

Метрики классификации и регрессии • Metrics

- как измерить, что алгоритм хорошо решает задачу?
- как сравнить несколько алгоритмов между собой и выбрать лучший?
- как выбрать правильный способ оценки и от чего это зависит?

		Actual Class	
		Cat	Non-Cat
Predicted Class	Cat	90 <i>TP</i>	60 <i>FP</i>
	Non-Cat	10 <i>FN</i>	940 <i>TN</i>

- *True Positive / TP / Истинно-положительное*: классификатор отнес объект к классу `cat` и объект действительно относится к классу `cat`
- *False Positive / FP / Ложно-положительное*: классификатор отнес объект к классу `cat`, но на самом деле объект не относится к этому классу
- *False Negative / FN / Ложно-отрицательное*: классификатор отнес объект к классу `non cat`, но на самом деле объект относится к классу `cat`
- *True Negative / TN / Истинно-отрицательное*: классификатор отнес объект к классу `non cat`, и объект действительно относится к этому классу

Accuracy

$$ACC = \frac{TP + TN}{TP + FP + FN + TN}$$

Если выборка несбалансирована, то эта метрика ничего не покажет

Precision

$$precision = \frac{TP}{TP + FP}$$

Recall

$$recall = \frac{TP}{TP + FN}$$

F1-score

$$f1 = 2 \times \frac{precision \times recall}{precision + recall}$$

гармоническое среднее между точностью и полнотой

Метрики классификации

- Часто метрики *precision*, *recall* и *f1* можно встретить с приставкой *micro* или *macro*
- Посмотрим на примере *precision*

		<i>gold labels</i>			
		urgent	normal	spam	
<i>system output</i>	urgent	8	10	1	$\text{precision}_u = \frac{8}{8+10+1}$
	normal	5	60	50	$\text{precision}_n = \frac{60}{5+60+50}$
	spam	3	30	200	$\text{precision}_s = \frac{200}{3+30+200}$
		$\text{recall}_u = \frac{8}{8+5+3}$	$\text{recall}_n = \frac{60}{10+60+30}$	$\text{recall}_s = \frac{200}{1+50+200}$	

Метрики классификации

Class 1: Urgent			Class 2: Normal			Class 3: Spam			Pooled		
	true urgent	true not		true normal	true not		true spam	true not		true yes	true no
system urgent	8	11	system normal	60	55	system spam	200	33	system yes	268	99
system not	8	340	system not	40	212	system not	51	83	system no	99	635
precision = $\frac{8}{8+11} = .42$			precision = $\frac{60}{60+55} = .52$			precision = $\frac{200}{200+33} = .86$			microaverage precision = $\frac{268}{268+99} = .73$		
			macroaverage precision = $\frac{.42+.52+.86}{3} = .60$								

ROC-AUC

Для вычисления ROC-AUC используются две дополнительные метрики:

- True positive rate: $TPR = \frac{TP}{TP + FN}$ (a.k.a Recall)

