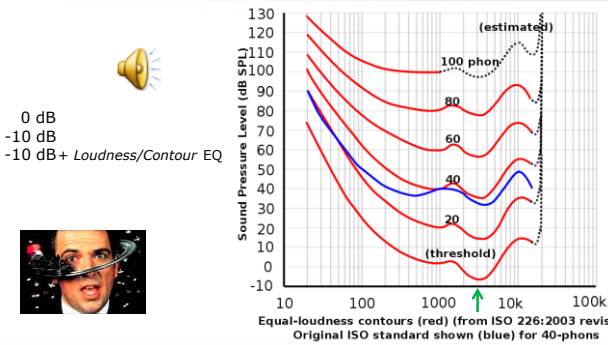
Pasma sygnału - przykłady



Liczba bitów - przykłady



Krzywe izofoniczne Fletchera-Munsona ISO 226:2003



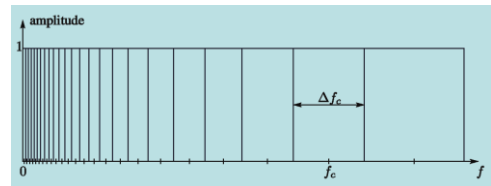
rys. Lindlöand, Wikimedia Commons



Rozdzielczość częstotliwościowa



- Jednostka JND (ang. *Just Noticeable Difference*)
 $2 \cdot \Delta f \approx f_{ref, Hz} \cdot 0,7\%$
- Pasma krytyczne (Barkhausen) – skala Bark (Zwicker, 1961)
 $\Delta f_{Bark}(f_{Hz}) = 25 + 75 \left(1 + 1,4 \left(\frac{f_{Hz}}{1000} \right)^{0,69} \right)$ $f_{Bark}(f_{Hz}) = 13 \arctan \left(\frac{0,76 f_{Hz}}{1000} \right) + 3,5 \arctan \left(\frac{f_{Hz}^2}{7500^2} \right)$



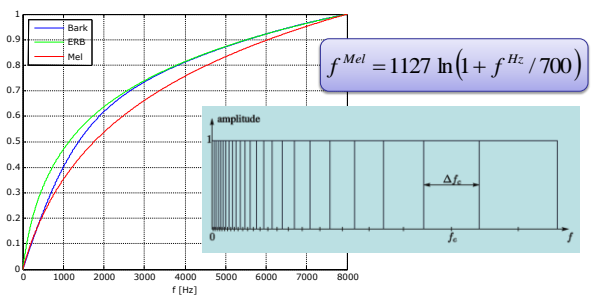
Rozdzielczość częstotliwościowa



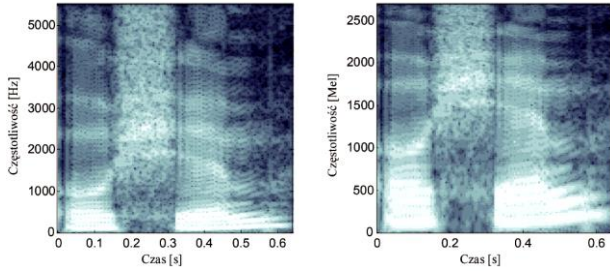
- Jednostka JND (ang. *Just Noticeable Difference*)
 $2 \cdot \Delta f \approx f_{ref, Hz} \cdot 0,7\%$
- Pasma krytyczne (Barkhausen) – skala Bark (Zwicker, 1961)
 $\Delta f_{Bark}(f_{Hz}) = 25 + 75 \left(1 + 1,4 \left(\frac{f_{Hz}}{1000} \right)^{0,69} \right)$ $f_{Bark}(f_{Hz}) = 13 \arctan \left(\frac{0,76 f_{Hz}}{1000} \right) + 3,5 \arctan \left(\frac{f_{Hz}^2}{7500^2} \right)$
- ERB (ang. *Equivalent Rectangular Bandwidth*)
 $ERB(f_{Hz}) = 21,4 \cdot \log_{10} \left(\frac{4,37 f_{Hz} + 1}{1000} \right)$
- Mel (Stevens, Volkman, Newman, 1937)
 $f^{Mel} = 1127,01048 \ln(1 + f / 700)$
- Więcej o transformacjach: <https://ccrma.stanford.edu/~jos/bbt/bbt.pdf>
- A perceptual space that can explain the robustness of bio-acoustic communication (53 min.):
<http://video.google.com/videoplay?docid=1131332207960335088>



Melowa psychoakustyczna skala częstotliwości



Spektrogram w skali melowej



Rys. M. Kepinski

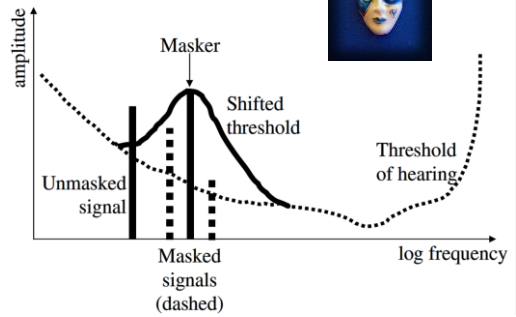
Rozdzielczość częstotliwościowa przykłady ($JND=f*0.7\%$)

f			
200 Hz			
500 Hz			
2000 Hz			

Rozdzielczość częstotliwościowa przykłady ($JND=f*0.7\%$)

f	df	0.5 JND	1 JND	2 JND
200 Hz				
500 Hz				
2000 Hz				

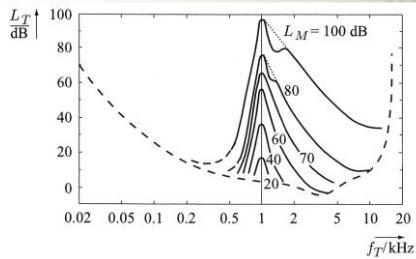
Maskowanie



66

Rys. L. Rabiner, Introduction to Digital Speech Processing, NOW 2007

Maskowanie tonów



Level L_T of a sinusoidal test tone of frequency f_T masked by
 — a sinusoidal masker with $f_M = 1$ kHz
 narrowband noise with level L_M and with critical bandwidth, centered at $f_M = f_c = 1$ kHz ([Zwicker, Fastl 1999])

Rys. P. Vary, R. Martin, Digital Speech Transmission, Wiley 2005

Maskowanie w czasie

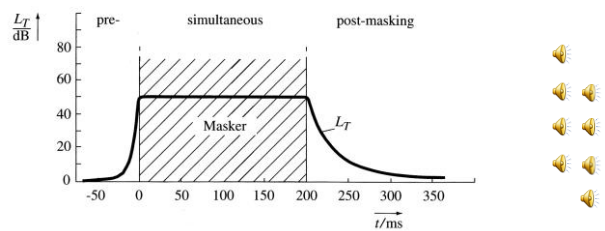


Figure 2.22: Pre- and post-masking: necessary level L_T for the audibility of a sinusoidal burst (test tone) masked by wideband noise (adapted from [Zwicker, Fastl 1999])

Rys. P. Vary, R. Martin, Digital Speech Transmission, Wiley 2005



Maskowanie w czasie - przykłady

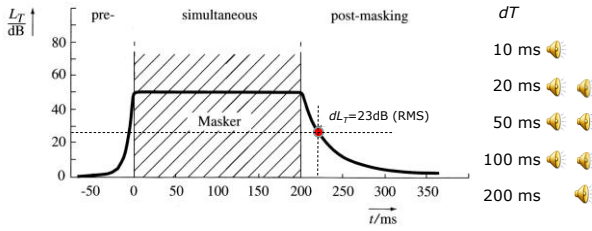
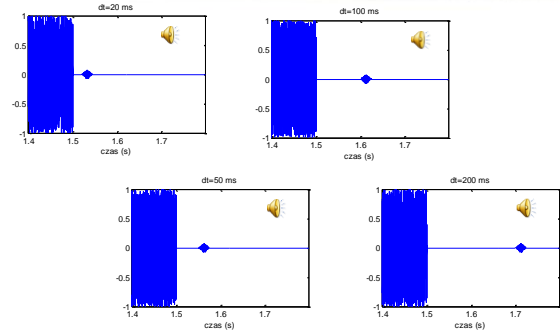


