



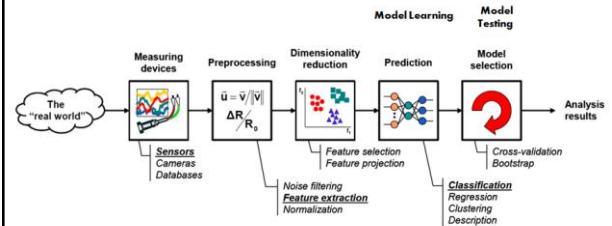
NHẬN DẠNG PATTERN RECOGNITION



ĐÁNH GIÁ MÔ HÌNH- EVALUATION

1

QUÁ TRÌNH HỌC MÁY CƠ BẢN



2

Nội dung

1. Tầm quan trọng của Evaluation?
2. Các tiêu chí đánh giá.
3. Các phương pháp đánh giá.
4. Một số độ đo tương ứng với bài toán.

3

Tại sao phải đánh giá ?

1. Biết được **khi nào huấn luyện mô hình thành công** ?
2. Biết được **mức độ thành công** của mô hình
3. Biết được **thời điểm dừng quá trình huấn luyện**
4. Biết được **khi nào cần cập nhật mô hình** ?

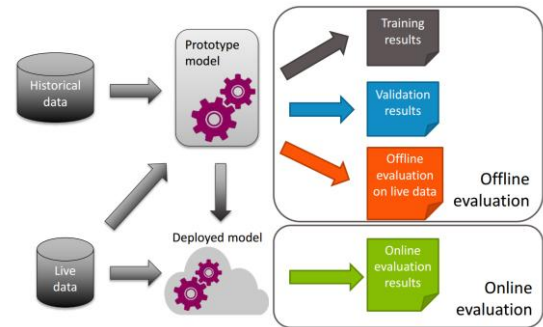
4

Một số câu hỏi căn bản khi evaluation

1. Đánh giá **khi nào** ?
2. **Các tiêu chí** đánh giá là gì ?
3. Dữ liệu – **Phương pháp đánh giá** ?
4. **Độ đo** nào được sử dụng ?

5

When to evaluation



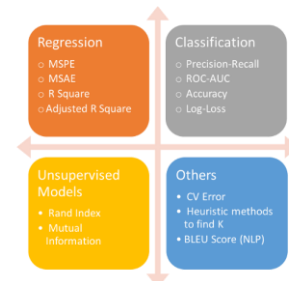
2. Các tiêu chí đánh giá

1. Tính chính xác (Accuracy)
2. Tính hiệu quả (Efficiency)
3. Khả năng xử lý nhiễu (Robustness).
4. Khả năng mở rộng (Scalability).
5. Khả năng diễn giải (Interpretability)
6. Mức độ phức tạp (complexity)

7

2. 1 Accuracy – chính xác

→ Tùy vào **bài toán, dữ liệu** sẽ có độ đo tương ứng.



8

2.2 Efficiency – hiệu quả

- Chi phí về **thời gian và tài nguyên** (bộ nhớ) cần thiết cho việc huấn luyện và kiểm thử hệ thống.



9

2.3 Robustness – xử lý nhiễu

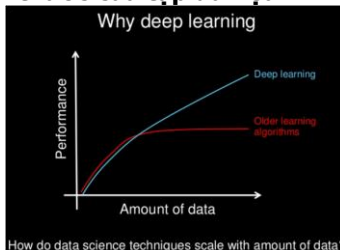
- Khả năng xử lý của hệ thống đối với các ví dụ **nhiều (lỗi)** hoặc **thiếu giá trị**.



10

2.4 Scalability – mở rộng

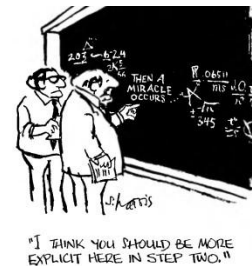
- **Hiệu năng** của hệ thống (ví dụ: tốc độ học, độ chính xác) **thay đổi** như thế nào đối với **kích thước của tập dữ liệu**



11

2.5 Interpretability – diễn giải

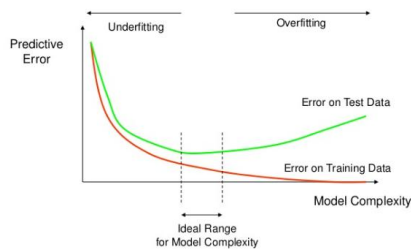
- **Mức độ dễ hiểu** (đối với người sử dụng) của các **kết quả** và **hoạt động** của hệ thống.



