

Desarrollo de un corrector ortográfico para aplicaciones de conversión texto-voz

Ana Armenta, José Gregorio Escalada, Juan María Garrido, Miguel Ángel Rodríguez

División de Tecnología del Habla

Telefónica I+D

Emilio Vargas 6, MADRID 28043

aalv@tid.es, jges@tid.es, jmgarri@tid.es, miguel@tid.es

Resumen: En Telefónica I+D hemos desarrollado un corrector ortográfico automático especialmente orientado a la adaptación de textos de Internet para su lectura por nuestro conversor texto-voz. El modelo de corrección detecta las palabras erróneas de acuerdo con el diccionario de formas correctas, y genera alternativas de corrección a través de transformaciones elementales de caracteres (inserción, borrado, sustitución o trasposición). La elección de la mejor de las alternativas se apoya en la probabilidad de ocurrencia (según su frecuencia en el corpus de entrenamiento), la probabilidad de la transformación elemental (las matrices de confusión), el contexto (según el modelo de lenguaje) y el grado de corrección del texto (factor de confianza).

Palabras clave: modelo de corrección, matrices de confusión, diccionario, modelo de lenguaje, abreviaturas con expansiones múltiples, factor de confianza

Abstract: Telefónica I+D has developed an automatic correction method specially focused on the adaptation of Internet text in order to be read by our text-to-speech system. The correction model detects errors by searching words in a dictionary, and then (if the word isn't found) correction candidates are generated by making elemental transformations (such as insertion, deletion, substitution or transposition). The selection of the best candidate takes into account the word probability (its frequency in the training corpus), the probability of the elemental transformation (confusion matrices), the context (language model) and the text correction level (reliability factor)

Keywords: correction model, confusion matrices, dictionary, language model, reliability factor

1 Introducción

Las aplicaciones en las que se utiliza un conversor texto-voz (CTV) para suministrar información al usuario o cliente final, a través del teléfono, son muy dependientes del grado de corrección en la escritura de los textos que reciben como entrada. Si el texto está mal escrito, mal puntuado, mal acentuado... la lectura que hará el CTV posiblemente será incorrecta.

El usuario puede sentirse defraudado por esa "mala lectura" que percibe; sobre todo porque él no ve el texto y no entiende la razón de esa lectura deficiente. Incluso puede que el usuario ni siquiera pueda entender de manera adecuada la

información recibida. Es decir, un servicio que funciona bien puede no resultar satisfactorio para el usuario por culpa de unos textos mal escritos.

En la división de Tecnología del Habla de Telefónica I+D hemos desarrollado una serie de técnicas de corrección que constituyen un ejemplo de aplicación de estrategias de procesamiento del lenguaje natural a las necesidades de productos y servicios industriales.

Aunque existen en el mercado multitud de correctores ortográficos de altas prestaciones, el tipo de aplicaciones en las que nuestra empresa demanda la

integración de un módulo corrector exige que el funcionamiento del mismo sea completamente automático, sin ninguna interacción con el usuario.

Además, el módulo corrector debe ser fácilmente adaptable a otros idiomas, o incluso a las características de los textos de unas y otras aplicaciones en las que se pueda integrar. Por este motivo se han desarrollado módulos que se encargan del preproceso y adaptación de los mensajes cortos y los mensajes de correo electrónico, para de esta forma conseguir una lectura más clara e inteligible.

Las tareas del corrector las hemos dividido en dos etapas. Una primera etapa de filtrado en la que se tratan aspectos determinísticos del texto de entrada (“emoticonos”, conversión de tipos y normalización de la escritura), con técnicas tradicionales de preproceso de texto. Y una segunda etapa de corrección ortográfica, basada en criterios estadísticos, con un modelado del lenguaje de los mensajes cortos y de los correos electrónicos y sus errores. En esta comunicación nos centraremos en la descripción del desarrollo de esta segunda etapa, esto es, en el modelo de corrección que hemos diseñado y desarrollado.

