Math for Machine Learning Week 1.2: Subspaces, Bases, and Orthogonality

By: Samuel Deng

Logistics and Announcements

Lesson Overview

Regression. Fill in gaps from last time: invertibility and Pythagorean theorem.

Subspaces. Subsets of $\mathcal{S} \subseteq \mathbb{R}^n$ where we "stay inside" when performing linear combinations of vectors.

Bases. A "language" to describe all vectors in a subspace.

Orthogonality. Orthonormal bases are "good" bases to work with.

becomes much simpler.

- **Projection.** Formal definition of projection and the relationship between projection and least squares.
- Least squares with orthonormal bases. If we have an orthonormal basis for span(col(X)), least squares

Lesson Overview

Big Picture: Least Squares



Lesson Overview

Big Picture: Gradient Descent

 $f(w)=w^2$





Least Squares A Quick Review

Matrices Review from linear algebra

A <u>matrix</u> is a box of numbers, or a list of vectors. We write $\mathbf{X} \in \mathbb{R}^{n \times d}$ as:

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \cdots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix}$$

Column definition: stack column vectors $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$ side-by-side next to each other. **Row definition:** take (by convention, column) vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, turn them into rows $\mathbf{x}_1^{\mathsf{T}}, \dots, \mathbf{x}_n^{\mathsf{T}} \in \mathbb{R}^{1 \times d}$, and stack them on top of each other.

or
$$\mathbf{X} = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow \end{bmatrix}$$

Multiplication Matrix-vector multiplication (column view)

To multiply a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and a vector $\mathbf{w} \in \mathbb{R}^d$, we can think of the column view:

$$\mathbf{X}\mathbf{w} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = w_1 \begin{bmatrix} \uparrow \\ \mathbf{x}_1 \\ \downarrow \end{bmatrix} + \dots + w_d \begin{bmatrix} \uparrow \\ \mathbf{x}_d \\ \downarrow \end{bmatrix}.$$

The result is $\mathbf{X}\mathbf{w} \in \mathbb{R}^n$.

Interpretation: Xw is a *linear combination* of the columns of X.

Multiplication Matrix-vector multiplication (equation view)

To multiply a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and a vector $\mathbf{w} \in \mathbb{R}^d$, we can think of the equation view:



The result is $\mathbf{X}\mathbf{w} \in \mathbb{R}^n$.

Interpretation: Xw compiles the "right-hand sides" of a system of linear equations.

$$\rightarrow \left[\begin{array}{c} \uparrow \\ \mathbf{w} \\ \downarrow \end{array} \right] = \left[\begin{array}{c} \mathbf{x}_1^{\mathsf{T}} \mathbf{w} \\ \vdots \\ \mathbf{x}_n^{\mathsf{T}} \mathbf{w} \end{array} \right]$$

Regression Setup (Example View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow \end{bmatrix} \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

<u>Goal</u>: For each $i \in [n]$, we predict: $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \ldots + w_d x_{id} \in \mathbb{R}$.

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Regression Setup (Feature View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} \mathbf{y} = \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Regression A note on intercepts

<u>**Goal:</u>** For each $i \in [n]$, what if we want to</u>

Solution: We modify add a "dummy" 1 to each example:

 $\mathbf{x}_i^{\mathsf{T}} = \begin{bmatrix} z \end{bmatrix}$

Same as transforming the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ into $\mathbf{X}' \in \mathbb{R}^{n \times (d+1)}$:

$$\mathbf{X} = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow \\ \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow \end{bmatrix}$$

Choose a weight vector that fits $\mathbf{X}': \mathbf{w} \in \mathbb{R}^{d+1}$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

 $\mathbf{X}'\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$. The last (d + 1) entry of \mathbf{w} is the intercept, w_0 .

We can always do this WLOG, so we'll focus on the "homogeneous" case.

predict:
$$\hat{y}_i = \mathbf{w}^{\mathsf{T}} \mathbf{x}_i + \mathbf{w}_0 = w_1 x_{i1} + \dots + w_d x_{id} + \mathbf{w}_0^{\mathsf{T}}$$

$$x_{i1} \ldots x_{id} 1$$
].

$$\implies \mathbf{X}' = \begin{bmatrix} \leftarrow \mathbf{x}_1^\top \rightarrow & 1 \\ \vdots & \vdots \\ \leftarrow \mathbf{x}_n^\top \rightarrow & 1 \end{bmatrix}$$

Least Squares Summary

Use the principle of *least squares* to find the $\hat{\mathbf{w}} \in \mathbb{R}^d$ that minimizes

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Using geometric intuition: $\hat{\mathbf{y}}$ is the vector for which $\hat{\mathbf{y}} - \mathbf{y}$ is perpendicular to span(col(X)).

