

# Ajuste y generalización

Que aprenda y que ese aprendizaje sea robusto  
frente a situaciones nuevas

# Fuentes de error

¿Cómo lo medimos?

Impredecibilidad del fenómeno

3) Limitaciones de la red

Método de ajuste

# ¿Cómo medimos el error?

- Por propiedades teóricas, muchas veces el error cuadrático
  - Simetría de signo
  - Crecimiento parabólico: poca importancia de errores pequeños, crecimiento fuerte no acotado
- Otras posibilidades: absoluto (igual crecimiento desde el principio), Tukey (satura), ...
- En clasificación, suele usarse el error entrópico
- Por conveniencia del algoritmo de ajuste: exponenciales
- Por conveniencia de la robustez: añadir parámetros del propio modelo, para recoger el apartado 3

¿Pardon me?

# Descripción de otras medidas de error

- Entropía cruzada simétrica (para caso binario):

$$\frac{-1}{N} \sum_1^N \sum_1^M (y_{ij} \ln(r_{ij}) + (1 - y_{ij}) \ln(1 - r_{ij})) = \frac{-1}{N} \left( \sum_{y_{ij}=1} \ln(r_{ij}) + \sum_{y_{ij}=0} \ln(1 - r_{ij}) \right)$$

- Área bajo ROC (clasificación binaria)
  - Las respuestas están entre 0 y 1 (no/sí)
  - Se va variando el umbral de decisión no-sí
  - Se obtiene una curva poniendo en abscisas la fracción de falsos positivos y en ordenadas la de positivos acertados
  - Se mide el área bajo esa curva

Hazlo para el modelo de la pizarra

# Descripción de otras medidas de error

- Coeficiente de correlación

$$\rho_s^2 = \frac{[\sum_{i=1}^N (r_i - \bar{r})(y_i - \bar{y})]^2}{\sum_{i=1}^N (r_i - \bar{r})^2 \sum_{i=1}^N (y_i - \bar{y})^2}$$

- En clasificación

- Sensibilidad **a una clase**:  $\frac{\text{RealesCorrectos}}{\text{Reales}}$
- Especificidad **de una clase**:  $\frac{\text{FalsosCorrectos}}{\text{Falsos}}$
- Valor predictivo positivo **en una clase**:  $\frac{\text{RealesCorrectos}}{\text{RealesCorrectos} + \text{FalsosFallados}}$
- Precisión:  $\frac{\text{Aciertos}}{\text{Total}}$

Otra vez, para el de la pizarra

# Impredecibilidad del problema

- Problemas

- 1) Aleatoriedad intrínseca

- 2) Tipo de función

- 3) Tamaño de la muestra, representatividad  $\Rightarrow$  extrapolación, generalización

- Soluciones

- 1)...

- 2) Adecuación de la red. Cálculos de procesador

- 3) Diseño previo. Aumento de muestra. Control estadístico. Remuestreo.

¿Cuál es ...

...por el problema en sí?

...por cómo cogemos nosotros los datos del problema?

# Muestra

- Problemas
  - Ruido, aleatoriedad
  - Cobertura, dominio
  - ¿Están todas las variables?
  - Representatividad de la distribución. Problema de las clases importantes mal representadas.
- La dependencia de la muestra crece con el número de pesos (¿ecuaciones>incógnitas?)
- Soluciones (nunca total, hay fórmulas de estimación de cota)
  - N° de pesos (grandes):
    - error
    - especializar pesos, que no jueguen todos siempre (n° efectivo menor)
  - Oscilaciones fuertes: error
  - Regresión bayesiana
  - Muestrear la influencia de la muestra: media (es la red), varianza (es la muestra)
  - Muestra artificial ¿?: interpolación, estimación probabilista

Programa un aumento artificial de la muestra de casas

# Limitaciones de la red

- Por ser red neuronal
- Por tener cierto número de capas y ciertas conexiones
- Por la cantidad de procesadores
- Por el cálculo realizado en cada uno
- Soluciones: selección de modelos
- Aproximación si es una capa oculta, ¿cuántos componentes principales representan el 70% de la varianza? En muchos casos normales, cerca de 5



# Métodos de ajuste

- Punto de partida
  - Aleatorio que no condicione mal los elementos utilizados en el algoritmo: puede que cada vez vaya a sitio distinto
  - Sacado de otro método: autocodificación
- Algoritmo de búsqueda/actualización. Puede influirle las características de la muestra
- Criterio de parada
  - Agotamiento, ¿?
  - Estabilización, ¿?

