

# Advanced Kernel Methods for Multi-Task Learning

Tesis dirigida por José Dorronsoro y Carlos Alaíz

**Carlos Ruiz Pastor**

January 9, 2023



Universidad Autónoma  
de Madrid

Acknowledgements 1

Acknowledgements 2

# Outline

0

## ► Introducción

Multi-Task Learning

Support Vector Machines

## ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

## ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

## ► Summary

# Table of Contents

## 1 Introducción

- ▶ **Introducción**
  - Multi-Task Learning**
  - Support Vector Machines**
- ▶ Una Formulación Convexa para Aprendizaje Multitarea
  - Convex Multi-Task Learning with Kernel Methods
  - Convex Multi-Task Learning with Neural Networks
  - Combinación Convexa de modelos Preentrenados
- ▶ Laplaciano Adaptativo para Aprendizaje Multitarea
  - Laplaciano de Grafo con Métodos de Kernel
  - Adaptive Graph Laplacian Algorithm
- ▶ Summary

# Introducción al Aprendizaje Automático

## 1 Introducción

- El Aprendizaje Automático intenta automatizar el proceso de aprendizaje
- En el aprendizaje supervisado tenemos:
  - un espacio de entrada  $\mathcal{X}$ ,
  - un espacio de salida  $\mathcal{Y}$ ,
  - y una distribución  $P(x, y)$  (desconocida) sobre  $\mathcal{X} \times \mathcal{Y}$
- Dada una función  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , definimos una función de pérdida como

$$\begin{aligned}\ell : \mathcal{Y} \times \mathcal{Y} &\rightarrow [0, \infty) \\ (\gamma, f(x)) &\rightarrow \ell(\gamma, f(x)),\end{aligned}$$

tal que  $\ell(\gamma, \gamma) = 0$  para todo  $\gamma \in \mathcal{Y}$

# Loss Functions

## 1 Introducción

- In classification, with the class labels  $y_i \in \{-1, 1\}$ , we can use:

$$\ell(y, f(x)) = [1 - yf(x)]_+ = \begin{cases} 0, & yf(x) \geq 1, \\ 1 - yf(x), & yf(x) < 1. \end{cases}$$

•

# Expected Risk

## 1 Introducción

- Given a space of hypothesis  $\mathcal{H} = \{h(\cdot, \alpha), \alpha \in A\}$
- Definition: Expected Risk

$$R_P(\alpha) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x, \alpha)) dP(x, y)$$

- Our goal is to find

$$\alpha^* = \arg \min_{\alpha \in A} \left\{ R_P(\alpha) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(x, \alpha)) dP(x, y) \right\},$$

however the distribution  $P(x, y)$  is unknown

# Empirical Risk

## 1 Introducción

- Instead, we have a set of  $n$  instances sampled from  $P(x, y)$ :

$$D_n = \{(x_i, y_i), i = 1, \dots, n\},$$

- Definition: Empirical Risk

$$\hat{R}_{D_n}(\alpha) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i, \alpha))$$

- Instead of the Expected Risk, we minimize this empirical risk:

$$\arg \min_{\alpha \in A} \left\{ \hat{R}_D(\alpha) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(x_i, \alpha)) \right\}$$



# Table of Contents

## 1 Introducción

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Multi-Task Learning

## 1 Introducción

# Table of Contents

## 1 Introducción

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Table of Contents

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- ▶ Introducción
  - Multi-Task Learning
  - Support Vector Machines
- ▶ Una Formulación Convexa para Aprendizaje Multitarea
  - Convex Multi-Task Learning with Kernel Methods
  - Convex Multi-Task Learning with Neural Networks
  - Combinación Convexa de modelos Preentrenados
- ▶ Laplaciano Adaptativo para Aprendizaje Multitarea
  - Laplaciano de Grafo con Métodos de Kernel
  - Adaptive Graph Laplacian Algorithm
- ▶ Summary

# Formulación Aditiva

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Una manera de implementar el MTL es combinar una parte común y otras específicas
- La formulación aditiva para el aprendizaje multitarea es

$$h_r(\cdot) = g(\cdot) + g_r(\cdot)$$

donde

- $g(\cdot)$  es la parte común
- $g_r(\cdot)$  es la parte específica
- Fue propuesta para SVM lineales con los modelos

$$h_r(\cdot) = \langle w + v_r, \cdot \rangle + b_r$$

# Formulación Convexa

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Proponemos la siguiente formulación convexa para el aprendizaje multitarea:

$$h_r(\cdot) = \lambda_r g(\cdot) + (1 - \lambda_r) g_r(\cdot),$$

con  $\lambda_r \in [0, 1]$ .

