

ECE 884 Deep Learning

Lecture 7: Regularization and Training

02/09/2021

Review of last lecture

- Loss Function
 - What it is, what is it for?
 - Example: cross-entropy loss function

Today's lecture

- Regularization
 - What it is, what is it for?
 - Examples
- Training

Problem of Loss Function

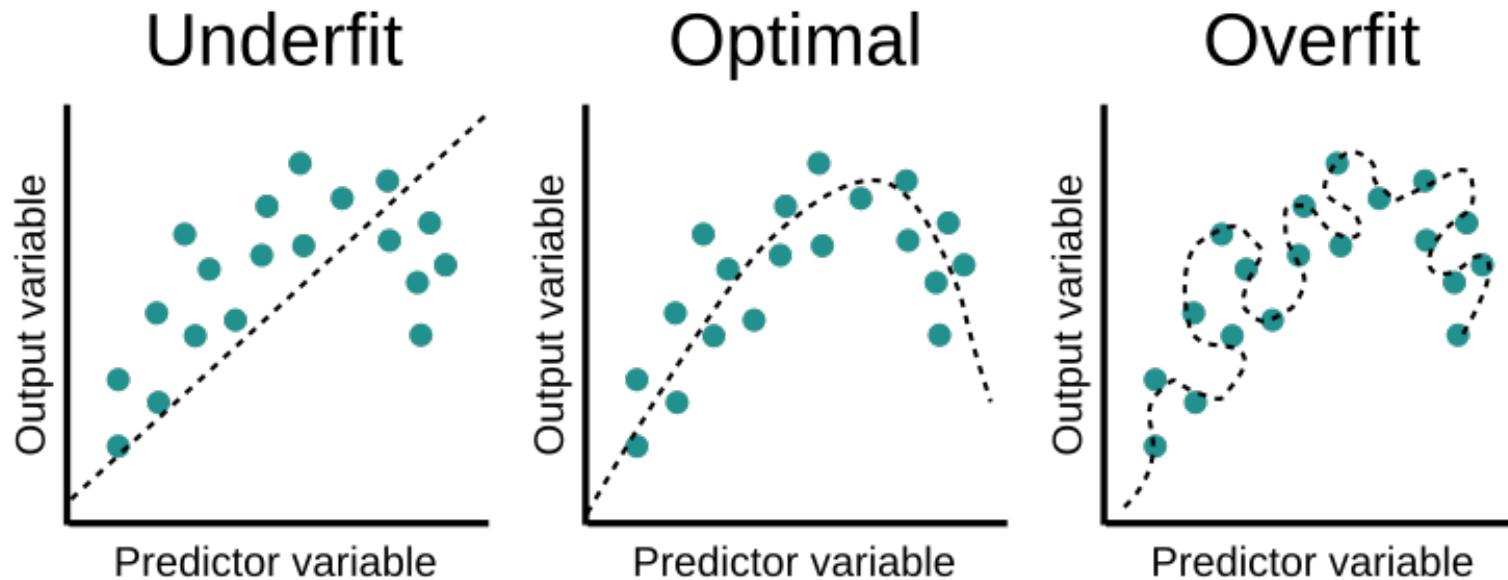
Loss functions encourage good performance on training data
but we really care about test data

Overfitting

A model is **overfit** when it performs too well on the training data, but has **poor** performance for unseen data

Overfitting

A model is **overfit** when it performs too well on the training data, but has **poor** performance for unseen data



Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)$$

Data loss: Model predictions
should match training data

Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too well* on training data

λ is a hyperparameter giving regularization strength

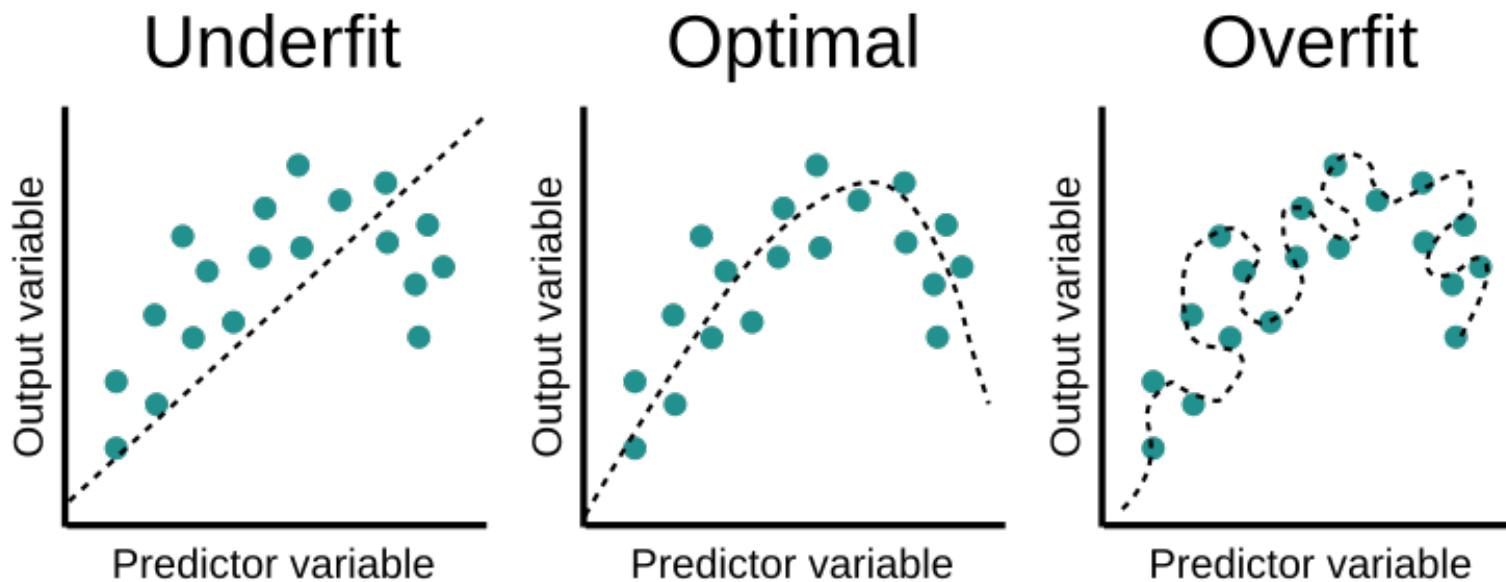
Example Regularization

L2 regularization: $R(W) = \sum_{k,l} W_{k,l}^2$

L1 regularization: $R(W) = \sum_{k,l} |W_{k,l}|$

Regularization Effect#1

Select simpler model with lower order



Regularization Effect#2

Select model with more features

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 Regularization

$$R(W) = \sum_{k,l} W_{k.l}^2$$

$$w_1^T x = w_2^T x = 1$$

Same predictions, so data loss will always be the same

Regularization Effect#2

Select model with more features

L2 Regularization

$$x = [1, 1, 1, 1]$$
$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$R(W) = \sum_{k,l} W_{k,l}^2$$

L2 regularization prefers weights to be “spread out”

$$w_1^T x = w_2^T x = 1$$

Same predictions, so data loss will always be the same

Summary

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

Loss function consists of **data loss** to fit the training data and **regularization** to prevent overfitting

Training

What is the purpose of training?

Find the model weights W that achieve the lowest total loss based on
the training dataset

Training

What is the purpose of training?

Training == Loss Optimization Problem

$$w^* = \arg \min_w L(w)$$

Training

What is the purpose of training?