By Pythagorean Theorem, any other vector $\tilde{y} \in \text{span}(\text{col}(X))$ gives a larger error:

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$$

Because $\hat{y} - y$ is perpendicular to span(col(X)), we obtain the normal equations:

$$\mathbf{X}^{\mathsf{T}}\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, and

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



Least Squares First missing item: invertibility of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$

"If there are no redundant features, then we can invert the normal equations"

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Regression Setup (Feature View)

<u>**Observed:**</u> Matrix of training samples $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector of training labels $\mathbf{y} \in \mathbb{R}^{n}$.

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} \mathbf{y} = \mathbf{y}$$

<u>**Unknown:**</u> Weight vector $\mathbf{w} \in \mathbb{R}^d$ with weights w_1, \ldots, w_d .

Choose a weight vector that "fits the training data": $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \hat{y}_i$ for $i \in [n]$, or:

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ where } \mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n.$$

 $\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}$.

Subspaces

Subspaces Idea

A <u>subspace</u> is a set of vectors that "stays within" the set under all linear combinations of the vectors.

Subspaces Definition

A <u>subspace</u> $\mathcal{S} \subseteq \mathbb{R}^n$ is a subset of vectors that satisfies the property: if $\mathbf{v}, \mathbf{w} \in \mathcal{S}$, then $\alpha \mathbf{v} + \beta \mathbf{w} \in \mathcal{S}$ for any $\alpha, \beta \in \mathbb{R}$.

Any subspace \mathcal{S} contains the zero vector: $\mathbf{0} \in \mathcal{S}$.

Example: $\mathcal{S}_0 := \mathbb{R}^2$

Example: $\mathcal{S}_1 := \{ \mathbf{v} \in \mathbb{R}^2 : v_1 = 0 \}$

Example: $\mathcal{S}_2 := \{ \mathbf{v} \in \mathbb{R}^3 : v_1 = v_2 \}$

Span Review

For a collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_d \in \mathbb{R}^n$, the span is the set of vectors we can attain through linear combinations of $\mathbf{a}_1, \dots, \mathbf{a}_d$:

$$\operatorname{span}(\mathbf{a}_1, \dots, \mathbf{a}_d) = \left\{ \mathbf{y} \in \mathbb{R}^n : \mathbf{y} = \sum_{i=1}^d \alpha_i \mathbf{a}_i, \alpha_i \in \mathbb{R} \right\}.$$

Recall that this is equivalent to all the $\mathbf{y} \in \mathbb{R}^{n \times d}$ we obtain from matrix vector multiplication!

$$y = \mathbf{A}\alpha$$
, i.e. $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

$$= \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ \mathbf{a}_1 & \dots & \mathbf{a}_d \\ \downarrow & \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_d \end{bmatrix}$$

Example: $\mathscr{S}_3 := \operatorname{span} \left(\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right)$

(Non)Example: $S_4 := \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 5 \}$

Subspaces Specific example: span(col(X))

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

 $span(col(\mathbf{X})) = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y} = w_1\mathbf{x}_1 + \dots + w_d\mathbf{x}_d\}$

We will refer to this, later, as CS(X), the <u>columnspace</u> of X.

Bases & Dimension

Basis Idea

For a subspace S, a <u>basis</u> is a *minimal* set of vectors that can "linearly describe" any vector in S. A "language" for vectors in S.

Basis Linear Independence and Span

Recall the following two notions.

A collection of vectors $\mathbf{a}_1, ..., \mathbf{a}_d \in \mathbb{R}^n$ is linear if $\alpha_i = 0$ for all $i \in [d]$.

For a collection of vectors $\mathbf{a}_1, ..., \mathbf{a}_d \in \mathbb{R}^n$, the linear combinations of $\mathbf{a}_1, ..., \mathbf{a}_d$:

 $\operatorname{span}(\mathbf{a}_1, \dots, \mathbf{a}_d) = \begin{cases} \mathbf{y} \\ \mathbf{y} \end{cases}$

A collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_d \in \mathbb{R}^n$ is <u>linearly independent</u> if $\alpha_1 \mathbf{a}_1 + \dots + \alpha_d \mathbf{a}_d = \mathbf{0}$ if and only

For a collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_d \in \mathbb{R}^n$, the <u>span</u> is the set of vectors we can attain through

$$\in \mathbb{R}^n : \mathbf{y} = \sum_{i=1}^d \alpha_i \mathbf{a}_i, \alpha_i \in \mathbb{R} \right\}.$$

Basis Definition

For a subspace $S \subseteq \mathbb{R}^n$, a set of vectors $\mathbf{a}_1, ..., \mathbf{a}_d \in S$ is a <u>basis</u> for S if: $S = \operatorname{span}(\mathbf{a}_1, ..., \mathbf{a}_d)$ and $\mathbf{a}_1, ..., \mathbf{a}_d$ are linearly independent.

