

V. Beyond Least Squares

Norm approximation

The bulk of this course has been spent studying variations of the least-squares problem: given observations $\mathbf{y} \in \mathbb{R}^M$ and a known $M \times N$ matrix \mathbf{A} , we solve

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

When \mathbf{A} has full column rank ($M \geq N$ and $\text{rank}(\mathbf{A}) = N$), then this problem has a unique solution that can be written in closed form: $\hat{\mathbf{x}} = \mathbf{A}^\dagger \mathbf{y}$, where $\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ is the pseudo-inverse of \mathbf{A} .

We did not make this explicit, but we also solved the “closest point” problem using the least-squares program above. Here we ask what is the closest point in an N dimensional subspace \mathcal{T} to a given point $\mathbf{y} \in \mathbb{R}^M$; that is, we want to solve

$$\underset{\mathbf{w} \in \mathcal{T}}{\text{minimize}} \|\mathbf{y} - \mathbf{w}\|_2^2.$$

If we have a basis $\mathbf{v}_1, \dots, \mathbf{v}_N$ for \mathcal{T} , then we are looking for the closest point in the span of $\{\mathbf{v}_n\}$ to \mathbf{y} . As every point in \mathcal{T} is a linear combination of these vectors, we can express any $\mathbf{w} \in \mathcal{T}$ as $\mathbf{w} = \mathbf{V}\mathbf{a}$ for some $\mathbf{a} \in \mathbb{R}^N$. So we can transform the problem above into

$$\underset{\mathbf{a} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{V}\mathbf{a}\|_2^2,$$

which of course has the solution $\hat{\mathbf{a}} = (\mathbf{V}^T \mathbf{V})^{-1} \mathbf{V}^T \mathbf{y}$. The corresponding optimal $\hat{\mathbf{w}} \in \mathbb{R}^M$ is then

$$\hat{\mathbf{w}} = \mathbf{V}\hat{\mathbf{a}} = \mathbf{V}(\mathbf{V}^T \mathbf{V})^{-1} \mathbf{V}^T \mathbf{y}.$$

We also know that this “approximation in a subspace” problem is equally easy if we use any norm which is induced by an inner product.

If $\langle \cdot, \cdot \rangle_H$ is a valid inner product, and $\|\cdot\|_H$ is its induced norm, then we can solve

$$\min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|_H$$

by forming

$$\mathbf{G} = \begin{bmatrix} \langle \mathbf{v}_1, \mathbf{v}_1 \rangle_H & \langle \mathbf{v}_1, \mathbf{v}_2 \rangle_H & \cdots & \langle \mathbf{v}_1, \mathbf{v}_N \rangle_H \\ \langle \mathbf{v}_2, \mathbf{v}_1 \rangle_H & \langle \mathbf{v}_2, \mathbf{v}_2 \rangle_H & & \langle \mathbf{v}_2, \mathbf{v}_N \rangle_H \\ \vdots & & & \\ \langle \mathbf{v}_N, \mathbf{v}_1 \rangle_H & \langle \mathbf{v}_N, \mathbf{v}_2 \rangle_H & \cdots & \langle \mathbf{v}_N, \mathbf{v}_N \rangle_H \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \langle \mathbf{y}, \mathbf{v}_1 \rangle_H \\ \langle \mathbf{y}, \mathbf{v}_2 \rangle_H \\ \vdots \\ \langle \mathbf{y}, \mathbf{v}_N \rangle_H \end{bmatrix}$$

and then taking $\hat{\mathbf{w}} = \mathbf{V}\mathbf{G}^{-1}\mathbf{b}$.

What if we are interested in finding the closest point in a subspace in a norm that is **not** induced by an inner product? For example¹, we may be interested in solving

$$\underset{\mathbf{w} \in \mathcal{T}}{\text{minimize}} \|\mathbf{y} - \mathbf{w}\|_1$$

or

$$\underset{\mathbf{w} \in \mathcal{T}}{\text{minimize}} \|\mathbf{y} - \mathbf{w}\|_\infty.$$

We can still (partially) characterize the solutions to these types of problems, but in general they do not have closed-form solutions. There is, however, a nice computational framework for solving them — in practice, they are not too much harder to solve than the least-squares problem.

¹It is a fact that the only ℓ_p norm that is induced by an inner product is the ℓ_2 norm.

