

Hacia una arquitectura flexible para sistemas de predicción de palabras: propuesta de diseño y evaluación

Sira E. Palazuelos Cagigas, José L. Martín Sánchez, Lisset Hierrezuelo Sabatela
Universidad de Alcalá
Escuela Politécnica Superior. Campus universitario s/n. 28805. Alcalá de Henares.
{sira, jlmartin, lisset}@depeca.uah.es

Javier Macías Guarasa
Universidad Politécnica de Madrid
Ciudad Universitaria s/n. 28040. Madrid.
macias@die.upm.es

Resumen: La predicción de palabras es un proceso mediante el cual se intenta predecir la palabra que un usuario está escribiendo o va a escribir, de forma simultánea a la propia escritura. Se utiliza fundamentalmente para disminuir el esfuerzo que requiere la entrada de textos en aplicaciones orientadas a personas con discapacidad. En el presente artículo se describe una propuesta para la arquitectura de un avanzado sistema de predicción de palabras. Esta arquitectura propuesta es modular y sumamente flexible, con bloques independientes e interfaces claras, de forma que permite utilizar fácilmente distintos algoritmos de predicción, e incluso predecir en varios idiomas. En la actualidad, se ha implementado un sistema de predicción de palabras según la arquitectura propuesta, que cuenta con un diccionario fijo para castellano con 160000 entradas, la posibilidad de generar diccionarios personales del usuario, y métodos de predicción basados en gramáticas probabilísticas de palabras y categorías y en el uso de una gramática formal. Dependiendo de la configuración utilizada, el sistema puede predecir correctamente más del 94% de las palabras, y reducir el número de pulsaciones necesarias para escribir el texto en más del 50%. Prueba de la flexibilidad de la arquitectura es que esta implementación se ha incluido en varios programas de predicción en distintos idiomas.

Palabras clave: Predicción de palabras, gramática probabilística independiente del contexto, gramáticas probabilísticas de palabras y categorías, ahorro de pulsaciones, personas con discapacidad, ayudas técnicas.

Abstract: Word prediction is a process that tries to guess the word a user is writing, at the same time he/she is doing it. It is mainly used to decrease the effort needed to write a text in applications devoted to people with disabilities. In this paper, we describe a proposal for the architecture of advanced word prediction systems. The proposed architecture is modular and flexible, with common interfaces between the modules to allow the use of different prediction algorithms or even the prediction in a different language. A word prediction system has been implemented according to this architecture. It consists of a 160000 entries lexicon for Spanish (that cannot be modified by the user), the possibility to create and store personal lexicons, prediction methods based on words and POS (parts of speech) probabilistic grammars and a formal grammar. The system can predict more than 94% of the words and reduce the keystrokes needed to write a text more than 50%, depending on the configuration of the prediction. The flexibility of the architecture allowed us to include the word prediction system in several applications for different languages.

Keywords: Word prediction, PCFG, probabilistic grammars based on words and POS, keystroke savings, people with disabilities, technical aids.

1 Introducción

La predicción de palabras es un proceso mediante el cual se intenta adivinar la palabra que un usuario va a escribir o está escribiendo, de forma simultánea a la propia escritura. Donde más intensivamente se utiliza la predicción es en los sistemas de ayuda a la escritura para personas con discapacidad, especialmente física. Estos usuarios no pueden manejar el teclado convencional y hay casos en que el acceso al ordenador se lleva a cabo mediante aplicaciones especiales, en las que se barren secuencialmente las opciones disponibles, para que puedan ser seleccionadas con un dispositivo que aproveche cualquier resto motor del usuario, en general, algún tipo de pulsador.

La lentitud de este método de escritura hace aconsejable el uso de técnicas de aceleración que permitan al usuario introducir texto a mayor velocidad. Además, el esfuerzo que necesitan estas personas para realizar cada pulsación puede ser muy grande, por lo que también es importante reducir el número de pulsaciones necesarias.