Bases are not unique – there are infinitely many bases for any subspace. However, all bases have the same number of elements.

Basis Examples

Example: $\mathscr{S}_0 := \mathbb{R}^2$

Basis Examples

Example: $\mathcal{S}_1 := \{ \mathbf{v} \in \mathbb{R}^2 : v_1 = 0 \}$

Basis Examples

Example: $\mathcal{S}_2 := \{ \mathbf{v} \in \mathbb{R}^3 : v_1 = v_2 \}$

Dimension of a Subspace Definition

The <u>dimension</u> of a subspace is the size of any of its bases. For a subspace S, write this as dim(S).

Matrices & Subspaces Every matrix comes with four subspaces

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix. Its <u>columnspace</u> is $CS(X) = \{y \in \mathbb{R}^n : y = Xw, \text{ for any } w \in \mathbb{R}^d\}$ (this was span(col(X))). Its <u>nullspace/kernel</u> is $NS(X) := \{ w \in \mathbb{R}^d : Xw = 0 \}.$ Its rowspace is $CS(X^{\top}) = \{ y \in \mathbb{R}^d : y = X^{\top}v, \text{ for any } v \in \mathbb{R}^n \}.$ Its left nullspace is $NS(X^{\top}) := \{ \mathbf{v} \in \mathbb{R}^n : X^{\top}\mathbf{v} = \mathbf{0} \}.$

Rank-nullity theorem: $d = \dim(CS(X)) + \dim(NS(X))$.

Matrices & Subspaces Columnspace of a matrix

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. We can think of its columnspace as:

$$CS(\mathbf{X}) := \{ \mathbf{y} \in \mathbb{R}^n : \mathbf{y} = \mathbf{X} \\ = \{ \mathbf{y} \in \mathbb{R}^n : \mathbf{y} = w \} \\ = \operatorname{span}(\mathbf{x}_1, \dots, \mathbf{x}_d) = 0 \}$$

This is a subspace that "comes with" any matrix.

 $\{\mathbf{w}, \text{ for any } \mathbf{w} \in \mathbb{R}^d\}$ $v_1 \mathbf{x}_1 + \ldots + w_d \mathbf{x}_d$, for any $w_i \in \mathbb{R}$ = span(col($\mathbf{x}_1, \ldots, \mathbf{x}_d$))

Matrices & Subspaces Rank of a matrix

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

The $\underline{\mathrm{rank}}$ of $\mathbf X$ is the number of linearly independent columns (which is the same as the number of linearly independent rows).

It is always the case that: rank(\mathbf{X}) $\leq \min\{n, d\}$. If rank(\mathbf{X}) = min $\{n, d\}$, then we say \mathbf{X} is full rank.


Matrices & Subspaces Rank & Invertibility

Let $\mathbf{X} \in \mathbb{R}^{d \times d}$ be a square matrix.

It is always the case that: $rank(\mathbf{X}) \leq d$. If $rank(\mathbf{X}) = d$, then we say **X** is full rank.

Basic fact from linear algebra:

X is invertible if and only if it is full rank.

Matrices & Subspaces Dimension of the columnspace

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

So, if $rank(\mathbf{X}) = d$, then $\mathbf{x}_1, \dots, \mathbf{x}_d$ form a basis for the columnspace!

- $CS(\mathbf{X}) = span(\mathbf{x}_1, \dots, \mathbf{x}_d)$
- $rank(\mathbf{X}) = how many of \mathbf{x}_1, \dots, \mathbf{x}_d$ are linearly independent

"If there are no redundant features, then we can invert the normal equations"

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

and rank(\mathbf{X}) = d, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Proof. To show that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, show $\operatorname{rank}(\mathbf{X}^{\mathsf{T}}\mathbf{X}) = d$.

Theorem (Invertibility of \mathbf{X}^{\mathsf{T}}\mathbf{X}). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$

and rank(\mathbf{X}) = d, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Proof. To show that $\mathbf{X}^{\top}\mathbf{X}$ is invertible, show $\mathbf{X}^{\top}\mathbf{X}$ has d linearly independent columns.

- Theorem (Invertibility of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$

 - $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0} \iff \mathbf{w} = \mathbf{0}.$

and rank(\mathbf{X}) = d, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Proof. To show that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, show $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ has d linearly independent columns.

Suppose $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0}$. Let $\mathbf{w} \in \mathbb{R}^d$ be any vector.

- Theorem (Invertibility of $\mathbf{X}^{\top}\mathbf{X}$). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$

 - $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0} \implies \mathbf{w} = \mathbf{0}.$

and rank(\mathbf{X}) = d, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Proof. To show that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, show $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ has d linearly independent columns.