Example: Let \mathcal{T} be the subspace of \mathbb{R}^2 defined by

$$\mathcal{T} = \{\mathbf{w} \in \mathbb{R}^2 : w[1] + 2w[2] = 0\}.$$

Given

$$\mathbf{y} = \begin{bmatrix} 3 \\ 2 \end{bmatrix},$$

find the minimizer \mathbf{w} of these three programs:

$$\min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|_2, \quad \min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|_1, \quad \min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|_\infty.$$

Example: Filter design

The standard “filter synthesis” problem is to find an finite-impulse response (FIR) filter whose discrete-time Fourier transform (DTFT) is as close to some target $H_*(\omega)$ as possible. When the deviation from the optimal response is measured using a uniform error, this is call “equiripple design”, since the error in the solution will tend to have ripples a uniform distance away from the ideal. That is, we would like to solve

$$\min_H \sup_{\omega \in [-\pi, \pi]} |H_*(\omega) - H(\omega)|, \quad \text{subject to } H(\omega) \text{ being FIR}$$

If we restrict ourselves to the case where $H_*(\omega)$ has linear phase (so the impulse response is symmetric around some time index)² we can recast this as a Chebyshev approximation problem.

A symmetric filter with $2K + 1$ taps has a real DTFT that can be written as a superposition of a DC term plus K cosines:

$$h_n = 0 \quad |n| > K \quad \Rightarrow \quad H(\omega) = \sum_{k=0}^K \tilde{h}_k \cos(k\omega), \quad \tilde{h}_k = \begin{cases} h_0, & k = 0 \\ 2h_k, & 1 \leq k \leq K. \end{cases}$$

So we are trying to solve

$$\min_{\mathbf{x}} \sup_{\omega \in [-\pi, \pi]} \left| H_*(\omega) - \sum_{k=0}^K x_k \cos(k\omega) \right|.$$

²The case with general phase can also be handled using convex optimization, but it is not naturally stated as a linear program.

We will approximate the supremum on the inside by measuring it at M equally spaced points $\omega_1, \dots, \omega_M$ between $-\pi$ and π . Then

$$\min_{\mathbf{x}} \max_{\omega_m} \left| H_*(\omega_m) - \sum_{k=0}^K x_k \cos(k\omega_m) \right| = \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|_{\infty},$$

where $\mathbf{y} \in \mathbb{R}^M$ and the $M \times (K + 1)$ matrix \mathbf{F} are defined as

$$\mathbf{y} = \begin{bmatrix} H_*(\omega_1) \\ H_*(\omega_2) \\ \vdots \\ H_*(\omega_M) \end{bmatrix} \quad \mathbf{F} = \begin{bmatrix} 1 & \cos(\omega_1) & \cos(2\omega_1) & \cdots & \cos(K\omega_1) \\ 1 & \cos(\omega_2) & \cos(2\omega_2) & \cdots & \cos(K\omega_2) \\ \vdots & & \ddots & & \\ 1 & \cos(\omega_M) & \cos(2\omega_M) & \cdots & \cos(K\omega_M) \end{bmatrix}$$

It should be noted that since the ω_m are equally spaced, the matrix \mathbf{F} (and its adjoint) can be applied efficiently using a fast discrete cosine transform. Just as with least-squares (recall steepest descent and conjugate gradients), this has a direct impact on the number of computations we need to solve the Chebyshev approximation problem above.

ℓ_1 and ℓ_∞ minimization

As we can see from the example above, minimizing $\|\mathbf{y} - \mathbf{Ax}\|_1$ or $\|\mathbf{y} - \mathbf{Ax}\|_\infty$ is doing something quite different than least squares. We can say the following qualitative things about these programs:

- Minimizing $\|\mathbf{y} - \mathbf{Ax}\|_1$ encourages the the residual $\mathbf{r} = \mathbf{y} - \mathbf{Ax}$ to be *sparse*. That is, many of the components of the optimal \mathbf{r} will be exactly equal to zero.
- Minimizing $\|\mathbf{y} - \mathbf{Ax}\|_\infty$ encourages the residual \mathbf{r} to be *uniform*. That is, many of the components of the optimal \mathbf{r} will have exactly the same magnitude.

In both of these cases, it is also easy to see that the solution need not be unique (although for “generic” problem instances, it will be — you need the column space of \mathbf{A} and \mathbf{y} to line up just right for the solution to be non-unique).

