

CHAPTERS 3 - 5

Inner Products and Norms

The most basic inner product is the dot product

$$\vec{v} \cdot \vec{w} = v_1 w_1 + v_2 w_2 + \dots + v_n w_n$$

On \mathbb{R}^n , the dot product obeys:

1) Bilinearity $(c\vec{u} + d\vec{v}) \cdot \vec{w} = c(\vec{u} \cdot \vec{w}) + d(\vec{v} \cdot \vec{w})$
 $\vec{u} \cdot (c\vec{v} + d\vec{w}) = c(\vec{u} \cdot \vec{v}) + d(\vec{u} \cdot \vec{w})$

2) Symmetry $\vec{u} \cdot \vec{v} = \vec{v} \cdot \vec{u}$

3) Positivity $\vec{v} \cdot \vec{v} \geq 0$ and $\vec{v} \cdot \vec{v} = 0$ if
and only if $\vec{v} = \vec{0}$.

From the dot product we get the 2-norm

$$\|\vec{v}\|_2 = \sqrt{\vec{v} \cdot \vec{v}} = \left(v_1^2 + v_2^2 + \dots + v_n^2 \right)^{1/2}$$

Other norms:

$$\|\vec{v}\|_p = \left(v_1^p + v_2^p + \dots + v_n^p \right)^{1/p}$$

Most common cases: $p = 1, 2, \infty$.

$$\|\vec{v}\|_\infty = \max_{1 \leq i \leq n} |v_i|.$$

Other norms and inner product come up naturally throughout mathematics and physics. The text book details this but we will largely avoid any generalizations. 2

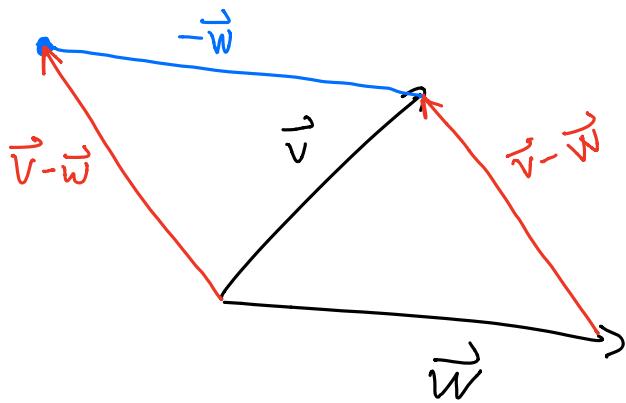
General properties of a norm:

- i) Positivity: $\|\vec{v}\| \geq 0$ and $\|\vec{v}\|=0$ if and only if $\vec{v}=\vec{0}$
- ii) Homogeneity: $\|c\vec{v}\|=|c|\|\vec{v}\|$.
- iii) Triangle inequality: $\|\vec{v}+\vec{w}\| \leq \|\vec{v}\| + \|\vec{w}\|$.

A norm measures the size of a vector.

$$\|\vec{v}-\vec{w}\|$$

measures the distance:



Matrices also have norms. We will encounter one later in the course.

We have three situations to consider:

i) Solving $A\vec{x} = \vec{b}$ when there is no solution.

ii) Solving $A\vec{x} = \vec{b}$ when there is no solution.

iii) Solving $A\vec{x} = \vec{b}$ when $\text{img } A \neq \mathbb{R}^m$ and $\ker A \neq \{\vec{0}\}$.

$$\begin{cases} \text{img } A \neq \mathbb{R}^m \\ \ker A = \{\vec{0}\} \end{cases}$$

$$\begin{cases} \ker A \neq \{\vec{0}\} \\ \text{img } A = \mathbb{R}^m \end{cases}$$

Some of these statements seem contradictory. But as we will see, we can make sense of them. But there are many ways to make sense of them.

Solving $A\vec{x} = \vec{b}$ when there is no solution

Since we cannot find \vec{x} such that $A\vec{x} = \vec{b}$, we find \vec{x} such that

$$\|A\vec{x} - \vec{b}\|_2$$

is as small as possible.

$$\text{i.e. } \|A\vec{x} - \vec{b}\|_2 < \|A\vec{y} - \vec{b}\| \text{ for } \vec{x} \neq \vec{y}.$$

Theorem Suppose A is $m \times n$ and $A^T A$ is invertible.⁴

For $\vec{x} = (A^T A)^{-1} A^T \vec{b}$ we have

$$\|A\vec{x} - \vec{b}\|_2 < \|A\vec{y} - \vec{b}\|_2 \quad \text{for every } \vec{y} \in \mathbb{R}^n, \vec{y} \neq \vec{x}.$$

This means that the solution of

$$A^T A \vec{x} = A^T \vec{b} \quad \leftarrow$$

is the unique minimum for $\|A\vec{y} - \vec{b}\|_2$.

