# Asymmetric Valleys: Beyond Sharp and Flat Local Minima

Haowei He | Gao Huang | Yang Yuan

IIIS | Department of Automation,
Tsinghua University

**October 26, 2019**

# Background: Local Landscape

A good understanding of local landscape is important for
1. Designing better optimization method
2. Explaining when and how a deep network achieves good generalization performance
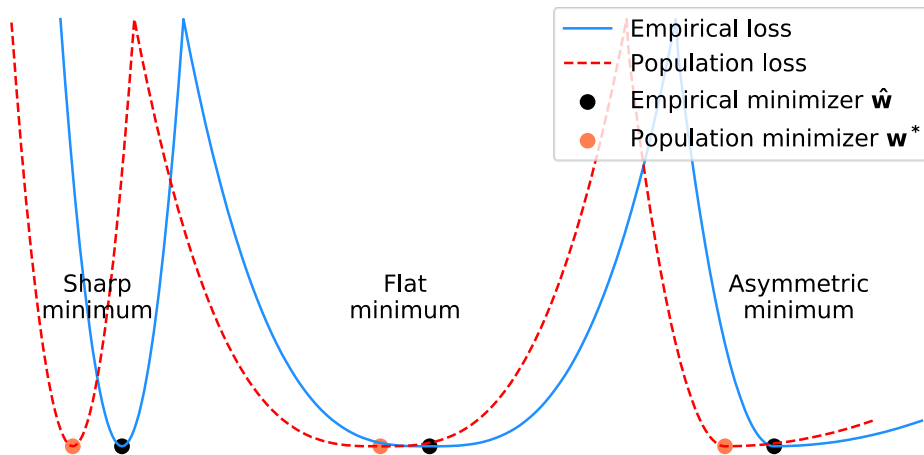
# Previous Work & Our Proposal

So what is a good local landscape for generalization?
## 1. Sharp or Flat
Flat minimum generalize better.*
## 2. Re-parameterization
## 3. Asymmetric Valley



*On large-batch training for deep learning: Generalization gap and sharp minima. ICLR, 2017.
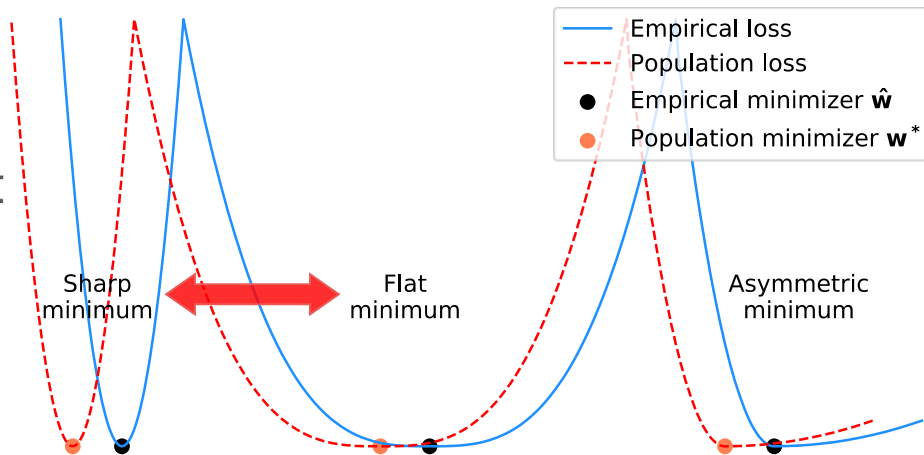
# Previous Work & Our Proposal

So what is a good local landscape for generalization?

