

# 1. MÁQUINAS DE APRENDIZAJE

Héctor Allende

Universidad Técnica Federico Santa María

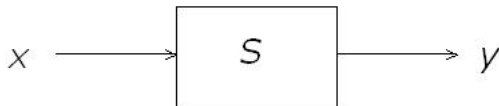
Febrero 2006

# AGENDA

- 1 MODELO ESTADÍSTICO DEL PROBLEMA
- 2 GENERALIZACIÓN Y TEORÍA VC
- 3 REGULARIZACIÓN
- 4 TEMAS PENDIENTES

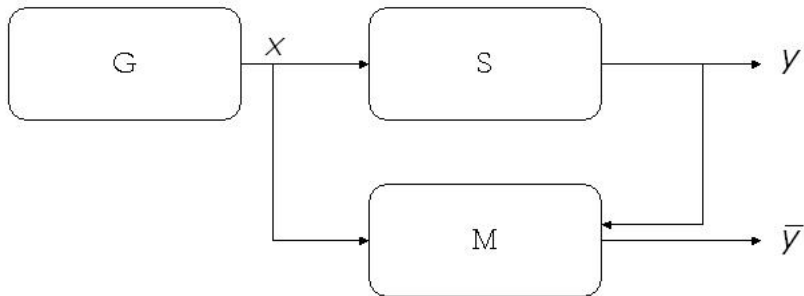
# PROBLEMA GENERAL

- **Materia Prima.** Datos, observaciones, mediciones.
- **Objetivo.** Obtener descripciones de alto nivel de esos datos, relaciones, modelos, patrones.
- Problema de Inducción. ¿Cuál es el alcance de estas descripciones? ¿Son válidas más allá de los datos que permitieron construirlas?
- Supuesto: Existe una regularidad subyacente a las observaciones: Sistema Generador.



# PROBLEMA GENERAL

- Identificar al Sistema: White-Box Models
- Imitar al Sistema: Black-Box Models



- ¿Identificar o imitar al sistema?

# MODELO ESTADÍSTICO DEL PROBLEMA

- Espacio de Observaciones  $Z = X \times Y$  ( ... con una medida  $P$  que implementa el supuesto de “regularidad” subyacente a las observaciones)
- Espacio de Hipótesis  $H$ : Colección de todos los modelos seleccionables para explicar los datos
- Función de Pérdida  $Q(f(x), y)$ : ¿Cuál es el costo de responder  $f(x)$  a una entrada  $x$  si el sistema responde con  $y$ ?

# MODELO ESTADÍSTICO DEL PROBLEMA

- **Aprendizaje**: Elegir  $F \in \Omega$  para minimizar el **Riesgo** asociado al modelo

$$E [Q(F(x), y)] = \int Q(F(x), y) dP(x, y) \quad (1)$$

No conocemos  $P$ !!, ¿Cómo elegir una hipótesis de  $H$  si no podemos computar (1)?

- No conocemos  $P$ , pero sí un conjunto de ejemplos  $D = \{(x_1, y_1), \dots, (x_m, y_m)\}$  que supondremos obtenidos i.i.d. de  $P$ .
- Funcional de Inducción: Criterio para elegir  $f$  sólo en base a la muestra  $D$ .

$$\begin{aligned}\hat{R} : H \times Z^n &\rightarrow \mathbb{R} \\ (f, D) &\mapsto R(f, D)\end{aligned}\tag{2}$$

- Elección Clásica: Funcional de Riesgo Empírico

$$\hat{R}^m(f, D) = \sum_{i=1}^m Q(f(x_i), y_i)\tag{3}$$

# PRINCIPIOS DE INDUCCIÓN

- Se obtiene al reemplazar  $P$  por

$$\hat{P} = \sum_{i=1}^m \delta(x - x_i) \delta(y - y_i) \quad (4)$$

- Alternativas a este principio básico se pueden obtener al considerar otros estimadores de  $P$ , por ejemplo

$$\hat{P} = \sum_{i=1}^m k(x - x_i, y - y_i) \quad (5)$$

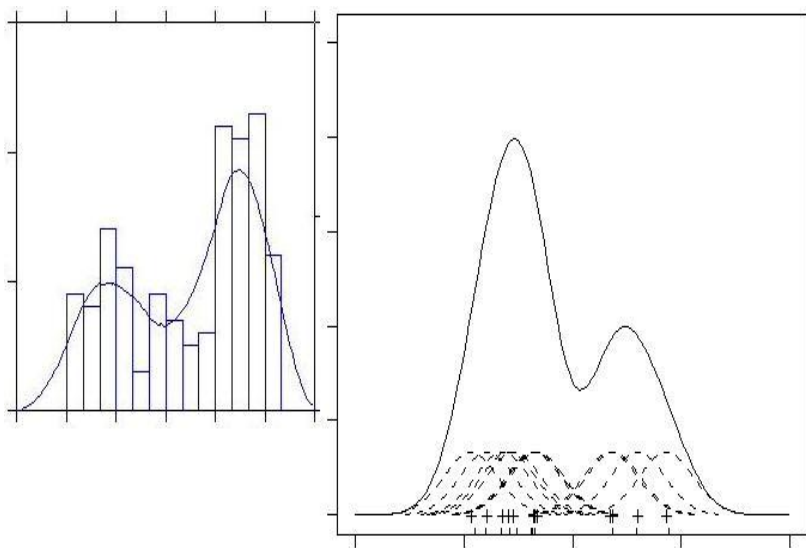
con  $k(\cdot, \cdot)$  un kernel centrado en  $(0, 0)$ .

