

## La inevitable estadística... en el mundo de la proteómica

Ana María Rodríguez Piñeiro<sup>1</sup>, María Páez de la Cadena, Francisco Javier Rodríguez Berrocal

Departamento de Bioquímica, Genética e Inmunología, Facultad de Biología, Campus As Lagoas-Marcosende, Universidad de Vigo, 36310 Vigo (Pontevedra)

### Resumen

En los estudios de proteómica, especialmente en los orientados a la búsqueda de biomarcadores, resulta imprescindible el uso de métodos estadísticos. La proteómica conlleva un tipo de experimentación que implica en la mayoría de los casos el manejo de un gran número de variables (los *spots* en los geles bidimensionales o los picos en los espectros de masas) pero, por el contrario, un escaso número de muestras. Los conjuntos de datos obtenidos en proteómica pueden tratarse mediante métodos univariantes, pero este tipo de técnicas incrementa la posibilidad de que aparezcan falsos positivos y no permite la detección de tendencias. Por el contrario, su análisis mediante métodos estadísticos multivariantes resulta mucho más eficiente. Estos métodos permiten la reducción de la dimensionalidad del conjunto de datos, o la utilización de estrategias de validación estadística como la validación cruzada o las pruebas de permutación siendo en general métodos más robustos y menos susceptibles a factores que induzcan confusión.

Por otro lado, un paso imprescindible en el establecimiento de marcadores es su validación biológica. Nuevamente, aquí, es necesario el uso de métodos estadísticos que permitan conocer el valor o la utilidad que tiene una determinada proteína, o varias, a la hora de diagnosticar la enfermedad, utilizarse como indicador pronóstico o para seguir la evolución de un paciente. Basándonos en nuestra propia experiencia, en este trabajo se revisan los métodos estadísticos, tanto univariantes como multivariantes, que se utilizan más frecuentemente en la búsqueda de biomarcadores mediante 2-DE y en su validación.

### Palabras clave:

Estadística, univariante, multivariante, electroforesis bidimensional, proteómica, biomarcadores.

### 1. Necesidad de la utilización de pruebas estadísticas en proteómica: consideraciones previas

La proteómica se considera, al cabo de más de 10 años de su definición por Wasinger y Wilkins (Wasinger *et al.*, 1995; Wilkins *et al.*, 1996), como el estudio de proteínas a gran escala, incluyendo su identificación, cuantificación, localización, estructura, modificaciones, interacciones, actividades y función. Aunque estos objetivos abarcan una gran variedad de técnicas, los métodos proteómicos por excelencia son la 2-DE<sup>2</sup> y la espectrometría de masas. Ambas técnicas son ampliamente utilizadas

debido a su gran resolución, ya que permiten la detección de cientos a miles de proteínas en un solo análisis. Sin embargo, esto se traduce en complicaciones a la hora de evaluar conjuntamente la gran cantidad de datos obtenidos y de ahí la importancia de la aplicación de las técnicas estadísticas correctas, que permitan analizar esa ingente cantidad de datos y obtener conclusiones significativas.

En particular, en los estudios de proteómica orientados a la búsqueda de biomarcadores, es decir, proteínas concretas diferenciadoras de un estado biológico (por ejemplo, una enfermedad), resulta imprescindible el uso de métodos estadísticos. Como ya han argumentado otros autores, aunque los conjuntos de datos obtenidos en proteómica pueden tratarse mediante métodos univariantes, este tipo de técnicas incrementa la posibilidad de que aparezcan

<sup>1</sup> Autor para la correspondencia

<sup>2</sup> 2-DE: electroforesis bidimensional.

falsos positivos y no permite la detección de tendencias (Dudoit *et al.*, 2003; Karp *et al.*, 2005; Wilkins *et al.*, 2006). Por el contrario, su análisis mediante métodos estadísticos multivariantes resulta mucho más eficiente (Karp *et al.*, 2005).

El desarrollo de este tipo de experimentación implica en la mayoría de los casos dos problemas estadísticos que se han dado en llamar maldiciones (de Noo *et al.*, 2006). La primera es la “maldición de la dimensionalidad” (*curse of dimensionality*), provocada porque el número de variables (los *spots* en los geles bidimensionales o los picos en los espectros de masas) está en el orden de millares a decenas de millares. La segunda es la “maldición de la dispersión del conjunto de datos” (*curse of dataset sparsity*), debida a un escaso número de muestras. Estos dos problemas afectan especialmente a los métodos univariantes, y se traducen en ocasiones en el descubrimiento de proteínas que discriminan las poblaciones a estudio, aun cuando estas no son en realidad diferentes. Este fenómeno se denomina “sobreajuste” (*overfitting*), y da lugar a que el modelo de discriminación encontrado no funcione adecuadamente cuando se analizan nuevas muestras. Para evitar este inconveniente existen varias soluciones, siendo una de ellas la reducción de la dimensionalidad del conjunto de datos, o la utilización de estrategias de validación estadística como la validación cruzada o las pruebas de permutación (Smit *et al.*, 2007). Aunque las aproximaciones multivariantes también se ven afectadas por estas dos maldiciones, estos métodos son en general más robustos y menos susceptibles a factores que induzcan a confusión.

Por otro lado, un paso imprescindible en el establecimiento de marcadores es su validación biológica. Estos estudios tienen, si cabe, una importancia mayor cuando se trata de encontrar buenos marcadores en relación con una determinada enfermedad. Nuevamente aquí, es imprescindible el uso de métodos estadísticos que permitan conocer el valor o la utilidad que tiene una determinada proteína, o varias, a la hora de diagnosticar la enfermedad, utilizarse como indicador pronóstico, o para seguir la evolución de un paciente. Estas capacidades se reflejan, en general, en la certeza con la que las proteínas seleccionadas clasifican correctamente las muestras en su verdadero grupo de origen.

Basándonos en nuestra propia experiencia (Ayude *et al.*, 2000; Rodríguez-Piñero *et al.*, 2005; Lemos-González *et al.*, 2007; Rodríguez-Piñero *et al.*, 2007) en este trabajo se revisan los métodos

estadísticos, tanto univariantes como multivariantes, que se utilizan más frecuentemente en la búsqueda de biomarcadores mediante 2-DE. En general, la estadística está presente en casi todas las fases del análisis de datos, incluso cuando el investigador no es consciente de ello (por ejemplo si está empleando programas especializados que “escondan” los procesos matemáticos del análisis). Sin embargo, en esta revisión nos centraremos en aquellos métodos que el investigador emplea de forma “consciente”. Por último, dedicaremos un apartado a describir las pruebas estadísticas que permiten determinar la sensibilidad y especificidad de un marcador o su capacidad de pronóstico.

