



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

گزارش تمرین سوم درس داده کاوی

«آشنایی با نرم افزار Weka»

گردآورنده: سعید دادخواه

استاد: دکتر ناظر فرد

بهمن ۱۳۹۵

بخش اول: ساخت فایل arff

برای تبدیل فایل txt به فایل weka به فرمت arff از نرم افزار weka استفاده خواهد شد. بعد از اجرای نرم افزار weka در Weka GUI Chooser در شاخه Tools از گزینه Arff Viewer برای این کار استفاده می شود. البته برای اجرای همه مراحل به جای استفاده از GUI از یک اسکریپت برای انجام همه مراحل استفاده خواهد شد و برای انجام این مرحله از کلاس `weka.core.converters.CSVLoader` استفاده می شود.

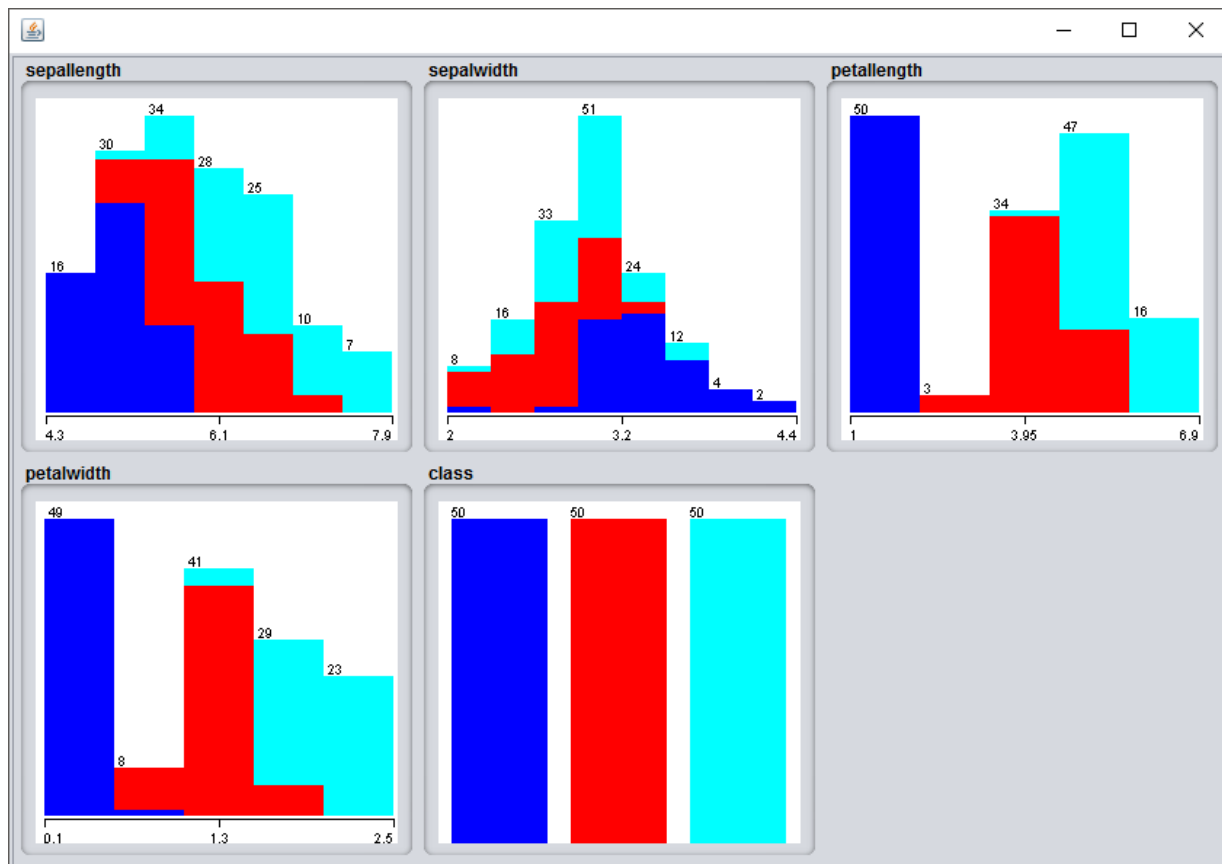
در صورتی که متد آغازین کلاس فوق اجرا شود برنامه به این شکل عمل خواهد کرد که آدرس یک فایل csv به عنوان ورودی می گیرد. ستون های این فایل csv باید دارای عنوان باشند. پس ابتدا یک برنامه به زبان پایتون اجرا می کنیم که یک فایل csv می سازد و عناوین ستون ها را نوشته و فایل `sample.txt` را به صورت csv می خواند و به ادامه فایل قبل اضافه می کند. پس از ساختن فایل csv کلاس بالا می تواند آن را خوانده و فایل arff را تولید کند.

