



روش‌های یادگیری ماشین در پردازش زبان طبیعی

Machine Learning Methods in Natural Language Processing

هادی ویسی

h.veisi@ut.ac.ir

دانشگاه تهران - دانشکده علوم و فنون نوین

نیمسال اول ۱۴۰۱-۱۴۰۰



معرفی درس ...

○ زمان و مکان

- شنبه و دوشنبه، ساعت ۸:۰۰ الی ۱۰:۰۰، دانشکده علوم و فنون نوین

○ وب سایت

- dsp.ut.ac.ir

○ هدف

- مرور روش‌های یادگیری ماشین در پردازش زبان طبیعی
 - مفاهیم یادگیری ماشین
 - مفاهیم پایه آمار و احتمال، نظریه اطلاعات و روش‌های تخمین
 - روش‌های تشابه یابی متن و دسته بندی متن
 - شبکه های عصبی مصنوعی و یادگیری عمیق
 - مرور نمونه کاربردها
- فعالیتهای تمرینی با رویکرد کاربردی



درس

منابع

- Christopher Bishop, Pattern Recognition and Machine Learning, Springer, 2006
- Raschka, Sebastian. *Python machine learning*. Packt Publishing Ltd, 2015.
- هادی ویسی، کبری مفاخری، سعید باقری شورکی، مبانی شبکه های عصبی: معماری، الگوریتمها و کاربردها، انتشارات نص، چاپ پنجم، زمستان ۱۳۹۹
- Laurene Fausette, Fundamentals of neural networks, architecture, algorithms and application, Prentice Hall, 1994 ترجمه
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, Deep Learning, MIT Press, 2016.
- هادی ویسی، مصطفی صالحی، وحید رنجبر، الما جعفری صدر، فرناز صادقی، محمد بحرانی، پردازش زبان و گفتار، انتشارات نویسه پارسی، پاییز ۱۴۰۰
- Daniel Jurafsky, James Martin, Speech and Language Processing, 2nd Edition, Prentice Hall, 2009.



معرفی درس ...

○ ارزیابی ...

• تمرین

- برای هر موضوع: ۴ یا ۵ تمرین
- همفکری و همکاری در یافتن پاسخ سوال‌ها توصیه می‌شود
- در صورت کپی بودن یکی یا چند مورد از پاسخ‌ها، کل نمره آن تمرین برای طرفین کپی در نظر گرفته نمی‌شود.
- تمرین‌های دارای پیاده‌سازی، باید هم شامل کدها و هم شامل گزارش مربوطه باشد
- تاخیر در تحویل
- ارسال پاسخ حداکثر تا ساعت ۲۳:۵۹ مهلت تعیین شده
- هر یک ساعت تاخیر در ارسال پاسخ‌ها (از یک ثانیه تا ۶۰ دقیقه!)، کسر یک درصد نمره آن تمرین به عنوان جریمه تاخیر
- ارسال پاسخ تمرین‌ها
 - تنها به صورت الکترونیکی و به ایمیل استاد درس است.
 - همه فایل‌های مرتبط با یک تمرین را در یک فایل فشرده شده
 - فرمت نام‌گذاری فایل ارسالی: ML4NLP_Family_StNo_HW#
 - وزن تمرین‌های مختلف با هم برابر نیست





معرفی درس ...

○ ارزیابی ...

- آزمونک (کوئیز)
 - نداریم (با توجه به غیرحضور بودن دوره)
- امتحان میان‌ترم
 - دوشنبه ۱۴۰۰/۰۹/۰۱ - ساعت ۸:۰۰
- امتحان پایان‌ترم
 - شامل کلیه مطالب تدریس شده: از جمله مطالب میان‌ترم
 - زمان: طبق برنامه دانشگاه



معرفی درس ...

○ ارزیابی ...

● پروژه: یک مورد

- پروژه کاربردی دارای پیاده‌سازی در Python یا سایر زبان‌های برنامه‌نویسی
- علاوه بر کد برنامه، گزارش مکتوب (به صورت تایپ شده)، داده‌ها و منابع هم تحویل گرفته می‌شود
- در صورت مساعد شدن شرایط کرونایی، تحویل به صورت حضوری است
- آخرین زمان تعیین موضوع پروژه: روز شنبه ۱۴۰۰/۰۸/۰۱
- تحویل پروژه: اولین هفته بعد از آخرین امتحان پایان‌ترم (دوشنبه ۱۴۰۰/۱۱/۱۱)
- موضوع الزاما مرتبط با یادگیری ماشین باشد
- برخی موضوع‌های پیشنهادی



- تشابه‌یابی متن با استفاده از نمایش‌های مبتنی بر یادگیری عمیق (مانند Bert)
- تشخیص احساس در متن با استفاده از یادگیری عمیق
- دسته‌بندی/خوشه‌بندی معنایی کلمات در یادگیری عمیق
- تشخیص گفتار برای تعداد کلمات محدود
- تبدیل متن به گفتار با استفاده از شبکه‌های عمیق مانند مبدل‌ها یا GAN
- تولید خودکار متن (مانند متن یا شعر) با شبکه‌های عصبی عمیق

معرفی درس ...

○ ارزیابی ...

- حضور و مشارکت در کلاس (نمره اضافی)

- مشارکت در بحث‌های کلاس

- مقاله (نمره اضافی)

- مقاله ارسال شده مورد قبول است

- به هر قیمتی مقاله ننویسید!

- همکاری با این درس در نوشتن مقاله را به اطلاع استاد خود برسانید

- بازنگری نمره‌ها و برگه‌ها

- در زمان تحویل پروژه درس (به صورت حضوری)

- دوشنبه ۱۱/۱۱/۱۴۰۰



معرفی درس ...



