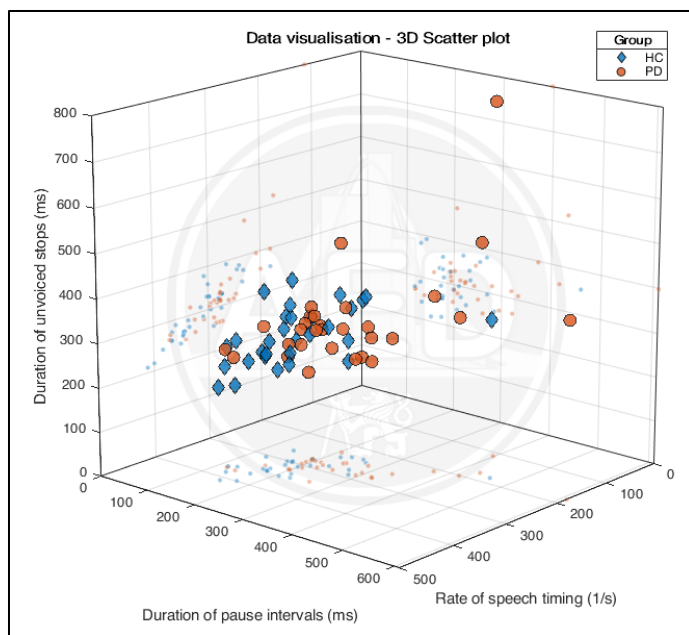


Zadání Cvičení #10

Popis dat: Pracovní data jsou uložena v souboru **data.csv**, který je k dispozici ke stažení na Moodle stránce tohoto předmětu, ve složce příslušného cvičení. Pro načtení dat do Matlabu využijte funkce `readtable`.

Data jsou ve formátu tabulky, která obsahuje data od pacientů s Parkinsonovou nemocí (label **PD**) a data od kontrolní skupiny zdravých lidí (label **HC**). Tabulka obsahuje ID kódy subjektů, identifikátory příslušnosti ke skupině (labels) a hodnoty těchto tří parametrů:

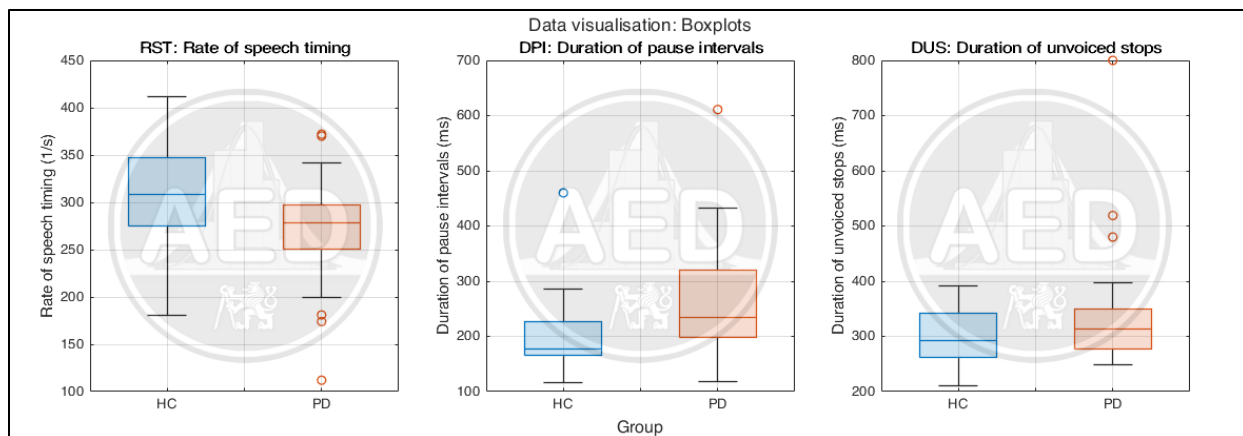
- **RST – Rate of Speech Timing**
 - Popisuje rychlost řeči pomocí identifikace znělých, neznělých a přerušovaných segmentů v záznamu. Jednotkou je s^{-1} .
- **DPI – Duration of pause intervals**
 - Popisuje kvalitu časování řeči (pauzy v řeči indikují problémy s iniciací), jednotkou jsou milisekundy, *ms*.
- **DUS – Duration of unvoiced stops**
 - Parametr popisující délku neznělých souhlásek, při kterých je náhle zastaven (a vypuštěn) proud vzduchu, jednotkou jsou opět milisekundy, *ms*.



