

Linear Predictions

Ju Sun*

March 23, 2021

Overview We first interpret supervised learning from the viewpoint of function approximation, and then survey classic linear prediction models and algorithms, including linear least squares for regression, and Perceptron, simple SVM, logistic regression for binary classification.

1 Function approximation view of supervised learning

Given: a data set $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ (called *training set*) so that $\mathbf{y}_i \approx f_*(\mathbf{x}_i) \forall i$. Here f_* is an unknown underlying function, and for all $i \in [N]$, the “ \approx ” sign in $\mathbf{y}_i \approx f_*(\mathbf{x}_i)$ is to allow noise or other errors over $\mathbf{y}_i = f_*(\mathbf{x}_i)$, e.g., $\mathbf{y}_i = f(\mathbf{x}_i) + \varepsilon_i$ for Gaussian noise ε_i . Some terminology:

\mathbf{x}_i is called the **input/predictor** (in statistics)/**features** (in pattern recognition),
 \mathbf{y}_i is called the **output/response** (in statistics)/**label** (in pattern recognition).

There are often three steps in a typical supervised learning workflow:

- **Step 1: Modeling (or Model selection).** Choose a family/set of functions \mathcal{H} , called the *hypothesis class* or *hypothesis set*, so that there exists an $f_{\mathcal{H}} \in \mathcal{H}$ that is “close” to f_* . Often, \mathcal{H} should be *reasonably large* to ensure there is indeed a good approximation $f_{\mathcal{H}}$ to f_* , and also *reasonably small/simple* so that such an $f_{\mathcal{H}}$ can be found in a computationally efficient manner.
- **Step 2: Computation (or Training).** Design an algorithm to find such an $f_{\mathcal{H}}$. In modern machine learning, one often first formulates the learning problem as an optimization problem, and then develops numerical optimization algorithms to solve the optimization problem—this is why optimization is a crucial component of modern machine learning¹.
 - **Optimization formulation.** A family of natural and popular formulation is *structural risk minimization, or SRM* (it is called this under certain additional probabilistic assumptions on the training set; we will talk more about this when introducing statistical learning theory later.):

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{y}_i, f(\mathbf{x}_i)) + \mathcal{R}(f). \quad (1.1)$$

*Department of Computer Science and Engineering and Department of Neurosurgery, University of Minnesota at Twin Cities. Email: jusun@umn.edu.

¹ISyE of UMN offers a new course *Optimization for Machine Learning* that covers popular scalable numerical optimization methods for solving large-scale (or practical-scale) machine learning problems. Other good resources include [Sra12, Sta].

Here ℓ represents the loss function to be chosen, which measures the difference between \mathbf{y}_i and $f(\mathbf{x}_i)$. Obviously, $\min_f \frac{1}{N} \cdot \sum_{i=1}^N \ell(\mathbf{y}_i, f(\mathbf{x}_i))$ tries to ensure that $\mathbf{y}_i \approx f(\mathbf{x}_i)$ for all i . The second term $\mathcal{R}(\mathbf{x})$ is typically called *regularizer* or *regularization term*, which puts certain preference on f that will be found: this is often needed when there are multiple or even infinitely many $f_{\mathcal{H}}$ that are good—perhaps because the \mathcal{H} we choose is larger than necessary—so that we have to restrict our search to certain f 's that are practically interesting.

- **Optimization algorithm.** For very simple problems, Eq. (1.1) may admit a closed-form analytic solution. But in modern machine learning, that is very rare and iterative numerical optimization methods are almost always needed to solve Eq. (1.1). Depending on the problem and formulation choice,

(i) Eq. (1.1) can be an unconstrained or constrained optimization problem. In general, unconstrained optimization problems are (much) easier to solve than constrained problems. So for applications with large-scale datasets, one often is willing to make reasonable compromise and tries to formulate unconstrained optimization problems or constrained optimization problems with very simple constraints so that scalable optimization algorithms can be developed;

(ii) Eq. (1.1) can be a convex or nonconvex problem. We will provide a quick review of convex analysis and optimization later. Analysis and optimization of convex problems are much more mature than nonconvex ones, and so there is an overall preference for convex formulations. But revival of deep learning after 2010 has substantially changed this—optimization problems in deep learning are always highly nonconvex.

- **Step 3: Generalization.** Measure how close the $f_{\mathcal{H}}$ found from **Step 2** is to f_* . Since we often do not know the true f_* , generalization can only be measured indirectly, e.g., by evaluating the mean of

$$d(f_*(\mathbf{x}_i), f_{\mathcal{H}}(\mathbf{x}_j)) \approx d(\mathbf{y}_j, f_{\mathcal{H}}(\mathbf{x}_j)) \quad (1.2)$$

over unseen (i.e., test) dataset $\{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^M$. Here $d(\cdot, \cdot)$ is a difference (error) function that may or may not be the same as the ℓ above. To study generalization in a rigorous manner, especially to quantify the relationship between the size of training set (i.e., *sample complexity*) and generalization, we need to put additional assumptions on the training set, e.g., the data points are sampled iid (i.e., independent and identically distributed) from an underlying probability distribution. We will talk about this in later lectures on statistical learning theory.

Below, we start with linear regression and linear classification problems, and illustrate Steps 1 & 2. We will derive their generalization properties in the learning theory lectures.

When \mathbf{y}_i 's in the training set are categorical, i.e., indicating the memberships of the inputs \mathbf{x}_i 's in a set of categories (e.g., {cat, dog, else} for image inputs, {COVID, Non-COVID} given patients' symptoms), the learning problem is often modeled as classification. Otherwise, it will be modeled as regression.