Es en este punto en el que la predicción de palabras puede aportar una notable mejora: al mostrar la palabra antes de que el usuario acabe de escribirla, es suficiente con seleccionarla, y se elimina la necesidad de completarla. De esta forma se reduce el número de pulsaciones necesarias para escribir el texto, y, por tanto, el esfuerzo necesario.

En la literatura hay una serie de trabajos sobre predicción de palabras para dispositivos móviles y PDAs, pero no son aplicables a nuestro caso ya que la fuente de la dificultad está en la prácticamente nula movilidad de los usuarios y no en las dificultades del dispositivo para la entrada de datos.

La estructura del artículo es la siguiente: en primer lugar se detalla la arquitectura del sistema, sus bloques y el flujo de datos dentro del sistema. A continuación se describen las características generales de las técnicas implementadas sobre la arquitectura propuesta, y que serán evaluadas en la sección 3, primero de forma aislada y luego en su conjunto. De esa evaluación se obtiene el flujo de trabajo, que se explica a continuación. Posteriormente, se describe brevemente algunas de las aplicaciones en las que está incluida la predicción y su papel en ellas, y, por último, las conclusiones.

2 Arquitectura y flujo de datos del sistema de predicción de palabras

En la Figura 1 se muestra nuestra propuesta para la arquitectura del sistema de predicción de palabras (en su versión simplificada), junto con flechas que indican el flujo de datos interno.

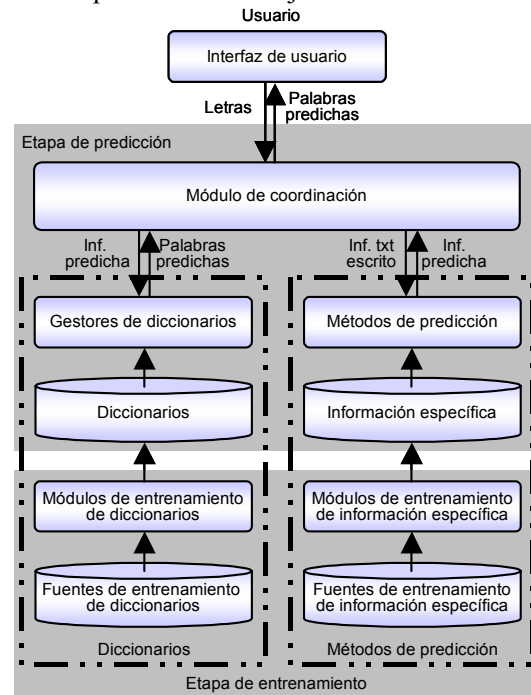


Figura 1: Arquitectura simplificada y flujo de datos del sistema de predicción de palabras.

Sus características fundamentales son la modularidad y flexibilidad, ya que ha sido diseñada para permitir utilizar distintos algoritmos de predicción, e incluso predecir en distintos idiomas simplemente cambiando las distintas fuentes de conocimiento implicadas.

Gran parte de su flexibilidad se debe a la existencia del módulo de coordinación, altamente configurable, que permite utilizar uno o varios métodos de predicción o diccionarios, dar prioridad a unos sobre otros, o incluso combinar sus resultados.

En los siguientes apartados se describe con mayor nivel de detalle cada uno de los módulos y las fuentes de conocimiento utilizadas, así como las técnicas implementadas.

2.1 Interfaz de usuario

Es la interfaz del programa en el que esté integrado el motor de predicción. Proporciona al módulo de coordinación el texto introducido por el usuario, muestra las palabras predichas y las inserta en el texto al ser seleccionadas, no

formando parte del sistema de predicción propiamente dicho.

2.2 Gestores de diccionarios

Existe uno por cada diccionario. Son los encargados de seleccionar la lista de las palabras más frecuentes del diccionario que cumplen las restricciones impuestas por los métodos de predicción.

2.3 Diccionarios

Los diccionarios son las fuentes de información básicas utilizadas por todos los métodos de predicción. Contienen tanto las palabras, como toda la información disponible asociada a ellas y que pueda ser necesaria en cualquier punto del proceso de predicción (si está disponible).

