

# Métodos de Regresión y Clasificación Lineales

Alvaro J. Riascos Villegas

Junio de 2019

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN
- 3 Ejemplo: Cáncer de próstata
- 4 Selección de variables
- 5 Métodos de contracción (shrinkage)
- 6 Modelos de lineales de clasificación

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathcal{F}$ .

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathcal{F}$ .

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathcal{F}$ .

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathcal{F}$ .

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathcal{F}$ .

## Modelo de Regresión Óptimo

- Supongamos que  $Y$  es una variable real y  $X$  es un  $p$ -vector:  
 $X \in \mathbb{R}^p$
- Denotamos la distribución conjunta por  $Pr(X, Y)$ .
- Nuestro objetivo es encontrar  $f : X \rightarrow Y$  para, dado  $X$ , predecir  $Y$  pero necesitamos una noción de pérdida.
- Vamos a utilizar el error al cuadrado (alternativamente podría ser del error absoluto):

$$R(f) = E[(Y - f(X))^2] \quad (1)$$

- $R(f)$  denota el error de prueba (o Riesgo de acuerdo a la terminología introducida en las clases anteriores).
- El problema es minimizar este error en un espacio de funciones  $\mathfrak{F}$ .

## Modelo de Regresión Óptimo

- $R(f)$  se puede escribir como:

$$R(f) = E_X E_{Y|X} [(Y - f(X))^2 | X]$$

- Ahora, si para cada  $X = x$  hacemos que  $f(x)$  sea tal que:

$$E_{Y|X} [(Y - f(x))^2 | X = x] \tag{2}$$

sea mínimo, entonces la función que a  $x$  le asocia  $f(x)$  resuelve el problema.

## Modelo de Regresión Óptimo

- $R(f)$  se puede escribir como:

$$R(f) = E_X E_{Y|X} [(Y - f(X))^2 | X]$$

- Ahora, si para cada  $X = x$  hacemos que  $f(x)$  sea tal que:

$$E_{Y|X} [(Y - f(x))^2 | X = x] \tag{2}$$

sea mínimo, entonces la función que a  $x$  le asocia  $f(x)$  resuelve el problema.

## Modelo de Regresión Óptimo

- Luego, formalmente el problema que queremos resolver es:

$$f(x) = \arg \min_c E_{Y|X} [(Y - c)^2 | X = x]$$

- Se puede demostrar que la solución a este problema es:

$$f(x) = E(Y|X = x)$$

## Modelo de Regresión Óptimo

- Luego, formalmente el problema que queremos resolver es:

$$f(x) = \arg \min_c E_{Y|X} [(Y - c)^2 | X = x]$$

- Se puede demostrar que la solución a este problema es:  
 $f(x) = E(Y|X = x)$

## Otras funciones de pérdida

- Si usamos como función de pérdida  $E[|Y - f(X)|]$  se puede demostrar que la solución óptima es:

$$\hat{f}(x) = \text{Median}(Y | X = x)$$

- Esta es menos sensible a datos atípicos pero mas difícil de trabajar analíticamente (no es diferenciable).

## Otras funciones de pérdida

- Si usamos como función de pérdida  $E[|Y - f(X)|]$  se puede demostrar que la solución óptima es:

$$\hat{f}(x) = \text{Median}(Y | X = x)$$

- Esta es menos sensible a datos atípicos pero más difícil de trabajar analíticamente (no es diferenciable).

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN**
- 3 Ejemplo: Cáncer de próstata
- 4 Selección de variables
- 5 Métodos de contracción (shrinkage)
- 6 Modelos de lineales de clasificación

# Regresión lineal

- El modelo de regresión lineal asume que:

$$f(x) \approx X^T \beta$$

- Si minimizamos el riesgo sujeto a que las funciones deben ser lineales obtenemos:

$$\beta = E (XX^T)^{-1} E(XY)$$

- El modelo de regresión lineal asume  $f(x)$  es globalmente lineal.

## Regresión lineal

- El modelo de regresión lineal asume que:

$$f(x) \approx X^T \beta$$

- Si minimizamos el riesgo sujeto a que las funciones deben ser lineales obtenemos:

$$\beta = E (XX^T)^{-1} E(XY)$$

- El modelo de regresión lineal asume  $f(x)$  es globalmente lineal.