2 Linear regression

For simplicity, we assume that $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ for all $i \in [N]$ (remember our convention is that scalars are non-bold small letters.). In linear regression, we model the relationship between \mathbf{x} and

y as linear—arguably the simplest possible:

$$\mathbf{y}_i \approx \langle \mathbf{w}, \mathbf{x}_i \rangle + b \quad \forall i \in [N], \quad \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}. \quad (2.1)$$

In words, the hypothesis class is the set of all linear functions in \mathbf{x} , which we can write as²

$$\mathcal{H}_L = \left\{ \mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \right\}. \quad (2.2)$$

2.1 Formulation

Once we decide the hypothesis class, we are ready to formulate the problem in an optimization problem. When employing the SRM framework, we need to choose an appropriate loss ℓ and regularizer \mathcal{R} . For simplicity, let's choose squared loss, i.e., $\ell(y_i, f(\mathbf{x}_i)) = (y_i - f(\mathbf{x}_i))^2$ for all i , and suppose that we do not need regularization now. This leads to a least-squares formulation:

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{N} \sum_{i=1}^N (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle - b)^2. \quad (2.3)$$

Now we are to turn this formulation into an equivalent yet compact form using matrix notations—this facilitates more direct translation of the mathematical expressions into modern numerical programming languages that are optimized for matrix computations. We append each \mathbf{x}_i with an additional coordinate 1, so that $\mathbf{x}'_i = \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \in \mathbb{R}^{d+1}$; correspondingly, we concatenate \mathbf{w} and b into $\mathbf{w}' = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} \in \mathbb{R}^{d+1}$. It is easy to verify that

$$\langle \mathbf{w}', \mathbf{x}'_i \rangle = \langle \mathbf{w}, \mathbf{x}_i \rangle + b \quad \forall i \in [N]. \quad (2.4)$$

We call this the *homogeneous form of linear functions*. This allows us to write Eq. (2.3) as

$$\min_{\mathbf{w}' \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{i=1}^N (y_i - \langle \mathbf{w}', \mathbf{x}'_i \rangle)^2. \quad (2.5)$$

The prime notation $(\cdot)'$ looks messy. Since we often use the homogeneous form, with slight abuse of notation, we will just remove the $(\cdot)'$ from \mathbf{w}' and \mathbf{x}'_i , with the understanding that the homogeneous notation will be either explicitly stated or can be inferred from the dimension of \mathbf{w} (i.e., $\mathbf{w} \in \mathbb{R}^{d+1}$). So we have

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \frac{1}{N} \sum_{i=1}^N (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2. \quad (2.6)$$

Last step, if we write

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad \text{and} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}, \quad (2.7)$$

and recall the definition of vector ℓ_2 norm, we arrive at a compact form of Eq. (2.6):

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} g(\mathbf{w}) \doteq \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2, \quad (2.8)$$

where we have omitted the factor $1/N$ in the objective, as it does not affect the solution. Our next job is to solve the optimization problem Eq. (2.8) to find a good linear model that fits the data well.

²Recall that we typically represent a set as: {generic form of elements in the set : constraints on elements in the set if any}. In Eq. (2.2), a generic linear function in \mathbf{x} can be represented as $\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle + b$, and the constraints are $\mathbf{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ —which are vacuous.

2.2 Solution via optimality condition

Our least-squares problem is unconstrained, and also the objective function $g(\mathbf{w})$ is relatively simple. So we shall try optimality conditions to see if they lead to somewhere.

We now quickly review optimality conditions for unconstrained optimization problems. Consider a minimization problem

$$\min_{\mathbf{z} \in \mathbb{R}^n} f(\mathbf{z}). \quad (2.9)$$

It is sufficient to consider minimization problems, as maximization problems of the form $\max_{\mathbf{z} \in \mathbb{R}^n} f(\mathbf{z})$ are equivalent to $\min_{\mathbf{z} \in \mathbb{R}^n} -f(\mathbf{z})$ —they have the same optimizer(s).

A point $\mathbf{z}_0 \in \mathbb{R}^n$ is a *local minimizer* of $f(\mathbf{z})$ if there exists a radius η so that $f(\mathbf{z}_0) \leq f(\mathbf{z})$ for all \mathbf{z} satisfying $\|\mathbf{z} - \mathbf{z}_0\|_2 \leq \eta$; in words, if $f(\mathbf{z}_0)$ is no larger than any other $f(\mathbf{z})$ in an η -ball around \mathbf{z}_0 . The value $f(\mathbf{z}_0)$ is called a *local minimum*. So minimizers concern the optimization variables, and minimums (or minima) concern the objective value.

For minimization problems, optimality conditions are mathematical conditions that any local minimizer must satisfy, and hence they are helpful for locating local minimizers either analytically, or numerically.

Theorem 2.1 (First-order necessary condition of optimality for unconstrained problems). *Assume f is first-order differentiable at \mathbf{z}_0 . If \mathbf{z}_0 is a local minimizer, then $\nabla f(\mathbf{z}_0) = \mathbf{0}$.*

The zero-gradient condition is a necessary condition for local minimizers, but not sufficient. A point where the gradient is zero (also called *first-order stationary point*, or FOSP) can be a local minimizer, local maximizer, or saddle point. It turns out the condition becomes sufficient for the family of convex functions—more on this when we talk about support vector machines and kernel methods. A salient feature of convex functions is that a local minimizer is also a global minimizer³. One way to tell convexity is through the Hessian.

Lemma 2.2 (Convexity through Hessian). *Assume f is second-order differentiable. Then f is convex if and only if $\nabla^2 f(\mathbf{z}) \succeq \mathbf{0}$ for all \mathbf{z} .*

Theorem 2.3 (First-order sufficient condition of optimality for unconstrained convex problems). *Assume f is convex and first-order differentiable at \mathbf{z}_0 . If $\nabla f(\mathbf{z}_0) = \mathbf{0}$, then \mathbf{z}_0 is a local and global minimizer of f .*

There is a more refined characterization of local minimizers using both gradient and Hessian.

Theorem 2.4 (Second-order necessary condition of optimality for unconstrained problems). *Assume f is second-order differentiable at \mathbf{z}_0 . If \mathbf{z}_0 is a local minimizer, then $\nabla f(\mathbf{z}_0) = \mathbf{0}$, and $\nabla^2 f(\mathbf{z}_0) \succeq \mathbf{0}$, i.e., Hessian at \mathbf{z}_0 is positive semidefinite.*

A point \mathbf{z}_0 satisfying $\nabla f(\mathbf{z}_0) = \mathbf{0}$ and $\nabla^2 f(\mathbf{z}_0) \succeq \mathbf{0}$ is called a *second-order stationary point*, or SOSP. Similar to FOSP, a SOSP can be a local minimizer, local maximizer, or saddle point. A stronger condition can ensure a local minimizer.

