

# Orthogonal Matching Pursuit

## A Machine Learning Perspective

Slides: Team RAAAK

A close-up photograph of a field of flowers. In the foreground, several white tulips are in full bloom, their petals slightly curved. A single, vibrant red tulip stands out prominently in the center. Behind the tulips, a dense field of small, white flowers with yellow centers, resembling daisies or baby's breath, stretches across the background.

Outliers

# Outliers in Linear Models

Samples that deviate strongly from the general pattern of our data

# Outliers in Linear Models

Samples that deviate strongly from the general pattern of our data

Not well explained by our linear model

# Outliers in Linear Models

Samples that deviate strongly from the general pattern of our data

Not well explained by our linear model

Sources of outliers include:

Errors in measurement or data entry

# Outliers in Linear Models

Samples that deviate strongly from the general pattern of our data

Not well explained by our linear model

Sources of outliers include:

Errors in measurement or data entry

Sampling errors

# Outliers in Linear Models

Samples that deviate strongly from the general pattern of our data

Not well explained by our linear model

Sources of outliers include:

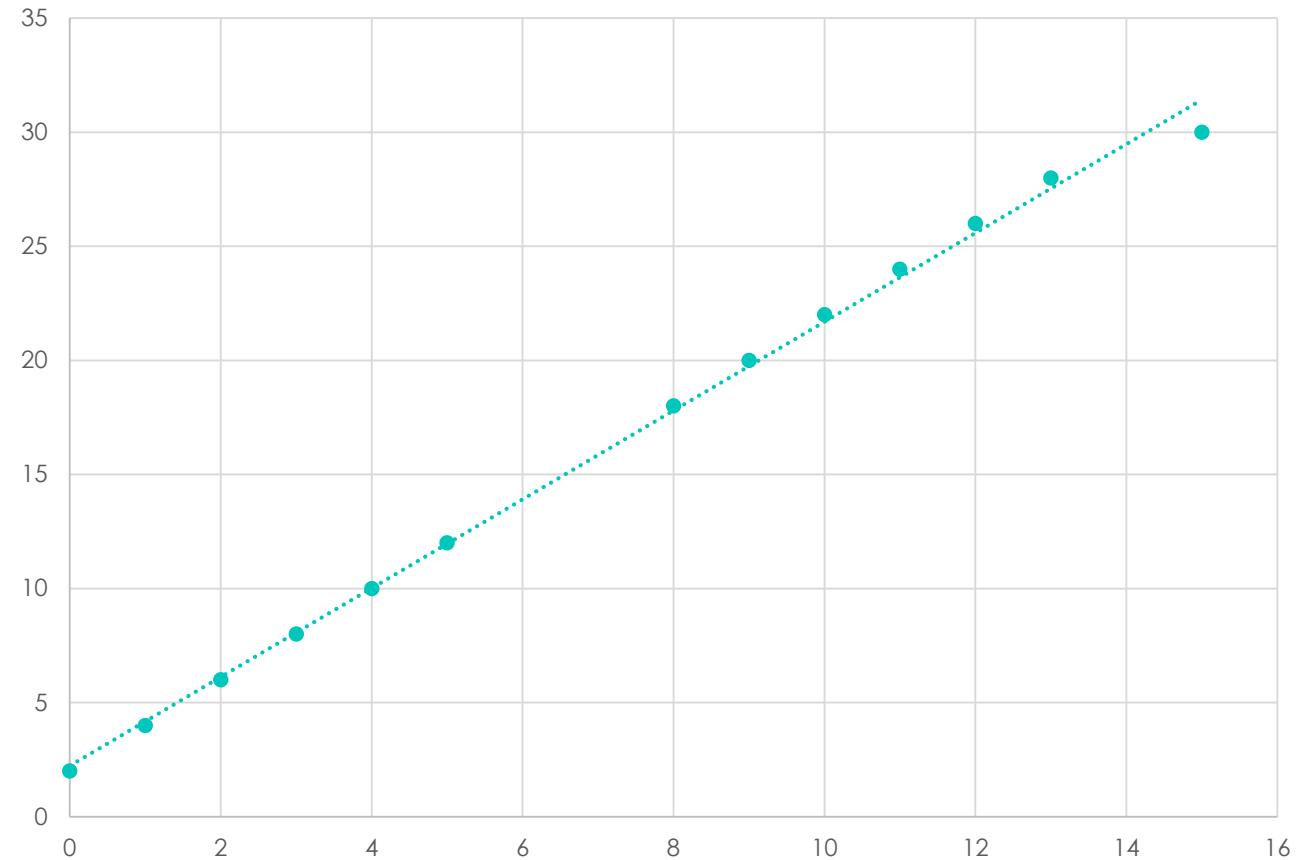
Errors in measurement or data entry

Sampling errors

Noise in our data that corrupts our true observation

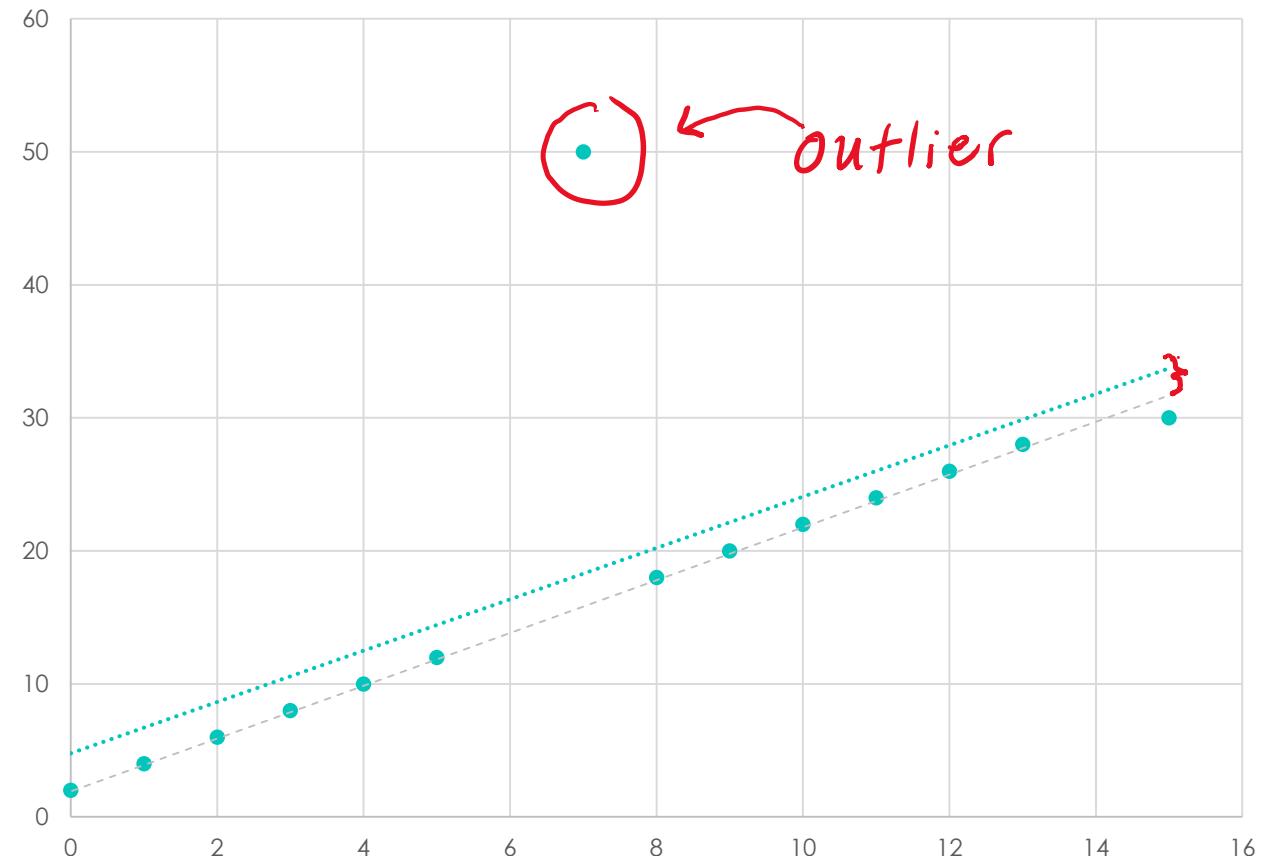
# The Impact of Outliers

- Outliers can drive up the error of our model!
- Cause us to not represent the true pattern of our data