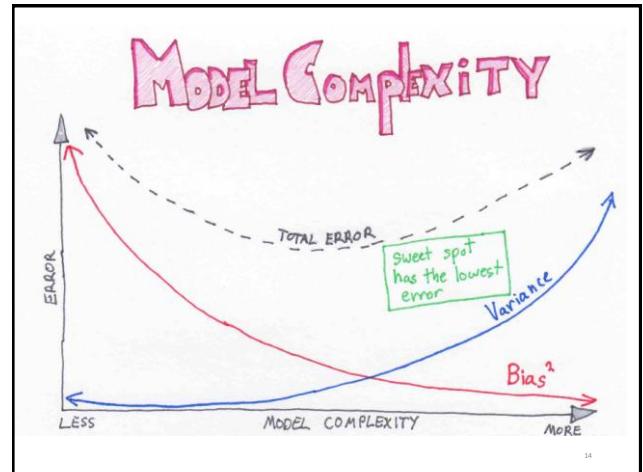
12

2.6 Complexity – mức độ phức tạp

→ **Mức độ phức tạp** của hệ thống (hàm hypothesis mục tiêu) học được.



13



14

3. Các phương pháp đánh giá

1. Hold-out / Repeated hold-out
2. Stratified sampling
3. Cross- validation
 1. K-Fold
 2. Leave one out
4. Bootstrap sampling.

15

3. 1 Hold-Out (Splitting)

→ Toàn bộ dữ liệu **D** được chia thành **2 tập con** không giao nhau

D_Train: dùng để huấn luyện hệ thống

D_Test: để đánh giá hiệu năng hệ thống sau khi học

Một số yêu cầu:

- Dữ liệu đã sử dụng ở D_Train thì không được ở trong D_Test
- Các tỉ lệ thường sử dụng D_Train = 2/3 D, D_Train 80% D_Test 20% , , D_Train 70% D_Test 30%

→ Thường phù hợp cho tập D có kích thước lớn.

16

3.1 Repeat Hold-Out

- Áp dụng hold-out nhiều lần
- Trong mỗi lần lặp một tỉ lệ nhất định của D được lựa chọn ngẫu nhiên để tạo tập dữ liệu D'
- Các giá trị lỗi (hoặc các giá trị đối với các tiêu chí đánh giá khác) được ghi nhận trong các bước lặp này được lấy trung bình cộng để xác định lỗi tổng thể

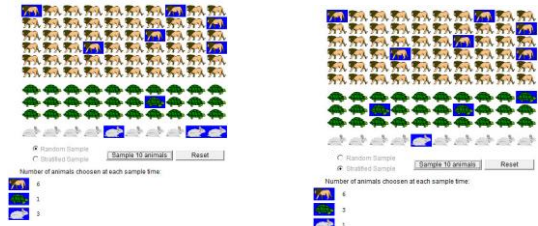
Hạn chế:

- Mỗi bước lặp sử dụng một tập kiểm thử khác nhau
- Có trường hợp một ví dụ trùng lặp trong các kiểm thử.

17

3.2 Statified sampling – lấy mẫu phân tầng

- Tập **dataset có kích thước nhỏ**.
- Dữ liệu không cân xứng (**unbalanced dataset**)
- **Mục tiêu:** Phân lớp (class distribution) trong tập huấn luyện và kiểm thử phải xấp xỉ như trong toàn tập Dataset.



3.2 Statified sampling – lấy mẫu phân tầng

- Lấy mẫu phân tầng có tác dụng:
 - Làm cân xứng (về phân bố lớp)
 - Đảm bảo tỉ lệ phân lớp (tỉ lệ các ví dụ giữa các lớp) trong tập huấn luyện và tập kiểm thử là xấp xỉ nhau.
- Phương pháp này **không** áp dụng cho bài toán hồi quy.

