

EE 127

Convex Optimization

Notes

Druv Pai

Contents

1	Logistics and Overview	3	5	Dealing With Noise	18
2	Least Squares	3	6	Convexity	22
3	Vector Calculus	5	7	Optimization Algorithms	26
4	Linear Algebra Review	6	8	Duality	29

1 Logistics and Overview

Office hours are 6:30-7:30 after class on Tuesday. Homeworks and self-grades are due on Friday at 11 PM.

Example 1. Say that you have 10^5 barrels of crude oil. You can process each barrel of crude oil into a barrel of jet fuel, which costs 10\$ or a barrel of gasoline, which costs 20\$. In an ideal world, you build all the gasoline possible. But we don't live in such a world, so we add some constraints, for example that we need to produce at least 10^3 barrels of jet fuel and $2 \cdot 10^3$ barrels of gasoline. Further, we can transport at most $1.8 \cdot 10^6$ barrel-miles. The gasoline distributor is 30 miles away, and the jet fuel distributor is 10 miles away. How do we maximize profit?

Solution. Let x_1 be the number of gallons of jet fuel and x_2 be the number of gallons of gasoline. The task resolves to

$$\begin{aligned} & \text{maximize } 10x_1 + 20x_2 \\ & \text{subject to } x_1 \geq 10^3 \\ & \quad x_2 \geq 2 \cdot 10^3 \\ & \quad 10x_1 + 30x_2 \leq 1.8 \cdot 10^6 \end{aligned}$$

This is a linear program, which is one of the problems that we will solve a lot during the latter part of the course. \square

A minimization program, in general, has an optimization variable $x \in \mathbb{R}^n$ with optimal value x^* among the feasible set, or set of values for x that satisfy the m constraints placed upon x . The program is of the form

$$\begin{aligned} & \text{minimize } f_0(x) \\ & \text{subject to } f_i(x) \leq b_i, \quad i \in \{1, \dots, m\} \end{aligned}$$

Example 2. Say that the i^{th} class has value α_i and workload β_i . Let $x_i = \mathbb{1}(\text{you take } i^{\text{th}} \text{ class})$. Then if the total workload you can handle is b , the integer linear program (linear program where the decision variable takes only integral values) for the classes you should take is given by

$$\begin{aligned} & \text{maximize } \alpha^\top x \\ & \text{subject to } \beta^\top x \leq b \end{aligned}$$

and solving this gives you a class schedule.

Example 3. Say that the i^{th} class has size x_i , credits c_i , and resources r_i . Then to maximize credit hours for students subject to a resource budget b , the integer linear program is

$$\begin{aligned} & \text{maximize } c^\top x \\ & \text{subject to } r^\top x \leq b \end{aligned}$$

These are examples of generic programs that we will solve by the end of the course.

2 Least Squares

The least squares problem is one of the most basic optimization problems. Given a matrix A and vector b (or appropriately sized tensors), the least squares problem finds the projection of b onto the span of the

columns of A :

$$x^* = \operatorname{argmin}_x \|Ax - b\|^2$$

Example 4. Say we have a set of data points $\{(x_i, y_i)\}_{i=1}^m$, where the noiseless relationship is $y = \beta_1 x + \beta_2$ and, we wish to find β . In matrix-vector form, we have

$$\begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

More simply, if x is a vector of the x_i , and y is a vector of the y_i ,

$$\begin{bmatrix} x & 1 \end{bmatrix} \beta = y$$

This leads to a quadratic cost, minimizable in β :

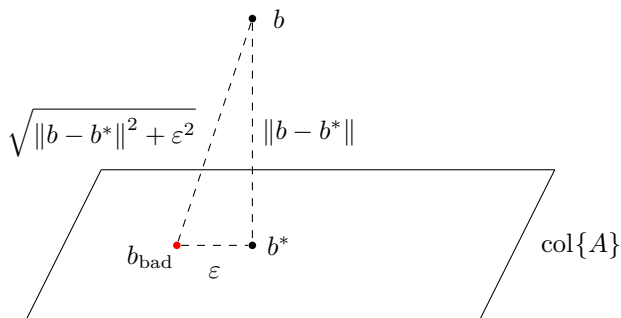
$$\beta^* = \operatorname{argmin}_{\beta_1, \beta_2} \sum_{i=1}^m [y_i - (\beta_1 x_i + \beta_2)]^2$$

which we can minimize in closed form, or numerical methods if m is large.

How do we actually find the quantity

$$x^* = \operatorname{argmin}_x \|Ax - b\|^2$$

One way is to solve it geometrically. Consider the subspace $\operatorname{col}(A)$; then b may be on $\operatorname{col}(A)$, or it may not be. If $b \in \operatorname{col}(A)$, then we are done; if not, then the closest point is $b^* = \operatorname{proj}_{\operatorname{col}(A)}(b)$, and other points are off by an appropriate amount.



By triangle geometry, we require $b - b^* = b - Ax^*$ to be orthogonal to the columns of A , so we require $Ax^* - b$ to be orthogonal to the columns of A , so $A^T(Ax^* - b) = 0$. Then

$$\begin{aligned} 0 &= A^T(Ax^* - b) \\ &= A^T Ax^* - A^T b \\ x^* &= (A^T A)^{-1} A^T b \end{aligned}$$

Of course, this requires $A^T A$ to be invertible, so A must be full rank. If not, then the relevant equation is $A^T Ax^* = A^T b$; one can show by rank-nullity that $A^T b \in \operatorname{range}(A^T A)$, so computing $A^T b$ and then row-reducing finds the solution.

The cost function for the least squares problem is quadratic; quadratic functions are convex, and thus all of their local minima are global minima. This is very nice because we immediately know how to minimize

such functions extremely well, often in closed form.

3 Vector Calculus

We begin with scalar-valued Taylor's theorem.

Theorem 5 (Taylor's Theorem in $\mathbb{R} \rightarrow \mathbb{R}$ Case)

Let $f: \mathbb{R} \rightarrow \mathbb{R}$. Then

$$f(x + \varepsilon) = \sum_{n=0}^{\infty} \frac{d^n f(x)}{dx^n} \frac{\varepsilon^n}{n!} = f(x) + \frac{df(x)}{dx} \varepsilon + \frac{1}{2} \frac{d^2 f(x)}{dx^2} \varepsilon^2 + \dots$$

in the case that f is smooth; there exists a remainder term $R_k(x)$ such that $|R_k(x)| \leq \left| \frac{\varepsilon^{k+1}}{(k+1)!} \frac{d^{k+1} f(x)}{dx^{k+1}} \right|$ otherwise.

Now we cover the multivariable case. First, some notation. Formally, if $f: \mathbb{R}^n \rightarrow \mathbb{R}$, then $\frac{\partial f(x)}{\partial x} = \left[\frac{\partial f(x)}{\partial x_1} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right]$ and $\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$. If $f(x)$ is scalar valued, then $\nabla_x f(x)$ has the same shape as x .

Now, we cover the second derivative. In the case that $f: \mathbb{R}^n \rightarrow \mathbb{R}$, the second derivative $\nabla_x^2 f(x)$ is called the Hessian. In particular, $\left[\nabla_x^2 f(x) \right]_{i,j} = \frac{\partial}{\partial x_j} (\nabla_x f(x))_i = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$.

Theorem 6 (Taylor's Theorem in $\mathbb{R}^n \rightarrow \mathbb{R}$ Case)

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$. Then

$$\begin{aligned} f(x + \varepsilon) &= f(x) + \frac{\partial f}{\partial x} \varepsilon + \frac{1}{2} \varepsilon^\top \nabla_x^2 f(x) \varepsilon + \dots \\ &= f(x) + (\nabla_x f(x))^\top \varepsilon + \frac{1}{2} \varepsilon^\top \nabla_x^2 f(x) \varepsilon + \dots \\ &= f(x) + \langle \nabla_x f(x), \varepsilon \rangle + \frac{1}{2} \varepsilon^\top \nabla_x^2 f(x) \varepsilon + \dots \end{aligned}$$

Derivatives at order $n \geq 3$ or more involve higher-order tensors and we won't cover those now.

Definition 7 (Level Set). Define the c -level set of f as $L_c(f) = \{x \mid f(x) = c\}$.

Since the gradient is sometimes interpreted as the direction of fastest movement or steepest ascent, we have that (by the definition of the directional derivative as $D_a f(x) = \langle \nabla_x f(x), a \rangle$, then for all c $[\nabla_x f(x)]_{x: f(x)=c} \perp L_c(f)$.

Example 8. Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be defined as $f(x) = x^\top x$. We have $\nabla_x f(x) = 2x$ manually, and $\nabla_x^2 f(x) = 2I$. Now we can compute it via Taylor's theorem:

$$\begin{aligned} f(x + \varepsilon) &= (x + \varepsilon)^\top (x + \varepsilon) \\ &= x^\top x + x^\top \varepsilon + \varepsilon^\top x + \varepsilon^\top \varepsilon \end{aligned}$$

$$= x^\top x + (2x)^\top \varepsilon + \frac{1}{2} (2\varepsilon^\top \varepsilon)$$

We have that $\nabla_x f(x) = 2x$ and $\nabla_x^2 = 2I$ by pattern matching.

Example 9. If $f(x) = x^\top a$, where $x, a \in \mathbb{R}^n$, then

$$f(x + \varepsilon) = (x + \varepsilon)^\top a = x^\top a + \varepsilon^\top a = f(x) + a^\top \varepsilon$$

so $\nabla_x f(x) = a$.

Example 10. If $f(x) = x^\top A x$ then

$$\begin{aligned} f(x + \varepsilon) &= (x + \varepsilon)^\top A (x + \varepsilon) \\ &= (x^\top + \varepsilon^\top) (A x + A \varepsilon) \\ &= x^\top A x + \varepsilon^\top A x + x^\top A \varepsilon + \varepsilon^\top A \varepsilon \\ &= x^\top A x + x^\top A^\top \varepsilon + x^\top A \varepsilon + \varepsilon^\top A \varepsilon \\ &= x^\top A x + x^\top (A^\top + A) \varepsilon + \varepsilon^\top A \varepsilon \end{aligned}$$

Therefore we have that $\nabla_x f(x) = (A + A^\top) x$ and $\nabla_x^2 f(x) = 2A$.

Example 11 (Least Squares). We want to find $x^* = \operatorname{argmin}_x \|Ax - b\|_2^2$. Let

$$f(x) = \|Ax - b\|_2^2 = (Ax - b)^\top (Ax - b) = x^\top A^\top A x - x^\top A^\top b - b^\top A x + b^\top b.$$

Using the previous example, $\nabla_x f(x) = 2(A^\top A x - A^\top b)$. Setting this to 0, we have $x^* = (A^\top A)^{-1} A^\top b$ which agrees.

4 Linear Algebra Review

Definition 12 (Vector). A vector x is an element of a vector space $X(\mathbb{F})$, which is a group that is closed under an addition operation, as well as another operation which multiplies a vector x by a member a of its component field \mathbb{F} .

Norms

We now discuss norms. Let X be a normed vector space, which is a vector space equipped with a norm.

Definition 13 (Norm). A function $f: X \rightarrow \mathbb{F}$ is a norm if

- $f(x) \geq 0$ for all $x \in X$
- $f(x) = 0$ if and only if $x = 0$
- $f(ax) = a f(x)$ for all $x \in X$ and $a \in \mathbb{F}$.
- $f(x + y) \leq f(x) + f(y)$ for all $x, y \in X$ (triangle inequality).

