

Adam AdamW

Adam??????

1.
 - $m_0=0$ $v_0=0$
 - η 0.001 β_1 0.9 β_2 0.999 ϵ 1e-8
2. $\theta_{\{t\}}$ $g_{\{t\}}$
3. $m_{\{t\}} = \beta_1 \cdot m_{\{t-1\}} + (1 - \beta_1) \cdot g_{\{t\}}$
4. $v_{\{t\}} = \beta_2 \cdot v_{\{t-1\}} + (1 - \beta_2) \cdot g_{\{t\}}^2$
5. **Bias Correction** m_t v_t 0 t
 $m_{\{t\}} = m_{\{t\}} / (1 - \beta_{\{1\}}^{\{t\}})$ $v_{\{t\}} = v_{\{t\}} / (1 - \beta_{\{2\}}^{\{t\}})$
6. $\theta_{\{t\}}$

```
class AdamOptimizer:
    def __init__(self, learning_rate=0.01, beta1=0.9, beta2=0.999, epsilon=1e-8):
        self.lr = learning_rate
        self.beta1 = beta1
        self.beta2 = beta2
        self.epsilon = epsilon
        self.m = 0
        self.v = 0
        self.t = 0

    def update(self, params, grads):
        self.t += 1
        for param, grad in zip(params, grads):
            self.m[param] = self.beta1 * self.m[param] + (1 - self.beta1) * grad
            self.v[param] = self.beta2 * self.v[param] + (1 - self.beta2) * grad**2
            m_hat = self.m[param] / (1 - self.beta1**self.t)
            v_hat = self.v[param] / (1 - self.beta2**self.t)
            param = param - self.lr * m_hat / (np.sqrt(v_hat) + self.epsilon)
        return params
```

weight decay

██

```
w = w - lr * w.grad - lr * weight_decay * w
```

L2

```
final_loss = loss + wd * all_weights.pow(2).sum() / 2
```

Adam L2

L2

```
grad_w = grad_w + weight_decay * w
```

AdamW

AdamW

L2

grad_w

```
grad_w = grad_w
# update step
w = w - learning_rate * grad_w - learning_rate * weight_decay * w
```

Revision #4

Created 2025-03-12 08:09:34 UTC by Colin

Updated 2026-04-29 07:34:10 UTC by Colin