- False positive rate: $FPR = \frac{FP}{FP + TN}$

ROC-AUC

id	оценка	класс
1	0.5	0
2	0.1	0
3	0.2	0
4	0.6	1
5	0.2	1
6	0.3	1
7	0.0	0

Табл. 1

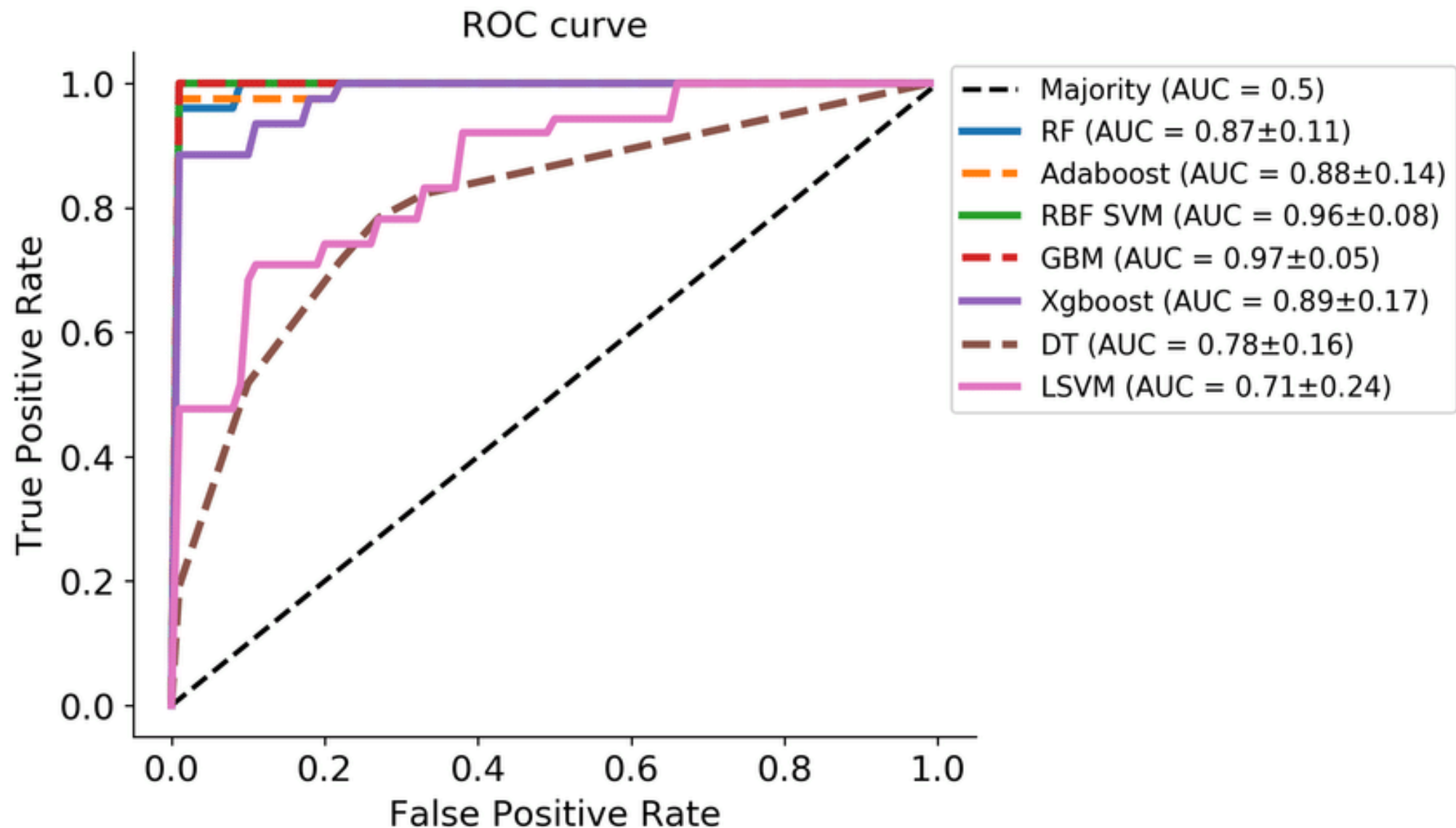
id	оценка	класс
4	0.6	1
1	0.5	0
6	0.3	1
3	0.2	0
5	0.2	1
2	0.1	0
7	0.0	0