One norm is the 2-norm on real sequences $x = (x_i)_{i=1}^n$, which is defined as

$$\|x\|_2 \stackrel{\text{def}}{=} \left(\sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}}$$

Definition 14. The ℓ^p norm on real sequences, defined as

$$\|x\|_p \stackrel{\text{def}}{=} \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

For $p \in [1, \infty)$, the p -norm is a proper norm (for $p \in [0, 1)$, the p -norm is what's called a seminorm).

When $p = 2$, we recover the Euclidean norm. When $p = 1$, we obtain $\|x\|_1 = \sum_{i=1}^n |x_i|$; when $p = 0$, we obtain $\|x\|_0 = n$; when $p \rightarrow \infty$, we obtain $\|x\|_\infty = \max_i |x_i|$.

Theorem 15 (Cauchy Schwarz Inequality)

Let X be a normed inner product space. Then for any $x, y \in X$, we have

$$|\langle x|y \rangle| \leq \|x\|_2 \|y\|_2$$

Proof. We have that

$$|\langle x|y \rangle| = \|x\|_2 \|y\|_2 \cos(\theta) = \|x\|_2 \|y\|_2 |\cos(\theta)| \leq \|x\|_2 \|y\|_2$$

as desired. □

There is a generalization of this for arbitrary norms:

Theorem 16 (Holder's Inequality)

Let X be a normed inner product space. Then for any p and q such that $\frac{1}{p} + \frac{1}{q} = 1$,

$$|\langle x|y \rangle| \leq \sum_{i=1}^n |x_i y_i| \leq \|x\|_p \|y\|_q$$

Proof. The proof will be covered once we have developed the tools of convexity.

Note that gradient proofs don't work because the functions on the left hand side are not differentiable. □

Example 17. Say that we have some $y \in \mathbb{R}^n$. We want to find

$$\max_{\|x\|_p \leq 1} \langle x|y \rangle$$

In the $p = 2$ case, we obtain $x^* = \frac{y}{\|y\|_2}$; by Theorem 15, we have that this is optimal.

In the $p \rightarrow \infty$ case, we pick $x^* = \text{sign}(y)$, so $\max_{\|x\|_p \leq 1} \langle x|y \rangle = \sum_{i=1}^n |y_i| = \|y\|_1$.

In the $p = 1$ case, we pick x^* such that $x_i^* = \mathbb{1}(|y_i| = \max_j |y_j|)$, so $\max_{\|x\|_p \leq 1} \langle x|y \rangle = \max_j |y_j| = \|y\|_\infty$.

It turns out that $\|\cdot\|_1$ and $\|\cdot\|_\infty$ are dual norms, and $\|\cdot\|_2$ is a self-dual norm.

Orthogonality and Orthonormalization

The Gram-Schmidt process converts a basis for a vector space into an orthonormal basis for the same vector space.

Algorithm 1 The Gram-Schmidt process.

Input: A basis $\{a_1, \dots, a_n\}$ for a vector space V .

Output: An orthonormal basis $\{q_1, \dots, q_n\}$ for V .

```

 $q_1 \leftarrow \frac{a_1}{\|a_1\|}$  ▷ Normalize  $a_1$ .
for  $i \in \{2, \dots, n-1\}$  do
     $p_i \leftarrow \sum_{j=1}^{i-1} q_j \langle a_i | q_j \rangle$  ▷ Project  $a_i$  onto  $\text{span}(q_1, \dots, q_{i-1})$ .
     $s_i \leftarrow a_i - p_i$  ▷ Subtract out the projection.
     $q_i \leftarrow \frac{s_i}{\|s_i\|_2}$  ▷ Normalize  $p_i$ .
return  $\{q_1, \dots, q_n\}$ 

```

We introduce the QR decomposition, where any full-rank matrix A can be expressed as $A = QR$, where Q is an orthogonal matrix and R is an upper-triangular matrix. The construction is

$$Q_{i,j} = q_{i,j}, \quad R_{i,j} = \langle q_i | a_j \rangle \text{ for } i \geq j \text{ or } \|s_i\|_2 \text{ for } i = j$$

One may also write $A = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$ since only a few components of the orthogonal matrix are used.

Claim 18. Let X be an inner product space and S be a subspace of V . Let $x \in X$. Then x can be written uniquely as the sum of $s \in S$ and $s' \in S^\perp$, where S^\perp is the orthogonal complement of S , i.e. $S^\perp = \{s' : \langle s' | s \rangle = 0 \forall s \in S\}$. Another way to write this is that $X = S \oplus S^\perp$.

Proof. We want to show that $S \oplus S^\perp = X$. Clearly $S \cap S^\perp = \{0\}$, and $S + S^\perp \subseteq X$. Assume that $W \subset X$ for the sake of contradiction. Let $\{w_n\}$ be an orthonormal basis for W , and extend this basis to X (say perhaps via Gram-Schmidt, or in the infinite dimensional case using Hahn-Banach extension theorem). Then there exists $x \notin W$ a new basis vector, and by Gram-Schmidt $x \perp W$, so $x \perp S$, and $x \perp S^\perp$, a contradiction.

It remains to show uniqueness. Consider $x_1, x_2 \in S$ and $y_1, y_2 \in S^\perp$, where $x_1 + y_1 = x_2 + y_2$ but $x_1 \neq x_2$ and $y_1 \neq y_2$. Then $x_1 - x_2 = y_2 - y_1$. Clearly $x_1 - x_2 \in S$ and $y_2 - y_1 \in S^\perp$, so $x_1 - x_2 = y_2 - y_1 = 0$, a contradiction. \square

Theorem 19 (Fundamental Theorem of Linear Algebra)

We have that

$$\text{range}(A) = \text{null}(A^\top)^\perp$$

and

$$\dim(\text{range}(A)) + \dim(\text{null}(A)) = \text{rank}(A) + \text{nullity}(A) = n$$

for $A \in \mathbb{R}^{m \times n}$.

Principal Component Analysis

Principal component analysis is a technique for dimensionality reduction. Say that we have unlabeled data points $\{x_i\}_{i=1}^n \in V$, where V is a normed inner product space, and has $\dim(V) = p$ which is considered “high-dimensional”. Further assume that the x_i are drawn from a distribution with 0 mean, that is, $x_1, \dots, x_n \sim p_x(\cdot)$ where $\text{Ex}[x] = 0$. Our goal is to find projections of x_i onto some subspace $W \subset V$, where

W is a normed inner product space, $\dim(W) = k$ where $k < p$, where $\text{proj}_W(x_i)$ is as close to x_i as possible, across all the x_i .

Let $k = 1$. The goal is to determine a w with $\|w\|_2^2 = 1$ such that the projections onto $\text{span}(w)$ are as close to the original vectors as possible. These projections are $\{\langle x_i | w \rangle\}_{i=1}^n$, with error norms $e_i^2 = \|x_i - \langle w | x_i \rangle w\|_2^2$, and average projection error $\frac{1}{n} \sum_{i=1}^n e_i^2$. Noting that $\text{Ex}[x] = 0$, we simplify the projection error to become

$$\begin{aligned} \|x_i - \langle w | x_i \rangle w\|_2^2 &= (x_i - \langle w | x_i \rangle w)^\top (x_i - \langle w | x_i \rangle w) \\ &= \|x_i\|_2^2 - 2 \langle w | x_i \rangle^2 + \langle w | x_i \rangle^2 \|w\|_2^2 \\ &= \|x_i\|_2^2 - \langle w | x_i \rangle^2 \\ e_i^2 &= \|x_i\|_2^2 - \langle w | x_i \rangle^2 \end{aligned}$$

The mean square projection error is

$$\text{MSE}(w) = \frac{1}{n} \sum_{i=1}^n \|x_i\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle w | x_i \rangle^2$$

The goal is to minimize the mean square projection error, and we can see that in expectation this is equivalent to the problem:

$$\begin{aligned} \text{Ex} \left[\min_w \text{MSE}(w) \right] &= \min_w \text{Ex}[\text{MSE}(w)] \\ &= \max_w \text{Ex} \left[\frac{1}{n} \sum_{i=1}^n \langle w | x_i \rangle^2 \right] \\ &= \max_w \text{Ex} [\langle w | x \rangle^2] \\ &= \max_w \text{Ex} [\langle w | x \rangle^2] + \text{Var}[\langle w | x \rangle] \\ &= \max_w \text{Var}[\langle w | x \rangle] \end{aligned}$$

Thus the principal component analysis maximizes the variance of w with the distribution p_X .

Another way to look at principal component analysis (and indeed, the way to consider an instantiation of principal component analysis, where we cannot work in expectation) is to note that, if $X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$ is the data matrix, we are looking for

$$\begin{aligned} \min_w \text{MSE}(w) &= \max_w \frac{1}{n} \sum_{i=1}^n \langle w | x_i \rangle^2 \\ &= \frac{1}{n} \|Xw\|_2^2 \\ &= \frac{1}{n} (Xw)^\top (Xw) \\ &= \frac{1}{n} w^\top X^\top X w \\ &= w^\top \left(\frac{X^\top X}{n} \right) w \end{aligned}$$

where $C = \left(\frac{X^\top X}{n} \right)$ is the covariance matrix of the data. Then the problem is to find $\max_{\|w\|_2=1} w^\top C w$.

This motivates the discussion of symmetric matrices.

Symmetric Matrices and Positive Semidefiniteness

Definition 20 (Symmetric Matrix). A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is a matrix A such that $A^\top = A$, or $a_{i,j} = a_{j,i}$ for each pair (i, j) .

The set of symmetric matrices is sometimes called \mathcal{S} .

Definition 21 (Diagonalizability). Let A be square and have characteristic polynomial $p(\lambda) = \det(A - \lambda I) = \prod_{i=1}^n (\lambda - \lambda_i)^{n_i}$. Then for each λ_i , we have $n_i = \dim(\text{null}(A - \lambda_i I))$.

Theorem 22 (Spectral Theorem for Real Symmetric Matrices)

Let $A \in \mathbb{R}^{n \times n}$ be a real symmetric matrix, with eigenvalues $\lambda_1, \dots, \lambda_k$ which themselves have algebraic multiplicities μ_1, \dots, μ_k . Further, define $\phi_i = \text{null}(A - \lambda_i I)$. Then for all i and $j \neq i$, we have

- $\lambda_i \in \mathbb{R}$,
- $\phi_i \perp \phi_j$,
- $\dim(\phi_i) = \mu_i$

Therefore, there exists orthogonal U and diagonal Λ such that $A = U\Lambda U^\top$.

Proof. The first two claims are trivial, so we prove the third. Let $U = \begin{bmatrix} u & U_1 \end{bmatrix}$, where we may obtain U by Gram-Schmidt. Then $U^\top A U = \begin{bmatrix} u^\top \\ U_1^\top \end{bmatrix} A \begin{bmatrix} u & U_1 \end{bmatrix} = \begin{bmatrix} u^\top \\ U_1^\top \end{bmatrix} \begin{bmatrix} \lambda u & A U_1 \end{bmatrix} = \begin{bmatrix} \lambda & 0 \\ 0 & U_1^\top A U_1 \end{bmatrix}$. If $B = U_1^\top A U_1$, then B is symmetric, and we are done by induction. \square

Definition 23 (Rayleigh Quotient). For a symmetric matrix $A \in \mathbb{R}^{n \times n}$, we define the Rayleigh quotient with respect to x as

$$\frac{x^\top A x}{x^\top x}$$

Theorem 24

We have that

$$\lambda_{\min}(A) \leq \frac{x^\top A x}{x^\top x} \leq \lambda_{\max}(A)$$

for all $x \in \mathbb{R}^n \setminus \{0\}$.