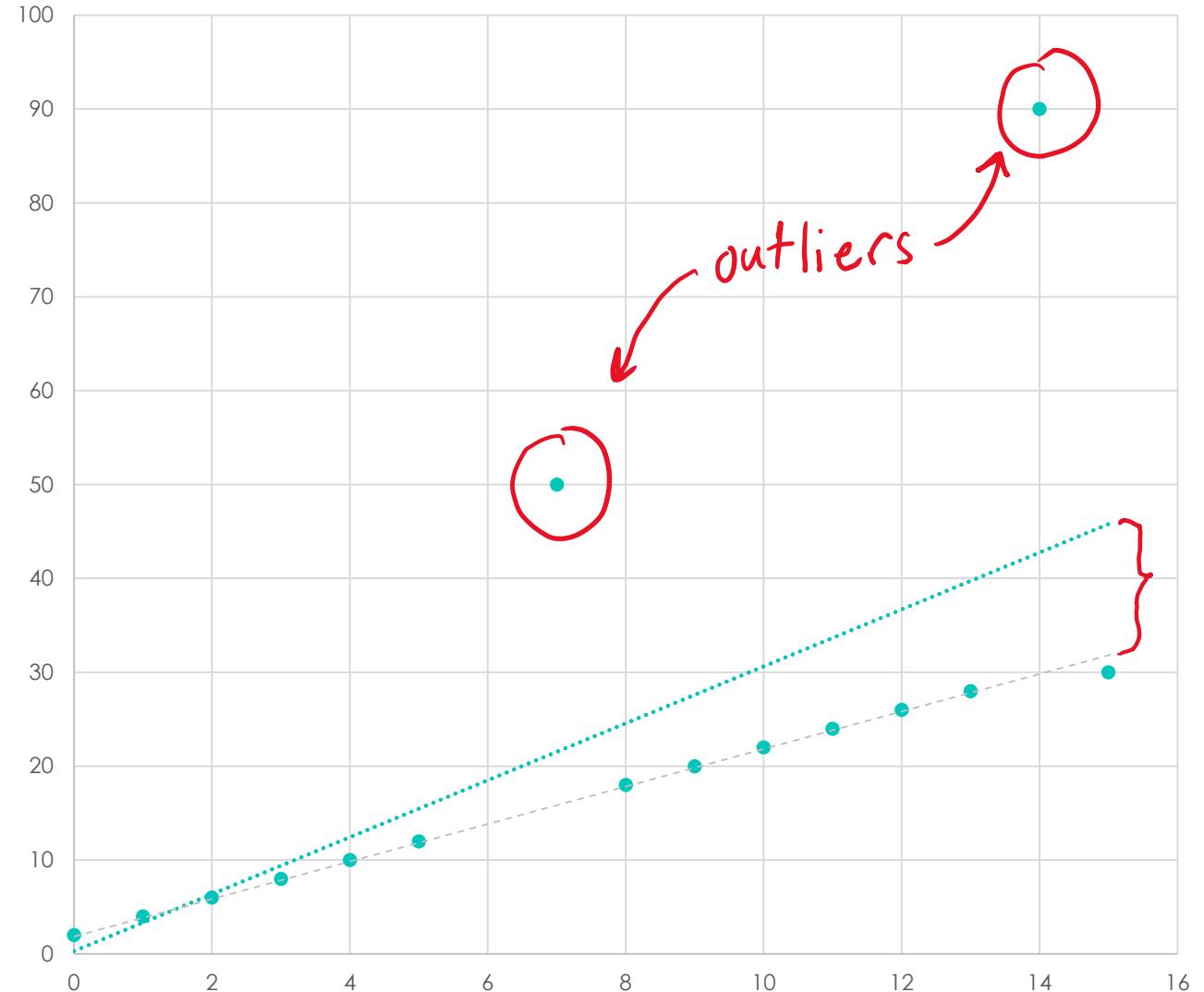
# The Impact of Outliers

- Outliers can drive up the error of our model!
- Cause us to not represent the true pattern of our data



# The Impact of Outliers

- Outliers can drive up the error of our model!
- Cause us to not represent the true pattern of our data



# Facts About OMP – A 16A Review

01

**Goal:** Solve a LS problem to get a sparse solution

02

**Sparse:** containing mostly zero entries

03

**K-sparse vector:** k controls the number of nonzero entries

04

**Assumption:** the columns of X (data matrix) are normalized

05

**Property:** OMP is a Greedy Algorithm

# OMP Review – Signal Processing Context

- You have seen this interpretation in EECS16A
- Context:
  - $m$  satellites that are potentially broadcasting unique signals
  - A received signal that is a linear combination of a subset of these signals



[This Photo](#) by Unknown Author is licensed under CC BY-SA-NC

# OMP Review – Signal Processing Context

- You have seen this interpretation in EECS16A
- What we're given:
  - $m$  satellites that are potentially broadcasting unique signals
  - A received signal that is a linear combination of a subset of these signals
- Goal
  - find which  $k$  satellites are transmitting signals that are present in our received signal

$$\mathbf{y} = \sum_{i \in K} \alpha_i \mathbf{s}_i^{(\tau_i)}$$

# OMP Review - Algorithm

---

**Algorithm 1:** OMP Algorithm from 16A

---

**Input:** Codes  $s_i$ , Received signal  $y$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** The indices of the used codes  $F$ , the sent messages  $x$

$e = y$ ,  $j = 1$ ,  $A = []$ ,  $F = \emptyset$ ;

**while**  $j \leq k$  and  $\|e\| \geq \epsilon$  **do**

$i, t = \text{FindMaxCrossCorrelation}();$

$F = F \cup \{i\};$

$A = [A | s_i^{(t)}];$

$x = (A^T A)^{-1} A^T y;$

$e = y - Ax;$

$j = j + 1;$

**end**

**return**  $F, x$

---

# OMP Review - Algorithm

---

**Algorithm 1:** OMP Algorithm from 16A

---

**Input:** Codes  $s_i$ , Received signal  $y$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** The indices of the used codes  $F$ , the sent messages  $x$

```
e = y, j = 1, A = [], F = ∅;  
while j ≤ k and ||e|| ≥ ε do  
    i, t = FindMaxCrossCorrelation();  
    F = F ∪ {i};  
    A = [A | si(t)];  
    x = (ATA)-1ATy;  
    e = y - Ax;  
    j = j + 1;  
end  
return F, x
```

---

**Recall:** We find the code and time shift with the maximum cross-correlation to predict one satellite that contributes to our received signal

# The Problem of Outlier Detection

We would like to **identify** and **remove** outlier data points to decrease model error



# Outlier Detection: An Analogous Context

## ○ Codes

○ In our new problem, codes correspond to our matrix  $X$  of data

# Outlier Detection: An Analogous Context

- Codes
  - In our new problem, codes correspond to our matrix  $X$  of data
- Message  $\alpha$ 
  - Rather than recover a message broadcasted by satellites, we aim to learn an appropriate set of weights  $w$

