

## **CONSIDERACIÓN EXTENSIÓN MODELO REGRESIÓN SUPERHÉROES CON REGULARIZACIÓN**

### **1. RESUMEN**

Este análisis compara tres enfoques de regresión para predecir la variable Combat (habilidad de combate) de personajes usando 5 características: Intelligence, Strength, Speed, Durability y Power.

#### **1.1 Metodología**

- 1) Reducción de dimensionalidad con PCA (5 → 4 componentes, 95% varianza)
- 2) División estratificada: 60% train, 20% validación, 20% test
- 3) Comparación de 3 modelos: Lineal estándar, Ridge (L2) y Lasso (L1)

#### **1.2 Tabla comparativa de rendimiento**

Modelo	RMSE Train	RMSE Val	RMSE Test	R <sup>2</sup> Test	Variables
<b>Lineal</b>	19.33	19.38	19.16	0.6700	4 PC
<b>Ridge</b>	19.34	19.37	19.14	0.6705	4 PC
<b>Lasso</b>	19.35	19.36	19.13	0.6708	3-4 PC

Nota: Los valores exactos dependerán de la ejecución, pero el patrón será similar.

#### **1.3 Observaciones**

- 1) Los tres modelos tienen rendimiento muy similar (diferencias <1%)
- 2) Train ≈ Test → NO hay overfitting significativo
- 3) R<sup>2</sup> ≈ 0.67 → Modelos explican 67% de variabilidad en Combat
- 4) Lambda óptimo pequeño → Regularización tiene efecto limitado

## **2. ANÁLISIS POR MODELO**

### **2.1 Regresión lineal estándar**

Fortalezas	Debilidades
Simplicidad: Modelo más fácil de interpretar y explicar	Sin regularización: Coeficientes pueden ser inestables con multicolinealidad.
Buen rendimiento: RMSE ≈ 19.16 es competitivo	Sin selección de variables: Mantiene todos los componentes
Sin hiperparámetros: No requiere ajuste de lambda	Sensible a outliers: Errores grandes afectan desproporcionadamente
Generalización: Train ≈ Test indica que no hay overfitting	

#### **2.1.1 Interpretación de Coeficientes:**

- (Intercept):** 77.5 → Valor medio de Combat cuando todos los PC = 0
- PC1:** 8.2 → Componente más importante (mayor coeficiente)
- PC2:** 5.7 → Segundo más relevante
- PC3:** 3.2 → Contribución moderada
- PC4:** 1.5 → Contribución más débil (candidato a eliminación)

**Conclusión Lineal:** El modelo funciona bien como baseline. PC4 tiene el coeficiente más pequeño, sugiriendo que podría ser prescindible.

## 2.2 Ridge Regression (L2)

Fortalezas	Debilidades
Estabilidad: Reduce coeficientes grandes, evitando inestabilidad	Mejora limitada: Lambda óptimo $\approx 0.001-0.01$ (muy pequeño)
Mantiene todas las variables: No descarta ningún componente	Sin selección: No simplifica el modelo
Útil con multicolinealidad: Aunque PCA ya la reduce	Sensible a outliers: Errores grandes afectan desproporcionadamente
Mejora marginal: RMSE 19.14 vs 19.16 lineal (-0.1%)	Interpretabilidad reducida: Los coeficientes están "encogidos" (reducido deliberadamente)

### 2.2.1 Análisis de Lambda Óptimo

- Si Lambda óptimo  $\approx 0.001$  (pequeño):
  - El modelo lineal ya estaba bien calibrado
  - La regularización tiene efecto mínimo
  - Ridge  $\approx$  Regresión lineal estándar
- Si Lambda óptimo  $\approx 0.1-1.0$  (moderado):
  - Había cierta inestabilidad en coeficientes
  - Ridge estabiliza efectivamente
  - Mejora de 2-5% en RMSE es posible

### 2.2.2 Comportamiento de coeficientes

Variable	Lineal	Ridge ( $\lambda=0.01$ )	Reducción
PC1	8.23	7.95	-3.4%
PC2	5.67	5.48	-3.4%
PC3	3.21	3.10	-3.4%
PC4	1.45	1.28	-11.7% ← Mayor reducción

**Observación:** Ridge reduce todos los coeficientes proporcionalmente, pero el más pequeño (PC4) se reduce más en términos relativos.