## Regresión lineal

- El modelo de regresión lineal asume que:

$$f(x) \approx X^T \beta$$

- Si minimizamos el riesgo sujeto a que las funciones deben ser lineales obtenemos:

$$\beta = E (XX^T)^{-1} E(XY)$$

- El modelo de regresión lineal asume  $f(x)$  es globalmente lineal.

## k-NN: k-vecinos más cercanos

- k-NN estima la esperanza condicional localmente como una función constante.

$$f(x) \approx \text{Ave}(y | x \in N_k(x))$$

## Comparación entre la regresión lineal y k-NN

- Ambos métodos aproximan  $E(Y|X = x)$  con promedios pero hacen supuestos muy distintos sobre la verdadera función de aprendizaje.
  - El modelo de regresión lineal asume  $f(x)$  es globalmente lineal.
  - k-NN asume que  $f(x)$  es localmente constante.

## Comparación entre la regresión lineal y k-NN

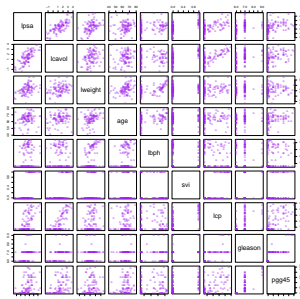
- Ambos métodos aproximan  $E(Y|X = x)$  con promedios pero hacen supuestos muy distintos sobre la verdadera función de aprendizaje.
  - El modelo de regresión lineal asume  $f(x)$  es globalmente lineal.
  - k-NN asume que  $f(x)$  es localmente constante.

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN
- 3 Ejemplo: Cáncer de próstata**
- 4 Selección de variables
- 5 Métodos de contracción (shrinkage)
- 6 Modelos de lineales de clasificación

## Ejemplo: Cáncer de próstata

- Stamey et.al (1989). Examina la relación entre el antígeno prostático (Ipsa) y marcadores clínicos.



**FIGURE 1.1.** Scatterplot matrix of the prostate cancer data. The first row shows the response against each of the predictors in turn. Two of the predictors, `svi` and `gleason`, are categorical.

- Correlación entre la variables predictoras.

**TABLE 3.1.** *Correlations of predictors in the prostate cancer data.*

	lcavol	lweight	age	lbph	svi	lcp	gleason
lweight	0.300						
age	0.286	0.317					
lbph	0.063	0.437	0.287				
svi	0.593	0.181	0.129	-0.139			
lcp	0.692	0.157	0.173	-0.089	0.671		
gleason	0.426	0.024	0.366	0.033	0.307	0.476	
pgg45	0.483	0.074	0.276	-0.030	0.481	0.663	0.757

- Regresión lineal entre lpsa y las variables predictoras.

**TABLE 3.2.** *Linear model fit to the prostate cancer data. The Z score is the coefficient divided by its standard error (3.12). Roughly a Z score larger than two in absolute value is significantly nonzero at the  $p = 0.05$  level.*

Term	Coefficient	Std. Error	Z Score
Intercept	2.46	0.09	27.60
lcavol	0.68	0.13	5.37
lweight	0.26	0.10	2.75
age	-0.14	0.10	-1.40
lbph	0.21	0.10	2.06
svi	0.31	0.12	2.47
lcp	-0.29	0.15	-1.87
gleason	-0.02	0.15	-0.15
pgg45	0.27	0.15	1.74

## Ejemplo: Cáncer de próstata

- **Primero se estandarizan variables para tener varianza unitaria.**
- Se elige muestra de entrenamiento de 67 observaciones y prueba de 30.
- Los Z Score son estadísticos de pruebas de hipótesis de coeficientes iguales a cero (valor absoluto mayor que 2 significa que se rechaza la hipótesis nula con una confianza del 95%).
- Si se usa como predictor de la *lpsa* el promedio de la variable en la muestra de entrenamiento, el error de prueba es 1.057 (i.e., error base). El modelo arroja un error de predicción de 0.521.

## Ejemplo: Cáncer de próstata

- Primero se estandarizan variables para tener varianza unitaria.
- Se elige muestra de entrenamiento de 67 observaciones y prueba de 30.
- Los Z Score son estadísticos de pruebas de hipótesis de coeficientes iguales a cero (valor absoluto mayor que 2 significa que se rechaza la hipótesis nula con una confianza del 95%).
- Si se usa como predictor de la *lpsa* el promedio de la variable en la muestra de entrenamiento, el error de prueba es 1.057 (i.e., error base). El modelo arroja un error de predicción de 0.521.