Proof. Write $A = U_A \Lambda_A U_A^\top$. Then

$$x^\top A x = x^\top U_A \Lambda_A U_A^\top x$$

Write $y = U_A^\top x$. Then

$$\begin{aligned} x^\top A x &= y^\top \Lambda_A y \\ &= \sum_{i=1}^n \lambda_i(A) y_i^2 \\ &\leq \lambda_{\max}(A) \sum_{i=1}^n y_i^2 \end{aligned}$$

$$\geq \lambda_{\min}(A) \sum_{i=1}^n y_i^2$$

Also note that $\sum_{i=1}^n y_i^2 = \|y\|_2^2 = \|U_A^\top x\|_2^2 = \|x\|_2^2$ since U_A is orthonormal. \square

Corollary 25. We have

$$\lambda_{\max}(A) = \max_{\|x\|_2=1} x^\top A x$$

and

$$\lambda_{\min}(A) = \min_{\|x\|_2=1} x^\top A x.$$

Proof. Let $x = v_{\max}$ (the eigenvector associated with $\lambda_{\max}(A)$) to get the first inequality, and $x = v_{\min}$ to get the second inequality; these clearly work when substituted into the proof. \square

Definition 26 (Positive Semidefinite). A matrix $A \in \mathbb{R}^{n \times n}$ is positive semidefinite (PSD) if for all $x \in \mathbb{R}^n$, we have $x^\top A x \geq 0$. We also say $A \succeq 0$ or $A > 0$.

Corollary 27. If A is PSD then $\lambda_i(A) \geq 0$.

Proof. Directly follows from Theorem 24. \square

Definition 28 (Matrix Square Root). If A is PSD then there exists B such that $A = B^\top B$. Indeed, if $A = U_A \Lambda_A U_A^\top$ then $B = U_A \Lambda_A^{1/2}$ (where $\Lambda_{A;i,j}^{1/2} = \sqrt{\Lambda_{A;i,j}}$) suffices, but is not symmetric; one can show that there is a symmetric matrix C such that $A = C^\top C$, but this itself is not unique in the case of A having repeated eigenvalues.

Back to PCA. We showed that our problem was equivalent to

$$w^* = \operatorname{argmax}_{\|w\|_2=1} w^\top C w$$

From 25, we have that w^* is the eigenvector corresponding to $\lambda_{\max}(C)$. Then w is the first principal component.

Singular Value Decomposition

Definition 29 (Singular Value Decomposition). Let $A \in \mathbb{R}^{m \times n}$ with $\operatorname{rank}(A) = r$. Then the singular value decomposition of A is an expression of A in the form

$$A = \sum_{i=1}^r \sigma_i u_i v_i^\top$$

with the convention $\sigma_1 \geq \dots \geq \sigma_r \geq 0$, $u_1, \dots, u_r \in \mathbb{R}^m$ are orthonormal, and $v_1, \dots, v_r \in \mathbb{R}^n$ are also orthonormal.

We also write $A = U_r \Sigma_r V_r^\top$, where $U_r = \begin{bmatrix} u_1 & \dots & u_r \end{bmatrix} \in \mathbb{R}^{m \times r}$, $\Sigma_r = \operatorname{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$, and

$V_r^\top = \begin{bmatrix} v_1^\top \\ \vdots \\ v_r^\top \end{bmatrix} \in \mathbb{R}^{r \times n}$. This is sometimes called the “compact” SVD.

We can compute the singular value decomposition via the following algorithm:

Algorithm 2 Computes the singular value decomposition.

Input: $A \in \mathbb{R}^{m \times n}$,

Output: $U_r \in \mathbb{R}^{m \times r}$ (orthonormal), $\Sigma_r \in \mathbb{R}^{r \times r}$ (diagonal), $V_r^\top \in \mathbb{R}^{r \times n}$ (orthonormal) such that $A = U_r \Sigma_r V_r^\top$.

Compute $A^\top A$.

▷ $A^\top A$ is symmetric.

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0 \leftarrow$ eigenvalues of $A^\top A$.

$v_1, \dots, v_r \leftarrow$ orthonormal eigenvectors of $A^\top A$ (such that $A^\top A v_i = \lambda_i v_i$ for each i).

$\sigma_i \leftarrow \sqrt{\lambda_i}$ for each i .

$u_1, \dots, u_r \leftarrow$ defined by $A v_i = \sigma_i u_i$ for all i .

return $[u_1 \ \dots \ u_r], \text{diag}(\sigma_1, \dots, \sigma_r), [v_1 \ \dots \ v_r]^\top$

Claim 30. The u_i are orthonormal.

Proof. First, we claim that $\|u_i\|_2^2 = 1$ for all i . We have

$$\begin{aligned}
 \|\sigma_i u_i\|_2^2 &= (\sigma_i u_i)^\top (\sigma_i u_i) \\
 &= (A v_i)^\top (A v_i) \\
 &= v_i^\top A^\top A v_i \\
 &= v_i^\top \lambda_i v_i \\
 &= \lambda_i v_i^\top v_i \\
 &= \lambda_i \|v_i\|_2^2 \\
 \sigma_i^2 \|u_i\|_2^2 &= \lambda_i \|v_i\|_2^2 \\
 \lambda_i \|u_i\|_2^2 &= \lambda_i \|v_i\|_2^2 \\
 \|u_i\|_2^2 &= \|v_i\|_2^2 \\
 &= 1
 \end{aligned}$$

so $\|u_i\|_2^2 = 1$ as desired. Now we claim that $\langle u_i | u_j \rangle = 0$ for all $i \neq j$. We have

$$\begin{aligned}
 \langle \sigma_i u_i | \sigma_j u_j \rangle &= (\sigma_i u_i)^\top (\sigma_j u_j) \\
 &= \sigma_i \sigma_j u_j^\top u_i \\
 \langle \sigma_i u_i | \sigma_j u_j \rangle &= \langle A v_i | A v_j \rangle \\
 &= v_j^\top A^\top A v_i \\
 &= v_j^\top \lambda_i v_i \\
 &= \lambda_i v_j^\top v_i \\
 \sigma_i \sigma_j u_j^\top u_i &= \lambda_i v_j^\top v_i \\
 &= 0 \\
 u_j^\top u_i &= 0 \\
 \langle u_i | u_j \rangle &= 0
 \end{aligned}$$

Therefore, the u_i are orthonormal. □

This works because $\lambda_i \geq 0$, so we prove this.

Claim 31. The $\lambda_i \geq 0$.

Proof. If $A^T A x = \lambda x$ then $x^T A^T A x = x^T \lambda x$ so $\|Ax\|_2^2 = \lambda \|x\|_2^2$, so $\lambda \geq 0$ since $\|\cdot\|_2^2 \geq 0$. \square

One may note that we have $AV_r = U_r \Sigma_r$ by construction, so then $AV_r V_r^T = U_r \Sigma_r V_r^T$. However, $\text{rank}(V_r V_r^T) = r \neq n$, so $AV_r V_r^T \neq I$. This is when we extend $V_r \in \mathbb{R}^{n \times r}$ to $V_n \in \mathbb{R}^{n \times n}$ via Gram Schmidt, and thereby develop the non-compact SVD.

Claim 32. A key fact we will use is that $\text{null}(A) = \text{null}(A^T A)$.

Proof. Take a vector $v \in \text{null}(A)$; then $A^T A v = A^T 0 = 0$, so $v \in \text{null}(A^T A)$, so $\text{null}(A) \subseteq \text{null}(A^T A)$. Now take a vector $v \in \text{null}(A^T A)$; then $0 = v^T A^T A v = (Av)^T (Av) = \|Av\|_2^2$, so $Av = 0$, so $v \in \text{null}(A)$, so $\text{null}(A^T A) \subseteq \text{null}(A)$. Therefore $\text{null}(A) = \text{null}(A^T A)$. \square

Now we are ready to construct the non-compact SVD.

Let $V_n = \begin{bmatrix} V_r & V_g \end{bmatrix}$, where $V_g \in \mathbb{R}^{n \times (n-r)}$ is the collection of orthonormal vectors that span $\mathbb{R}^n / \text{range}(V_r)$, generated by Gram-Schmidt. Then

$$V_n V_n^T = \begin{bmatrix} V_r & V_g \end{bmatrix} \begin{bmatrix} V_r^T \\ V_g^T \end{bmatrix} = V_r V_r^T + V_g V_g^T$$

Therefore,

$$AV_n V_n^T = AV_r V_r^T + AV_g V_g^T$$

Claim 33. $\text{range}(V_g) \subseteq \text{null}(A^T A)$.

Proof. First, we claim that $\text{range}(V_r) \oplus \text{null}(A^T A) = \mathbb{R}^n$. We start with the almost tautological equation $\text{range}(A^T A) \oplus \text{range}(A^T A)^\perp = \mathbb{R}^n$. Then, we rewrite $\text{range}(A^T A) = \text{range}(V_r)$ since V_r has columns that are orthonormal basis vectors for $\text{range}(A^T A)$. Then, we write that $\text{range}(A^T A)^\perp = \text{null}((A^T A)^T) = \text{null}(A^T A)$. This yields that $\text{range}(V_r) \oplus \text{null}(A^T A) = \mathbb{R}^n$.

Now, we have that $\text{range}(V_g) \cap \text{range}(V_r) = \emptyset$. This follows from construction of V_r . To see why, assume otherwise; then there exists nonzero $v \in \text{range}(V_g) \cap \text{range}(V_r)$, so there exists a linear combination of basis vectors, without all coefficients from one space being 0, in $\text{range}(V_g)$ and $\text{range}(V_r)$ that equals v . But since V_r has a full complement of r basis vectors for its rank r subspace, this means that if $v \in \text{range}(V_r)$, then there would be no need to extend $\text{range}(V_r)$ to $\text{span}(v)$ by Gram-Schmidt so $v \notin \text{range}(V_g)$, a contradiction. \square

Since $\text{range}(V_g) \subseteq \text{null}(A^T A)$, $\text{range}(V_g) \in \text{null}(A)$ by Claim 32. Hence $AV_g V_g^T = 0 V_g^T = 0$, so $AV_r V_r^T = AV_n V_n^T$. But then V_n has rank n and is orthonormal, so $V_n V_n^T = I$.

Extending U_r to U_n in a similar manner gives the non-compact SVD.

We consider the geometry of the SVD. In general, if S is the unit circle, then if U is an orthonormal linear transformation, $US = S$. However, if T is a general linear transformation, TS is an ellipse, or a generalization of an ellipse to $\text{range}(A)$ dimensions. We attempt to find the longest axis direction:

$$x^* = \operatorname{argmax}_{\|x\|_2=1} \|Ax\|_2 = \operatorname{argmax}_{\|x\|_2=1} x^T A^T A x = v_{\max}(A^T A)$$

with the longest squared axis direction being $\lambda_{\max}(A^T A)$, so the longest axis direction is $\sigma_1(A^T A)$.

