



به نام خدا



کارگاه علم داده با پایتون پیشرفته

جلسه اول: مقدمه ای بر یادگیری ماشین و مروری بر کتابخانه های Numpy
و Pandas

مدرس :

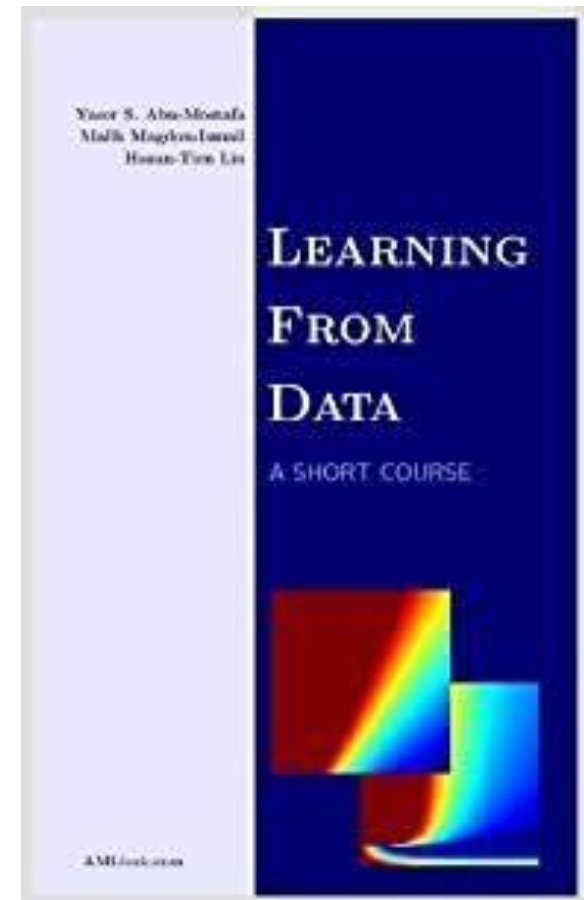
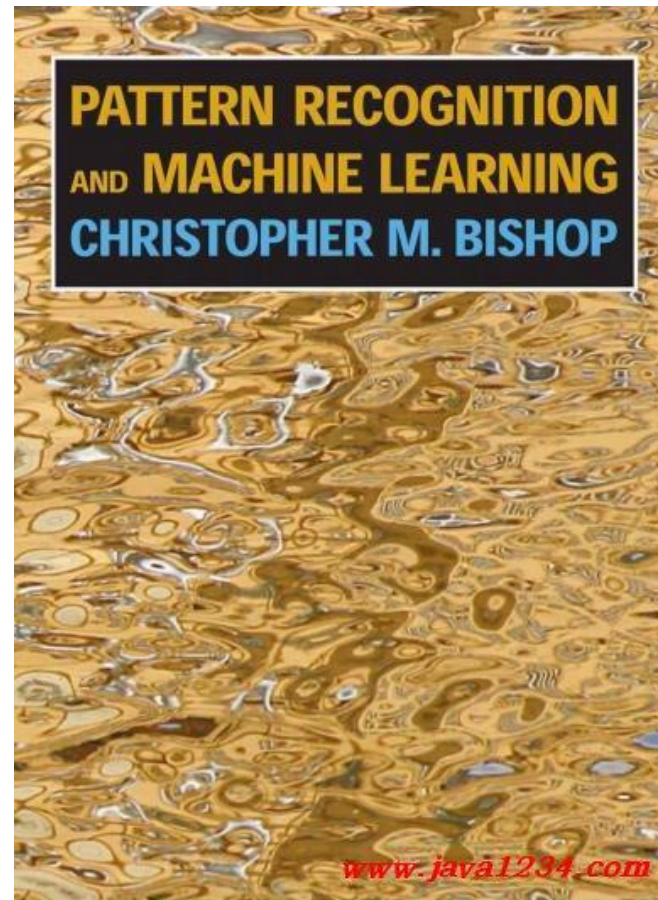
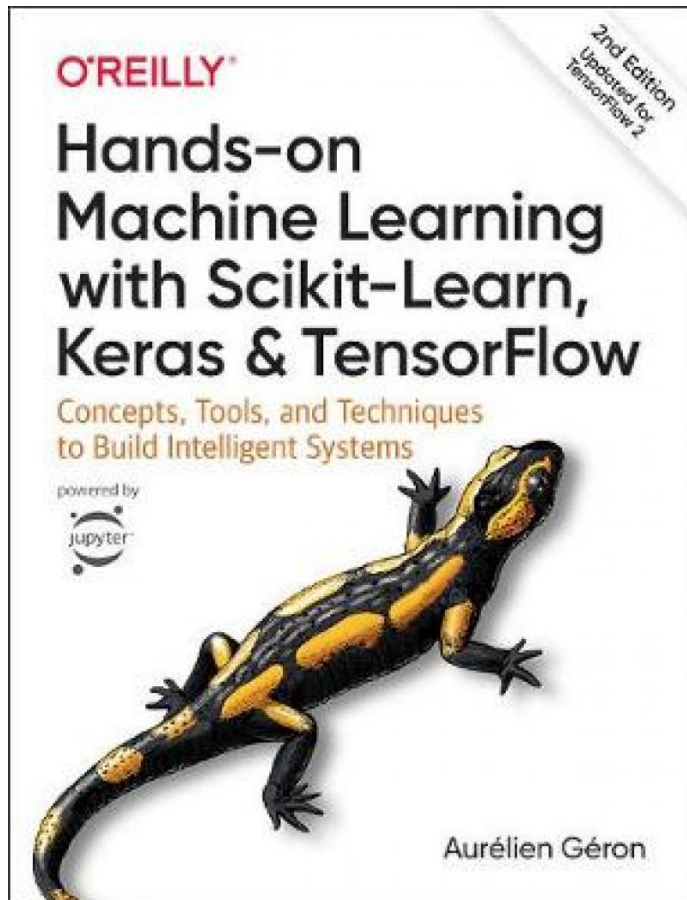
مهرناز جلیلی