Theorem 2.5 (Second-order sufficient condition of optimality for unconstrained problems). *Assume f is second-order differentiable at \mathbf{z}_0 . If $\nabla f(\mathbf{z}_0) = \mathbf{0}$, and $\nabla^2 f(\mathbf{z}_0) \succ \mathbf{0}$, i.e., Hessian at \mathbf{z}_0 is **positive definite**, then \mathbf{z}_0 is a local minimizer.*

³A convex function has a unique local minimum—which is also the global minimum, but could have multiple local minimizers that are also global minimizers.

For our problem Eq. (2.8), the least squares objective is a quadratic polynomial and hence is second-order differentiable. The Hessian is $2\mathbf{X}^\top\mathbf{X}$, which is positive semidefinite. So $g(\mathbf{w})$ is a convex function.

Invoking Theorem 2.1, if $\mathbf{w}_0 \in \mathbb{R}^{d+1}$ is a local minimizer, then

$$\nabla g(\mathbf{w}_0) = 2\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\mathbf{w}_0) = \mathbf{0} \implies \mathbf{X}^\top\mathbf{X}\mathbf{w}_0 = \mathbf{X}^\top\mathbf{y}. \quad (2.10)$$

- If $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$ has full column rank, i.e., with linearly independent columns, or equivalently the N data points in \mathbf{X} span the $(d+1)$ -dimensional input space, then $\mathbf{X}^\top\mathbf{X}$ has full rank and hence is invertible. Then

$$\mathbf{w}_0 = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}. \quad (2.11)$$

Obviously, \mathbf{w}_0 is uniquely defined by the right side of Eq. (2.11) and so is the unique global minimizer of $g(\mathbf{w})$.

- Otherwise, $\mathbf{X}^\top\mathbf{X}$ is not invertible and there are multiple (in fact, infinitely many) global minimizers. A particular global minimizer can be found through the *pseudo-inverse*. Recall that for any matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$ of rank $r \leq \min(m, n)$, its *compact SVD* can be written as $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ where $\mathbf{U} \in \mathbb{R}^{m \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ is diagonal, and $\mathbf{V} \in \mathbb{R}^{n \times r}$. The pseudo-inverse of \mathbf{M} is then $\mathbf{M}^\dagger = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^\top \in \mathbb{R}^{n \times m}$. Obviously, pseudo-inverse is defined for any matrix, square or not. A solution to Eq. (2.10), and hence a global minimizer to $g(\mathbf{w})$, is then

$$\mathbf{w}_0 = (\mathbf{X}^\top\mathbf{X})^\dagger\mathbf{X}^\top\mathbf{y}. \quad (2.12)$$

For a square invertible matrix, its pseudo-inverse coincides with its inverse. So actually this provides a generic form of global minimizer to $g(\mathbf{w})$, whether $\mathbf{X}^\top\mathbf{X}$ is invertible or not.

The closed-form solution we get here seems nice, but probably only for small-scale problems. The cost of calculating $\mathbf{X}^\top\mathbf{X}$ is $O(d^2N)$, and inverting $\mathbf{X}^\top\mathbf{X} \in \mathbb{R}^{(d+1) \times (d+1)}$ costs $O(d^3)$. This is daunting when d and N are large. Can we do better?

2.3 Solution via iterative optimization

In numerical optimization, iterative methods start with an initial guess and produce a sequence of points that gradually approach a solution. Gradient descent (GD) is a basic yet powerful iterative method. For an unconstrained problem $\min_{\mathbf{z} \in \mathbb{R}^n} f(\mathbf{z})$, GD runs like this:

Algorithm 1 Gradient descent for minimizing $f(\mathbf{z})$

Input: initialization $\mathbf{z}^{(0)}$, $k = 1$, stopping precision $\varepsilon > 0$

- 1: **while** $\|\nabla f(\mathbf{z}^{(k-1)})\|_2 > \varepsilon$ **do**
 - 2: choose a step size $t^{(k)}$
 - 3: update the estimate: $\mathbf{z}^{(k)} = \mathbf{z}^{(k-1)} - t^{(k)}\nabla f(\mathbf{z}^{(k-1)})$
 - 4: update the counter: $k = k + 1$
 - 5: **end while**
-

Intuitively, to find a local minimizer, one can try to construct a sequence of iterates that carry monotonically decreasing function values. If we make a small movement \mathbf{d} from \mathbf{z} , Taylor's theorem says

$$f(\mathbf{z} + \mathbf{d}) \approx f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), \mathbf{d} \rangle \implies f(\mathbf{z} + \mathbf{d}) - f(\mathbf{z}) \approx \langle \nabla f(\mathbf{z}), \mathbf{d} \rangle. \quad (2.13)$$

We hope to make $f(\mathbf{z} + \mathbf{d}) - f(\mathbf{z})$ as negative as possible to make rapid progress toward a local minimizer, and so we can try to make $\langle \nabla f(\mathbf{z}), \mathbf{d} \rangle$ as negative as possible. Obviously aligning \mathbf{d} with $-\nabla f(\mathbf{z})$ ensures the fastest possible progress, as

$$-\frac{\nabla f(\mathbf{z})}{\|\nabla f(\mathbf{z})\|_2} = \arg \min_{\mathbf{d}: \|\mathbf{d}\|_2=1} \langle \nabla f(\mathbf{z}), \mathbf{d} \rangle. \quad (2.14)$$

This is where GD comes from.

Of course, the step sizes $t^{(k)}$'s are controlling the magnitudes of movement. On one hand, they should be made as large as possible to allow fast progress. On the other, they should be reasonably small to make the first-order Taylor approximation in Eq. (2.13) reasonably accurate. There are two popular strategies for choosing the step sizes:

- **Fixed step size** Choose a sufficiently small constant as the step size for all iterations. If the value is not sufficiently small, typically the objective value will blow up after a while. **Pros:** simple; **Cons:** could be conservative for most iterations.
- **Adaptive step size via backtracking line search** Search for an appropriate (large) step size adapted to local landscape of the function.