19

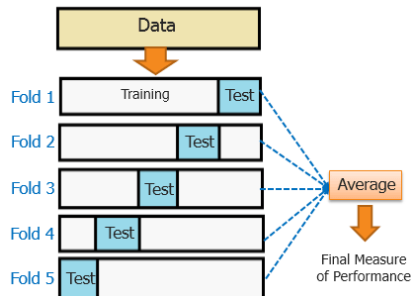
3.3 Cross-Validation

Cross validation - Phương pháp phù hợp khi tập dữ liệu **vừa và nhỏ**

- **K-Fold Cross Validation**
- **Stratified K-fold Cross Validation**
- **Leave One Out Cross Validation**

20

3.3 Cross-Validation



21

3.3 Cross-Validation

- Tập dữ liệu D được **chia thành k tập con không giao nhau** (gọi là fold) có **kích thước xấp xỉ nhau**
- Mỗi lần (trong k lần lặp), **một tập con được sử dụng làm tập kiểm thử**, và **k-1 tập con còn lại được dùng làm tập huấn luyện**.
- K giá trị lỗi (mỗi giá trị tương ứng với fold) được **tính trung bình cộng để thu được độ lỗi tổng thể**.
- Các lựa chọn cho K thường là **5 hoặc 10**

22

K-Fold Cross-Validation

k-fold is better. LOOCV will have less bias, but more variance. K-fold # of folds is a bias-variance trade-off.

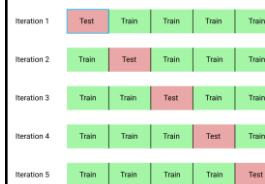
$$CV(k) = \frac{1}{k} \sum_{i=1}^k MSE_i$$

mean square error
of folds

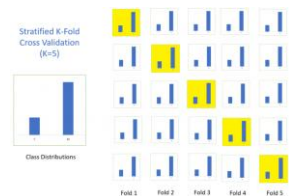
K=10 is typical

23

3.3 Cross-Validation



K-Fold Cross Validation



Stratified K-Fold Cross Validation

Các lớp trong tập có phân bố lớp xấp xỉ bằng nhau

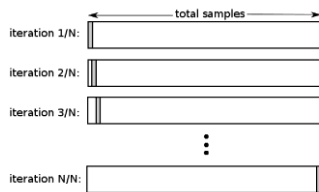
24

3.3 Cross-Validation

Leave-one-out cross-validation – Phù hợp cho tập D (rất) nhỏ

Là một dạng của cross-validation trong đó:

- Số lượng các folds **bằng kích thước của tập dữ liệu**.
- **Mỗi fold chỉ bao gồm một ví dụ.**



25

3.4 Bootstrap sampling

- **Cross-validation** – lấy mẫu không lặp lại (sampling without replacement)
- **Bootstrap sampling** – lấy mẫu có lặp lại (sampling with replacement)



26

3.4 Bootstrap sampling

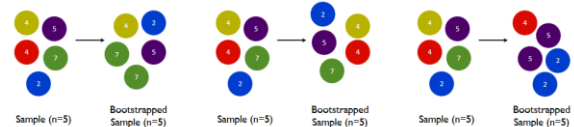
- **Cross-validation** – lấy mẫu không lặp lại (sampling without replacement)
- **Bootstrap sampling** – lấy mẫu có lặp lại (sampling with replacement)



27

3.4 Bootstrap sampling

- **Phù hợp với tập dữ liệu có kích thước (rất) nhỏ**



28

BOOTSTRAP

Note: with replacement

i	X1	X2
1	1	1.0
2	2	2.0
3	3	3.0

→

i	X1	X2
1	1	1.0
2	2	2.0
3	3	3.0
...
1	1	1.0
2	2	2.0

Note: normally we would create many more than two bootstrapped data sets.

Bootstrap allows us to simulate obtaining many new datasets by repeated sampling with replacement from the original dataset.