Suppose $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0}$. Let $\mathbf{w} \in \mathbb{R}^d$ be any vector. Take a dot product of both sides with \mathbf{w} :

 $\mathbf{W}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}$

 $\mathbf{W}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{W}$

- Theorem (Invertibility of $\mathbf{X}^{\mathsf{T}}\mathbf{X}$). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$

$\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0} \implies \mathbf{w} = \mathbf{0}.$

$$\mathbf{w} = \mathbf{w}^{\mathsf{T}} \mathbf{0} = \mathbf{0}.$$

$$\mathbf{v} = \|\mathbf{X}\mathbf{w}\|^2 = 0$$

Least Squares First missing item: invertibility of X¹X

Theorem (Invertibility of $\mathbf{X}^{\top}\mathbf{X}$). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Proof. To show that $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, show $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ has d linearly independent columns. $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} =$

Suppose $\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{0}$. Let $\mathbf{w} \in \mathbb{R}^d$ be any vector. Take a dot product of both sides with \mathbf{w} : $\mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \mathbf{w}^{\mathsf{T}}\mathbf{0} = 0.$

 $\mathbf{w}^{\mathsf{T}}\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} = \|\mathbf{X}\mathbf{w}\|^2 = 0 \implies \mathbf{X}\mathbf{w} = \mathbf{0}.$

$$= 0 \implies w = 0.$$

But $rank(\mathbf{X}) = d$, so \mathbf{X} has d linearly independent columns. Therefore, $\mathbf{w} = \mathbf{0}$.

and rank(\mathbf{X}) = d_i , then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible.

Theorem (Invertibility of $\mathbf{X}^{\top}\mathbf{X}$). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix, with columns $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. If $n \ge d$

Least Squares Summary

Use the principle of *least squares* to find the $\hat{\mathbf{w}} \in \mathbb{R}^d$ that minimizes

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Using geometric intuition: $\hat{\mathbf{y}}$ is the vector for which $\hat{\mathbf{y}} - \mathbf{y}$ is perpendicular to span(col(X)).

By Pythagorean Theorem, any other vector $\tilde{y} \in \text{span}(\text{col}(X))$ gives a larger error:

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$$

Because $\hat{y} - y$ is perpendicular to span(col(X)), we obtain the normal equations:

$$\mathbf{X}^{\mathsf{T}}\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, and

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



Least Squares Summary

Use the principle of *least squares* to find the $\hat{\mathbf{w}} \in \mathbb{R}^d$ that minimizes

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Using geometric intuition: $\hat{\mathbf{y}}$ is the vector for which $\hat{\mathbf{y}} - \mathbf{y}$ is perpendicular to $\mathrm{CS}(\mathbf{X})$.

By Pythagorean Theorem, any other vector $\tilde{y} \in CS(X)$ gives a larger error:

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$$

Because $\hat{y} - y$ is perpendicular to CS(X), we obtain the normal equations:

$$\mathbf{X}^{\mathsf{T}}\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, and

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



By Pythagorean Theorem, any other vector $\tilde{y} \in CS(X)$ gives a larger error:

"The vector closest to y in the subspace is perpendicular."

 $\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$

Orthogonality Definition and Orthonormal Bases

Norms and Inner Products **Euclidean Norm**

Recall the notion of "length" from \mathbb{R}^2 . For a v $\|\mathbf{x}\|_{2}$:

Generalizing this, for $\mathbf{x} \in \mathbb{R}^n$, the <u>Euclidean norm (ℓ_2 -norm)</u> is:

$$\|\mathbf{x}\|_2 := \sqrt{x_1^2 + \ldots + x_n^2} = \sqrt{\mathbf{x}^\top \mathbf{x}}.$$

In this course, dropping the "2" and just writing $\|\mathbf{x}\|$ denotes the Euclidean norm.

vector
$$\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$$
,

$$x = \sqrt{x_1^2 + x_2^2}.$$

 $\|\mathbf{x}\|_2^2 = \mathbf{x}^{\mathsf{T}}\mathbf{x}.$

Orthogonality Definition

Two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ are <u>orthogonal</u> if $\langle \mathbf{v}, \mathbf{w} \rangle$ geometric notion of "perpendicular."

A set of vectors is orthogonal if every pair of distinct vectors in the set is orthogonal.

