

Unit 6

Optimization & Gradient Descent

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 4300: Introduction to Machine Learning, Sp20

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**
- 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

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 $w_0 = b$ 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)
```

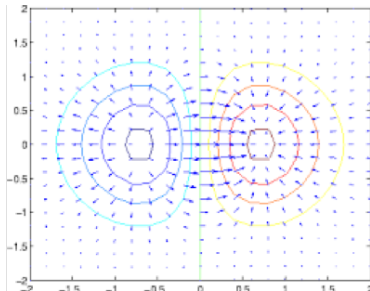
- Can we solve this problem ourselves?
- Yes! And the tools we will learn will be very useful later

Outline

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

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 sufficiently smooth 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 $f(\cdot)$
- If the input is vector-valued, then the gradient is vector-valued:

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

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

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

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

Example 1

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

```
def feval(w):

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

    # Gradient
    df0 = 2*w[0]+2*(w[1]**3)
    df1 = 6*w[0]*(w[1]**2)
    fgrad = np.array([df0, df1])

    return f, fgrad

# Point to evaluate
w = np.array([2,4])
f, fgrad = feval(w)

f      = 260.000000
fgrad  = [132 192]
```


Example 2

- Loss function:

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

- Fits an exponential model to data
- Partial derivatives:

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

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

- Gradient:

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

```
def Jeval(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
```

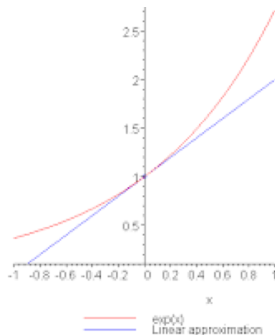
First-order approximation of scalar-input functions

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

$$f(w) = f(w_0) + \frac{df(w_0)}{dw}(w - w_0) + O((w - w_0)^2)$$
 - $O(\epsilon^2)$: “grows no faster than $C\epsilon^2$ for some $C > 0$ ”
- The first-order Taylor approximation of $f(w)$ at w_0 :

$$f(w) \approx f(w_0) + \frac{df(w_0)}{dw}(w - w_0)$$

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



First-order approximation of vector-input functions

- Consider function $f(\mathbf{w})$ with vector-valued 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 $f(\mathbf{w})$ at \mathbf{w}_0** is

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

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

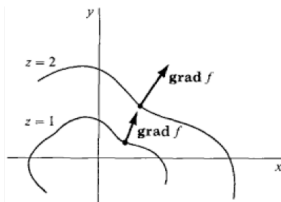
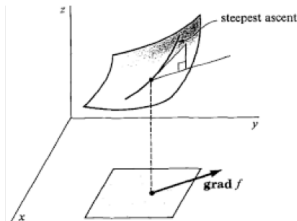
Understanding the gradient

- The gradient tells both the **direction** of maximum increase and the **slope**. Why?
- Choose a reference point \mathbf{w}_0 , and look at slope in direction \mathbf{u} (where $\|\mathbf{u}\| = 1$)

$$\frac{f(\mathbf{w}_0 + \epsilon \mathbf{u}) - f(\mathbf{w}_0)}{\epsilon} = \frac{\nabla f(\mathbf{w}_0)^T (\epsilon \mathbf{u}) + O(\epsilon^2)}{\epsilon}$$

$$\stackrel{\epsilon \rightarrow 0}{=} \nabla f(\mathbf{w}_0)^T \mathbf{u} = \underbrace{\|\nabla f(\mathbf{w}_0)\| \left(\frac{\nabla f(\mathbf{w}_0)}{\|\nabla f(\mathbf{w}_0)\|} \right)^T}_{\in [-1, 1]} \frac{\mathbf{u}}{\|\mathbf{u}\|}$$

- Cauchy-Schwarz says $\frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \in [-1, 1]$ for any \mathbf{a}, \mathbf{b} , and that $\frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} = 1$ when \mathbf{a}, \mathbf{b} are colinear
- Thus $\|\nabla f(\mathbf{w}_0)\|$ is the maximum slope, and it occurs in the direction of $\nabla f(\mathbf{w}_0)$



First-order approximations of matrix-input functions

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

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

- $\|\mathbf{A}\|_F^2 = \sum_i \sum_j a_{ij}^2$ 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 $f(\mathbf{W}) = \mathbf{a}^\top \mathbf{W} \mathbf{b} = \sum_i \sum_j a_i w_{ij} b_j$ for fixed vectors \mathbf{a} and \mathbf{b}
- The partial derivatives are $\frac{\partial f(\mathbf{W})}{\partial w_{ij}} = a_i b_j = [\nabla f(\mathbf{W})]_{ij}$ at any \mathbf{W}
- The gradient matrix is $\nabla f(\mathbf{W}) = \mathbf{a} \mathbf{b}^\top$ at any \mathbf{W}
- The linear approximation of $f(\mathbf{W})$ at \mathbf{W}_0 is

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

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

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

perfect approx since f is linear!

Example 3 in Python

- Function $f(\mathbf{W}) = \mathbf{a}^T \mathbf{W} \mathbf{b}$
 - In Python, use `.dot` for matrix multiplication
- Gradient $\nabla f(\mathbf{W}) = \mathbf{a} \mathbf{b}^T$
 - Want to set `fgrad[i,j]=a[i]b[j]`
 - But want to avoid for-loops
 - Use **Python broadcasting**
 - `a[:,None]` is $m \times 1$
 - `b[None,:]` is $1 \times n$