## 2. Análisis estadístico de mapas bidimensionales de proteínas

Una de las ventajas de la utilización de la 2-DE es que el proteoma puede visualizarse mediante mapas en donde las proteínas aparecen como manchas o *spots*, separadas en función de su punto isoeléctrico y de su masa molecular. Un proceso fundamental antes de su obtención es un cuidadoso diseño experimental, lo que ayuda a reducir los errores sistemáticos, mejorando la precisión de los análisis estadísticos posteriores y disminuyendo el número de falsos positivos. A este respecto, Chich *et al.* (2007) han propuesto recientemente la utilización de un diseño por bloques de muestras que sean robustos y la aleatorización de las muestras en cada bloque.

Dada la complejidad de estos mapas, su análisis ha de realizarse a través de un programa específico para imágenes 2D<sup>3</sup>. En general, estos programas se basan en asignar, a cada uno de los puntos (píxeles) detectados en un gel, unas coordenadas  $X$  e  $Y$ , que corresponden a su posición horizontal y vertical en la imagen, y un valor  $Z$  que determina la densidad óptica. De esta forma, las coordenadas  $X$ ,  $Y$  y  $Z$  de cada mancha proteica vienen determinadas por las correspondientes a los puntos que la forman. Esto proporciona el área de cada *spot*, así como su valor de intensidad óptica, y por lo tanto una medida cuantitativa (volumen) de cada proteína.

Cuando las manchas proteicas son detectadas, la imagen original del gel es filtrada y el volumen de los *spots* se ajusta a un modelo gaussiano, sobre el que se puede realizar la comparación. En caso

<sup>3</sup> 2D: bidimensional/es.

necesario, existen varios métodos con los que el investigador puede estimar los datos perdidos, como reemplazar los valores por 0, por un valor de densidad mínimo detectable en el gel, por la media de las medidas del mismo *spot* en otras muestras, etc. (Chang *et al.*, 2004; Chich *et al.*, 2007). A continuación se suelen normalizar los valores de los *spots* en todos los geles analizados, en la mayoría de los casos obteniendo su volumen relativo, que es comparable entre geles y que permite la aplicación de los métodos estadísticos. En relación a este último paso, cabe destacar el estudio de Randic *et al.* (2006), quienes recientemente han demostrado que existe un número máximo de *spots*, que ellos proponen ser los 300 *spots* más intensos, a partir del cual no es posible obtener más información estadística, es decir, que los datos derivados de las restantes proteínas no aportan mayor información sobre el mapa proteómico que se quiere caracterizar. Esta afirmación no es aplicable en el caso de la búsqueda de variaciones en proteínas concretas, especialmente en la búsqueda de biomarcadores, que suelen ser proteínas de baja abundancia. Sin embargo, sí debe ser tenido en cuenta cuando se pretende utilizar el propio mapa proteico para caracterizar a un grupo de población.

Por otra parte, antes de realizar la comparación entre mapas 2D de diversas muestras poblacionales, el investigador debe comprobar la reproducibilidad de la técnica, proceso en el que también está presente la estadística ya que, por ejemplo, es necesario emplear métodos de regresión para estimar el coeficiente de variación entre réplicas. Desde el punto de vista estadístico cabe destacar que, a pesar de la conocida variabilidad de la 2-DE, las réplicas biológicas tienen más valor que las réplicas técnicas (Chich *et al.*, 2007), y por ello son más recomendables, especialmente en el caso de muestras complejas como son las humanas. Con el fin de analizar la utilidad de un diseño experimental, Hunt *et al.* (2005) proponen, antes de comenzar un estudio, llevar a cabo estudios piloto con 3 muestras por grupo y 3 geles por muestra, y una serie de herramientas estadísticas para evaluar la potencia del diseño y determinar el número mínimo de muestras y geles necesarios para el experimento definitivo.

Por último, hay que recordar que los procesos de escaneado de la imagen, el pre-procesado (normalización, filtrado, etc.) y la propia comparación (alineamiento de imágenes y *spots*) también conllevan una serie de errores que, en gran medida, dependen del tipo de densitómetro empleado y del programa

de análisis 2D, y que afectarán a los resultados de los análisis estadísticos (Meleth *et al.*, 2005). Los distintos programas disponibles para 2-DE han sido comparados por varios autores, y así por ejemplo Garrels *et al.* (1989) estimaron la reproducibilidad del sistema QUEST en un 97,6%, y Rosengren *et al.* (2003) determinaron que el PD-Quest derivado de este sistema es más fiable que otros programas en el alineamiento de *spots*, con una menor tasa de falsos positivos.

### 3. Tratamientos estadísticos para comparar la expresión proteica

#### 3.1. Métodos univariantes: comparación de la expresión de proteínas individuales en mapas 2D

De modo general, en la búsqueda proteómica de biomarcadores se consideran proteínas candidatas a aquellas cuya expresión (bien en términos de cantidad absoluta o de volumen relativo) varía en los mapas problema en relación a los controles. Las pruebas estadísticas univariantes que se emplean en estos análisis de mapas 2D se pueden clasificar como métodos de filtrado, que se aplican sobre los datos pre-procesados (tratados como se ha descrito en el apartado anterior), y que posteriormente se pueden combinar con métodos de reducción de la dimensionalidad o de clasificación, como se comentará más adelante. Las pruebas univariantes más comunes son las de significación, generalmente con el fin de comparar medias/medianas de las variables. La prueba *t* de Student compara diferencias entre las medias del grupo de casos y del grupo control, asumiendo la normalidad de los datos (prueba paramétrica). Aunque este test se ha aplicado en numerosos estudios proteómicos, los datos obtenidos a partir de mapas 2D no suelen cumplir los requisitos paramétricos, es decir, que las variables medidas sean continuas y presenten una distribución normal, con varianzas homogéneas en las dos muestras poblacionales, y que los datos sean independientes. Por ello es más correcto emplear el test *U* de Mann-Whitney, que no asume estos requisitos (prueba no paramétrica). Para ambos tipos de test, se considera que una variable difiere en los grupos significativamente cuando el valor de probabilidad *P* asociado al estadístico es menor que un valor  $\alpha$  (generalmente 0,05 para un 95% de confianza en el resultado, o 0,01 para un 99% de confianza). Estas dos pruebas son las más empleadas en 2-DE, ya que normalmente las muestras a comparar no están relacionadas. Sin embargo, en algunos casos se trabaja con muestras

relacionadas (como es el caso del análisis de tejidos sanos y enfermos de un mismo paciente) y para ellos existen pruebas específicas: la prueba *t* para muestras relacionadas (prueba *t* dependiente) en el caso paramétrico, y la prueba *T* de Wilcoxon en el caso no paramétrico. Por otra parte, existen otros tipos de análisis univariantes aplicables a los datos de 2-DE, como por ejemplo el ANOVA<sup>4</sup> cuando lo que se desea analizar es la varianza muestral.