## Ejemplo: Cáncer de próstata

- Primero se estandarizan variables para tener varianza unitaria.
- Se elige muestra de entrenamiento de 67 observaciones y prueba de 30.
- Los Z Score son estadísticos de pruebas de hipótesis de coeficientes iguales a cero (valor absoluto mayor que 2 significa que se rechaza la hipótesis nula con una confianza del 95%).
- Si se usa como predictor de la *lpsa* el promedio de la variable en la muestra de entrenamiento, el error de prueba es 1.057 (i.e., error base). El modelo arroja un error de predicción de 0.521.

## Ejemplo: Cáncer de próstata

- Primero se estandarizan variables para tener varianza unitaria.
- Se elige muestra de entrenamiento de 67 observaciones y prueba de 30.
- Los Z Score son estadísticos de pruebas de hipótesis de coeficientes iguales a cero (valor absoluto mayor que 2 significa que se rechaza la hipótesis nula con una confianza del 95%).
- Si se usa como predictor de la *lpsa* el promedio de la variable en la muestra de entrenamiento, el error de prueba es 1.057 (i.e., error base). El modelo arroja un error de predicción de 0.521.

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN
- 3 Ejemplo: Cáncer de próstata
- 4 Selección de variables**
- 5 Métodos de contracción (shrinkage)
- 6 Modelos de lineales de clasificación

## Selección de variables

- Dos problemas típicos:
  - 1 Error de prueba: Es posible reducir el error de prueba disminuyendo el número de variables (reduciendo complejidad y varianza) aunque aumente el sesgo: reducir error de estimación por más error de aproximación.
  - 2 Interpretación: Un menor número de variables usualmente permite una mejor interpretación.
- Vamos a discutir diferentes formas de reducir el número de variables, Todas estas técnicas son ejemplos de técnicas de selección de modelos que más adelante volveremos a discutir.

## Selección de variables

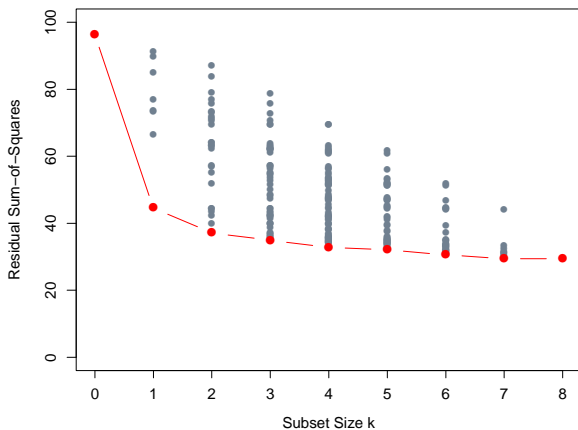
- Dos problemas típicos:
  - 1 Error de prueba: Es posible reducir el error de prueba disminuyendo el número de variables (reduciendo complejidad y varianza) aunque aumente el sesgo: reducir error de estimación por más error de aproximación.
  - 2 Interpretación: Un menor número de variables usualmente permite una mejor interpretación.
- Vamos a discutir diferentes formas de reducir el número de variables, Todas estas técnicas son ejemplos de técnicas de selección de modelos que más adelante volveremos a discutir.

- Mejor subconjunto de variables.
  - Se elige el subconjunto de variables que minimiza el error de prueba.
  - Computacionalmente intensivo. Es computacionalmente viable solo para casos con menos de 40 variables predictoras.

# Selección de variables: Mejor subconjunto de variables

- Mejor subconjunto de variables.
  - Se elige el subconjunto de variables que minimiza el error de prueba.
  - Computacionalmente intensivo. Es computacionalmente viable solo para casos con menos de 40 variables predictoras.

# Selección de variables: Mejor subconjunto de variables



**FIGURE 3.5.** All possible subset models for the prostate cancer example. At each subset size is shown the residual sum-of-squares for each model of that size.

# Selección de variables: Forward Backward selection

- Forward: Comenzar con un modelo que solo tiene una constante y añadir secuencialmente la variable que mas ayude a reducir el error de predicción. En cada etapa se reestima.
- Es computacionalmente menos demandante.
- En la medida que es un búsqueda restringida (a la forma específica como se recorren las variables) tiene mayor sesgo pero potencialmente menos varianza que el mejor subconjunto.
- Backward: Comienza con el modelo que tiene todas la variables y va eliminando la variable que menos contribuye al poder predictivo (una forma de implementarlo es eliminando la variable con el menor Z-score). En cada etapa se reestima.

