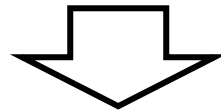
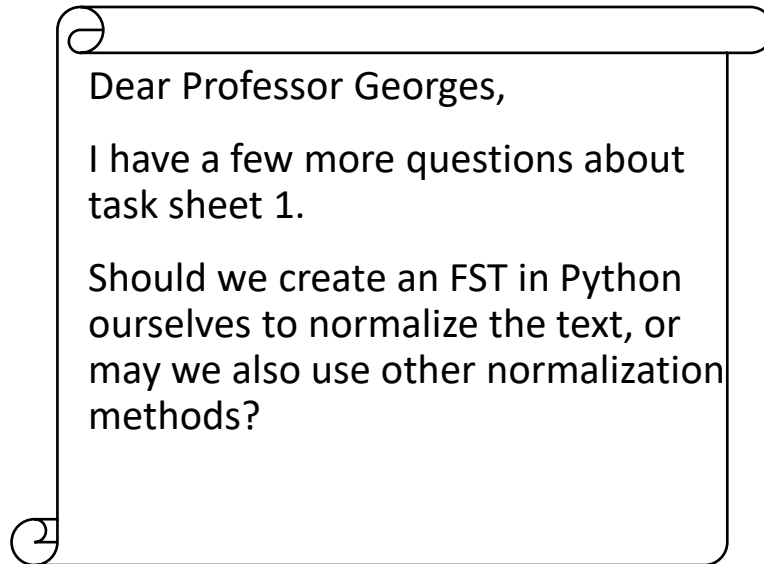
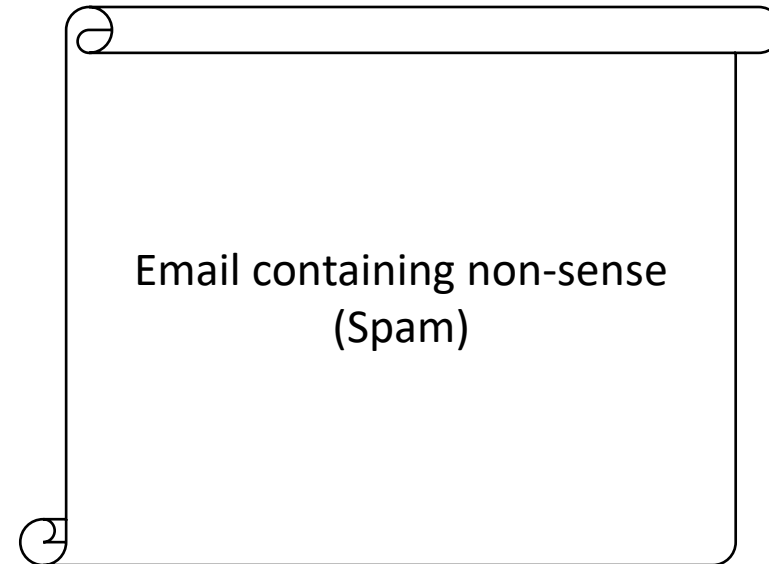


Spam Detection

Spam vs. non-Spam



Spam or ham?



Spam or ham?

Spam filter with Bayes classifier

$$SF(\underline{w}) = \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w})$$



Rings a
bell?

Spam filter with Bayes classifier

Let w be the word sequence.

Question: How to model a spam filter using Bayes classifier?

Spam filter with Bayes classifier

Let \underline{w} be the word sequence.

Question: How to model a spam filter using Bayes classifier?

Answer: A **Spam Filter**, denoted by $\mathbf{SF}(\underline{w})$, taking \underline{w} as input, can be modeled as:

$$\mathbf{SF}(\underline{w}) = \dots$$

Spam filter with Bayes classifier

$$SF(\underline{w}) = \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} \dots?$$

Spam filter with Bayes classifier

$$SF(\underline{w}) = \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w})$$

Spam filter with Bayes classifier

$$\begin{aligned} SF(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \end{aligned}$$

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

Question: What needs to be estimated?

Spam filter with Bayes classifier

$$\begin{aligned} SF(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{spam}) &= P_{\text{spam}}(\underline{w}) \\ P(\underline{w} | \text{no-spam}) &= P_{\text{no-spam}}(\underline{w}) \end{aligned} \right\} \text{Question: How to compute that?}$$

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{spam}) &= P_{\text{spam}}(\underline{w}) \\ P(\underline{w} | \text{no-spam}) &= P_{\text{no-spam}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{spam}) &= P_{\text{spam}}(\underline{w}) \\ P(\underline{w} | \text{no-spam}) &= P_{\text{no-spam}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

Question:

Which quantity is missing?