برای انجام این مراحل باید پوشه های `data` و پوشه حاوی `make_weka.bat` و `add_header.py` در کنار یکدیگر قرار گیرند و فایل `sample.txt` نیز در پوشه `data` باشد. با اجرای `make_weka.bat` ابتدا کد پایتون `add_header.py` اجرا می شود و فایل `sample.csv` در پوشه `data` ایجاد می شود که همان فایل txt است با این تفاوت که عناوین ستون ها نیز به آن ها اضافه شده است. پس از آن کلاس بالا اجرا می شود و فایل `sample.arff` از روی فایل csv در همان پوشه `data` ساخته می شود.

بخش دوم: Visualize All

با انتخاب گزینه Visualize All پنجره زیر باز می شود. در این نمودارها به هر کلاس یک رنگ اختصاص داده می شود و چگونگی توزیع هر یک از کلاس ها در هر کدام از ویژگی ها نمایش داده می شود. برای نمایش این موارد در ویژگی های عددی از نمودار هیستوگرام و برای مقادیر غیر عددی از نمودار میله ای استفاده می شود.

نمودار آخر که کلاس هر کدام از نمونه ها را نمایش می دهد کاملاً مشخص می کند که از هر کلاس پنجاه نمونه داریم. از این پس کلاس های مربوط به رنگ آبی، قرمز و آسمانی به ترتیب کلاس اول تا سوم نامیده خواهند شد. نمودار اول نشان می دهد که به طوری کلی کلاس اول در این ویژگی کمتر از کلاس دوم و مخصوصاً کلاس سوم است. نمودار دوم نشان می دهد که این ویژگی نمی تواند نقش موثری در جداسازی کلاس سوم از دیگر کلاس ها داشته باشد ولی می تواند به تشخیص کلاس اول و دوم کمک کند. نمودار سوم و چهارم نشان می دهند که این ویژگی ها در تشخیص کلاس اول می توانند بسیار خوب عمل کنند. در مورد دو کلاس دیگر نیز این ویژگی ها می توانند موثر واقع شوند.