Algorithm 2 Gradient descent for minimizing $f(\mathbf{z})$ with backtracking line search

Input: initialization $\mathbf{z}^{(0)}$, $k = 1$, stopping precision $\varepsilon > 0$ close to 0

- 1: **while** $\|\nabla f(\mathbf{z}^{(k-1)})\|_2 > \varepsilon$ **do**
 - 2: choose initial step size t , $\rho \in (0, 1)$, and $\eta \in (0, 1)$
 - 3: **while** $f(\mathbf{z}^{(k-1)} - t\nabla f(\mathbf{z}^{(k-1)})) - f(\mathbf{z}^{(k-1)}) > -\eta t \|\nabla f(\mathbf{z}^{(k-1)})\|_2^2$ **do**
 - 4: decrease the step size: $t = \rho t$
 - 5: **end while**
 - 6: set the step size: $t^{(k)} = t$
 - 7: update the estimate: $\mathbf{z}^{(k)} = \mathbf{z}^{(k-1)} - t^{(k)} \nabla f(\mathbf{z}^{(k-1)})$
 - 8: update the counter: $k = k + 1$
 - 9: **end while**
-

Pros: relatively large step size, and hence fast movement and rapid convergence

Cons: slightly more computation each iteration for searching the good step size

The backtracking line-search strategy is highly recommended for practical implementation of GD. See Appendix B for the intuition behind the backtracking line search rule.

Since $\nabla f(\mathbf{z}) = \mathbf{0}$ is the necessary condition for \mathbf{z} being a local minimizer, we set the stopping criterion as checking if $\|\nabla f(\mathbf{z})\|_2$ is sufficiently close to 0.

When we apply GD to our least-squares problem, the gradient update step is:

$$\mathbf{w}^{(k)} = \mathbf{w}^{(k-1)} - t^{(k)} \nabla g(\mathbf{w}^{(k-1)}) = \mathbf{w}^{(k-1)} - 2t^{(k)} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}^{(k-1)}), \quad (2.15)$$

which costs $O(dN)$. Easy to check computing $\|\nabla f(\mathbf{z}^{(k-1)})\|_2$, $f(\mathbf{z}^{(k-1)} - t\nabla f(\mathbf{z}^{(k-1)}))$ and alike costs $O(dN)$ also. So the cost for each iteration is $O(dN)$, and so the total cost is $O(dNT)$ if T is the total number of iteration taken to find an approximate minimizer. When $T \ll \min(d, N)$, GD is computationally favorable for our problem compared to computing the solution using the closed-form formula Eq. (2.11) or Eq. (2.12).

2.4 Popular variants

When $\mathbf{X}^\top \mathbf{X}$ is not invertible, there are infinitely many global minimizers to our least squares problem. Particularly, this happens when $N < d + 1$, i.e., the number of data points is smaller than the input dimension. In statistics, this belongs to the family of so-called *high-dimensional problems* where typically a regularization term is added.

- **Ridge regression** takes the form

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad (2.16)$$

for a certain $\lambda > 0$. Since our linear model is $y \approx \langle \mathbf{w}, \mathbf{x} \rangle = \sum_j w_j x_j$, regularizing $\|\mathbf{w}\|_2^2$ ensures that entries in \mathbf{w} are all reasonably small so that any change to the input \mathbf{x} only causes a small change in the predicted value. In other words, the learned model is stable. Since the Hessian $2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \succeq \mathbf{0}$ everywhere, the objective is convex. In fact, $2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \succeq 2\mathbf{I}$ ⁴ and so the objective is *strongly* convex. Thus, ridge regression has a unique global minimizer.

- **Lasso** takes the form

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \quad (2.17)$$

for a certain $\lambda > 0$. Compared to $\|\mathbf{w}\|_2^2 = \sum_j w_j^2$, $\|\mathbf{w}\|_1 = \sum_j |w_j|$ penalizes large entries of \mathbf{w} much less and small entries much more. The net effect is that regularizing using $\|\mathbf{w}\|_1$ tends to lead a solution that is sparse—containing very few large entries and the rest negligible in magnitude. This is useful when one is interested to select only few most important features (i.e., columns) from \mathbf{X} . Both $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$ and $\lambda \|\mathbf{w}\|_1$ are convex, and so the positive combination is convex. However, in general, Lasso does not have a unique global minimizer either. **Elastic net**

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2^2 \quad (2.18)$$

which integrates Lasso and ridge regression is a fix to this and has a unique global minimizer, alongside other benefits over Lasso, such as stability when selecting correlated features.

Comparisons of various popular linear models for linear regression and classification can be found here

https://scikit-learn.org/stable/modules/linear_model.html.

3 Review of subspaces and hyperplanes

Consider a line L in \mathbb{R}^2 , as illustrated in Fig. 1 (left).

- If L passes through the origin, there are two ways to represent L . One possibility is to find a vector $\mathbf{v} \in \mathbb{R}^2$ aligned with the L , then $L = \{\lambda \mathbf{v} : \lambda \in \mathbb{R}\}$, called *basis representation*. The other possibility is to find a vector \mathbf{w} that is orthogonal (i.e., normal) to L , i.e., orthogonal to all vector in L , and then $L = \{\mathbf{x} \in \mathbb{R}^2 : \langle \mathbf{w}, \mathbf{x} \rangle = 0\}$, called *normal representation*. Obviously, $\langle \mathbf{w}, \mathbf{v} \rangle = 0$, i.e., \mathbf{w} is orthogonal to \mathbf{v} .

⁴For two symmetric matrices $\mathbf{M}_1, \mathbf{M}_2 \in \mathbb{R}^{n \times n}$, $\mathbf{M}_1 \succeq \mathbf{M}_2$ means $\mathbf{M}_1 - \mathbf{M}_2 \succeq \mathbf{0}$.