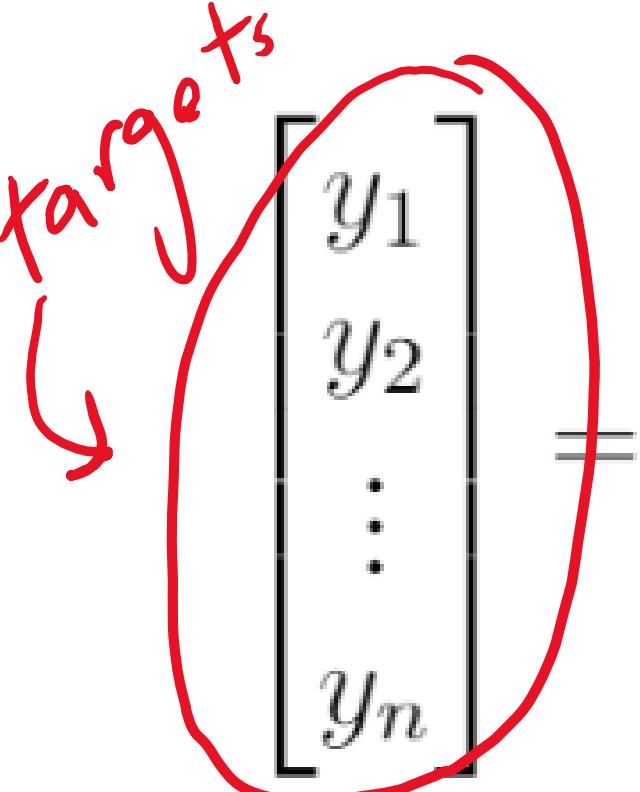
# Outlier Detection: An Analogous Context

- Codes
  - In our new problem, codes correspond to our matrix  $\mathbf{X}$  of data
- Message  $\alpha$ 
  - Rather than recover a message broadcasted by satellites, we aim to learn an appropriate set of weights  $\mathbf{w}$
- Time shifts  $\tau_i$ 
  - In Outlier Detection, we no longer have the problem of time delays that we saw in the signal processing context
  - This means we no longer have to compute any cross-correlation values

# Outlier Detection: An Analogous Context

- Codes
  - In our new problem, codes correspond to our matrix  $\mathbf{X}$  of data
- Message  $\alpha$ 
  - Rather than recover a message broadcasted by satellites, we aim to learn an appropriate set of weights  $\mathbf{w}$
- Time shifts  $\tau_i$ 
  - In Outlier Detection, we no longer have the problem of time delays that we saw in the signal processing context
  - This means we no longer have to compute any cross-correlation values
- Meaning of k
  - Previously, we thought of k as related to the number of satellites transmitting unique codes
  - Now, k refers to the **number of outliers we aim to detect**

# OMP for Outlier Detection – Initial Setup


$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} w + \epsilon.$$

# OMP for Outlier Detection – Initial Setup

Recall:  
we are  
given  
data  
in  
 $(\vec{x}_i, y_i)$   
pairs

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix} w + \epsilon.$$

A red circle highlights the matrix  $\begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix}$ . A red arrow points from this circle to the text "data matrix". A red X is drawn over the term  $w + \epsilon.$ .

# OMP for Outlier Detection – Initial Setup

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} {x_1}^T \\ {x_2}^T \\ \vdots \\ {x_n}^T \end{bmatrix} w + \epsilon.$$

weights  
we aim  
to  
learn  
through  
least  
squares

# OMP for Outlier Detection – Initial Setup

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} w + \epsilon.$$

*noise.*

$\epsilon$  is labeled as noise.

$\mathbf{x}_i^T$  is labeled as  $\mathbf{x}_i$ .

# OMP for Outlier Detection – Error & Residuals

- Consider an error term  $\epsilon_i$  added to each predicted output

$$y_i = \mathbf{x}_i^T \mathbf{w} + \epsilon_i$$

# OMP for Outlier Detection – Error & Residuals

- Consider an error term  $\epsilon_i$  added to each predicted output

$$y_i = \mathbf{x}_i^T \mathbf{w} + \epsilon_i$$

- Constructing the error vector  $\epsilon$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \sum \epsilon_i \mathbf{e}_i$$

# OMP for Outlier Detection – Error & Residuals

- Consider an error term  $\epsilon_i$  added to each predicted output

$$y_i = \mathbf{x}_i^T \mathbf{w} + \epsilon_i$$

- Constructing the error vector  $\epsilon$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \sum \epsilon_i \mathbf{e}_i$$

- Matrix-vector product form

$$[\mathbf{X} \ \mathbf{I}] \begin{bmatrix} \mathbf{w} \\ \epsilon \end{bmatrix}$$

Recall:

$$\hat{\epsilon}_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{ith index}$$

# OMP for Outlier Detection – Error & Residuals

$$[X \ I] \begin{bmatrix} w \\ \epsilon \end{bmatrix}$$

Note: We assume there are only a few outliers, so  $\epsilon$  contains only a few large values

# OMP for Outlier Detection – Inputs

- Data
  - Given in the form of a data matrix  $\mathbf{X}$
  - And targets  $\mathbf{y}$
- Desired sparsity  $k$  of our solution
- Potentially also a desired threshold  
(more on this later)

# OMP for Outlier Detection – Outputs

- A set of indices in our data that correspond to outlier data points
- We'll call this  $F$

# The Algorithm: Initialization

---

Residual  $r$

---

Matrix  $A$

---

Outlier Set  $F$

# OMP: Initial Values

- **Residual:** difference of target and predicted value (using OLS)
  - $r = \mathbf{y} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- **A:** matrix that is updated after each iteration based on the standard basis vector we have considered
  - $A = X$
- **Outlier Set:** We haven't found any outlier data points yet, so this will be initialized to the empty set
- We also set our index  $j$  to be 1 initially, and add 1 at each iteration

# OMP: The Algorithm So Far

---

**Algorithm 2:** OMP to Detect Outliers

---

**Input:** Data matrix  $X$ , Predictions  $y$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** Set of data indices  $F$  where the outliers are

