

# Data mining & Machine Learning

CS 373

Purdue University

Dan Goldwasser

[dgoldwas@purdue.edu](mailto:dgoldwas@purdue.edu)

# Today's Lecture

## Learning as Optimization

- *Our discussion of learning so far focused on specific algorithms*
- *Today, we'll discuss a new way of thinking about learning – optimizing an objective function, consisting of both a data-fitting term (performance on training data) and generalization term (e.g., margin).*
- *We will start with finding the best margin classifier*
- *Then, we will extend this framework to include more general terms for controlling overfitting*

# Perceptron In Practice

- The perceptron algorithm is actually a pretty good practical algorithm.
  - It's also very simple to implement and modify.

*In the next slides we will introduce some design choices and understand their implication on how Perceptron works.*

# (Realistic) Perceptron

---

- We learn  $\mathbf{f}:\mathbf{x}\rightarrow \{-1,+1\}$  represented as  $\mathbf{f} = \text{sgn}\{\mathbf{w}\bullet\mathbf{x}\}$
- Where  $\mathbf{x}=\{0,1\}^n$
- Given Labeled examples:  $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots (\mathbf{x}_m, \mathbf{y}_m)\}$

1. Initialize  $\mathbf{w}=0 \in \mathbf{R}^n$

**2. For Iter = 0,...,T**

3. Iterate over all the examples

a. **Predict** the label of instance  $\mathbf{x}$  to be  $\mathbf{y}' = \text{sgn}\{\mathbf{w}\bullet\mathbf{x}\}$

b. If  $\mathbf{y}' \neq \mathbf{y}$ , **update** the weight vector:

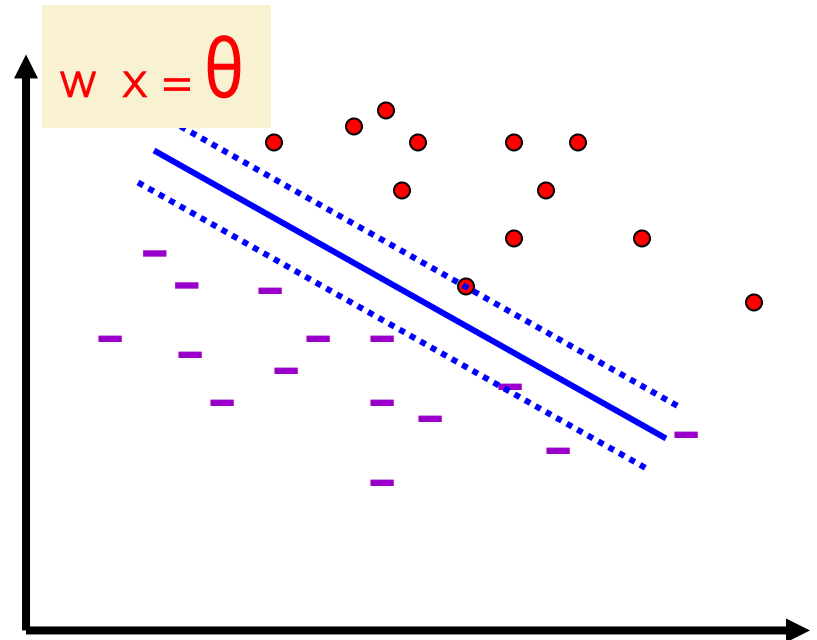
$$\mathbf{w} = \mathbf{w} + r \mathbf{y} \mathbf{x} \quad (r - \text{a constant, learning rate})$$

Otherwise, if  $\mathbf{y}' = \mathbf{y}$ , leave weights unchanged.

# Regularization: Perceptron with Margin

- *Weights with better margin generalize better*
  - Perceptron finds *any* separating hyperplane
- **Thick Separator (aka as Perceptron with Margin)**

- **Predict positive**  
bias term  $w x - \theta > \gamma$  margin
- **Predict negative**  
 $w x - \theta < -\gamma$
- **Mistake:**  
 $\gamma > (w x - \theta) > -\gamma$



# Regularization: Perceptron with Margin

---

- **Perceptron margin** : hyperparameter that has to be tuned using the validation set (try different values)
  - *In the future: the data will decide the margin*
- The impact of margin regularization in perceptron becomes smaller as  $w$  grows

$$w x - \theta > \gamma$$

*what happens if we multiply the weights by 2? by 2000?*

- ***Can we control the growth of  $w$ ?***
- In practice, **very effective**

# Perceptron: Robust Variation

---

- *The perceptron algorithm counts later points more than earlier points*

1: (0,1,...,1,0,1)

2: (0,1,...,1,0,1)

...

100: (0,1,...,1,0,1)

...

10000: (0,1,...,1,0,1)

Makes some mistakes, update..

After 100 examples, learner stops making mistakes

**We keep going...**

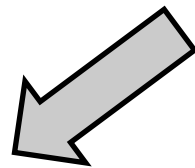
**BUT** then at the 10,000 example the learner makes a mistake!

***Is this a problem?***

# Voted Perceptron

---

- *Training:*
  - *Learner remembers how long each hypothesis survived (no mistakes on  $w$ )*
- Test:
  - Weighted vote of all participating hypotheses



**BIG Idea in ML:**  
**reduce variance using classifier ensemble**

$$\hat{y} = \text{sign} \left( \sum_{k=1}^K c^{(k)} \text{sign} \left( w^{(k)} \cdot \hat{x} + b^{(k)} \right) \right)$$

- Heavy: (1) Computational effort (2) **Storage**



# Averaged Perceptron

---

- **Training:** *Maintain a running weighted average of survived hypotheses*
- **Test:** *Predict according to the averaged weight vector*

**Voted Perceptron**  $\rightarrow \hat{y} = \text{sign} \left( \sum_{k=1}^K c^{(k)} \text{sign} \left( \mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right)$

$\hat{y} = \text{sign} \left( \sum_{k=1}^K c^{(k)} \left( \mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right)$   $\leftarrow$  **Averaged Perceptron**

- *An efficient approximation of voted perceptron*
- *Almost always better than regular perceptron!*

# Understanding Weighted Perceptron

```
D = {xi, yi}i=1,...,n
w0 = (0, ..., 0)
t=0
repeat T times
  for (xi, yi) in D
    y' = sign(wtx) prediction based on
    the current model
    if (y' != y)
      wt+1 = wt + x y r Update Rule
    else
      wt+1 = wt
  t = t + 1
return (w1 + ... + wnT) / (nT) return the
averaged result
```

**key idea:** represent the averaged function **explicitly** by keeping in memory all the functions visited by the algorithm.

