Teoría de la información Trabajo Integrador N°2

- Integrantes: Gutierrez Paredes José María
- Email: GutierrezJoseMaria01@hotmail.com
- Fecha de entrega: 06/12
- Repositorio: https://github.com/guteeeeeeeee/teoriaInformacion.git

Índice

Resumen	3
Introducción	4
Desarrollo	
Parte 1	5
Procesamiento	
Codificación	
Compresión	
Descompresión	
Conclusiones	
Parte 2	8
Conclusión	
Anexo	11

Resumen

La realización de este trabajo busca asimilar los conceptos vistos durante la teoría y poder aplicarlos de manera práctica para así poder sacar nuestras propias conclusiones.

Los temas que se tratan son la codificación y compresión de archivos, y los canales de comunicación.

La codificación se estudia porque resulta fundamental ya que permite, mediante el uso de algún algoritmo, que un mensaje transmita la información que contiene en las menores longitudes posibles. De esta manera logrando reducir su tamaño lo que resulta en una compresión de su tamaño original.

Además se estudian los canales de comunicación porque cumplen un papel fundamental en el esquema de la comunicación ya que es el medio por el que se transmite el mensaje enviado desde la fuente a su destino.

Introducción

El problema tratado es poder codificar el mensaje enviado mediante los algoritmos ya estudiados. Lo que se busca es lograr que el mensaje que se quiera transmitir pueda ser enviado de la manera más rápida y ocupando la menor cantidad de espacio posible (compresión). Esto se logra mediante un código compacto. Se estudiaron dos métodos para lograr estos códigos: el algoritmo de Huffman y el de Shannon-Fano. Pero eso no es todo ya que se necesita también poder descomprimir ese archivo para recuperar la información contenida por lo que se emplea una tabla de codificación.

El otro problema tratado es el de analizar las propiedades de ciertos canales de comunicación para así poder compararlos entre sí y obtener conclusiones acerca de sus propiedades. En ellos lo más común es que haya ruido y ocurran pérdidas de la información enviada ya que los símbolos de entrada X pueden convertirse en símbolos de salida Y con ciertas probabilidades propias de cada canal. Provocando de esta manera que el receptor conociendo el símbolo de salida no pueda saber con certeza cuál símbolo fue el que se le envió provocando así problemas en la comunicación del mensaje.

Desarrollo

Parte 1

La realización de esta parte se puede dividir en 4: procesamiento, codificación, compresión y descompresión.

Procesamiento:

En el procesamiento lo que se hizo fue leer el archivo de texto dado e ir almacenando cada palabra con sus datos importantes como su frecuencia de aparición, largo de la palabra y probabilidad de aparición.

Codificación:

En la etapa de codificación lo que se hizo fue aplicar los algoritmos de Huffman y Shannon-Fano para así poder obtener la codificación binaria de cada palabra.

El primero consiste en ir agrupando en un nodo padre los dos símbolos (también son nodos) con menor probabilidad de aparición logrando de esta manera ir formando un árbol binario. Una vez construido, mediante su recorrido se podrá saber la codificación correspondiente de cada símbolo.

En cambio el algoritmo de Shannon-Fano lo que hace es, a partir de ordenar todas las probabilidades de emisión, ir subdividiendo el vector que las contiene en dos partes iguales buscando que tengan las probabilidades de cada parte sean las más cercanas posibles. De esta manera a cada subdivisión se le otorga un bit distinto de codificación.

Una vez obtenidas las codificaciones de cada palabra leída luego de haber aplicado ambos métodos sigue la etapa de compresión.

Compresión:

Primero se escribe la tabla con las palabras y sus correspondientes codificaciones. Lo que primero se indica es en 1 byte la cantidad de palabras codificadas, luego para cada palabra se sigue el siguiente formato: 1 byte para indicar cantidad de caracteres que conforman la palabra, 1 byte por cada carácter, 1 byte para indicar la cantidad de bits de la codificación y 2 bytes para

la codificación. Esto debido a que en C no se puede emplear unidades más pequeñas de escritura que el byte.

1 Byte	1 Byte * (Cant de caracteres)	1 Byte	2 Bytes
Cant caracteres	Palabra	Cantidad bits code	Code

Luego se vuelve a leer el archivo de texto original y se reemplaza, en un archivo binario, cada palabra por su codificación.

Algoritmo	Largo Medio	Rendimiento	Redundancia
Shannon-Fano	9.458161	0.994820	0.005180
Huffman	9.437098	0.997041	0.002959

De esta manera se obtiene el archivo binario con el texto original comprimido en bits junto con la tabla de codificaciones para poder realizar la descompresión.

