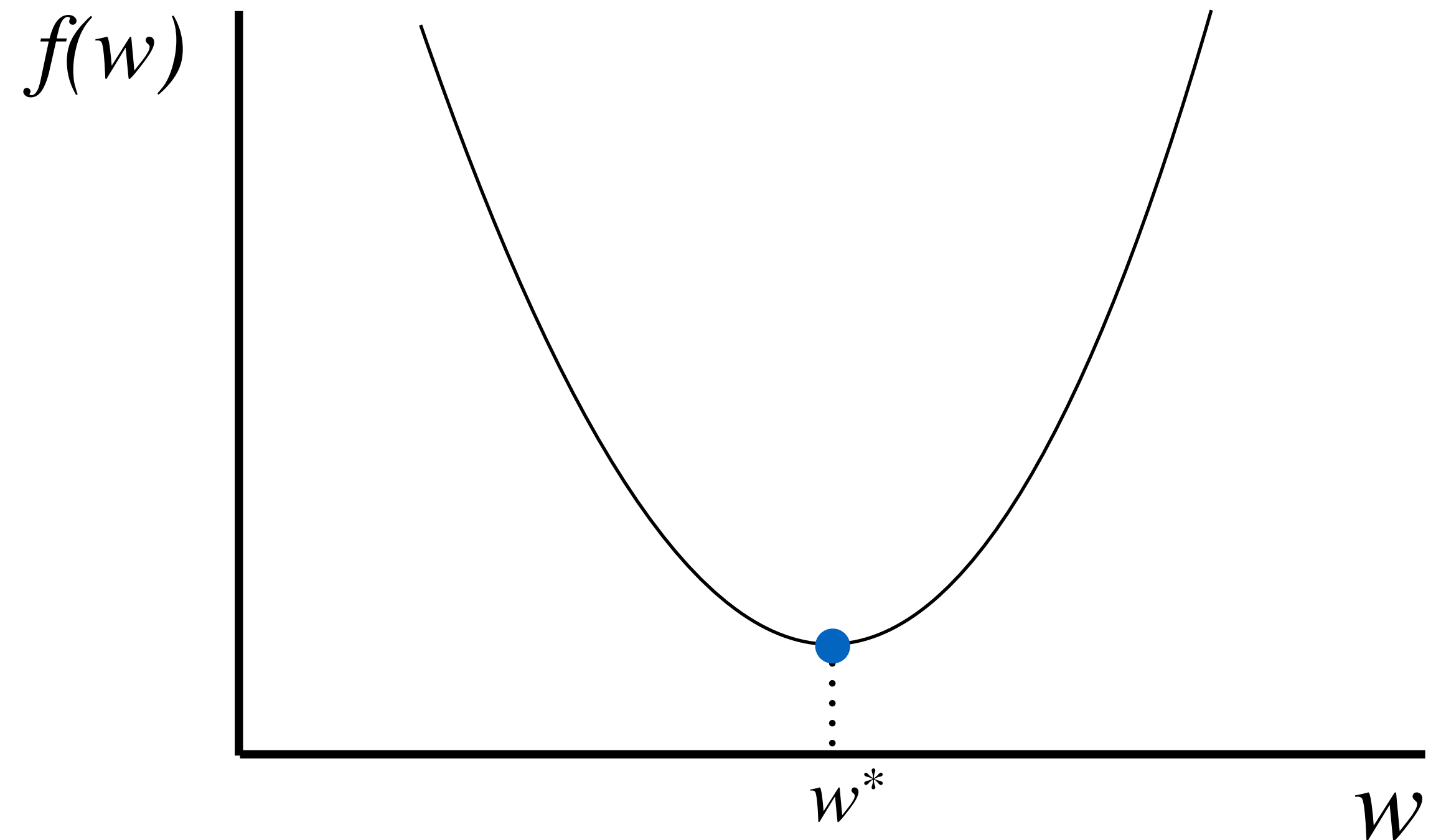# Gradient Descent

# Linear Regression Optimization

**Goal**: Find $\mathbf{w}^*$ that minimizes

$$f(\mathbf{w}) = ||\mathbf{X}\mathbf{w} - \mathbf{y}||_2^2$$
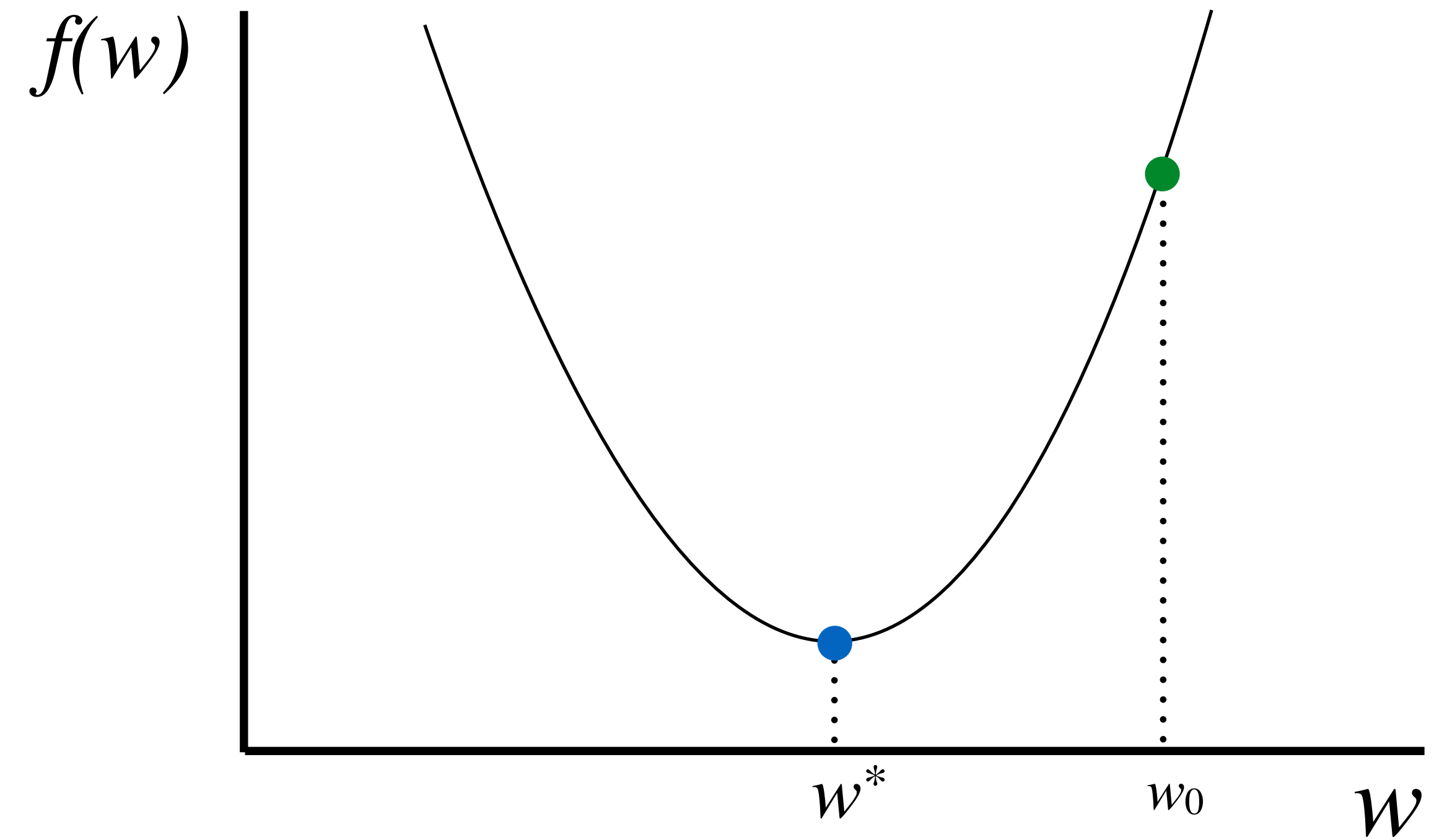
- Closed form solution exists
- Gradient Descent is iterative (Intuition: go downhill!)



Scalar objective: $f(w) = ||w\mathbf{x} - \mathbf{y}||_2^2 = \sum_{j=1}^{n}(wx^{(j)} - y^{(j)})^2$
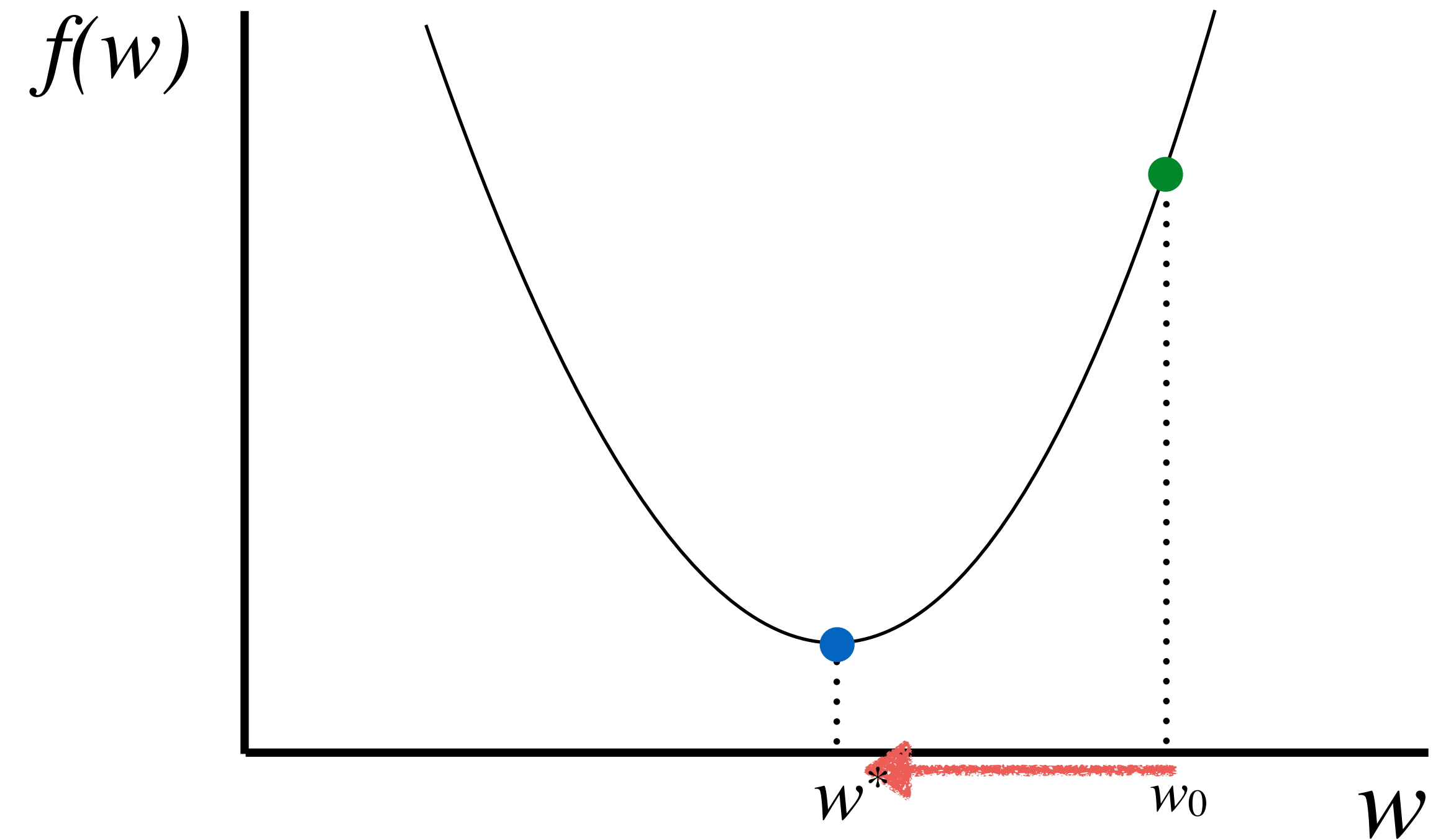
# Gradient Descent