- Obtenemos el denominado *Vicinal Risk*

$$\hat{R}(f, D) = \sum_{i=1}^m Q(f(x_i), y_i) k(x - x_i, y - y_i) \quad (6)$$



# PRINCIPIOS DE INDUCCIÓN



- ¿Qué principio de inducción es mejor?
- Antes de responder esta pregunta veamos como el modelo estadístico del aprendizaje puede servir para responder a problemas clásicos en estadística.
- **Regresión**: La relación entre  $X$  e  $Y$  es estocástica, i.e., dado un  $x$  existe un conjunto probable de respuestas  $y$ . Nos interesa estimar

$$r(x) = \int y dP(y|x) \quad (7)$$

denominada *función de regresión*

- Consideremos el funcional de riesgo asociado a la función de pérdida cuadrática,

$$R(f) = \int (y - f(x))^2 dP(x, y) \quad (8)$$

- Entonces es posible mostrar que

$$R(f) = \int (y - r(x))^2 dP(x, y) + \int (f(x) - r(x))^2 dP(x)$$

- La primera parte cuantifica la varianza de  $y$  y la segunda la diferencia entre el modelo  $f$  y la curva de regresión  $r(x)$ .

- De esta forma, si

$$\int y^2 dP(x, y) < \infty \quad \int r^2(x) dP(x, y) < \infty$$

el mínimo de  $R(f)$  es la curva de regresión cuando  $r(x) \in H$ .

- Si  $r(x) \notin H$ , entonces  $f$  es función de  $H$  que resulta más cercana a  $r(x)$  bajo la métrica  $L_2(P)$ .
- Además, si  $\tilde{f}$  es tal que  $R(\tilde{f}) - R(f) < \epsilon$ , entonces la distancia entre  $f$  y  $\tilde{f}$  es menor que  $\sqrt{\epsilon}$ .
- De esta forma, el problema de regresión (clásico) es un problema de aprendizaje con la función de pérdida cuadrática.

- Consideremos ahora el problema de distinguir entre un conjunto de categorías  $w_1, w_2, \dots, w_p$ .
- Dado un conjunto de características  $x$ , supongamos que nuestro modelo  $f$  clasifica a  $x$  en la clase  $w_i$ .
- La probabilidad de error es entonces

$$e(f|x) = 1 - P(y = w_i|x) = \sum_{k \neq i} P(y = w_k|x) \quad (9)$$

- Definamos la función

$$L(y, f(x)) = \begin{cases} 0 & y = f(x) \\ 1 & y \neq f(x) \end{cases} \quad (10)$$

- Entonces la suma anterior es equivalente a

$$e(f|x) = \sum_k L(y, f(x))P(y|x) \quad (11)$$

- El error de clasificación esperado se obtiene integrando en el espacio de las características,

$$\begin{aligned} e(f) &= \int_X \sum_k L(y, \phi) P(y|x) P(x) dx \\ &= \int_X \sum_k L(y, \phi) P(y, x) dx \\ &= \int_{X \times Y} L(y, \phi) dP(x, y) \end{aligned}$$

- Esta expresión es exactamente el riesgo asociado a la función de pérdida  $Q = L$  denominada *Misclassification Loss Function*.

- Consideremos ahora el problema de estimar una función de densidad de probabilidad  $p_0$  desde un conjunto  $H$ .
- La *entropía* para una función  $p \in H$  se define como

$$R(p) = - \int p_0(x) \log p(x) dx = - \int \log p(x) dP(x) \quad (12)$$

Es decir, el riesgo asociado a la función de pérdida

$$Q(p) = - \log p(x) \quad (13)$$

- El mínimo del funcional anterior es el mismo que obtendríamos si sumamos a  $R$  una constante como  $R(p_0)$ ,

$$R(p) = - \int p_0(x) (\log p(x) - \log p_0(x)) dx = - \int p_0(x) \log \frac{p(x)}{p_0(x)} dx$$

- La minimización del funcional anterior es la base para la estimación de densidades en el approach estadístico clásico.
- Si aplicamos la desigualdad de Jensen y un poco de álgebra obtenemos la siguiente desigualdad para  $R(p)$

$$R(p) \leq \log \left( 1 - \left( \frac{1}{2} \int |p(x) - p_0(x)| dx \right)^2 \right)$$

- Esta inecuación nos permite ver que el mínimo del funcional de riesgo se alcanza cuando  $p(x)$  es la distribución verdadera  $p_0(x)$ .



# GENERALIZACIÓN

- Hemos visto que problemas clásicos pueden analizarse como problemas de minimización del funcional de riesgo.
  - Como este funcional no puede ser computado en la práctica necesitamos un funcional empírico basado en los datos.
  - Sin embargo necesitamos que el funcional permita elegir funciones que se comporten bien más allá de los ejemplos particulares que se utilizan para evaluarlo.
- 
- *El funcional debe permitir generalizar*, i.e. el riesgo alcanzado con la hipótesis seleccionada debiera ser similar al riesgo mínimo que se puede obtener
  - *El valor del funcional debe ser indicativo del riesgo real*, i.e., el valor del funcional empírico obtenido debiera ser también similar al mínimo que se puede obtener.

# GENERALIZACIÓN Y CONSISTENCIA

- Para formalizar esta idea, sea  $f_m$  la hipótesis que minimiza  $\hat{R}$  en una muestra de tamaño  $m$  y  $f_0$  la hipótesis que minimiza el riesgo real  $R$ .

## CONSISTENCIA (DÉBIL)

Diremos que el principio de inducción  $\tilde{R} : \Omega \rightarrow \mathbb{R}$  es *consistente* si y sólo si se satisfacen las siguientes condiciones

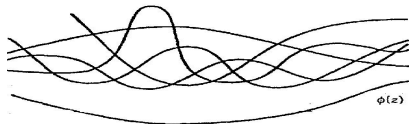
$$R(f_m) \xrightarrow[m]{P} R(f_0) \quad (14)$$

$$\tilde{R}(f_m) \xrightarrow[m]{P} R(f_0) \quad (15)$$

- Hablaremos de “consistencia universal” si  $\hat{R}$  es consistente independiente de la distribución  $P(x, y)$  de los ejemplos.

