

Gradient Methods

Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no¹

Department of Physics and Center for Computing in Science Education, University
of Oslo, Oslo, Norway¹

March 7, 2025

© 1999-2025, Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no. Released under CC

Attribution-NonCommercial 4.0 license

Overview of week March 3-7

Gradient methods:

1. Semi-Newton methods (Broyden's algorithm and Broyden-Farberg-Goldberg-Shanno algorithm)
2. Steepest descent and conjugate gradient descent
3. Stochastic gradient descent and variants thereof
4. Automatic differentiation

Teaching Material, videos and written material

1. These lecture notes
2. Video on the Conjugate Gradient methods
3. Recommended background literature, [Convex Optimization](#) by Boyd and Vandenberghe. Their [lecture slides](#) are very useful (warning, these are some 300 pages).

Brief reminder on Newton-Raphson's method

Let us quickly remind ourselves how we derive the above method. Perhaps the most celebrated of all one-dimensional root-finding routines is Newton's method, also called the Newton-Raphson method. This method requires the evaluation of both the function f and its derivative f' at arbitrary points. If you can only calculate the derivative numerically and/or your function is not of the smooth type, we normally discourage the use of this method.

The equations

The Newton-Raphson formula consists geometrically of extending the tangent line at a current point until it crosses zero, then setting the next guess to the abscissa of that zero-crossing. The mathematics behind this method is rather simple. Employing a Taylor expansion for x sufficiently close to the solution s , we have

$$f(s) = 0 = f(x) + (s - x)f'(x) + \frac{(s - x)^2}{2}f''(x) + \dots$$

Small values

For small enough values of the function and for well-behaved functions, the terms beyond linear are unimportant, hence we obtain

$$f(x) + (s - x)f'(x) \approx 0,$$

yielding

$$s \approx x - \frac{f(x)}{f'(x)}.$$

Having in mind an iterative procedure, it is natural to start iterating with

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

Simple geometric interpretation

The above is Newton-Raphson's method. It has a simple geometric interpretation, namely x_{n+1} is the point where the tangent from $(x_n, f(x_n))$ crosses the x -axis. Close to the solution, Newton-Raphson converges fast to the desired result. However, if we are far from a root, where the higher-order terms in the series are important, the Newton-Raphson formula can give grossly inaccurate results. For instance, the initial guess for the root might be so far from the true root as to let the search interval include a local maximum or minimum of the function. If an iteration places a trial guess near such a local extremum, so that the first derivative nearly vanishes, then Newton-Raphson may fail totally

Extending to more than one variable

Newton's method can be generalized to systems of several non-linear equations and variables. Consider the case with two equations

$$\begin{aligned}f_1(x_1, x_2) &= 0 \\ f_2(x_1, x_2) &= 0,\end{aligned}$$

which we Taylor expand to obtain

$$\begin{aligned}0 = f_1(x_1 + h_1, x_2 + h_2) &= f_1(x_1, x_2) + h_1 \partial f_1 / \partial x_1 + h_2 \partial f_1 / \partial x_2 + \dots \\ 0 = f_2(x_1 + h_1, x_2 + h_2) &= f_2(x_1, x_2) + h_1 \partial f_2 / \partial x_1 + h_2 \partial f_2 / \partial x_2 + \dots\end{aligned}$$

Jacobian

Defining the Jacobian matrix \hat{J} we have

$$\hat{J} = \begin{pmatrix} \partial f_1 / \partial x_1 & \partial f_1 / \partial x_2 \\ \partial f_2 / \partial x_1 & \partial f_2 / \partial x_2 \end{pmatrix},$$

we can rephrase Newton's method as

$$\begin{pmatrix} x_1^{n+1} \\ x_2^{n+1} \end{pmatrix} = \begin{pmatrix} x_1^n \\ x_2^n \end{pmatrix} + \begin{pmatrix} h_1^n \\ h_2^n \end{pmatrix},$$

where we have defined

$$\begin{pmatrix} h_1^n \\ h_2^n \end{pmatrix} = -\hat{J}^{-1} \begin{pmatrix} f_1(x_1^n, x_2^n) \\ f_2(x_1^n, x_2^n) \end{pmatrix}.$$

Inverse of the Jacobian

We need thus to compute the inverse of the Jacobian matrix and it is to understand that difficulties may arise in case \hat{J} is nearly singular.

It is rather straightforward to extend the above scheme to systems of more than two non-linear equations. In our case, the Jacobian matrix is given by the Hessian that represents the second derivative of the energy.

Steepest descent

The basic idea of gradient descent is that a function $F(\mathbf{x})$, $\mathbf{x} \equiv (x_1, \dots, x_n)$, decreases fastest if one goes from \mathbf{x} in the direction of the negative gradient $-\nabla F(\mathbf{x})$.

It can be shown that if

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma_k \nabla F(\mathbf{x}_k),$$

with $\gamma_k > 0$.

For γ_k small enough, then $F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k)$. This means that for a sufficiently small γ_k we are always moving towards smaller function values, i.e a minimum.

More on Steepest descent

The previous observation is the basis of the method of steepest descent, which is also referred to as just gradient descent (GD). One starts with an initial guess x_0 for a minimum of F and computes new approximations according to

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k), \quad k \geq 0.$$

The parameter γ_k is often referred to as the step length or the learning rate within the context of Machine Learning.

The ideal

Ideally the sequence $\{x_k\}_{k=0}$ converges to a global minimum of the function F . In general we do not know if we are in a global or local minimum. In the special case when F is a convex function, all local minima are also global minima, so in this case gradient descent can converge to the global solution. The advantage of this scheme is that it is conceptually simple and straightforward to implement. However the method in this form has some severe limitations: In machine learning we are often faced with non-convex high dimensional cost functions with many local minima. Since GD is deterministic we will get stuck in a local minimum, if the method converges, unless we have a very good initial guess. This also implies that the scheme is sensitive to the chosen initial condition. Note that the gradient is a function of $x = (x_1, \dots, x_n)$ which makes it expensive to compute numerically.

The sensitiveness of the gradient descent

The gradient descent method is sensitive to the choice of learning rate γ_k . This is due to the fact that we are only guaranteed that $F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_k)$ for sufficiently small γ_k . The problem is to determine an optimal learning rate. If the learning rate is chosen too small the method will take a long time to converge and if it is too large we can experience erratic behavior.

Many of these shortcomings can be alleviated by introducing randomness. One such method is that of Stochastic Gradient Descent (SGD), see below.

Convex function

Convex function: Let $X \subset \mathbb{R}^n$ be a convex set. Assume that the function $f : X \rightarrow \mathbb{R}$ is continuous, then f is said to be convex if

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

for all $x_1, x_2 \in X$ and for all $t \in [0, 1]$. If \leq is replaced with a strict inequality in the definition, we demand $x_1 \neq x_2$ and $t \in (0, 1)$ then f is said to be strictly convex. For a single variable function, convexity means that if you draw a straight line connecting $f(x_1)$ and $f(x_2)$, the value of the function on the interval $[x_1, x_2]$ is always below the line as illustrated below.