2 Modelo de corrección

2.1 El modelo inicial

Se diseñó un primer modelo de corrección ortográfica basado en la detección y tratamiento de palabras aisladas. El funcionamiento del sistema es, a grandes rasgos, el siguiente. Cada palabra se busca en un diccionario que contiene todas las palabras que no precisan corrección ortográfica. Esto incluye tanto las consideradas válidas en el idioma (válidas como entradas de un diccionario) como las formas conjugadas, los diminutivos, los neologismos, palabras extranjeras frecuentes, etc.

Para completar nuestro diccionario con todas estas formas que no suelen aparecer en los diccionarios convencionales, se revisaron semiautomáticamente las

palabras de varios corpus, incorporando al diccionario toda forma que se considerara correcta. Para ello utilizamos un corpus de texto periodístico en formato electrónico, un corpus de mensajes de correo electrónico obtenido de diversas listas de distribución y un corpus de mensajes cortos recogido entre el personal de la empresa.

Cuando una palabra no aparece en el diccionario (es lo que llamaremos en adelante un ‘error no ambiguo’ porque no hay ambigüedad en su incorrección), el modelo de corrección ortográfica establece que se realicen transformaciones simples en la palabra (inserción, borrado, sustitución o trasposición de caracteres), y cada una de las formas generadas se comprueba frente al diccionario. De todas las formas generadas que aparecen en el diccionario, el sistema elige como correcta aquella que tenga mayor probabilidad, obtenida como el resultado de multiplicar la “probabilidad a priori” [1] de la forma generada (probabilidad de la palabra en el idioma), por la “probabilidad de canal” de la transformación (probabilidad de que un carácter se sustituya por otro, desaparezca, etc, en función de la identidad del carácter y de su contexto, reflejada en las matrices de confusión).

De acuerdo con el teorema de Bayes:

$$P(\text{corr}_i / \text{error}) = \frac{P(\text{corr}_i) * P(\text{err} / \text{corr}_i)}{P(\text{err})} \quad (1)$$

donde $P(\text{corr}_i / \text{error})$ es la probabilidad de que la forma correcta sea corr_i , condicionada a que aparezca la forma err ; $P(\text{corr}_i)$ es la probabilidad a priori de la forma corr_i ; $P(\text{err} / \text{corr}_i)$ es la probabilidad de que se haya producido la forma err , condicionada a la forma correcta corr_i ; y $P(\text{err})$ es la probabilidad de que aparezca err , que será igual para todas las alternativas de corrección. Por tanto, para maximizar la probabilidad de encontrar la forma correcta, suponiendo

que se ha detectado el error de manera fiable, basta con encontrar el máximo de la ecuación (2):

$$P(\text{corr}_i / \text{error}) = P(\text{corr}_i) * P(\text{err} / \text{corr}_i) \quad (2)$$

En una primera versión, la “probabilidad a priori” se estimaba directamente como la frecuencia de aparición de la palabra en un corpus, y se guardaba asociada a la entrada de la palabra en el diccionario. Posteriormente se incorporó un modelo de lenguaje que permite incluir información contextual en la estimación de la probabilidad “a priori”, y así poder afrontar la corrección de lo que en adelante llamaremos errores ambiguos (palabras que vistas de forma aislada son correctas, pero que no lo son si se tiene en cuenta el contexto). Las “probabilidades de canal” se calculan a partir de las matrices de confusión, que, a su vez, se estiman como el número de veces que se ha producido cada uno de los errores básicos (inserción, borrado, sustitución o trasposición) en el corpus de entrenamiento [2].

El entrenamiento se realiza con técnicas de *boot-strapping*, utilizando el propio corrector ortográfico para detectar formas erróneas y su correspondiente forma correcta. Se obtiene el alineamiento óptimo [3] (con menor número de errores elementales) entre las dos formas, y se incrementan las cuentas correspondientes de las matrices de confusión. Con estas nuevas matrices se vuelve a procesar el corpus para identificar de nuevo los errores. Inicialmente, estas matrices se entrenaron con pares “forma errónea-forma correcta”, repasadas manualmente, para acelerar la convergencia del método.