Para el texto sin codificar cada palabra tiene un largo medio de 4.75 caracteres donde cada uno ocupa 1 byte (8 bits) mientras que para la codificación de Huffman el largo medio de cada una es de 9.437 bits. De esta manera se hace evidente el nivel de compresión de tamaño que se logra a partir de la codificación de las palabras.

Los valores de rendimiento y redundancia son los esperables ya que aunque tienen diferentes métodos su objetivo es el mismo: lograr que cada palabra tenga la longitud de su codificación de acuerdo a las probabilidades de aparición dándoles menores longitudes a las más probables.

Descompresión:

Para que se pueda descomprimir el archivo primero se lee la tabla de tabla de codificaciones y se va almacenando cada palabra con su codificación.

Luego se van leyendo los bits que conforman al archivo binario y se busca dentro de todas las codificaciones almacenadas cual es la palabra a la que le corresponde y en base a eso se van reemplazando, en un archivo de texto, por la palabra a la que representan.

Conclusiones:

La tabla de codificación impacta fuertemente en el tamaño del archivo comprimido. Esto es debido a que para poder notificar al receptor la información necesaria para que pueda descomprimirlo éste debe conocer la codificación de cada palabra.

Almacenar la codificación de cada palabra ocupa 2 bytes pero el verdadero problema ocurre al momento de almacenar cada palabra (hay 4155 en la tabla) ya que se ocupará 1 byte por cada carácter y en promedio está formada por 7.44 caracteres por lo que la tabla solo con las palabras ocupará 30 Kb, y si se le agregan las codificaciones de cada una quedaría de un tamaño de 47 Kb.

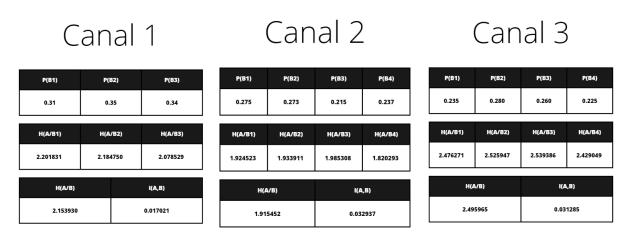
De esta manera se puede concluir que si bien se logró una reducción de tamaño del 20,25% (de 79 Kb a 63 Kb) se le debe dar otro manejo a la tabla de codificaciones ya que su gran tamaño opaca el buen trabajo de la compresión del archivo ya que redujo considerablemente el tamaño del texto de 79 Kb a 16 Kb (sin contar la tabla). Esto fue exitoso debido a que se aprovechó que los algoritmos codificaron las palabras en códigos a nivel de bits y se los fue escribiendo en el archivo uno tras otro logrando de manera que no se desperdicie espacio.

Aunque si bien la tabla es pesada es necesaria ya que sin ella no sería posible la descompresión del archivo ya que solo serían bits que no tienen ningún significado aparente. Ella es la que le da valor a cada conjunto de bits dotándolos de significado que solo el que la posee podrá entender.

Una posible solución sería encontrar la forma de comprimir la representación de las palabras dentro de la tabla ya que es lo que más espacio ocupa dentro de ella.

Segunda parte

Para caracterizar los canales utilizados en la transmisión de información se utiliza para cada uno la matriz de canal la cual nos indica las probabilidades de que salga la salida Y si se emite la entrada X. Éstas se encuentran en el apéndice de este informe.



En el análisis de los canales dados se encuentra que todos contienen ruido.

Las probabilidades de los símbolos de salida son cercanas a 1/k siendo k cantidad de símbolos de salida lo que demuestra que son estadísticamente independientes.

Se puede verificar en cada canal estudiado que las entropías a posteriori son valores cercanos entre sí. Recordando que estos valores son la cantidad media de incertidumbre de conocer el valor de entrada sabiendo el símbolo de salida. Esto indica que el receptor no sabría cual símbolo puede ser el de entrada aun conociendo el de salida. Por lo que los canales estudiados no serían los ideales para transmitir la información dada ya que no habría una correlación entre la entrada y la salida lo que dificultaría la interpretación por parte del receptor.

Otro variable a analizar es la equivocación del canal la cual indica la cantidad de información sobre la entrada que no deja pasar el canal. Mientras mayor sea este número menor fiabilidad tendrá el canal. En el caso ideal de un canal sin ruido este valor sería nulo ya que esto significa que si conozco la salida, conozco perfectamente la entrada por lo que no habría problemas de interpretar

el mensaje enviado ya que durante su transmisión por el canal no hubo ruido que perturbara la información transmitida.

Esto es malo ya que la información obtenida al conocer el símbolo de salida no servirá para identificar al símbolo de entrada por lo que se tiene una pérdida de la información.

