

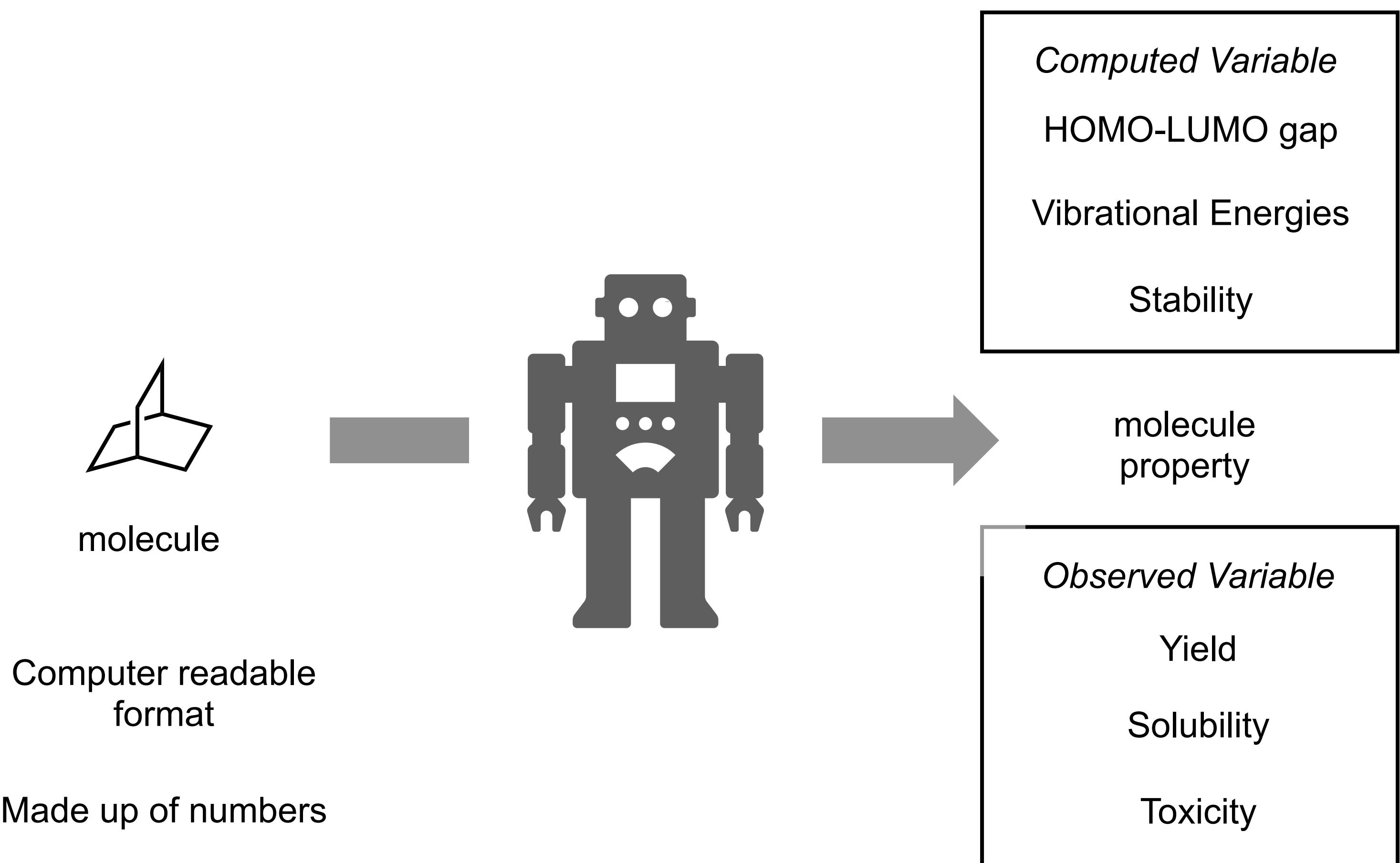
# *Machine Learning for Chemists*

– *Principal Component Analysis*  
  &  
*Support Vector Machines* –

---

12 February 2024

# Recall: What is Machine Learning, Really?





The “big idea” behind PCA

An example of PCA used in the wild

Introduction to Classifications

The “big idea” behind SVMs

An example of SVMs used in the wild

Live Demo!

# A Brief Introduction to Dimensions

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"

## 1 Variable

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"

# A Brief Introduction to Dimensions

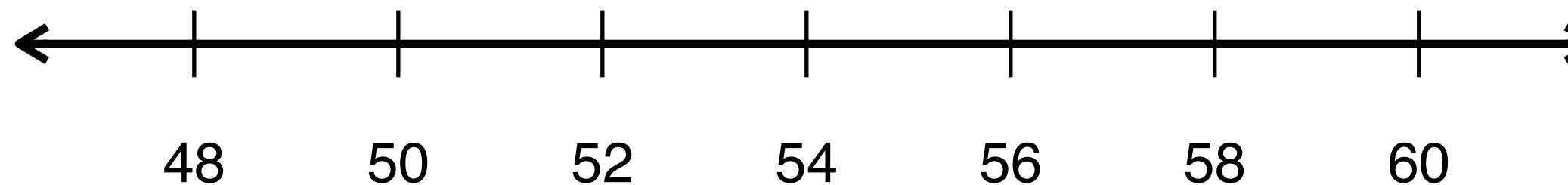


UNIVERSITY OF  
CAMBRIDGE

1 Variable

1 dimension

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"



Height

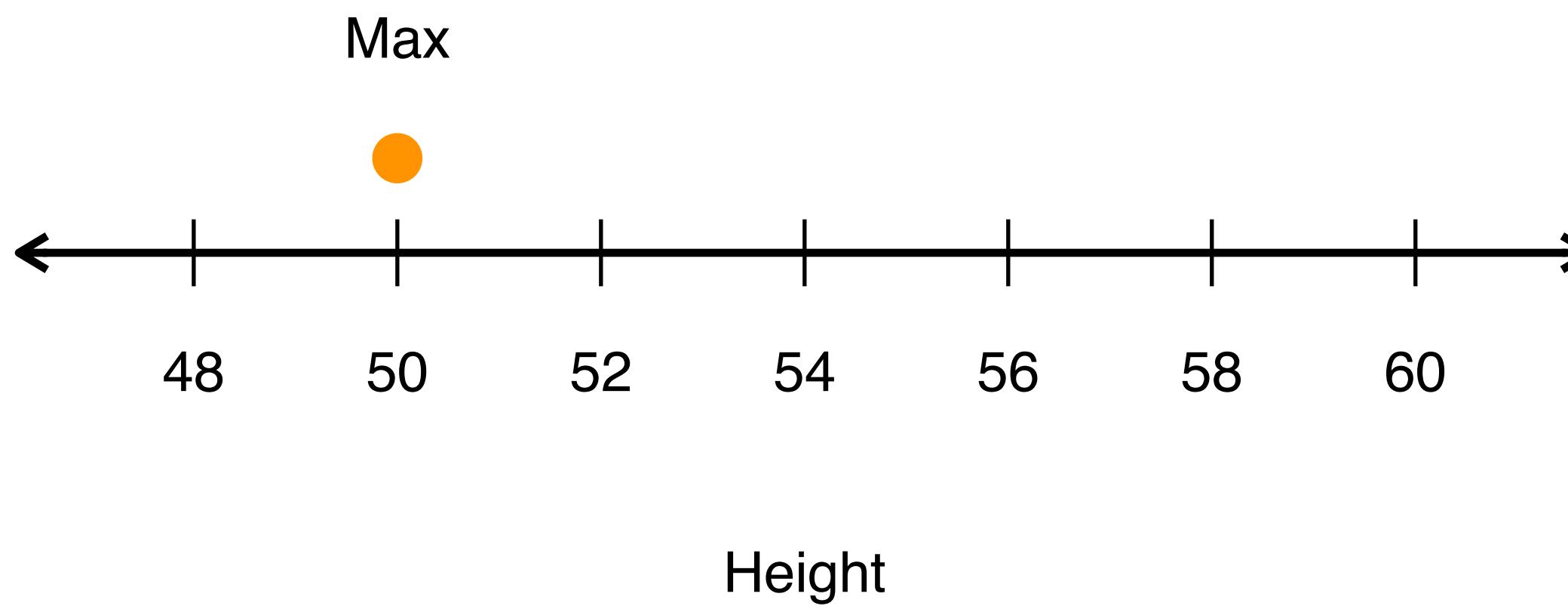
# A Brief Introduction to Dimensions



1 Variable

1 dimension

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"



# A Brief Introduction to Dimensions

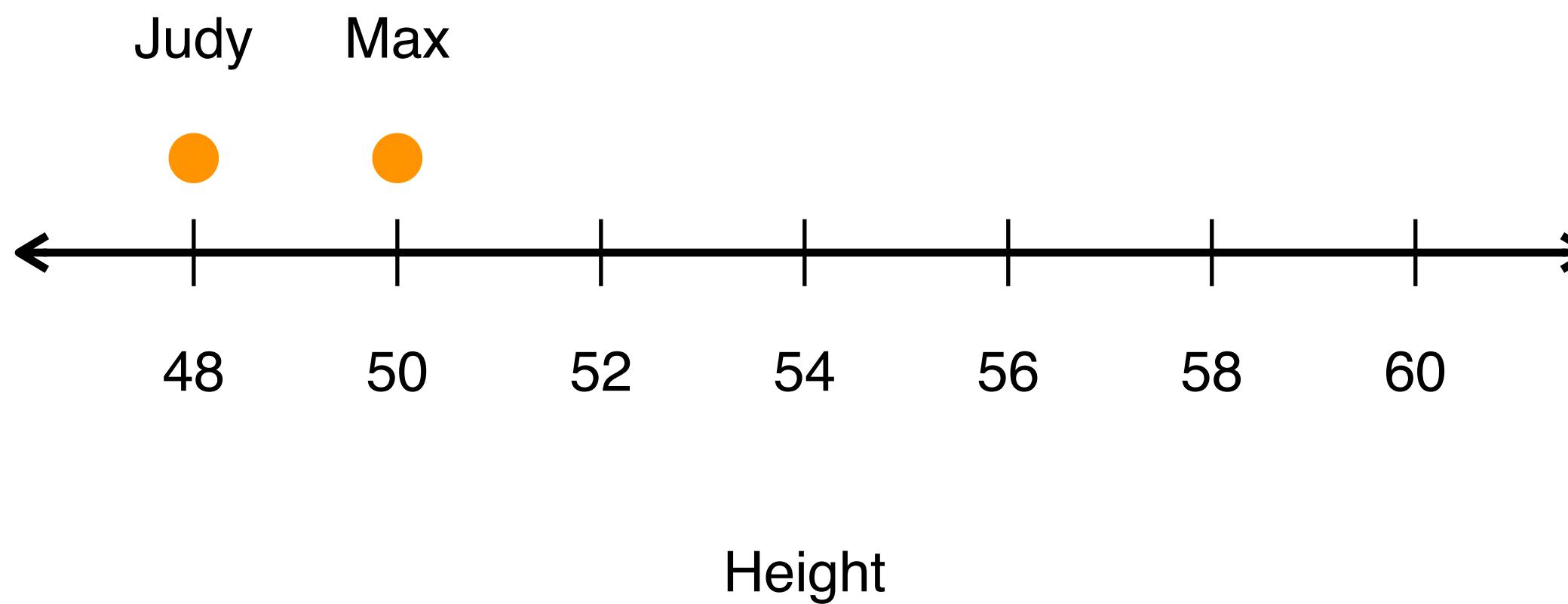


UNIVERSITY OF  
CAMBRIDGE

1 Variable

1 dimension

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"



# A Brief Introduction to Dimensions

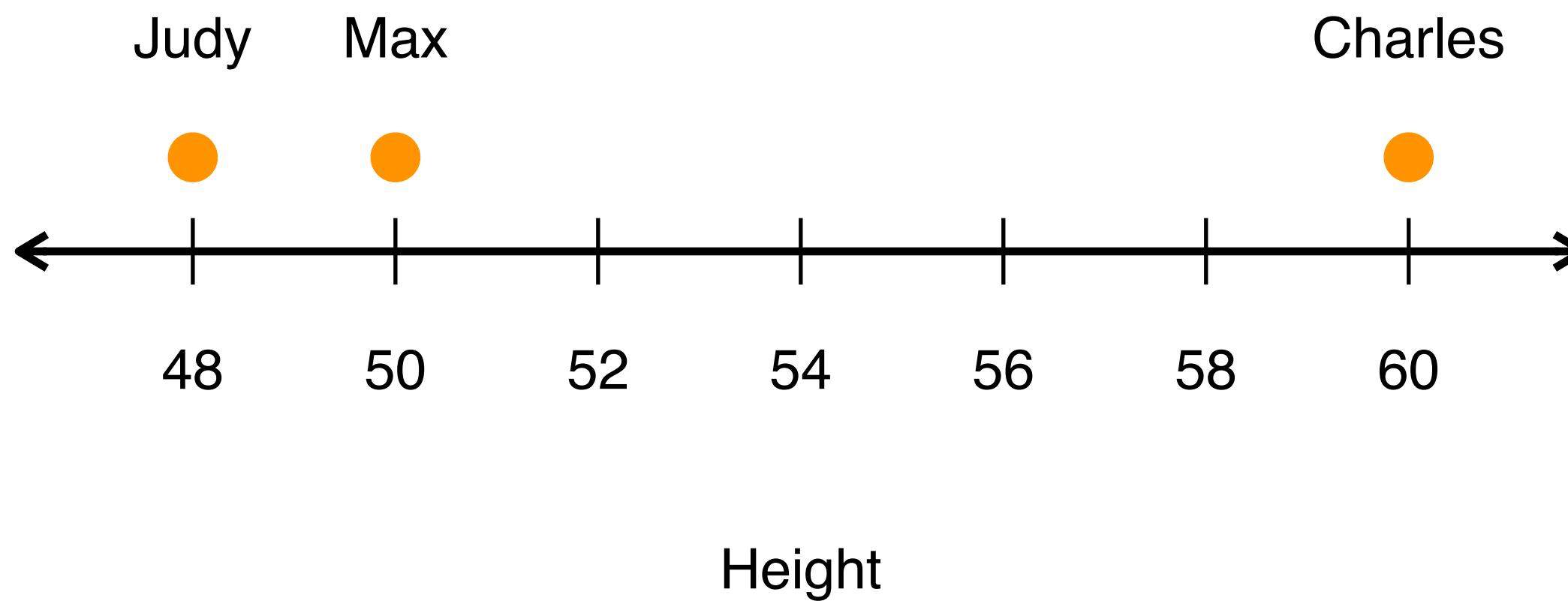


UNIVERSITY OF  
CAMBRIDGE

1 Variable

1 dimension

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"



# A Brief Introduction to Dimensions

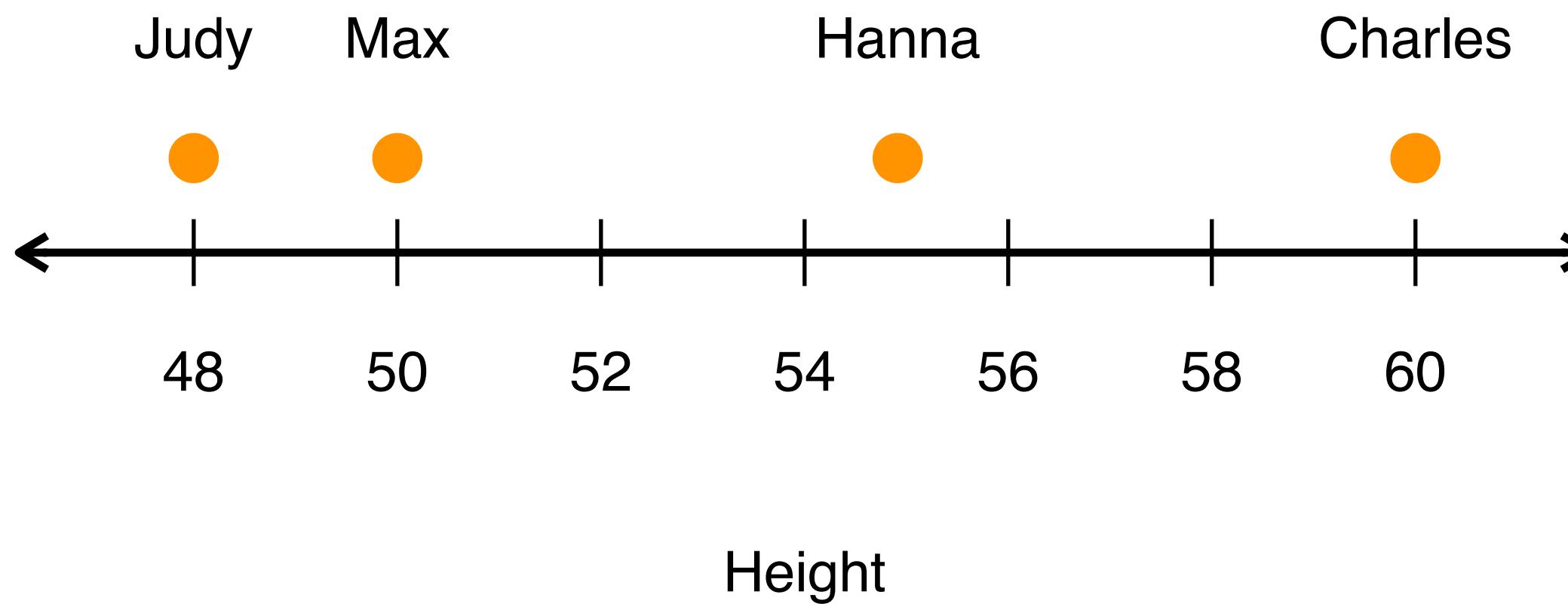


UNIVERSITY OF  
CAMBRIDGE

1 Variable

1 dimension

Name	Height
Max	50"
Judy	55"
Charles	60"
Hanna	48"



2 Variables

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7

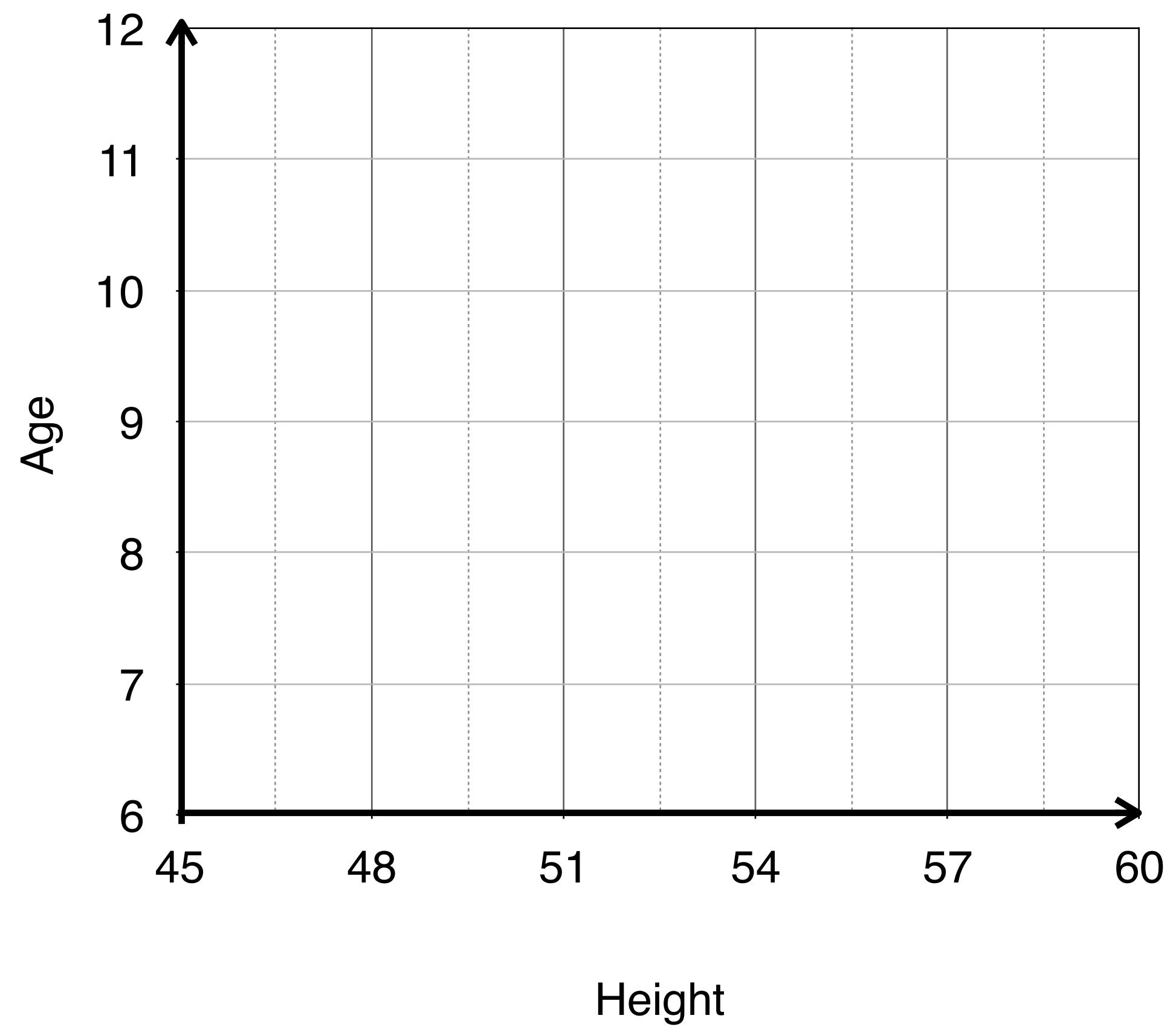
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



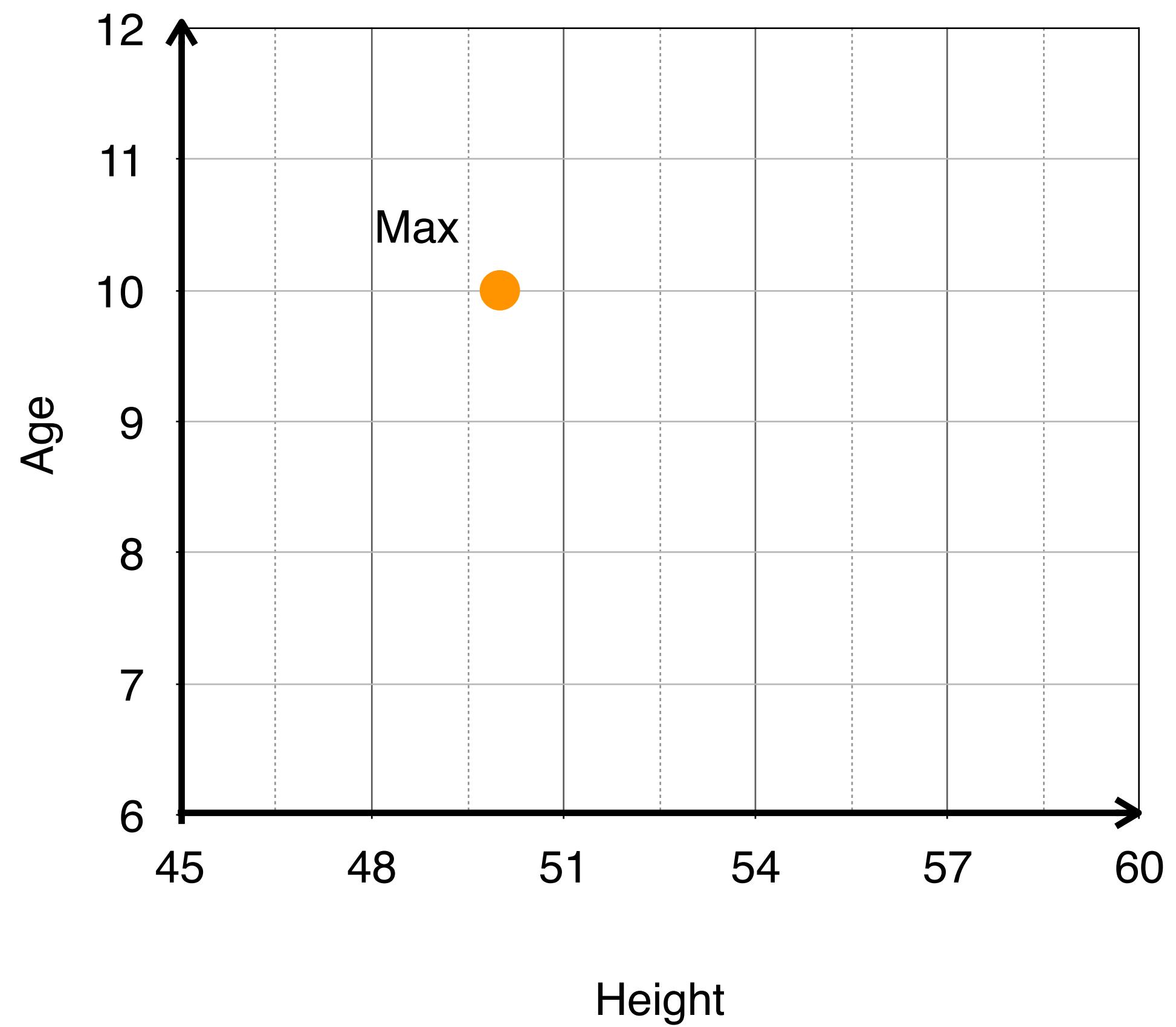
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