En nuestra arquitectura proponemos la inclusión de tres tipos de diccionarios cuya combinación abarca todas las posibilidades planteables en un sistema de predicción de palabras:

Diccionario general. Es el que debe proporcionar la información frecuencial y gramatical necesaria para todos los métodos de predicción. Es estático, es decir, no evoluciona con el uso del programa, ni es ampliable por el usuario. En nuestra implementación actual contiene más de 160.000 entradas (significante, lema, frecuencia absoluta en el texto de entrenamiento, categoría gramatical y conjunto de rasgos asociados) correspondientes a 130.000 significantes diferentes.

Diccionario personal. Se genera en cada sesión, empezando vacío al principio de la misma, y borrándose cuando acaba. Su objetivo es introducir el vocabulario nuevo a medida que se escribe y ajustar la frecuencia de cada palabra a la que tiene en el texto concreto que está siendo escrito. Por su propio funcionamiento su rendimiento mejora a medida que la sesión avanza.

Diccionarios temáticos. Diseñados para no perder la eficacia del diccionario personal al acabar cada sesión. Se pueden generar diccionarios adaptados a distintos contextos semánticos. De esta manera el usuario puede disponer de diccionarios mucho más eficientes que el generar al escribir textos sobre un determinado tema.

2.4 Módulos de entrenamiento de diccionarios

Este bloque engloba los procedimientos necesarios para generar los diccionarios a partir

de las fuentes de entrenamiento de los mismos (normalmente, aunque no siempre, en una etapa previa a la predicción). Para la generación se usan como fuentes de conocimiento los textos disponibles, categorizados o no, y se someten a procesos automáticos y/o manuales (dependiendo del diccionario).

2.5 Métodos de predicción de palabras

Son los algoritmos encargados de decidir las restricciones que deben cumplir las palabras predichas (conjunto de categorías, frecuencias, rasgos, significante o lema que deben tener) a partir de la información proporcionada por el módulo de coordinación y de las fuentes de información específica de cada uno. Los métodos de predicción incluidos en la implementación actual sobre la arquitectura propuesta están basados en las siguientes gramáticas:

Gramáticas probabilísticas de palabras.

El método más sencillo es el basado en la frecuencia absoluta de las palabras (*unigramas*). La única información que necesita es la parte ya escrita de la palabra en curso, y la lista de palabras que genera (cuando la restricción se aplica a un diccionario) es el conjunto de palabras más frecuentes de dicho diccionario que comienzan por dicha cadena de letras. Se puede aplicar a todos los diccionarios.

Para aumentar la cantidad de información a incluir en la predicción, se incrementa la longitud de la historia de palabras a considerar, pasando a los métodos basado en *bigramas*, *trigramas*, *etc.* que incluyen información sobre secuencias de dos, tres o más palabras. En la implementación actual, los *bigramas* y *trigramas* se aplican al diccionario personal y a los temáticos, siendo los que mejores resultados proporcionan. Su principal problema es la cantidad de texto que necesitan para entrenarse de forma robusta, especialmente cuando se consideran secuencias largas.

Gramáticas probabilísticas de categorías.

Determinan la probabilidad de las categorías a las que debe pertenecer la palabra siguiente, en base a las categorías de las palabras anteriores. Solucionan parte del problema de los métodos estadísticos básicos, reduciendo la cantidad de texto necesario para ser entrenados. Además, introducen información gramatical de corto alcance en el proceso de predicción, al basarse en el cálculo de las probabilidades de cada secuencia de categorías (dos categorías, *bipos* o

tres categorías, *trijos*), e introducir esa información en el proceso de predicción.

Estos métodos han sido mejorados con un sistema de comprobación de concordancia de rasgos que rechaza las palabras cuyos rasgos no son adecuados según una serie concreta de reglas. Por ejemplo, comprueban la concordancia en género y número en secuencias de artículo, nombre y adjetivo.

Para incrementar la robustez de las gramáticas probabilísticas se han introducido mecanismos de suavizado, tanto los clásicos, basados en descuento como los de retroceso directo a modelos de menos complejidad. (Goodman, 2001).