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{spam}) &= P_{\text{spam}}(\underline{w}) \\ P(\underline{w} | \text{no-spam}) &= P_{\text{no-spam}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

$P(S) \Rightarrow$ Question: How to compute that?

Spam filter with Bayes classifier

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{spam}) &= P_{\text{spam}}(\underline{w}) \\ P(\underline{w} | \text{no-spam}) &= P_{\text{no-spam}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

$$P(\text{spam}) + P(\text{no-spam}) = 1.0 \quad \text{Question: How to compute that?}$$

Spam filter with Bayes classifier

Predictor Prior Probability

$$\begin{aligned} \text{SF}(\underline{w}) &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(S | \underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) / P(\underline{w}) \\ &= \operatorname{argmax}_{S=\{\text{spam}, \text{no-spam}\}} P(\underline{w} | S) P(S) \end{aligned}$$

Likelihood

$$P(\underline{w} | \text{spam}) = P_{\text{spam}}(\underline{w})$$

$$P(\underline{w} | \text{no-spam}) = P_{\text{no-spam}}(\underline{w})$$

How about an n-gram language model?

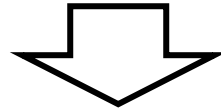
Class Prior

$$P(\text{spam}) + P(\text{no-spam}) = 1.0 \quad \text{can be estimated}$$

Language Identification

Language Identification (LID)

„Wenn wir das nehmen und darauf aufbauen, dann können wir einen Schritt weiter gehen: Wenn das Meer nicht glücklich ist, ist keiner glücklich. “



English or German?

„And if we just take that and we build from there, then we can go to the next step, which is that if the ocean ain't happy, ain't nobody happy. “



English or German?

LID with Bayes Classifier

$$\text{LID}(\underline{w}) = \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w})$$

LID with Bayes Classifier

$$\begin{aligned} \text{LID}(\underline{w}) &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) / P(\underline{w}) \end{aligned}$$

Bayes Theorem

$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)}$$

LID with Bayes Classifier

$$\begin{aligned} \text{LID}(\underline{w}) &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) / \cancel{P(\underline{w})} \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) \end{aligned}$$

Independent of
 $P(\underline{w})$

LID with Bayes Classifier

$$\begin{aligned} \text{LID}(\underline{w}) &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) / P(\underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) \end{aligned}$$

$P(\underline{w} | L) = ?$ $P(L) = ?$

LID with Bayes Classifier

$$\begin{aligned} \text{LID}(\underline{w}) &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) / P(\underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{de}) &= P_{\text{German}}(\underline{w}) \\ P(\underline{w} | \text{en}) &= P_{\text{English}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

LID with Bayes Classifier

$$\begin{aligned} \text{LID}(\underline{w}) &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(L | \underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) / P(\underline{w}) \\ &= \operatorname{argmax}_{L=\{\text{de}, \text{en}\}} P(\underline{w} | L) P(L) \end{aligned}$$

$$\left. \begin{aligned} P(\underline{w} | \text{de}) &= P_{\text{German}}(\underline{w}) \\ P(\underline{w} | \text{en}) &= P_{\text{English}}(\underline{w}) \end{aligned} \right\} \text{How about an n-gram language model?}$$

$$P(\text{de}) + P(\text{en}) = 1.0$$

Usually uniformly distributed,
depends on application

LID with Bayes Classifier

$$P(\underline{w} | \text{de}) = P_{\text{German}}(\underline{w}) = P_{\text{German}}(F(\underline{w}))$$

$$P(\underline{w} | \text{en}) = P_{\text{English}}(\underline{w}) = P_{\text{English}}(F(\underline{w}))$$