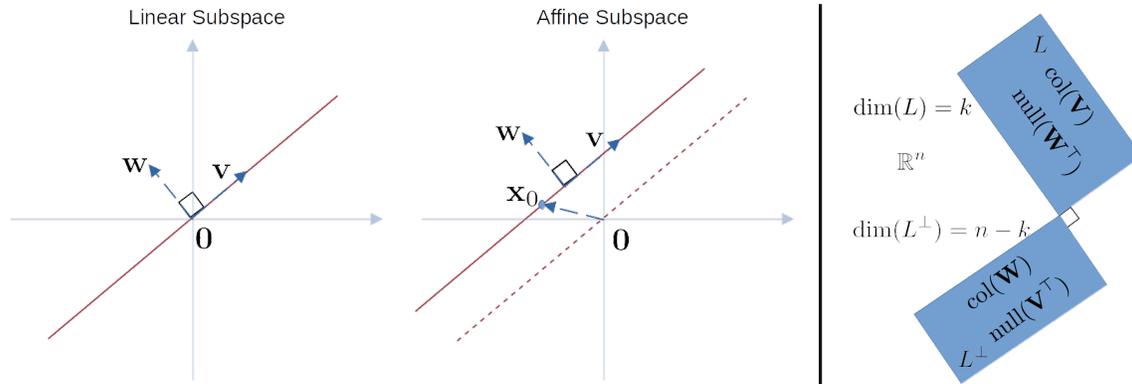


Figure 1: (left) Illustration of subspaces and hyperplanes; (right) Geometric picture of basis and normal representations. The picture is adapted from Sec 4.1 of the famous linear algebra book [Str16].

- If L does not pass through the origin, we can find an arbitrary point $x_0 \in L$ and write $L = x_0 + L' = \{x_0 + x : x \in L'\}$ for an L' that passes through the origin. So we can find v aligned with L' and w orthogonal to L' , so that

$$L = \underbrace{\{x_0 + \lambda v : \lambda \in \mathbb{R}\}}_{\text{basis representation}} = \underbrace{\{x \in \mathbb{R}^2 : \langle w, x \rangle = \langle w, x_0 \rangle\}}_{\text{normal representation}}. \quad (3.1)$$

Recall that a set $S \subset \mathbb{R}^n$ is called a subspace if for any $u, v \in S$, $\alpha u + \beta v \in S$ for all $\alpha, \beta \in \mathbb{R}$, i.e., any linear combination of u and v stays in the set. Geometrically, subspaces can be thought of as high-dimensional “flats” in \mathbb{R}^n , and they are natural generalizations of lines. Subspaces also admit both basis and normal representations that generalize the corresponding representation for lines: for any k -dimensional subspace $L \subset \mathbb{R}^n$,

- **basis representation:** for any k linearly independent vectors $\{v_1, \dots, v_k\}$ that span L , i.e., $\{v_1, \dots, v_k\}$ is a basis for L ,

$$L = \left\{ \sum_{i=1}^k \alpha_i v_i : \alpha_i \in \mathbb{R} \forall i \right\} = \{V\alpha : \alpha \in \mathbb{R}^k\} = \text{col}(V), \quad (3.2)$$

where $V \doteq [v_1 \dots v_k]$ and $\text{col}(\cdot)$ indicates the column space.

- **normal representation:** for any $n - k$ linearly independent vectors $\{w_1, \dots, w_{n-k}\}$ that are orthogonal to L , i.e., $\langle w_j, x \rangle = 0$ for all $j \in [n - k]$ and all $x \in L$,

$$L = \{x \in \mathbb{R}^n : \langle x, w_j \rangle = 0 \forall j \in [n - k]\} = \{x \in \mathbb{R}^n : W^T x = 0\} = \text{null}(W^T), \quad (3.3)$$

where $W \doteq [w_1 \dots w_{n-k}]$ and $\text{null}(\cdot)$ denotes the null space. Moreover, W spans the unique $(n - k)$ -dimensional orthogonal subspace L^\perp of L .⁵

The geometric aspect of the discussion is summarized in Fig. 1 (right).

⁵Two subspaces L and L' are said to be orthogonal to each other if z and z' are orthogonal to each other, i.e., $\langle z, z' \rangle = 0$, for all $z \in L$ and $z' \in L'$.

When the subspace L has dimension $n - 1$, it deserves a special name—*hyperplane*, which is a critical element of machine learning; in the next section, we need this object for linear classification.

All subspaces we talk of contain the origin; if we want to emphasize this fact, we prepend the adjective *linear*, i.e., call them *linear subspaces*. This is also to distinguish them with flats that do not necessarily pass through the origin—as generalization of lines that do not necessarily; we call these flats *affine subspaces*.

Similar to how we represent “affine” lines, any affine subspace L is a shifted linear subspace, i.e., $L = \mathbf{x}_0 + L'$ for certain $\mathbf{x}_0 \in L$ and linear subspace L' . This implies the following basis and normal representations for L :

- **basis representation:** assume $\mathbf{V} \in \mathbb{R}^{n \times k}$ spans L' and $\mathbf{x}_0 \in L$, then

$$L = \{\mathbf{x}_0 + \mathbf{x} : \mathbf{x} \in L'\} = \{\mathbf{x}_0 + \mathbf{V}\boldsymbol{\alpha} : \boldsymbol{\alpha} \in \mathbb{R}^k\}. \quad (3.4)$$

Moreover, $\dim(L) = \dim(L') = k$.

- **normal representation:** assume \mathbf{W} spans $(L')^\perp$ and $\mathbf{x}_0 \in L$, then

$$L = \{\mathbf{x}_0 + \mathbf{x} : \mathbf{x} \in L'\} = \{\mathbf{x}_0 + \mathbf{x} : \mathbf{W}^\top \mathbf{x} = \mathbf{0}\} = \{\mathbf{z} \in \mathbb{R}^n : \mathbf{W}^\top \mathbf{z} = \mathbf{W}^\top \mathbf{x}_0\}. \quad (3.5)$$

A natural question is whether $\mathbf{W}^\top \mathbf{x}_0$ is unique given that \mathbf{x}_0 is an arbitrary point on L . The answer is yes, as for any two points $\mathbf{x}_0, \mathbf{x}'_0 \in L$, $\mathbf{W}^\top \mathbf{x}_0 - \mathbf{W}^\top \mathbf{x}'_0 = \mathbf{W}^\top (\mathbf{x}_0 - \mathbf{x}'_0) = \mathbf{W}^\top \mathbf{V}\boldsymbol{\alpha}$ for a certain $\boldsymbol{\alpha} \in \mathbb{R}^{n-k}$. But $\mathbf{W}^\top \mathbf{V} = \mathbf{0}$, implying that $\mathbf{W}^\top \mathbf{V}\boldsymbol{\alpha} = \mathbf{0}$ whatever the $\boldsymbol{\alpha}$ is.