4. Một số độ đo

1. Accuracy/ Error
2. Precision/Recall
3. F-Score
4. AP/MAP

Confusion matrix (ma trận nhầm lẫn)

- **TP_i** (true positive): Số lượng các ví dụ thuộc lớp c_i được phân loại chính xác vào lớp c_i
- **FP_i** (false positive): Số lượng các ví dụ không thuộc lớp c_i bị phân loại nhầm vào lớp c_i
- **TN_i** (true negative): Số lượng các ví dụ không thuộc lớp c_i được phân loại (chính xác)
- **FN_i** (false negative): Số lượng các ví dụ thuộc lớp c_i bị phân loại nhầm (vào các lớp khác c_i)

Lớp c_i		Được phân lớp bởi hệ thống	
		Thuộc	Ko thuộc
Phân lớp thực sự (đúng)	Thuộc	TP _i	FN _i
	Ko thuộc	FP _i	TN _i

Confusion matrix (ma trận nhầm lẫn)

		Actual Values	
		1	0
Predicted Values	1	TRUE POSITIVE 	FALSE POSITIVE TYPE 1 ERROR
	0	FALSE NEGATIVE TYPE 2 ERROR	TRUE NEGATIVE

Error Types

TYPE I error
(false positive)

NHẦM

TYPE II error
(false negative)

BỎ SÓT

33

Error Types

		Decision Based on test	
		Accept	Reject
In Reality	TRUE	✓	✗ Type I error
	FALSE	✗ Type II error	✓

Type 1: **Loại bỏ** ví dụ mà đúng ra **không nên loại bỏ**

Type 2: **Chấp nhận** ví dụ mà đúng ra **không nên chấp nhận**

34

4. 1 Accuracy – độ chính xác

→ Mức độ dự đoán (phân lớp) **chính xác** của hệ thống (đã được huấn luyện) đối với ví dụ kiểm chứng (test data).

		Actual	
		Positives(1)	Negatives(0)
Predicted	Positives(1)	TP	FP
	Negatives(0)	FN	TN

Error = 1 - accuracy

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

35

ACCURACY

$$Acc = \frac{1}{n} \sum 1(\hat{y}_i = y_i)$$

\hat{y}_i : Predicted y
 y_i : True y
 n : number of observations
 $1(\hat{y}_i = y_i)$: Indicator function

A common metric in classification. Fails when we have highly imbalanced classes. In those cases F1 is more appropriate.

ChrisAlbon

36

4. 1 Accuracy – độ chính xác

- Là độ đo tính toán đơn giản nhất.
- Phù hợp cho các bài toán bộ dữ liệu cân bằng trong đó tỉ lệ FP (nhầm) và FN (bỏ sót) cân bằng nhau.

Hạn chế:

- Chỉ thể hiện độ chính xác không thể hiện loại lỗi trong mô hình.

37

4. 2 Precision/Recall

		Thực tế (Actual)	
		1	0
(Predicted)	1	True Positive	False Positive
	0	False Negative	True Negative

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} = \frac{\text{True Positive}}{\text{Total Predicted Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \frac{\text{True Positive}}{\text{Total Actual Positive}}$$

38

Precision

"Precision is about the predicted positives"

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

$$\text{Precision} = \frac{\text{True Positives}}{\text{Everything Predicted Positive}}$$

Precision is the ability of a classifier not to label as positive an observation that is negative. It measures the purity of positive predictions.

BY CHRIS ALBON

39

Recall

"Recall is about the real positives"

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

Recall is the ability of the classifier to find all positive examples. If we wanted to be certain to find all positive examples, we'd maximize recall!

BY CHRIS ALBON

40

4. 2 Precision/Recall

■ Precision đối với lớp c_i

→ Tổng số các ví dụ thuộc lớp c_i được phân loại chính xác chia cho tổng số các ví dụ được phân loại vào lớp c_i

$$Precision(c_i) = \frac{TP_i}{TP_i + FP_i}$$

■ Recall đối với lớp c_i

→ Tổng số các ví dụ thuộc lớp c_i được phân loại chính xác chia cho tổng số các ví dụ thuộc lớp c_i

$$Recall(c_i) = \frac{TP_i}{TP_i + FN_i}$$

41

4. 2 Precision/Recall

- Làm thế nào để tính toán được giá trị Precision và Recall (một cách tổng thể) cho toàn bộ các lớp $C=\{c_i\}$?

■ Trung bình vi mô (Micro-averaging)

$$Precision = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FP_i)} \quad Recall = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FN_i)}$$

■ Trung bình vĩ mô (Macro-averaging)

$$Precision = \frac{\sum_{i=1}^{|C|} Precision(c_i)}{|C|} \quad Recall = \frac{\sum_{i=1}^{|C|} Recall(c_i)}{|C|}$$

42

4. 2 Precision/Recall

		Actual	
		Spam	Not Spam
Predict	Spam	8	32
	Not Spam	2	8

• $Prec = 8/(8+32) = 20\%$

• $Rec = 8/10 = 80\%$

→ Tỷ lệ xác suất bộ lọc chính xác khi **xác định 1 mail là thư rác** là 20%.