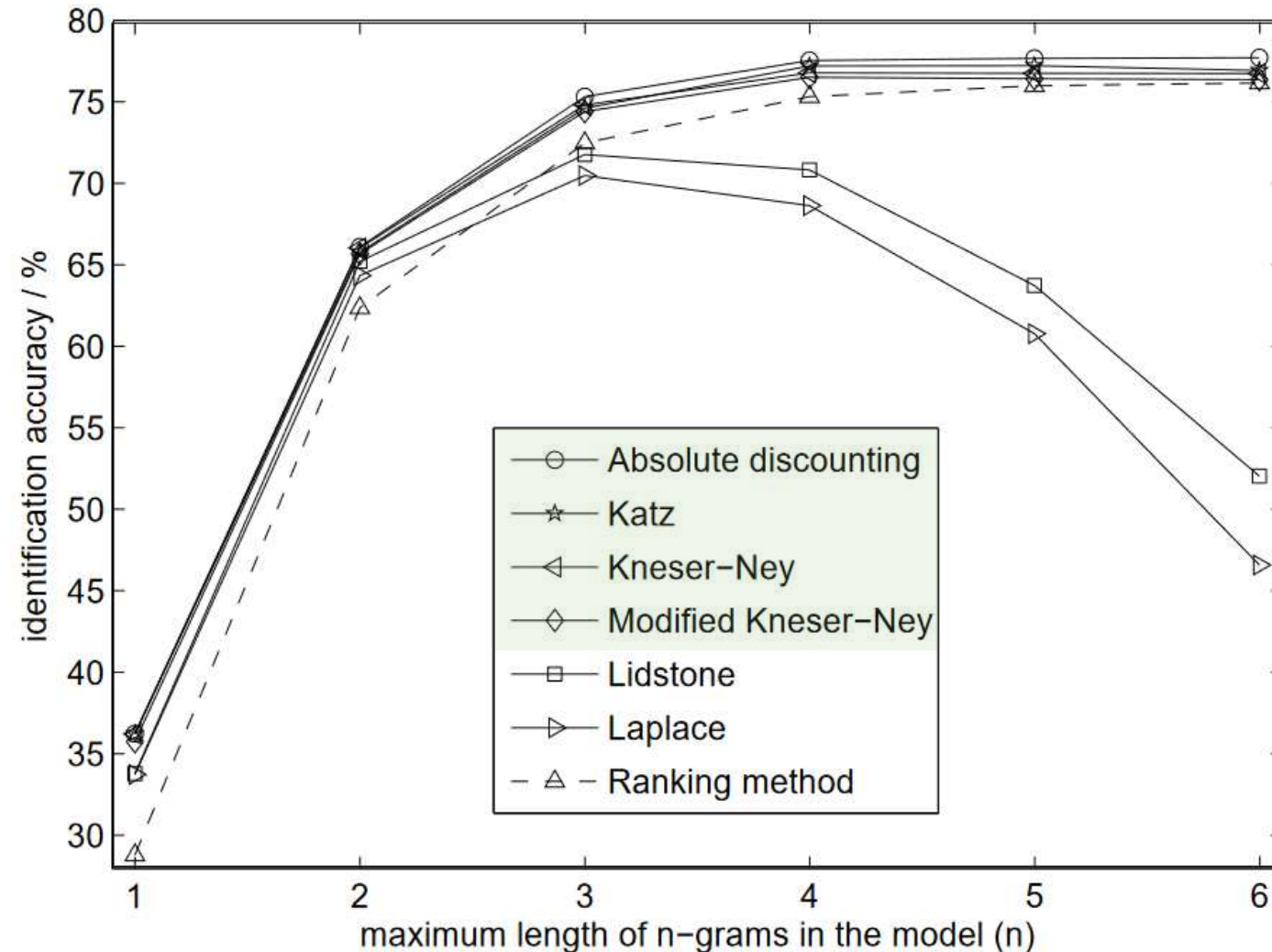
$F(\underline{w})$:= Text preprocessing

- Tokenization
- Character Set
- Punctuation
-

LID with n-gram Language Models

Different character-based n-gram models

Language	Character set size (characters)
English	57
French	60
German	70
Spanish	61
Portuguese	64
Finnish	57
Greek	68
Russian	67
Tagalog	57
Ashéninka	60
Cashinahua	59
Japanese	506
Chinese	539
Kikongo	44
⋮	⋮
Average	66.0



LID too booring?



Multi-lingual ...

- **Automatic Speech Recognition**
- **Spoken Language Understanding**

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

„Play Music“ was correctly
recognized in 80% of all
cases

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

„Increase“ was correctly
recognized in 50% of all
cases

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

„Decrease“ was confused with
„Increase“ in 40% of all cases.

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

→ Ideally:
each diagonal
element in matrix
equal to 1.0

Confusion Matrix (accuracy)

	Play Music	Increase	Decrease	Stop	Reference
Play Music	0.8	0	0	0.2	
Increase	0	0.5	0.4	0.1	
Decrease	0	0.4	0.5	0.1	
Stop	0.2	0.1	0.1	0.6	
Prediction					

Wanted:
Metric telling
„how diagonal is
the matrix“?