Obrázek 1: 3D Scatter graf pro obě vyšetřované skupiny

Obě skupiny podstoupily vyšetření řeči pomocí volného monologu. Řečové nahrávky byly analyzovány a ze záznamu byly vypočteny hodnoty parametrů. Další informace viz [Hlavnička et al. 2017](#) (Reference).

Pro zjednodušení práce **jsme za vás provedli vizualizaci dat** (viz grafy na této stránce). Na **obrázku 1** můžete vidět vykreslená data pomocí 3D grafu typu scatter, včetně projekcí dat do rovin x-y, y-z a x-z. Tyto projekce pak přirozeně odpovídají 2D scatter grafům pro jednotlivé dvojice parametrů. **Obrázek 2** pak obsahuje klasické boxploty pro všechny tři parametry.



Obrázek 2: Boxploty pro analyzované parametry

Zadání úlohy	body
<p>Z poskytnutých dat natrénujte dva klasifikační modely <u>na standardizovaných datech</u> (u obou modelů ponechte defaultní nastavení Matlabu pro funkci <code>fitcsvm</code>):</p> <ul style="list-style-type: none"> SVM klasifikátor s Lineárním Kernelem SVM klasifikátor s Radial Basis Function (RBF) Kernelem <p>Oba modely validujte na stejných datech, ze kterých jste model natrénovali:</p> <ul style="list-style-type: none"> Srovnajte skutečné (<i>true</i>) labels a predikované labels pomocí <i>Confusion Matrix</i> <ul style="list-style-type: none"> Využijte funkce <code>confusionmat</code> pro výpočet a <code>confusionchart</code> pro zobrazení výsledků. Nezapomeňte uvést správné označení skupin. V tabulce si také můžete zobrazit přesnosti klasifikace pro jednotlivé skupiny pomocí argumentu <code>RowSummary</code>. Vykreslete si ROC (<i>Receiver operating characteristics</i>) křivku pro každý z modelů <u>do jednoho obrázku</u>. <ul style="list-style-type: none"> Použijte funkci <code>perfcurve</code> (Jako pozitivní skupinu berte pacienty s PN). Uveďte také hodnotu AUC pro oba modely. <p>Vykreslete obrázek s vypočtenými maticemi zmatení (Confusion Matrices, <i>zmaticemi</i>) a dobře ho popište. Také vykreslete obrázek s ROC křivkami pro oba modely. Slovy запиšte hodnoty AUC a vašimi slovy zhodnoťte výsledky.</p>	1.5
<p>Provedte 5-Fold Cross-Validaci SVM klasifikátoru s lineárním kernelem.</p> <p>Cross-validujte SVM s lineárním kernelem, model vždy uče na <u>standardizovaných datech</u>.</p> <p>Určete přesnost modelu na základě těchto veličin (hodnoty <i>True/False Positive/Negative</i> můžete určit manuálně nebo například pomocí funkce <code>confusionmat</code>):</p> <div style="display: flex; align-items: center;"> <ul style="list-style-type: none"> • Senzitivita • Specificita • Accuracy <div style="margin-left: 20px;"> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; display: inline-block; text-align: center;"> <p style="color: green; margin: 0;">Specificity</p> $SPC = \frac{TN}{TN + FP}$ </div> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; display: inline-block; text-align: center; margin-left: 10px;"> <p style="color: red; margin: 0;">Sensitivity</p> $SEN = \frac{TP}{TP + FN}$ </div> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; display: inline-block; text-align: center; margin-left: 10px;"> <p style="color: black; margin: 0;">Accuracy</p> $ACC = \frac{TP + TN}{TP + FP + FN + TN}$ </div> </div> </div> <p>Způsob výpočtu 5-fold cross-validace lze provést (libovolně vyberte):</p> <ul style="list-style-type: none"> Manuálním rozdělením na testovací a trénovací množiny pro 5-fold schéma, iterativním výpočtem SVM modelu a průměrováním výsledných hodnot senzitivity, specificity a accuracy z jednotlivých iterací. 😞 Rozdělením dat pro 5-fold schéma pomocí funkce <code>cvpartition</code> a opět iterativním výpočtem modelu SVM a průměrováním výsledků. 😊 Přímým použitím vestavěné cross-validace ve funkci pro výpočet SVM modelu, <code>fitcsvm</code>, a klasifikace pomocí funkce <code>kfoldPredict</code>. 😊 <p>Jedním z výše uvedených způsobů naimplementujte 5-fold cross-validaci. Uveďte výsledné vypočtené hodnoty senzitivity, specificity a accuracy.</p>	1