Least Norm Solution to Linear System, Pseudoinverses

We now consider the minimum-norm problem, which runs sort of counter to the least squares problem. The least squares problem has a matrix $A \in \mathbb{R}^{m \times n}$, where $m > n$, and we want to approximate the best inverse to A , to solve in x the overdetermined system $Ax = b$. The minimum norm problem has a matrix $A \in \mathbb{R}^{m \times n}$, where $m < n$, and there can be infinitely many solutions to $Ax = b$. We wish to find $\min_x \|x\|_2$ subject to the constraint $Ax = b$. Recall that from Theorem 19, that $\text{null}(A) \oplus \text{range}(A^T) = \mathbb{R}^n$. Consider a decomposition of a solution x into $x = x_n + x_r$, where $x_n \in \text{null}(A)$, so $Ax_n = 0$, and $x_r \in \text{range}(A^T)$, so $A^T y = x_r$. Then

$$Ax = A(x_n + x_r) = A(x_n + A^T y) = AA^T z.$$

Therefore

$$\begin{aligned} \|x\|_2^2 &= \|x_n + x_r\|_2^2 \\ &= \|x_n + A^T y\|_2^2 \\ &= (x_n + A^T y)^T (x_n + A^T y) \\ &= \|x_n\|_2^2 + 2 \langle x_n | A^T y \rangle + \|A^T y\|_2^2 \\ &= \|x_n\|_2^2 + \|A^T y\|_2^2 \end{aligned} \quad (\text{By Theorem 19})$$

To minimize $\|x\|_2$, we can pick $x_n = 0$ since $Ax = AA^T y$ is not a function of x_n . Therefore the chosen solution is

$$\begin{aligned} Ax &= AA^T y \\ &\stackrel{\text{set}}{=} b \\ y &= (AA^T)^{-1} b \end{aligned}$$

We have $x = x_n + x_r = x_n + A^T y$. Since in the minimum norm case $x_n = 0$, then we choose $x = A^T z = A^T (AA^T)^{-1} b$. Note that we want A to be full rank, or else we can reduce A until it has linearly independent columns.

Consider the compact singular value decomposition of $A = U_r \Sigma_r V_r^T$, with U_r and V_r orthonormal and Σ_r square and diagonal. Then the formula for x is

$$\begin{aligned} x &= A^T (AA^T)^{-1} b \\ &= (V_r \Sigma_r^T U_r^T) \left[(U_r \Sigma_r V_r^T) (V_r \Sigma_r^T U_r^T) \right]^{-1} b \end{aligned}$$

$$\begin{aligned}
&= \left(V_r \Sigma_r U_r^\top \right) \left(U_r \Sigma_r V_r^\top V_r \Sigma_r U_r^\top \right)^{-1} b \\
&= V_r \Sigma_r U_r^\top \left(U_r \Sigma_r^2 U_r^\top \right)^{-1} b \\
&= V_r \Sigma_r U_r^\top U_r \Sigma_r^{-2} U_r^\top \\
&= V_r \Sigma_r^{-1} U_r^\top b
\end{aligned}$$

We can take this inverse because

$$\begin{aligned}
\left(U_r \Sigma_r^2 U_r^\top \right)^{-1} \left(U_r \Sigma_r^2 U_r^\top \right) &= U_r \Sigma_r^{-2} U_r^\top U_r \Sigma_r^2 U_r^\top \\
&= U_r U_r^\top \\
&= \begin{bmatrix} I_r & 0_{m-r} \\ 0_{m-r} & 0_{(m-r)^2} \end{bmatrix}
\end{aligned}$$

If A is full row rank (as assumed), then $m = r$ and so $U_r U_r^\top = I_r$.

Definition 34 (Pseudoinverse). If $A = U_r \Sigma_r V_r^\top$ is the compact SVD of A , then $A^\dagger = V_r \Sigma_r^{-1} U_r^\top$ is the (Moore-Penrose) pseudoinverse of A .

Some properties of the pseudoinverse are:

- $AA^\dagger A = A$
- $AA^\dagger = U_r U_r^\top$
- $A^\dagger A = V_r V_r^\top$
- $A^\dagger AA^\dagger = A^\dagger$

Writing A^\dagger in terms of A gives

- If A is invertible, then $A^{-1} = A^\dagger$.
- If A is full row rank, then $A^\dagger = A^\top (AA^\top)^{-1}$ (the minimum-norm solution, also the right inverse).
- If A is full column rank, then $A^\dagger = (A^\top A)^{-1} A^\top$ (the least-squares solution, also the left inverse).

Matrix Norms, Low Rank Approximation

We now discuss matrix norms.

Definition 35 (Frobenius Norm). The Frobenius norm is the induced ℓ_2 -norm:

$$\|A\|_F = \|A\|_2 = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{i,j}|^2 \right)^{1/2} = \text{trace}(A^\top A)^{1/2}$$

Since it is the induced ℓ_2 -norm, it has a nice property.

Claim 36. If U and V are orthonormal matrices, then $\|AV\|_F = \|UA\|_F = \|A\|_F$.

Proof. Pick U and V orthonormal. Then

$$\|UA\|_F^2 = \text{trace}\left((UA)^\top (UA)\right) = \text{trace}\left(A^\top U^\top UA\right) = \text{trace}\left(A^\top A\right) = \|A\|_F^2$$

Also,

$$\|AV\|_F^2 = \text{trace}\left((AV)^\top (AV)\right) = \text{trace}\left(V^\top A^\top AV\right) = \text{trace}\left(A^\top AVV^\top\right) = \text{trace}\left(A^\top A\right) = \|A\|_F^2$$

as desired. \square

Definition 37 (Operator Norm). The spectral norm, or operator norm, also called the ℓ_2 -norm, gives

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \left(x^\top A^\top A x\right)^{1/2} = \lambda_{\max}\left(A^\top A\right)^{1/2} = \sigma_{\max}(A)$$

We now discuss the low rank approximation of matrices.

Theorem 38 (Eckart-Young-Mirsky Theorem)

Let $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = r$ and $A = U_A \Sigma_A V_A^\top$ be the full SVD of A . Let $A_k = \sum_{i=1}^k \sigma_i(A) u_i(A) v_i(A)^\top$, with $\sigma_1(A) \geq \dots \geq \sigma_n(A) > 0$. Then

$$A_k = \underset{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B)=k}}{\text{argmin}} \|A - B\|_2 = \underset{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B)=k}}{\text{argmin}} \|A - B\|_F.$$

Proof for Spectral Norm. We want to show that $A_k = \sum_{i=1}^k \sigma_i(A) u_i(A) v_i(A)^\top = \underset{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B)=k}}{\text{argmin}} \|A - B\|_2$.

Consider

$$\begin{aligned} \|A - A_k\|_2 &= \left\| \sum_{i=k+1}^n \sigma_i(A) u_i(A) v_i(A)^\top \right\|_2 \\ &= \sigma_{k+1}(A) \quad \left(\text{Spectral re-composition of } \sum_{i=k+1}^n \sigma_i(A) u_i(A) v_i(A)^\top, \text{ then Definition 37.} \right) \end{aligned}$$

We want to show that for any other matrix B with $\text{rank}(B) = k$, that $\|A - B\|_2 \geq \sigma_{k+1}(A)$. Since $\text{rank}(B) = k$, then $\dim(\text{null}(B)) = n - k > 0$. We have that since $\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}$, then $\|A\|_2 \geq \frac{\|Ax\|_2}{\|x\|_2}$ for any x , so $\|A\|_2 \|x\|_2 \geq \|Ax\|_2$. Also, define $V_{k+1} = [v_1 \ \dots \ v_{k+1}]$, and $\text{range}(V_{k+1}) = k + 1$. Therefore by Pidgeonhole principle $\text{null}(B) \cap \text{range}(V_{k+1}) \neq \emptyset$. Therefore

$$\begin{aligned} \|A - B\|_2^2 &\geq \|(A - B)w\|_2^2 && (\|w\|_2 = 1.) \\ &= \|Aw\|_2^2 && (\text{Pick } w \in \text{null}(B).) \\ &= \left\| A \sum_{i=1}^{k+1} \alpha_i v_i(A) \right\|_2^2 && (\text{Pick } w \in \text{null}(B) \cap \text{range}(V_{k+1}).) \\ &= \left\| U_A \Sigma_A V_A^\top \sum_{i=1}^{k+1} \alpha_i v_i(A) \right\|_2^2 \\ &= \left\| U_A \Sigma_A V_A^\top V \alpha \right\|_2^2 \\ &= \|U_A \Sigma_A \alpha\|_2^2 \\ &= \|\Sigma_A \alpha\|_2^2 \\ &= \sum_{i=1}^{k+1} \alpha_i^2 \sigma_i(A)^2 \end{aligned}$$

$$\begin{aligned}
&\geq \sigma_{k+1}(A)^2 \left(\sum_{i=1}^n \alpha_i^2 \right) \\
&\geq \sigma_{k+1}(A)^2 \quad (\text{Since } \|w\|_2 = 1, \sum_{i=1}^{k+1} \alpha_i^2 = 1.)
\end{aligned}$$

as desired. \square

Proof for Frobenius Norm. We want to show that $A_k = \sum_{i=1}^k \sigma_i(A) u_i(A) v_i(A)^\top = \operatorname{argmin}_{\substack{B \in \mathbb{R}^{m \times n} \\ \operatorname{rank}(B)=k}} \|A - B\|_F$.

Consider

$$\begin{aligned}
\|A - A_k\|_F^2 &= \left\| \sum_{i=1}^n \sigma_i(A) u_i(A) v_i(A)^\top - \sum_{i=1}^k \sigma_i(A) u_i(A) v_i(A)^\top \right\|_F^2 \\
&= \left\| \sum_{i=k+1}^n \sigma_i(A) u_i(A) v_i(A)^\top \right\|_F^2
\end{aligned}$$

Recall that the Frobenius norm is invariant to an orthonormal transformation. In particular, if U and W are orthonormal matrices, $\|A\|_F = \|UA\|_F = \|AW\|_F$. So then if $A = U_A \Sigma_A V_A^\top$ is the singular value decomposition of A ,

$$\|A\|_F^2 = \|U_A \Sigma_A V_A^\top\|_F^2 = \|\Sigma_A\|_F^2 = \sum_{i=1}^n \sigma_i(A)^2.$$

Therefore

$$\begin{aligned}
\|A - A_k\|_F^2 &= \left\| \sum_{i=k+1}^n \sigma_i(A) u_i(A) v_i(A)^\top \right\|_F^2 \\
&= \sum_{i=k+1}^n \sigma_i(A)^2
\end{aligned}$$

Now we want to show that for any other matrix B with $\operatorname{rank}(B) = k$, that $\|A - B\|_F^2 \geq \sum_{i=k+1}^n \sigma_i(A)^2$. In particular, we want to show that

$$\|A - B\|_F^2 = \sum_{i=1}^n \sigma_i(A - B)^2 \geq \sum_{i=k+1}^n \sigma_i(A)^2 = \|A - A_k\|_F^2$$

In particular, it's sufficient to show that $\sigma_i(A - B) \geq \sigma_{k+i}(A)$.

Define $A - B = C$ and $C_j = \sum_{i=1}^j \sigma_i(C) u_i(C) v_i(C)^\top$, then $\sigma_i(C) = \|C - C_{i-1}\|_2$. Also we write $\sigma_{k+i}(A) = \|A - A_{k+i-1}\|_2 = \|B + C - A_{k+i-1}\|_2$. Then

$$\begin{aligned}
\sigma_i(A - B) &\geq \|C - C_i\|_2 \\
&\geq \|A - B - C_i\|_2 \\
&\geq \|A - A_{k+i}\|_2 & (\operatorname{rank}(B + C_i) \leq \operatorname{rank}(B) + \operatorname{rank}(C_i) = k + i.) \\
&\geq \sigma_{k+i}(A)
\end{aligned}$$

which completes the result. \square

5 Dealing With Noise

We've already covered least squares as a way to deal with noise. The methods we discuss will cover additional methods of dealing with noise.