# GENERALIZACIÓN Y CONSISTENCIA

- Supongamos que en el espacio  $\Omega$  existe una función  $\phi$  tal que  $\forall(x, y), Q(y, \phi(x)) < Q(y, f(x)), \forall f \in \Omega$ .
- Es evidente que cualquier principio de inducción  $\tilde{R}$  basado en cualquier muestra  $D$  encontrará  $\phi(x)$  como solución al problema.
- Esto nos indica que la noción de consistencia es demasiado débil y necesita un poco de precisión.



# CONSISTENCIA FUERTE

## CONSISTENCIA (ESTRICTA)

Diremos que un principio de inducción  $\tilde{R}$  es estrictamente consistente en  $H$  si en cualquier  $\lambda(c) \neq \emptyset$

$$\lambda(c) := \{g \in H : R(g) \geq c\} \quad (16)$$

se satisface

$$\tilde{R}(f_m^c) \xrightarrow[m]{P} R(f_0^c) \quad (17)$$

donde  $f_m^c$  es la función que minimiza  $\tilde{R}$  en el subconjunto  $\lambda(c)$  y  $f_0^c$  es la función que minimiza  $R$  en el mismo subconjunto.

- Notemos que claramente la consistencia estricta implica la consistencia débil

# TEOREMA CLAVE DE LA TEORÍA VC

- Uno de los resultados claves de la *Teoría Estadística del Aprendizaje* es el que relaciona la consistencia del principio del riesgo empírico con el estudio de un proceso empírico del “peor caso”.

## PROCESO EMPÍRICO DE UNA COLA

Se define como la secuencia  $(\xi_+^m)_m$ ,  $m \in \mathbb{N}$  donde

$$\xi_+^m = \sup_{f \in H} \left( R(f) - \hat{R}^m(f) \right)$$

$$\text{y } \hat{R}^m = 1/m \sum_{i=1}^m Q(f(x_i), y_i)$$

- Hablamos del “peor caso” pues para cada tamaño muestral  $m$  tomamos el supremo de las diferencias entre  $R$  y  $\hat{R}^m$ .

# TEOREMA CLAVE DE LA TEORÍA VC

Y ahora el resultado clave ...

## TEOREMA DE EQUIVALENCIA (VAPNIK, CHERVONENKIS)

Sea  $R(f)$  absolutamente acotado en  $H$ . Entonces las siguientes condiciones son equivalentes:

- 1  $\hat{R}^m$  es estrictamente consistente en  $H$ .
- 2 El proceso empírico de una cola converge a cero en probabilidad, es decir

$$P(\xi_+^m > \epsilon) \xrightarrow{m \rightarrow \infty} 0 \quad (18)$$

- Si una de las condiciones es válida independiente de la distribución  $P(x, y)$  de los ejemplos, la otra también es universal.

# TEOREMA CLAVE DE LA TEORÍA VC

- ¿Bajo que condiciones el proceso de una cola converge a cero?. Consideremos un caso simple en que  $Q$  toma valores en  $\{0, 1\}$  (clasificación) y  $H$  es finito.
- La desigualdad de Chernoff puede aplicarse sobre cada función de  $H$  para concluir que

$$\begin{aligned} P(\xi_+^m > \epsilon) &\leq 2Ne^{-2\epsilon^2 m} \\ &\leq 2e^{(\frac{\ln N}{I} - 2\epsilon^2)m} \end{aligned}$$

- Esto sugiere que si

$$\frac{\ln N}{m} \xrightarrow{m \rightarrow \infty} 0$$

el proceso de una cola convergerá a cero en probabilidad.

# CAPACIDAD Y DILEMA BIAS-VARIANZA

- El resultado anterior sugiere que el tamaño del espacio de soluciones  $H$  debe mantenerse bajo control para obtener consistencia y una buena generalización.
- Conceptualmente esta idea es la misma que expresa el clásico dilema Sesgo-Varianza en estadística.
- El Sesgo de un estimador  $f_m$  que se obtiene de una muestra de ejemplos  $D$  se define como

$$Bias^2(f_m) = E_P (E_T [f_m(x)] - f_0(x))^2$$

donde  $T$  es la distribución conjunta de los ejemplos.

- La varianza del estimador se define como

$$Var(f_m) = E_P E_T (f_m(x) - E_T [f_m(x)])^2$$



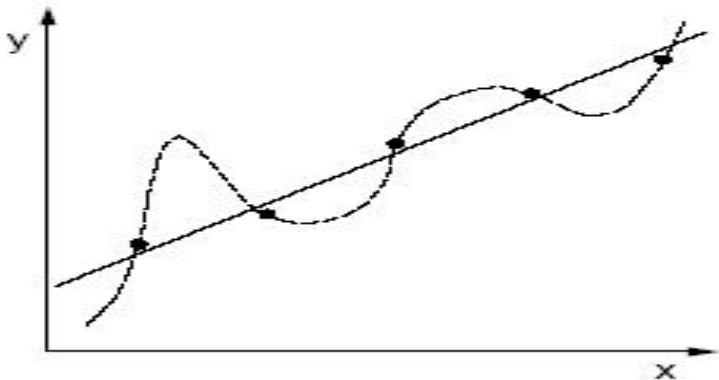
# CAPACIDAD Y DILEMA BIAS-VARIANZA

- A medida que aumentamos el tamaño del espacio de soluciones es más probable encontrar una función que modele perfectamente el conjunto de ejemplos.
- Por otro lado, si tenemos muchas soluciones, pequeños cambios en el conjunto de ejemplos darán origen a una solución diferente. Es decir, la varianza esperada tiende a aumentar si el espacio se hace más grande.
- Intuitivamente estimadores inestables tienden a no generalizar adecuadamente i.e. sobreajustan.

Al modificar la complejidad estructural del espacio de soluciones existe un compromiso entre sesgo y varianza ó bien entre ajuste y generalización.

# CAPACIDAD Y DILEMA BIAS-VARIANZA

Consideremos el siguiente ejemplo:



- ¿Qué aproximación es preferible?

