

# روش‌های آماری در پردازش زبان طبیعی

(۱۶۲-۰۵-۸۳)

## Statistical Methods in Natural Language Processing

هادی ویسی

[h.veisi@ut.ac.ir](mailto:h.veisi@ut.ac.ir)

دانشگاه تهران - دانشکده علوم و فنون نوین



## معرفی درس ...

### ○ زمان و مکان

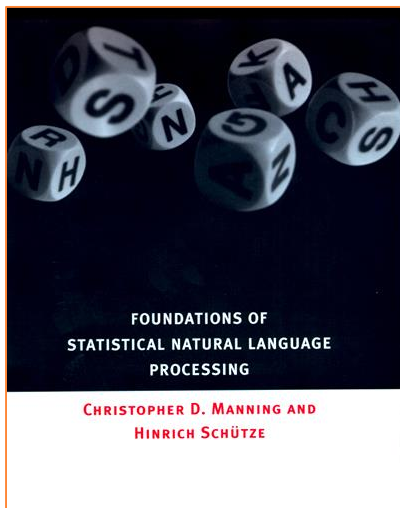
- شنبه و دوشنبه، ساعت ۱۱:۰۰ الی ۱۲:۳۰
- دانشکده علوم و فنون نوین - ساختمان ۲۳

### ○ وب سایت

- <http://dsp.ut.ac.ir/courses/statistical-nlp-fall1394>

### ○ هدف

- مرور روش‌های یادگیری ماشین در پردازش زبان طبیعی (با تاکید بر روش‌های آماری)
  - مفاهیم پایه آمار و احتمال، نظریه اطلاعات و روش‌های تخمین
  - مفاهیم یادگیری ماشین
  - روش‌های مدل‌سازی آماری
  - مرور نمونه کاربردها
- فعالیت‌های تمرینی با رویکرد کاربردی



- Chris Manning and Hinrich Schütze, Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.
- Daniel Jurafsky, James Martin, Speech and Language Processing, 2nd Edition, Prentice Hall, 2009.
- Bishop, Christopher M. Pattern recognition and machine learning. springer, 2006.
- Igor Bolshakov, Alexander Gelbukh, Computational Linguistics, Models, Resources, Applications, 2004.

## معرفی درس ...

### ○ ارزیابی ...

#### • تمرین

- برای هر موضوع
- همفکری و همکاری در یافتن پاسخ سوال‌ها توصیه می‌شود
- در صورت کپی بودن یکی یا چند مورد از پاسخ سوالات یک تمرین، کل نمره آن تمرین 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
- در صورت داشتن یک روز تاخیر در ارسال پاسخ‌ها (از یک ثانیه تا ۲۴ ساعت!)، ۲۵٪ نمره آن تمرین به عنوان جریمه تاخیر
- در صورت تاخیر دو روزه ۵۰٪ نمره مربوطه به عنوان جریمه تاخیر
- پس از دو روز نمره‌ای در نظر گرفته نمی‌شود



- وزن تمرین‌های مختلف با هم برابر نیست

## معرفی درس ...

### ○ ارزیابی ...

- آزمونک (کويز)
  - از مطالب هر بخش
  - ممکن است بدون اطلاع قبلی باشد

### • امتحان میان‌ترم

- نداریم!

### • امتحان پایان‌ترم

- شامل کلیه مطالب تدریس شده
- زمان: ??



## معرفی درس ...

### ○ ارزیابی ...

#### ● پروژه: یک مورد

- پروژه کاربردی دارای پیاده‌سازی در MATLAB، Python یا سایر زبان‌های برنامه‌نویسی
- علاوه بر کد برنامه، گزارش مکتوب (به صورت تایپ شده) تحویل گرفته می‌شود
- تحویل حضوری

- آخرین زمان تعیین موضوع پروژه: روز دوشنبه ۱۳۹۴/۰۷/۲۷
- تحویل پروژه: ۵ روز بعد از آخرین امتحان پایان‌ترم

#### ● مقاله (نمره اضافی)

- مقاله آماده انتشار مورد قبول است
- به هر قیمتی مقاله ننویسید!

