

Machine Learning

ITCS 5356

Gradient Descent

Least Mean Squares

Razvan C. Bunescu

Department of Computer Science @ CCI

razvan.bunescu@charlotte.edu

ML is Optimization

- Try to find the value for w that minimizes:

$$J(w) = \frac{1}{2}w^2 - 4w + 9$$

$$J(w) = \frac{1}{2}(w - 4)^2 + 1$$

- Set $\nabla J(w) = 0$
 - $\Rightarrow w - 4 = 0$
 - $\Rightarrow w = 4$

Machine Learning is Optimization

- Parametric ML involves minimizing an **objective function** $J(\mathbf{w})$:
 - Also called **cost function** or **loss function**.
 - Want to find $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$
- Numerical optimization procedure:
 1. Start with some guess for \mathbf{w}^0 , set $\tau = 0$.
 2. Update \mathbf{w}^τ to $\mathbf{w}^{\tau+1}$ such that $J(\mathbf{w}^{\tau+1}) \leq J(\mathbf{w}^\tau)$.
 3. Increment $\tau = \tau + 1$.
 4. Repeat from 2 until J cannot be improved anymore.

Gradient-based Optimization

- How to update \mathbf{w}^τ to $\mathbf{w}^{\tau+1}$ such that $J(\mathbf{w}^{\tau+1}) \leq J(\mathbf{w}^\tau)$?

- Move \mathbf{w} in the direction of **steepest descent**:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau + \eta \Delta$$

- Δ is the direction of steepest descent, i.e. direction along which J decreases the most.
- η is the learning rate and controls the magnitude of the change.

Gradient-based Optimization

- Move \mathbf{w} in the direction of **steepest descent**:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} + \eta \Delta$$

- What is the direction of **steepest descent** of $J(\mathbf{w})$ at \mathbf{w}^{τ} ?
 - The gradient $\nabla J(\mathbf{w})$ is in the direction of **steepest ascent**.
 - Set $\Delta = -\nabla J(\mathbf{w}) \Rightarrow$ the **gradient descent** update:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$$

Gradient Descent Algorithm

- Want to minimize a function $J : R^n \rightarrow R$.
 - J is differentiable and convex.
 - compute gradient of J i.e. *direction of steepest increase*:

$$\nabla J(\mathbf{w}) = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n} \right]$$

1. Set learning rate $\eta = 0.001$ (or other small value).
2. Start with some guess for \mathbf{w}^0 , set $\tau = 0$.
3. Repeat for epochs E or until J does not improve:
4. $\tau = \tau + 1$.
5. $\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$

What if objective is not differentiable?

- **Subgradient methods.**
 - Minimize convex functions that are not necessarily differentiable.
- **Gradient free methods:**
 - **Evolutionary Programming.**
 - **Bayesian Optimization.**
 - <https://arxiv.org/abs/1807.02811>
 - **Particle swarm optimization.**
 - **Surrogate optimization**
 - **Simulated annealing.**
 - ...

Gradient Descent Algorithm

- Want to minimize a function $J : R^n \rightarrow R$.
 - J is differentiable and convex.
 - compute gradient of J i.e. *direction of steepest increase*:

$$\nabla J(\mathbf{w}) = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n} \right]$$

1. Set learning rate $\eta = 0.001$ (or other small value).
2. Start with some guess for \mathbf{w}^0 , set $\tau = 0$.
3. Repeat for epochs E or until J does not improve:
4. $\tau = \tau + 1$.
5. $\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$

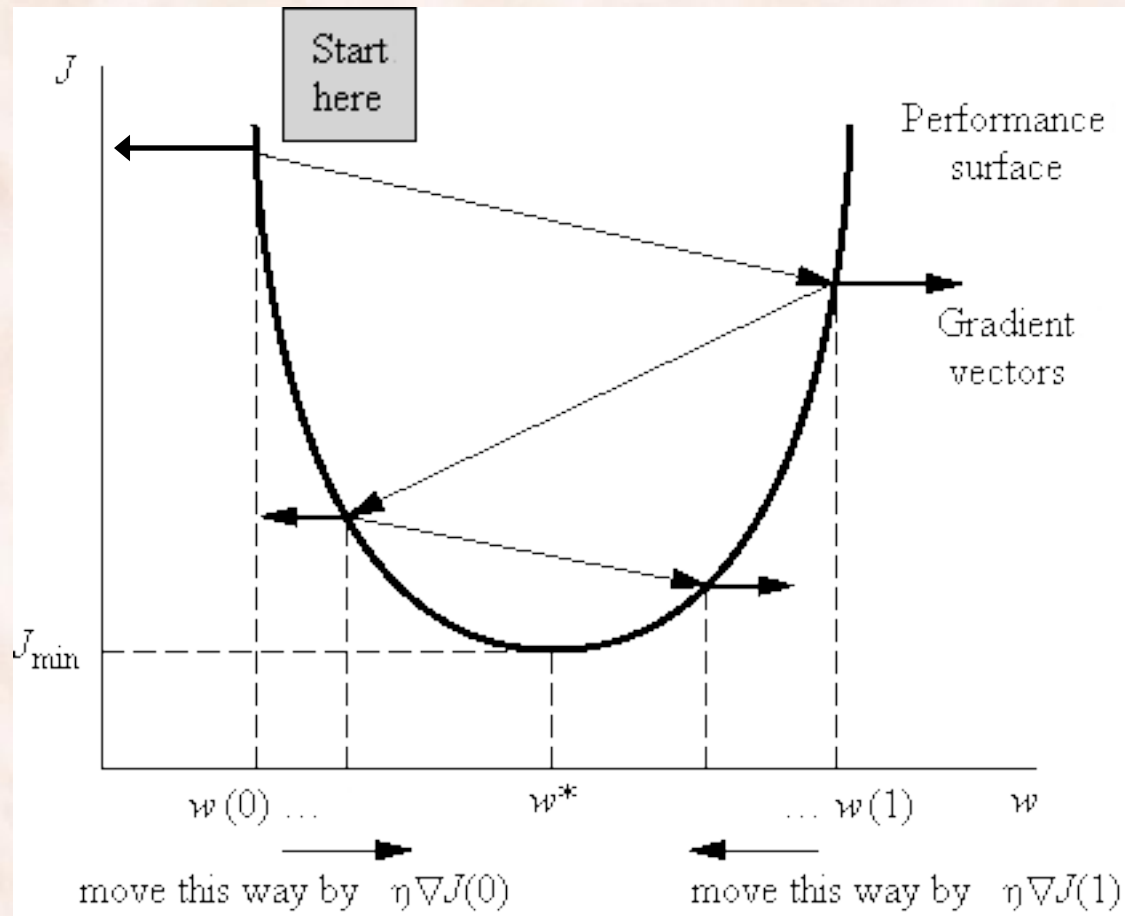
Gradient Descent: Example

- Let's use gradient descent to minimize:

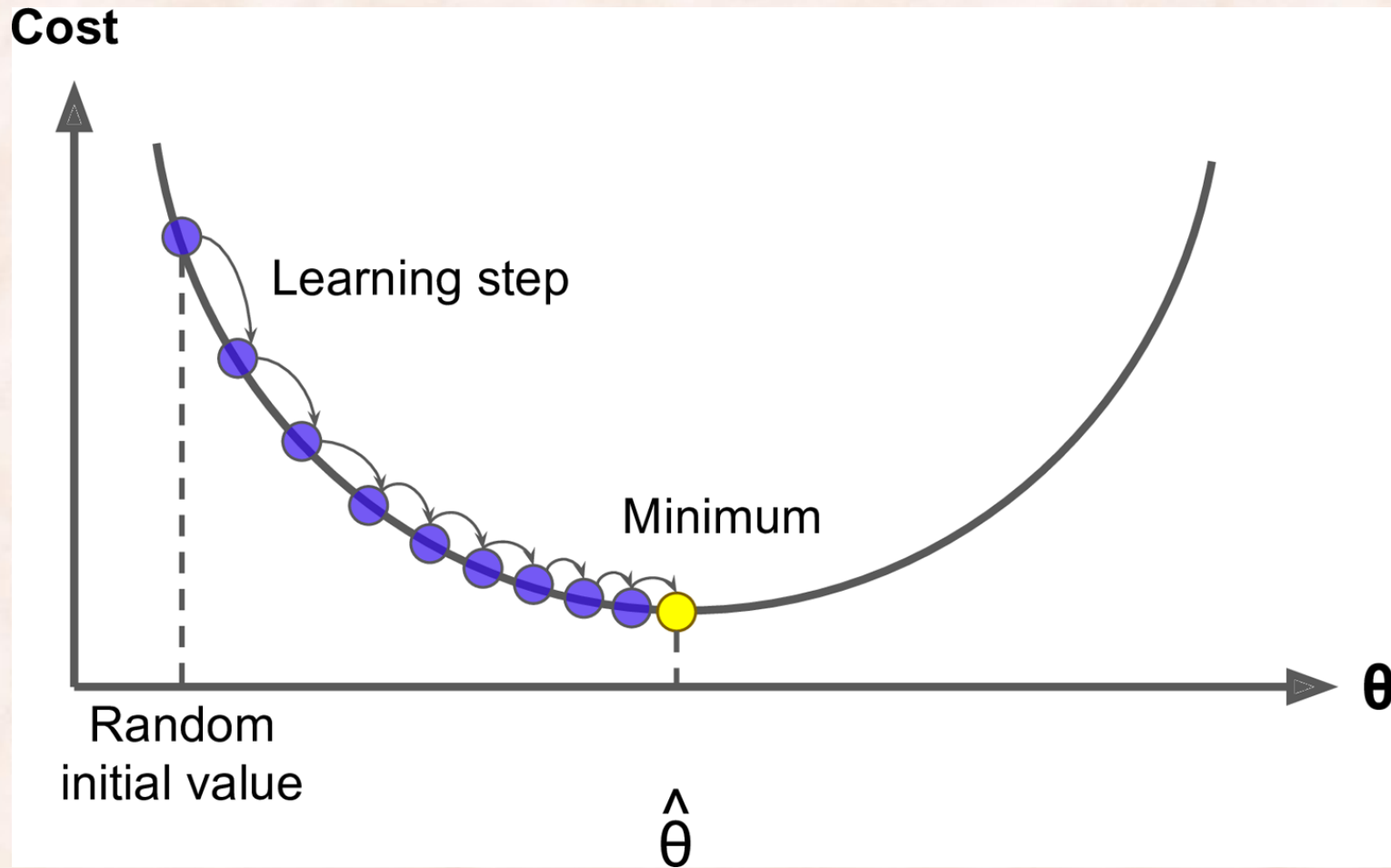
$$J(w) = \frac{1}{2}w^2 - 4w + 9$$

- Start from $w^0 = 0$, use learning rate $\eta = 0.5$
- What if $\eta = 1.0$?
- What if $\eta = 2.0$?

Gradient Descent: Large Updates



Gradient Descent: Small Updates

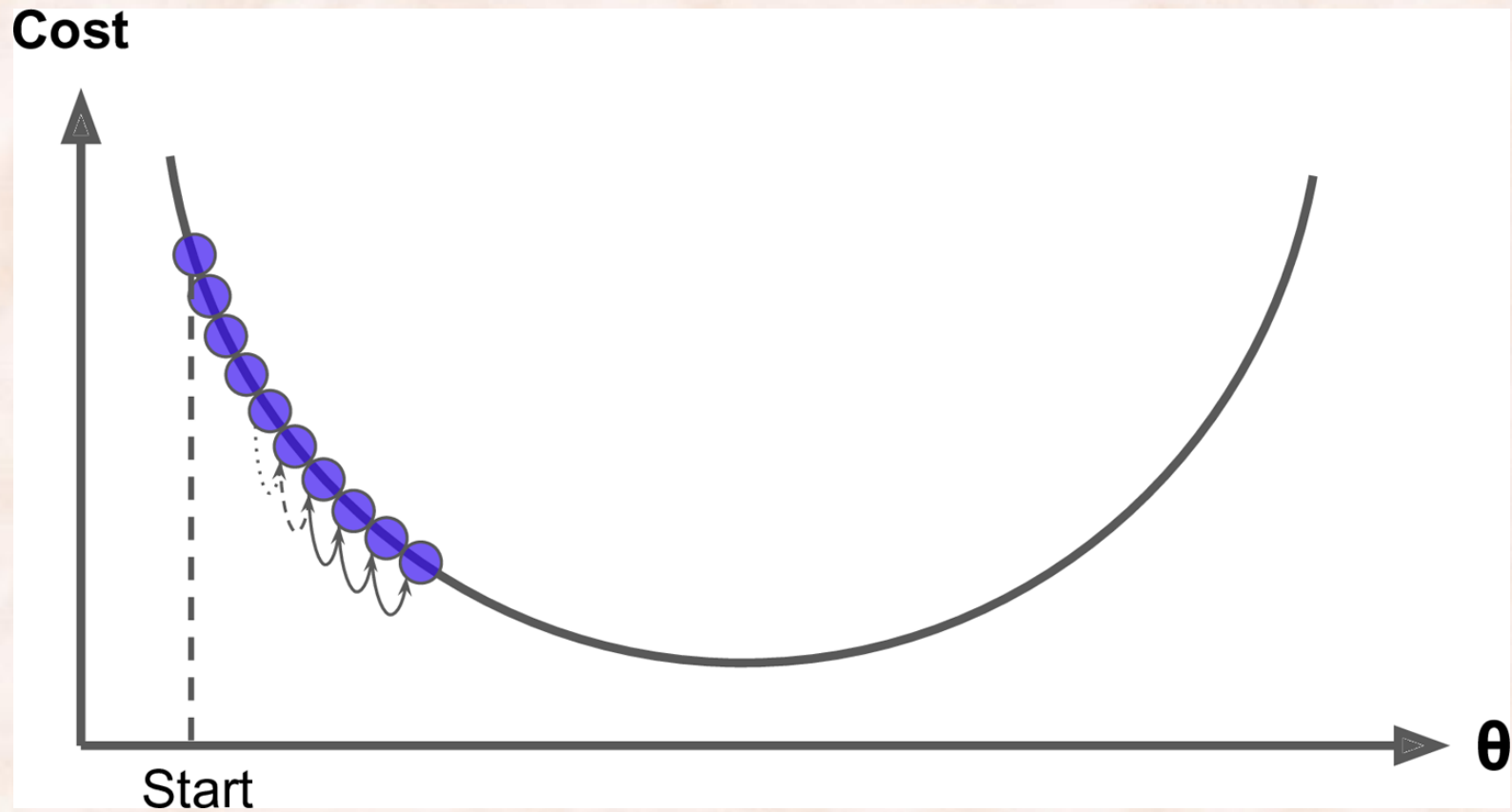


The Learning Rate

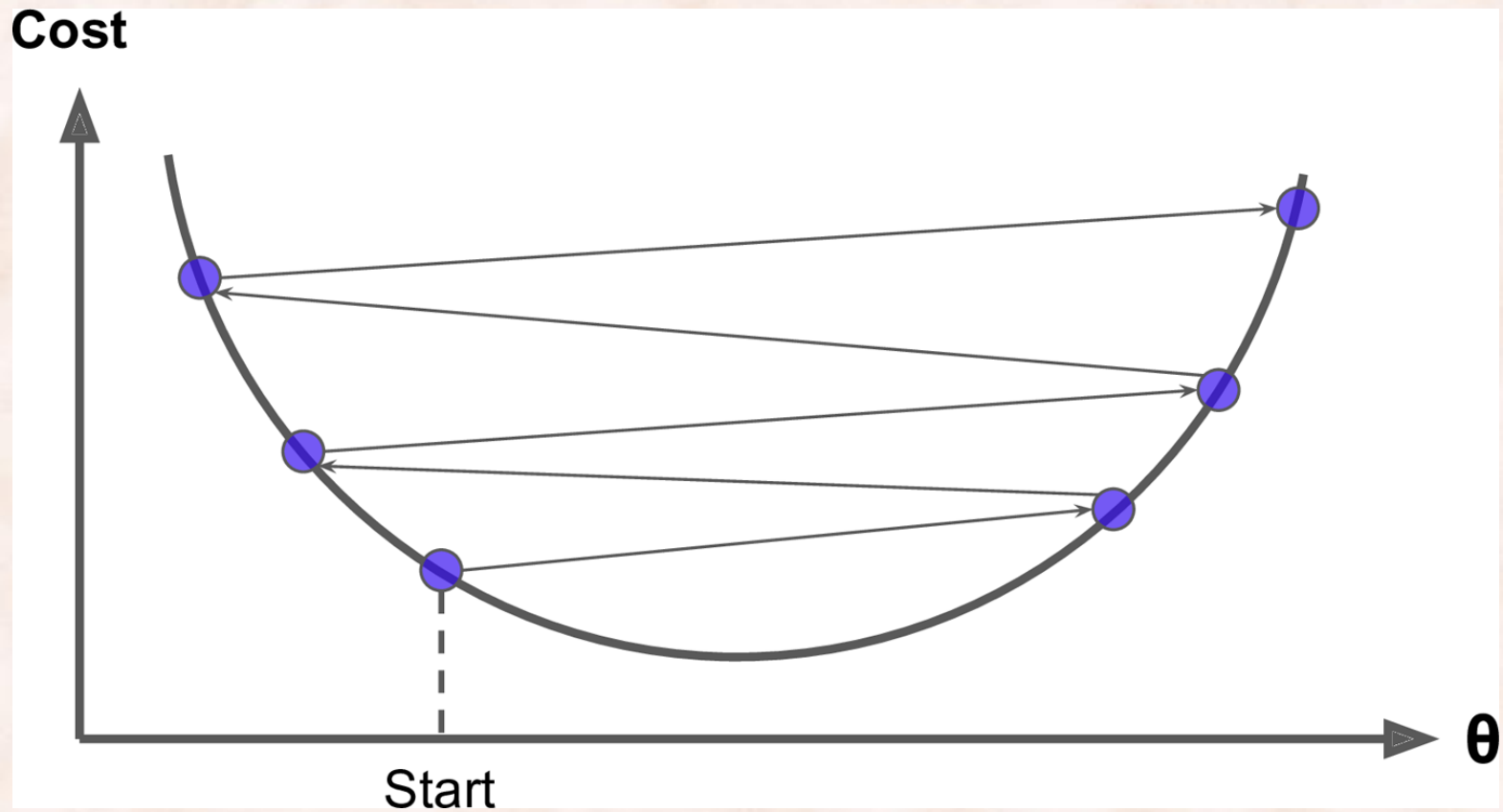
1. Set **learning rate** $\eta = 0.001$ (or other small value).
2. Start with some guess for \mathbf{w}^0 , set $\tau = 0$.
3. Repeat for epochs E or until J does not improve:
4. $\tau = \tau + 1$.
5. $\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$

- How big should the **learning rate** be?
 - If learning rate too small => slow convergence.
 - If learning rate too big => oscillating behavior => may not even converge.