→ Tỷ lệ xác suất một **thư rác** bị bộ lọc phát hiện là 80%.

43

4. 2 Precision/Recall

- Một mô hình tốt mong muốn khi Precision và Recall **đều cao**.

- Chọn Precision hay Recall tùy thuộc vào bài toán.

Hạn chế:

- Precision và Recall **thường mất cân bằng nhau**.

44

4.3 F- Score

$$F_{\beta} = (1 + \beta^2) \frac{\text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}$$

- Khi $\beta > 1$, recall được coi trọng hơn precision
- Khi $\beta < 1$, precision được coi trọng hơn.
- Khi $\beta = 1$, precision và recall coi trọng như nhau.
- β thường được sử dụng là $\beta = 2$ và $\beta = 0.5$

45

F1 Score

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F1 score can be interpreted as the harmonic mean of precision and recall. Values range from 0 (bad) to 1 (good).

46

4.3 F1 -Score

- F là một **trung bình điều hòa (harmonic mean)** của các tiêu chí Precision và Recall. Nó có xu hướng **lấy giá trị gần với giá trị nào nhỏ hơn giữa 2 tiêu chí** này.
- F1 có giá trị lớn nếu cả 2 giá trị Precision và Recall đều lớn \rightarrow F1 càng cao độ phân lớp càng tốt.

47

4.4 Average Precision

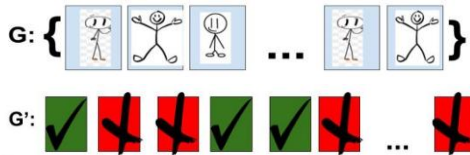
Information Retrieval



- Q to be the user query
- G to be a set of labeled data in the database

48

4. 4 Average Precision



- Ground truth positives (GTP) – number of **True of query Q**.
- $d(i,j)$ to be a score function to show how similar object i is to j
- G' which an ordered set of G according to score function $d(,)$

49

4. 4 Average Precision

$$AP@k = \frac{1}{GTP} \sum_{i=1}^k \frac{TP \text{ seen}}{i}$$

- K to be the index of G'
- GTP refers to the total number of ground truth positives for the query
- TP seen refers to the number of true positives seen till k

50

4. 4 Average Precision

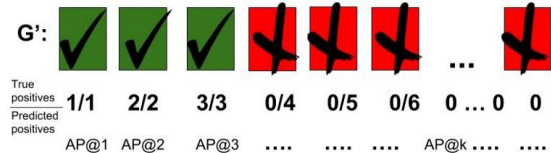


$$\text{Overall AP} = \frac{1}{3} (1/1 + 0/2 + 0/3 + 2/4 + 3/5 + 0/6 + 0 \dots + 0) = 0.7$$

Calculation of a AP for a given query, Q, with a GTP=3

51

4. 4 Average Precision



$$\text{Overall AP} = \frac{1}{3} (1/1 + 2/2 + 3/3 + 0/4 + 0/5 + 0/6 + 0 \dots + 0) = 1.0$$

Calculation of a perfect AP for a given query, Q, with a GTP=3

52

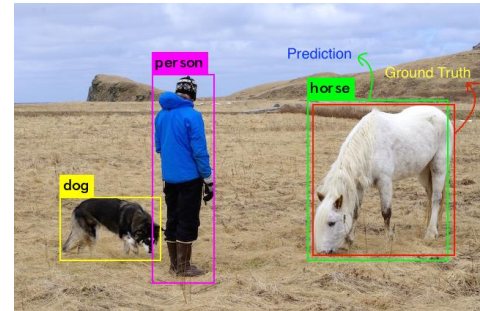
4.4 Mean Average Precision - mAP

$$\text{MAP} = \frac{\sum_{q=1}^Q \text{AveP}(q)}{Q}$$

- Q is the number of queries
- $\text{AveP}(q)$ is the average precision (AP) for a given query, q