```
def feval(W,a,b):
    # Function
    f = a.dot(W.dot(b))

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

    return f, fgrad

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

f, fgrad = feval(W,a,b)
```

Outline

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

Stationary points

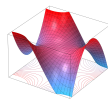
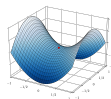
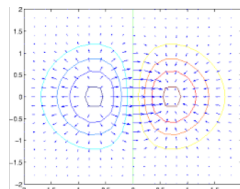
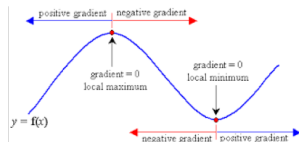
- A **stationary point** of J is any w such that

$$\nabla J(w) = 0$$
- The unconstrained minimizer of a smooth cost 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}



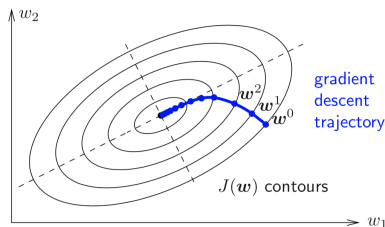
Gradient descent

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

- **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
- $\alpha_k > 0$ is the **stepsize** or **learning rate**
 - Often \mathbf{w}^0 is chosen randomly
 - Basically, we take downhill steps until we reach the bottom



Analysis of gradient descent

- The Taylor series of $J(\mathbf{w})$ 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

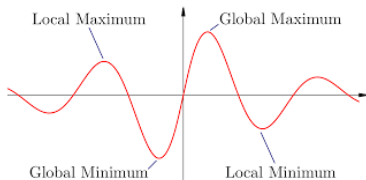
$$\begin{aligned} J(\mathbf{w}^{k+1}) &= J(\mathbf{w}^k) - \alpha_k \nabla J(\mathbf{w}^k)^\top \nabla J(\mathbf{w}^k) + O(\alpha_k^2 \|\nabla J(\mathbf{w}^k)\|^2) \\ &= J(\mathbf{w}^k) - \alpha_k \|\nabla J(\mathbf{w}^k)\|^2 + O(\alpha_k^2 \|\nabla J(\mathbf{w}^k)\|^2) \end{aligned}$$

- Thus, if α_k is sufficiently small, then $J(\mathbf{w}^k)$ will not increase
 - Why? Can always make α_k small enough so that $C\alpha_k^2 < \alpha_k$ for any $C > 0$, in which case the middle term will dominate the last term

Local vs. global minimizers

■ Definitions

- \hat{w} is a **global minimizer** if $J(\hat{w}) \leq J(w) \forall w$
- \hat{w} is a **local minimizer** if $J(\hat{w}) \leq J(w) \forall w$ in some open neighborhood of \hat{w}
- In most cases, gradient descent only guarantees convergence to a **local minimum**
- For a convex function, any **local minimum** is a **global minimum** (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 $J(\mathbf{w})$ in two stages:

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 $f_i(z_i) = \ln[1 + e^{z_i}] - y_i z_i$
 where $J(\mathbf{w}) = f(\mathbf{z}(\mathbf{w}))$

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

$$\frac{\partial f(\mathbf{z}(\mathbf{w}))}{\partial w_j} = \sum_i \frac{\partial f(\mathbf{z})}{\partial z_i} \frac{\partial z_i(\mathbf{w})}{\partial w_j}$$

■ Here, $\frac{\partial f(\mathbf{z})}{\partial z_i} = \frac{1}{1 + e^{z_i}} e^{z_i} - y_i = \frac{1}{e^{-z_i} + 1} - y_i$ and $\frac{\partial z_i(\mathbf{w})}{\partial w_j} = a_{ij}$

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 $f_i(z_i) = \ln[1 + e^{z_i}] - y_i z_i$
 $= -\ln\left(\frac{e^{z_i}}{1+e^{z_i}}\right) + (1 - y_i)z_i$
- Finally $J(\mathbf{w}) = f(\mathbf{z}) = \sum_i f_i(z_i)$

- **Backward pass:** Compute gradient:

- Recall $\frac{\partial J(\mathbf{w})}{\partial w_j} = \sum_i a_{ij} [\nabla f(\mathbf{z})]_i$

$$\text{where } [\nabla f(\mathbf{z})]_i = \frac{1}{1+e^{-z_i}} - y_i$$

- So $\nabla J(\mathbf{w}) = \begin{bmatrix} \frac{\partial J(\mathbf{w})}{\partial w_0} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_d} \end{bmatrix} = \mathbf{A}^T \nabla f(\mathbf{z})$