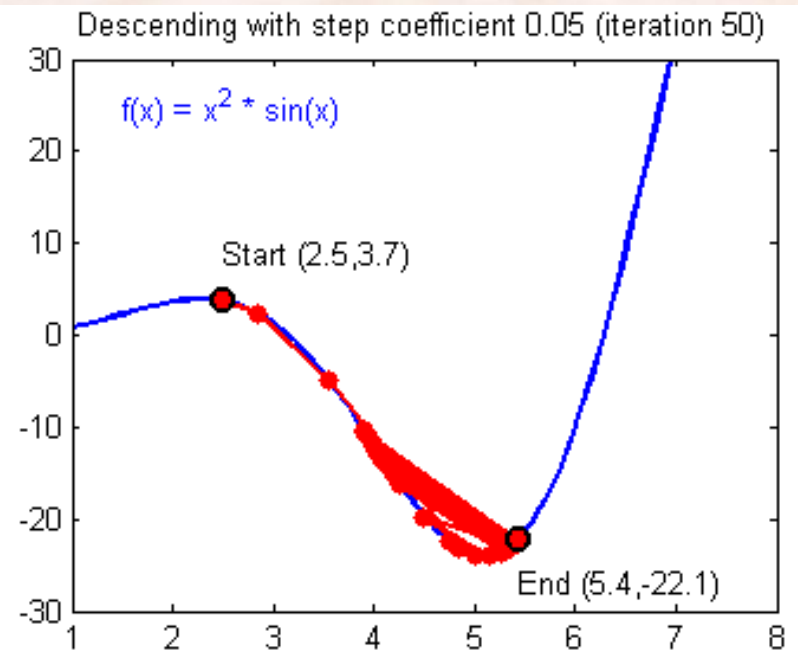
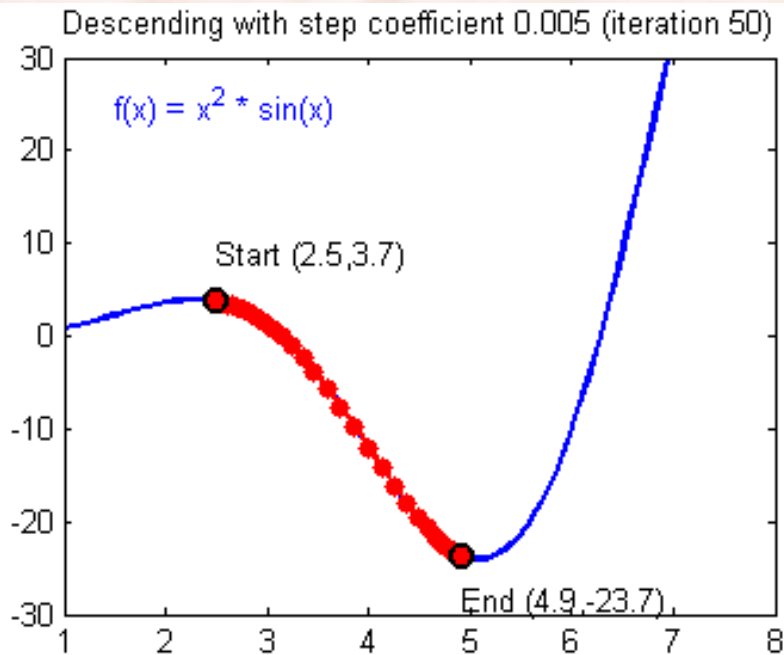
Learning Rate too Small



