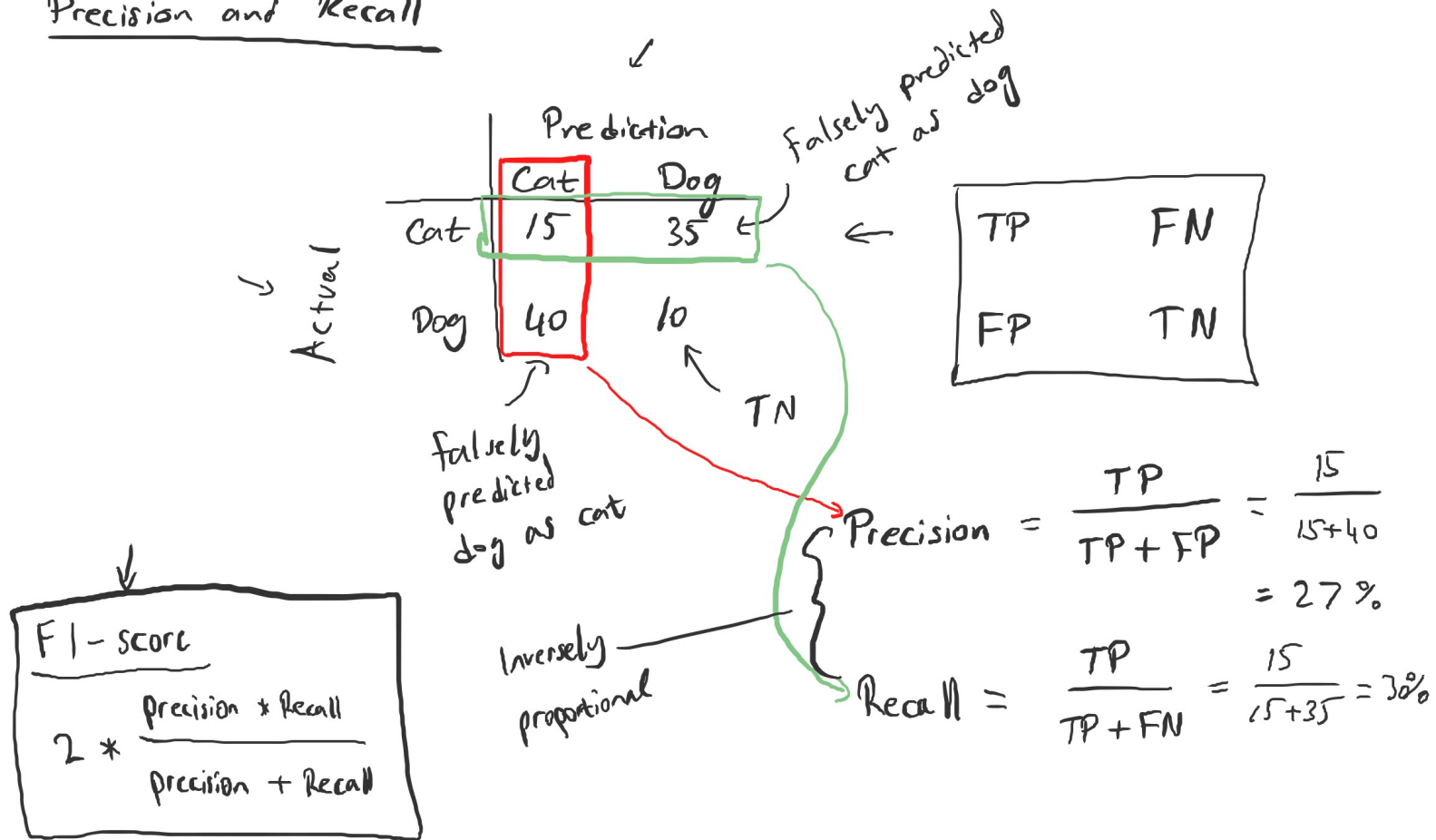## Evaluation Metrics in NLP

- Accuracy
- Precision
- Recall
- F1 - score
→ - BLEU score (ex. Machine Translation)

NLP

## Accuracy ←

- How many samples did your model predict w/ the correct label out of all samples?

# Precision and Recall

Prediction

| | Cat | Dog |
|---|---|---|
| Cat | 15 | 35 |
| Dog | 40 | 10 |

Actual

Falsely predicted cat as dog

TN

Falsely predicted dog as cat

| | |
|---|---|
| TP | FN |
| FP | TN |

Inversely proportional

$$Precision = \frac{TP}{TP+FP} = \frac{15}{15+40}$$

$$= 27\%$$

$$Recall = \frac{TP}{TP+FN} = \frac{15}{15+35} = 30\%$$

F1 - score

$$2 * \frac{precision * Recall}{precision + Recall}$$

# BLEU (bilingual evaluation understudy)

ROUGE

French : Le chat est sur le tapis

Human {
Reference 1 : The cat is on the mat. ←

Reference 2 : There is a cat on the mat. ←

Model Output : the the the the the the the ←

Modified Precision (BLEU) = $\dfrac{2}{7}$

Count-clip

$= \dfrac{count}{\;}\; \dfrac{4}{6} = \dfrac{2}{3}$

sum    /    sum

## BLEU w/ bigrams

Model Output : The cat the cat on the mat

| bigrams | Count | count-clip |
|---|---|---|
| The cat | 2 | 1 |
| cat the | 1 | 0 |
| cat on | 1 | 1 |
| on the | 1 | 1 |
| the mat | 1 | 1 |