We can actually characterize these properties a little more carefully. Here is concrete result for ℓ_1 :

Theorem: Given $\mathbf{y} \in \mathbb{R}^M$ and an $M \times N$ matrix \mathbf{A} with full column rank, there is a solution $\hat{\mathbf{x}}_1$ to

$$\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{Ax}\|_1 \tag{1}$$

such that $\hat{\mathbf{r}}_1 = \mathbf{y} - \mathbf{A}\hat{\mathbf{x}}_1$ has at most $M - N$ non-zero terms (and so at least N terms which are exactly equal to zero).

Proof: Suppose \mathbf{x} is a vector such that $\mathbf{r} = \mathbf{y} - \mathbf{Ax}$ has more than $M - N$ non-zero terms. We will show that there is another vector $\mathbf{x}' \neq \mathbf{x}$ such that one of two things is true:

1. $\|\mathbf{y} - \mathbf{Ax}'\|_1 < \|\mathbf{y} - \mathbf{Ax}\|_1$, and so \mathbf{x} can not be the minimizer;

2. $\|\mathbf{y} - \mathbf{A}\mathbf{x}'\|_1 = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1$, and $\mathbf{r}' = \mathbf{y} - \mathbf{A}\mathbf{x}'$ has at least one fewer non-zero term than \mathbf{r} .

Let $\Gamma \subset \{1, \dots, M\}$ be the subset of locations where $\mathbf{y} - \mathbf{A}\mathbf{x}$ is equal to zero. For now we will assume that Γ is not empty, we will remove this restriction below. By assumption, Γ has fewer than N members. We will use $|\Gamma|$ the number of terms in Γ , so $|\Gamma| < N$.

Let \mathbf{A}_Γ be the $|\Gamma| \times N$ matrix formed by extracting the rows of \mathbf{A} corresponding to Γ . Since $|\Gamma| < N$, this matrix has fewer rows than columns, and so it has a non-trivial null space. Let $\mathbf{d} \in \mathbb{R}^N$ be any non-zero vector in this null space. By construction, $\mathbf{y} - \mathbf{A}(\mathbf{x} + \delta\mathbf{d})$ will also be 0 on the set Γ for all δ .

Let Γ^c be the complement of Γ (the set of all indices that are not in Γ). Define $\mathbf{a}_m \in \mathbb{R}^N$ so that \mathbf{a}_m^T is the m th row of \mathbf{A} . Then for any δ , we can write

$$\begin{aligned} \|\mathbf{y} - \mathbf{A}(\mathbf{x} + \delta\mathbf{d})\|_1 &= \sum_{m \in \Gamma^c} |y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}| \\ &= \sum_{m \in \Gamma^c} \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d})(y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}). \end{aligned}$$

Since δ is arbitrary, we can make its magnitude small enough so that

$$\text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}) = \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x}) \quad \text{for all } m \in \Gamma^c, \quad (2)$$

and so for appropriately small δ ,

$$\begin{aligned} \|\mathbf{y} - \mathbf{A}(\mathbf{x} + \delta\mathbf{d})\|_1 &= \sum_{m \in \Gamma^c} \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x})(y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}) \\ &= \sum_{m=1}^M \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x})(y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}), \end{aligned}$$

where we have added back in the terms on Γ since by construction $y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d} = 0$ for $m \in \Gamma$. We can now break that sum

into two parts, yielding

$$\|\mathbf{y} - \mathbf{A}(\mathbf{x} + \delta\mathbf{d})\|_1 = \|\mathbf{y} - \mathbf{Ax}\|_1 - \delta \sum_{m=1}^M \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x}) \mathbf{a}_m^T \mathbf{d}. \quad (3)$$

If the sum on the right above is exactly equal to zero, we can choose any δ that maintains (2). This equality is maintained exactly until an additional term in $y[m] - \mathbf{a}_m^T \mathbf{x} - \delta \mathbf{a}_m^T \mathbf{d}$ is driven to zero (and this is guaranteed to happen for some $\delta = \delta_0$), meaning that condition (2) above holds for $\mathbf{x}' = \mathbf{x} + \delta_0 \mathbf{d}$: $\|\mathbf{y} - \mathbf{Ax}'\|_1 = \|\mathbf{y} - \mathbf{Ax}\|_1$ and there is at least one fewer nonzero term in $\mathbf{r}' = \mathbf{y} - \mathbf{Ax}'$ than in \mathbf{r} .

