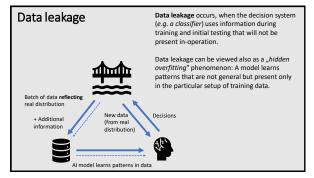
Mechatronic Engineering program: Python for machine learning and data science

Data and models interpretation

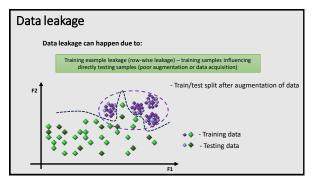
Ziemowit Dworakowski

AGH University of Krakow

1

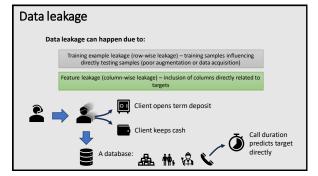


2

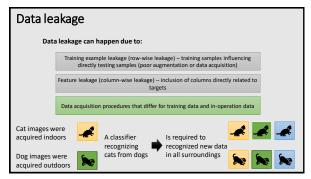


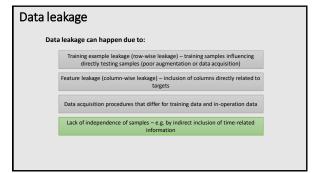


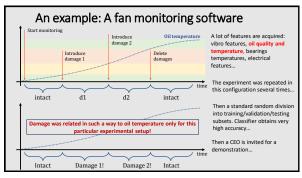
Δ



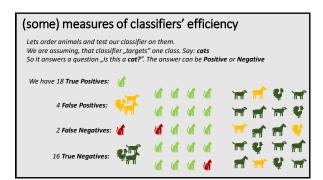
5

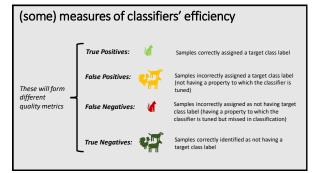


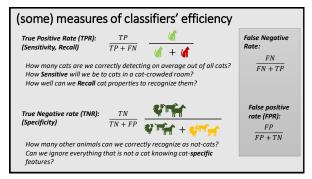


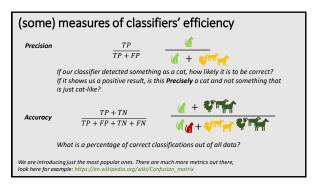


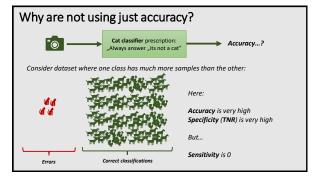
Data diversity So how to check if our dataset may contain/cause data leakage? Unfortunately, there are no straightforward methods for that... But we can: Check features and think if they can be measured with no prior knowledge of targets (if no possible leakage) Check data distribution and think if it is close to one expected in-operation Check for experimental conditions and think if they look plausible (if all the expected in-operation conditions are covered) Check for experimental conditions and think if a covered in they look plausible (if all the expected in-operation conditions are covered) Wherever possible, refrain from random train/test splits in favor of informed splits ensuring independence of samples (mimicking possible future experiment setting)

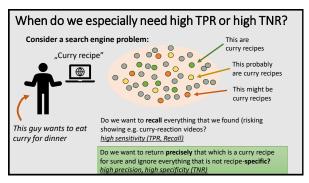


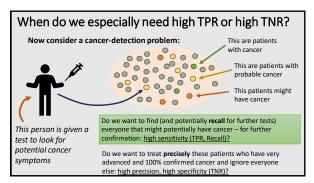


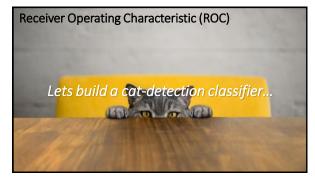


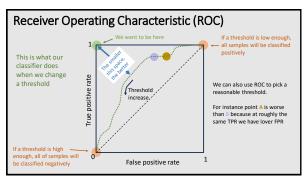


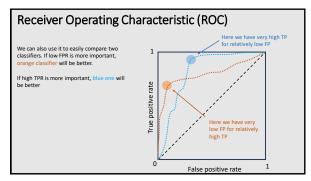


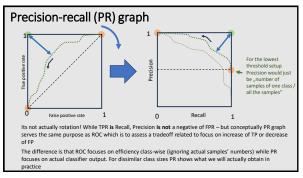












Precision-recall (PR) graph Receiver Operating Characteristic (ROC) How to use them? 1) We can compare classifiers independently from threshold setup (calculation of area-under-curve) 2) We can use these graphs to configure final output with respect to our actual needs