Find the model weights W that achieve the lowest total loss based on the training dataset

$$w^* = \arg \min_w L(w)$$

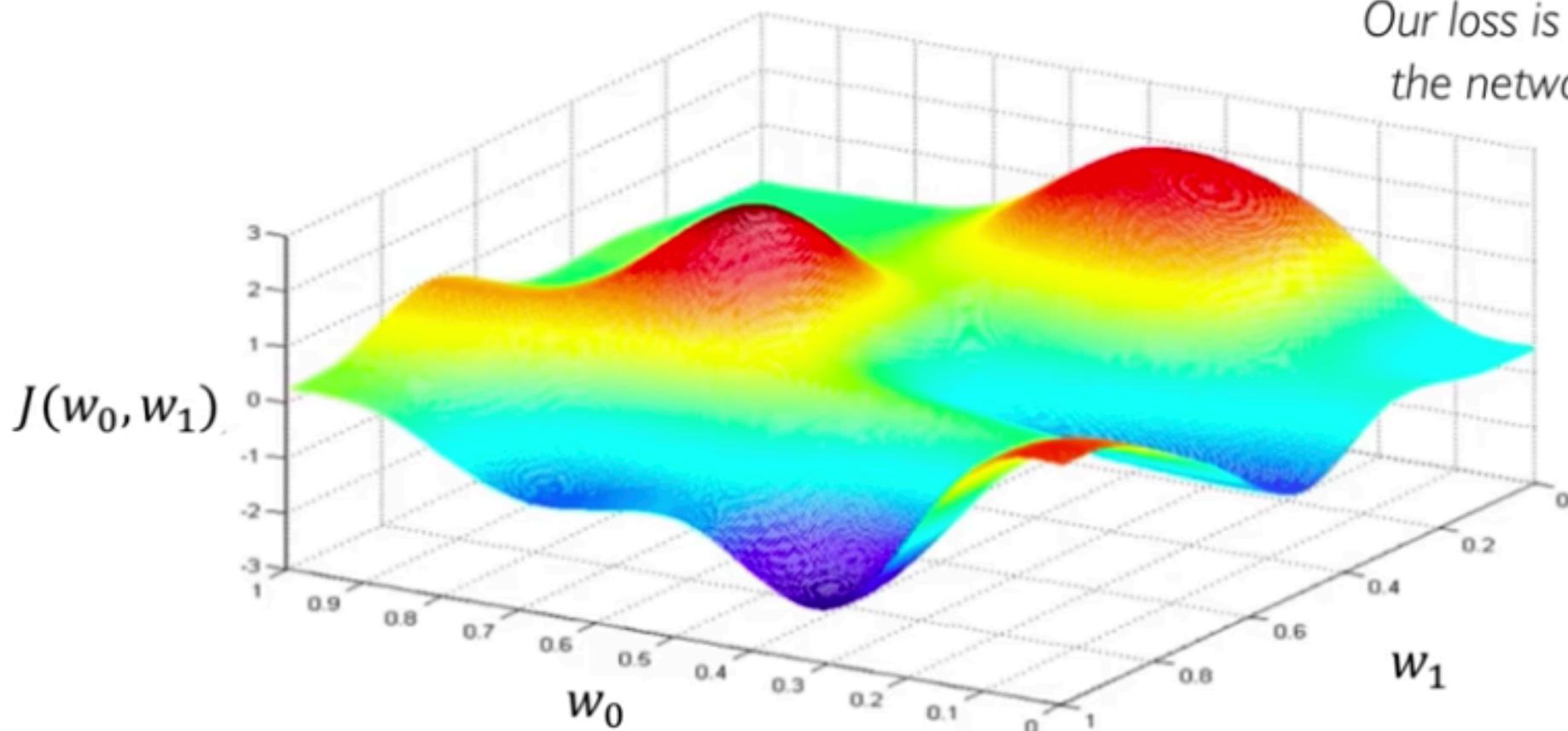
Remember:

$$W = \{W^{(0)}, W^{(1)}, \dots\}$$

Optimization Landscape

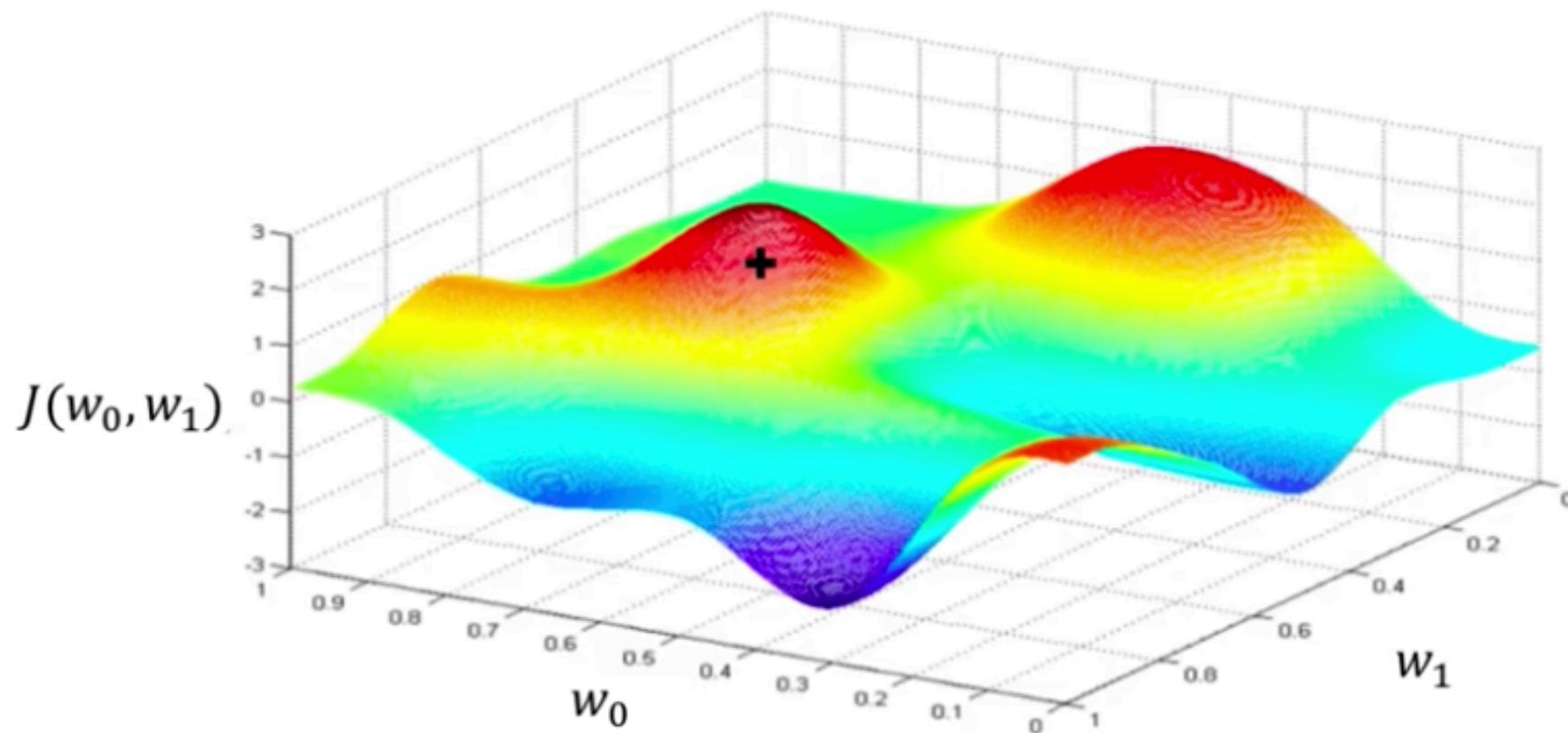
$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$

Remember:
Our loss is a function of
the network weights!



Optimization Method: Gradient Decent

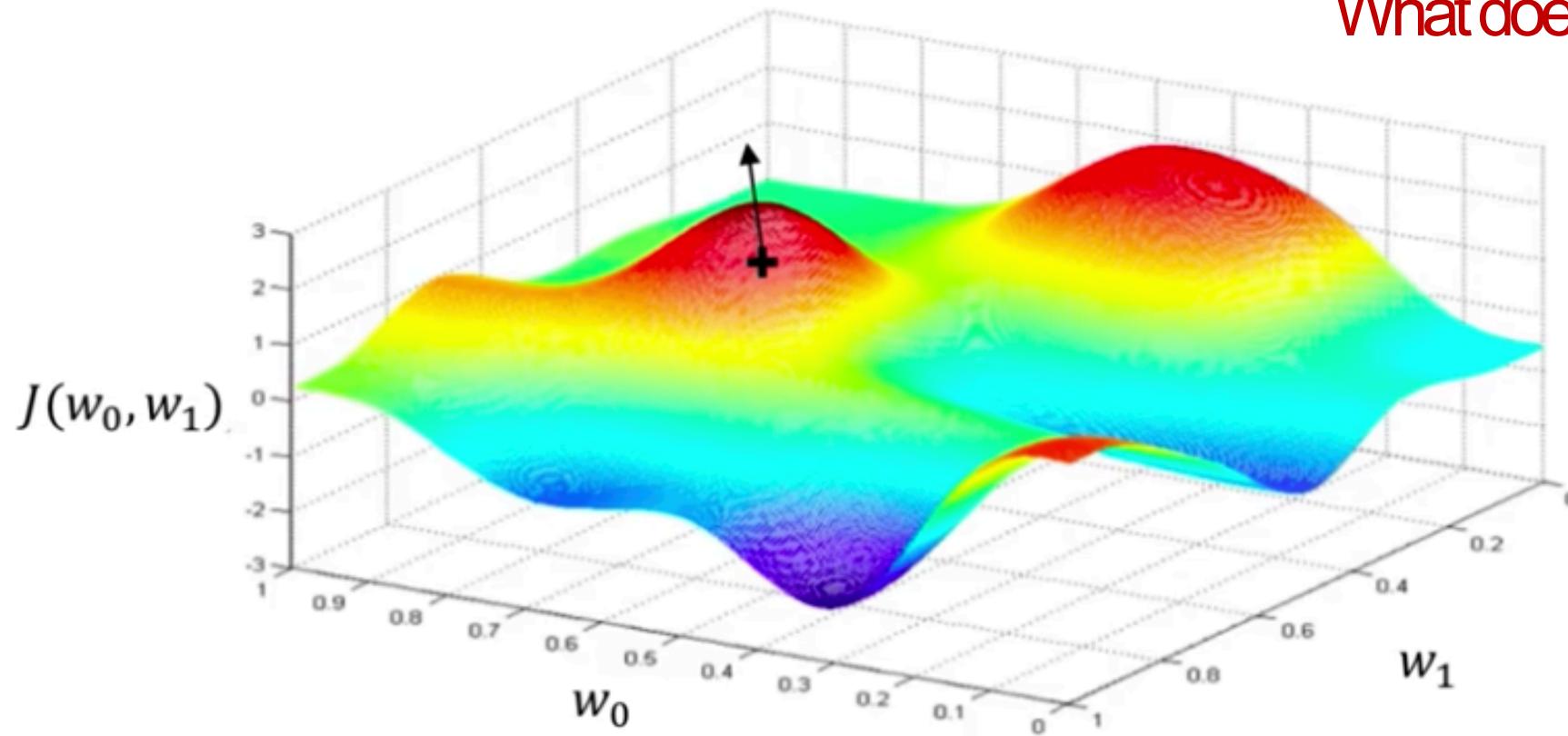
Randomly pick an initial (w_0, w_1)