Figure 2.22: Pre- and post-masking; necessary level L_T for the audibility of a sinusoidal burst (test tone) masked by wideband noise (adapted from [Zwicker, Fastl 1999])

Rys. P. Vary, B. Martin, Digital Speech Transmission, Wiley 2005

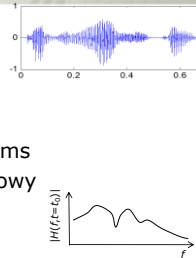


Maskowanie szumem – przykłady $dL_T=23\text{dB (RMS)}$

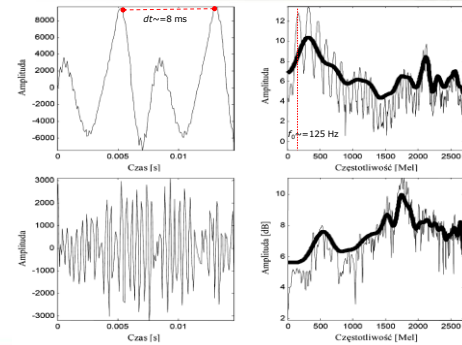


Sygnal mowy

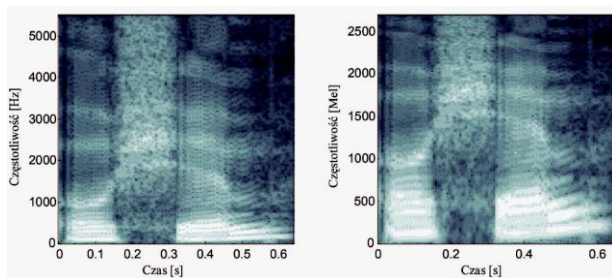
- Sygnal czasowo-częstotliwościowy
- Zmienna moc i widmo chwilowe
- Zakres (*) częstotliwości $50 < f < 8\text{kHz}$
- Niestacjonarny globalnie
- *Quasi-stacjonarny* w porcjach 10-20ms
- Dźwięczny – bezdźwięczny – impulsowy
- Częstotliwość podstawowa – f_0
- Formanty: $f_1, f_2, f_3, (f_4)$
- Częstotliwościowa zmienność własna sygnału (np. pozycje formantów) odpowiada zdolności rozdzielczej słuchu (Mel) !!!



Widmo - ton krtaniowy i formanty */i/, /s/*



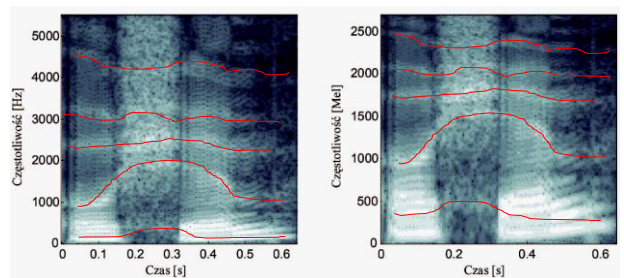
Spektrogramy i formanty



Rys. M. Kepiński



Spektrogramy i formanty



Rys. M. Kepiński



Formanty – średnie częstotliwości (Hz)

w kontekstach

VOWEL	FEMALE FORMANT				
	F1	F2	F3	F4	F5
[i]	361	2732	3431	4306	5178
[ɨ]	484	2807	2933	4423	5136
[e]	399	2140	2966	4349	5250
[ɛ]	821	1627	2710	4232	5240
[ɔ]	619	1194	2738	4091	5165
[o]	418	963	2810	4278	5045

VOWEL	MALE FORMANT				
	F1	F2	F3	F4	F5
[i]	282	2096	2760	3338	3887
[ɨ]	388	1742	2288	3355	3889
[e]	474	1750	2416	3357	3923
[ɛ]	420	1315	2243	3379	3842
[ɔ]	491	1035	2208	3323	3812
[o]	354	918	2156	3216	3708

izolowane

VOWEL	FEMALE FORMANTS			
	F1	F2	F3	F4
[i]	247	2775	3110	4135
[ɨ]	329	2335	3041	4214
[e]	623	2210	2979	3964
[ɛ]	999	1545	2756	3951
[ɔ]	641	1080	2793	3883
[o]	334	731	3140	4091

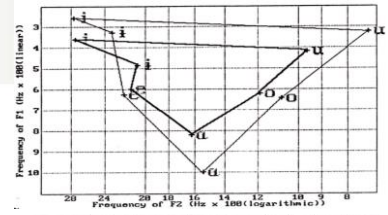
VOWEL	MALE FORMANTS			
	F1	F2	F3	F4
[i]	213	2216	2971	3542
[ɨ]	325	1996	2821	3422
[e]	553	1794	2861	3608
[ɛ]	723	1261	2513	3422
[ɔ]	538	883	2571	3326
[o]	280	602	2459	3332

Tab. 3. Kieśla, Alphabetic variation of Polish vowels, w: Studia Phonetica Posnaniensia, vol. 6, Poznań 2000

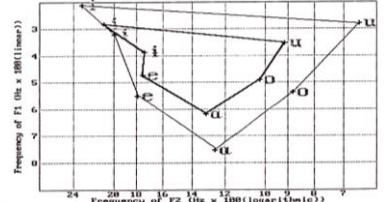


Średnie pozycje formantów

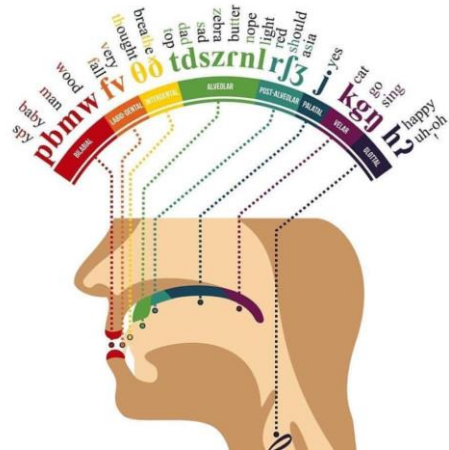
Kobiety



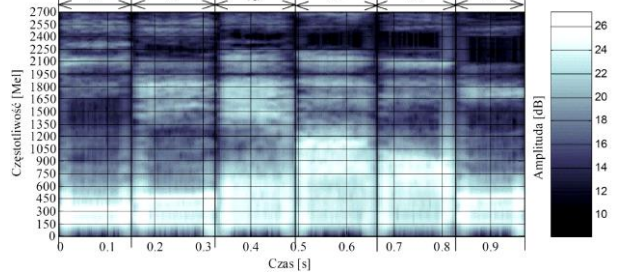
Mężczyźni



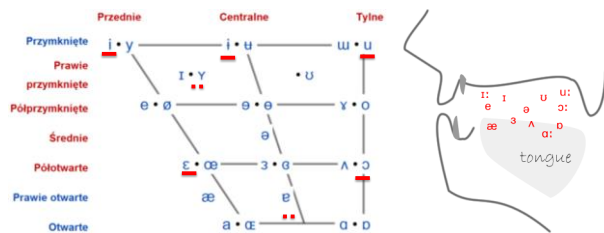
Tab. 3. Kieśla, Alphabetic variation of Polish vowels, w: Studia Phonetica Posnaniensia, vol. 6, Poznań 2000



Przykłady formantów widma szeregu samogłosek polskich



Alfabet IPA – samogłoski



z: Wikipeda [rev.]



