

# Confusion matrices

|      |         | Predicted |     |         |         |
|------|---------|-----------|-----|---------|---------|
|      |         | pos       | neg | neutral | Support |
| Gold | pos     | 15        | 10  | 100     | 125     |
|      | neg     | 10        | 15  | 10      | 35      |
|      | neutral | 10        | 100 | 1000    | 1110    |

A threshold was imposed for these categorical predictions.

# Accuracy

The correct predictions divided by the total number of examples.

|      |         | Predicted |     |         |
|------|---------|-----------|-----|---------|
|      |         | pos       | neg | neutral |
| Gold | pos     | 15        | 10  | 100     |
|      | neg     | 10        | 15  | 10      |
|      | neutral | 10        | 100 | 1000    |

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best.
- Value encoded: how often is the system correct?
- Weaknesses:
  - ▶ No per-class metrics.
  - ▶ Failure to control for class size.

# Accuracy and the cross-entropy loss

Accuracy is inversely proportional to the negative log-loss (a.k.a. cross entropy loss; [sklearn link](#)):

$$-\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{i,k} \log(p_{i,k})$$

## Precision

For class  $k$ : the correct predictions for  $k$  divided by the sum of all guesses for  $k$ .

|      |           | Predicted |      |         |
|------|-----------|-----------|------|---------|
|      |           | pos       | neg  | neutral |
| Gold | pos       | <b>15</b> | 10   | 100     |
|      | neg       | <b>10</b> | 15   | 10      |
|      | neutral   | <b>10</b> | 100  | 1000    |
|      | Precision | 0.43      | 0.12 | 0.90    |

Precision for pos:  $15 / (15 + 10 + 10) = 0.43$

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best. (Caveat: undefined values resulting from dividing by 0 need to be mapped to 0.)
- Value encoded: penalize incorrect guesses.
- Weakness: Achieve high precision for  $k$  simply by rarely guessing  $k$ .

# Recall

For class  $k$ : the correct predictions for  $k$  divided by the sum of all true members of  $k$ .

|      |         | Predicted |           |            | Recall |
|------|---------|-----------|-----------|------------|--------|
|      |         | pos       | neg       | neutral    |        |
| Gold | pos     | <b>15</b> | <b>10</b> | <b>100</b> | 0.12   |
|      | neg     | 10        | 15        | 10         | 0.43   |
|      | neutral | 10        | 100       | 1000       | 0.90   |

Recall for pos:  $15 / (15 + 10 + 100) = 0.12$

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best.
- Value encoded: penalize missed true cases.
- Weakness: Achieve high recall for  $k$  simply by always guessing  $k$ .

## F scores

$$F_{\beta}(k) = (\beta^2 + 1) \cdot \frac{\text{Precision}(k) \cdot \text{Recall}(k)}{(\beta^2 \cdot \text{Precision}(k)) + \text{Recall}(k)}$$

|      |         | Predicted |     |         | F <sub>1</sub> |
|------|---------|-----------|-----|---------|----------------|
|      |         | pos       | neg | neutral |                |
| Gold | pos     | 15        | 10  | 100     | 0.19           |
|      | neg     | 10        | 15  | 10      | 0.19           |
|      | neutral | 10        | 100 | 1000    | 0.90           |

- Bounds: [0, 1], with 0 the worst and 1 the best; always between precision and recall.
- Value encoded: how much do predictions for  $k$  align with true instances of  $k$ , with  $\beta$  controlling the weight places on precision vs. recall
- Weaknesses:
  - ▶ No normalization for the size of the dataset.
  - ▶ Ignores the values off the row and column for  $k$ .

# Averaging F scores

- Macro-averaging
- Weighted averaging
- Micro-averaging

## Macro-averaged F scores

|      |         | Predicted |     |         | F <sub>1</sub> |
|------|---------|-----------|-----|---------|----------------|
|      |         | pos       | neg | neutral |                |
| Gold | pos     | 15        | 10  | 100     | <b>0.19</b>    |
|      | neg     | 10        | 15  | 10      | <b>0.19</b>    |
|      | neutral | 10        | 100 | 1000    | <b>0.90</b>    |
|      |         |           |     |         | 0.43           |

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best.
- Value encoded: same values as F scores plus the assumption that all classes are equal.
- Weaknesses:
  - ▶ A classifier that does well only on small classes might not do well in the real world.
  - ▶ A classifier that does well only on large classes might do poorly on small but vital smaller ones.



## Weighted average F scores

|      |         | Predicted |     |         | Support | F <sub>1</sub> |
|------|---------|-----------|-----|---------|---------|----------------|
|      |         | pos       | neg | neutral |         |                |
| Gold | pos     | 15        | 10  | 100     | 125     | <b>0.19</b>    |
|      | neg     | 10        | 15  | 10      | 35      | <b>0.19</b>    |
|      | neutral | 10        | 100 | 1000    | 1110    | <b>0.90</b>    |
|      |         |           |     |         |         | <b>0.43</b>    |

$$\frac{0.19 \cdot 125 + 0.19 \cdot 35 + 0.90 \cdot 1110}{125 + 35 + 1110}$$

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best.
- Value encoded: same values as  $F_\beta$  plus the assumption that class size matters.
- Weaknesses: Large classes will dominate.

# Micro-averaged F scores

|      |         | Predicted |     |         |
|------|---------|-----------|-----|---------|
|      |         | pos       | neg | neutral |
| Gold | pos     | 15        | 10  | 100     |
|      | neg     | 10        | 15  | 10      |
|      | neutral | 10        | 100 | 1000    |

|  |     | yes | no   |  |     | yes | no   |  |     |      |     |
|--|-----|-----|------|--|-----|-----|------|--|-----|------|-----|
|  | yes | 15  | 110  |  | yes | 15  | 20   |  | yes | 1000 | 110 |
|  | no  | 20  | 1125 |  | no  | 110 | 1125 |  | no  | 110  | 50  |

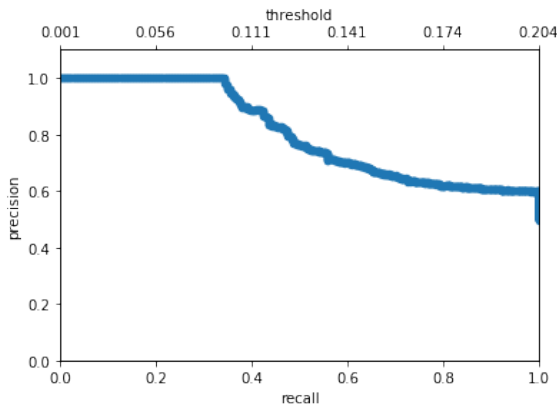
|     | yes  | no   | F <sub>1</sub> |
|-----|------|------|----------------|
| yes | 1030 | 240  | 0.81           |
| no  | 240  | 2300 | 0.91           |

# Micro-averaged F scores

- Bounds:  $[0, 1]$ , with 0 the worst and 1 the best.
- Value encoded: Micro-averaged  $F_1$  for “yes” = accuracy.
- Weaknesses:
  - ▶ Same as for weighted F scores, plus
  - ▶ a score for “yes” and “no”, hence no single summary number.

# Precision–recall curves

Summarizes the relationship between precision and recall by using each predicted probability as a potential threshold:



Average precision provides a summary of the curve.