Como se indicó anteriormente, aplicar un valor de  $\alpha$  para un conjunto de muchas variables, las cuales son analizadas de manera univariada, conduce a un incremento notable de la probabilidad de obtener falsos positivos. Por ejemplo, analizando 1000 *spots* con una prueba como la *t* de Student, y con un nivel de confianza del 95%, se pueden seleccionar hasta 50 marcadores por azar. Por ello es necesario aplicar *a posteriori* una corrección (Horgan, 2007). Uno de los métodos de corrección más empleados es el de Bonferroni, que fija un valor de  $\alpha$  para el conjunto de datos, y compara el estadístico obtenido para cada variable con el valor [ $\alpha$ /número de variables], disminuyendo así la tasa de falsos positivos (Smit *et al.*, 2007). Sin embargo, es necesario tener en cuenta que la corrección de Bonferroni, considerada como una “prueba familiar” (*familywise*) ya que fija un valor de  $\alpha$  idéntico para todas las variables, puede ser contraproducente en el caso de grandes conjuntos de datos, ya que puede llegar a anular la significación estadística de variables que realmente son significativas (Perco *et al.*, 2006). Para estos casos existen otras pruebas que consideran distintos valores de  $\alpha$  para cada variable, y que pueden ser más apropiadas para el cálculo de la tasa de falsos positivos (Smit *et al.*, 2007).

Una consideración necesaria en este tipo de comparaciones es que, desde un punto de vista estadístico, el análisis de un número limitado de muestras de pacientes conlleva el riesgo de no detectar algunas proteínas alteradas (lo que se denomina “error de tipo II”), pero eso no significa que, si se incluyeran más muestras en el estudio, proteínas que se hubieran detectado mostrando una variación significativa en su expresión, no fueran nuevamente detectadas (lo que sería un “error de tipo I”) (Byrjalsen *et al.*, 1999).

En nuestro grupo, hemos aplicado diversas técnicas estadísticas univariantes para el análisis de

mapas 2D. Por ejemplo, para la comparación de proteínas séricas de individuos sanos y de pacientes con CCR<sup>5</sup> hemos empleado la prueba de Mann-Whitney (Rodríguez-Piñeiro *et al.*, 2004), mientras que para la comparación de tejido sano y tumoral de pacientes se empleó la prueba de Wilcoxon (Álvarez-Chaver *et al.*, 2007).

### 3.2. Métodos multivariantes: comparación de los niveles de expresión de conjuntos de proteínas en mapas 2D

Una estrategia diferente consiste en buscar no una, sino un patrón o conjunto de proteínas que varíen o que contribuyan a la diferenciación entre el grupo de casos y el grupo de controles. Fundamentalmente, en el análisis de mapas 2D se utilizan dos tipos de métodos multivariantes: los que permiten la reducción de la dimensionalidad del conjunto de datos (apartado 3.2.1), y aquellos que permiten la clasificación de las muestras en diferentes grupos (apartado 3.2.2). Por otra parte, los métodos multivariantes pueden ser empleados con dos finalidades: en primer lugar, permiten seleccionar un conjunto de proteínas cuya variación combinada discrimina los grupos analizados (apartado 4), proteínas que más tarde se valoran en la práctica clínica mediante otras técnicas distintas de la 2-DE; en segundo lugar, las pruebas multivariantes pueden en sí mismas emplearse para extraer toda la información que contiene el subproteoma visualizado en el mapa 2D, de forma que el propio mapa es la herramienta discriminatoria de los grupos (apartado 5).

#### 3.2.1. Métodos de reducción de datos

##### 3.2.1.1. Análisis de componentes principales (ACP)

El ACP<sup>6</sup> es una clase de análisis factorial que intenta identificar una serie de factores o variables subyacentes a los datos obtenidos (variables predictoras), de forma que expliquen la mayor parte de la varianza observada para el mayor número posible de variables medidas. Cada uno de estos nuevos factores, denominados componentes principales (CP<sup>7</sup>), es una combinación lineal de las variables originales, cumpliendo además la condición de ser ortonor-

<sup>4</sup> ANOVA: análisis de varianza.

<sup>5</sup> CCR: cáncer colorrectal.

<sup>6</sup> ACP: análisis de componentes principales.

<sup>7</sup> CP: componentes principales.

males entre sí, de modo que no existe correlación entre las CP, es decir, recogen distinta “parte” de la información contenida en las variables originales. El ACP genera un número de CP igual al número de variables originales, y cada componente expresa en orden decreciente un porcentaje de la varianza total de los datos, de modo que a menudo bastan las tres primeras CP para explicar una proporción suficiente de la varianza, reduciendo así de forma efectiva el número de variables que es necesario manejar. Por lo tanto, el ACP permite reducir un número elevado de variables (por ejemplo, el volumen relativo de todos los *spots* comunes a controles y casos) a unos pocos factores que permiten establecer diferencias entre los dos grupos de población. Esta prueba requiere que los datos sean variables cuantitativas (no categóricas), con una distribución normal bivariada para cada pareja de variables, y que las observaciones sean independientes. El modelo de análisis factorial especifica que las variables vienen determinadas por los factores comunes (los factores estimados por el modelo) y por factores únicos (los cuales no se superponen entre las distintas variables observadas); las estimaciones calculadas se basan en el supuesto de que ningún factor único está correlacionado con los demás, ni con los factores comunes.

El ACP es uno de los análisis multivariantes más populares en 2-DE, y ha sido aplicado en numerosos estudios, tanto en nuestro grupo (Rodríguez-Piñeiro *et al.*, 2007) como en otros: Kovarova *et al.* (2000) lo aplicaron para caracterizar cambios en líneas celulares de leucemia en respuesta a tratamiento; Roblick *et al.* (2004) estudiaron las diferencias relativas a la progresión del CCR analizando el proteoma de tejido normal, pólipos, adenomas, tejidos tumorales y metastáticos; recientemente, Verhoecx *et al.* (2005) combinaron DIGE con ACP para corroborar las diferencias encontradas en mapas analizados mediante técnicas univariantes. Por otra parte, Marengo *et al.* (2003) describieron un nuevo método denominado ACP de tres vías, que incluye información tanto sobre la intensidad como sobre la posición de los *spots*.

### 3.2.1.2. Mínimos cuadrados parciales (MCP)

La técnica de MCP<sup>8</sup>, más conocida como PLS por sus siglas en inglés (*partial least squares*) es similar al ACP, distinguiéndose de éste en que tam-

bién tiene en cuenta la covarianza de los datos a la hora de buscar las relaciones con los grupos. De esta forma, se pretende encontrar una variable predictora Y en términos de un conjunto de variables X, es decir, las variables no se consideran de forma independiente, sino que se dividen en variables dependientes e independientes, y por ello es que se incluye el estudio de covarianzas.

