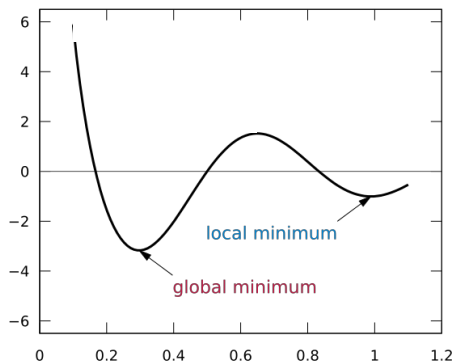# Convex Function vs. Nonconvex Function: A Little Bit Theory

**Shusen Wang**

# Global Extremum vs. Local Extremum

**Local Minimum** of a function $f(\mathbf{w})$
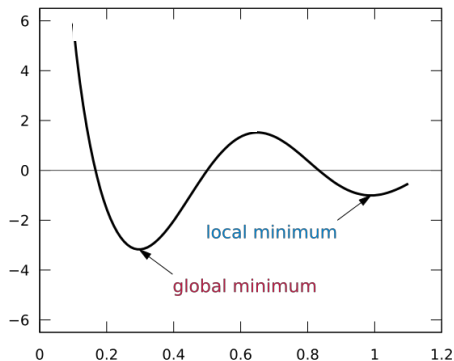
If $f(\mathbf{w}^{\star}) \leq f(\mathbf{w})$ for all $\mathbf{w}$ in a neighborhood of $\mathbf{w}^{\star}$, then $\mathbf{w}^{\star}$ is a **local minimum** of $f$.

**Global Minimum** of a function $f(\mathbf{w})$

If $f(\mathbf{w}^{\star}) \leq f(\mathbf{w})$ for all $\mathbf{w}$ in the domain of $f$, then $\mathbf{w}^{\star}$ is a **global minimum** of $f$.

- A global minimum is a local minimum.
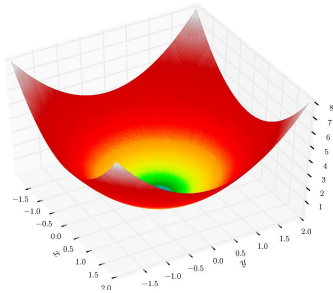- Global minimum may not be unique.

# Properties of Local Minimum

Properties of local minimum $\mathbf{w}^\star$:

1. The gradient at $\mathbf{w}^\star$, $\nabla f(\mathbf{w}^\star) \in \mathbb{R}^d$, is all-zeros.

2. The Hessian matrix at $\mathbf{w}^\star$, $\nabla^2 f(\mathbf{w}^\star) \in \mathbb{R}^{d \times d}$, is positive semidefinite (i.e., all of its $d$ eigenvalues are nonnegative.)
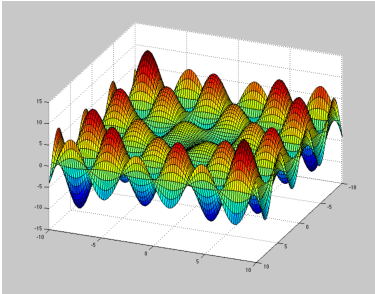
# Convex Function



**Graph of a convex function**

- **Convex function**: The line segment between any two points on the graph of the function lies above or on the graph

Properties of a convex function $f$:

1. Local minimum = global minimum.
2. $\nabla f(\mathbf{w}^\star) = \mathbf{0} \iff \mathbf{w}^\star$ is a global minimum.
3. The Hessian matrix $\nabla^2 f(\mathbf{w})$ is positive semi-definite everywhere.
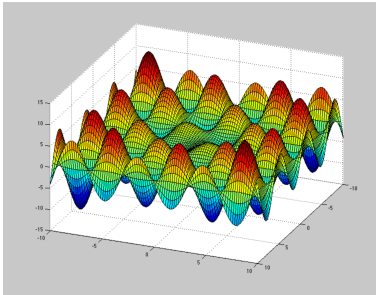
# Nonconvex Function



**Graph of a nonconvex function**

Properties:

1. Local minimum ~~=~~ global minimum.
2. $\nabla f(\mathbf{w}^\star) = \mathbf{0}$ ~~⬌~~ $\mathbf{w}^\star$ is a global minimum.
3. The Hessian matrix $\nabla^2 f(\mathbf{w})$ is positive semi-definite everywhere.