Normal equations.

\vec{x} is called the least-squares solution

Proof It suffices to consider

$$\begin{aligned} \|A\vec{y} - \vec{b}\|_2^2 &= (A\vec{y} - \vec{b})^T (A\vec{y} - \vec{b}) \\ &= (\vec{y}^T A^T - \vec{b}^T) (A\vec{y} - \vec{b}) \\ &= \vec{y}^T A^T A \vec{y} - \vec{y}^T A^T \vec{b} - \vec{b}^T A \vec{y} + \vec{b}^T \vec{b} =: F(\vec{y}) \end{aligned}$$

$$\frac{\partial F}{\partial y_i} ?$$

$$\frac{\partial \vec{y}}{\partial y_i} = \vec{e}_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i\text{th entry.}$$

$$\begin{aligned}\frac{\partial F}{\partial y_i} &= \vec{e}_i^T A^T A \vec{y} + \vec{y}^T A^T A \vec{e}_i - \vec{b}^T A \vec{e}_i - \vec{e}_i^T A^T \vec{b} \\ &= 2 \vec{e}_i^T (A^T A \vec{y} - A^T \vec{b})\end{aligned}$$

$$\Rightarrow \nabla F(\vec{y}) = 2 (A^T A \vec{y} - A^T \vec{b})$$

The minimum occurs where the gradient is zero:

Solve $A^T A \vec{y} - A^T \vec{b}$.

One then needs to show that this critical point is indeed a minimum.



Ex Find the least-squares solution of

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \vec{x} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

$$A^T A = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 2 & 3 \end{pmatrix}$$

$$A^T \vec{b} = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \vec{x} = \begin{pmatrix} -\frac{1}{2} \\ 1 \end{pmatrix}$$

$$\text{Solve } \begin{pmatrix} 2 & 2 & | & 1 \\ 2 & 3 & | & 2 \end{pmatrix} \rightarrow \begin{pmatrix} 2 & 2 & | & 1 \\ 0 & 1 & | & 1 \end{pmatrix}$$

Block Matrix Notation

Consider a matrix

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{pmatrix} \quad \left(\begin{array}{c} X_1 \\ X_2 \\ X_3 \\ X_4 \end{array} \right)$$

||

$$\begin{pmatrix} A_{11} & A_{21} \\ A_{12} & A_{22} \end{pmatrix} \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \end{pmatrix} = \begin{pmatrix} A_{11} \vec{x}_1 + A_{21} \vec{x}_2 \\ A_{12} \vec{x}_1 + A_{22} \vec{x}_2 \end{pmatrix}$$

\bigcirc = block of all zeros.

Orthogonality

- Two vectors \vec{u}, \vec{v} are orthogonal if $\vec{u} \cdot \vec{v} = 0$.

Ex $\vec{u} = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, \vec{v} = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}$

- A square matrix Q is orthogonal if

$$Q^T Q = I$$

Ex

Permutation matrices!

From a computational point of view there is a much better way to solve the normal equations.

Suppose $A = \begin{pmatrix} Q \\ \emptyset \end{pmatrix}$

$m \times n$ $m \times m$

$m \geq n$

$n \times n$
 \hat{R}
 \emptyset
 $(m-n) \times n$
zero matrix

when $Q^T Q = I$ and $\hat{R}_{ii} \neq 0$.

$$\begin{aligned} A^T A &= (R^T \emptyset) Q^T Q \begin{pmatrix} \hat{R} \\ \emptyset \end{pmatrix} \\ &= (\hat{R}^T \emptyset) \begin{pmatrix} \hat{R} \\ \emptyset \end{pmatrix} \\ &= \hat{R}^T \hat{R}. \end{aligned}$$

$$\begin{aligned} A^T \vec{b} &= (\hat{R}^T \emptyset) Q^T \vec{b} \\ &= \hat{R}^T (I \emptyset) Q^T \vec{b} \end{aligned}$$

$$A^T A \vec{x} = A^T \vec{b}$$

$$\hat{R}^T R \vec{x} = \hat{R}^T (I \emptyset) Q^T \vec{b}$$

$$\hat{R}^T R \vec{x} = (I \emptyset) Q^T \vec{b}$$

Recall: For the LU factorization, we used elementary matrices 8

$$A = \begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \end{bmatrix}$$

$$E_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ x & x & x \end{bmatrix}, E_2 E_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix}$$

$$E_3 E_2 E_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \end{bmatrix}$$

Then it turned out that because each of E_3, E_2, E_1 is special upper triangular,

$$(E_3 E_2 E_1)^{-1} = L$$

is also special upper triangular.

