1> (a) Pf $I(X,Y) = H(X) - H(X|Y)$

$$= -E_p(\ln p(x)) + E_p(\ln p(x|y))$$

$$= -E_p(\ln p(x) - \ln p(x|y)) \qquad [\text{from linearity of expectations}].$$

$$= -E_p \ln \left( \frac{p(x)}{p(x|y)} \right)$$

$$= -E_p \ln \left( \frac{p(x) p(y)}{p(x,y)} \right) \checkmark$$

(b) $H(X,Y) = -E_p(\ln p(x,y)) \quad$ ⤳ $E_p$

$$= -E_p(\ln p(x|y) p(y))$$

$$= -E_p(\ln p(x|y) + \ln p(y))$$

$$= -E_p \ln p(x|y) \quad - E_p \ln p(y)$$

$$= H(X|Y) + H(Y) \checkmark$$

(c) $I(X,Y|z) = H(X|z) - H(X|Y,z)$

$$= -\{ E_p \ln p(x|z) - E_p \ln p(x|Y,z) \}$$

$$= -\{ E_p \ln p(x|z) - E_p \ln \frac{p(x,Y,z)}{p(Y,z)} \}$$

$$= -\{ \not{p} E_p \ln p(x|z) - E_p \ln \frac{p(x,Y,z) p(z)}{p(Y,z) P(z)} \}$$

①

$$= -\{E_p \ln p(x, z) - E_p \ln \frac{p(x, y \mid z)}{p(y \mid z)}\}$$

$$= - E_p \ln \frac{p(x \mid z) \, p(y \mid z)}{p(x, y \mid z)} \checkmark$$

The conditional independence assumption, that guarantees $I(X, Y \mid Z) = 0$ is:

$$p(x \mid z) \, p(y \mid z) = p(x, y \mid z) \checkmark$$

2> (a)

$$\sigma_n^2 = E((X - EX)^2)$$
$$\sigma_y^2 = E((Y - EY)^2)$$

From Cauchy - Schwarz inequality,

$$\{E[(X - EX)^{\textcircled{2}} (Y - EY)^{\textcircled{0}}]\}^2 \leq E[X - EX]^2 \; E[(Y - EY)^2]$$

$$\{cov(X, Y)\}^2 \leq \sigma_n^2 \sigma_y^2$$

$$\left\{\frac{cov(X, Y)}{\sigma_n \sigma_y}\right\}^2 \leq 1$$

$$\rho_{xy}^2 \leq 1$$

$$|\rho_{xy}| \leq 1 . \checkmark$$

(b)     $X = aY$     $\Rightarrow$     $E(X) = a \beta E(Y)$     [linearity of expectation]

$$\textcircled{i}$$

$$\text{cov}(X, Y) = E((X - EX)(Y - EY))$$

$$= E[(X - EX) \, a(X - EX)]$$

$$= a \, E(a(X - EX)^2)$$

$$= a \, E[(X - EX)^2]$$

$$= a \, \sigma_x^2$$

$$\sigma_x^2 = a^2 \sigma_y^2$$

$$\sigma_x = |a| \, \sigma_y.$$

$$\therefore \quad \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} \neq 1$$

$$\rho_{xy} = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} = \frac{a \sigma_x^2}{|a| \sigma_x \cdot \sigma_x} = \frac{a}{|a|}.$$

$$\rho_{xy} = 1 \qquad \text{when } a > 0 \quad \checkmark$$

$$\rho_{xy} = -1 \qquad \text{when } a < 0. \quad \checkmark$$

(c) $\quad I(X, Y) = -E_p \ln \dfrac{p(x) \, p(Y)}{p(X, Y)}.$

$I(X, Y) = 0$ implies that

$$p(X) \, p(Y) = p(X, Y)$$

We can then write

$$E(XY) = \iint XY \, p(X, Y) \, dx \, dy$$

$$= \iint XY \, p(x) \, p(y) \, dx \, dy$$

$$= \int x \, p(x) \, dx \int y \, p(y) \, dy = E(X) E(Y)$$

③

$$\therefore \quad E(XY) - E(X)E(Y) = 0$$

$$\Rightarrow \quad cov(X,Y) = 0$$

~~Hence $\rho(X,Y)=0$~~

Hence, $\rho_{XY} = 0$ ✓

(d) No, even of $\rho_{XY} = 0$, $I(X,Y)$ can be non zero. Counterexample:
Let X, Y take values, $(1,0), (0,1), (-1,0), (0,-1)$ with probability 0.25.