Now consider the case where the sum on the right in (3) is nonzero. Since the sign of δ is arbitrary, so we can choose so that $\delta \sum_{m=1}^M \text{sign}(y[m] - \mathbf{a}_m^T \mathbf{x}) \mathbf{a}_m^T \mathbf{d} > 0$, and so

$$\|\mathbf{y} - \mathbf{A}(\mathbf{x} + \delta\mathbf{d})\|_1 < \|\mathbf{y} - \mathbf{Ax}\|_1.$$

Hence, \mathbf{x} cannot be a solution to (1) since there is another vector $\mathbf{x}' = \mathbf{x} + \delta\mathbf{d}$ that results in a smaller value of the functional. This establishes the theorem.

The key part of finding \mathbf{x}' was the fact that \mathbf{A}_Γ was guaranteed to have a non-trivial null space. This was true simply because it had more columns than rows — if $\Gamma \geq N$ (meaning the solution has at least N terms exactly equal to zero), no such guarantee exists.

A similar statement can be made for ℓ_∞ norm approximation — we may explore this more on the homework.

ℓ_1 and ℓ_∞ as linear programs

Unlike the ℓ_2 case, the solutions to the problems

$$\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1 \quad \text{and} \quad \min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_\infty,$$

do not have closed form solutions, and the mapping from the data \mathbf{y} to the solution $\hat{\mathbf{x}}$ is not linear at all. In addition, these functionals are not differentiable, and so the gradient is not defined at every point, meaning that we cannot use a standard gradient descent algorithms.

There is, however, a very well-defined computational framework for solving these two problems. They can be recast as **linear programs**. This is a type of (convex) optimization program that appears all over applied mathematics (operations research, statistics, machine learning, etc), and so there has been a lot of work in the last 30 on efficient methods for solving them.

A **linear program** (LP) is an optimization program of the form

$$\min_{\mathbf{z} \in \mathbb{R}^Q} \mathbf{c}^T \mathbf{z} \quad \text{subject to} \quad \mathbf{M}\mathbf{z} \leq \mathbf{b}. \quad (4)$$

Above, $\mathbf{c} \in \mathbb{R}^Q$, \mathbf{M} is a $P \times Q$ matrix, and $\mathbf{b} \in \mathbb{R}^P$. The nomenclature comes from the fact that both the functional $\mathbf{c}^T \mathbf{z}$ and the constraints $\mathbf{M}\mathbf{z}$ are linear in \mathbf{z} . Notice that the presence of the constraints are necessary for this even to make sense, as $\mathbf{c}^T \mathbf{z}$ is unbounded below if left unconstrained.

There are several ways in which the ℓ_1 (and ℓ_∞) norm approximation problem can be recast as a linear program; we will illustrate one of them here. The trick we will use is to introduce a dummy variable

$\mathbf{u} \in \mathbb{R}^M$ whose entries bound the magnitudes of the residuals:

$$|y[m] - \mathbf{a}_m^T \mathbf{x}| \leq u[m], \quad \text{for all } m = 1, \dots, M.$$

We can write this as

$$-u[m] \leq y[m] - \mathbf{a}_m^T \mathbf{x} \leq u[m],$$

or combining all the constraints together,

$$-\mathbf{u} \leq \mathbf{y} - \mathbf{A}\mathbf{x} \leq \mathbf{u}$$

which indicates that for a pair \mathbf{x}, \mathbf{u} to be feasible, we will need $\mathbf{u} \geq 0$. Then

$$\min_{\substack{\mathbf{x} \in \mathbb{R}^N \\ \mathbf{u} \in \mathbb{R}^M}} \sum_{m=1}^M u[m] \quad \text{subject to} \quad \begin{array}{l} -\mathbf{A}\mathbf{x} - \mathbf{u} \leq -\mathbf{y} \\ \mathbf{A}\mathbf{x} - \mathbf{u} \leq \mathbf{y} \end{array} \quad (5)$$

will have $|y[m] - \mathbf{a}_m^T \hat{\mathbf{x}}| = \hat{u}[m]$ at its minimum, which means the $\hat{\mathbf{x}}$ in the solution above will be the same as the $\hat{\mathbf{x}}$ which solves (1). (The constraint $\mathbf{u} \geq 0$ is implicit above, and so does not need to be incorporated explicitly.) We can also see that (5) is a linear program of the form (4) with $Q = M + N$ variables and $P = 2M$ constraints, and

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} -\mathbf{A} & -\mathbf{I} \\ \mathbf{A} & -\mathbf{I} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -\mathbf{y} \\ \mathbf{y} \end{bmatrix}.$$

Recasting $\min \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_\infty$ as a linear program is even easier. Now we just need a single dummy variable $t \in \mathbb{R}$, which constrains

$$|y[m] - \mathbf{a}_m^T \mathbf{x}| \leq t \quad \text{for all } m = 1, \dots, M.$$