¿De dónde podemos sacar un criterio para decidir que ya lo hace suficientemente bien?  
¿Qué podemos mirar para decidir que ya no va a progresar más?

# Algoritmo de actualización

- Búsqueda pura
  - Algoritmos genéticos
  - Cristalización
  - “Extremo”: capa oculta grande aleatoria, sólo se ajusta la de salida
- Uso de derivadas
  - Descenso gradiente
  - Gradiente conjugado
  - Cuasi-Newton, Métrica variable, Levenberg-Marquardt

Muestra tus conocimientos de optimización, por favor

# Problemas

- Espacio de enorme dimensión
  - Grave problema especialmente para búsqueda pura
- Mínimos locales
  - Grave especialmente con las derivadas
  - La exponencial del error cuadrático tiene menos mínimos locales
  - Evitar regiones sin gradiente: función de activación que tenga fluctuaciones en zonas “planas” (Ej: RRELU)
  - Un punto inicial distinto puede evitar ese mínimo local

¿Qué dimensión tiene el espacio de búsqueda del perceptrón de la pizarra?

# Enfoque probabilista bayesiano

Buscamos un buen modelo, que sea lo más probablemente posible el real

# Hipótesis

- La red aproxima la forma funcional perfectamente
- Los conjuntos de pesos distintos producen modelos distintos
- La parte no predecible tiene media nula, varianza constante y no está autocorrelada, ni correlada con las variables
- Esta parte tiene una estructura de probabilidad
- La muestra es representativa de la población total

Valora de 0 a 100 la fiabilidad de cada hipótesis

# Regresión bayesiana

$\text{Max}(\text{Prob}(\text{pesos}/\text{muestra})) = \text{Max}(\text{Prob}(\text{muestra}/\text{pesos})\text{Prob}(\text{pesos})/\text{Prob}(\text{muestra}))$

$\text{Max}(\log(\text{Prob}(\text{pesos}/\text{muestra}))) \iff$

$\text{Max}(\log(\text{Prob}(\text{muestra}/\text{pesos}) + \log(\text{Prob}(\text{pesos}))))$

Revisa los paréntesis

# Regresión bayesiana

$$\log(\text{Prob}(\text{muestra/pesos})) = \dots$$

$$\log\left(\prod_i \frac{e^{-\frac{(r_i - y_i)^2}{2\sigma_d^2}}}{\sqrt{2\pi\sigma_d}}\right) = \dots$$

$$\sum_i \left(\log \frac{e^{-\frac{(r_i - y_i)^2}{2\sigma_d^2}}}{\sqrt{2\pi\sigma_d}}\right) = \dots$$

$$-\left(N_d \log \sqrt{2\pi\sigma_d} + \sum_i \frac{(r_i - y_i)^2}{2\sigma_d^2}\right) = \dots$$

$$-\left(\text{Cte} + \frac{1}{2\sigma_d^2} \sum_i (r_i - y_i)^2\right)$$

¿Por qué nos pulimos el valor concreto de la constante?

# Regresión bayesiana

$$\log(\text{Prob}(\text{pesos})) = \dots$$

$$\log\left(\prod_j \frac{e^{-\frac{p_j^2}{2\sigma_p^2}}}{\sqrt{2\pi\sigma_p}}\right) = \dots$$

$$-\left(N_p \log \sqrt{2\pi\sigma_p} + \sum_j \frac{p_j^2}{2\sigma_p^2}\right) = \dots$$

$$-\left(K + \frac{1}{2\sigma_p^2} \sum_j p_j^2\right)$$

Repite con una función de densidad que creas mejor



# Regresión bayesiana

Max(Prob(pesos/muestra))

$$\text{Max}\left(-\left(\text{Cte} + \frac{1}{2\sigma_d^2} \sum_i (r_i - y_i)^2\right) - \left(K + \frac{1}{2\sigma_p^2} \sum_j p_j^2\right)\right)$$

$$\text{Min}\left(\text{Cte} + \frac{1}{2\sigma_d^2} \sum_i (r_i - y_i)^2 + K + \frac{1}{2\sigma_p^2} \sum_j p_j^2\right)$$

$$\text{Min}\left(\sum_i (r_i - y_i)^2 + \frac{\sigma_d^2}{\sigma_p^2} \sum_j p_j^2\right)$$

¿Y con otras funciones de densidad?