

## **Resumen del artículo *An overview of gradient descent optimization algorithms***

**Alumno: Andrés Padrón Quintana**

El artículo de Sebastian Ruder ofrece una visión clara y práctica sobre algoritmos de optimización para el descenso por gradiente en aprendizaje automático. Inicia diferenciando las tres variantes básicas—batch, estocástico y mini-batch—y explica el compromiso entre exactitud de la actualización y costo computacional.

Después, se presentan los principales retos del entrenamiento: selección del learning rate, sensibilidad a superficies no convexas con múltiples mínimos locales y la presencia de saddle points, así como la necesidad de tasas de aprendizaje que se adapten a la frecuencia de características escasas. Sobre esta base, el texto desarrolla los optimizadores más usados: Momentum y Nesterov (aceleración anticipada) para reducir oscilaciones y ganar inercia en valles; Adagrad para tasas adaptativas por parámetro (especialmente útil con datos escasos), y sus mejoras Adadelta y RMSprop que evitan la disminución monótona de la tasa; Adam, que combina momento y promedios exponenciales con corrección de sesgo; así como extensiones como AdaMax (norma infinito) y Nadam (Nesterov + Adam). El artículo también compara trayectorias sobre funciones de prueba y muestra que los métodos adaptativos suelen escapar con mayor rapidez de puntos de silla y converger de forma estable. Se complementa con estrategias transversales—barajado de datos o curriculum learning, normalización por lotes, early stopping y ruido en gradientes—y ofrece un vistazo a esquemas paralelos y distribuidos (Hogwild!, Downpour SGD, EASGD) para acelerar el entrenamiento.

En conclusión, Ruder sugiere que Adam y métodos afines suelen ser elecciones seguras por su robustez y velocidad, si bien SGD con un buen calendario de learning rate puede alcanzar resultados competitivos cuando el tiempo de cómputo y la sintonía fina no son restricciones críticas.

### **Aplicación en el área laboral**

1. Modelos de riesgo de crédito y scoring: En portafolios con variables categóricas de alta cardinalidad y datos escasos (p. ej., historiales cortos, codificaciones one-hot), optimizadores adaptativos como Adam o RMSprop agilizan la convergencia y reducen la sensibilidad al escalado de entradas. Ejemplo: entrenamiento de una regresión logística o red feed-forward para probabilidad de incumplimiento, usando early stopping y reducción de tasa cuando la pérdida de validación se estanca. Actualmente este tema me apasiona mucho debido a que estoy haciendo mi tesis de licenciatura sobre este tema.
2. Predicción operativa en seguros y pricing: Para redes que estiman frecuencias o severidades, RMSprop/Adam estabilizan el aprendizaje cuando la varianza de gradientes es alta; el barajado por épocas y batch normalization acortan tiempos de entrenamiento.

3. Modelos de alta dimensión (series de tiempo con embeddings o múltiples features): En arquitecturas profundas, Momentum/Nesterov ayudan a avanzar por ravines; Adam/Nadam permite tamaños de paso eficaces sin ajustar manualmente cada capa. Ejemplo: red para nowcasting de indicadores financieros con regularización y validación temporal.
4. Entrenamiento distribuido en pipelines empresariales: En escenarios con grandes volúmenes, estrategias como Hogwild! (para esparsidad) o esquemas asíncronos tipo Downpour pueden reducir tiempos de ciclo cuando existen recursos de cómputo paralelos.

## Referencias

Ruder, S. (2017). *An overview of gradient descent optimization algorithms*. arXiv preprint arXiv:1609.04747.