دانشجو کارشناسی ارشد علم داده ها

دانشگاه شهید بهشتی



منابع



سرفصل ها

درس ۱: مقدمه ای بر یادگیری ماشین (تئوری)

مروری بر کتابخانه های Numpy و Pandas

درس ۲: کار با دیتا فریم ها در پکیج Pandas (عملی)
پاک سازی داده ها

درس ۳: رگرسیون خطی ساده و چندگانه (تئوری و عملی)

درس ۴: طبقه بندی یا classification (تئوری)
K نزدیک ترین همسایه (KNN) (تئوری و عملی)

درس ۵: رگرسیون لجستیک (تئوری و عملی)

درس ۶: درخت تصمیم (تئوری و عملی)

درس ۷: ماشین بردار پشتیبان (SVM) (تئوری و عملی)

درس ۸: تقلیل ابعاد (تئوری و عملی)

درس ۹: کشف داده های پرت (تئوری و عملی)

درس ۱۰: Case study

مقدمه ای بر یادگیری ماشین

یادگیری ماشین در ابتدا زیر شاخه ای از هوش مصنوعی محسوب میشد ولی در حال حاضر بهترین رویکرد برای حل مسائل هوش مصنوعی است.



تشخیص صوت

(Speech Recognition)

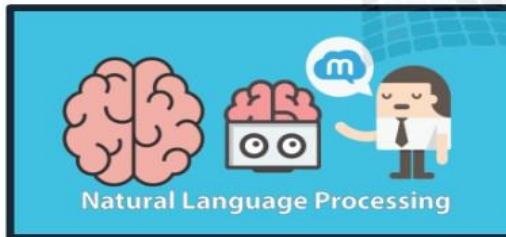
تشخیص کلمات سیگنال صوتی



بینایی ماشین

(computer vision)

تشخیص اشیاء داخل تصویر



پردازش زبان طبیعی



Robotics

برای این که بتوان یک مساله را با تکنیک های یادگیری ماشین حل کرد:

۱- الگویی باید وجود داشته باشد

۲- این الگو توسط ریاضی شناخته شده باشد

۳- داده ای داشته باشیم

Clump	UnifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class
5	1	1	1	2	1	3	1	1	benign
5	4	4	5	7	10	3	2	1	benign
3	1	1	1	2	2	3	1	1	malignant
6	8	8	1	3	4	3	7	1	benign
4	1	1	3	2	1	3	1	1	benign
8	10	10	8	7	10		7	1	malignant
1	1	1	1	2	10	3	1	1	benign
2	1	2	H	2	1	3	1	1	benign
2	1	1	1	2	1	1	1	5	benign
4	2	1	1	2	1	2	1	1	benign

