



۱ classification

۱.۱ ارزیابی

برای ارزیابی از Cross-Validation استفاده میکنیم به این صورت که داده ها را به k بخش قسمت کرده

در هر که در هر بخش یک قسمت را به داده های test و مابقی را به داده های train بخش بندی میکنیم

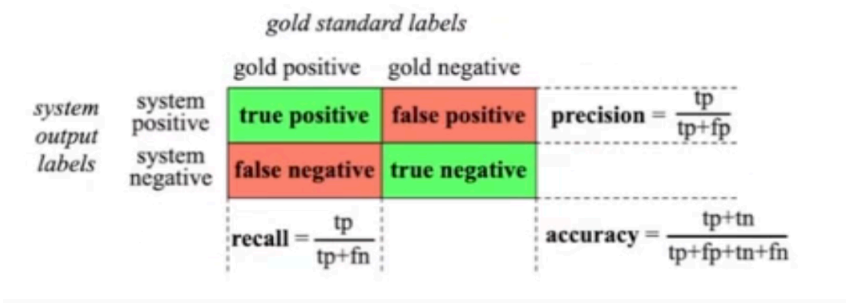
در نهایت مدلی که بیشترین میانگین و کمترین واریانس را دارد به عنوان بهترین مدل در نظر میگیریم

*معیار های ارزیابی عبارتند از :

$$\text{Micro } f_1 = \sum_c \alpha_i f_{1c}, \text{ Macro } f_1 = \frac{1}{|c|} \sum_c f_{1c}$$

** α_i بر اساس تعداد داده ها

*سایر معیار های ارزیابی در شکل ۱ آمده اند.



شکل ۱:معیارهای ارزیابی

۲.۱ GenerativeClassifiers – NaveBayes

$$\text{Buyes: Naive } \hat{y} = \operatorname{argmax}_{y \in Y} P(y \mid x) = \operatorname{argmax}_{y \in Y} \frac{P(x|y)P(y)}{P(n) \leadsto \text{ignore}}$$

اگر حالت unigram در نظر بگیریم $P(n \mid y) = \prod_{i=1}^n P(w_{v_i} \mid y) P(y)$

و اگر تعداد کلمات $|V|$ باشد برای یادگیری آن $۱ + ۲|V|$ پارامتر لازم داریم

۳.۱ DiscriminativeClassifier – logisticregression

$$w.x = \sum_{i=1} w_i x_i \quad \rho(y = 1) = \sigma(w \cdot x + b) = \hat{y}$$

$$\rho(y = 0) = 1 - \rho(y = 1) = 1 - \hat{y}$$

$$p\left(y(x) = \hat{y}^y \times (1 - \hat{y})^{1-y}\right) \rightarrow -\log p(y \mid x) = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y})) = L_{CE}(y, \hat{y}')$$

برای جلوگیری از overfitting:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^m \log p\left(y^i \mid x^i\right) - \underbrace{\alpha R(\theta)}_{\text{regularization}} \quad R(\theta) = \begin{cases} \|\theta\|_2^2 = \sum \theta_i^2 \\ \|\theta\|_1 = \sum |\theta_i| \end{cases}$$

۲ تحلیل sentiment

**تحلیل احساسات جنبه های زیادی دارد و میتواند در سطح doc ، token ، ... بررسی شود

* بحث sentiment و aspect خیلی باهم آمیخته هستند ولی معمولاً این دوتا رو باهم برچسب زده نداریم؛

به عبارت دیگر مثلاً ما یک کلمه را نداریم که از یک منظر برچسب مثبت و از یک منظر برچسب منفی داشته باشه

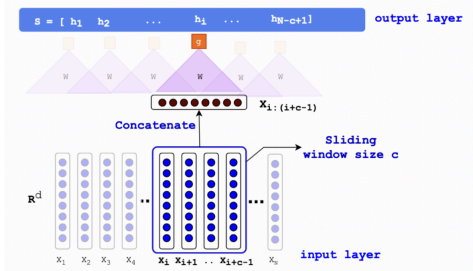
*راه های شناسایی تشخیص مثبت یا منفی بودن جمله عبارتند از :

۱.تعداد کلمات مثبت و منفی را می شماریم اگر تعداد مثبت ها بیشتر بود جمله مثبت و اگر تعداد منفی ها بیشتر بود جمله منفی .

مشکل: احتمال دارد تعداد مثبت و منفی یکی باشد در نتیجه جمله خنثی در نظر گرفته میشود

۲. استفاده از Bayes Naive که فقط کلمات مثبت منفی را در نظر بگیریم یا کل کلمات رو در نظر بگیریم

۳.استفاده از Regression Logistic یا استفاده از network neural ها



شکل ۲:استفاده از cnn برای حل مسئله طبقه بندی

**اگر Multiclass بود از Softmax در تابع فعال سازی لایه اخر استفاده میکنم اگر multilabel ، Multiclass بود از sigmoid .