Start at a random point

# Gradient Descent

Start at a random point

Determine a descent direction

# Gradient Descent
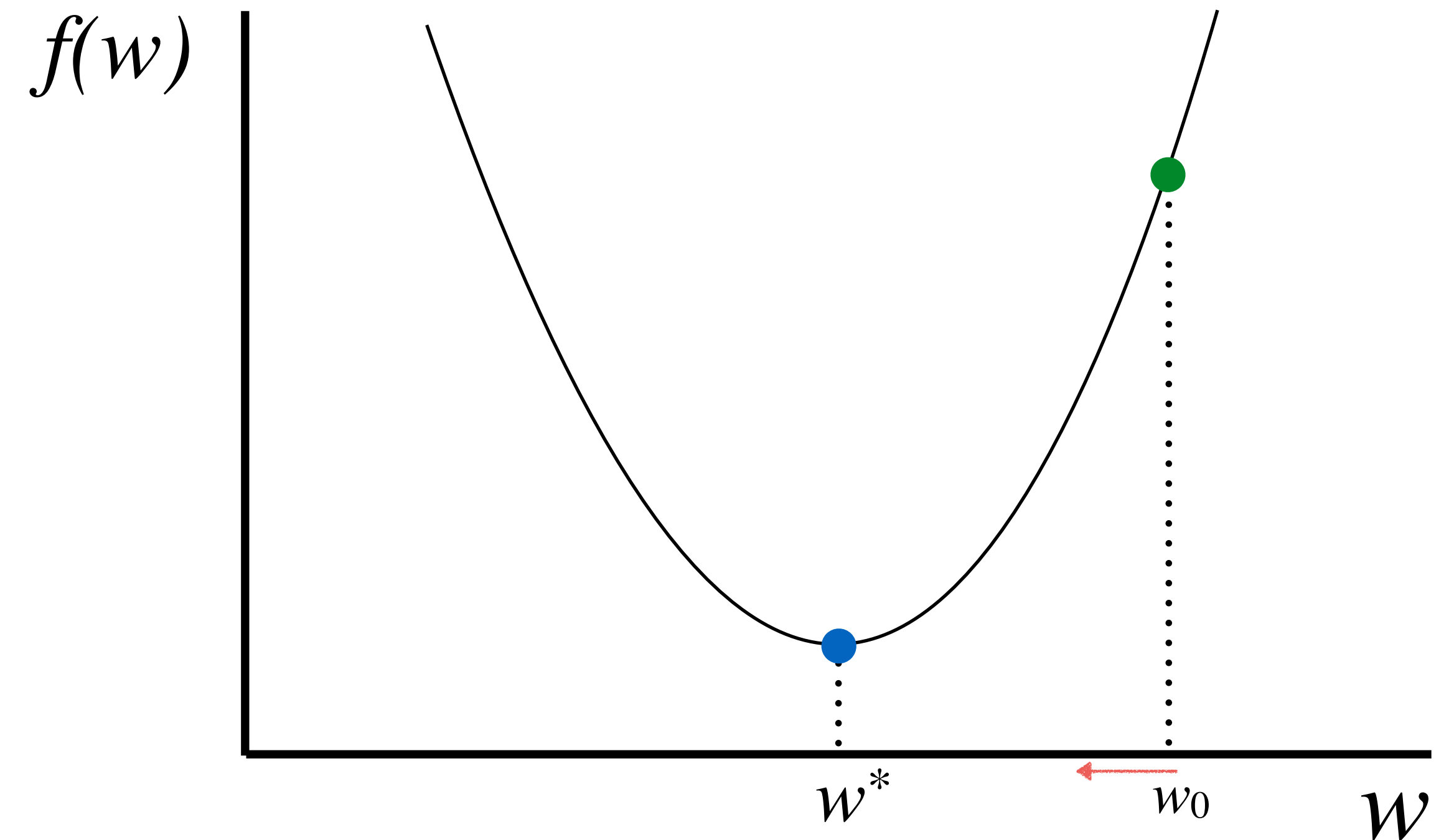
Start at a random point

Determine a descent direction
Choose a step size

$f(w)$



$w^*$  $w_0$  $w$

# Gradient Descent

Start at a random point

Determine a descent direction
Choose a step size
Update

$f(w)$

$w^*$  $w_1$  $w_0$  $w$

# Gradient Descent

Start at a random point
**Repeat**
    Determine a descent direction
    Choose a step size
    Update
**Until** stopping criterion is satisfied

# Gradient Descent

Start at a random point
**Repeat**
| Determine a descent direction
  Choose a step size
  Update
**Until** stopping criterion is satisfied

$f(w)$

$w^*$    $w_1$  $w_0$

$w$

# Gradient Descent

Start at a random point
**Repeat**
   | Determine a descent direction
    Choose a step size
    Update
**Until** stopping criterion is satisfied

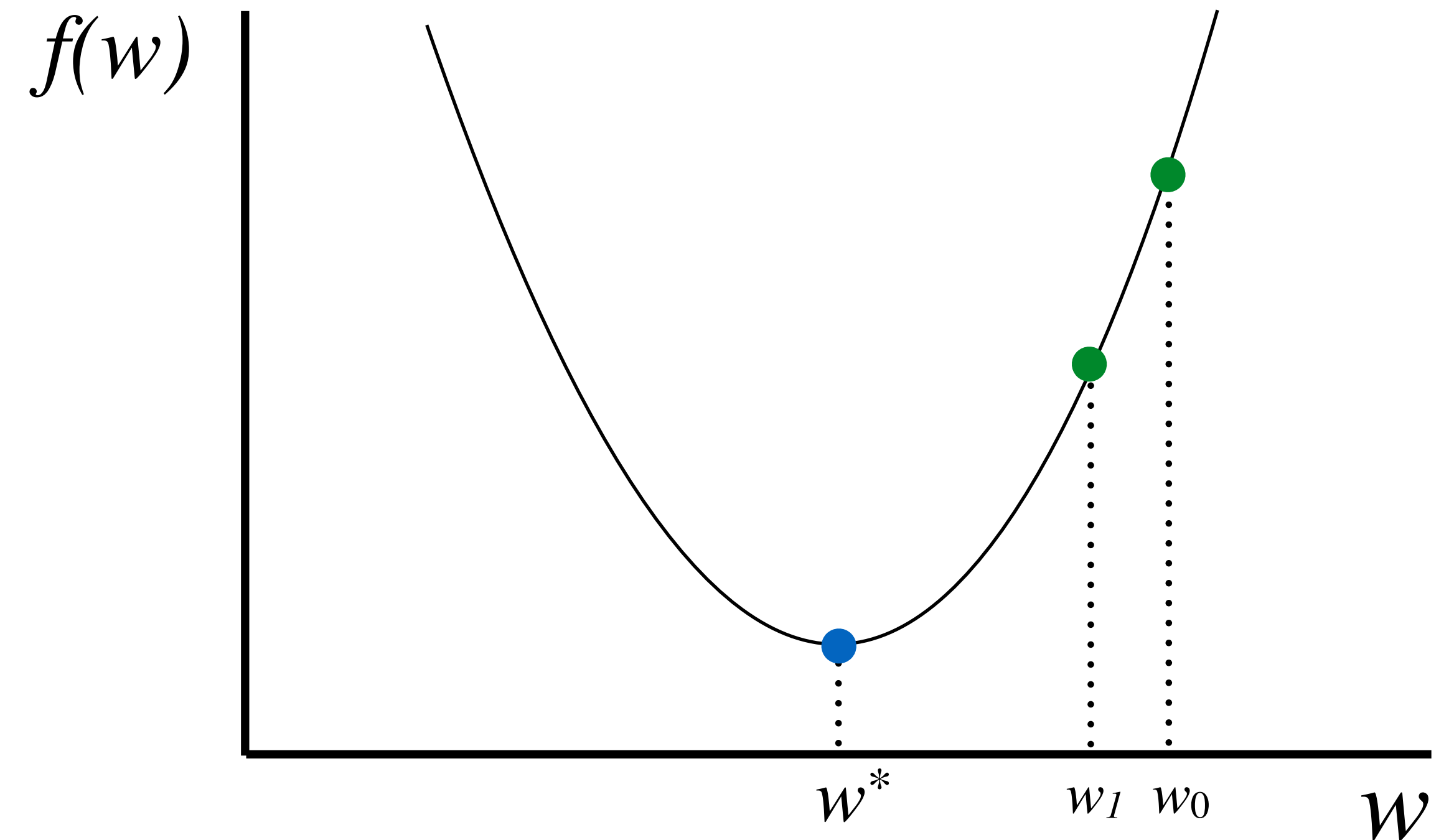# Gradient Descent

Start at a random point
**Repeat**
    Determine a descent direction
    Choose a step size
    Update
**Until** stopping criterion is satisfied

# Gradient Descent

Start at a random point
**Repeat**
    Determine a descent direction
    Choose a step size
    Update
**Until** stopping criterion is satisfied