- Un ejemplo de control de capacidad se da en la interpolación con splines cúbicos. En este caso se favorecen los modelos “simples” utilizando un término inversamente proporcional a la segunda derivada del estimador,

$$\hat{R}(f) = \sum_{i=1}^m (y_i - f(x_i))^2 + \lambda \int [f''(x)]^2 dx$$

- Otro ejemplo clásico son las medidas para selección de modelos denominadas *BIC* ó *AIK* (Criterio de Akaike) que son inversamente proporcionales a la varianza muestral.
- El tema es cómo medir adecuadamente la “capacidad” del espacio de soluciones. Claramente la cardinalidad no es suficiente pues en general trabajaremos con espacios infinitos.

- El funcional de riesgo empírico genérico no toma en cuenta ninguna medida estructural. Por lo tanto su capacidad para generalizar debe depender de características del espacio de hipótesis  $H$ .
- El tema es cómo medir adecuadamente la “capacidad” del espacio de soluciones. Claramente la cardinalidad no es suficiente pues en general trabajaremos con espacios infinito-dimensionales.
- **Observación Clave:** No interesa directamente el **tamaño** del espacio sino las diferentes configuraciones que se alcanzan.
- Dos funciones que siempre reportan el mismo riesgo son idénticas aunque paramétricamente no lo sean.

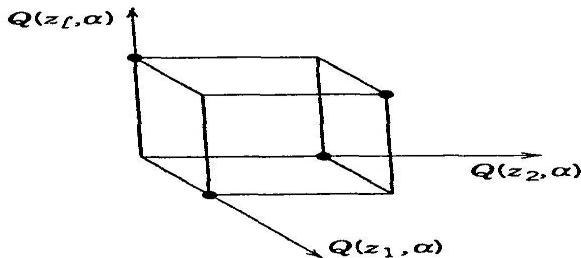
- La observación anterior, motiva la definición de una medida de capacidad que tome en cuenta los valores que puede tomar la función  $Q(f(x), y)$  en el espacio  $H$ .
- Dado un conjunto de ejemplos  $D_m = \{(x_i, y_i), i = 1, \dots, m\}$  y una función  $f \in H$  definamos,

$$q(f) = (Q(f(x_1), y_1), \dots, Q(f(x_m), y_m))$$

- De esta forma,  $q(H) = \{q(f), f \in H\}$ .
- En el caso más simple en que  $Q$  toma valores en  $\{0, 1\}$ ,  $\#q(H)$  es finita ( $\leq 2^m$ ) y parece una medida apropiada de capacidad.

# MEDIDAS DE CAPACIDAD

- Gráficamente  $\#q(H)$  es el número de vértices diferentes que se inducen sobre el hipercubo  $q(H)$  que se forma tomando en cada eje un ejemplo  $z_i = (x_i, y_i)$ .



- Sobre  $\#q(H)$  podemos introducir las siguientes definiciones de capacidad.

## ENTROPÍA DE UN ESPACIO DE HIPÓTESIS

Sea  $P(x, y)$  la distribución del conjunto de ejemplos  $D$ . La *entropía aleatoria* de un espacio de hipótesis  $H$  se define como

$$\mathcal{E}^H(D) = \ln(\#q(H)) \quad (19)$$

Se define además la *Entropía* de  $H$

$$\mathcal{E}^H(m) = E_P(\mathcal{E}^H(D)) \quad (20)$$

## TEOREMA

Para que el proceso de una cola asociado a un funcional de inducción  $\hat{R}$  converja en probabilidad a cero es suficiente que se verifique

$$\frac{\mathcal{E}^H(m)}{m} \xrightarrow{m \rightarrow \infty} 0 \quad (21)$$

- Este teorema relaciona la habilidad de generalización de una máquina de aprendizaje -que implemente el riesgo empírico- con una medida específica de capacidad.



# GENERALIZACIÓN DE LAS MEDIDAS

- La búsqueda de condiciones que sean además de suficientes sean necesarias y que sean válidas para funciones  $Q$  con valores en  $\mathbb{R}$  nos llevan a explorar otros conceptos de capacidad.

## $\epsilon$ -NET Y NÚMERO DE CUBRIMIENTO

Sea  $A$  un conjunto dotado de una métrica  $\rho$ . Una  $\epsilon$ -net de  $A$  es un conjunto  $\Lambda$  tal que  $\forall f \in A$ , existe  $\tilde{f} \in \Lambda$  que satisface

$$\rho(f, \tilde{f}) < \epsilon$$

Denominamos  $\epsilon$ -net minimal a aquella  $\epsilon$ -net  $\tilde{\Lambda}$  tal que  $\tilde{\Lambda} \subset \Lambda$  para toda  $\epsilon$ -net  $\Lambda$ . Llamaremos **número de cubrimiento** (covering number)  $\mathcal{N}(\epsilon, A, \rho)$  a la cardinalidad de este conjunto minimal.

# $\epsilon$ -ENTROPÍA

- Consideremos el conjunto  $q(H)$  formado por los vectores  $q(f) = (Q(f(x_1), y_1), \dots, Q(f(x_m), y_m)), \forall f \in H$ .

## $\epsilon$ -ENTROPÍA

Sea  $P(x, y)$  la distribución del conjunto de ejemplos  $D$ . La  $\epsilon$ -entropía aleatoria de un espacio de hipótesis  $H$  se define como

$$\mathcal{E}^H(\epsilon, D) = \ln(\mathcal{N}(\epsilon, q(H), l^\infty)) \quad (22)$$

donde  $l^\infty$  es la métrica del supremo. Se define además la  $\epsilon$ -entropía de  $H$

$$\mathcal{E}^H(\epsilon, m) = E_P(\mathcal{E}^H(\epsilon, D)) \quad (23)$$

- El siguiente teorema cierra el problema de consistencia del riesgo empírico utilizando las medidas introducidas.