### 3.2.2. Métodos de clasificación

#### 3.2.2.1. Análisis discriminante (AD)

El AD<sup>9</sup> permite clasificar las muestras en una serie de grupos previamente establecidos. Así se obtiene un modelo predictivo para pronosticar el grupo de pertenencia de una muestra particular, a partir de las características observadas. Existen diversos tipos de AD: lineal, cuadrático, regularizado y diagonal. En el modelo lineal para dos grupos, usualmente empleado en proteómica, se genera una función discriminante basada en la combinación lineal de las variables predictoras que proporcionan la mejor discriminación entre los dos grupos. La función se genera a partir de una muestra de casos independientes y mutuamente exclusivos, cuyo grupo de pertenencia (único) es conocido, y posteriormente se puede aplicar a nuevos casos no clasificados (siempre y cuando se disponga de las medidas de las variables predictoras). Concretamente, cuando se trata de mapas 2D se suele aplicar el AD lineal a los valores de cantidad o de volumen relativo de los *spots* comunes a todos los mapas de los grupos analizados.

Una variante del AD que se está haciendo cada vez más popular en proteómica se denomina **AD de componentes principales** (ADCP, más conocido como PC-DA por sus siglas en inglés), y consiste en la reducción de datos mediante ACP, y la realización de un AD posterior sobre las puntuaciones así obtenidas. Esta combinación de métodos ha sido empleada en proteómica bajo diversos nombres (revisado en Smit *et al.*, 2007).

Otro método de reducción de datos combinable con el AD es el método de MCP, que en ese caso se denomina **AD-MCP** (PLS-DA en inglés). Recientemente, Karp *et al.* (2005) demostraron su utilidad y complementariedad con métodos univariantes aplicados a la proteómica de expresión en bacterias, aunque también existen numerosos ejemplos de su

<sup>8</sup> MCP: mínimos cuadrados parciales.

<sup>9</sup> AD: análisis discriminante.

utilización en estudios de metabolómica (revisado en Smit *et al.* 2007). En relación a la búsqueda de biomarcadores, Verhoeckx *et al.* (2005) emplearon este método en combinación con técnicas univariantes para estudiar efectos antiinflamatorios.

### 3.2.2.2. Otros métodos de clasificación

El método denominado **máquinas de vectores de soporte (SVM en inglés)** realiza la separación de los grupos construyendo un hiperplano que separa los casos de los controles, asignando el grupo de pertenencia de cada muestra según en qué lado del hiperplano se encuentre. También es combinable con el ACP y el MCP. Esta técnica es bastante popular en proteómica clínica, y especialmente se ha empleado para discriminar varios tipos de cáncer (ver ejemplos en Smit *et al.*, 2007).

La **regresión logística** es un método de clasificación basado en aplicar una regresión lineal para ajustar los datos al logaritmo natural de la probabilidad de que una muestra de una clase esté fuera de esa clase. Intuitivamente, este análisis realiza un proceso similar al AD, y por ello también puede combinarse con métodos previos de reducción de datos. La regresión logística ha sido recientemente empleada en proteómica del cáncer (Smit *et al.*, 2007).

Las **redes neuronales artificiales** son métodos no lineales, en los que se construye un modelo con una capa neuronal interna formada por los datos, una o más capas neuronales intermedias y ocultas, que consisten en las transformaciones de los datos de origen, y una capa neuronal externa que constituye la respuesta predicha. Este método también está adquiriendo popularidad en proteómica clínica (ver Smit *et al.*, 2007).

El método de los **centroides reducidos más cercanos** (*nearest shrunken centroids*) asigna las muestras a aquella clase cuya mediana de clase es más cercana al valor de la muestra. Aunque su fundamento parece *a priori* demasiado simple para producir resultados significativos, esta prueba ha sido empleada con éxito por Kemperman *et al.* (2007) para distinguir pacientes con y sin proteinuria.

Los **árboles de clasificación** se basan en el uso de un algoritmo que sucesivamente va dividiendo los datos de cada "nódulo paterno" en subconjuntos que constituyen los "nódulos hijos", maximizando la homogeneidad en los hijos mediante la elección de la/las proteína/s que más disminuyen la heteroge-

neidad. En proteómica clínica, se pueden encontrar ejemplos como la detección de pacientes con cáncer de páncreas (Bhattacharyya *et al.*, 2004), o del efecto del tratamiento de pacientes con leucemia (Albitar *et al.*, 2006).

Los **ensambles o conjuntos de clasificadores** (*ensemble classifiers*) combinan varios tipos de reglas de clasificación (clasificadores básicos) para construir un clasificador con mayor resolución. Cada muestra es clasificada primero por las reglas básicas, y la predicción se realiza en forma de "voto por mayoría". Un ejemplo es el bosque de clasificación, donde se construyen múltiples árboles de clasificación cuyos resultados se ensamblan posteriormente (Tong *et al.*, 2003).

## 4. Tratamiento estadístico para la selección de proteínas candidatas como biomarcadores

Como se ha comentado anteriormente, una utilidad importante de los métodos multivariantes en el análisis de mapas 2D es poder seleccionar proteínas con potencial valor como biomarcadores, antes de realizar los estudios de validación biológica. Por ejemplo, es muy frecuente que al comparar el proteoma de individuos afectados (por enfermedad, tratamiento, etc.) con el de individuos control se detecte un número elevado de proteínas cuya expresión individual está alterada. Estas alteraciones, aun teniendo relación con la enfermedad o situación, pueden no tener mucho valor como marcadores para distinguir los dos grupos considerados, o incluso pueden deberse al azar (error tipo I). En estos casos, la aplicación de métodos multivariantes sobre el conjunto de *spots* cuya alteración se ha detectado mediante métodos univariantes permite agrupar aquellas proteínas con mayor influencia en la distinción de grupos, y seleccionar las que teóricamente serían más útiles en la práctica clínica.

Nuestro grupo se ha beneficiado de esta aplicación de las técnicas multivariantes para la selección de potenciales marcadores (Rodríguez-Piñeiro *et al.*, 2007). En una comparación de mapas 2D de glicoproteínas séricas de individuos sanos y pacientes con CCR, se aplicó primero la prueba *U* de Mann-Whitney para seleccionar aquellas proteínas alteradas por la patología, encontrando 28 proteínas significativamente diferentes entre los grupos comparados. Sobre ellas se aplicaron métodos de reducción de datos (ACP) y clasificación (AD), re-

duciendo el grupo de proteínas a los mejores candidatos a biomarcadores en base a su capacidad de discriminación de los grupos. Así, finalmente se pudo concluir que 3 de ellas (clusterina, factor I del complemento y  $\beta$ -2-glicoproteína I) constituían un potencial panel de biomarcadores para el CCR.

