

# Lecture 10: Systems Evaluation

**Chenshu Wu**

Department of Computer Science

2025 Spring

April 15, 2025



香港大學

THE UNIVERSITY OF HONG KONG



# Contents

- Learning Outcome: How to evaluate a system/method?
- Metrics and Charts
  - Accuracy, Precision, Recall, F1-Score
  - Confusion Matrix
  - Detection Rate and False Alarm Rate
  - ROC Curve
  - Boxplots and Violin plots
  - CDF
- Methodology



# Why to evaluate?

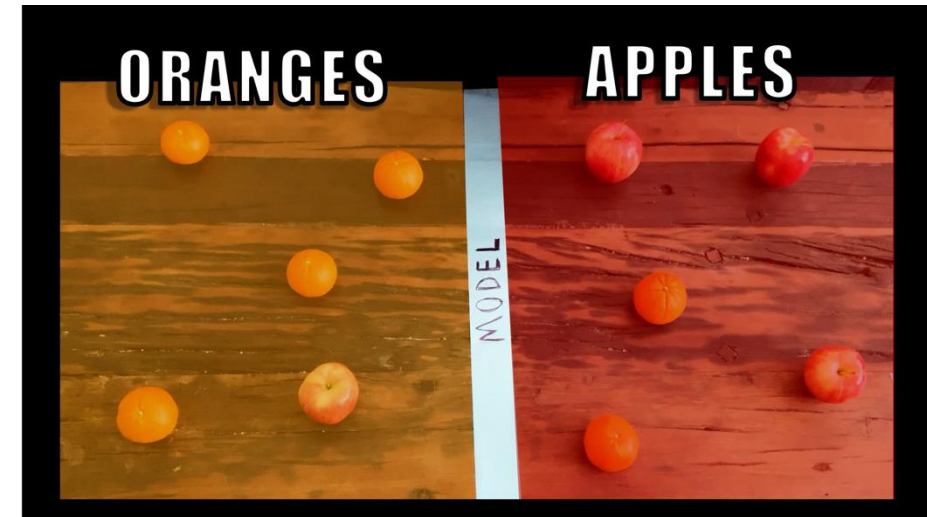
- Ensure system requirements are met
- Understand the boundaries (including the limitations)
- Compare with benchmarks or SOTA
- Explore trade-offs

# How did you evaluate your results?

- In labs and projects

# Classification Metrics

- consider building a model to classify apples and oranges on a flat surface



Source: <https://kimberlyfessel.com/mathematics/data/accuracy-precision-recall/>

# Accuracy

- Accuracy: all of the correctly classified observations and divide by the total number of observations
- one of the most popular classification metrics
- Imbalanced dataset; that is, what if we have 990 oranges and only 10 apples?



# Confusion matrix (binary classification)

		Predicted	
		Spam	Not spam
Actual	Spam	<b>True positive</b> ✓	<b>False negative</b> ✗
	Not spam	<b>False positive</b> ✗	<b>True negative</b> ✓