## TEOREMA DE CONSISTENCIA DEL RIESGO EMPÍRICO

Para que obtener convergencia en probabilidad del proceso de una cola asociado a un funcional de inducción  $\hat{R}$  es suficiente y necesario que para todo  $\epsilon, \delta, \sigma$  exista un  $\sigma$ -cubrimiento  $\Theta$  de  $H$  tal que

$$\lim_{m \rightarrow \infty} \frac{\mathcal{E}^\Theta(\epsilon, m)}{m} < \delta \quad (24)$$

$\Theta$  es un  $\sigma$ -cubrimiento de  $H$  si  $\forall f \in H \exists \tilde{f} \in \Theta$  tal que

$$\begin{aligned} Q(f(x), y) &\geq Q(\tilde{f}(x), y) \\ R(f) - R(\tilde{f}) &< \sigma \end{aligned}$$

# CONSISTENCIA UNIVERSALMENTE VÁLIDA

- La noción de entropía depende de la distribución de probabilidad  $P(x, y)$  asociada a los ejemplos en  $D$ .
- Para asegurar la consistencia del riesgo empírico de manera independiente de la medida de probabilidad es necesario introducir una cota para la entropía, que denominaremos *función de crecimiento* (*growth function*) del espacio de hipótesis  $H$ ,

$$\mathcal{G}(\epsilon, H, m) = \ln \sup_D \mathcal{N}(\epsilon, q(H), l^\infty)$$

- Para hacer universales los teoremas de consistencia es suficiente reemplazar la entropía  $\mathcal{E}^H(\epsilon, m)$  por  $\mathcal{G}(\epsilon, H, m)$ .

## TEOREMA

Supongamos que  $q(H)$  esta formada por funciones que toman valores en  $\{0, 1\}$ . Entonces

$$\mathcal{G}(H, m) = m \ln m$$

ó bien existe un  $h \in \mathbb{N}$  tal que

$$\mathcal{G}(H, m) \begin{cases} = m \ln m & m \leq h \\ \leq h \left(1 + \ln \frac{m}{h}\right) & m > h \end{cases}$$

- Si ocurre el primer caso diremos que  $h$  es infinita.
- El número  $h$  se conoce como **Dimensión VC**

- Recordemos que el principio del riesgo empírico es universalmente consistente en el conjunto de indicatrices  $H$  si se verifica

$$\lim_{m \rightarrow \infty} \frac{\mathcal{G}(H, m)}{m} = 0$$

- Notemos que si  $H$  tiene dimensión VC finita  $h$ , entonces tenemos,

$$\lim_{m \rightarrow \infty} \frac{\mathcal{G}(H, m)}{m} < \lim_{m \rightarrow \infty} \frac{h \left(1 + \ln \frac{m}{h}\right)}{m} = 0$$

Si el espacio de hipótesis tiene dimensión VC finita se verifica universalmente la consistencia del riesgo empírico. Sin embargo no se trata de una condición necesaria.

## TEOREMA: COTA ADITIVA PARA EL RIESGO EMPÍRICO

Supongamos que la dimensión VC  $h$  asociada a un espacio de hipótesis  $H$  es finita. Entonces, independiente de la distribución de los ejemplos,

$$P\left(\sup_{f \in H} |R(f) - \hat{R}^m(f)| > \epsilon\right) < 4 \exp\left(\frac{h(1 + \ln(2m/h))}{m} - (\epsilon - 1/m)^2\right) m$$

Equivalente podemos decir, que con probabilidad  $1 - \eta$

$$R(f_m) < \hat{R}^m(f_m) + \sqrt{\mathcal{L}(m)} + \frac{1}{m}$$

donde  $f_m$  es el mínimo del riesgo empírico  $\hat{R}^m$

$$\mathcal{L}(m) = 4 \frac{h(\ln(2m/h) + 1) - \ln \eta/4}{m}$$

## TEOREMA: COTA RELATIVA PARA EL RIESGO EMPÍRICO

Supongamos que la dimensión VC  $h$  asociada a un espacio de hipótesis  $H$  es finita. Entonces,

$$P\left(\sup_{f \in H} \frac{R(f) - \hat{R}^m(f)}{\sqrt{R(f)}} > \epsilon\right) < 4 \exp\left(\frac{h(1 + \ln(2m/h))}{m} - \frac{\epsilon^2}{4}\right) m$$

Equivalente podemos decir, que con probabilidad  $1 - \eta$

$$R(f_m) < \hat{R}^m(f_m) + \frac{\mathcal{L}(m)}{2} \left(1 + \sqrt{1 + \frac{4\hat{R}^m(f_m)}{\mathcal{L}(m)}}\right)$$

donde  $f_m$  es el mínimo del riesgo empírico  $\hat{R}^m$ .



# DIMENSIÓN VC PARA FUNCIONES REALES

- La dimensión VC para el caso de un espacio  $q(H)$  toma valores reales se puede definir introduciendo un conjunto auxiliar  $\Theta(H)$  de funciones que tomen valores en  $\{0, 1\}$ .

$$\Theta(H) = \{\theta(Q(f(x), y) - \beta), f \in H, \beta \in \mathbb{R}\}$$

donde  $\theta$  es una indicatriz centrada en 0.

- La función de crecimiento de  $H$  se define en este caso como en el caso de funciones valuadas en  $\{0, 1\}$  pero sobre el conjunto  $\Theta(H)$ . Lo mismo con la dimensión VC.
- Este truco permite extender todos los teoremas al caso general de funciones  $Q$  reales.

# OBSERVACIONES ACERCA DE LA TEORÍA VC

- Si bien el truco anterior permite obtener resultados elegantes para explicar las condiciones en que el riesgo empírico es consistente, éste no incorpora ninguna información acerca de la escala de las observaciones.
- En el caso de clasificación dicha información de escala no es importante pues los valores que toman las hipótesis  $f(x)$  son sólo una referencia para distinguir entre clases.
- En el caso de regresión sin embargo la pérdida de la información de escala hace que los resultados de la teoría VC no sean útiles en la práctica. Por ejemplo, las cotas para el riesgo empírico suelen ser muy poco *tight* como para ser informativas.