# Efficient Weighted Perceptron

---

- We can express the learned function as a sum of all its updates:

$$w_0 = (0, \dots, 0)$$

$$w_1 = w_0 + \Delta_1 = \Delta_1$$

$$w_2 = w_1 + \Delta_2 = \Delta_1 + \Delta_2$$

$$w_3 = w_2 + \Delta_3 = \Delta_1 + \Delta_2 + \Delta_3$$

...

- This allows us to express the average of the classifiers in terms of updates:

$$(w_1 + w_2 + w_3)/3$$

$$= (\Delta_1)/3 + (\Delta_1 + \Delta_2)/3 + (\Delta_1 + \Delta_2 + \Delta_3)/3$$

$$= (3/3)\Delta_1 + (2/3)\Delta_2 + (1/3)\Delta_3$$

# Efficient Weighted Perceptron

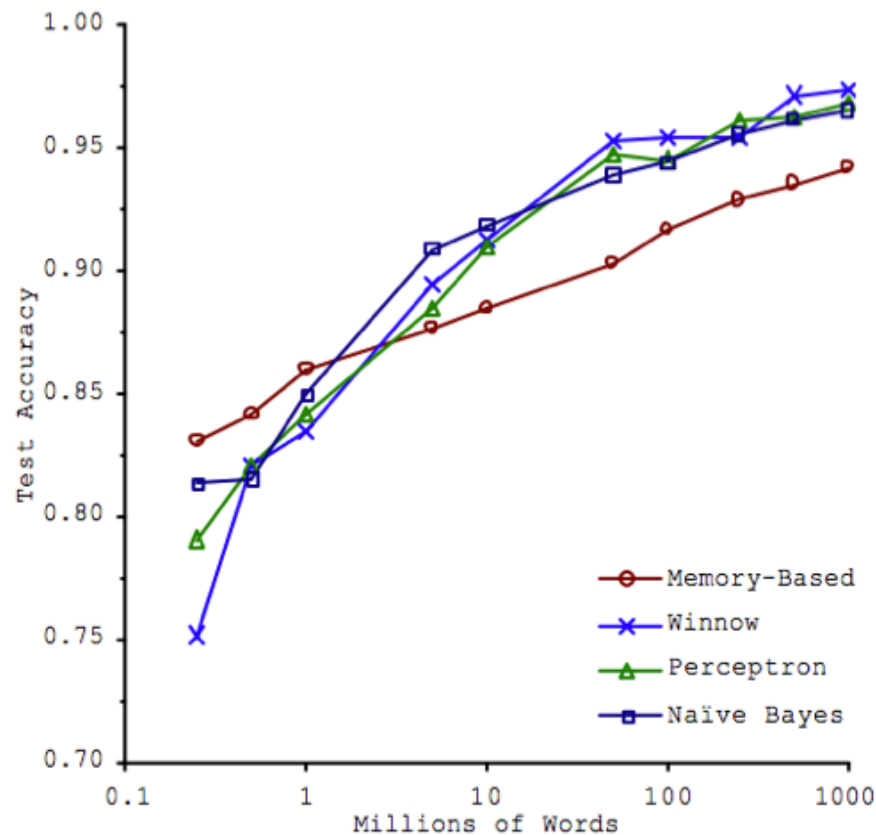
```
D = {xi, yi}i=1,...,n
w = (0,...,0)      current function weights
a = (0,...,0)      counter of all the updates seen so far
step = nT
repeat T times
  for (xi, yi) in D
    y' = sign(wx)    prediction based
                     on current model
    if (y' != y)
      w = w + xiyi  update Rule
      a = a + (step/nT)(xiyi)  update the weight
                                counter
    step = step - 1
return a  return the averaged result
```

**key idea:** represent the averaged function as averaged updates, instead of explicitly keeping in memory all the functions visited by the algorithm.

# Practical Example

Task: **context sensitive spelling**

*"I didn't know {weather,whether} to laugh or cry"*



**Source:** Scaling to very very large corpora for natural language disambiguation Michele Banko, Eric Brill. MSR, 2001.

# **Learning as Optimization**

# Learning as Optimization

---

- The perceptron algorithm is a mistake driven algorithm.
  - The model is updated in response to an error
- In that sense, it is minimizing the number of mistake on the training data.
  - Does it provide a globally optimal result?
- In the next couple of lectures we will make the optimization process explicit, and essentially equate learning with an optimization objective.
  - **Key advantage:** "one size fits all": if you can optimize a function, the difference between the algorithms is minimal!

# Searching over models/patterns

---

- Consider a **space** of possible models  $M=\{M_1, M_2, ..., M_k\}$  with parameters  $\theta$
- Search could be over model structures or parameters, e.g.:
  - **Parameters**: In a linear regression model, what are regression coefficients ( $\beta$ ) that minimize squared loss on the training data?
  - **Model structure**: In a decision trees, what is the tree structure that minimizes 0/1 loss on the training data?



# Optimization

---

- **Non-smooth** functions:

- If the function is *discrete*, then traditional optimization methods that rely on smoothness are not applicable. Instead we need to use **combinatorial optimization**
- **Example:** *Choosing what features (structure) to add to a decision tree*
- *You need to define the search space and a search procedure.*
  - **Formally:** *State Space (including Initial State, Final States), Transition Function (i.e., actions), Scoring Function.*

# Search algorithms for discrete spaces

---

- **Conduct the search by:**

- Considering a particular state (*model*)
- Testing to see if it is the goal state (*model with maximum score*)
- And if not, expand the current state to generate successor states by applying all possible actions  
(*determine alternative models to consider next*)

**Search strategies differ in their choice of how to expand states**

# Optimization

---

- **Smooth** functions:

- If a function is ***smooth***, it is differentiable and the derivatives are continuous, then we can use gradient-based optimization
  - If function is ***convex***, we can solve the minimization problem in closed form:  $\min_{\theta} S(\theta)$  using **convex optimization**
  - If function is smooth but **non-convex**, we can use iterative search over the surface of  $S$  to find a local minimum (e.g., hill-climbing)

# Convex optimization problems

---

$$\begin{array}{ll}\text{minimize} & f(x) \\ \text{subject to} & x \in C\end{array}$$

