

# Unit 6

## Optimization & Gradient Descent

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 5307: Introduction to Machine Learning, Sp23

# Learning objectives

- Identify the cost function, parameters, and constraints in an **optimization** problem
- Compute the **gradient** of a cost function for scalar, vector, or matrix parameters
- Efficiently compute a gradient in Python
- Write the **gradient-descent** update
- Understand the effect of the **stepsize** on convergence
- Be familiar with adaptive stepsize schemes like the **Armijo rule**
- Be familiar with **constrained** optimization
- Understand the implications of **convexity** for gradient descent
- Determine if a loss function is convex

# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- Gradients of Multi-Variable Functions
- Gradient Descent
- Adaptive Stepsize via the Armijo Rule
- Convexity
- Constrained Optimization

# Motivation: Build an Optimizer for Logistic Regression

- Recall the optimization problem for binary logistic regression with  $y_i \in \{0, 1\}$ :

$$\mathbf{w}_{\text{ml}} \triangleq \arg \min_{\mathbf{w}} \sum_{i=1}^n (\ln[1 + e^{z_i}] - y_i z_i) \quad \text{for } z_i = [1 \ \mathbf{x}_i^T] \mathbf{w}$$

which has no closed-form solution. (For brevity, we use  $w_0 \triangleq b$  here.)

- Previously, we used the `LogisticRegression` method in `sklearn` to solve it:

```
logreg = linear_model.LogisticRegression(C=1e5)
```

```
Xs = preprocessing.scale(X)
logreg.fit(Xs, y)
```

```
LogisticRegression(C=100000.0, class_weight=None, dual=False,
    fit_intercept=True, intercept_scaling=1, max_iter=100,
    multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
    solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
```

- Can we solve this problem ourselves?
- Yes! And the tools we will learn will be very useful later

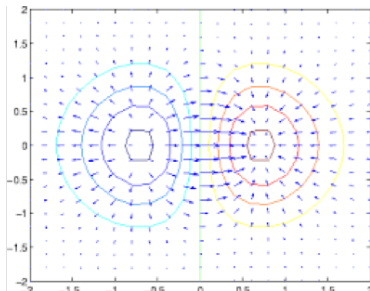
# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- **Gradients of Multi-Variable Functions**
- Gradient Descent
- Adaptive Stepsize via the Armijo Rule
- Convexity
- Constrained Optimization

# Gradients and optimization

- Often, we need to find the **minimizer** of a **cost**, i.e.,  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} J(\mathbf{w})$
- The **gradient**  $\nabla J(\mathbf{w})$  is very useful in this case

- $\nabla J(\hat{\mathbf{w}}) = \mathbf{0}$  at the minimizer  $\hat{\mathbf{w}}$
- $\nabla J(\mathbf{w})$  gives the **direction** of maximum increase and **slope** at  $\mathbf{w}$
- $\nabla J(\mathbf{w})$  exists if  $J(\cdot)$  is differentiable at  $\mathbf{w}$ .
  - We will assume this is the case



- The gradient can also be used to **linearly approximate** a function (details later)

# Definition of the gradient

- Consider a scalar-valued function  $J(\cdot)$
- If the input is vector-valued, then the gradient is vector-valued:

$$\nabla J(\mathbf{w}) = \begin{bmatrix} \partial J(\mathbf{w})/\partial w_1 \\ \vdots \\ \partial J(\mathbf{w})/\partial w_d \end{bmatrix}$$

- If the input is matrix-valued, then the gradient is matrix-valued:

$$\nabla J(\mathbf{W}) = \begin{bmatrix} \partial J(\mathbf{W})/\partial w_{11} & \cdots & \partial J(\mathbf{W})/\partial w_{1k} \\ \vdots & & \vdots \\ \partial J(\mathbf{W})/\partial w_{d1} & \cdots & \partial J(\mathbf{W})/\partial w_{dk} \end{bmatrix}$$

- The gradient always has the same dimensions as the input!

# Example 1

- $J(\mathbf{w}) = w_1^2 + 2w_1w_2^3, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$
- Partial derivatives:
  - $\partial J(\mathbf{w})/\partial w_1 = 2w_1 + 2w_2^3$
  - $\partial J(\mathbf{w})/\partial w_2 = 6w_1w_2^2$
- Gradient:  $\nabla J(\mathbf{w}) = \begin{bmatrix} 2w_1 + 2w_2^3 \\ 6w_1w_2^2 \end{bmatrix}$
- Example on right:
  - computes  $J(\mathbf{w})$  &  $\nabla J(\mathbf{w})$  at  $\mathbf{w} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$
  - gradient is a numpy array

```
def Jeval(w):

    # Function
    J = w[0]**2 + 2*w[0]*(w[1]**3)

    # Gradient
    dJ0 = 2*w[0]+2*(w[1]**3)
    dJ1 = 6*w[0]*(w[1]**2)
    Jgrad = np.array([dJ0, dJ1])

    return J, Jgrad

# Point to evaluate
w = np.array([2,4])
J, Jgrad = Jeval(w)

print('J={0}, Jgrad={1}'.format(J,Jgrad))

J=260, Jgrad=[132 192]
```



## Example 2

- $J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (y_i - ae^{-bx_i})^2, \quad \mathbf{w} = \begin{bmatrix} a \\ b \end{bmatrix}$

- Fits an exponential model to data

- Partial derivatives:

$$\frac{\partial J(\mathbf{w})}{\partial a} = \sum_{i=1}^n (y_i - ae^{-bx_i})(-e^{-bx_i})$$

$$\frac{\partial J(\mathbf{w})}{\partial b} = \sum_{i=1}^n (y_i - ae^{-bx_i})(ax_ie^{-bx_i})$$