Let $Ax = y$, where $A \in \mathbb{R}^{n \times n}$ is invertible. If $y \rightarrow y + \Delta_y$, and because of this $x \rightarrow x + \Delta_x$, we want to find the relative change in Δ_x i.e. $\frac{\|\Delta_x\|_2}{\|x\|_2}$. We have

$$\begin{aligned} A(x + \Delta_x) &= y + \Delta_y \\ A\Delta_x &= \Delta_y \\ \Delta_x &= A^{-1}\Delta_y \\ \|\Delta_x\|_2 &= \|A^{-1}\Delta_y\|_2 \\ &\leq \|A^{-1}\|_2 \|\Delta_y\|_2 \end{aligned} \quad (\text{Definition of spectral norm.})$$

We also desire to bound $\|x\|_2$.

$$\begin{aligned} Ax &= y \\ \|Ax\|_2 &= \|y\|_2 \\ &\leq \|A\|_2 \|x\|_2 \end{aligned}$$

Therefore

$$\begin{aligned} \frac{\|\Delta_x\|_2}{\|x\|_2} &= \frac{\|A^{-1}\|_2 \|\Delta_y\|_2}{\|x\|_2} \\ &\leq \frac{\|A^{-1}\|_2 \|\Delta_y\|_2}{\|y\|_2 / \|A\|_2} \\ &\leq \|A\|_2 \|A^{-1}\|_2 \left(\frac{\|\Delta_y\|_2}{\|y\|_2} \right) \end{aligned}$$

We know that $\|A\|_2 = \sigma_{\max}(A)$ and $\|A^{-1}\|_2 = \sigma_{\min}(A)^{-1}$, so

$$\frac{\|\Delta_x\|_2}{\|x\|_2} \leq \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} \cdot \frac{\|\Delta_y\|_2}{\|y\|_2}$$

Definition 39 (Condition Number). The *condition number* of A is the quantity $K(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$.

The idea is that a lower condition number means a system given by $Ax = b$ is more resistant to noise.

Example 40. Let's look at least squares in terms of noise resistance. The full solution is that $x^* = (A^T A)^{-1} A^T b$, and the normal equation is $A^T A x^* = A^T b$. We care about the stability of this system. We have

$$K(A^T A) = \frac{\sigma_{\max}(A^T A)}{\sigma_{\min}(A^T A)} = \frac{\sigma_{\max}(A)^2}{\sigma_{\min}(A)^2} = \left(\frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} \right)^2 = \kappa(A)^2$$

due to the definition of singular values. We observe that this can create instability due to noise very fast (quadratically) compared to the underlying system $Ax = b$.

This is important because it motivates *ridge regression*. We wish to change the singular values of $A^T A$ so that we get a better condition vector.

Claim 41. If $A \in \mathbb{R}^{n \times n}$, then $\lambda_i(A + \lambda I) = \lambda_i(A) + \lambda$.

Proof. Let $v_i(A)$ be such that $Av_i(A) = \lambda_i(A)v_i(A)$. Then

$$(A + \lambda I)v_i(A) = Av_i(A) + \lambda Iv_i(A) = \lambda_i(A)v_i(A) + \lambda v_i(A) = (\lambda_i(A) + \lambda)v_i(A)$$

as desired. \square

Notice that the least squares problem is the optimization problem

$$\min_x \|Ax - b\|_2^2$$

We consider the new *regularized* optimization problem

$$\min_x \|Ax - b\|_2^2 + \underbrace{\lambda^2 \|x\|_2^2}_{\text{penalty term}}$$

The λ parameter trades off how much we care about $\|x\|_2^2$ being small against our accuracy. The penalty term is also called a *regularizer*.

This is a quadratic convex problem. Let $f(x) = \|Ax - b\|_2^2 + \lambda^2 \|x\|_2^2$. Then expanding we have

$$\begin{aligned} f(x) &= \|Ax - b\|_2^2 + \lambda^2 \|x\|_2^2 \\ &= (Ax - b)^\top (Ax - b) + \lambda^2 x^\top x \\ &= x^\top A^\top Ax - x^\top A^\top b - b^\top Ax + b^\top b + \lambda^2 x^\top x \\ \nabla_x f(x) &= \nabla_x (x^\top A^\top Ax) - 2\nabla_x (b^\top Ax) + \nabla_x (b^\top b) + \nabla_x (\lambda^2 x^\top x) \\ &= 2A^\top Ax - 2A^\top b + 2\lambda^2 x \\ &\stackrel{\text{set}}{=} 0 \\ 2A^\top Ax + 2\lambda^2 x &= 2A^\top b \\ (A^\top A + \lambda^2 I)x^* &= A^\top b \\ x^* &= (A^\top A + \lambda^2 I)^{-1} A^\top b \end{aligned}$$

This is the solution to the ridge regression problem, or ℓ_2 regularized least squares. One can see that $K(A^\top A + \lambda^2 I) < K(A^\top A)$, showing the resistance to noise of this system.

Another interpretation of the ridge regression is that we can encode prior information on the norm of x . If we have information that the i^{th} coordinate of x has $\lambda x_i \approx 0$ for some $\lambda \neq 0$, then we can augment A to obtain

$$\underbrace{\begin{bmatrix} A \\ \lambda I \end{bmatrix}}_{\tilde{A}} x = \underbrace{\begin{bmatrix} b \\ 0 \end{bmatrix}}_{\tilde{b}}$$

The least squares formulation of this matrix has

$$x^* = (\tilde{A}^\top \tilde{A})^{-1} \tilde{A}^\top \tilde{b} = \begin{bmatrix} A^\top & \lambda I \end{bmatrix} \begin{bmatrix} A \\ \lambda I \end{bmatrix}^{-1} \begin{bmatrix} b \\ 0 \end{bmatrix} = (A^\top A + \lambda^2 I)^{-1} A^\top b,$$

What if we had some information that the i^{th} coordinate of x has $\lambda x_i \approx x'_i$? Then we can generalize, and

obtain the matrix equation

$$\begin{bmatrix} A \\ \lambda I \end{bmatrix} x = \begin{bmatrix} b \\ x_0 \end{bmatrix}$$

We can also weight the different rows of the data matrix.

Definition 42 (Tikhonov Regularization). Let W_1, W_2 be weight matrices, usually diagonal. Then the *Tikhonov regularization* problem is

$$\min_x \|W_1 (Ax - b)\|_2^2 + \|W_2 (x - x')\|_2^2$$

for the side information vector x' .

These two different perspectives are first, to improve the condition number of A , and second, to induce prior information into our system.

The next goal is to use a probabilistic model to choose our hyperparameter for i.e. ridge regression. The motivation for this is that we want to include probabilistic information into the model.

Let's say we have $\{(x_i, y_i)\}_{i=1}^m$ be the data points, where $y_i = g(x_i) + Z_i$, where $Z_i \sim \text{Normal}(0, \sigma_i^2)$, and $p_{Z_i}(z) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{z^2}{2\sigma_i^2}\right)$. Dealing with general $g \in L^2[x]$ is difficult, so we restrict g to be linear. In particular, let $g(x) = w^\top x$, and w is the set of model parameters (what we want to learn), so that $y_i = w^\top x_i + Z_i$.

Least squares finds w such that $\sum_{i=1}^m (w^\top x_i - y_i)^2$ is minimized. We incorporate the probabilistic information into this model using *maximum likelihood estimation*. We find the w that maximizes the probability density of the data:

$$\text{MLE}(w | x, y) = \underset{w}{\operatorname{argmax}} p_{X,Y|w}(x, y | w)$$

In this particular case, X is deterministic, so

$$\begin{aligned} \text{MLE}(w | x, y) &= \text{MLE}(w | y) \\ &= \underset{w}{\operatorname{argmax}} p_{Y|w}(y | w) \\ &= \underset{w}{\operatorname{argmax}} \prod_{i=1}^m p_{Y_i|w}(y_i | w) \\ &= \underset{w}{\operatorname{argmax}} \prod_{i=1}^m p_{Z_i|w}(y_i - w^\top x_i | w) \\ &= \underset{w}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - w^\top x_i)^2}{2\sigma_i^2}\right) \\ &= \underset{w}{\operatorname{argmax}} \frac{1}{(\sqrt{2\pi})^n} \frac{1}{\prod_{i=1}^m \sigma_i} \exp\left(-\sum_{i=1}^m \frac{(y_i - w^\top x_i)^2}{2\sigma_i^2}\right) \\ &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^m \frac{(y_i - w^\top x_i)^2}{2\sigma_i^2} \\ &= \underset{w}{\operatorname{argmin}} \|S(y - xw)\|_2^2 \end{aligned}$$

where $S = \text{diag}\left(\frac{1}{\sqrt{2\pi\sigma_1^2}}, \dots, \frac{1}{\sqrt{2\pi\sigma_m^2}}\right)$. This is related to Tikhonov regularization and weighted least squares, and suggests that a way of interpreting least squares is to see which components or features matter more to the final prediction.

If we have a prior on w , we can use *maximum a posteriori estimation*:

$$\text{MAP}(w | x, y) = \underset{w}{\text{argmax}} p_{w|X,Y}(w | x, y) = \underset{w}{\text{argmax}} p_{X,Y|w}(x, y | w) p_w(w)$$

In this case, assume that $w \sim \text{Normal}(\mu, \Sigma_w)$ is the prior distribution on w . Then

$$\begin{aligned} \text{MAP}(w | x, y) &= \underset{w}{\text{argmax}} p_{Y|w}(y | w) p_w(w) \\ &= \underset{w}{\text{argmax}} p_w \prod_{i=1}^m p_{Y|w}(w | w) \\ &= \underset{w}{\text{argmax}} \left[\prod_{i=1}^m \frac{\exp\left(-\frac{(y_i - w^\top x_i)^2}{2\sigma_i^2}\right)}{\sqrt{2\pi\sigma_i^2}} \right] \left[\frac{\exp\left(-(w - \mu)^\top \Sigma_w^{-1} (w - \mu)\right)}{\sqrt{(2\pi)^m \det(\Sigma_w)}} \right] \\ &= \underset{w}{\text{argmin}} \exp\left(\sum_{i=1}^m \frac{(y_i - w^\top x_i)^2}{2\sigma_i^2} + (w - \mu)^\top \Sigma_w^{-1} (w - \mu)\right) \\ &= \underset{w}{\text{argmin}} \|S(xw - y)\|_2^2 + \left\| \Sigma_w^{-1/2} (w - \mu) \right\|_2^2 \end{aligned}$$

The $\Sigma_w^{-1/2}$ is similar to the regularization term, and enforces the viewpoint that the regularizer enforces prior beliefs and certainty of those beliefs.

We use this to introduce *principal component regression*, which is the technique of projecting onto low dimensional principal components and doing regression from there.

