

FACULTAD DE INGENIERÍA DE LA UNIVERSIDAD DE LA REPÚBLICA

INTRODUCCIÓN AL RECONOCIMIENTO DE PATRONES

CURSO 2013

Proyecto final

Autores:

José Luis NUNES
Matías TAILANIÁN

Tutores:

Federico LECUMBERRY
Ignacio RAMIREZ



11 de diciembre de 2013

Índice general

1. Introducción	2
2. Base de datos	3
2.1. Características fenotípicas	3
2.2. Características genotípicas - Determinación de polimorfismos (SNPs)	4
2.3. Resumen	4
3. Técnicas utilizadas	5
3.1. Primera etapa	6
3.2. Segunda etapa - Clases balanceadas	7
3.2.1. Naive Bayes	8
3.2.2. C 4.5	8
3.2.3. K-nn	8
3.3. PCA	8
3.4. LDA?	8
3.5. NN?	8
3.6. SVM?	8
4. Conclusiones	9

Introducción

El objetivo principal del proyecto es la investigación en técnicas que permitan contribuir con la predicción de fertilidad de rodeo lechero y la calidad de la carne integrando técnicas de reconocimiento de patrones sobre datos de alta dimensión.

A lo largo de los años se ha intentado relacionar ciertas características fenotípicas del ganado bovino con algunos indicadores genéticos. Es un problema actualmente abierto de gran interés mundial y Uruguay, siendo un país esencialmente ganadero, no puede estar ajeno. Desde hace varios años se viene desarrollando una línea de investigación relacionada con la genética molecular de la calidad de la carne bovina en el área de Genética de la Facultad de Veterinaria. Se ha realizado un abordaje desde diferentes puntos de vista, como la caracterización de genes conocidos en rodeos vacunos de nuestro país, o la búsqueda y análisis de nuevos genes asociados a la calidad de carne tanto bovina como ovina [1].

Obtener una correlación entre la información genética y la capacidad reproductiva de bovinos, el impacto en la calidad de la leche y de la carne de los animales de nuestros rodeos tendría un alto impacto en la producción. A lo largo de los últimos años en Facultad de Veterinaria se han creado bases de datos con información fenotípica y genotípica relacionada con la calidad cárnica y la fertilidad bovina aplicada a la producción lechera. En esta oportunidad se trabajará con una base de datos confeccionada en nuestro país durante el año 2009, provista por la veterinaria Dra. Ana Meikle.

Base de datos

Para lo que sigue de este trabajo es importante conocer la base de datos con la que se trabajará, y realizar una breve descripción de las características relevantes.

La base de datos consta de información de un seguimiento realizado durante 9 meses sobre 891 vacas de 7 tambos diferentes y consiste en las características detalladas a continuación.

2.1. Características fenotípicas

- Edad.
- Condición corporal (BCS) desde 30 días previos al parto hasta 120 días posparto, con una frecuencia de al menos una vez cada 30 días. Se utilizó la escala de 5 puntos (1=flaca, 5= gorda).
- Cantidad de partos (Cantidad de lactancias). Es un indicador importante que tiene que ver con la historia y desgaste de la vaca.
- Anestro: es la cantidad de días que pasaron desde el parto hasta el reinicio. El reinicio se define como el día en que la progesterona aumenta a un nivel determinado, indicando que la vaca volvió a ciclar.
- Intervalo entre partos. Es la cantidad de días que pasan entre 2 partos consecutivos (válido solamente para vacas multíparas). Para maximizar la producción de leche lo que se busca es que la vaca quede peñada una vez al año, es decir un intervalo entre partos de 365 aproximadamente.
- Secado: la cantidad de días que pasan entre el último día que se le saca leche a la vaca y el parto. Es válido solo para vacas multíparas. Es un indicador de cuanto tiempo se dejó descansar a la vaca antes del parto. Cuanto más tiempo tiene, se prepara el cuerpo y llega en mejor forma.
- Servicios: Cantidad de inceminaciones realizadas para lograr la preñez.
- Concentración de progesterona hasta los 60 días posparto.

- Promedio de cantidad de grasa en la leche durante los primeros 100 días posparto.
- Promedio de cantidad de leche durante los primeros 100 días posparto.

2.2. Características genotípicas - Determinación de polimorfismos (SNPs)

Se extrajo ADN y se obtuvieron muestras de buena calidad. En la figura 2.1 se muestra el procedimiento realizado para caracterizar el gen IGF-I bovino en 3 clases: “AA”, “AB” y “BB”.