- Where  $f$  is a **convex function** (*score function*)  
 $C$  is a **convex set** (*constraints on model parameters or structure*)  
 $x$  is the optimization variable (*includes data and parameters*)
- For convex optimization problems, **all locally optimal points are globally optimal**
- **Example algorithms:** Quadratic programming (SVMs), least squares estimation, *maximum likelihood estimation*

# Loss functions

---

- To formalize performance let's define a **loss function**:  $loss(y, \hat{y})$ 
  - Where  $\hat{y}$  is the gold label
- The loss function measures the error on a single instance
  - Specific definition depends on the learning task

## *Regression*

$$y = w^T x$$

$$loss(y, \hat{y}) = (y - \hat{y})^2$$

## *Binary classification*

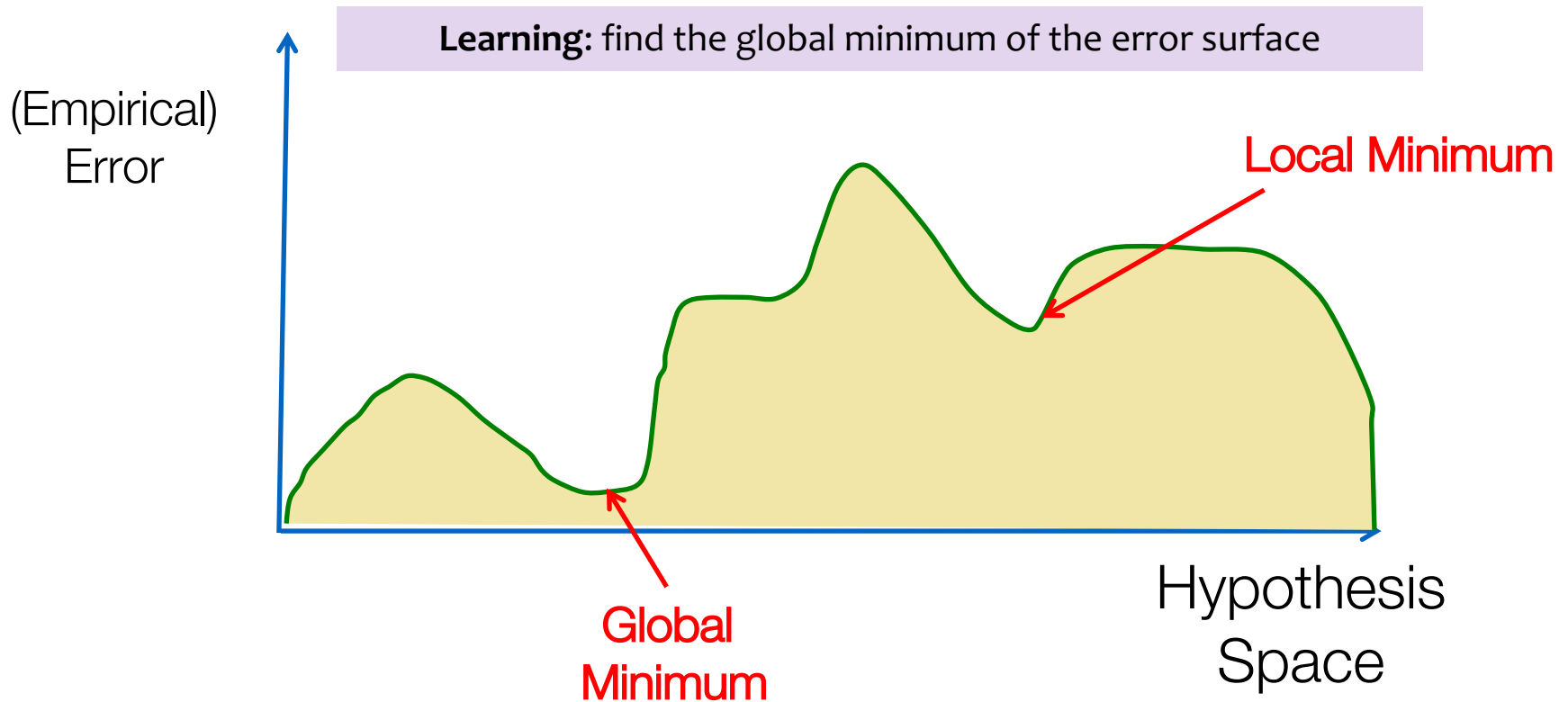
$$y = \text{sign}(w^T x)$$

$$loss(y, \hat{y}) = \begin{cases} 0 & y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$$

# Error Surface

---

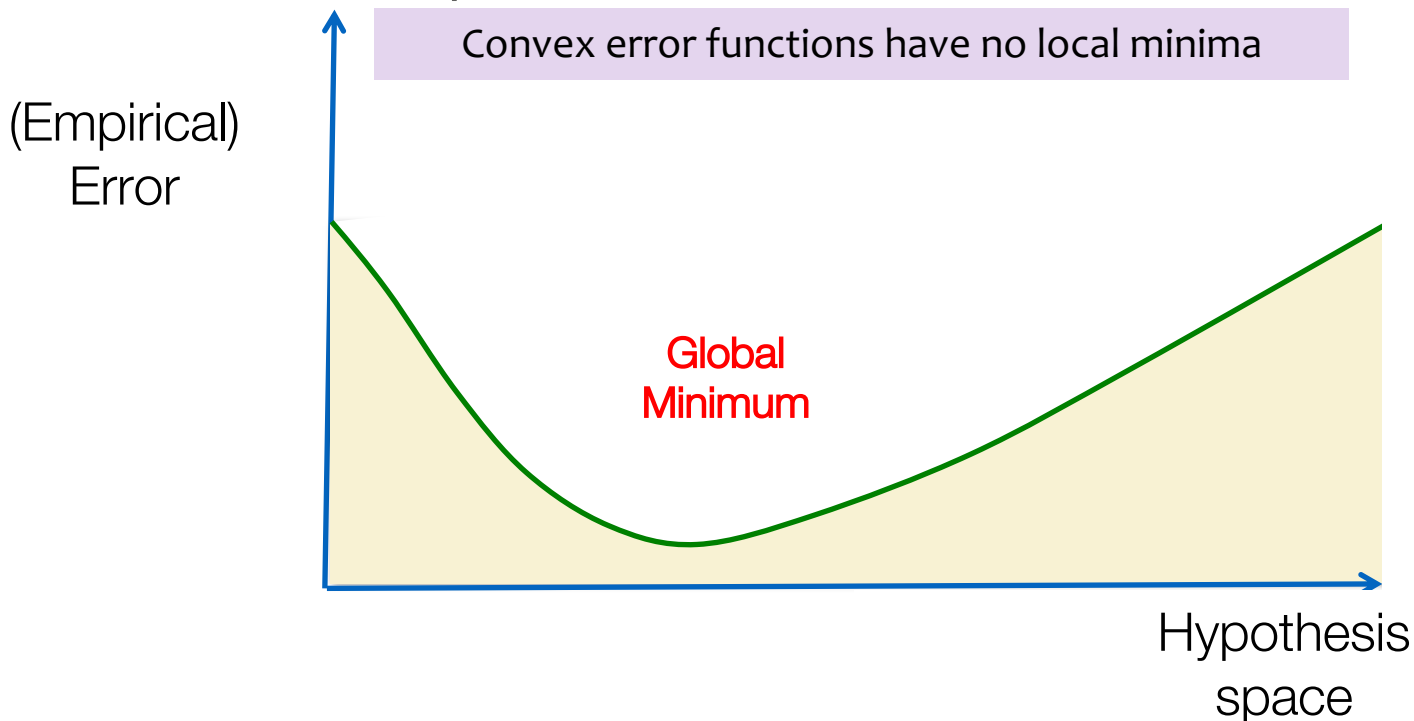
- **Linear classifiers:** hypothesis space parameterized by  $w$
- *Error/Loss/Risk* are all functions of  $w$