- Los hiperparámetros  $\lambda_r$  regulan la influencia de cada parte:
  - $\lambda_1, \dots, \lambda_T = 0$ : modelos independientes (ITL)
  - $\lambda_1, \dots, \lambda_T = 1$ : modelo común (CTL)
- La interpretación de los hiperparámetros es más sencilla

# Table of Contents

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Formulacion Convexa con Métodos de Kernel

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- La formulación aditiva con métodos de kernel puede expresarse con los modelos:

$$h_r(\cdot) = \{\langle \mathbf{w}, \phi(\cdot) \rangle + \mathbf{b}\} + \{\langle \mathbf{v}_r, \phi_r(\cdot) \rangle + \mathbf{d}_r\}$$

- Con nuestra formulación convexa los modelos son:

$$h_r(\cdot) = \lambda_r \{\langle \mathbf{w}, \phi(\cdot) \rangle + \mathbf{b}\} + (1 - \lambda_r) \{\langle \mathbf{v}_r, \phi_r(\cdot) \rangle + \mathbf{d}_r\}$$

- Desarrollamos tres variantes de SVM:
  - L1-SVM
  - L2-SVM
  - LS-SVM



# Formulación Aditiva para MTL L1-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Primal - L1-SVM Aditiva

$$\begin{aligned} \arg \min_{\mathbf{w}, \mathbf{v}, \mathbf{b}, \boldsymbol{\xi}} \quad & J(\mathbf{w}, \mathbf{v}, \mathbf{b}, \boldsymbol{\xi}) = C \sum_{r=1}^T \sum_{i=1}^{m_r} \xi_i^r + \frac{1}{2} \sum_{r=1}^T \|\mathbf{v}_r\|^2 + \frac{\mu}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \gamma_i^r (\langle \mathbf{w}, \phi(\mathbf{x}_i^r) \rangle + \langle \mathbf{v}_r, \phi_r(\mathbf{x}_i^r) \rangle + b_r) \geq p_i^r - \xi_i^r, \\ & \xi_i^r \geq 0; \quad i = 1, \dots, m_r, \quad r = 1, \dots, T. \end{aligned}$$

- El parámetro  $\mu$  (junto con  $C$ ) regula la influencia de cada parte:
  - $\mu \rightarrow \infty$ : modelos independientes (ITL)
  - $C \rightarrow 0, \mu \rightarrow 0$ : modelo común (CTL)

# Formulación Aditiva para MTL L1-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Dual - L1-SVM Aditiva

$$\begin{aligned} \min_{\alpha} \quad & \Theta(\alpha) = \frac{1}{2} \alpha^\top \left( \frac{1}{\mu} Q + K \right) \alpha - \mathbf{p} \alpha \\ \text{s.t.} \quad & 0 \leq \alpha_i^r \leq C; \quad i = 1, \dots, m_r, \quad r = 1, \dots, T, \\ & \sum_{i=1}^{m_r} \alpha_i^r y_i^r = 0; \quad r = 1, \dots, T. \end{aligned}$$

- El parámetro  $\mu$  (junto con  $C$ ) regula la influencia de cada parte:
  - $\mu \rightarrow \infty$ : modelos independientes (ITL)
  - $C \rightarrow 0, \mu \rightarrow 0$ : modelo común (CTL)

# Formulación Convexa para MT L1-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Primal - L1-SVM Convexa

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}} \quad & J(\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}) = C \sum_{r=1}^T \sum_{i=1}^{m_r} \xi_i^r + \frac{1}{2} \sum_{r=1}^T \|\mathbf{v}_r\|^2 + \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i^r (\lambda_r \{ \langle \mathbf{w}, \phi(\mathbf{x}_i^r) \rangle + b \} + (1 - \lambda_r) \{ \langle \mathbf{v}_r, \phi_r(\mathbf{x}_i^r) \rangle + d_r \}) \geq p_i^r - \xi_i^r, \\ & \xi_i^r \geq 0, \quad i = 1, \dots, m_r, \quad r = 1, \dots, T. \end{aligned}$$

- Los hiperparámetros  $\lambda_r$  regulan la influencia de cada parte:
  - $\lambda_1, \dots, \lambda_T = 0$ : modelos independientes (ITL)
  - $\lambda_1, \dots, \lambda_T = 1$ : modelo común (CTL)
- El hiperparámetro  $C$  no interviene en la definición de los modelos