# REGULARIZACIÓN Y CONTROL DE CAPACIDAD

- Los resultados de la teoría VC afirman la idea de que una máquina de aprendizaje que generaliza, debe mantener controlada la capacidad del espacio de soluciones  $H$ .

Las cotas para el riesgo obtenidas de la teoría VC tienen la forma genérica

$$R(f) \leq \hat{R}^m(f) + \Lambda(H)$$

donde  $\hat{R}^m$  es el riesgo empírico y  $\Lambda(H)$  es una medida de la complejidad estructural del espacio de soluciones.

- $\Lambda(H)$  es una medida global en el espacio de soluciones. No depende de una hipótesis candidata en particular. Más aún es independiente de los datos.
- El desarrollo de un principio adaptivo de aprendizaje a partir de estos resultados es difícil.

- La idea clave en regularización es la restricción del espacio de soluciones mediante la incorporación de un funcional  $\Omega(f)$  que penalice hipótesis complejas.
- El funcional de riesgo empírico regularizado tiene la forma

$$\tilde{R}^{reg}(f) = \hat{R}^m(f) + \lambda \Omega(f) \quad (25)$$

donde  $\Omega(\cdot)$  se denomina *regularizador* y el correspondiente  $\lambda$  *parámetro de regularización*.

- El parámetro  $\lambda$  controla el tradeoff entre ajuste a los datos ( $\hat{R}^m(f)$ ) y complejidad estructural ( $\Omega(f)$ ).
- Hipótesis complejas son admisibles siempre que los datos lo requieran.

- Un ejemplo clásico (en estadística bayesiana) consiste en considerar un regularizador proporcional a la norma de la primera derivada

$$\hat{R}^{reg}(f) = \hat{R}^m(f) + \lambda \|\delta_x f\|^2$$

- De entre dos funciones con el mismo riesgo empírico, se prefiere la que sea más suave.
- Otro ejemplo se da en redes neuronales (*weight decay neural networks*)

$$\hat{R}^{reg}(f) = \hat{R}^m(f) + \lambda \|w\|^2$$

- En este caso  $w$  es el vector de pesos de la red.  $\|w\|^2$  intenta restringir la norma de la hipótesis  $\|f\|^2$ .

# PROBLEMAS INVERSOS

- Consideremos dos espacios de funciones  $M$  y  $N$  y un operador  $A : M \rightarrow N$  que asigna a cada  $f \in M$  un único  $F \in N$ .
- Supongamos que queremos estimar  $f$ , pero sólo podemos observar  $F$ . Si conocemos  $A$ , el problema se reduce a resolver la ecuación,

$$Af = F \quad (26)$$

- Denominaremos este tipo de problemas **Problemas Inversos**.
- Por ejemplo, la estimación de densidades  $f(t)$  a partir de un conjunto de ejemplos es un problema inverso

$$\int_{-\infty}^x f(t)dt = F(x)$$

donde  $F(x)$  es la distribución de los ejemplos.

# PROBLEMAS INVERSOS

- Para que el proceso de estimación tenga sentido, nos interesa que dos soluciones próximas en el espacio observable  $N$  tengan preimágenes cercanas en el espacio de las soluciones  $M$ .
- Esta idea se captura técnicamente en la noción de continuidad la cual requiere que  $M$  y  $N$  sean espacios métricos. El operador  $A$  es continuo

## CONTINUIDAD DE OPERADORES

Sean  $(X, \rho_1)$  y  $(Y, \rho_2)$  dos espacios métricos y  $C : X \rightarrow Y$  un operador.  $C$  se dice continuo si y sólo si  $\forall \epsilon, \forall x \in A \exists \sigma > 0$  tal que

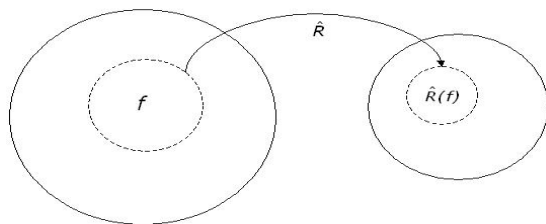
$$CB_{\rho_1}(x, \sigma) \subset B_{\rho_2}(Ax, \epsilon)$$

donde  $B_{\rho_1}(x, \sigma)$  es la bola (abierta) de radio  $\sigma$  en la métrica  $\rho_1$  y análogamente  $B_{\rho_2}(Ax, \epsilon)$ .

- Claramente nos interesa la continuidad del operador inverso  $A^{-1}$ .

# APRENDIZAJE COMO PROBLEMA INVERSO

- El problema de estimar una función desde un conjunto de ejemplos puede verse como un problema inverso.
- El operador  $A$  es en este caso un funcional de inducción  $\hat{R}$  que mapea cada función de un espacio de hipótesis  $H$  a un valor único  $\hat{R}(f)$  en  $\mathbb{R}$ .



Nos interesa que si dos hipótesis tienen valores similares en el funcional de inducción, éstas estén próximas en el espacio de hipótesis pues esto garantiza que las funciones son similares más allá de los ejemplos particulares.



# APRENDIZAJE COMO PROBLEMA INVERSO

- Una visión alternativa se obtiene al considerar  $A$  como el operador que toma una función  $f$  del espacio de hipótesis y genera una muestra de ejemplos  $D(f) = \{(x_i, f(x_i)), i = 1, \dots, m)\}$ .
- En este caso  $X_m = \{x_i; i = 1, \dots, m\}$  es un parámetro del operador  $A$ , encargado de evaluar  $f$  en  $X_m$ .
- El problema es estimar  $f$  desde una muestra particular.