1. Sharp or Flat

2. Re-parameterization

Flat and sharp minimum can convert to each other.*

3. Asymmetric Valley



Empirical loss
Population loss
Empirical minimizer $\hat{\mathbf{w}}$
Population minimizer $\mathbf{w}^*$

Sharp minimum　　Flat minimum　　Asymmetric minimum

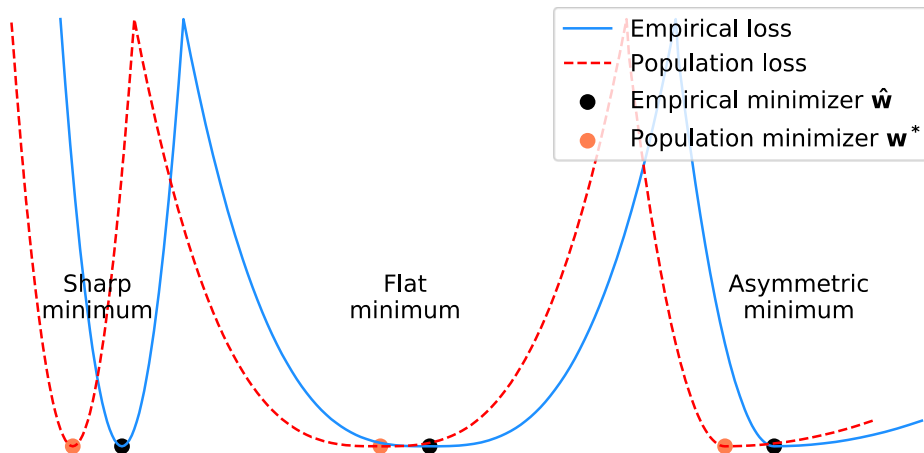*Sharp minima can generalize for deep nets. ICML, 2017.

# Previous Work & Our Proposal

So what is a good local landscape for generalization?

1. Sharp or Flat
2. Re-parameterization
3. Asymmetric Valley

**Our work.***
Minimum with asymmetric direction.
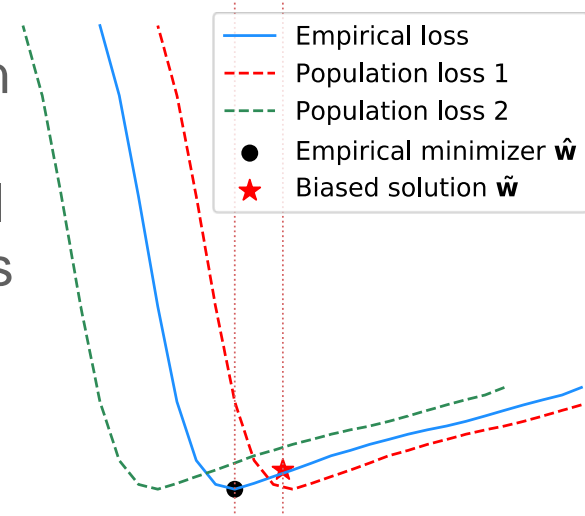Biased solution generalize better.

* Asymmetric Valleys: Beyond sharp and flat local minimum. NeurIPS, 2019.

# What is asymmetric valley?
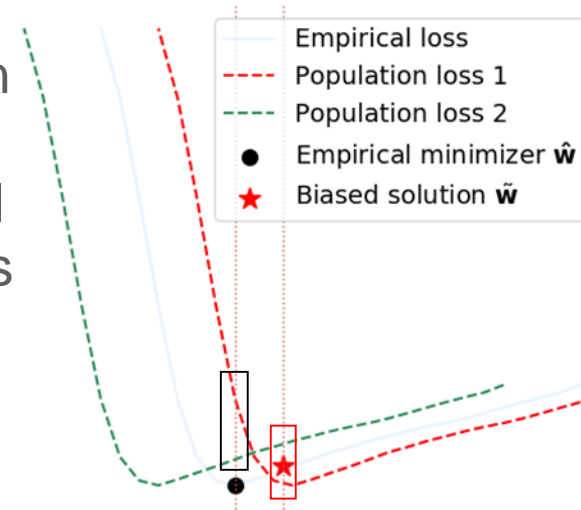
Intuitively illustration:

1. Loss grows fast on one side and slowly on the other side.

2. With a random shift between the empirical and population loss, the red star solution has lower population loss in expectation.