# Convex Error Surfaces

---

- **Convex functions** have a single minimum point
  - Local minimum = global minimum
  - *Easier to optimize*



Convex optimization

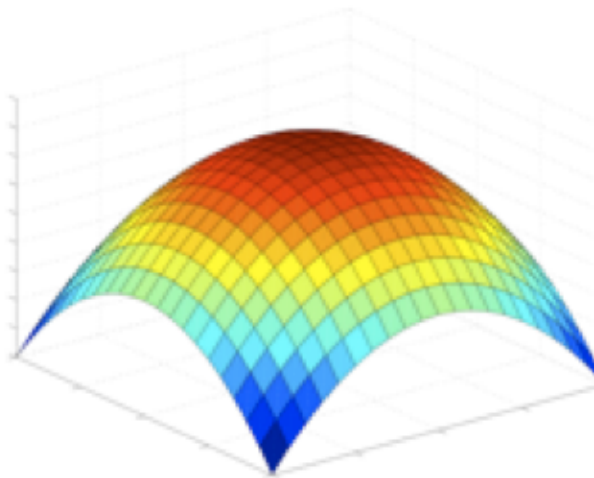


# Concave vs convex

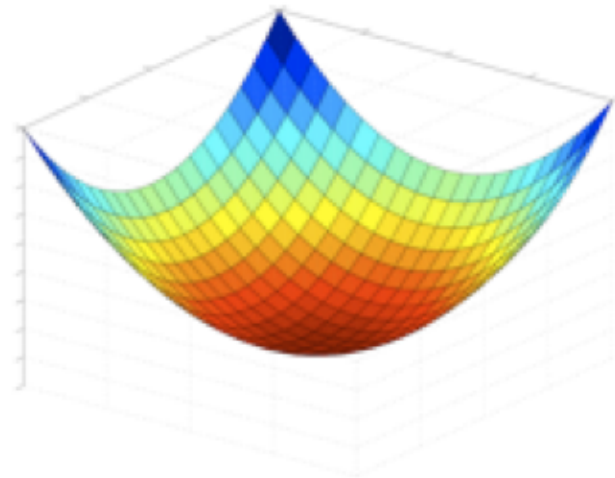
---

- Maximizing a concave function is equivalent to minimizing a convex function

**concave**



**convex**



# Score function: Likelihood

---

- Let  $D = \{x(1), \dots, x(n)\}$
- Assume the data  $D$  are independently sampled from the same distribution:  
$$p(X|\theta)$$
- The likelihood function represents the probability of the data as a function of the model parameters:

$$\begin{aligned} L(\theta|D) &= L(\theta|x(1), \dots, x(n)) \\ &= p(x(1), \dots, x(n)|\theta) \\ &= \prod_{i=1}^n p(x(i)|\theta) \end{aligned}$$

**If instances are independent,  
likelihood is product of probs**

# Maximum likelihood estimation

---

- Most widely used method of parameter estimation
- “Learn” the best parameters by finding the values of  $\theta$  that maximizes likelihood:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(\theta)$$

- Often easier to work with loglikelihood:

$$\begin{aligned} l(\theta|D) &= \log L(\theta|D) \\ &= \log \prod_{i=1}^n p(x(i)|\theta) \\ &= \sum_{i=1}^n \log p(x(i)|\theta) \end{aligned}$$

# Maximum likelihood estimation

---

- Define likelihood, take derivative, set to 0, and solve
- **Example**
  - Toss a weighted coin 100 times, observe 30 heads
  - What is the MLE estimate for the  $p$  parameter of the Binomial distribution that generated the data?
  - First define likelihood

$$\begin{aligned} L(p|H=30, n=100) &= P(H=30|n=100, p) \\ &= \binom{100}{30} p^{30} (1-p)^{70} \end{aligned}$$

**Take derivative, set to 0, solve**

$$\begin{aligned}0 &= \frac{d}{dp} \left( \binom{100}{30} p^{30} (1-p)^{70} \right) \\&\propto 30p^{29}(1-p)^{70} - 70p^{30}(1-p)^{69} \\&= p^{29}(1-p)^{69} [30(1-p) - 70p] \\&= p^{29}(1-p)^{69} [30 - 100p] \\p &= 0, 1, \frac{30}{100}\end{aligned}$$

# Gradient descent

---

- For some convex functions, we may be able to take the derivative, but it may be difficult to directly solve for parameter values
- **Solution:**
  - Start at some value of the parameters
  - Take derivative and use it to move the parameters in the direction of the solution
  - Repeat