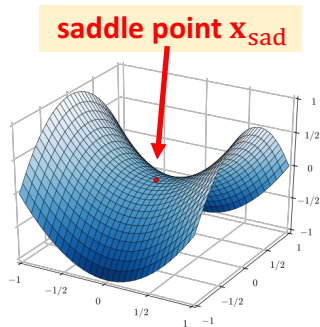# Global Minimum is Unlikely to Reach



**Graph of a nonconvex function**

- Number of local minima $\gg$ number of global minima.
- The final solution depends on the initialization.
- Reaching one of the global minima is very unlikely.

# Saddle Point



**saddle point $\mathbf{x}_{\text{sad}}$**

Graph of a nonconvex function

**Definition of saddle point:**

1. The gradient of $f$ at a saddle point is all-zeros: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$.
2. The Hessian matrix $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues.**.

# Saddle Point vs. Local Minimum

| **saddle point $\mathbf{w}_{\text{sad}}$** | **local minimum $\mathbf{w}^{\star}$** |
| --- | --- |
| • Gradient: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues**. | • Gradient: $\nabla f(\mathbf{w}^{\star}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}^{\star})$ does **not** have **negative eigenvalues**. |

# Saddle Point vs. Local Minimum

| saddle point $\mathbf{w}_{\text{sad}}$ | local minimum $\mathbf{w}^\star$ |
|---|---|
| • Gradient: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues**. | • Gradient: $\nabla f(\mathbf{w}^\star) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}^\star)$ does **not** have **negative eigenvalues**. |

• Full gradient descent stops at either a saddle point or a local minimum.
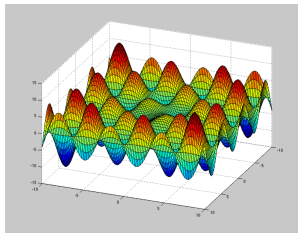
# Saddle Point vs. Local Minimum

| saddle point $\mathbf{w}_{\text{sad}}$ | local minimum $\mathbf{w}^\star$ |
|---|---|
| • Gradient: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues**. | • Gradient: $\nabla f(\mathbf{w}^\star) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}^\star)$ does **not** have **negative eigenvalues**. |

• Full gradient descent stops at either a saddle point or a local minimum.



• In 2D, numbers of saddle points and local minimum are comparable.
• It is not true in high-dim.

# Saddle Point vs. Local Minimum

| **saddle point $\mathbf{w}_{\text{sad}}$** | **local minimum $\mathbf{w}^{\star}$** |
|---|---|
| • Gradient: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$.<br>• Hessian: $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues**. | • Gradient: $\nabla f(\mathbf{w}^{\star}) = \mathbf{0}$.<br>• Hessian: $\nabla^2 f(\mathbf{w}^{\star})$ does **not** have **negative eigenvalues**. |

- Full gradient descent stops at either a saddle point or a local minimum.

- In high dim, the number of saddle points is much larger than local minima.
  - The Hessian has $d$ eigenvalues, each of which can be positive or negative.
  - ➡ $2^d$ combinations.
  - One out of the $2^d$ combinations corresponds to local minima.
  - $2^d - 2$ combinations corresponds to saddle points.

# Saddle Point vs. Local Minimum

| **saddle point $\mathbf{w}_{\text{sad}}$** | **local minimum $\mathbf{w}^{\star}$** |
|---|---|
| • Gradient: $\nabla f(\mathbf{w}_{\text{sad}}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}_{\text{sad}})$ has **both positive** and **negative eigenvalues**. | • Gradient: $\nabla f(\mathbf{w}^{\star}) = \mathbf{0}$. <br> • Hessian: $\nabla^2 f(\mathbf{w}^{\star})$ does **not** have **negative eigenvalues**. |

- Full gradient descent stops at either a saddle point or a local minimum.