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

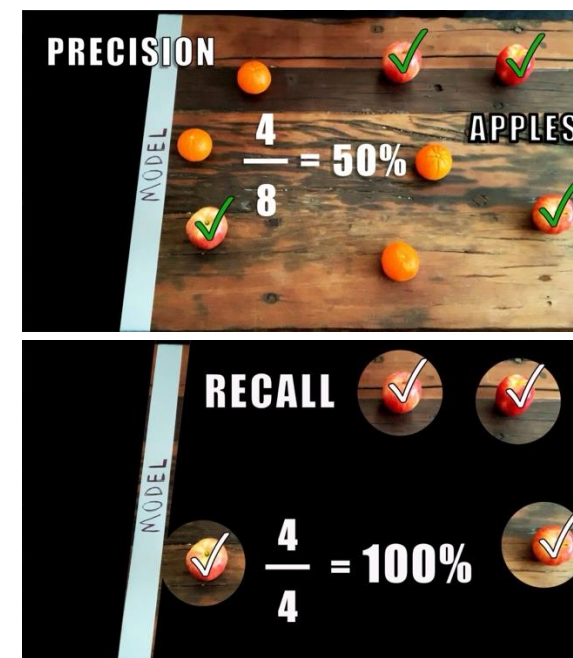
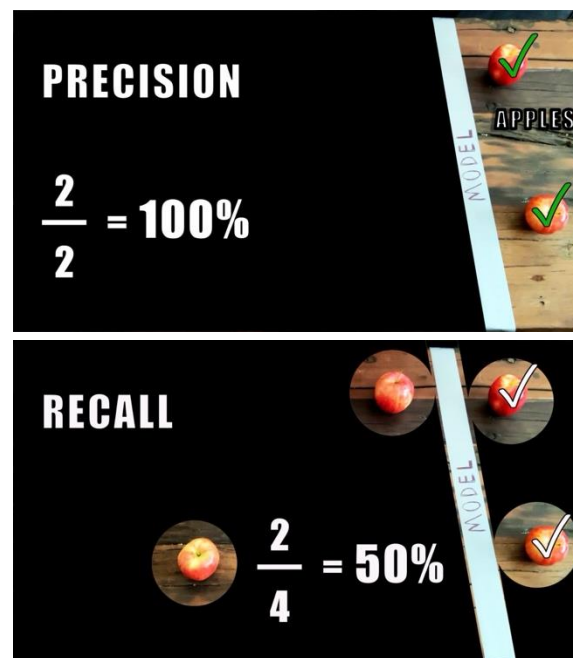
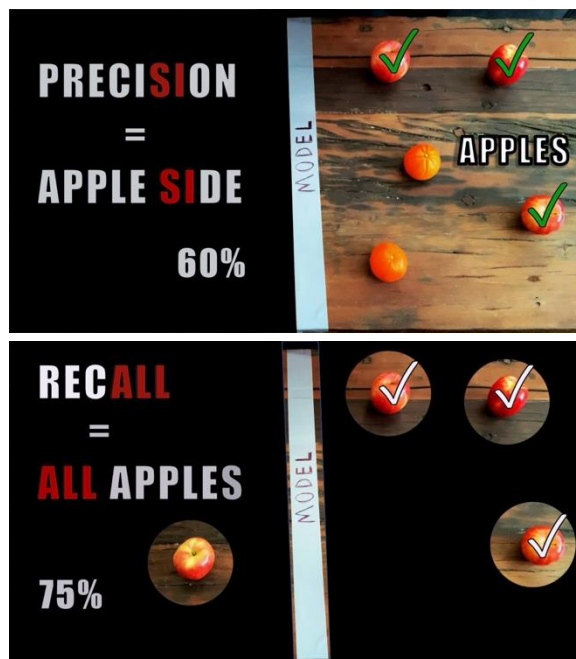
$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

- precision measures the accuracy of positive predictions
- recall measures the ability to find all relevant positive instances

# Precision and Recall

- Both precision and recall are defined in terms of just one class, oftentimes the positive—or minority—class.
- Precision-Recall Tradeoff

$$\text{F1-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$





# Confusion matrix

True material	wood	100.0%				
	plastic		99.0%	1.0%		
	ceramic		14.5%	74.8%	10.7%	
	water		0.1%	4.6%	94.3%	1.0%
	aluminum				3.8%	96.2%
		wood	plastic	ceramic	water	aluminum

Predicted material

Actual	push	85.0	0.0	4.0	5.0	3.0	3.0
	sweep	1.1	95.8	0.0	0.0	2.1	1.1
	clap	0.0	0.0	99.0	0.0	1.1	0.0
	slide	0.0	0.0	0.9	97.2	0.0	1.9
	circle	1.0	5.1	0.0	3.0	85.9	5.1
	zigzag	2.9	1.9	0.0	1.9	0.0	93.3
		push	sweep	clap	slide	circle	zigzag

Predicted

True material	A: wood desk	.83			.17																	
	B: side table	.03	.14	.36	.47																	
	C: book	.26	.15	.55						.04												
	D: dining table	.55			.45																	
	E: yoga ball				.39	.40	.21															
	F: play mat					.52	.48															
	G: plastic chair							.79	.21													
	H: water filter							.05	.84	.11												
	I: yoga mat				.01			.09	.42	.48												
	J: chop board										.99	.01										
	K: teatable										.03	.97										
	L: painted door										.45		.50	.05								
	M: ipad front												.14	.86								
	N: drywall												.17	.66	.17							
	O: acivity desk													.04	.96							
	P: bread toaster													.01	.76	.14	.09					
	Q: metal plate														.08	.25	.38	.30				
	R: candy box															.41	.43	.16				
	S: ipad back															.06	.03	.55	.36			
	T: metal door																		1.0			
	U: nut box																		.08	.92		
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U

Predicted material



# When to use which?

- Precision: Use when false positives are more costly than false negatives.
- Recall: Use when false negatives are more costly than false positives.
- F1-score: Use when you want to balance precision and recall.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{F1-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