$$E(X) = 0, \quad E(Y) = 0 \quad E(XY) = 0$$

$$\therefore \quad E(XY) - E(X) \cdot E(Y) = 0$$

$$cov(X,Y) = 0$$

$$I(X,Y) = \sum p(x,y) \ln \frac{p(x)p(y)}{p(x,y)}$$

$$= p(x=1, y=0) \ln \frac{p(x=1)p(y=0)}{p(x=1, y=0)} + p(x=0, y=1) \ln \frac{p(x=0)p(y=1)}{p(x=0, y=1)}$$

$$+ p(x=-1, y=0) \ln \frac{p(x=-1)p(y=0)}{p(x=-1, y=0)} + p(x=0, y=-1) \ln \frac{p(x=0)p(y=-1)}{p(x=0, y=-}$$

$$= -\frac{1}{4} \ln \frac{\frac{1}{4} \cdot \frac{1}{2}}{\frac{1}{4}} - \frac{1}{4} \ln \frac{\frac{1}{2} \cdot \frac{1}{4}}{\frac{1}{4}} - \frac{1}{4} \ln \frac{\frac{1}{4} \cdot \frac{1}{2}}{\frac{1}{4}}$$

$$- \frac{1}{4} \ln \frac{\frac{1}{2} \cdot \frac{1}{4}}{\frac{1}{4}}$$  ✓

$$= 4 \ln 2 \neq 0$$

④

Problem 2

$\left(\dfrac{30}{30}\right)$

2.1>

$$P(Y=1|X) = \frac{P(X|Y=1)\,P(Y=1)}{P(X|Y=1)\,P(Y=1) + P(X|Y=0)\,P(Y=0)}$$

$$= \frac{1}{1 + \dfrac{P(Y=0)}{P(Y=1)}\,\dfrac{P(X|Y=0)}{P(X|Y=1)}}$$

$$= \frac{1}{1 + \exp \ln\left[\dfrac{P(Y=0)}{P(Y=1)} \cdot \dfrac{P(X|Y=0)}{P(X|Y=1)}\right]}$$

$$= \frac{1}{1 + \exp\left[\ln \dfrac{(1-\pi)}{\pi} + \sum_i \ln \dfrac{P(x_i|Y=0)}{P(x_i|Y=1)}\right]}$$

Now, we know that:

$$P(x_i|y=k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}}\,\exp \frac{-(x_i - \mu_{ik})^2}{2\sigma_{ik}^2}$$

$$\therefore \ln P(x_i|y=k) = -\ln \sigma_{ik}\sqrt{2\pi}$$
$$- \frac{(x_i - \mu_{ik})^2}{2\sigma_{ik}^2}$$

$$= -\ln \sigma_{ik}\sqrt{2\pi} - \frac{(x_i^2 - 2x_i\mu_{ik} + \mu_{ik}^2)}{2\sigma_{ik}^2}$$

Hence $k = 0, 1,$

We can write

$$\sum_i [\ln P(x_i | Y=0) - \ln P(x_i | Y=1)]$$

$$= \ln(\sigma_{i1}\sqrt{2\pi}) + \frac{x_i^2}{2\sigma_{i1}^2} - \frac{x_i \mu_{i1}}{\sigma_{i1}^2} + \frac{\mu_{i1}^2}{2\sigma_{i1}^2}$$

$$- \ln(\sigma_{i0}\sqrt{2\pi}) - \frac{x_i^2}{2\sigma_{i0}^2} + \frac{x_i \mu_{i0}}{\sigma_{i0}^2} + \frac{\mu_{i0}^2}{2\sigma_{i0}^2}$$