Learning Rate too Large



Learning Rates vs. GD Behavior

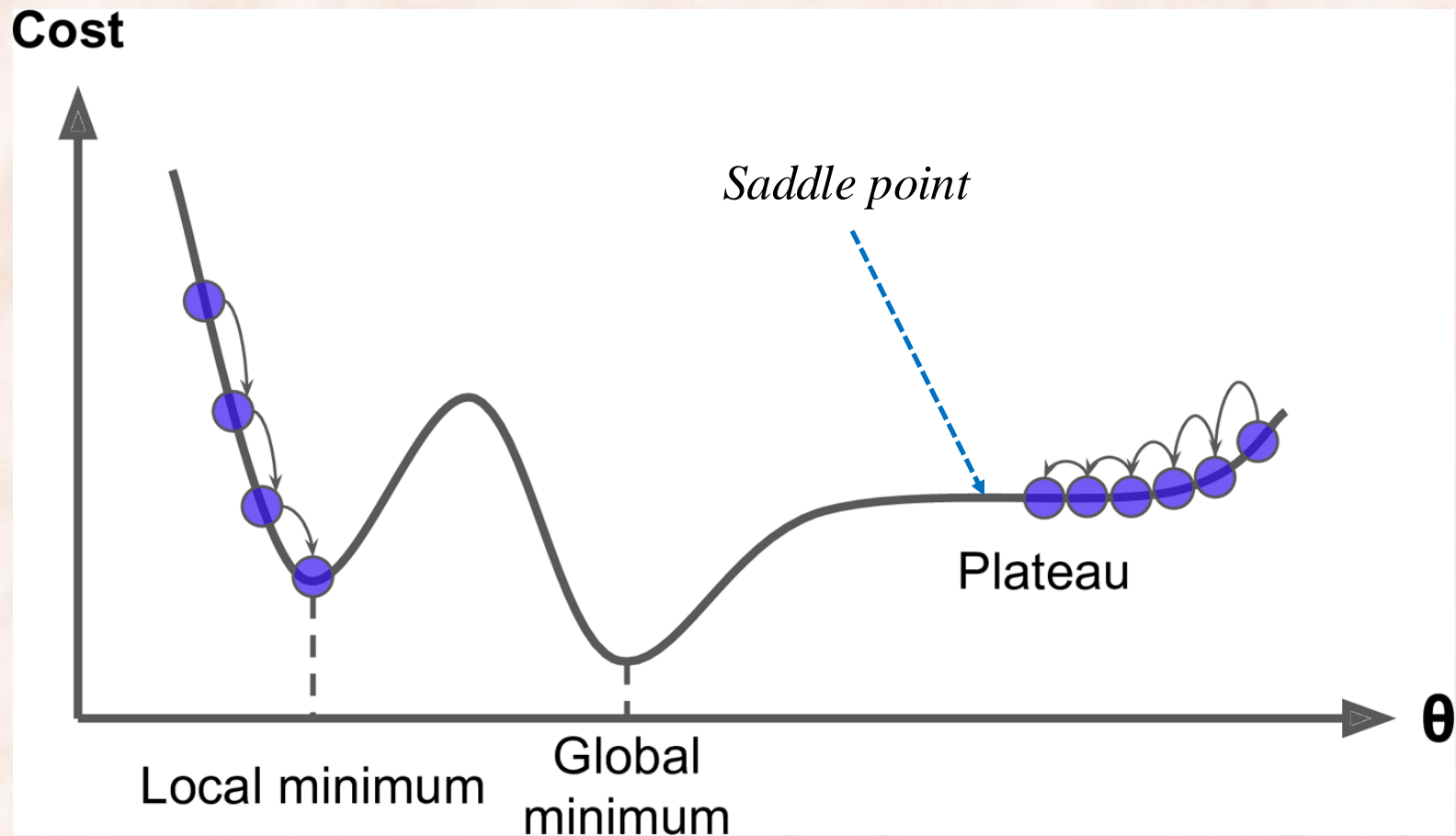


<http://scs.ryerson.ca/~aharley/neural-networks/>

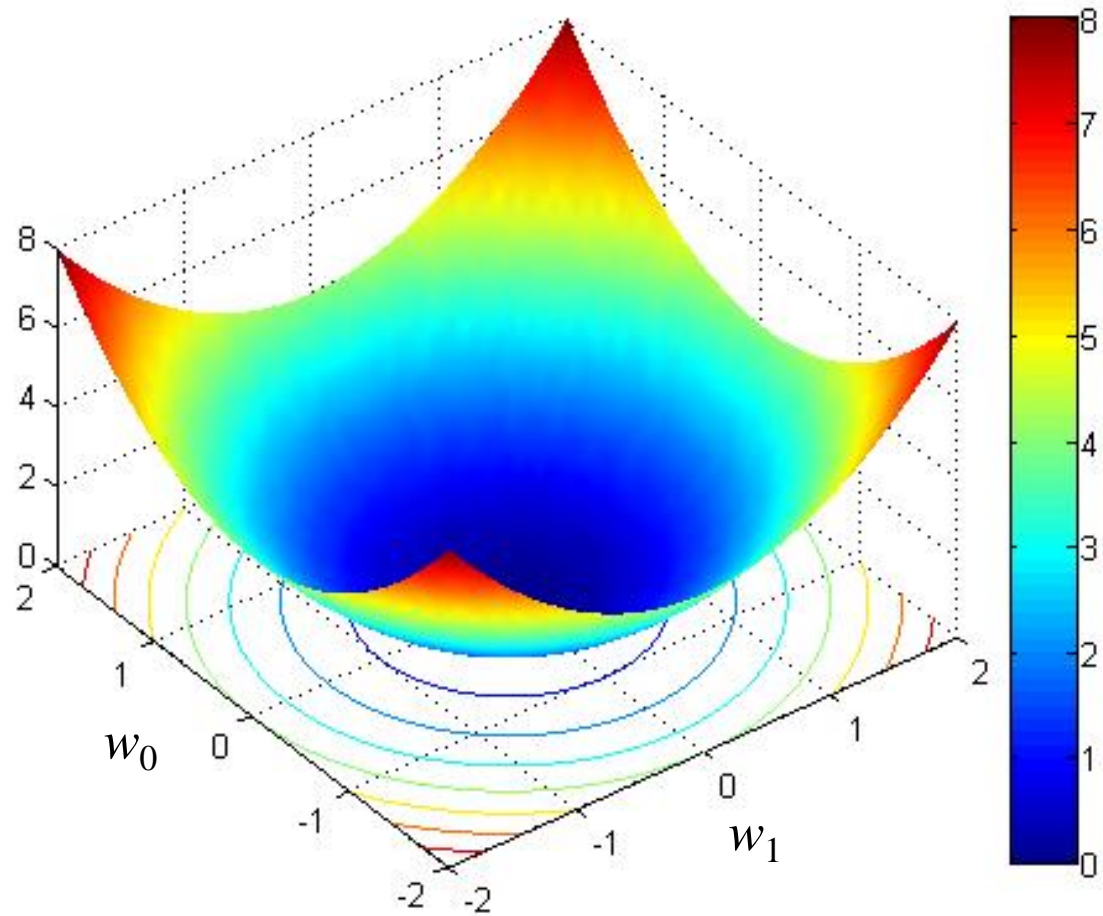
The Learning Rate

- How big should the **learning rate** be?
 - If learning rate too big => oscillating behavior.
 - If learning rate too small => hinders convergence.
- Use **line search** (backtracking line search, conjugate gradient, ...).
- Use **second order methods** (Newton's method, L-BFGS, ...).
 - Requires computing or estimating the Hessian.
- Use a simple learning rate **annealing schedule**:
 - Start with a relatively large value for the learning rate.
 - Decrease the learning rate as a function of the number of epochs or as a function of the improvement in the objective.
- Use **adaptive learning rates**:
 - Adagrad, Adadelta, RMSProp, Adam.