Then

$$\min_{\substack{\mathbf{x} \in \mathbb{R}^N \\ t \in \mathbb{R}}} t \quad \text{subject to} \quad \begin{aligned} -\mathbf{A}\mathbf{x} - t\mathbf{1} &\leq -\mathbf{y} \\ \mathbf{A}\mathbf{x} - t\mathbf{1} &\leq \mathbf{y}, \end{aligned}$$

will have $\hat{t} = \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_\infty$ at its minimum, and so the $\hat{\mathbf{x}}$ in the solution above will match the solution to $\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_\infty$. This is a linear program with $Q = N + 1$ variables and $P = 2M$ constraints — we can plug into (4) with

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ t \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \quad \mathbf{M} = \begin{bmatrix} -\mathbf{A} & -\mathbf{1} \\ \mathbf{A} & -\mathbf{1} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -\mathbf{y} \\ \mathbf{y} \end{bmatrix}.$$

Convexity and uniqueness

From our geometric reasoning above, we can see that the ℓ_1 and ℓ_∞ norm approximation problems do not always have unique solutions. For example:

$$\min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|_\infty, \quad \mathcal{T} = \text{span} \left(\begin{bmatrix} 0 \\ 1 \end{bmatrix} \right), \quad \mathbf{y} = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

has solutions

$$\hat{\mathbf{w}} = \begin{bmatrix} 0 \\ \beta \end{bmatrix} \quad \text{for all } \beta \in [-1, 5].$$

But it is true that all of the solutions for both of these problems will be clustered together in a connected set.

In general, for any valid norm $\|\cdot\|$, the approximation problem

$$\min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|, \quad \mathcal{T} = \text{subspace}$$

is **convex**. Here are the require formal definitions:

Convexity: A subset \mathcal{U} of a vector space is called convex if

$$\mathbf{x}, \mathbf{y} \in \mathcal{U} \quad \Rightarrow \quad (1 - \lambda)\mathbf{x} + \lambda\mathbf{y} \in \mathcal{U} \quad \text{for all } 0 \leq \lambda \leq 1.$$

A real-valued functional on a vector space $f(\cdot) : \mathcal{S} \rightarrow \mathbb{R}$ is called convex if

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) \quad \text{for all } 0 \leq \lambda \leq 1.$$

If we replace the \leq above with a strict inequality,

$$f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) < (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y}) \quad \text{for all } 0 < \lambda < 1,$$

for all $\mathbf{x} \neq \mathbf{y}$, then $f(\cdot)$ is called strictly convex. There is a similar definition³ for sets: $\mathcal{U} \subset \mathcal{S}$ is called strictly convex if for all $\mathbf{x}, \mathbf{y} \in \mathcal{U}$

$$(1 - \lambda)\mathbf{x} + \lambda\mathbf{y} \in \text{interior}(\mathcal{U}), \quad 0 < \lambda < 1.$$

An optimization program in a vector space \mathcal{S} is called convex (or a “convex program”) if it has the form

$$\min f(\mathbf{x}), \quad \mathbf{x} \in \mathcal{U},$$

where $\mathcal{U} \subset \mathcal{S}$ is a convex set, and $f(\cdot)$ is a convex function.

The main consequence of an optimization program being convex is that all local minima are also global minima. This fact is extremely useful computationally, since we know a lot of simple ways that are guaranteed to find local minima (e.g. steepest descent).

Fact: Any valid norm $\|\cdot\|$ on a vector space \mathcal{S} is a convex function. You can prove this at home; it is a direct consequence of the triangle inequality.

Fact: For any valid norm $\|\cdot\|$ on a vector space \mathcal{S} , the unit ball

$$\mathcal{B} = \{\mathbf{x} \in \mathcal{S} : \|\mathbf{x}\| \leq 1\}$$

is convex. This is also easy enough to prove that you can do it at home in your spare time.

³We will need to recall the definition of the interior of a set. We say that $\mathbf{z} \in \text{interior}(\mathcal{U})$ if there exists any $\epsilon > 0$ such that the neighborhood around \mathbf{z} given by $\mathcal{B}_\epsilon(\mathbf{z}) = \{\mathbf{z}' \in \mathcal{S} : \|\mathbf{z} - \mathbf{z}'\| \leq \epsilon\}$ is also included in \mathcal{U} , $\mathcal{B}_\epsilon(\mathbf{z}) \subset \mathcal{U}$.