- Empirical loss
- Population loss 1
- Population loss 2
- ● Empirical minimizer $\hat{\mathbf{w}}$
- ★ Biased solution $\tilde{\mathbf{w}}$

# What is asymmetric valley?
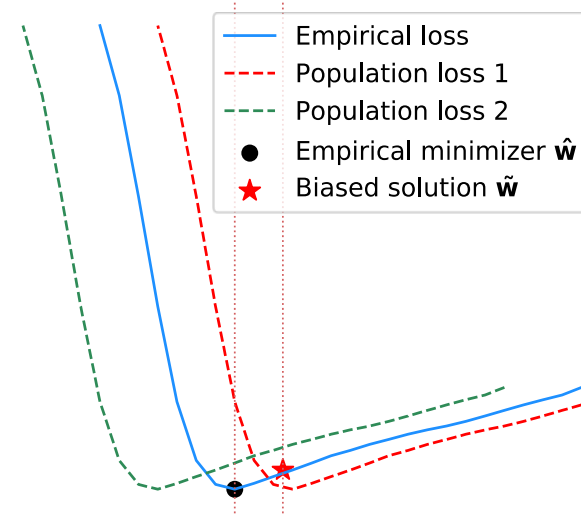
Intuitively illustration:

1. Loss grows fast on one side and slowly on the other side.

2. With a random shift between the empirical and population loss, the <span style="color:red">red star</span> solution has lower population loss in expectation.



Empirical loss
Population loss 1
Population loss 2
● Empirical minimizer $\hat{w}$
★ Biased solution $\tilde{w}$
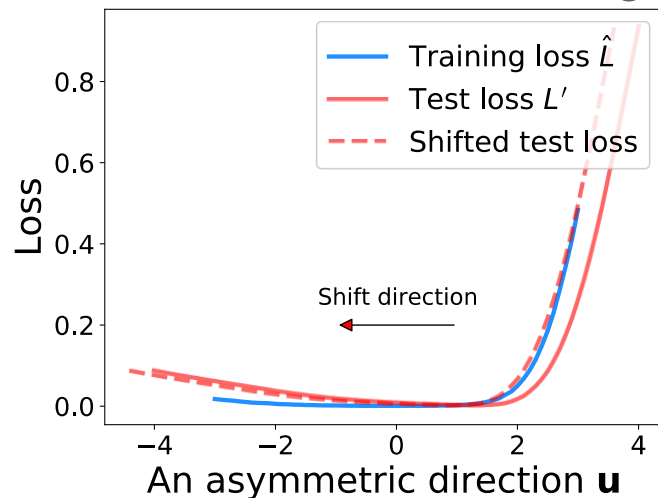
# What is asymmetric valley?

Two interesting implications:
1. converging to *which* local minimum may not be critical. However, it matters a lot *where* the solution locates.
2. the solution with lowest *a priori* generalization error is not necessarily the minimizer of the training loss.



Empirical loss
Population loss 1
Population loss 2
● Empirical minimizer $\hat{w}$
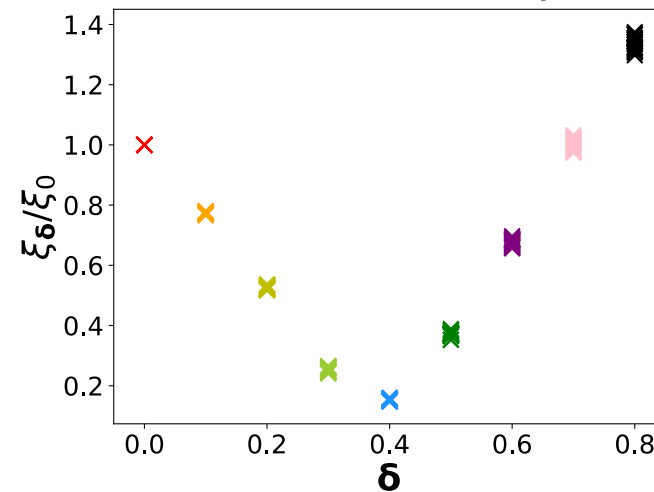★ Biased solution $\tilde{w}$

# Assumptions and Verification

Random shift assumption: assume a random symmetric horizontal shift between training loss and test loss with 0 expectation.