- Gradient:

$$\nabla J(\mathbf{w}) = \begin{bmatrix} -\sum_{i=1}^n (y_i - ae^{-bx_i})e^{-bx_i} \\ a\sum_{i=1}^n (y_i - ae^{-bx_i})x_ie^{-bx_i} \end{bmatrix}$$

```
def Jval(y,x,w):
    # Unpack vector
    a = w[0]
    b = w[1]

    # Compute the loss function
    yerr = y-a*np.exp(-b*x)
    J = 0.5*np.sum(yerr**2)

    # Compute the gradient
    dJ_da = -np.sum( yerr*np.exp(-b*x))
    dJ_db = np.sum( yerr*a*x*np.exp(-b*x))
    Jgrad = np.array([dJ_da, dJ_db])

    return J,Jgrad

# Generate some random data
ny = 100
y = np.random.randn(ny)
x = np.random.rand(ny)

# Set some arbitrary parameters
w = np.array([1,2])

# Evaluate cost and gradient
J, Jgrad = Jval(x,y,w)
print('J\t={0}\nJgrad\t={1}'.format(J,Jgrad))

J          =3835.653983662096
Jgrad      =[ 7849.48274896 13608.30853013]
```

# First-order approximation of scalar-input functions

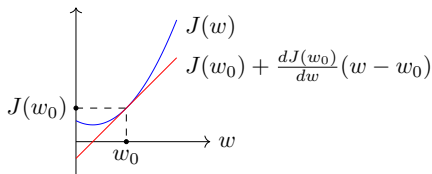
- Consider function  $J(\cdot)$  with scalar input  $w$
- The Taylor series of  $J(\cdot)$  at  $w_0$  can be written as

$$J(w) = J(w_0) + \frac{dJ(w_0)}{dw}(w - w_0) + O((w - w_0)^2)$$

- Informally,  $=O(\epsilon^2)$  means “is at most a constant times  $\epsilon^2$  when  $\epsilon$  is close 0.”  
Formally,  $g(\epsilon) = O(\epsilon^2)$  means  $\exists C, \delta > 0$  s.t.  $g(\epsilon) \leq C\epsilon^2 \forall |\epsilon| \in (0, \delta)$

- The first-order Taylor approximation of  $J(\cdot)$  at  $w_0$ :

$$J(w) \approx J(w_0) + \frac{dJ(w_0)}{dw}(w - w_0)$$



- This approximates  $J(\cdot)$  by a linear function in the neighborhood of  $w_0$
- Note that  $\frac{dJ(w_0)}{dw} = J'(w_0)$  is the slope at  $w_0$
- What if the function has a *vector-valued* input?

# First-order approximation of vector-input functions

- Consider scalar-valued function  $J(\cdot)$  with input  $\mathbf{w} = [w_1, \dots, w_d]^\top$
- Fix a point  $\mathbf{w}_0 = [w_{01}, \dots, w_{0d}]^\top$
- Then the **first-order Taylor approximation of  $J(\cdot)$  at  $\mathbf{w}_0$**  is

$$\begin{aligned}
 J(\mathbf{w}) &= J(\mathbf{w}_0) + \sum_{j=1}^d \frac{\partial J(\mathbf{w}_0)}{\partial w_j} (w_j - w_{0j}) + O(\|\mathbf{w} - \mathbf{w}_0\|^2) \\
 &= J(\mathbf{w}_0) + \sum_{j=1}^d [\nabla J(\mathbf{w}_0)]_j [\mathbf{w} - \mathbf{w}_0]_j + O(\|\mathbf{w} - \mathbf{w}_0\|^2) \\
 &= J(\mathbf{w}_0) + \nabla J(\mathbf{w}_0)^\top (\mathbf{w} - \mathbf{w}_0) + O(\|\mathbf{w} - \mathbf{w}_0\|^2) \\
 &\approx J(\mathbf{w}_0) + \underbrace{\nabla J(\mathbf{w}_0)^\top (\mathbf{w} - \mathbf{w}_0)}_{\langle \nabla J(\mathbf{w}_0), \mathbf{w} - \mathbf{w}_0 \rangle} \quad \text{for } \mathbf{w} \text{ near } \mathbf{w}_0
 \end{aligned}$$

- $\langle \mathbf{a}, \mathbf{b} \rangle \triangleq \mathbf{a}^\top \mathbf{b}$  is the **inner product** for real-valued vectors
- This approximates  $J(\cdot)$  by a linear function in the neighborhood of  $\mathbf{w}_0$

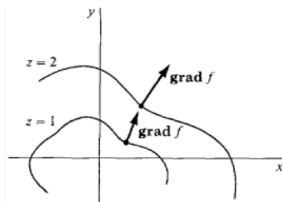
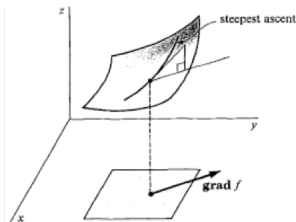
# Understanding the gradient

- The gradient tells both the **direction** of maximum increase and the **slope**. Can we prove this?
- Choose a reference point  $w_0$ , and look at slope in direction  $u$  (where  $\|u\| = 1$ )

$$\frac{J(w_0 + \epsilon u) - J(w_0)}{\epsilon} = \frac{\nabla J(w_0)^T (\epsilon u) + O(\epsilon^2)}{\epsilon}$$

$$\stackrel{\epsilon \rightarrow 0}{=} \nabla J(w_0)^T u = \underbrace{\|\nabla J(w_0)\|}_{\in [-1, 1]} \frac{\nabla J(w_0)^T u}{\|\nabla J(w_0)\| \|u\|}$$