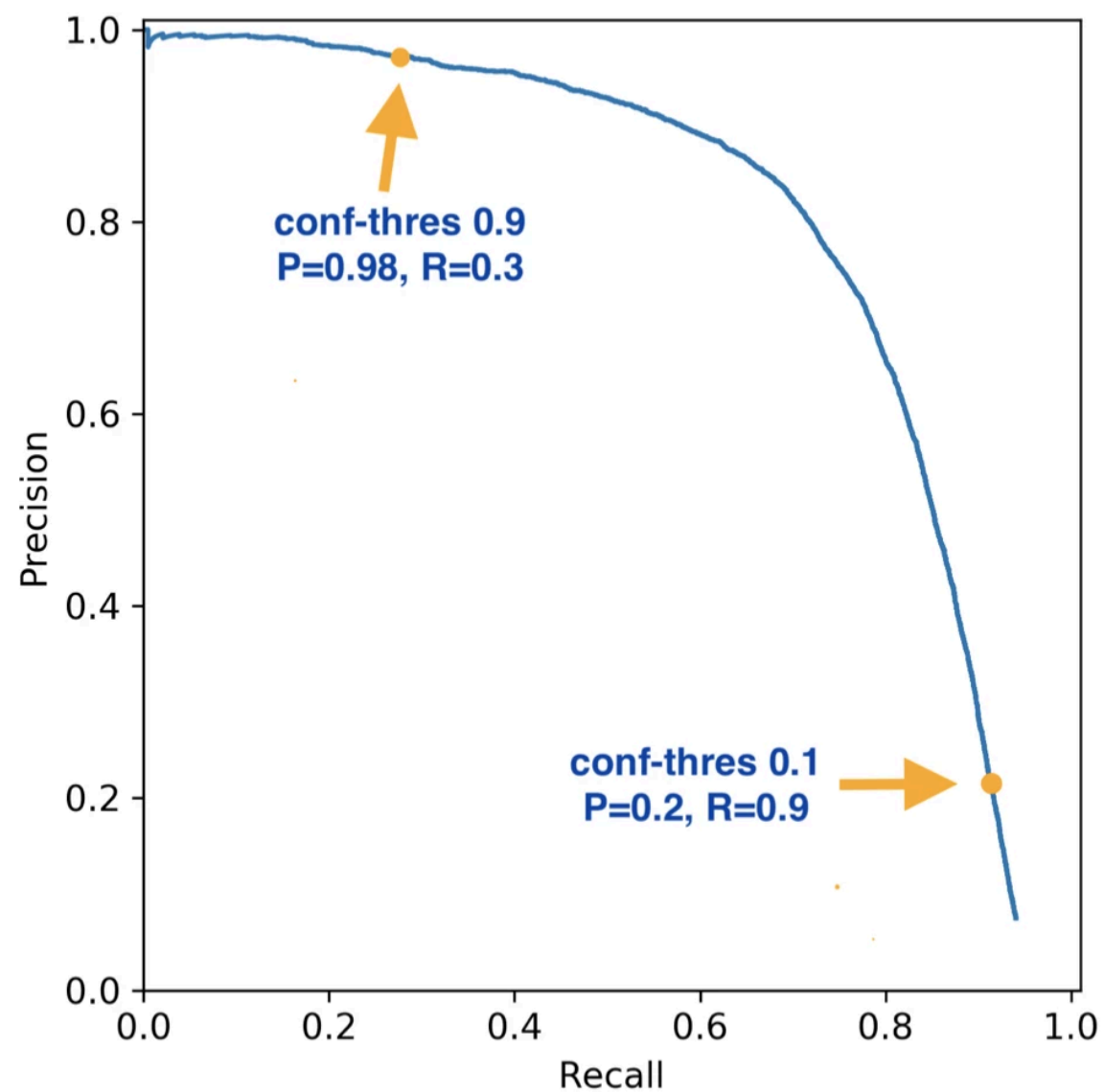
Табл. 2

id	> 0.25	класс
4	1	1
1	1	0
6	1	1
3	0	0
5	0	1
2	0	0
7	0	0

Табл. 3

ROC-AUC пример графика





Метрики регрессии

Mean Squared Error

y – истинные значение целевой переменной

\hat{y} – предикт алгоритма

N – число объектов в выборке

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

MSE применяется в ситуациях, когда нам надо подчеркнуть большие ошибки и выбрать модель, которая дает меньше больших ошибок прогноза.

Root Mean Squared Error

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

Легко интерпретировать, поскольку он имеет те же единицы, что и исходные значения (в отличие от MSE).

Mean Absolute Error

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Среднеквадратичный функционал сильнее штрафует за большие отклонения по сравнению со среднеабсолютным, и поэтому более чувствителен к выбросам.

Mean Absolute Percentage Error

$$MAPE = \frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100$$

Этот коэффициент можно интерпретировать в долях или процентах. Если получилось, например, что $MAPE=11.4\%$, то это говорит о том, что ошибка составила 11,4% от фактических значений.

Symmetric mean absolute percentage error

$$SMAPE = \frac{100\%}{N} \sum_{i=1}^N \frac{|\hat{y}_i - y_i|}{|y_i| + |\hat{y}_i|}$$

R2 score

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

- выбор метрики **всегда** зависит от задачи
- метрика \neq функция потерь
 - функцию потерь мы минимизируем по параметрам модели
 - метрикой мы измеряем, насколько качественно работает модель
- `sklearn:` **Metrics and scoring: quantifying the quality of predictions**