Przetwarzanie sygnału mowy

- Rejestracja sygnału
 - mikrofon, mikrofony
 - wzmacniacz, kompresja dynamiki ($L > 30\text{dB}$)
 - przetwornik A/C z filtrem przeciwaliasingowym ($F_s > 8\text{kHz}$, $Q > 8\text{bit}$) $s(n) = s_a(t_0 + n \cdot dt)$, $n = 0, 1, 2, \dots$
- Normalizacja ...
- Redukcja szumu ...
- Kształtowanie charakterystyki
 - kompresja dynamiki, kompandacja
 - equalizacja, preemfaza
- Detekcja – energia, charakterystyka, ...
- Kodowanie, kompresja, parametryzacja, analiza, rozpoznawanie, synteza, ...



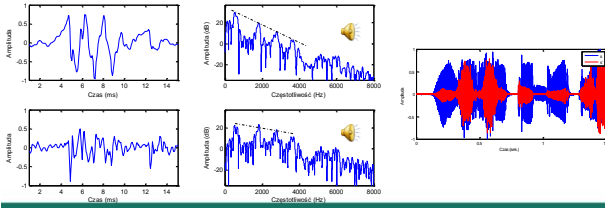
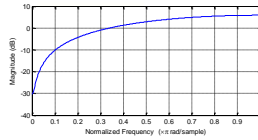
Preemfaza

- Filtracja górnoprzepustowa FIR

$$H(z) = 1 - bz^{-1}$$

$$0.92 < b < 0.98$$

- Kompensacja charakterystyki emisji głosu
- Zwiększenie istotności wyższych formantów



Parametryzacja



- Zapis rozpoznawanych obiektów w postaci ich liczbowych identyfikatorów
- Przetwarzanie wstępne sygnału (*pre-processing*, filtracja, normalizacja, redukcja szumu itp...)
- Wybór i zdefiniowanie wielkości charakterystycznych dla danej klasy obiektów
- Ekstrakcja cech i ich agregacja w postaci wektora
- Przetwarzanie wektorów cech (*post-processing*)
 - *VarNorm*, *SoftMax*, *PCA*, *SVD*, ...
- Transformacja, zmniejszenie wymiarowości, bez straty informacji



Ekstrakcja cech

- Redukcja redundancji informacyjnej zapisu danego obiektu, przy zachowaniu jego cech istotnych z punktu widzenia możliwości dyskryminacyjnych.
- Przekleństwo wymiarowości: duża liczba cech wymaga dużej ilości danych (N^2). Proporcja $N/L > 30$
- Wybór dokonany przez eksperta
- Metody automatyczne
 - Maksymalizacja zdolności dyskryminacyjnej
 - Minimalizacja błędów klasyfikacji
 - Maksymalizacja odległości międzyklasowej
 - Minimalizacja rozmiaru klas w przestrzeni
 - Analiza wariancji, statystyczne testy istotności, *test t*, ...



Cechy sygnału mowy

- Prozodia
- Energia
- Widmowe
- Cepstralne
- Formantowe
- Artykulacyjne
- Czasowe
- Inne temporalne
- Supersegmentalne
- Entropowe
- Statystyczne

Metody parametryzacji:

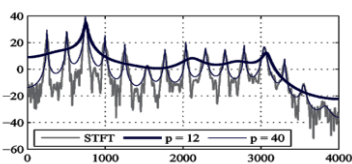
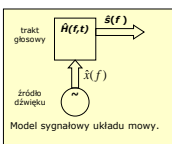
MFC, **MFCC**, LPC, **PLP**, **RASTA**, TRAPS, PWT, WPT, VQ, ...



Ekstrakcja cech sygnału mowy



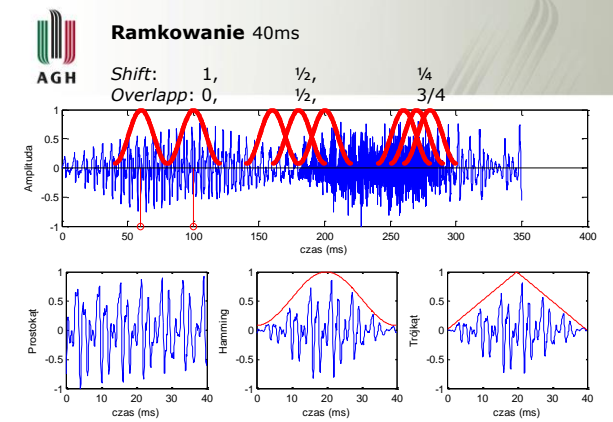
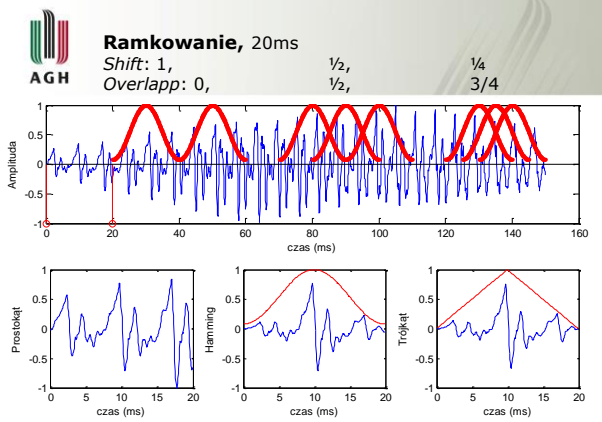
- Zastosowanie: rozpoznawanie mowy, rozpoznawanie mówców (PLP, MFCC), kodowanie i kompresja mowy (CELP), diagnostyka, przetwarzanie (*vocoder*), ...
- Mowa to zjawisko czasowo-częstotliwościowe, więc
- Analiza częstotliwościowa (STFT, FFT, Falki, LP) w krótkich (ok. 20-30 ms) okresach czasu (ramkach).
- Sygnał mowy:



Ramkowanie

- Konieczne zgromadzenie danych przed ich analizą
- Długość bufora / ramki
 - Zbyt długa – sygnał przestaje być stacjonarny
 - Zbyt krótka
 - problem w analizie niskich częstotliwości
 - brak okresu krtaniowego w ramce
 - niemożność ekstrakcji cech toru głosowego
 - silne efekty brzegowe
- Długość zazwyczaj 15-25 ms
- Zakładka (ang. *overlap*, *frame-shift*): 1, 1/2, 3/4, 1/3, ... ramki
- Okienkowanie (zmniejsza efekty brzegowe):

Hamming, Gauss, trójkąt



Ekstrakcja tonu krztaniowego F_0

AGH

- Zakres częstotliwości 80-1000 Hz, (częściej 100-500 Hz)

Ekstrakcja tonu krztaniowego F_0

AGH

- Zakres częstotliwości 80-1000 Hz, (częściej 100-500 Hz)
- Autokorelacja

$$s_w(m) = w(n-m)s(m)$$

$$r(l) = \sum_m s_w(m+l)s_w(m)$$

Ekstrakcja tonu krztaniowego F_0

AGH

- Cepstrum

$$s_w^c = F^{-1}(\log(|F(s_w)|))$$

$$\hat{s}(f,t) = \hat{x}(f,t) \cdot \hat{H}(f,t)$$

$$\log \hat{s}^c(f,t) = \log \hat{x}^c(f,t) + \log \hat{H}^c(f,t)$$

Ekstrakcja tonu krztaniowego F_0

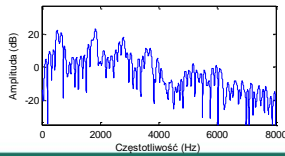
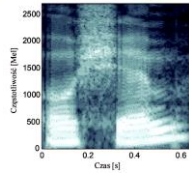
AGH

- Cepstrum – brak pobudzenia krztaniowego



Ekstrakcja formantów

- Zakres częstotliwości: 200Hz – 4000Hz
- skala Mel
- FFT – wygładzanie widma amplitudowego
- Predykcja liniowa – modelowanie traktu głosowego
- Inne modele
 - Statystyczne
 - Filtracyjne (Z)
 - AI
 - ...