Optimization Method: Gradient Decent

Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$

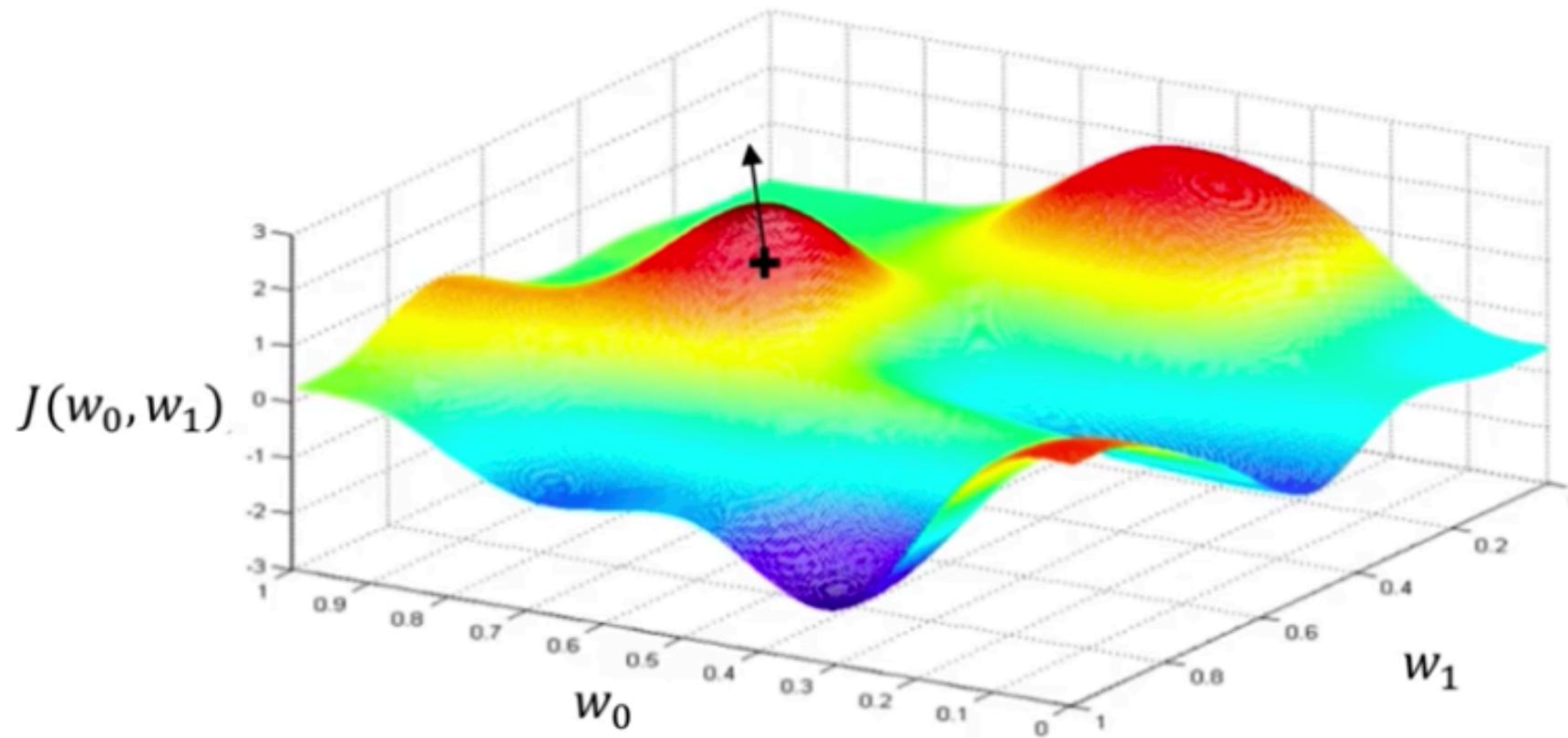
Why computing gradient?
What does gradient represent?



Optimization Method: Gradient Decent

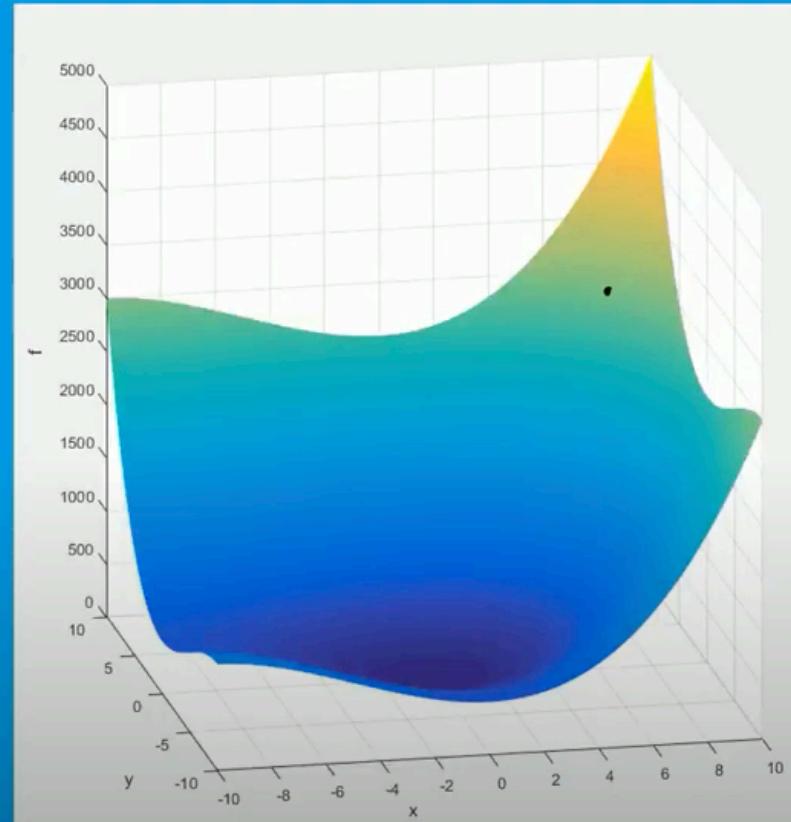
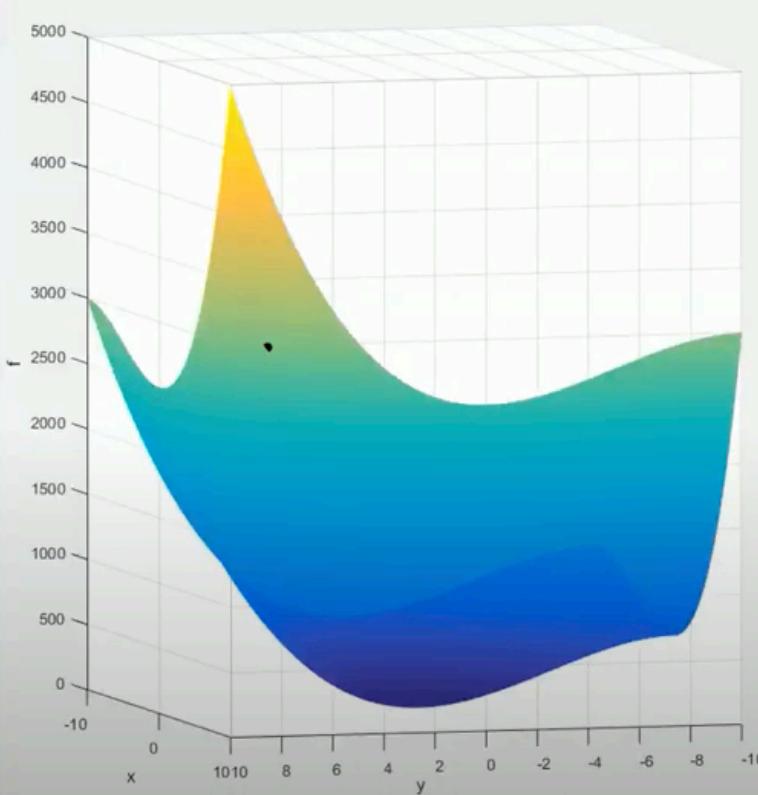
Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$