Conditions on convex functions

In the following we state first and second-order conditions which ensures convexity of a function f . We write D_f to denote the domain of f , i.e the subset of R^n where f is defined. For more details and proofs we refer to: [S. Boyd and L. Vandenberghe. Convex Optimization. Cambridge University Press.](#)

First order condition

Suppose f is differentiable (i.e $\nabla f(x)$ is well defined for all x in the domain of f). Then f is convex if and only if D_f is a convex set and

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

holds for all $x, y \in D_f$. This condition means that for a convex function the first order Taylor expansion (right hand side above) at any point a global under estimator of the function. To convince yourself you can make a drawing of $f(x) = x^2 + 1$ and draw the tangent line to $f(x)$ and note that it is always below the graph.

Second condition

Second order condition

Assume that f is twice differentiable, i.e the Hessian matrix exists at each point in D_f . Then f is convex if and only if D_f is a convex set and its Hessian is positive semi-definite for all $x \in D_f$. For a single-variable function this reduces to $f''(x) \geq 0$. Geometrically this means that f has nonnegative curvature everywhere.

This condition is particularly useful since it gives us an procedure for determining if the function under consideration is convex, apart from using the definition.

More on convex functions

The next result is of great importance to us and the reason why we are going on about convex functions. In machine learning we frequently have to minimize a loss/cost function in order to find the best parameters for the model we are considering.

Ideally we want the global minimum (for high-dimensional models it is hard to know if we have local or global minimum). However, if the cost/loss function is convex the following result provides invaluable information:

Any minimum is global for convex functions

Consider the problem of finding $x \in \mathbb{R}^n$ such that $f(x)$ is minimal, where f is convex and differentiable. Then, any point x^* that satisfies $\nabla f(x^*) = 0$ is a global minimum.

This result means that if we know that the cost/loss function is convex and we are able to find a minimum, we are guaranteed that it is a global minimum.

Some simple problems

1. Show that $f(x) = x^2$ is convex for $x \in \mathbb{R}$ using the definition of convexity. Hint: If you re-write the definition, f is convex if the following holds for all $x, y \in D_f$ and any $\lambda \in [0, 1]$
$$\lambda f(x) + (1 - \lambda)f(y) - f(\lambda x + (1 - \lambda)y) \geq 0.$$
2. Using the second order condition show that the following functions are convex on the specified domain.
 - ▶ $f(x) = e^x$ is convex for $x \in \mathbb{R}$.
 - ▶ $g(x) = -\ln(x)$ is convex for $x \in (0, \infty)$.
3. Let $f(x) = x^2$ and $g(x) = e^x$. Show that $f(g(x))$ and $g(f(x))$ is convex for $x \in \mathbb{R}$. Also show that if $f(x)$ is any convex function then $h(x) = e^{f(x)}$ is convex.
4. A norm is any function that satisfy the following properties
 - ▶ $f(\alpha x) = |\alpha|f(x)$ for all $\alpha \in \mathbb{R}$.
 - ▶ $f(x + y) \leq f(x) + f(y)$
 - ▶ $f(x) \leq 0$ for all $x \in \mathbb{R}^n$ with equality if and only if $x = 0$

Using the definition of convexity, try to show that a function satisfying the properties above is convex.

Broyden's Algorithm for Solving Nonlinear Equations

Broyden's algorithm is a quasi-Newton method used to solve systems of nonlinear equations. Unlike Newton's method, which requires the computation of the Jacobian matrix at each iteration, Broyden's method approximates the Jacobian (or its inverse) to reduce computational cost. This makes it particularly useful for high-dimensional problems where computing the exact Jacobian is expensive.

Problem Formulation

Consider a system of n nonlinear equations:

$$F(x) = 0,$$

where $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a vector-valued function, and $x \in \mathbb{R}^n$ is the vector of unknowns. The goal is to find x^* such that $F(x^*) = 0$.

Just a short reminder of Newton's Method

Newton's method iteratively updates the solution as:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}_k^{-1} \mathbf{F}(\mathbf{x}_k),$$

where $\mathbf{J}_k = \mathbf{J}(\mathbf{x}_k)$ is the Jacobian matrix of \mathbf{F} evaluated at \mathbf{x}_k . However, computing \mathbf{J}_k and its inverse at each iteration can be computationally expensive.

Broyden's Method

Broyden's method approximates the Jacobian (or its inverse) to avoid recomputing it at every iteration. There are two variants of Broyden's method:

1. Broyden's Good Method: Updates an approximation of the Jacobian matrix.
2. Broyden's Bad Method: Updates an approximation of the inverse Jacobian matrix.

Broyden's Good Method

Let B_k be the approximation of the Jacobian J_k at iteration k . The update rule for B_k is:

$$B_{k+1} = B_k + \frac{(F(x_{k+1}) - F(x_k) - B_k s_k) s_k^\top}{s_k^\top s_k},$$

where $s_k = x_{k+1} - x_k$ is the step vector. The new approximation B_{k+1} satisfies the **secant equation:

$$B_{k+1} s_k = F(x_{k+1}) - F(x_k).$$

The solution is updated as:

$$x_{k+1} = x_k - B_k^{-1} F(x_k).$$

Broyden's Bad Method

Let H_k be the approximation of the inverse Jacobian J_k^{-1} at iteration k . The update rule for H_k is:

$$H_{k+1} = H_k + \frac{(s_k - H_k y_k) s_k^\top H_k}{s_k^\top H_k y_k},$$

where $y_k = F(x_{k+1}) - F(x_k)$. The new approximation H_{k+1} satisfies the ****inverse secant equation****:

$$H_{k+1} y_k = s_k.$$

The solution is updated as:

$$x_{k+1} = x_k - H_k F(x_k).$$

Algorithm Steps

The steps of Broyden's algorithm (Good Method) are as follows:

1. Initialize x_0 and B_0 (e.g., $B_0 = I$, the identity matrix).
2. For $k = 0, 1, 2, \dots$:
 - ▶ Compute the step: $s_k = -B_k^{-1}F(x_k)$.
 - ▶ Update the solution: $x_{k+1} = x_k + s_k$.
 - ▶ Compute $y_k = F(x_{k+1}) - F(x_k)$.
 - ▶ Update the Jacobian approximation:

$$B_{k+1} = B_k + \frac{(y_k - B_k s_k)s_k^\top}{s_k^\top s_k}.$$

- ▶ Check for convergence: If $\|F(x_{k+1})\| < \epsilon$, stop.

Advantages and Limitations

Advantages

1. Avoids the need to compute the exact Jacobian at each iteration.
2. Suitable for high-dimensional problems where computing the Jacobian is expensive.

Limitations

1. Convergence is not guaranteed for all problems.
- oThe approximation of the Jacobian may become inaccurate over time, requiring periodic reinitialization.

Applications

Broyden's algorithm is widely used in:

1. Optimization problems.
2. Solving systems of nonlinear equations in scientific computing.
3. Machine learning for training certain types of models.

Broyden-Fletcher-Goldfarb-Shanno algorithm

The Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm is an iterative method for solving unconstrained nonlinear optimization problems. It belongs to the family of quasi-Newton methods, which aim to approximate the inverse Hessian matrix used in Newton's method for optimization. BFGS is widely used due to its efficiency and the fact that it does not require the computation of second derivatives.