Predykcja Liniowa

LP, ang. Linear Prediction

- Metoda modelowania cech traktu głosowego
- Umożliwia ekstrakcję formantów
- Zakładamy model sygnału $s(n) = \sum_{k=1}^K \alpha(k)s(n-k) + Ge(n)$
- Idea predykcji
 - ze znanych próbek sygnału wyznacz kolejne $\tilde{s}(n) = \sum_{k=1}^K \alpha(k)s(n-k)$
 - błąd predykcji $d(n) = s(n) - \tilde{s}(n) = s(n) - \sum_{k=1}^K \alpha(k)s(n-k)$
 - to odpowiedź FIR o transmitancji

$$A(z) = 1 - \sum_{k=1}^K \alpha(k)z^{-k} \quad \text{więc} \quad H_{10}(z) = \frac{G}{1 - \sum_{k=1}^K \alpha(k)z^{-k}}$$



Predykcja liniowa

wyznaczanie współczynników (LPC)

- Szukamy $d \equiv Ge \Leftrightarrow a \equiv a$
 - Kryterium najmniejszych kwadratów
- $$J = \left\| \left\langle s_w(m) - \sum_{k=1}^K \alpha(k)s_w(m-k) \right\rangle_m \right\|^2 \xrightarrow{\min_a} \Phi_{i,k} \alpha = \psi$$
- Metoda kowariancji (niewygodna, rozkład macierzy kowariancji)
 - Metoda autokorelacji
 - Stosujemy okienkowanie (określone granice estymacji)

- Ustalone K $J = \sum_{m=m_0-M}^{m_0+M+K} d^2(m)$



Predykcja liniowa

wyznaczanie współczynników

Macierz

$$\Phi_{i,k} : \varphi(i,k) = \sum_{m=m_0-M-i}^{m_0+M-i} s_w(m)s_w(m+i-k) = \sum_{m=m_0-M-k}^{m_0+M-k} s_w(m)s_w(m+k-i) \Rightarrow \varphi^T = \varphi$$

Więc jest to wektor autokorelacji

$$\varphi(i,k) = \varphi(i-k) \Rightarrow \varphi(k) = \varphi(-k) = \sum_m s_w(m)s_w(m+k)$$

Po wstawieniu do warunku $\min\{J\}$, otrzymujemy:

$$\sum_{k=1}^K \alpha(k)\varphi(i-k) = \varphi(i), \quad i = 1, 2, \dots, K$$

czyli

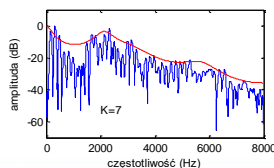
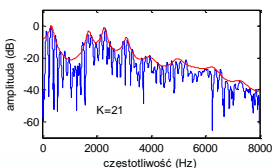
$$\begin{bmatrix} \varphi(0) & \varphi(1) & \dots & \varphi(K-1) \\ \varphi(1) & \varphi(0) & \dots & \varphi(K-2) \\ \dots & \dots & \dots & \dots \\ \varphi(K-1) & \varphi(K-2) & \dots & \varphi(0) \end{bmatrix} \begin{bmatrix} \alpha(1) \\ \dots \\ \alpha(K) \end{bmatrix} = \begin{bmatrix} \varphi(1) \\ \dots \\ \varphi(K) \end{bmatrix}$$



Charakterystyka filtra predycyjnego

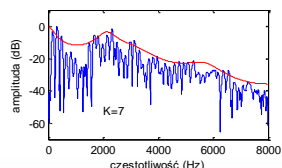
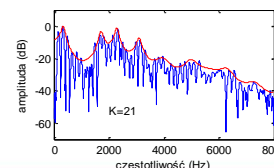
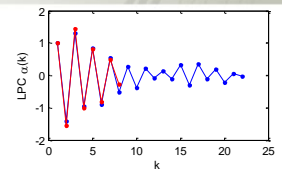
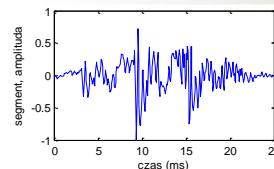
$$|H(f/f_s)| = \left| \frac{G}{1 - \sum_{k=1}^K \alpha(k)e^{-2\pi j k f / f_s}} \right|$$

- Pozwala wyznaczyć obwiednię widma segmentu
- Stała obwiedni (dobroć) zależy od K - rzędu predyktora
- Maksima lokalne obwiedni – formanty



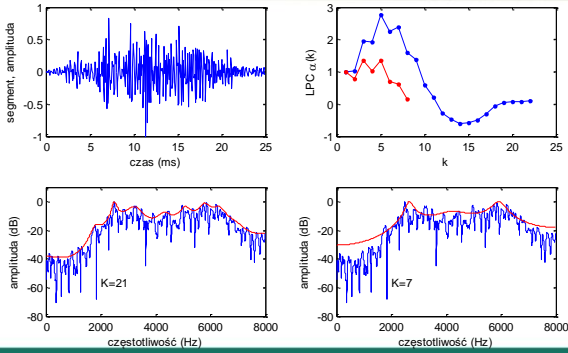
LPC

przykłady – głoska dźwięczna

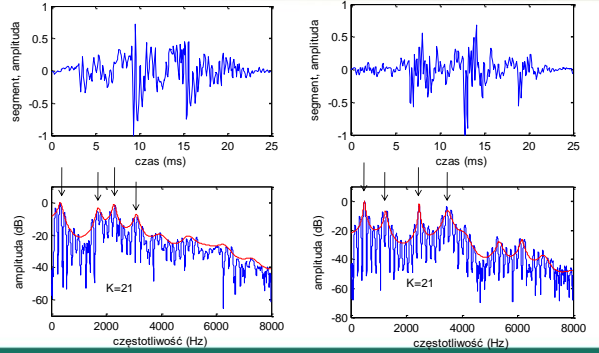




LPC przykłady - głoska bezdźwięczna

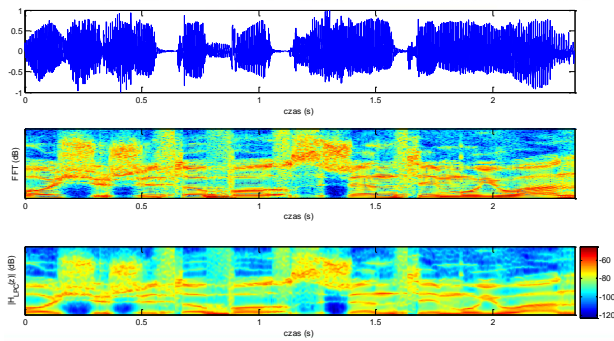


LPC przykłady - formanty

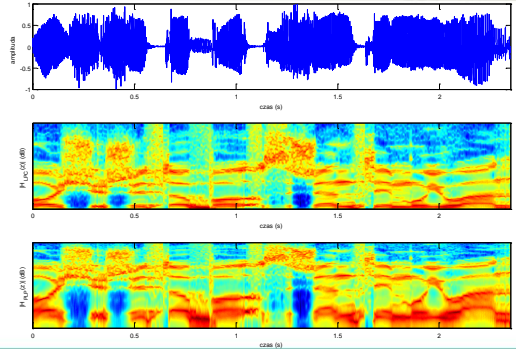


FFT vs LPC

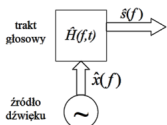
„Woź żyto bo zżęty mój lan, (...)”



PLP (ang. Perceptual Linear Prediction) przykład – widmo LP w skali melowej „Woź żyto bo zżęty mój lan, (...)”