```
# Create a function with all the parameters
def feval_param(w,X,y):
    """
    Compute the loss and gradient 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))
    f = np.sum((1-y)*z - np.log(py))

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

Python approach #1: Use 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:

```
feval = @(w) feval_param(w,Xtr,ytr)
```

```
feval = lambda w: feval_param(w,Xtr,ytr)

# You can now pass a parameter like w0
f0, fgrad0 = feval(w0)
```

```
# Create a function with all the parameters
def feval_param(w,X,y):
    """
    Compute the loss and gradient 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))
    f = np.sum((1-y)*z - np.log(py))

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

Python approach #2: Create a “class”

Another approach is to create a “class” (i.e., object-oriented programming)

- Includes an **constructor** that

- loads data (X, y)
- does pre-computations
- called during **instantiation**:

```
log_fun = LogisticFun(Xtr,ytr)
```

- Includes an **feval** function to

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

```
# Call the function
f, fgrad = log_fun.feval(w0)
```

```
class LogisticFun(object):
    def __init__(self,X,y):
        """
        Class to compute loss & gradient for binary logistic regression

        The constructor takes the training features `X` and responses `y`
        """
        self.X = X
        self.y = y
        n = X.shape[0]
        self.A = np.column_stack((np.ones(n), X))

    def feval(self,w):
        """
        Compute the loss and gradient for a given weight vector `w`

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

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


Always check your gradient implementation!

- 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
p = X.shape[1]+1
w0 = np.random.randn(p)

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

# Measure the function and gradient at w0 and w1
f0, fgrad0 = log_fun.feval(w0)
f1, fgrad1 = log_fun.feval(w1)

# Predict the amount the function should have changed based on the gradient
df_est = fgrad0.dot(w1-w0)

# Print the two values to see if they are close
print("Actual f1-f0      = %12.4e" % (f1-f0))
print("Predicted f1-f0 = %12.4e" % df_est)