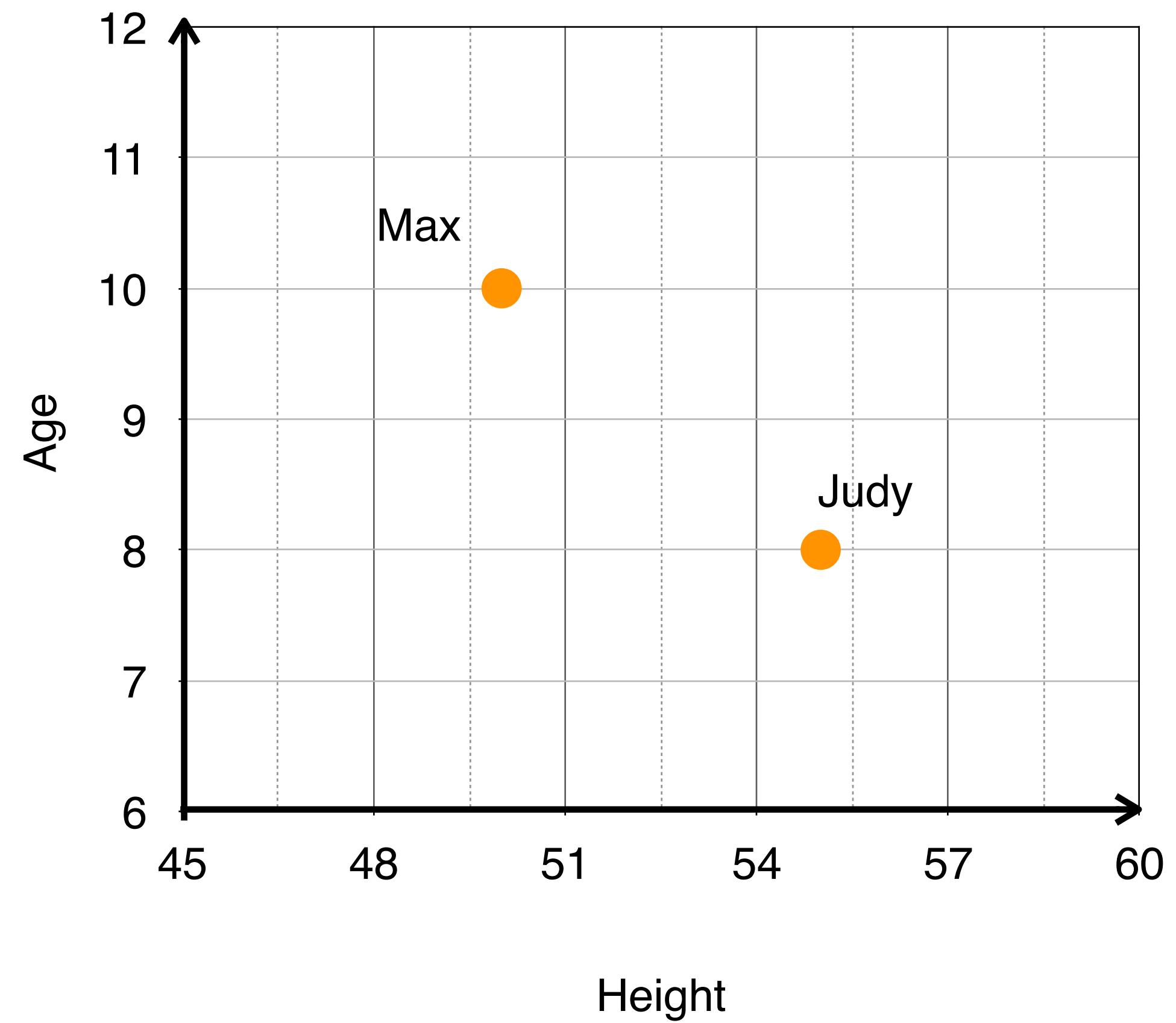
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



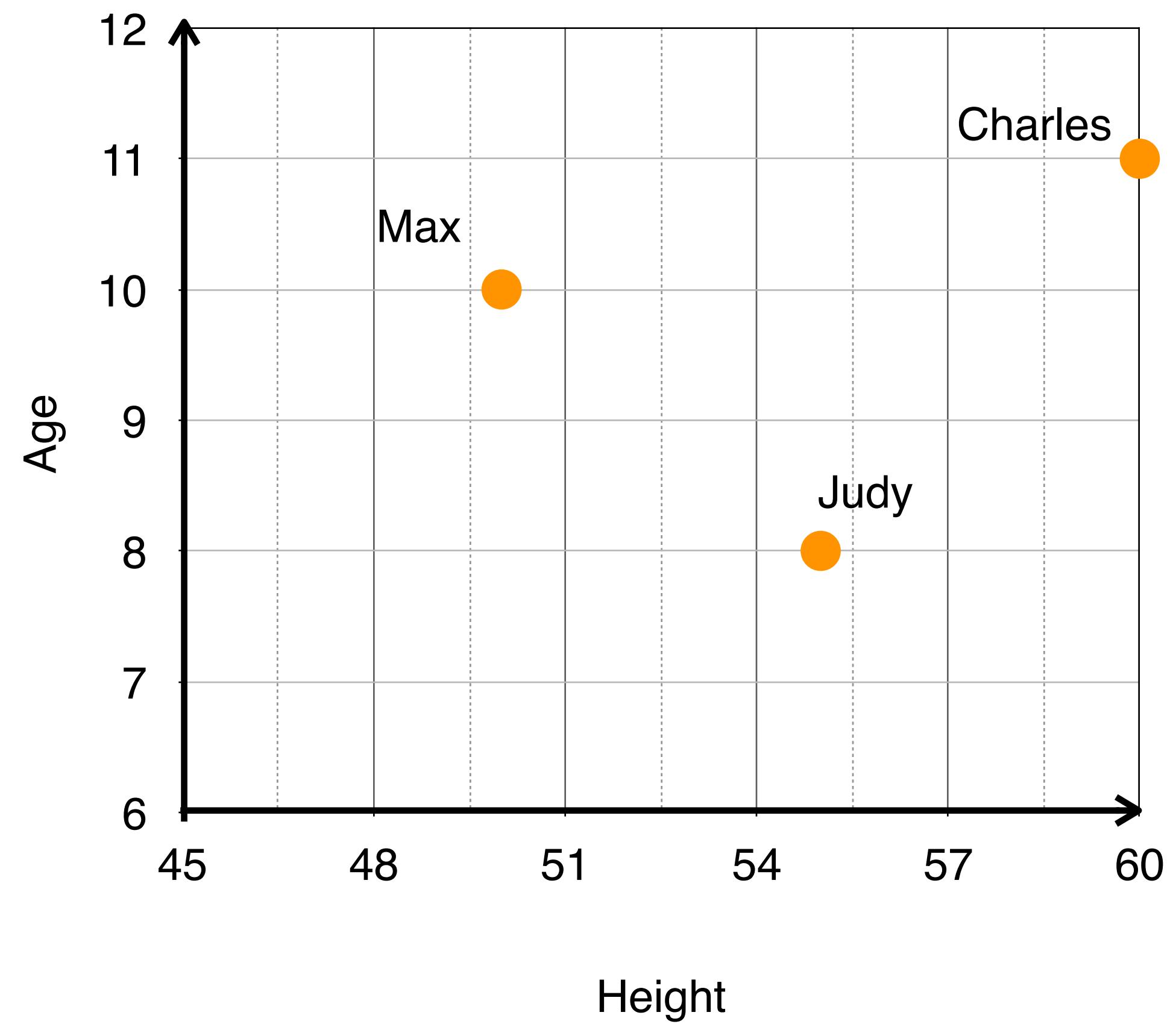
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



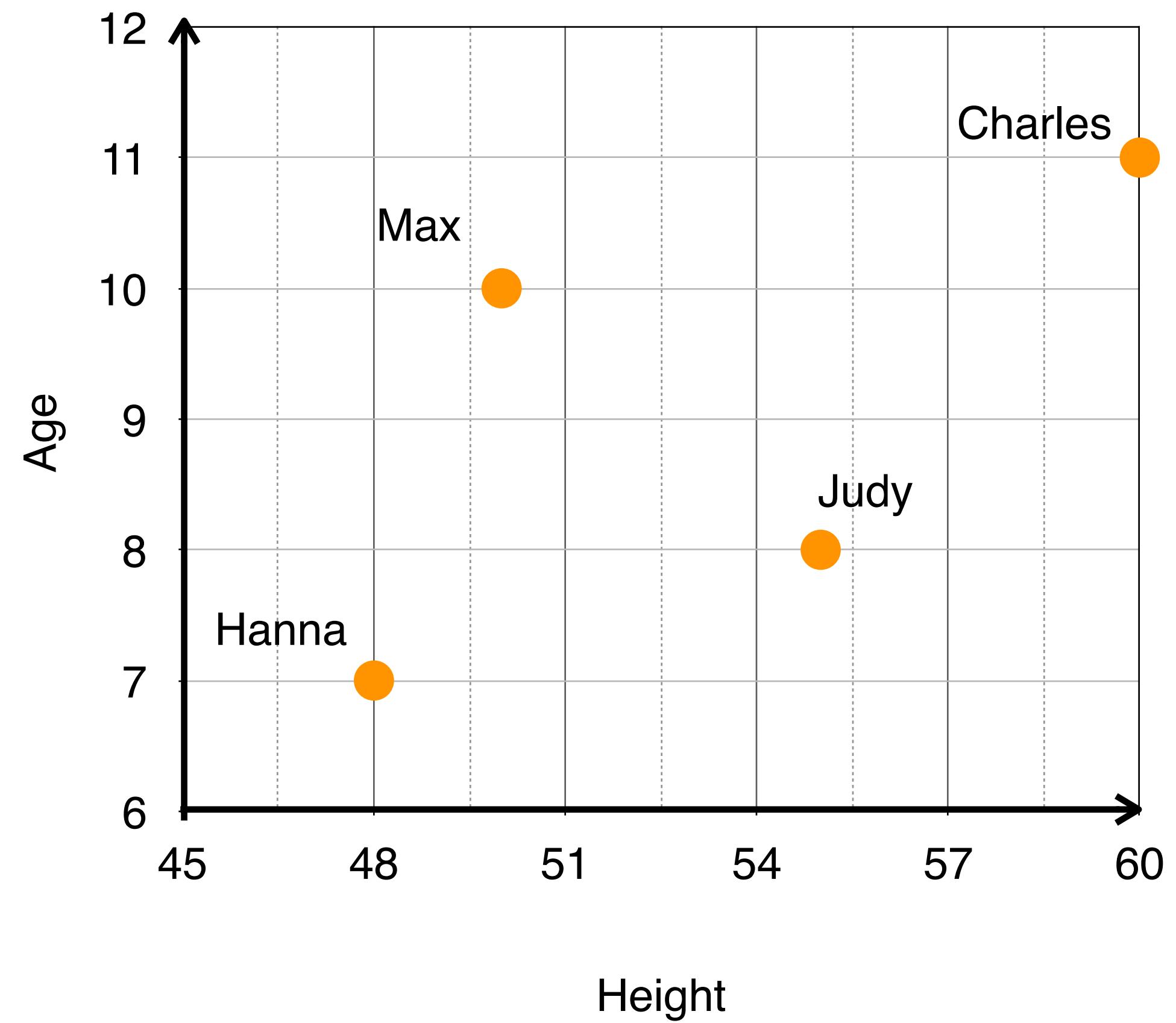
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



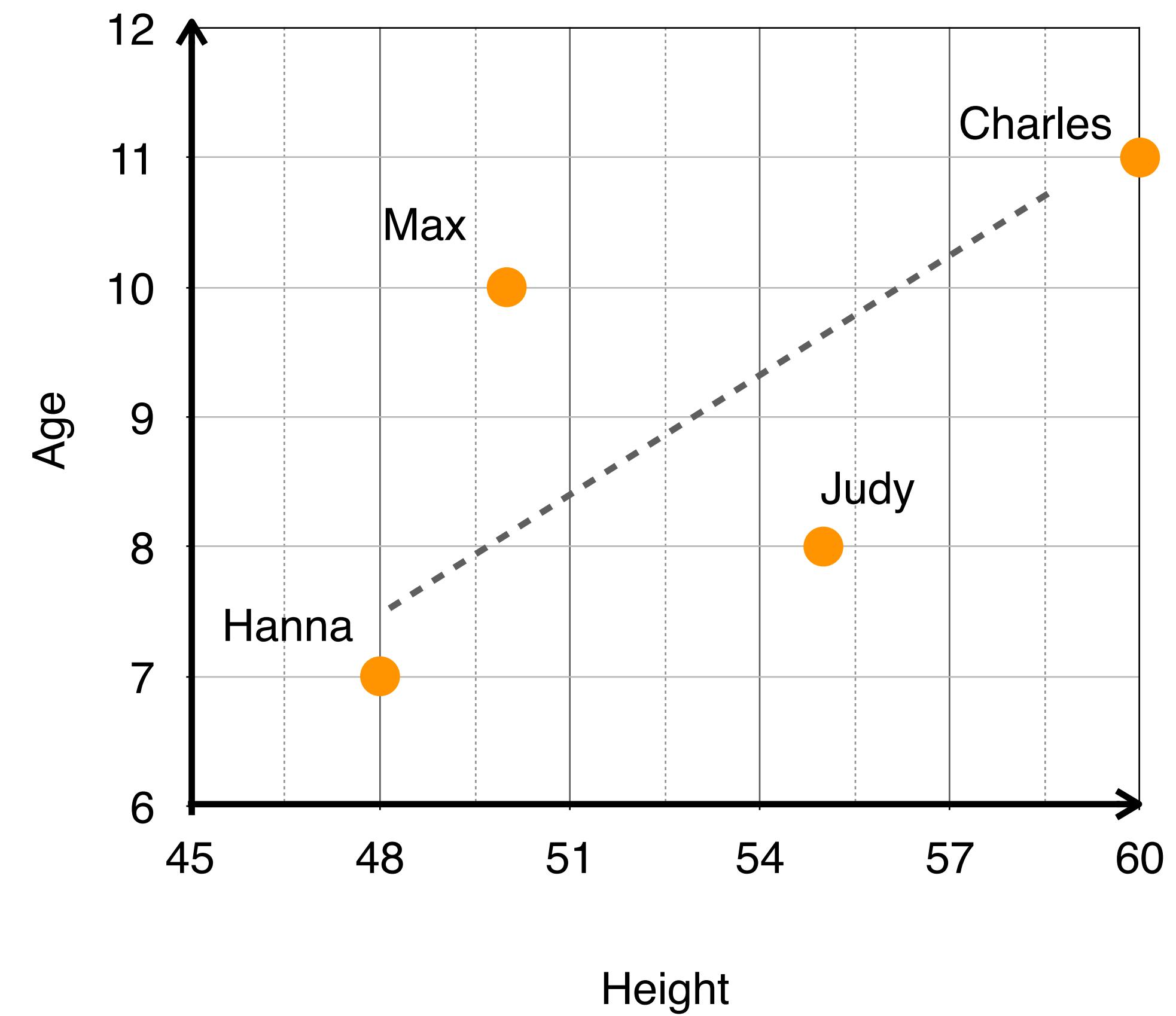
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



**positive** correlation  
between height and age

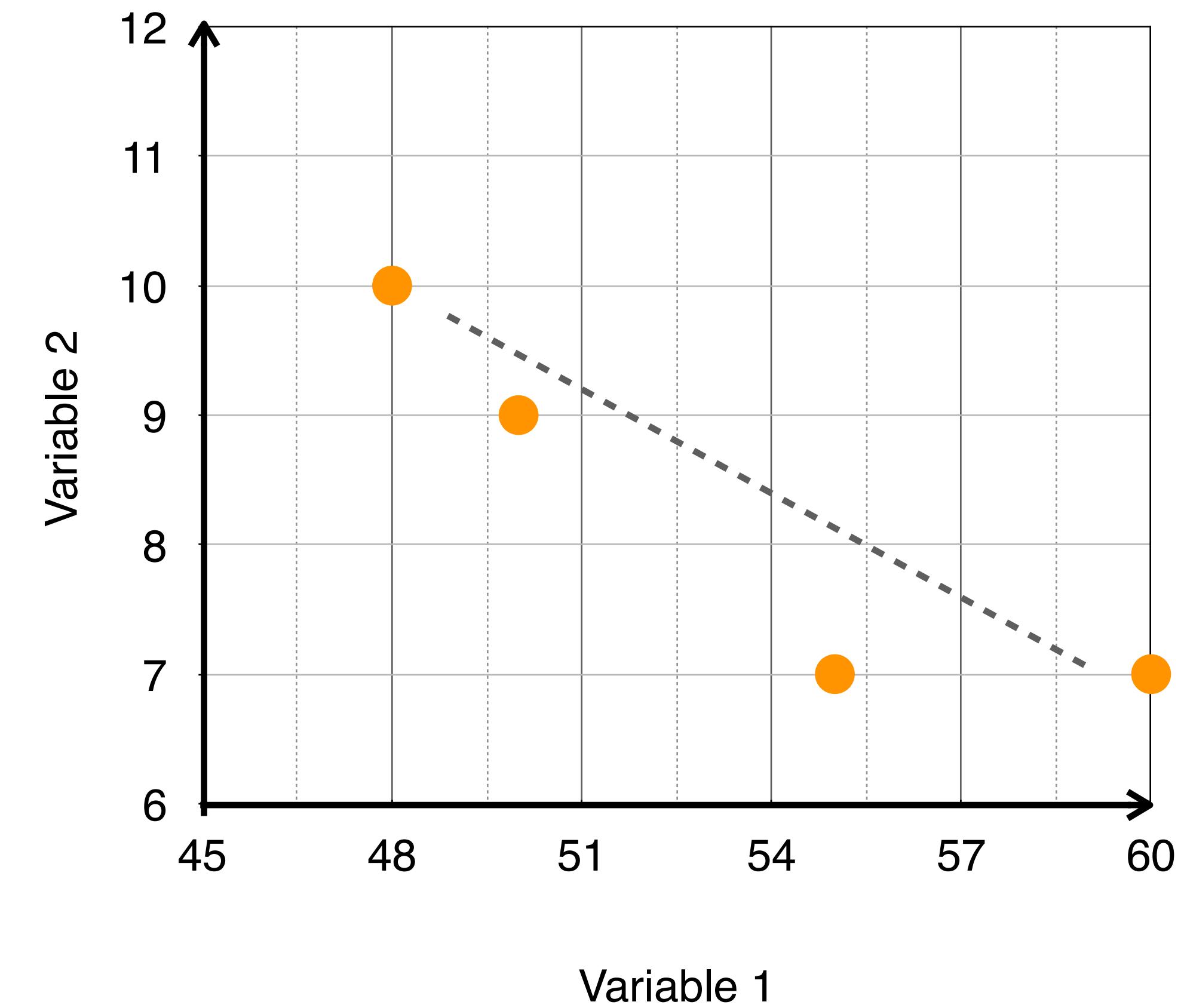
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



**negative** correlation between  
variable 1 and variable 2

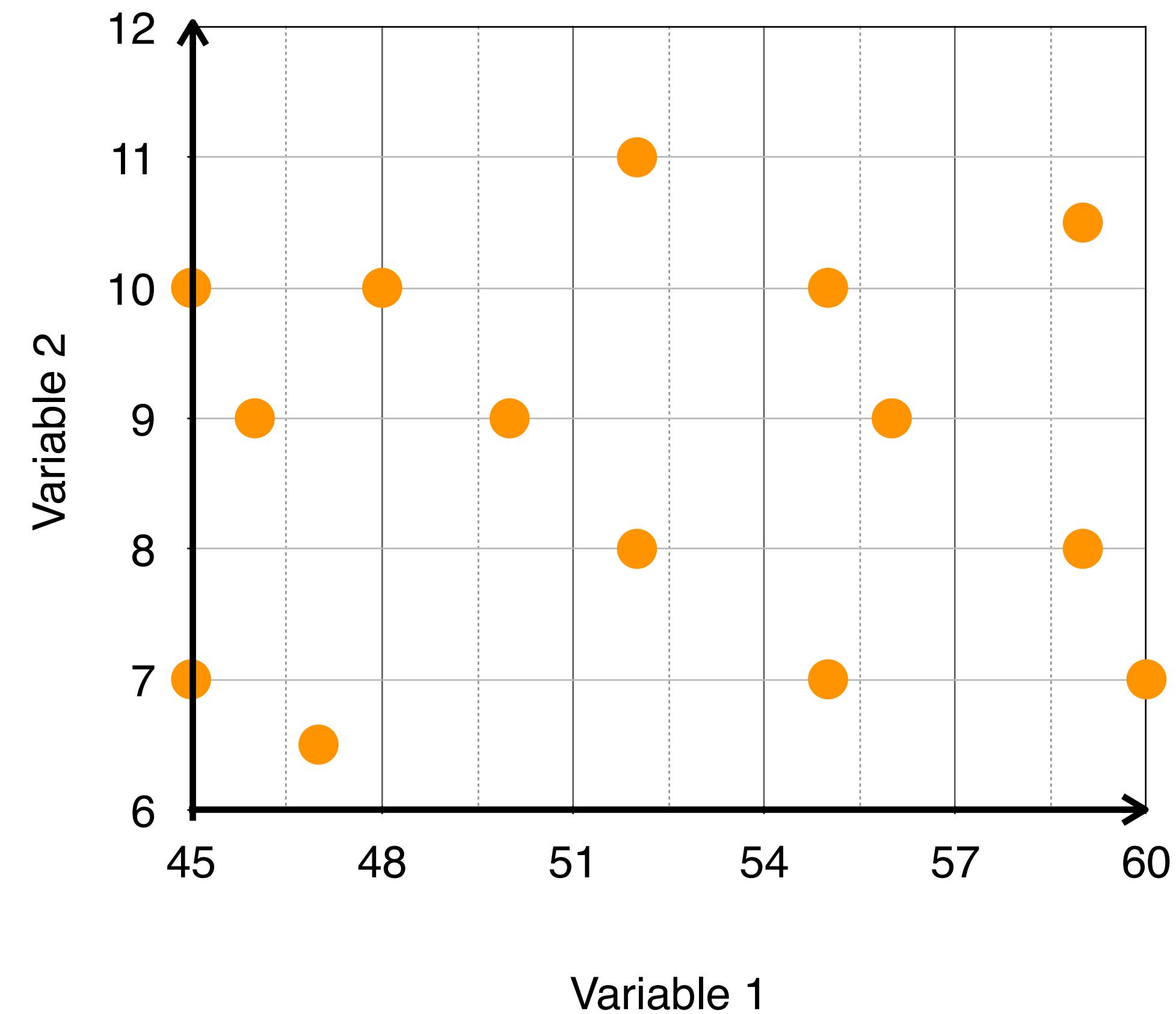
# A Brief Introduction to Dimensions



2 Variables

2 dimensions

Name	Height	Age
Max	50"	10
Judy	55"	8
Charles	60"	11
Hanna	48"	7



**no** correlation between variable  
1 and variable 2

3 Variables

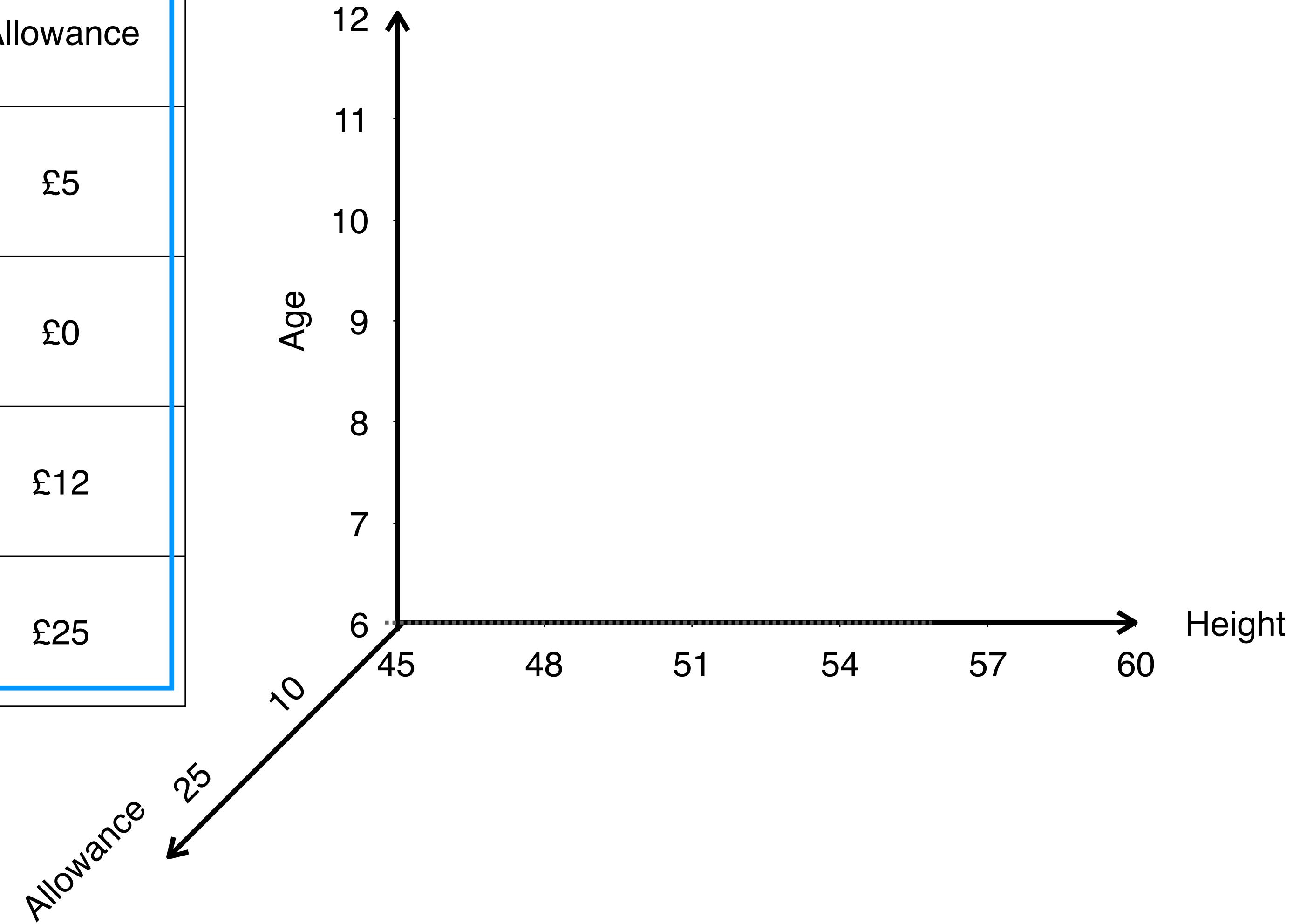
Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25