$r = y - X(X^T X)^{-1}X^T y$ ,  $j = 1$ ,  $A = X$ ,  $F = \emptyset$ ;

**while**                           **do**

|

j = j + 1;

**end**

**return**  $F$

---

# Identifying an index

- We'd like to pinpoint the index  $i$  with the largest error value (captured by our residual)
- How would we do this?
- Check for understanding: what should replace the red question mark below?

$$i = \operatorname{argmax}_i |\langle \mathbf{r}, ? \rangle|$$

# Identifying an index

- We'd like to pinpoint the index  $i$  with the largest error value (captured by our residual)
- How would we do this?
- **We take the inner product of the residual with the  $i$ 'th standard basis vector**

$$i = \operatorname{argmax}_i |\langle \mathbf{r}, \mathbf{e}_i \rangle|$$

$$i = \operatorname{argmax}_i |\langle r, e_i \rangle| = \operatorname{argmax}_i |r[i]|$$

This simplifies to the i'th index of the residual,  
which we add to our set of outlier indices

# OMP: The Algorithm So Far

---

**Algorithm 2:** OMP to Detect Outliers

---

**Input:** Data matrix  $X$ , Predictions  $y$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** Set of data indices  $F$  where the outliers are

$r = y - X(X^T X)^{-1}X^T y$ ,  $j = 1$ ,  $A = X$ ,  $F = \emptyset$ ;

**while**  $j \leq k$  and  $\|r\| \geq \epsilon$  **do**

|  $i = \operatorname{argmax}_i |r[i]|$ ;

|  $F = F \cup \{i\}$ ;

|  $j = j + 1$ ;

**end**

**return**  $F$

---

# Augmenting A

- We've pinpointed the  $i$ 'th standard basis vector to have the maximum inner product with our residual
- So, we augment  $A$  by appending the  $i$ 'th standard basis vector as follows:

$$A = [A | e_i];$$

# Recomputing the Residual

- The next step to generate a new value of the residual, taking into account the fact that we've identified one new outlier data point
- To do this, we first perform OLS to generate a new prediction with our augmented version of  $A$ 
  - Define  $\hat{x} = (A^T A)^{-1} A^T y$
  - Our updated prediction of the targets is  $A\hat{x}$

# Recomputing the Residual

- The next step to generate a new value of the residual, taking into account the fact that we've identified one new outlier data point
- To do this, we first perform OLS to generate a new prediction with our augmented version of  $A$

$$\text{Define } \hat{x} = (A^T A)^{-1} A^T y$$

- Our updated prediction of the targets is  $A\hat{x}$
- We then redefine our residual exactly as we did in the initial step, as the difference between the target vector and our prediction

$$r = y - A\hat{x}$$

# OMP: The Algorithm So Far

---

**Algorithm 2:** OMP to Detect Outliers

---

**Input:** Data matrix  $X$ , Predictions  $\mathbf{y}$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** Set of data indices  $F$  where the outliers are

$$\mathbf{r} = \mathbf{y} - X(X^T X)^{-1} X^T \mathbf{y}, j = 1, A = X, F = \emptyset;$$

