

دانشکده علوم و فنون نوین

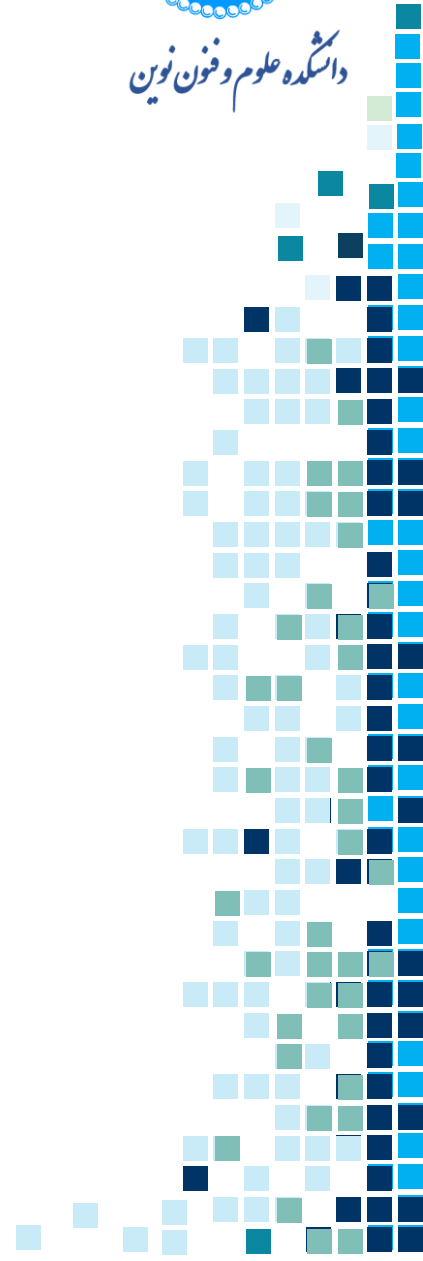
روش‌های یادگیری ماشین در پردازش زبان طبیعی

Machine Learning Methods in Natural Language Processing

هادی ویسی

h.veisi@ut.ac.ir

نیم‌سال اول ۱۴۰۲-۱۴۰۳



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

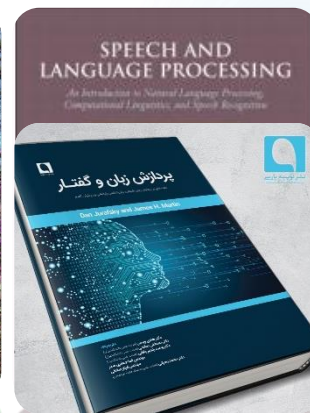
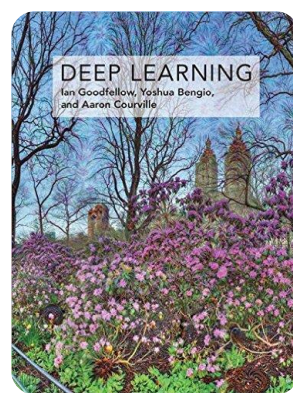
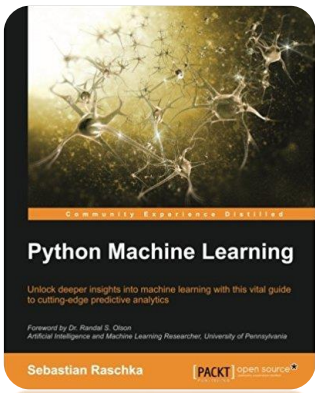
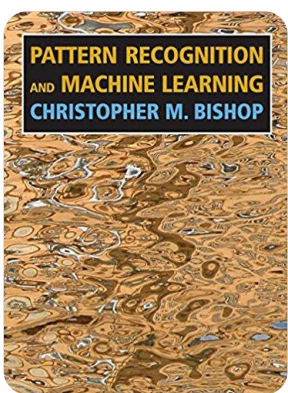
◀ مفاهیم پایه آمار و احتمال، نظریه اطلاعات و روش‌های تخمین

◀ روش‌های سنتی یادگیری ماشین

◀ شبکه‌های عصبی مصنوعی و یادگیری عمیق

◀ مرور نمونه کاربردها

◀ فعالیت‌های تمرینی با رویکرد کاربردی



دانشگاه علم و فنون نون

منابع

◀ Christopher Bishop, Pattern Recognition and Machine Learning, Springer, ۲۰۰۶

◀ Raschka, Sebastian. *Python machine learning*. Packt Publishing Ltd, ۲۰۱۵.