Najděte optimální hodnoty hyperparametrů SVM klasifikátoru s RBF kernelem - *Box constraint* a *Kernel scale* (v přednášce značené jako *C* a *σ*) pomocí techniky *Grid Search*.

- Mřížku pro techniku *Grid Search* nastavte v rozsahu hodnot **0.01 až 3.01** pro oba parametry, přičemž testovací krok zvolte jako **0.1**.
- Na každém bodu mřížky hyperparametrů vypočtete hodnoty **specificity, senzitivity a accuracy cross-validovaného modelu** pomocí techniky *Leave-One-Out*.
 - Ať už jste v předchozím bodě realizovali cross-validaci jakkoli, zde využijte vstupního argumentu `Leaveout` funkce `fitcsvm`.
 - Výpočet modelů ve všech bodech mřížky může zabrat několik minut v závislosti na rychlosti vašeho počítače. Pokud máte nainstalovaný *Parallel Computing Toolbox*, můžete pro zrychlení výpočtu využít cyklu `parfor`, který automaticky paralelizuje výpočty na všechna procesorová jádra vašeho počítače.
- **Výsledné optimalizační povrchy vykreslete** pomocí funkce `plotGridSearch`, která je dostupná ke stažení na Moodlu.
 - Tato funkce vám také vrátí tabulku s optimálními hodnotami hyperparametrů z hlediska specificity, senzitivity a accuracy.
 - Popis toho, jak do funkce `plotGridSearch` vložit vstupní data, naleznete přímo v ní.

Implementujte optimalizaci parametrů pomocí techniky *Grid Search*. Vytvořte pomocí funkce `plotGridSearch` obrázek s vypočtenou mřížkou a zapište výsledné optimální body z mřížky (optimální hodnoty parametrů *Box Constraint* a *Kernel Scale* pro nejlepší hodnoty senzitivity, specificity a accuracy).

1.5

Nepovinný bonus:

Naprogramujte v Matlabu nejjednodušší neuronovou síť – [Perceptron](#). 🤖

Načtěte ze souboru **bonus.mat** data pomocí funkce `load`. Data v proměnných *A* a *B* jsou vektory o velikost 2x100, tedy dva soubory stovky 2D bodů. Tyto body mohou reprezentovat např. hodnoty dvou parametrů pro dvě rozdílné skupiny lidí.

Nejprve trochu teorie... 😊

Uvažujte rozhodovací funkci $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, která rozděluje prostor na dvě poloviny. Pokud $\mathbf{w} = [w_0, w_1, w_2]$ je normálový vektor k rozhodovací hranici a $\mathbf{x} = [1, x_1, x_2]$ je jeden datový bod, pak pokud:

- $g(\mathbf{x}) > 0$, \mathbf{x} patří do třídy 1
- $g(\mathbf{x}) < 0$, \mathbf{x} patří do třídy 2
- $g(\mathbf{x}) = 0$, \mathbf{x} leží na rozhodovací hranici.

A navíc absolutní hodnota funkce $g(\mathbf{x})$ udává vzdálenost bodu od rozhodovací hranice.

