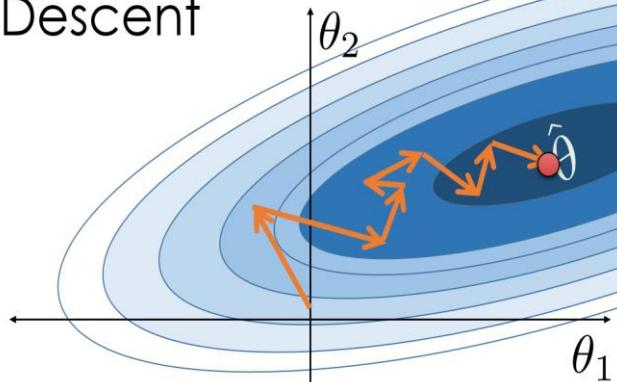


دورهی آموزشی «علم داده»
Data Science Course

Stochastic Gradient
Descent



جلسه سیام (بخش دوم)
گرادیان کاهشی تصادفی
و تکانه
Stochastic Gradient Descent
& Momentum

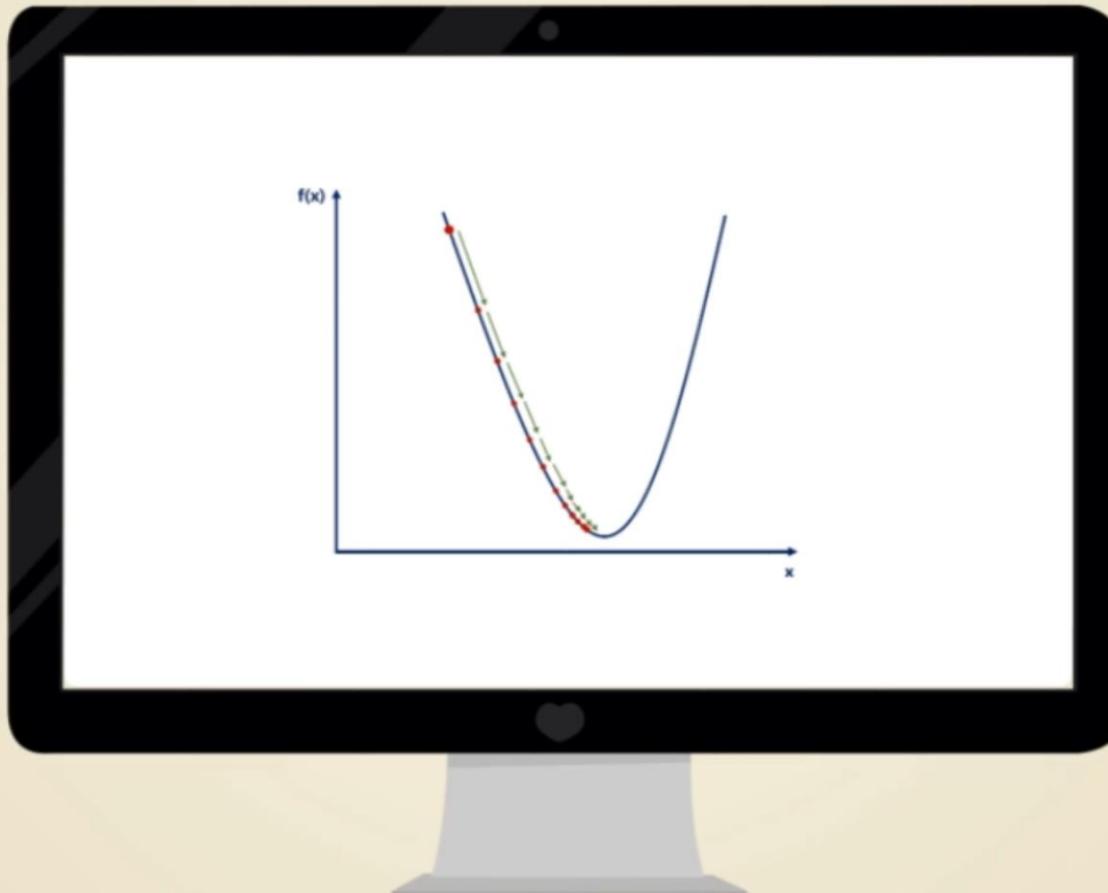
مدرس: محمد فزونی
عضو هیئت علمی دانشگاه گنبدکاووس

OPTIMIZATION

الگوریتم‌هایی که بکار می‌گیریم تا پارامترهای مدل‌مون رو در جهت بهبود عملکرد، تغییر بدیم

Algorithms we will use to vary our model's parameters

OPTIMIZATION



OPTIMIZATION

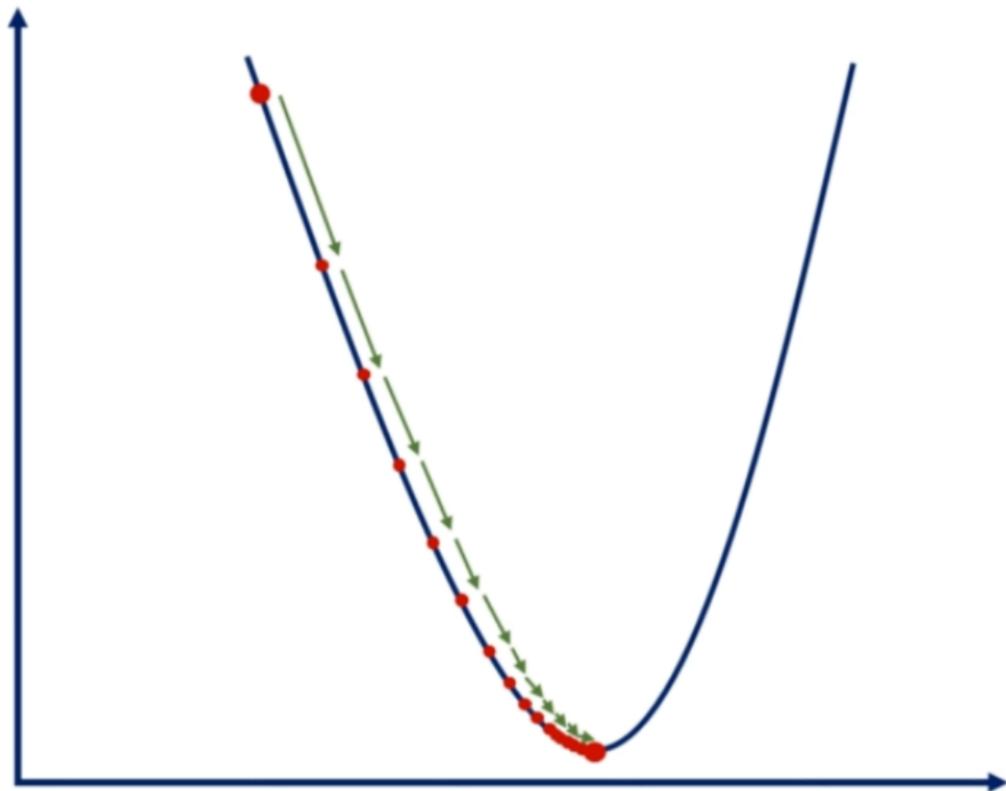
IN



گرادیان
کاهشی،
روی کاغذ،
نامبر وان
هست، ولی
در هنگام
بکارگیری
كمی دست و
پا شکسته
عمل می‌کنه



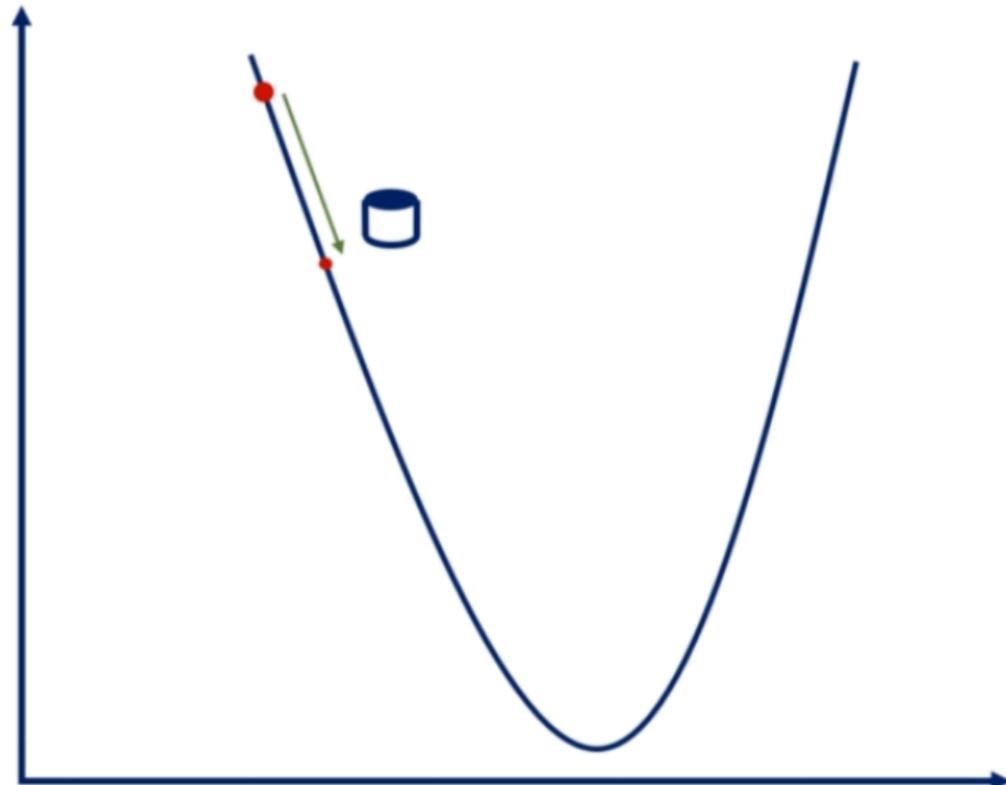
Gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