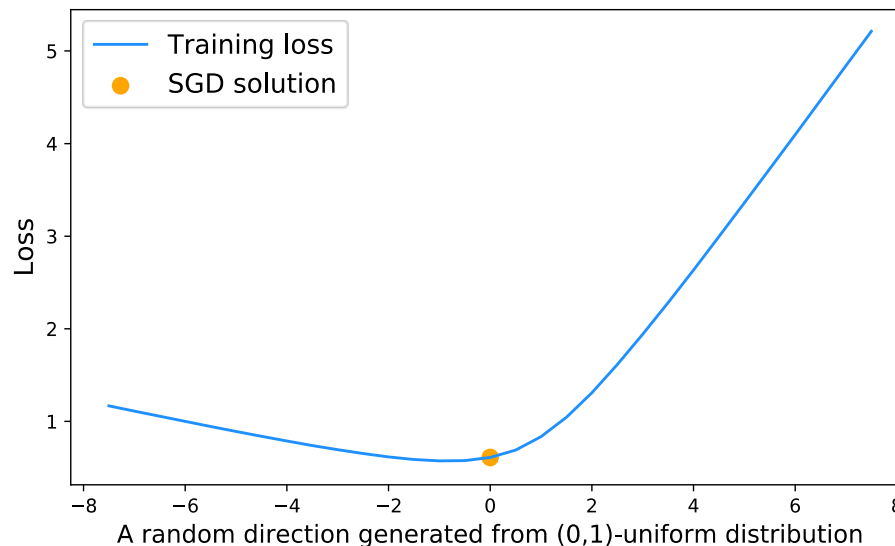red dashed line has the best match

$\delta$ : relative shift

$\xi_\delta$: loss difference with $\delta$ shift

# Assumptions and Verification

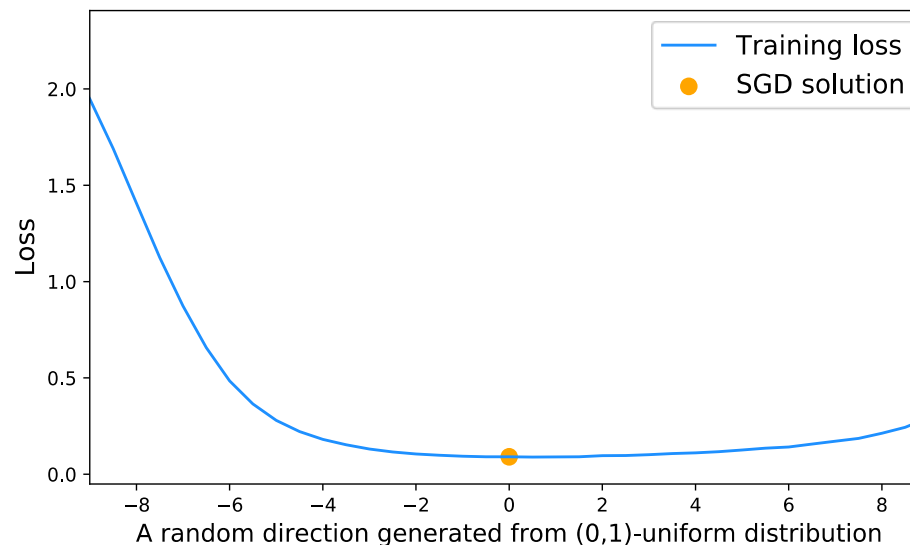Locally asymmetric assumption: assume locally asymmetric property

# Assumptions and Verification

Locally asymmetric assumption: asymmetric valley in a 2D case
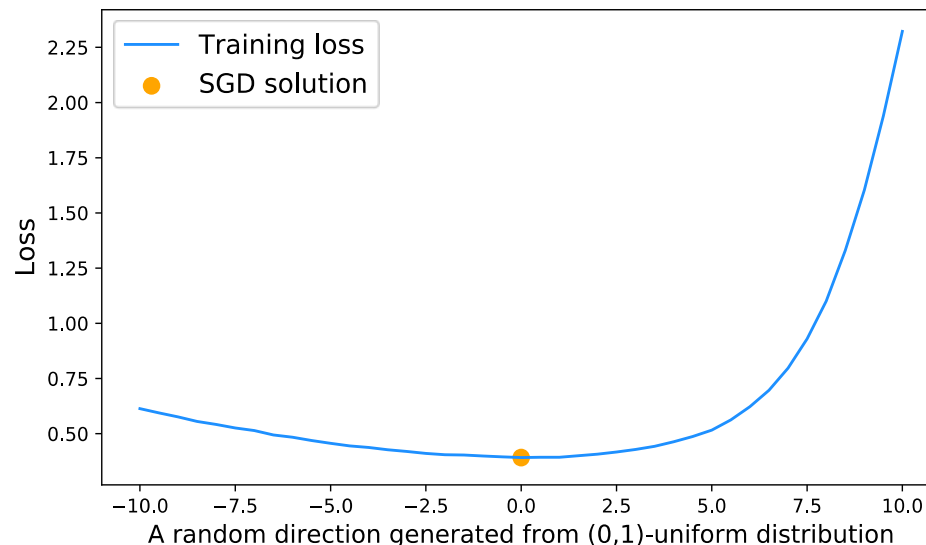(1 single neuron with its weight, bias and sigmoid activation)