◀ هادی ویسی، کبری مفاخری، سعید باقری شورکی، مبانی شبکه های عصبی: معماری، الگوریتمها و کاربردها، انتشارات نص، چاپ پنجم، زمستان ۱۳۹۹

◀ ترجمه Laurene Fausette, Fundamentals of neural networks, architecture, algorithms and application, Prentice Hall, 1994

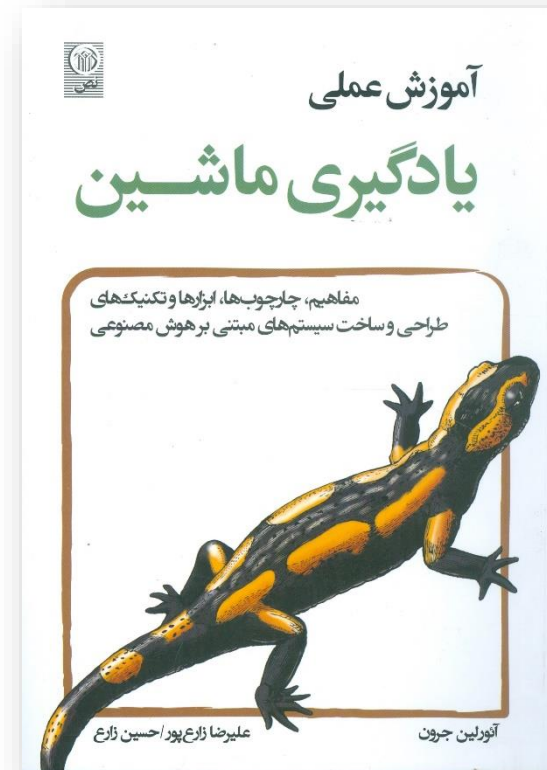
◀ Ian Goodfellow, Yoshua Bengio, Aaron Courville, Deep Learning, MIT Press, ۲۰۱۶.

◀ 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 ...

منابع □



◀ Géron, Aurélien. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems.* " O'Reilly Media, Inc.", ۲۰۱۹

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
- ◀ دیرکرد در تحویل

- ارسال پاسخ حداکثر تا ساعت ۲۳:۵۹ مهلت تعیین شده
- هر یک ساعت دیرکرد در ارسال پاسخها (از یک ثانیه تا ۶۰ دقیقه!)، کسر یک درصد نمره آن تمرین به عنوان جریمه دیرکرد
- امکان بخشودگی یک مورد دیرکرد (برای یک تمرین) حداکثر به اندازه یک روز (۲۴ ساعت) به انتخاب دانشجو

◀ 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 (کويز)

◀ یکی یا دو سوال در هر بار

◀ ممکن است بدون اطلاع قبلی باشد

◀ امتحان میان‌ترم

◀ دوشنبه ۱۴۰۲/۰۸/۲۹ ساعت ۸:۰۰ (حضوری)

◀ امتحان پایان‌ترم

◀ شامل کلیه مطالب تدریس شده: از جمله مطالب میان‌ترم

◀ زمان: طبق برنامه دانشگاه

◀ بازنگری نمره‌ها و برگه‌ها

◀ در زمان تحویل پروژه درس: اولین هفته بعد از آخرین امتحان پایان‌ترم



معرفی درس ...

□ ارزیابی . . .

◀ پروژه درس (اختیاری، نمره اضافی)

◀ پروژه کاربردی دارای پیاده‌سازی در Python یا سایر زبان‌های برنامه‌نویسی

◀ موضوع الزاما مرتبط با مطالب درس باشد

◀ آخرین مهلت انتخاب موضوع: ۱۴۰۲/۰۹/۰۱

◀ تحویل به صورت حضوری است

◀ موارد تحویل دادنی:

● کلیه کدهای پروژه

● گزارش مکتوب (به صورت تایپ شده) شامل توضیح روش و جزئیات پیاده‌سازی و نتایج بدست آمده و تحلیل‌های مربوطه

● داده‌های مورد استفاده در پروژه

● مقاله‌ها و منابع مورد استفاده

◀ بارم‌بندی نمرات:

● انجام درست پیاده‌سازی و مرتب بودن کدها: ۵۰٪

● کامل بودن گزارش (شامل نحوه استفاده از کد و مبانی علمی کار) و رعایت اصول نگارشی در آن: ۲۵٪

● ارائه نتایج و تحلیل آن (در گزارش): ۲۵٪

◀ تحویل پروژه: اولین هفته بعد از آخرین امتحان پایان‌ترم



معرفی درس ...

□ ارزیابی ...

◀ ارائه شفاهی

- ◀ هدف: آشنایی با مطالب بهروز در حوزه درس
- ◀ انتخاب یک موضوع مرتبط با مطالب درس
- ◀ مطالعه منابع لازم و انجام یک ارائه کوتاه در کلاس
 - منابع متعلق به سه سال اخیر
- ◀ زمان هر ارائه ۲۰ تا ۲۵ دقیقه
- ◀ انتخاب موضوع با هماهنگی استاد
- ◀ برخی موضوع های پیشنهادی
 - هوش مصنوعی اعتمادپذیر و لزوم آن در پردازش زبان طبیعی
 - مدل های زبانی بزرگ (LLM): روش ساخت و کاربردها
 - مروری بر Bard و ChatGPT و نحوه ساخت آنها
 - مروری بر کتابخانه ها و ابزارهای مدرن در پردازش گفتار



معرفی درس ...