Actual f1-f0      =  -4.3270e-04
Predicted f1-f0 =  -4.3271e-04
```

A simple implementation of gradient descent

- We now implement gradient descent
- Inputs:
 - feval: function that computes cost & gradient
 - initial parameters w^0
 - learning rate α
 - number of iterations m
- Outputs:
 - final parameters w^m
 - history of cost & parameters (for debugging)

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

    feval:  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': [], 'f': []}

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

        # Evaluate the function and gradient
        f0, fgrad0 = feval(w0)

        # Take a gradient step
        w0 = w0 - lr*fgrad0

        # Save history
        hist['f'].append(f0)
        hist['w'].append(w0)

    # Convert to numpy arrays
    for elem in ('f', 'w'):
        hist[elem] = np.array(hist[elem])
    return w0, f0, hist
```

Gradient descent for logistic regression

- random initialization
- 1000 iterations
- convergence is slow!
- test accuracy is poor!
 - 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)

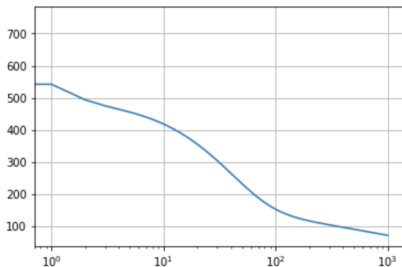
Test accuracy = 0.968198
```

```
# Initial condition
winit = np.random.randn(p)

# Parameters
feval = log_fun.feval
nit = 1000
lr = 1e-4

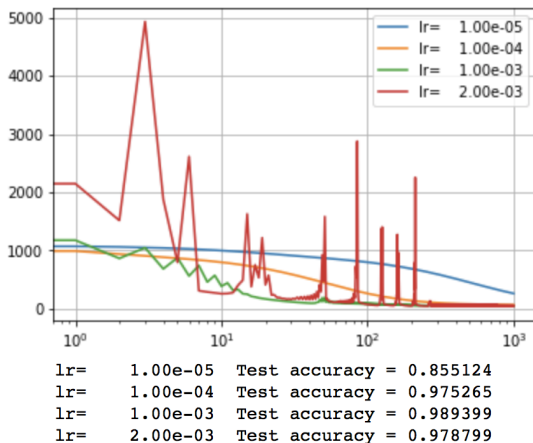
# Run the gradient descent
w, f0, hist = grad_opt_simp(feval, winit, lr=lr, nit=nit)

# Plot the training loss
t = np.arange(nit)
plt.semilogx(t, hist['f'])
plt.grid()
```



Effect of stepsize (or learning rate)

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



Outline

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

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:**

- Fix some $c \in (0, 1)$, usually $c = \frac{1}{2}$
- At each k , choose some stepsize $\alpha_k > 0$ such that

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

- Cost is guaranteed to decrease at each iteration (unless $\nabla J(\mathbf{w}^k) = \mathbf{0}$)
- Decreases by some fraction c of that predicted by linear approximation of $J(\mathbf{w}^{k+1})$

- A **simple Armijo-based α_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 $\beta_{\text{decr}} < 1$

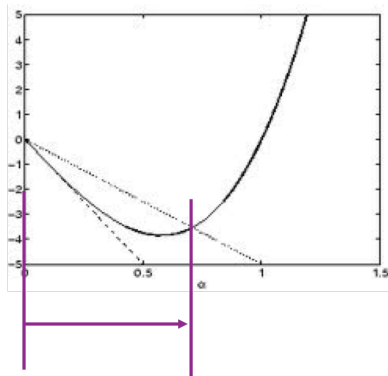
- Or use **line-search**, meaning test several α_k at each k and choose the best
 - More accurate but more expensive than simple Armijo

Armijo rule illustrated

- Recall Armijo rule: fix $c \in (0, 1)$ and accept any stepsize $\alpha_k > 0$ satisfying

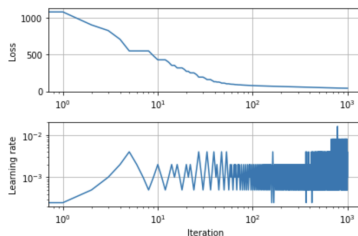
$$J(\mathbf{w}^{k+1}; \alpha_k) \leq \underbrace{J(\mathbf{w}^k) - c \alpha_k \|\nabla J(\mathbf{w}^k)\|^2}_{\triangleq y_c(\alpha_k)}$$

- Curve shows $J(\mathbf{w}^{k+1}; \alpha_k)$ versus α_k
 - line-search would give samples of this
- Dashed line shows $y_1(\alpha_k)$ versus α_k
 - this is the $J(\mathbf{w}^{k+1})$ predicted by linear approximation of $J(\mathbf{w}^k)$
- Dotted line shows $y_{\frac{1}{2}}(\alpha_k)$ versus α_k
 - purple bars show set of $\{\alpha_k\}$ satisfying the Armijo rule with $c = \frac{1}{2}$
 - what about other values of $c \in (0, 1)$?



Armijo example in Python

- The simple Armijo method applied to logistic regression



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

Test accuracy = 0.989399