Given an objective function $f(x)$, the goal is to minimize $f(x)$ with respect to the vector x . The BFGS method iteratively updates an estimate of the inverse Hessian matrix and a search direction to find the minimum of $f(x)$.

BFGS optimization problem

The optimization problem is to minimize $f(\mathbf{x})$ where \mathbf{x} is a vector in \mathbb{R}^n , and f is a differentiable scalar function. There are no constraints on the values that \mathbf{x} can take.

The algorithm begins at an initial estimate for the optimal value \mathbf{x}_0 and proceeds iteratively to get a better estimate at each stage.

The search direction \mathbf{p}_k at stage k is given by the solution of the analogue of the Newton equation

$$B_k \mathbf{p}_k = -\nabla f(\mathbf{x}_k),$$

where B_k is an approximation to the Hessian matrix, which is updated iteratively at each stage, and $\nabla f(\mathbf{x}_k)$ is the gradient of the function evaluated at \mathbf{x}_k . A line search in the direction \mathbf{p}_k is then used to find the next point \mathbf{x}_{k+1} by minimising

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k),$$

over the scalar $\alpha > 0$.

BFGS optimization problem, setting up the equations

We are given the following unconstrained optimization problem:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable objective function, and $\mathbf{x} \in \mathbb{R}^n$ is the vector of decision variables.

The first-order necessary conditions for optimality are given by:

$$\nabla f(\mathbf{x}^*) = 0$$

where $\nabla f(\mathbf{x})$ denotes the gradient of $f(\mathbf{x})$.

BFGS Algorithm Overview

The BFGS method is an iterative procedure that approximates the inverse Hessian matrix H_k at each iteration. The update of the current solution x_k involves computing a search direction and step length. The general steps of the BFGS algorithm are as follows:

- ▶ Initialize x_0 and choose an initial guess for the inverse Hessian approximation H_0 , typically $H_0 = I$ (the identity matrix).
- ▶ For each iteration k , do the following:
 - ▶ Compute the gradient at the current point: $\nabla f(x_k)$.
 - ▶ Compute the search direction:

$$p_k = -H_k \nabla f(x_k)$$

- ▶ Perform a line search to find an appropriate step size α_k , which minimizes $f(x_k + \alpha_k p_k)$.
- ▶ Update the current point:

$$x_{k+1} = x_k + \alpha_k p_k$$

- ▶ Compute the new gradient $\nabla f(x_{k+1})$.

Final steps

- Update the inverse Hessian approximation using the following formula:

$$s_k = x_{k+1} - x_k, \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

$$H_{k+1} = H_k + \frac{s_k s_k^T}{s_k^T y_k} - \frac{H_k y_k y_k^T H_k}{y_k^T H_k y_k}$$

- Repeat the process until convergence, i.e., $\|\nabla f(x_k)\|$ is sufficiently small.

Convergence and Termination Criteria

The BFGS algorithm is guaranteed to converge to a local minimum under certain conditions, such as the objective function being smooth and convex. The algorithm terminates when the gradient $\nabla f(\mathbf{x}_k)$ is sufficiently close to zero, or when a maximum number of iterations is reached.

A typical convergence criterion is:

$$\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$$

where ϵ is a small tolerance value.

Properties of BFGS

The BFGS algorithm has several key properties that make it widely used:

1. **No need for second derivatives:** The BFGS method approximates the Hessian matrix and avoids the direct computation of second derivatives.
2. **Superlinear convergence:** The BFGS method typically converges faster than gradient descent methods, especially near the optimum.
3. **Memory efficiency:** The BFGS algorithm maintains a low-rank approximation of the inverse Hessian, making it computationally efficient.
4. **Positive definiteness:** If the initial Hessian approximation is positive definite, the BFGS update preserves this property.

Final words on the BFGS

The Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm is an effective and widely used method for solving unconstrained optimization problems. By iteratively approximating the inverse Hessian matrix, BFGS achieves efficient convergence without requiring the explicit computation of second derivatives. Its performance is highly dependent on the choice of the initial guess and the line search strategy.

Standard steepest descent

Before we proceed, we would like to discuss the approach called the **standard Steepest descent**, which again leads to us having to be able to compute a matrix. It belongs to the class of Conjugate Gradient methods (CG).

The success of the CG method for finding solutions of non-linear problems is based on the theory of conjugate gradients for linear systems of equations. It belongs to the class of iterative methods for solving problems from linear algebra of the type

$$\hat{A}\hat{x} = \hat{b}.$$

In the iterative process we end up with a problem like

$$\hat{r} = \hat{b} - \hat{A}\hat{x},$$

where \hat{r} is the so-called residual or error in the iterative process. When we have found the exact solution, $\hat{r} = 0$.

Gradient method

The residual is zero when we reach the minimum of the quadratic equation

$$P(\hat{x}) = \frac{1}{2} \hat{x}^T \hat{A} \hat{x} - \hat{x}^T \hat{b},$$

with the constraint that the matrix \hat{A} is positive definite and symmetric. This defines also the Hessian and we want it to be positive definite.

Steepest descent method

We denote the initial guess for \hat{x} as \hat{x}_0 . We can assume without loss of generality that

$$\hat{x}_0 = 0,$$

or consider the system

$$\hat{A}\hat{z} = \hat{b} - \hat{A}\hat{x}_0,$$

instead.

Steepest descent method

One can show that the solution \hat{x} is also the unique minimizer of the quadratic form

$$f(\hat{x}) = \frac{1}{2} \hat{x}^T \hat{A} \hat{x} - \hat{x}^T \hat{b}, \quad \hat{x} \in \mathbb{R}^n.$$

This suggests taking the first basis vector \hat{r}_1 (see below for definition) to be the gradient of f at $\hat{x} = \hat{x}_0$, which equals

$$\hat{A} \hat{x}_0 - \hat{b},$$

and $\hat{x}_0 = 0$ it is equal $-\hat{b}$.

Final expressions

We can compute the residual iteratively as

$$\hat{r}_{k+1} = \hat{b} - \hat{A}\hat{x}_{k+1},$$

which equals

$$\hat{b} - \hat{A}(\hat{x}_k + \alpha_k \hat{r}_k),$$

or

$$(\hat{b} - \hat{A}\hat{x}_k) - \alpha_k \hat{A}\hat{r}_k,$$

which gives

$$\alpha_k = \frac{\hat{r}_k^T \hat{r}_k}{\hat{r}_k^T \hat{A} \hat{r}_k}$$

leading to the iterative scheme

$$\hat{x}_{k+1} = \hat{x}_k - \alpha_k \hat{r}_k,$$

Conjugate gradient method

In the CG method we define so-called conjugate directions and two vectors \hat{s} and \hat{t} are said to be conjugate if

$$\hat{s}^T \hat{A} \hat{t} = 0.$$

The philosophy of the CG method is to perform searches in various conjugate directions of our vectors \hat{x}_i obeying the above criterion, namely

$$\hat{x}_i^T \hat{A} \hat{x}_j = 0.$$

Two vectors are conjugate if they are orthogonal with respect to this inner product. Being conjugate is a symmetric relation: if \hat{s} is conjugate to \hat{t} , then \hat{t} is conjugate to \hat{s} .