# Gradient Descent

Start at a random point
**Repeat**
 Determine a descent direction
 Choose a step size
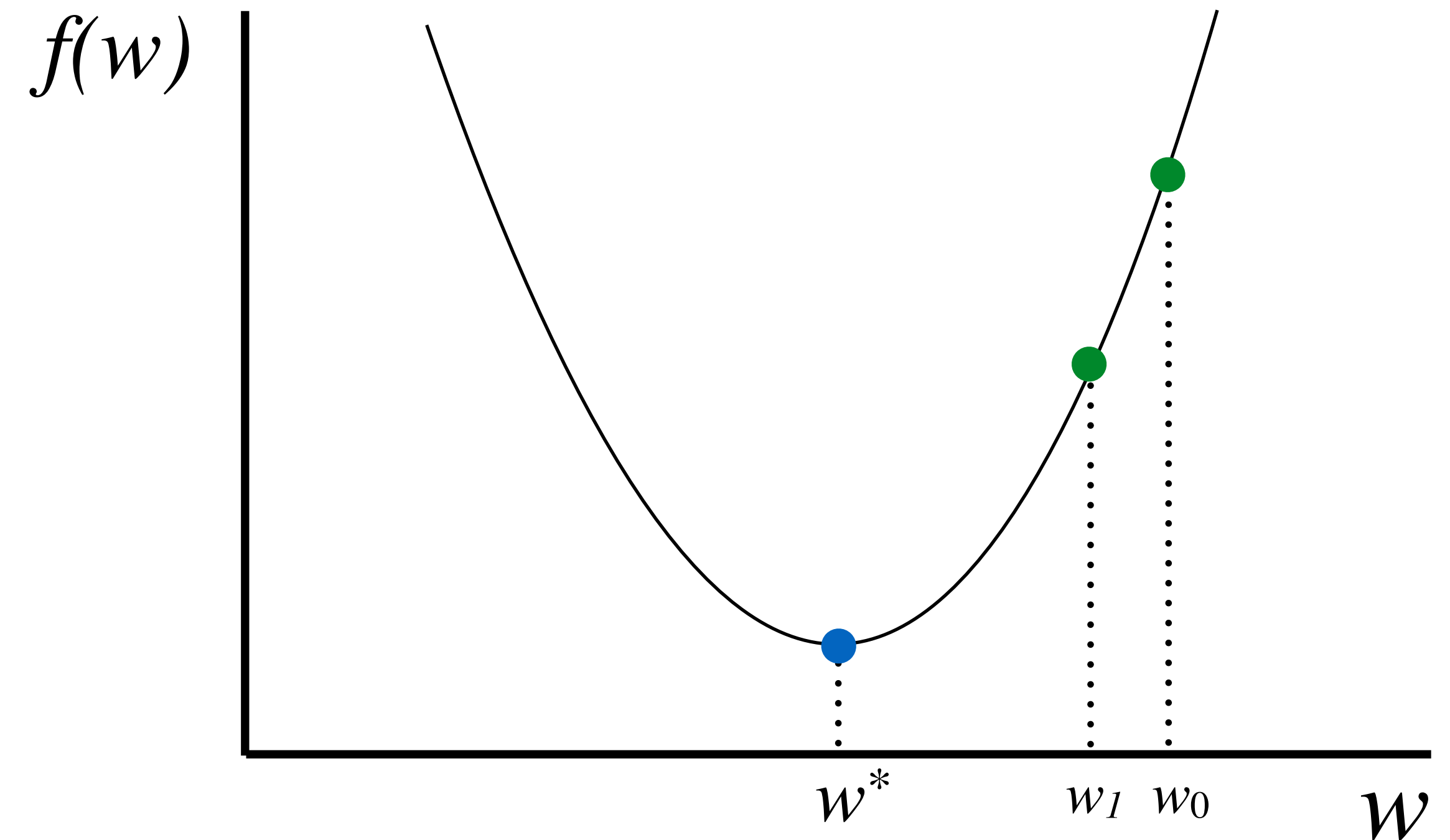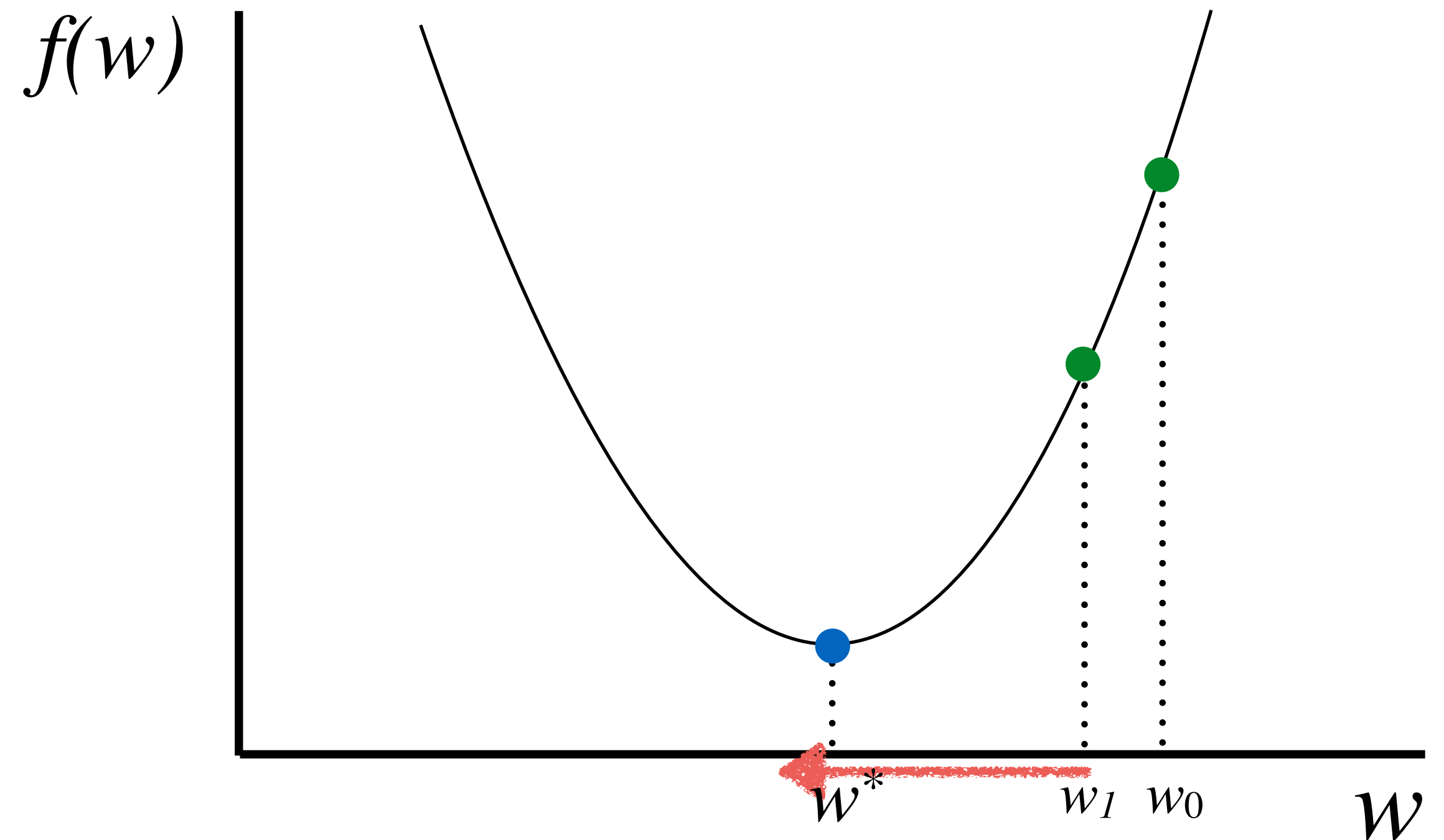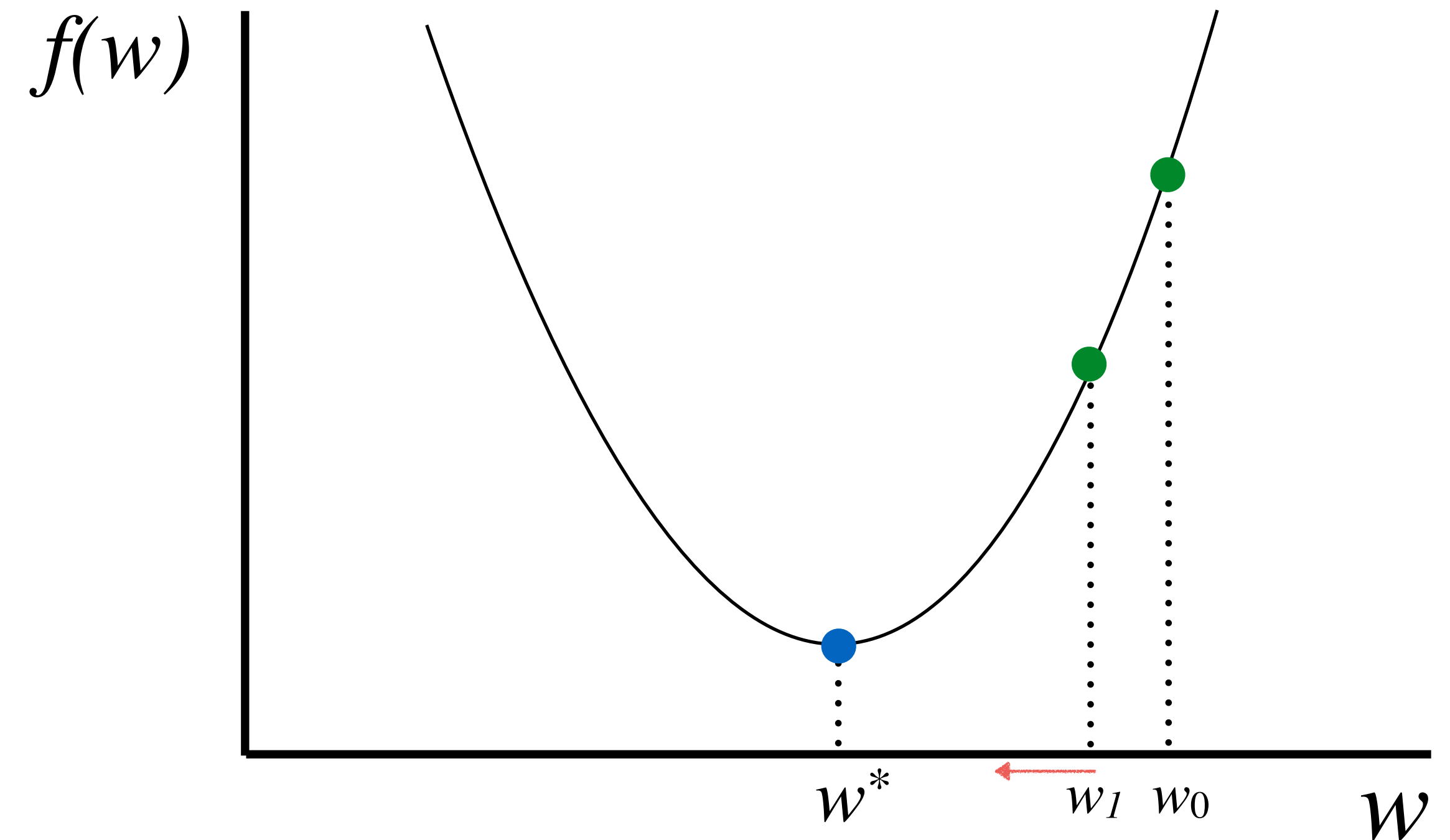 Update
**Until** stopping criterion is satisfied
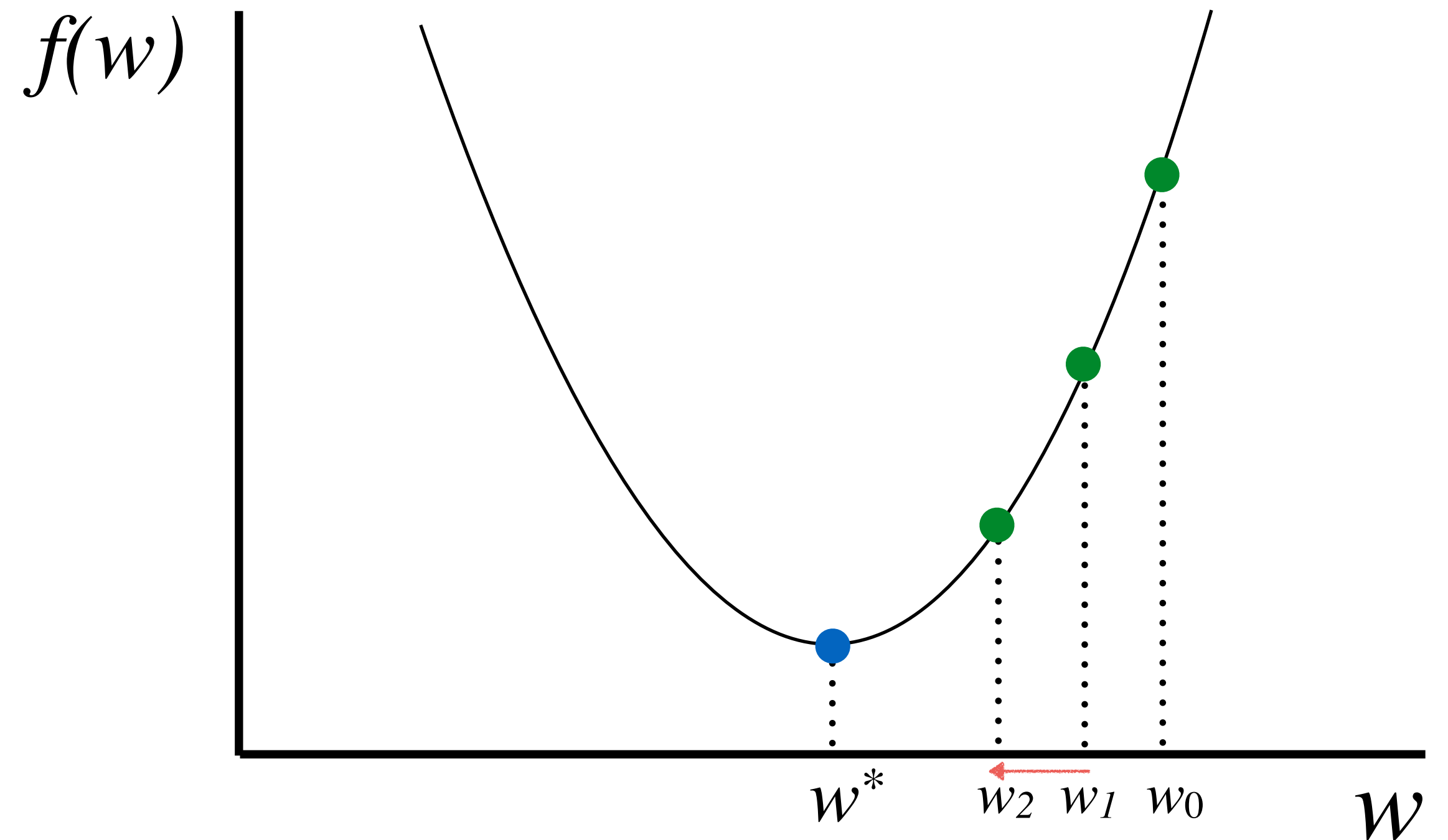
# Gradient Descent

Start at a random point
**Repeat**
    Determine a descent direction
    Choose a step size
    Update
**Until** stopping criterion is satisfied

$f(w)$

$w^*$ $\cdots$ $w_2$ $w_1$ $w_0$

$w$

# Where Will We Converge?



Convex
$f(w)$
$w^*$
$w$

Any local minimum is a global minimum

Non-convex
$g(w)$
$w'$
$w^*$
$w$

Multiple local minima may exist

**Least Squares, Ridge Regression and
Logistic Regression are all convex!**

# Choosing Descent Direction (1D)

$f(w)$

positive ⇒ go left!

zero ⇒ done!

$w^*$   $w_0$   $w$

$f(w)$

negative ⇒ go right!

$w_0$   $w^*$   $w$

We can only move in two directions
Negative slope is direction of descent!

Step Size

**Update Rule:** $w_{i+1} = w_i - \alpha_i \dfrac{df}{dw}(w_i)$

Negative Slope

# Choosing Descent Direction



"Gradient2" by Sarang. Licensed under CC BY-SA 2.5 via Wikimedia Commons
http://commons.wikimedia.org/wiki/File:Gradient2.svg#/media/File:Gradient2.svg

2D Example:

- Function values are in black/white and black represents higher values
- Arrows are gradients

We can move anywhere in $\mathbb{R}^d$

Negative gradient is direction of *steepest* descent!