En el Anexo 1 mostramos nuestras matrices de confusión con la intención de aclarar tanto el método de construcción de las matrices como su significado. En la matriz de inserción $Ins[x, y]$ se computa el número de veces que ‘x’ fue escrito como ‘xy’; en la matriz de borrado $Borr[x, y]$, las

veces que los caracteres ‘xy’ fueron escritos como ‘x’; en la matriz de sustituciones $Sust[x, y]$, las veces que ‘x’ fue escrito como ‘y’; y finalmente, en la matriz de trasposiciones $Tras[x, y]$, el número de veces que ‘xy’ fue escrito como ‘yx’. Las matrices de confusión constituyen un modelado del tipo de errores ortográficos que se dan en el corpus de entrenamiento, y son la base para el cálculo de las probabilidades de canal.

De esta forma, la probabilidad de canal $P(\text{err} / \text{corr}_i)$, se calcula como:

- $Ins[x, y] / \text{cuenta}[x]$ en caso de inserción
- $Borr[x, y] / \text{cuenta}[y]$ en caso de borrado
- $Sust[x, y] / \text{cuenta}[x]$ en caso de sustitución
- $Tras[x, y] / \text{cuenta}[x, y]$ en caso de trasposición

donde $\text{cuenta}[x]$, $\text{cuenta}[y]$ y $\text{cuenta}[x, y]$ representan el número de apariciones de x , y , y xy respectivamente en el corpus de entrenamiento.

Esta aproximación permite adaptar el comportamiento del módulo corrector a distintos tipos de contenidos, e incluso a distintos idiomas. Entrenando el diccionario de formas y las matrices de confusión con textos procedentes de un tipo de contenidos determinado, tendremos un modelo de corrección adaptado al entorno en cuestión.

2.2 Corrección de uno o dos errores simples

En un primer momento, el módulo corrector intentaba tratar, para cada palabra considerada incorrecta (no presente en el diccionario), un solo error simple de los recogidos en las matrices de confusión. Esta limitación se impuso para mantener unas buenas prestaciones de tiempo de respuesta. Por cada letra de la forma

“errónea”, se generan unas 50 alternativas, por lo que en promedio se generan por cada palabra unas 200 formas que hay que consultar en el diccionario. Si se intentan corregir dos errores, se pasa a unas 40.000 consultas por palabra, y para tres errores, unos 8.000.000.

Para poder obtener mejores tiempos de respuesta, y por tanto aumentar el número de errores tratados, se implementaron, en una segunda versión, técnicas de acceso al diccionario basadas en estructuras *trie* [4][5]. También se utilizaron técnicas de poda que reducen el número de consultas. De esta forma, se consiguieron unos tiempos de respuesta aceptables para la corrección de dos errores simples (similares a los de la versión anterior para un único error).

Otro inconveniente de incrementar el número de errores tratados es un incremento de la probabilidad “de falsa alarma” (modificar palabras correctas que no han sido incluidas en el diccionario). Si se intenta una sola transformación elemental, es posible que no se genere ninguna forma recogida en el diccionario, y por tanto se mantiene la forma original. Pero si se acumulan más transformaciones, es seguro que en algún momento se generarán formas recogidas en el diccionario, y por tanto se modificará la palabra.

	1 error simple	2 errores simples
Corrección de errores no ambiguos	86,80(%)	87,63(%)
Corrupción de palabras correctas	0,18(%)	0,31(%)
Expansión correcta de abreviaturas	76,71(%)	77,48(%)

Tabla 1. Comportamiento del corrector ante una y dos transformaciones elementales, en ambos casos se utiliza el modelo de lenguaje

En la Tabla 1 se presentan resultados del corrector cuando corrige con una o dos transformaciones elementales. Como se ve,

al tratar dos errores simples aumenta la tasa de corrección de errores no ambiguos, pero aumenta también la tasa de “corrupción” de palabras correctas.