$$= \ln\frac{\sigma_{i1}}{\sigma_{i0}} + \frac{x_i^2}{2}\left(\frac{1}{\sigma_{i1}^2} - \frac{1}{\sigma_{i0}^2}\right) + x_i\left(\frac{\mu_{i0}}{\sigma_{i0}^2} - \frac{\mu_{i1}}{\sigma_{i1}^2}\right)$$

$$+ \left(\frac{\mu_{i1}^2}{2\sigma_{i1}^2} - \frac{\mu_{i0}^2}{2\sigma_{i0}^2}\right).$$

$$P(Y=1|X) = \frac{1}{1 + \frac{(1-\pi)}{\pi}\exp\left[\sum_i \left\{ \ln\frac{\sigma_{i1}}{\sigma_{i0}} + \frac{x_i^2}{2}\left(\frac{1}{\sigma_{i1}^2} - \frac{1}{\sigma_{i0}^2}\right) + x_i\left(\frac{\mu_{i0}}{\sigma_{i0}^2} - \frac{\mu_{i1}}{\sigma_{i1}^2}\right) + \left(\frac{\mu_{i1}^2}{2\sigma_{i1}^2} - \frac{\mu_{i0}^2}{2\sigma_{i0}^2}\right)\right\}\right]}$$

Therefore this model is NOT the form used by logistic regression. This is because the exponential term in the logistic function has an exponent linear in x, whereas, in this case, the exponent is a <u>quadratic</u> function of x.

2.2) similar to previous case, let us calculate

$$P(Y=1|X)$$

$$P(Y=1|x) = \frac{P(Y=1)\,P(x|Y=1)}{P(Y=1)\,P(x|Y=1) + P(Y=0)\,P(x|Y=0)}$$

$$= \frac{1}{1 + \dfrac{P(Y=0)\,P(x|Y=0)}{P(Y=1)\,P(x|Y=1)}}$$

$$= \frac{1}{1 + \dfrac{(1-\pi)}{\pi}\,\exp\,\ln\,\dfrac{P(x|Y=0)}{P(x|Y=1)}}$$

now

$$\ln P(x|Y=k) = -\ln 2\pi\,\sigma_1\,\sigma_2\sqrt{1-\rho^2}$$

$$-\left[\frac{\sigma_2^2(x_1-\mu_{1k})^2 + \sigma_1^2(x_2-\mu_{2k})^2 - 2\sigma_1\sigma_2(x_1-\mu_{1k})(x_2-\mu_2\rho_k)}{2(1-\rho^2)\sigma_1^2\sigma_2^2}\right]$$

$$= -\ln 2\pi\,\sigma_1\,\sigma_2\sqrt{1-\rho^2}$$

$$-\left[\frac{x_1^2\sigma_2^2 + \sigma_1^2 x_2^2 - 2\sigma_1\sigma_2 x_1 x_2}{2(1-\rho^2)\sigma_1^2\sigma_2^2}\right] \quad \underleftarrow{\text{does not depend on } k}$$

$$-\left[\frac{\sigma_2^2(-2x_1\mu_{1k} + \mu_{1k}^2) + \sigma_1^2(-2x_2\mu_{2k} - \mu_{2k}^2)}{} \right.$$
$$\left. \underleftarrow{\text{depends on } k} \right.$$
$$\frac{-2\sigma_1\sigma_2(\mu_{1k}\mu_{2k} - x_1\mu_{2k} - x_2\mu_{1k})}{2(1-\rho^2)\sigma_1^2\sigma_2^2}$$

The ~~os~~ first and second term does not depend on $k$. Hence when we write.

$\ln \frac{p}{k} (X|Y=0) - \ln P(X|Y=1)$, the first and second terms cancel. We are left with:

$\ln P(X|Y=0) - \ln P(X|Y=1)$

$$= \frac{\sigma_2^2(-2X_1\mu_{11} + \mu_{11}^2) + \sigma_1^2(-2X_2\mu_{21} + \mu_{21}^2) - 2\sigma_1\sigma_2(\mu_{11}\mu_{21} - X_1\mu_{21} - X_2\mu_{11})}{2(1-\rho^2)\sigma_1^2\sigma_2^2}$$

$$+ \frac{\sigma_2^2(-2X_1\mu_{10} + \mu_{10}^2) + \sigma_1^2(-2X_2\mu_{20} + \mu_{20})^2 - 2\sigma_1\sigma_2(\mu_{10}\sigma_{20} - X_1\mu_{20} - X_2\mu_{20})}{2(1-\rho^2)\sigma_1^2\sigma_2^2}$$

$$= \frac{X_1(-2\sigma_2^2\mu_{11} + 2\sigma_2^2\mu_{10} \cancel{\phantom{xxxx}} + 2\sigma_1\sigma_2(\mu_{21} - \mu_{20})) + X_2(-2\sigma_1^2\mu_{21} + 2\sigma_1^2\mu_{20}) + 2\sigma_1\sigma_2(\mu_{21} - \mu_{20}) - 2\sigma_1\sigma_2(\mu_{11}\mu_{21} - \mu_{10}\mu_{20}) + \sigma_2^2(\mu_{11}^2 - \mu_{10}^2) + \sigma_1^2(\mu_{21}^2 - \mu_{20}^2)}{2(1-\rho^2)\sigma_1^2\sigma_2^2}$$

$$= w_0' + w_1'X_1 + w_2'X_2$$

Therefore, we can finally write

$$P(Y=1|x) = \frac{1}{1 + \exp \left( \cdots \right)}$$

$$P(Y=1|x) = \frac{1}{1 + \exp \ln\left(\frac{1-\pi}{\pi}\right) \left( w_0' + w_1' x_1 + w_2' x_2 \right)}$$

$$= \frac{1}{1 + \exp \left( w_0 + w_1 x_1 + w_2 x_2 \right)}$$

This IS the form used by the logistic regression, since the exponent is linear in $(x_1, x_2)$.

3.1) For each category, the no. of parameters we will need to estimate is

$$|V| \times |L| = 5 \times 10^4 \times 10^3$$
$$= 5 \times 10^7.$$

Even if there are 10 categories, and each category has a 100 documents, we will only have $100 \times 1000 = 10^5$ words.

Hence, we have more $P(X_i / Y)$ to evaluate than the no. of words, and so, most of these estimates will be either 0, or inaccurate.

✓

3.2) The overall testing accuracy is 78.521%. The confusion matrix is printed in the next page.

✓

3.3) Adding up individual columns of the confusion matrix suggests some code groups are more miscategorized than others.

'comp' groups are more miscategorized because I think they would contain similar content. (ibm vs mac vs windows)

Similarly electronics, since its a very

broad category, and has similarities to hardware groups in "comp".

"talk. politics. misc" and "talk. religion. misc" are also miscategorized a lot, since they contain broad miscellaneous topics.

3.4> The plot is attached in the next page. I took 10 points between $10^{-5}$ and 1 and the reported values are:

$10^{-5}$: 78.3%

$3.6 \times 10^{-5}$: 78.6%

$1.29 \times 10^{-4}$: 78.97%

$4.64 \times 10^{-4}$: 79.6%

$1.67 \times 10^{-3}$: 79.73%

$5.99 \times 10^{-3}$: 80.06%

$2.15 \times 10^{-2}$: 80.59%

$7.74 \times 10^{-2}$: 80.61%

$2.78 \times 10^{-1}$: 80.39%

$1.00$: 78.11%

At low values of $\alpha$, accuracy drops, because some $P(x_i/Y)$ are taken to be very small, simply because there are no training examples for those words and categories.

At high values of $\alpha$, the prior dominates over the evidence, hence likelihood estimates from training examples are washed out by the prior, which is much bigger. So, accuracy drops.

3.5) The important metric here is $I(Y, X)$
For each word $X_j$ in the vocabulary, we
can define $I(Y, X_j)$ as

$$I(Y, X_j) = H(Y) - H(Y|X_j)$$

To calculate this metric, we can calculate

$$H(Y) = \sum_{k} P(Y=Y_k) \ln P(Y=Y_k).$$

$$H(Y|X_j) = P(X_j=1) H(Y|X_j=1)$$
$$+ P(X_j=0) H(Y|X_j=0).$$

First we calculate
$$P(X_j=1) = \sum_{k} P(X_j=1|Y=Y_k) P(Y=Y_k)$$
$$P(X_j=0) = 1 - P(X_j=1)$$