انواع الگوریتم های یادگیری

۱- یادگیری با ناظر (Supervised Learning)

۲- یادگیری بدون ناظر (Unsupervised Learning)

۳- یادگیری نیمه نظارتی (Semi-supervised Learning)

۴- یادگیری تقویتی (Reinforcement Learning)

Supervised vs. Unsupervised

Supervised Learning

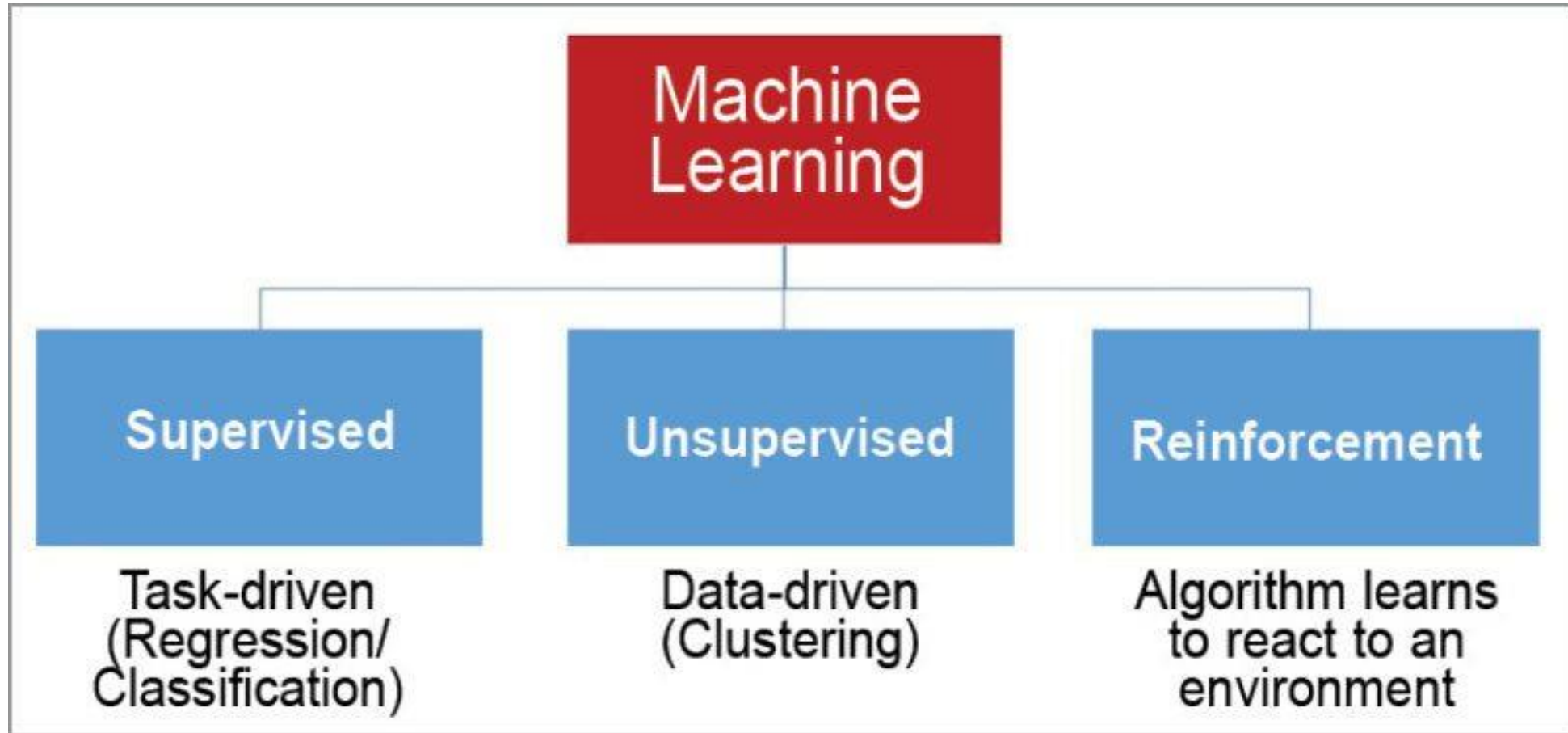
- **Classification:**
Classifies labeled data
- **Regression:**
Predicts trends using previous labeled data
- Has more evaluation methods than unsupervised learning
- Controlled environment