Table 1: Summary of representations for linear and affine subspaces

	basis representation	normal representation
linear subspace L with basis \mathbf{V} and normal basis \mathbf{W} ($\dim(L) = k$)	$\{\mathbf{V}\boldsymbol{\alpha} : \boldsymbol{\alpha} \in \mathbb{R}^k\} = \text{col}(\mathbf{V})$	$\{\mathbf{x} : \mathbf{W}^\top \mathbf{x} = \mathbf{0}\} = \text{null}(\mathbf{W}^\top)$
affine subspace L with basis \mathbf{V} , point $\mathbf{x}_0 \in L$, and normal basis \mathbf{W} ($\dim(L) = k$)	$\{\mathbf{x}_0 + \mathbf{V}\boldsymbol{\alpha} : \boldsymbol{\alpha} \in \mathbb{R}^k\} = \mathbf{x}_0 + \text{col}(\mathbf{V})$	$\{\mathbf{x} : \mathbf{W}^\top \mathbf{x} = \mathbf{W}^\top \mathbf{x}_0\} = \mathbf{x}_0 + \text{null}(\mathbf{W}^\top)$

As expected, *affine hyperplanes* are affine subspaces with dimension one less the ambient dimension. Of special interest in this case is the normal representation

$$\{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{w}, \mathbf{x} \rangle = \langle \mathbf{w}, \mathbf{x}_0 \rangle\}. \quad (3.6)$$

Of course, linear hyperplane takes the form $\{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{w}, \mathbf{x} \rangle = 0\}$.

Obviously, an affine subspaces can be a linear subspace in our definition. So henceforth, *subspaces are defaulted to affine subspaces, unless the word “linear” is appended; similarly for hyperplanes*. We are now ready for studying linear classification.

4 Linear classification

We focus on binary classification: $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{1, -1\}$ for $i \in [N]$. A first idea is to use a linear function to map any input \mathbf{x} to either 1 or -1 , but this is unrealistic as the output lies in a continuum

of values. To get around, we pass the output through the simple sign (\cdot) function which naturally takes value in $\{1, -1\}$.⁶ This leads to the hypothesis class

$$\mathcal{H}_H = \left\{ \mathbf{x} \mapsto \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b) : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \right\}, \quad (4.1)$$

which is the set of all hyperplanes.

A natural objective is to find a tuple (\mathbf{w}_0, b_0) so that

$$\text{sign}(\langle \mathbf{w}_0, \mathbf{x}_i \rangle + b_0) = y_i \iff y_i(\langle \mathbf{w}_0, \mathbf{x}_i \rangle + b_0) > 0 \forall i \in [N]. \quad (4.2)$$

For convenience, we use the homogeneous representation again with abuse of the notation, and the goal is to:

$$\text{find } \mathbf{w} \in \mathbb{R}^{d+1} \text{ s. t. } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 0 \forall i \in [N]. \quad (4.3)$$

This is a *feasibility problem* in optimization. The training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ is said to be *linearly separable* if there exists a $\mathbf{w}_0 \in \mathbb{R}^{d+1}$ that solves problem (4.3), i.e., satisfies $y_i(\langle \mathbf{w}_0, \mathbf{x}_i \rangle) > 0$ for all $i \in [N]$.

4.1 Classic solution: Perceptron

Perceptron is a classic algorithm designed for solving problem (4.3). Invented by Frank Rosenblatt in 1958, it helped to fuel the first wave of excitement about neural networks around 60's, but later on also helped to kill the excitement and cause major setbacks for neural networks research due to the famous 1969 book [MM17] that elucidates the limitations of Perceptron. Nonetheless, Perceptron is a critical milestone in the development of binary classification as well as online learning algorithms. We describe only the binary classification aspect here.

Algorithm 3 The Perceptron algorithm for binary classification

Input: training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, initialization $\mathbf{w}^{(0)} = \mathbf{0}$, $k = 1$

- 1: **while** $\exists i$ s. t. $y_i \langle \mathbf{w}^{(k-1)}, \mathbf{x}_i \rangle \leq 0$ **do**
 - 2: update the estimate: $\mathbf{w}^{(k)} = \mathbf{w}^{(k-1)} + y_i \mathbf{x}_i$
 - 3: update the counter: $k = k + 1$
 - 4: **end while**
-

To get a sense why the update step is sensible, note that for an $i \in N$ with $y_i \langle \mathbf{w}^{(k-1)}, \mathbf{x}_i \rangle \leq 0$, after the update,

$$y_i \langle \mathbf{w}^{(k)}, \mathbf{x}_i \rangle = y_i \langle \mathbf{w}^{(k-1)} + y_i \mathbf{x}_i, \mathbf{x}_i \rangle \quad (4.4)$$

$$= y_i \langle \mathbf{w}^{(k-1)}, \mathbf{x}_i \rangle + \|\mathbf{x}_i\|_2^2 \quad (4.5)$$

$$> y_i \langle \mathbf{w}^{(k-1)}, \mathbf{x}_i \rangle, \quad (4.6)$$

i.e., the value moves toward positive and we are making progress. Here, we assume $\|\mathbf{x}_i\|_2 > 0$ for all $i \in [N]$. The overall convergence behavior of the Perceptron algorithm is captured by the following theorem.

⁶We define the sign function as $\text{sign}(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$.

Theorem 4.1. Assume the training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ is linearly separable, and define the data radius $R \doteq \max_{i \in [N]} \|\mathbf{x}_i\|_2$ and the margin parameter⁷ $M = \min \{\|\mathbf{w}\|_2 : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \forall i \in [N]\}$. Then, the Perceptron algorithm will take at most $(RM)^2$ steps to find a feasible \mathbf{w} , when it stops.