2.3 Abreviaturas con expansiones múltiples

Una característica muy frecuente de las abreviaturas o códigos que se utilizan en los mensajes cortos de móviles (y cada vez más en los mensajes de correo electrónico) es que a menudo tienen varias lecturas posibles, de manera que queda a la inteligencia del lector el encontrar la expansión correcta.

Para el tratamiento automático de estos mensajes, tan pernicioso puede resultar el no expandir estas abreviaturas como expandirlas con una opción incorrecta.

Por esta razón, uno de los aspectos abordados en nuestro desarrollo ha sido el tratamiento de abreviaturas con varias posibilidades de expansión.

Para ello, manteniendo como referencia el modelo seguido en el corrector, se optó por una representación muy sencilla para cada abreviatura:

$abr \rightarrow exp1(freq1) \rightarrow exp2(freq2) \rightarrow exp3(freq3) \dots$

Con esta representación, una abreviatura *abr* puede ser sustituida por *exp1*, o *exp2*, o *exp3*, o... El número entre paréntesis representa la frecuencia con la que se da esa expansión. Si no se indica ningún número, esta frecuencia se estima a partir de la recogida en el diccionario para cada una de las formas *exp1*, *exp2*, etc.

Cuando se encuentra en el texto una forma que aparece en la tabla de abreviaturas, se proponen todas las posibles expansiones, y será la evaluación proporcionada por el modelo de lenguaje la que determine cuál es la forma correcta.

3 Modelo de lenguaje

3.1 Un primer modelo de lenguaje

Como ya hemos mencionado, en una segunda fase de desarrollo del corrector nos planteamos la inclusión en el modelo

de corrección de información contextual a través de un modelo de lenguaje. Puesto que se puede considerar que el alfabeto de símbolos coincide con el de las palabras recogidas en el diccionario de formas correctas (unas 100.000), no parecía adecuado utilizar un modelo de lenguaje basado directamente en la forma de las palabras. Se optó entonces por la utilización de categorías gramaticales, enriquecidas con información adicional que recogiera las relaciones entre las palabras (información de género, de número, verbo auxiliar, etc). Se diseñó así un conjunto de 75 etiquetas.

Se empleó el módulo de proceso lingüístico del CTV para procesar los distintos corpus manejados (texto periodístico, mensajes de correo electrónico y mensajes cortos corregidos) y generar los pares “palabra<>etiqueta”, utilizando el módulo de categorización gramatical del conversor. A partir de esta secuencia de pares se obtuvo tanto la cuenta de trigramas que constituye el modelo de lenguaje como las etiquetas correspondientes a cada forma del diccionario (una misma forma puede funcionar con distintas etiquetas en distintos contextos; por ejemplo, ‘ría’ puede ser verbo o nombre).

Sobre la cuenta de trigramas (aparecieron 59.245 secuencias, de las 421.875 posibles) no se utilizó ninguna técnica de *discounting*, pues se observó que al conservarse la decisión local, palabra a palabra, no se perdía ninguna precisión al rechazar las formas con una etiqueta que nunca se dio en el contexto correspondiente. La probabilidad asociada a cada trigrama fue incorporada entonces al modelo de corrección.

En la Tabla 2 se presentan: en A, los resultados del corrector antes de la incorporación del modelo de lenguaje de categorías (cuando solamente se tenían en cuenta la probabilidad a priori y las matrices de confusión); en B, los resultados del corrector cuando funciona

con la probabilidad correspondiente al modelo de lenguaje de categorías.

	A	B
Corrección de errores no ambiguos	83,30(%)	86,80(%)
Corrupción de palabras correctas	0,50(%)	0,18(%)
Expansión correcta de abreviaturas	74,34(%)	76,71(%)

Tabla 2. Comportamiento del corrector sin y con modelo de lenguaje de categorías (A y B, respectivamente)

3.2 Mejoras y ajustes del modelo de lenguaje

Cuando se ha pretendido extender el uso del modelo de lenguaje para decidir entre las distintas expansiones de una abreviatura (y para decidir frente a formas ambiguas, potencialmente correctas), se ha comprobado que muchas expansiones (las de las abreviaturas más frecuentes), se corresponden a la misma categoría gramatical, o se utilizan en contextos gramaticales muy similares. Por ejemplo, algunas de las expansiones más habituales de ‘t’ son ‘te’ y ‘tú’, y ambas aparecen frecuentemente ante verbo.