Examples:

- The ℓ_1 norm (and unit ball) is convex, but not strictly convex.
- The ℓ_2 norm (and unit ball) is strictly convex.
- The ℓ_∞ norm (and unit ball) is convex, but not strictly convex.

In fact, the ℓ_p norm (and unit ball) is strictly convex for all $1 < p < \infty$. This fact follows directly from the **Minkowski Inequality**⁴. If $\{x_n\}$ and $\{y_n\}$ are sequences of real numbers, then for any $1 < p < \infty$,

$$\left(\sum_{n=1}^N |x_n + y_n|^p \right)^{1/p} \leq \left(\sum_{n=1}^N |x_n|^p \right)^{1/p} + \left(\sum_{n=1}^N |y_n|^p \right)^{1/p},$$

with equality holding only when there exists an $\alpha \in \mathbb{R}$ such that

$$x_n = \alpha y_n \quad \text{for all } n = 1, \dots, N.$$

Uniqueness of ℓ_p norm minimization

In a vector space \mathcal{S} with a strictly convex norm $\|\cdot\|$, the program

$$\min_{\mathbf{w} \in \mathcal{T}} \|\mathbf{y} - \mathbf{w}\|, \quad \mathcal{T} = \text{subspace},$$

has a unique solution for all $\mathbf{y} \in \mathcal{S}$.

Proof: Consider two vectors $\mathbf{w}_1, \mathbf{w}_2 \in \mathcal{T}$ that have equal distance to \mathbf{y} :

$$\|\mathbf{y} - \mathbf{w}_1\| = \|\mathbf{y} - \mathbf{w}_2\| = \rho.$$

⁴Proving this would make another good homework problem, but if you are impatient you can look at http://en.wikipedia.org/wiki/Minkowski_inequality.

We first show that if $\mathbf{w}_1 \neq \mathbf{w}_2$, we can always find another point in \mathcal{T} that is even closer to \mathbf{y} . By strict convexity,

$$\|0.5(\mathbf{y} - \mathbf{w}_1) + 0.5(\mathbf{y} - \mathbf{w}_2)\| < \rho,$$

and so re-arranging the terms inside the norm above,

$$\left\| \mathbf{y} - \frac{\mathbf{w}_1 + \mathbf{w}_2}{2} \right\| < \rho,$$

and so $\mathbf{w}_3 = (\mathbf{w}_1 + \mathbf{w}_2)/2$, which of course is also a member of \mathcal{T} , is closer to \mathbf{y} than both \mathbf{w}_1 and \mathbf{w}_2 . Thus the closest point to \mathbf{y} in \mathcal{T} must be unique.

Closest point in a convex set

Let $\mathbf{x}_0 \in \mathbb{R}^N$ be given, and let \mathcal{C} be a non-empty, closed, convex set. The **projection** of \mathbf{x}_0 onto \mathcal{C} is the closest point (in the standard Euclidean distance, for now) in \mathcal{C} to \mathbf{x}_0 :

$$P_{\mathcal{C}}(\mathbf{x}_0) = \arg \min_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x}_0 - \mathbf{y}\|_2$$

We will see below that there is a unique minimizer to this problem, and that the solution has geometric properties that are analogous to the case where \mathcal{C} is a subspace.

Projection onto a subspace

Let's recall how we solve this problem in the special case where $\mathcal{C} := \mathcal{T}$ is a K -dimensional *subspace*. In this case, the solution $\hat{\mathbf{x}} = P_{\mathcal{T}}(\mathbf{x}_0)$ is unique, and is characterized by the **orthogonality principle**:

$$\mathbf{x}_0 - \hat{\mathbf{x}} \perp \mathcal{T}$$

meaning that $\langle \mathbf{y}, \mathbf{x}_0 - \hat{\mathbf{x}} \rangle = 0$ for all $\mathbf{y} \in \mathcal{T}$.

Affine sets are not fundamentally different than subspaces. Any affine set \mathcal{C} can be written as a subspace \mathcal{T} plus an offset \mathbf{v}_0 :

$$\mathcal{C} = \mathcal{T} + \mathbf{v}_0 = \{\mathbf{x} : \mathbf{x} = \mathbf{y} + \mathbf{v}_0, \mathbf{y} \in \mathcal{T}\}.$$

It is easy to translate the results for subspaces above to say that the projection onto an affine set is unique, and obeys the orthogonality principle

$$\langle \mathbf{y} - \hat{\mathbf{x}}, \mathbf{x}_0 - \hat{\mathbf{x}} \rangle = 0, \quad \text{for all } \mathbf{y} \in \mathcal{C}. \quad (6)$$

You can solve this problem by shifting \mathbf{x}_0 and \mathcal{C} by negative \mathbf{v}_0 , projecting $\mathbf{x}_0 - \mathbf{v}_0$ onto the subspace $\mathcal{C} - \mathbf{v}_0$, and then shifting the answer back.