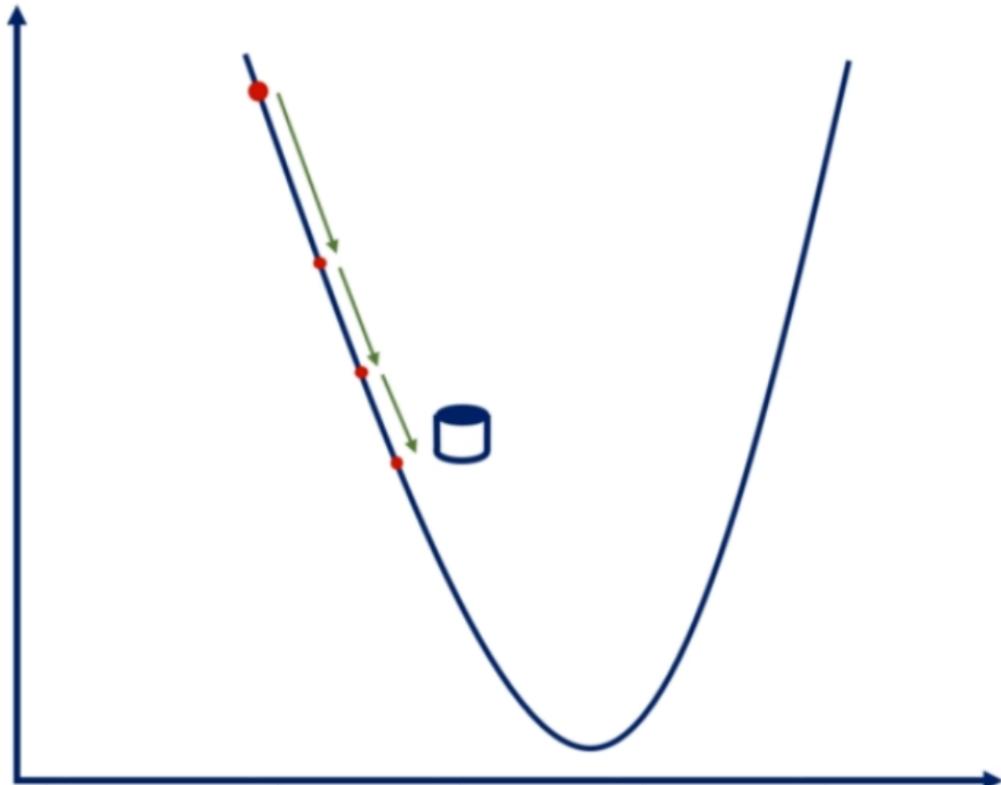
Gradient descent



x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
.
.
x_{n1}	x_{n2}	x_{n3}	x_{n4}	x_{n5}	x_{n6}	x_{n7}	x_{n8}

n x k

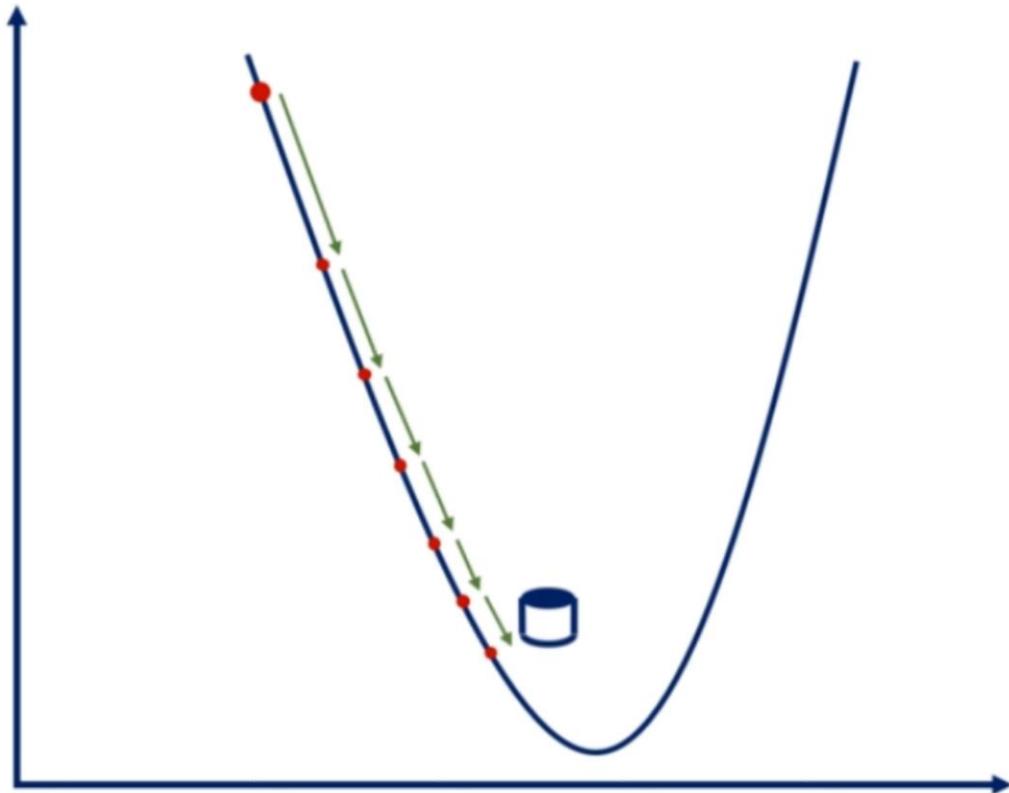
Gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

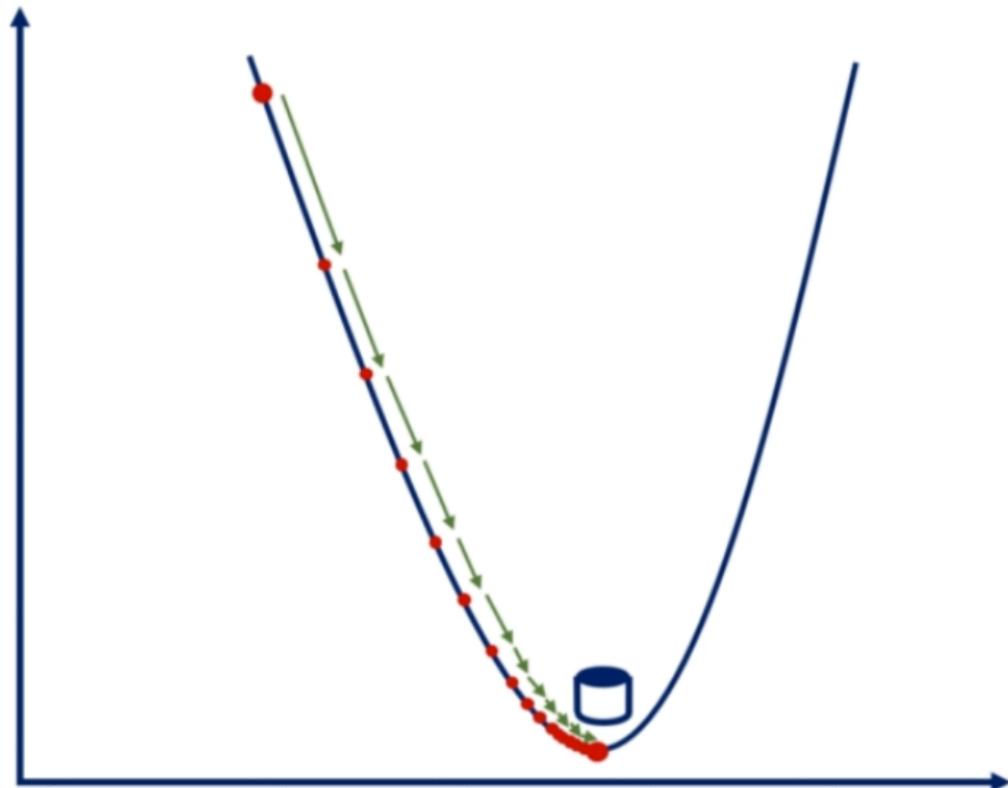
Gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

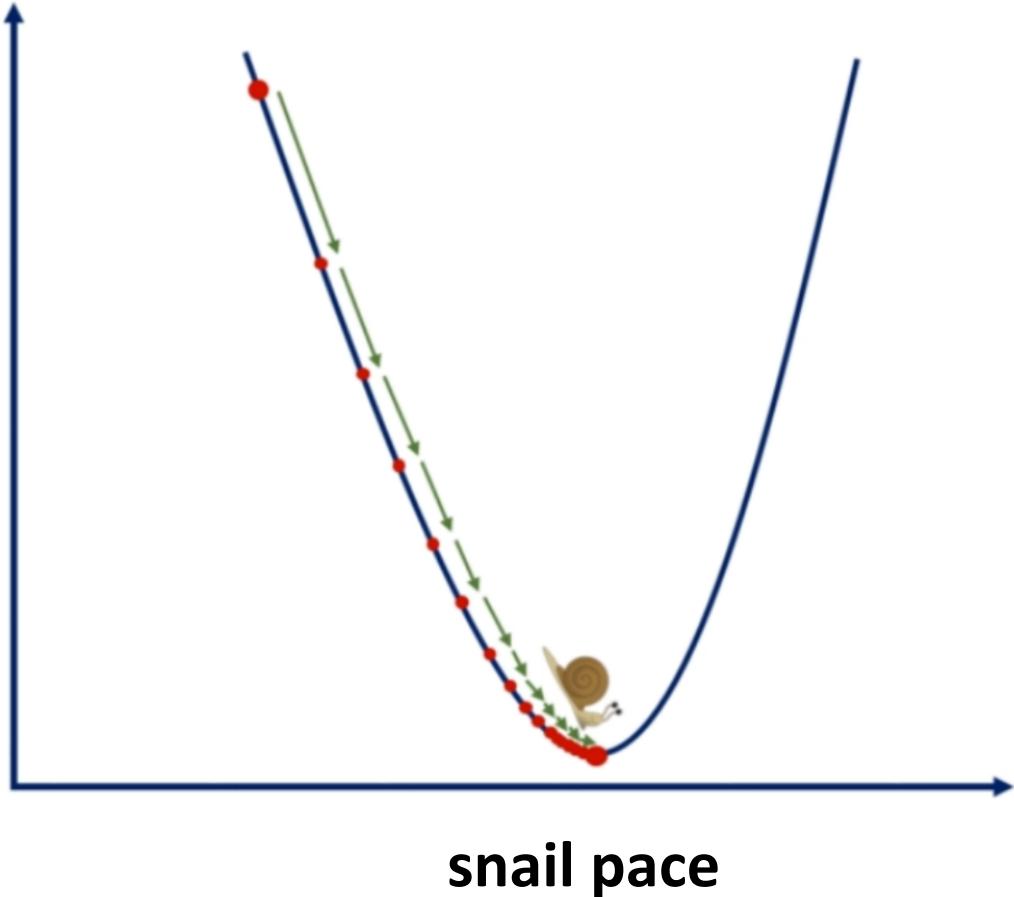
Gradient descent



x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
.
x_{n1}	x_{n2}	x_{n3}	x_{n4}	x_{n5}	x_{n6}	x_{n7}	x_{n8}