## 5. Utilización de los mapas proteicos como herramientas diagnósticas

Las pruebas multivariantes también sirven para determinar si los mapas 2D tienen en sí mismos utilidad diagnóstica, es decir, si el proteoma de un individuo enfermo contiene la información necesaria para su diagnóstico, del mismo modo que los patrones de expresión génica observados mediante microarrays. En este caso, las pruebas multivariantes se aplican a todo el conjunto de proteínas de un mapa 2D para determinar si la información global que contienen permite discriminar los grupos de población estudiados (si el fin es diagnóstico, individuos sanos y pacientes).

Las diferencias entre los mapas de casos y controles se pueden determinar basándose en un rasgo cuantitativo como es la cantidad de proteína (volumen relativo), pero también puede hacerse utilizando un criterio cualitativo como es la posición de las proteínas en los mapas. Dentro de los estudios realizados por nuestro grupo, hemos comprobado la utilidad diagnóstica de los mapas proteicos utilizando ambos criterios, cuantitativo (Rodríguez-Piñeiro *et al.*, 2007) y cualitativo (Rodríguez-Piñeiro *et al.*, 2005). Resumimos, a continuación, los resultados de estos trabajos.

### 5.1. Comparando cantidad: ACP y AD

Cuando las pruebas multivariantes se aplican a los valores de todas las proteínas comunes a los mapas de casos y controles, sin tener en cuenta los resultados en los análisis univariantes, se mejora la detección de estados patológicos (en general, de los casos) ya que puede ocurrir que proteínas que no presentan por sí mismas una alteración significativa contribuyan al estado patológico en combinación con otras proteínas (LaBaer, 2005).

En nuestro grupo, se han empleado el ACP y el AD para determinar si el glicoproteoma sérico completo permitía distinguir a personas sanas de pacientes con CCR (Rodríguez-Piñeiro *et al.*, 2007). En primer lugar, se aplicó ACP sobre todos los datos de

volumen relativo (cantidad) de 363 *spots* comunes a controles y pacientes. Esta prueba permitió reducir el 100% de la varianza entre los grupos, dada inicialmente por las 363 variables, a 9 CP. Sobre ellos se aplicó una prueba de homogeneidad de varianza (mediante el estadístico de Levene), comprobándose que no existían diferencias significativas entre las varianzas del grupo de controles y del grupo de pacientes. Por ello, se aplicó una prueba de análisis de varianza de tipo ANOVA, encontrando que el segundo CP daba cuenta de las diferencias entre los grupos. Esto se pudo corroborar gráficamente.

Por otra parte, la aplicación del AD a los 363 *spots* comunes permitió encontrar una función discriminante (de clasificación o diagnóstica en este caso) para los grupos control y paciente, que explicaba el 100% de la varianza muestral, obteniéndose una puntuación para cada individuo que permite clasificarlo en uno u otro grupo con un nivel de confianza del 99,2%.

Estos resultados demuestran que, al igual que se había observado mediante ACP, es posible determinar la condición de un individuo en función de su patrón 2D de expresión de glicoproteínas séricas. Sin embargo, cabe reconocer que los valores de confianza obtenidos son mucho mejores cuando el análisis se restringe a aquellas proteínas detectadas previamente como diferenciales mediante técnicas univariantes.

### 5.2. Comparando posición: análisis de deformaciones relativas

Aunque clásicamente los análisis de mapas 2D se centran en la comparación de intensidad, volumen o cantidad de proteína, o en la presencia o ausencia de determinados *spots*, también es posible incluir variables de posición (datos de masa y punto isoeléctrico), como en el ACP de tres vías propuesto por Marengo *et al.* (2003) y comentado anteriormente. Dentro del estudio de la utilidad de los mapas como herramienta diagnóstica, nuestro grupo ha utilizado otro tipo de análisis multivariante que no tiene en cuenta la cantidad de proteína, sino sólo sus coordenadas de posición en los ejes *X* e *Y* de los mapas 2D. Esta técnica podría considerarse similar al ACP, y se denomina análisis de deformaciones relativas (RWA<sup>10</sup> por sus siglas en inglés) (Rodríguez-Piñeiro *et al.*, 2005). Este método es de

<sup>10</sup> RWA: análisis de deformaciones relativas.

uso generalizado y ha sido ampliamente comprobado en morfometría aplicada a sistemas biológicos (Bookstein, 1991; Rohlf, 1993).

En el estudio que realizamos comparando mapas de suero de donantes y pacientes con CCR, la aplicación de la morfometría geométrica nos permitió la detección de una variación local significativa en punto isoelectrico y masa molecular, que claramente separaba el grupo control del grupo de pacientes con CCR (Rodríguez-Piñeiro *et al.*, 2005). Esta utilidad podría ser una futura aplicación diagnóstica del método, de forma análoga al ACP y AD, ya que se podrían almacenar los patrones de posición de *spots* de un grupo numeroso de muestras de individuos sanos y de casos, y posteriormente comprobar estadísticamente si una muestra desconocida, con un patrón proteico concreto, corresponde al fenotipo de una persona sana o afectada por una patología o condición fisiológica.

Aún más, el método permite la detección gráfica de aquellas manchas proteicas que contribuyen en mayor medida a los cambios detectados entre los grupos, teóricamente relacionados con la condición o enfermedad estudiada. Existen varios fenómenos que explican las variaciones en posición: modificaciones post-traduccionales tales como fosforilación, glicosilación, acetilación, y otras; cambios en la secuencia aminoacídica de una proteína (polimorfismos de aminoácidos, por ejemplo); la existencia de formas proteicas truncadas; etc. Tanto la magnitud como la dirección de estos cambios se pueden mostrar en forma de vectores en los mapas 2D, y así en nuestro estudio se pudieron identificar aquellas proteínas que mostraban grandes desviaciones (alteraciones), junto con proteínas que mostraban cambios menores o incluso proteínas que no variaban en relación a la condición estudiada.

En comparación con la estrategia de Marengo *et al.* (2003), quienes aplicaron el ACP de tres vías para diferenciar grupos en dos conjuntos de muestras de controles y casos, el RWA separa la variación local de los efectos globales (Bookstein, 1991; Rohlf, 1993; Zelditch *et al.*, 2004). Aplicado a mapas 2D, corrige así el sesgo debido a variaciones que afecten a todas las manchas proteicas simultáneamente en cualquiera de las dos dimensiones de separación. Esto tiene una especial importancia en el caso de la 2-DE, ya que los errores técnicos y otras variaciones aleatorias tendrían probablemente un efecto en la localización de todos los *spots* resueltos (al menos en una de las dos dimensiones).