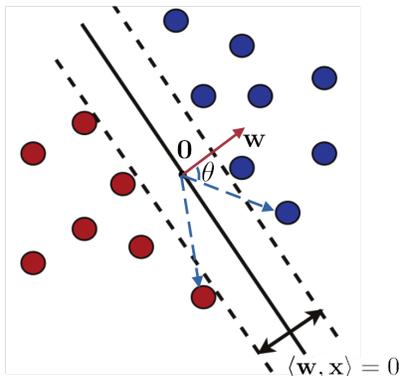


Figure 2: Separating hyperplane in binary classification and interpretation of the Perceptron convergence theorem.

So there are at least two limitations of Perceptron: 1) it cannot deal with linearly non-separable training data, and the algorithm will not even stop in those scenarios; 2) even if the training data are separable, the running time can be exponential in the worst case. Below, we introduce two modern methods that address the limitations of Perceptron.

4.2 Modern solution I: plain SVM

One solution is to reformulate problem (4.3) into a form that is amenable to numerical optimization methods. Iterative methods often work with *compact* constraint sets so that convergence could be established, where none of the inequalities is strict. Note that for any \mathbf{w}_0 satisfying $\langle \mathbf{w}_0, \mathbf{x}_i \rangle > 0 \forall i \in [N]$, there exists an $\eta > 0$ so that $\langle \mathbf{w}_0, \mathbf{x}_i \rangle \geq \eta \forall i \in [N]$. So there exists a $\lambda > 0$ so that $\langle \lambda \mathbf{w}_0, \mathbf{x}_i \rangle \geq 1 \forall i \in [N]$. So problem (4.3) is equivalent to

$$\text{find } \mathbf{w} \in \mathbb{R}^{d+1} \text{ s. t. } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \forall i \in [N]. \quad (4.8)$$

For any feasible \mathbf{w}_0 for problem (4.8), obviously $\lambda \mathbf{w}_0$ for all $\lambda > 1$ is also feasible. We can further refine the formulation by controlling the magnitude of \mathbf{w} , say using

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2 \text{ s. t. } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \forall i \in [N]. \quad (4.9)$$

This is the homogeneous form of *hard-margin SVM* that we will discuss later. Obviously, the ℓ_2 norm we use here is arbitrary; in principle, one can use any function that is monotonically increasing in the “magnitude” of \mathbf{w} , e.g., all vector norms.

When the training set is not linearly separable, there is no feasible solution for problem (4.8). One can add in controlled slackness to allow slight constraint violation, e.g., via

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2 + C \sum_{i=1}^N \xi_i \text{ s. t. } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i, \forall i \in [N]. \quad (4.10)$$

⁷This will become clear when we talk about the max-margin aspect of SVMs.

Here the changed lower bounds $1 - \xi_i$ in the constraints introduce slackness, and the term $\sum_{i=1}^N \xi_i$ in the objective controls the size of the slackness. This is the homogeneous form of *soft-margin SVM*.

Both problems in Eqs. (4.9) and (4.10) are convex (quadratic) optimization problems and can be solved efficiently.

4.3 Modern solution II: logistic regression

First note that the *logistic function* (see Fig. 3)

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (4.11)$$

maps any input naturally into the $[0, 1]$ interval. So

$$\psi(z) \doteq 2\phi(z) - 1 \in [-1, 1], \quad (4.12)$$

and $\psi(z) \rightarrow 1$ when $z \rightarrow \infty$ and $\psi(z) \rightarrow -1$ when $z \rightarrow -\infty$. We can therefore choose the hypothesis class

$$\mathcal{H}_{\text{logistic}} = \left\{ \mathbf{x} \mapsto 2\phi(\langle \mathbf{w}, \mathbf{x} \rangle) - 1 : \mathbf{w} \in \mathbb{R}^{d+1} \right\}. \quad (4.13)$$

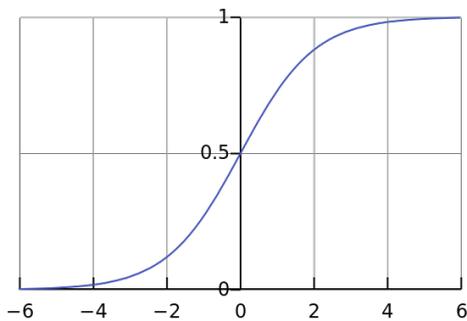


Figure 3: Logistic function. Figure courtesy of Wikipedia.

To pick an appropriate loss function, note that: when $y_i = 1$, we hope that $\psi(\langle \mathbf{w}, \mathbf{x}_i \rangle)$ is as large as possible, or $1 + \exp(-\langle \mathbf{w}, \mathbf{x}_i \rangle)$ is as small as possible; when $y_i = -1$, we hope that $\psi(\langle \mathbf{w}, \mathbf{x}_i \rangle)$ is as small as possible, or $1 + \exp(-\langle \mathbf{w}, \mathbf{x}_i \rangle)$ is as large as possible. Combining the two, we hope to make

$$1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \quad (4.14)$$

as small as possible for all $i \in [N]$. So one can formulate an optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)). \quad (4.15)$$

Although this is a convex problem, the exponential term may cause some numerical issues as the exponent $-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$ may get large for certain i 's. Because making $1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$ small is equivalent to making $\log(1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle))$ small, we arrive at our *logistic regression* formulation:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)), \quad (4.16)$$

which is numerically more stable, as the logarithm function counters the possible blowup of the exponential term. One can easily verify Eq. (4.16) is also a convex problem by checking the Hessian.

Logistic regression can also be derived from the maximum likelihood principle, as we will explore in homework.

Further reading

Main reference is Chapter 9 of [SSS14]. Chapters 2–4 of [HTF09] are good supplements. [NW06] is a highly recommended reference for numerical optimization.

Disclaimer

This set of notes is preliminary and has not been thoroughly proofread. Typos and factual errors are well expected and hence use it with caution. Bug reports are very welcome and should be sent to Prof. Ju Sun via jusun@umn.edu.