Step Size

**Update Rule:** $\mathbf{w}_{i+1} = \mathbf{w}_i - \alpha_i \nabla f(\mathbf{w}_i)$

Negative Slope

# Gradient Descent for Least Squares

**Update Rule:** $w_{i+1} = w_i - \alpha_i \dfrac{df}{dw}(w_i)$

Scalar objective: $f(w) = ||w\mathbf{x} - \mathbf{y}||_2^2 = \displaystyle\sum_{j=1}^{n}(wx^{(j)} - y^{(j)})^2$

Derivative: $\dfrac{df}{dw}(w) = 2\displaystyle\sum_{j=1}^{n}(wx^{(j)} - y^{(j)})x^{(j)}$
(chain rule)

Scalar Update: $w_{i+1} = w_i - \alpha_i \displaystyle\sum_{j=1}^{n}(w_i x^{(j)} - y^{(j)})x^{(j)}$
(2 absorbed in $\alpha$ )

Vector Update: $\mathbf{w}_{i+1} = \mathbf{w}_i - \alpha_i \displaystyle\sum_{j=1}^{n}(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

# Choosing Step Size

$f(w)$

$f(w)$

$f(w)$

$w^*$    $w$

$w^*$    $w$

$w^*$    $w$

Too small: converge
very slowly

Too big: overshoot and
even diverge

Reduce size over time

Theoretical convergence results for various step sizes

A common step size is $\alpha_i = \dfrac{\alpha \quad \text{---- Constant}}{n\sqrt{i} \quad \text{---- Iteration \#}}$

# Training Points

# Parallel Gradient Descent for Least Squares

Vector Update: $\mathbf{w}_{i+1} = \mathbf{w}_i - \alpha_i \sum_{j=1}^{n} (\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

Compute summands in parallel!
note: workers must all have $\mathbf{w}_i$

Example: $n = 6$; 3 workers

$\mathbf{w}_{i+1}$

workers:

| $-\mathbf{x}^{(1)}-$ | $-\mathbf{x}^{(3)}-$ | $-\mathbf{x}^{(2)}-$ |
| $-\mathbf{x}^{(5)}-$ | $-\mathbf{x}^{(4)}-$ | $-\mathbf{x}^{(6)}-$ |

O($nd$) Distributed Storage

map:

$(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$    $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$    $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

O($nd$) Distributed Computation    O($d$) Local Storage

reduce:

$\sum_{j=1}^{n} (\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

O($d$) Local Computation    O($d$) Local Storage

```python
> for i in range(numIters):
      alpha_i = alpha / (n * np.sqrt(i+1))
      gradient = train.map(lambda lp: gradientSummand(w, lp))
                      .sum()
      w -= alpha_i * gradient
  return w
```

workers:

$- \mathbf{x}^{(1)} -$
$- \mathbf{x}^{(5)} -$

$- \mathbf{x}^{(3)} -$
$- \mathbf{x}^{(4)} -$

$- \mathbf{x}^{(2)} -$
$- \mathbf{x}^{(6)} -$

O($nd$) Distributed Storage

$\mathbf{w}_{i+1}$

map:

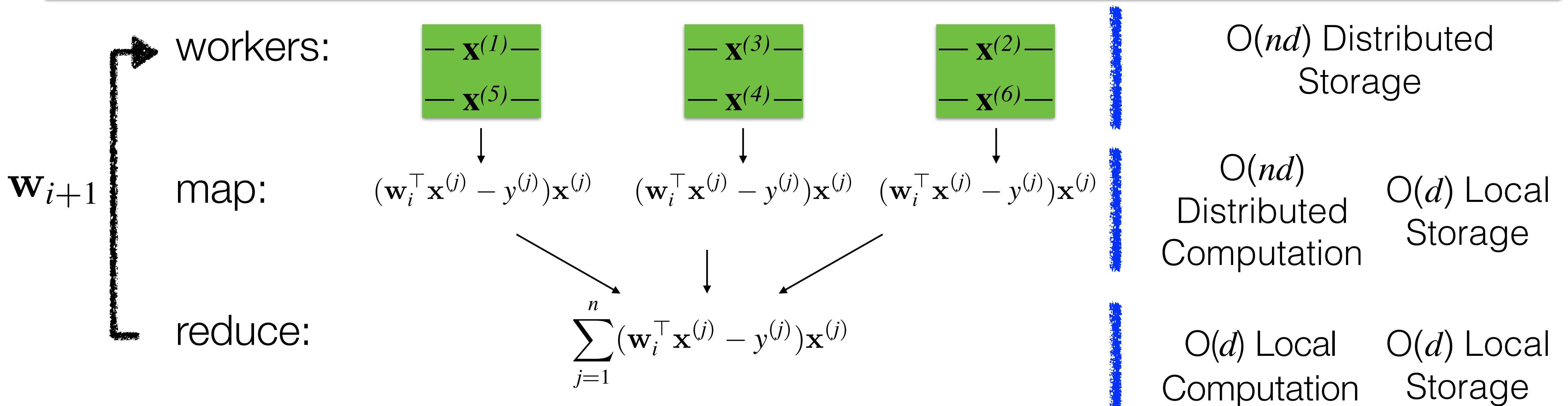$(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$ $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$ $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

O($nd$) Distributed Computation    O($d$) Local Storage

reduce:

$$\sum_{j=1}^{n}(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$$
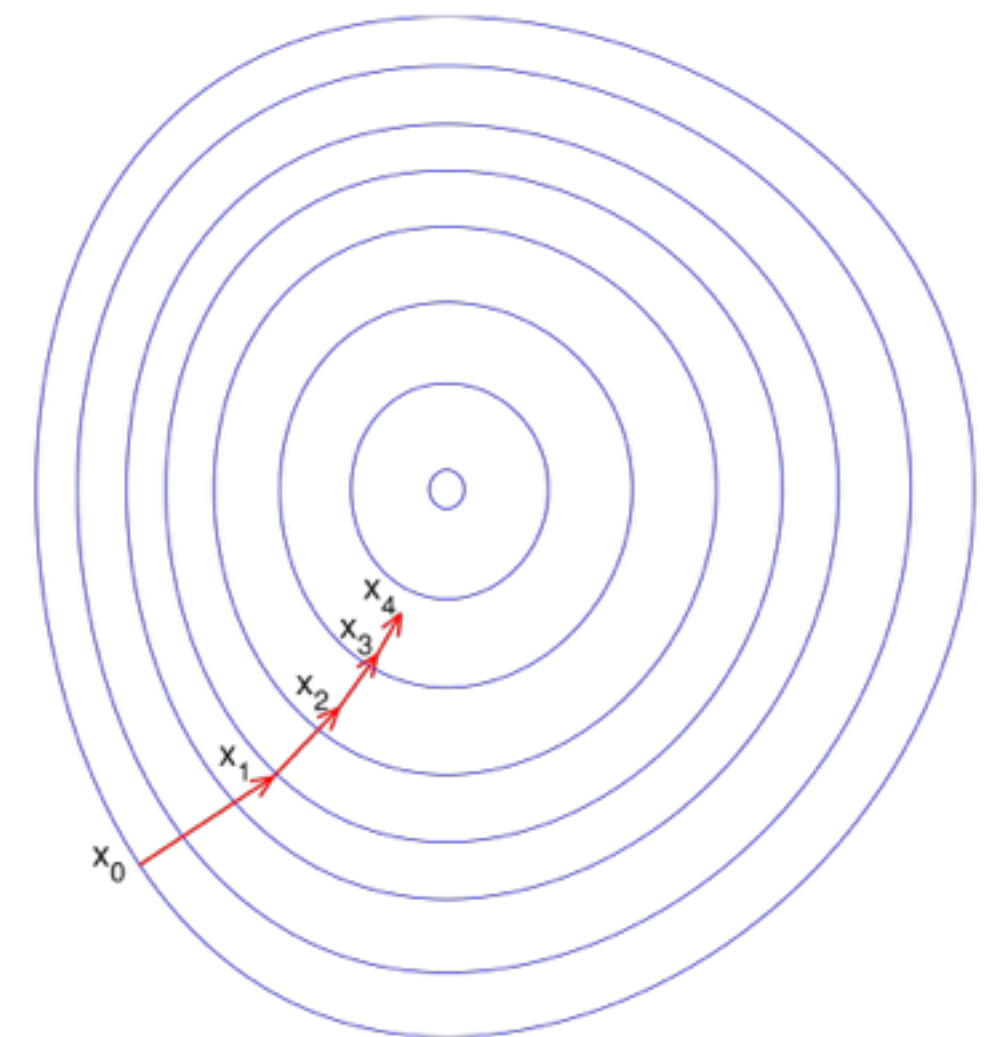
O($d$) Local Computation    O($d$) Local Storage

# Gradient Descent Summary

Pros:
- Easily parallelized
- Cheap at each iteration
- Stochastic variants can make things even cheaper

Cons:
- Slow convergence (especially compared with closed-form)
- **Requires communication across nodes!**
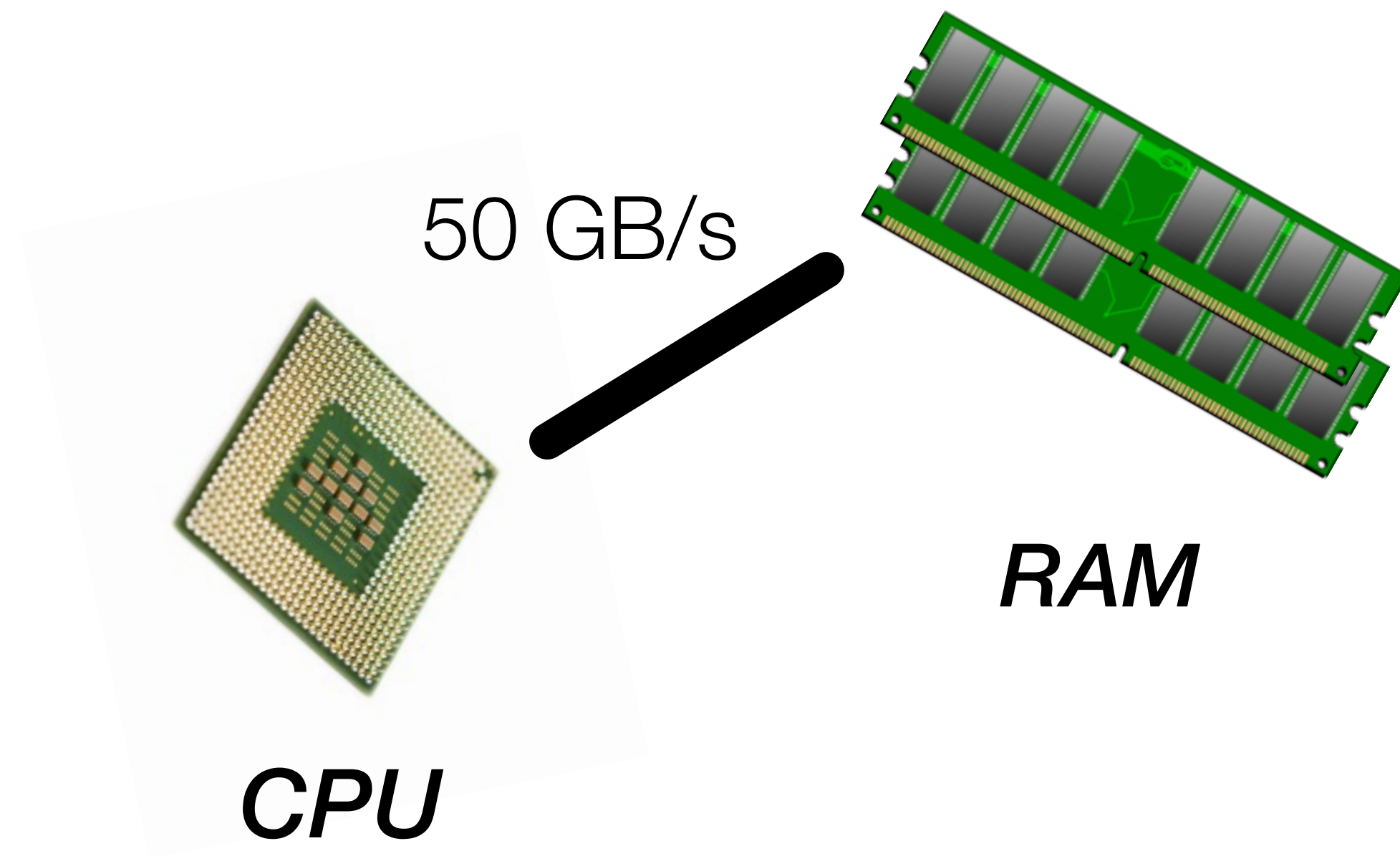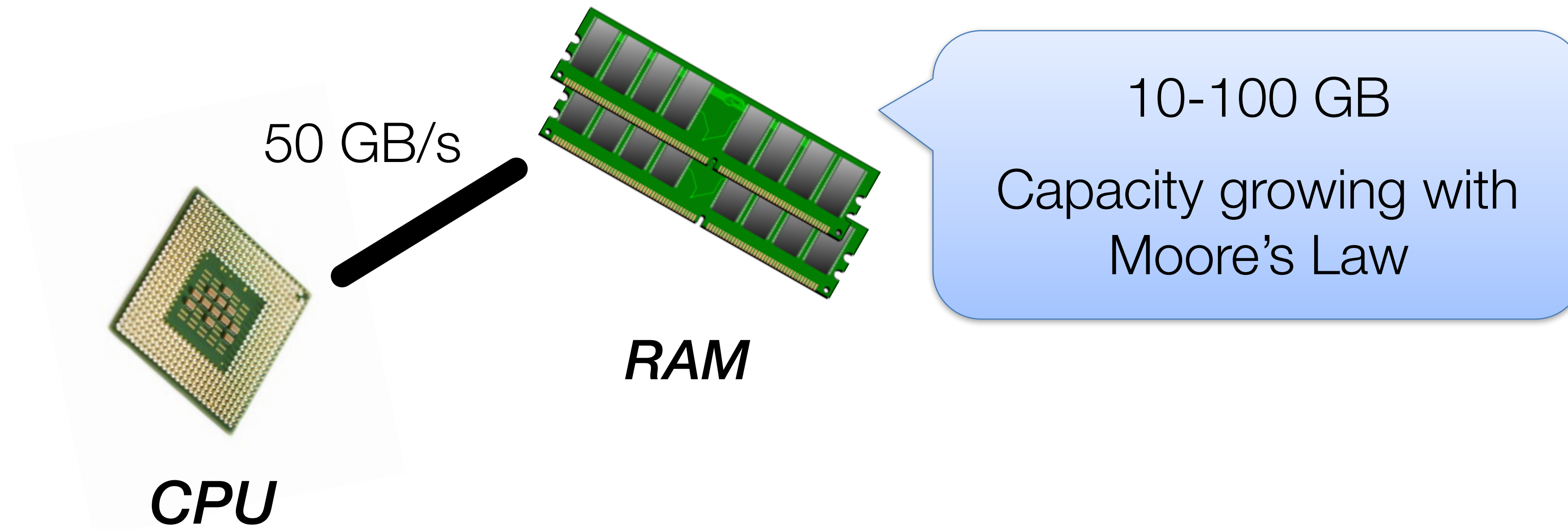
# Communication Hierarchy