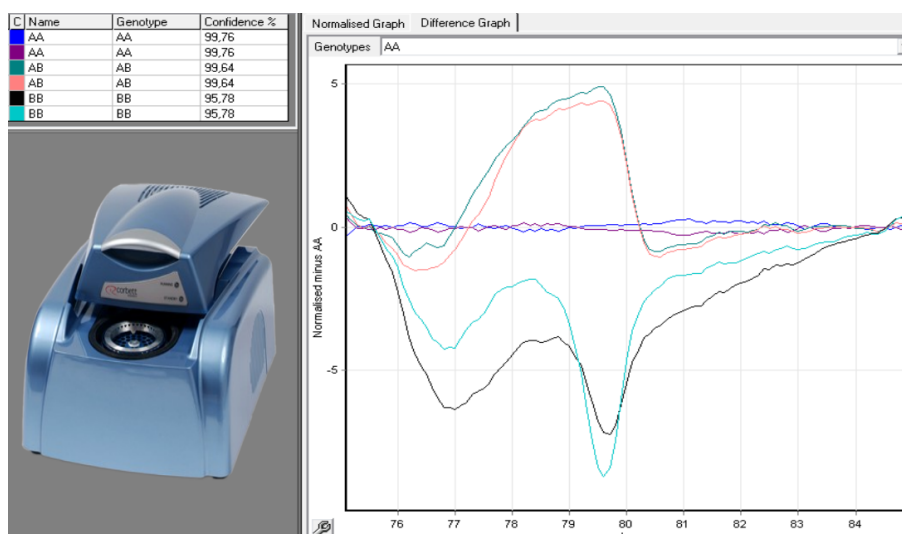


Figura 2.1: Determinación del genotipo

El eje de las ordenadas muestra la fluorescencia normalizada al genotipo “AA”. Se puede observar que el genotipo “AB” se encuentra por encima del genotipo normalizado, mientras que el genotipo “BB” por debajo.

2.3. Resumen

En primer lugar fue necesario entender e interiorizarse con la base de datos, los conceptos y la terminología específica del problema. La lectura e interpretación de los datos no resulta para nada sencilla, por lo que se realizaron dos reuniones con la Dr. Ana Meikle, encargada de la investigación que llevó a cabo la confección de la base asociada a la producción de leche y fertilidad del ganado bovino. En ambas reuniones se buscó depurar la base con el fin de quitar información repetida y no relevante para el estudio. También se estableció un orden de relevancia en las características fenotípicas.

En resumen contamos con una base de datos acotada y “limpia” con varias características fenotípicas que se quieren correlacionar con los genotipos de cada individuo. En particular se abordará el problema como un trabajo de clasificación, tomando los genotipos como clases.

Técnicas utilizadas

La mayoría de los algoritmos de *machine learning* seleccionan los atributos apropiados para realizar sus decisiones. Por ejemplo los métodos de árbol de decisión eligen en cada paso la característica que mejor separa en clases. Cuantas más características tengamos, en teoría tendremos más poder de discriminación, pero en la práctica no ocurre así. Agregar información irrelevante o características distractivas, confunde a los algoritmos de machine learning. Está probado que el número de instancias necesarias para entrenamiento para producir un cierto nivel de performance crece exponencialmente con la cantidad de atributos irrelevantes.

Dado el efecto negativo de las características irrelevantes es muy común en muchos algoritmos de *machine learning* que se realice una etapa de selección de características donde se eliminen las irrelevantes. La mejor manera de realizar una selección de características es manual, adquiriendo un amplio conocimiento del problema en cuestión y logrando interpretar cada una de las características. Sin embargo es posible realizar una selección con algoritmos automáticos que resultan muy útiles. Reducir la dimensionalidad de los datos puede mejorar la performance de los algoritmos, además de reducir la complejidad del problema y bajar los requerimientos de capacidad de cómputo. Por esta razón se realizará una etapa de selección automática de características previa a la clasificación, como se explicará más adelante.

En una primera etapa de análisis sobre la base de datos se ataca el problema con los clasificadores **C4.5**, **Naive Bayes** y **K-NN**.

El algoritmo **C4.5** se utiliza para generar un árbol de decisión que puede ser utilizado para clasificación. Por la naturaleza de los árboles, es un algoritmo no paramétrico, por lo que resulta robusto ante *outliers*.

Naive Bayes es un clasificador basado en la aplicación del teorema de Bayes. No particiona el espacio de instancias e ignora de forma robusta a las características irrelevantes. Asume por diseño que todas las características son independientes entre si, y paga un precio muy alto cuando hay características redundantes.

El algoritmo **k-NN** (k-nearest neighbors) es un algoritmo de clasificación super-

vizado no paramétrico que pedice la clase de las instancias basado en los k vecinos más cercanos. Es un algoritmo muy sensible a la propia estructura de los datos. Cuando $k \rightarrow \infty$, el algoritmo asegura una tasa de error no superior al doble de la tasa de error de Bayes (mínimo alcanzable dada la distribución de los datos).

Con esta batería de clasificadores se cubre un amplio espectro y se utilizan algunos de los algoritmos más utilizados para problemas de reconocimiento de patrones.