Gramática probabilística independiente del contexto modificada. Se ha introducido un analizador sintáctico de Early (Allen, 1994) (Early, 1970) para mejorar la calidad de la lista de palabras predichas. Permite predecir la categoría, rasgos y, en algunos casos, lema o significante, de la palabra siguiente de acuerdo con el conjunto de reglas incluido en el sistema.

Se han introducido mecanismos para soportar el tratamiento de rasgos (un total de 63 en la implementación actual, sobre los se puede comprobar concordancia, imponer o prohibir su aparición en una determinada palabra dependiendo de su entorno, etc.), lo que, entre otras consecuencias, implica un tremendo incremento de potencia expresiva, de modo que algunas de sus reglas actuales equivaldrían a más de 200 reglas tradicionales. Además, soporta que podamos imponer o prohibir lemas o significantes concretos en una posición cualquiera de una regla. La inclusión de todas estas características implica una enorme optimización, reduciendo en gran medida el número de reglas (118 en la implementación actual) y de categorías (21 en la versión actual) total del sistema, sin perder potencia expresiva.

La ambigüedad de las palabras y la posibilidad de que el fragmento escrito de la frase en curso se corresponda con varias reglas se gestiona manteniendo activas todas las reglas que el fragmento pueda estar siguiendo y añadiendo información frecuencial a las reglas, que se combina con la frecuencia de las palabras de los distintos diccionarios.

Para mejorar su eficacia se han modificado las categorías gramaticales asignadas a cada palabra, y, en vez de utilizar el conjunto tradicional de categorías, se agrupan las palabras según su comportamiento en relación con las demás: por ejemplo, las palabras “a” y

“de” se han separado del resto de las preposiciones, por su comportamiento diferente cuando son seguidas por un artículo. Así mismo, “al” y “del” también tienen su propia categoría. Del mismo modo se ha procedido con todas las palabras hasta conseguir el conjunto de categorías más adecuado.

2.6 Información específica

La necesaria para completar la información que requiere cada método de predicción (matrices de n-gramas, reglas de contexto libre, etc.).

2.7 Módulos de entrenamiento de información específica

Engloban los procedimientos necesarios para generar la información que necesita cada método de predicción a partir de sus fuentes de entrenamiento.

2.8 Métodos auxiliares de aceleración

Además de la predicción de palabras, en la implementación actual se incluye una serie de técnicas auxiliares con una base pragmática que ayudan a mejorar el rendimiento del sistema global. En concreto:

- Eliminación de las palabras previamente rechazadas por el usuario: si se le ha mostrado una determinada palabra y, tras verla, no la ha escogido, el sistema supone que no era la adecuada y no volverá a mostrarla (hasta que proceda).
- Predicción de las terminaciones más frecuentes que comiencen por la última letra escrita por el usuario.
- Inserción automática de los espacios que siguen a los signos de puntuación o las mayúsculas después de punto.

2.9 Módulo de coordinación

Es el elemento clave de la arquitectura de predicción propuesta, y que permite dotarla de gran flexibilidad. Establece el flujo de trabajo del sistema de predicción. Dicho flujo se determina a partir de la experimentación sobre cada implementación concreta. En el apartado 4 detallaremos el adoptado en la implementación actual.

3 Evaluación

La evaluación de los sistemas de predicción de palabras es una tarea que todavía no está estandarizada, por lo que cada centro de

investigación ha proporcionado los resultados que ha considerado interesantes para su tarea concreta. Esta situación dificulta la comparación entre sistemas a partir de resultados publicados sobre ellos.

Además de los parámetros de la propia estrategia general de presentación de propuestas (por ejemplo, la longitud de la lista de palabras candidatas que se presenta al usuario), influyen otros muchos factores relacionados tanto con la interfaz del programa, como con el idioma o el texto en el que se esté escribiendo, que pueden producir variaciones muy significativas en los resultados. En (Palazuelos et al., 1999) se puede encontrar una exposición y discusión de los factores que afectan a los resultados de la predicción de un sistema determinado.