Conjugate gradient method

An example is given by the eigenvectors of the matrix

$$\hat{v}_i^T \hat{A} \hat{v}_j = \lambda \hat{v}_i^T \hat{v}_j,$$

which is zero unless $i = j$.

Conjugate gradient method

Assume now that we have a symmetric positive-definite matrix \hat{A} of size $n \times n$. At each iteration $i + 1$ we obtain the conjugate direction of a vector

$$\hat{x}_{i+1} = \hat{x}_i + \alpha_i \hat{p}_i.$$

We assume that \hat{p}_i is a sequence of n mutually conjugate directions. Then the \hat{p}_i form a basis of R^n and we can expand the solution $\hat{A}\hat{x} = \hat{b}$ in this basis, namely

$$\hat{x} = \sum_{i=1}^n \alpha_i \hat{p}_i.$$

Conjugate gradient method

The coefficients are given by

$$A\mathbf{x} = \sum_{i=1}^n \alpha_i A\mathbf{p}_i = \mathbf{b}.$$

Multiplying with $\hat{\mathbf{p}}_k^T$ from the left gives

$$\hat{\mathbf{p}}_k^T A\hat{\mathbf{x}} = \sum_{i=1}^n \alpha_i \hat{\mathbf{p}}_k^T A\hat{\mathbf{p}}_i = \hat{\mathbf{p}}_k^T \hat{\mathbf{b}},$$

and we can define the coefficients α_k as

$$\alpha_k = \frac{\hat{\mathbf{p}}_k^T \hat{\mathbf{b}}}{\hat{\mathbf{p}}_k^T A\hat{\mathbf{p}}_k}$$

Conjugate gradient method and iterations

If we choose the conjugate vectors \hat{p}_k carefully, then we may not need all of them to obtain a good approximation to the solution \hat{x} . We want to regard the conjugate gradient method as an iterative method. This will us to solve systems where n is so large that the direct method would take too much time.

We denote the initial guess for \hat{x} as \hat{x}_0 . We can assume without loss of generality that

$$\hat{x}_0 = 0,$$

or consider the system

$$\hat{A}\hat{z} = \hat{b} - \hat{A}\hat{x}_0,$$

instead.

Conjugate gradient method

One can show that the solution \hat{x} is also the unique minimizer of the quadratic form

$$f(\hat{x}) = \frac{1}{2} \hat{x}^T \hat{A} \hat{x} - \hat{x}^T \hat{b}, \quad \hat{x} \in \mathbb{R}^n.$$

This suggests taking the first basis vector \hat{p}_1 to be the gradient of f at $\hat{x} = \hat{x}_0$, which equals

$$\hat{A} \hat{x}_0 - \hat{b},$$

and $\hat{x}_0 = 0$ it is equal $-\hat{b}$. The other vectors in the basis will be conjugate to the gradient, hence the name conjugate gradient method.

Conjugate gradient method

Let \hat{r}_k be the residual at the k -th step:

$$\hat{r}_k = \hat{b} - \hat{A}\hat{x}_k.$$

Note that \hat{r}_k is the negative gradient of f at $\hat{x} = \hat{x}_k$, so the gradient descent method would be to move in the direction \hat{r}_k . Here, we insist that the directions \hat{p}_k are conjugate to each other, so we take the direction closest to the gradient \hat{r}_k under the conjugacy constraint. This gives the following expression

$$\hat{p}_{k+1} = \hat{r}_k - \frac{\hat{p}_k^T \hat{A} \hat{r}_k}{\hat{p}_k^T \hat{A} \hat{p}_k} \hat{p}_k.$$

Conjugate gradient method

We can also compute the residual iteratively as

$$\hat{r}_{k+1} = \hat{b} - \hat{A}\hat{x}_{k+1},$$

which equals

$$\hat{b} - \hat{A}(\hat{x}_k + \alpha_k \hat{p}_k),$$

or

$$(\hat{b} - \hat{A}\hat{x}_k) - \alpha_k \hat{A}\hat{p}_k,$$

which gives

$$\hat{r}_{k+1} = \hat{r}_k - \hat{A}\hat{p}_k,$$

Using gradient descent methods, limitations

- ▶ **Gradient descent (GD) finds local minima of our function.** Since the GD algorithm is deterministic, if it converges, it will converge to a local minimum of our energy function. Because in ML we are often dealing with extremely rugged landscapes with many local minima, this can lead to poor performance.
- ▶ **GD is sensitive to initial conditions.** One consequence of the local nature of GD is that initial conditions matter. Depending on where one starts, one will end up at a different local minima. Therefore, it is very important to think about how one initializes the training process. This is true for GD as well as more complicated variants of GD.
- ▶ **Gradients are computationally expensive to calculate for large datasets.** In many cases the energy function is a sum of terms, with one term for each data point. For example, in linear regression, $E \propto \sum_{i=1}^n (y_i - \mathbf{w}^T \cdot \mathbf{x}_i)^2$; for logistic regression, the square error is replaced by the cross entropy. To calculate the gradient we have to sum over *all* n data points. Doing this at every GD step becomes extremely

Improving gradient descent with momentum

We discuss here some simple examples where we introduce what is called 'memory' about previous steps, or what is normally called momentum gradient descent. The mathematics is explained below in connection with Stochastic gradient descent.

```
from numpy import asarray
from numpy import arange
from numpy.random import rand
from numpy.random import seed
from matplotlib import pyplot

# objective function
def objective(x):
    return x**2.0

# derivative of objective function
def derivative(x):
    return x * 2.0

# gradient descent algorithm
def gradient_descent(objective, derivative, bounds, n_iter, step_size):
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = bounds[:, 0] + rand(len(bounds)) * (bounds[:, 1] -
    # run the gradient descent
    for i in range(n_iter):
```

Same code but now with momentum gradient descent

```
from numpy import asarray
from numpy import arange
from numpy.random import rand
from numpy.random import seed
from matplotlib import pyplot

# objective function
def objective(x):
    return x**2.0

# derivative of objective function
def derivative(x):
    return x * 2.0

# gradient descent algorithm
def gradient_descent(objective, derivative, bounds, n_iter, step_size,
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = bounds[:, 0] + rand(len(bounds)) * (bounds[:, 1] -
    # keep track of the change
    change = 0.0
    # run the gradient descent
    for i in range(n_iter):
        # calculate gradient
        gradient = derivative(solution)
        # calculate update
        new_change = step_size * gradient + momentum * change
```

Overview video on Stochastic Gradient Descent

What is Stochastic Gradient Descent

Batches and mini-batches

In gradient descent we compute the cost function and its gradient for all data points we have.

In large-scale applications such as the [ILSVRC challenge](#), the training data can have on order of millions of examples. Hence, it seems wasteful to compute the full cost function over the entire training set in order to perform only a single parameter update. A very common approach to addressing this challenge is to compute the gradient over batches of the training data. For example, a typical batch could contain some thousand examples from an entire training set of several millions. This batch is then used to perform a parameter update.

Stochastic Gradient Descent (SGD)

In stochastic gradient descent, the extreme case is the case where we have only one batch, that is we include the whole data set.