Unsupervised Learning

- **Clustering:**
Finds patterns and groupings from unlabeled data
- Has fewer evaluation methods than supervised learning
- Less controlled environment

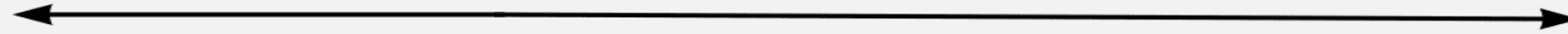


اصطلاحات و مفاهیم آن ها



Train & Test

Dataset

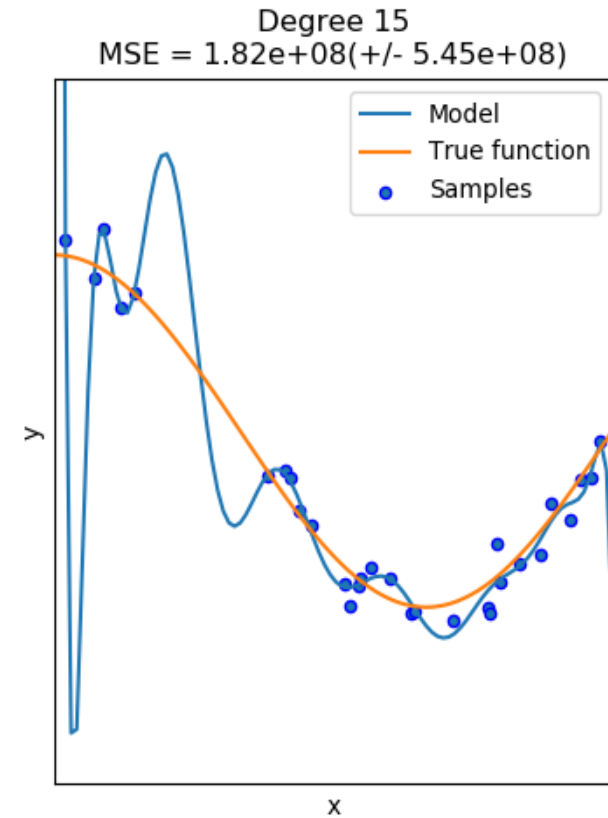
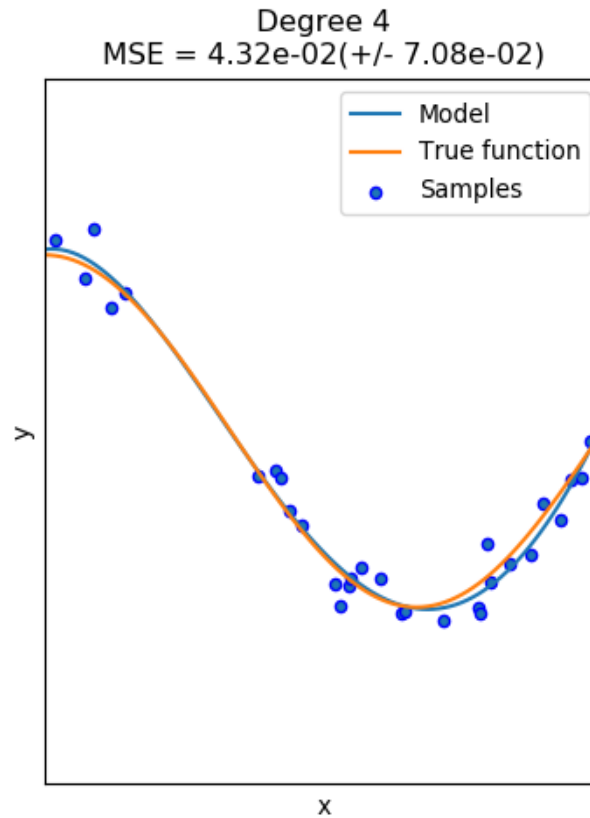
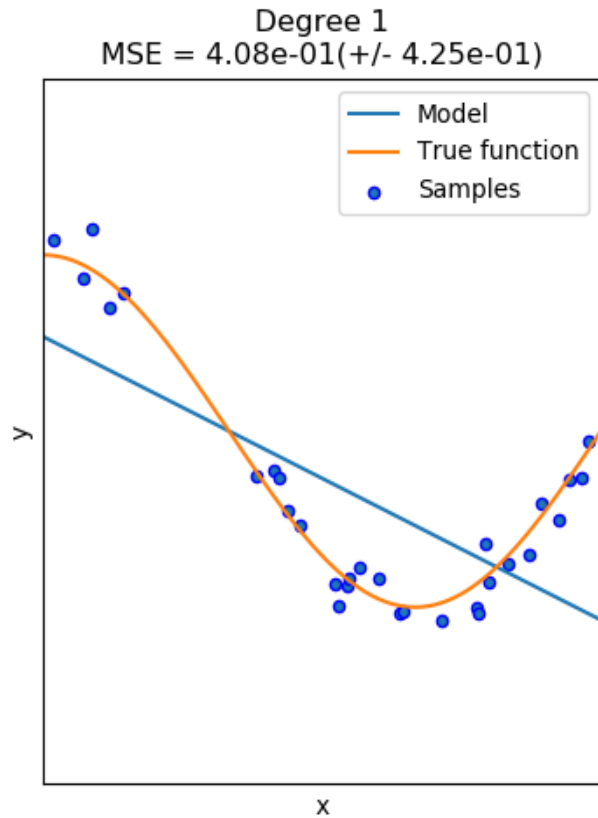