- In high dim, the number of saddle points is much larger than local minima.

- If a neural net is optimized by the full gradient descent, it will converge to a saddle point.

# Be Careful When Optimizing a Nonconvex Function

**Be careful about the initialization!**

- Bad initialization results in convergence to bad regions.
  - Because of the nonconvexity, global minimum cannot be attained.
- Heuristics:
  - The trainable parameters (e.g., the filters of ConvNet) are randomly initialized with proper scaling.
  - Bad scaling leads to terrible results.
  - All-zero and all-one initializations are bad ideas.

# Be Careful When Optimizing a Nonconvex Function

**Be careful about the initialization!**

**Be careful about the optimization algorithm!**

- Full gradient descent will be stuck in a saddle point.
  - Because the gradient is near zero when approaching the saddle point.
- Stochastic gradient descent (SGD) can escape the saddle points.
  - Because it is random and noisy.

# Be Careful When Optimizing a Nonconvex Function
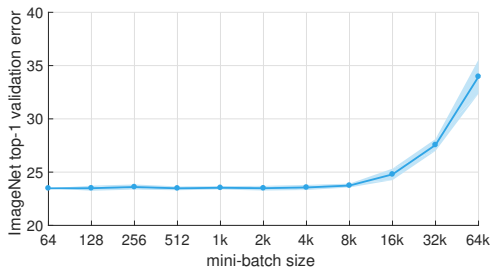
**Be careful about the initialization!**

**Be careful about the optimization algorithm!**

**Be careful about the batch size!**

- Small batch size does not make full use of GPUs.
- For parallel computing with multiple GPUs, larger batch size ➜ lower per-epoch runtime.
- Large batch size, e.g., $10K$, may result in bad generalization.
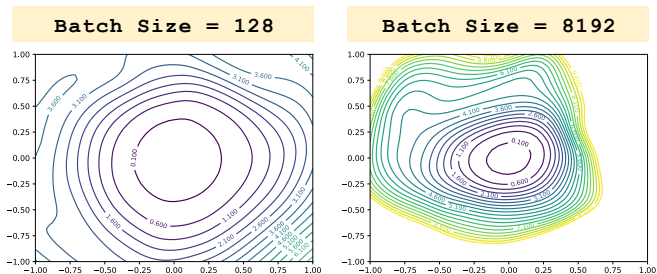
# ... More about the Batch Size

- Batch size larger than $8K$ results in poor generalization.
- Large batch size is good for computation.
- Lots of tricks are required in *large batch training*.



The figure is from the paper "*Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour*"

# … **More about the Batch Size**

- Researchers' conjecture:
  - Small batch size ➜ flat local minima; Big batch size ➜ shape local minima.
  - Flat local minima generalizes better (on the test set).



The figure is from paper https://arxiv.org/abs/1712.09913

# … **More about the Batch Size**

- There are papers supportive of small batch training, e.g.,
  https://arxiv.org/pdf/1804.07612.pdf

The presented results confirm that using small batch sizes achieves the best training stability and generalization performance, for a given computational cost, across a wide range of experiments. In all cases the best results have been obtained with batch sizes $m = 32$ or smaller, often as small as $m = 2$ or $m = 4$.

**Yann LeCun**
@ylecun

Follow

Training with large minibatches is bad for your health.
More importantly, it's bad for your test error.
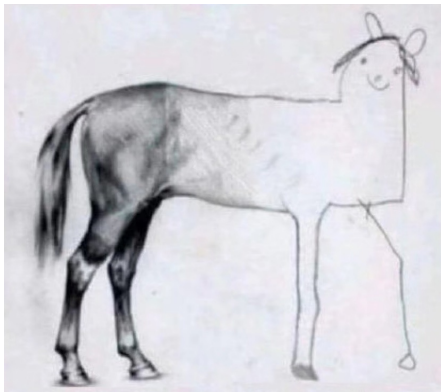Friends dont let friends use minibatches larger than 32. arxiv.org/abs/1804.07612

1:00 PM - 26 Apr 2018

**447** Retweets **1,173** Likes

# Do Not Believe Deep Learning Theories Blindly



Empirical study

Explanations