now, from Bayes's rule,

$$P(Y=Y_k|X_j=1) = \frac{P(X_j=1|Y=Y_k) P(Y=Y_k)}{P(X_j=1)}$$

we know the numerator from estimates made
from training examples. The denominator
is calculated above.
now, from these values we can calculate
$$H(Y|X_j=1) = \sum_{k} P(Y=Y_k|X_j=1) \ln P(Y=Y_k|X_j=1).$$

similarly, for $H(Y|X_j = 0)$, we can write

$$P(X_j = 0 | Y = Y_k) = 1 - P(X_j = 1 | Y = Y_k)$$
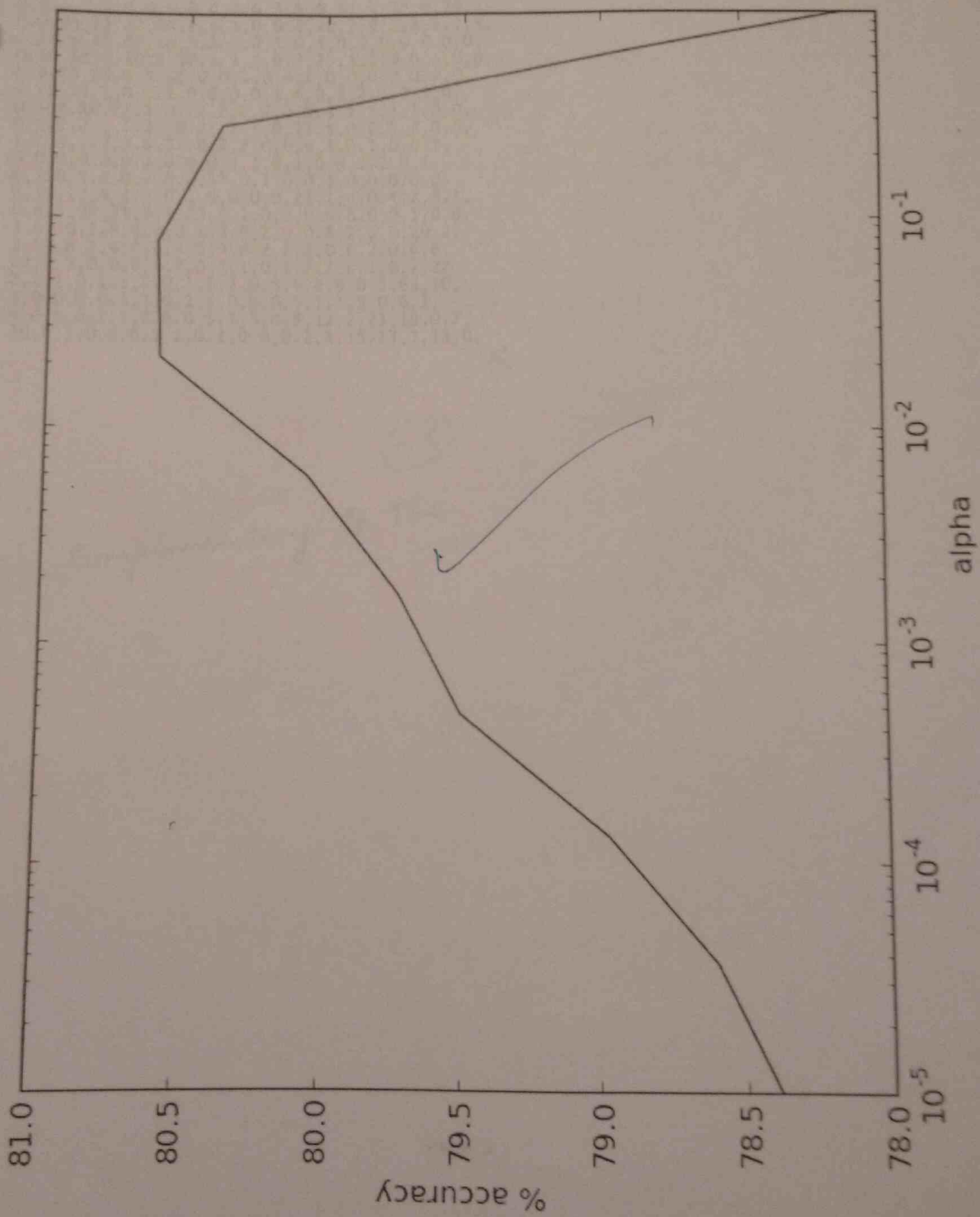
$$\therefore \quad P(Y = Y_k | X_j = 0) = \frac{P(X_j = 0 | Y = Y_k) \, P(Y = Y_k)}{P(X_j = 0)}.$$