# Communication Hierarchy



**CPU**

2 billion cycles/sec per core

Clock speeds not changing,
but number of cores growing
with Moore's Law

# Communication Hierarchy

50 GB/s

*RAM*

*CPU*

# Communication Hierarchy

**CPU**

50 GB/s

**RAM**

10-100 GB

Capacity growing with Moore's Law

# Communication Hierarchy

50 GB/s

*RAM*

*CPU*

100 MB/s

*Disk*

# Communication Hierarchy

50 GB/s

**RAM**

**CPU**

100 MB/s

**Disk**

1-2 TB

Capacity growing exponentially, but not speed

# Communication Hierarchy



50 GB/s

*RAM*

*CPU*

100 MB/s

*Disk* ×10

# Communication Hierarchy



50 GB/s

RAM

CPU

100 MB/s

Disk

×10

# Communication Hierarchy

50 GB/s

**RAM**

**CPU**

100 MB/s

**Disk**  ×10

**Network**

Top-of-rack switch

# Communication Hierarchy



CPU — 50 GB/s — RAM

100 MB/s — Disk ×10

Network — 10Gbps (1 GB/s) — Top-of-rack switch

10Gbps — Nodes in same rack

# Communication Hierarchy

RAM

50 GB/s

CPU

100
MB/s

Disk

×10

Network

10Gbps
(1 GB/s)

Top-of-rack
switch

10Gbps

3Gbps

Nodes in
same rack

Nodes in
other racks

# Summary

Access rates fall sharply with distance

50× gap between memory and network!



| CPU | 50 GB/s | RAM | 1 GB/s | Local disks | 1 GB/s | Rack | 0.3 GB/s | Different Racks |

*Must be mindful of this hierarchy when developing parallel algorithms!*

# Distributed ML:
# Communication Principles

# Communication Hierarchy

Access rates fall sharply with distance
- Parallelism makes computation fast
- Network makes communication slow



| CPU | 50 GB/s | RAM | 1 GB/s | Local disks | 1 GB/s | Rack | 0.3 GB/s | Different Racks |

*Must be mindful of this hierarchy when developing parallel algorithms!*

## 2nd Rule of thumb
Perform parallel and in-memory computation

Persisting in memory reduces communication
- Especially for iterative computation (gradient descent)

Scale-up (powerful multicore machine)
- No network communication
- Expensive hardware, eventually hit a wall

CPU

RAM

Disk

CPU

RAM

Disk

## 2nd Rule of thumb
Perform parallel and in-memory computation

Persisting in memory reduces communication
- Especially for iterative computation (gradient descent)

Scale-out (distributed, e.g., cloud-based)
- Need to deal with network communication
- Commodity hardware, scales to massive problems
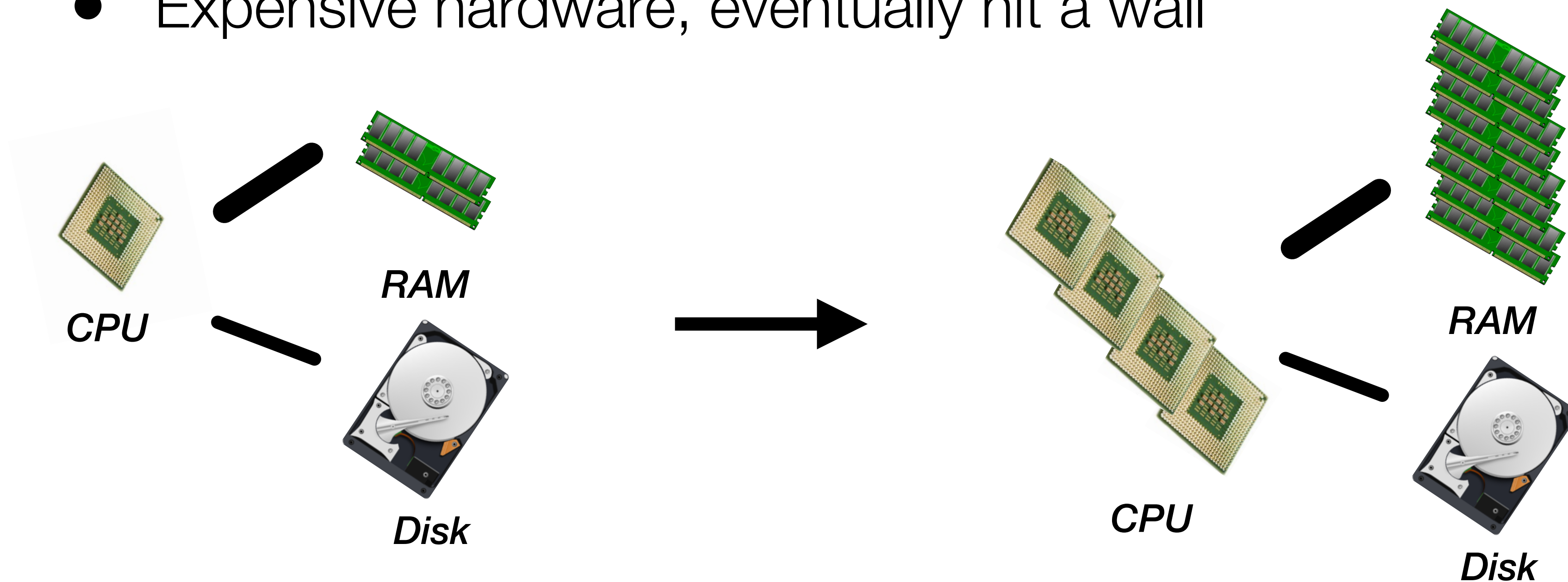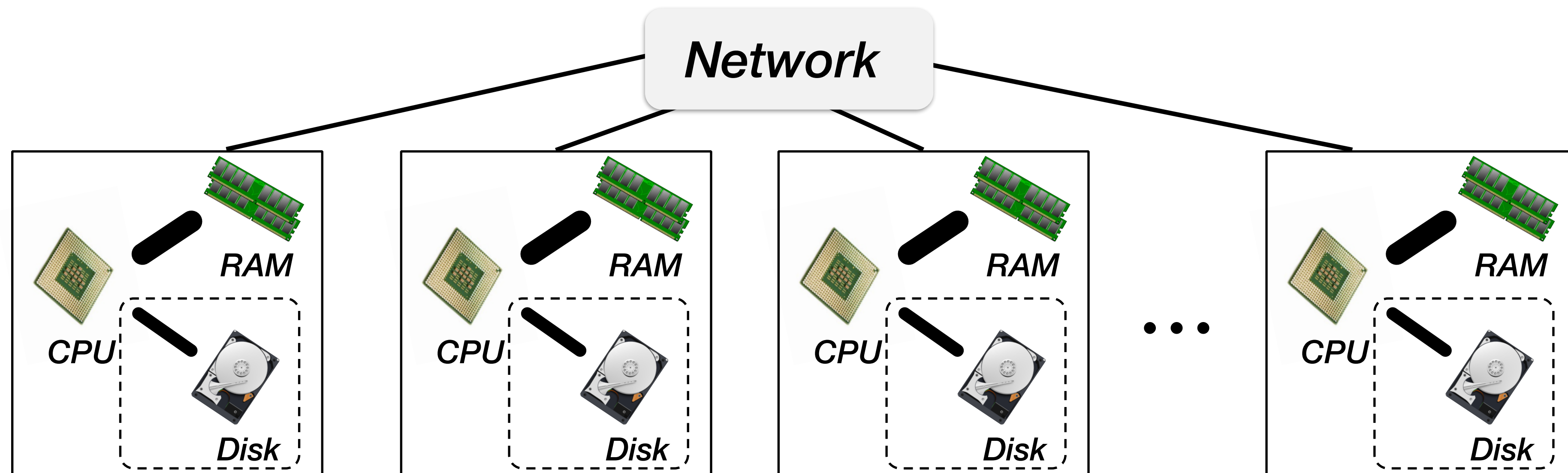
## 2nd Rule of thumb
Perform parallel and in-memory computation

Persisting in memory reduces communication
- Especially for iterative computation (gradient descent)

Scale-out (distributed, e.g., cloud-based)
- Need to deal with network communication
- Commodity hardware, scales to massive problems

```
> train.cache()   ← Persist training data across iterations
  for i in range(numIters):
      alpha_i = alpha / (n * np.sqrt(i+1))
      gradient = train.map(lambda lp: gradientSummand(w, lp)).sum()
      w -= alpha_i * gradient
```

**Q**: How should we leverage distributed computing while mitigating network communication?

First Observation: We need to store and potentially communicate Data, Model and Intermediate objects
- **A**: Keep large objects local

## 3rd Rule of thumb
Minimize Network Communication - Stay Local

**Example:** Linear regression, big $n$ and small $d$
- Solve via closed form (not iterative!)
- Communicate $O(d^2)$ intermediate data

workers:

$\mathbf{x}^{(1)}$  $\mathbf{x}^{(3)}$  $\mathbf{x}^{(2)}$

$\mathbf{x}^{(5)}$  $\mathbf{x}^{(4)}$  $\mathbf{x}^{(6)}$

map:

$\mathbf{x}^{(i)}$ $\mathbf{x}^{(i)}$    $\mathbf{x}^{(i)}$ $\mathbf{x}^{(i)}$    $\mathbf{x}^{(i)}$ $\mathbf{x}^{(i)}$

reduce:

$$\left( \sum \mathbf{x}^{(i)} \mathbf{x}^{(i)} \right)^{-1}$$

## 3rd Rule of thumb
## Minimize Network Communication - Stay Local

$\mathbf{W}_{i+1}$

workers:

$\boxed{— \mathbf{x}^{(1)}— \\ — \mathbf{x}^{(5)}—}$   $\boxed{— \mathbf{x}^{(3)}— \\ — \mathbf{x}^{(4)}—}$   $\boxed{— \mathbf{x}^{(2)}— \\ — \mathbf{x}^{(6)}—}$

map:

$(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$   $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$   $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

reduce:

$$\sum_{j=1}^{n}(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$$

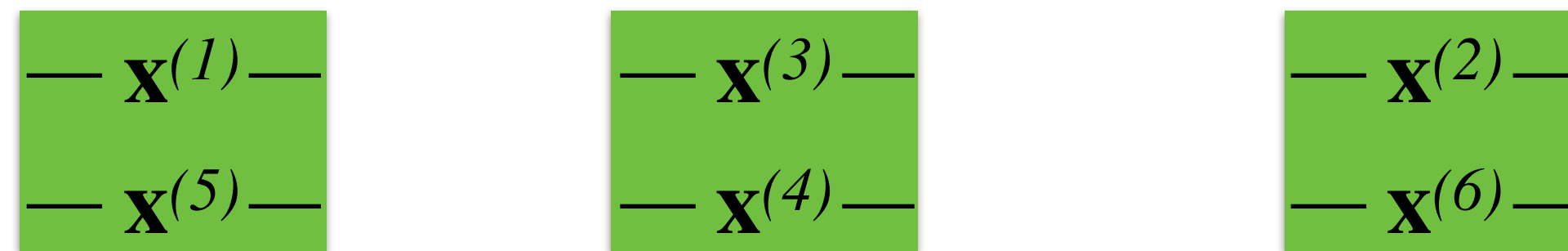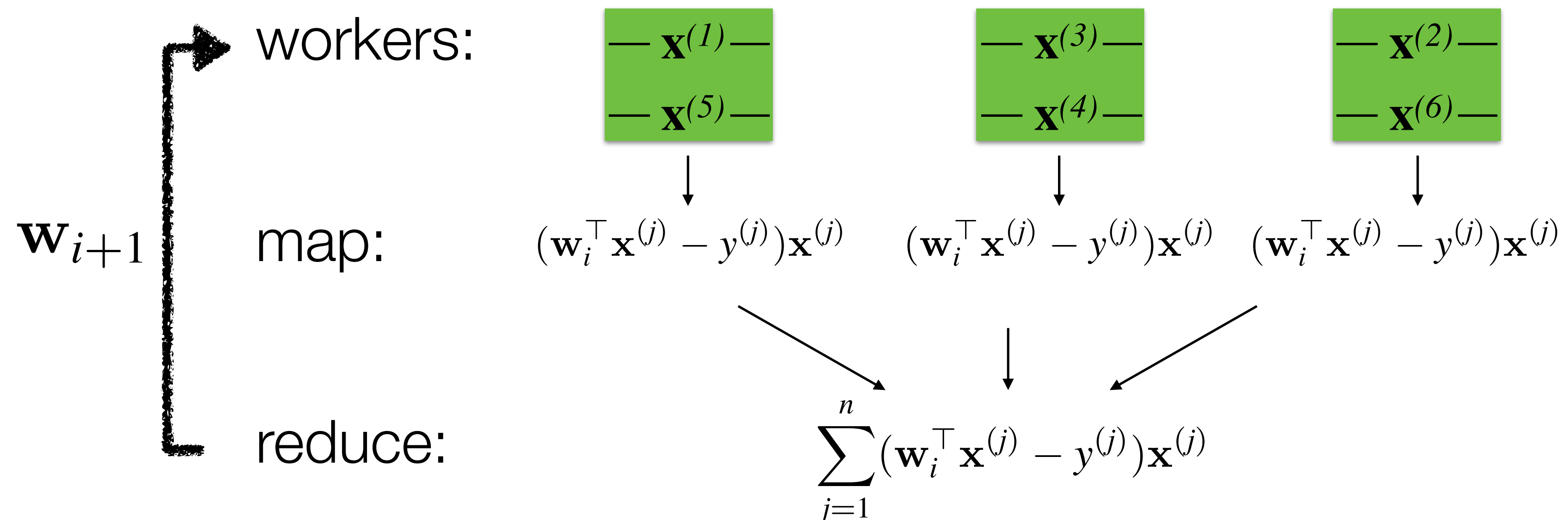**3rd Rule of thumb**
Minimize Network Communication - Stay Local