## Gradient Descent Rule:

$$\underline{\mathbf{w}}_{\text{new}} = \underline{\mathbf{w}}_{\text{old}} - \eta \Delta(\underline{\mathbf{w}})$$

where

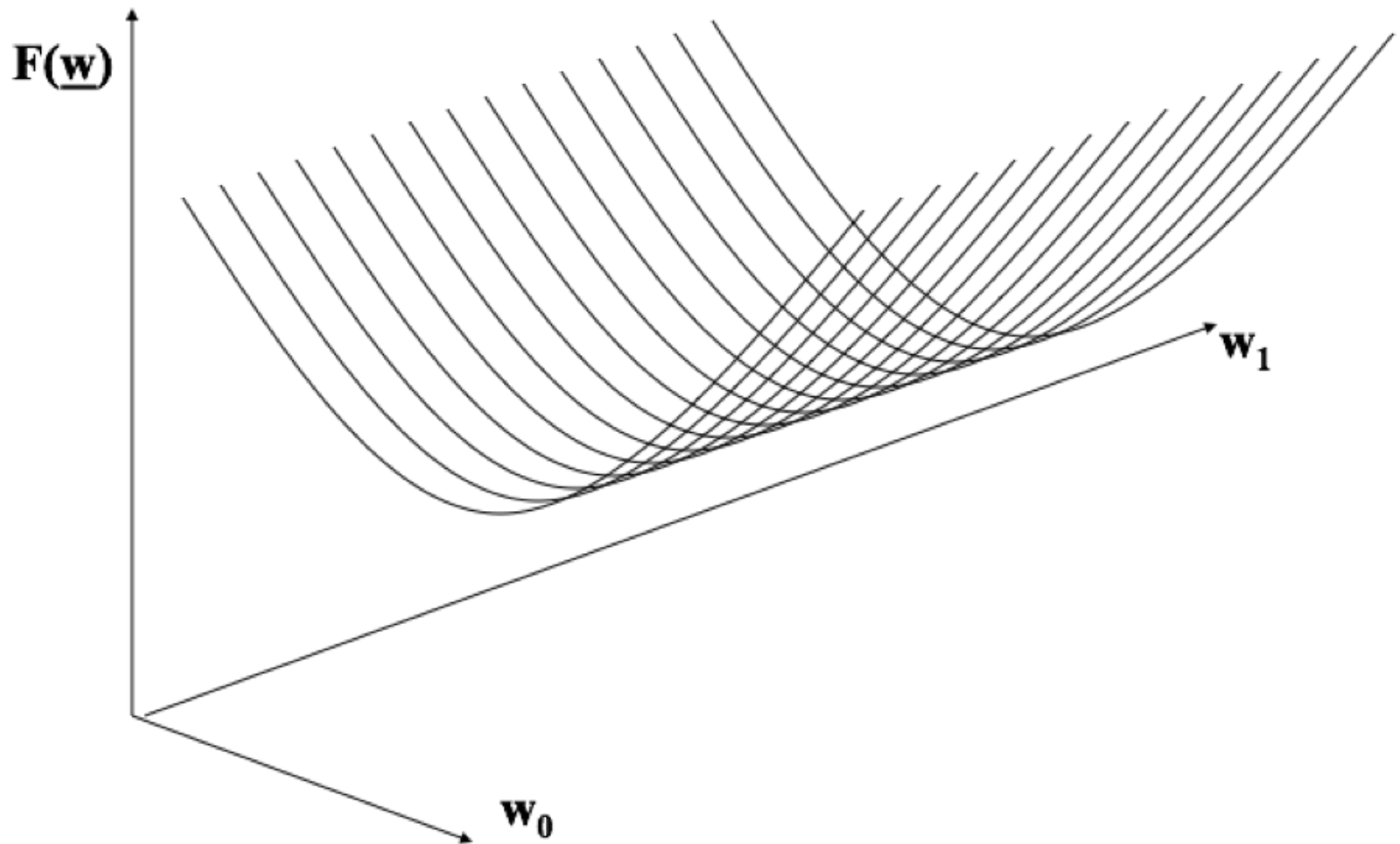
$\Delta(\underline{\mathbf{w}})$  is the gradient and

$\eta$  is the learning rate (small, positive)

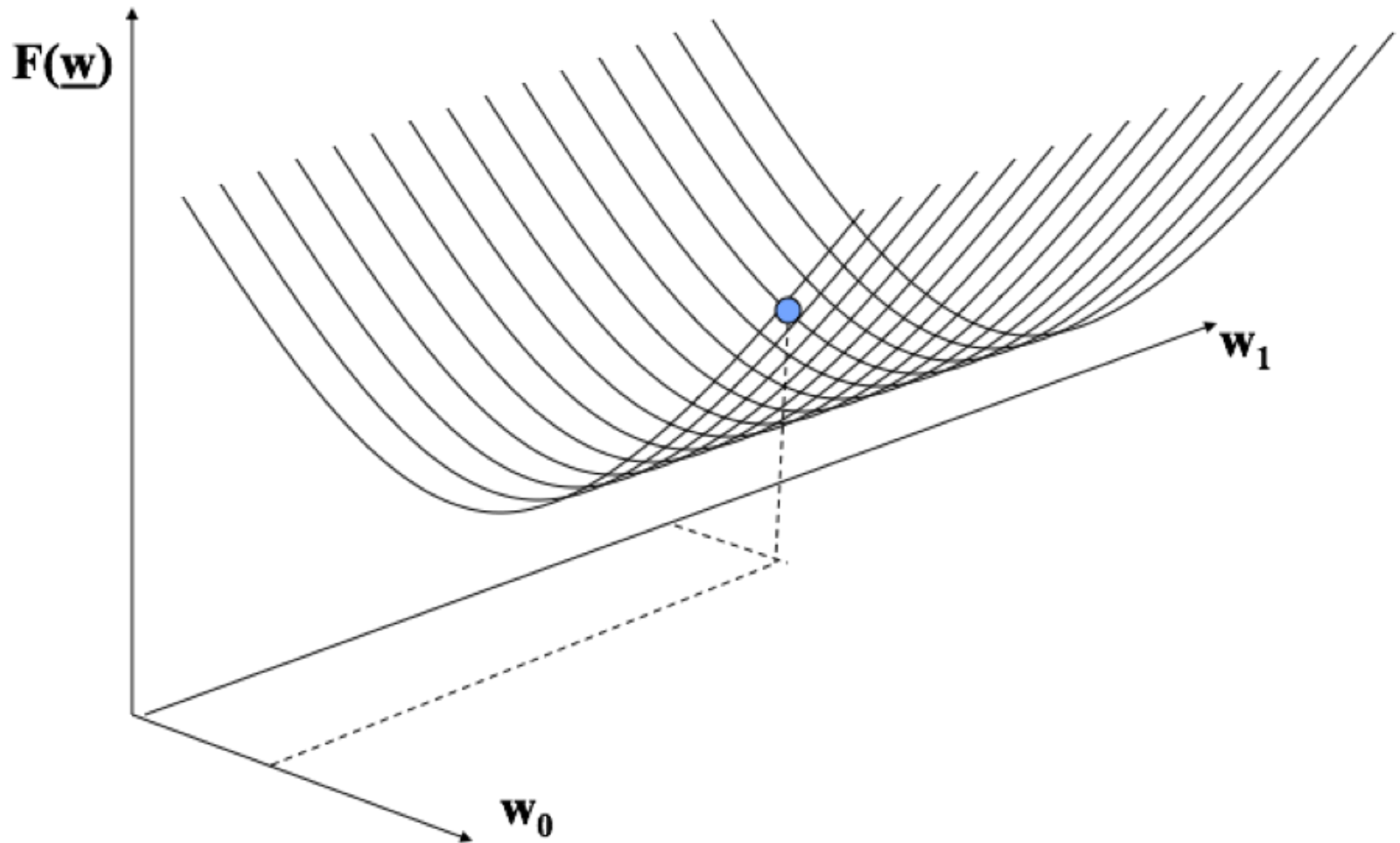
Notes:

1. This moves us downhill in direction  $\Delta(\underline{\mathbf{w}})$  (steepest downhill direction)
2. How far we go is determined by the value of  $\eta$

# Illustration of gradient descent

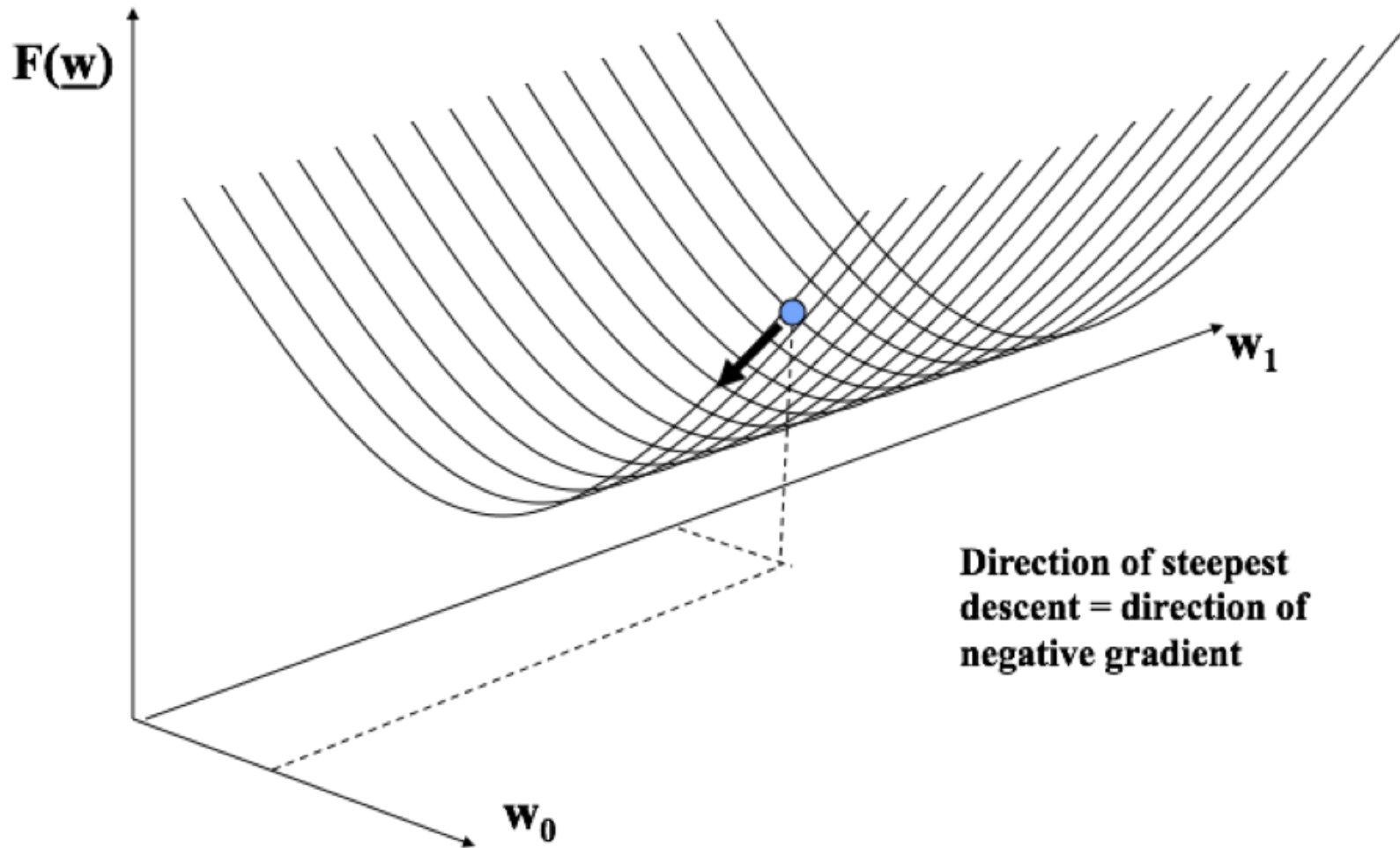


# Illustration of gradient descent

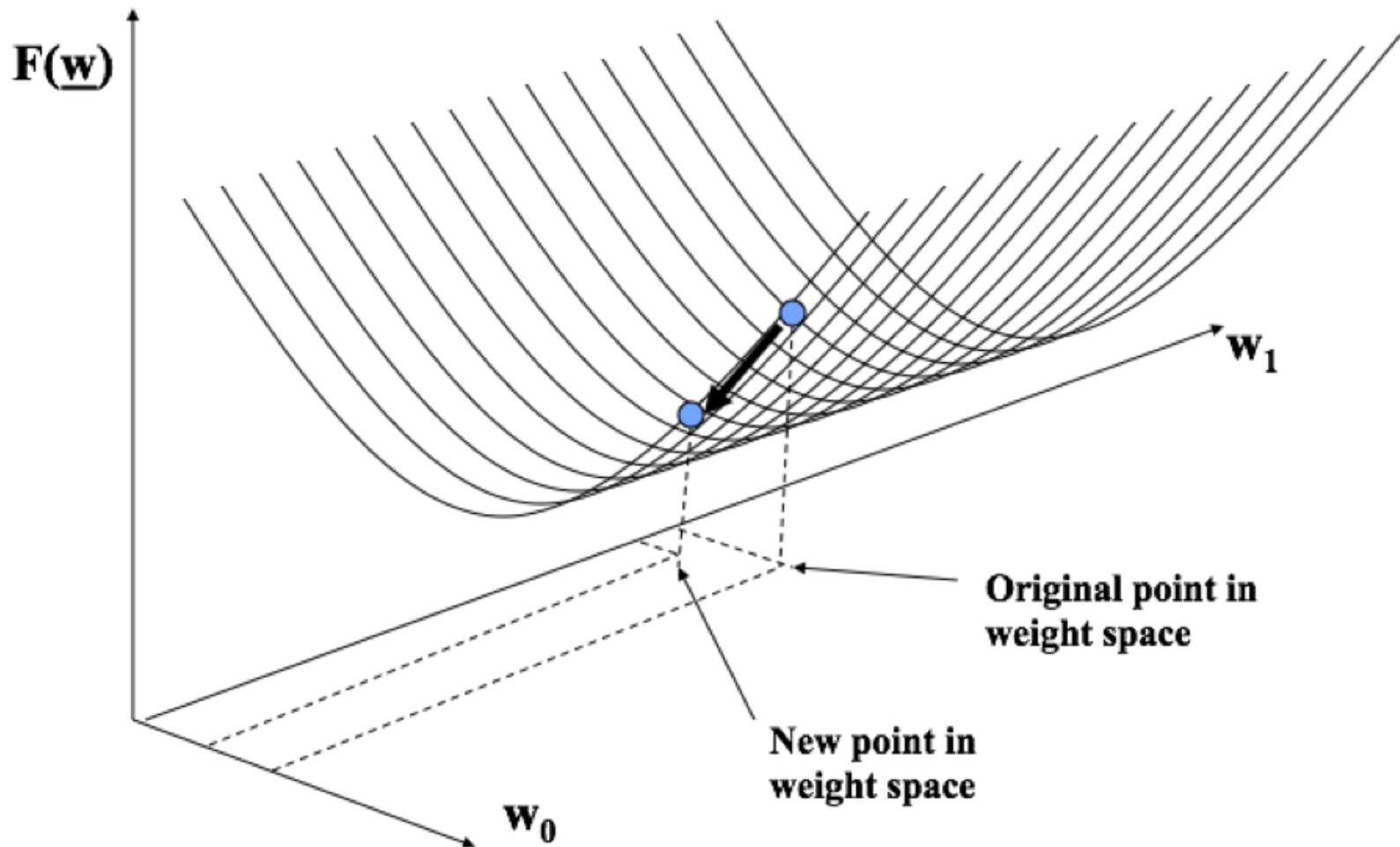




# Illustration of gradient descent



# Illustration of gradient descent

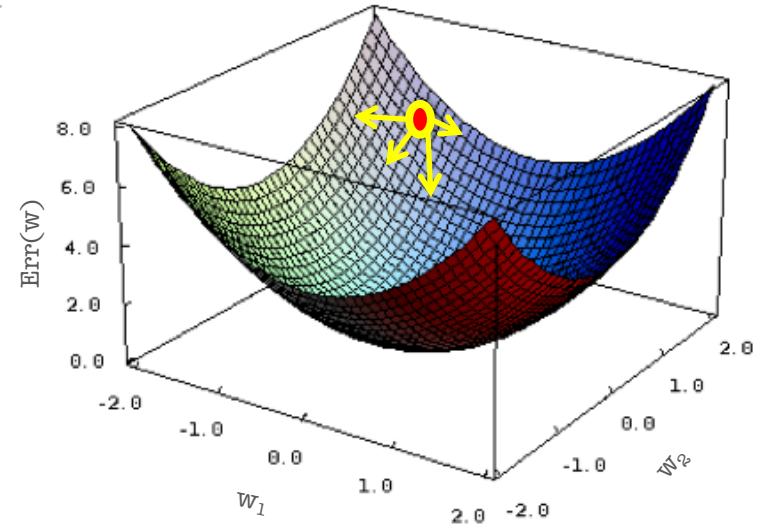