## 6. Métodos estadísticos empleados para la validación biológica de los marcadores.

### 6.1. Tratamiento estadístico para la validación de los biomarcadores

Otro campo dentro de la búsqueda de biomarcadores que requiere la utilización de métodos estadísticos es la validación de los marcadores. Un marcador para una enfermedad puede ser útil en el diagnóstico, para determinar el pronóstico del enfermo, o en el seguimiento de la evolución del paciente. Para evaluar la utilidad del marcador en cada uno de estos aspectos se requiere cuantificar la proteína (comúnmente con métodos tipo ELISA<sup>11</sup>) en un número amplio de pacientes y controles, o de pacientes en diferentes situaciones (por ejemplo, con tratamiento y sin tratamiento) y analizar los datos con diferentes pruebas estadísticas. Sin embargo, antes de realizar la comparación debe observarse la distribución de los datos, por ejemplo aplicando la prueba de Kolmogorov-Smirnov para cada muestra poblacional, lo que permite determinar si la población se ajusta o no a una distribución normal. A continuación se aplica la prueba de Levene para determinar si existe homogeneidad de varianza. Como se mencionó anteriormente, estas dos condiciones permiten aplicar pruebas paramétricas, como la *t* de Student de contraste de hipótesis para dos muestras independientes.

En general, puede decirse que existen dos aproximaciones estadísticas para evaluar biomarcadores. La primera modeliza el riesgo de enfermedad (o el resultado de la enfermedad) con técnicas como la regresión logística; así, un marcador es considerado útil si tiene un efecto importante en el riesgo. La segunda evalúa la efectividad de la clasificación dada por el marcador, utilizando parámetros tales como la sensibilidad, especificidad, valores predictivos, y las curvas ROC<sup>12</sup> (Pepe *et al.*, 2007).

### 6.2. Estimación de los parámetros diagnósticos

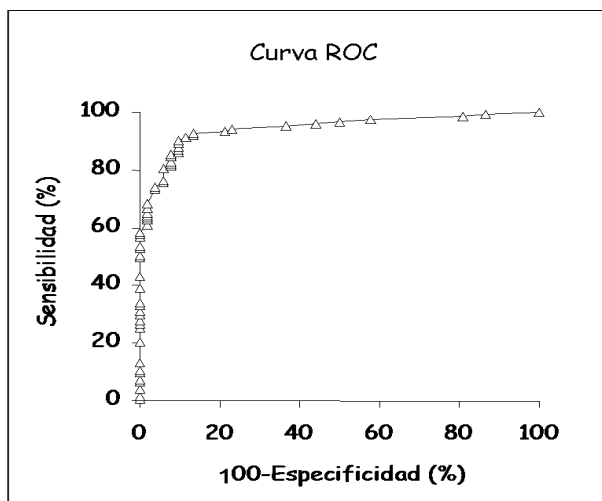
La validez de la determinación de los niveles de una proteína sérica como prueba diagnóstica para una patología se estudia mediante la elaboración de una curva ROC, siguiendo el procedimiento descrito por Beck y Shultz (1986). En dicha curva se representa, para cada uno de los posibles puntos de corte tomados en función de los niveles del marcador en

<sup>11</sup> ELISA: ensayo inmunoabsorbente ligado a enzima.

<sup>12</sup> ROC: características operativas relativas o para el receptor.



la muestra poblacional de pacientes, el porcentaje de verdaderos positivos (sensibilidad) en el eje de ordenadas, frente al porcentaje de falsos positivos (100-especificidad) en el eje de abscisas (Figura 1). El porcentaje de verdaderos positivos se calcula en función del número de individuos correctamente clasificados como enfermos por la prueba, con respecto al total de individuos de la muestra. El porcentaje de falsos positivos se calcula en función del número de individuos que, estando sanos, son erróneamente clasificados como enfermos por la prueba, con respecto al total de individuos de la muestra.



**Figura 1.** Ejemplo de una curva ROC obtenida para un marcador tumoral sérico. La gráfica muestra en ordenadas el porcentaje de verdaderos positivos (sensibilidad), y en abscisas el porcentaje de falsos positivos (100-especificidad). El área bajo la curva indica la eficacia diagnóstica del marcador.

Un test diagnóstico perfecto tendría un 100% de especificidad y sensibilidad. Sin embargo, esto es una utopía. Lo que ocurre en la realidad es que si, utilizando la curva ROC, cambiamos el punto de corte para aumentar la sensibilidad, disminuiríamos la especificidad y viceversa. La decisión depende de lo que se quiera conseguir. Por ejemplo, si queremos confirmar la presencia de la enfermedad, el test debería tener una alta especificidad (es decir, pocos falsos positivos); por el contrario, si quisiéramos descartar la enfermedad, entonces el test debería ser muy sensible (es decir, presentar pocos falsos negativos).

Tanto la sensibilidad como la especificidad proporcionan información acerca de la probabilidad de obtener un resultado concreto (positivo o negativo) en función de la verdadera condición del individuo con respecto a la enfermedad. Sin embargo, el clínico normalmente desconoce si el individuo está enfermo o no, y el test debe ayudarle a discernirlo.

En este sentido son útiles los parámetros “valor predictivo positivo” y “valor predictivo negativo”. El VPP<sup>13</sup> es la probabilidad de padecer la enfermedad si se obtiene un resultado positivo en el test, mientras que el VPN<sup>14</sup> es la probabilidad de que un sujeto con un resultado negativo en la prueba esté realmente sano. No obstante, a pesar de su utilidad clínica hay que tener en cuenta que estos parámetros dependen en gran medida de la prevalencia de la enfermedad (Peng, 2006).

Finalmente, a la hora de validar marcadores debe aprovecharse la circunstancia que ofrece la proteómica al detectar simultáneamente varias proteínas alteradas y por lo tanto candidatas. La determinación de parámetros diagnósticos puede hacerse combinando los valores de varios marcadores y aplicando pruebas multivariantes. A modo de ejemplo, en trabajos realizados en nuestro laboratorio hemos observado que combinando los valores determinados en suero para el receptor del factor de crecimiento epidérmico (sEGFR<sup>15</sup>) y tres de sus ligandos (EGF, TGF- $\alpha$  y anfirregulina) se mejora enormemente la sensibilidad y la especificidad obtenidas para cada marcador de forma independiente (Lemos-González *et al.*, 2007). Por otro lado, mediante la aplicación del AD, y con la combinación de los niveles séricos del sEGFR y de sus ligandos EGF y TGF- $\alpha$ , se clasificaron correctamente un 90,2% de individuos sanos y un 100% de pacientes de cáncer no microcítico de pulmón, mientras que la aplicación del modelo de regresión logística, con la misma combinación de sEGFR, EGF y TGF- $\alpha$ , permitió clasificar correctamente un 100% de donantes sanos y un 100% de pacientes de cáncer de cabeza y cuello.

## 7. Análisis estadístico en espectrometría de masas

Aunque no haremos una revisión detallada de la estadística en espectrometría de masas, queremos recordar que los métodos estadísticos son también relevantes en este campo de la proteómica. En general, las pruebas descritas para 2-DE son también aplicables sobre los espectros de masas, donde cada pico de masas es una variable similar a los *spots* de los mapas 2D.