Two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ are <u>orthogonal</u> if $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^\top \mathbf{w} = 0$. In \mathbb{R}^2 and \mathbb{R}^3 , this corresponds to our

Orthogonality Pythagorean Theorem

<u>Theorem (Pythagorean Theorem).</u> If vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ are orthogonal, then

Proof. Let $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ be orthogonal vectors. Expand the square $\|\mathbf{v} + \mathbf{w}\|^2$. $\|\mathbf{v} + \mathbf{w}\|^2 = \langle \mathbf{v} + \mathbf{w} \rangle$ $= \langle \mathbf{v}, \mathbf{v} \rangle +$ $= \langle \mathbf{v}, \mathbf{v} \rangle +$ $= \|\mathbf{v}\|^2 +$

 $\|\mathbf{v} + \mathbf{w}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2.$

$$\langle \mathbf{v}, \mathbf{w} + \mathbf{w} \rangle$$

$$+ \langle \mathbf{v}, \mathbf{w} \rangle + \langle \mathbf{w}, \mathbf{v} \rangle + \langle \mathbf{w}, \mathbf{w} \rangle$$

$$+ 2 \langle \mathbf{v}, \mathbf{w} \rangle + \langle \mathbf{w}, \mathbf{w} \rangle$$

$$+ \|\mathbf{w}\|^{2}$$

By Pythagorean Theorem, any other vector $\tilde{y} \in CS(X)$ gives a larger error: $\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$



Theorem (Projection minimizes distance). Let $\hat{y} \in CS(X)$ be the vector where $\hat{y} - y$ is orthogonal to any vector in CS(X) and let $\tilde{y} \in CS(X)$ be any other vector. Then



Theorem (Projection minimizes distance). Let $\hat{y} \in CS(X)$ be the vector where $\hat{y} - y$ is orthogonal to any vector in CS(X) and let $\tilde{y} \in CS(X)$ be any other vector. Then $\|\hat{y} - y\|^2 \le \|\tilde{y} - y\|^2$.

Proof. Because $\hat{y}\in CS(X)$ and $\tilde{y}\in CS(X)$ and CS(X) is a subspace, $\tilde{y}-\hat{y}\in CS(X).$

The vector $\hat{y} - y$ is orthogonal to any vector in span(col(X)), so $\hat{y} - y$ is orthogonal to $\tilde{y} - \hat{y}$.

By the Pythagorean Theorem:

 $\|\hat{\mathbf{y}} - \mathbf{y}\|^2 + \|\tilde{\mathbf{y}} - \hat{\mathbf{y}}\|^2 = \|\hat{\mathbf{y}} - \mathbf{y} + \tilde{\mathbf{y}} - \hat{\mathbf{y}}\|^2 = \|\tilde{\mathbf{y}} - \mathbf{y}\|^2$

But because norms are always nonnegative,

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$$





Theorem (Projection minimizes distance). Let $\hat{y} \in CS(X)$ be the vector where $\hat{y} - y$ is orthogonal to any vector in CS(X) and let $\tilde{y} \in CS(X)$ be any other vector. Then



Least Squares Summary

Use the principle of *least squares* to find the $\hat{\mathbf{w}} \in \mathbb{R}^d$ that minimizes

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Using geometric intuition: $\hat{\mathbf{y}}$ is the vector for which $\hat{\mathbf{y}} - \mathbf{y}$ is perpendicular to $\mathrm{CS}(\mathbf{X})$.

By Pythagorean Theorem, any other vector $\tilde{y} \in CS(X)$ gives a larger error:

$$\|\hat{\mathbf{y}} - \mathbf{y}\|^2 \le \|\tilde{\mathbf{y}} - \mathbf{y}\|^2.$$

Because $\hat{\mathbf{y}} - \mathbf{y}$ is perpendicular to $CS(\mathbf{X})$, we obtain the normal equations:

$$\mathbf{X}^{\mathsf{T}}\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^{\mathsf{T}}\mathbf{y}.$$

If $n \ge d$ and $rank(\mathbf{X}) = d$, then $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is invertible, and

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



Least Squares Summary

Goal: Find the $\hat{\mathbf{w}} \in \mathbb{R}^d$ that minimizes $\|\mathbf{X}\mathbf{w}-\mathbf{y}\|^2.$ <u>Theorem (OLS).</u> If $n \ge d$ and $rank(\mathbf{X}) = d$, then: $\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

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



Orthogonality Projections

Projection Idea: A vector's "shadow" on another set

 $\hat{\mathbf{y}}$ in S to \mathbf{y} .





For an arbitrary set $S \subseteq \mathbb{R}^n$, the <u>projection</u> of a vector $\mathbf{y} \in \mathbb{R}^n$ onto the set S is the closest vector

Denote this vector $\Pi_{S}(\mathbf{y}) := \hat{\mathbf{y}}$.

Projection Projection of a vector onto an arbitrary set

 $\hat{\mathbf{y}}$ in S to \mathbf{y} .