# A Brief Introduction to Dimensions

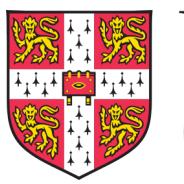
3 Variables

3 dimensions

Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25



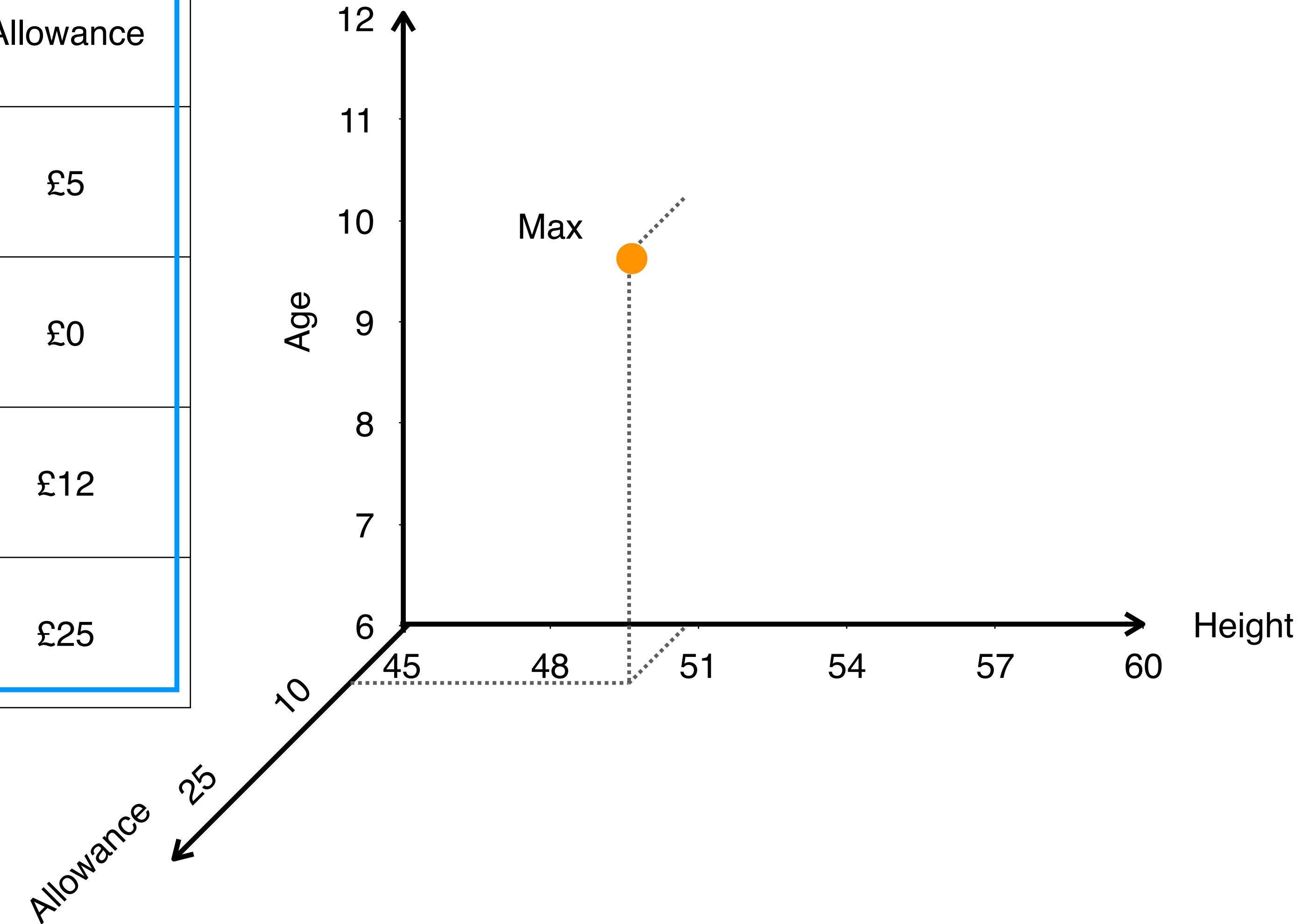
# A Brief Introduction to Dimensions



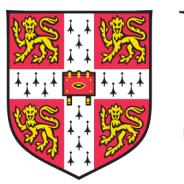
3 Variables

3 dimensions

Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25



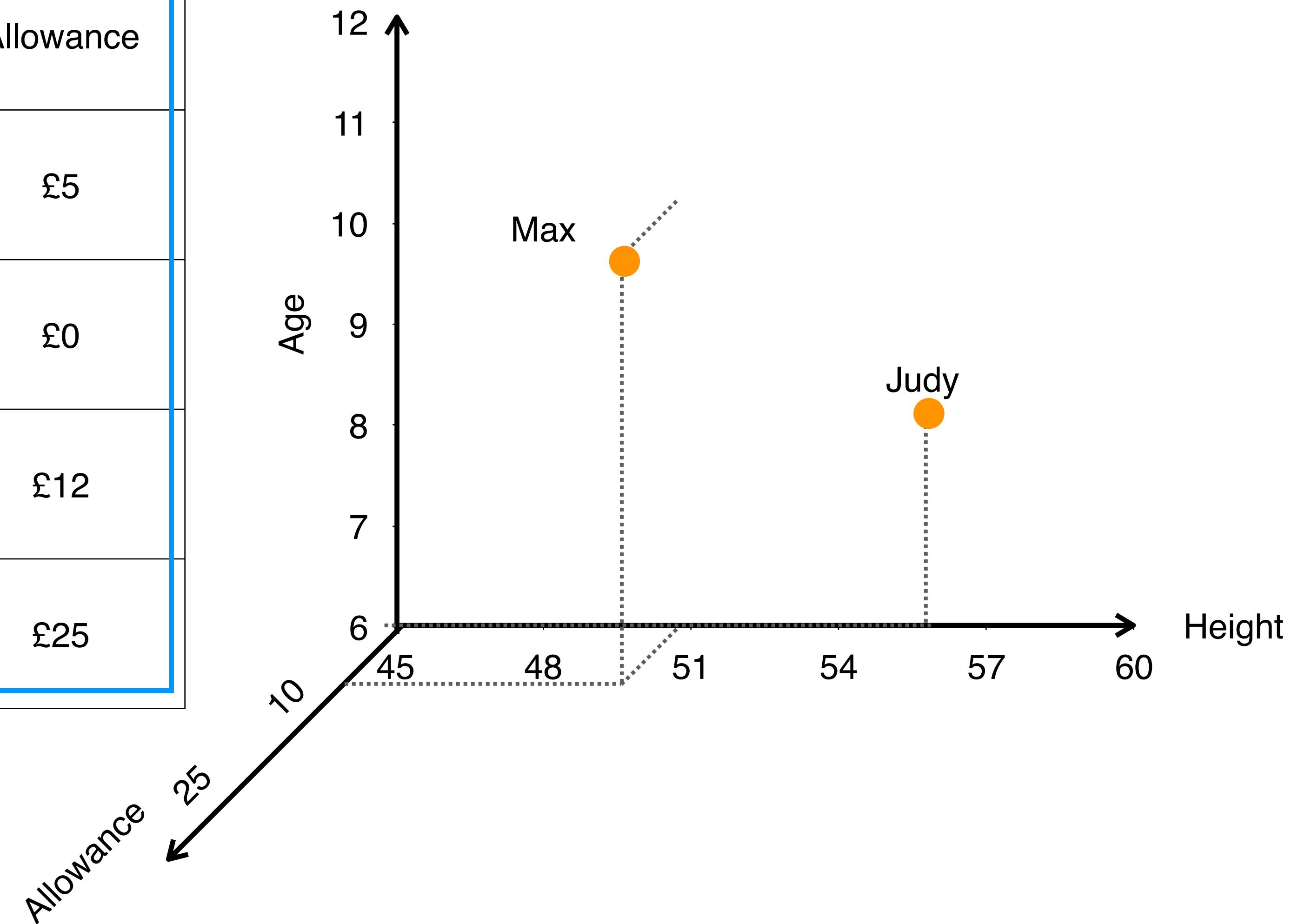
# A Brief Introduction to Dimensions



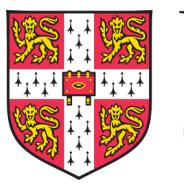
3 Variables

3 dimensions

Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25



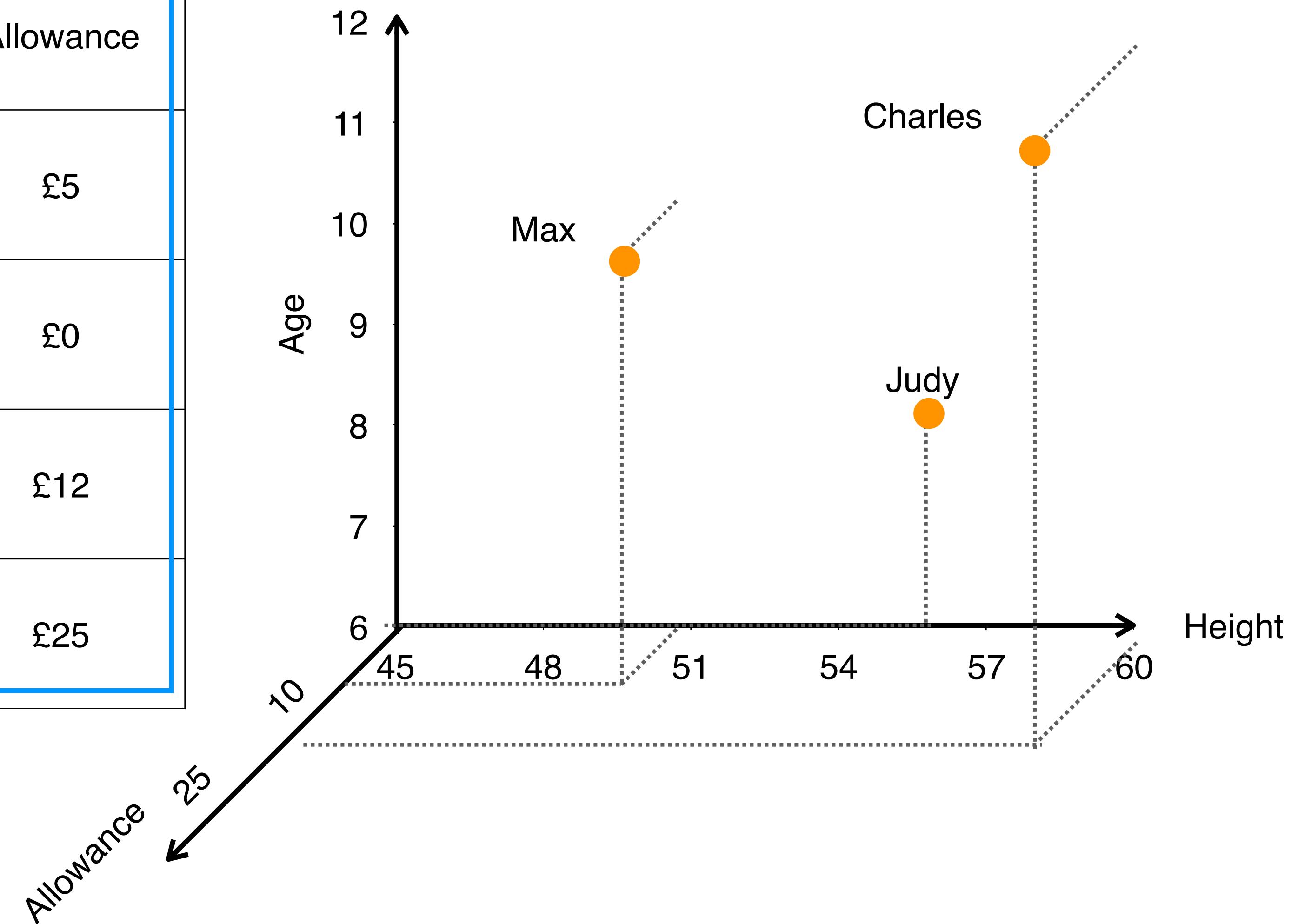
# A Brief Introduction to Dimensions



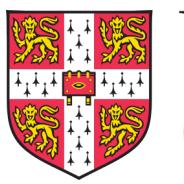
3 Variables

3 dimensions

Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25



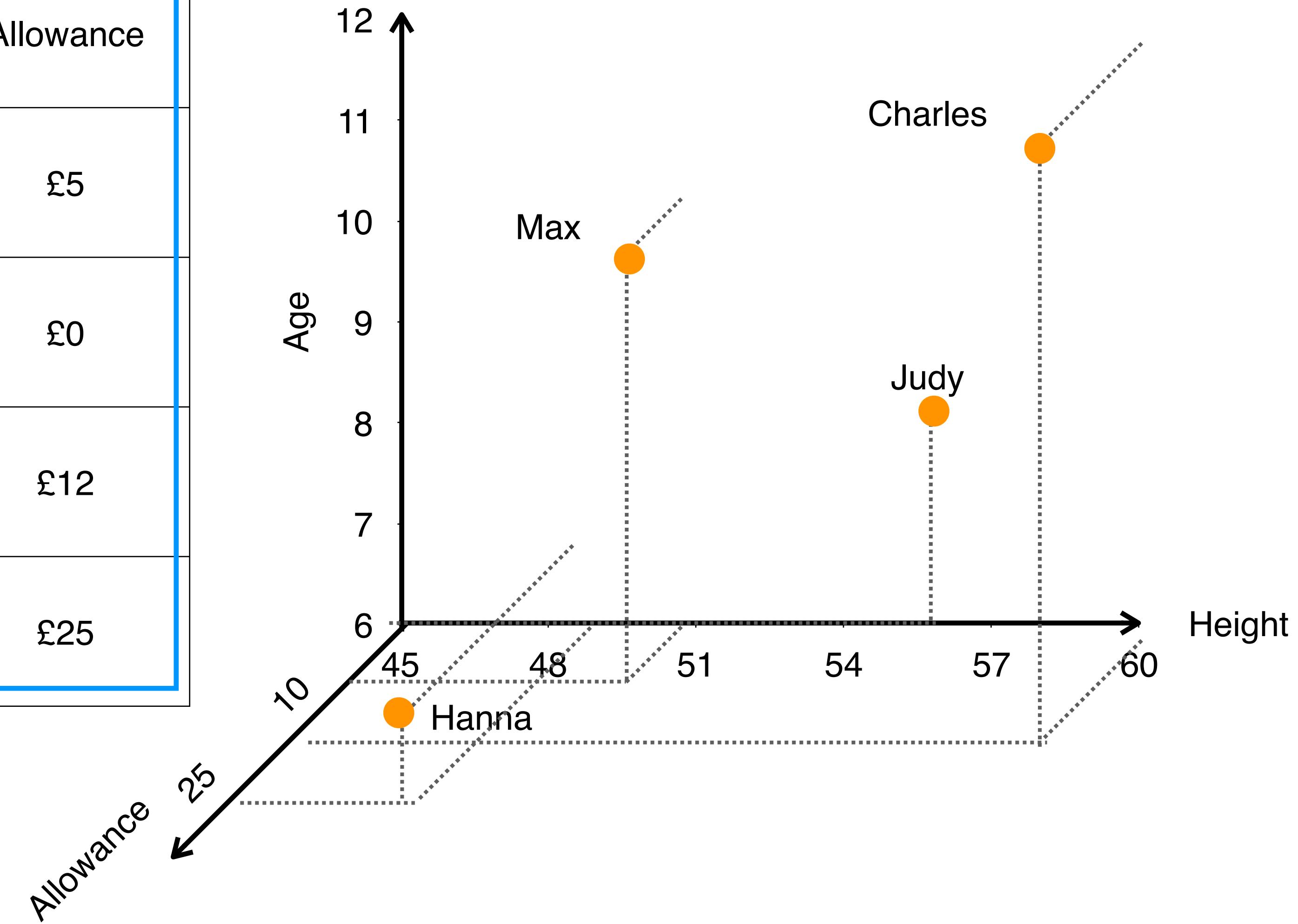
# A Brief Introduction to Dimensions



3 Variables

3 dimensions

Name	Height	Age	Allowance
Max	50"	10	£5
Judy	55"	8	£0
Charles	60"	11	£12
Hanna	48"	7	£25



# Why Would You Ever Want to Reduce Your Dimensions?

N Variables

N dimensions

Name	Height	Age	Allowance	Fav. #	...	...	...
Max	50"	10	£5	2	...	...	...
Judy	55"	8	£0	7	...	...	...
Charles	60"	11	£12	12	...	...	...
Hanna	48"	7	£25	55	...	...	...

# Why Would You Ever Want to Reduce Your Dimensions?

N Variables

N dimensions

Name	Height	Age	Allowance	Fav. #	...	...	...
Max	50"	10	£5	2	...	...	...
Judy	55"	8	£0	7	...	...	...
Charles	60"	11	£12	12	...	...	...
Hanna	48"	7	£25	55	...	...	...

You cannot draw N > 3 dimensions

# Are All These Dimensions Really That Important?

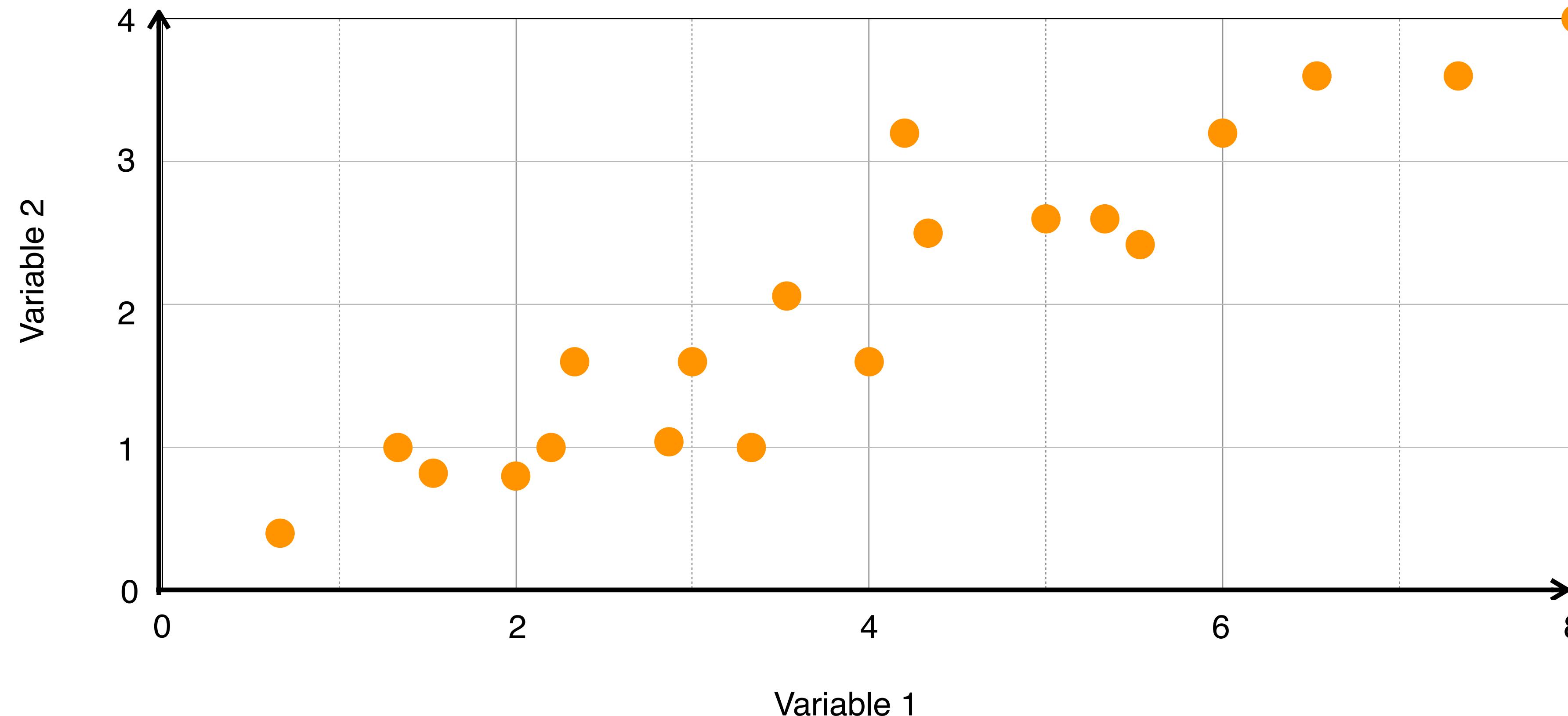


## Principal Component Analysis (PCA)

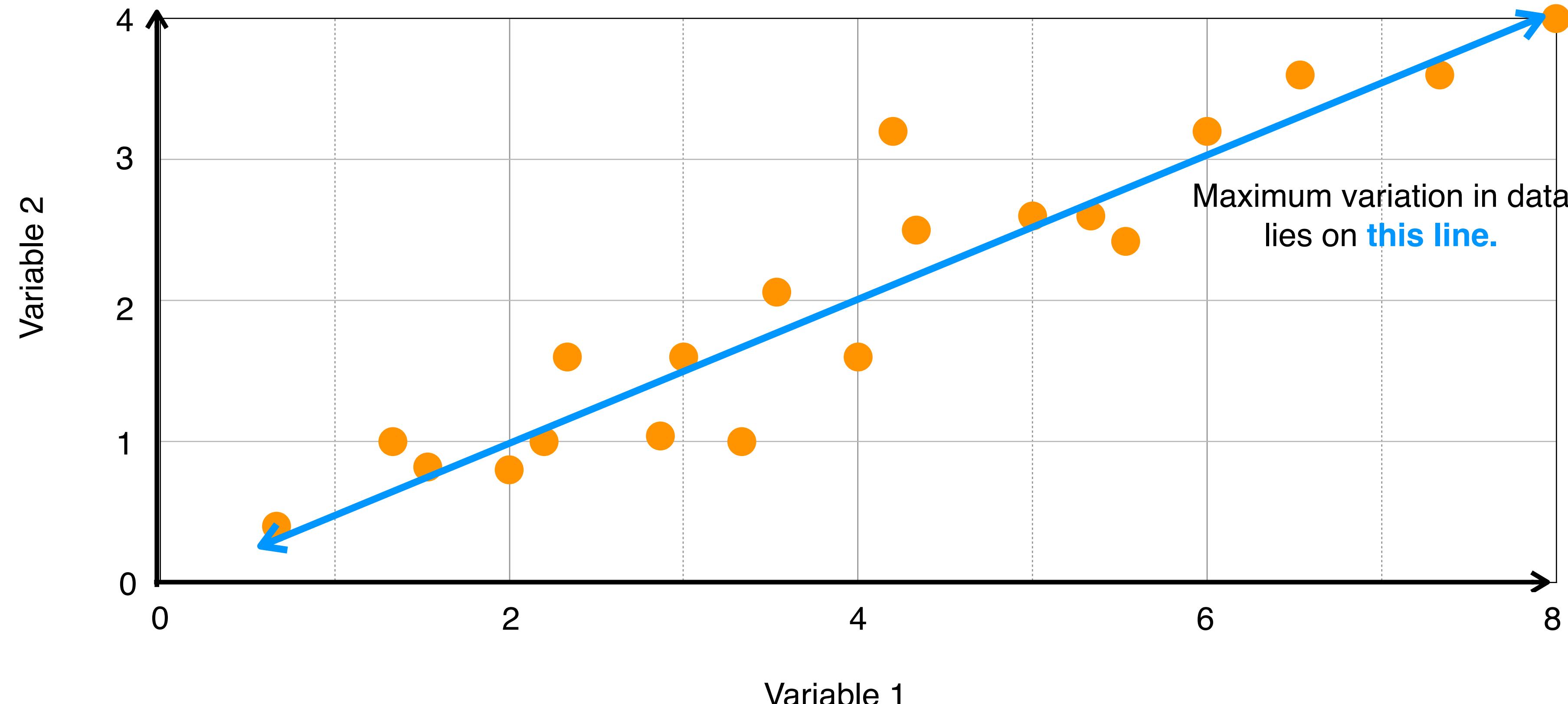
Finds a way to squeeze multidimensional data to more manageable dimensional data