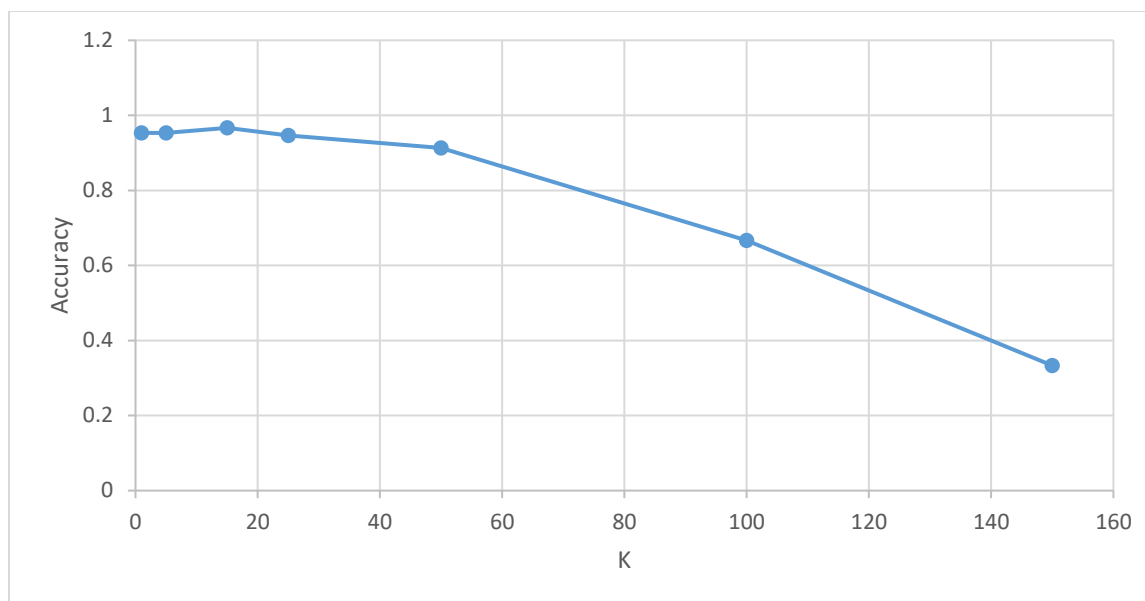


بخش سوم: نرمال سازی و گسسته سازی

با توجه به دستورات آورده شده در صورت پروژه مراحل زیر انجام شدند. برای نرمال سازی در تنظیمات فیلتر مقدار scale برابر دو و مقدار translation برابر منفی یک در نظر گرفته شدند. برای گسسته سازی نیز مقدار attribute indices برابر سه تا چهار و bins برابر پنج در نظر گرفته شدند. داده های به دست آمده در فایل iris-norm-disc.arff در پوشه data در دسترس هستند.

بخش چهارم: تاثیر k در الگوریتم K-NN

همانطور که در نمودار زیر مشاهده می شود با افزایش k دقت کاهش می یابد که این کاهش دقت به خاطر افزایش بایاس است. نکته دیگری که از این نمودار می توان دریافت کرد کاهش دقت با نزدیک شدن k به مقادیر کوچک است که به دلیل افزایش واریانس است. البته به دلیل ویژگی های مجموعه داده ها این مقدار کاهش دقت چشمگیر نیست و کاهش دقت کمی را می بینیم.



بخش پنجم: مقایسه درخت یک سطحی و چند سطحی

نتایج به شکل زیر هستند.

الگوریتم	دقت	TPR	FPR	Precision
DecisionStump	۶۶,۶۷٪	۰,۶۶۷	۰,۱۶۷	۰,۵۰۰
J48	۹۶٪	۰,۹۶۰	۰,۰۲۰	۰,۹۶۰

ماتریس درهم ریختگی برای الگوریتم‌ها به شکل زیر است.

کلاس	DecisionStump			J48		
	a	b	c	a	b	c
a	۵۰	۰	۰	۴۹	۱	۰
b	۰	۵۰	۰	۰	۴۷	۳
c	۰	۵۰	۰	۰	۲	۴۸

اگر معیارها را به معیارهای خوبی و بدی تقسیم کنیم از معیارهای گفته شده دقت، TPR و Precision معیارهای خوبی هستند یعنی با افزایش این مقادیر یعنی الگوریتم بهتر عمل می‌کند و برعکس FPR معیار بدی است و با افزایش مقدار آن یعنی الگوریتم بدتر عمل کرده است. DecisionStump درختی با یک سطح است. انتظار می‌رود برای مجموعه داده‌های پیچیده درخت‌های چند سطحی بهتر از درخت‌های یک سطحی عمل کنند چون می‌توان

گفت اگر الگوریتم ساخت درخت برای درخت‌های یک‌سطحی و چندسطحی یکسان باشد در بدترین حالت درخت‌های چندسطحی مانند درخت‌های یک‌سطحی عمل خواهند کرد. پس می‌توان انتظار داشت معیارهای خوبی برای J48 که درخت چندسطحی است و معیارهای بدی برای DecisionStump بیشتر باشد. همانطور که مشاهده می‌شود دقت، TPR و Precision الگوریتم J48 بسیار بیشتر از DecisionStump است و مقدار FPR نیز کمتر است. ماتریس درهم‌ریختگی نیز نشان‌گر همین نکته‌ها است.