# Formulación Convexa para MT L1-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Dual - SVM Convexa

$$\min_{\alpha} \quad \Theta(\alpha) = \frac{1}{2} \alpha^T (\Lambda Q \Lambda + (I_n - \Lambda) K (I_n - \Lambda)) \alpha - \mathbf{p} \alpha$$

$$\text{s.t.} \quad 0 \leq \alpha_i^r \leq C; \quad i = 1, \dots, m_r, \quad r = 1, \dots, T,$$

$$\sum_{i=1}^{m_r} \alpha_i^r y_i^r = 0; \quad r = 1, \dots, T,$$

donde

$$\Lambda = \text{diag}(\overbrace{\lambda_1, \dots, \lambda_1}^{m_1}, \dots, \overbrace{\lambda_T, \dots, \lambda_T}^{m_T})$$

# Formulación Convexa para MT L1-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Dual - SVM Convexa ( $\lambda$ común)

$$\begin{aligned} \min_{\alpha} \quad & \Theta(\alpha) = \frac{1}{2} \alpha^\top \left( \lambda^2 Q + (1 - \lambda)^2 K \right) \alpha - \mathbf{p} \alpha \\ \text{s.t.} \quad & 0 \leq \alpha_i^r \leq C; \quad i = 1, \dots, m_r, \quad r = 1, \dots, T, \\ & \sum_{i=1}^{m_r} \alpha_i^r y_i^r = 0; \quad r = 1, \dots, T, \end{aligned}$$

- El hiperparámetro  $\lambda$  regula la influencia de cada parte:
  - $\lambda = 0$ : modelos independientes (ITL)
  - $\lambda = 1$ : modelo común (CTL)
- El hiperparámetro  $C$  no interviene en la definición de los modelos

# Proposiciones

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Proposicion (Equivalencia entre formulaciones para SVM)

*Para valores  $\lambda \in (0, 1)$ , la formulación aditiva con hiperparámetros  $C_{add}, \mu$  y la formulación convexa con  $C_{conv}$  y un  $\lambda$  común,  $\lambda_1, \dots, \lambda_T = \lambda$ , son equivalentes cuando*

$$C_{add} = (1 - \lambda)^2 C_{conv}, \mu = (1 - \lambda)^2 / \lambda^2.$$

### Proposicion (Equivalencia con CTL e ITL)

- *Para  $\lambda = 0$ , la formulación convexa con un  $\lambda$  común es equivalente a modelos independientes (ITL).*
- *Para  $\lambda = 1$  la formulación convexa con un  $\lambda$  común es equivalente a un modelo común (CTL).*

# Formulación convexa para MT L2-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Primal - MTL L2-SVM Convexa

$$\begin{aligned} \arg \min_{\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}} \quad & J(\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}) = \frac{C}{2} \sum_{r=1}^T \sum_{i=1}^{m_r} (\xi_i^r)^2 + \frac{1}{2} \sum_{r=1}^T \|\mathbf{v}_r\|^2 + \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \gamma_i^r (\lambda_r \{ \langle \mathbf{w}, \phi(\mathbf{x}_i^r) \rangle + b \} + (1 - \lambda_r) \{ \langle \mathbf{v}_r, \phi_r(\mathbf{x}_i^r) \rangle + d_r \}) \geq p_i^r - \xi_i^r, \end{aligned}$$

# Formulación Convexa para MT L2-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Dual - MTL L2-SVM Convexa

$$\begin{aligned} \min_{\alpha} \quad & \Theta(\alpha) = \frac{1}{2} \alpha^T \left( \{ \Lambda Q \Lambda + (I_n - \Lambda) K (I_n - \Lambda) \} + \frac{1}{C} I \right) \alpha - p \alpha \\ \text{s.t.} \quad & 0 \leq \alpha_i^r, \quad i = 1, \dots, m_r, \quad r = 1, \dots, T, \\ & \sum_{i=1}^{m_r} \alpha_i^r y_i^r = 0, \quad r = 1, \dots, T. \end{aligned}$$