Directions

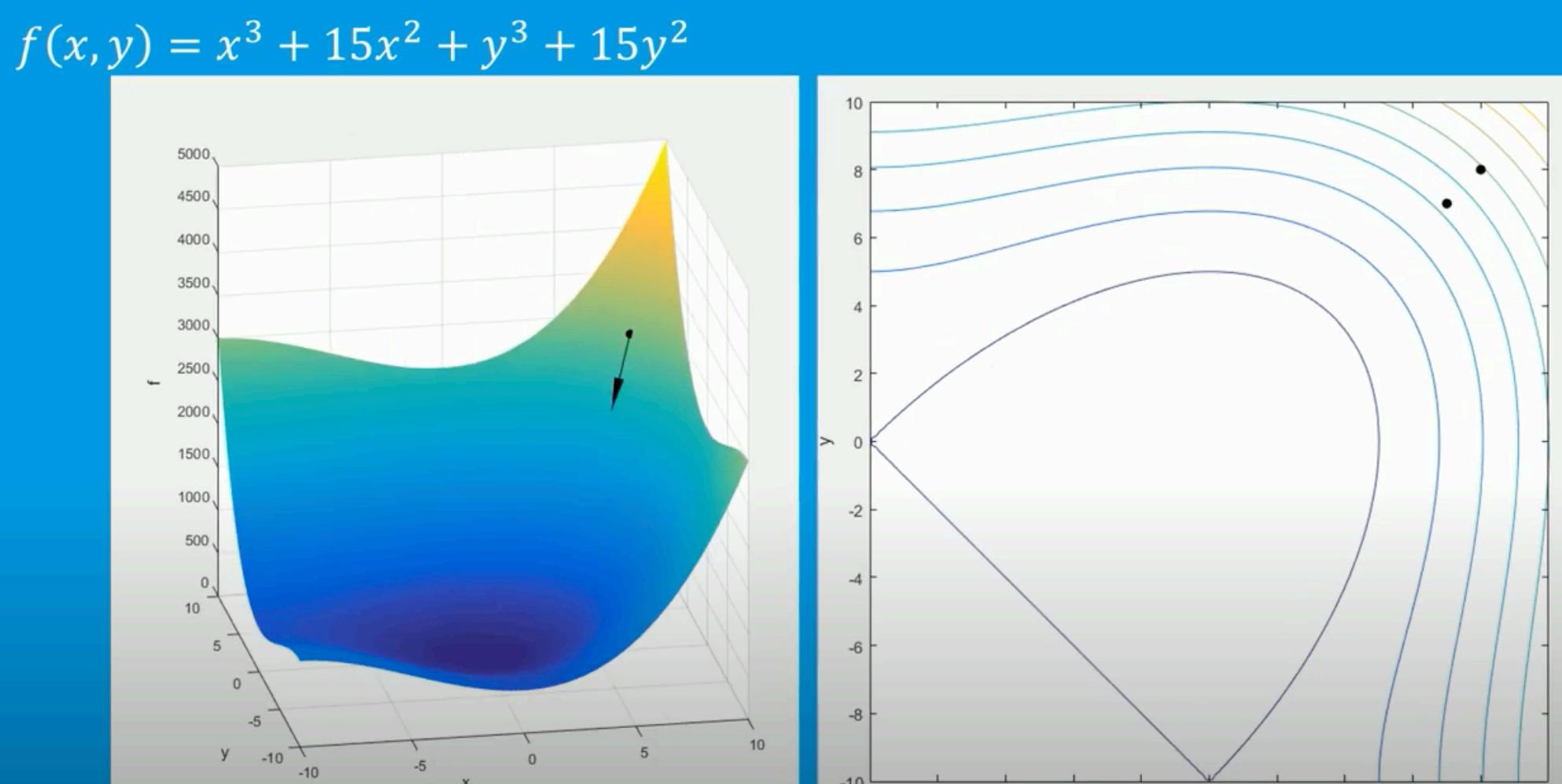


Optimization Method: Gradient Decent

$$f(x, y) = x^3 + 15x^2 + y^3 + 15y^2$$

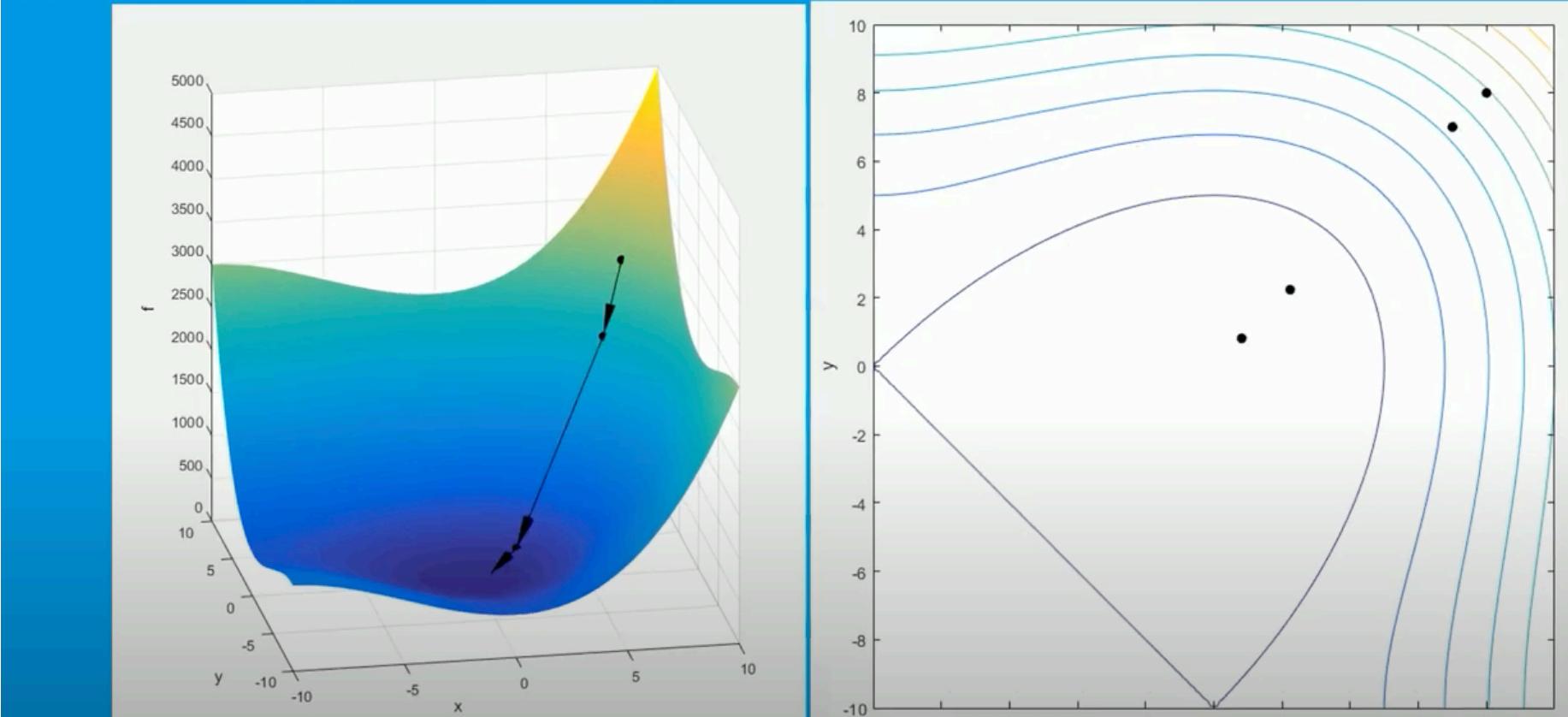


Optimization Method: Gradient Decent



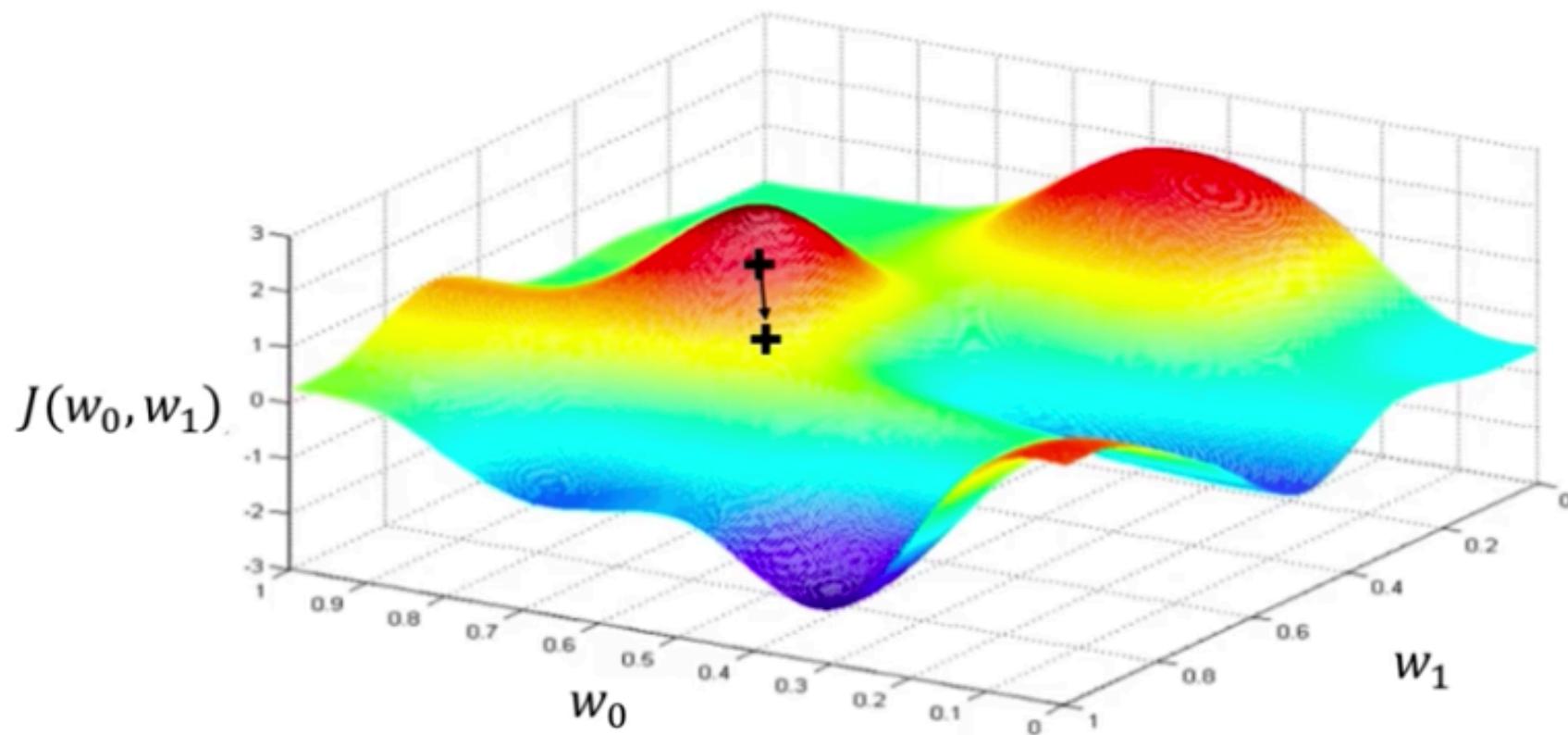
Optimization Method: Gradient Decent

$$f(x, y) = x^3 + 15x^2 + y^3 + 15y^2$$



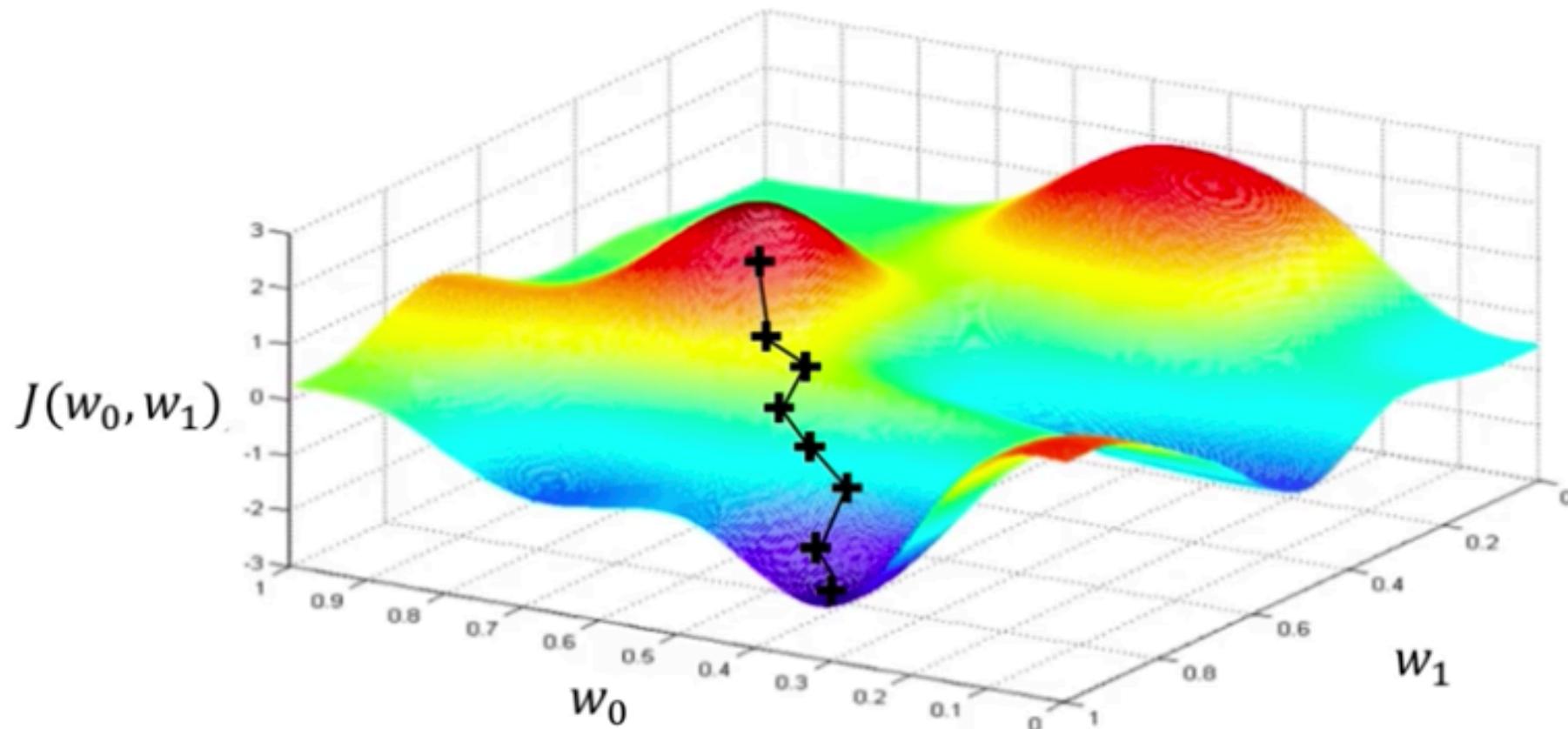
Optimization Method: Gradient Decent