Por las razones expuestas anteriormente, para poder verificar realmente la eficacia de un determinado método o diccionario, y, sobre todo, para poder hacer comparaciones entre ellos, es necesario, realizar series de experimentos entre los que lo único que varíe sea el o los factores concretos a evaluar, ya que, en caso contrario, la influencia de las demás modificaciones puede cambiar totalmente los resultados.

De cada experimento de la serie se proporciona el número de palabras predichas correctamente y el número de pulsaciones ahorradas, siendo éste último el parámetro en el que estamos más interesados, ya que es indicativo directo del esfuerzo que se reduce al usuario. Ambos resultados se proporcionan en valor absoluto y relativo con respecto a la escritura del texto sin predicción, y con bandas de fiabilidad para una confianza del 95%.

Los experimentos se han realizado con textos cortos de distintos tipos, ya que en algunos casos la longitud influye en el rendimiento de los métodos de predicción y los diccionarios, y se ha decidido evaluar en las condiciones más cercanas a las que tendría un usuario típico del sistema.

3.1 Evaluación de los modelos basados en gramáticas probabilísticas de categorías

En la Tabla 1 se muestra la comparación de los resultados de ciertas variaciones del método basado en *trijos* con el uso de los unigramas. Como puede verse, el uso de los *trijos* sin ningún tipo de retroceso ni suavizado empeora los resultados, mejorando con la utilización de técnicas de suavizado, pero sin conseguir

superar el método sin información gramatical. El número de pulsaciones ahorradas sólo mejora al utilizar retroceso a *bipos* y *unigramas* si los *trijos* no ofrecen predicciones válidas. Se puede observar que el número de palabras predichas es menor que en el caso de predicción con unigramas. Esto es así porque sin información gramatical las palabras pequeñas aparecen siempre como primera opciones y se predicen en todos los casos. Al añadir información gramatical se predice un mayor número de palabras largas (por eso mejora el número de pulsaciones ahorradas), pero menos palabras cortas, que son menos útiles para el usuario.

| Método de predicción | Núm. palabras predichas | Núm. pulsac. ahorradas |
|---------------------------------------------------------------|-------------------------|------------------------|
| Unigramas del dicc. general | 89.35%±0.31 | 36.37%±0.19 |
| <i>Trijos</i> sin suavizado | 75.84%±0.43 | 32.83%±0.18 |
| <i>Trijos</i> con suavizado | 78.82%±0.41 | 33.81%±0.18 |
| <i>Trijos</i> con retroceso a <i>bipos</i> y <i>unigramas</i> | 88.51%±0.32 | 37.87%±0.19 |

Tabla 1: Efecto del uso de modelos basados en gramáticas probabilísticas de categorías.

3.2 Evaluación del modelo basado en una gramática probabilística independiente del contexto

El modelo formal fue probado con una gran cantidad de textos de diferentes temáticas. Ofrecemos los resultados de dos de ellos, por considerarlos especialmente significativos.

| Método de predicción | Núm. palabras predichas | Núm. pulsac. ahorradas |
|----------------------------------------------------------------------------------|-------------------------|------------------------|
| Unigramas del dicc. general | 91.37%±0.42 | 37.83%±0.3 |
| <i>Trijos</i> con retroceso a <i>bipos</i> y <i>unigramas</i> | 90.59%±0.43 | 40.18%±0.3 |
| Gramática formal con retroceso a <i>trijos</i> , <i>bipos</i> y <i>unigramas</i> | 88.78%±0.19 | 39.93%±0.3 |

Tabla 2: Comparación entre la gramática formal y los *trijos* con retroceso a *bipos* y *unigramas* sobre un texto técnico.

En la Tabla 2 se muestran los del experimento llevado a cabo sobre un texto de carácter técnico de 17.500 palabras que necesita 122.425 pulsaciones para ser escrito sin predicción. Como se puede observar, el método de predicción basado en la gramática independiente de contexto aumentada no proporciona mejores resultados que los métodos basados en *tripos* con retroceso, en parte porque la gramática formal ha sido entrenada para generar textos de diferente estilo: sin listas, paréntesis, y menor número de términos extranjeros, tablas y números, que no están considerados en las reglas del sistema actual.