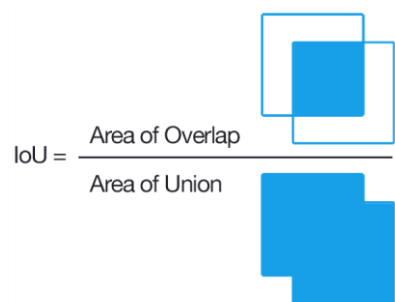
53

IoU - Intersection over union



54

IoU - Intersection over union



55

IoU - Intersection over union

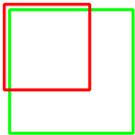


$\text{IoU} = \{0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\}$

56

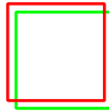
IoU - Intersection over union

IoU: 0.4034



Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

Passcal VOC Challenge : based on 50% IOU

COCO Challenge: ranging from 5% to 95%

57

IoU - Intersection over union



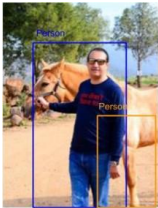
True Positive

Example (IoU > 0.5)

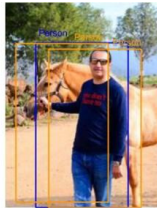
58

IoU - Intersection over union

False Positive



IoU < 0.5



Duplicate BB are considered as FP



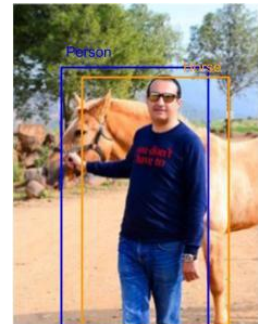
No IoU

59

IoU - Intersection over union

False Negative

IoU > 0.5 but has the wrong classification



60

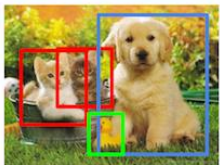
Detection

Classification



CAT

Object Detection



CAT, DOG, DUCK

61

Detection



Class	X coordinate	Y coordinate	Box Width	Box Height
Dog	100	600	150	100
Horse	700	300	200	250
Person	400	400	100	500

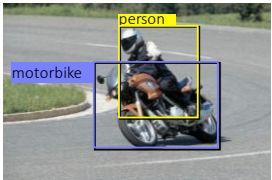
62

Problem formulation

{ airplane, bird, motorbike, person, sofa }



Input



Desired output

Evaluating a detector



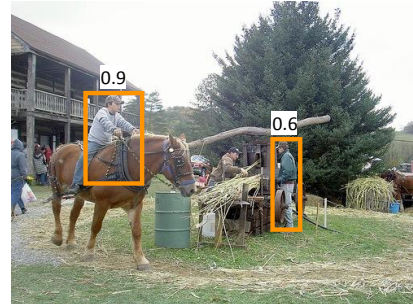
Test image (previously unseen)

First detection ...



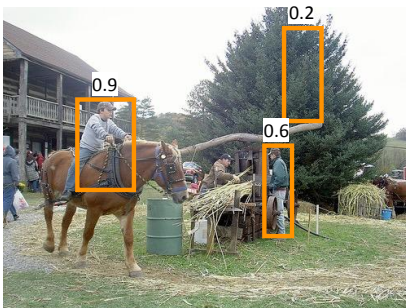
□ 'person' detector predictions

Second detection ...



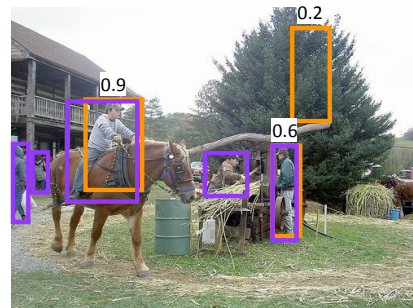
□ 'person' detector predictions

Third detection ...