La información mutua se puede definir como la cantidad de información que se obtiene de la entrada conociendo la salida. Para un correcto uso del canal este valor tiene que ser máximo.

En los canales estudiados este valor es cercano a cero lo que significa que no son los óptimos para ser utilizados en la transmisión de las fuentes dadas.

Dado que los canales ya están dados y no se pueden modificar lo que se podría hacer es modificar las probabilidades de emisión buscando minimizar la equivocación del canal y maximizar la información mutua.

Si se desea enviar un mensaje el canal elegido sería el segundo ya que este tiene el menor ruido aunque tengo un bajo valor en la cantidad de información.

Conclusión

Como conclusión se puede decir que fue un trabajo interesante ya que se aplicaron conceptos vistos en la teoría que luego de ser implementados quedaron mucho más claros.

Se evidenció la importancia de la codificación ya que los textos que utilizamos generalmente tienen poco rendimiento y mucha redundancia. Lo que significa que para transmitir la información que contienen ocupan mucho más espacio del necesario en el almacenamiento y en su transmisión. Además las palabras utilizadas resultan ser más largas que la información que transmiten. Pero a partir de otorgar longitudes de acuerdo a la probabilidad de aparición se logra una lectura más rápida de las palabras más utilizadas y por lo tanto a la hora de almacenarlas una reducción importante de su tamaño respecto al texto original.

Lo que también se dio cuenta es que a la hora de comprimir se le debe dar otro tratamiento a la tabla de codificaciones ya que, para indicar su codificación, cada palabra ocupará mucho más espacio del necesario por lo que lo recomendable sería reducir su tamaño de alguna manera porque sino toda la compresión del texto original se ve opacada por ella. Sin embargo es muy importante ya que es la responsable de darle las herramientas al receptor para que pueda descomprimir y recuperar el mensaje originalmente enviado. Sino ¿cuál sería el sentido de comprimir algo que no se puede recuperar?

También se aprendió sobre la compresión sin pérdidas ya que se redujo el tamaño del archivo propuesto y luego de descomprimirlo se obtuvo el mismo texto originalmente enviado.

También hay que tener cuidado con los canales de comunicación utilizados ya que pueden provocar problemas a la hora de transmitir el mensaje que deseamos que llegue y entienda el receptor. En los casos estudiados los canales tenían ruido y pérdida de la información por lo que el receptor al recibir el mensaje luego de pasar por el canal no podrá lo que le fue enviado ya que las salidas que le llegan no tienen correspondencia con los símbolos enviados.

Resulta fundamental el análisis para detectar canales que se adecuen a la fuente que queremos utilizar ya que aunque se hagan correctamente los pasos para armar y codificar un mensaje finalmente al receptor le llegará algo que no podrá entender.

Anexo

Canal 1:

Símbolo	P(i)
S1	0,2
S2	0,1
S3	0,3
S4	0,3
S5	0,1

Matriz del canal				
	B1	B2	В3	
S1	0,3	(a) 0,3	1-B1-B2 0,4	
S2	(b) 0,4	0,4	1-B1-B2 0,2	
S3	0,3	(a) 0,3	1-B1-B2 0,4	
S4	(a) 0,3	0,4	1-B1-B2 0,3	
S5	0,3	(b) 0,4	1-B1-B2 0,3	

Canal 2:

Símbolo	P(i)
S1	0,25
S2	0,33
S3	0,27
S4	0,15

Matriz del canal					
	B1	B2	В3	B4	
S1	0,2	(a) 0,3	(b) 0,2	1-B1-B2-B3	0,3
S2	(a) 0,3	0,3	(b) 0,2	1-B1-B2-B3	0,2
S3	(c) 0,3	(b) 0,2	0,2	1-B1-B2-B3	0,3
S4	(c) 0,3	0,3	(a) 0,3	1-B1-B2-B3	0,1

Canal 3:

Símbolo	P(i)
S1	0,15
S2	0,1
S3	0,20
S4	0,25
S5	0,14
S6	0,16

Matriz del canal					
Prob(S1/B1)	B1	B2	B3	B4	
S1	0,2	(a) 0,3	(b) 0,2	1-B1-B2-B3	0,3
S2	(c) 0,3	(a) 0,3	0,3	1-B1-B2-B3	0,1
S3	(b) 0,2	0,2	(c) 0,3	1-B1-B2-B3	0,3
S4	(a) 0,3	0,3	(b) 0,2	1-B1-B2-B3	0,2
S5	0,2	(c) 0,3	(a) 0,3	1-B1-B2-B3	0,2
S6	(b) 0,2	(c) 0,3	0,3	1-B1-B2-B3	0,2