 $\hat{\mathbf{y}} \in S$

- For an arbitrary set $S \subseteq \mathbb{R}^n$, the <u>projection</u> of a vector $\mathbf{y} \in \mathbb{R}^n$ onto the set S is the closest vector
 - Denote this vector $\Pi_{S}(\mathbf{y}) := \hat{\mathbf{y}}$.
 - "Closest" in a Euclidean ("least squares") distance sense:
 - $\Pi_{S}(\mathbf{y}) = \arg \min \|\hat{\mathbf{y}} \mathbf{y}\| = \|\hat{\mathbf{y}} \mathbf{y}\|^{2}.$

Projection Projection of a vector onto a subspace

Let $\mathscr{X} \subseteq \mathbb{R}^n$ be a subspace, with the basis $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the matrix with $\mathbf{x}_1, \dots, \mathbf{x}_d$ as its columns. Any point $\hat{\mathbf{y}} \in \mathcal{X}$ is a linear combination:

> $\hat{\mathbf{y}} = w_1 \mathbf{x}$ $= \mathbf{X}\mathbf{w}$

 $\Pi_{\mathcal{X}}(\mathbf{y}) = \arg \min \|\hat{\mathbf{y}} - \mathbf{y}\|^2$ ŷ∈ℒ

$$\mathbf{x}_1 + \ldots + w_d \mathbf{x}_d$$

The projection of **y** onto \mathcal{X} is:

Projection Projection of a vector onto a subspace

Let $\mathscr{X} \subseteq \mathbb{R}^n$ be a subspace, with the basis $\mathbf{x}_1, ..., \mathbf{x}_d \in \mathbb{R}^n$. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the matrix with $\mathbf{x}_1, ..., \mathbf{x}_d$ as its columns. Any point $\hat{\mathbf{y}} \in \mathscr{X}$ is a linear combination:

 $\hat{\mathbf{y}} = w_1 \mathbf{x}$ = $\mathbf{X} \mathbf{w}$

This is equivalent to finding:

 $\hat{\mathbf{w}} = \arg \mathbf{w}$ $\hat{\mathbf{w}} \in \mathbf{w}$

$$\mathbf{x}_1 + \dots + w_d \mathbf{x}_d$$

$$\min_{\mathbb{R}^d} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2$$

Least Squares as Projection Projection Matrix

$$\hat{\mathbf{w}} = \arg \min_{\hat{\mathbf{w}} \in \mathbb{R}^d} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2$$

This is just least squares! By what we've learned...

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

$$\Pi_{\mathcal{X}}(\mathbf{y}) = \hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

Let $P_{\mathbf{X}} := \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}} \in \mathbb{R}^{n \times n}$ be the projection matrix for span(col(\mathbf{X})).



Linearity is the central property in linear algebra. Cooking is typically linear.

Bacon, egg, cheese (on bagel) <u>Bacon, egg, cheese (on roll)</u> Lox sandwich

1 egg	1 egg
1 slice of cheese	1 slice
1 slice bacon	1 slice
1 Kaiser roll	0 Kaise
0 cream cheese	0 crear
0 slices of lox	0 slices
0 bagel	1 bage

0 egg

0 slice of cheese

0 slice bacon

0 Kaiser roll

1 cream cheese

2 slices of lox

of cheese

bacon

er roll

m cheese

s of lox

1 bagel

Linearity is the central property in linear algebra.

A function ("transformation") $T: \mathbb{R}^d \to \mathbb{R}^n$ is <u>linear</u> if T satisfies these two properties for any two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$:

 $T(\mathbf{a} + \mathbf{b}) = T(\mathbf{a}) + T(\mathbf{b})$

 $T(c\mathbf{a}) = cT(\mathbf{a})$ for any $c \in \mathbb{R}$.

Example. Consider the function $T : \mathbb{R}^3 \to \mathbb{R}$, defined by:

 $T(\mathbf{x}) = 2x_1 + 3x_3.$

- Matrices also play by these rules. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{w}, \mathbf{v} \in \mathbb{R}^{d}$ be vectors.
 - $\mathbf{X}(\mathbf{w} + \mathbf{v}) = \mathbf{X}\mathbf{w} + \mathbf{X}\mathbf{v}$
 - $\mathbf{X}(c\mathbf{w}) = c(\mathbf{X}\mathbf{w})$ for any $c \in \mathbb{R}$.

<u>Theorem (Equivalence of linear transformations and matrices).</u> Any linear transformation $T : \mathbb{R}^d \to \mathbb{R}^n$ has a corresponding matrix $\mathbf{A}_T \in \mathbb{R}^{n \times d}$ such that:

Any matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ has a corresponding linear transformation $T_{\mathbf{A}} : \mathbb{R}^d \to \mathbb{R}^n$ such that:

 $T_{\mathbf{A}}($

$T(\mathbf{x}) = \mathbf{A}_T \mathbf{x}.$

$$(\mathbf{x}) = \mathbf{A}\mathbf{x}$$

This means that matrix-vector multiplication is the same as applying a linear transformation.

So one way of thinking of a matrix is an "action" applied to vectors.