Dále uvažujme funkci $f(\mathbf{w}, \mathbf{x}_i) = l_i \cdot g(\mathbf{w}, \mathbf{x}_i)$. Tato funkce „kontroluje“ správnou klasifikaci tím, že kombinuje údaj o skutečné příslušnosti (v labelu l_i) a výsledek klasifikace z funkce $g(\mathbf{w}, \mathbf{x}_i)$. Správně klasifikované body budou mít výsledek funkce $f(\mathbf{w}, \mathbf{x}_i)$ vždy kladný, nesprávně klasifikované pak vždy záporný.

0.5

Můžeme pak vytvořit chybovou funkci J_p , která se nazývá **Perceptron criterion function** a která bude akumulovat vzdálenosti všech špatně klasifikovaných bodů od rozhodovací hranice (mínus je tam protože špatně klasifikované body budou mít zápornou hodnotu funkce a my pro ně chceme kladnou chybu):

$$J_p(\mathbf{w}) = \sum_{x_i \in X} -f(\mathbf{w}, x_i) = \sum_{x_i \in X} -l_i \cdot \mathbf{w}^T x_i$$

Pokud tato funkce bude mít nulovou hodnotu, znamená to, že všechny vzorky byly správně klasifikovány – tím pádem je cílem jí minimalizovat.

Najdeme minimum – provedeme derivaci této funkce vzhledem k \mathbf{w} :

$$\nabla J_p = \sum_{x_i \in X} -l_i x_i$$

Tím pádem poté, co náhodně „nastřelíme“ rozhodovací hranici iniciací vektoru \mathbf{w} , se podíváme na momentální klasifikace všech vzorků. Vybere ty vzorky, které jsou špatně klasifikované a s jejich pomocí updatujeme naši rozhodovací hranici ve směru nižší hodnoty chybové funkce.

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \cdot \nabla J_p = \mathbf{w}(t) + \eta \sum_{x_i \in X} l_i x_i$$

Kde t značí číslo „epochy“, po kterou je rozhodovací hranice stejná a pracuje se na aktualizaci, a η značí „velikost kroku“, tedy jak rychle se bude hranice moci měnit skrze epochy. η je dobré nastavit např. na hodnotu 0.1, 0.01, 0.001 apod.

Prakticky... 🐒

- Obě skupiny **sjednoťte do vektoru** o velikost 2x200.
- **Rozšiřte sjednocený vektor dat vektorem samých jedniček.** Toto je důležité proto, aby algoritmus mohl vypočítat *bias* (váhu, která posouvá rozhodovací hranici).
- Připravte si vektor labelů $L = l_1, l_2, \dots, l_N$ o velikosti 1x200, který bude odpovídat vektoru dat, s **hodnotami 1 pro skupinu A a -1 pro skupinu B.**
- **Inicializujte si vektor vah (weights) \mathbf{w}** o velikosti 3x1 náhodnými malými čísly (např. 0.1). Vektor vah má jeden prvek pro každou z dimenzí dat a jeden prvek pro výpočet *biasu* w_0 . Také si nastavte hodnotu velikosti kroku η .
- **Implementujte algoritmus batch perceptronu** teoreticky popsany výše a v pseudokódu shrnutý v rámečku.

Do Moodle odevzdejte obrázek s vaší kódovou implementací Perceptronového klasifikátoru a grafické zobrazení dat a vypočtené rozhodovací hranice.

Algoritmus Hromadného (Batch) Perceptronu

- 1: Inicializujte hodnoty: \mathbf{w}, η
- 2: Iterujte epochy:
- 3: $X = \{ \}$
- 4: **Pro:** $i = 1: N$
- 5: **Pokud:** $f(\mathbf{w}, x_i) < 0$, pak $x_i \rightarrow X$
- 6: $\mathbf{w} = \mathbf{w} + \eta \cdot \sum_{x_i \in X} l_i x_i$
- 7: **Dokud** $X = \{ \}$

Reference

Hlavnička, J., Čmejla, R., Tykalová, T., Šonka, K., Růžička, E., and Rusz, J. (2017). *Automated analysis of connected speech reveals early biomarkers of Parkinson's disease in patients with rapid eye movement sleep behaviour disorder*. Scientific reports, 7, 12, <https://doi.org/10.1038/s41598-017-00047-5>.