نکته دیگری که می‌توان از ماتریس درهم‌ریختگی استدلال کرد این است که کلاس a تفاوت بیشتری با کلاس‌های b و c دارد و راحت‌تر می‌توان نمونه‌های مربوط به آن را تشخیص داد. در بین دسته‌بندی‌های اشتباه بیشتر اشتباه‌ها در جداسازی کلاس b و c رخ می‌دهد به‌طوری که در J48 تمامی تشخیص‌های غلط کلاس b به کلاس c و تمامی تشخیص‌های غلط کلاس c به کلاس b نسبت داده شده‌اند. در DecisionStump این قضیه به‌طور جدی‌تری خود را نمایش می‌دهد به این صورت که الگوریتم به خوبی کلاس a را تشخیص می‌دهد ولی قدرت جداسازی کلاس‌های b و c را ندارد و در صورتی که یک نمونه را متعلق به کلاس a نداند آن را به عنوان نمونه‌ای از کلاس b تشخیص می‌دهد.

بخش ششم

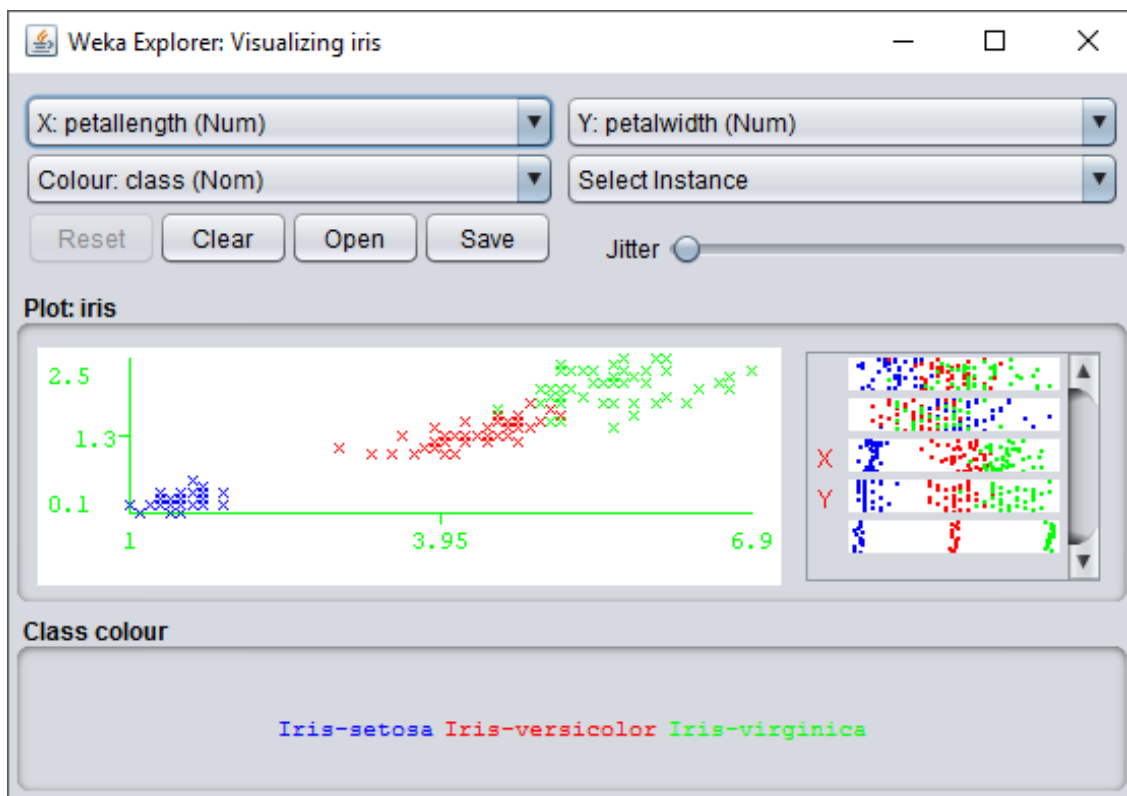
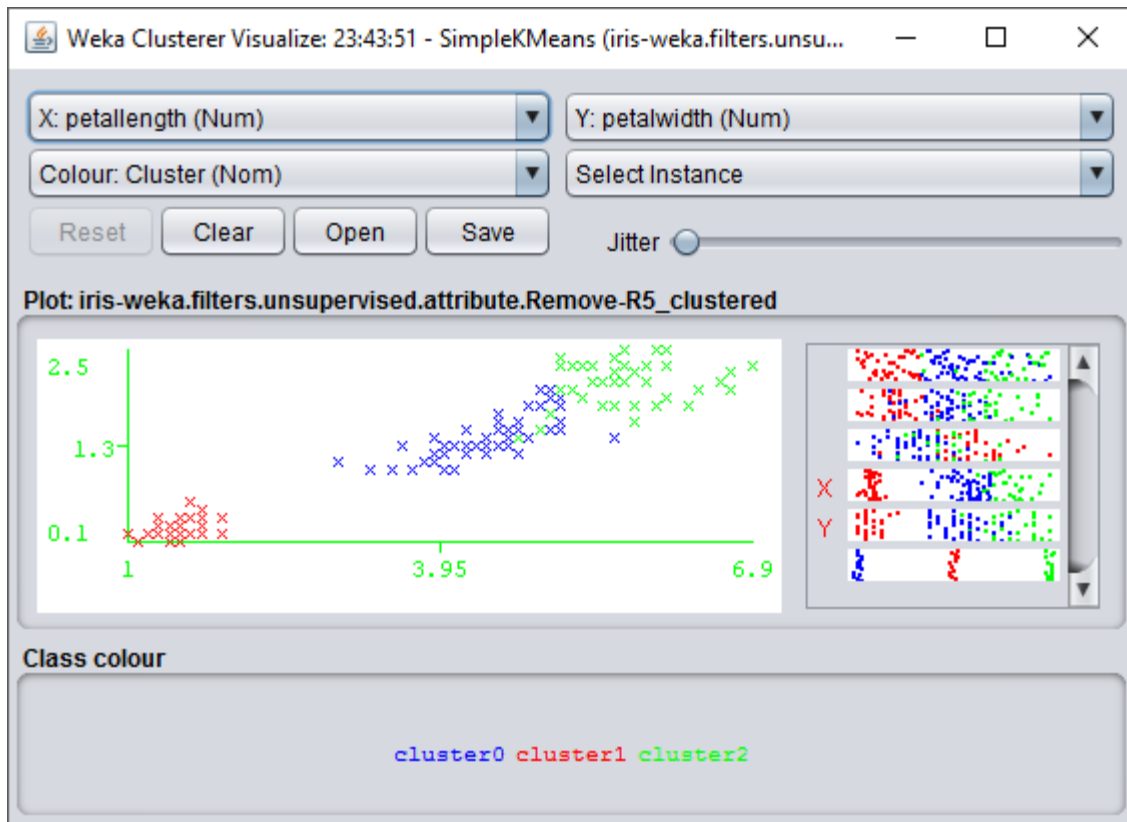
رگرسیون