En la Tabla 3 mostramos los resultados sobre un texto que se corresponde en mayor medida con el estilo de la gramática. Se trata de un fragmento de un texto de carácter técnico de 2.058 palabras, que necesita 13.957 pulsaciones para ser escrito sin ayuda de la predicción.

| Método de predicción | % palabras predichas | % pulsac. ahorradas |
|--------------------------------|----------------------|---------------------|
| Unigramas del dicc. general | 93.49%±1.07 | 41.07%±0.82 |
| <i>Tripes</i> con retroceso | 93.25%±1.08 | 43.90%±0.82 |
| Gramática formal con retroceso | 92.27%±1.15 | 44.37%±0.82 |

Tabla 3: Comparación entre la gramática formal y los *tripos* sobre un fragmento del mismo estilo que la gramática.

En este caso se puede apreciar un incremento de eficacia con respecto al modelo basado en *tripos* con retroceso a *bipos* y *unigramas*, en número de pulsaciones ahorradas. De todas formas, los resultados no se pueden considerar concluyentes (no son estadísticamente significativos porque los intervalos de confianza se solapan) pero podemos apreciar la tendencia a mejorar los resultados cuando el estilo de la gramática y del texto son similares.

3.3 Evaluación del diccionario personal y de los temáticos

El efecto del diccionario personal es más significativo para textos grandes, pero la evaluación debe hacerse tanto con textos grandes como pequeños (que serán los que escriba el usuario).

En la Tabla 4 se pueden observar los excelentes resultados obtenidos debido al

aprendizaje del vocabulario y la adaptación de las frecuencias del diccionario personal sobre el proceso de predicción, como mejora de los dos factores que se evalúan. En este caso usamos un texto de carácter técnico con 37.496 palabras en total, y que sin predicción requiere 253.932 pulsaciones para escribirlo. En esta serie de experimentos se considera que cuando no haya predicciones del diccionario personal las palabras provendrán del diccionario general.

| Núm. Palabras escritas | % palabras predichas | % pulsac. ahorradas |
|------------------------|----------------------|---------------------|
| 100 | 86.00%±6.80 | 39.32%±3.60 |
| 200 | 89.00%±4.34 | 42.67%±2.57 |
| 500 | 90.00%±2.63 | 45.56%±1.65 |
| 1000 | 91.40%±1.74 | 46.35%±1.19 |
| 5000 | 92.26%±0.74 | 49.49%±0.53 |
| 10000 | 92.32%±0.52 | 50.55%±0.37 |
| 37496 | 92.32%±0.27 | 51.52%±0.19 |

Tabla 4: Evolución de los resultados del diccionario personal al escribir un texto.

En el caso de los diccionarios temáticos, su uso proporcionará mejores resultados, siempre que se correspondan realmente con el tema del que se está escribiendo, como se muestra en la Tabla 5, en que los textos de prueba y entrenamiento pertenecen a un fragmento de El Quijote, de 7.813 palabras y 46.008 pulsaciones. En caso contrario, la predicción que proporcionan es contraproducente.

| Diccionarios utilizados | % palabras predichas | % pulsac. ahorradas |
|-------------------------------------|----------------------|---------------------|
| Dicc. general y personal | 88.74%±0.7 | 39.00%±0.4 |
| Dicc. general, personal y temático. | 92.45%±0.6 | 42.12%±0.4 |

Tabla 5: Resultados de la predicción utilizando un diccionario temático del mismo tema sobre el que se está escribiendo.