Schemat wyznaczenia MFCC (Mel-frequency cepstral coefficients)

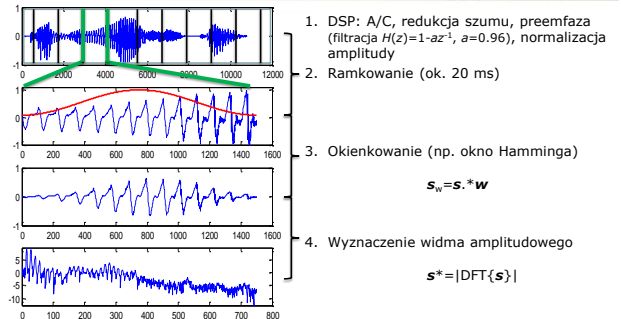


$$c[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |\hat{s}(e^{j\omega})| e^{j\omega n} d\omega$$

$$\hat{s}(f) = \hat{x}(f) \cdot \hat{H}(f, t)$$

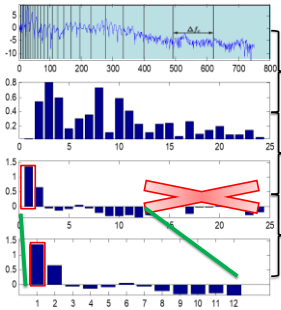


Schemat wyznaczenia MFCC (Mel-frequency cepstral coefficients) kroki 1-4

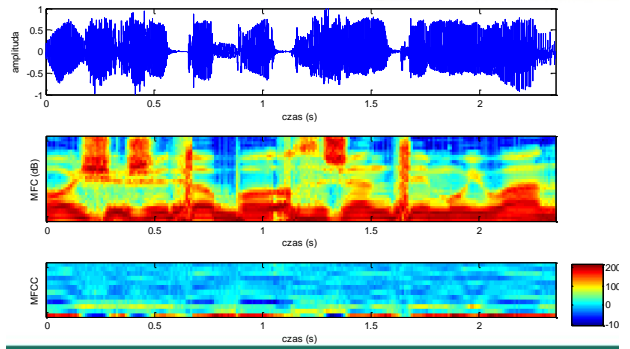




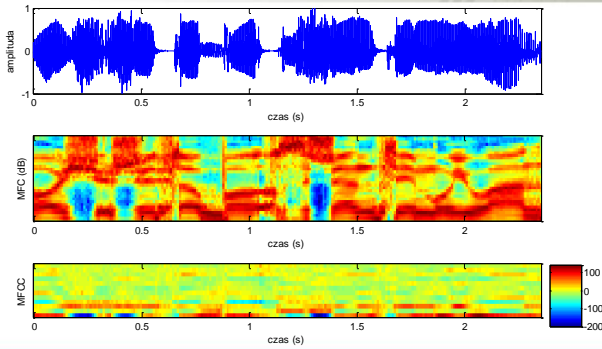
Schemat wyznaczania MFCC (Mel-frequency cepstral coefficients) kroki 5-8



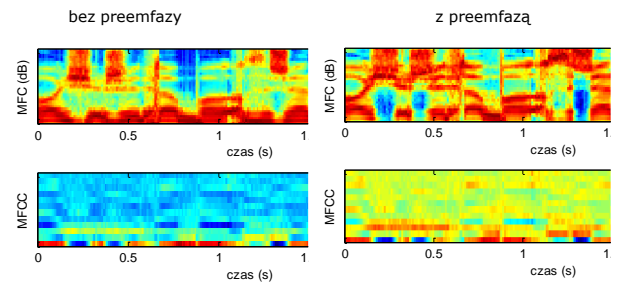
Przykład parametryzacji MFCC (32 MFC, 12 MFCC \ c₀, bez preemfazy) „Woź żyto bo zżęty mój lan, (...)”



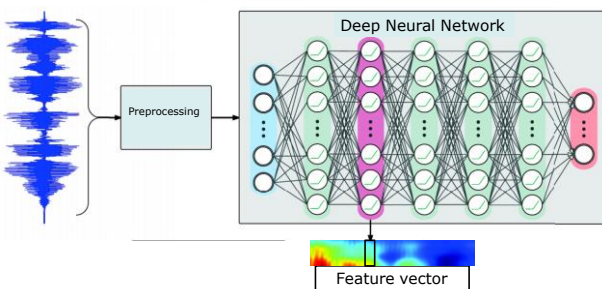
Przykład parametryzacji MFCC (32 MFC, 12 MFCC \ c₀, z preemfazą) „Woź żyto bo zżęty mój lan, (...)”



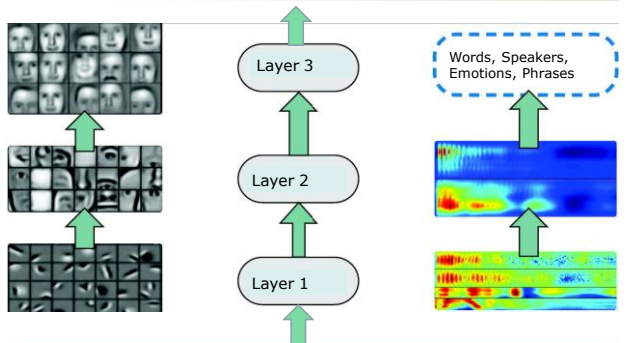
Przykład parametryzacji MFCC „Woź żyto bo zżęty mój lan, (...)”



Bottleneck Features – Deep Neural Networks



Bottleneck Features – Deep Neural Networks



Parametryzacja	Model	Ewaluacja
<ul style="list-style-type: none"> Ekstrakcja cech Normalizacja Redukcja wymiarowości 	<ul style="list-style-type: none"> Architektura Trening modelu Predykcja 	<ul style="list-style-type: none"> Analiza Statystyki Testowanie Interpretacja

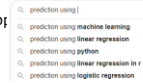
Classification

Predicts category (class, label)
Known labelled input $x : \{c_1, c_2, \dots\}$



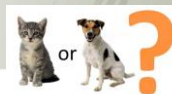
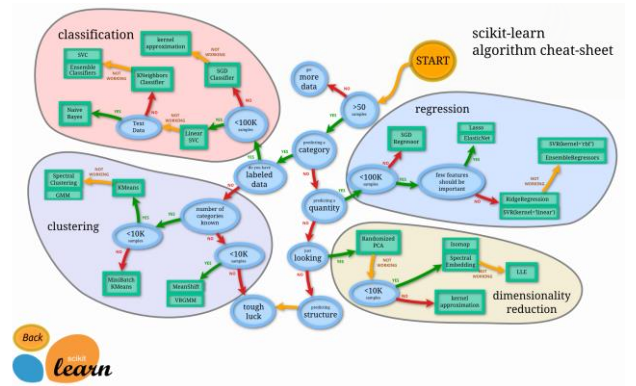
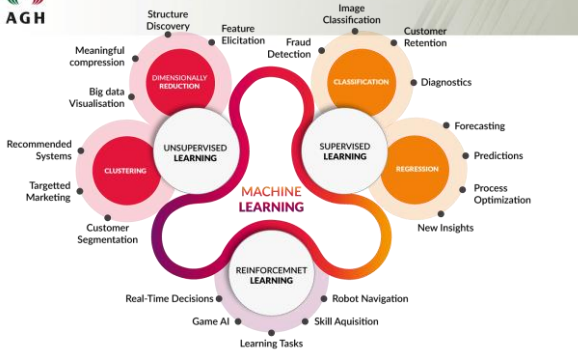
Regression

Predicts value
Known I/O value map $y=f(x)$



Clustering

Groups similar elements (no labels)



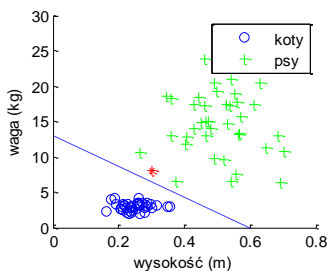
Predict known label for unseen data

\mathbf{x} - feature vector

$$\mathbf{x} = [x_1, \dots, x_L]^T$$

c_i - class, label from the set of M classes, $i=1, \dots, M$

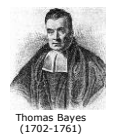
$$\mathbf{x} \rightarrow c_i, \bigcup_{i=1}^M c_i = \Omega$$



MAP (*Maximum a Posteriori prob.*): we are looking for a specific label c_i^*

$$c_i^* = \arg \max_{c_i, i=1, \dots, M} P(c_i | \mathbf{x})$$

That would maximize the observed probability of the class given the input \mathbf{x} .