# Error Surface for Squared Loss

We add the  $\frac{1}{2}$  for convenience

$$Err(w) = \frac{1}{2} \sum_{d \in D} (\hat{y}_d - y_d)^2$$

$y = w^T x$



Since  $\hat{y}$  is a constant (for a given dataset), the Error function is a *quadratic function* of  $W$  (paraboloid)

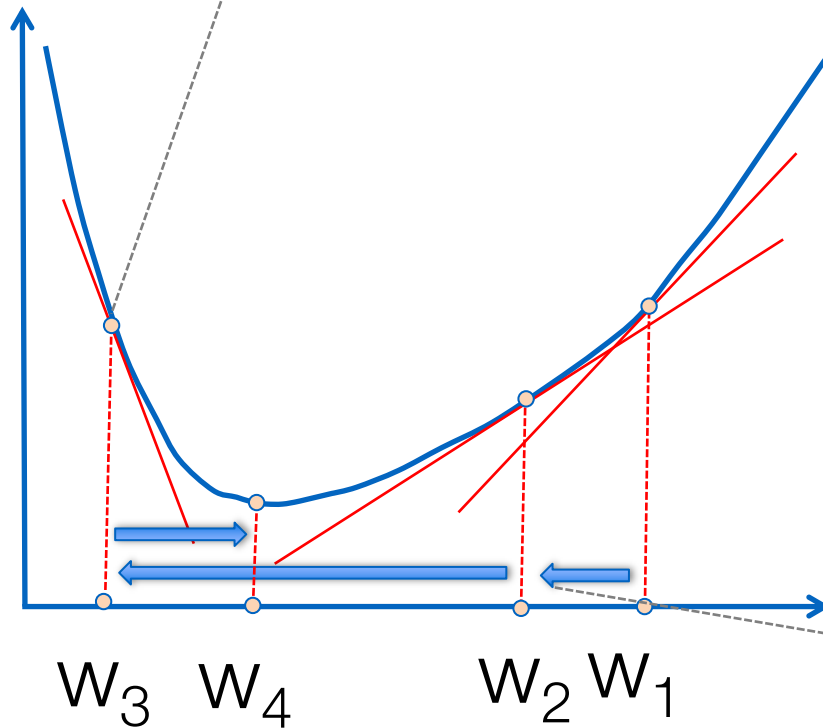
→ Squared Loss function is convex!

How can we find the global minimum?

# Gradient Descent Intuition

(3) We also need to determine the step size (aka learning rate).