- Cauchy-Schwarz says  $\frac{a^T b}{\|a\| \|b\|} \in [-1, 1]$  for any  $a, b$ , and that  $\frac{a^T b}{\|a\| \|b\|} = 1$  when  $a, b$  are colinear
- Thus  $\|\nabla J(w_0)\|$  is the maximum slope, and it occurs in the direction of  $\nabla J(w_0)$



# First-order approximations of matrix-input functions

- Consider function  $J(\cdot)$  with matrix-valued input  $\mathbf{W} = [w_{ij}]$
- Fix a point  $\mathbf{W}_0$
- The **first-order Taylor approximation** of  $J(\cdot)$  at  $\mathbf{W}_0$  is

$$\begin{aligned}
 J(\mathbf{W}) &= J(\mathbf{W}_0) + \sum_{i=1}^d \sum_{j=1}^k \frac{\partial J(\mathbf{W}_0)}{\partial w_{ij}} (w_{ij} - w_{0,ij}) + O(\|\mathbf{W} - \mathbf{W}_0\|_F^2) \\
 &= J(\mathbf{W}_0) + \sum_{j=1}^k \sum_{i=1}^d [\nabla J(\mathbf{W}_0)^\top]_{ji} [\mathbf{W} - \mathbf{W}_0]_{ij} + O(\|\mathbf{W} - \mathbf{W}_0\|_F^2) \\
 &= J(\mathbf{W}_0) + \text{tr} \{ \nabla J(\mathbf{W}_0)^\top (\mathbf{W} - \mathbf{W}_0) \} + O(\|\mathbf{W} - \mathbf{W}_0\|_F^2) \\
 &\approx J(\mathbf{W}_0) + \underbrace{\text{tr} \{ \nabla J(\mathbf{W}_0)^\top (\mathbf{W} - \mathbf{W}_0) \}}_{\langle \nabla J(\mathbf{W}_0), \mathbf{W} - \mathbf{W}_0 \rangle} \quad \text{for } \mathbf{W} \text{ near } \mathbf{W}_0
 \end{aligned}$$

- $\|\mathbf{A}\|_F^2 = \sum_i \sum_j a_{ij}^2$  is the squared **Frobenius norm**
- $\text{tr}\{\mathbf{A}\} \triangleq \sum_j a_{jj}$  is the **trace** (i.e., sum of diagonal elements)
- $\langle \mathbf{A}, \mathbf{B} \rangle \triangleq \text{tr}\{\mathbf{A}^\top \mathbf{B}\} = \sum_i \sum_j a_{ij} b_{ij}$  is the **inner product** for real-valued matrices

## Example 3

- Suppose that  $J(\mathbf{W}) = \mathbf{a}^\top \mathbf{W} \mathbf{b} = \sum_j \sum_k a_j w_{jk} b_k$  for fixed vectors  $\mathbf{a}$  and  $\mathbf{b}$
- Then  $[\nabla J(\mathbf{W})]_{j'k'} = \frac{\partial J(\mathbf{W})}{\partial w_{j'k'}} = a_{j'} b_{k'}$  for any  $\mathbf{W}$
- Thus the gradient matrix is  $\nabla J(\mathbf{W}) = \mathbf{a} \mathbf{b}^\top$  for any  $\mathbf{W}$
- The linear approximation of  $J(\cdot)$  at  $\mathbf{W}_0$  is

$$\begin{aligned}
 J(\mathbf{W}) &\approx J(\mathbf{W}_0) + \text{tr}\{\nabla J(\mathbf{W}_0)^\top (\mathbf{W} - \mathbf{W}_0)\} \\
 &= \mathbf{a}^\top \mathbf{W}_0 \mathbf{b} + \text{tr}\{(\mathbf{a} \mathbf{b}^\top)^\top (\mathbf{W} - \mathbf{W}_0)\} \\
 &= \mathbf{a}^\top \mathbf{W}_0 \mathbf{b} + \text{tr}\{\mathbf{b} \mathbf{a}^\top (\mathbf{W} - \mathbf{W}_0)\} \\
 &= \mathbf{a}^\top \mathbf{W}_0 \mathbf{b} + \text{tr}\{\mathbf{a}^\top (\mathbf{W} - \mathbf{W}_0) \mathbf{b}\} \\
 &= \mathbf{a}^\top \mathbf{W}_0 \mathbf{b} + \mathbf{a}^\top (\mathbf{W} - \mathbf{W}_0) \mathbf{b} \\
 &= \mathbf{a}^\top \mathbf{W} \mathbf{b}
 \end{aligned}$$

since  $\text{tr}\{\mathbf{B} \mathbf{A}\} = \text{tr}\{\mathbf{A} \mathbf{B}\}$

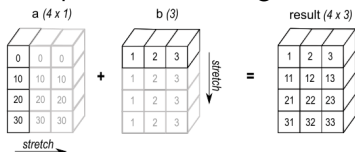
since  $\text{tr}\{c\} = c$  for scalar  $c$

perfect approx since  $J$  is linear!