Binary classification: $M=2$

Class priors: $P(c_1), P(c_2)$

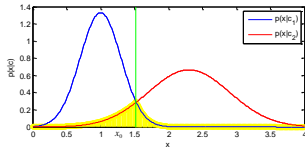
If we knew likelihood $p(\mathbf{x}|c_i)$:

$$P(c_i | \mathbf{x}) = \frac{p(\mathbf{x} | c_i)P(c_i)}{p(\mathbf{x})}$$



Binary Naive Bayes

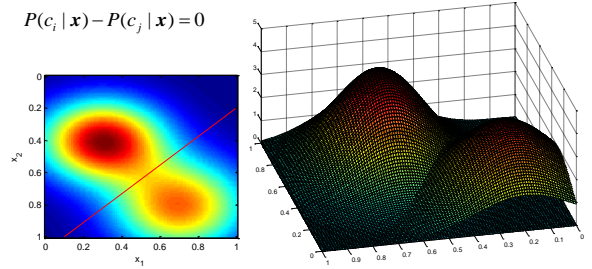
$$p(x | c_1)P(c_1) \stackrel{?}{>} p(x | c_2)P(c_2)$$



Naive Bayes – 2D

Dyscrimination Hyperplane

$$P(c_i | x) - P(c_j | x) = 0$$



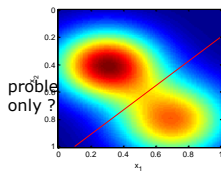
Decision function (Fischer discriminator)

Monotonic function of MAP

$$g_i(x) \equiv f(P(c_i | x)) \quad \forall j \neq i \quad g_i(x) > g_j(x) \Rightarrow x \in c_i$$

Decision plane

$$g_{ij}(x) \equiv g_i(x) - g_j(x) = 0, \quad i, j = 1, \dots, M, \quad i \neq j$$



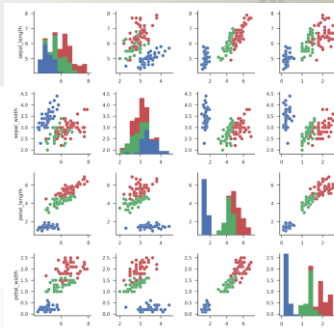
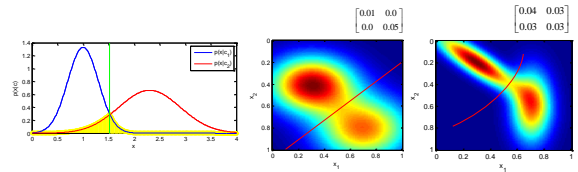
How to tackle an N-class using binary classifiers



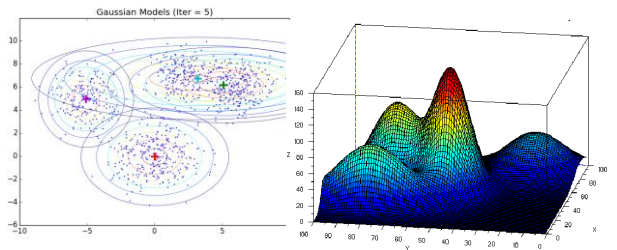
Naive Bayes – Normal PDF

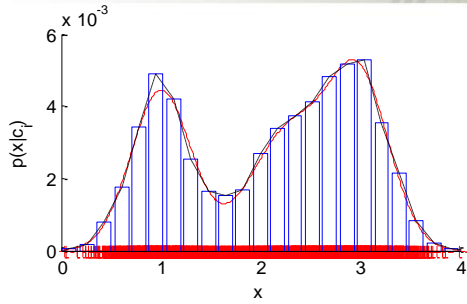
$$\text{Multivariate N-PDF } p(x | c_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{d/2}} \exp\left(-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right)$$

$$\mu_i = E[x] \quad \Sigma_i = E[(x - \mu_i)(x - \mu_i)^T] \quad \Theta_i = \{\mu_i, \Sigma_i\}$$



```
>>> from sklearn import datasets
>>> iris = datasets.load_iris()
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
>>> y_pred = gnb.fit(iris.data, iris.target).predict(iris.data)
>>> print("Number of mislabeled points out of a total %d points : %d"
... % (iris.data.shape[0], (iris.target != y_pred).sum()))
Number of mislabeled points out of a total 150 points : 6
```





Gaussian Mixture Model



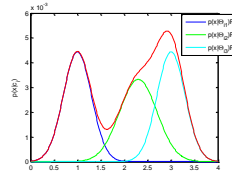
GMM – statistical model of complex distributions

$$p(x | c_i) = \sum_{j=1}^J p(x | \theta_{ij}) P_j$$

$$\sum_{j=1}^J P_j = 1, \int_{-\infty}^{\infty} p(x | \theta_j) dx = 1$$

$$\max \left\{ p^{ML} = \prod_k p(x_k, \theta, P_1, \dots, P_J) \right\}$$

mixture.GaussianMixture

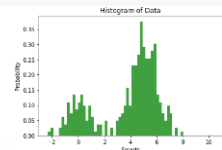


```
import numpy as np
from sklearn import mixture
np.random.seed(1)
g = mixture.GaussianMixture(n_components=3)
# Generate random observations with two modes centered on 0
# and 10 to use for training.
obs = np.concatenate((np.random.randn(100, 1), 10 + np.random.randn(300, 1), 5 + np.random.randn(200, 1)))
g.fit(obs)
p = g.predict([[0], [2], [6], [9], [10]])
print(p)
g.predict_proba([[0], [2], [6], [9], [10]])
print(p)
```

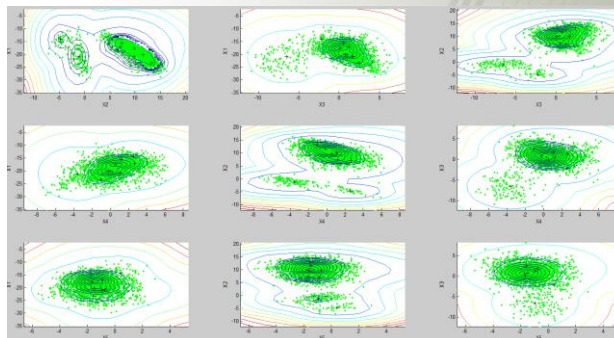
```
[ 1 1 2 0]
[[1.74450697e-23 9.99985504e-01 1.44962481e-05]
 [2.7870519e-14 8.18856673e-01 1.81143327e-01]
 [3.62624405e-04 3.43844650e-10 9.99637375e-01]
 [9.98599654e-01 3.64032460e-22 1.40034569e-03]
 [9.99984542e-01 1.86890101e-27 1.54577593e-05]]
```

```
import numpy as np
from sklearn import mixture
np.random.seed(1)
g1 = mixture.GaussianMixture(n_components=3)
g2 = mixture.GaussianMixture(n_components=3)
# Generate random observations with two modes centered on 0
# and 10 to use for training.
obs1 = np.concatenate((np.random.randn(100, 1), 5 + np.random.randn(300, 1))) # 0 & 5
obs2 = np.concatenate((2.5 + np.random.randn(100, 1), 6 + np.random.randn(200, 1))) # 2.5 & 6
g1.fit(obs1)
g2.fit(obs2)
p1g1_score = g1.score([4])
print("log p(x=4|g1)=", p1g1_score)
p2g2_score = g2.score([4])
print("log p(x=4|g2)=", p2g2_score)
p1g1_score = g1.score([2])
print("log p(x=2|g1)=", p1g1_score)
p2g2_score = g2.score([2])
print("log p(x=2|g2)=", p2g2_score)
print("What is P(g|x)=? , Bayes")
print("P(g1|x=2)=", np.exp(p1)*0.5/(np.exp(p1)*0.5+np.exp(p2)*0.5))
print("P(g2|x=2)=", np.exp(p2)*0.5/(np.exp(p1)*0.5+np.exp(p2)*0.5))
```