# Selección de variables: Forward Backward selection

- Forward: Comenzar con un modelo que solo tiene una constante y añadir secuencialmente la variable que mas ayude a reducir el error de predicción. En cada etapa se reestima.
- Es computacionalmente menos demandante.
- En la medida que es un búsqueda restringida (a la forma específica como se recorren las variables) tiene mayor sesgo pero potencialmente menos varianza que el mejor subconjunto.
- Backward: Comienza con el modelo que tiene todas la variables y va eliminando la variable que menos contribuye al poder predictivo (una forma de implementarlo es eliminando la variable con el menor Z-score). En cada etapa se reestima.

# Selección de variables: Forward Backward selection

- Forward: Comenzar con un modelo que solo tiene una constante y añadir secuencialmente la variable que mas ayude a reducir el error de predicción. En cada etapa se reestima.
- Es computacionalmente menos demandante.
- En la medida que es un búsqueda restringida (a la forma específica como se recorren las variables) tiene mayor sesgo pero potencialmente menos varianza que el mejor subconjunto.
- Backward: Comienza con el modelo que tiene todas la variables y va eliminando la variable que menos contribuye al poder predictivo (una forma de implementarlo es eliminando la variable con el menor Z-score). En cada etapa se reestima.

# Selección de variables: Forward Backward selection

- Forward: Comenzar con un modelo que solo tiene una constante y añadir secuencialmente la variable que mas ayude a reducir el error de predicción. En cada etapa se reestima.
- Es computacionalmente menos demandante.
- En la medida que es un búsqueda restringida (a la forma específica como se recorren las variables) tiene mayor sesgo pero potencialmente menos varianza que el mejor subconjunto.
- Backward: Comienza con el modelo que tiene todas la variables y va eliminando la variable que menos contribuye al poder predictivo (una forma de implementarlo es eliminando la variable con el menor Z-score). En cada etapa se reestima.

# Selección de variables: Forward stagewise

- Comienza con el intercepto y va añadiendo variables según la variable que más correlación tenga con el residuo del modelo anterior.
- No se reestima, se añade la variable si aun no esta en la regresión, o se le añe de el coeficiente de correlación al coeficiente que en ese momento tiene la variable (i.e., puede ser cero en ese momento).
- Es una búsqueda más restringida que forward stepwise.
- Genera un conjunto anidado de variables (i.e. solo aumentan y no cambian en cada iteración).
- Es un método competitivo cuando existen muchas variables.

# Selección de variables: Forward stagewise

- Comienza con el intercepto y va añadiendo variables según la variable que más correlación tenga con el residuo del modelo anterior.
- No se reestima, se añade la variable si aun no esta en la regresión, o se le añe de el coeficiente de correlación al coeficiente que en ese momento tiene la variable (i.e., puede ser cero en ese momento).
- Es una búsqueda más restringida que forward stepwise.
- Genera un conjunto anidado de variables (i.e. solo aumentan y no cambian en cada iteración).
- Es un método competitivo cuando existen muchas variables.

# Selección de variables: Forward stagewise

- Comienza con el intercepto y va añadiendo variables según la variable que más correlación tenga con el residuo del modelo anterior.
- No se reestima, se añade la variable si aun no esta en la regresión, o se le añe de el coeficiente de correlación al coeficiente que en ese momento tiene la variable (i.e., puede ser cero en ese momento).
- Es una búsqueda más restringida que forward stepwise.
- Genera un conjunto anidado de variables (i.e. solo aumentan y no cambian en cada iteración).
- Es un método competitivo cuando existen muchas variables.

## Selección de variables: Forward stagewise

- Comienza con el intercepto y va añadiendo variables según la variable que más correlación tenga con el residuo del modelo anterior.
- No se reestima, se añade la variable si aun no esta en la regresión, o se le añe de el coeficiente de correlación al coeficiente que en ese momento tiene la variable (i.e., puede ser cero en ese momento).
- Es una búsqueda más restringida que forward stepwise.
- Genera un conjunto anidado de variables (i.e. solo aumentan y no cambian en cada iteración).
- Es un método competitivo cuando existen muchas variables.