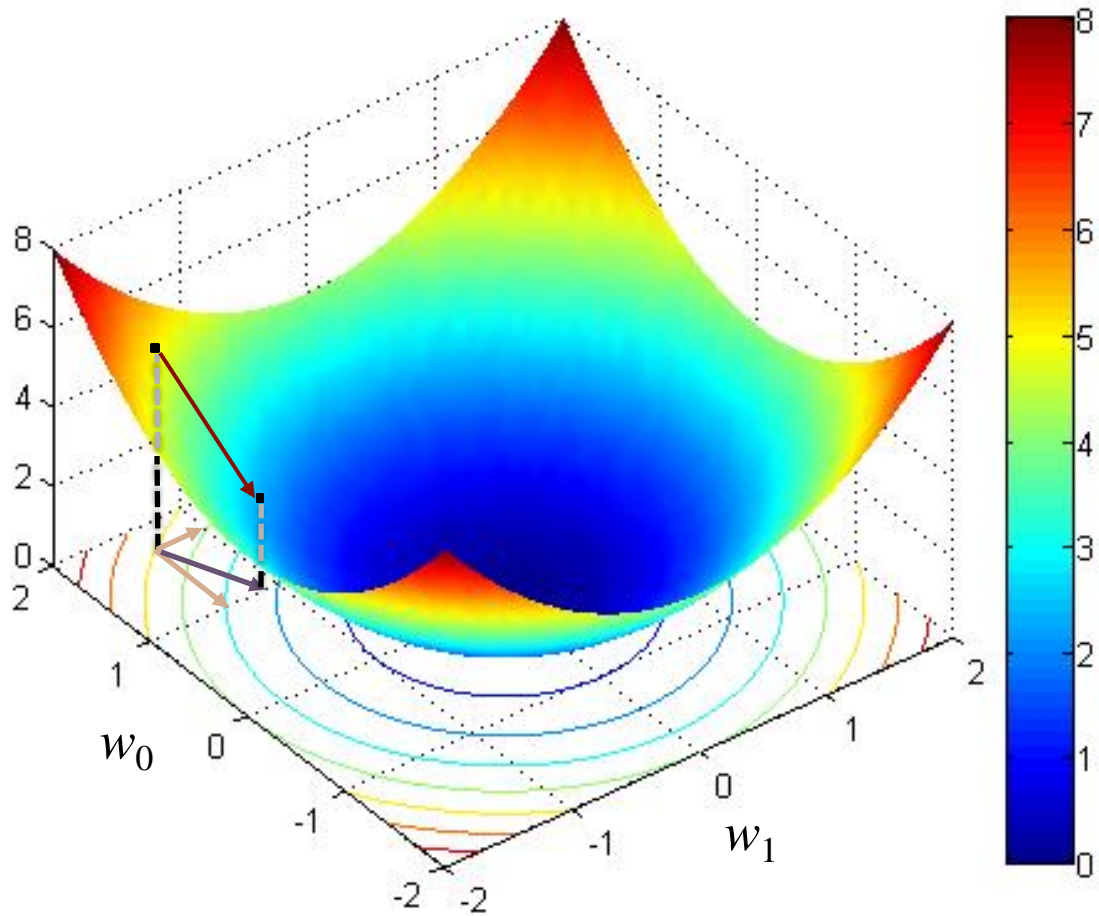
Gradient Descent: Nonconvex Objective



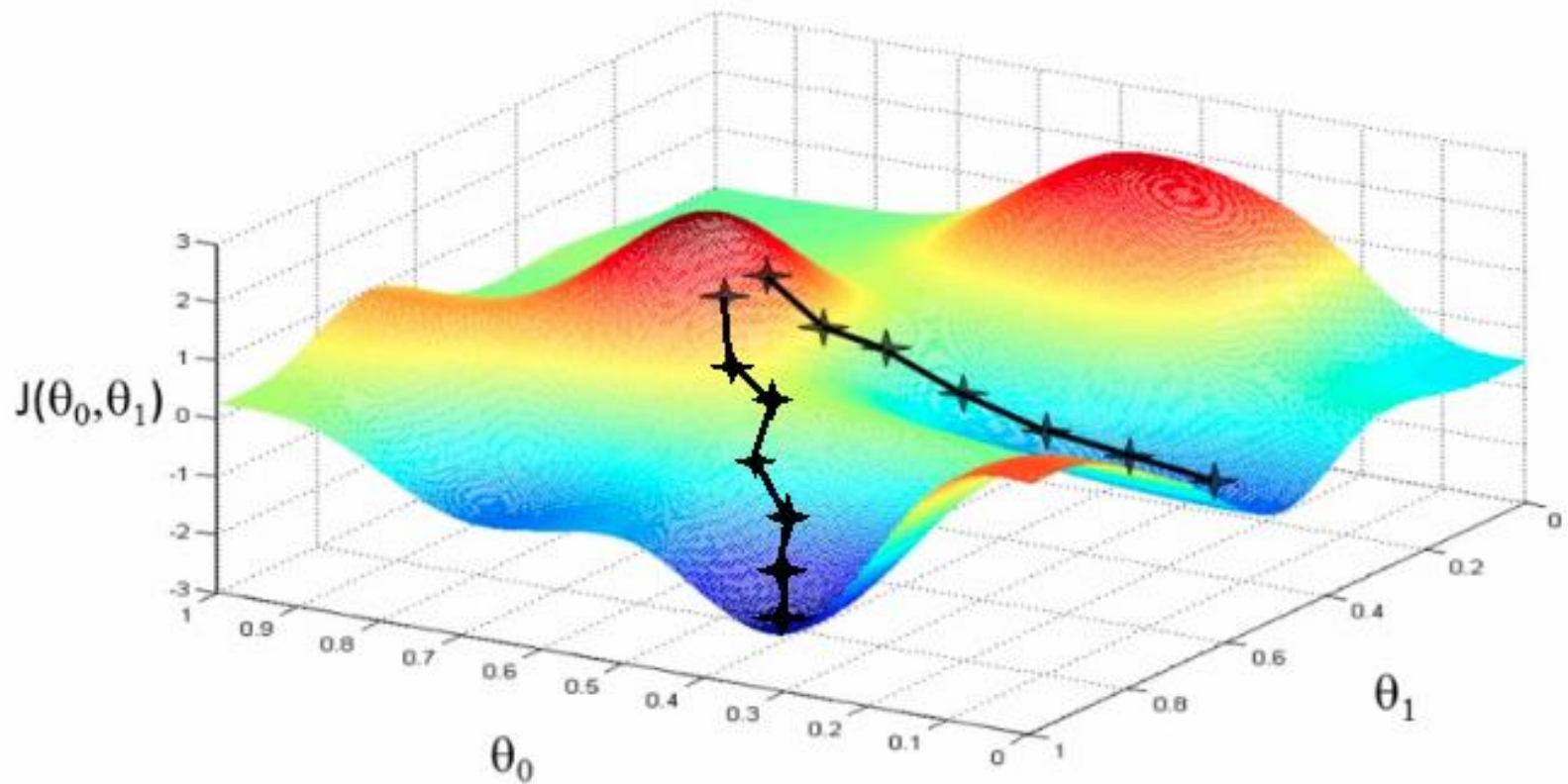
Convex Multivariate Objective



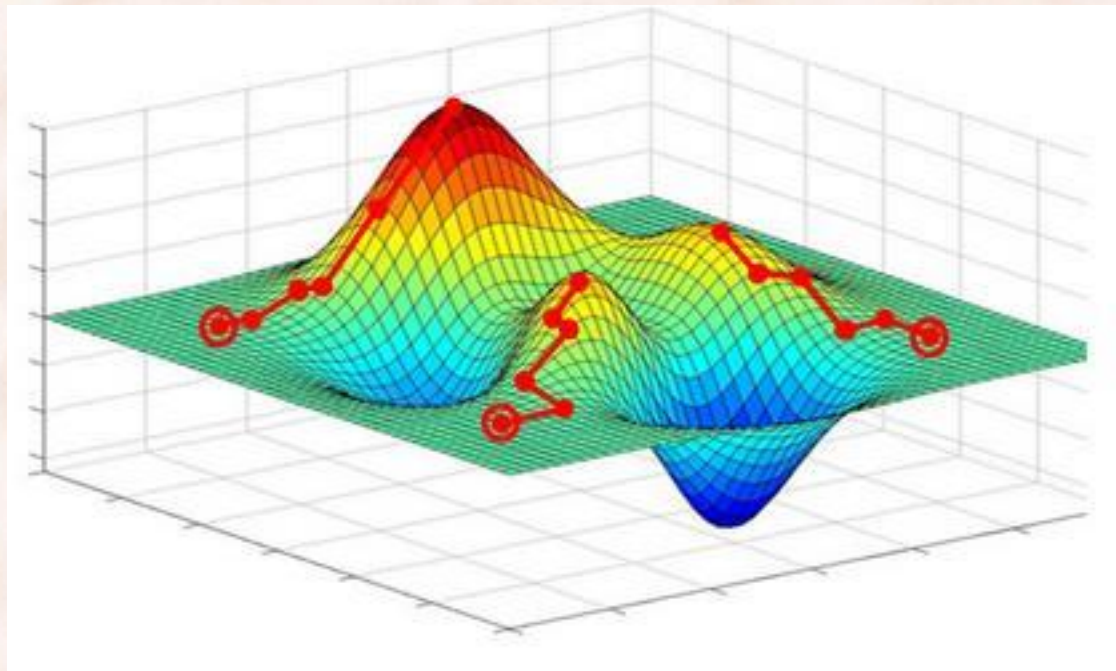
Gradient Step and Contour Lines



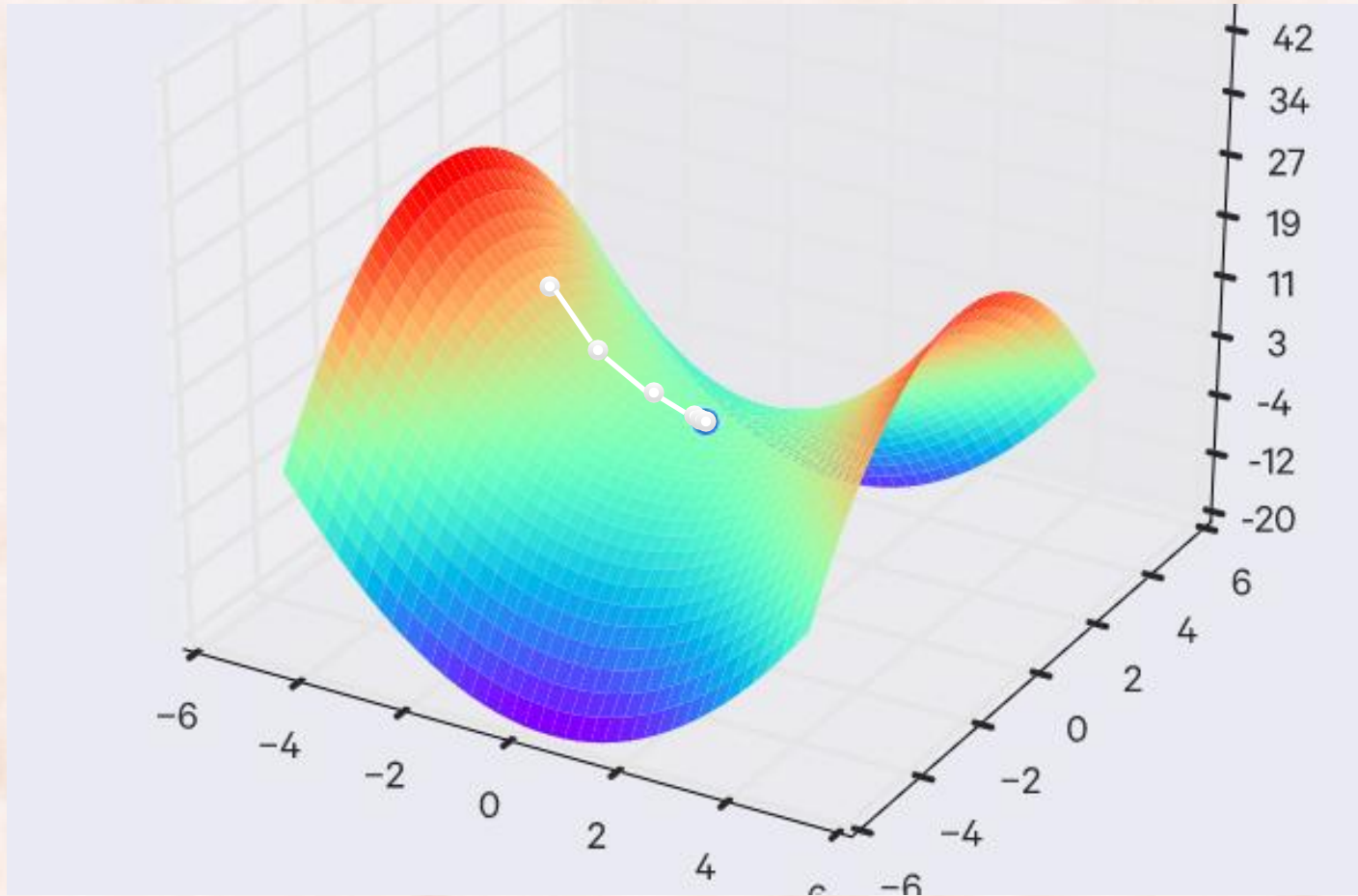
Gradient Descent: Nonconvex Objectives



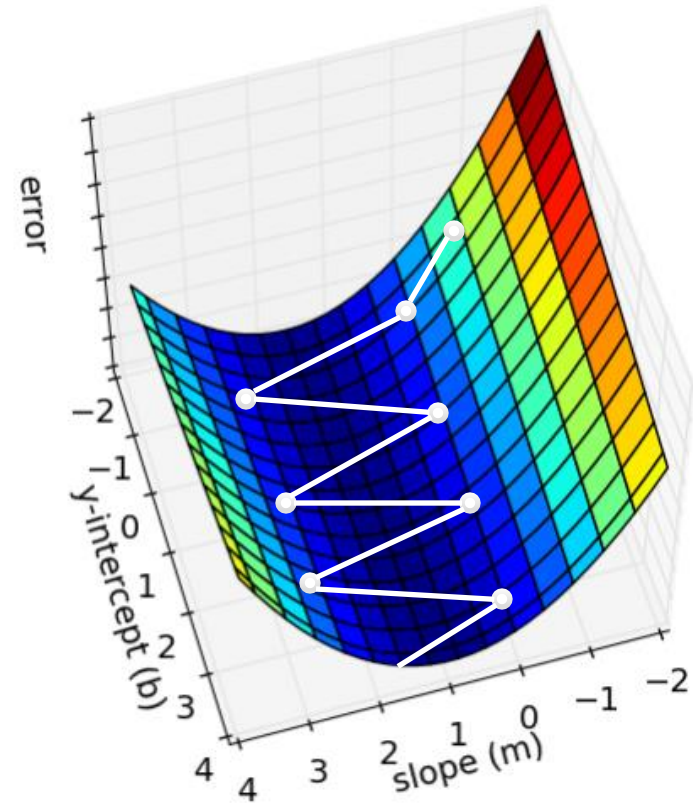
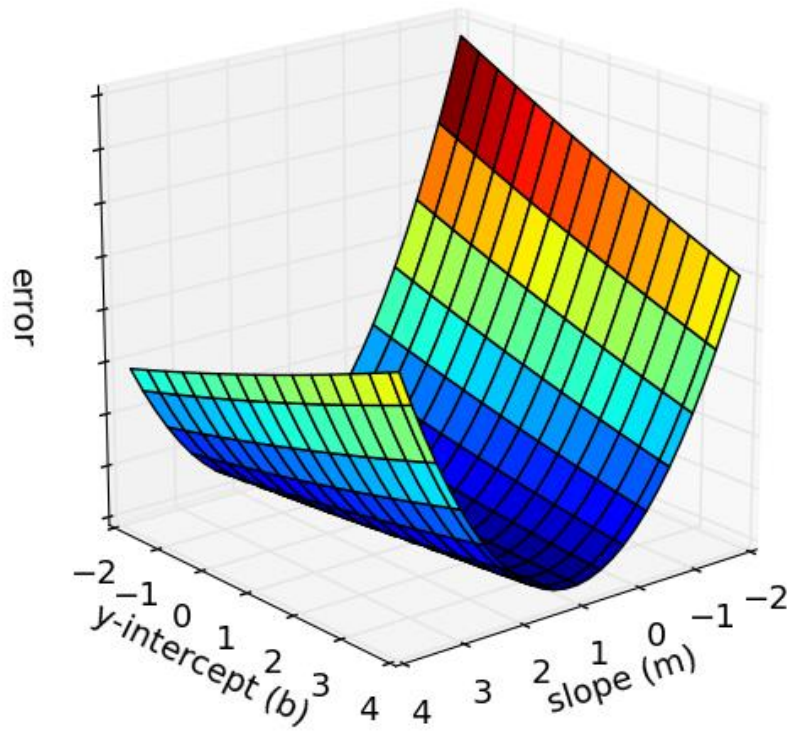
Gradient Descent & Plateaus



Gradient Descent & Saddle Points



Gradient Descent & Ravines



Gradient Descent & Ravines

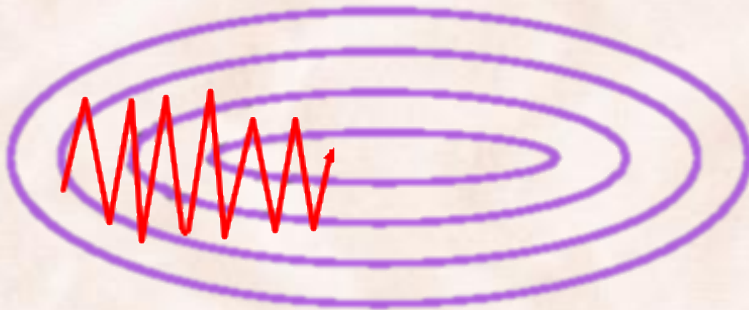
- **Ravines** are areas where the cost surface curves much more steeply in one dimension than another.
 - Common around local optima.
 - GD oscillates across the slopes of the ravines, making slow progress towards the local optimum along the bottom.
- Use **momentum** to help accelerate GD in the relevant directions and dampen oscillations:
 - Add a fraction of the past **update vector** to the current update vector.
 - The momentum term increases for dimensions whose previous gradients point in the same direction.
 - It reduces updates for dimensions whose gradients change sign.
 - Also reduces the risk of getting stuck in local minima.

Gradient Descent & Momentum

Vanilla Gradient Descent:

$$\mathbf{v}^{\tau+1} = \eta \nabla J(\mathbf{w}^{\tau})$$

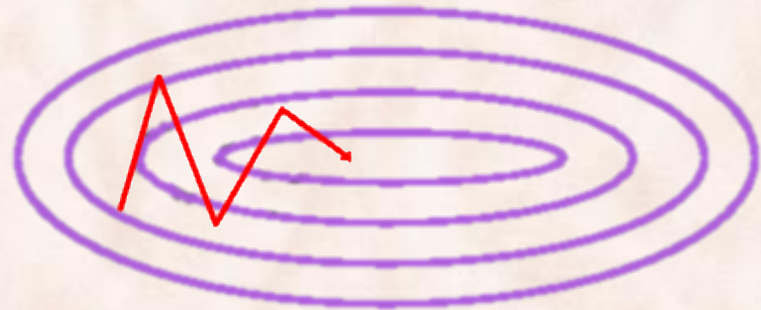
$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$



Gradient Descent w/ Momentum:

$$\mathbf{v}^{\tau+1} = \gamma \mathbf{v}^{\tau} + \eta \nabla J(\mathbf{w}^{\tau})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$



γ is usually set to 0.9 or similar.

The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions.

Batch vs. Stochastic Gradient Descent

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$$

- Depending on how much data is used to compute the gradient at each step:
 - **Batch gradient descent:**
 - Use all the training examples.
 - **Stochastic gradient descent (SGD).**
 - Use one training example, update after each.
 - **Minibatch gradient descent.**
 - Use a constant number of training examples (minibatch).

Batch Gradient Descent for Linear Regression

- Sum-of-squares error:

$$\hat{y}_n = \mathbf{w}^T \mathbf{x}^{(n)}$$

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y_n)^2$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla J(\mathbf{w}^{\tau})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^{\tau T} \mathbf{x}^{(n)} - t_n) \mathbf{x}^{(n)}$$

Stochastic Gradient Descent for Linear Regression

$$\hat{y}_n = \mathbf{w}^T \mathbf{x}^{(n)}$$

- Sum-of-squares error:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y_n)^2 = \frac{1}{N} \sum_{n=1}^N \text{loss}(\mathbf{w}^\tau, \mathbf{x}^{(n)})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau - \eta \nabla \text{loss}(\mathbf{w}^\tau, \mathbf{x}^{(n)})$$

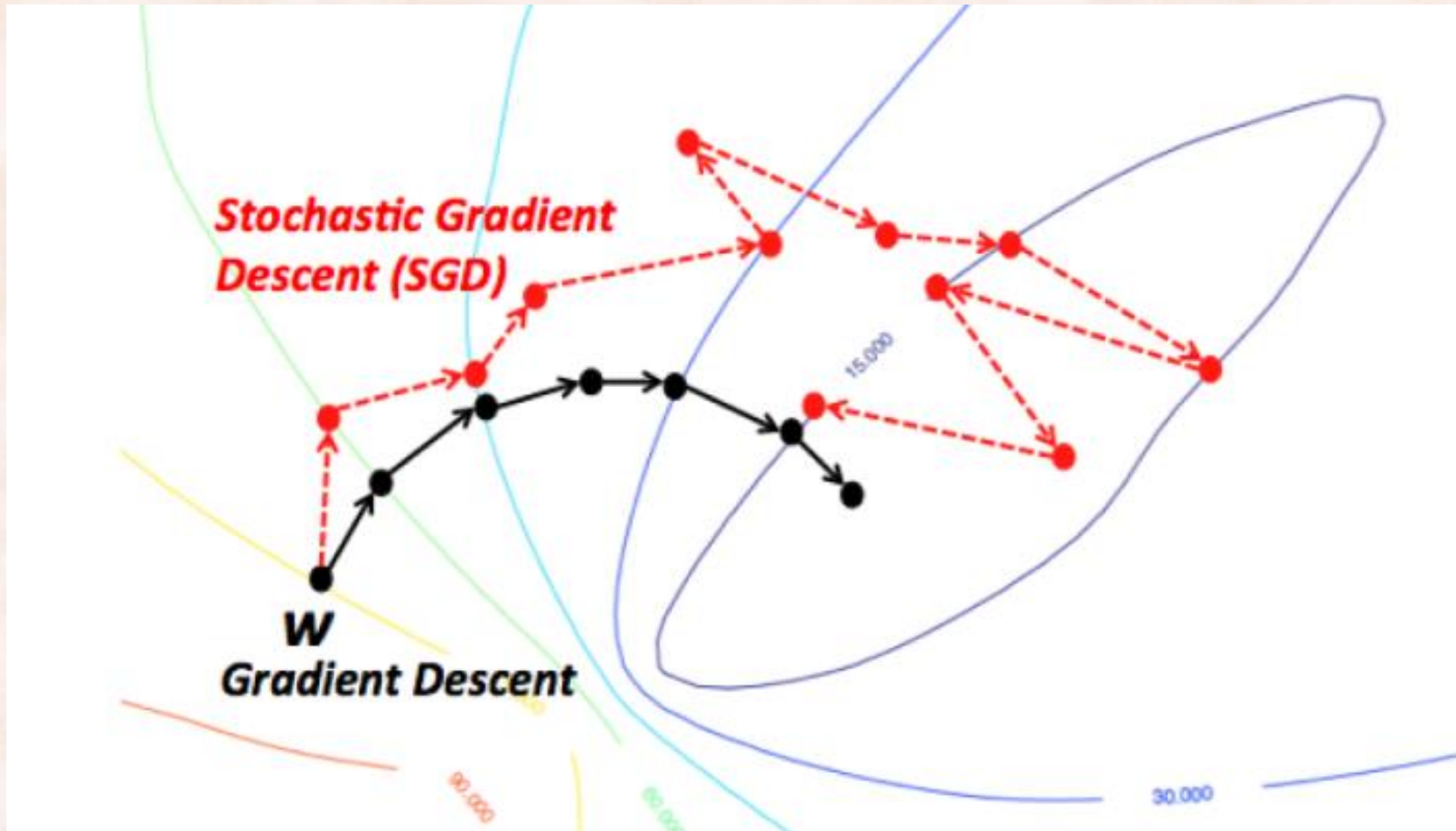
$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau - \eta (\mathbf{w}^T \mathbf{x}^{(n)} - y_n) \mathbf{x}^{(n)}$$

