

Zadání Cvičení #9

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 funkci `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), údaje o pohlaví, věku, a hodnoty těchto dvou akustických parametrů:

- **stdPWR** – Směrodatná odchylka intenzity řeči, jednotka dB.
- **stdF0** – Směrodatná odchylka F0, jednotka půltóny – *semitones* (st).

Oba parametry byly vyhodnoceny z nahrávek souvislé řeči (monologu) pomocí programu PRAAT (viz [Boersma and Weenink 2017](#)).

Zadání úlohy	body
<p>Vizualizujte si data pomocí 2D-scatter grafu. Obě skupiny od sebe barevně odlište a použijte libovolné značky a velikost tak, aby byly data dobře viditelná. Navíc si můžete případně vykreslit rozdělení dat pomocí boxplotů.</p> <p>Odpovězte svými slovy na tyto otázky:</p> <ol style="list-style-type: none"> Pohledem na vaši vizualizaci zhodnoťte, jak dobře lze oddělit skupiny v 2D prostoru pomocí lineární rozhodovací hranice. Do vizualizace můžete zakreslit (od ruky) kam byste rozhodovací hranici umístili. Zhodnoťte, zda není pro klasifikaci některý z parametrů nadbytečný, a zvažte, jak by se obecně daly nadbytečné (redundantní) parametry automaticky identifikovat ve vícedimenzionálních datasetech. 	1
<p>Proveďte klasifikaci dat pomocí následujících metod a pro každou z nich vykreslete do obrázku s daty příslušnou rozhodovací hranici (<i>decision boundary</i>), abyste je mohli mezi sebou snadno vizuálně porovnat:</p> <ol style="list-style-type: none"> Logistická regrese (Logistic regression) <ul style="list-style-type: none"> Použijte funkce <code>glmfit</code> a <code>glmval</code>, se standardním nastavením pro binomiální rozdělení výstupní funkce. Rozhodovací hranici můžete vypočítat z definice logistické funkce nebo můžete využít funkci <code>contour</code> na úrovni 0.5. Lineární diskriminantní analýza (LDA) <ul style="list-style-type: none"> Pro výpočet LDA modelu použijte funkci <code>fitdiscr</code>. Klasifikaci dat můžete provést pomocí funkce <code>predict</code>. Rozhodovací hranici vykreslete za použití koeficientů z vypočteného LDA modelu. 	1.5

<p>3. Support Vector Machine (SVM) s lineárním jádrem</p> <ul style="list-style-type: none"> • Výpočet SVM modelu proveďte pomocí funkce <code>fitcsvm</code>. • Pokud chcete, můžete si vykreslit vhodným způsobem do obrázku jaké datové body použil model jako <i>Support Vectors</i>. • Koeficienty pro vykreslení rozhodovací hranice SVM modelu získáte z Matlabovské proměnné (<i>structure</i>) modelu (pole <code>Beta</code> a <code>Bias</code>). 	
<p>Porovnejte jednotlivé klasifikátory z hlediska výpočetní náročnosti a přesnosti klasifikace.</p> <p>a) Pro měření výpočetní náročnosti využijte časoměrné funkce <code>tic</code> a <code>toc</code>.</p> <p>b) Pro každý z klasifikátorů vypočtete a zobrazte matici záměn (<i>Confusion matrix</i>) – využijte funkci <code>confusionmat</code>, resp. <code>confusionchart</code>. Přesnost klasifikace určete jako procento dat, které by model přiřadil do správné skupiny.</p> <ul style="list-style-type: none"> • Zapište časy výpočtu pro všechny modely a procentuální přesnosti klasifikace. • Svémi slovy zhodnoťte výsledky a porovnejte modely. 	1.5
<p>Nepovinný bonus:</p> <p>Implementujte manuálně algoritmus pro výpočet lineární diskriminační analýzy LDA. Pro otestování funkce použijte data uložená v souboru bonus.csv, který je dostupný ke stažení na Moodle (proměnné jsou značeny x a y).</p> <p>Odevzdejte obrázek s vaší kódovou implementací a grafické zobrazení dat a rozhodovací hranice.</p>	0.5

Reference

Boersma, P. and Weenink, D. (2017). *Praat: doing phonetics by computer* [Computer program]. Version 6.0.30, retrieved 22 July 2017 from <http://www.praat.org/>