# Formulación convexa para MT LS-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Primal - MTL LS-SVM Convexa

$$\begin{aligned} \arg \min_{\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}} \quad & J(\mathbf{w}, \mathbf{v}, b, \mathbf{d}, \boldsymbol{\xi}) = \frac{C}{2} \sum_{r=1}^T \sum_{i=1}^{m_r} (\xi_i^r)^2 + \frac{1}{2} \sum_{r=1}^T \|\mathbf{v}_r\|^2 + \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \gamma_i^r (\lambda_r \{ \langle \mathbf{w}, \phi(\mathbf{x}_i^r) \rangle + b \} + (1 - \lambda_r) \{ \langle \mathbf{v}_r, \phi_r(\mathbf{x}_i^r) \rangle + d_r \}) = p_i^r - \xi_i^r, \end{aligned}$$

# Formulación Convexa para MT LS-SVM

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Problema Dual - MTL LS-SVM Convexa

$$\left[ \begin{array}{c|c|c} 0 & \mathbf{0}_T^\top & \mathbf{y}^\top \Lambda \\ \hline \mathbf{0}_T & \mathbf{0}_{T \times T} & A^\top Y (I_n - \Lambda) \\ \hline \mathbf{y} & YA & \widehat{Q} + \frac{1}{c} I_n \end{array} \right] \begin{pmatrix} b \\ d_1 \\ \vdots \\ d_T \\ \alpha \end{pmatrix} = \begin{pmatrix} 0 \\ \mathbf{0}_T \\ \mathbf{p} \end{pmatrix},$$

# Table of Contents

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Redes Neuronales MT

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- La manera más común de adaptar las redes neuronales es el *hard sharing*
  - Capas ocultas compartidas por todas las tareas
  - Capas de salida específicas para cada tarea
- El modelo se puede expresar como:

$$h_r(\cdot) = g_r(\cdot; w_r, \Theta) = \{\langle w_r, f(\cdot; \Theta) \rangle\} + d_r$$

- $w_r, d_r$  son los parámetros de las capas de salida específicas
- $\Theta$  son los parámetros de las capas ocultas compartidas

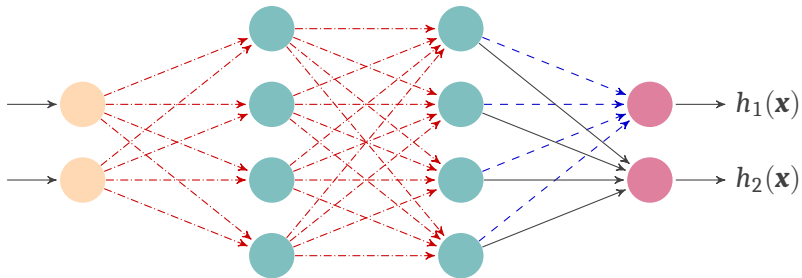


Figure: Ejemplo de *Hard Sharing* para dos tareas .

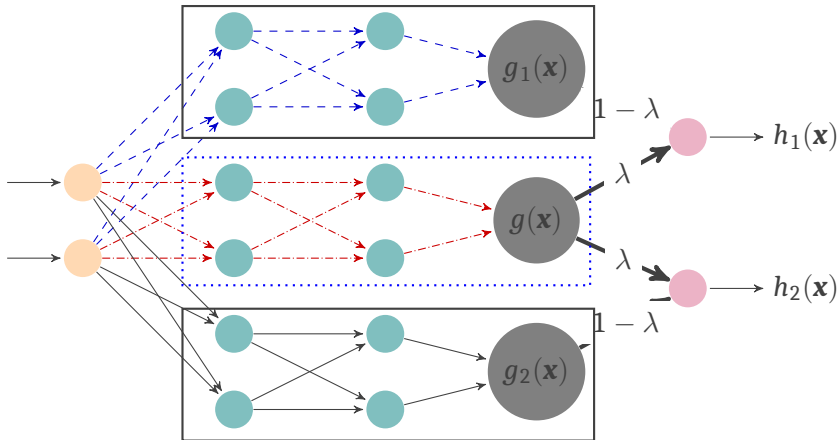
# Formulación Convexa para Redes Neuronales MT

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Proponemos la formulación convexa para redes neuronales MT, combinando:
  - Una parte común  $g(\cdot; w, \Theta)$
  - Una parte específica  $g_r(\cdot; w_r, \Theta_r)$
- Los modelos son:

$$\begin{aligned} h_r(\cdot) &= \lambda_r g(\cdot; w, \Theta) + (1 - \lambda_r) g_r(\cdot; w_r, \Theta_r) \\ &= \lambda_r \{ \langle w, f(\cdot; \Theta) \rangle + b \} + (1 - \lambda_r) \{ \langle w_r, f_r(\cdot; \Theta_r) \rangle + d_r \}. \end{aligned}$$