$$H(Y|X_j = 0) = \sum_k P(Y = Y_k | X_j = 0) \, \ln P(Y = Y_k | X_j = 0)$$

3.6> The program calculates $I(Y, X_j)$ for each word $X_j$ in the vocabulary, and calculates the words with 100 highest $I(Y, X_j)$. These are printed in the next page.

3.7> 1) The dataset still contains often used words like 'of', 'is', 'we', etc. probably because longer documents containing lots of such words were chosen.

2) It also contains too many words related to computers, such as windows, monitor, ram. etc. This might introduce ~ some bias.

```
0,0,1,0,0,0,0,0,0,0,2,0,3,5,0,11,1,12,6,39,
0,0,33,11,17,54,7,3,1,0,0,3,20,7,8,2,1,1,1,3,
0,13,0,30,13,16,5,1,0,0,1,0,4,0,0,0,0,0,0,0,
0,14,57,0,30,6,32,2,1,1,0,3,25,3,1,0,0,1,0,0,
0,9,19,20,0,3,16,0,0,1,0,4,7,0,0,0,0,0,1,0,
1,22,21,1,0,0,1,0,0,0,0,1,4,0,3,2,1,0,1,0,
0,4,4,10,12,1,0,14,2,2,2,0,8,3,1,1,1,1,0,0,
0,1,2,2,1,17,0,27,1,1,0,11,5,0,0,2,2,0,0,
1,1,3,1,2,3,8,17,0,2,2,0,6,4,1,0,1,0,0,1,
0,0,0,0,0,0,1,0,0,0,4,1,0,1,0,0,1,2,0,1,
0,1,0,1,0,0,2,0,0,17,0,1,0,0,1,0,0,0,0,0,
2,11,12,4,3,5,0,1,0,0,0,0,21,1,4,0,4,2,5,1,
0,8,5,32,21,3,7,13,3,1,0,2,0,8,6,0,0,1,0,0,
3,6,10,1,8,6,4,0,1,3,0,2,9,0,5,2,5,0,10,2,
3,10,8,2,4,4,6,4,0,3,0,2,7,8,0,0,2,0,6,6,
24,1,3,0,0,0,0,2,0,5,1,0,1,7,3,0,1,6,2,27,
2,2,1,0,1,1,2,0,1,2,2,9,3,6,2,0,0,3,63,10,
3,0,0,0,0,1,1,0,1,1,0,0,0,5,1,1,5,0,6,3,
4,0,5,0,1,1,2,6,0,5,1,5,0,8,12,2,23,18,0,7,
26,0,3,0,0,0,1,1,0,1,0,0,0,2,1,15,13,1,13,0,
```

X

(-5)

complimentary to this

100words

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| windows | god | he | scsi | car | drive | space | team | dos | bike | file |
| of | that | mb | game | key | mac | jesus | window | dod | hockey | the |
| graphics | | card | image | his | gun | encryption | | sale | apple | |
| government | | season | we | games | israel | disk | files | ide | controller | |
| players | shipping | | chip | program | was | cars | nasa | win | year | were |
| they | turkish | motif | people | armenian | | play | drives | bible | use | |
| widget | pc | clipper | offer | jpeg | baseball | | bus | my | nhl | |
| software | | is | db | server | jews | os | | israeli | output | data |
| system | who | league | armenians | | for | christian | | christians | | |
| entry | mhz | ftp | price | christ | guns | thanks | church | color | teams | |
| privacy | condition | | launch | him | com | monitor | ram | | | |

for metric $f(x_i, y)$.