**Example:** Linear regression, big $n$ and big $d$
- Gradient descent, communicate $\mathbf{w}_i$
- O($d$) communication OK for fairly large $d$
- Compute locally on data (*Data Parallel*)

workers:

$— \mathbf{x}^{(1)} —$    $— \mathbf{x}^{(3)} —$    $— \mathbf{x}^{(2)} —$
$— \mathbf{x}^{(5)} —$    $— \mathbf{x}^{(4)} —$    $— \mathbf{x}^{(6)} —$

$\mathbf{w}_{i+1}$

map:

$(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$    $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$    $(\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$

reduce:

$$\sum_{j=1}^{n} (\mathbf{w}_i^\top \mathbf{x}^{(j)} - y^{(j)})\mathbf{x}^{(j)}$$

## 3rd Rule of thumb
### Minimize Network Communication - Stay Local

**Example:** Hyperparameter tuning for ridge regression with small $n$ and small $d$

- Data is small, so can communicate it
- 'Model' is collection of regression models corresponding to different hyperparameters
- Train each model locally (*Model Parallel*)

## 3rd Rule of thumb
## Minimize Network Communication - Stay Local

**Example:** Linear regression, big $n$ and huge $d$
- Gradient descent
- O($d$) communication slow with hundreds of millions parameters
- Distribute data and model (*Data and Model Parallel*)
- Often rely on sparsity to reduce communication

## 3rd Rule of thumb
### Minimize Network Communication

**Q**: How should we leverage distributed computing while mitigating network communication?

First Observation: We need to store and potentially communicate Data, Model and Intermediate objects
- **A**: Keep large objects local

Second Observation: ML methods are typically iterative
- **A**: Reduce # iterations

## 3rd Rule of thumb
Minimize Network Communication - Reduce Iterations

Distributed iterative algorithms must compute and communicate
- In Bulk Synchronous Parallel (BSP) systems, e.g., Apache Spark, we strictly alternate between the two

Distributed Computing Properties
- Parallelism makes computation fast
- Network makes communication slow

Idea: Design algorithms that **compute more, communicate less**
- Do more computation at each iteration
- Reduce total number of iterations

## 3rd Rule of thumb
## Minimize Network Communication - Reduce Iterations

Extreme: **Divide-and-conquer**
- Fully process each partition locally, communicate final result
- Single iteration; minimal communication
- Approximate results

```
> w = train.mapPartitions(localLinearRegression)
                  .reduce(combineLocalRegressionResults)
```

```
> for i in range(numIters):
      alpha_i = alpha / (n * np.sqrt(i+1))
      gradient = train.map(lambda lp: gradientSummand(w, lp)).sum()
      w -= alpha_i * gradient
```

## 3rd Rule of thumb
Minimize Network Communication - Reduce Iterations

Less extreme: **Mini-batch**
- Do more work locally than gradient descent before communicating
- Exact solution, but diminishing returns with larger batch sizes

```
> for i in range(fewerIters):
      update = train.mapPartitions(doSomeLocalGradientUpdates)
                      .reduce(combineLocalUpdates)
      w += update
```

```
> for i in range(numIters):
      alpha_i = alpha / (n * np.sqrt(i+1))
      gradient = train.map(lambda lp: gradientSummand(w, lp)).sum()
      w -= alpha_i * gradient
```

## 3rd Rule of thumb
Minimize Network Communication - Reduce Iterations

**Throughput**: How many bytes per second can be read

**Latency**: Cost to send message (independent of size)

| Latency | |
|---|---|
| **Memory** | 1e-4 ms |
| **Hard Disk** | 10 ms |
| **Network (same datacenter)** | .25 ms |
| **Network (US to Europe)** | >5 ms |

We can amortize latency!
- Send larger messages
- *Batch* their communication
- E.g., Train multiple models together

**1st Rule of thumb**
Computation and storage should be linear (in $n, d$)

**2nd Rule of thumb**
Perform parallel and in-memory computation

**3rd Rule of thumb**
Minimize Network Communication

# Lab Preview
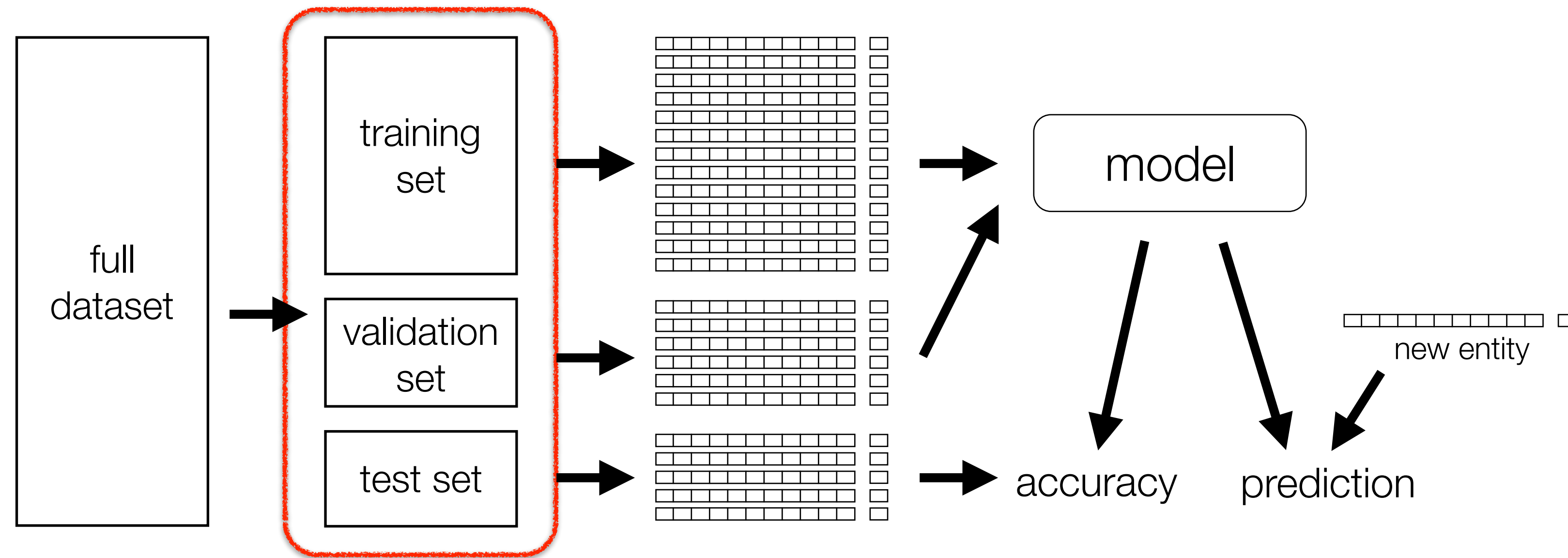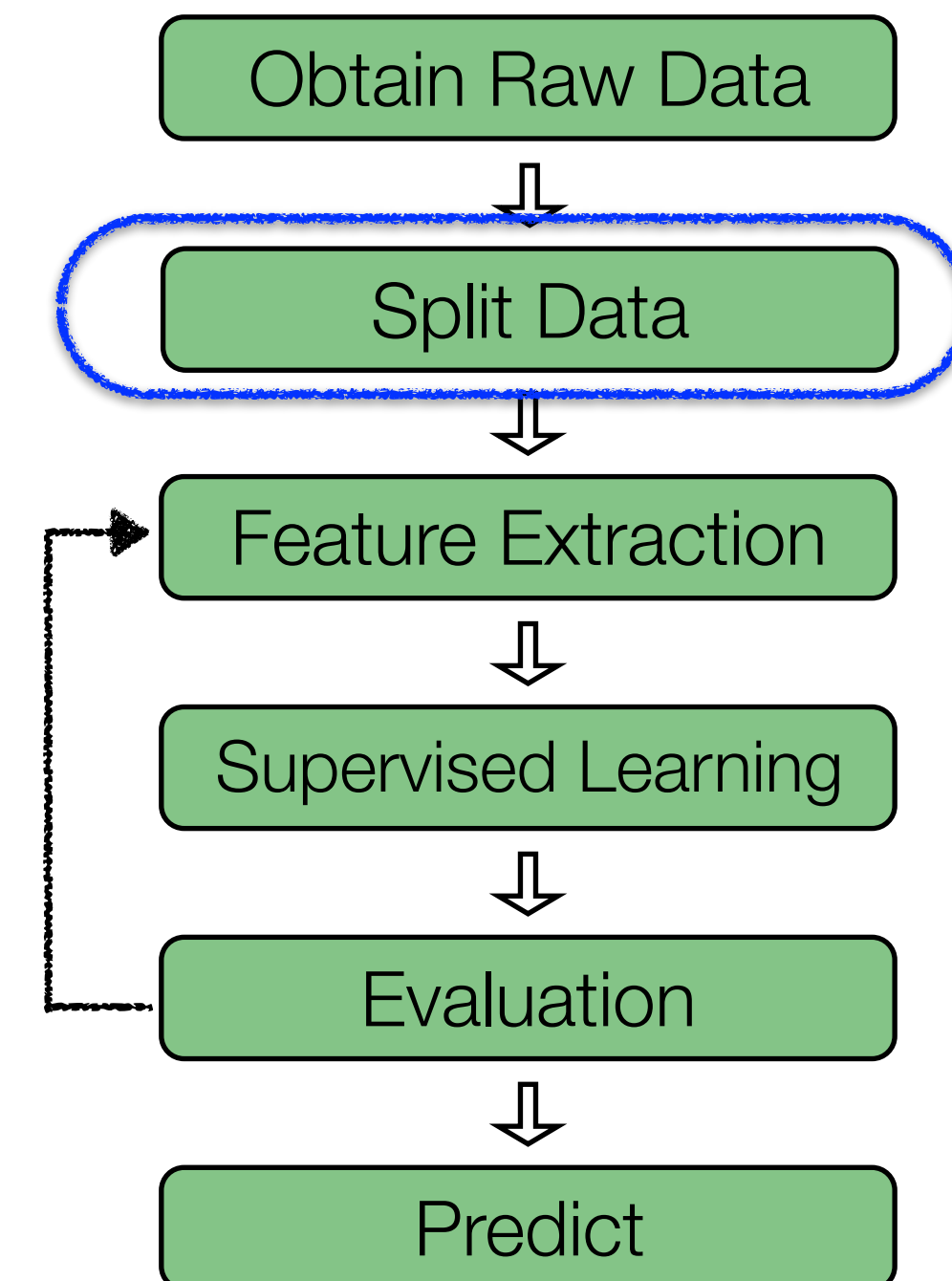
**Goal**: Predict song's release year from audio features

Obtain Raw Data
⇩
Split Data
⇩
Feature Extraction
⇩
Supervised Learning
⇩
Evaluation
⇩
Predict