Now for

$$A = \begin{bmatrix} x & x & x \\ x & x & y \\ x & x & x \\ x & x & x \end{bmatrix} \quad (Q_3 Q_2 Q_1)^{-1}$$

We find Q_1 , $Q_1 Q_1^T = I$:

$$= (Q_3 Q_2 Q_1)^T$$

$$Q_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix}$$

$$Q_2 Q_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \end{bmatrix}$$

$$Q_3 Q_2 Q_1 A = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{bmatrix}$$

The QR factorization There are many ways to compute the $A = QR$ factorization of a matrix. We will use a method that is closer to row reduction:

Theorem Let $\vec{x} \in \mathbb{R}^n$. Define

$$\vec{w} = \vec{x} + \begin{pmatrix} \text{sign}(x_1) \| \vec{x} \|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Then $H = I - 2 \frac{\vec{w} \vec{w}^T}{\| \vec{w} \|_2^2}$ satisfies

- $H^T H = I$

$$\text{Sign}(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases}$$

- $H \vec{x} = \begin{pmatrix} -\text{sign}(x_1) \| \vec{x} \|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

This H is called a Householder reflector.

We also use $H(\vec{x})$ to refer to this matrix

Define

$$H(\alpha) = I \quad \alpha \in \mathbb{R}.$$

6

Important: Consider a vector $\begin{pmatrix} a \\ \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{pmatrix} = \begin{pmatrix} a \\ \vec{x} \end{pmatrix}$

$$\begin{pmatrix} I & 0 \\ Q & H(\vec{x}) \end{pmatrix} \begin{pmatrix} a \\ \vec{x} \end{pmatrix} = \begin{pmatrix} a \\ \| \vec{x} \|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} I & 0 \\ Q & H(\vec{x}) \end{pmatrix} \begin{pmatrix} \vec{a} \\ \vec{x} \end{pmatrix} = \begin{pmatrix} \vec{a} \\ \| \vec{x} \|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Algorithm (QR factorization)

Input: An $m \times n$ matrix A

Output: An $m \times m$ matrix Q , an $m \times n$ matrix R

Step 1: Initialize $Q = I$, $R = A$

Step 2: For $j=1, 2, \dots, n$ do steps 3-5

Step 3:

$$\text{Set } H_0 = H(R(j:end, j))$$

Step 4:

$$\text{Set } R(j:end, j:end) = I - R(j:end, j:end)$$

Step 5:

$$\text{Set } Q(:, j:end) = Q(:, j:end) H_0^T$$

Output (Q, R).

As we will see in the next coding project, the details of coding this make it look a bit different.

How the algorithm actually proceeds:

$$A = R = \begin{pmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{pmatrix} \rightarrow \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \\ 0 & x & x \end{pmatrix} \rightarrow \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & x \end{pmatrix}$$

$\rightarrow \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} \hat{R} \\ 0 \end{pmatrix}$

How to use the QR algorithm to effectively solve the normal equations

$$\begin{pmatrix} A \end{pmatrix} \vec{x} = \begin{pmatrix} b \end{pmatrix}$$

$$\begin{pmatrix} A & | & b \end{pmatrix} = \begin{pmatrix} x & x & x & | & x \\ x & x & x & | & x \\ x & x & x & | & x \\ x & x & x & | & x \end{pmatrix} \rightarrow \begin{pmatrix} x & x & x & | & x \\ 0 & x & x & | & x \\ 0 & x & x & | & x \\ 0 & x & x & | & x \end{pmatrix}$$

$\rightarrow \begin{pmatrix} x & x & x & | & x \\ 0 & x & x & | & x \\ 0 & 0 & x & | & x \\ 0 & 0 & x & | & x \end{pmatrix} \rightarrow \begin{pmatrix} x & x & x & | & x \\ 0 & x & x & | & x \\ 0 & 0 & x & | & x \\ 0 & 0 & 0 & | & x \end{pmatrix}$

$$\hat{R} \quad (\begin{pmatrix} I & 0 \end{pmatrix} Q^T \vec{b})$$

→

$$\left(\begin{array}{ccc|c} X & X & X & X \\ 0 & X & X & X \\ 0 & 0 & X & X \\ 0 & 0 & 0 & X \\ 0 & 0 & 0 & 0 \end{array} \right)$$

The absolute value of this entry gives

$$\|A\vec{x} - \vec{b}\|_2 \text{ where } A^T A \vec{x} = A^T \vec{b}.$$

To summarize: Suppose A is $m \times n$, $m > n$. Form $(A; b)$ reduce to upper triangular form using the Householder algorithm.