Training Dataset

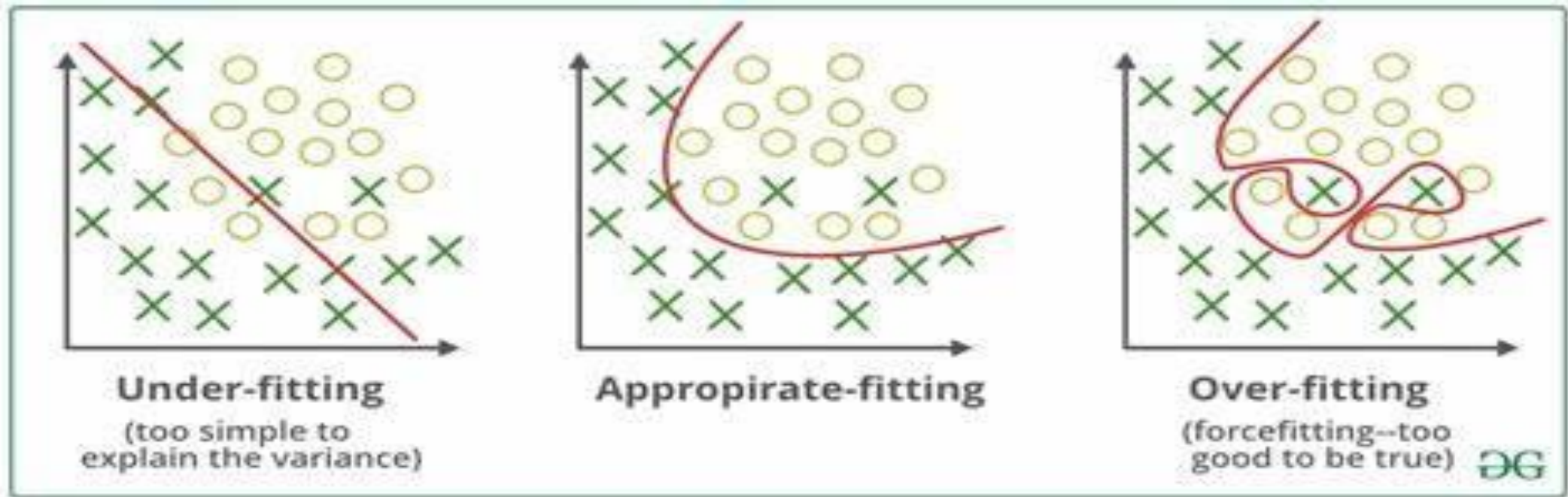
Test Dataset

dataaspirant.com

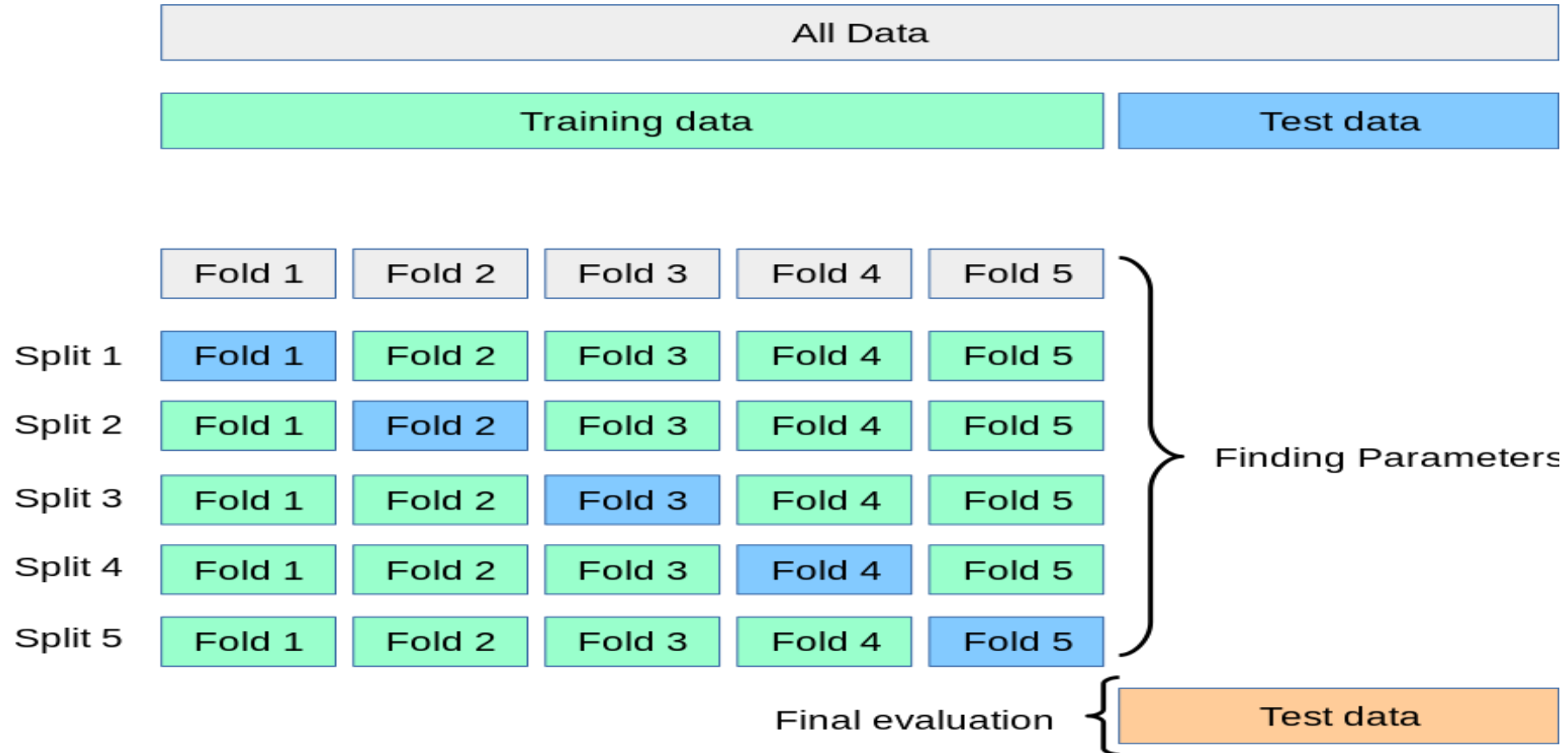
Overfit vs. Underfit



Overfit vs. Underfit



Cross-validation





ارزیابی سیستم های یادگیری

ماتریس درهم ریختگی (confusion matrix)

		PREDICTED	
		Positive	Negative
ACTUAL	Positive	TRUE POSITIVE	FALSE NEGATIVE
	Negative	FALSE POSITIVE	TRUE NEGATIVE