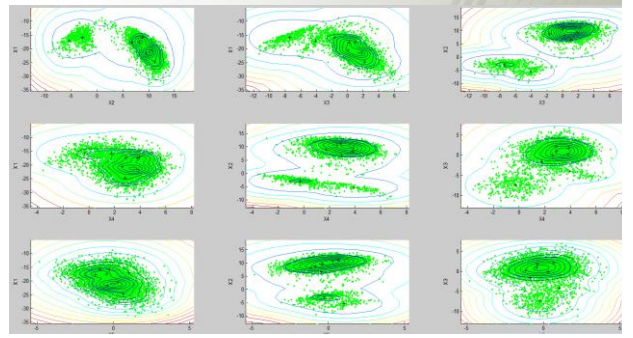
```
log p(x=4|g1) = -1.7072020753278263
log p(x=4|g2) = -2.2636967356701074
log p(x=2|g1) = -4.153287400767302
log p(x=2|g2) = -2.21274778571532
What is P(g|x)=? , Bayes
P(g1|x=2) = 0.1255860595922101
P(g2|x=2) = 0.8744139404077899
```



GMM w działaniu Estymacja K-Means EM na DCT(log(X)), wzrost 3sigma, /u/, 8GMM

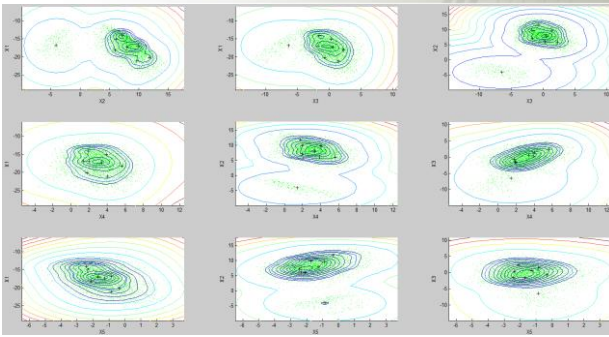


GMM w działaniu Estymacja K-Means EM na DCT(log(X)), wzrost 3sigma, /m/, 8GMM

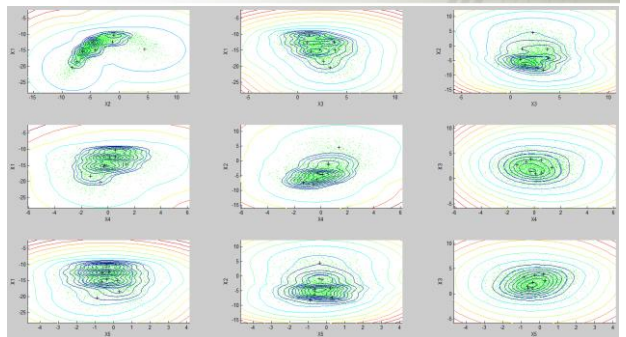




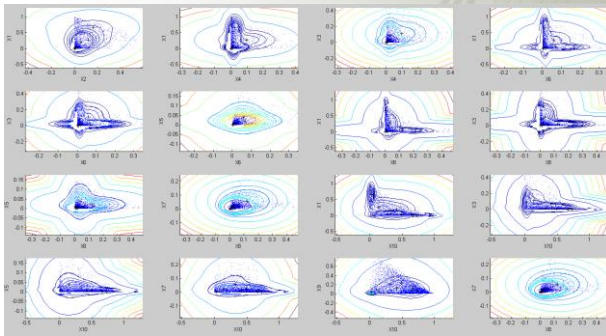
GMM w działaniu
Estymacja K-Means EM na $DCT(\log(X))$,
widok 3sigma, /l/, 8GMM



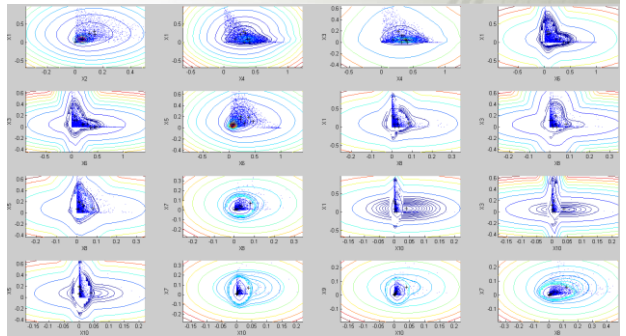
GMM w działaniu
Estymacja K-Means EM na $DCT(\log(X))$,
widok 3sigma, /s/, 8GMM



GMM w działaniu
Estymacja K-Means EM na X ,
widok 3sigma, /s/, 8GMM



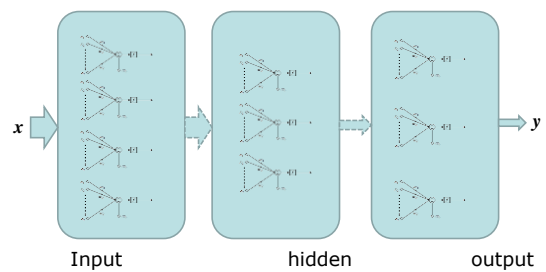
GMM w działaniu
Estymacja K-Means EM na X ,
widok 3sigma, /a/, 8GMM



Artificial Neural Networks



Multilayer Perceptron



Rosenblatt perceptron



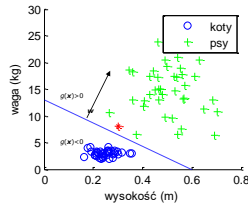
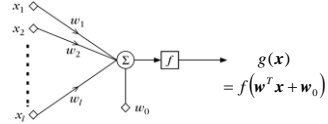
Frank Rosenblatt (1928-1971)

Algebraic function which defines a discriminating hyperplane

$$g(x) = w^T x + w_0, \quad w = [w(1), \dots, w(L)]$$

$$g_{ij} = w^T (x_i - x_j) = 0$$

w - hyperplane normal vector



<https://meta.stackoverflow.com/questions/30514/advertising-in-a-ny-a-platform-to-train-machine>

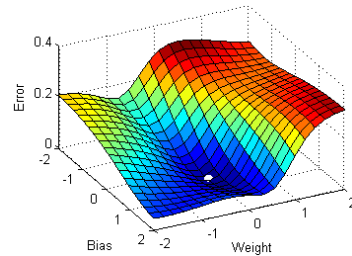
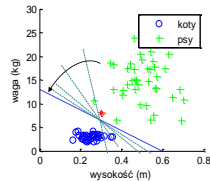
Neural weights estimation

- Iterative perceptron algorithm

$$w_{i+1} = w_i - \rho_i \sum_{x \in X_{err}} \delta_x x, \quad \rho_i \in \mathbb{R}^+, \quad \delta_x = \begin{cases} -1, & x \in C_1 \\ 1, & x \in C_2 \end{cases}$$