```
for it in range(nit):

    # Take a gradient step
    w1 = w0 - lr*fgrad0

    # Evaluate the test point by computing the objective function, f1,
    # at the test point and the predicted decrease, df_est
    f1, fgrad1 = feval(w1)
    df_est = fgrad0.dot(w1-w0)

    # Check if test point passes the Armijo condition
    alpha = 0.5
    if (f1-f0 < alpha*df_est) and (f1 < f0):
        # If descent is sufficient, accept the point and increase the
        # learning rate
        lr = lr*2
        f0 = f1
        fgrad0 = fgrad1
        w0 = w1
    else:
        # Otherwise, decrease the learning rate
        lr = lr/2
```

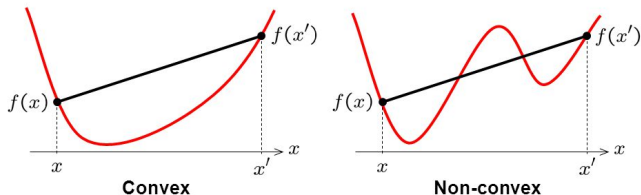
What values of β_{incr} , β_{decr} are being used?

Outline

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

Local minimizers of convex functions

- Theorem: if $f(w)$ is **convex** and w is a local minimizer, then w is a global minimizer
- 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



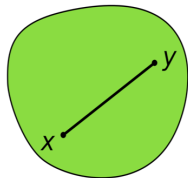
Convex sets

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

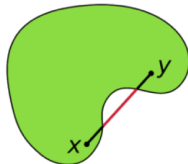
$$tx + (1 - t)y \in D$$

- The line between any two points in D remains in D
- Examples:
 - Circle, square, ellipse
 - \mathbb{R}^n
 - a hyperplane in \mathbb{R}^n
 - the half-space $\{x : Ax \leq b\}$ for any matrix A and vector b

a convex set



a nonconvex set



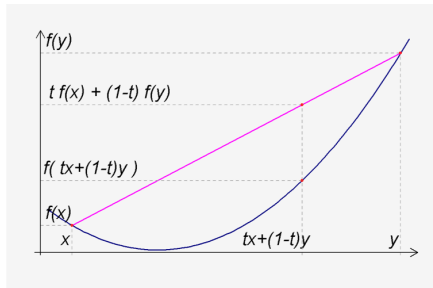
Convex functions

- A function $f(x) \in \mathbb{R}$ 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:

- $f(x) = ax + b$
- $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + b$
- $f(x) = ax^2 + bx + c$ iff $a \geq 0$
- norms are convex (e.g., $\|\mathbf{x}\|$, $\|\mathbf{x}\|_1$)
- $f(x)$ is convex if $f''(x)$ exists everywhere and $f''(x) \geq 0$
 - vector case: Hessian must exist everywhere and be positive semidefinite
- if f and g are convex, then so are $f + g$ and $f(g(\cdot))$
- RSS, logistic loss, and their L1 or L2 regularized versions are all convex



Optimization topics that we did not cover

- Our Armijo-based optimizer is okay, but not nearly as fast as methods in `sklearn`
- There are many topics that we did not cover, e.g.,
 - Newton's method and quasi-Newton methods (i.e., matrix-valued α_k)
 - non-smooth optimization (i.e., gradient does not exist everywhere)
 - constrained optimization
- Take an optimization class and learn more!
 - ECE-5759 (Autumn)

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**
- Understand the implications of **convexity** for gradient descent
- Determine if a loss function is convex