[lr:tokenization]Tokenization
[lr:wordrep]Words Rep-re-sen-taion
[lr:ml]Machine Learn-ing
[lr:oversampling]Oversampling
[lr:lm]Language Mod-els
[lr:evalmetrics]Evaluation Met-rics

1

2

Hazm
NLTK

3

?
?
?
?
?

4

?

**Trem fre-quency (TF):**
*TF*
*TF*
??

$$tf\left(t,D\right)=\frac{f_{t,d}}{\Sigma_{t`\in d}f_{t`,d}}$$

(1)

**Inverse Doc-u-ment Fre-quency (iDF):**
??

$$idf\left(t,D\right)=\log\frac{N}{\left|\{d\in D,t\in d\}\right|}$$

(2)

*TF*
*iDF*
??

$$tf-idf\left(t,d,D\right)=tf\left(t,D\right).idf\left(t,D\right)$$

(3)

?

5

?

?

*Gaus-sian*
*Naive*
*Bayes*
*SVC*
*Lin-earSVC*
*Ran-dom*
*For-est*
*Lo-gis-tic*
*Re-gres-sion*
?

6

$$P\left(y|X\right)=\frac{P(X|y).P(X)}{P(y)}$$

*X*
*y*
$P\left(X|y\right)$
*y*

$$p\left(x=v|C_{k}\right)=\frac{1}{\sqrt{2\pi\sigma_{k}^{2}}}e^{-\frac{(\epsilon-\mu_{k})^{2}}{2\mu_{k}^{2}}}$$