<sup>13</sup> VPP: valor predictivo positivo.

<sup>14</sup> VPN: valor predictivo negativo.

<sup>15</sup> sEGFR: receptor del factor de crecimiento epidérmico.

Un proceso fundamental en el que la estadística está presente es la identificación de proteínas a partir de los espectros obtenidos. Así, motores de búsqueda como ProFound están basados en métodos Bayesianos y utilizan filtros “finos” para establecer las asociaciones entre los datos empíricos y las bases de datos. Otros motores, como SEQUEST, utilizan filtros más “burdos” pero que requieren menos cálculos y son por ello más rápidos aunque sin una pérdida de fiabilidad. PeptideProphet es un programa que emplea pruebas como el AD lineal y la regla de Bayes para convertir los resultados de SEQUEST en una probabilidad (Alterovitz *et al.*, 2007).

En ocasiones los datos de espectrometría de masas no se emplean para la identificación de proteínas, sino para la búsqueda de patrones diferenciales, como en los experimentos de SELDI-TOF<sup>16</sup>. En estos casos, se suelen emplear métodos multivariantes de reducción de datos y clasificación como los anteriormente descritos, o análisis de correlaciones y distancias (Cox *et al.*, 2005). Un ejemplo de la aplicación de métodos de reducción de datos se puede encontrar en el análisis de la localización subcelular de proteínas de membrana presentado por Sadowski *et al.* (2006). Por otra parte, las técnicas de clasificación han sido recientemente empleadas por Meuwis *et al.* (2007), quienes han utilizado SELDI-TOF para la detección de biomarcadores séricos de enfermedad inflamatoria intestinal.

## 8. Conclusiones

La combinación de técnicas estadísticas univariantes y multivariantes es de gran utilidad en la búsqueda de biomarcadores mediante técnicas proteómicas. No sólo maximiza la información que puede obtenerse a partir de mapas 2D, y mejora la detección de cambios proteómicos verdaderos y relevantes como marcadores, sino que también ayuda a su validación biológica.

## Agradecimientos

Los trabajos mencionados han sido parcialmente financiados por la Xunta de Galicia (PGIDIT05PXI B31002PR). A.M. Rodríguez Piñeiro disfruta de un contrato Ángeles Alvariño (Xunta de Galicia).

## 9. Referencias

- Albitar M, Potts S J, Giles F J, O'Brien S *et al.* 2006. Proteomic-based prediction of clinical behavior in adult acute lymphoblastic leukemia. *Cancer* 106: 1587-1594.
- Alterovitz G, Liu J, Afkhami E y Ramoni M F. 2007. Bayesian methods for proteomics. *Proteomics* 7: 2843-2855.
- Álvarez-Chaver P, Rodríguez-Piñeiro A M, Rodríguez-Berrocal F J, Martínez-Zorzano V S *et al.* 2007. Identification of hydrophobic proteins as biomarker candidates for colorectal cancer. *International Journal of Biochemistry and Cell Biology* 39: 529-540.
- Ayude D, Fernández-Rodríguez J, Rodríguez-Berrocal F J, Martínez-Zorzano V S *et al.* 2000. Value of the serum alpha-L-fucosidase activity in the diagnosis of colorectal cancer. *Oncology* 59: 310-316.
- Beck J R y Schultz E K. 1986. The use of relative operating characteristic (ROC) curves in test performance evaluation. *Archives of Pathology and Laboratory Medicine* 110: 13-20.
- Bhattacharyya S, Siegel E R, Petersen G M, Chari S T *et al.* 2004. Diagnosis of pancreatic cancer using serum proteomic profiling. *Neoplasia* 6: 674-686.
- Bookstein F L. 1991. Morphometric tools for landmark data. Cambridge University Press, Nueva York, Estados Unidos.
- Byrjalsen I, Mose Larsen P, Fey S J, Nilas L *et al.* 1999. Two-dimensional gel analysis of human endometrial proteins: characterization of proteins with increased expression in hyperplasia and adenocarcinoma. *Molecular Human Reproduction* 5: 748-756.
- Chang J, van Remmen H, Ward W F, Regnier F E *et al.* 2004. Processing of data generated by 2-dimensional gel electrophoresis for statistical analysis: missing data, normalization, and statistics. *Journal of Proteome Research* 3: 1210-1218.
- Chich J-F, David O, Villers F, Schaeffer B *et al.* 2007. Statistics for proteomics: Experimental design and 2-DE differential analysis. *Journal of Chromatography B* 849: 261-272.
- Cox B, Kislinger T y Emili A. 2005. Integrating gene and protein expression data: pattern analysis and profile mining. *Methods* 35: 303-314.
- De Noo M E, Tollenaar R A E M, Deelder A M y Bouwman L H. 2006. Current status and prospects

<sup>16</sup> SELDI-TOF: ionización/desorción por láser potenciada por superficie, acoplada a tiempo de vuelo.