- $w, \Theta$  son los parámetros de la red común (capa de salida y ocultas)
- $w_r, \Theta_r$  son los parámetros de las redes específicas (capa de salida y ocultas)



**Figure:** Ejemplo de formulación convexa con redes neuronales para dos tareas.

# Formulación Convexa para Redes Neuronales MT

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- El riesgo a minimizar en este caso es

$$\hat{R}_D = \sum_{r=1}^T \sum_{i=1}^{m_r} \ell(h_r(x_i^r), y_i^r) + \frac{\mu}{2} \left( \|w\|^2 + \sum_{r=1}^T \|w_r\|^2 + \Omega(\Theta) + \Omega(\Theta_r) \right).$$

- Se puede aplicar el descenso por gradiente con

$$\begin{aligned} \nabla_w h_t(x_i^t) &= \lambda_t f(x_i^t, \Theta), & \nabla_{\Theta} h_t(x_i^t) &= \lambda_t \langle w, \nabla_{\Theta} f(x_i^t, \Theta) \rangle; \\ \nabla_{w_t} h_t(x_i^t) &= (1 - \lambda_t) f_t(x_i^t, \Theta), & \nabla_{\Theta_t} h_t(x_i^t) &= (1 - \lambda_t) \langle w, \nabla_{\Theta_t} f_t(x_i^t, \Theta_t) \rangle; \\ \nabla_{w_r} h_t(x_i^t) &= 0, & \nabla_{\Theta_r} h_t(x_i^t) &= 0, \text{ for } r \neq t. \end{aligned}$$

- Los gradientes se escalan adecuadamente con  $\lambda_t$  y  $(1 - \lambda_t)$



# Formulación Convexa para Redes Neuronales MT

## 2 Una Formulación Convexa para Aprendizaje Multitarea

---

### Algorithm 1: Pase “forward”

---

```

Input:  $X_{mb}, t_{mb}$                                 // Minibatch data and task labels
Output:  $f$                                           // Forward pass for the minibatch
Data:  $\lambda$                                        // Parameter of convex combination
Data:  $g; g_1, \dots, g_T$                          // Modules of the common and specific networks
for  $x_i, t_i \in (X_{mb}, t_{mb})$  do
  |  $f_i \leftarrow \lambda g(x_i) + (1 - \lambda) g_{t_i}(x_i)$  // Convex combination
end

```

---

- El pase “backward” se hace con la diferenciación automática de PyTorch

# Table of Contents

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Combinación Convexa de modelos Preentrenados

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Alternativa a la formulación convexa para aprendizaje MT
- Consideramos la combinación convexa de
  - modelo común  $g(\cdot)$  entrenado
  - modelos específicos  $g_r(\cdot)$  entrenados
- Minimizamos el riesgo eligiendo los hiperparámetros  $\lambda_1, \dots, \lambda_T$  óptimos

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} \ell(\lambda_r g(x_i^r) + (1 - \lambda_r) g_r(x_i^r), y_i^r),$$

# Formulación Unificada Clasificación

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Hinge loss (classification):

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} [1 - y_i^r \{ \lambda_r g(x_i^r) + (1 - \lambda_r) g_r(x_i^r) \}]_+.$$

- Squared hinge loss (classification):

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} [1 - y_i^r \{ \lambda_r g(x_i^r) + (1 - \lambda_r) g_r(x_i^r) \}]_+^2.$$

- Ambas se pueden expresar como:

$$\sum_{r=1}^T \sum_{i=1}^{m_r} u(\lambda_r c_i^r + d_i^r), \text{ donde } c_i^r = y_i^r (g_r(x_i^r) - g(x_i^r)), \text{ } d_i^r = 1 - y_i^r g_r(x_i^r)$$

# Formulación Unificada Regresión

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- Absolute loss (regression):

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} |y_i^r - \{\lambda_r g(x_i^r) + (1 - \lambda_r) g_r(x_i^r)\}|.$$

- Squared loss (regression):

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} (y_i^r - \{\lambda_r g(x_i^r) + (1 - \lambda_r) g_r(x_i^r)\})^2.$$

- Ambas se pueden expresar como:

$$\sum_{r=1}^T \sum_{i=1}^{m_r} u(\lambda_r c_i^r + d_i^r), \text{ donde } c_i^r = g(x_i^r) - g_r(x_i^r), \text{ } d_i^r = g_r(x_i^r) - y_i^r$$