dataaspirant.com

$$Accuracy = \frac{T_p + T_n}{T_p + T_n + F_p + F_n}$$

$$Precision = \frac{T_p}{T_p + F_p}$$

$$Recall = \frac{T_p}{T_p + T_n}$$

$$F_1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

محاسبه کنید

Actual \ Predicted	Spam +	Real -
Spam +	1558	255
Real -	125	2633

Accuracy = 92%

Precision = 92.5%

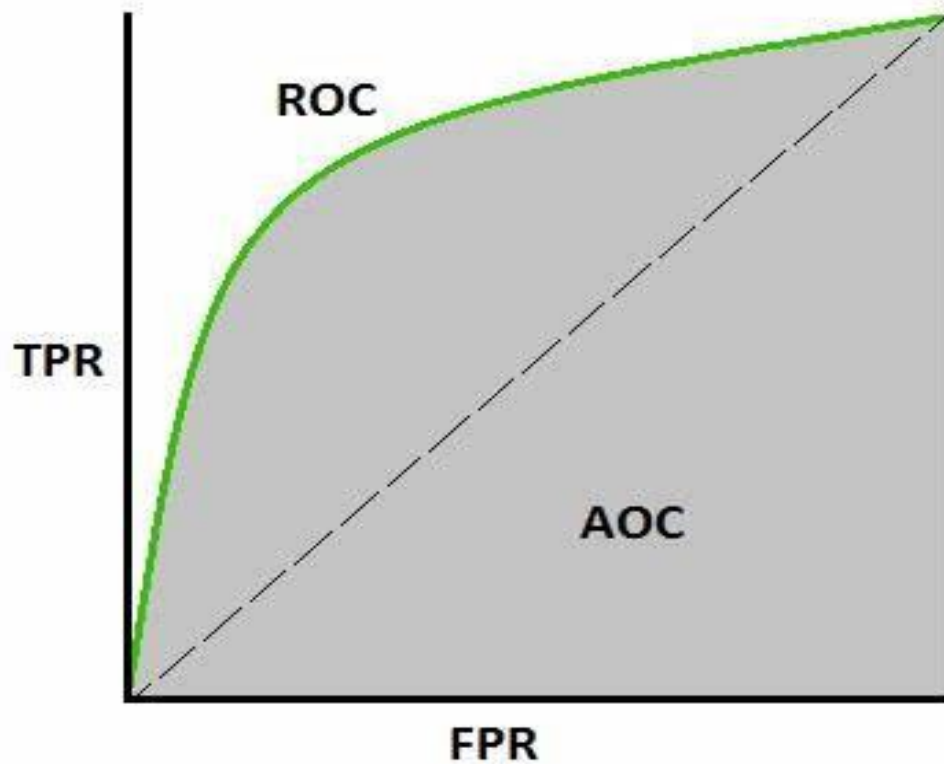
Actual \ Predicted	Cancer = yes	Cancer = no
Cancer = yes	90	210
Cancer = no	140	9560