Projection onto a general convex set

In general, there is no closed-form expression for the projector onto a given convex set. However, the concepts above (orthogonality, projection onto a subspace) can help us understand the solution for an arbitrary convex set.

Uniqueness If \mathcal{C} is closed and convex, then for any \mathbf{x}_0 , the program

$$\min_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x}_0 - \mathbf{y}\|_2 \quad (7)$$

has a **unique** solution.

First of all, that there is at least one point in \mathcal{C} such that the minimum is obtained is a consequence of \mathcal{C} being closed and of $\|\mathbf{x}_0 - \mathbf{y}\|_2$ being a continuous function of \mathbf{y} . Let $\hat{\mathbf{x}}$ be one such minimizer; we will show that $\|\mathbf{x}_0 - \mathbf{y}\|_2 > \|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2$ for all $\mathbf{y} \in \mathcal{C}$, $\mathbf{y} \neq \hat{\mathbf{x}}$.

Consider first all the points in \mathcal{C} which are co-aligned with $\hat{\mathbf{x}}$. Let

$$\mathcal{I} = \{\alpha \in \mathbb{R} : \alpha \hat{\mathbf{x}} \in \mathcal{C}\}.$$

Since \mathcal{C} is convex and closed, this is a closed interval of the real line (that contains at least the point $\alpha = 1$). The function

$$g(\alpha) = \|\mathbf{x}_0 - \alpha \hat{\mathbf{x}}\|_2^2 = \alpha^2 \|\hat{\mathbf{x}}\|_2^2 - 2\alpha \langle \hat{\mathbf{x}}, \mathbf{x}_0 \rangle + \|\mathbf{x}_0\|_2^2,$$

is minimized at $\alpha = 1$ by construction, and since its second derivative is strictly positive, we know that this is the unique minima. So if there is another minimizer of (7), then it is not co-aligned with $\hat{\mathbf{x}}$.

Now let \mathbf{y} be any point in \mathcal{C} that is not co-aligned with $\hat{\mathbf{x}}$. We will show that \mathbf{y} cannot minimize (7) because the point $\hat{\mathbf{x}}/2 + \mathbf{y}/2 \in \mathcal{C}$

is definitively closer to \mathbf{x}_0 . We have

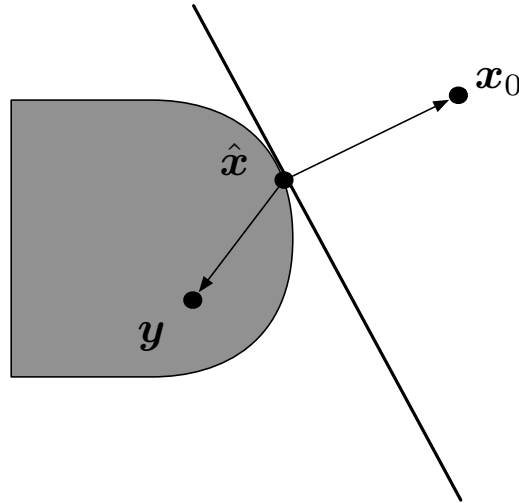
$$\begin{aligned}
 \left\| \mathbf{x}_0 - \frac{\hat{\mathbf{x}}}{2} - \frac{\mathbf{y}}{2} \right\|_2^2 &= \left\| \frac{\mathbf{x}_0 - \hat{\mathbf{x}}}{2} + \frac{\mathbf{x}_0 - \mathbf{y}}{2} \right\|_2^2 \\
 &= \frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2^2}{4} + \frac{\|\mathbf{x}_0 - \mathbf{y}\|_2^2}{4} + \frac{\langle \mathbf{x}_0 - \hat{\mathbf{x}}, \mathbf{x}_0 - \mathbf{y} \rangle}{2} \\
 &< \frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2^2}{4} + \frac{\|\mathbf{x}_0 - \mathbf{y}\|_2^2}{4} + \frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2 \|\mathbf{x}_0 - \mathbf{y}\|_2}{2} \\
 &= \left(\frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2}{2} + \frac{\|\mathbf{x}_0 - \mathbf{y}\|_2}{2} \right)^2 \\
 &\leq \|\mathbf{x}_0 - \mathbf{y}\|_2^2.
 \end{aligned}$$

The strict inequality above follows from Cauchy-Schwarz, while the last inequality follows from the fact that $\hat{\mathbf{x}}$ is a minimizer. This shows that no $\mathbf{y} \neq \hat{\mathbf{x}}$ can also minimize (7), and so $\hat{\mathbf{x}}$ is unique.