**while**                           **do**

$$i = \operatorname{argmax}_i |\mathbf{r}[i]|;$$

$$F = F \cup \{i\};$$

$$A = [A | \mathbf{e}_i];$$

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{y};$$

$$\mathbf{r} = \mathbf{y} - Ax;$$

$$j = j + 1;$$

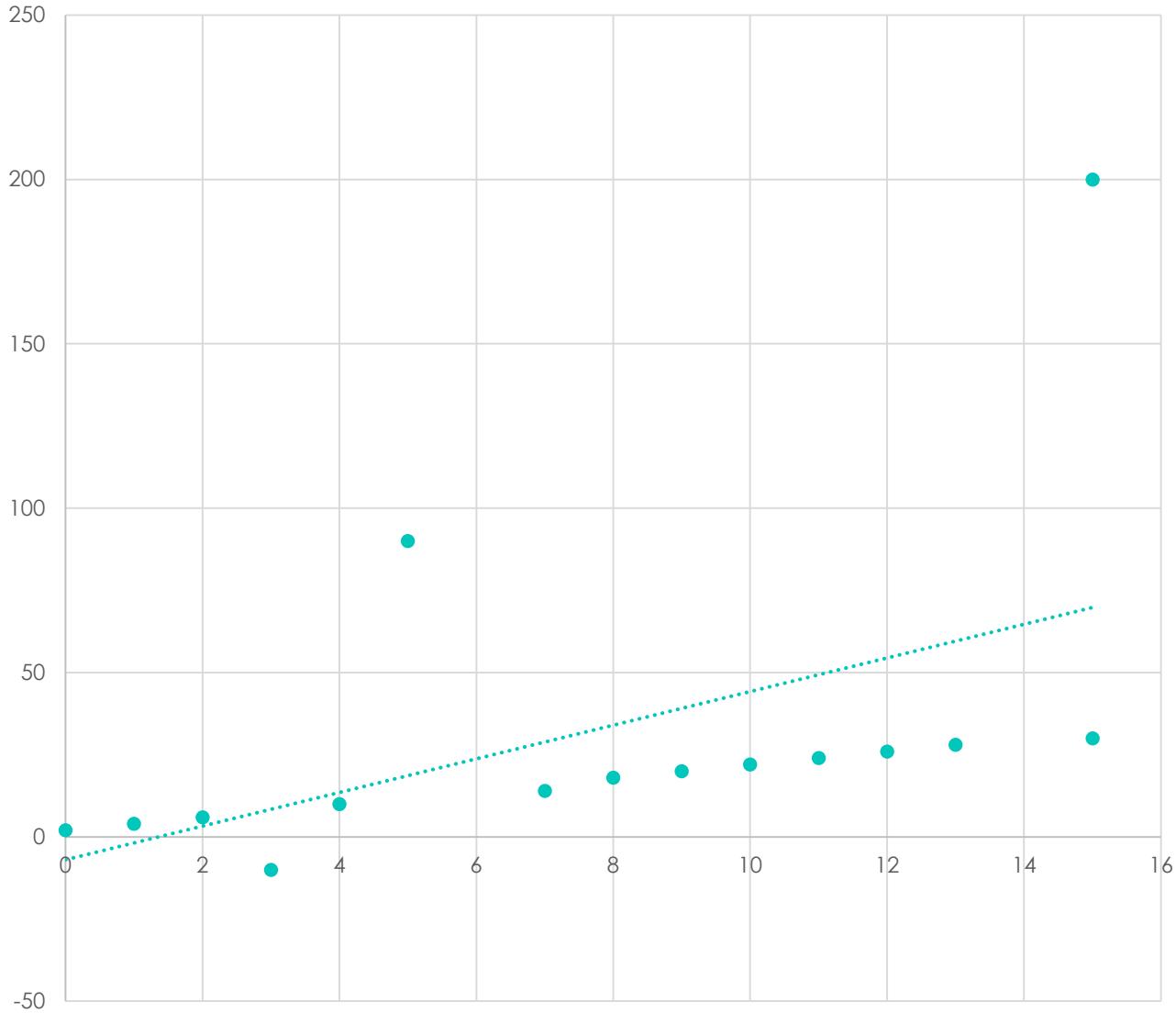
**end**

**return**  $F$

---

# OMP for Outlier Detection – Big Picture

- Each iteration pinpoints one outlier to be removed
- Goal: finding the vectors in a given basis (e.g. the standard basis consisting of vectors  $e_i$ ) that can best capture the pattern in our data, while excluding those directions that increase error in our predictions