Accuracy = 96%

recall = 30%

منحنی Roc

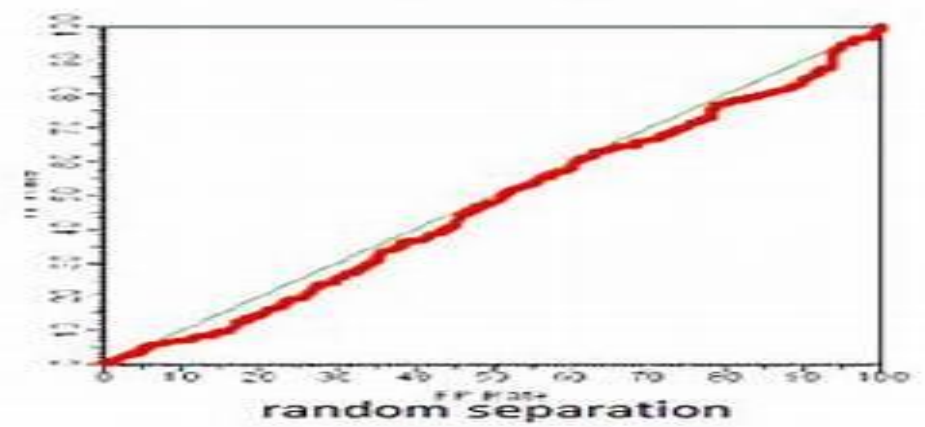
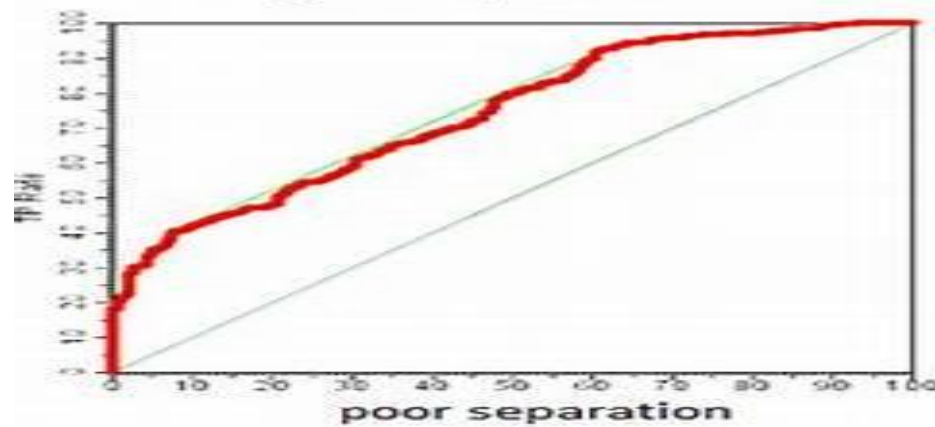
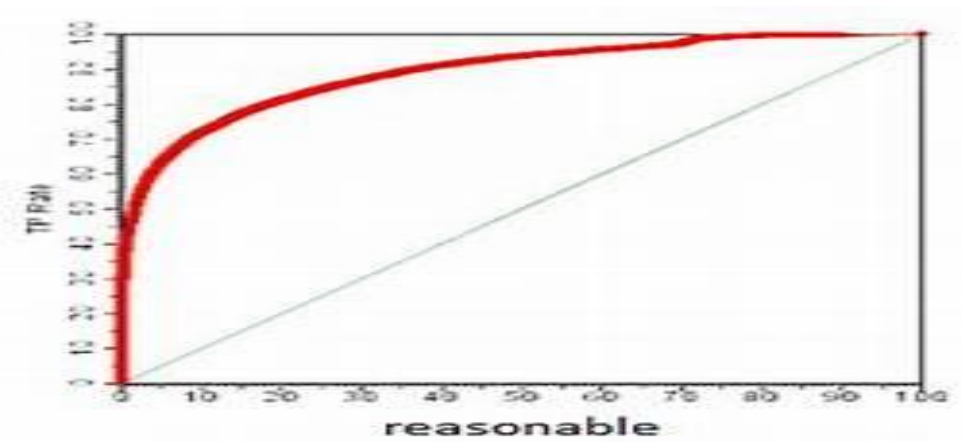
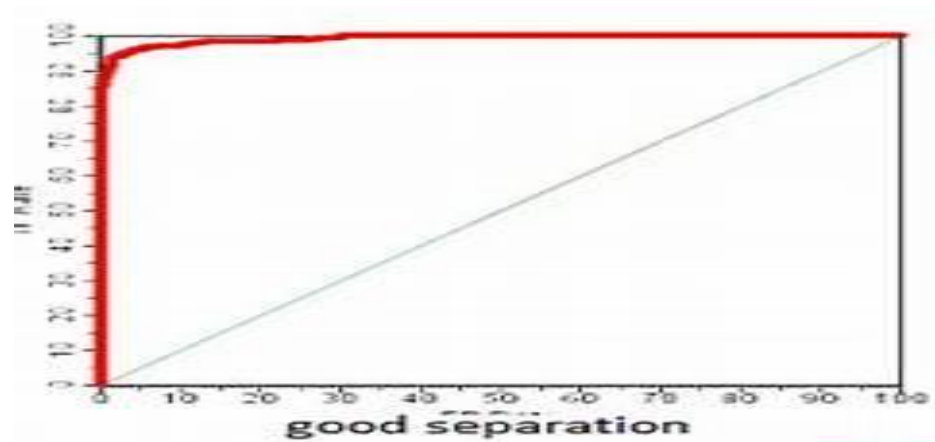
Receiver Operating Characteristic Curve



در نمودار ROC هر دوشاخص «حساسیت» (Sensitivity) یا «بازیابی» (Recall) ترکیب شده و به صورت یک منحنی نمایش داده می‌شوند.

اغلب برای بررسی کارایی الگوریتم‌های دسته‌بندی از منحنی ROC استفاده می‌کنند. این موضوع در شاخه یادگیری ماشین با نظارت بیشتر مورد توجه قرار گرفته است

مقایسه





Numpy & Pandas



برای دانلود و مشاهده کدها و اسلایدها به لینک گیت‌هاب زیر مراجعه کنید:

[mehrjy/Advanced-python-course \(github.com\)](https://github.com/mehrjy/Advanced-python-course)