# Detection

- Detection Rate

- How many events (X) are correctly detected out of a total of N of true events?  $DR = X/N$

- False Alarm Rate

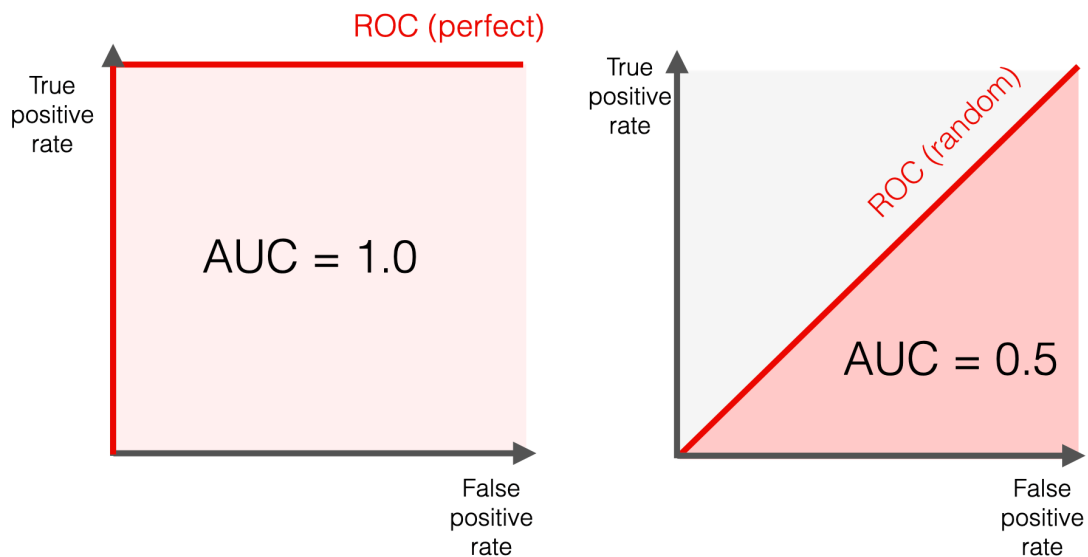
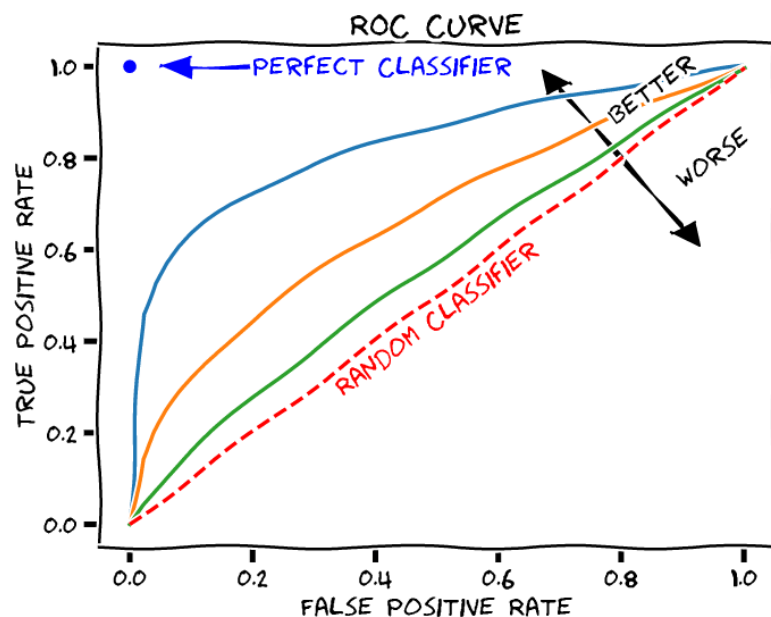
- How many events are detected when there are actually no events of interests?
- Can depend on the definition of the period/number of samples for “no events”.