## Example 3 in Python

- Function  $J(W) = a^T W b$ 
  - In Python, can use `.dot` for matrix multiplication
- Gradient  $\nabla J(W) = a b^T$ 
  - Want to set  $Jgrad[j,k] = a[j]b[k]$
  - But want to avoid for-loops
  - Use **Python broadcasting**!
    - `a[:,None]` is  $m \times 1$
    - `b[None,:]` is  $1 \times n$
    - `"*"` stretches as needed

- Example of broadcasting addition:



```
def Jeval(W,a,b):
    # Function
    J = a.dot(W.dot(b))

    # Gradient -- Use python broadcasting
    Jgrad = a[:,None]*b[None,:]

    return J,Jgrad

# Some random data
m = 4
n = 3
W = np.random.randn(m,n)
a = np.random.randn(m)
b = np.random.randn(n)

J,Jgrad = Jeval(W,a,b)

print('J\t={0}\nJgrad\t={1}'.format(J,Jgrad))
```

```
J          =-0.18205418344529506
Jgrad      =[[ -0.19137257 -0.09198831 -0.1026666 ]
 [ 1.66114969  0.79847574  0.89116532]
 [-0.93328792 -0.44860964 -0.50068566]
 [ 1.00634285  0.48372543  0.5398778 ]]
```

# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- Gradients of Multi-Variable Functions
- **Gradient Descent**
- Adaptive Stepsize via the Armijo Rule
- Convexity
- Constrained Optimization



# Stationary points

- A **stationary point** of  $J(\cdot)$  is any  $w$  such that

$$\nabla J(w) = 0$$

- The unconstrained minimizer of differentiable  $J$

$$\hat{w} = \arg \min_w J(w)$$

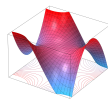
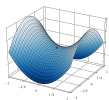
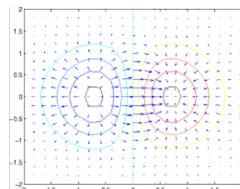
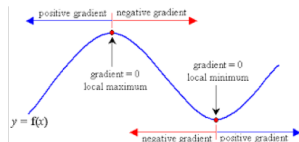
is one such stationary point

- In general, stationary points can be either

- **minimizers**,
- **maximizers**, or
- **saddle points**

- But often we cannot explicitly solve  $\nabla J(w) = 0$

- Instead, use a numerical approach to find  $\hat{w}$

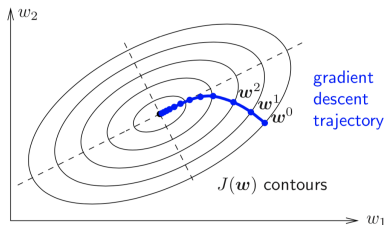


# Gradient descent

- Goal: Find the minimizer of  $J(\cdot)$ , i.e.,  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} J(\mathbf{w})$ 
  - $J(\cdot)$  is called the **objective** or **cost** or **loss** function
  - The optimization is “**unconstrained**” since there are no constraints on  $\mathbf{w}$
  - Will assume that  $\nabla J(\mathbf{w})$  exists (i.e.,  $J(\mathbf{w})$  is differentiable) at all  $\mathbf{w}$

- **Gradient descent (GD) algorithm:**

- Choose an initial  $\mathbf{w}^0$ , then iterate  $\mathbf{w}^{k+1} = \mathbf{w}^k - \alpha_k \nabla J(\mathbf{w}^k)$  until convergence
- In words: take small downhill steps until you reach the bottom
  - $\alpha_k > 0$  is the **stepsize** or **learning rate**
  - Often  $\mathbf{w}^0$  is chosen randomly



# Analysis of gradient descent

- The Taylor series expansion of  $J(\cdot)$  at  $\mathbf{w}^k$  is

$$J(\mathbf{w}) = J(\mathbf{w}^k) + \nabla J(\mathbf{w}^k)^\top (\mathbf{w} - \mathbf{w}^k) + O(\|\mathbf{w} - \mathbf{w}^k\|^2)$$

- Evaluating this at  $\mathbf{w} = \mathbf{w}^{k+1}$  yields

$$J(\mathbf{w}^{k+1}) = J(\mathbf{w}^k) + \nabla J(\mathbf{w}^k)^\top (\mathbf{w}^{k+1} - \mathbf{w}^k) + O(\|\mathbf{w}^{k+1} - \mathbf{w}^k\|^2)$$

- From the GD update, we know  $\mathbf{w}^{k+1} - \mathbf{w}^k = -\alpha_k \nabla J(\mathbf{w}^k)$ , and so

$$\begin{aligned} J(\mathbf{w}^{k+1}) &= J(\mathbf{w}^k) - \alpha_k \nabla J(\mathbf{w}^k)^\top \nabla J(\mathbf{w}^k) + O(\alpha_k^2 \|\nabla J(\mathbf{w}^k)\|^2) \\ &= J(\mathbf{w}^k) - \alpha_k \|\nabla J(\mathbf{w}^k)\|^2 + \underbrace{O(\alpha_k^2)}_{\leq C\alpha_k^2} \|\nabla J(\mathbf{w}^k)\|^2 \\ &\leq C\alpha_k^2 \end{aligned}$$

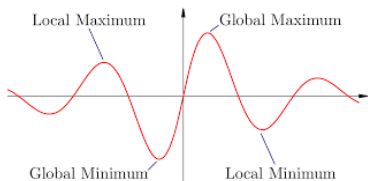
- Thus, if stepsize  $\alpha_k$  is sufficiently small, then  $J(\mathbf{w}^{k+1}) \leq J(\mathbf{w}^k)$ 
  - Why? For any  $C > 0$ , the choice  $\alpha_k < 1/C$  will make  $C\alpha_k^2 < \alpha_k$ , in which case the middle term will dominate the last term and GD will make progress

# Local versus global minima

## ■ Definitions

- $J(\hat{\mathbf{w}})$  is a **global minima** if  $J(\hat{\mathbf{w}}) \leq J(\mathbf{w}) \forall \mathbf{w}$
- $J(\hat{\mathbf{w}})$  is a **local minima** if  $J(\hat{\mathbf{w}}) \leq J(\mathbf{w}) \forall \mathbf{w}$  in some open neighborhood of  $\hat{\mathbf{w}}$

- In most cases, gradient descent only guarantees convergence to a **local minimum**
- For a **convex function**, any local minimum is a global minimum! (Will discuss more later ...)