□ ارزیابی ...

◀ پروژه درس (اختیاری، نمره اضافی)

◀ 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
- درک معنا با استفاده از مدل‌های زبانی بزرگ (LLM)



معرفی درس ...



□ ارزیابی

عنوان	وزن	توضیح
تمرین	۵۰%	بعد از هر موضوع (وزن تمرین‌ها برابر نیست)
آزمونک (کوین)	۵%	ممکن است بدون اعلام قبلی باشد.
آزمون میان‌ترم	۲۰%	دوشنبه ۲۹/۰۸/۱۴۰۲ ساعت ۱۰:۰۰ (حضور)
آزمون پایان‌ترم	۲۰%	از کل مطالب درس، مطابق برنامه دانشگاه
ارائه شفاهی	۵%	ارائه کلاسی از یک موضوع به‌روز
پروژه (نمره اضافی)	۱۰%	موضوع اختیاری، مرتبط با مطالب درس (آخرین مهلت انتخاب موضوع: ۰۱/۰۹/۱۴۰۲) تحویل پروژه: اولین هفته بعد از آخرین امتحان پایان‌ترم

معرفی درس . . .



دستیار آموزشی

◀ امیرمحمد کویشپور

a.m.kouyeshpour@ut.ac.ir ◀



◀ سیاوش حسینپور صفاریان

siavash.saffaria@ut.ac.ir ◀

◀ برگزاری کلاس‌های حل تمرین و رفع اشکال

◀ به ویژه راهنمایی پیاده‌سازی و کدنویسی

معرفی درس ...

□ سرفصل‌ها . . .

◀ مروری بر مفاهیم و اصول یادگیری ماشین

◀ مروری بر مبانی آمار و احتمال

◀ احتمال (توأم، 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

◀ روش‌های سنتی یادگیری ماشین

◀ بیز ساده

◀ نزدیک‌ترین همسایه

◀ رگرسیون

◀ درخت تصمیم و جنگل تصادفی

◀ ماشین بردار پشتیبان (SVM)

معرفی درس ...

□ سرفصل‌ها ...

◀ شبکه عصبی مصنوعی و یادگیری عمیق

◀ مبانی و مفاهیم

◀ شبکه عصبی پرسپترون و آدالین

◀ 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)

● سازوکار توجه (Attention)

◀ مدل‌ها (Transformer)

● BERT

● GPT

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

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

توضیح	موضوع	تاریخ	هفته
	معرفی درس	۱۴۰۲/۰۷/۰۳ و ۰۱	۱
	مروری بر مفاهیم یادگیری ماشین	۱۴۰۲/۰۷/۱۰ و ۰۸	۲
	مروری بر مبانی آمار و احتمال	۱۴۰۲/۰۷/۱۷ و ۱۵	۳
آزمونک	مروری بر نظریه اطلاعات و روش های تخمین	۱۴۰۲/۰۷/۲۴ و ۲۲	۴
	روش های سنتی: بیز ساده و KNN و رگرسیون	۱۴۰۲/۰۷/۲۹	۵
تمرین	روش های سنتی: درخت تصمیم و SVM	۱۴۰۲/۰۸/۰۱	۵
	روش های نمایش کلمات و اسناد: بازیابی	۱۴۰۲/۰۸/۰۸ و ۰۶	۶
آزمونک	اطلاعات و تشابه یابی متون	۱۴۰۲/۰۸/۱۵ و ۱۳	۷
	شبکه عصبی مصنوعی	۱۴۰۲/۰۸/۲۲ و ۲۰	۸
تمرین	مبانی، پرسپترون و آدالاین	۱۴۰۲/۰۸/۲۹ و ۲۷	۹
	شبکه عصبی پرسپترون چندلایه	۱۴۰۲/۰۹/۰۶ و ۰۴	۱۰
	کاربردها و Wav2Vec	۱۴۰۲/۰۹/۱۳ و ۱۱	۱۱
	شبکه عصبی عمیق: خودرمزگذارها	۱۴۰۲/۰۹/۲۰ و ۱۸	۱۲
	شبکه عصبی عمیق: GAN و CNN	۱۴۰۲/۰۹/۲۷ و ۲۵	۱۳
	شبکه عصبی بازگشتی	۱۴۰۲/۱۰/۰۴ و ۰۲	۱۴
تمرین	LSTM و سازوکار توجه	۱۴۰۲/۱۰/۱۱ و ۰۹	۱۵
	مبدل ها، Bert و GPT	۱۴۰۲/۱۰/۱۸ و ۱۶	۱۶
آزمونک	ارائه دانشجویان	۱۴۰۲/۱۰/۲۵ و ۲۳	۱۷
تمرین	مدل مخفی مارکوف (HMM)		
	ارائه دانشجویان		
	مدل مخفی مارکوف (HMM): کاربردها		
	روش های خوشه بندی		
	جبرانی (در صورت نیاز)		

اعلام موضوع پروژه

آزمون میان ترم