○ ارزیابی

عنوان	وزن	توضیح
تمرین	۵۰٪	بعد از هر موضوع (وزن تمرین‌ها برابر نیست)
آزمونک (کويز)	-	به دلیل غیرحضورى کلاس‌ها بودن ندارم
امتحان میان‌ترم	۱۵٪	دوشنبه ۱۴۰۰/۰۹/۰۱ ساعت ۸:۰۰
امتحان پایان‌ترم	۲۰٪	از کل مطالب درس، مطابق برنامه دانشگاه
پروژه	۱۵٪	موضوع اختیاری، تعیین موضوع تا شنبه ۱۴۰۰/۰۸/۰۱ تحويل پروژه: اولین هفته بعد از آخرین امتحان پایان‌ترم (دوشنبه ۱۴۰۰/۱۱/۱۱)
حضور و مشارکت کلاس (نمره اضافی)	۵٪	پاسخ دادن به سوالات حین تدریس و مشارکت در بحث‌های کلاس
مقاله (نمره اضافی)	۱۵٪	مقاله ارسال شده به مجله/کنفرانس مورد قبول است



معرفی درس ...

○ سرفصل‌ها ...

- مروری بر مفاهیم و اصول یادگیری ماشین

- مروری بر مبانی آمار و احتمال

- احتمال (توام، شرطی)، امید ریاضی

- قانون بیز

- متغیر تصادفی

- توابع توزیع

- مروری بر نظریه اطلاعات و آنتروپی

- مروری بر روش‌های تخمین

- کمینه میانگین مربعات خطا (MMSE)

- تخمین بیشینه شباهت (MLE)

- تخمین بیز (Bayesian)



معرفی درس ...

○ سرفصل‌ها ...

- بازیابی اطلاعات و تشابه‌یابی متون
 - نمایش کلمات و متن: تبدیل متن به بردار ویژگی
- دسته‌بندی متون با روش بیز (ساده)
- شبکه عصبی مصنوعی
 - مبانی و مفاهیم
 - شبکه عصبی پرسپترون
 - شبکه عصبی پرسپترون چندلایه (MLP)
 - نمایش کلمات/جمله/سند با بردار کلمات
 - یادگیری عمیق
 - شبکه خودرمزگذار، پیچشی (CNN) و شبکه مولد مقابله‌ای (GAN)
 - شبکه‌های عصبی بازگشتی (RNN)
 - شبکه حافظه کوتاه مدت ماندگار (LSTM)
 - مفهوم توجه (Attention)
 - مبدل‌ها (Transformer)



معرفی درس ...

○ سرفصل‌ها

- مدل مخفی مارکوف (HMM)
 - کاربرد در برچسپ‌زنی اجزای کلام (POS: Part-of-Speech tagging)
 - کاربرد در تشخیص گفتار (Speech Recognition)
- روش‌های خوشه‌بندی
 - روش k-میانگین
 - الگوریتم امید-بیشینه (EM)



معرفی درس

○ زمان بندی

هفته	تاریخ	موضوع	توضیحات
۱	۱۴۰۰/۰۷/۰۵ و ۰۳	معرفی درس	
		مروری بر مفاهیم یادگیری ماشین	
۲	۱۴۰۰/۰۷/۱۲ و ۱۰	مروری بر مفاهیم یادگیری ماشین	
۳	۱۴۰۰/۰۷/۱۹ و ۱۷	مروری بر مبانی آمار و احتمال	تمرین
۴	۱۴۰۰/۰۷/۲۶ و ۲۴	مروری بر نظریه اطلاعات و روش‌های تخمین	
۵	۱۴۰۰/۰۸/۰۳ و ۰۱	بازیابی اطلاعات و تشابه‌یابی متون	
۶	۱۴۰۰/۰۸/۱۰ و ۰۸	بازیابی اطلاعات و تشابه‌یابی متون	اعلام موضوع پروژه
۷	۱۴۰۰/۰۸/۱۷ و ۱۵	دسته‌بندی متون	تمرین
۸	۱۴۰۰/۰۸/۲۴ و ۲۲	شبکه عصبی مصنوعی	
۹	۱۴۰۰/۰۸/۲۹ و ۱۴۰۰/۰۹/۰۱	شبکه عصبی پرسپترون چندلایه	آزمون میان ترم
۱۰	۱۴۰۰/۰۹/۰۸ و ۰۶	شبکه عصبی عمیق	تمرین
۱۱	۱۴۰۰/۰۹/۱۵ و ۱۳	شبکه عصبی عمیق: CNN	
۱۲	۱۴۰۰/۰۹/۲۲ و ۲۰	شبکه عصبی عمیق: GAN	
۱۳	۱۴۰۰/۰۹/۲۹ و ۲۷	شبکه عصبی بازگشتی و LSTM	
۱۴	۱۴۰۰/۱۰/۰۶ و ۰۴	مفهوم توجه در شبکه عصبی و مبدل‌ها	تمرین
۱۵	۱۴۰۰/۱۰/۱۳ و ۱۱	مدل مخفی مارکوف (HMM)	
۱۶	۱۴۰۰/۱۰/۲۰ و ۱۸	مدل مخفی مارکوف (HMM): کاربردها	

۱۴۰۰/۰۷/۰۵: اربعین

متناسب با شرایط و سطح کلاس، و همچنین تغییرات پیش‌بینی نشده در زمان بندی، ممکن است سرفصل مطالب و یا زمان بندی‌های کلاس مقداری تغییر داشته باشد