# Gradient of cross-entropy loss

- Recall the binary logistic regression problem when  $y_i \in \{0, 1\}$ :

$$\mathbf{w}_{\text{ml}} \triangleq \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n (\ln[1 + e^{z_i}] - y_i z_i)}_{\text{"cross entropy loss"} J(\mathbf{w})} \quad \text{for } z_i = [\mathbf{1} \ \mathbf{x}_i^T] \mathbf{w}$$

- To solve this problem, think of cost  $J(\mathbf{w})$  as a composition of two functions:

1) **linear transformation**:  $\mathbf{z}(\mathbf{w}) = \mathbf{A}\mathbf{w}$  for  $\mathbf{A} = [\mathbf{1} \ \mathbf{X}]$

2) **separable** function:  $f(\mathbf{z}) = \sum_{i=1}^n f_i(z_i)$  for  $f_i(z_i) = \ln[1 + e^{z_i}] - y_i z_i$   
so that  $J(\mathbf{w}) = f(\mathbf{z}(\mathbf{w}))$

- Then apply the **multivariable chain rule**:

$$\frac{\partial J}{\partial w_j} = \sum_i \frac{\partial f}{\partial z_i} \frac{\partial z_i}{\partial w_j}$$

- Here,  $\frac{\partial f}{\partial z_i} = \frac{1}{1 + e^{z_i}} e^{z_i} - y_i = \frac{1}{e^{-z_i} + 1} - y_i$  and  $\frac{\partial z_i}{\partial w_j} = a_{ij}$

# Computing cost *and* gradient

- Usually we want to compute *both*  $J(\mathbf{w})$  and  $\nabla J(\mathbf{w})$

- Forward pass:** Compute cost:

- First compute  $\mathbf{z} = \mathbf{A}\mathbf{w}$

- Then  $f_i(z_i) = \ln[1 + e^{z_i}] - y_i z_i$   

$$= -\ln\left(\frac{e^{z_i}}{1+e^{z_i}}\right) + (1 - y_i)z_i$$
  

$$= -\ln\left(\frac{1}{1+e^{-z_i}}\right) + (1 - y_i)z_i$$

- Finally  $J(\mathbf{w}) = f(\mathbf{z}) = \sum_i f_i(z_i)$

- Backward pass:** Compute gradient:

- Multivariate chain rule:

$$\frac{\partial J}{\partial w_j} = \sum_i \frac{\partial f}{\partial z_i} \frac{\partial z_i}{\partial w_j} = \sum_i \frac{\partial f}{\partial z_i} a_{ij}$$

$$\text{with } \frac{\partial f}{\partial z_i} = \frac{1}{1+e^{-z_i}} - y_i = [\nabla f]_i$$

$$\text{So } \nabla J(\mathbf{w}) = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \vdots \\ \frac{\partial J}{\partial w_d} \end{bmatrix} = \begin{bmatrix} \sum_i a_{i0} [\nabla f]_i \\ \vdots \\ \sum_i a_{id} [\nabla f]_i \end{bmatrix} = \begin{bmatrix} a_{10} & \cdots & a_{n0} \\ \vdots & & \vdots \\ a_{1d} & \cdots & a_{nd} \end{bmatrix} \begin{bmatrix} [\nabla f]_1 \\ \vdots \\ [\nabla f]_n \end{bmatrix} = \mathbf{A}^T \nabla f$$

```
def Jeval_param(w,X,y):
    """
    Compute the loss and gradient of w given w,X,y
    """
    # Construct transform matrix
    n = X.shape[0]
    A = np.column_stack((np.ones(n,), X))

    # The loss is the binary cross entropy
    z = A.dot(w)
    py = 1/(1+np.exp(-z))
    J = np.sum((1-y)*z - np.log(py))

    # Gradient
    df_dz = py-y
    Jgrad = A.T.dot(df_dz)
    return J, Jgrad
```

# Only- $w$ -dependence using a lambda function

When implementing gradient descent in **Python**, we want a cost/gradient evaluation function that **only depends on  $w$**

- First create a function that
  - computes cost & gradient
  - given  $w, X, y$
- Then create a “**lambda function**” that
  - fixes  $X, y$  at training values
  - same as “anonymous function” in Matlab:

`Jeval = @(w) Jeval_param(w,Xtr,ytr)`

```
Jeval = lambda w: Jeval_param(w,Xtr,ytr)

# Evaluate J and Jgrad at w0
J0, Jgrad0 = Jeval(w0)
```

```
def Jeval_param(w,X,y):
    """
    Compute the loss and gradient of w given w,X,y
    """
    # Construct transform matrix
    n = X.shape[0]
    A = np.column_stack((np.ones(n), X))

    # The loss is the binary cross entropy
    z = A.dot(w)
    py = 1/(1+np.exp(-z))
    J = np.sum((1-y)*z - np.log(py))

    # Gradient
    df_dz = py-y
    Jgrad = A.T.dot(df_dz)
    return J, Jgrad
```

# Only- $w$ -dependence by creating a Python class

A more powerful approach to building an only- $w$ -dependent function is to create a “class” (i.e., object-oriented programming)