□ 'person' detector predictions

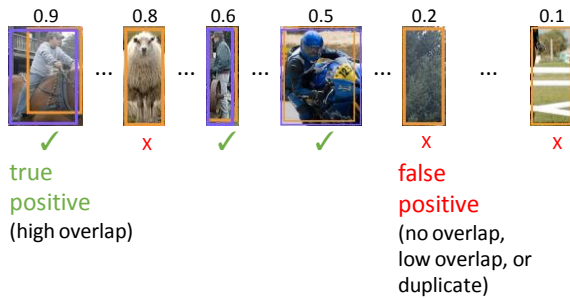
Compare to ground truth



□ 'person' detector predictions

□ ground truth 'person' boxes

Sort by confidence



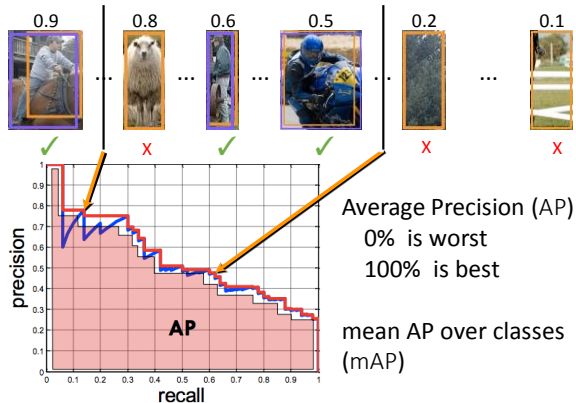
Evaluation metric



$$\text{precision@t} = \frac{\# \text{true positives@t}}{\# \text{true positives@t} + \# \text{false positives@t}} \quad \frac{\checkmark}{\checkmark + \text{X}}$$

$$\text{recall@t} = \frac{\# \text{true positives@t}}{\# \text{ground truth objects}}$$

Evaluation metric



COO- Dataset

Average Precision (AP):

- AP % AP at IoU=.50:.95 (primary challenge metric)
- AP_{IoU=.50} % AP at IoU=.50 (PASCAL VOC metric)
- AP_{IoU=.75} % AP at IoU=.75 (strict metric)

AP Across Scales:

- AP_{small} % AP for small objects: area < 32²
- AP_{medium} % AP for medium objects: 32² < area < 96²
- AP_{large} % AP for large objects: area > 96²

Average Recall (AR):

- AR_{max=1} % AR given 1 detection per image
- AR_{max=10} % AR given 10 detections per image
- AR_{max=100} % AR given 100 detections per image

AR Across Scales:

- AR_{small} % AR for small objects: area < 32²
- AR_{medium} % AR for medium objects: 32² < area < 96²
- AR_{large} % AR for large objects: area > 96²

Normalized Confusion matrix (chuẩn hóa ma trận nhầm lẫn)

	Predicted as Positive	Predicted as Negative
Actual : Positive	True Positive (TP)	False Negative (FN)
Actual : Negative	False Positive (FP)	True Negative (TN)

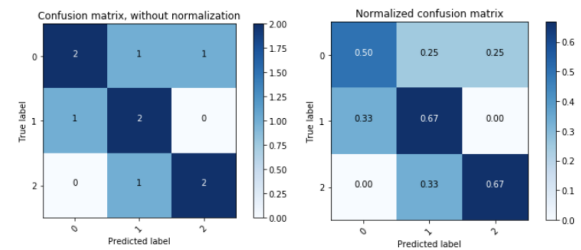


	Predicted as Positive	Predicted as Negative
Actual : Positive	$TPR = TP / (TP + FN)$	$FNR = FN / (TP + FN)$
Actual : Negative	$FPR = FP / (FP + TN)$	$TNR = TN / (FP + TN)$

→ False Positive Rate (FNR) + True Positive Rate = 1

73

Normalized Confusion matrix (chuẩn hóa ma trận nhầm lẫn)



74

Normalized Confusion matrix (chuẩn hóa ma trận nhầm lẫn)

- **False Positive Rate** còn được gọi là False Alarm Rate (tỉ lệ **báo động nhầm**).
- **False Negative Rate** còn được gọi là Miss Detection Rate (tỉ lệ **bỏ sót**)

	Predicted as Positive	Predicted as Negative
Actual : Positive	$TPR = TP / (TP + FN)$	$FNR = FN / (TP + FN)$
Actual : Negative	$FPR = FP / (FP + TN)$	$TNR = TN / (FP + TN)$

→ Trong một số bài toán việc tăng hay giảm **FNR, FPR** phụ thuộc vào ngưỡng nào đó

75

Tài liệu tham khảo

Slide được tham khảo từ:

- <http://www.cs.virginia.edu/~hw5x/Course/IR2015/site/lectures/>
- <https://nlp.stanford.edu/IR-book/newslides.html>
- <https://course.ccs.neu.edu/cs6200s14/slides.html>