Take the first n rows and solve using backward substitution.

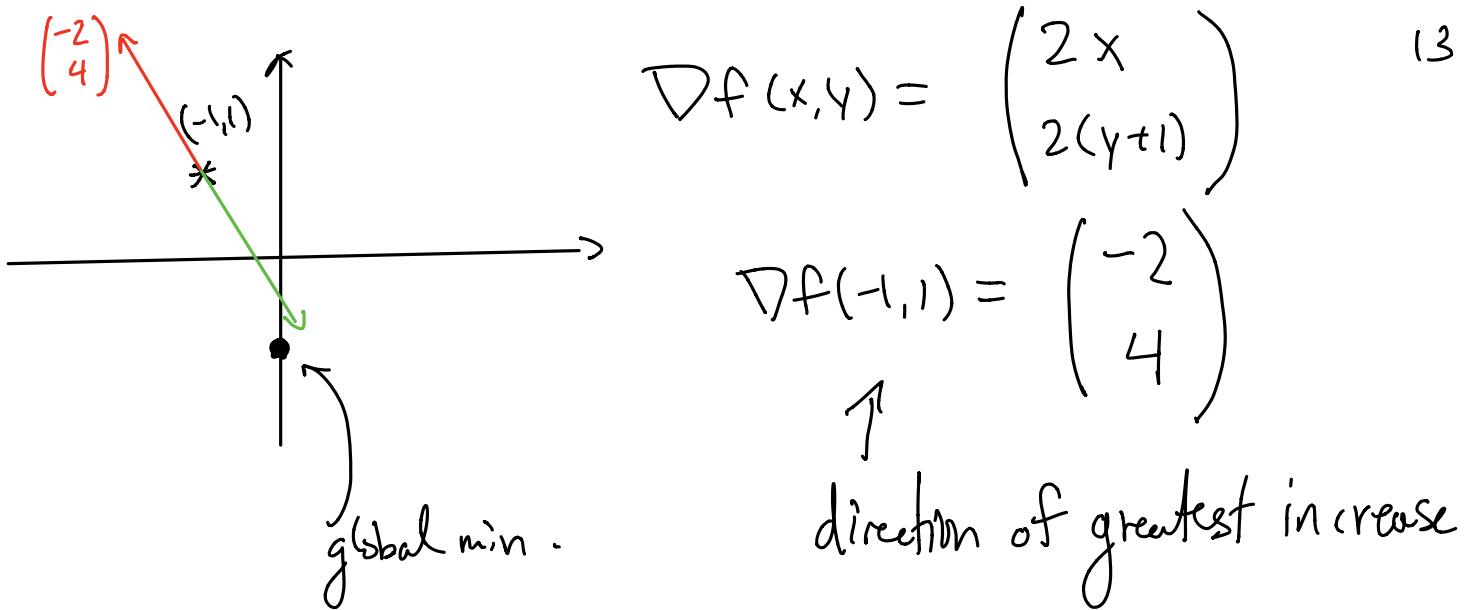
Gradient descent for the normal equations

Consider

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}, \text{ say } f(x, y) = x^2 + (y+1)^2.$$

The global minimum for f occurs at $(x, y) = (0, -1)$.

If we are standing at a point $(-1, 1)$. How do we figure out what direction to move in?



So, you want to take a step in the direction opposite that of the gradient. But how big of a step? That's the hard part...

Our goal: Minimize $F(\vec{y}) = \| A\vec{y} - \vec{b} \|_2^2$

$$\nabla F(\vec{y}) = 2 \left(A^T A \vec{y} - A^T \vec{b} \right)$$

\vec{y}_0 initial guess.

$$\vec{y}_{k+1} = \vec{y}_k - h_k (A^T A \vec{y}_k - A^T \vec{b})$$

$h_k = ?$

Choose: $h_k = \frac{(A^T \vec{b} - A^T A \vec{y}_k)^T (A^T \vec{b} - A^T A \vec{y}_k)}{(A A^T \vec{b} - A A^T A \vec{y}_k)^T (A A^T \vec{b} - A A^T A \vec{y}_k)}$

This looks complicated, so break it down:

define

$$\vec{r}_k = A^T A \vec{y}_k - A^T \vec{b}$$

$$h_k = \frac{\vec{r}_k^T \vec{r}_k}{(\vec{A} \vec{r}_k)^T (\vec{A} \vec{r}_k)}$$

$$\vec{y}_{k+1} = \vec{y}_k - h_k \vec{r}_k$$

You could choose $h_k = h > 0$ small, but this is, in general, much slower. For problems outside of least-squares such as those from machine learning, an optimal choice of h_k is elusive.