This process is called Stochastic Gradient Descent (SGD) (or also sometimes on-line gradient descent). This is relatively less common to see because in practice due to vectorized code optimizations it can be computationally much more efficient to evaluate the gradient for 100 examples, than the gradient for one example 100 times. Even though SGD technically refers to using a single example at a time to evaluate the gradient, you will hear people use the term SGD even when referring to mini-batch gradient descent (i.e. mentions of MGD for “Minibatch Gradient Descent”, or BGD for “Batch gradient descent” are rare to see), where it is usually assumed that mini-batches are used. The size of the mini-batch is a hyperparameter but it is not very common to cross-validate or bootstrap it. It is usually based on memory constraints (if any), or set to some value, e.g. 32, 64 or 128. We use powers of 2 in practice because many vectorized operation implementations work faster when their inputs are sized in powers of 2.

In our notes with SGD we mean stochastic gradient descent with

Stochastic Gradient Descent

Stochastic gradient descent (SGD) and variants thereof address some of the shortcomings of the Gradient descent method discussed above.

The underlying idea of SGD comes from the observation that the cost function, which we want to minimize, can almost always be written as a sum over n data points $\{\mathbf{x}_i\}_{i=1}^n$,

$$C(\beta) = \sum_{i=1}^n c_i(\mathbf{x}_i, \beta).$$

Computation of gradients

This in turn means that the gradient can be computed as a sum over i -gradients

$$\nabla_{\beta} C(\beta) = \sum_i^n \nabla_{\beta} c_i(\mathbf{x}_i, \beta).$$

Stochasticity/randomness is introduced by only taking the gradient on a subset of the data called minibatches. If there are n data points and the size of each minibatch is M , there will be n/M minibatches. We denote these minibatches by B_k where $k = 1, \dots, n/M$.

SGD example

As an example, suppose we have 10 data points (x_1, \dots, x_{10}) and we choose to have $M = 5$ minibatches, then each minibatch contains two data points. In particular we have

$B_1 = (x_1, x_2), \dots, B_5 = (x_9, x_{10})$. Note that if you choose $M = 1$ you have only a single batch with all data points and on the other extreme, you may choose $M = n$ resulting in a minibatch for each datapoint, i.e $B_k = x_k$.

The idea is now to approximate the gradient by replacing the sum over all data points with a sum over the data points in one the minibatches picked at random in each gradient descent step

$$\nabla_{\beta} C(\beta) = \sum_{i=1}^n \nabla_{\beta} c_i(x_i, \beta) \rightarrow \sum_{i \in B_k}^n \nabla_{\beta} c_i(x_i, \beta).$$

The gradient step

Thus a gradient descent step now looks like

$$\beta_{j+1} = \beta_j - \gamma_j \sum_{i \in B_k}^n \nabla_{\beta} c_i(\mathbf{x}_i, \beta)$$

where k is picked at random with equal probability from $[1, n/M]$. An iteration over the number of minibatches (n/M) is commonly referred to as an epoch. Thus it is typical to choose a number of epochs and for each epoch iterate over the number of minibatches, as exemplified in the code below.

Simple example code

```
import numpy as np

n = 100 #100 datapoints
M = 5 #size of each minibatch
m = int(n/M) #number of minibatches
n_epochs = 10 #number of epochs

j = 0
for epoch in range(1,n_epochs+1):
    for i in range(m):
        k = np.random.randint(m) #Pick the k-th minibatch at random
        #Compute the gradient using the data in minibatch Bk
        #Compute new suggestion for
        j += 1
```

Taking the gradient only on a subset of the data has two important benefits. First, it introduces randomness which decreases the chance that our optimization scheme gets stuck in a local minima. Second, if the size of the minibatches are small relative to the number of datapoints ($M < n$), the computation of the gradient is much cheaper since we sum over the datapoints in the k -th minibatch and not all n datapoints.

When do we stop?

A natural question is when do we stop the search for a new minimum? One possibility is to compute the full gradient after a given number of epochs and check if the norm of the gradient is smaller than some threshold and stop if true. However, the condition that the gradient is zero is valid also for local minima, so this would only tell us that we are close to a local/global minimum. However, we could also evaluate the cost function at this point, store the result and continue the search. If the test kicks in at a later stage we can compare the values of the cost function and keep the β that gave the lowest value.

Slightly different approach

Another approach is to let the step length γ_j depend on the number of epochs in such a way that it becomes very small after a reasonable time such that we do not move at all. Such approaches are also called scaling. There are many such ways to [scale the learning rate](#) and [discussions here](#). See also <https://towardsdatascience.com/learning-rate-schedules-and-adaptive-learning-rate-methods-> for a discussion of different scaling functions for the learning rate.

Time decay rate

As an example, let $e = 0, 1, 2, 3, \dots$ denote the current epoch and let $t_0, t_1 > 0$ be two fixed numbers. Furthermore, let $t = e \cdot m + i$ where m is the number of minibatches and $i = 0, \dots, m - 1$. Then the function

$$\gamma_j(t; t_0, t_1) = \frac{t_0}{t + t_1}$$

goes to zero as the number of epochs gets large. I.e. we start with a step length $\gamma_j(0; t_0, t_1) = t_0/t_1$ which decays in *time* t .

In this way we can fix the number of epochs, compute β and evaluate the cost function at the end. Repeating the computation will give a different result since the scheme is random by design. Then we pick the final β that gives the lowest value of the cost function.

```
import numpy as np

def step_length(t,t0,t1):
    return t0/(t+t1)

n = 100 #100 datapoints
M = 5   #size of each minibatch
m = int(n/M) #number of minibatches
n_epochs = 500 #number of epochs
```

Code with a Number of Minibatches which varies

In the code here we vary the number of mini-batches.

```
# Importing various packages
from math import exp, sqrt
from random import random, seed
import numpy as np
import matplotlib.pyplot as plt

n = 100
x = 2*np.random.rand(n,1)
y = 4+3*x+np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x]
XT_X = X.T @ X
theta_linreg = np.linalg.inv(X.T @ X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)
# Hessian matrix
H = (2.0/n)* XT_X
EigValues, EigVectors = np.linalg.eig(H)
print(f"Eigenvalues of Hessian Matrix:{EigValues}")

theta = np.random.randn(2,1)
eta = 1.0/np.max(EigValues)
Niterations = 1000

for iter in range(Niterations):
    gradients = 2.0/n*X.T @ ((X @ theta)-y)
    theta -= eta*gradients
```

Replace or not

In the above code, we have use replacement in setting up the mini-batches. The discussion [here](#) may be useful.