$n \times k$

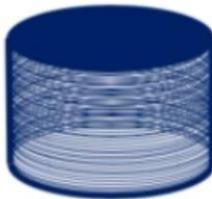
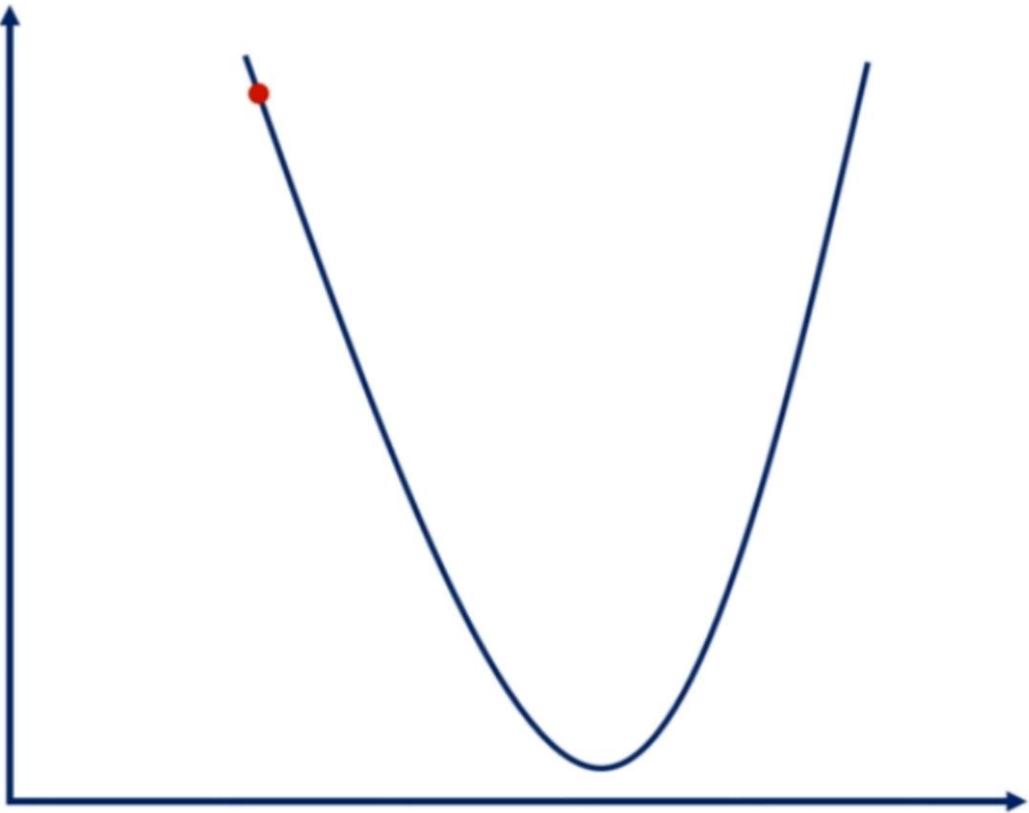
Gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

Stochastic gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

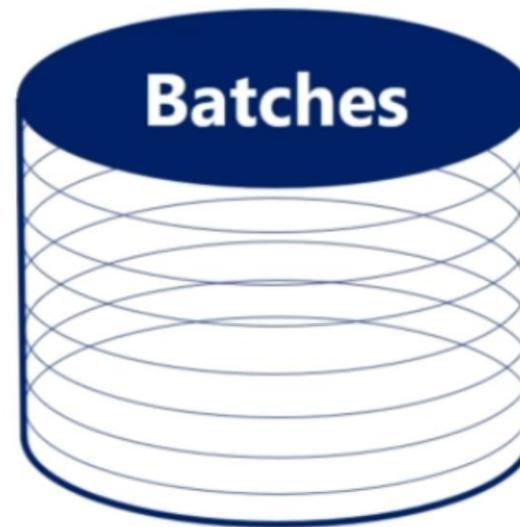


It works in the same way, but updates the weights many times inside a single epoch



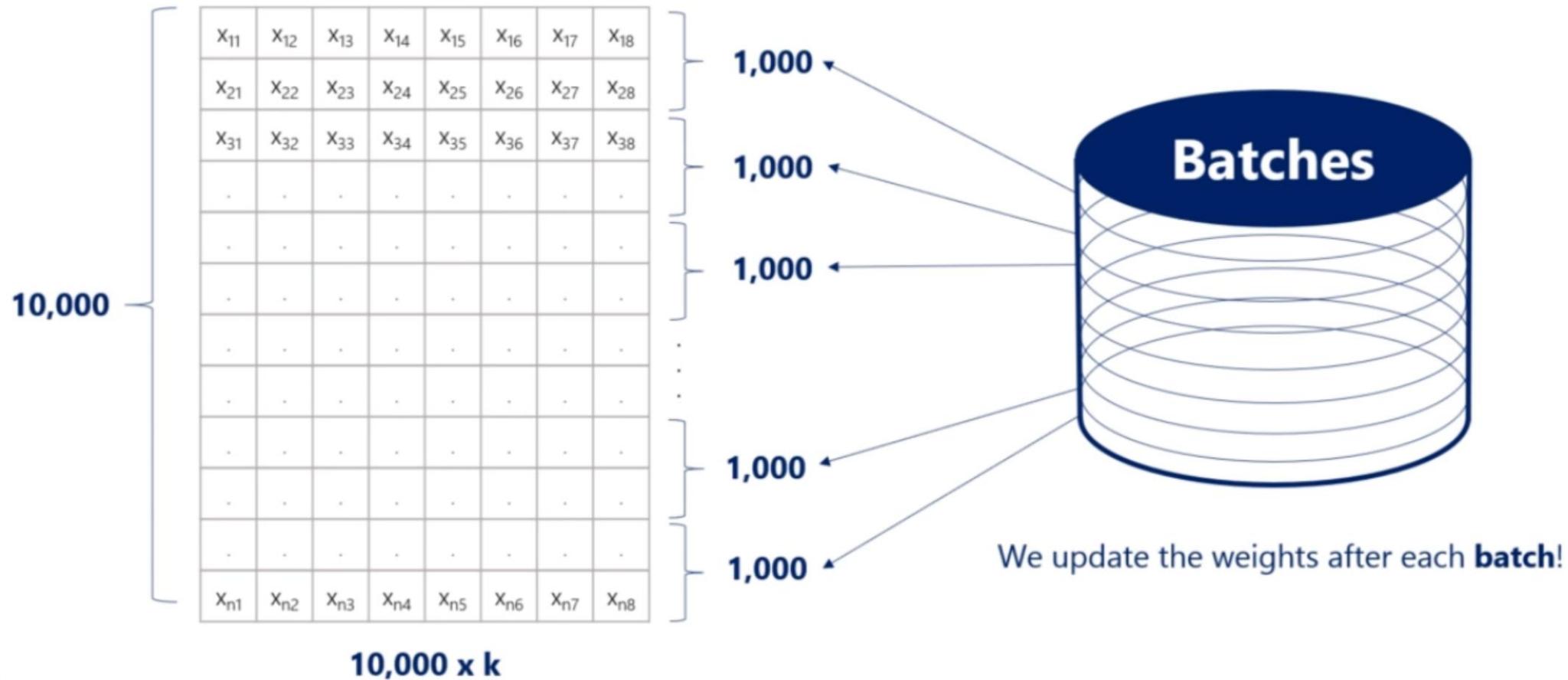
Batching

x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
.
.
.
.
.
.
.
x_{n1}	x_{n2}	x_{n3}	x_{n4}	x_{n5}	x_{n6}	x_{n7}	x_{n8}



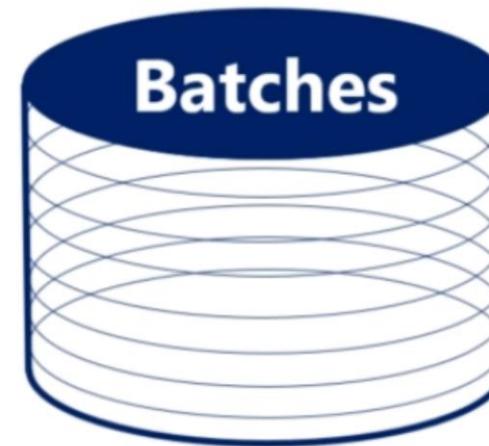
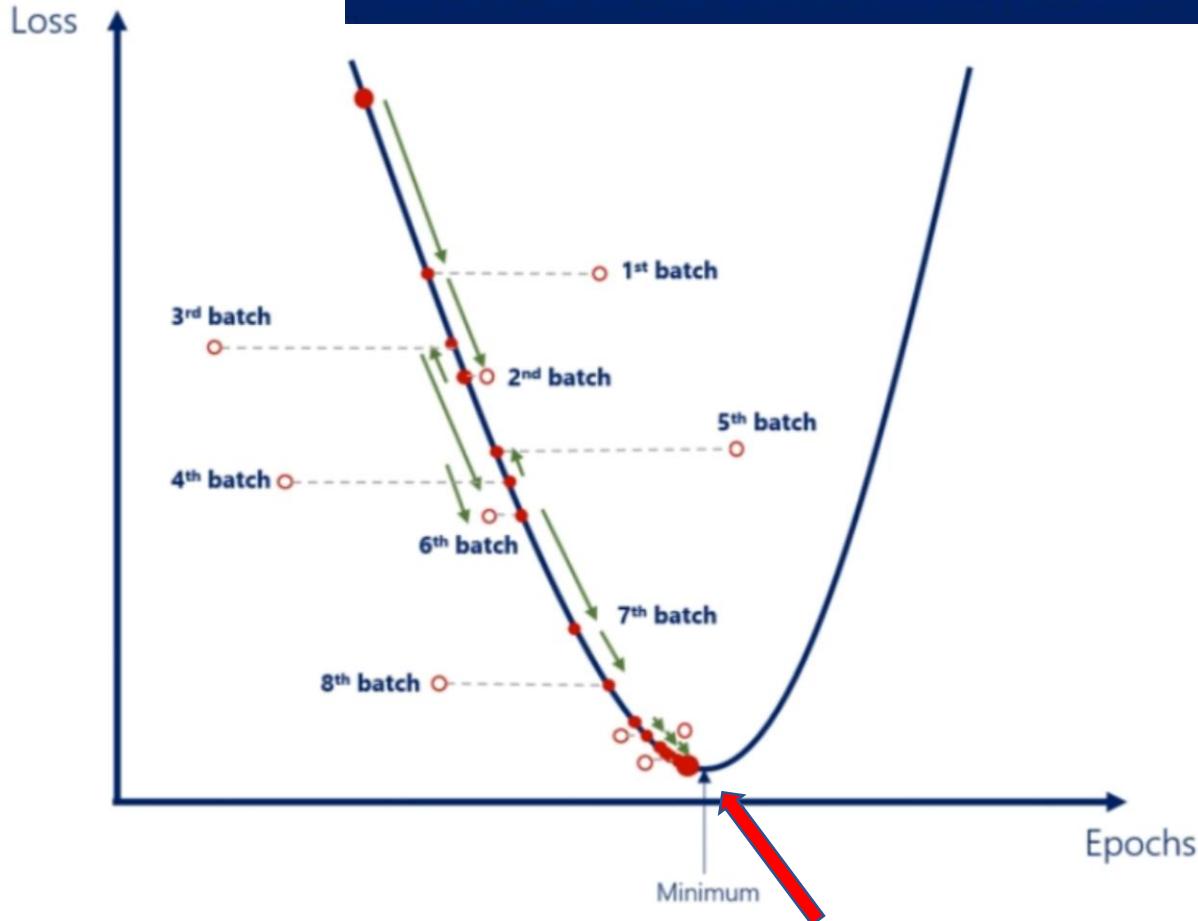
Batching: the process of splitting the dataset in n batches (mini-batches)