Take small step in opposite direction of gradient



Optimization Method: Gradient Decent

Repeat until convergence



Optimization Method: Gradient Decent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Batch Gradient Descent

Problem: This is very expensive when N is large!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

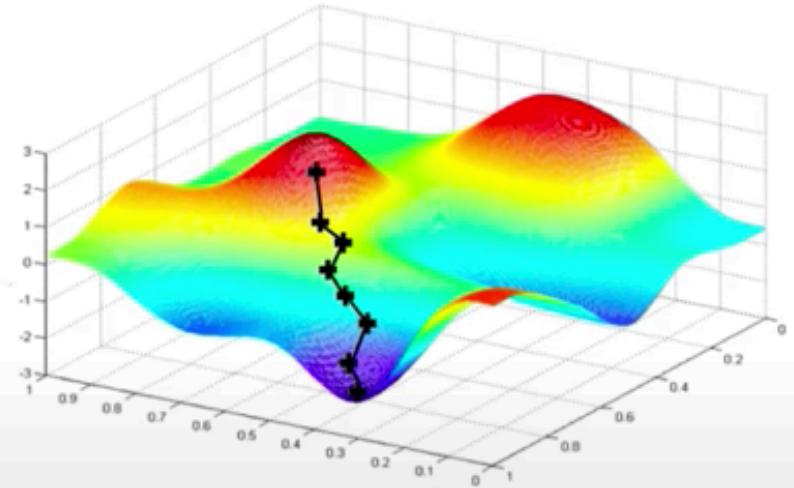
$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Stochastic Gradient Descent (SGD)