- Includes an **constructor** that

- loads data ( $X, y$ )
- does pre-computations

```
# Instantiate the class
log_fun = LogisticFun(Xtr,ytr)
```

- Includes an **Jeval** function to

- compute cost & gradient
- using data stored in the class

```
# Call the method
J0, Jgrad0 = log_fun.Jeval(w0)
```

```
class LogisticFun(object):
    def __init__(self,X,y):
        """
        The constructor takes in the training features `X` and
        """
        self.X = X
        self.y = y
        n = X.shape[0]
        self.A = np.column_stack((np.ones(n,), X))

    def Jeval(self,w):
        """
        The Jeval method computes the loss and gradient at weight `w`
        """
        # The loss is the binary cross entropy
        z = self.A.dot(w)
        py = 1/(1+np.exp(-z))
        J = np.sum((1-self.y)*z - np.log(py))

        # Gradient
        df_dz = py-self.y
        Jgrad = self.A.T.dot(df_dz)
        return J, Jgrad
```



# Always check your gradient implementation!

- Randomly pick  $w_0$  and  $w_1$  that are close together
- Check that  $J(w_1) - J(w_0) \approx \nabla J(w_0)^T(w_1 - w_0)$

```
# Take a random initial point
d = X.shape[1]
w0 = np.random.randn(d+1)

# Perturb the point
step = 1e-6
w1 = w0 + step*np.random.randn(d+1)

# Measure the function and gradient at w0 and w1
J0, Jgrad0 = log_fun.Jeval(w0)
J1, Jgrad1 = log_fun.Jeval(w1)

# Predict the amount the function should have changed based on the gradient
dJ_pred = Jgrad0.dot(w1-w0)

# Print the two values to see if they are close
print("Actual J1-J0    = %12.4e" % (J1-J0))
print("Predicted J1-J0 = %12.4e" % dJ_pred)

Actual J1-J0    =  8.7296e-04
Predicted J1-J0 =  8.7296e-04
```

- The above follows from  $J(w_1) \approx J(w_0) + \nabla J(w_0)^T(w_1 - w_0)$

# A simple implementation of gradient descent

- We now implement gradient descent
- Inputs:
  - $J_{\text{eval}}(w)$ : function that computes cost & gradient
  - initial parameters  $w^0$
  - learning rate  $\alpha$
  - number of iterations  $m$
- Outputs:
  - final cost & parameters  $w^m$
  - history of cost & parameters (for debugging)

```
def grad_opt_simp(Jeval, winit, lr=1e-3, nit=1000):
    """
    Simple gradient descent optimization

    Jeval: A function that returns f, fgrad, the objective
           function and its gradient
    winit: Initial estimate
    lr:    learning rate
    nit:   Number of iterations
    """
    # Initialize
    w0 = winit

    # Create history dictionary for tracking progress per iteration.
    # This isn't necessary if you just want the final answer, but it
    # is useful for debugging
    hist = {'w': [], 'J': []}

    # Loop over iterations
    for it in range(nit):

        # Evaluate the function and gradient
        J0, Jgrad0 = Jeval(w0)

        # Take a gradient step
        w0 = w0 - lr * Jgrad0

        # Save history
        hist['J'].append(J0)
        hist['w'].append(w0)

    # Convert to numpy arrays
    for elem in ('J', 'w'):
        hist[elem] = np.array(hist[elem])

    return w0, J0, hist
```

# Gradient descent for logistic regression

- random initialization
- 2000 iterations
- convergence is slow!
- test accuracy is not great!
  - weights not converged

```
def predict(X,w):
    z = X.dot(w[1:]) + w[0]
    yhat = (z > 0)
    return yhat

yhat = predict(Xts,w)
acc = np.mean(yhat == yts)
print("Test accuracy = %f" % acc)

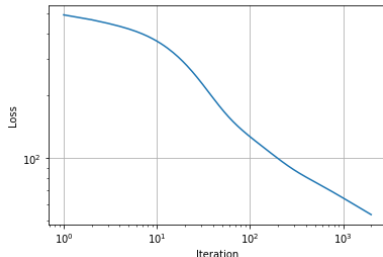
Test accuracy = 0.968198
```

```
# Initial condition
winit = rng.standard_normal(d+1)

# Parameters
Jeval = log_fun.Jeval
nit = 2000
lr = 1e-4

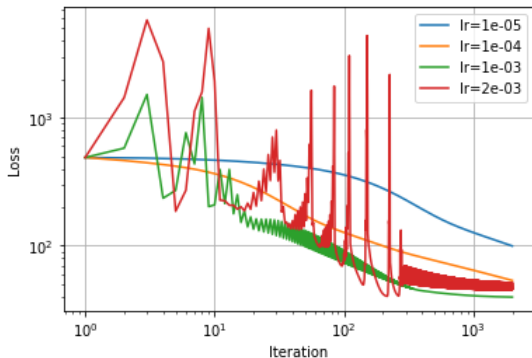
# Run the gradient descent
w, J0, hist = grad_opt_simp(Jeval, winit, lr=lr, nit=nit)

# Plot the training loss
t = 1+np.arange(nit)
plt.loglog(t, hist['J'])
plt.grid()
plt.ylabel('Loss')
plt.xlabel('Iteration');
```



# Effect of stepsize (or learning rate)

- stepsize too small  $\Rightarrow$  slow convergence
- stepsize too large  $\Rightarrow$  instability (overshoots optimal solution)



$lr= 1.00e-05$ , Loss = 99.58, Test accuracy = 0.9187  
 $lr= 1.00e-04$ , Loss = 53.56, Test accuracy = 0.9859  
 $lr= 1.00e-03$ , Loss = 39.68, Test accuracy = 0.9894  
 $lr= 2.00e-03$ , Loss = 51.46, Test accuracy = 0.9859

# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- Gradients of Multi-Variable Functions
- Gradient Descent
- Adaptive Stepsize via the Armijo Rule
- Convexity
- Constrained Optimization

# The Armijo approach to adaptive stepsize

- Recall our gradient-descent analysis result:

$$J(\mathbf{w}^{k+1}; \alpha_k) = J(\mathbf{w}^k) - \alpha_k \|\nabla J(\mathbf{w}^k)\|^2 + O(\alpha_k^2 \|\nabla J(\mathbf{w}^k)\|^2)$$

- Armijo rule:**

- At each  $k$ , choose some stepsize  $\alpha_k > 0$  and check if

$$J(\mathbf{w}^{k+1}; \alpha_k) \leq J(\mathbf{w}^k) - c \alpha_k \|\nabla J(\mathbf{w}^k)\|^2$$

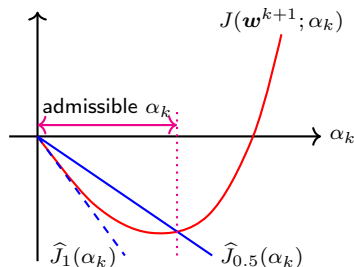
- Here  $c \in (0, 1)$  is a design parameter, for example  $c = 0.5$
- If yes, cost guaranteed to decrease (unless  $\nabla J(\mathbf{w}^k) = \mathbf{0}$ , when it stays the same)
  - Decreases by fraction  $c$  of that predicted by linear approximation of  $J(\mathbf{w}^{k+1})$
- A **simple Armijo-based  $\alpha_k$ -update:**
  - If Armijo rule passes: accept  $\mathbf{w}^{(k+1)}$  and set  $\alpha_{k+1} = \beta_{\text{incr}} \alpha_k$  for some  $\beta_{\text{incr}} > 1$
  - If Armijo rule fails: reject  $\mathbf{w}^{(k+1)}$  and set  $\alpha_{k+1} = \beta_{\text{decr}} \alpha_k$  for some  $\beta_{\text{decr}} < 1$
- Alternative: try several  $\alpha_k$  and choose the one with smallest  $J(\mathbf{w}^{k+1}; \alpha_k)$ 
  - Can be more accurate than Armijo, but much more expensive

# Armijo rule illustrated

- Recall  $J(\mathbf{w}^{k+1}; \alpha_k) = J(\mathbf{w}^k) - \alpha_k \|\nabla J(\mathbf{w}^k)\|^2 + O(\alpha_k^2 \|\nabla J(\mathbf{w}^k)\|^2)$
- Recall Armijo rule: fix  $c \in (0, 1)$  and accept any stepsize  $\alpha_k > 0$  satisfying

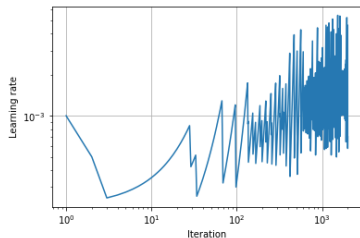
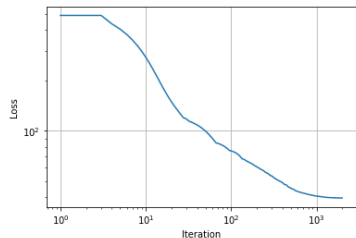
$$J(\mathbf{w}^{k+1}; \alpha_k) \leq J(\mathbf{w}^k) - c \alpha_k \|\nabla J(\mathbf{w}^k)\|^2 \triangleq \hat{J}_c(\alpha_k)$$

- Red curve shows  $J(\mathbf{w}^{k+1}; \alpha_k)$  versus  $\alpha_k$ 
  - line-search would give samples of this
- Dashed line shows  $\hat{J}_1(\alpha_k)$  versus  $\alpha_k$ 
  - this is the  $J(\mathbf{w}^{k+1}; \alpha_k)$  predicted by linear approximation
- Blue line shows  $\hat{J}_{0.5}(\alpha_k)$  versus  $\alpha_k$ 
  - purple arrows show range of  $\alpha_k$  satisfying the Armijo rule with  $c = 0.5$
  - what about other values of  $c \in (0, 1)$ ?



# Armijo example in Python

- The Armijo method applied to logistic regression:



```
# Loop over GD iterations
for it in range(nit):

    # Take a gradient step
    w1 = w0 - lr*Jgrad0
    J1, Jgrad1 = Jeval(w1)

    # Check if update passes the Armijo condition
    if (J1 < J0 - c*lr*np.linalg.norm(Jgrad0)**2):
        # If passes...
        w0 = w1 # accept the update
        J0 = J1 # accept the update
        Jgrad0 = Jgrad1 # accept the update
        lr = lr*beta_incr # increase learning rate
    else:
        # If fails...
        lr = lr*beta_decr # decrease learning rate
```

Loss is better than best fixed-stepsize method:

```
yhat = predict(Xts,w)
acc = np.mean(yhat == yts)
print("Loss = %6.2f, Test accuracy = %6.4f" % (J0,acc))
```

Loss = 39.55, Test accuracy = 0.9894



# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- Gradients of Multi-Variable Functions
- Gradient Descent
- Adaptive Stepsize via the Armijo Rule
- Convexity
- Constrained Optimization