## Formulación Unificada

### 2 Una Formulación Convexa para Aprendizaje Multitarea

- En todos los casos tenemos que minimizar

$$\hat{R}_D(\lambda_1, \dots, \lambda_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} u(\lambda_r c_i^r + d_i^r)$$

- Como es separable, tenemos en cada tarea el problema

$$\arg \min_{\lambda_r \in [0,1]} \mathcal{J}(\lambda_r) = \sum_{i=1}^{m_r} u(\lambda_r c_i^r + d_i^r),$$

- Usando el Teorema de Fermat

$$\lambda^* = \arg \min_{0 \leq \lambda \leq 1} \mathcal{J}(\lambda) \iff (0 \in \partial \mathcal{J}(\lambda^*) \text{ and } \lambda^* \in (0, 1)) \text{ or } \lambda^* = 0 \text{ or } \lambda^* = 1.$$

# Combinación Convexa con Error Cuadrático

## 2 Una Formulación Convexa para Aprendizaje Multitarea

- La función a minimizar es

$$\arg \min_{\lambda \in [0,1]} \mathcal{J}(\lambda) = \sum_{i=1}^m (\lambda c_i + d_i)^2.$$

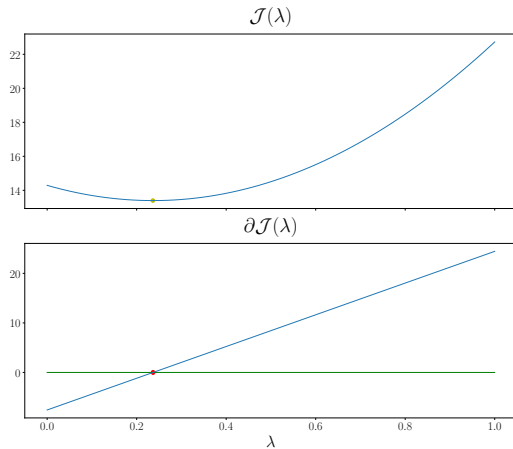
- La derivada es

$$\mathcal{J}'(\lambda) = \sum_{i=1}^m 2c_i(\lambda c_i + d_i).$$

- Como es derivable, resolviendo  $\mathcal{J}'(\lambda) = 0$  obtenemos

$$\lambda' = -\frac{\sum_{i=1}^m d_i c_i}{\sum_{i=1}^m (c_i)^2}.$$

- La solución es entonces  $\lambda^* = \max(\min(\lambda', 1), 0)$





# Combinación Convexa con Error Absoluto

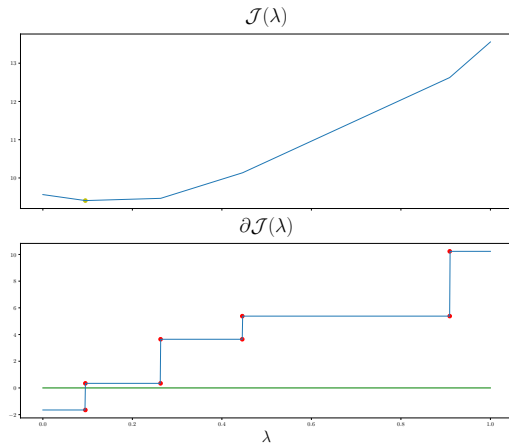
## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Proposición ( $\lambda^*$ óptimo para el problema con valor absoluto)

- $\lambda^* = 0$  es óptimo si y solo si:  $-\sum_{i: 0 > \lambda_{(i)}} |c_{(i)}| + \sum_{i: 0 < \lambda_{(i)}} |c_{(i)}| \leq 0$
- $\lambda^* \in (0, 1)$  es óptimo si y solo si  $0 < \lambda^* = \lambda_{(k)} < 1$  para algún  $k = 1, \dots, m$ , y

$$-\sum_{i: \lambda_{(k)} > \lambda_{(i)}} |c_{(i)}| + \sum_{i: \lambda_{(k)} < \lambda_{(i)}} |c_{(i)}| \in [-|c_{(k)}|, |c_{(k)}|]$$