Minimizes information lost in “squeezing” process

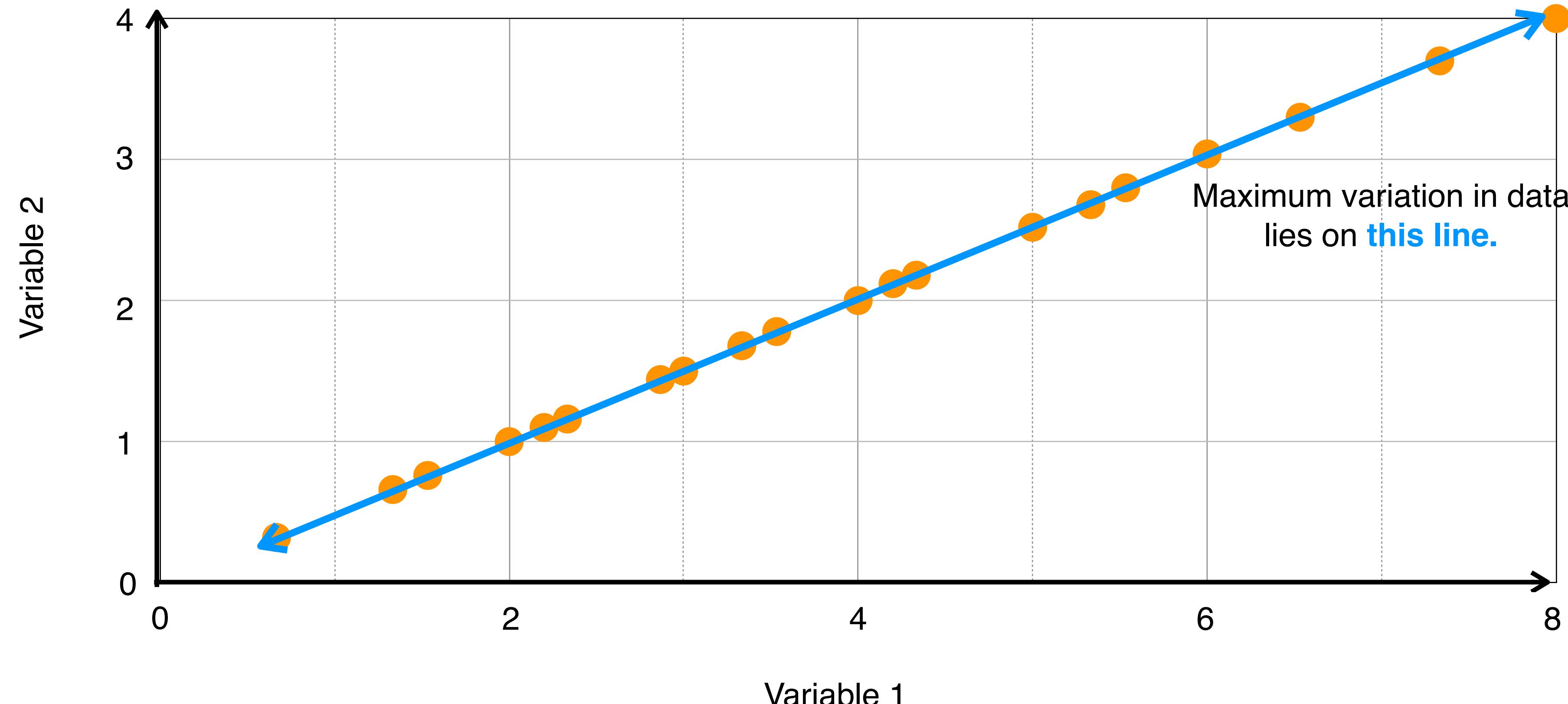
Allows us to visualize our data



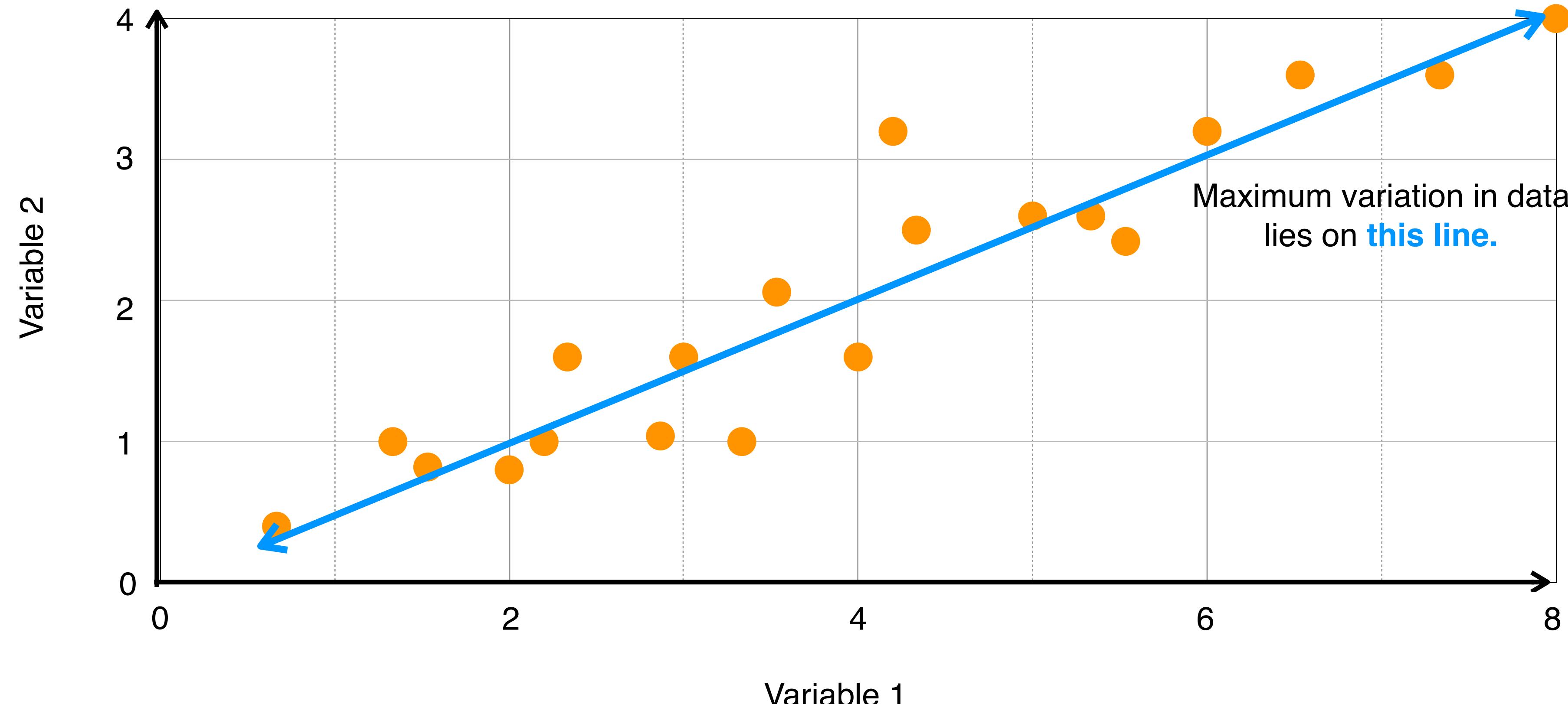
We want to find 2 perpendicular axes that encapsulate the most information (variation) possible.



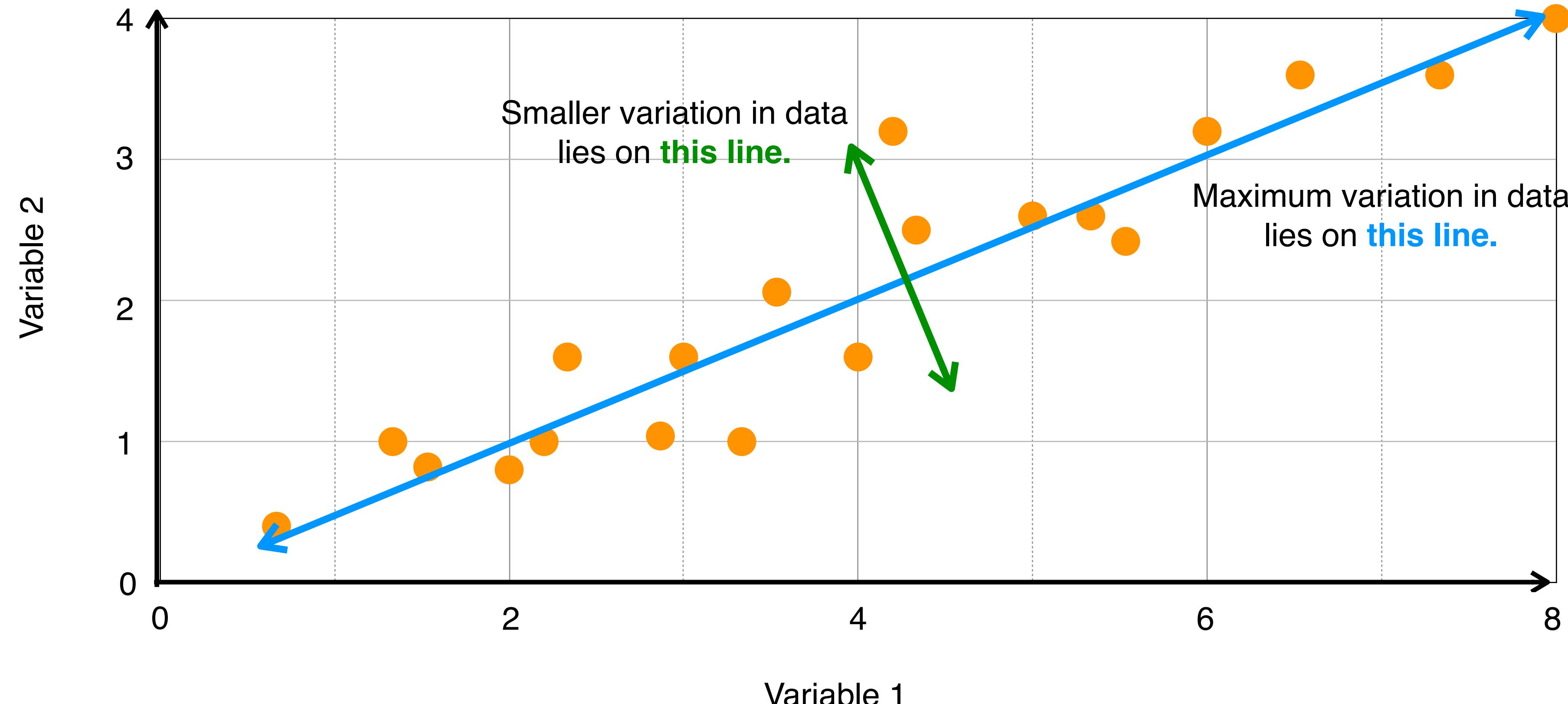
We want to find 2 perpendicular axes that encapsulate the most information (variation) possible.



We want to find 2 perpendicular axes that encapsulate the most information (variation) possible.

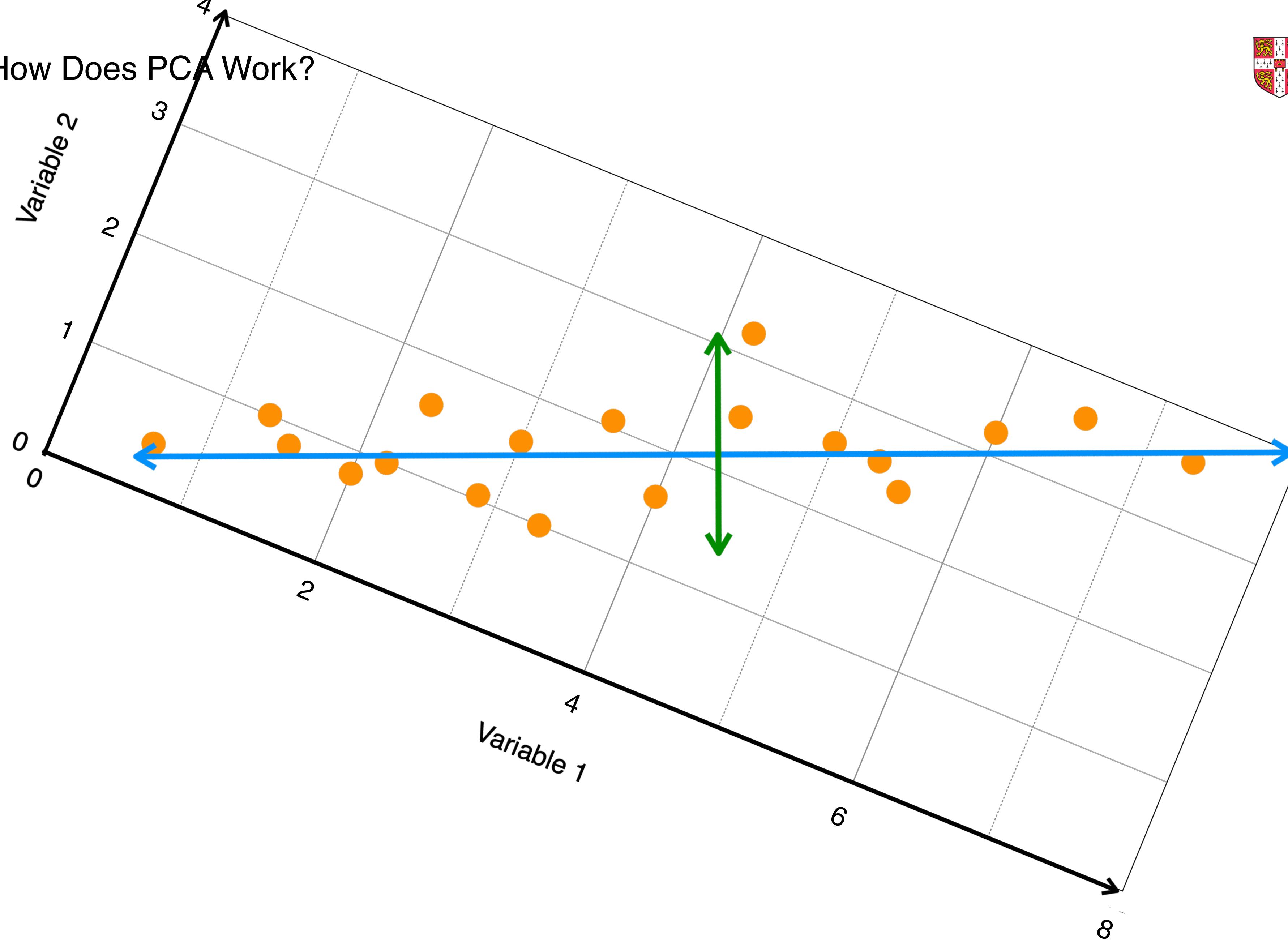


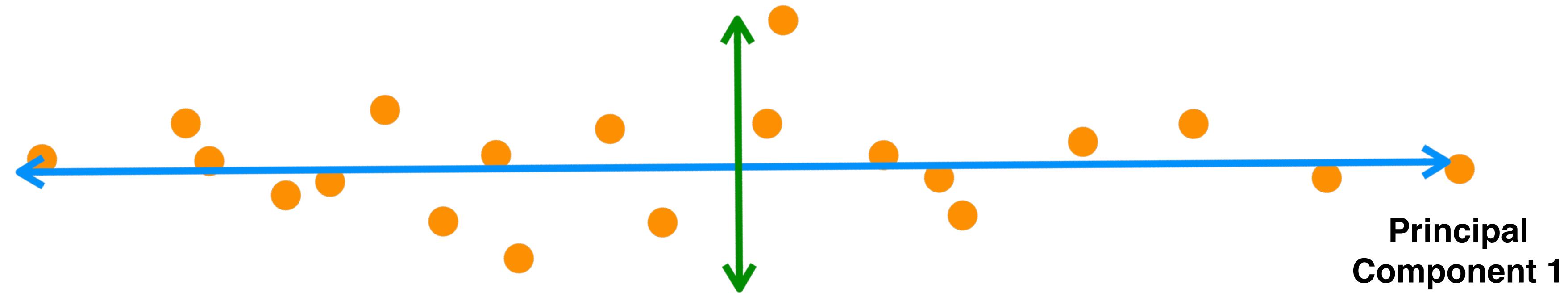
We want to find 2 perpendicular axes that encapsulate the most information (variation) possible.



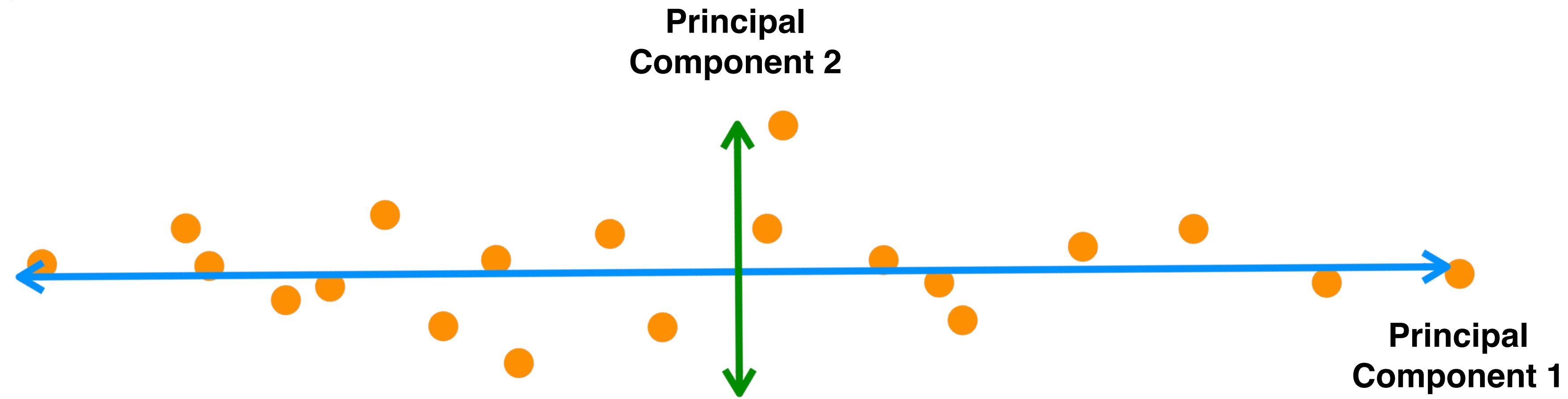
We want to find 2 perpendicular axes that encapsulate the most information (variation) possible.

## The Big Idea: How Does PCA Work?



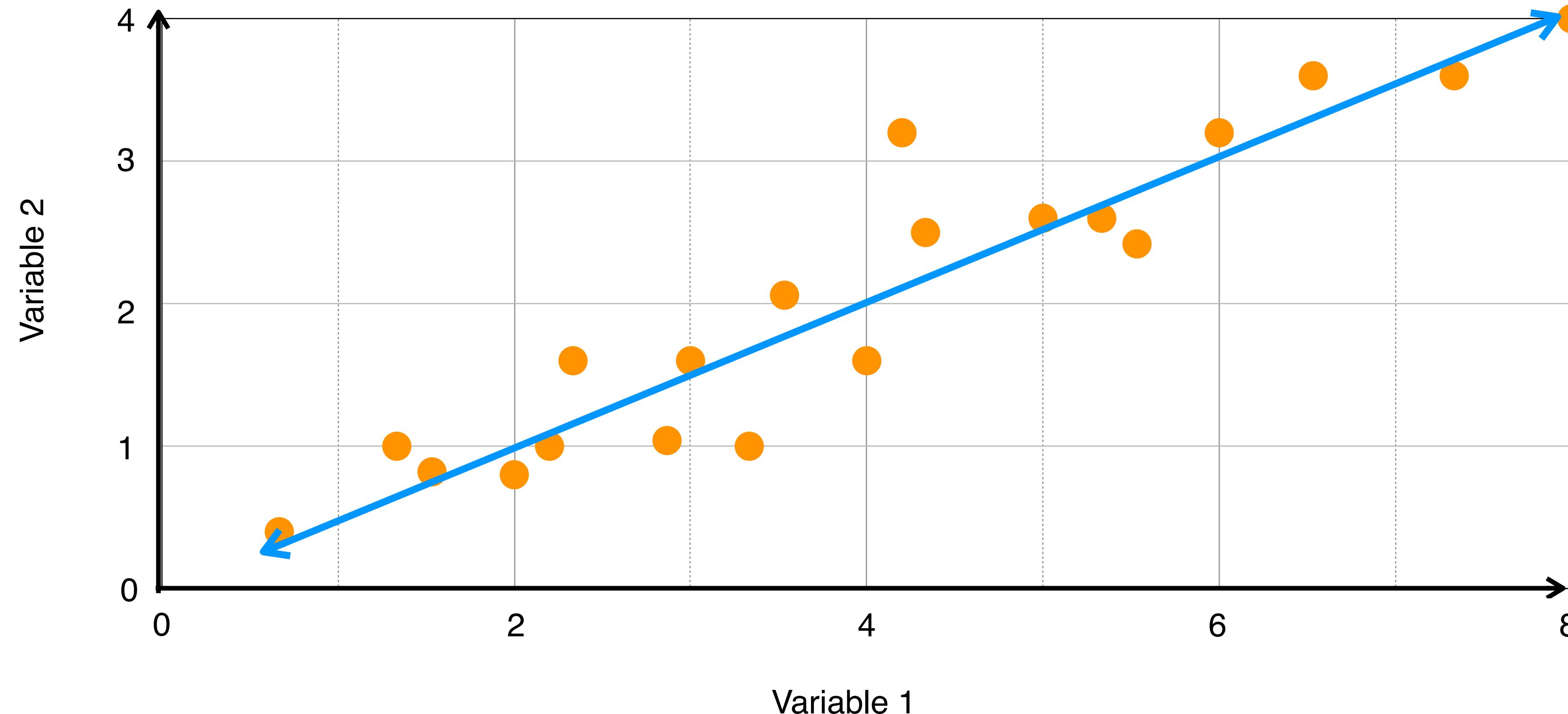


Principal Component 1 (PC1) captures the direction where most of the data variation occurs.