$T(\mathbf{x}) = \mathbf{A}_T \mathbf{x}$ and $T_{\mathbf{A}}(\mathbf{x}) = \mathbf{A} \mathbf{x}$

Least Squares as Projection **Projection Matrix**

making up the matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$,



Encodes an *action* on vectors!

is the projection matrix onto \mathscr{X} .

To project a vector $\mathbf{y} \in \mathbb{R}^n$ onto \mathcal{X} , compute:

Let $\mathcal{X} \subseteq \mathbb{R}^d$ be a subspace with basis $\mathbf{x}_1, ..., \mathbf{x}_d \in \mathbb{R}^n$. If $\mathbf{x}_1, ..., \mathbf{x}_d$ are linearly independent,

 $\Pi_{\mathscr{X}}(\mathbf{y}) = \hat{\mathbf{y}} = P_{\mathbf{x}}\mathbf{y} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}.$

Least Squares Orthonormal Bases and Projection
Norms and Inner Products Unit Vectors

A vector $\mathbf{v} \in \mathbb{R}^d$ is a <u>unit vector</u> if $\|\mathbf{v}\| = 1$.

We can convert any vector into a unit vector by dividing itself by its norm:

V ||**V**||

Orthonormal Basis "Good" Bases

How should we represent a subspace?

Take, for example, the subspace $\mathcal{S} = \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 0 \}.$

Orthonormal Basis "Good" Bases

Attempt 1: Use the span of a set of vectors: span $\begin{pmatrix} 2 & 0 & 2 \\ 1 & 0 & 3 \\ 0 & 0 & 0 \end{pmatrix}$.

Attempt 2: Use the span of a set of linearly independent vectors (a basis):

Attempt 3: Use the span of an orthonormal set of vectors (an orthonormal basis):

$$\mathcal{S} = \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 0 \}$$

- $\operatorname{span}\left(\begin{array}{cc|c}2 & 0\\1 & 1\end{array}\right).$
- $\left|\begin{array}{c|c}0\\0\end{array}, 1\\0\end{array}\right).$ span

Orthonormal Basis "Good" Bases

$\mathcal{S} = \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 0 \}$ $\operatorname{span}\left(\begin{bmatrix}2\\1\\0\end{bmatrix},\begin{bmatrix}0\\1\\0\end{bmatrix},\begin{bmatrix}2\\3\\0\end{bmatrix}\right) \qquad \operatorname{span}\left(\begin{bmatrix}2\\1\\0\end{bmatrix},\begin{bmatrix}0\\1\\0\end{bmatrix}\right) \qquad \operatorname{span}\left(\begin{bmatrix}1\\0\\0\end{bmatrix},\begin{bmatrix}0\\1\\0\end{bmatrix}\right)$

Orthonormal Basis Definition

A set of vectors $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathcal{S}$ is an <u>orthonormal basis</u> for the subspace \mathcal{S} if they are a basis for \mathcal{S} and, additionally:

 $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$ for $i \neq j$.

 $\|\mathbf{u}_{i}\| = 1$ for $i \in [n]$.

Orthonormal Basis Orthogonal Matrices

A square matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$ is an <u>orthogonal matrix</u> if its columns $\mathbf{u}_1, \dots, \mathbf{u}_d \in \mathbb{R}^d$ are orthogonal unit vectors:

These form an orthonormal basis for span(col(U)).

Its rows are also orthogonal.

 $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$ for $i \neq j$.

 $\|\mathbf{u}_{i}\| = 1$ for $i \in [d]$.

Orthonormal Basis Orthogonal Matrices

A matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$ is an <u>semi-orthogonal matrix</u> if its columns $\mathbf{u}_1, \dots, \mathbf{u}_d \in \mathbb{R}^n$ are orthogonal unit vectors:

These form an orthonormal basis for span(col(U)).

 $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$ for $i \neq j$.

 $\|\mathbf{u}_{i}\| = 1$ for $i \in [d]$.

Orthonormal Basis Properties of Orthogonal Matrices

Let a square matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$ be an <u>orthogonal matrix</u>. Then:

U is its own inverse: $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{U}\mathbf{U}^{\mathsf{T}} = \mathbf{I}$.

U is length-preserving: ||Uv|| = ||v||.

Orthonormal Basis Properties of Orthogonal Matrices

Let matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$ be an <u>semi-orthogonal matrix</u>. Then:

U is its own left inverse: $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{I}$.

U is length-preserving: ||Uv|| = ||v||.

subspace $\mathcal{S} = \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 0 \}$ and the vector



A basis is just a "language" for representing vectors in a subspace. For example, consider the

 $\hat{\mathbf{y}} = \begin{bmatrix} \mathbf{1} \\ 1 \\ 0 \end{bmatrix}$

subspace $\mathcal{S} = \{ \mathbf{v} \in \mathbb{R}^3 : v_3 = 0 \}$ and the vector



A basis is just a "language" for representing vectors in a subspace. For example, consider the

 $\hat{\mathbf{y}} = \begin{bmatrix} \mathbf{1} \\ 1 \\ 0 \end{bmatrix}$

Every subspace $\mathscr{X} \subseteq \mathbb{R}^n$ has many choices of bases.