Momentum based GD

The stochastic gradient descent (SGD) is almost always used with a *momentum* or inertia term that serves as a memory of the direction we are moving in parameter space. This is typically implemented as follows

$$\begin{aligned} \mathbf{v}_t &= \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} E(\boldsymbol{\theta}_t) \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \mathbf{v}_t, \end{aligned} \tag{1}$$

where we have introduced a momentum parameter γ , with $0 \leq \gamma \leq 1$, and for brevity we dropped the explicit notation to indicate the gradient is to be taken over a different mini-batch at each step. We call this algorithm gradient descent with momentum (GDM). From these equations, it is clear that \mathbf{v}_t is a running average of recently encountered gradients and $(1 - \gamma)^{-1}$ sets the characteristic time scale for the memory used in the averaging procedure. Consistent with this, when $\gamma = 0$, this just reduces down to ordinary SGD as discussed earlier. An equivalent way of writing the updates is

More on momentum based approaches

Let us try to get more intuition from these equations. It is helpful to consider a simple physical analogy with a particle of mass m moving in a viscous medium with drag coefficient μ and potential $E(w)$. If we denote the particle's position by w , then its motion is described by

$$m \frac{d^2 w}{dt^2} + \mu \frac{dw}{dt} = -\nabla_w E(w).$$

We can discretize this equation in the usual way to get

$$m \frac{w_{t+\Delta t} - 2w_t + w_{t-\Delta t}}{(\Delta t)^2} + \mu \frac{w_{t+\Delta t} - w_t}{\Delta t} = -\nabla_w E(w).$$

Rearranging this equation, we can rewrite this as

$$\Delta w_{t+\Delta t} = -\frac{(\Delta t)^2}{m + \mu \Delta t} \nabla_w E(w) + \frac{m}{m + \mu \Delta t} \Delta w_t.$$

Momentum parameter

Notice that this equation is identical to previous one if we identify the position of the particle, w , with the parameters θ . This allows us to identify the momentum parameter and learning rate with the mass of the particle and the viscous drag as:

$$\gamma = \frac{m}{m + \mu\Delta t}, \quad \eta = \frac{(\Delta t)^2}{m + \mu\Delta t}.$$

Thus, as the name suggests, the momentum parameter is proportional to the mass of the particle and effectively provides inertia. Furthermore, in the large viscosity/small learning rate limit, our memory time scales as $(1 - \gamma)^{-1} \approx m/(\mu\Delta t)$.

Why is momentum useful? SGD momentum helps the gradient descent algorithm gain speed in directions with persistent but small gradients even in the presence of stochasticity, while suppressing oscillations in high-curvature directions. This becomes especially important in situations where the landscape is shallow and flat in some directions and narrow and steep in others. It has been argued that first-order methods (with appropriate initial conditions) can perform comparable to more expensive second order methods

Second moment of the gradient

In stochastic gradient descent, with and without momentum, we still have to specify a schedule for tuning the learning rates η_t as a function of time. As discussed in the context of Newton's method, this presents a number of dilemmas. The learning rate is limited by the steepest direction which can change depending on the current position in the landscape. To circumvent this problem, ideally our algorithm would keep track of curvature and take large steps in shallow, flat directions and small steps in steep, narrow directions. Second-order methods accomplish this by calculating or approximating the Hessian and normalizing the learning rate by the curvature. However, this is very computationally expensive for extremely large models. Ideally, we would like to be able to adaptively change the step size to match the landscape without paying the steep computational price of calculating or approximating Hessians.

Recently, a number of methods have been introduced that accomplish this by tracking not only the gradient, but also the second moment of the gradient. These methods include AdaGrad, AdaDelta, Root Mean Squared Propagation (RMS-Prop), and

RMS prop

In RMS prop, in addition to keeping a running average of the first moment of the gradient, we also keep track of the second moment denoted by $s_t = \mathbb{E}[g_t^2]$. The update rule for RMS prop is given by

$$\begin{aligned} g_t &= \nabla_{\theta} E(\theta) \\ s_t &= \beta s_{t-1} + (1 - \beta) g_t^2 \\ \theta_{t+1} &= \theta_t - \eta_t \frac{g_t}{\sqrt{s_t + \epsilon}}, \end{aligned} \tag{3}$$

where β controls the averaging time of the second moment and is typically taken to be about $\beta = 0.9$, η_t is a learning rate typically chosen to be 10^{-3} , and $\epsilon \sim 10^{-8}$ is a small regularization constant to prevent divergences. Multiplication and division by vectors is understood as an element-wise operation. It is clear from this formula that the learning rate is reduced in directions where the norm of the gradient is consistently large. This greatly speeds up the convergence by allowing us to use a larger learning rate for flat directions.

ADAM optimizer

A related algorithm is the ADAM optimizer. In [ADAM](#), we keep a running average of both the first and second moment of the gradient and use this information to adaptively change the learning rate for different parameters. The method is efficient when working with large problems involving lots of data and/or parameters. It is a combination of the gradient descent with momentum algorithm and the RMSprop algorithm discussed above.

In addition to keeping a running average of the first and second moments of the gradient (i.e. $\mathbf{m}_t = \mathbb{E}[\mathbf{g}_t]$ and $\mathbf{s}_t = \mathbb{E}[\mathbf{g}_t^2]$, respectively), ADAM performs an additional bias correction to account for the fact that we are estimating the first two moments of the gradient using a running average (denoted by the hats in the update rule below). The update rule for ADAM is given by (where multiplication and division are once again understood to be element-wise operations below)

$$\begin{aligned}\mathbf{g}_t &= \nabla_{\theta} E(\theta) \\ \mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t\end{aligned}\tag{4}$$

Algorithms and codes for Adagrad, RMSprop and Adam

The algorithms we have implemented are well described in the text by [Goodfellow, Bengio and Courville](#), chapter 8.

The codes which implement these algorithms are discussed after our presentation of automatic differentiation.

Practical tips

- ▶ **Randomize the data when making mini-batches.** It is always important to randomly shuffle the data when forming mini-batches. Otherwise, the gradient descent method can fit spurious correlations resulting from the order in which data is presented.
- ▶ **Transform your inputs.** Learning becomes difficult when our landscape has a mixture of steep and flat directions. One simple trick for minimizing these situations is to standardize the data by subtracting the mean and normalizing the variance of input variables. Whenever possible, also decorrelate the inputs. To understand why this is helpful, consider the case of linear regression. It is easy to show that for the squared error cost function, the Hessian of the cost function is just the correlation matrix between the inputs. Thus, by standardizing the inputs, we are ensuring that the landscape looks homogeneous in all directions in parameter space. Since most deep networks can be viewed as linear transformations followed by a non-linearity at each layer, we expect this intuition to hold beyond the linear case.

Automatic differentiation

Automatic differentiation (AD), also called algorithmic differentiation or computational differentiation, is a set of techniques to numerically evaluate the derivative of a function specified by a computer program. AD exploits the fact that every computer program, no matter how complicated, executes a sequence of elementary arithmetic operations (addition, subtraction, multiplication, division, etc.) and elementary functions (exp, log, sin, cos, etc.). By applying the chain rule repeatedly to these operations, derivatives of arbitrary order can be computed automatically, accurately to working precision, and using at most a small constant factor more arithmetic operations than the original program.

Automatic differentiation is neither:

- ▶ Symbolic differentiation, nor
- ▶ Numerical differentiation (the method of finite differences).