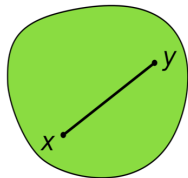
# Convex sets

- A set  $D$  is **convex** if, for any  $x, y \in D$  and  $t \in [0, 1]$

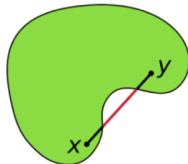
$$\underbrace{tx + (1 - t)y}_{\text{convex combination of } x \text{ and } y} \in D$$

- The line between any two points in  $D$  remains in  $D$
- Examples of convex sets:
  - Circle, square, ellipse
  - $\mathbb{R}^n$
  - a hyperplane in  $\mathbb{R}^n$
  - a half-space in  $\mathbb{R}^n$

a convex set



a nonconvex set

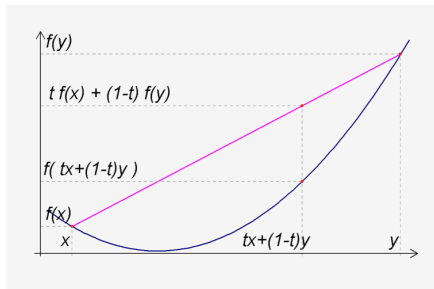


# Convex functions

- A function  $f(\cdot) \in \mathbb{R}$  is **convex** if
  - 1) its domain  $D$  is a convex set, and
  - 2) for any  $x, y \in D$  and  $t \in [0, 1]$ ,
 
$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y)$$

- Examples of convex functions:

- $f(x) = ax + b$
- $f(x) = a^T x + b$
- $f(x) = ax^2 + bx + c$  iff  $a \geq 0$
- $f(x)$  is convex if  $f''(x)$  exists everywhere and  $f''(x) \geq 0$ 
  - vector case: Hessian must exist everywhere and be positive semidefinite
- norms are convex (e.g.,  $\|x\|$ ,  $\|x\|_1$ )
- if  $f$  and  $g$  are convex, then so is  $f + g$
- if  $f$  and  $g$  are convex and  $f$  is non-decreasing, then  $f(g(\cdot))$  is convex
- RSS, logistic loss, and their L1 or L2 regularized versions are all convex



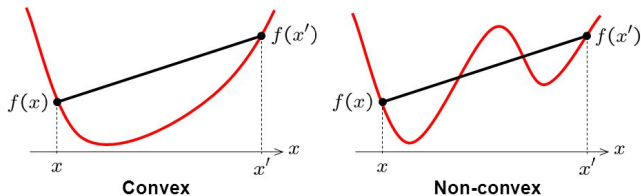
# Local minimizers of convex functions

## Theorem

If  $f(\cdot)$  is **convex** and  $x$  is a local minimizer of  $f$ , then  $x$  is a global minimizer of  $f$

### ■ Implications for optimization:

- Recall: with proper stepsize, gradient descent converges to a local minimizer
- But local minimizers are not always global minimizers!
- With a **convex** function, gradient descent converges to a *global* minimizer



# Outline

- Motivating Example: Build an Optimizer for Logistic Regression
- Gradients of Multi-Variable Functions
- Gradient Descent
- Adaptive Stepsize via the Armijo Rule
- Convexity
- **Constrained Optimization**

# Constrained Optimization

- In some cases, we have constraints on the design variables  $\mathbf{w}$ :
  - **equality constraints**:  $h_l(\mathbf{w}) = 0$  for  $l = 1 \dots L$
  - **inequality constraints**:  $g_m(\mathbf{w}) \leq 0$  for  $m = 1 \dots M$
- To handle these constraints, we can reformulate the constrained problem

$$\arg \min_{\mathbf{w}} J(\mathbf{w}) \text{ such that } \begin{cases} h_l(\mathbf{w}) = 0, & l = 1 \dots L \\ g_m(\mathbf{w}) \leq 0, & m = 1 \dots M \end{cases}$$

as an *unconstrained* one involving additional design variables!

- With only equality constraints, we'd add **Lagrange multipliers**  $\{\lambda_l\}_{l=1}^L$ :

$$\arg \min_{\mathbf{w}, \boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}) \text{ with } \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}) \triangleq J(\mathbf{w}) + \sum_{l=1}^L \lambda_l h_l(\mathbf{w}),$$

which exploits the fact  $0 = \frac{\partial \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda})}{\partial \lambda_l} = h_l(\mathbf{w}) \forall l = 1 \dots L$  at the optimum  $(\mathbf{w}, \boldsymbol{\lambda})$

- With **inequality constraints**:  $\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = J(\mathbf{w}) + \sum_l \lambda_l h_l(\mathbf{w}) + \sum_m \mu_m g_m(\mathbf{w})$

# Other optimization topics

- There are many topics that we did not cover, e.g.,
  - Newton's method and quasi-Newton methods (i.e., matrix-valued  $\alpha_k$ )
  - non-differentiable optimization (i.e., gradient does not exist everywhere: LASSO)
- Also, our Armijo-based optimizer works well with convex functions, but it may not work well with non-convex functions, as encountered in deep learning
- Take an optimization class and learn more!
  - ECE-5500 Intro to Optimization (Autumn every year) (was 5759)
  - ECE-6500 Convex Optimization (Spring odd years) (was 7100)
  - ECE-8101 Nonconvex Optimization for Machine Learning (Autumn even years)

# Learning objectives

- Identify the cost function, parameters, and constraints in an **optimization** problem
- Compute the **gradient** of a cost function for scalar, vector, or matrix parameters
- Efficiently compute a gradient in Python
- Write the **gradient-descent** update
- Understand the effect of the **stepsize** on convergence
- Be familiar with adaptive stepsize schemes like the **Armijo rule**
- Be familiar with **constrained** optimization
- Understand the implications of **convexity** for gradient descent
- Determine if a loss function is convex