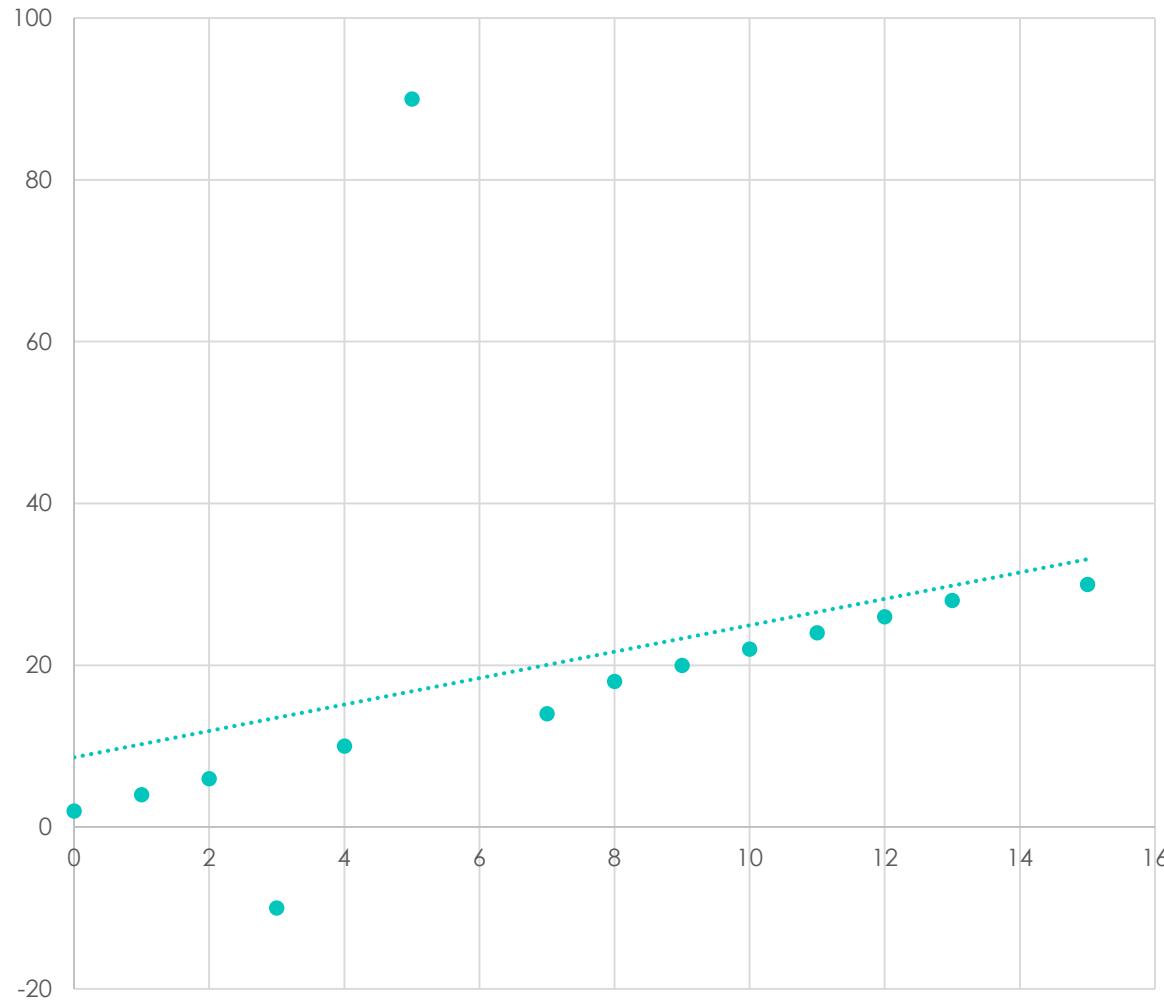


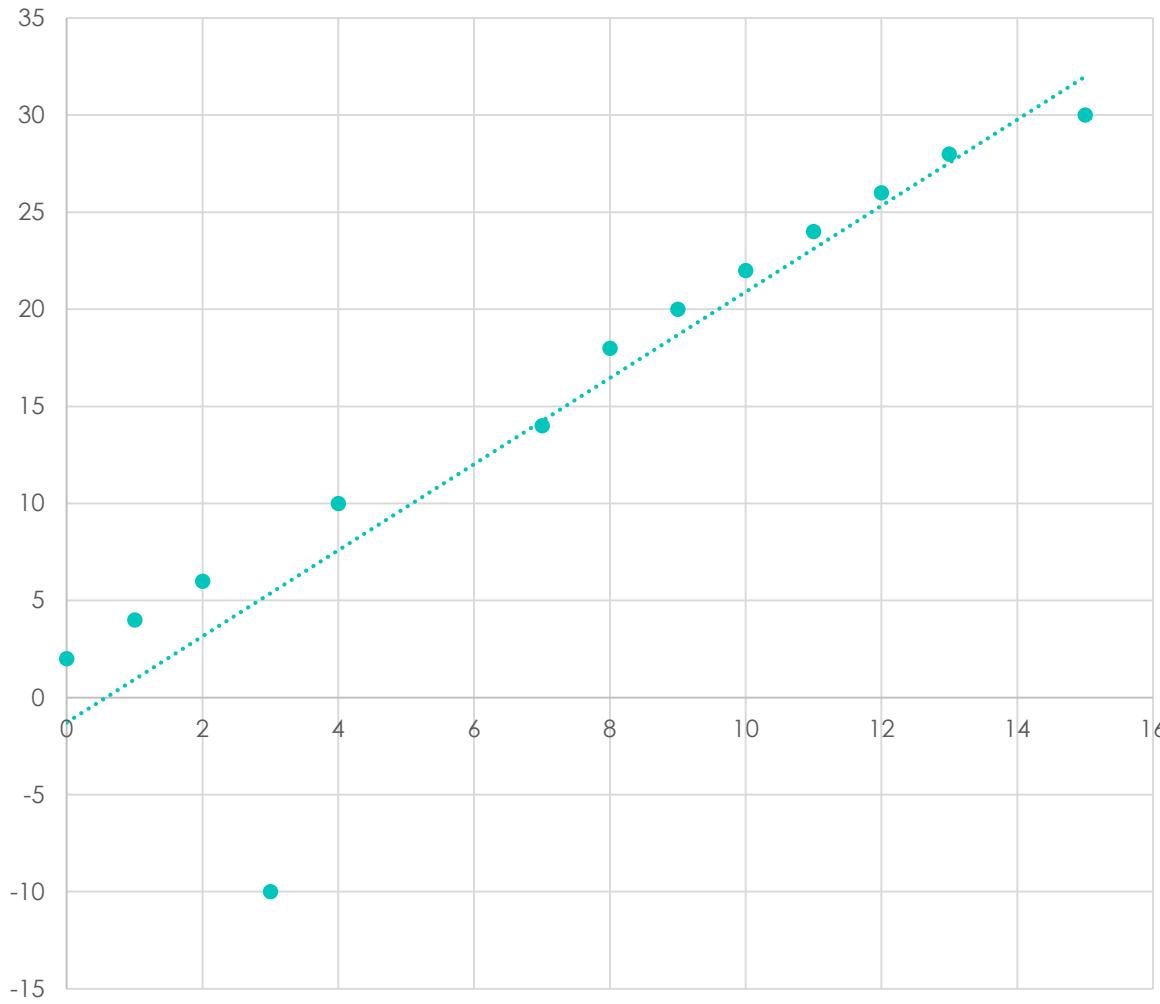
# OMP for Outlier Detection and Removal - Visualization