Obtuseness. Relative to the solution $\hat{\mathbf{x}}$, every vector in \mathcal{C} is at an obtuse angle to the error $\mathbf{x}_0 - \hat{\mathbf{x}}$. More precisely, $P_{\mathcal{C}}(\mathbf{x}_0) = \hat{\mathbf{x}}$ if and only if

$$\langle \mathbf{y} - \hat{\mathbf{x}}, \mathbf{x}_0 - \hat{\mathbf{x}} \rangle \leq 0. \quad (8)$$

Compare with (6) above. Here is a picture:



We first prove that (8) \Rightarrow $P_{\mathcal{C}}(\mathbf{x}_0) = \hat{\mathbf{x}}$. For any $\mathbf{y} \in \mathcal{C}$, we have

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}_0\|_2^2 &= \|\mathbf{y} - \hat{\mathbf{x}} + \hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 \\ &= \|\mathbf{y} - \hat{\mathbf{x}}\|_2^2 + \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 + 2\langle \mathbf{y} - \hat{\mathbf{x}}, \hat{\mathbf{x}} - \mathbf{x}_0 \rangle \\ &\geq \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 + 2\langle \mathbf{y} - \hat{\mathbf{x}}, \hat{\mathbf{x}} - \mathbf{x}_0 \rangle. \end{aligned}$$

Note that the inner product term above is the same as (8), but with the sign of the second argument flipped, so we know this term must be non-negative. Thus

$$\|\mathbf{y} - \mathbf{x}_0\|_2^2 \geq \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2.$$

Since this holds uniformly over all $\mathbf{y} \in \mathcal{C}$, $\hat{\mathbf{x}}$ must be the closest point in \mathcal{C} to \mathbf{x}_0 .

We now show that $P_{\mathcal{C}}(\mathbf{x}_0) = \hat{\mathbf{x}} \Rightarrow$ (8). Let \mathbf{y} be an arbitrary point in \mathcal{C} . Since \mathcal{C} is convex, the point $\hat{\mathbf{x}} + \theta(\mathbf{y} - \hat{\mathbf{x}})$ must also be in \mathcal{C} for all $\theta \in [0, 1]$. Since $\hat{\mathbf{x}} = P_{\mathcal{C}}(\mathbf{x}_0)$,

$$\|\hat{\mathbf{x}} + \theta(\mathbf{y} - \hat{\mathbf{x}}) - \mathbf{x}_0\|_2 \geq \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2 \quad \text{for all } \theta \in [0, 1]. \quad (9)$$

It is a basic geometric fact that that for any two vectors \mathbf{x}, \mathbf{y} ,

$$\text{if } \|\mathbf{x} + \theta\mathbf{y}\|_2 \geq \|\mathbf{x}\|_2 \text{ for all } \theta \in [0, 1] \text{ then } \langle \mathbf{x}, \mathbf{y} \rangle \geq 0. \quad (10)$$

Applying this to (9) means

$$\langle \mathbf{y} - \hat{\mathbf{x}}, \hat{\mathbf{x}} - \mathbf{x}_0 \rangle \geq 0 \quad \Rightarrow \quad \langle \mathbf{y} - \hat{\mathbf{x}}, \mathbf{x}_0 - \hat{\mathbf{x}} \rangle \leq 0.$$

Finally, to prove (10), we expand the norm as

$$\|\mathbf{x} + \theta\mathbf{y}\|_2^2 = \|\mathbf{x}\|_2^2 + \theta^2\|\mathbf{y}\|_2^2 + 2\theta\langle \mathbf{x}, \mathbf{y} \rangle,$$

from which we can immediately deduce that

$$\begin{aligned} \frac{\theta}{2}\|\mathbf{y}\|_2^2 + \langle \mathbf{x}, \mathbf{y} \rangle &\geq 0 \quad \text{for all } \theta \in [0, 1] \\ &\Rightarrow \langle \mathbf{x}, \mathbf{y} \rangle \geq 0. \end{aligned}$$