## Selección de variables: Forward stagewise

- Comienza con el intercepto y va añadiendo variables según la variable que más correlación tenga con el residuo del modelo anterior.
- No se reestima, se añade la variable si aun no esta en la regresión, o se le añe de el coeficiente de correlación al coeficiente que en ese momento tiene la variable (i.e., puede ser cero en ese momento).
- Es una búsqueda más restringida que forward stepwise.
- Genera un conjunto anidado de variables (i.e. solo aumentan y no cambian en cada iteración).
- Es un método competitivo cuando existen muchas variables.

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN
- 3 Ejemplo: Cáncer de próstata
- 4 Selección de variables
- 5 Métodos de contracción (shrinkage)**
- 6 Modelos de lineales de clasificación

## Métodos de contracción

- El mejor subconjunto puede tener el menor error de prueba. Pero es un método discreto. Las variables o se incluyen o se descartan.
- El método puede tener una varianza alta.
- Los métodos de contracción son más continuos y usualmente tienen menor varianza.
- Vamos a considerar la regresión de Ridge y Lasso.

## Métodos de contracción

- El mejor subconjunto puede tener el menor error de prueba. Pero es un método discreto. Las variables o se incluyen o se descartan.
- El método puede tener una varianza alta.
- Los métodos de contracción son más continuos y usualmente tienen menor varianza.
- Vamos a considerar la regresión de Ridge y Lasso.

## Métodos de contracción

- El mejor subconjunto puede tener el menor error de prueba. Pero es un método discreto. Las variables o se incluyen o se descartan.
- El método puede tener una varianza alta.
- Los métodos de contracción son más continuos y usualmente tienen menor varianza.
- Vamos a considerar la regresión de Ridge y Lasso.

## Métodos de contracción

- El mejor subconjunto puede tener el menor error de prueba. Pero es un método discreto. Las variables o se incluyen o se descartan.
- El método puede tener una varianza alta.
- Los métodos de contracción son más continuos y usualmente tienen menor varianza.
- Vamos a considerar la regresión de Ridge y Lasso.

- Resuelve el problema:

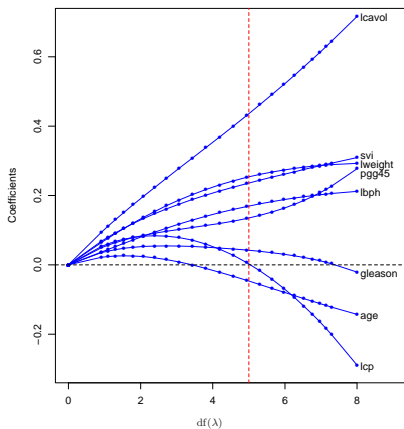
$$\min\left\{\sum_{i=1}^N (y_i - \beta_0 - x_i^T \beta)^2 + \lambda \|\beta\|^2\right\} \quad (3)$$

# Métodos de contracción: Ridge

**TABLE 3.3.** *Estimated coefficients and test error results, for different subset and shrinkage methods applied to the prostate data. The blank entries correspond to variables omitted.*

Term	LS	Best Subset	Ridge	Lasso	PCR	PLS
Intercept	2.465	2.477	2.452	2.468	2.497	2.452
lcavol	0.680	0.740	0.420	0.533	0.543	0.419
lweight	0.263	0.316	0.238	0.169	0.289	0.344
age	-0.141		-0.046		-0.152	-0.026
lbph	0.210		0.162	0.002	0.214	0.220
svi	0.305		0.227	0.094	0.315	0.243
lcp	-0.288		0.000		-0.051	0.079
gleason	-0.021		0.040		0.232	0.011
pgg45	0.267		0.133		-0.056	0.084
Test Error	0.521	0.492	0.492	0.479	0.449	0.528
Std Error	0.179	0.143	0.165	0.164	0.105	0.152

# Métodos de contracción: Ridge



**FIGURE 3.8.** Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter  $\lambda$  is varied. Coefficients are plotted versus  $df(\lambda)$ , the effective degrees of freedom. A vertical line is drawn at  $df = 5.0$ , the value chosen by cross-validation.