The problem is to minimize $\|Xw - y\|_2^2$. If $X = U_X \Sigma_X V_X^\top$, then the least squares solution is $w^* = V \Sigma^\dagger U^\top y$, where $\Sigma_{i,j}^\dagger = \frac{\delta_{i,j}}{\Sigma_{i,i}}$. For principal components regression, we do the regression with only the top k principal components, i.e., $w^* = V \Sigma^{\text{PCR}} U^\top y$ where $\Sigma^{\text{PCR}} = \frac{\mathbf{1}(i=j \wedge i \leq k)}{\Sigma_{i,i}}$.

There's yet another perspective on ridge regression, which is motivated by this. In particular, ridge regression is a soft form of principal components analysis. Starting from the ridge regression standpoint, we have

$$\begin{aligned} w^* &= \underset{w}{\text{argmin}} \|Xw - y\|_2^2 + \lambda \|w\|_2^2 \\ &= \underset{w=V_X z}{\text{argmin}} \|XV_X z - y\|_2^2 + \lambda \|V_X z\|_2^2 & (X = U_X \Sigma_X V_X^\top) \\ &= \underset{z}{\text{argmin}} \|XV z - y\|_2^2 + \lambda \|z\|_2^2 \\ z^* &= \left[(XV_X)^\top (XV_X) + \lambda I \right]^{-1} (XV_X)^\top y \\ &= \left(V_X^\top X^\top X V_X + \lambda I \right)^{-1} (XV_X)^\top y \\ &= \left[V_X^\top (U_X \Sigma_X V_X^\top) V_X + \lambda I \right]^{-1} (XV_X)^\top y \\ &= \left(\Sigma_X^\top \Sigma_X + \lambda I \right)^{-1} \Sigma_X^\top U^\top y \end{aligned}$$

$$= \begin{bmatrix} \frac{\sigma_1}{\sigma_1^2 + \lambda} & & \\ & \ddots & \\ & & \frac{\sigma_n}{\sigma_n^2 + \lambda} \end{bmatrix} U^\top y$$

and one can read off the asymptotic behavior of z^* as a function of λ .

6 Convexity

Today we discuss convex sets and convex functions. Convex problems are particularly nice because there exists a unique optimum point, which can be computed in closed form if the objective function is differentiable.

Definition 43 (Convex Combination). A linear combination $\sum_{i=1}^n a_i x_i$ is a *convex combination* if $\sum_{i=1}^n a_i = 1$ and $a_i \geq 0$ for all i .

Definition 44 (Convex Set). A set C is *convex* if for any set of points $S = \{x_1, \dots, x_n\}$ any convex combination of S is in C . Equivalently, C is convex if for any $x, y \in C$, then $\alpha x + (1 - \alpha)y \in C$, for $\alpha \in [0, 1]$.

Pictorially, C is convex if any line between two points in C consists only of points in C .

Definition 45 (Convex Hull). The *convex hull* of a set of points $S = \{x_1, \dots, x_n\}$ is the minimal-measure convex set that has S as a subset.

Example 46. Let $C = \{x \mid a^\top x = b\}$. Then we claim C is convex. Take $x_1, x_2 \in C$. Let $\alpha \in [0, 1]$ and $x_3 = \alpha x_1 + (1 - \alpha)x_2$. Then $a^\top x_3 = \alpha a^\top x_1 + a^\top x_2 - \alpha a^\top x_2 = \alpha a^\top x_1 - a^\top x_2 + a^\top x_2 = a^\top x_2 = b$, so $x_3 \in C$, so C is convex.

Definition 47 (Half Space). A *half space* corresponding to vectors a and b is the set $\{x \mid a^\top x \geq b\}$ or $\{x \mid a^\top x \leq b\}$.

Example 48. We claim $\mathbb{S}_{\geq 0}^n$ (the set of positive semidefinite matrices) is convex. Take $\alpha \in [0, 1]$. Let $A_1, A_2 \in \mathbb{S}_{\geq 0}^n$ and $A_3 = \alpha A_1 + (1 - \alpha)A_2$. Then for any x we have $\alpha x^\top A_1 x + (1 - \alpha)x^\top A_2 x \geq \alpha \cdot 0 + (1 - \alpha) \cdot 0 = 0$, so $A_3 \in \mathbb{S}_{\geq 0}^n$, so $\mathbb{S}_{\geq 0}^n$ is convex.

Theorem 49 (Separating Hyperplane Theorem)

Let $C, D \subseteq \mathbb{R}^n$ be two convex, compact sets with $C \cap D = \emptyset$. Then there exists a hyperplane $a^\top x = b$, such that

- for all $x \in C$ we have $a^\top x \geq b$.
- for all $x \in D$ we have $a^\top x \leq b$.

Proof. Define $d(C, D) = \inf \{\|c - d\|_2 \mid c \in C, d \in D\}$. Let c and d be the minimizing points; we know that such a pair exists since C and D are compact. Then the hyperplane with $a = d - c$ as the normal vector and that passes through $x_0 = \text{midpoint}(c, d) = \frac{c+d}{2}$ is given by

$$f(x) = a^\top (x - x_0)$$

$$= (d - c)^T \left(x - \frac{c + d}{2} \right)$$

We have that $f(d) = (d - c)^T \left(d - \frac{c + d}{2} \right) = \frac{1}{2} \|d - c\|_2^2 \geq 0$. Similarly $f(c) = -\frac{1}{2} \|d - c\|_2^2$.

We want to show that for all $x \in D$ that $f(x) \geq 0$. By symmetry this shows that for all $x \in C$ that $f(x) \leq 0$.

Assume for the sake of contradiction that $u \in D$ such that $f(u) < 0$. Then

$$\begin{aligned} f(u) &= (d - c)^T \left(u - \frac{c + d}{2} \right) \\ &= (d - c)^T \left(u - d + d - \frac{c + d}{2} \right) \\ &= (d - c)^T \left[(u - d) + \frac{d - c}{2} \right] \\ &= (d - c)^T (u - d) + \frac{1}{2} \|d - c\|_2^2 \\ &= \langle d - c | u - d \rangle + \frac{1}{2} \|d - c\|_2^2 \\ &< 0 \end{aligned}$$

Therefore $\langle d - c | u - d \rangle < 0$.

Let $t \in [0, 1]$; then define $p = d + t(u - d)$. Clearly $p \in D$. Then

$$\begin{aligned} \|c - p\|_2^2 &= \|c - d - t(u - d)\|_2^2 \\ &= [(c - d) - t(u - d)]^T [(c - d) - t(u - d)] \\ &= \|c - d\|_2^2 + t^2 \|u - d\|_2^2 - 2 \langle c - d | t(u - d) \rangle \\ &= \|c - d\|_2^2 + t^2 \|u - d\|_2^2 - 2t \langle c - d | u - d \rangle \end{aligned}$$

We want $t^2 \|u - d\|_2^2 - 2t \langle c - d | u - d \rangle < 0$, so we want $2t \langle d - c | u - d \rangle + t^2 \|u - d\|_2^2 < 0$, so we want $2 \langle d - c | u - d \rangle + t \|u - d\|_2^2 < 0$. Since we can choose t and $\langle d - c | u - d \rangle < 0$, choosing $t < 2 \frac{\langle d - c | u - d \rangle}{\|u - d\|_2^2}$ shows that $\|c - p\|_2^2 \leq 0$. But then $\|c - p\|_2^2 < \|c - d\|_2^2$, which is a contradiction with the definition of c and d . \square

Definition 50 (Convex Function). A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for all $\alpha \in [0, 1]$, $x, y \in \mathbb{R}^n$, we have

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

Definition 51 (Epigraph). The epigraph of f is the set

$$\text{epi}(f) = \{(x, t) : x \in \text{domain}(f), f(x) \leq t\}$$

f is a convex function if and only if $\text{epi}(f)$ is a convex set.

Theorem 52 (First-Order Conditions)

let $f \in C^1(\mathbb{R}^n, \mathbb{R})$. Then f is convex if and only if $\text{domain}(f)$ is convex and for each $x, y \in \mathbb{R}^n$

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

This is what gives convexity its power. The derivative (gradient) of f is a local property at x ; however, we can find global properties of f if f is convex. In particular, if f is convex and $\nabla_x f(x^*) = 0$ then $f(y) \geq f(x^*) + 0$ for all x , so $f(x^*)$ is a global minimum.

Theorem 53 (Second-Order Conditions)

Let $f \in C^2(\mathbb{R}^n, \mathbb{R})$. Then f is convex if and only if $\text{domain}(f)$ is convex and $\nabla_x^2 f(x) \succeq 0$ (is positive semidefinite).

Theorem 54 (Jensen's Inequality for Probability)

Let f be a convex function and X a random variable. Then

$$f(\text{Ex}[X]) \leq \text{Ex}[f(X)]$$

In the discrete non-probabilistic case, we have

Theorem 55 (Jensen's Inequality)

Let f be a convex function and $\lambda_1, \dots, \lambda_n$ be such that $\sum_{i=1}^n \lambda_i = 1$. Then for all $x_1, \dots, x_n \in \mathbb{R}$, we have

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i)$$

The proof is by definitions of convexity.

A general optimization problem is

$$\begin{aligned} x^* &= \underset{x}{\text{argmin}} f(x) \\ \text{s.t. } f_i(x) &\leq 0, \quad i \in [1, m] \\ g_i(x) &= 0, \quad i \in [1, \ell] \end{aligned}$$

Definition 56 (Convex Problem). A *convex problem* is an optimization problem where $f(x)$ and $f_i(x)$ and $g_i(x)$ are all convex, and the feasible region $D = \text{domain}(f) \cap \bigcap_{i=1}^m \text{domain}(f_i) \cap \bigcap_{i=1}^{\ell} \text{domain}(g_i)$ is convex.

Example 57 (Linear Program). A linear program

$$\begin{aligned} x^* &= \underset{x}{\text{argmin}} c^\top x \\ \text{s.t. } Ax &= b \end{aligned}$$

is a convex problem.

Example 58. The optimization problem

$$\begin{aligned} x^* &= \underset{x}{\text{argmin}} x^2 \\ \text{s.t. } x &\geq 1 \\ x &\leq 2 \end{aligned}$$

is a convex problem; in this case the constraint $x \geq 1$ is called a *active constraint*.

In general, if $f(x)$ is convex, then to find the optimum point x^* we set $[\nabla_x f(x)]_{x=x^*} = 0$.

Example 59. The optimization problem

$$\begin{aligned} x^* &= \operatorname{argmin}_x x^2 \\ \text{s.t. } x &\geq 1 \\ x &\leq -2 \end{aligned}$$

is *infeasible*; the feasible region $D = \emptyset$.

Example 60. The optimization problem

$$\begin{aligned} x^* &= \operatorname{argmin}_x x_1 + x_2 \\ \text{s.t. } x_1^2 &\geq 2 \\ x_2^2 &\leq 1 \end{aligned}$$

is a convex problem; taking the gradient gives $\nabla_x f(x) = 1$ (the ones vector), but by geometry or observation one sees that $x^* = \begin{bmatrix} -\sqrt{2} \\ -1 \end{bmatrix}$ works. Both constraints are active.

Example 61. The optimization problem

$$\begin{aligned} x^* &= \operatorname{argmin}_x x_1 \\ \text{s.t. } x_1 + x_2 &\geq 0 \end{aligned}$$

is *unbounded*; we can take x_1 arbitrarily low and consider $x_2 = -x_1$.

At this point it's useful to consider inf vs min; $\inf S$ need not be attained by any value in S .

Theorem 62

Let X be a convex and closed set. Then if $x^* = \operatorname{argmin}_{x \in X} c^\top x$ then $x \in \overline{X} - X^o$, the set of limit points of X .