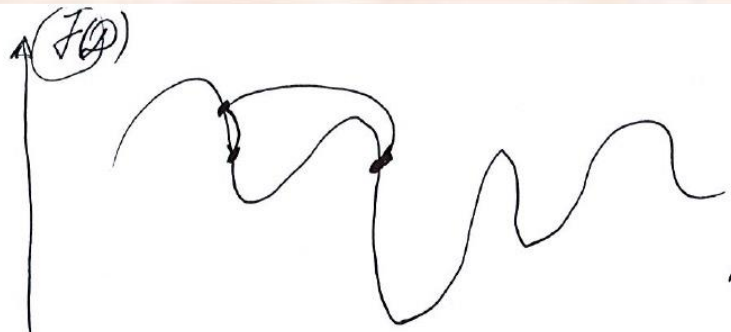
- Update parameters \mathbf{w} after each example, sequentially:
 \Rightarrow the *least-mean-square* (LMS) algorithm.

Batch GD vs. Stochastic GD

- Accuracy:
- Time complexity:
- Memory complexity:
- Online learning:

Batch GD vs. Stochastic GD

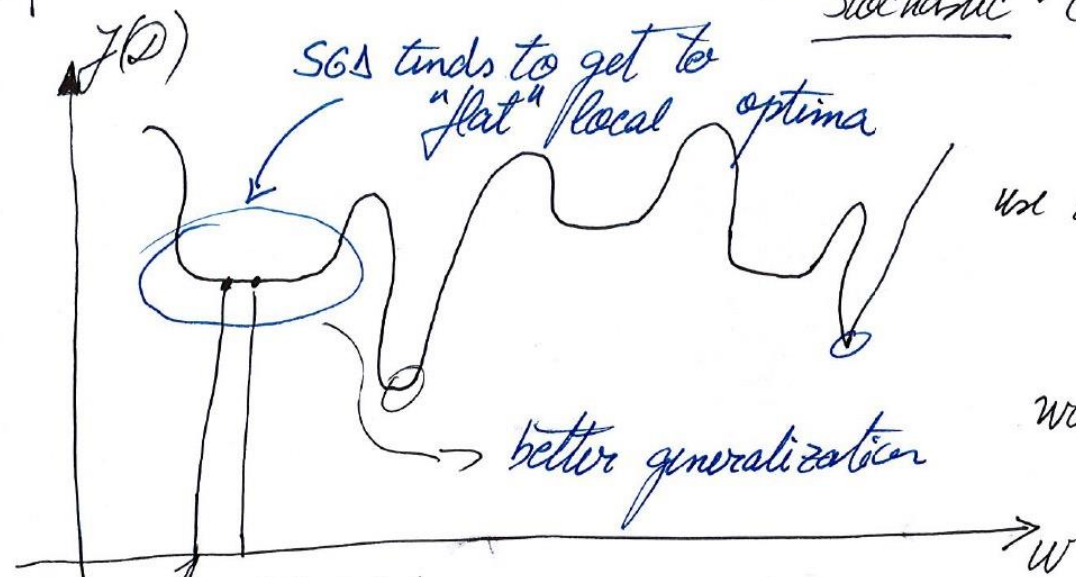




Stochastic w/ mini-batch B of entire dataset D , $B \subseteq D$

batch GA: Compute $J(D)$ and use $\frac{\partial J(D)}{\partial w}$

Stochastic: Compute $J(B)$, use $\frac{\partial J(B)}{\partial w}$

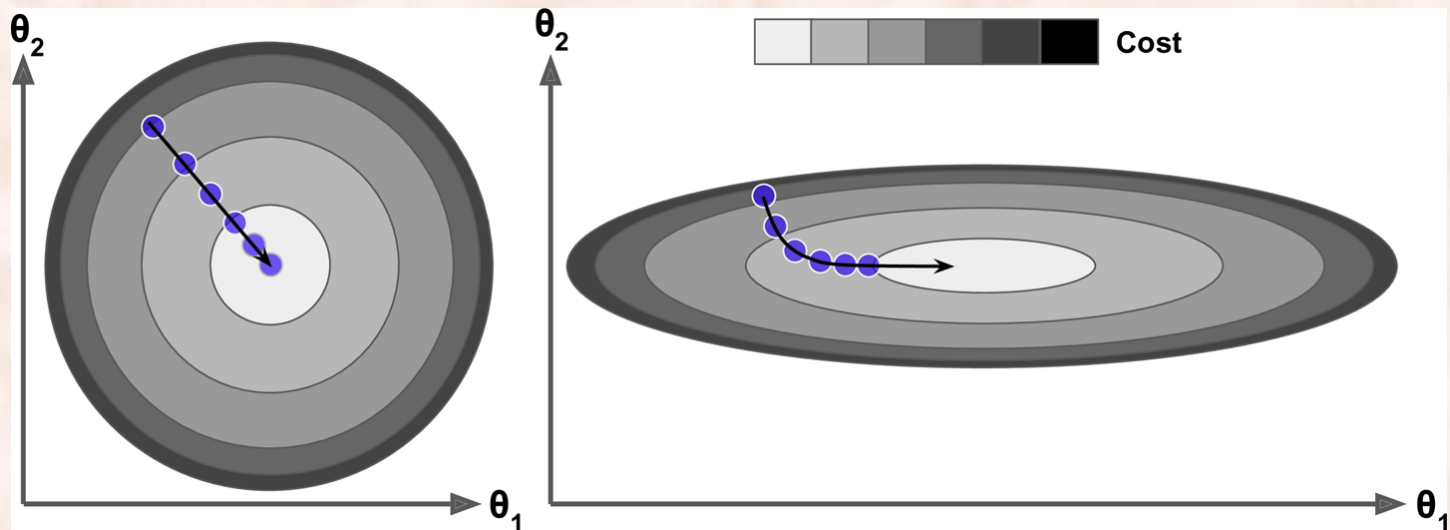


use it as a substitute for $\frac{\partial J(D)}{\partial w}$, want to minimize $J(D)$.

0.9873
0.98

Pre-processing Features

- Features may have very different scales, e.g. $x_1 = \text{rooms}$ vs. $x_2 = \text{size in sq ft}$.
 - **Right** (*different scales*): GD goes first towards the bottom of the bowl, then slowly along an almost flat valley.
 - **Left** (*scaled features*): GD goes straight towards the minimum.



Feature Scaling

- **Scaling between $[0, 1]$ or $[-1, +1]$:**
 - For each feature x_j , compute min_j and max_j **over the training examples.**
 - Scale x_j as follows: $\hat{x}_j = \frac{x_j - min_j}{max_j - min_j}$
- **Scaling to standard normal distribution:**
 - For each feature x_j , compute sample μ_j and sample σ_j **over the training examples.**
 - Scale x_j as follows: $\hat{x}_j = \frac{x_j - \mu_j}{\sigma_j}$
- **Use the same scaling factors at test time:**
 - Clip to min_j and max_j .

n	x_j	\hat{x}_j
1	2	0.2
2	1	0
3	6	1
<hr/>		
4	8	?

Train (rows 1-3)
Test (row 4)

$\min_j = 1$ $\max_j = 6$ Scaling to $[0, 1]$

$$n=1: \hat{x}_j = \frac{x_j - \min_j}{\max_j - \min_j} = \frac{2-1}{6-1} = \frac{1}{5}$$

$$n=2: \hat{x}_j = \frac{1-1}{6-1} = 0$$

$$n=3: \hat{x}_j = \frac{6-1}{6-1} = 1$$

$$n=4: \hat{x}_j = \frac{8-1}{6-1} = \frac{7}{5} = \underline{\underline{1.4}} \quad (1)$$

clipping: clip to $[0, 1]$ (2)
 $\hat{x}_j = 1$

Standardization

• sample mean μ_j :

$$\mu_j = \frac{1}{N} \sum_{n=1}^N x_j^{(n)}$$

• sample std. dev σ_j :

$$\sigma_j = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_j^{(n)} - \mu_j)^2}$$

$$\hat{x}_j = \frac{x_j - \mu_j}{\sigma_j}$$

Gradient Descent vs. Normal Equations

- **Gradient Descent:**

- Need to select learning rate η .
- May need many iterations:
 - Can do *Early Stopping* on validation data for regularization.
- Scalable when number of training examples N is large.

- **Normal Equations:**

- No iterations \Rightarrow easy to code.
- Computing $(X^T X)^{-1}$ has cubic time complexity \Rightarrow slow for large N .
- $X^T X$ may be singular:
 1. Redundant (linearly dependent) features.
 2. #features $>$ #examples \Rightarrow do *feature selection* or *regularization*.