Nos interesa que dos conjuntos de ejemplos similares generen hipótesis similares. El operador  $A^{-1}$  que mapea un conjunto de ejemplos a una función en el espacio de soluciones es nuestro algoritmo. Nos interesa que el algoritmo sea estable.

# BUEN Y MAL CONDICIONAMIENTO

## BUEN CONDICIONAMIENTO EN EL SENTIDO DE HADAMARD

Un problema inverso  $Af = F$  con  $f \in M$  y  $F \in N$  se dice *bien condicionado en el sentido de Hadamard* (well-posed) ssi

- 1 Para todo  $F$ ,  $f$  existe y es única
- 2 El problema es estable, i.e.,  $A^{-1}$  es continuo

## BUEN CONDICIONAMIENTO EN EL SENTIDO DE TIKHONOV

Un problema inverso  $Af = F$  con  $f \in M$  y  $F \in N$  se dice *bien condicionado en el sentido de Tikhonov* ssi  $\exists M' \subset M$  tal que

- 1 Para todo  $F \in N' = AM'$ ,  $f$  existe, es única y esta en el *espacio de correctitud*  $M'$
- 2 El problema es estable, i.e.,  $A^{-1}$  es continuo en  $M'$  y  $N'$ .

- Notemos que si el problema está bien condicionado, entonces para cualquier secuencia  $F_1, F_2, \dots, F_m$  convergente a  $F$  la secuencia de las preimágenes  $f_1, f_2, \dots, f_m$ ,  $f_i = A^{-1}F_i$  converge a  $f = A^{-1}F$ .
- Condiciones suficientes bajo las cuales tenemos un problema bien condicionado son las siguientes:

## TEOREMA CLAVE DE REGULARIZACIÓN

Sea  $M' \subset M$  un espacio métrico compacto y  $A : M \rightarrow N$  un operador continuo. Entonces  $A^{-1} : N' \rightarrow M'$  es continuo en  $N' = AM'$ .

# MÉTODO DE REGULARIZACIÓN

- Introducido por Tikhonov hacia 1963 para resolver el problema inverso  $Af_0 = F_0$ ,  $f_0 \in M$ ,  $F_0 \in N$ .
- La idea consiste en restringir las soluciones candidatas  $f$  a un subespacio compacto de  $M$ .
- Para ello se considera un funcional semicontinuo  $W : M \rightarrow \mathbb{R}$  denominado **regularizador**.

## REGULARIZADOR

$W : M \rightarrow \mathbb{R}$  se denomina regularizador si es semicontinuo y satisface las propiedades:

- 1  $W$  está definido para  $f_0$ .
- 2 En su dominio de definición,  $W(f) \geq 0$
- 3 Los conjuntos

$$M_c := \{f \in M : W(f) \leq c\}$$

son todos compactos.

- Para resolver el problema inverso, lo tradicional es minimizar en el espacio  $M$ , la norma de la diferencia

$$\tilde{R}(f, F_0) = \rho_2(Af - F_0)$$

donde  $\rho_2$  es la norma del espacio observable  $N$ .

- El problema es que si en vez de observar  $F_0$  observamos  $F_\delta$  muy cercano a  $F_0$ , digamos

$$\rho_2(F_\delta - F) \leq \delta$$

la solución  $f_\delta$  que minimiza  $\tilde{R}(f, F_\delta)$  no resulta cercana a  $f_0$ , a menos que el problema sea bien condicionado.

- Tikhonov introduce en cambio el funcional

$$\hat{R}_\lambda^{reg}(f, F_0) = \rho_2^2(Af - F_0) + \gamma W(f)$$

con  $\gamma > 0$ .

- Si minimizamos ahora,  $\hat{R}_\lambda^{reg}(f, F_\delta)$  ¿La solución que se obtiene  $f_\delta^\gamma$ , es cercana a  $f_0$ ?
- Más aún, supongamos que tenemos secuencias  $F_\delta$  convergentes a  $F_0$  ¿Las correspondientes secuencias  $f_\delta^\gamma$  son convergentes a  $f_0$ ?
- La respuesta dependerá de que tan cercano sea  $F_\delta$  a  $F_0$  y de cómo se elija el parámetro  $\gamma$ .

# MÉTODO DE REGULARIZACIÓN

## TEOREMA

Sean  $(M, \rho_1), (N, \rho_2)$  espacios métricos. Supongamos que en vez de  $F_0$  observamos aproximaciones  $F_\delta \in N$  tal que  $\rho_2(F_0 - F_\delta) \leq \delta$ . Supongamos que el parámetro de regularización es tal que

$$\gamma(\delta) \xrightarrow{\delta \rightarrow 0} 0$$

y además,

$$\lim_{\delta \rightarrow 0} \frac{\delta^2}{\gamma(\delta)} \leq r < \infty$$

Entonces, los elementos  $f_\delta^{\gamma(\delta)}$  que minimizan los funcionales  $\tilde{R}_{\gamma(\delta)}(f, F_\delta)$  convergen a la solución exacta  $f_0 = A^{-1}F_0$  a medida que  $\delta \rightarrow 0$ .

- La demostración del teorema anterior se basa (en parte) en la siguiente desigualdad

$$W(f_\delta^{\gamma(\delta)}) \geq W(f) + \frac{\delta^2}{\gamma(\delta)}$$

- Como los conjuntos  $W(f) \leq c$  son compactos  $\forall c$ , la minimización de  $R_\gamma^{reg}$  considerando un  $F_\delta$  cercano a  $F_0$  genera  $f_\delta^{\gamma(\delta)}$  cercanos a  $f_0$ .
- Las condiciones impuestas sobre  $\gamma(\delta)$  permiten concluir el resultado más fuerte de que secuencias de  $F_\delta$ 's convergentes en  $N$  generan secuencias  $f_\delta^{\gamma(\delta)}$ 's convergentes en  $M$ .