- $\lambda^* = 1$  es óptimo en otro caso

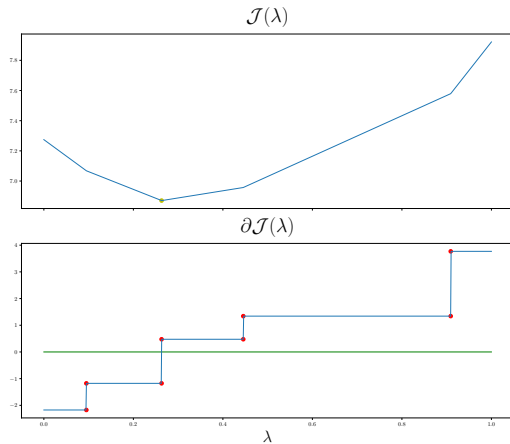


# Combinación Convexa con Error Hinge

## 2 Una Formulación Convexa para Aprendizaje Multitarea

### Proposición ( $\lambda^*$ óptimo para el problema con error hinge)

- $\lambda^* = 0$  es óptimo si y solo si:  $-\sum_{i: 0 > \lambda_{(i)}} \max(0, c_{(i)}) - \sum_{0 < \lambda_{(i)}} \min(0, c_{(i)}) \leq 0$
- $\lambda^* \in (0, 1)$  es óptimo si y solo si  $0 < \lambda^* = \lambda_{(k)} < 1$  para algún  $k = 1, \dots, m$ , y
 
$$-\sum_{i: \lambda_{(k)} > \lambda_{(i)}} \max(0, c_{(i)}) - \sum_{i: \lambda_{(k)} < \lambda_{(i)}} \min(0, c_{(i)}) \in [\min(0, c_{(k)}), \max(0, c_{(k)})]$$
- $\lambda^* = 1$  es óptimo en otro caso



# Combinación Convexa con Error Hinge Cuadrático

## 2 Una Formulación Convexa para Aprendizaje Multitarea

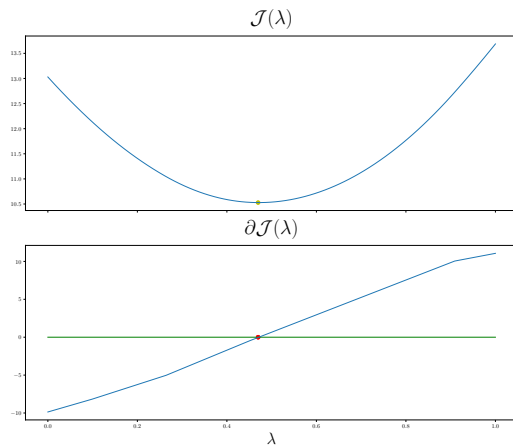
### Proposición ( $\lambda^*$ óptimo para el problema con error hinge cuadrático)

- $\lambda^* = 0$  es óptimo si y solo si:  $-\sum_{i: 0 > c_{(i)}, 0 < \lambda_{(i)}} 2c_i d_i - \sum_{i: 0 < c_{(i)}, 0 > \lambda_{(i)}} 2c_i d_i \leq 0$
- $\lambda^* \in (0, 1)$  es óptimo si y solo si  $0 < \lambda^* = \hat{\lambda}_{(k)} < 1$  para algún  $k = 1, \dots, m$ , donde

$$\hat{\lambda}_{(k)} = - \frac{\sum_{i: \lambda_{(k+1)} \geq \lambda_{(i)}} \max(0, c_{(i)}) d_{(i)} + \sum_{i: \lambda_{(k)} \leq \lambda_{(i)}} \min(0, c_{(i)}) d_{(i)}}{\sum_{i: \lambda_{(k+1)} \geq \lambda_{(i)}} \max(0, c_{(i)})^2 + \sum_{i: \lambda_{(k)} \leq \lambda_{(i)}} \min(0, c_{(i)})^2},$$

y además  $\lambda_{(k)} \leq \hat{\lambda}_k \leq \lambda_{(k+1)}$

- $\lambda^* = 1$  es óptimo en otro caso



# Table of Contents

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

- ▶ Introducción
  - Multi-Task Learning
  - Support Vector Machines
- ▶ Una Formulación Convexa para Aprendizaje Multitarea
  - Convex Multi-Task Learning with Kernel Methods
  - Convex Multi-Task Learning with Neural Networks
  - Combinación Convexa de modelos Preentrenados
- ▶ Laplaciano Adaptativo para Aprendizaje Multitarea
  - Laplaciano de Grafo con Métodos de Kernel
  - Adaptive Graph Laplacian Algorithm
- ▶ Summary