#### ● بازنگری نمره‌ها و برگه‌ها

- در زمان تحویل پروژه درس (به صورت حضوری)



## معرفی درس ...



### ○ ارزیابی

عنوان	وزن	توضیح
تمرین	۴۰٪	بعد از هر موضوع (وزن تمرین‌ها برابر نیست)
امتحان کوتاه (کويز)	۱۵٪	ممکن است بدون اطلاع قبلی باشد
امتحان پایان ترم	۳۰٪	از کل مطالب درس روز ??/۱۳۹۴/۱۰ ساعت ??:۰۰
پروژه	۱۵٪	موضوع اختیاری تعیین موضوع تا دوشنبه ۲۷/۰۷/۱۳۹۴ تحويل پروژه: ۵ روز بعد از آخرین امتحان پایان ترم
مقاله (نمره اضافی)	۱۵٪	مقاله آماده انتشار مورد قبول است



## معرفی درس ...

### ○ سرفصل‌ها ...

- 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)
- مروری بر مفاهیم و اصول یادگیری ماشین
- مدل مخفی مارکوف (HMM)
  - کاربرد در برچسپ‌زنی اجزای کلام (POS: Part-of-Speech tagging)
- تبدیل متن به بردار ویژگی





## معرفی درس ...

### ○ سرفصل‌ها

- روش‌های خوشه‌بندی
  - روش  $k$ -میانگین و الگوریتم امید-بیشینه (EM)
- روش ییز (ساده)
- شبکه عصبی مصنوعی
  - مبانی و مفاهیم
  - شبکه عصبی پرسپترون چندلایه (MLP)
  - شبکه عصبی احتمالاتی (PNN)
  - شبکه عصبی خودسازمان‌ده (SOM)
- گرامر مستقل از متن احتمالاتی (PCFG)
- نمونه کاربردها
  - دسته‌بندی متون
  - مروری بر بازیابی اطلاعات در متن
  - مروری بر بازشناسی گفتار
  - مروری بر ترجمه ماشینی



## معرفی درس

### ○ زمان‌بندی

توضیحات	موضوع	تاریخ	هفته
	معرفی درس و مروری بر مبانی آمار و احتمال	۲۳/۰۶/۱۳۹۴ و ۲۱	۱
	مروری بر نظریه اطلاعات و روش‌های تخمین	۳۰/۰۶/۱۳۹۴ و ۲۸	۲
تمرین، کویز	مروری بر مفاهیم یادگیری ماشین	۰۶/۰۷/۱۳۹۴ و ۰۴	۳
	مدل مخفی مارکوف (HMM)	۱۳/۰۷/۱۳۹۴ و ۱۱	۴
تمرین، کویز	کاربرد مدل مخفی مارکوف در POS Tagging	۲۰/۰۷/۱۳۹۴ و ۱۸	۵
اعلام موضوع پروژه	تبدیل متن به بردار ویژگی	۲۷/۰۷/۱۳۹۴ و ۲۵	۶
	خوشه‌بندی	۰۴/۰۸/۱۳۹۴ و ۰۲	۷
	روش E-M و روش بیز ساده	۱۱/۰۸/۱۳۹۴ و ۰۹	۸
تمرین، کویز	شبکه عصبی مصنوعی	۱۸/۰۸/۱۳۹۴ و ۱۶	۹
	شبکه عصبی پرسپترون چندلایه	۲۵/۰۸/۱۳۹۴ و ۲۳	۱۰
	شبکه عصبی احتمالاتی	۳۰/۰۸/۱۳۹۴	۱۱
	شبکه عصبی خودسازمان‌ده	۰۲/۰۹/۱۳۹۴	
تمرین، کویز	گرامر مستقل از متن احتمالاتی (PCFG)	۰۹/۰۹/۱۳۹۴ و ۰۷	۱۲
	دسته‌بندی متون	۱۶/۰۹/۱۳۹۴ و ۱۴	۱۳
	بازیابی اطلاعات در متن	۲۳/۰۹/۱۳۹۴ و ۲۱	۱۴
تمرین، کویز	مروری بر بازشناسی گفتار	۳۰/۰۹/۱۳۹۴ و ۲۸	۱۵
	مروری بر ترجمه ماشینی	۰۷/۱۰/۱۳۹۴ و ۰۵	۱۶