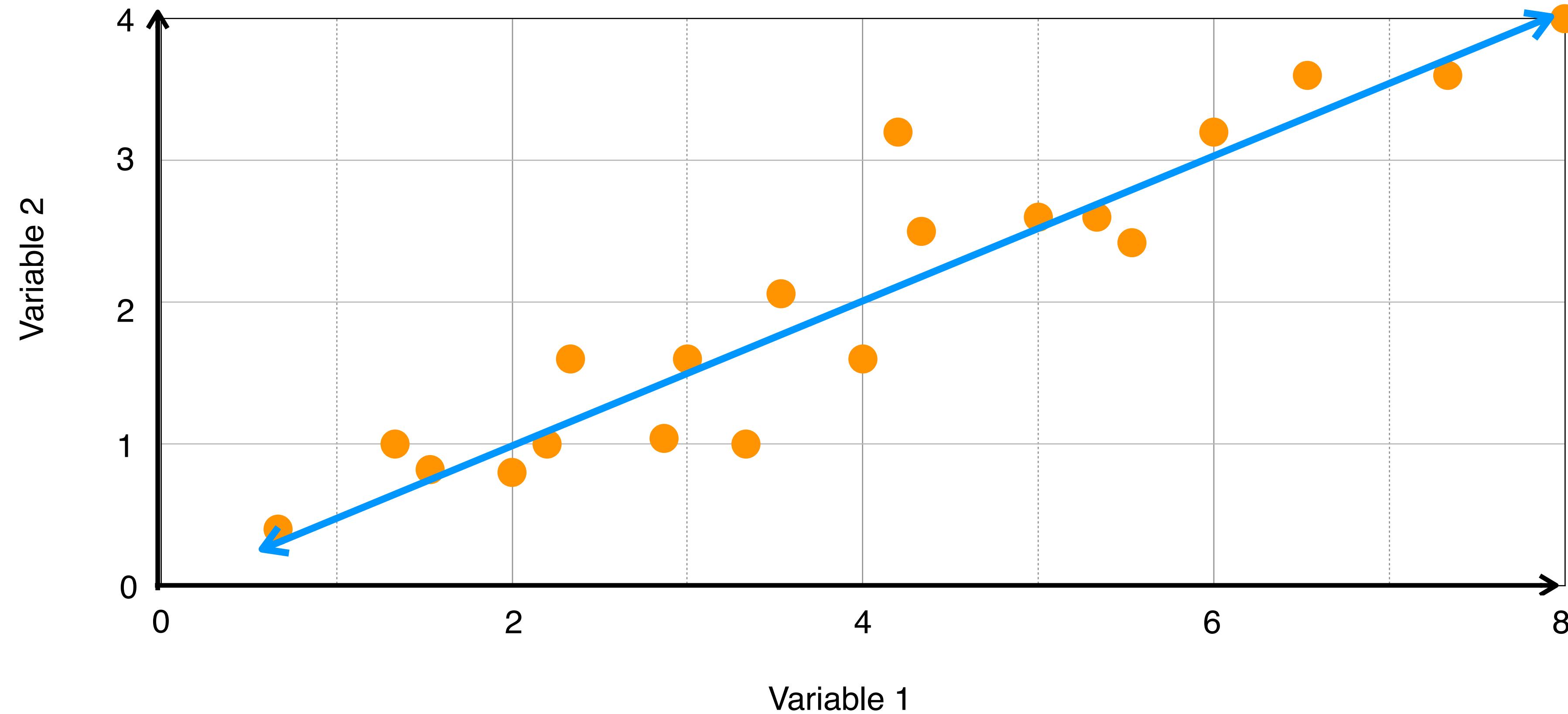


Principal Component 1 (PC1) captures the direction where most of the data variation occurs.

Principal Component 2 (PC2) captures the direction where the 2nd most data variation occurs.



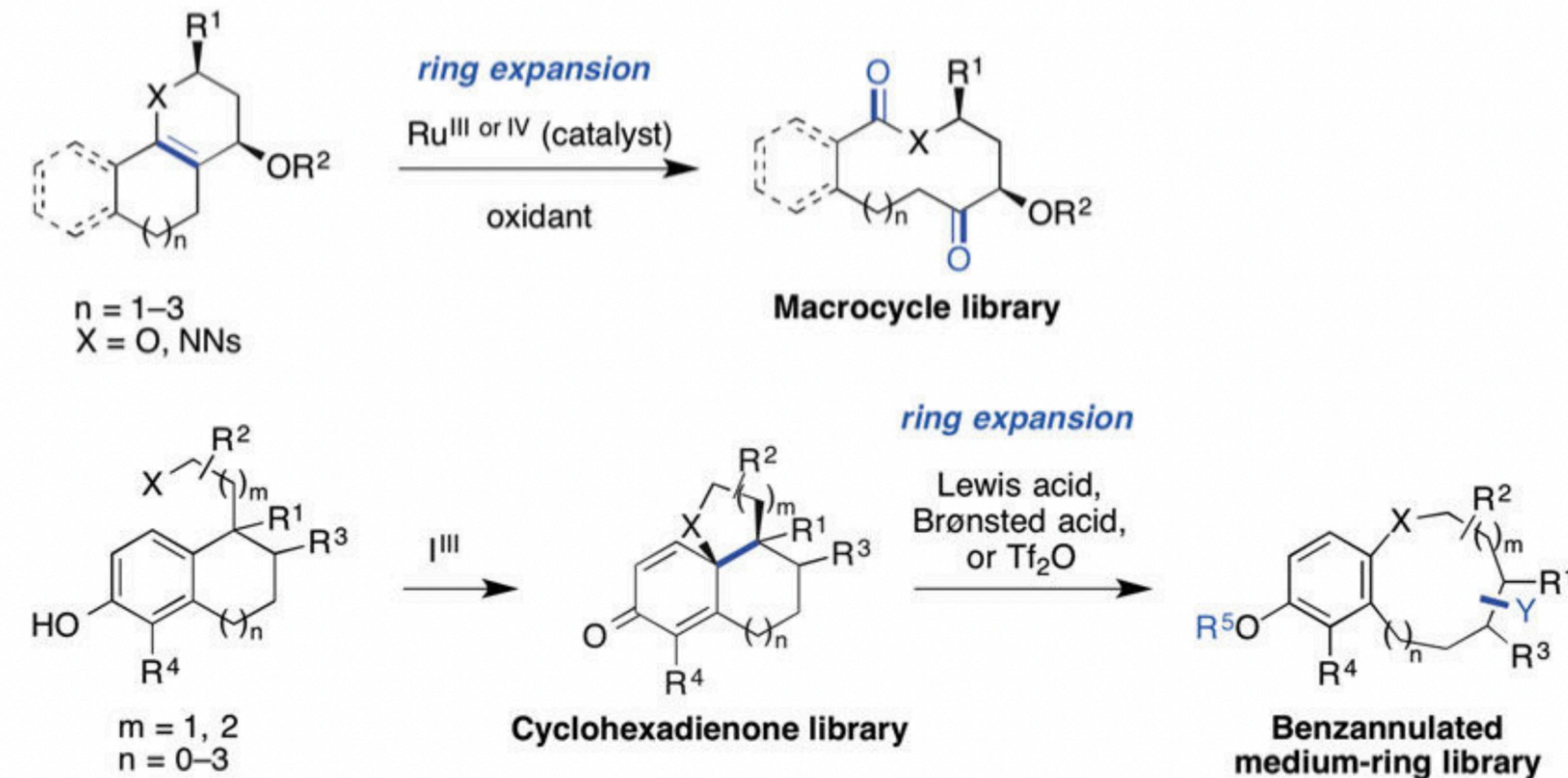
Principal Component 1 =  $(2 \cdot \text{Variable 1}) + (1 \cdot \text{Variable 2})$



Principal Component 1  $\approx (2 \cdot \text{Variable 1}) + (1 \cdot \text{Variable 2})^*$

\*NOTE: This is a slight simplification!  
In reality PC1 would be expressed as  
a linear combination of variable 1 and  
2 that has this 2:1 relationship, but has  
been scaled up or down.

## Guiding library design

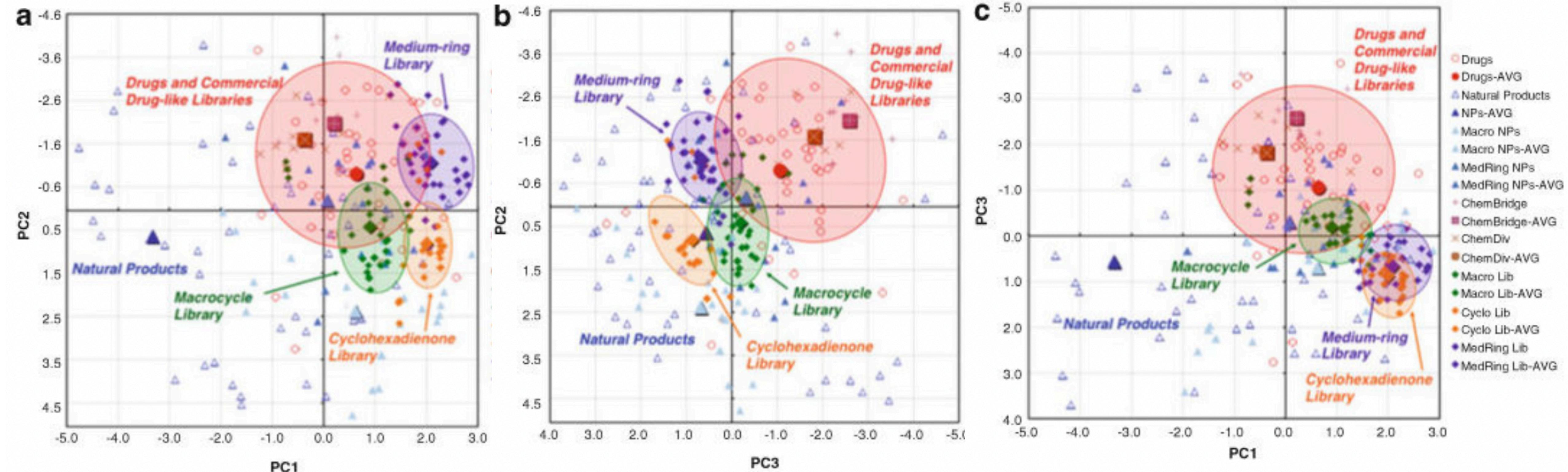
**Fig. 1.**

Recently developed routes to natural product-like macrocycle and medium-ring libraries using ring expansion strategies

PCA to take 20-dimensional molecular featurization down to 3D.

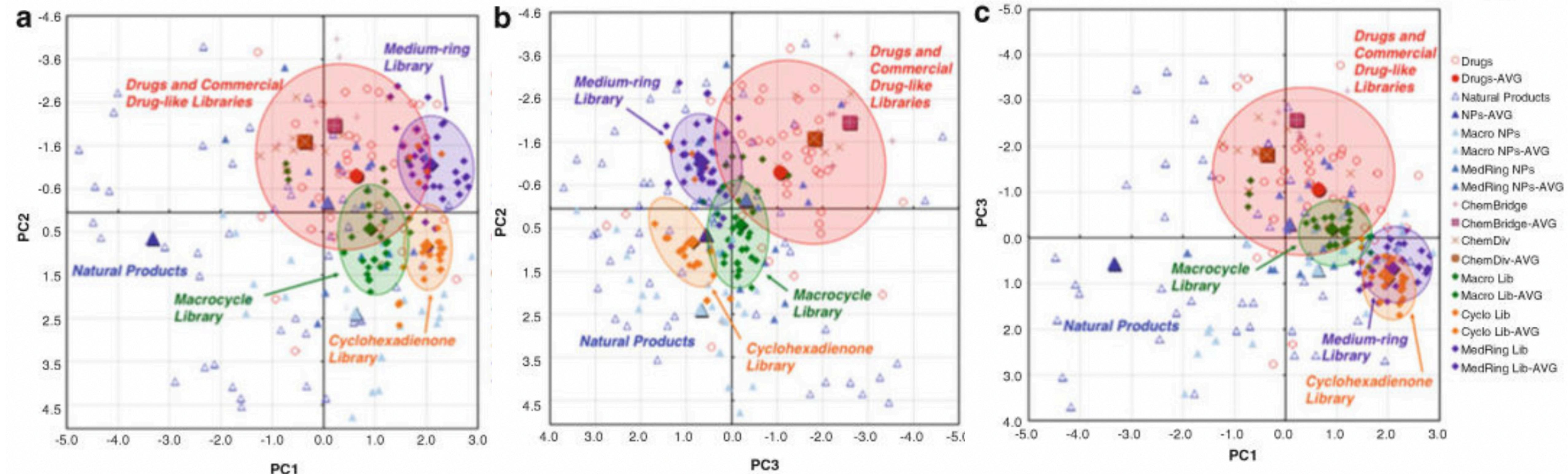
	PC1	PC2	PC3
MW	-0.3207	-0.1482	0.1008
N	-0.2080	-0.1055	-0.3331
O	-0.3180	0.0179	0.1159
HBD	-0.3087	0.0677	-0.1180
HBA	-0.3346	-0.0057	-0.0256
RotB	-0.2587	-0.1209	-0.0540
tPSA	-0.3336	0.0109	-0.0818
nStereo	-0.2789	0.0936	0.2882
nStMW	-0.1612	0.2880	0.2964
Rings	-0.1669	-0.2384	0.0698
RngAr	-0.0475	-0.3660	-0.3227
RngSys	-0.1761	-0.2315	-0.1506
RngLg	-0.1459	-0.0128	0.2521
RRSys	-0.0015	-0.0362	0.1920
ALOGPs	0.0859	-0.4127	0.2677
ALOGpS	0.0378	0.4319	-0.2060
Fsp3	-0.0910	0.3464	0.3224
LogD	0.1850	-0.2569	0.3114
relPSA	-0.1797	0.2262	-0.3112
VWSA	-0.3122	-0.1253	0.1626

# PCA: A Real World Example



Modify library molecules so that they fall into natural product space  
using principal component factors.

# PCA: A Real World Example



Modify library molecules so that they fall into natural product space  
using principal component factors.

Move PC1 left and PC2 down

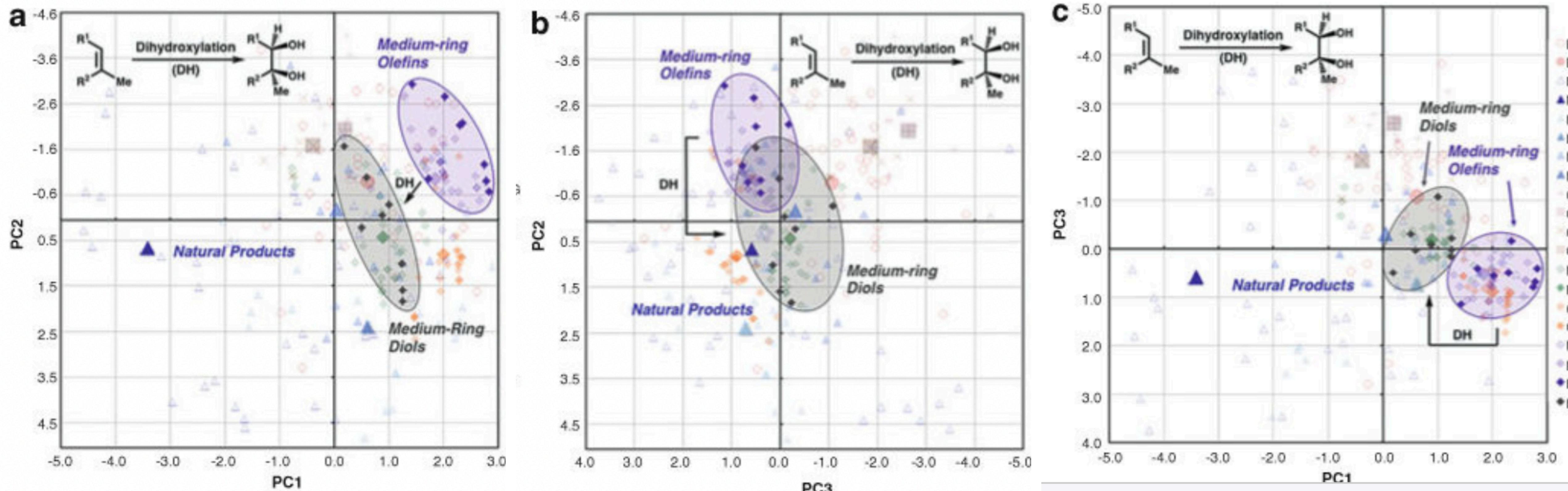
Important (negatively correlated) features for PC1 are molecular weight, number of oxygen atoms, and number of hydrogen bond donors and acceptors.

	PC1	PC2	PC3
MW	-0.3207	-0.1482	0.1008
N	-0.2080	-0.1055	-0.3331
O	-0.3180	0.0179	0.1159
HBD	-0.3087	0.0677	-0.1180
HBA	-0.3346	-0.0057	-0.0256
RotB	-0.2587	-0.1209	-0.0540
tPSA	-0.3336	0.0109	-0.0818
nStereo	-0.2789	0.0936	0.2882
nStMW	-0.1612	0.2880	0.2964
Rings	-0.1669	-0.2384	0.0698
RngAr	-0.0475	-0.3660	-0.3227
RngSys	-0.1761	-0.2315	-0.1506
RngLg	-0.1459	-0.0128	0.2521
RRSys	-0.0015	-0.0362	0.1920
ALOGPs	0.0859	-0.4127	0.2677
ALOGpS	0.0378	0.4319	-0.2060
Fsp3	-0.0910	0.3464	0.3224
LogD	0.1850	-0.2569	0.3114
relPSA	-0.1797	0.2262	-0.3112
VWSA	-0.3122	-0.1253	0.1626

Important features for PC2 are density of stereocenters and the number of sp<sup>3</sup> bonds.

	PC1	PC2	PC3
MW	-0.3207	-0.1482	0.1008
N	-0.2080	-0.1055	-0.3331
O	-0.3180	0.0179	0.1159
HBD	-0.3087	0.0677	-0.1180
HBA	-0.3346	-0.0057	-0.0256
RotB	-0.2587	-0.1209	-0.0540
tPSA	-0.3336	0.0109	-0.0818
nStereo	-0.2789	0.0936	0.2882
nStMW	-0.1612	0.2880	0.2964
Rings	-0.1669	-0.2384	0.0698
RngAr	-0.0475	-0.3660	-0.3227
RngSys	-0.1761	-0.2315	-0.1506
RngLg	-0.1459	-0.0128	0.2521
RRSys	-0.0015	-0.0362	0.1920
ALOGPs	0.0859	-0.4127	0.2677
ALOGpS	0.0378	0.4319	-0.2060
Fsp3	-0.0910	0.3464	0.3224
LogD	0.1850	-0.2569	0.3114
relPSA	-0.1797	0.2262	-0.3112
VWSA	-0.3122	-0.1253	0.1626

# PCA: A Real World Example



Answers typically a “yes/no” question.

Is this a cat?

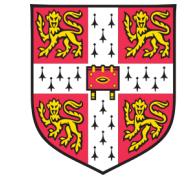


Answers typically a “yes/no” question.

Is this a cat?



Yes, this is a **cat**



Answers typically a “yes/no” question.

Is this a cat?

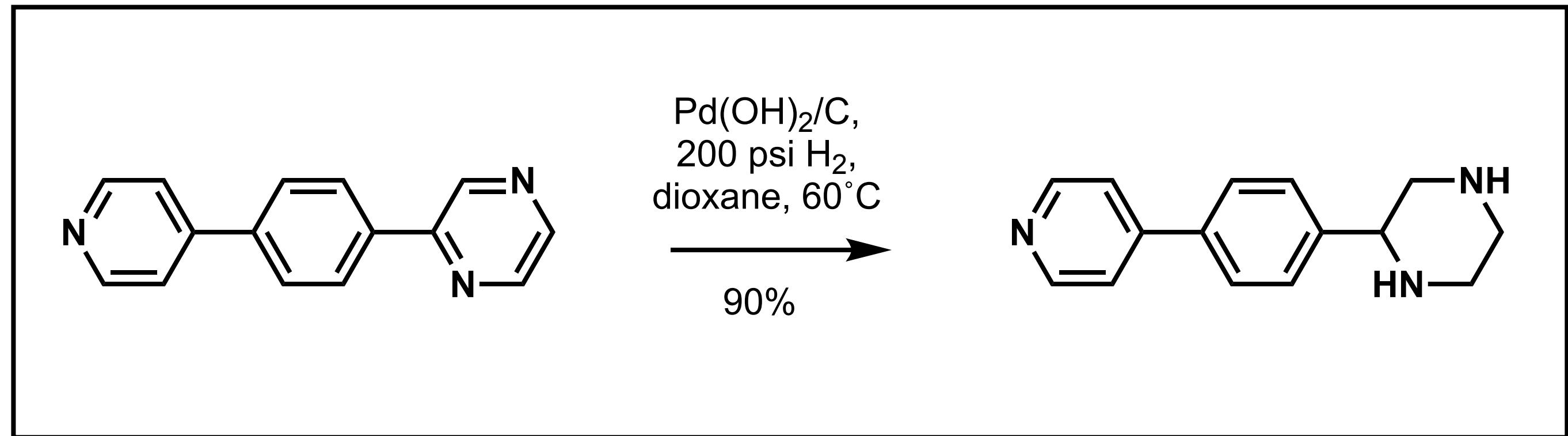


No, this is a **not a cat**

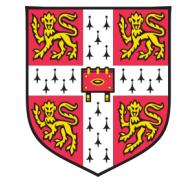


Answers typically a “yes/no” question.

Is this reaction high yielding?

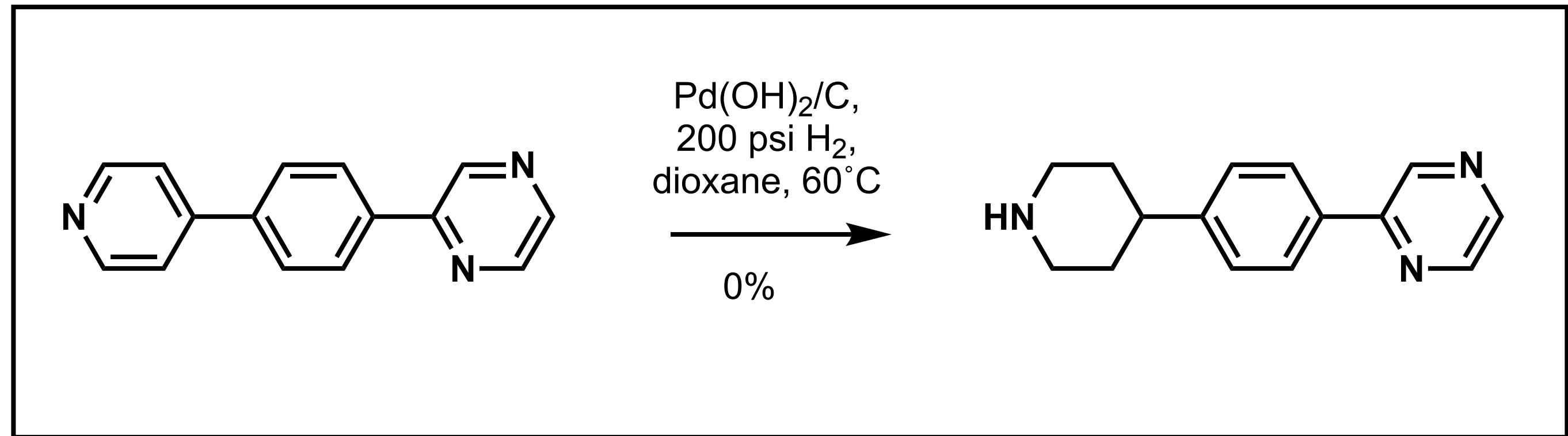


Yes, is a **high yielding rxn**



Answers typically a “yes/no” question.

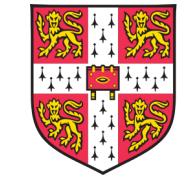
Is this reaction high yielding?



No, this is **not a high yielding rxn**

Can also assign labels from a predetermined list

What is this? A cat? A dog? An owl?



Can also assign labels from a predetermined list

What is this? A cat? A dog? An owl?



**This is a cat**



Can also assign labels from a predetermined list

What is this? A cat? A dog? An owl?



**This is a dog**



Can also assign labels from a predetermined list

What is this? A cat? A dog? An owl?