Batching

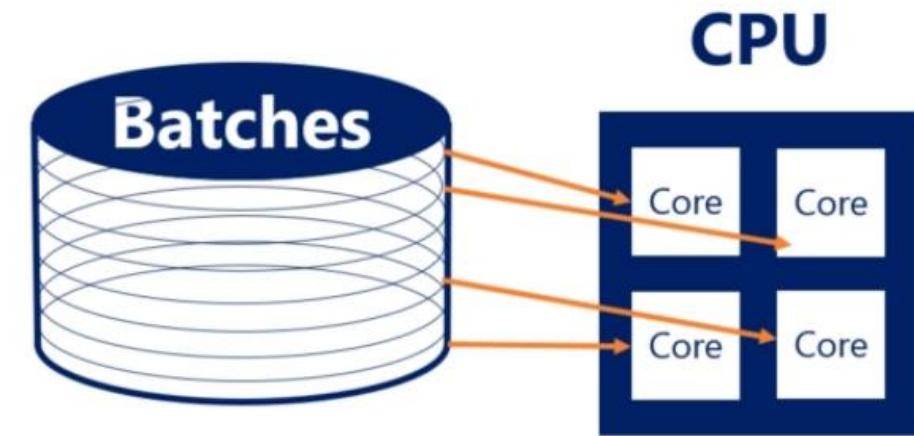
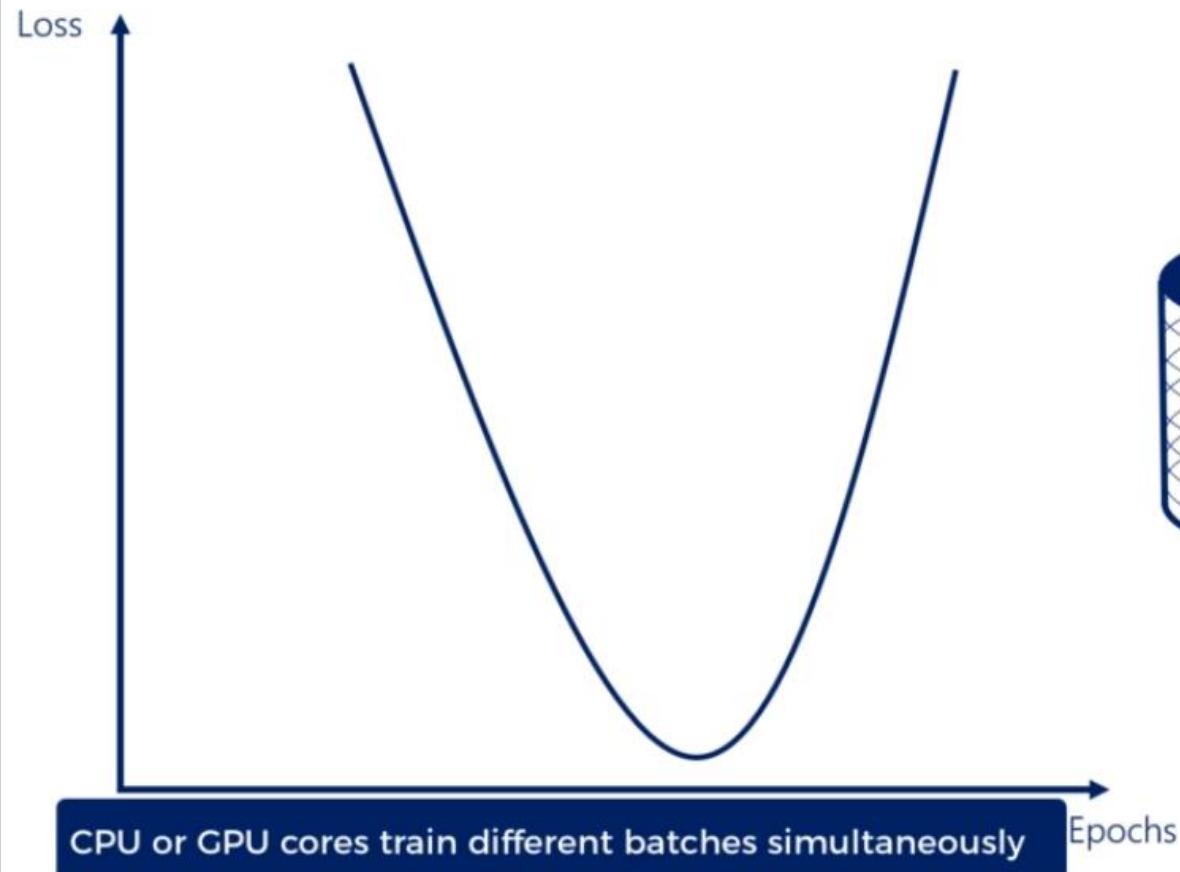


Batching

The SGD comes at a cost: it approximates things a bit



Batching





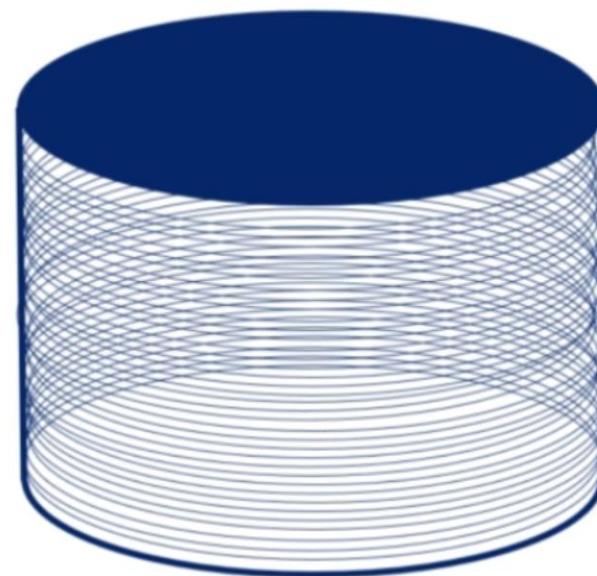
10,000

	X ₁₂	X ₁₃	X ₁₄	X ₁₅	X ₁₆	X ₁₇	X ₁₈
X ₂₁	X ₂₂	X ₂₃	X ₂₄	X ₂₅	X ₂₆	X ₂₇	X ₂₈
X ₃₁	X ₃₂	X ₃₃	X ₃₄	X ₃₅	X ₃₆	X ₃₇	X ₃₈
.
.
.
.
.
.
.
X _{n1}	X _{n2}	X _{n3}	X _{n4}	X _{n5}	X _{n6}	X _{n7}	X _{n8}

10,000 x k

1
1
1
.
.
.
.
.
.
.
.
.
.
.
.
1

We update the weights after each **sample!**

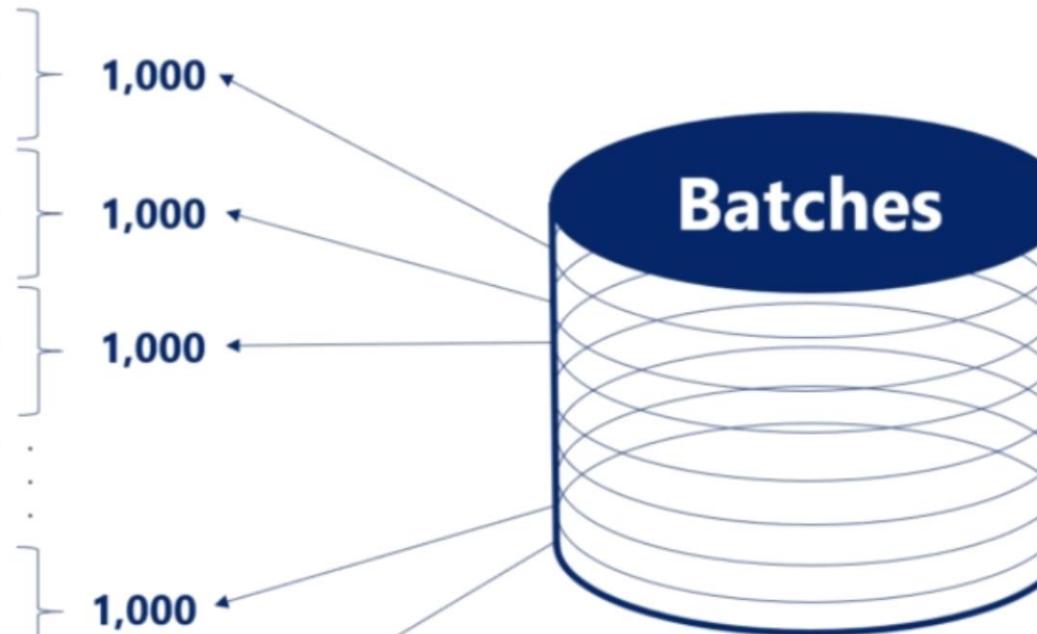




Mini-batch GD

10,000

X ₁₁	X ₁₂	X ₁₃	X ₁₄	X ₁₅	X ₁₆	X ₁₇	X ₁₈
X ₂₁	X ₂₂	X ₂₃	X ₂₄	X ₂₅	X ₂₆	X ₂₇	X ₂₈
X ₃₁	X ₃₂	X ₃₃	X ₃₄	X ₃₅	X ₃₆	X ₃₇	X ₃₈
.
.
.
.
.
X _{n1}	X _{n2}	X _{n3}	X _{n4}	X _{n5}	X _{n6}	X _{n7}	X _{n8}



We update the weights after each **batch**!

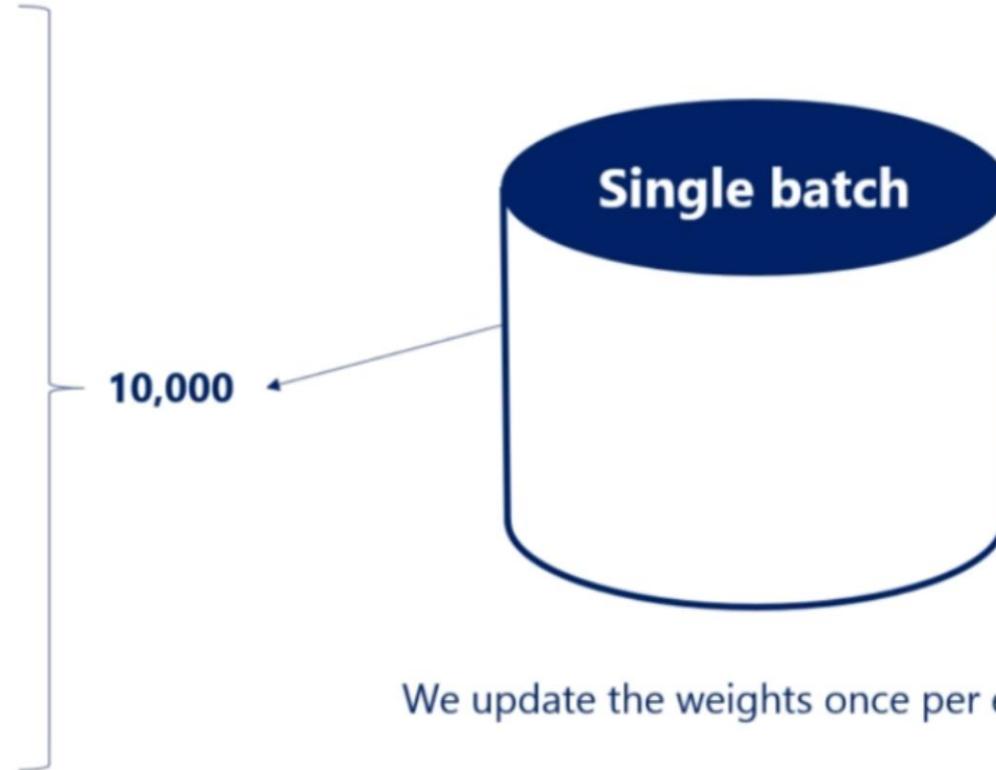
Practitioners refer to the mini-batch GD as SGD