Evaluation

“A systematic determination of a subject's merit, worth and significance, using criteria governed by a set of standards.”

True/False P ositive/Negative

True Positive: A (spam) E-Mail is correctly recognized as spam

True Negative: A (non-spam) E-Mail is correctly recognized as non-spam

False Positive: A (non-spam) e-mail is incorrectly recognized as spam

False Negative: A (spam) e-mail is incorrectly recognized as non-spam

Accuracy

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

condition positive (P) := TP + FN

condition negative (N) := TN + FP

True Positive:

A (spam) E-Mail is correctly recognized as spam

True Negative:

A (non-spam) E-Mail is correctly recognized as non-spam

False Positive:

A (non-spam) e-mail is incorrectly recognized as spam

False Negative:

A (spam) e-mail is incorrectly recognized as non-spam

Precision

“Rate of relevant instances
among the retrieved instances.”

Positive Predictive Value (PPV)

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

True Positive:

A (spam) E-Mail is correctly recognized as spam

True Negative:

A (non-spam) E-Mail is correctly recognized as non-spam

False Positive:

A (non-spam) e-mail is incorrectly recognized as spam

False Negative:

A (spam) e-mail is incorrectly recognized as non-spam

Recall

“Rate of the total amount of relevant instances that were actually retrieved.”

Sensitivity, Hit Rate, True Positive Rate (TPR)

$$\text{TPR} = \frac{\text{TP}}{P} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

condition positive (P) := TP + FN

True Positive:

A (spam) E-Mail is correctly recognized as spam

True Negative:

A (non-spam) E-Mail is correctly recognized as non-spam

False Positive:

A (non-spam) e-mail is incorrectly recognized as spam

False Negative:

A (spam) e-mail is incorrectly recognized as non-spam

Fall-Out

False Positive Rate (FPR)

$$\text{FPR} = \frac{\text{FP}}{N} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

condition negative (N) $:= \text{TN} + \text{FP}$

True Positive:

True Negative:

False Positive:

False Negative:

A (spam) E-Mail is correctly recognized as spam

A (non-spam) E-Mail is correctly recognized as non-spam

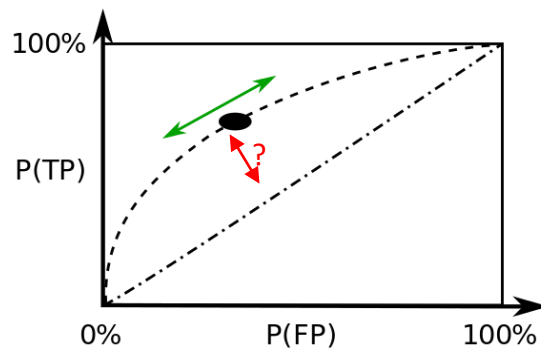
A (non-spam) e-mail is incorrectly recognized as spam

A (spam) e-mail is incorrectly recognized as non-spam

F_β - Score

F1 Score for $\beta = 1$ (or Sørensen–Dice Coefficient)

$$F_\beta = (1 + \beta^2) \cdot \frac{P \cdot R}{\beta^2 \cdot P + R}$$



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

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

P : Precision

R : Recall

F_{β} -Score

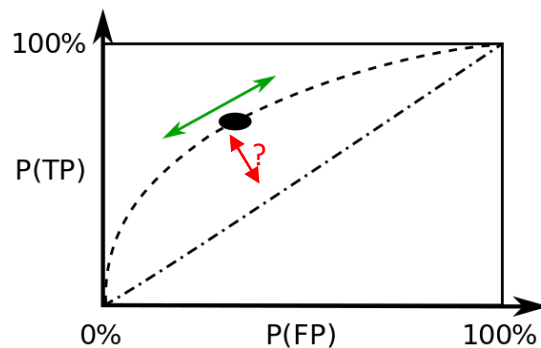
Use it only if your test set contains an equal amount
of samples per class!

If your test set is imbalanced w.r.t. class sample size,
the score is misleading!

Matthews Correlation Coefficient

Robustness despite imbalanced data set (w.r.t. #samplesPerClass)

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}$$



True Positive:

True Negative:

False Positive:

False Negative:

A (spam) E-Mail is correctly recognized as spam

A (non-spam) E-Mail is correctly recognized as non-spam

A (non-spam) e-mail is incorrectly recognized as spam

A (spam) e-mail is incorrectly recognized as non-spam