# Aprendizaje Multitarea con Regularización Laplaciana

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

- Otra manera de acoplar distintas tareas es usar una regularización Laplaciana
- Consideramos un grafo donde
  - Los nodos representan tareas
  - Las aristas y sus pesos representan las relaciones entre las tareas
- La matriz de adyacencia  $A$  tiene los pesos de las aristas
- La matriz de grados  $D$  es una matriz diagonal donde

$$(D)_{rr} = \sum_{s=1}^T (A)_{rs}$$

- La matriz Laplaciana se define como  $L = D - A$



# Aprendizaje Multitarea con Regularización Laplaciana

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

- Dados los modelos para cada tarea definidos como

$$h_r(\cdot) = \langle w_r, \cdot \rangle + b_r$$

- Definimos la regularización

$$\sum_{r=1}^T \sum_{s=1}^T (A)_{rs} \|w_r - w_s\|^2,$$

- Esta regularización se puede expresar como

$$\sum_{r=1}^T \sum_{s=1}^T (A)_{rs} \|w_r - w_s\|^2 = \sum_{r=1}^T \sum_{s=1}^T (L)_{rs} \langle w_r, w_s \rangle,$$

# Table of Contents

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Laplaciano de Grafo con Métodos de Kernel

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

- Consideramos el problema de minimización

$$R(u_1, \dots, u_T) = \sum_{r=1}^T \sum_{i=1}^{m_r} \ell(y_i^r, \langle u_r, \phi(x_i^r) \rangle) + \mu \sum_r \sum_s (E)_{rs} \langle u_r, u_s \rangle, \quad (0)$$

- Si usamos el vector  $\mathbf{u}^\top = (u_1^\top, \dots, u_T^\top)$  lo expresamos como

$$R(\mathbf{u}) = \sum_{r=1}^T \sum_{i=1}^{m_r} \ell(y_i^r, \langle \mathbf{u}, \mathbf{e}_r \otimes \phi(x_i^r) \rangle) + \mu (\mathbf{u}^\top (E \otimes I) \mathbf{u}). \quad (1)$$

# Laplaciano de Grafo con Métodos de Kernel

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

### Lemma

Las soluciones  $u_1^*, \dots, u_T^*$  de (0), o equivalentemente la solución  $\mathbf{u}^*$  de (1), se pueden obtener minimizando

$$S(\mathbf{w}) = \sum_{r=1}^T \sum_{i=1}^{m_r} \ell(y_i^r, \langle \mathbf{w}, (B_r \otimes \phi(x_i^r)) \rangle) + \mu \mathbf{w}^\top \mathbf{w}, \quad (2)$$

donde  $\mathbf{w} \in \mathbb{R}^p \otimes \mathcal{H}$  con  $p \geq T$  y  $B_r$  son las columnas de  $B \in \mathbb{R}^{p \times T}$ , una matriz de rango máximo tal que  $E^{-1} = B^\top B$ .

El kernel reproductor correspondiente es:

$$\langle B_r \otimes \phi(x_i^r), B_s \otimes \phi(x_j^s) \rangle = (E^{-1})_{rs} k(x_i^r, x_j^s)$$

# Table of Contents

## 3 Laplaciano Adaptativo para Aprendizaje Multitarea

### ► Introducción

Multi-Task Learning

Support Vector Machines

### ► Una Formulación Convexa para Aprendizaje Multitarea

Convex Multi-Task Learning with Kernel Methods

Convex Multi-Task Learning with Neural Networks

Combinación Convexa de modelos Preentrenados

### ► Laplaciano Adaptativo para Aprendizaje Multitarea

Laplaciano de Grafo con Métodos de Kernel

Adaptive Graph Laplacian Algorithm

### ► Summary

# Table of Contents

## 4 Summary

- ▶ Introducción
  - Multi-Task Learning
  - Support Vector Machines
- ▶ Una Formulación Convexa para Aprendizaje Multitarea
  - Convex Multi-Task Learning with Kernel Methods
  - Convex Multi-Task Learning with Neural Networks
  - Combinación Convexa de modelos Preentrenados
- ▶ Laplaciano Adaptativo para Aprendizaje Multitarea
  - Laplaciano de Grafo con Métodos de Kernel
  - Adaptive Graph Laplacian Algorithm
- ▶ **Summary**

# Good Luck!

## 4 Summary

- Enough for an introduction! You should know enough by now

# Advanced Kernel Methods for Multi-Task Learning

*Thank you for listening!*

*Any questions?*