Batch GD

10,000

x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
.
.
.
.
.
.
.
x_{n1}	x_{n2}	x_{n3}	x_{n4}	x_{n5}	x_{n6}	x_{n7}	x_{n8}

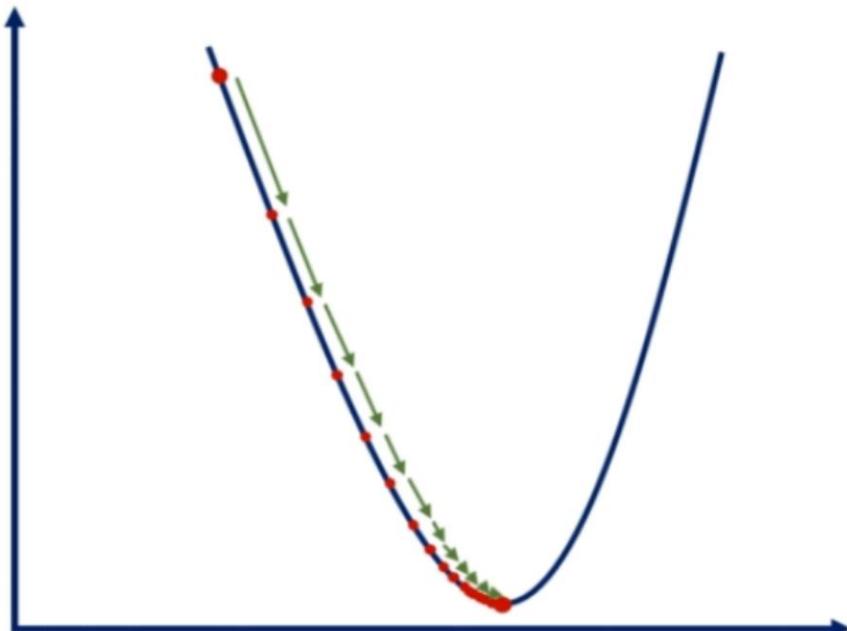
10,000 x k



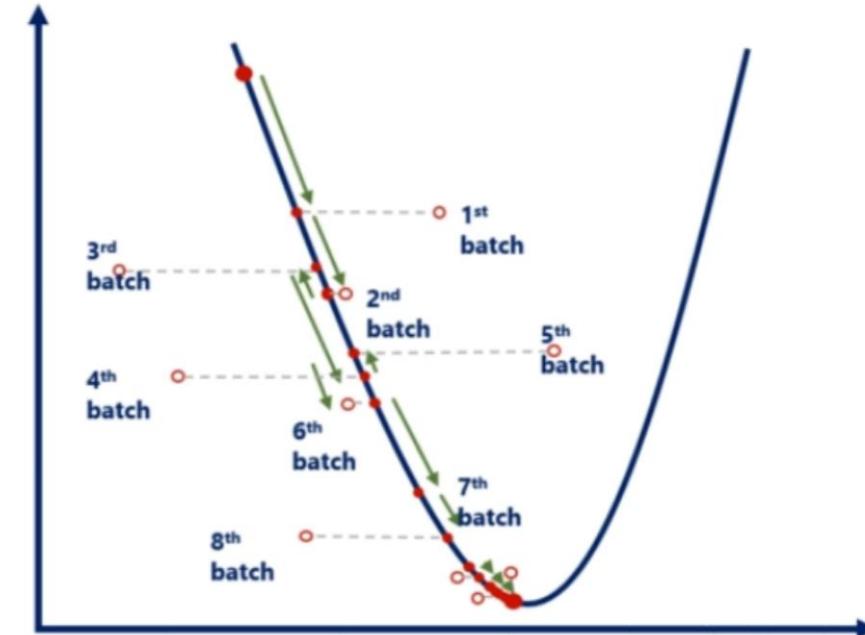
معایب گرایان کاهشی و نحوه‌ی رفع آنها

Gradient descent pitfalls

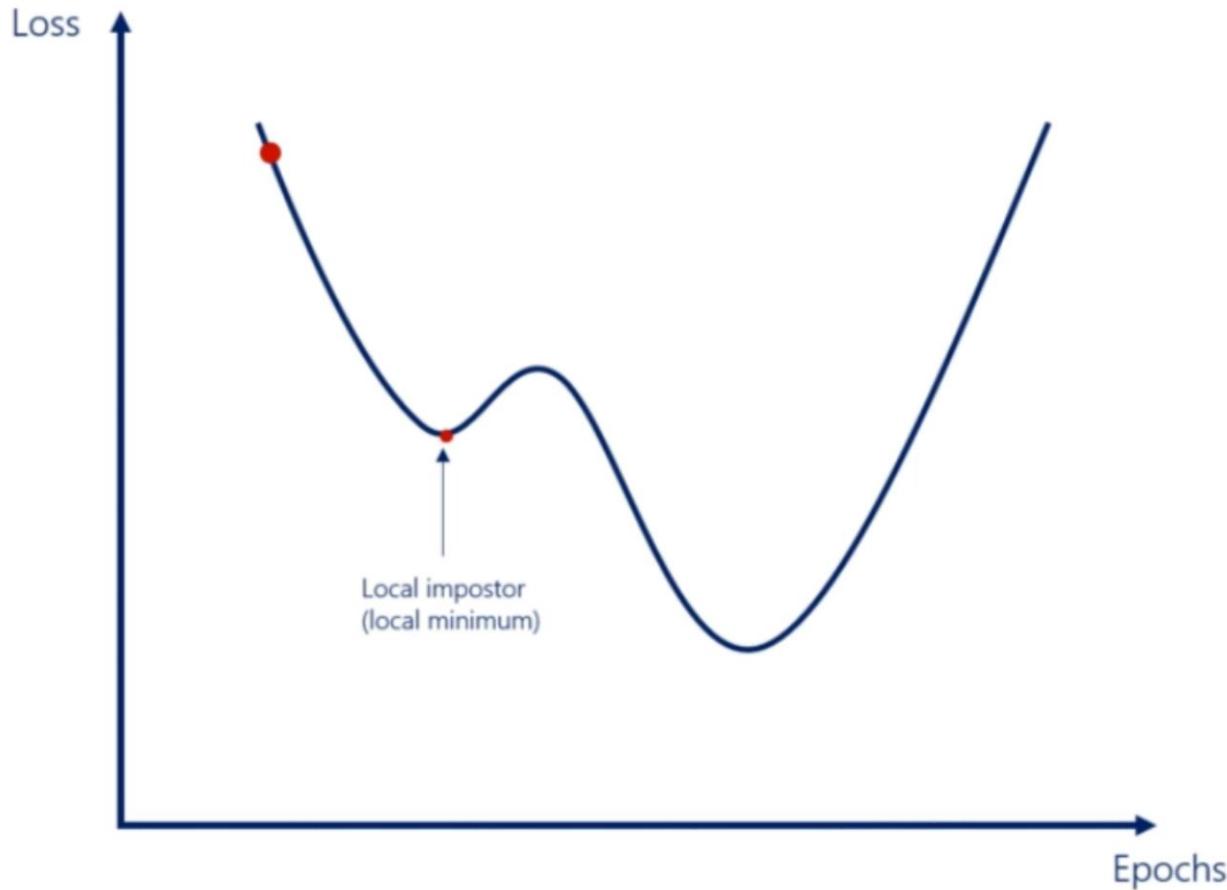