Single point to calculate GD

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



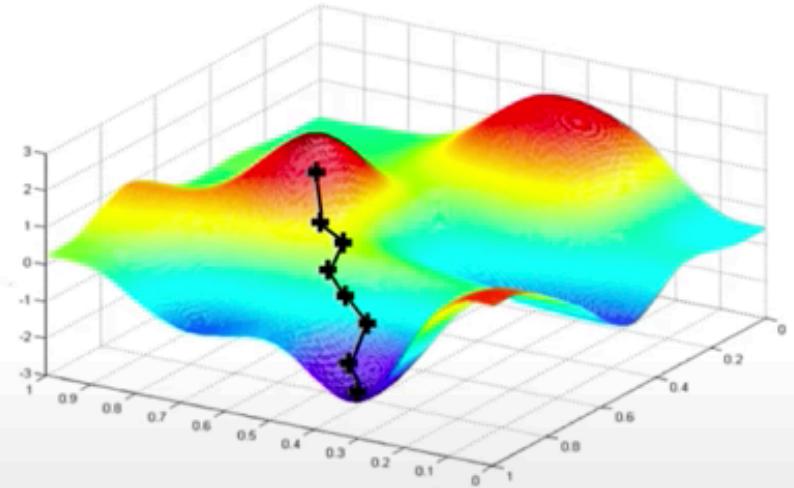
Easy to compute but
very noisy (stochastic)!

Stochastic Gradient Descent (SGD)

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

Problem: calculated gradient decent is noisy!



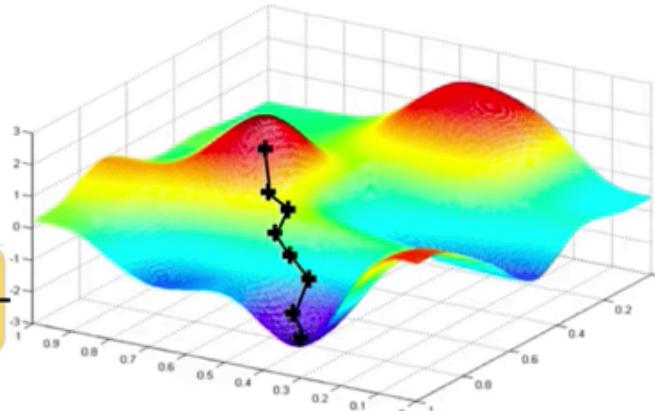
Easy to compute but
very noisy (stochastic)!

Stochastic Gradient Descent (SGD)

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$
5. Update weights,
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$
6. Return weights

Fast to compute and a much better estimate of the true gradient!



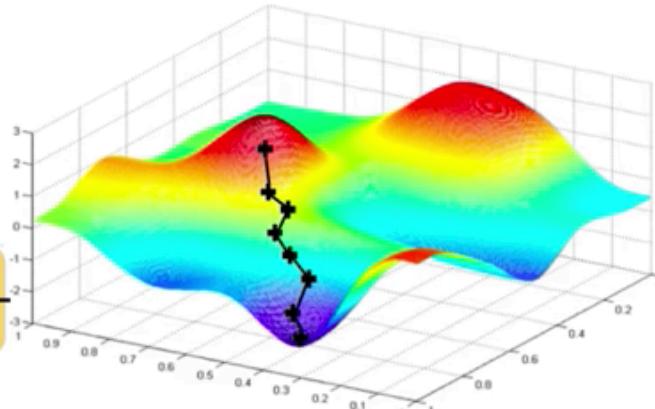
Approximate sum using a **minibatch** of 32 / 64 / 128 examples

Hyperparameters of Stochastic Gradient Descent (SGD)

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$
5. Update weights,
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$
6. Return weights

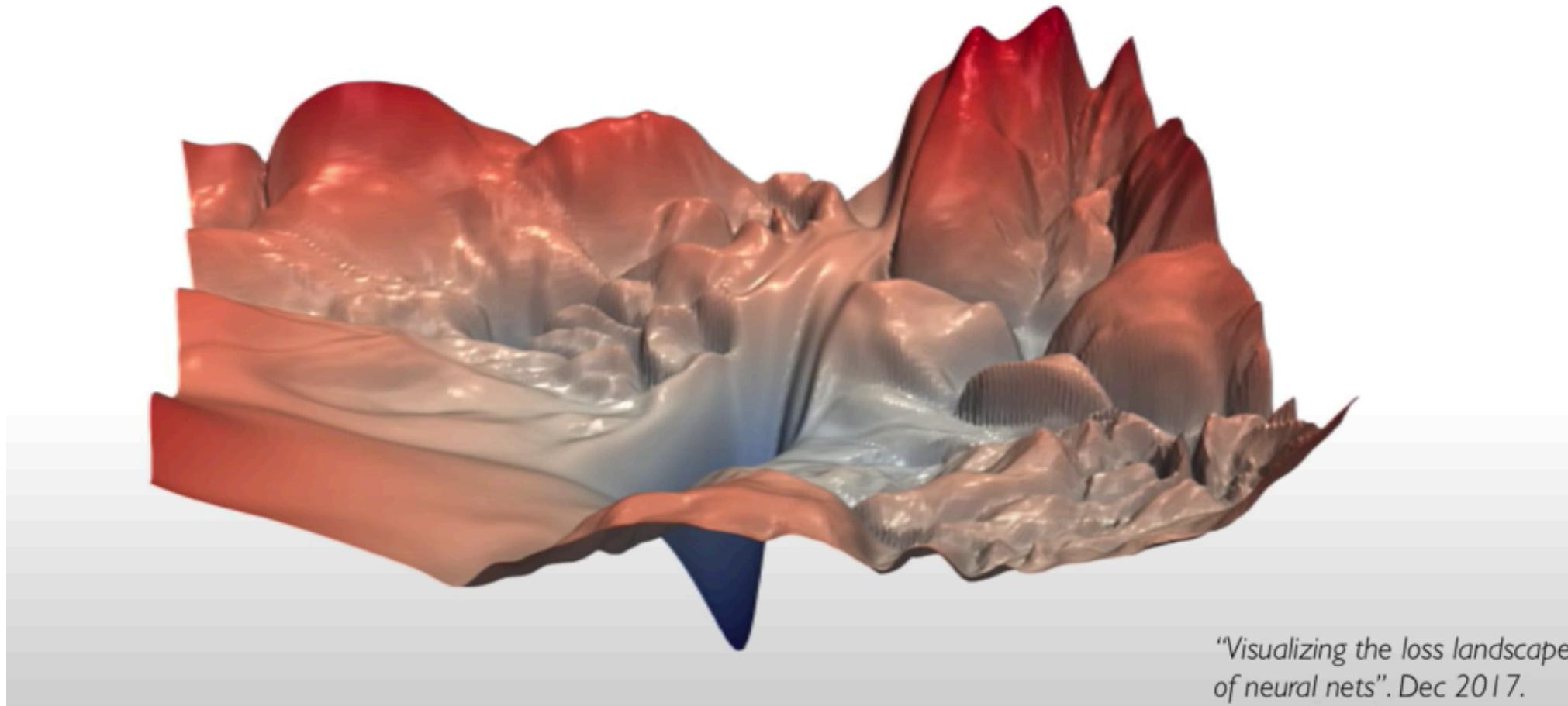
Fast to compute and a much better estimate of the true gradient!



Hyperparameters:

- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling

Initialization Matters a lot!



"Visualizing the loss landscape
of neural nets". Dec 2017.