Symbolic differentiation can lead to inefficient code and faces the difficulty of converting a computer program into a single expression, while numerical differentiation can introduce round-off errors in the discretization process and cancellation

Using autograd

Here we experiment with what kind of functions Autograd is capable of finding the gradient of. The following Python functions are just meant to illustrate what Autograd can do, but please feel free to experiment with other, possibly more complicated, functions as well.

```
import autograd.numpy as np
from autograd import grad
```

```
def f1(x):
    return x**3 + 1
```

```
f1_grad = grad(f1)
```

```
# Remember to send in float as argument to the computed gradient from
a = 1.0
```

```
# See the evaluated gradient at a using autograd:
```

```
print("The gradient of f1 evaluated at a = %g using autograd is: %g"%(
```

```
# Compare with the analytical derivative, that is  $f1'(x) = 3x^2$ 
```

```
grad_analytical = 3*a**2
```

```
print("The gradient of f1 evaluated at a = %g by finding the analytic
```

Autograd with more complicated functions

To differentiate with respect to two (or more) arguments of a Python function, Autograd need to know at which variable the function is being differentiated with respect to.

```
import autograd.numpy as np
from autograd import grad
def f2(x1,x2):
    return 3*x1**3 + x2*(x1 - 5) + 1
```

```
# By sending the argument 0, Autograd will compute the derivative w.r.
f2_grad_x1 = grad(f2,0)
```

```
# ... and differentiate w.r.t x2 by sending 1 as an additional argument
f2_grad_x2 = grad(f2,1)
```

```
x1 = 1.0
x2 = 3.0
```

```
print("Evaluating at x1 = %g, x2 = %g"%(x1,x2))
print("-"*30)
```

```
# Compare with the analytical derivatives:
```

```
# Derivative of f2 w.r.t x1 is: 9*x1**2 + x2:
f2_grad_x1_analytical = 9*x1**2 + x2
```

```
# Derivative of f2 w.r.t x2 is: x1 - 5:
f2_grad_x2_analytical = x1 - 5
```

More complicated functions using the elements of their arguments directly

```
import autograd.numpy as np
from autograd import grad
def f3(x): # Assumes x is an array of length 5 or higher
    return 2*x[0] + 3*x[1] + 5*x[2] + 7*x[3] + 11*x[4]**2

f3_grad = grad(f3)

x = np.linspace(0,4,5)

# Print the computed gradient:
print("The computed gradient of f3 is: ", f3_grad(x))

# The analytical gradient is: (2, 3, 5, 7, 22*x[4])
f3_grad_analytical = np.array([2, 3, 5, 7, 22*x[4]])

# Print the analytical gradient:
print("The analytical gradient of f3 is: ", f3_grad_analytical)
```

Note that in this case, when sending an array as input argument, the output from Autograd is another array. This is the true gradient of the function, as opposed to the function in the previous example. By using arrays to represent the variables, the output from Autograd might be easier to work with, as the output is closer

Functions using mathematical functions from Numpy

```
import autograd.numpy as np
from autograd import grad
def f4(x):
    return np.sqrt(1+x**2) + np.exp(x) + np.sin(2*np.pi*x)

f4_grad = grad(f4)

x = 2.7

# Print the computed derivative:
print("The computed derivative of f4 at x = %g is: %g"%(x,f4_grad(x)))

# The analytical derivative is:  $x/\sqrt{1 + x^2} + \exp(x) + \cos(2\pi x)$ 
f4_grad_analytical = x/np.sqrt(1 + x**2) + np.exp(x) + np.cos(2*np.pi*x)

# Print the analytical gradient:
print("The analytical gradient of f4 at x = %g is: %g"%(x,f4_grad_anal
```

More autograd

```
import autograd.numpy as np
from autograd import grad
def f5(x):
    if x >= 0:
        return x**2
    else:
        return -3*x + 1

f5_grad = grad(f5)

x = 2.7

# Print the computed derivative:
print("The computed derivative of f5 at x = %g is: %g"%(x,f5_grad(x)))
```

And with loops

```
import autograd.numpy as np
from autograd import grad
def f6_for(x):
    val = 0
    for i in range(10):
        val = val + x**i
    return val
```

```
def f6_while(x):
    val = 0
    i = 0
    while i < 10:
        val = val + x**i
        i = i + 1
    return val
```

```
f6_for_grad = grad(f6_for)
f6_while_grad = grad(f6_while)
```

```
x = 0.5
```

```
# Print the computed derivaties of f6_for and f6_while
```

```
print("The computed derivative of f6_for at x = %g is: %g"%(x,f6_for_grad(x)))
```

```
print("The computed derivative of f6_while at x = %g is: %g"%(x,f6_while_grad(x)))
```

```
import autograd.numpy as np
from autograd import grad
```

```
# Both of the functions are implementation of the sum: sum(x**i) for i in range(10)
```

Using recursion

```
import autograd.numpy as np
from autograd import grad
```

```
def f7(n): # Assume that n is an integer
    if n == 1 or n == 0:
        return 1
    else:
        return n*f7(n-1)
```

```
f7_grad = grad(f7)
```

```
n = 2.0
```

```
print("The computed derivative of f7 at n = %d is: %g"%(n,f7_grad(n)))
```

```
# The function f7 is an implementation of the factorial of n.
# By using the product rule, one can find that the derivative is:
```

```
f7_grad_analytical = 0
for i in range(int(n)-1):
    tmp = 1
    for k in range(int(n)-1):
        if k != i:
            tmp *= (n - k)
    f7_grad_analytical += tmp
```

```
print("The analytical derivative of f7 at n = %d is: %g"%(n,f7_grad_an
```


Unsupported functions

Autograd supports many features. However, there are some functions that is not supported (yet) by Autograd.

Assigning a value to the variable being differentiated with respect to

```
import autograd.numpy as np
from autograd import grad
def f8(x): # Assume x is an array
    x[2] = 3
    return x*2

f8_grad = grad(f8)

x = 8.4

print("The derivative of f8 is:",f8_grad(x))
```

Here, Autograd tells us that an 'ArrayBox' does not support item assignment. The item assignment is done when the program tries to assign `x[2]` to the value 3. However, Autograd has implemented the computation of the derivative such that this assignment is not possible.

The syntax `a.dot(b)` when finding the dot product

```
import autograd.numpy as np
from autograd import grad
def f9(a): # Assume a is an array with 2 elements
    b = np.array([1.0,2.0])
    return a.dot(b)

f9_grad = grad(f9)

x = np.array([1.0,0.0])

print("The derivative of f9 is:",f9_grad(x))
```

Here we are told that the 'dot' function does not belong to Autograd's version of a Numpy array. To overcome this, an alternative syntax which also computed the dot product can be used:

```
import autograd.numpy as np
from autograd import grad
def f9_alternative(x): # Assume a is an array with 2 elements
    b = np.array([1.0,2.0])
    return np.dot(x,b) # The same as  $x_1*b_1 + x_2*b_2$ 

f9_alternative_grad = grad(f9_alternative)

x = np.array([3.0,0.0])
```

Recommended to avoid

The documentation recommends to avoid inplace operations such as

`a += b`

`a -= b`

`a *= b`

`a /=b`

Using Autograd with OLS

We conclude the part on optimization by showing how we can make codes for linear regression and logistic regression using **autograd**.