**Raw Data**: Millionsong Dataset from UCI ML Repository
- Explore features
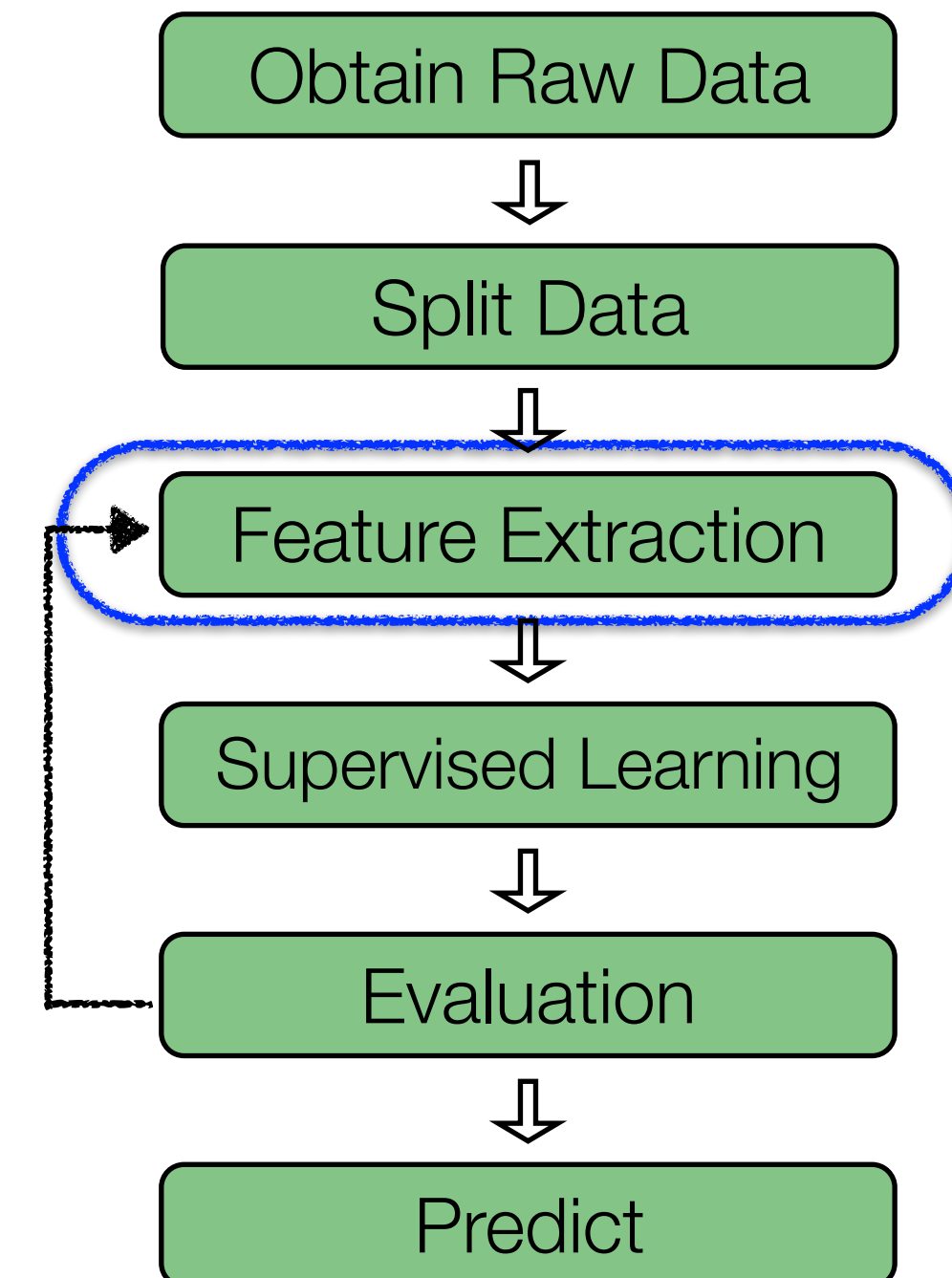- Shift labels so that they start at 0 (for interpretability)
- Visualize data

Obtain Raw Data

Split Data

Feature Extraction

Supervised Learning

Evaluation

Predict

**Split Data**: Create training, validation, and test sets

**Feature Extraction**:
- Initially use raw features
- Subsequently compare with quadratic features

Obtain Raw Data
⇩
Split Data
⇩
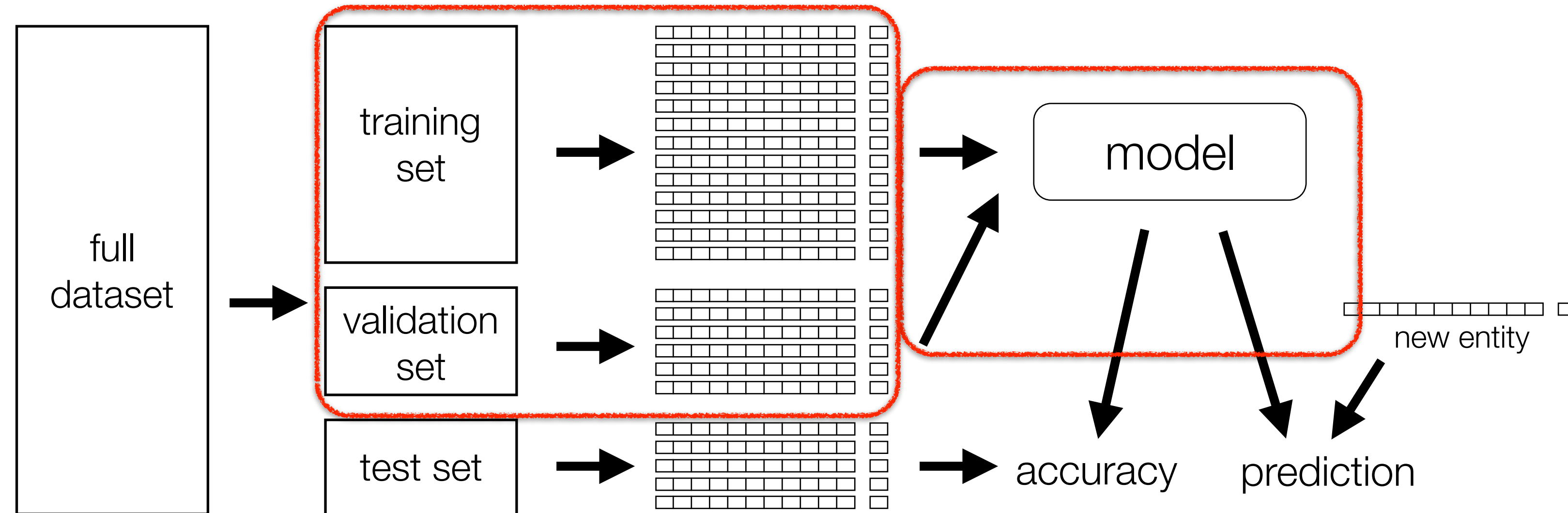Feature Extraction
⇩
Supervised Learning
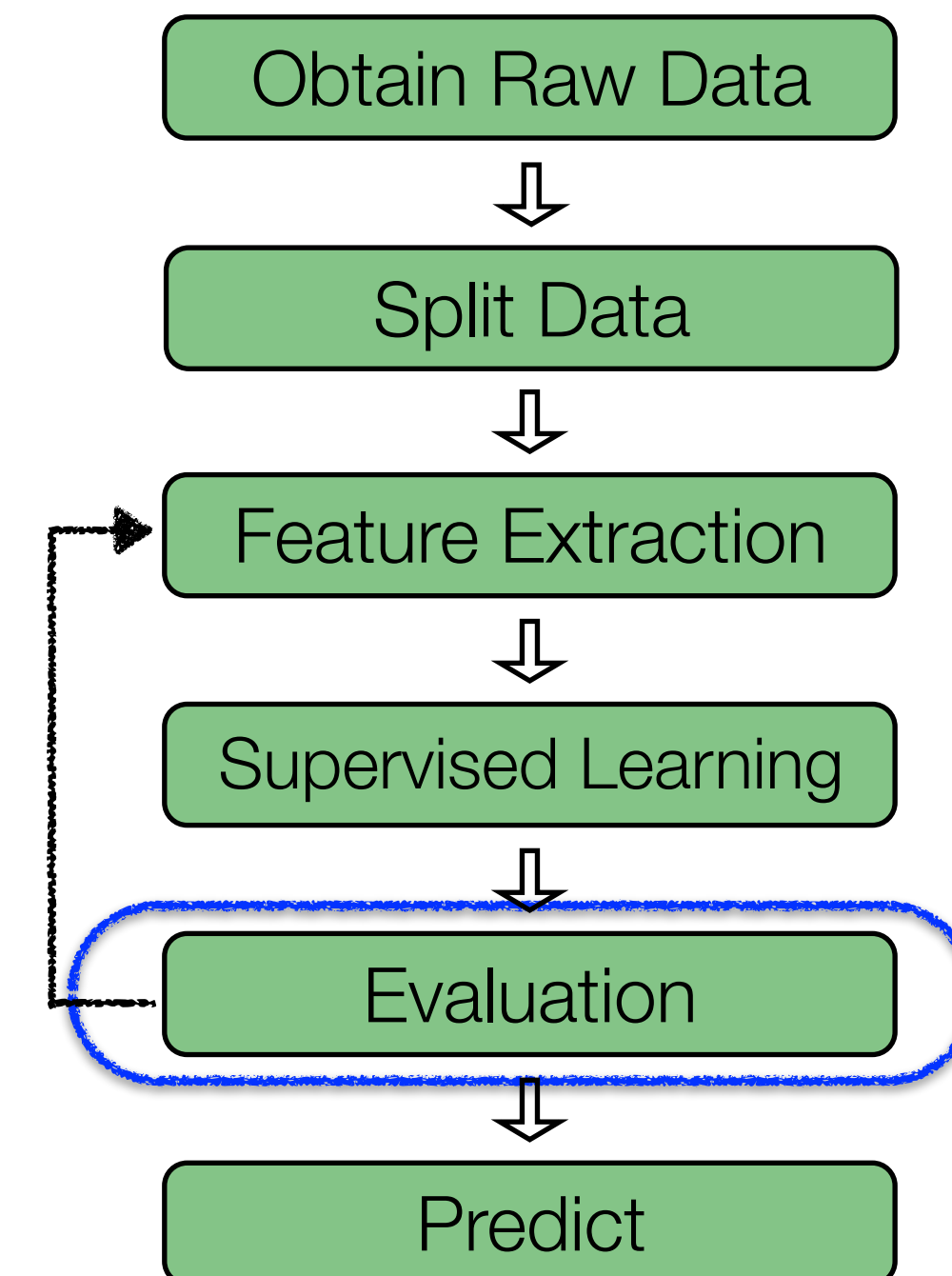⇩
Evaluation
⇩
Predict

**Supervised Learning**: Least Squares Regression
- First implement gradient descent from scratch
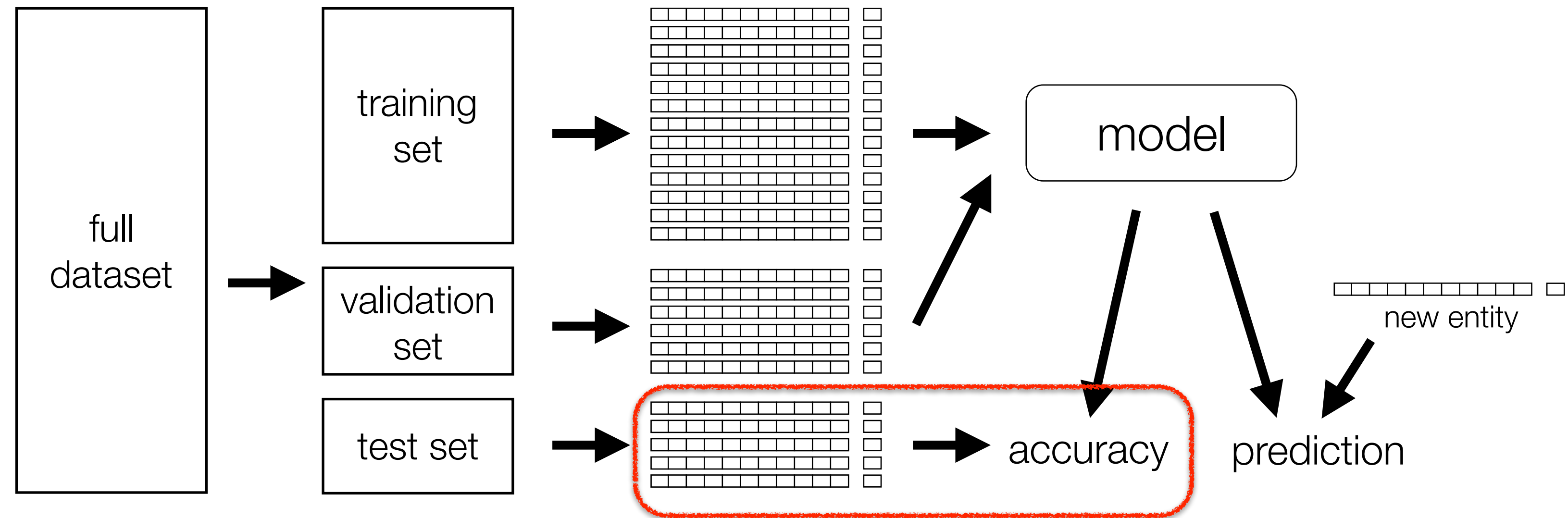- Then use MLlib implementation
- Visualize performance by iteration