Sin embargo, ‘tú’ aparecerá casi exclusivamente frente a formas en segunda persona del singular, mientras que ‘te’ admite una distribución más amplia. Una posibilidad era incluir la persona de la conjugación de los verbos como un elemento más en la etiqueta de categoría, pero este camino conducía a una explosión del número de etiquetas, que difuminaba la capacidad aglutinante del modelo, y la posibilidad de su correcto entrenamiento.

Por esta razón se optó por complementar el modelo de etiquetas de categorías con un modelo de formas de palabras. Además de los trigramas de categorías ya descritos, el modelo incluye ahora trigramas de palabras. Para evitar la explosión del número de trigramas posibles

(para unas 100.000 palabras en el diccionario, se tendrían 1.000.000.000.000.000 de posibles trigramas), se limitó el modelo a los trigramas en los que los tres elementos tienen una frecuencia muy alta (entre las 5.000 palabras más frecuentes). Aunque esta aproximación es muy poco ortodoxa, nos facilita un mecanismo suficiente para resolver los contextos más frecuentes (que es donde más suelen utilizarse las abreviaturas).

Las probabilidades de los dos submodelos se combinan con un peso que se ha ajustado para que en los casos correctos ambos submodelos queden equilibrados.

3.3 Tratamiento de errores en formas ambiguas

Un número muy elevado de los errores presentes en los mensajes analizados transformaban la palabra correcta en otra forma que coincidía con otra palabra correcta. Se trata de formas ortográficamente correctas cuyo uso es incorrecto en ese contexto (aproximadamente en la mitad de los errores analizados se daba esta circunstancia). Este fenómeno es el que se ha denominado “errores en formas ambiguas”, y no eran tratables en la primera versión del corrector, pues sólo se intentaba corregir una palabra si no aparecía en el diccionario de formas correctas. Puesto que este tipo de errores generan palabras que sí van a aparecer en el diccionario de formas correctas, estas palabras se daban por buenas, y no se intentaban transformar.

Dada la importancia de este volumen de errores, y confiando en las posibilidades del modelo de lenguaje para detectar contextos anómalos (que marcarán posibles errores en formas ambiguas), se decidió afrontar su tratamiento.

Un elemento nuevo es el del grado de confianza del texto original y de las formas propuestas. Puesto que se trata de modificar una palabra que en principio

podría ser correcta por otra distinta, se incluyó un procedimiento en la interfaz, que permite especificar un nivel de confianza (de 0 a 100) del texto que se está tratando. Si la confianza en la corrección inicial del texto es alta, sólo se aceptarán transformaciones que supongan un claro incremento de la probabilidad resultante. Si la confianza es baja (el usuario supone que es muy probable que el texto contenga formas erróneas), se aceptarán las transformaciones propuestas con umbrales más bajos de probabilidad. Hay valores especiales del factor de confianza: el 110 se emplea para corregir únicamente mediante expansión de abreviaturas e inserciones simples; el 120 se emplea para solamente expandir abreviaturas.

4 Resultados

Hemos presentado un módulo corrector, desarrollado en Telefónica I+D, completamente automático, cuyos componentes (modelo de lenguaje, matrices de confusión, diccionario, tabla de abreviaturas) se pueden adaptar fácilmente al idioma y tipo de texto según sea la aplicación a la que se quiera incorporar.