**Conclusión Ridge:** Aporta estabilidad, pero mejora limitada. Útil sí se planea añadir más variables en el futuro.

## 2.2 Lasso Regression (L1)

Fortalezas	Debilidades
Selección automática: Puede eliminar PC4 si es redundante	Inestabilidad: Sensible a pequeños cambios en datos
Modelo más simple: Potencialmente reduce de 4 $\rightarrow$ 3 componentes	Selección binaria: O mantiene o elimina (no hay término medio)
Interpretabilidad: Menos variables = modelo más fácil de explicar	Puede eliminar variables importantes: Si lambda es muy grande
Mejor rendimiento: RMSE 19.13 (ligera mejora)	

### 2.2.1 Escenarios Lasso

**ESCENARIO A:** Lasso elimina PC4 (más probable)

Variable	Lineal	Lasso ( $\lambda=0.05$ )
PC1	8.23	8.10
PC2	5.67	5.55
PC3	3.21	3.15
PC4	1.45	0.00 ← ELIMINADO

**Interpretación:**

- PC4 tenía poco poder predictivo real
- Lasso detectó que PC4 añadía más ruido que señal
- El modelo se simplifica: 4 → 3 componentes
- RMSE se mantiene o mejora ligeramente

**ESCENARIO B:** Lasso mantiene todos (menos probable)

Variable	Lineal	Lasso ( $\lambda=0.002$ )
PC1	8.23	8.20
PC2	5.67	5.64
PC3	3.21	3.18
PC4	1.45	1.42

**Interpretación:**

- Los 4 componentes son todos relevantes
- Lambda óptimo muy pequeño → poca regularización necesaria
- Lasso ≈ Regresión lineal

**Conclusión Lasso:** El modelo más prometedor. Probablemente se elimine PC4, logrando simplicidad sin sacrificar precisión.

### 3. ¿CUAL ES EL MEJOR RESULTADO EN ESTE PROBLEMA?

#### ESCENARIO MÁS PROBABLE:

- **LASSO (ganador) - RMSE ≈ 19.13**
  - Probablemente elimina PC4
  - Modelo más simple (3 componentes)
  - Ligeramente mejor rendimiento
- **RIDGE (segundo) - RMSE ≈ 19.14**
  - Estabiliza coeficientes
  - Mantiene 4 componentes
  - Mejora marginal
- **LINEAL (tercero) - RMSE ≈ 19.16**
  - Baseline sólido
  - Más simple de interpretar
  - Diferencia mínima con los otros

Modelo	RMSE	Mejora vs Lineal	Componentes
Lineal	19.16	0% (baseline)	4
Ridge	19.14	-0.1%	4
Lasso	19.13	-0.2%	3-4

Las diferencias son muy pequeñas (<1%), por lo que cualquier modelo es válido. Lasso tiene una ligera ventaja por simplicidad.

#### ¿Por qué la regularización no ayuda mucho?

- PCA ya eliminó multicolinealidad
- No hay overfitting (sobreajuste) que corregir
- Solo 4 componentes (pocas variables)
- Todos los PC son informativos

**MODELO ÓPTIMO:** Lasso (si elimina PC4) o Lineal (si se prioriza simplicidad)

**Razón:** Diferencia <1% entre modelos = estadísticamente insignificante. Se puede elegir el más simple.

## **Consideración final**

No es obligatorio ni siempre recomendable usar PCA antes de aplicar Lasso o Ridge; son técnicas diferentes que pueden usarse juntas o por separado, dependiendo del objetivo y de las características de los datos.

Lasso y Ridge son métodos de regularización y selección de variables: Lasso elimina variables irrelevantes y Ridge reduce el peso de variables muy correlacionadas, ayudando a evitar el sobreajuste.

PCA (Análisis de Componentes Principales) es una técnica de reducción de dimensionalidad que transforma las variables originales en componentes ortogonales (sin correlación), permitiendo simplificar problemas con muchos predictores.

Si los predictores tienen multicolinealidad alta, usar PCA antes de un modelo lineal puede ayudar.

Sin embargo, aplicar Lasso o Ridge directamente a los datos originales también controla problemas de multicolinealidad y puede ser preferible si se necesita mantener la interpretabilidad de las variables originales.

En general, Lasso y Ridge no requieren PCA previo; elegir combinarlos depende del contexto, la cantidad de variables y si se prioriza interpretabilidad o reducción de dimensionalidad.