3.4 Efecto de las técnicas auxiliares de aceleración

En la Tabla 6 podemos ver los resultados al ir aplicando sucesivamente los mecanismos auxiliares de aceleración sobre el mismo texto usado en el apartado anterior, donde los ahorros en pulsaciones son realmente importantes. Estos mecanismos siempre ahorran, al menos, pulsaciones. Alguno, como la adición de

espacios y mayúsculas automáticamente tras los puntos, ahorra más pulsaciones que algunos métodos de predicción, aunque no implica mayor número de palabras predichas.

| Método auxiliar | % palabras predichas | % pulsac. ahorradas |
|------------------------------------------|----------------------|---------------------|
| Experimento sin mecanismos aux. | 92.45%±0.6 | 42.12%±0.4 |
| Añadición espacios y mayúsc. tras punto | 92.45%±0.6 | 45.15%±0.4 |
| Predicción de terminaciones | 93.86%±0.5 | 45.87%±0.5 |
| Eliminando palabras que aparecen 2 veces | 94.10%±0.5 | 46.18%±0.5 |
| Eliminando palabras que aparecen 1 vez | 94.73%±0.5 | 47.09%±0.5 |

Tabla 6: Efecto de los mecanismos auxiliares de aceleración.

En (Palazuelos, 2001) se describe en detalle toda la arquitectura del sistema de predicción, y los distintos módulos, particularizando en la implementación actual, dedicando especial atención al método basado en la gramática formal aumentada, incluyendo el conjunto de categorías, reglas y sus probabilidades. Además, se proporciona mayor información sobre el proceso de evaluación, tanto objetiva como subjetiva.

4 Flujo de trabajo en la implementación actual

Fruto de la evaluación realizada, se ha determinado que el flujo óptimo del módulo de coordinación, que se describe a continuación:

1. Recibe de la interfaz del programa el último carácter escrito por el usuario:
 - 1.1. Si es letra, la añade a la palabra actual que está escribiendo y genera una lista con las terminaciones más frecuentes que comienzan por dicha letra.
 - 1.2. Si es un signo, da por finalizada la palabra actual, añade automáticamente un espacio, y, si el signo era un punto, la siguiente letra será mayúscula.
2. Con la palabra anterior a la actual, consulta los *bigramas* del diccionario personal, y obtiene la lista de palabras más frecuentes que la han seguido y comienzan por el fragmento escrito de la palabra actual.
3. Repite el proceso para obtener los *bigramas* del diccionario temático activo.

4. Con las palabras anteriores a la actual y sus categorías y rasgos, consulta al método de predicción más potente, para que le proporcione la lista de restricciones que debe cumplir la palabra que está escribiendo en la actualidad: conjunto de categorías a la que debe pertenecer y su probabilidad, rasgos que debe cumplir para cada una, y lema/significante si aplica.
5. Si la lista de restricciones está vacía (por ejemplo, en el caso del analizador, si la estructura de la frase que está escribiendo el usuario no está incluida en la gramática), consulta a otros métodos menos potentes, hasta que obtiene una lista de condiciones.
6. Consulta todos los diccionarios para obtener las palabras más frecuentes de cada uno que cumplen las restricciones.
7. Elimina las palabras rechazadas por el usuario de todas las listas.
8. Combina las listas de palabras para generar una única lista dando prioridad a los bigramas del diccionario temático y personal, por ese orden, y, si la lista no está completa, se añaden las más frecuentes del diccionario personal, temático y general, por ese orden. También es posible combinar las listas ponderando la importancia de cada diccionario, pero el factor de ponderación ideal varía en gran medida con la longitud del texto escrito, por lo que en la actualidad no se utiliza esta opción.

5 Aplicaciones

El motor de predicción implementado sobre la arquitectura propuesta en este artículo ha sido incluido en varias aplicaciones:

- *MELE*, sistema de aprendizaje del castellano como segunda lengua. (Palazuelos et al., 1999a).
- *PROFET*, teclado virtual con predicción de palabras multi-idioma donde el sistema descrito en este artículo realiza la predicción para castellano. (Carlberger et al., 1997).
- *PredWin*, editor de texto y operaciones matemáticas y comunicador para personas con problemas físicos. (Palazuelos et al., 1999b).
- *Comunicador* es un sistema de comunicación para personas con graves discapacidades físicas basado en un entorno gráfico de acceso a mensajes. (Domínguez et al., 2004)

En cada una de ellas tiene un objetivo diferente: En *MELE* aumenta la calidad del texto escrito por la persona extranjera, ya que las palabras elegidas de la lista de predicciones están correctamente escritas. En *PROFET*, *PredWin* y *Comunicador* ayuda a la escritura de textos y la composición de mensajes para comunicación con menor esfuerzo.