- Trade-off between DR and FAR

- Applications can be more DR-sensitive or FAR-sensitive
- ROC curve

# ROC curve

- Receiver-Operating Characteristics (ROC)
- Useful to show sensitivity to threshold, or search for the best threshold



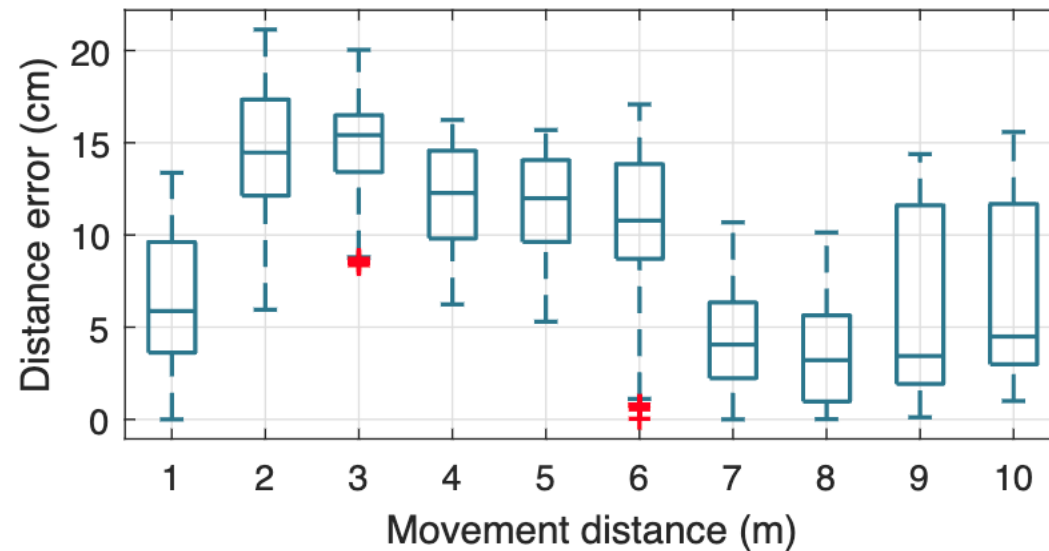
AUC-ROC

# Regression Tasks

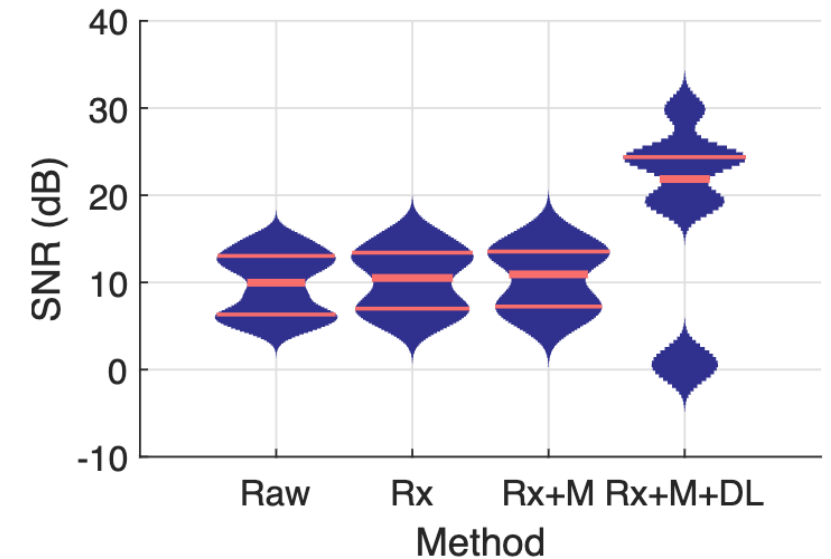
- Example tasks:
  - Localization
  - Heart/Breathing rate estimation
- Calculation
  - L1/L2 norm:  $e = ||\text{estimate} - \text{ground\_truth}||$
  - Consider n samples:  $1/n * \text{sum}(e_i)$
  - mean squared error (MSE)
  - root mean squared error (RMSE)
- Statistics
  - Mean, median, percentiles
  - More informative formats?

$$\text{MSE} = \overset{\text{Mean}}{\frac{1}{n}} \sum_{i=1}^n \left( \overset{\text{Error}}{Y_i - \hat{Y}_i} \right) \overset{\text{Squared}}{^2}$$