GD



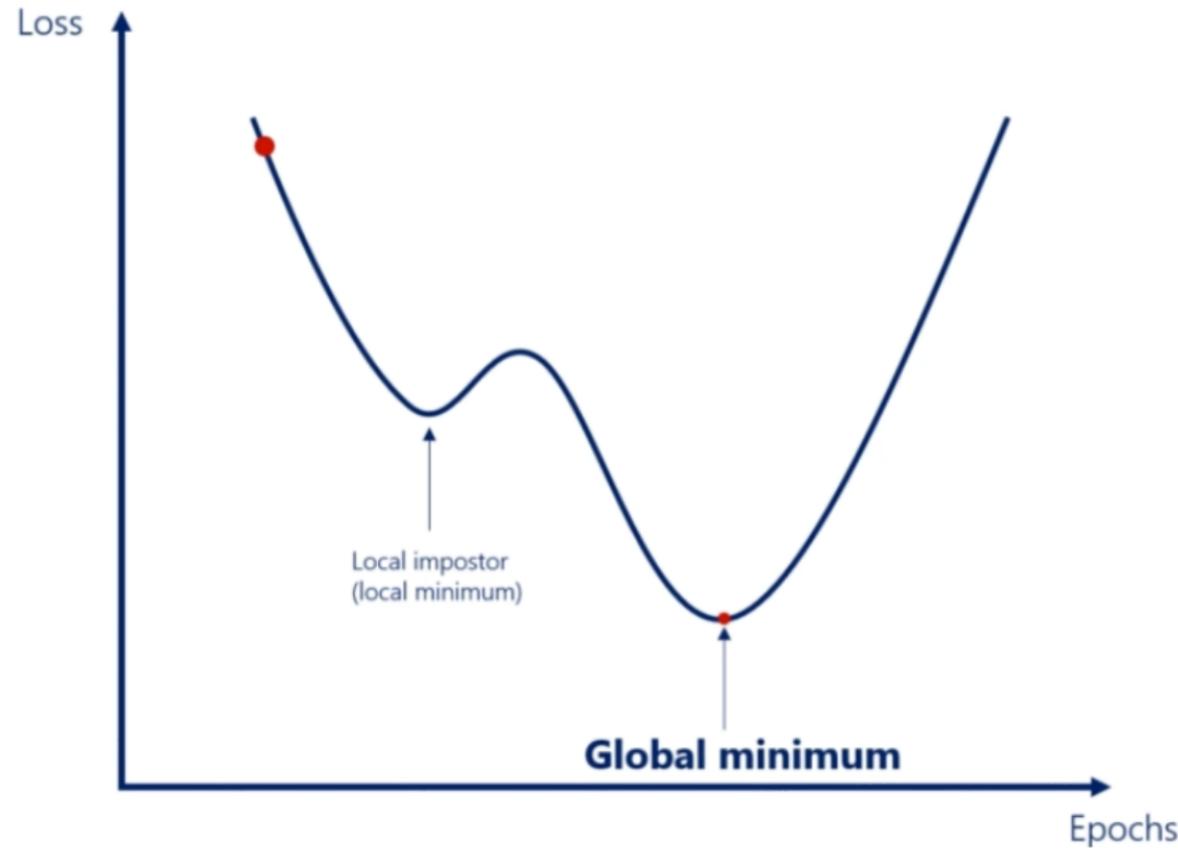
SGD



Gradient descent pitfalls

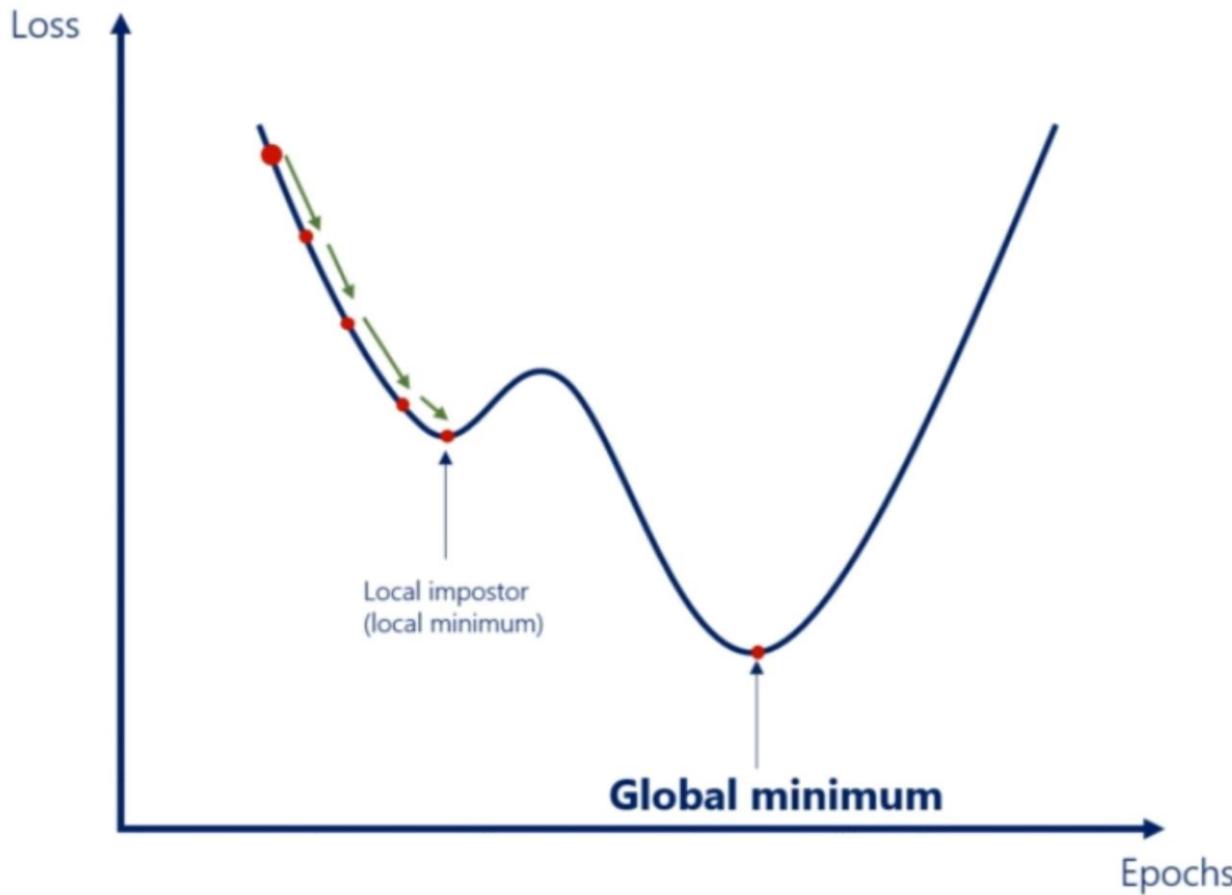


Gradient descent pitfalls

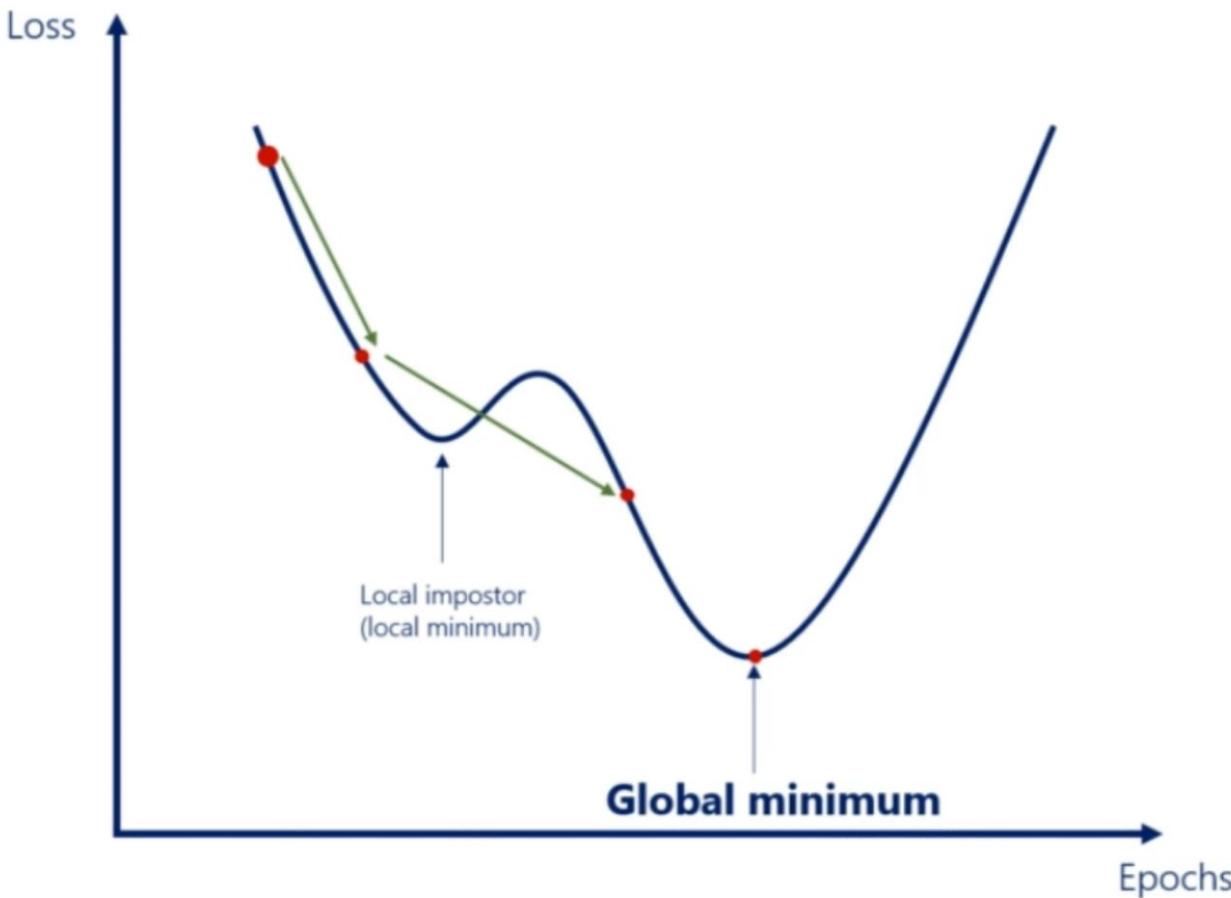


Each local minimum is a suboptimal solution to the optimization problem

Gradient descent pitfalls