# OMP for Outlier Detection and Removal

- 

## Visualization

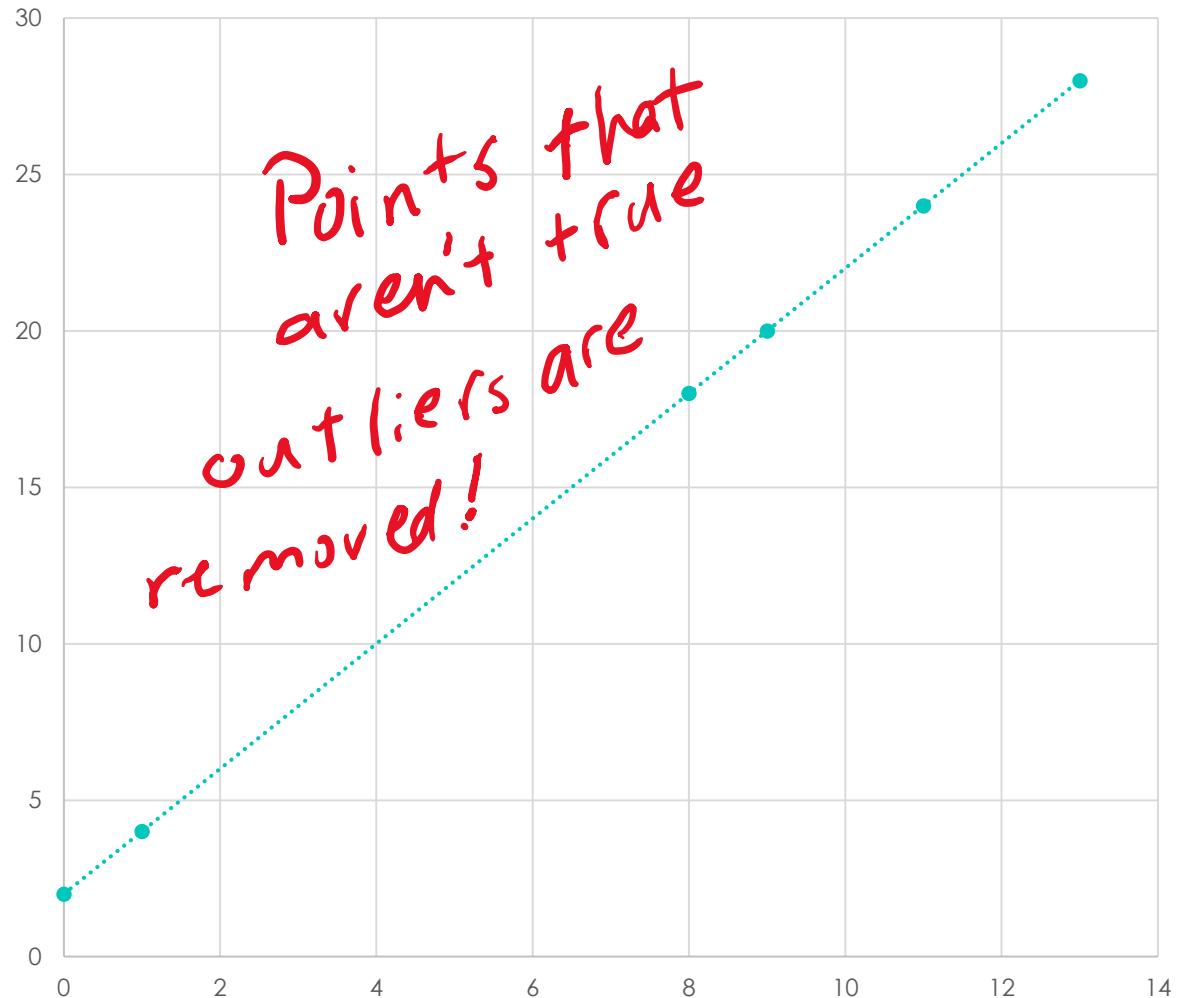




# OMP for Outlier Detection and Removal

-

## Visualization



After  
Several  
Iterations...

## Question

How do we know when  
to stop running  
Orthogonal Matching  
Pursuit for outlier  
detection?



# Stopping Conditions

- Option 1: Stop after a certain number of outliers are removed
  - This is especially helpful when we have an idea of how many outliers exist in our data
  - How does this affect our algorithm?
    - For  $k$  outliers, iterate while  $j \leq k$

# Stopping Conditions

- Option 1: Stop after a certain number of outliers are removed
  - This is especially helpful when we have an idea of how many outliers exist in our data
  - How does this affect our algorithm?
    - For  $k$  outliers, iterate while  $j \leq k$
- Option 2: stop once we've reached a certain residual threshold
  - Stop iterating when the norm of our residual  $r$  falls below a predefined value  $\epsilon$
  - How does this affect our algorithm?
    - Iterate while  $\|r\| \geq \epsilon$



Note: the values of k and epsilon are set **by us**, and are not a quantity determined by the OMP algorithm.



This means we can tune them!

# Tuning our Stopping Condition

- Last week, you learned about validation (and more specifically, cross-validation) as a method to turn hyperparameters so that our model performs well on our test dataset
- We can use this for finding an OMP stopping condition as well!
- Option 1
  - Cross-validation to set our residual threshold  $\epsilon$
  - We choose the value of  $\epsilon$  for which the error on our validation set is lowest
- Option 2
  - Tuning our number of removed points  $k$
  - We choose the value of  $k$  for which the error on our validation set is lowest
- Option 3
  - We stop when there are only two points left in our dataset (the minimum needed to create our linear model)

# OMP for Outlier Detection – Putting it all together

---

**Algorithm 2:** OMP to Detect Outliers

---

**Input:** Data matrix  $X$ , Predictions  $\mathbf{y}$ , Sparsity  $k$ , Residual Threshold  $\epsilon$

**Output:** Set of data indices  $F$  where the outliers are

$$\mathbf{r} = \mathbf{y} - X(X^T X)^{-1} X^T \mathbf{y}, j = 1, A = X, F = \emptyset;$$

**while**  $j \leq k$  and  $\|\mathbf{r}\| \geq \epsilon$  **do**

$$i = \operatorname{argmax}_i |\mathbf{r}[i]|;$$

$$F = F \cup \{i\};$$

$$A = [A | \mathbf{e}_i];$$

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{y};$$

$$\mathbf{r} = \mathbf{y} - Ax;$$

$$j = j + 1;$$

**end**

**return**  $F$

---