# Boxplots & Violin plots



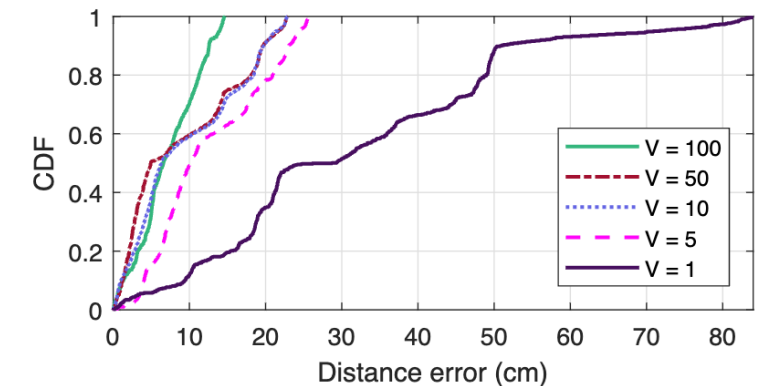
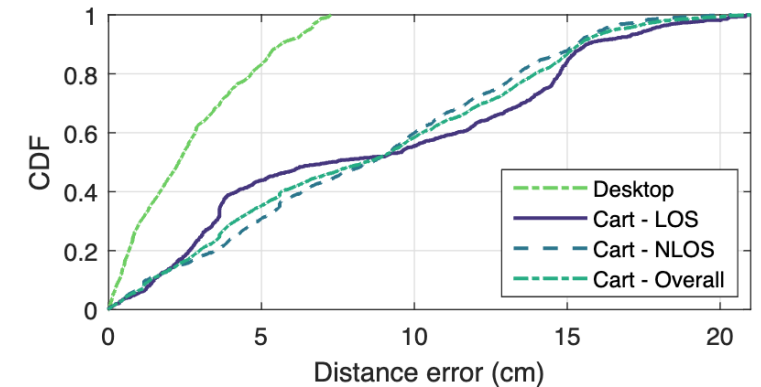
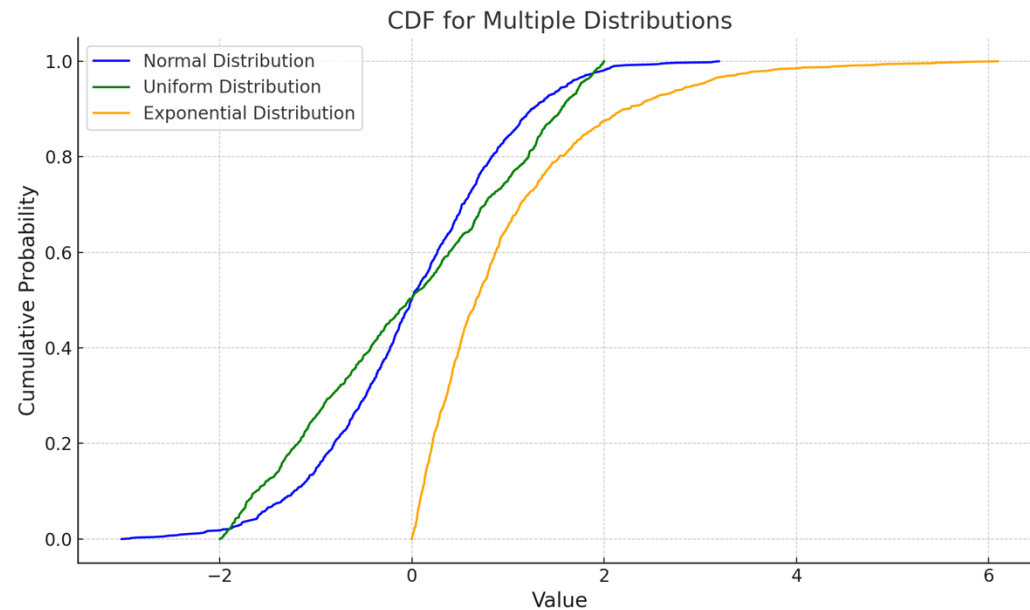
showing the median, quartiles, and outliers



show the full distribution, including the median, quartiles, and the overall distribution shape

# Cumulative Distribution Function (CDF)

- Show the distribution of a variable over its entire/possible range
- Give a sense of percentile performance
- Default metric for localization



# Beyond Accuracy

- Computation Complexity
  - Latency
  - Throughput
  - Memory/cpu usage
  - Energy/Power consumption
- 
- Scalability
  - Security & Privacy metrics
- 
- Usually more important for a (real-time) systems
  - Essential to clearly state what platforms are used for evaluation



# What to evaluate

- Overall performance
  - One or two key metrics
- Ablation study
  - How different components contribute to the overall performance?
- Parameter study
  - How does the system work under all impacting factors?
- Comparison study
  - How does your method perform compared with baseline approaches?
- Cross-validation
- Field trials/Real-world evaluation

# Case Study: Breathing Rate Estimation



# Summary

- There are many more metrics to explore and use.
- Some are specific to certain applications/fields (eg, PESQ, STOI for speech)
- Need to decide the most appropriate metric(s) depending on your applications and the performance you care about.
- The key is that, it is important to conduct extensive evaluation in a logical and scientific manner, and report the results in a clear way.



# Questions?

- Thank you!