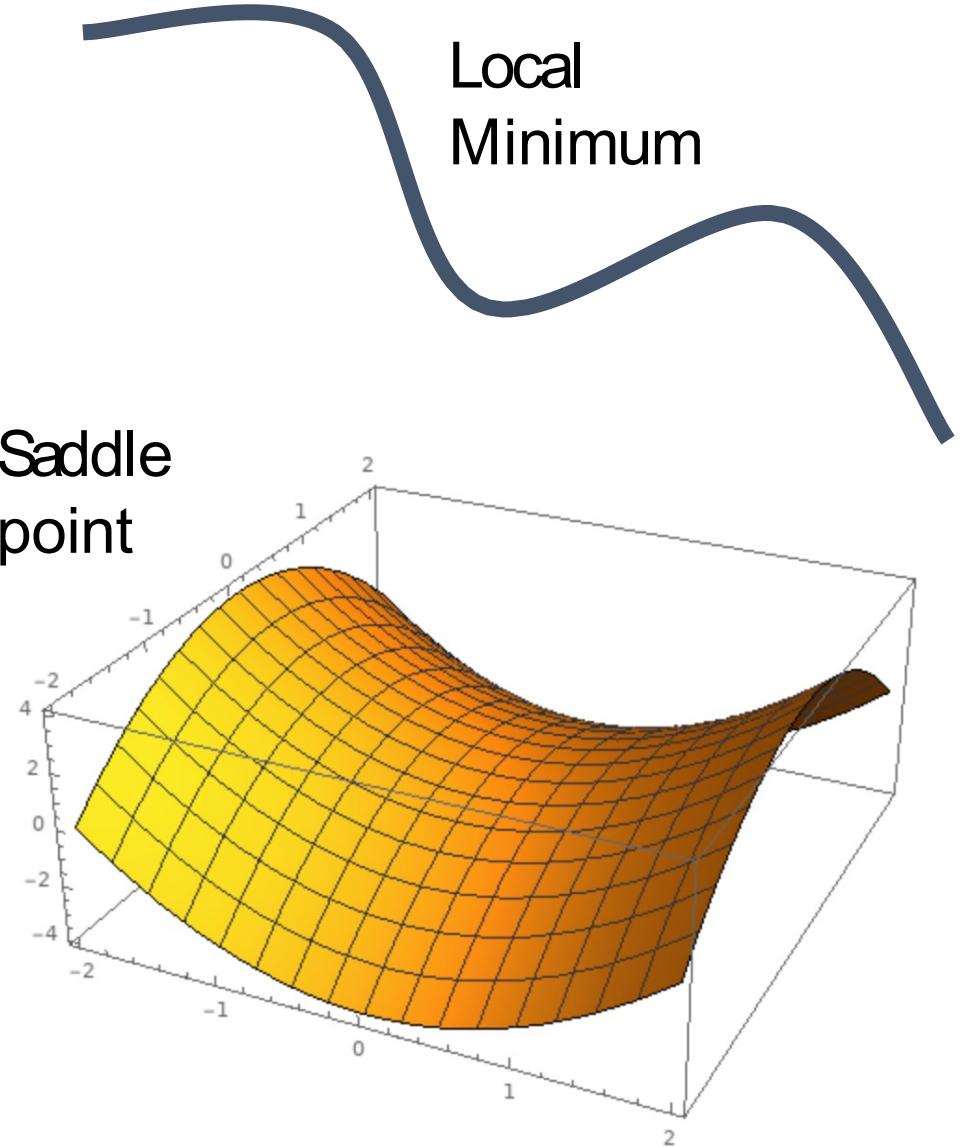
Pros and Cons

GRADIENT BASED OPTIMIZATION

- Pros
 - Widely Used
 - Fast
 - Scales Well
- Cons
 - Requires smooth gradients
 - Cost of computing gradients
 - Local minima

Problem#1 with SGD

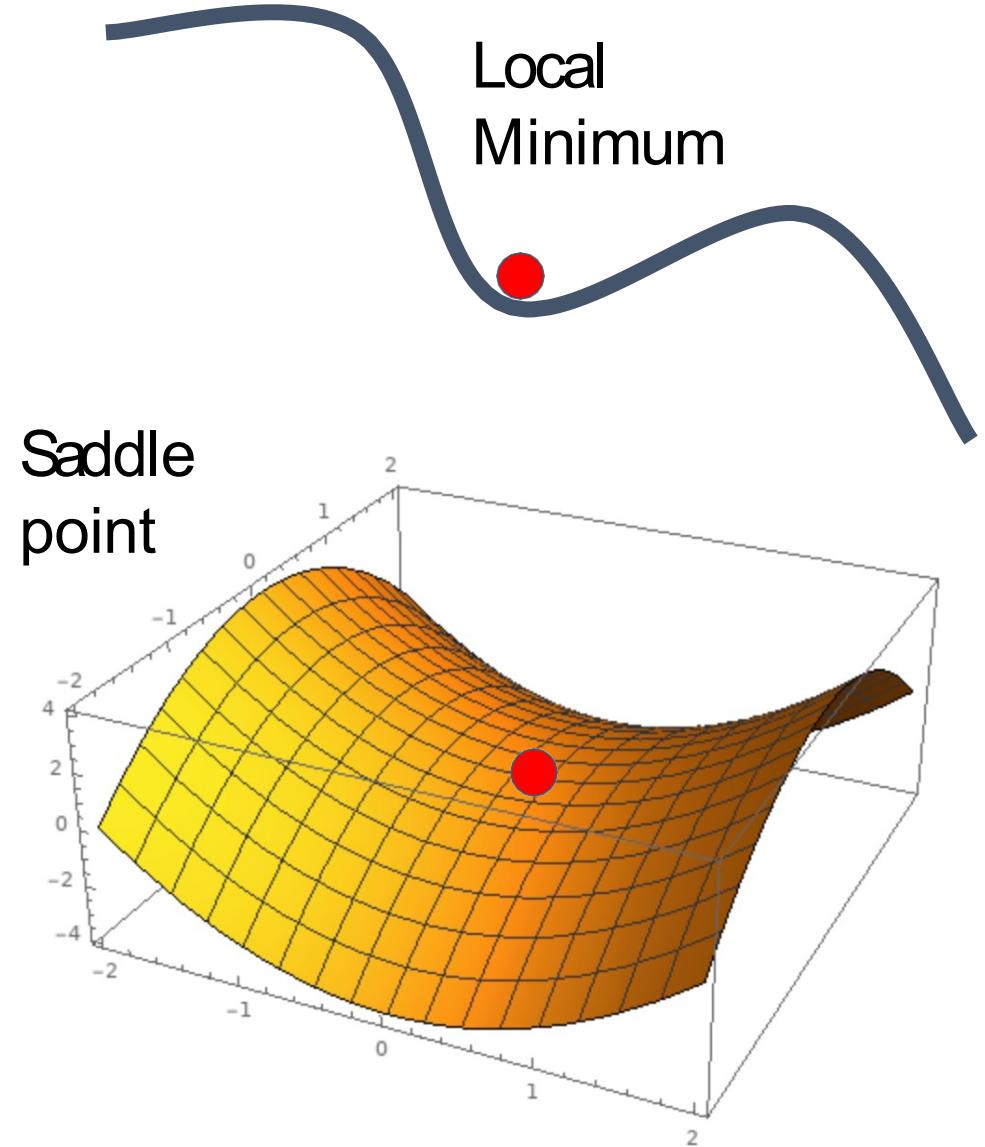
What if the loss
function has a **local
minimum** or **saddle
point**?



Problem#1 with SGD

What if the loss
function has a **local
minimum** or **saddle
point**?

Zero gradient, gradient
descent gets stuck



SGD+Momentum

SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

SGD+Momentum

$$\begin{aligned}v_{t+1} &= \rho v_t + \nabla f(x_t) \\x_{t+1} &= x_t - \alpha v_{t+1}\end{aligned}$$

V is velocity as a moving average
of gradients.

SGD+Momentum

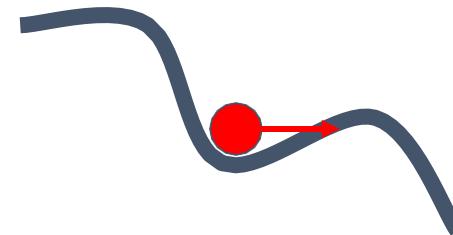
SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

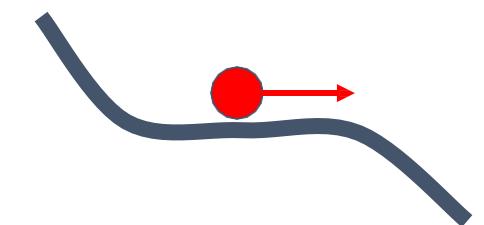
SGD+Momentum

$$\begin{aligned}v_{t+1} &= \rho v_t + \nabla f(x_t) \\x_{t+1} &= x_t - \alpha v_{t+1}\end{aligned}$$