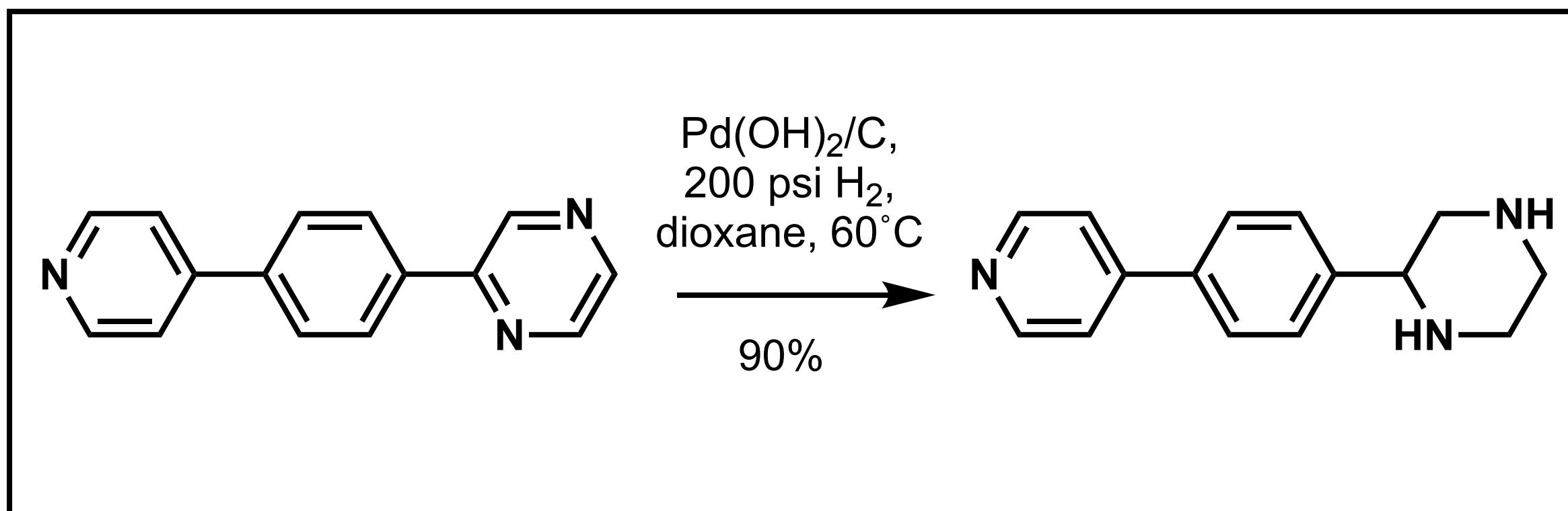
This is **an owl**

True Positive (TP) = Correctly identified as **X**

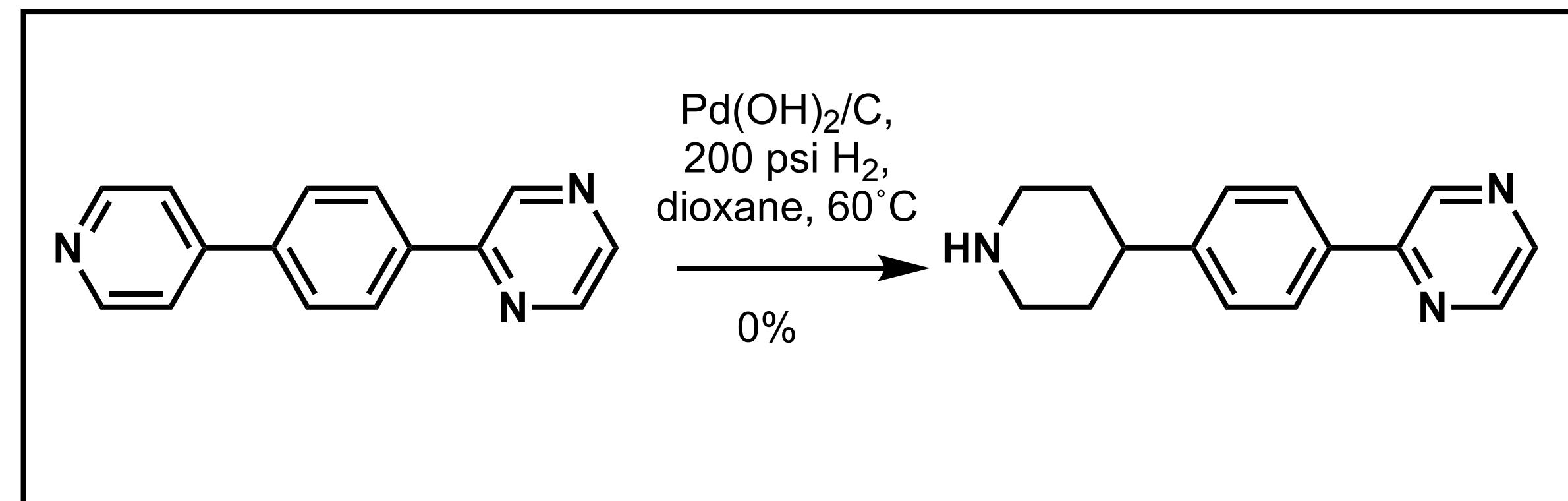
True Negative (TN) = Correctly identified as **not X**

False Positive (FP) = Incorrectly identified as **X**

False Negative (FN) = Incorrectly identified as **not X**

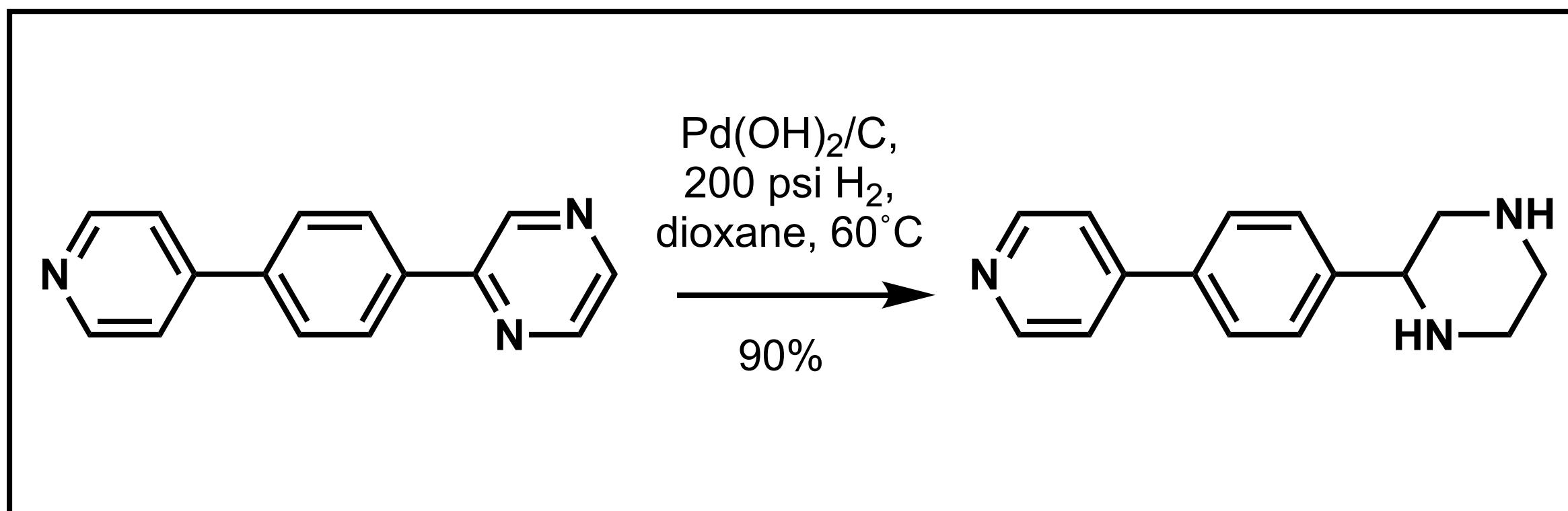


Rxn A

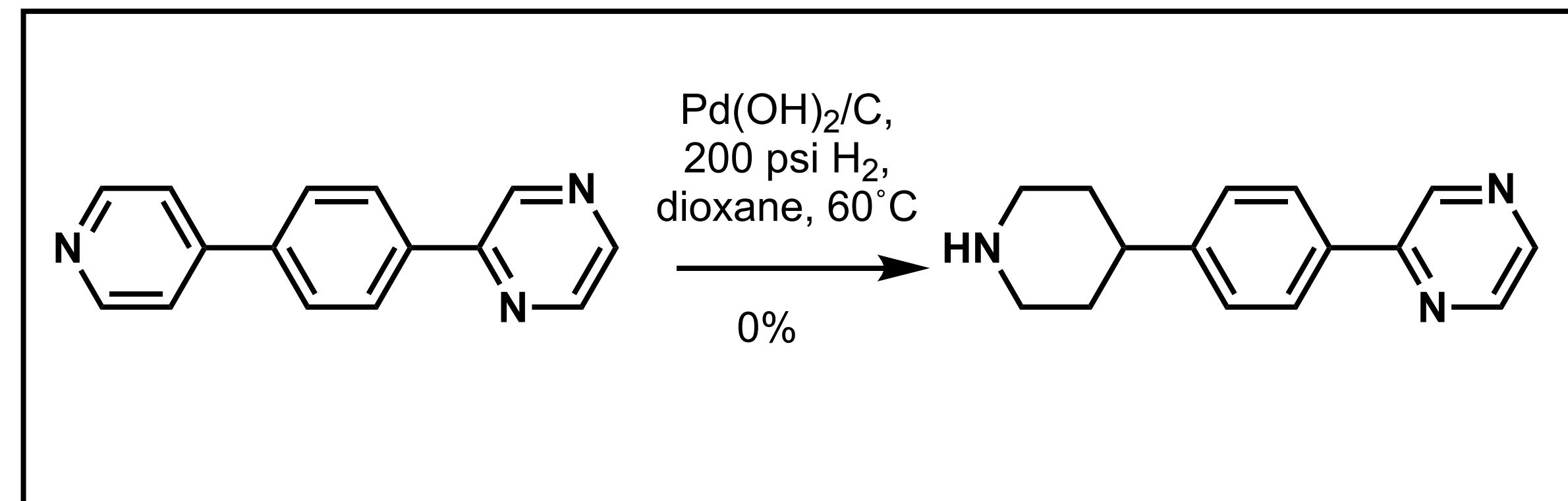


Rxn B

Classification	Truth	Result

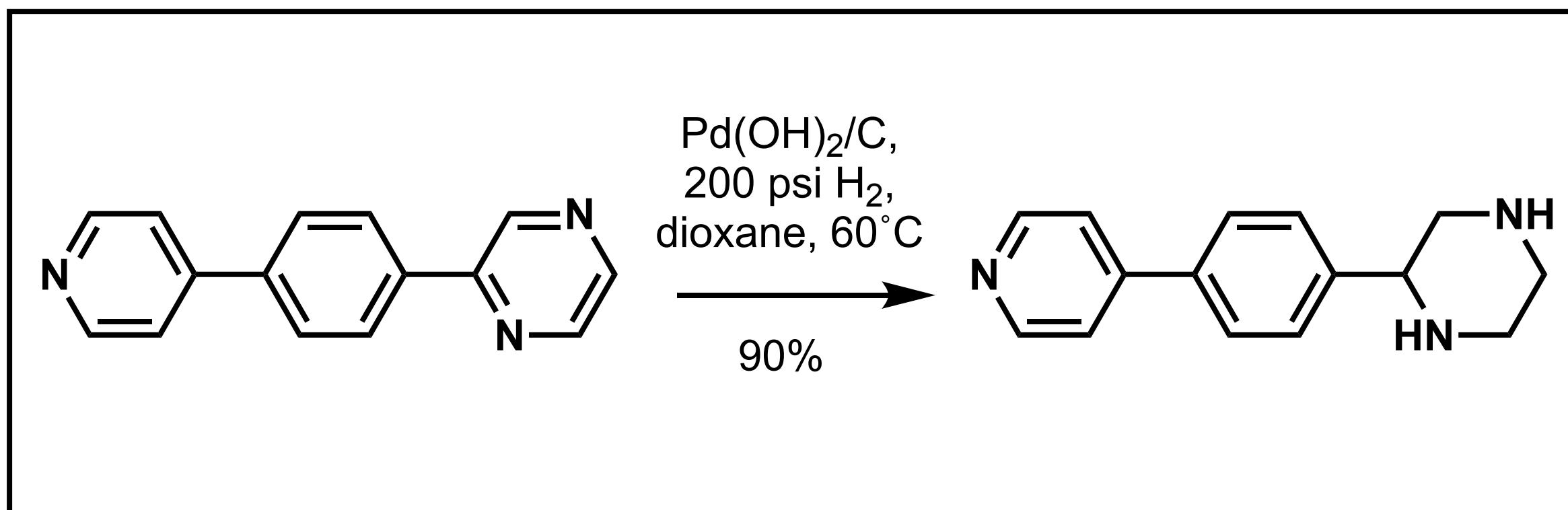


Rxn A

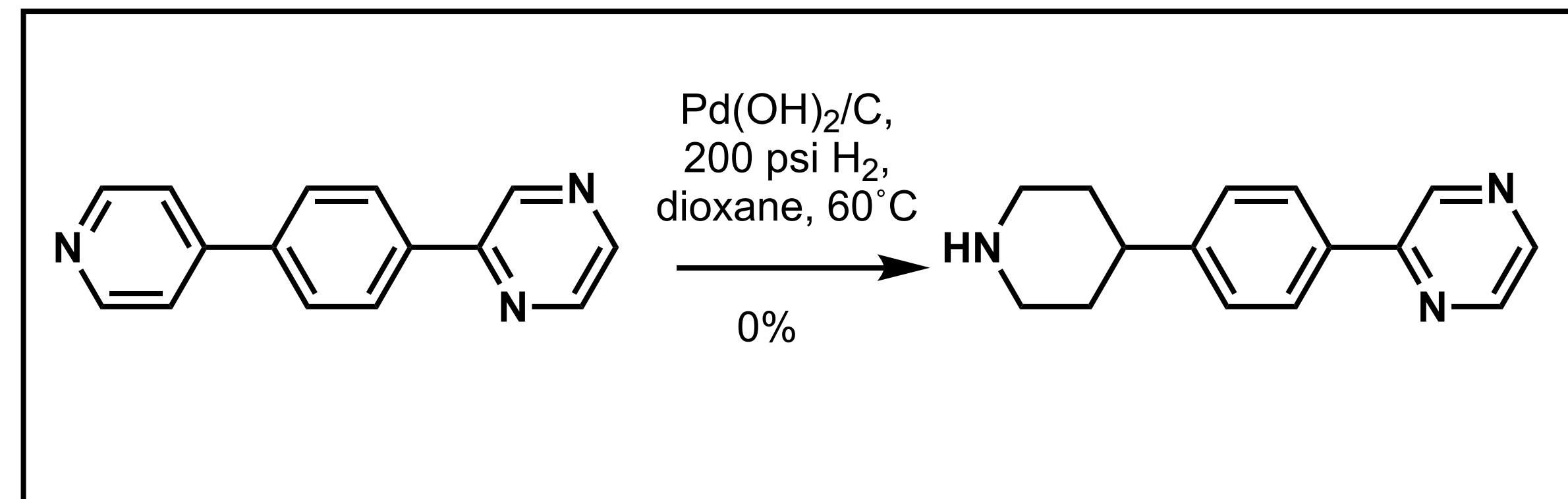


Rxn B

Classification	Truth	Result
	High Yielding	
	High Yielding	

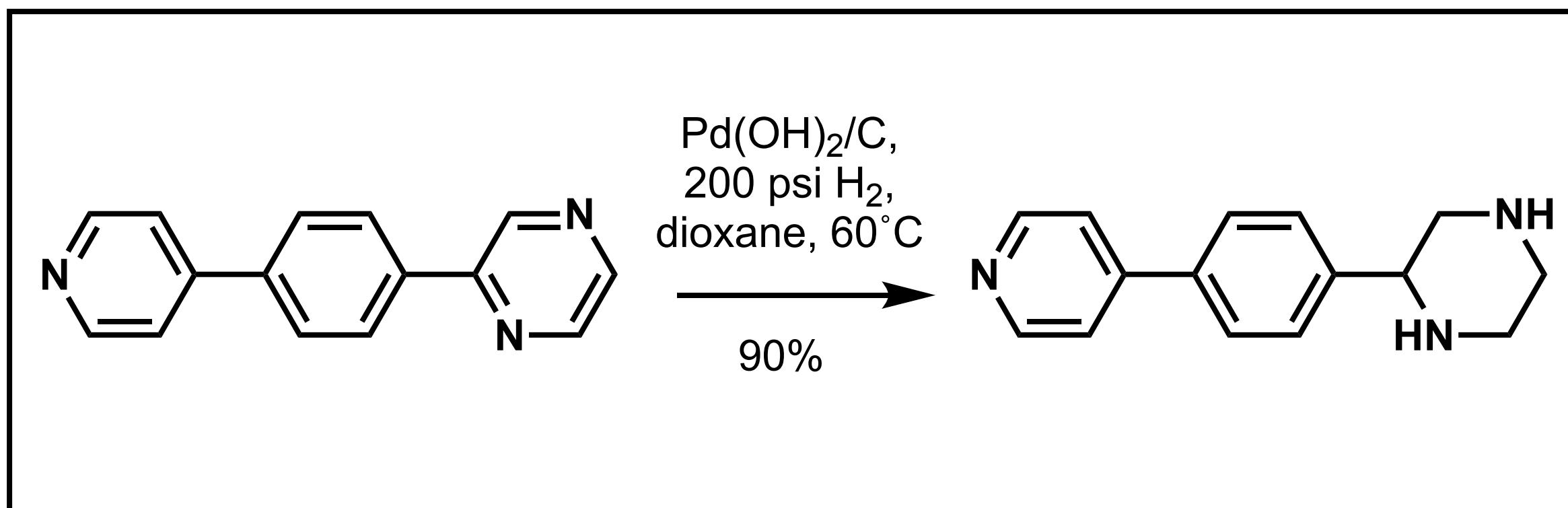


Rxn A

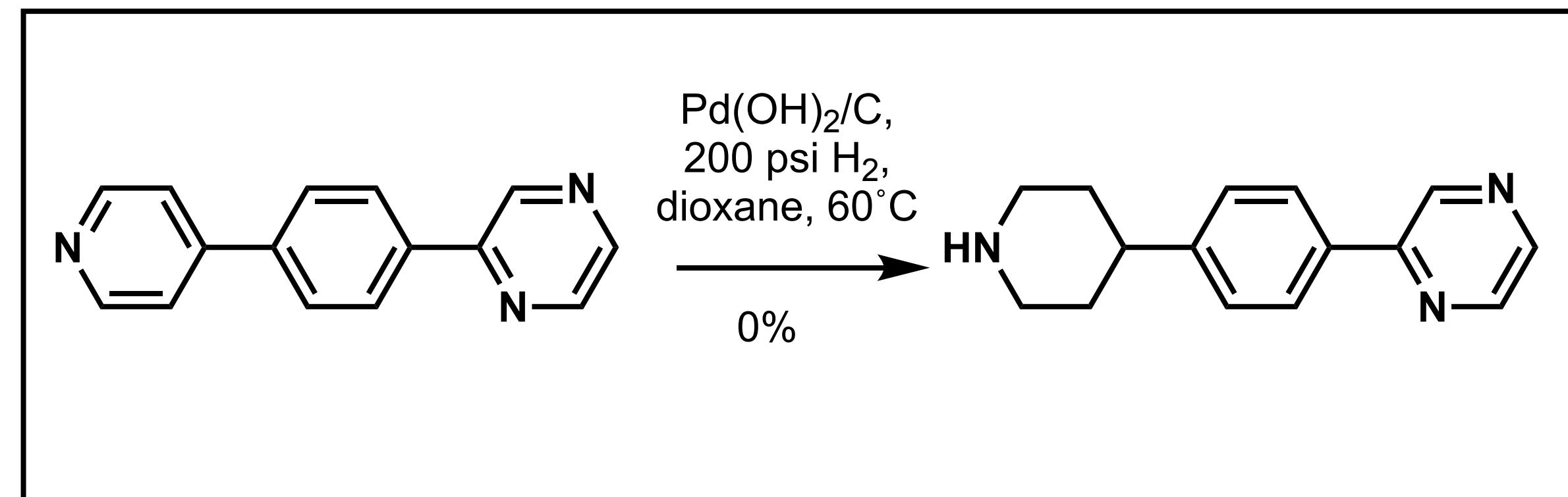


Rxn B

Classification	Truth	Result
High Yielding	High Yielding	True Positive
	High Yielding	

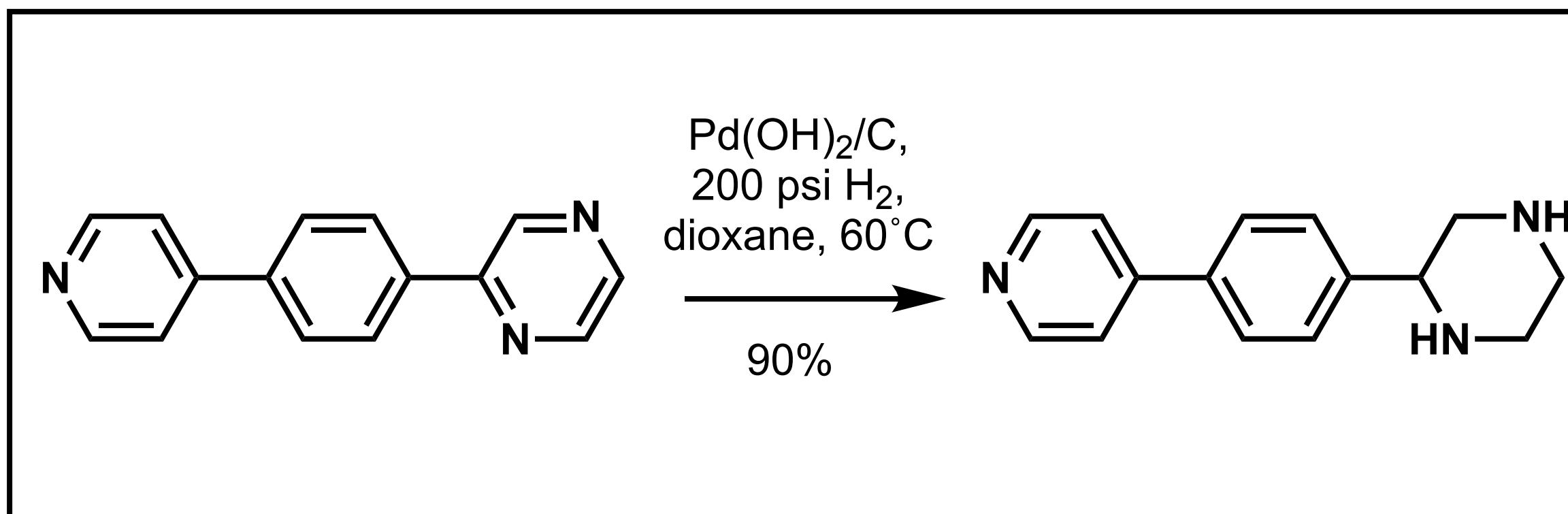


Rxn A

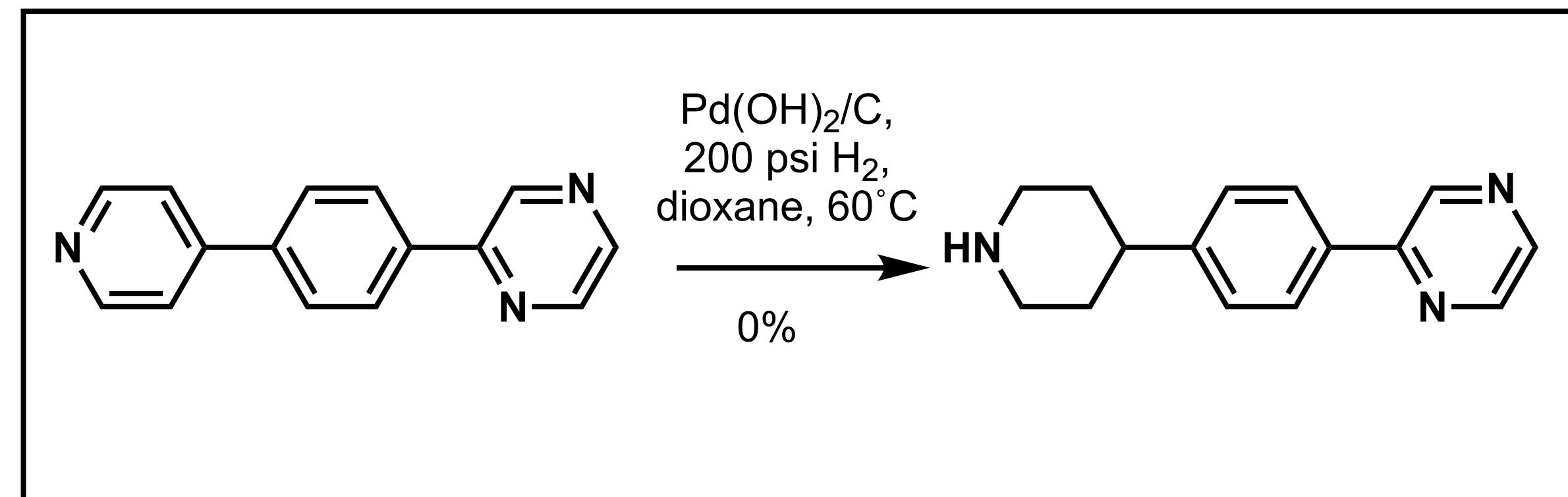


Rxn B

Classification	Truth	Result
High Yielding	High Yielding	True Positive
Low Yielding	High Yielding	False Negative