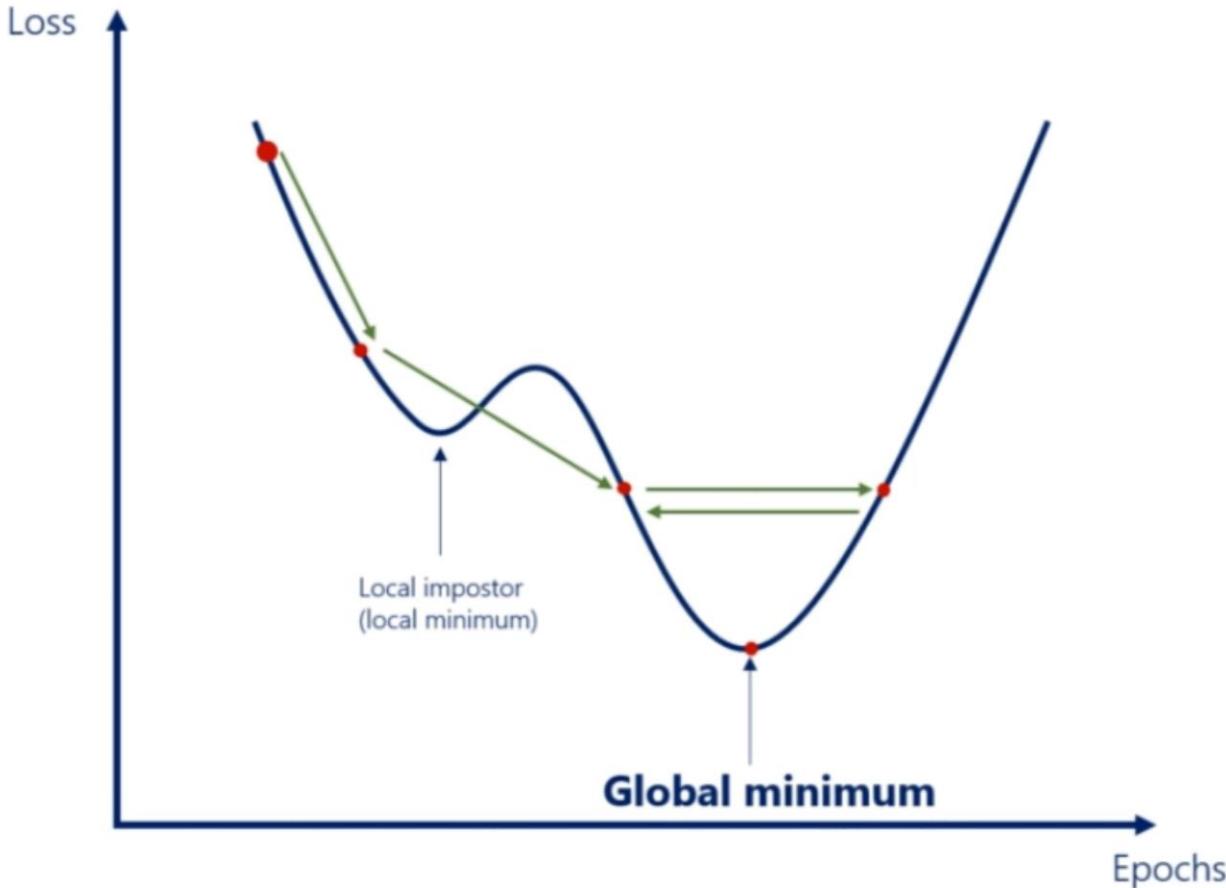


Gradient descent pitfalls



ممکن است
با افزایش
نرخ
یادگیری،
به سمت
 نقطه‌ی
 مینیمم
 حرکت کنیم

Gradient descent pitfalls

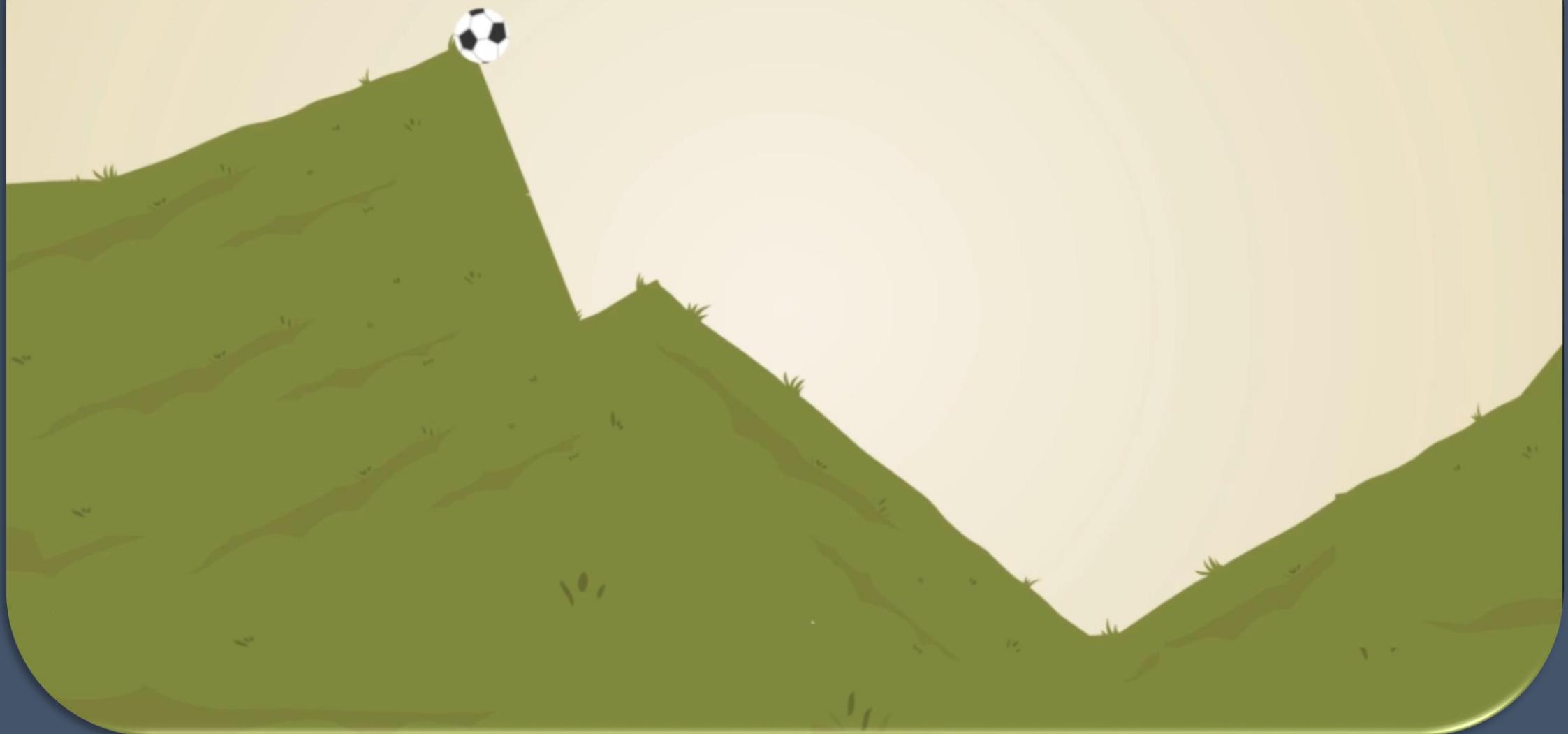


اما بزرگ
انتخاب
کردن نرخ
یادگیری،
خطر در
نوسان
بودن را
ایجاد
خواهد کرد

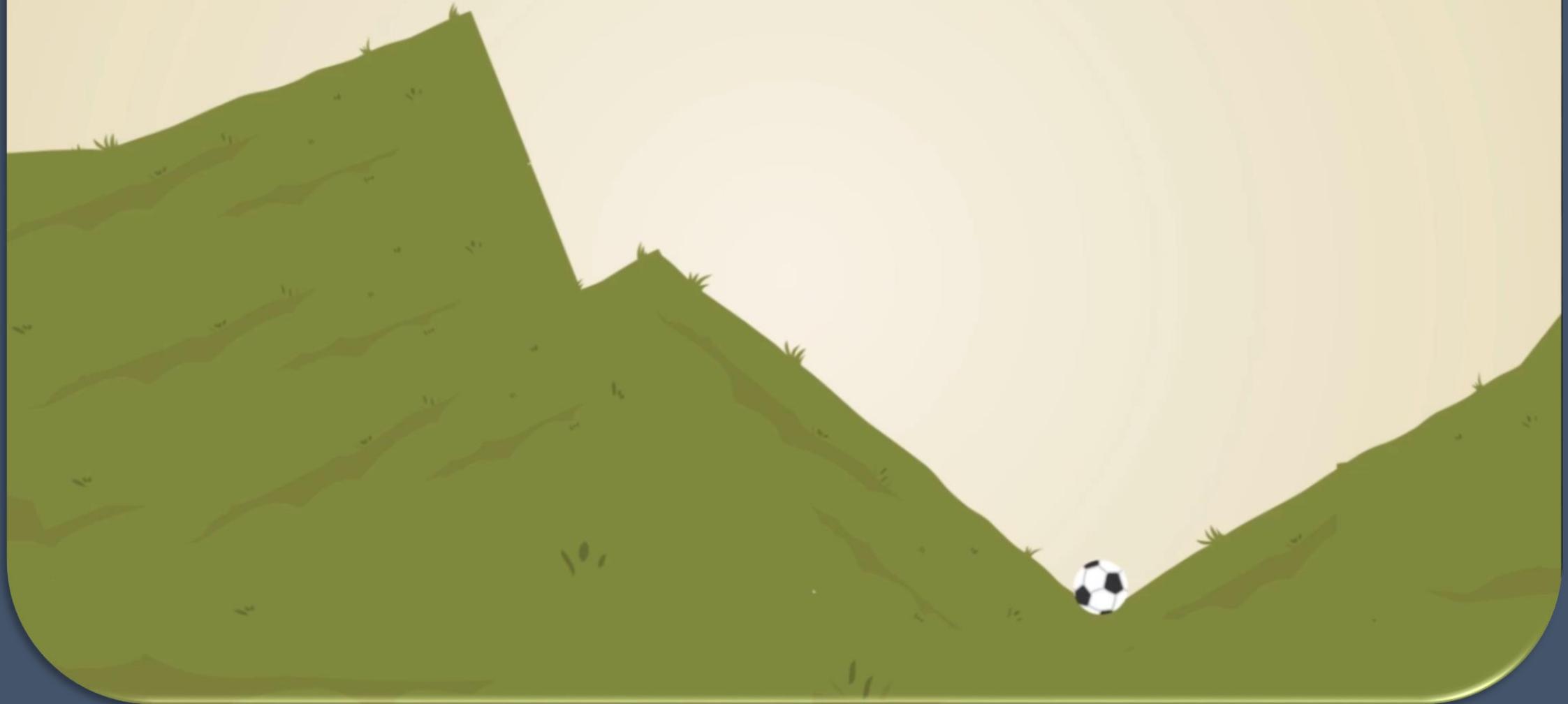
What is the remedy?