Solving $A\vec{x} = \vec{b}$ when there are infinitely many solutions

Suppose: A is $m \times n$ $m < n$
 $\text{im } A = \mathbb{R}^m$
 $\ker A \neq \{\vec{0}\}$
 $m < n$ implies $\ker A \neq \{\vec{0}\}$

By assumption, $A\vec{x} = \vec{b}$ has infinitely many solutions.

Pick one solution \vec{x}_p . Now let $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_k$ be a basis for the kernel of A .

$$A(\vec{x}_p + c_1\vec{x}_1 + c_2\vec{x}_2 + \dots + c_k\vec{x}_k) = \vec{b}$$

for any choice of c_1, c_2, \dots, c_k . For example,

$$\begin{cases} c_1 = 1,000,000,000 \\ c_2 = c_3 = \dots = c_k. \end{cases} \quad \text{works.}$$

It seems unlikely that this is a relevant solution.

One approach: Among all possible solutions of $\vec{A}\vec{x} = \vec{b}$ find the one with the smallest 2-norm.

We won't solve this, but it is
not difficult



Our first solution:

Find the minimum of

$$\|\vec{A}\vec{y} - \vec{b}\|_2^2 + \|\vec{y}\|_2^2$$

over all choices of $\vec{y} \in \mathbb{R}^n$. This is different than the above question.

I claim that this a problem we already know how to solve:

$$\left\| \begin{pmatrix} I \\ A \end{pmatrix} \vec{y} - \begin{pmatrix} 0 \\ \vec{b} \end{pmatrix} \right\|_2^2 = \|\vec{y}\|_2^2 + \|\vec{A}\vec{y} - \vec{b}\|_2^2$$

Replace A with $\begin{pmatrix} I \\ A \end{pmatrix}$ and \vec{b} with $\begin{pmatrix} 0 \\ \vec{b} \end{pmatrix}$ and solve the normal equations!

But we need not solve this:

$$\|A\vec{y} - \vec{b}\|_2^2 + \gamma \|\vec{y}\|_2^2 \quad \gamma \geq 0.$$

Effects of γ :

If we want to prioritize $A\vec{x} \approx \vec{b}$ and we care less about how large $\|\vec{x}\|_2$ is, we choose γ small, say $\gamma \approx 10^{-16}$

If we care less about how close $A\vec{x}$ is to \vec{b} but more about

So, we need to find the minimum of

$$\left\| \begin{bmatrix} \sqrt{\gamma} I \\ A \end{bmatrix} \vec{y} - \begin{bmatrix} \vec{0} \\ \vec{b} \end{bmatrix} \right\|_2.$$

Application Suppose your rock company produces
 20 tons of road base (RB)
 10 tons of sand (S)
 10 tons of gravel (G)

each day. You have 4 categories of jobs and you are wondering if you can increase production by 50%. But you don't want to restructure the company too much..

Equipment : As a simplification, you estimate

$$\text{Truck} : 1 \text{ hour} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \begin{matrix} \text{RB} \\ \text{S} \\ \text{G} \end{matrix}$$

$$\text{Shovel A} : 1 \text{ hour} \rightarrow \begin{pmatrix} 0 \\ 2 \\ 1 \end{pmatrix}$$

$$\text{Mixer} : 1 \text{ hour} \rightarrow \begin{pmatrix} 1 \\ -\frac{1}{2} \\ -\frac{1}{2} \end{pmatrix}$$

$$\text{Shovel B} : 1 \text{ hour} \rightarrow \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

$$A = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 2 & -\frac{1}{2} & 1 \\ 1 & 1 & -\frac{1}{2} & 1 \end{pmatrix}$$

$$\text{Currently, } A \vec{x} = \begin{pmatrix} 20 \\ 10 \\ 10 \end{pmatrix}$$

We want

$$A \vec{x}_{\text{new}} = \begin{pmatrix} 30 \\ 15 \\ 15 \end{pmatrix}$$

$$A(\vec{x}_{\text{new}} - \vec{x}) = \begin{pmatrix} 10 \\ 5 \\ 5 \end{pmatrix} = \vec{b}$$

One approach

$$\text{minimize } \left(\|A\vec{y} - \vec{b}\|_2 + \underbrace{\|\vec{y}\|_2}_\text{keeps change small!} \right)$$