Para analizar los resultados se utilizan las implementaciones provistas en el software **Weka** [2].

3.1. Primera etapa

Se utiliza el clasificador compuesto **AttributeSelectedClassifier** que aplica una técnica de selección de características antes de entrenar al clasificador. Se logra entonces una reducción de la dimensionalidad. La estrategia elegida para la selección de características se realiza utilizando el enfoque **wrapper**, que evalúa el set de atributos utilizando un esquema de aprendizaje y utiliza validación cruzada para estimar la precisión del esquema de aprendizaje. El clasificador utilizado para estimar esta precisión es un árbol de decisión **C 4.5**.

En todos los casos se utilizó validación cruzada con 10 subconjuntos.

Una vez hecha la extracción de características, y así la reducción de dimensionalidad, se estudia el desempeño de diferentes clasificadores. En la tabla 3.1 se muestran los resultados para los algoritmos C4.5, Naive Bayes y k-NN.

	Tiempo [s]	$\sqrt{\text{MSE}}$	F-Measure	κ	Bien Clasif [%]
C4.5	0.84	0.45	0.331	0	49.83
Bayes	0.37	0.45	0.345	0.0044	49.60
k-NN	0.37	0.53	0.402	0.0043	45.23

Tabla 3.1: Resultados etapa 1

Se puede observar en la tabla anterior que se obtienen resultados muy similares para los 3 algoritmos. Para C4.5 y Bayes se obtiene un porcentaje de aciertos un poco menor al 50 %, mientras que para k-NN los resultados son un poco inferiores.

La medida de error *kappa-statistic* (κ), es un indicador de la performance del algoritmo que tiene en cuenta las coincidencias por azar. Se calcula como

$$\kappa = \frac{P_0 - P_e}{1 - P_e}$$

donde P_0 es la proporción de coincidencias observadas y P_e la proporción de coincidencias esperadas en las hipótesis de independencia, es decir, coincidencias por azar. Se puede ver en la tabla 3.1 que se obtuvieron valores de κ realmente bajísimos, siendo este un indicador más de la mala performance alcanzada por los algoritmos.

Por otro lado resulta interesante analizar las matrices de confusión que resultan de estos algoritmos. Para el C4.5 se obtiene la siguiente matriz:

1	a	b	c	<— classified as	C4.5
2	0	309	0	a = AA	
3	0	444	0	b = AB	
4	0	138	0	c = BB	

Claramente el resultado obtenido no es el esperado. En este caso clasifica todos los patrones como pertenecientes a la clase “AB”, y como esta clase representa casi el 50% de todas las muestras, el porcentaje de aciertos coincide. Este es un resultado determinístico, que más allá del porcentaje de aciertos, significa que el clasificador no funcionó adecuadamente. A su vez, analizando el árbol de decisión se puede ver que tiene una sola hoja.

Por otro lado las matrices de confusión para los algoritmos *Naive Bayes* y *k-NN* son las siguientes:

1	a	b	c	<— classified as	Bayes
2	6	300	3	a = AA	
3	8	435	1	b = AB	
4	1	136	1	c = BB	

1	a	b	c	<— classified as	k-NN
2	63	233	13	a = AA	
3	81	331	32	b = AB	
4	18	111	9	c = BB	

Si bien estos dos casos no se obtuvo un resultado determinístico como con C4.5, igualmente los resultados tienen un fuerte sesgo hacia la clasificación de los patrones como pertenecientes a la clase “AB”.

Para mitigar el fenómeno de la salida determinística (y el sesgo) mencionado, el siguiente paso es atacar el problema del desbalance de clases. Para ello se realiza un sorteo aleatorio de las muestras pertenecientes a las 2 clases mayoritarias, de forma que las 3 clases tengan la misma cantidad de patrones. Los resultados se presentan en la siguiente sección.

3.2. Segunda etapa - Clases balanceadas

	Tiempo [s]	# Caract.	$\sqrt{\text{MSE}}$	F-Measure	Mal Clasif [%]
Bayes	0.62	23	0.166	0.957	4.33
C4.5	0.28	23	0.315	0.839	15.88

Tabla 3.2: CfsSubset con búsqueda BestFirst

3.2.1. Naive Bayes

3.2.2. C 4.5

3.2.3. K-nn

3.3. PCA

3.4. LDA?

3.5. NN?

3.6. SVM?

Conclusiones

Bibliografía

- [1] Eileen Armstrong Reborati, *Detección y análisis de genes asociados a la calidad de la carne en bovinos*, Tesis de doctorado, Madrid, 2011.
- [2] Mark Hall, Eibe Frank, Geoffrey Holmes, Bernhard Pfahringer, Peter Reutemann, Ian H. Witten (2009); The WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11.
- [3] Meyer, K. (2007). WOMBAT – A tool for mixed model analyses in quantitative genetics by REML, J. Zhejiang Uni.