The first example shows results with ordinary least squares.

```
# Using Autograd to calculate gradients for OLS
```

```
from random import random, seed
```

```
import numpy as np
```

```
import autograd.numpy as np
```

```
import matplotlib.pyplot as plt
```

```
from autograd import grad
```

```
def CostOLS(beta):
```

```
    return (1.0/n)*np.sum((y-X @ beta)**2)
```

```
n = 100
```

```
x = 2*np.random.rand(n,1)
```

```
y = 4+3*x+np.random.randn(n,1)
```

```
X = np.c_[np.ones((n,1)), x]
```

```
XT_X = X.T @ X
```

```
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
```

```
print("Own inversion")
```

```
print(theta_linreg)
```

```
# Hessian matrix
```

```
H = (2.0/n)* XT_X
```

```
EigValues, EigVectors = np.linalg.eig(H)
```

```
print(f"Eigenvalues of Hessian Matrix:{EigValues}")
```

Same code but now with momentum gradient descent

```
# Using Autograd to calculate gradients for OLS
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

def CostOLS(beta):
    return (1.0/n)*np.sum((y-X @ beta)**2)

n = 100
x = 2*np.random.rand(n,1)
y = 4+3*x#+np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x]
XT_X = X.T @ X
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)
# Hessian matrix
H = (2.0/n)* XT_X
EigValues, EigVectors = np.linalg.eig(H)
print(f"Eigenvalues of Hessian Matrix:{EigValues}")

theta = np.random.randn(2,1)
eta = 1.0/np.max(EigValues)
Niterations = 30
```

But none of these can compete with Newton's method

```
# Using Newton's method
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

def CostOLS(beta):
    return (1.0/n)*np.sum((y-X @ beta)**2)

n = 100
x = 2*np.random.rand(n,1)
y = 4+3*x+np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x]
XT_X = X.T @ X
beta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(beta_linreg)
# Hessian matrix
H = (2.0/n)* XT_X
# Note that here the Hessian does not depend on the parameters beta
invH = np.linalg.pinv(H)
EigValues, EigVectors = np.linalg.eig(H)
print(f"Eigenvalues of Hessian Matrix:{EigValues}")

beta = np.random.randn(2,1)
Niterations = 5
```

Including Stochastic Gradient Descent with Autograd

In this code we include the stochastic gradient descent approach discussed above. Note here that we specify which argument we are taking the derivative with respect to when using **autograd**.

```
# Using Autograd to calculate gradients using SGD
```

```
# OLS example
```

```
from random import random, seed
```

```
import numpy as np
```

```
import autograd.numpy as np
```

```
import matplotlib.pyplot as plt
```

```
from autograd import grad
```

```
# Note change from previous example
```

```
def CostOLS(y,X,theta):
```

```
    return np.sum((y-X @ theta)**2)
```

```
n = 100
```

```
x = 2*np.random.rand(n,1)
```

```
y = 4+3*x+np.random.randn(n,1)
```

```
X = np.c_[np.ones((n,1)), x]
```

```
XT_X = X.T @ X
```

```
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
```

```
print("Own inversion")
```

```
print(theta_linreg)
```

```
# Hessian matrix
```

```
H = (2.0/n)* XT_X
```

```
EigValues, EigVectors = np.linalg.eig(H)
```

Same code but now with momentum gradient descent

```
# Using Autograd to calculate gradients using SGD
# OLS example
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

# Note change from previous example
def CostOLS(y,X,theta):
    return np.sum((y-X @ theta)**2)

n = 100
x = 2*np.random.rand(n,1)
y = 4+3*x+np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x]
XT_X = X.T @ X
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)
# Hessian matrix
H = (2.0/n)* XT_X
EigValues, EigVectors = np.linalg.eig(H)
print(f"Eigenvalues of Hessian Matrix:{EigValues}")

theta = np.random.randn(2,1)
eta = 1.0/np.max(EigValues)
```


Similar (second order function now) problem but now with AdaGrad

```
# Using Autograd to calculate gradients using AdaGrad and Stochastic Gradient Descent
# OLS example
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

# Note change from previous example
def CostOLS(y,X,theta):
    return np.sum((y-X @ theta)**2)

n = 1000
x = np.random.rand(n,1)
y = 2.0+3*x +4*x*x

X = np.c_[np.ones((n,1)), x, x*x]
XT_X = X.T @ X
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)

# Note that we request the derivative wrt third argument (theta, 2 here)
training_gradient = grad(CostOLS,2)
# Define parameters for Stochastic Gradient Descent
```

RMSprop for adaptive learning rate with Stochastic Gradient Descent

```
# Using Autograd to calculate gradients using RMSprop and Stochastic
# OLS example
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

# Note change from previous example
def CostOLS(y,X,theta):
    return np.sum((y-X @ theta)**2)

n = 1000
x = np.random.rand(n,1)
y = 2.0+3*x +4*x*x# +np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x, x*x]
XT_X = X.T @ X
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)

# Note that we request the derivative wrt third argument (theta, 2 here)
training_gradient = grad(CostOLS,2)
# Define parameters for Stochastic Gradient Descent
```

And finally ADAM

```
# Using Autograd to calculate gradients using RMSprop and Stochastic
# OLS example
from random import random, seed
import numpy as np
import autograd.numpy as np
import matplotlib.pyplot as plt
from autograd import grad

# Note change from previous example
def CostOLS(y,X,theta):
    return np.sum((y-X @ theta)**2)

n = 1000
x = np.random.rand(n,1)
y = 2.0+3*x +4*x*x# +np.random.randn(n,1)

X = np.c_[np.ones((n,1)), x, x*x]
XT_X = X.T @ X
theta_linreg = np.linalg.pinv(XT_X) @ (X.T @ y)
print("Own inversion")
print(theta_linreg)

# Note that we request the derivative wrt third argument (theta, 2 here)
training_gradient = grad(CostOLS,2)
# Define parameters for Stochastic Gradient Descent
n_epochs = 50
M = 5 #size of each minibatch
```

And Logistic Regression

```
import autograd.numpy as np
from autograd import grad
```

```
def sigmoid(x):
    return 0.5 * (np.tanh(x / 2.) + 1)
```

```
def logistic_predictions(weights, inputs):
    # Outputs probability of a label being true according to logistic
    return sigmoid(np.dot(inputs, weights))
```

```
def training_loss(weights):
    # Training loss is the negative log-likelihood of the training labels
    preds = logistic_predictions(weights, inputs)
    label_probabilities = preds * targets + (1 - preds) * (1 - targets)
    return -np.sum(np.log(label_probabilities))
```

```
# Build a toy dataset.
```

```
inputs = np.array([[0.52, 1.12, 0.77],
                   [0.88, -1.08, 0.15],
                   [0.52, 0.06, -1.30],
                   [0.74, -2.49, 1.39]])
```

```
targets = np.array([True, True, False, True])
```

```
# Define a function that returns gradients of training loss using Auto
training_gradient_fun = grad(training_loss)
```

```
# Optimize weights using gradient descent.
```

```
weights = np.array([0.0, 0.0, 0.0])
```

Introducing JAX

Presently, instead of using **autograd**, we recommend using **JAX**. **JAX** is Autograd and **XLA (Accelerated Linear Algebra)**, brought together for high-performance numerical computing and machine learning research. It provides composable transformations of Python+NumPy programs: differentiate, vectorize, parallelize, Just-In-Time compile to GPU/TPU, and more.

Here's a simple example on how you can use **JAX** to compute the derivative of the logistic function.

```
import jax.numpy as jnp
from jax import grad, jit, vmap

def sum_logistic(x):
    return jnp.sum(1.0 / (1.0 + jnp.exp(-x)))

x_small = jnp.arange(3.)
derivative_fn = grad(sum_logistic)
print(derivative_fn(x_small))
```