*What happens if we overshoot?*



(1) The derivative of the function at  $w_1$  is the slope of the tangent line  
→ Positive slope (increasing)

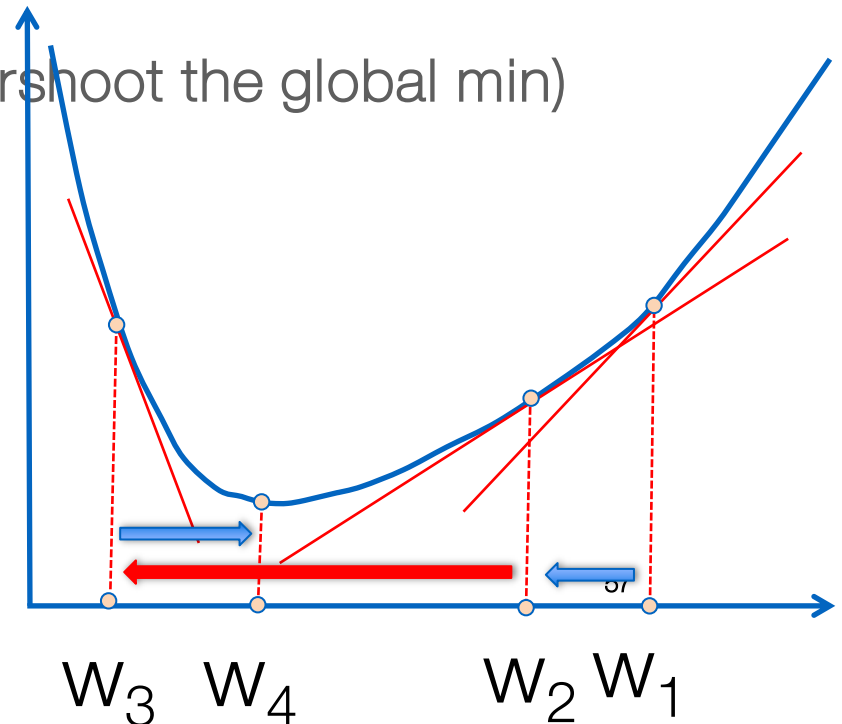
*Which direction should we move to decrease the value of  $Err(w)$  ?*

(2) The gradient determines the direction of steepest increase of  $Err(w)$  *(go in the opposite direction)*

*What is the gradient of  $Error(w)$  at this point?*

# Note about GD step size

- Setting the step size to a very small value
  - Slow convergence rate
- Setting the step size to a large value
  - May oscillate (consistently overshoot the global min)
- Tune experimentally
  - More sophisticated algorithm  
set the value automatically  
(conjugate gradient)



# The Gradient of Error(w)

---

The gradient is a generalization of the derivative

$$\nabla Err(w) = \left( \frac{\partial Err(w)}{\partial w_0}, \frac{\partial Err(w)}{\partial w_1}, \dots, \frac{\partial Err(w)}{\partial w_n} \right)$$

The gradient is a vector of partial derivatives.

It Indicates the *direction of steepest increase* in  $Err(w)$ ,  
for each one of  $w$ 's coordinates

# Gradient Descent Updates

---

- **Compute** the gradient of the training error at each iteration
  - Batch mode: compute the gradient over all training examples

$$\nabla Err(w) = \left( \frac{\partial Err(w)}{\partial w_0}, \frac{\partial Err(w)}{\partial w_1}, \dots, \frac{\partial Err(w)}{\partial w_n} \right)$$

- **Update w:**

Learning rate ( $>0$ )

$$w^{i+1} = w^i - \alpha \nabla Err(w^i)$$

## Computing $\nabla \text{Err}(\mathbf{w}^i)$ for Squared Loss

---

$$\text{Err}(w) = \frac{1}{2} \sum_{d \in D} (y_d - f(x_d))^2$$

$$\begin{aligned} \frac{\partial \text{Err}(\mathbf{w})}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (y_d - f(\mathbf{x}_d))^2 \\ &= \frac{1}{2} \frac{\partial}{\partial w_i} \sum_{d \in D} (y_d - f(\mathbf{x}_d))^2 \\ &= \frac{1}{2} \sum_{d \in D} 2(y_d - f(\mathbf{x}_d)) \frac{\partial}{\partial w_i} (y_d - \mathbf{w} \cdot \mathbf{x}_d) \\ &= - \sum_{d \in D} (y_d - f(\mathbf{x}_d)) x_{di} \end{aligned}$$



## Batch Update Rule for Each $w_i$

---

- **Implementing gradient descent:**

*As you go through the training data, accumulate the change in each  $w_i$  of  $W$*

$$\Delta w_i = \alpha \sum_{d=1}^D (y_d - \mathbf{w}^i \cdot \mathbf{x}_d) x_{di}$$

# Gradient Descent for Squared Loss

---

Initialize  $\mathbf{w}^0$  randomly

for  $i = 0 \dots T$ :

$\Delta \mathbf{w} = (0, \dots, 0)$

for every training item  $d = 1 \dots D$ :

$f(\mathbf{x}_d) = \mathbf{w}^i \cdot \mathbf{x}_d$

for every component of  $\mathbf{w}$   $j = 0 \dots N$ :

$\Delta w_j += \alpha(y_d - f(\mathbf{x}_d)) \cdot x_{dj}$

$\mathbf{w}^{i+1} = \mathbf{w}^i + \Delta \mathbf{w}$

return  $\mathbf{w}^{i+1}$  when it has converged

# Batch vs. Online Learning

---

- The Gradient Descent algorithm makes updates after going over the entire data set
  - Data set can be huge
  - Streaming mode (we cannot assume we saw all the data)
  - Online learning allows “adapting” to changes in the target function
- Stochastic Gradient Descent
  - Similar to GD, updates after each example
  - Can we make the same convergence assumptions as in GD?
- Variations: update after a subset of examples (mini-batch)

# Stochastic Gradient Descent

---

Initialize  $\mathbf{w}^0$  randomly

for  $m = 0 \dots M$ :

$$f(\mathbf{x}_m) = \mathbf{w}^i \cdot \mathbf{x}_m$$

$$\Delta w_j = \alpha(y_d - f(\mathbf{x}_m)) \cdot x_{mj}$$

$$\mathbf{w}^{i+1} = \mathbf{w}^i + \Delta \mathbf{w}$$

return  $\mathbf{w}^{i+1}$  when it has converged