Proof. Assume for the sake of contradiction that $x^* \in X^o$, the interior of X . Then there exists r such that there exists some ball $\{z \in X : \|x^* - z\|_2 < r\} = B_r(x^*) \subseteq X^o$. We have that $\nabla_x f(x) = c$. Consider $z = \frac{rc}{\|c\|_2^2} \in B_r(x^*)$, so $z \in X^o$. Then

$$\begin{aligned} f(x^* - z) &= c^\top x^* - c^\top z \\ &= c^\top x^* - \frac{r}{\|c\|_2} c^\top c \\ &= c^\top x^* - r\|c\|_2 \\ &\leq c^\top x^* \end{aligned}$$

which has a better optimum, a contradiction. □

We define the *epigraph reformulation* as

$$\min_{x \in X} f(x) = \min_{f(x) \leq t} t$$

where t is known as a *slack variable*.

Definition 63 (LASSO Regression). We now introduce LASSO regression, or ℓ^1 normalization. We have the optimization problem

$$x^* = \operatorname{argmin}_x \|Ax - y\|_2^2 + \|x\|_1$$

Adding the epigraph reformulation, we have

$$\begin{aligned} x^* &= \operatorname{argmin}_{x,t} \|Ax - y\|_2^2 + \sum_{i=1}^n t_i \\ \text{s.t. } & -t_i \leq x_i \leq t_i, \quad i \in [1, n] \\ & t_i \geq 0, \quad i \in [1, n] \end{aligned}$$

Definition 64 (Monotone Transformation). If $\phi(x)$ is continuous and strictly increasing, then over a domain D we have that $\operatorname{argmin}_{x \in D} f(x) = \operatorname{argmin}_{x \in D} \phi(f(x))$, and similarly for monotone decreasing functions.

Definition 65 (Logistic Regression). Say that we have the data points X_1, \dots, X_m and corresponding labels $y_1, \dots, y_m \in \{-1, 1\}$. We wish to predict $p_{Y|X}(1 | X_i)$. We achieve this by learning w and β such that $p_{Y|X}(1 | X_i) = \sigma(w^\top X_i + \beta)$ where $\sigma(x) = \frac{1}{1+e^{-x}}$. In particular, $p_{Y|X}(y_i | X_i) = \sigma(-y_i (w^\top X_i + \beta))$.

How do we solve for w and β ? We use maximum likelihood estimation. We have

$$\begin{aligned} (w^*, \beta^*) &= \operatorname{argmax}_{w, \beta} \prod_{i=1}^m p_{Y|X}(y_i | X_i) \\ &= \operatorname{argmax}_{w, \beta} \prod_{i=1}^m \sigma(-y_i (w^\top X_i + \beta)) \\ &= \operatorname{argmax}_{w, \beta} \sum_{i=1}^m \log(\sigma(-y_i (w^\top X_i + \beta))) \end{aligned}$$

and the rest is computation.

7 Optimization Algorithms

Gradient Descent

For the purposes of this lecture, we consider unconstrained optimization problems of the form

$$x^* = \min_{x \in \mathbb{R}^n} f(x)$$

For most functions f , even if f is convex, there is no closed-form solution. To optimize globally these functions, we use gradient descent. This is a consequence of Taylor's theorem, which states that

$$\begin{aligned} f(x + \Delta_x) &= f(x) + \nabla_x f(x)^\top \Delta_x + o(\|\Delta_x\|_2^3) \\ f(x + su) &= f(x) + \nabla_x f(x)^\top su \\ &= f(x) + s \langle \nabla_x f(x) | u \rangle \end{aligned}$$

We want to find the minimum, so we want $\langle \nabla_x f(x) | u \rangle < 0$. By Cauchy-Schwarz, we see that the minimizing direction is $u^* = -\nabla_x f(x)$.

This yields an algorithm:

Algorithm 3 Gradient descent algorithm.

Input: Initial vector $x^{(0)}$, differentiable function $f \in C^1(\mathbb{R}^n, \mathbb{R})$, learning rate η .

Output: Local minimum of f .

```

 $k \leftarrow 0$ 
while  $\nabla_{x^{(k)}} f(x^{(k)}) \neq 0$  do
     $x^{(k+1)} \leftarrow x^{(k)} - \eta \nabla_{x^{(k)}} f(x^{(k)})$ 
     $k \leftarrow k + 1$ 
return  $x^{(k)}$ 

```

The interesting questions are how and when gradient descent converges to the minimal point, i.e. when $x^* = \text{GRADIENTDESCENT}(x^{(0)}, f, \eta)$.

Example 66 (Least Squares Gradient Descent). Let $f(x) = \|Ax - b\|_2^2$, the least squares objective. Then $\nabla_x f(x) = 2A^\top Ax - 2A^\top b$. Then the gradient update rule is

$$\begin{aligned}
 x^{(k+1)} &= x^{(k)} - \eta \nabla_{x^{(k)}} f(x^{(k)}) \\
 &= x^{(k)} - \eta (2A^\top Ax^{(k)} - 2A^\top b) \\
 &= (I - 2\eta A^\top A) x^{(k)} + 2\eta A^\top b
 \end{aligned}$$

Recalling the analysis of stability, we want $|\lambda_i(I - 2\eta A^\top A)| < 1$.

If A is full column rank, then

$$\begin{aligned}
 x^{k+1} - x^* &= x^{(k+1)} - (A^\top A)^{-1} A^\top b \\
 &= (I - 2\eta A^\top A) x^{(k)} + 2\eta A^\top b - (A^\top A)^{-1} A^\top b \\
 &= (I - 2\eta A^\top A) x^{(k)} + 2\eta (A^\top A) (A^\top A)^{-1} A^\top b - (A^\top A)^{-1} A^\top b \\
 &= (I - 2\eta A^\top A) x^{(k)} + (2\eta A^\top A - I) (A^\top A)^{-1} A^\top b \\
 &= (I - 2\eta A^\top A) \left[x^{(k)} - (A^\top A)^{-1} A^\top b \right] \\
 &= (I - 2\eta A^\top A)^{k+1} \left[x^{(0)} - (A^\top A)^{-1} A^\top b \right]
 \end{aligned}$$

We want to figure out for which η that $|\lambda_i(I - 2\eta A^\top A)| \leq 1$; we'll do this later.

For very large A , computing gradient descent is lower cost computationally than computing the closed form iterates.

We want to generalize the idea of least squares to smooth, strongly convex functions.

Definition 67 (μ -Strongly Convex). A function f is μ -strongly convex on a domain D if for all $x, y \in D$

$$f(y) \geq f(x) + \nabla_x f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|_2^2.$$

Definition 68 (L -Smooth). A function f is L -smooth on a domain D if for all $x, y \in D$

$$f(y) \leq f(x) + \nabla_x f(x)^\top (y - x) + \frac{L}{2} \|y - x\|_2^2.$$

Theorem 69

Let $(x^{(k)})_{k \in \mathbb{N}}$ be a sequence of gradient descent iterates for a μ -strongly convex, L -smooth function f . Then there exists C such that

$$\|x^{(k+1)} - x^*\|_2^2 \leq C^{k+1} \|x^{(0)} - x^*\|_2^2.$$

Proof. We require a small claim on smoothness:

Claim. If f is L -smooth, then

$$\|\nabla_x f(x)\|_2^2 \leq 2L(f(x) - f(x^*)).$$

Proof. We have that

$$\begin{aligned} f(x^*) &\leq f(x) \\ &\leq f\left(x - \frac{\nabla_x f(x)}{L}\right) \\ f\left(x - \frac{\nabla_x f(x)}{L}\right) &\leq f(x) + [\nabla_{x'} f(x')]_{x'=x}^\top (-L[\nabla_{x'} f(x')]_{x'=x}) + \frac{L}{2} \left\| \frac{[\nabla_{x'} f(x')]_{x'=x}}{L} \right\|_2^2 \\ &= f(x) - \frac{1}{L} \left\| [\nabla_{x'} f(x')]_{x'=x} \right\|_2^2 + \frac{1}{2L} \left\| [\nabla_{x'} f(x')]_{x'=x} \right\|_2^2 \\ &= f(x) - \frac{\left\| [\nabla_{x'} f(x')]_{x'=x} \right\|_2^2}{2L} \\ f(x^*) &\leq f(x) - \frac{\left\| [\nabla_{x'} f(x')]_{x'=x} \right\|_2^2}{2L} \\ \left\| [\nabla_{x'} f(x')]_{x'=x} \right\|_2^2 &\leq 2L(f(x) - f(x^*)) \end{aligned}$$

as desired. ■

By the definition of strong convexity,

$$\begin{aligned} f(x^*) &\geq f(x) + \nabla_x f(x)^\top (x^* - x) + \frac{\mu}{2} \|x^* - x\|_2^2 \\ \nabla_x f(x)^\top (x - x^*) &\geq f(x) - f(x^*) = \frac{\mu}{2} \|x^* - x\|_2^2 \end{aligned}$$

Now we want to bound the difference $\|x^{(k+1)} - x^*\|_2^2$:

$$\begin{aligned} \|x^{(k+1)} - x^*\|_2^2 &= \|x^{(k)} - \eta \nabla_{x^{(k)}} f(x^{(k)}) - x^*\|_2^2 \\ &= \left\| (x^{(k)} - x^*) - \eta \nabla_{x^{(k)}} f(x^{(k)}) \right\|_2^2 \\ &= \|x^{(k)} - x^*\|_2^2 + \eta^2 \left\| \nabla_{x^{(k)}} f(x^{(k)}) \right\|_2^2 - 2\eta \left\langle \nabla_{x^{(k)}} f(x^{(k)}) \mid x^{(k)} - x^* \right\rangle \\ &\leq \|x^{(k)} - x^*\|_2^2 + 2\eta^2 L (f(x^{(k)}) - f(x^*)) - 2\eta \left(f(x^{(k)}) - f(x^*) + \frac{\mu}{2} \|x^{(k)} - x^*\|_2^2 \right) \\ &\leq (1 - \eta\mu) \|x^{(k)} - x^*\|_2^2 + (2\eta^2 L - 2\eta) (f(x^{(k)}) - f(x^*)) \end{aligned}$$

$$\begin{aligned} &\leq \left(1 - \frac{\mu}{L}\right) \|x^{(k)} - x^*\|_2^2 && (\text{Choose } \eta = \frac{1}{L}.) \\ \|x^{(k+1)} - x^*\|_2^2 &\leq \left(1 - \frac{\mu}{L}\right)^{k+1} \|x^{(0)} - x^*\|_2^2 \end{aligned}$$

Therefore, one may write $C = 1 - \frac{\mu}{L}$ and prove the claim. \square

This is exponential convergence, but since C is a linear function, sometimes it's called linear convergence.

Let f be differentiable.

- If f is μ -strongly convex and L -smooth, then convergence is $\mathcal{O}(\exp(-t))$.
- If f is μ -strongly convex and Lipschitz-continuous, or convex and smooth, then convergence is $\mathcal{O}\left(\frac{1}{t}\right)$.
- If f is convex and Lipschitz continuous, then convergence is $\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$.

A function f is K -Lipschitz-continuous over D if for all $x, y \in D$, $\|f(x) - f(y)\|_2 \leq K\|x - y\|_2$, where $\|\cdot\|$ can be replaced in a general metric space by a distance function.