Some are better than others.

Let $\mathscr{X} \subseteq \mathbb{R}^n$ be a subspace, with $\dim(\mathscr{X}) = d$. One basis: $\mathbf{x}_1, ..., \mathbf{x}_d \in \mathbb{R}^n$, with matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$. Another basis: $\mathbf{u}_1, ..., \mathbf{u}_d \in \mathbb{R}^n$, with matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$. Then,

$$\mathscr{X} = \mathrm{CS}(\mathbf{U}) = \mathrm{CS}(\mathbf{X}).$$





Let $\mathscr{X} \subseteq \mathbb{R}^n$ be a subspace, with $\dim(\mathscr{X}) = d$.

 $\mathscr{X} = \mathrm{CS}(\mathbf{U}) = \mathrm{CS}(\mathbf{X}).$

Therefore, for any $\hat{\mathbf{y}} \in \mathcal{X}$, we can write:

 $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{U}\hat{\mathbf{w}}_{onb}$

Both $\hat{\mathbf{w}}, \hat{\mathbf{w}}_{onb} \in \mathbb{R}^d$ are valid ways to "represent" $\hat{\mathbf{y}}$.





How do we find $\hat{\mathbf{w}}_{onb} \in \mathbb{R}^d$ in $\hat{\mathbf{y}} = \mathbf{U}\hat{\mathbf{w}}_{onb}$?

Least squares!

$$\hat{\mathbf{w}}_{onb} = \underset{\hat{\mathbf{w}}_{onb} \in \mathbb{R}^d}{\arg \min} \|\mathbf{y} - \mathbf{U}\hat{\mathbf{w}}_{onb}\|^2$$

The columns of ${f U}$ give an ONB for ${\mathcal X}...$

$$\hat{\mathbf{w}}_{onb} = (\mathbf{U}^{\mathsf{T}}\mathbf{U})^{-1}\mathbf{U}^{\mathsf{T}}\mathbf{y}$$
$$= \mathbf{U}^{\mathsf{T}}\mathbf{y}$$





Orthonormal Basis Why do we like an orthogonal basis?

- Let \mathscr{X} be a subspace. Let $\Pi_{\mathscr{X}}(\mathbf{y}) = \arg \min \|\|\hat{\mathbf{y}} \mathbf{y}\|^2$ be the projection of \mathbf{y} onto \mathscr{X} . ŷ∈𝒴
- For an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $\mathbf{CS}(\mathbf{X}) = \mathcal{X}$,

$$\hat{\mathbf{w}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

For a semi-orthogonal matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$ with $CS(\mathbf{U}) = \mathcal{X}$,

$$\hat{\mathbf{w}}_{onb} = \mathbf{U}^{\mathsf{T}}$$

Much simpler – no inverse operations!

y and $\hat{\mathbf{y}} = \mathbf{U}\mathbf{U}^{\mathsf{T}}\mathbf{y}$.

y and $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$.

Orthonormal Basis Why do we like an orthogonal basis?

<u>Theorem (Projection with orthogonal matrices).</u> Let $\mathscr{X} \subseteq \mathbb{R}^n$ be a subspace and let $\mathbf{y} \in \mathbb{R}^n$, the <u>projection</u> of \mathbf{y} onto \mathcal{X} , i.e.

 $\Pi_{\mathcal{Y}}(\mathbf{y}) = \mathbf{a}$

is given by

$\mathbf{u}_1, \dots, \mathbf{u}_d \in \mathbb{R}^n$ be an orthonormal basis for \mathcal{X} , with semi-orthogonal matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$. For any

$$\underset{\hat{\mathbf{y}} \in \mathcal{X}}{\operatorname{rg min}} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$$

 $\Pi_{\mathscr{X}}(\mathbf{y}) = \mathbf{U}\mathbf{U}^{\mathsf{T}}\mathbf{y}.$

Recap

Lesson Overview

Regression. Fill in gaps from last time: invertibility and Pythagorean theorem.

Subspaces. Subsets of $\mathcal{S} \subseteq \mathbb{R}^n$ where we "stay inside" when performing linear combinations of vectors.

Bases. A "language" to describe all vectors in a subspace.

Orthogonality. Orthonormal bases are "good" bases to work with.

becomes much simpler.

- **Projection.** Formal definition of projection and the relationship between projection and least squares.
- Least squares with orthonormal bases. If we have an orthonormal basis for CS(X), least squares

Lesson Overview

Big Picture: Least Squares



Lesson Overview

Big Picture: Gradient Descent

 $f(w)=w^2$