Obtain Raw Data
⇩
Split Data
⇩
Feature Extraction
⇩
Supervised Learning
⇩
Evaluation
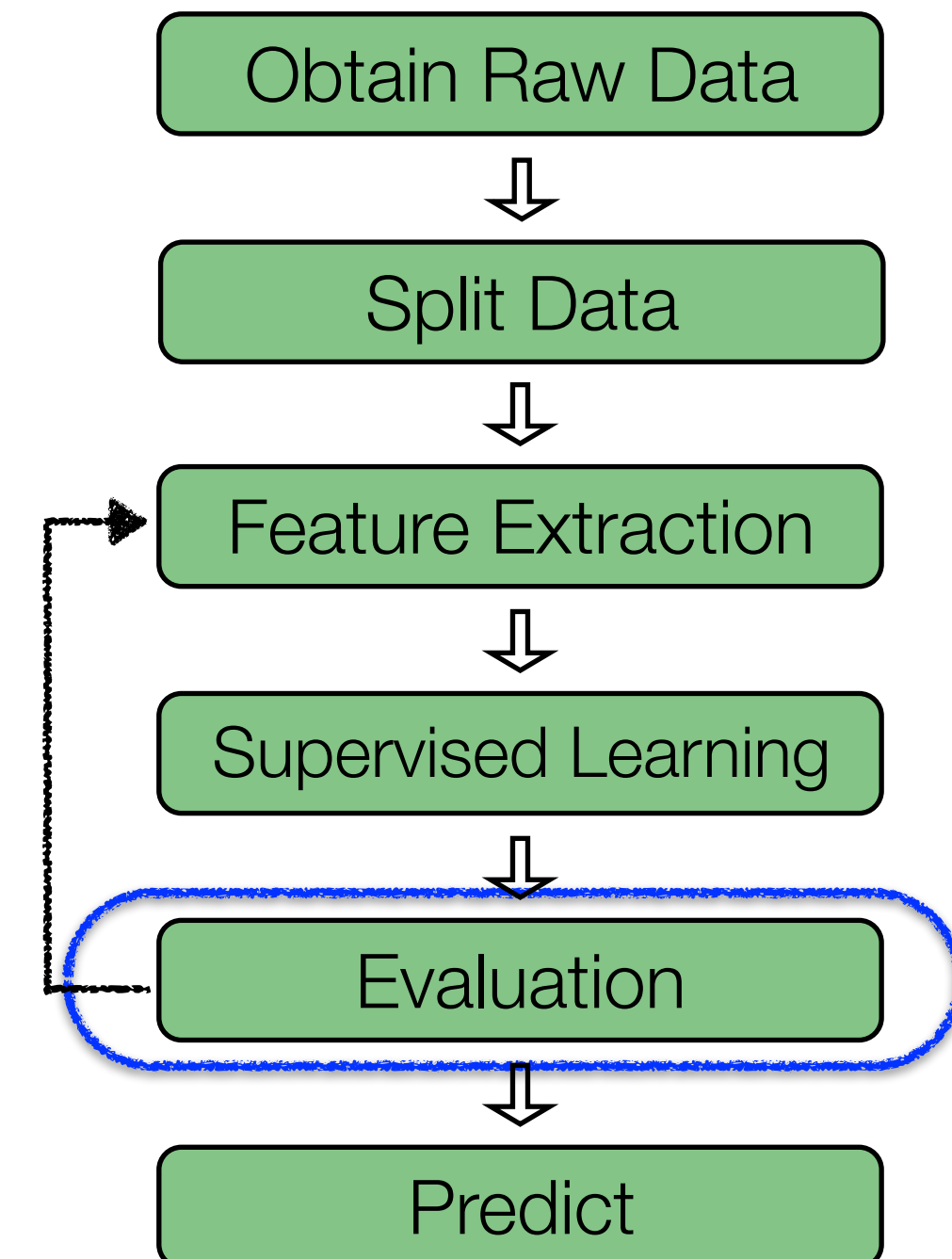⇩
Predict

**Evaluation (Part 1)**: Hyperparameter tuning

- Use grid search to find good values for regularization and step size hyperparameters
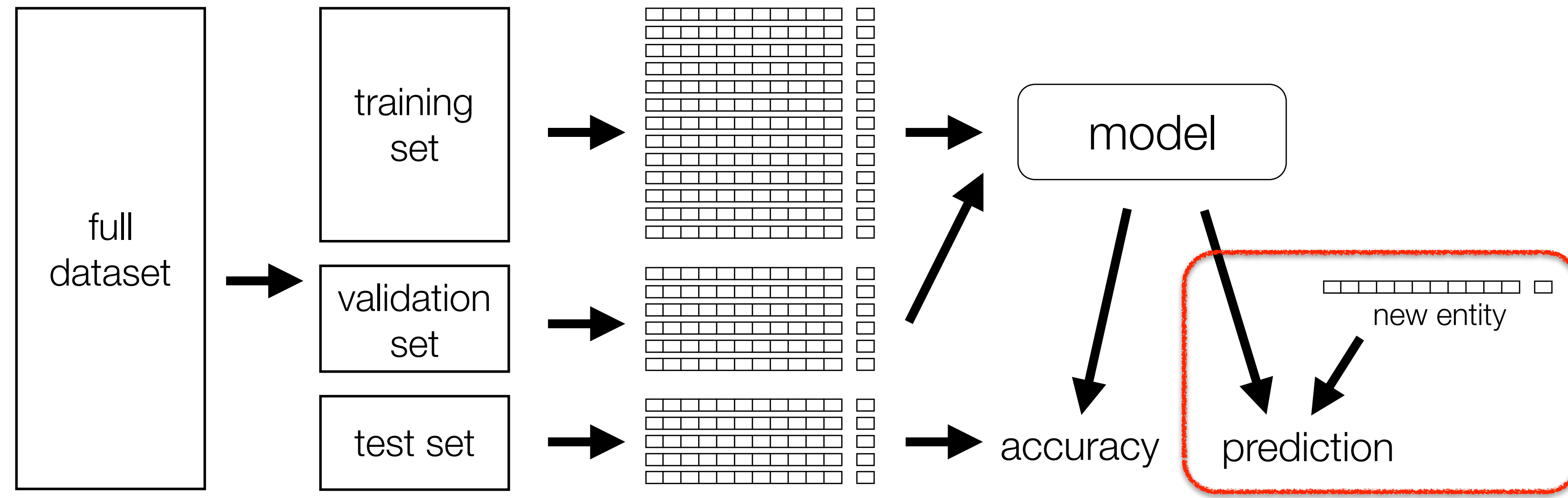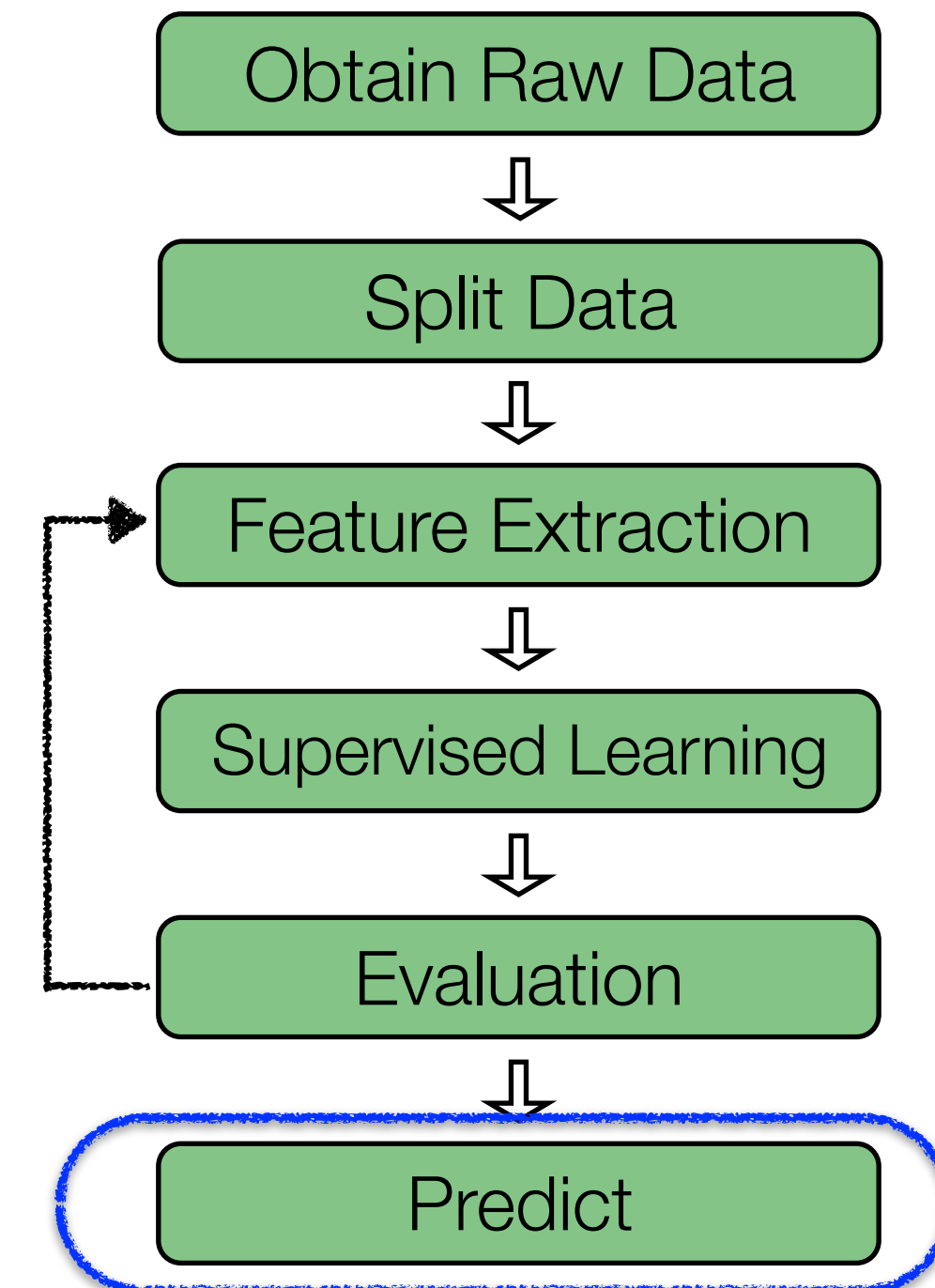- Evaluate using RMSE
- Visualize grid search

**Evaluation (Part 2)**: Evaluate final model
- Evaluate using RMSE
- Compare to baseline model that returns average song year in training data

**Predict**: Final model could be used to predict song year for new songs (we won't do this though)

Obtain Raw Data
⇩
Split Data
⇩
Feature Extraction
⇩
Supervised Learning
⇩
Evaluation
⇩
Predict

# MLlib and Pipelines

# Spark's Machine Learning Library (MLlib)

- Consists of common learning algorithms and utilities
    - Classification
    - Regression
    - Clustering
    - Collaborative filtering
    - Dimensionality reduction

- Two packages:
    - spark.mllib
    - spark.ml

# ML: Transformer

- A *Transformer* is a class which can transform one DataFrame into another DataFrame
- A Transformer implements `transform()`

- Examples
  - HashingTF
  - LogisticRegressionModel
  - Binarizer

# ML: Estimator

- An *Estimator* is a class which can take a DataFrame and produce a Transformer
- An Estimator implements `fit()`

- Examples
  - LogisticRegression
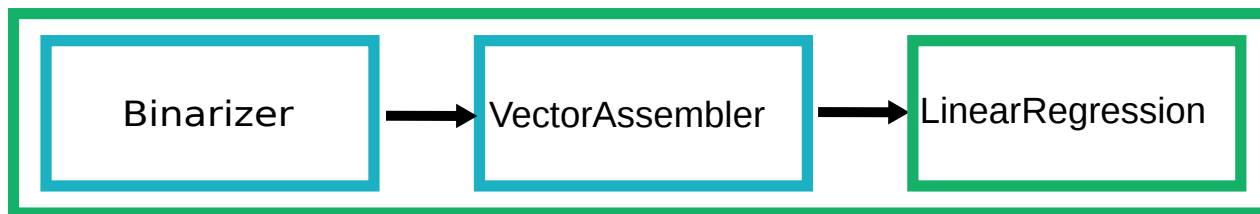  - StandardScaler
  - Pipeline

# ML: Pipelines

A *Pipeline* is an estimator that contains stages representing a resusable workflow. Pipeline stages can be either estimators or transformers.
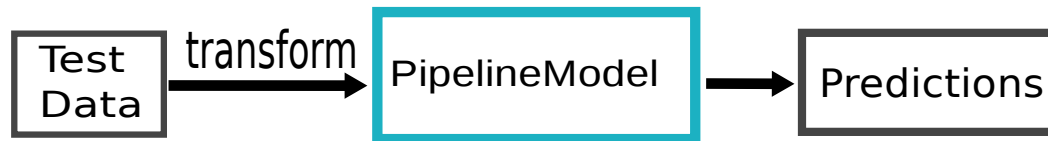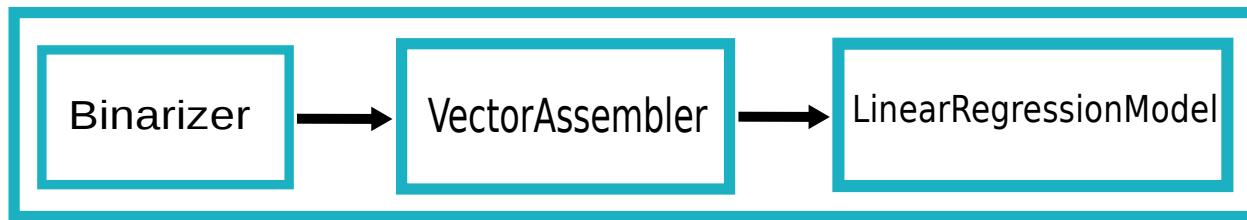
## Pipeline

```
Transformer  →  Transformer  →  Estimator
```

## Pipeline

```
Binarizer  →  VectorAssembler  →  LinearRegression
```

# ML: PipelineModel

Train Data → **fit** → Pipeline → PipelineModel

## PipelineModel

Binarizer → VectorAssembler → LinearRegressionModel

Test Data → **transform** → PipelineModel → Predictions

# ML: Standard Scaler Pipeline

## Pipeline

StandardScaler → LinearRegression

Train Data → **fit** → Pipeline → PipelineModel

## PipelineModel

StandardScalerModel → LinearRegressionModel

Test Data → **transform** → PipelineModel → Predictions