8 Duality

Unconstrained optimization of a differentiable function f can be solved by gradient descent, or solving for x in the differential equation $\nabla_x f(x) = 0$. Constrained optimization problems can be turned into unconstrained optimization problems in the following way. Let $f(x)$ be the objective, $g(x) \leq 0$ be the set of inequality constraints put into vector form, and $h(x) = 0$ be the set of equality constraints put into vector form:

$$\begin{aligned} p^* &= \min_x f(x) \\ \text{s.t. } &g_i(x) \leq 0, i \in \{1, \dots, n_{\leq}\} \\ &h_i(x) = 0, i \in \{1, \dots, n_{=}\} \end{aligned}$$

We define the *Lagrangian* as

$$L(x, \lambda, \mu) = f(x) + \lambda^\top g(x) + \mu^\top h(x)$$

where $\lambda \geq 0$ (as in, $\lambda_i \geq 0$). The λ and μ are *dual variables* or *Lagrange multipliers*. Define $\min_x \mathcal{L}(x, \lambda, \mu) = \ell(\lambda, \mu)$.

It's easy to show that the pointwise maximum of convex functions is convex; a corollary is that the pointwise minimum of concave functions is concave. Therefore $\ell(\lambda, \mu)$ is concave, and the concavity does not depend on $f(x)$.

Claim 70. $\ell(\lambda, \mu) \leq p^*$, for all $\lambda \geq 0$ (elementwise) and μ .

Proof. Say that \tilde{x} is a feasible point (fulfills constraints) for the primal problem (the one that defines p^*). Then $g(\tilde{x}) \leq 0$ elementwise, and $h(\tilde{x}) = 0$ elementwise. Therefore for any positive $\lambda \geq 0$ we have $\lambda^\top g(\tilde{x}) \leq 0$ and $\mu^\top h(\tilde{x}) = 0$, since $\lambda \geq 0$ componentwise. Then

$$\begin{aligned} \mathcal{L}(\tilde{x}, \lambda, \mu) &= f(\tilde{x}) + \lambda^\top g(\tilde{x}) + \mu^\top h(\tilde{x}) \\ &\leq f(\tilde{x}) \end{aligned}$$

Since $\ell(\lambda, \mu) = \inf_x \mathcal{L}(x, \lambda, \mu)$, it must be true that $\ell(\lambda, \mu) \leq \mathcal{L}(\tilde{x}, \lambda, \mu) = p^*$ as desired. \square

An interpretation of this result is that we can do an unconstrained optimization of the objective function $f(x) + \sum_{i=1}^{n_{\leq}} g_i(x) (\infty \cdot \mathbb{1}(g_i(x) > 0)) + \sum_{i=1}^{n_{=}} h_i(x) (\infty \cdot \mathbb{1}(h_i(x) \neq 0))$ and obtain the minimum every time if we have a minimizer that can deal with such penalties. But this too harsh of a penalty to be practical, and it's not smooth; the Lagrange multipliers give a linear penalty on violating the constraints.

Example 71 (Minimum Norm Solution). We use Lagrange multipliers to solve the optimization problem

$$\begin{aligned} p^* &= \min_x x^\top x \\ \text{s.t. } Ax &= b \end{aligned}$$

The Lagrangian is $\mathcal{L}(x) = x^\top x + \mu^\top (Ax - b)$; define $\ell(\mu) = \min_x \mathcal{L}(x, \mu)$. We have $\nabla_x \mathcal{L}(x) = 2x + A^\top \mu \stackrel{\text{set}}{=} 0$, so $x^* = -\frac{1}{2} A^\top \mu$. Then

$$\begin{aligned} \ell(\mu) &= \mathcal{L}(x^*, \mu) \\ &= \mathcal{L}\left(-\frac{1}{2} A^\top \mu, \mu\right) \\ &= \frac{1}{4} \mu^\top A A^\top \mu + \mu^\top \left[A \left(-\frac{1}{2} A^\top \mu\right) - b\right] \\ &= -\frac{1}{4} \mu^\top A A^\top \mu - \mu^\top b \end{aligned}$$

Since $p^* \geq \ell(\mu)$, we wish to maximize $\ell(\mu)$ over all μ . We have

$$\begin{aligned} \nabla_\mu \ell(\mu) &= \nabla_\mu \left(-\frac{1}{4} \mu^\top A A^\top \mu - \mu^\top b\right) \\ &= -\frac{1}{4} (2 A A^\top) \mu - b \\ &\stackrel{\text{set}}{=} 0 \\ \mu^* &= -2 (A A^\top)^{-1} b \\ x^* &= -\frac{1}{2} A^\top \left(-2 (A A^\top)^{-1} b\right) \\ &= A^\top (A A^\top)^{-1} b \end{aligned}$$

This motivates the construction of the *Lagrangian dual problem*:

$$d^* = \max_{\lambda \geq 0} \ell(\lambda, \mu)$$

This is the case even when the problem is not convex.

Example 72 (Partitioning Problem). Consider the problem

$$p^* = \min_{x_i^2=1} x^\top W x$$

where $W \in \mathbb{S}^n$. The Lagrangian is

$$\begin{aligned} \mathcal{L}(x, \mu) &= x^\top W x + \sum_{i=1}^n \mu_i (x_i^2 - 1) \\ &= x^\top (W + \text{diag}(\mu)) x - \sum_{i=1}^n \mu_i \end{aligned}$$

If $W + \text{diag}(\mu)$ is not positive semidefinite, then $\inf_x \mathcal{L}(x, \mu) = -\infty$. Now assume that $W + \text{diag}(\mu)$ is positive semidefinite, so we can pick x such that $x^\top (W + \text{diag}(\mu)) x = 0$, so $\ell(\mu) = \inf_x \mathcal{L}(x, \mu) = -\sum_{i=1}^n \mu_i$. The dual problem is $\max_{W + \text{diag}(\mu) \succeq 0} \ell(\mu)$. A lower bound is $\mu = -\lambda_{\min}(W)1$, giving $d^* \geq n = \lambda_{\min}(W)$.

Formally, for a Lagrangian $\mathcal{L}(x, \lambda, \mu)$ and $\ell(\lambda, \mu) = \min_x \mathcal{L}(x, \lambda, \mu)$, the primal problem is

$$\begin{aligned} p^* &= \min_x f(x) \\ \text{s.t. } g(x) &\leq 0 \\ h(x) &= 0 \end{aligned}$$

and the dual problem is

$$\begin{aligned} d^* &= \max_{\lambda, \mu} \ell(\lambda, \mu) \\ \text{s.t. } \lambda &\geq 0 \end{aligned}$$

We know that $d^* \leq p^*$; this is called *weak duality*. Under some conditions $d^* = p^*$; this is called *strong duality*. The difference $p^* - d^*$ is the *duality gap*.

Theorem 73 (Min-Max Theorem)

For any sets X, Y , and any function $f: X \times Y \rightarrow \mathbb{R}$,

$$\min_{x \in X} \max_{y \in Y} f(x, y) \geq \max_{y \in Y} \min_{x \in X} f(x, y).$$

Proof. Fix $x_0 \in X$ and $y_0 \in Y$. Define $h(y) = \min_{x \in X} f(x, y)$, and similarly define $g(x) = \max_{y \in Y} f(x, y)$. Then

$$\begin{aligned} h(y_0) &= \min_{x \in X} f(x, y_0) \\ &\leq f(x_0, y_0) \\ &\leq \max_{y \in Y} f(x_0, y) \\ &\leq g(x_0) \end{aligned}$$

Therefore for each x and y , $h(y) \leq g(x)$, so

$$\begin{aligned} \max_{y \in Y} h(y) &\leq \min_{x \in X} g(x) \\ \max_{y \in Y} \min_{x \in X} f(x, y) &\leq \min_{x \in X} \max_{y \in Y} f(x, y) \end{aligned}$$

as desired. □

Clearly this is relevant in terms of duality. Consider for simplicity the problem

$$\begin{aligned} p^* &= \min_{\substack{x \\ g(x) \leq 0}} f(x) \\ \mathcal{L}(x, \lambda) &= f(x) + \lambda^\top g(x) \\ \max_{\lambda \geq 0} \mathcal{L}(x, \lambda) &= \max_{\lambda \geq 0} (f(x) + \lambda^\top g(x)) \\ &= \begin{cases} \infty, & g(x) \neq 0 \\ f(x), & g(x) \leq 0 \end{cases} \end{aligned}$$

$$\begin{aligned}
p^* &= \min_x \max_{\lambda \geq 0} \mathcal{L}(x, \lambda) \\
d^* &= \max_{\lambda \geq 0} \ell(\lambda) \\
&= \max_{\lambda \geq 0} \min_x \mathcal{L}(x, \lambda) \\
p^* &\geq d^*
\end{aligned}$$

which proves weak duality from the min-max theorem.

We want to find where strong duality holds, i.e. where $p^* = d^*$, which is true sometimes.

Theorem 74 (Slader's Condition)

For a convex problem

$$\begin{aligned}
p^* &= \min_x f(x) \\
\text{s.t. } g(x) &\leq 0 \\
h(x) &= 0
\end{aligned}$$

where there exists a point x_0 such that $g(x) < 0$ (a strictly feasible region), then strong duality holds.

A comment on strict feasibility: this means that, if the feasible region is X , the condition implies that $\overline{X} - X \neq \emptyset$, that is, the set is not only its limit points.

Theorem 75 (Refined Slader's Condition)

For a convex problem

$$\begin{aligned}
p^* &= \min_x f(x) \\
\text{s.t. } g(x) &\leq 0 \\
h(x) &= 0
\end{aligned}$$

where there exists x_0 such that for all affine constraints g_i we have $g_i(x) \leq 0$ and all other constraints we have $g_i(x) < 0$, strong duality holds.

Duality is of special interest in linear programs. Suppose we have the linear program

$$\begin{aligned}
p^* &= \min_x c^\top x \\
\text{s.t. } Ax &\leq b
\end{aligned}$$

Clearly, each constraint is linear or affine, so strong duality always holds (as long as the linear program has a feasible region).

The Lagrangian is

$$\begin{aligned}
\mathcal{L}(x, \lambda) &= c^\top x + \lambda^\top (Ax - b) \\
&= (A^\top \lambda + c)^\top x - \lambda^\top b \\
\ell(\lambda) &= \min_x \mathcal{L}(x, \lambda) \\
&= \min_x (A^\top \lambda + c)^\top x - \lambda^\top b
\end{aligned}$$

$$= \begin{cases} -\infty, & A^T \lambda + c \neq 0 \\ -b^T \lambda, & A^T \lambda + c = 0 \end{cases}$$

The dual is then

$$\begin{aligned} d^* &= \max_{\lambda \geq 0} \ell(\lambda) \\ &= \max_{\substack{\lambda \geq 0 \\ A^T \lambda + c = 0}} -b^T \lambda \end{aligned}$$

One can see that strong duality holds here, as long as the feasible region is nonempty.

We finish with some intuition about Slater's condition. Consider the optimization

$$\begin{aligned} p^* &= \min_x f(x) \\ \text{s.t. } g(x) &\leq 0 \\ h(x) &= 0 \end{aligned}$$

and consider the set of ordered tuples $G = \{g(x), f(x)\}$. Then $p^* = \min \{t \mid (u, t) \in G, u \leq 0\}$. The Lagrangian is $\mathcal{L}(u, t, \lambda) = t + \lambda^T u$, an affine hyperplane. In fact these hyperplanes are *supporting hyperplanes* which have at least one point on the curve, and separate it from another half-plane.