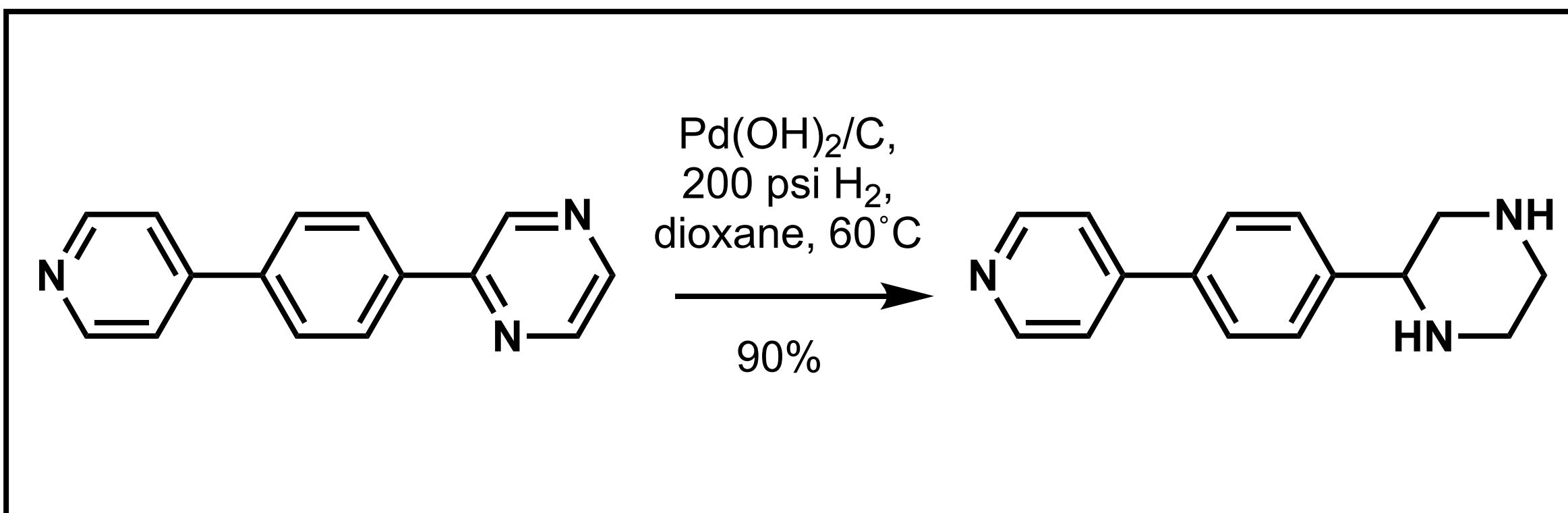
Rxn A



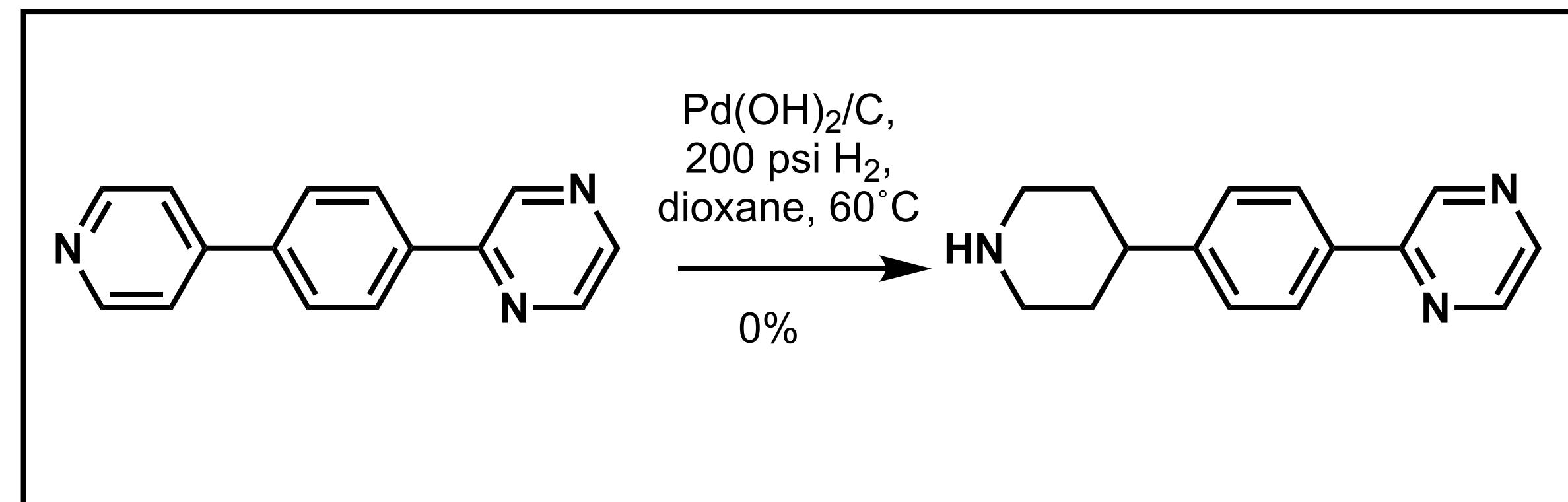
Rxn B

Classification	Truth	Result
High Yielding	High Yielding	True Positive
Low Yielding	High Yielding	False Negative

Classification	Truth	Result
	Low Yielding	
	Low Yielding	



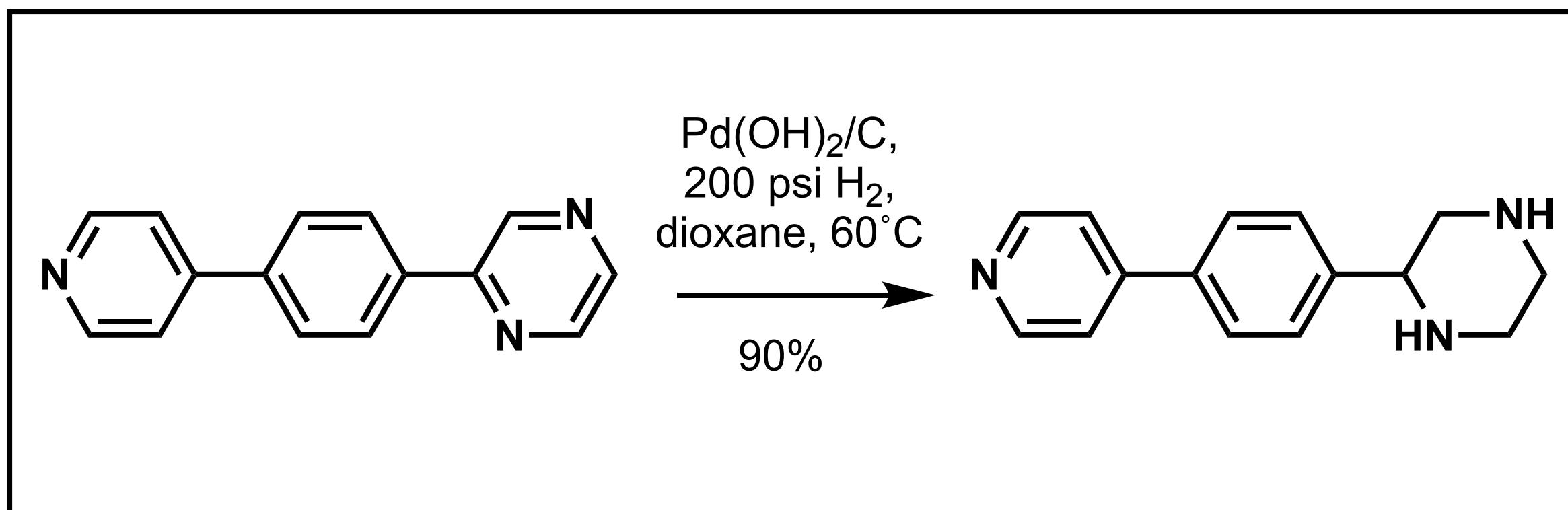
Rxn A



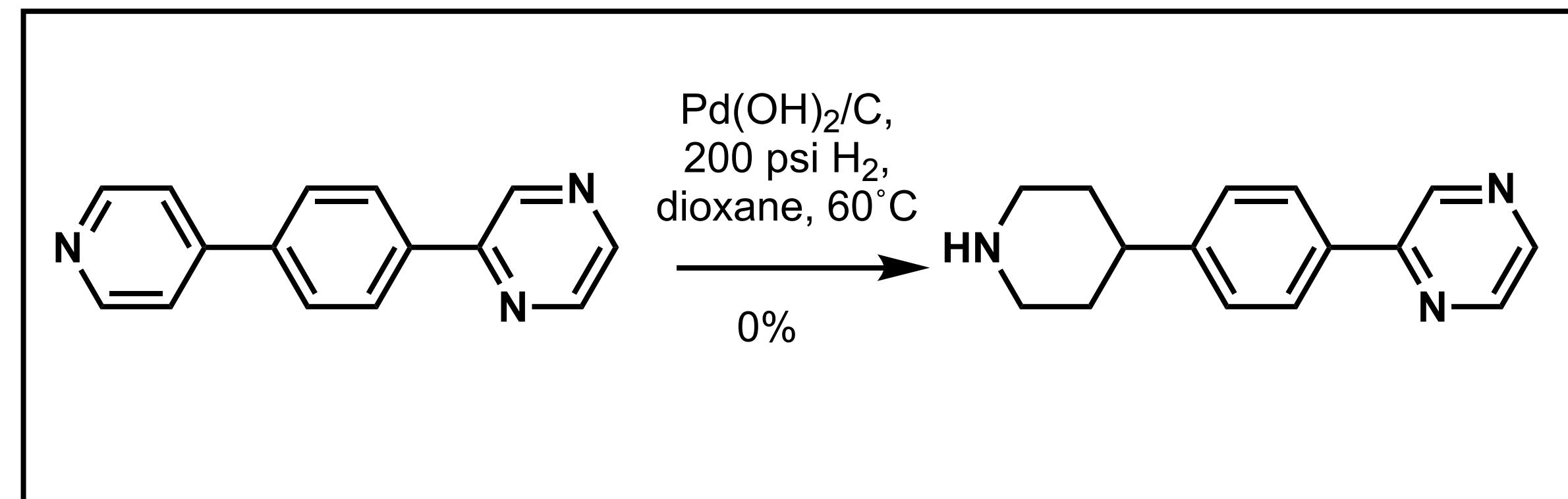
Rxn B

Classification	Truth	Result
High Yielding	High Yielding	True Positive
Low Yielding	High Yielding	False Negative

Classification	Truth	Result
High Yielding	Low Yielding	False Positive
Low Yielding	Low Yielding	



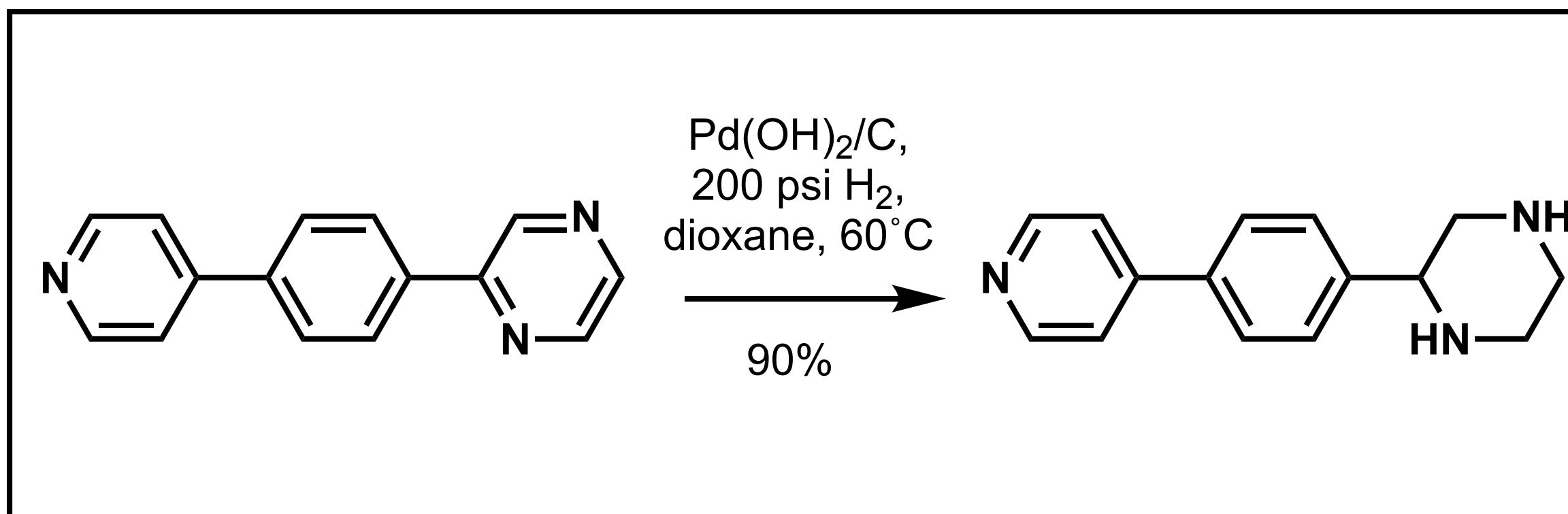
Rxn A



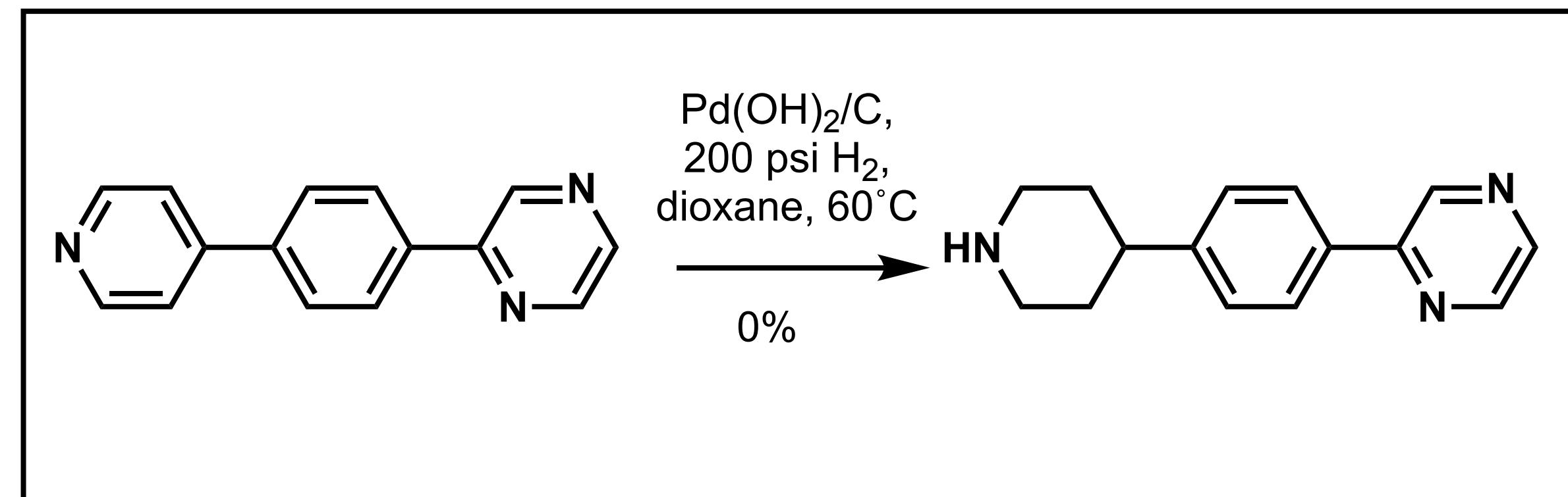
Rxn B

Classification	Truth	Result
High Yielding	High Yielding	True Positive
Low Yielding	High Yielding	False Negative

Classification	Truth	Result
High Yielding	Low Yielding	False Positive
Low Yielding	Low Yielding	True Negative



Rxn A



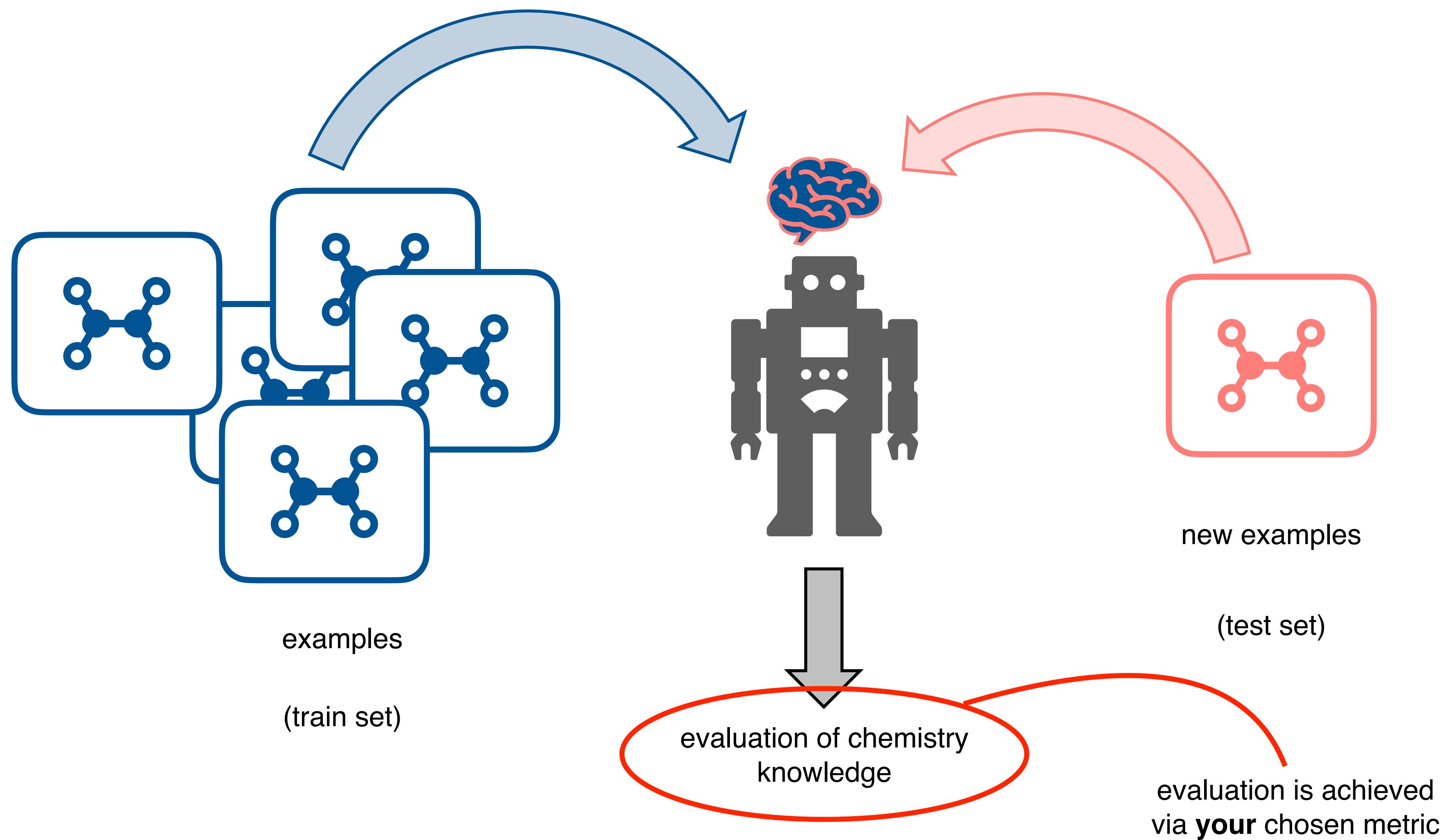
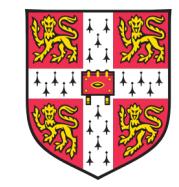
Rxn B

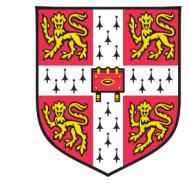
**Confusion Matrix**

What actually happened (ground truth)

		What we predicted (predictions)	
		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

# Why Are Metrics **SO** Important? Recall Our Robot



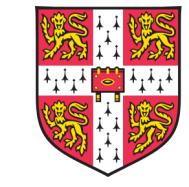


## Confusion Matrix

What actually happened (ground truth)

What we predicted (predictions)	What actually happened (ground truth)	
	High Yielding	Not High Yielding
High Yielding	True Positive	False Positive
Not High Yielding	False Negative	True Negative

High accuracy model = often gets it right



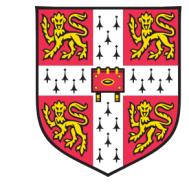
## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

High accuracy model = often gets it right

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$



## Confusion Matrix

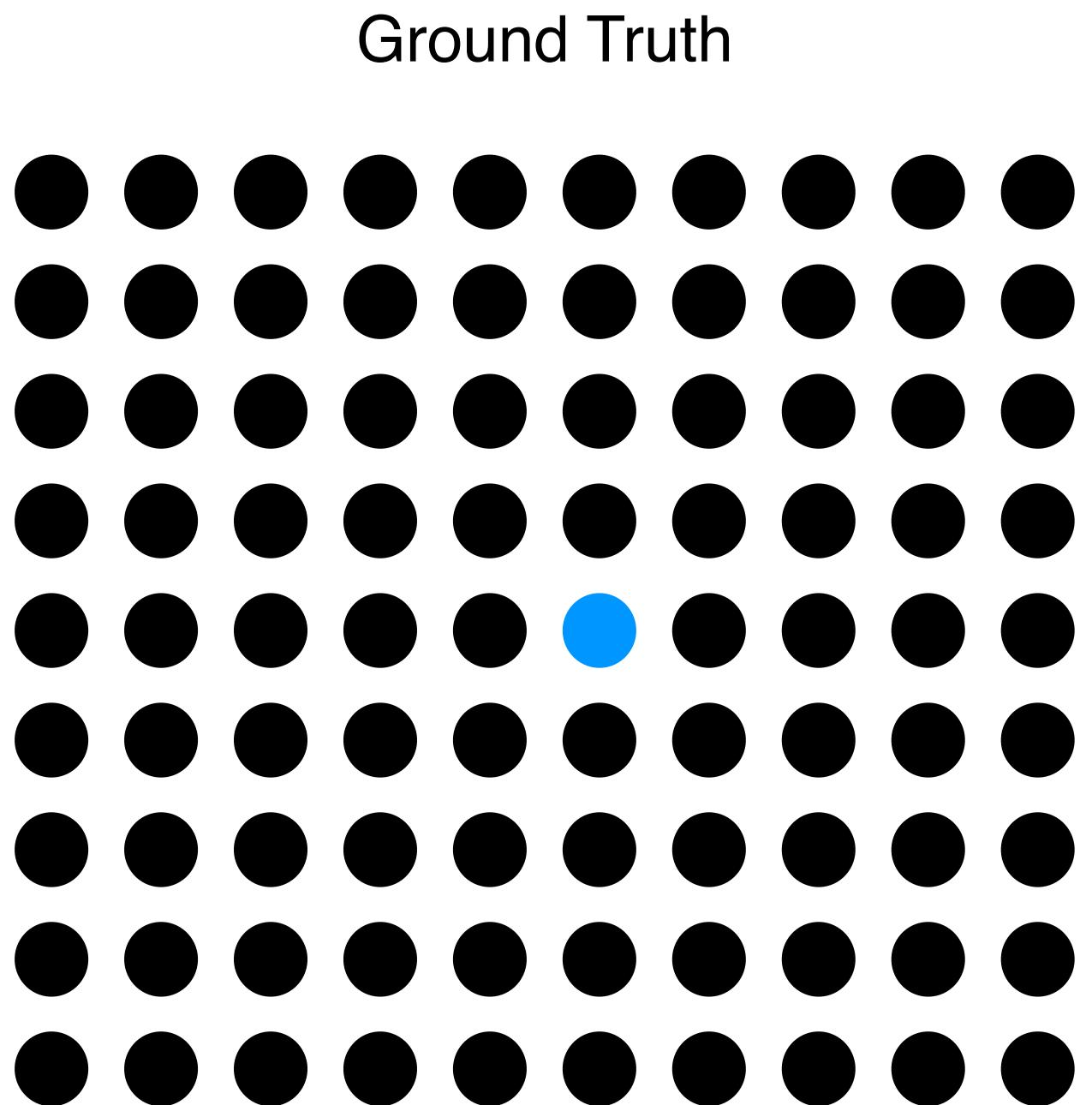
What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