- Gradient descent

$$w' = w - L * dF(w, X) / dw$$



FFN/MLP, one-hot classifier

```
import keras
from keras.models import Sequential
from keras.layers import Dense, Dropout, Activation
from keras.optimizers import SGD

# Generate dummy data
import numpy as np
x_train = np.random.random((1000, 20))
y_train = keras.utils.to_categorical(np.random.randint(10, size=(1000, 1)), num_classes=10)
x_test = np.random.random((100, 20))
y_test = keras.utils.to_categorical(np.random.randint(10, size=(100, 1)), num_classes=10)

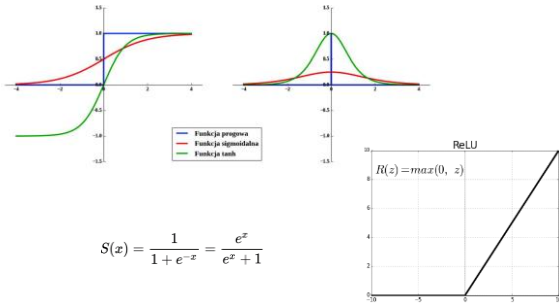
model = Sequential()
# Dense(64) is a fully-connected layer with 64 hidden units.
# in the first layer, you must specify the expected input data shape:
# here, 20-dimensional vectors.
model.add(Dense(64, activation='relu', input_dim=20))
model.add(Dropout(0.5))
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(10, activation='softmax'))

sgd = SGD(lr=0.01, decay=1e-6, momentum=0.9, nesterov=True)
model.compile(loss='categorical_crossentropy',
              optimizer=sgd,
              metrics=['accuracy'])

model.fit(x_train, y_train,
          epochs=20,
          batch_size=128, verbose=0)
score = model.evaluate(x_test, y_test, batch_size=128)
print(score)
```

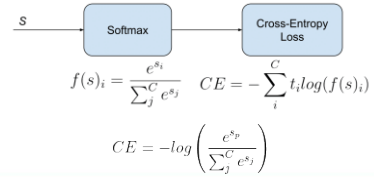
100/100 [=====] - 0s 704us/step
[2.2926833629608154, 0.1299999523162842]

Activation functions and derivatives

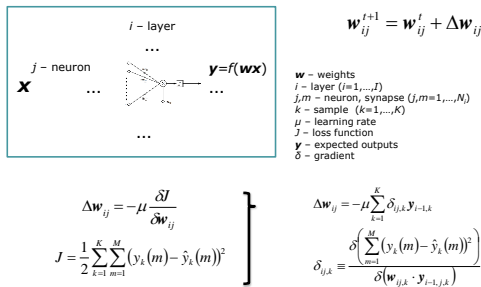


Softmax network outputs

- Usually used as the network output activation
- Trained to produce probability estimation
- 1 example – one class only (no multivariate class problems, one-hot approach)
- Commonly trained under a log loss (or cross-entropy) criterion



Backpropagation - assumptions



Backpropagation - solution

$\mathbf{w}_{ij}^{t+1} = \mathbf{w}_{ij}^t - \mu \sum_{k=1}^K \delta_{i,j,k} \mathbf{y}_{i-1,k}$

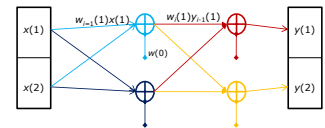
$f(x) = \frac{1}{1 + \exp(-\alpha x)}$
 $f'(x) = \text{af}(x)(1 - f(x))$

Solution for $J=L^*MSE(\mathbf{y}, \mathbf{y}^{\wedge})$

$i = I: \delta_{i=L,j,k} = [f(\mathbf{w}_{ij} \mathbf{y}_{i-1,k}) - \hat{y}_k(j)] f'(\mathbf{w}_{ij} \mathbf{y}_{i-1,k})$

$i < I: \delta_{i-1,j,k} = \left(\sum_{m=1}^{N_i} \delta_{im} w_{im}(j) \right) f'(\mathbf{w}_{i-1,j} \mathbf{y}_{i-2,k})$

w - weights
i - layer ($i=1, \dots, I$)
j, m - neuron, synapse ($j, m=1, \dots, N$)
k - sample ($k=1, \dots, K$)
 μ - learning rate
J - loss function
 \mathbf{y} - expected outputs
 δ - gradient



Backpropagation - implementation

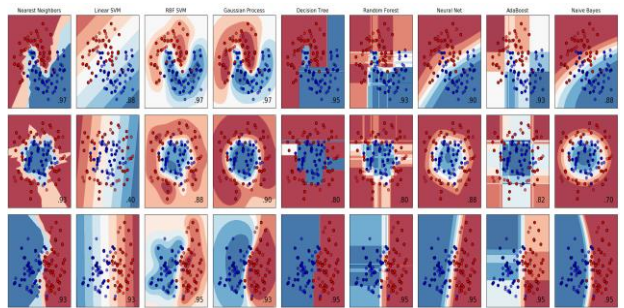
1. Inicjalizacja – wylosuj wszystkie wagi $\mathbf{w}^t=0$
2. Wyznacz odpowiedzi $\mathbf{y}_{i,k}=f(\mathbf{w}\mathbf{y}_{i-1,k})$ wszystkich neuronów sieci dla każdego wektora treningowego (*forward step*)
3. Oblicz funkcję kosztu *J* (krok pośredni)
4. $i=I$, Dla wszystkich wzorców (*k*) wyznacz wartość $\delta_{i,j,k}$ każdego neuronu wyjściowego (*backward step*)
5. Dla wszystkich warstw (kolejno wstecz) $i=I, I-1, \dots, 2$ wyznacz $\delta_{i-1,j,k}$ wszystkich neuronów dla każdego wzorca (*k*)
6. Ustaw wszystkie wagi $\mathbf{w}_{ij}^{t+1} = \mathbf{w}_{ij}^t - \mu \sum_{k=1}^K \delta_{i,j,k} \mathbf{y}_{i-1,k}$
7. $t=t+1$, goto 2, until $J < J_{tr}$ OR $dJ/dw < \epsilon$ OR $t > T$

Stochastic gradient descent

- GD Problems
 - RAM
 - Computational power
- SGD Solutions
 - Batch sample
 - Batch-size
 - Epochs



Comparison of classifiers

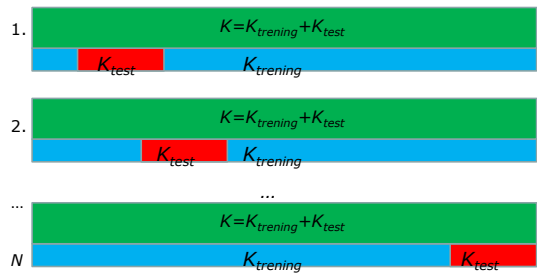


Model evaluation

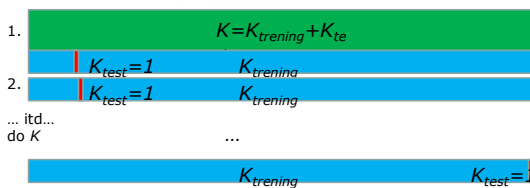
- How good is our model ?
- Generalization, Robustness
- Loss Function
- Performance: Accuracy, Error rate (%)

N-fold Cross-validation

$$ERR = \frac{1}{N} \sum_{n=1}^N ERR_n$$



Leave-one-out



Maksymalizacja mocy zbioru testowego i treningowego jednocześnie

Overfitting

- Train/Eval mismatch
- Data problem
- Model problem
- Optimizer problem



How to avoid overfitting

- Better Data (quantity, quality, representativeness)
 - Data augmentation
 - Feature engineering
 - Hyper-parameter optimization (x-val)
 - Model (architecture and size change)
 - Regularization L1, L2
 - Batch-Normalization
 - Dropout
 - Shorter training
-