# MÉTODO DE REGULARIZACIÓN

- Supongamos que  $M$  es un espacio de Hilbert (e.g., en métodos de kernel). Una elección típica es en este caso

$$W^H(f) = \|f\|^2 = \langle f, f \rangle$$

- ¿Es  $W^H$  un regularizador válido?. La respuesta es negativa dado que los conjuntos  $M_c = \{f, W^H(f) \leq c\}$  son sólo débilmente compactos.
- Sin embargo las propiedades de un espacio de Hilbert permiten establecer

## TEOREMA

Sea  $M$  un espacio de Hilbert y  $W(f) = \langle f, f \rangle$ . Entonces si  $\gamma_0$  satisface las condiciones del teorema anterior con  $r = 0$ , los elementos  $f_\delta^{\gamma(\delta)}$  convergen a  $f_0$ .

# REGULARIZACIÓN Y MÉTODOS BAYESIANOS

- Desde el punto de vista bayesiano (McKay, 1996) el problema de aprendizaje se entiende como la estimación de la distribución *a-posteriori* de las hipótesis  $p(h|(X, Y))$  dado un conjunto de observaciones u ejemplos  $(X, Y) = \{(x_1, y_1), i = 1, \dots, m\}$ .
- Usando la regla de Bayes podemos encontrar que

$$p(h|X, Y) = \frac{p(X, Y|h)p(h)}{p(Y|X)}$$

- La característica central de la estimación bayesiana es que asume cierto conocimiento previo de los datos y del tipo de relaciones funcionales que podemos encontrar.
- Este conocimiento previo se recoge imponiendo una estructura particular a  $p(h)$  y a  $p(X, Y|h)$ .

# REGULARIZACIÓN Y MÉTODOS BAYESIANOS

- Desde el punto de vista bayesiano (McKay, 1996) el problema de aprendizaje se entiende como la estimación de la distribución *a-posteriori* de las hipótesis  $p(h|(X, Y))$  dado un conjunto de observaciones u ejemplos  $(X, Y) = \{(x_1, y_1), i = 1, \dots, m\}$ .
- Usando la regla de Bayes podemos encontrar que

$$p(h|X, Y) = \frac{p(X, Y|h)p(h)}{p(Y|X)}$$

- La característica central de la estimación bayesiana es que asume cierto conocimiento previo de los datos y del tipo de relaciones funcionales que podemos encontrar.

Este conocimiento previo se recoge imponiendo una estructura particular a  $p(h)$  denominado *a-priori* y a  $p(X, Y|h)$  denominada *verosimilitud*.

# REGULARIZACIÓN Y MÉTODOS BAYESIANOS

- El *a-posteriori* puede verse como una corrección del *a-priori* distribucional tomando en cuenta la información que nos entregan las observaciones.
- Supongamos que los datos se generan según

$$y_i = f(x_i) + \xi_i \Leftrightarrow y - f(x_i) = \xi_i$$

- En este caso, es razonable suponer que

$$p(x_i, y_i) = p(y_i - f(x_i)) = p(\xi_i)$$

- Si el conjunto de datos se genera de manera independiente e idénticamente distribuida, obtenemos que

$$p(X, Y|h) = \prod_{i=1}^m p(\xi_i)$$

- Dado el *a-posteriori*  $p(h|X, Y)$  y una determinada entrada  $x$  la predicción de  $y$  corresponde a estimar

$$p(y|X, Y, x) = \int_H p(y|h, x)p(h|X, Y)dh$$

- **Aproximación MAP (*Maximum a Posteriori*)**. Muchas veces aproximación de este tipo de integrales es intratable. Una aproximación razonable consiste en elegir la hipótesis de  $H$  que maximiza  $p(h|X, Y)$  y utilizar esta hipótesis para generar la salida  $y = f(x)$ .
- Maximizar  $p(h|X, Y)$  es equivalente a minimizar  $-\ln p(h|X, Y)$ , que a su vez resulta equivalente a minimizar

$$\tilde{R}^{bay}(h) = -\ln p(X, Y|h) - \ln p(h) = \sum_{i=1}^m -\ln p(y - h(x_i)) - \ln p(h)$$

- Notemos la similitud de la función objetivo anterior y el riesgo empírico regularizado

$$\tilde{R}^{bay}(h) \propto R_{emp}(h) + \lambda \Omega(h)$$

- Consideremos el siguiente *a-priori* gaussiano sobre  $H$  con parámetro de escala  $\sigma_h$

$$p(h) = \frac{1}{K_h} e^{-\frac{\|h\|^2}{\sigma_h}}$$

- Supongamos además que el ruido asociado al modelo es gaussiano con parámetro de escala  $\sigma_\xi$ ,

$$p(\xi) = \frac{1}{K_\xi} e^{-\frac{\xi^2}{\sigma_\xi}}$$

- Entonces  $\tilde{R}^{bay}(h)$  tiene la forma,

$$\begin{aligned}\tilde{R}^{bay}(h) &= \sum_{i=1}^m -\ln p(y - h(x_i)) - \ln p(h) \\ &= \sum_{i=1}^m \frac{1}{\sigma_\xi} (y_i - f(x_i))^2 + \frac{1}{\sigma_h} \|h\|^2 \\ &\propto \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 + \frac{\sigma_\xi}{m\sigma_h} \|h\|^2\end{aligned}$$

- Notemos que  $\tilde{R}^{bay}(h)$  coincide con el riesgo regularizado  $R_\gamma^{reg}$  si  $\gamma = \frac{\sigma_\xi}{m\sigma_h}$ .

- El *a-priori*  $p(h)$  sobre el espacio de las soluciones tiene el rol de regularizar la información acerca del error que aportan las observaciones.
- Seleccionar un regularizador es equivalente a seleccionar un *a-priori* sobre  $H$ .
- Esta “visión bayesiana” también nos previene de que el término de error depende del modelo de ruido que estemos asumiendo, quizá implícitamente.



- Teoría Algorítmica del Aprendizaje
- Generalización, Estabilidad y Robustez
- Funciones de Pérdida y Modelos de Ruido