$$PetalWidth = -0.2103 \times SepalLength + 0.2288 \times SepalWidth + 0.5261 \times PetalLength - 0.2487$$

RMSE برای این مدل در صورت استفاده از روش 10 Fold مقدار ۰,۱۹۶۴ است.

خوشه‌بندی

تصویر اول از تصاویر زیر تصویر به‌دست آمده از خوشه‌بندی است و تصویر دوم نیز تصویر به دست آمده از نمایش داده‌های اصلی است. شباهت واضحی بین دو تصویر وجود دارد. هرچند که برای دو کلاسی که در مرکز و بالا سمت راست نمودار قرار گرفته‌اند در برخی موارد کلاس به درستی تشخیص داده نشده است ولی برای کلاسی که در نزدیکی مبدا قرار گرفته است خوشه‌بندی کاملاً صحیح عمل کرده است.



بخش هفتم: ویژگی‌های ارزشمند

با استفاده از InfoGainAttributeEval اگر از همه مجموعه داده‌ها استفاده کنیم به ترتیب petalwidth و petallength از ۱,۴۱۸ و ۱,۳۷۸ را به خود اختصاص می‌دهند و اگر از 10 Fold استفاده کنیم هر دو رنگ ۱,۵ با خطای ۰,۵ را به دست می‌آورند. با استفاده از CorrelationAttributeEval اگر از همه مجموعه داده‌ها استفاده کنیم همان دو ویژگی مقادیر ۰,۶۱۵ و ۰,۵۹۲ و اگر از 10 Fold استفاده کنیم مقادیر ۰,۶۱۵ و ۰,۵۹۲ را به خود اختصاص می‌دهند.