Estos programas se están utilizando en la actualidad, y los usuarios y terapeutas consideran que la predicción es beneficiosa en términos de reducción de esfuerzo al escribir los textos y de aumento de la calidad de los mismos. Sin embargo, el aumento de velocidad no siempre llega a producirse, ya que la aceleración que se produce al reducir el número de pulsaciones puede verse contrarrestada por el tiempo necesario para localizar las palabras en la lista de candidatas predichas (Magnuson y Hunnicutt, 2002).

6 Conclusiones

En este artículo se ha descrito y evaluado una propuesta de arquitectura para un sistema avanzado de predicción de palabras que actualmente está incluido en varios programas de ayuda a la escritura y comunicación.

La arquitectura del sistema es modular, de gran flexibilidad, con módulos independientes e interfaces bien definidas entre bloques, de forma que permite de fácilmente cambiar los algoritmos de predicción, e incluso predecir en diferentes idiomas. En la actualidad se dispone de versiones del motor de predicción para castellano, inglés y sueco, y el trabajo de adaptación se ha visto tremendamente minimizado por las ventajas de la arquitectura propuesta.

La implementación actual cuenta con un diccionario para castellano con 160000 entradas fijo, y la posibilidad de generar diccionarios personales del usuario. Acerca de los métodos de predicción, incluye algoritmos basados en gramáticas probabilísticas basadas en palabras y en categorías y en una gramática formal.

Dependiendo de la configuración usada, el sistema puede predecir correctamente más del 94% de las palabras del texto, y reducir más del 50% el número de pulsaciones necesarias para escribir el texto.

Agradecimientos

Parte de este trabajo está siendo financiado por la Comunidad de Madrid, en el proyecto

GR/SAL/0814/2004 “Sistema de comunicación global para personas con discapacidades físicas y/o intelectuales”.

Bibliografía

- Allen J. 1994. “Natural language Understanding”. Benjamin/Cummings Publishing Company Inc 2ª Ed.
- Carlberger A., Carlberger J., Magnuson T., Palazuelos S.E., Hunnicutt M.S. & Aguilera S. 1997. Profet, a New Generation of Word Prediction: An Evaluation Study. Proceedings of the Workshop on NLP for Communication Aids, ACL/EACL'97 (pp. 23-28). Ed. ACL. Madrid.
- Domínguez L.M., Palazuelos S.E., Martín J.L., Gómez M., García J.C. 2004. Comunicador gráfico para usuarios de pulsador con acceso óptimo a mensajes. Actas Iberdiscap 2004. Páginas: 184-188. San José. Costa Rica.
- Earley J. 1970. An Efficient Context-Free Parsing Algorithm. Comm. of the ACM. Volumen 13. No 2. Páginas: 94-102.
- Goodman J. T. 2001. A bit of progress in language modelling. Extended version. Microsoft technical report. MSR-TR-2001-72.
- Magnuson T., Hunnicutt S. 2002. Measuring effectiveness of word-prediction: the advantage of long term use. Speech, Music and Hearing, KTH, Stockholm, Sweden. TMH-QPSR Volume 43: 57-67.
- Palazuelos S. E., Aguilera S., Rodrigo J. L., Godino J., Martín J. 1999. Considerations on the Automatic Evaluation of word prediction systems. Augmentative and Alternative Communication: New Directions in Research and Practice. Pags: 92-104. Whurr Publishers. Londres.
- Palazuelos S.E., Rodrigo J. L., Godino J. I., Aliaga F., Martín J. L., Aguilera S. 1999. Predicción de palabras en castellano. SEPLN Procesamiento del lenguaje natural Volumen: 25. Páginas: 151-158.
- Palazuelos Cagigas, S. E. 2001. Aportación a la predicción de palabras en castellano y su integración en sistemas de ayuda a personas con discapacidad física. Tesis Doctoral. UPM. <http://www.depeda.uah.es/personal/sira/TesisSiraEspanol.pdf>