- of clinical proteomics studies on detection of colorectal cancer: Hopes and fears. *World Journal of Gastroenterology* 12: 6594-6601.
- Dudoit S, Shaffer J P y Boldrick J C. 2003. Multiple hypothesis testing in microarray experiments. *Statistical Science* 18: 71-103.
- Garrels J I. 1989. The QUEST system for quantitative analysis of two-dimensional gels. *Journal of Biological Chemistry* 264: 5269-5282.
- Horgan G W. 2007. Sample size and replication in 2D gel electrophoresis studies. *Journal of Proteome Research* 6: 2884-2887.
- Hunt S M N, Thomas M R, Sebastian L T, Pedersen S K *et al.* 2005. Optimal replication and the importance of experimental design for gel-based quantitative proteomics. *Journal of Proteome Research* 4: 809-819.
- Karp N A, Griffin J L y Lilley K S. 2005. Application of partial least squares discriminant analysis to two-dimensional difference gel studies in expression proteomics. *Proteomics* 5: 81-90.
- Kemperman R F, Horvatovich P L, Hoekman B, Reijmers T H *et al.* 2007. Comparative urine analysis by liquid chromatography-mass spectrometry and multivariate statistics: method development, evaluation, and application to proteinuria. *Journal of Proteome Research* 6: 194-206.
- Kovárová H, Hajdúch M, Korínková G, Halada P *et al.* 2000. Proteomics approach in classifying the biochemical basis of the anticancer activity of the new olomoucine-derived synthetic cyclin-dependent kinase inhibitor, bohémine. *Electrophoresis* 21: 3757-3764.
- LaBaer J. 2005. So, you want to look for biomarkers (Introduction to the special biomarkers issue). *Journal of Proteome Research* 4: 1053-1059.
- Lemos-González Y, Rodríguez-Berrocal F J, Cordero O J, Gómez C *et al.* 2007. Alteration of the serum levels of the epidermal growth factor and its ligands in patients with non-small cell lung cancer and head and neck carcinoma. *British Journal of Cancer* 96: 1569-1578.
- Marengo E, Leardi R, Robotti E, Righetti P G *et al.* 2003. Application of three-way principal component analysis to the evaluation of two-dimensional maps in proteomics. *Journal of Proteome Research* 2: 351-360.
- Meleth S, Deshane J y Kim H. 2005. The case for well-conducted experiments to validate statistical protocols for 2D gels: different pre-processing = different lists of significant proteins. *BMC Biotechnology* 5: 7-21.
- Meuwis M A, Fillet M, Geurts P, de Seny D *et al.* 2007. Biomarker discovery for inflammatory bowel disease, using proteomic serum profiling. *Biochemical Pharmacology* 73: 1422-1433.
- Peng X. 2006. Developing and evaluating genomics- or proteomics-based diagnostic tests: statistical perspectives. *Methods in Molecular Medicine* 129: 27-39.
- Pepe M S, Feng Z, Huang Y, Longton G *et al.* 2007. Integrating the predictiveness of a marker with its performance as a classifier. *American Journal of Epidemiology* (en prensa, DOI 10.1093/aje/kwm305).
- Perco P, Rapberger R, Siehs C, Lukas A *et al.* 2006. Transforming omics data into context: Bioinformatics on genomics and proteomics raw data. *Electrophoresis* 27: 2659-2675.
- Randic M. 2006. Quantitative characterizations of proteome: Dependence on the number of proteins considered. *Journal of Proteome Research* 5: 1575-1579.
- Roblick U J, Hirschberg D, Habermann J K, Palmberg C *et al.* 2004. Sequential proteome alterations during genesis and progression of colon cancer. *Cellular and Molecular Life Sciences* 61: 1246-1255.
- Rodríguez-Piñeiro A M, Ayude D, Rodríguez-Berrocal F J y Páez de la Cadena M. 2004. Concanavalin A chromatography coupled to two-dimensional gel electrophoresis improves protein expression studies of the serum proteome. *Journal of Chromatography B* 803: 337-343.
- Rodríguez-Piñeiro A M, Carvajal-Rodríguez A, Rolán-Alvarez E, Rodríguez-Berrocal F J *et al.* 2005. Application of relative warp analysis to the evaluation of two-dimensional gels in proteomics: Studying isoelectric point and relative molecular mass variation. *Journal of Proteome Research* 4: 1318-1323.
- Rodríguez-Piñeiro A M, Rodríguez-Berrocal F J y Páez de la Cadena M. 2007. Improvements in the search for potential biomarkers by proteomics: Application of principal component and discriminant analyses for two-dimensional maps evaluation. *Journal of Chromatography B* 849: 251-260.
- Rohlf, F. J. 1993. *Contributions to Morphometrics*. Museo Nacional de Ciencias Naturales, Madrid, España.

- Rosengren A T, Salmi J M, Aittokallio T, Westerholm J *et al.* 2003. Comparison of PDQuest and Progenesis software packages in the analysis of two-dimensional electrophoresis gels. *Proteomics* 3: 1936-1946.
- Sadowski P, Dunkley T P J, Shadforth I P, Dupree P *et al.* 2006. Quantitative proteomic approach to study subcellular localization of membrane proteins. *Nature Protocols* 1: 1778-1789.
- Smit S, Hoefsloot H C J y Smilde A K. 2007. Statistical data processing in clinical proteomics. *Journal of Chromatography B* (en prensa, DOI 10.1016/j.jchromb.2007.10.042).
- Tong W, Hong H, Fang H, Xie Q *et al.* 2003. Decision forest: combining the predictions of multiple independent decision tree models. *Journal of Chemical Information and Computer Sciences* 43: 525-531.
- Verhoeckx K C, Gaspari M, Bijlsma S, van der Greef J *et al.* 2005. In search of secreted protein biomarkers for the anti-inflammatory effect of beta2-adrenergic receptor agonists: application of DIGE technology in combination with multivariate and univariate data analysis tools. *Journal of Proteome Research* 4: 2015-2023.
- Wasinger V, Cordwell S, Cerpa-Poljak A, Yan J *et al.* 1995. Progress with gene-product mapping of the Mollicutes: *Mycoplasma genitalium*. *Electrophoresis* 16: 1090-1094.
- Wilkins M R, Pasquali C, Appel R D, Ou K *et al.* 1996. From proteins to proteomes: large scale protein identification by two-dimensional electrophoresis and amino acid analysis. *BioTechnology* 14: 61-65.
- Wilkins M R, Appel R D, Van Eyk J E, Chung M C M *et al.* 2006. Guidelines for the next 10 years of proteomics. *Proteomics* 6: 4-8.
- Zelditch M L; Swiderski D L; Sheets H D; Fink W L. 2004. *Geometric Morphometrics for Biologists: A Primer*. Academic Press, Nueva York, Estados Unidos.

## El plasma ICP para la ionización elemental en Espectrometría de Masas: su aplicación en Proteómica Cuantitativa”

*Alfredo Sanz-Medel*

Departamento de Química Física y Analítica, Universidad de Oviedo, C/Julián Clavería 8, 33006 Oviedo (Asturias). E-mail: asm@uniovi.es

### Introducción

A estas alturas ya nadie discute que la Espectrometría de Masas (EM) se ha convertido en una de las técnicas clave, probablemente la más poderosa por su sensibilidad y aplicabilidad a problemas reales, para la caracterización de proteínas. No es de extrañar, pues, que las técnicas de EM ya clásicas (e.g. MALDI-MS y ESI-(MS)<sup>n</sup>) sean hoy imprescindibles en estudios de proteómica estructural y funcional. Asimismo, la importancia y aplicaciones de la EM analítica crece enormemente en los últimos años también en el área emergente de la biología de sistemas.

En cualquier caso, tales técnicas “convencionales” de EM proporcionan una información química

molecular, es decir, de las moléculas, de los aminoácidos, péptidos y proteínas buscados. Mucho menos conocida en el mundo de la Proteómica (mi asistencia a congresos de Espectrometría de Masas en los simposios dedicados a su aplicación en proteínas así lo corrobora de forma irrefutable) es la Espectrometría de Masas Elemental, que utiliza como fuente de ionización un plasma inducido de radiofrecuencias (el ICP-MS). En este caso, la información directa obtenida es atómica, es decir, de los elementos constitutivos de la biomolécula. Por eso, para conocer la naturaleza de la especie química concreta (biomolécula) donde se halla el elemento medido por ICP-MS es preciso acoplar ese detector a un sistema de separación previa y potente de la biomolécula estudiada (p.e. HPLC, CE o electroforesis 2D) en la