Implementation: Vectorization

- **Version 1:** Compute gradient component-wise.

$$\nabla J(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y_n) \mathbf{x}^{(n)}$$

$$\hat{y}_n = \mathbf{w}^T \mathbf{x}^{(n)}$$

```
grad = np.zeros(K)
```

```
for n in range(N):
```

```
    h = w.dot(X[n]) // This NumPy code assumes examples stored in rows of X.
```

```
    temp = h - y[n]
```

```
    for k in range(K):
```

```
        grad[k] = grad[k] + temp * X[n,k]
```

```
for k in range(K):
```

```
    grad[k] = grad[k] / N
```

Implementation: Vectorization

- **Version 2:** Compute gradient, partially vectorized.

$$\nabla J(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y_n) \mathbf{x}^{(n)}$$

$$\hat{y}_n = \mathbf{w}^T \mathbf{x}^{(n)}$$

```
grad = np.zeros(K)
```

```
for n in range(N): // This NumPy code assumes examples stored in rows of X.
```

```
    grad = grad + (w.dot(X[n])) - y[n] * X[n]
```

```
grad = grad / N
```

Implementation: Vectorization

- **Version 3:** Compute gradient, vectorized.

$$\nabla J(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y_n) \mathbf{x}^{(n)}$$

$$\hat{y}_n = \mathbf{w}^T \mathbf{x}^{(n)}$$

$$\text{grad} = \mathbf{X.T.dot}(\mathbf{X.dot}(\mathbf{w}) - \mathbf{y}) / \mathbf{N}$$

NumPy code above assumes examples stored in columns of X

Implementation: Gradient Checking

- Want to minimize $J(\theta)$, where θ is a scalar.

- Mathematical definition of derivative:

$$\frac{d}{dq} J(q) = \lim_{e \rightarrow 0} \frac{J(q+e) - J(q-e)}{2e}$$

- Numerical approximation of derivative:

$$\frac{d}{dq} J(q) \approx \frac{J(q+e) - J(q-e)}{2e} \quad \text{where } \epsilon = 0.0001$$

Implementation: Gradient Checking

- If θ is a vector of parameters θ_i ,
 - Compute numerical derivative with respect to each θ_i .
 - Aggregate all derivatives into numerical gradient $G_{\text{num}}(\theta)$.
- Compare numerical gradient $G_{\text{num}}(\theta)$ with implementation of gradient $G_{\text{imp}}(\theta)$:

$$\frac{\|G_{\text{num}}(\theta) - G_{\text{imp}}(\theta)\|}{\|G_{\text{num}}(\theta) + G_{\text{imp}}(\theta)\|} \leq 10^{-6}$$

Supplemental Topics

Gradient Descent Optimization Algorithms

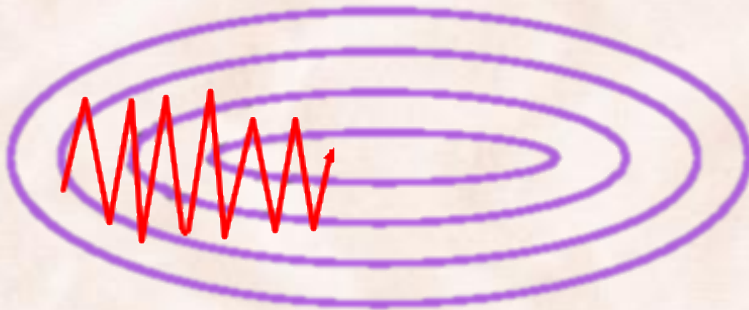
- **Momentum.**
- **Nesterov Accelerated Gradient (NAG).**
- Adaptive learning rates methods:
 - Idea is to perform larger updates for infrequent params and smaller updates for frequent params, by accumulating previous gradient values for each parameter.
 - **Adagrad:**
 - Divide update by sqrt of sum of squares of past gradients.
 - **Adadelta.**
 - **RMSProp.**
 - **Adaptive Moment Estimation (Adam)**

Gradient Descent & Momentum

Vanilla Gradient Descent:

$$\mathbf{v}^{\tau+1} = \eta \nabla J(\mathbf{w}^{\tau})$$

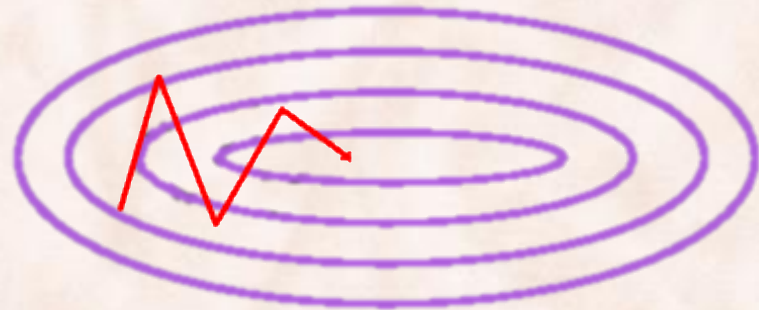
$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$



Gradient Descent w/ Momentum:

$$\mathbf{v}^{\tau+1} = \gamma \mathbf{v}^{\tau} + \eta \nabla J(\mathbf{w}^{\tau})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$



γ is usually set to 0.9 or similar.

The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions.

Momentum & Nesterov Accelerated Gradient

GD with Momentum:

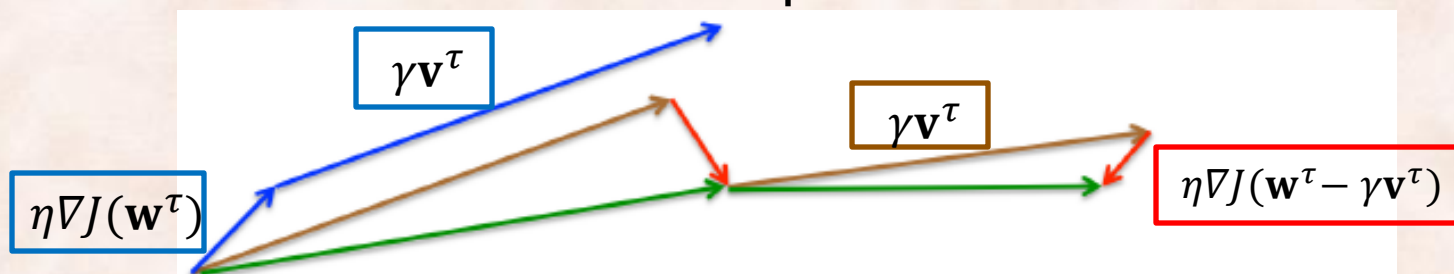
$$\mathbf{v}^{\tau+1} = \gamma \mathbf{v}^{\tau} + \eta \nabla J(\mathbf{w}^{\tau})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$

Nesterov Accelerated Gradient:

$$\mathbf{v}^{\tau+1} = \gamma \mathbf{v}^{\tau} + \eta \nabla J(\mathbf{w}^{\tau} - \gamma \mathbf{v}^{\tau})$$

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \mathbf{v}^{\tau+1}$$



Nesterov update (Source: G. Hinton's lecture 6c)

By making an anticipatory update, NAGs prevents GD from going too fast
 \Rightarrow significant improvements when training RNNs.

AdaGrad

- Optimized for problems with sparse features.
- Per-parameter learning rate: make smaller updates for params that are updated more frequently:

$$w_i = w_i - \eta \frac{g_{t,i}}{\sqrt{\epsilon + G_{t,i}}} \quad \text{where } G_{t,i} = \sum_{\tau=1}^t g_{\tau,i}^2$$
$$g_{t,i} = \frac{\partial J(\mathbf{w})}{\partial w_i}$$

- Require less tuning of the learning rate compared with SGD.

RMSProp

- Element-wise gradient: $g_i^t = \nabla_{w_i} J(\mathbf{w}_t)$
- Gradient is $\mathbf{g}_t = [g_1^t, g_2^t, \dots, g_K^t]$
- Element-wise square gradient: $\mathbf{g}_t^2 = \mathbf{g}_t \circ \mathbf{g}_t$

RMSProp:

$$\mathbf{E}_t[\mathbf{g}^2] = \gamma \mathbf{E}_{t-1}[\mathbf{g}^2] + (1 - \gamma) \mathbf{g}_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\mathbf{E}_t[\mathbf{g}^2] + \epsilon}} \mathbf{g}_t$$

γ is usually set to 0.9, η is set to 0.001

Adam: Adaptive Moment Estimation

- Maintain an exponentially decaying average of past gradients (1st m.) and past squared gradients (2nd m.):
 - 1) $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$
 - 2) $\mathbf{v}_t = \beta_1 \mathbf{v}_{t-1} + (1 - \beta_1) \mathbf{g}_t^2$
- Biased towards 0 during initial steps, use bias-corrected first and second order estimates:

- 1) $\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t}$

- 2) $\hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t}$

Adam: Adaptive Moment Estimation

- First and second moment:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

- Bias-correction:

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t} \text{ and } \hat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t}$$

Adam:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}} \hat{\mathbf{m}}_t$$

Visualization

- Adagrad, RMSprop, Adadelta, and Adam are very similar algorithms that do well in similar circumstances.
 - Insofar, **Adam** might be the best overall choice.