# Contenido

- 1 Teoría de la decisión
- 2 Regresión lineal y k-NN
- 3 Ejemplo: Cáncer de próstata
- 4 Selección de variables
- 5 Métodos de contracción (shrinkage)
- 6 Modelos de lineales de clasificación**

## Regresión lineal

- Se ajusta una regresión lineal por cada clase

$$\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

- Un punto  $x$  se clasifica  $k$  si  $\hat{f}_k(x) \geq \hat{f}_l(x) \forall l$ .
- De esta forma la frontera de decisión entre la clase  $k$  y  $l$  sería cuando  $\hat{f}_k(x) = \hat{f}_l(x)$  que corresponde a un plano.
- Esto pasaría para cualquier 2 clases.
- El espacio se divide en regiones separadas por fronteras de decisión hiperplanas (no necesariamente pasan por el origen).

## Regresión lineal

- Se ajusta una regresión lineal por cada clase

$$\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

- Un punto  $x$  se clasifica  $k$  si  $\hat{f}_k(x) \geq \hat{f}_l(x) \forall l$ .
- De esta forma la frontera de decisión entre la clase  $k$  y  $l$  sería cuando  $\hat{f}_k(x) = \hat{f}_l(x)$  que corresponde a un plano.
- Esto pasaría para cualquier 2 clases.
- El espacio se divide en regiones separadas por fronteras de decisión hiperplanas (no necesariamente pasan por el origen).

## Regresión lineal

- Se ajusta una regresión lineal por cada clase

$$\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

- Un punto  $x$  se clasifica  $k$  si  $\hat{f}_k(x) \geq \hat{f}_l(x) \forall l$ .
- De esta forma la frontera de decisión entre la clase  $k$  y  $l$  sería cuando  $\hat{f}_k(x) = \hat{f}_l(x)$  que corresponde a un plano.
- Esto pasaría para cualquier 2 clases.
- El espacio se divide en regiones separadas por fronteras de decisión hiperplanas (no necesariamente pasan por el origen).

## Regresión lineal

- Se ajusta una regresión lineal por cada clase

$$\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

- Un punto  $x$  se clasifica  $k$  si  $\hat{f}_k(x) \geq \hat{f}_l(x) \forall l$ .
- De esta forma la frontera de decisión entre la clase  $k$  y  $l$  sería cuando  $\hat{f}_k(x) = \hat{f}_l(x)$  que corresponde a un plano.
- Esto pasaría para cualquier 2 clases.
- El espacio se divide en regiones separadas por fronteras de decisión hiperplanas (no necesariamente pasan por el origen).

## Regresión lineal

- Se ajusta una regresión lineal por cada clase

$$\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$$

- Un punto  $x$  se clasifica  $k$  si  $\hat{f}_k(x) \geq \hat{f}_l(x) \forall l$ .
- De esta forma la frontera de decisión entre la clase  $k$  y  $l$  sería cuando  $\hat{f}_k(x) = \hat{f}_l(x)$  que corresponde a un plano.
- Esto pasaría para cualquier 2 clases.
- El espacio se divide en regiones separadas por fronteras de decisión hiperplanas (no necesariamente pasan por el origen).

## Regresión logística

- Este método de clasificación es lineal porque las fronteras de decisión son lineales. Luego transformaciones monótonas de  $\hat{f}_k(x)$  también van a definir modelos de clasificación lineales. Por ejemplo: el modelo logístico.

$$Pr(G = 1|X = x) = \frac{e^{\beta_{1,0} + \beta_1^T x}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

$$Pr(G = k|X = x) = \frac{e^{\beta_{k,0} + \beta_k^T x}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

$$Pr(G = K|X = x) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

## Regresión logística

- Este método de clasificación es lineal porque las fronteras de decisión son lineales. Luego transformaciones monótonas de  $\hat{f}_k(x)$  también van a definir modelos de clasificación lineales. Por ejemplo: el modelo logístico.

$$Pr(G = 1|X = x) = \frac{e^{\beta_{1,0} + \beta_1^T x}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

$$Pr(G = k|X = x) = \frac{e^{\beta_{k,0} + \beta_k^T x}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

$$Pr(G = K|X = x) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_l^T x}}$$

## Regresión logística

- Las fronteras de decisión se pueden calcular cómo:

$$\log\left(\frac{\Pr(G = 1|X = x)}{\Pr(G = K|X = x)}\right) = \beta_{1,0} + \beta_1^T x$$
$$\log\left(\frac{\Pr(G = k-1|X = x)}{\Pr(G = K|X = x)}\right) = \beta_{k-1,0} + \beta_{k-1}^T x$$