MOMENTUM

MOMENTUM



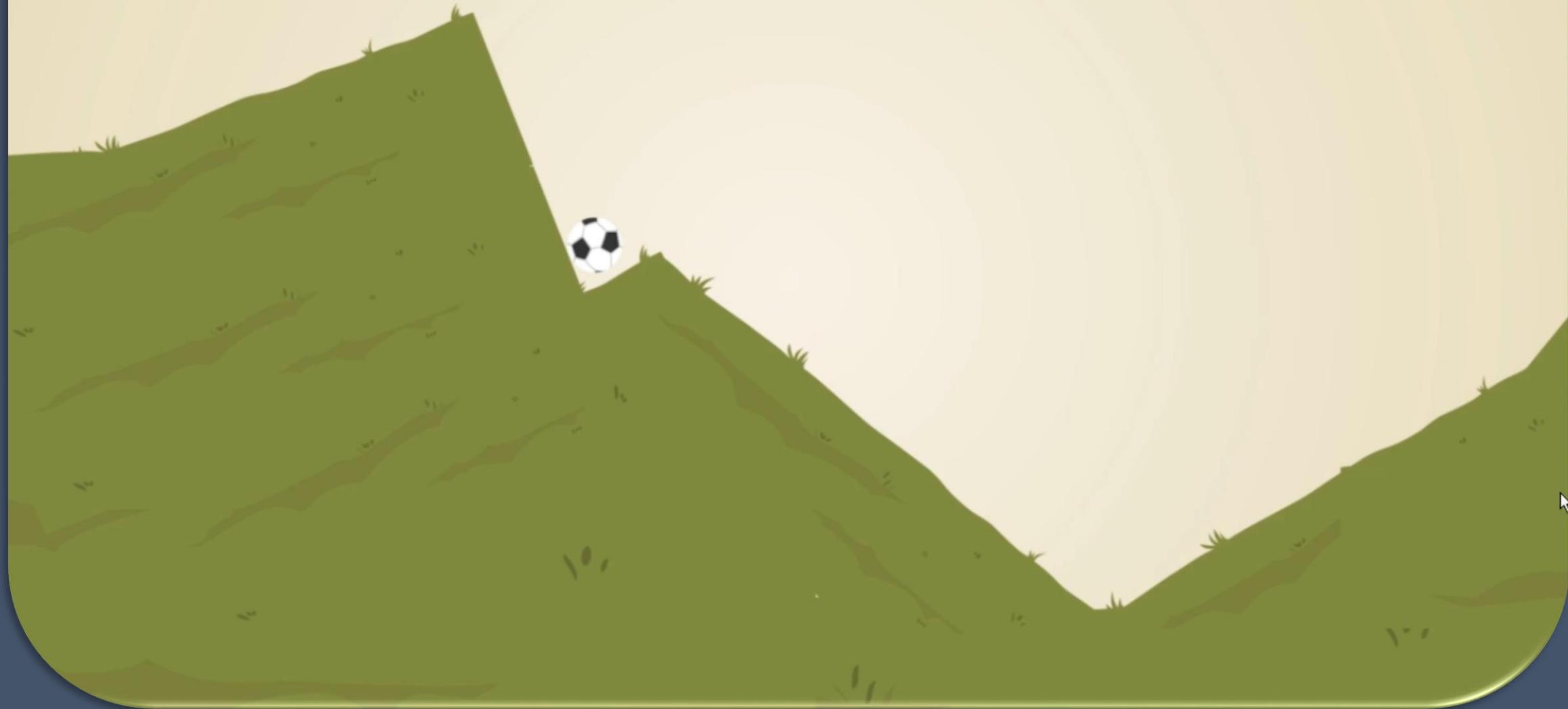
MOMENTUM



MOMENTUM



MOMENTUM

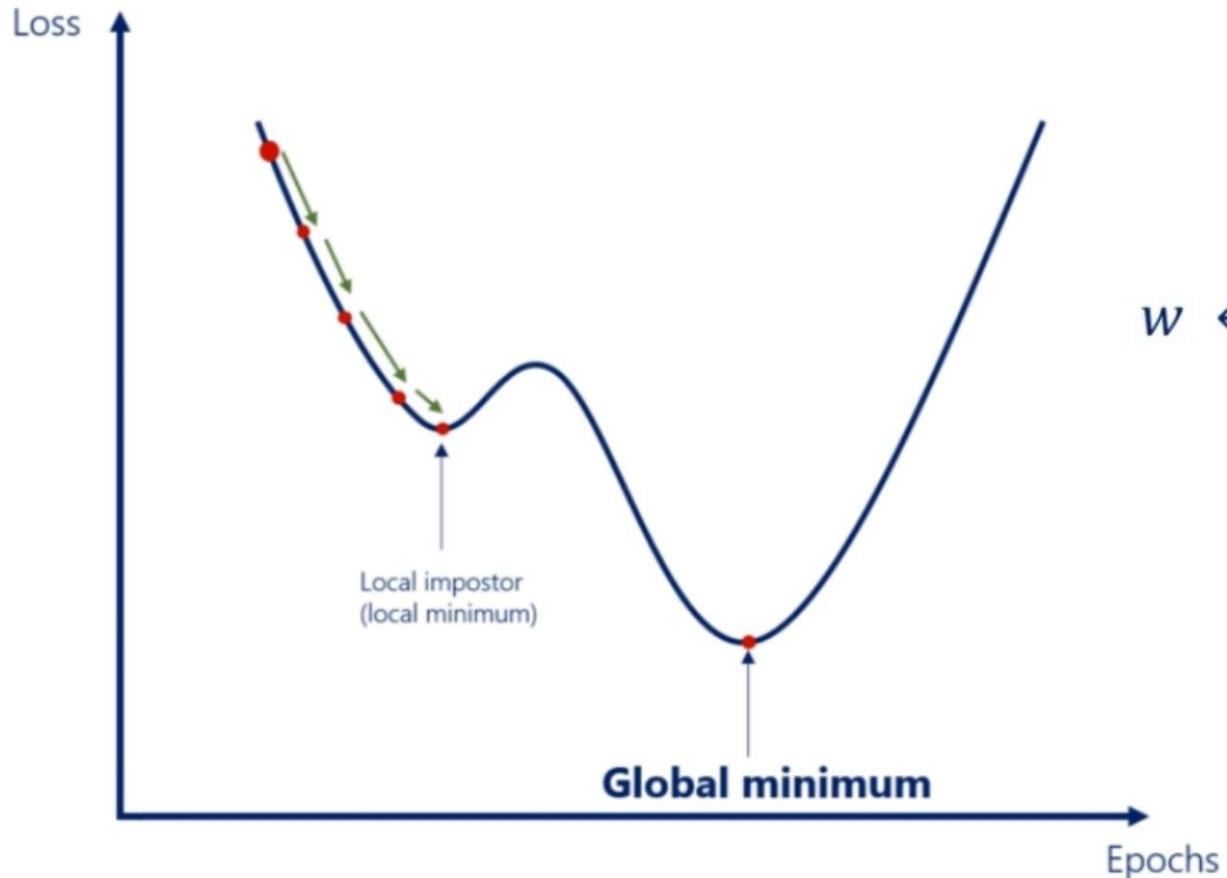


MOMENTUM



The momentum accounts for the fact that the ball actually going downhill

Momentum



$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

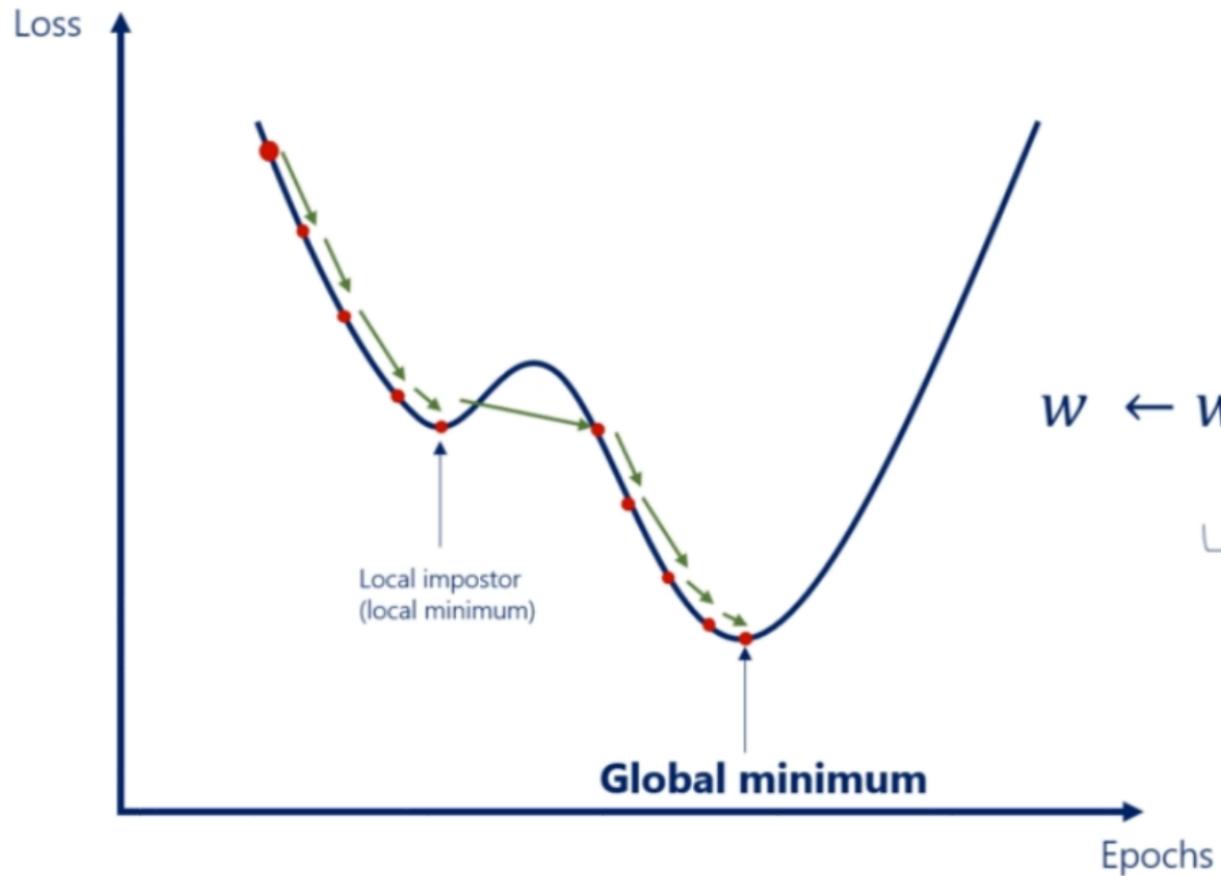
$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \eta \frac{\partial L}{\partial w}(t-1)$$

Current update

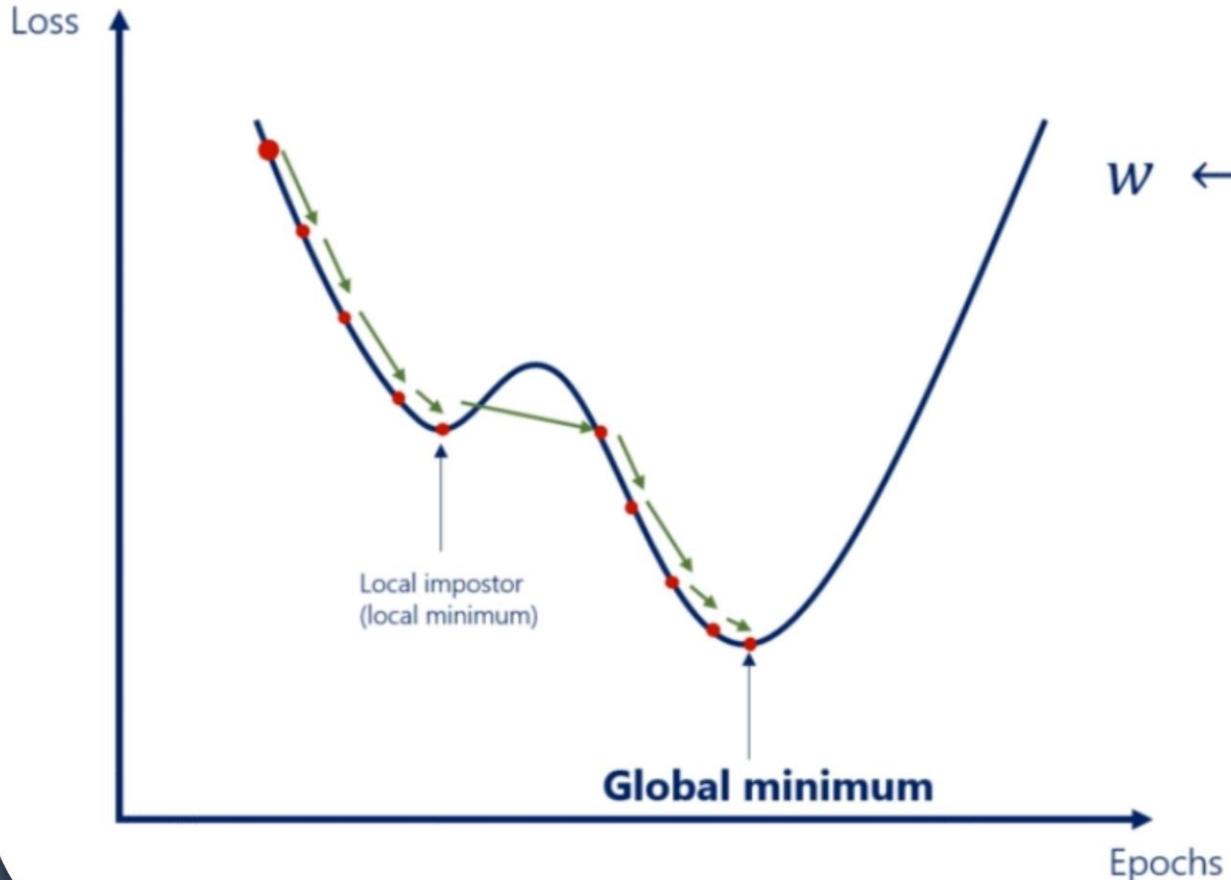
Update a moment ago

The best way to check how fast the ball rolls, is to check how fast it rolled a **moment ago**

Momentum



Momentum



$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \alpha \eta \frac{\partial L}{\partial w}(t-1)$$

Current update

Update a moment ago

$\alpha = 0.9$ is conventional

α is a hyperparameter

در ویدیوی بعدی نشون می‌دهیم که چطور
نرخ پادگیری بهینه‌ای را انتخاب کنیم تا مدل
طبق میل و خواسته‌ی ما جلو بره

Just Stay Tuned