A Proof of Theorem 4.1

Proof. Since the requirement $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 0$ for all $i \in [N]$ is homogenous in \mathbf{w} , we care only about the direction $\frac{\mathbf{w}}{\|\mathbf{w}\|_2}$ instead of \mathbf{w} itself. Suppose \mathbf{w}_* satisfies $\langle \mathbf{w}_*, \mathbf{x}_i \rangle \geq 1$ for all $i \in [N]$ with $\|\mathbf{w}_*\|_2 = B$. We want to show that $\mathbf{w}^{(T)}$ aligns with \mathbf{w}_* , i.e.,

$$\frac{\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle}{\|\mathbf{w}_*\|_2 \|\mathbf{w}^{(T)}\|_2} = 1 \quad (\text{A.1})$$

when T is large enough.

First, for any k ,

$$\begin{aligned} \|\mathbf{w}^{(k)}\|_2^2 &= \|\mathbf{w}^{(k-1)} + y_i \mathbf{x}_i\|_2^2 = \|\mathbf{w}^{(k-1)}\|_2^2 + \underbrace{\|\mathbf{x}_i\|_2^2}_{\leq R^2} + \underbrace{2y_i \langle \mathbf{w}^{(k-1)}, \mathbf{x}_i \rangle}_{\leq 0} \\ &\leq \|\mathbf{w}^{(k-1)}\|_2^2 + R^2 \implies \|\mathbf{w}^{(k)}\|_2^2 - \|\mathbf{w}^{(k-1)}\|_2^2 \leq R^2. \end{aligned} \quad (\text{A.2})$$

Using telescoping summation, we obtain

$$\|\mathbf{w}^{(T)}\|_2^2 = \sum_{k=1}^T \left(\|\mathbf{w}^{(k)}\|_2^2 - \|\mathbf{w}^{(k-1)}\|_2^2 \right) \leq TR^2. \quad (\text{A.3})$$

On the other hand, for any k ,

$$\langle \mathbf{w}_*, \mathbf{w}^k \rangle - \langle \mathbf{w}_*, \mathbf{w}^{k-1} \rangle = \langle \mathbf{w}_*, \mathbf{w}^k - \mathbf{w}^{k-1} \rangle = y_i \langle \mathbf{w}_*, \mathbf{x}_i \rangle \geq 1. \quad (\text{A.4})$$

Applying telescoping summation again, we obtain

$$\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle = \sum_{k=1}^T \left(\langle \mathbf{w}_*, \mathbf{w}^{(k)} \rangle - \langle \mathbf{w}_*, \mathbf{w}^{(k-1)} \rangle \right) \geq T. \quad (\text{A.5})$$

Combining Eqs. (A.3) and (A.5), we finally obtain

$$\frac{\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle}{\|\mathbf{w}_*\|_2 \|\mathbf{w}^{(T)}\|_2} \geq \frac{T}{B\sqrt{TR}} = \frac{\sqrt{T}}{BR}. \quad (\text{A.6})$$

When $T \geq (RB)^2$,

$$\frac{\sqrt{T}}{BR} \geq 1 \implies \frac{\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle}{\|\mathbf{w}_*\|_2 \|\mathbf{w}^{(T)}\|_2} \geq 1. \quad (\text{A.7})$$

But $\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle \leq \|\mathbf{w}_*\|_2 \|\mathbf{w}^{(T)}\|_2$ due to the Cauchy-Schwarz inequality, implying that $\frac{\langle \mathbf{w}_*, \mathbf{w}^{(T)} \rangle}{\|\mathbf{w}_*\|_2 \|\mathbf{w}^{(T)}\|_2} \leq 1$. Thus, it takes at most $T = (RB)^2$ iterations to attain Eq. (A.1), completing the proof. ■

B Why backtracking line search?

To simplify the notation, suppose that the current iterate is \mathbf{z} and so the gradient is $\nabla f(\mathbf{z})$. We want to choose a step size t so that $f(\mathbf{z} - t\nabla f(\mathbf{z})) - f(\mathbf{z})$ is as negative as possible in order to minimize the objective fast. Since we assume that f is first-order differentiable, Taylor's theorem tells us

$$f(\mathbf{z} - t\nabla f(\mathbf{z})) - f(\mathbf{z}) = -t\|\nabla f(\mathbf{z})\|_2^2 + o(t\|\nabla f(\mathbf{z})\|_2) \quad (\text{B.1})$$

as $t \rightarrow 0$. Now the linear term $-t\|\nabla f(\mathbf{z})\|_2^2$ is negative, and the lower-order term $o(t\|\nabla f(\mathbf{z})\|_2)$ may be either positive or negative. In any case, when $t > 0$ is sufficiently small, $-t\|\nabla f(\mathbf{z})\|_2^2$ will dominate the right side of Eq. (B.1), and we can reach a level so that

$$-t\|\nabla f(\mathbf{z})\|_2^2 + o(t\|\nabla f(\mathbf{z})\|_2) \leq -\eta t\|\nabla f(\mathbf{z})\|_2^2 \quad (\text{B.2})$$

for a pre-fixed $\eta \in (0, 1)$. Of course, any smaller t still satisfies this. But since we hope to set t to be the largest possible, our line search is backward: we start with a large t and gradually decreases it whenever Eq. (B.2) is violated.

References

- [HTF09] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The elements of statistical learning*, Springer New York, 2009.
- [MM17] Seymour A. Papert Marvin Minsky, *Perceptrons, reissue of the 1988 expanded edition with a new foreword by léon bottou*, MIT Press Ltd, 2017.
- [NW06] Jorge Nocedal and Stephen J. Wright, *Numerical optimization*, 2 ed., Springer New York, 2006.
- [Sra12] Suvrit Sra, *Optimization for machine learning*, MIT Press, Cambridge, Mass, 2012.
- [SSS14] Shai Ben-David Shai Shalev-Shwartz, *Understanding machine learning*, Cambridge University Press, 2014.
- [Sta] *Stanford EE364b - Convex Optimization II*, <http://stanford.edu/class/ee364b/>.
- [Str16] Gilbert Strang, *Introduction to linear algebra*, 5 ed., Cambridge Press, Wellesley, MA, 2016.