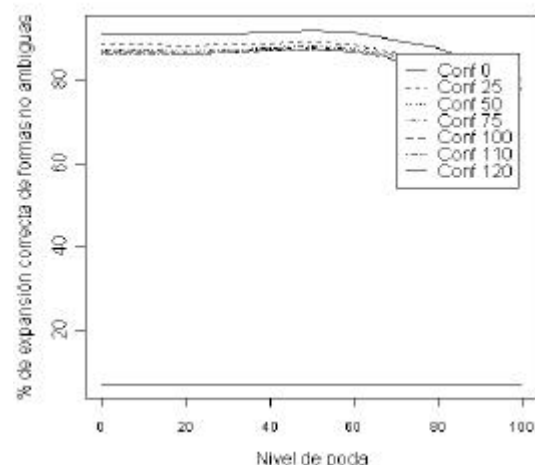


Figura 1. Variación del porcentaje de expansión correcta de formas no ambiguas en función del factor de confianza

En concreto, nosotros lo hemos aplicado a textos de Internet en español

(mensajes cortos y mensajes de correo electrónico), ya que estos van a ser leídos por el CTV, que exige como entrada texto correctamente escrito.

Las figuras 1, 2 y 3 describen el comportamiento de nuestro corrector de mensajes cortos según se trate de errores no ambiguos, errores ambiguos o abreviaturas, para distintos niveles de poda y del factor de confianza.

Como vemos, el comportamiento del corrector se mantiene constante hasta que se aplican podas del 50-60%, punto en el que comienza a decaer la tasa de correcciones acertadas. También queda patente el distinto comportamiento del corrector ante variaciones del factor de confianza según se trate de errores no ambiguos, errores ambiguos o abreviaturas. El corrector debe ser configurado con las opciones más apropiadas para cada servicio en el que se vaya a utilizar.

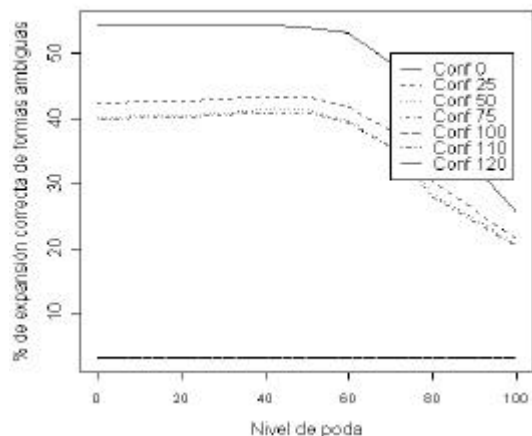


Figura 2. Variación del porcentaje de expansión correcta de formas ambiguas en función del factor de confianza

El módulo corrector que hemos descrito está siendo utilizado en diversos servicios

de lectura de mensajes cortos en castellano que utilizan el CTV.

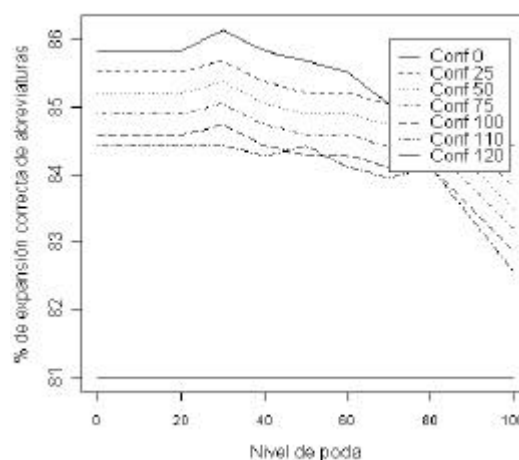


Figura 3. Variación del porcentaje de expansión correcta de abreviaturas en función del factor de confianza

Bibliografía

- [1] Shannon, C.E. 1948. A mathematical theory of communication, *Bell System Technical Journal*, vol.27:379-423 y 623-656. .
- [2] Church, K. y W. Gale. 1991. Probability scoring for spelling correction. *Statistics and Computing I*: 93-100.
- [3] Wagner, R.A. y M.J. Fisher. 1974. The String-to-String Correction Problem . *Journal of the asociation for Computing Machinery* , vol 21, No 1: 168-173.
- [4] Muth, F. y A.L. Tharp, . 1977. Correcting human error in alphanumeric terminal input. *Information Processing and Management*, 13(6): 329-337.
- [5] Knuth, D.E. 1968 . *The Art of Computer Programming*, Volume 3 / Sorting and Searching