# **Assumptions and Verification**

Locally asymmetric assumption: asymmetric valley in Deep neural networks (DenseNet-100 on CIFAR10)

# Assumptions and Verification

Locally asymmetric assumption: asymmetric valley in Deep neural networks (ResNet-164 on CIFAR100)
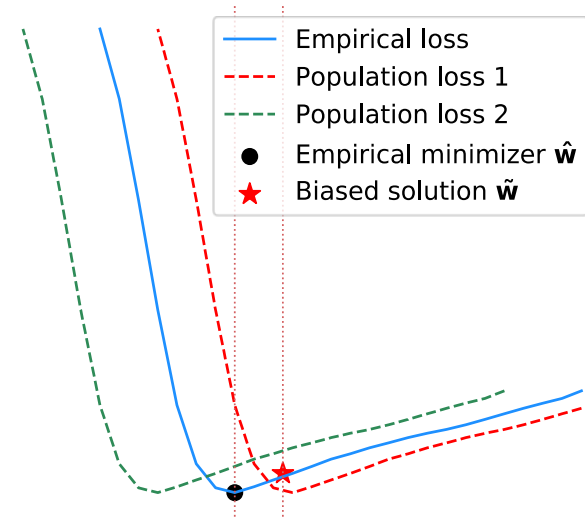
# Theorem (informal)

Bias leads to better generalization

$$E_\delta L(\hat{w}^*) - E_\delta L(\hat{w}^* + c_0) > 0$$

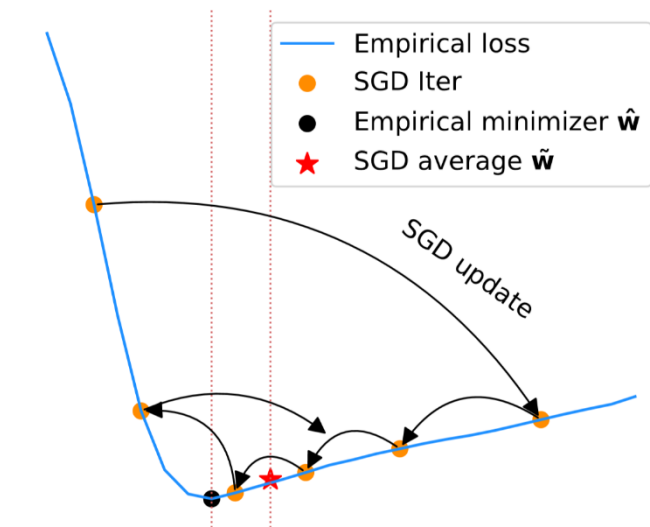where $c_0$ is a bias towards the flat side, $\hat{w}^*$ is an empirical solution



Legend:
— Empirical loss
--- Population loss 1
--- Population loss 2
● Empirical minimizer $\hat{\mathbf{w}}$
★ Biased solution $\tilde{\mathbf{w}}$

# **Theorem (informal)**

SGD averaging generates a bias

$$E[\overline{w}] > c_0 > 0$$

where $c_0$ is a bias towards the flat side, $\overline{w}$ is SGD average

# Thanks

## Asymmetric Valleys:
## Beyond Sharp and Flat Local Minima

Haowei He | Gao Huang | Yang Yuan

**October 26, 2019**