What about	t regression?					
Regression is much simpler – you just measure how far your model's prediction is from the actual target:						
MSE (Mean squared error)	Easy interpretation (same unit)					
MAE (Mean absolute error)	 Focuses on "important errors" 					
RMSE (Root mean squared error)	 Focuses on "important errors" but is also easy to interpret 					
R2 (R Squared, coefficient of determination) $R^2 = 1 - \frac{SS_{res}}{SS_{tot}} * * * * * * * * * * * * * * * * * * $	Allows comparison of models using different datasets					

Confus	ion matrix			
		Actual positive	Actual negative	
	Predicted positive	TP	₽P P	
	Predicted negative	FN	TN TN	

Confus	ion m	natrix	:	Now: if we have B sample, how	
	Actual A (100)	Actual B (200)	Actual C (100)	Actual D (20)	likely we are to correctly classify
Predicted A	100	10	30	0	We can calculate actual probability of B output given B:
Predicted B	0	150	10	0	
Predicted C	0	0	60	0	$P(y_B B) = \frac{\# y_B B}{\# y_A B + \# y_B B + \# y_C B + \# y_D B},$
Predicted D	0	40	0	20	Number of B outputs for B class (TP) # v _e B
		Ц.	ı	$P(y_B B) = \frac{\# y_B B}{\#B}$	
Sum of numbers in column is equal to					≠
number of samples in class					Number of samples in B (TP + FN)

Confusion matrix							
	Actual A (100)	Actual B (200)	Actual C (100)	Actual D (20)	Now: classifier says "A" – what does that mean?		
Predicted A	100	10	10 30 0 If distribution of testing se		If distribution of testing set reflects reality, we can calculate actual		
Predicted B	0	150	10	0	probability of A given A output:		
Predicted C	0	0	60	0	$P(A y_A) = \frac{\# y_A A}{\# y_A A + \# y_A B + \# y_A C + \# y_A D},$		
Predicted D	0	40	0	20	Number of A outputs for A class (TP) # $\gamma_A A$		
	um of nui		olumn is in class	$P(A y_A) = \frac{\# y_A A}{\# y_A}$ Number of A outputs (TP + FP)			

Confusion matrix							
	Actual A (100)	Actual B (200)	Actual C (100)	Actual D (20)	We can use confusion matrix for: - Estimation of probability of		
Predicted A	100	10	30	0	correct classification		
Predicted B	0	150	10	0	- Estimation of probability of class presence		
Predicted C	0	0	60	0	- Experiment planning (which		
Predicted D	0	40	0	20	classes require more samples) - Model tuning (which classes		
					require higher accuracy)		

26

Gentle introduction to bayes rule

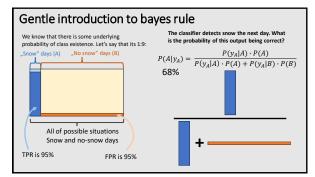
Imagine that we have a weather forecast classifier that calculates probability of snow the next day

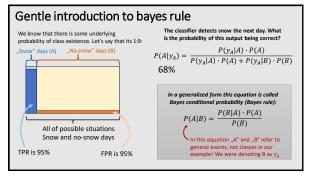
The classifier is really good, it has overall 95% accuracy, Snow prediction is as well 95% specific and 95% sensitive (snow is predicted for 5% of not-snow days and 5% of snow days is not predicted)

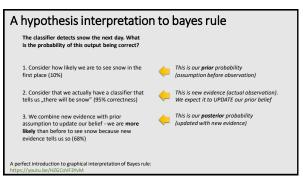
The classifier detects snow the next day. What is the probability of this output being correct?

What is the answer if we know that only 10% of days are snowy?

What if we have this answer in the middle of June?







Things to remember:
mings to remember

- 1. Data leakage: explain the problem, four main sources for it, provide at least two different practical examples
 2. Explain risks related to assessment of data using only visual or only statistical means, explain what are important aspects to consider when looking at a new dataset
 3. Define what is a false positive, true positive, false negative and true negative indication, provide a practical example
 4. Define classifier metrics: FPR, FPR, FNR, TNR, Precision, Accuracy, Recall, Sensitivity, Specificity
 5. Draw examples of ROC and PR diagrams, describe elements and show how threshold setup affects the curves, explain how two different classifiers can be compared on one diagram
 6. Explain how ROC and PR diagrams are different from usage perspective and how are they similar (what purpose do they serve)
 7. Draw an example of a confusion matrix, explain how can it be used to improve experiment or understand classification outcomes
 8. Draw a graphical interpretation of a Bayes rule, explain probability of correct classification given TPR, FPR and class prevalence

•			