Local Minima

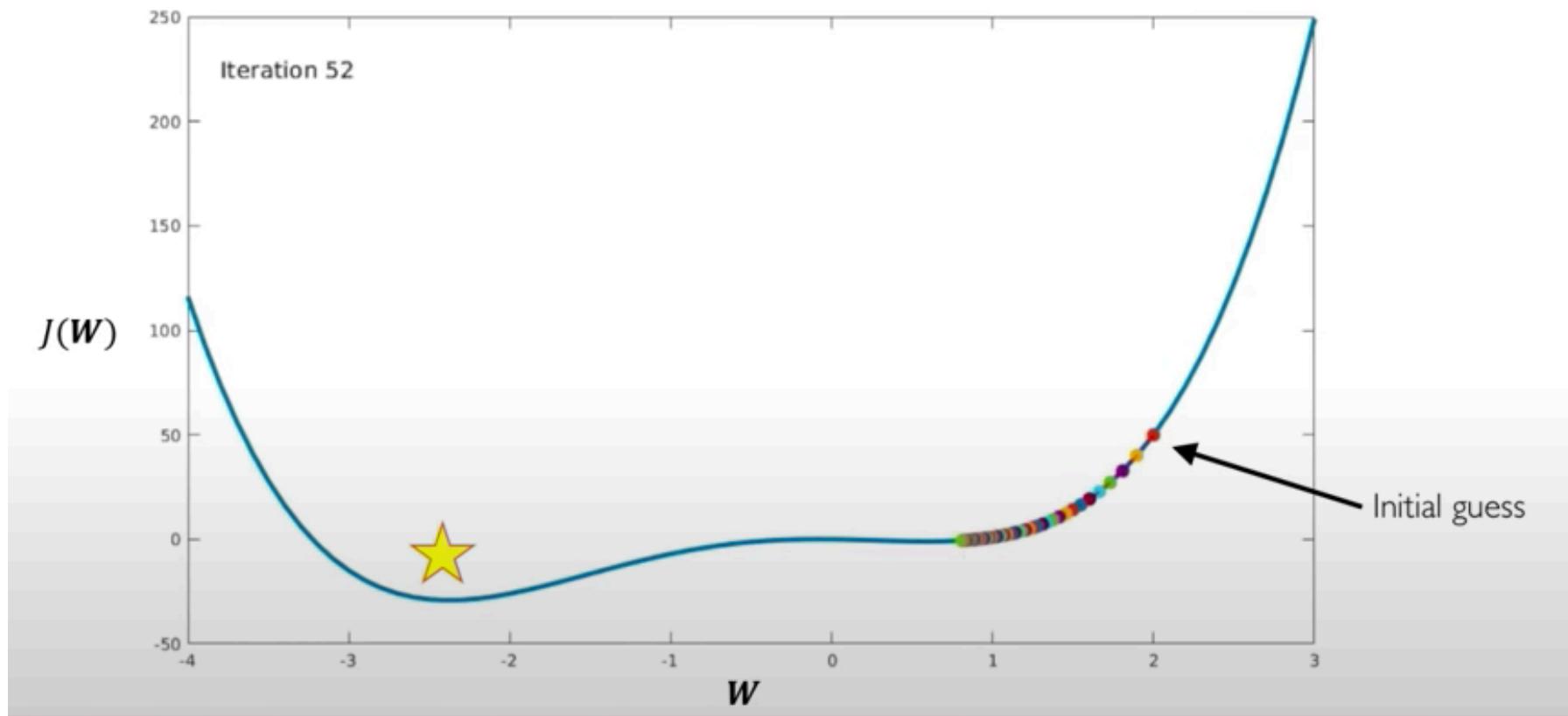


Saddle points



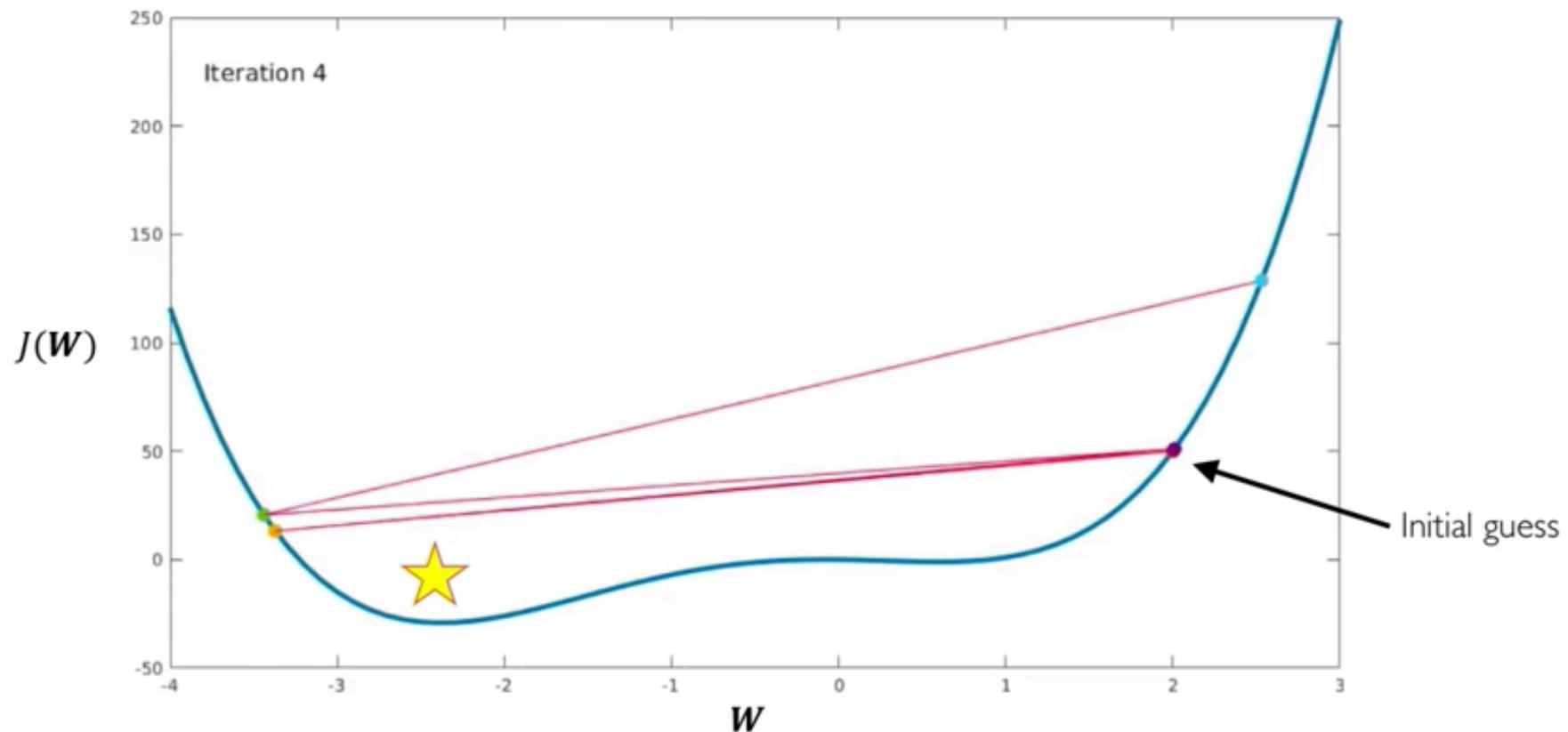
Problem#2 with SGD: Learning Rate

Small learning rate converges slowly and gets stuck in false local minima



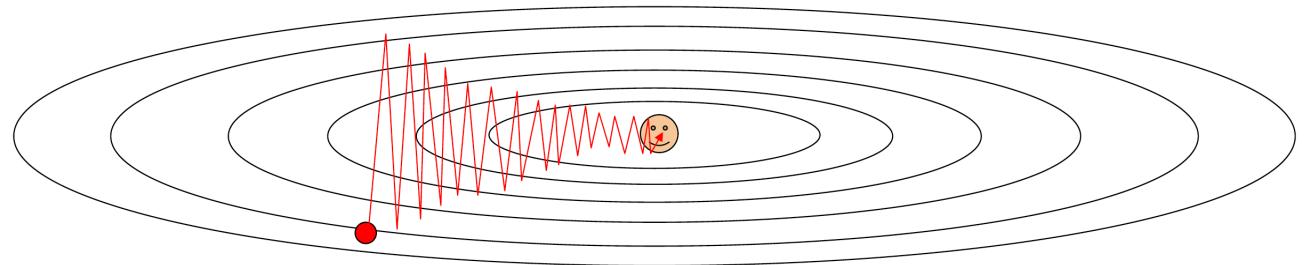
Problem#2 with SGD: Learning Rate

Large learning rates overshoot, become unstable and diverge



Optimizers with Adaptive Learning Rate

- RMSProp
- AdaGrad
- Adam == Momentum + RMSProp
- AdamW



Per-parameter learning rate

Learning rate decay

In practice:

- **Adam / AdamW** is a good default choice in many cases **SGD+Momentum** can outperform Adam but may require more tuning

Any Question ?