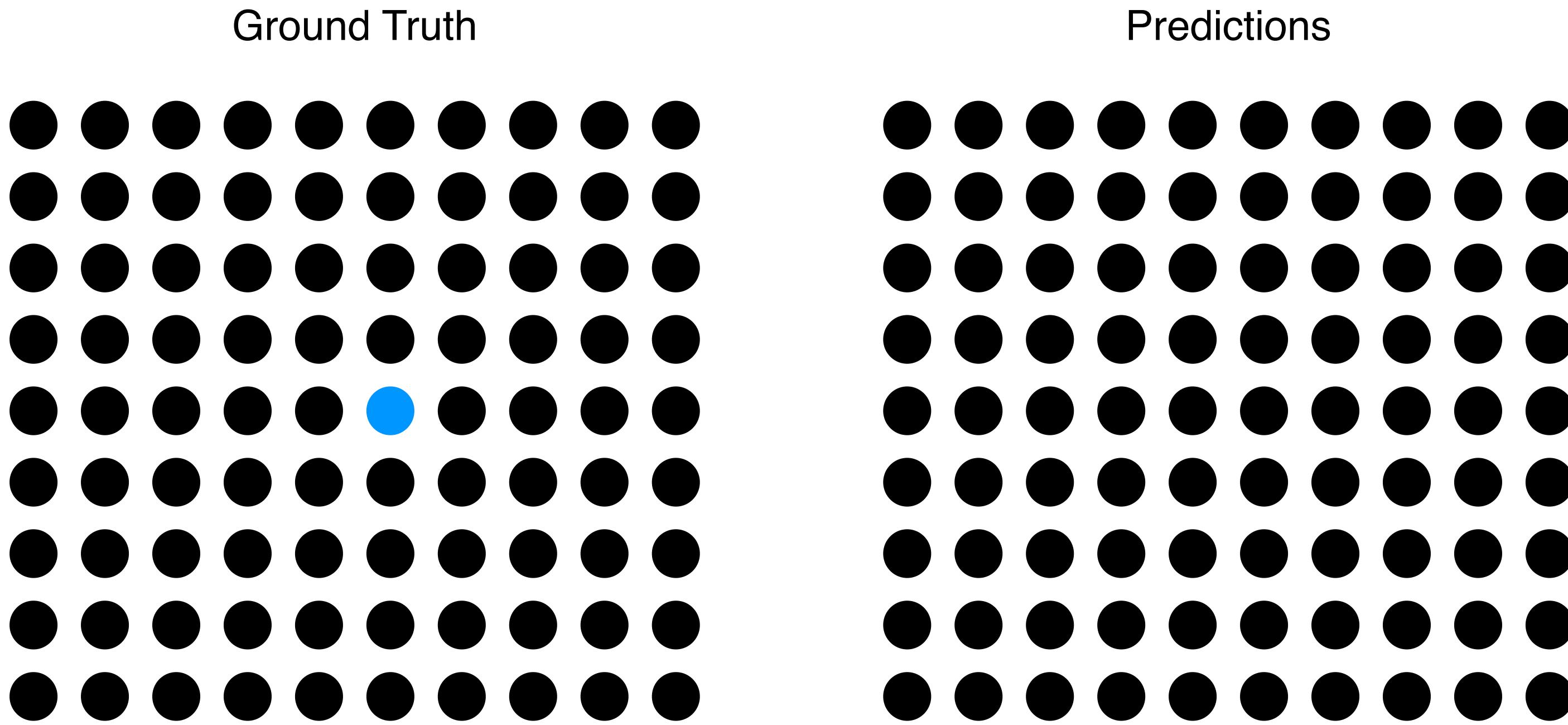
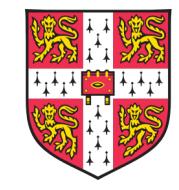
High accuracy model = often gets it right

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Not always a good metric in unbalanced data

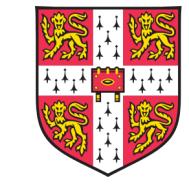


$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$



$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

**Accuracy = 99%**

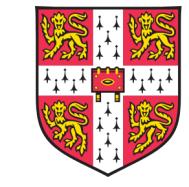


## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

High precision model = when it does predict a TP, it gets it right.



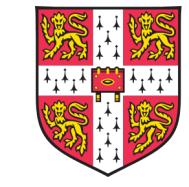
## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

High precision model = when it does predict a TP, it gets it right.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$



## Confusion Matrix

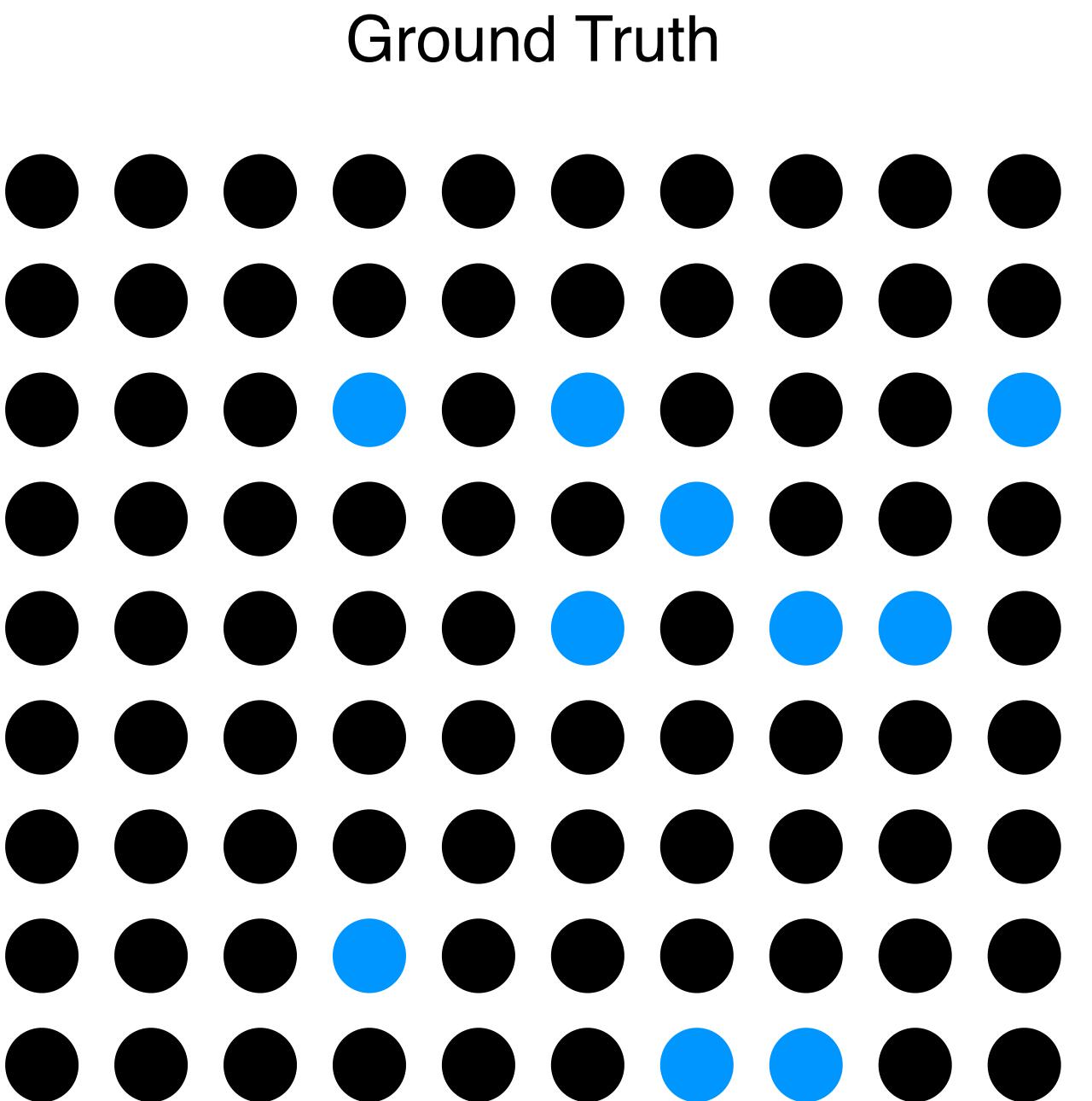
What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

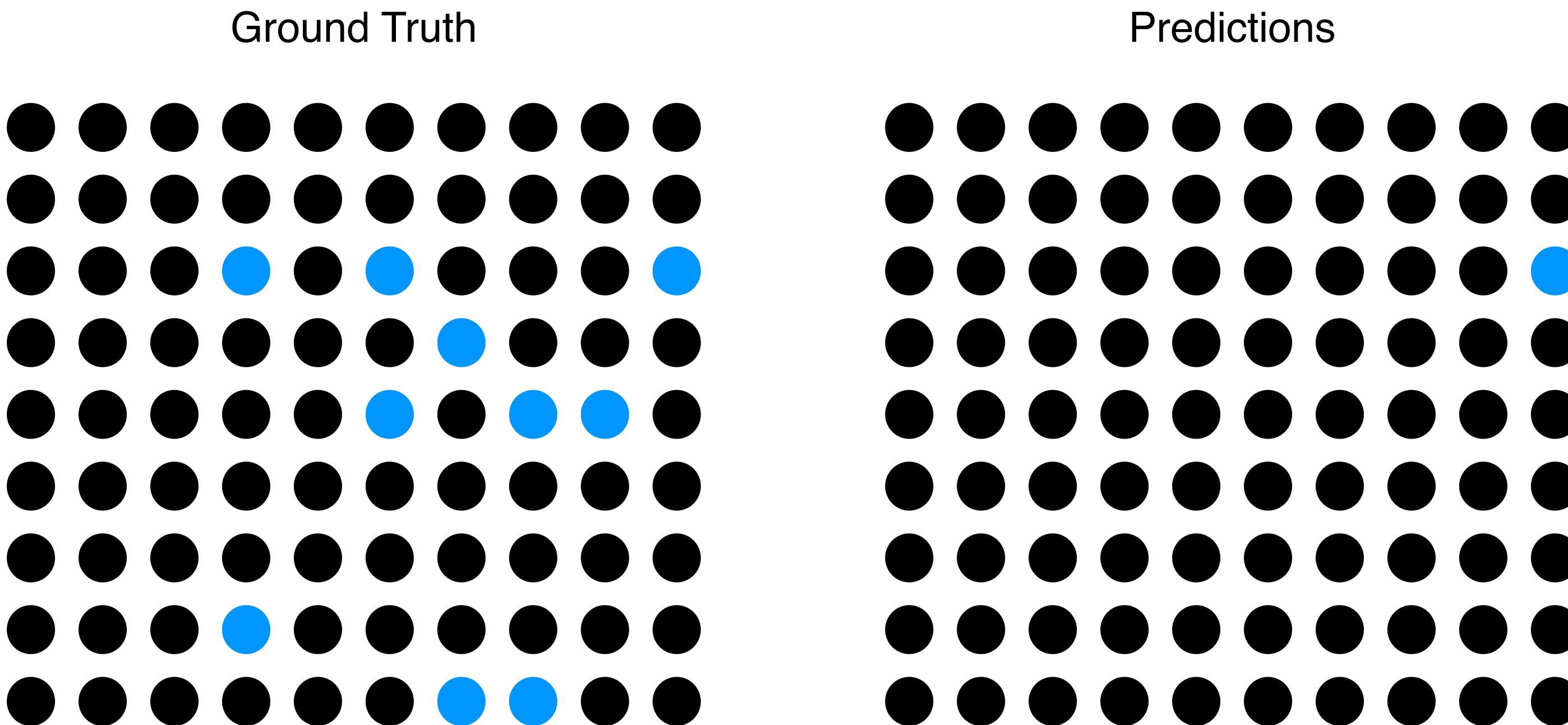
High precision model = when it does predict a TP, it gets it right.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Precision says nothing about how good the model is with finding all the positives.

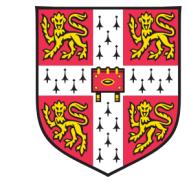


$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

**Precision = 100%**

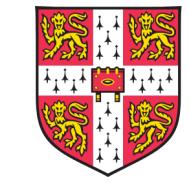


### Confusion Matrix

What actually happened (ground truth)

What we predicted (predictions)	What actually happened (ground truth)	
	High Yielding	Not High Yielding
High Yielding	True Positive	False Positive
Not High Yielding	False Negative	True Negative

High recall model = finds most of the TP.



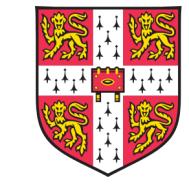
### Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

High recall model = finds most of the TP.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$



## Confusion Matrix

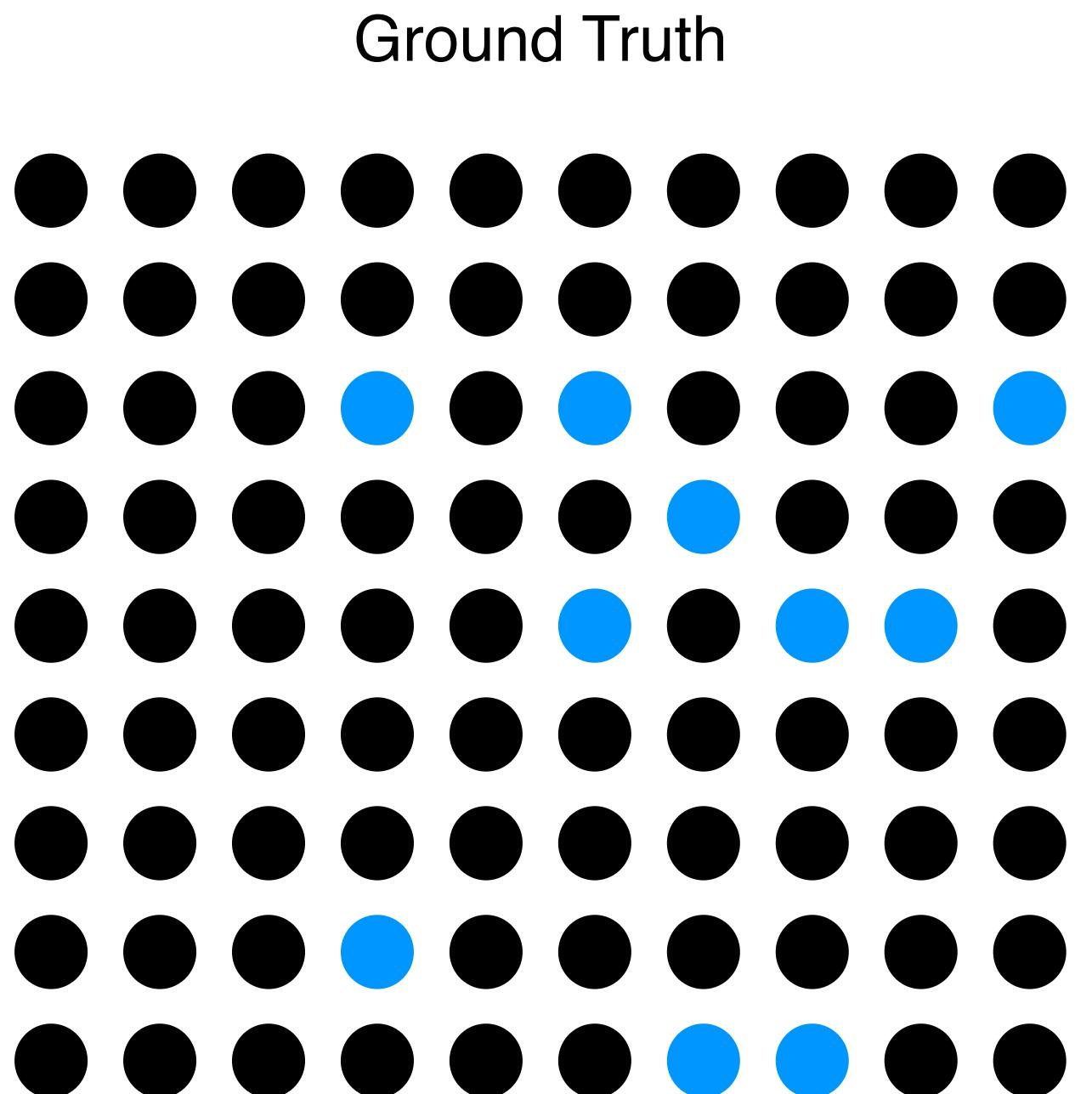
What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

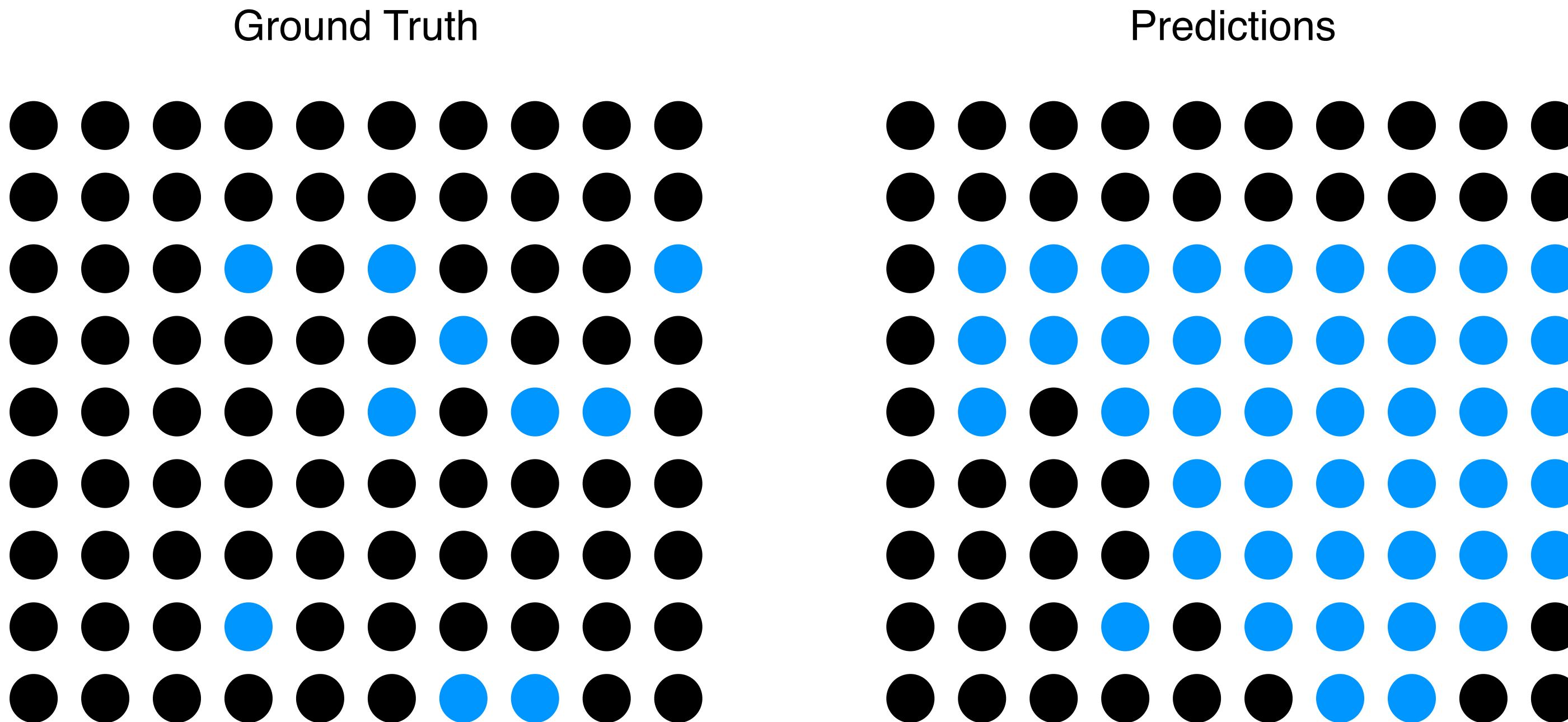
High recall model = finds most of the TP.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Recall says nothing about how good the model is at distinguishing between true positives and false positives.



$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$



$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

**Recall = 100%**

## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

It is up to **you** if you'd like to prioritize high precision or high recall.

## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

It is up to **you** if you'd like to prioritize high precision or high recall.

Oftentimes, we'd like both.

## Confusion Matrix

What actually happened (ground truth)

		High Yielding	Not High Yielding
What we predicted (predictions)	High Yielding	True Positive	False Positive
	Not High Yielding	False Negative	True Negative

It is up to **you** if you'd like to prioritize high precision or high recall.

Oftentimes, we'd like both.

$$\text{F-Score} = \frac{2 \cdot \text{TP}}{\text{TP} + \text{TN} + \text{FP}}$$

Less about visualization, more about prediction

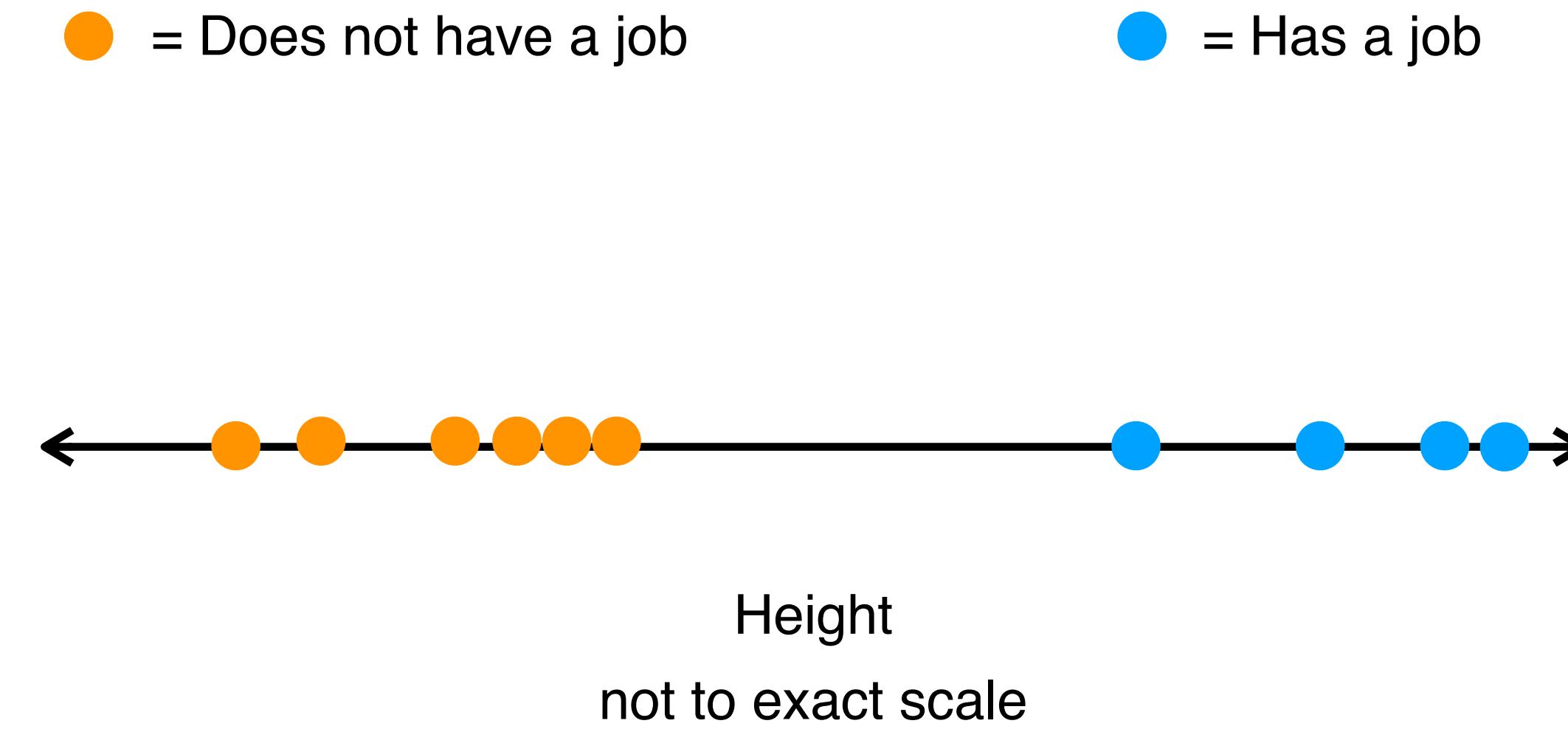
Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes

Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes

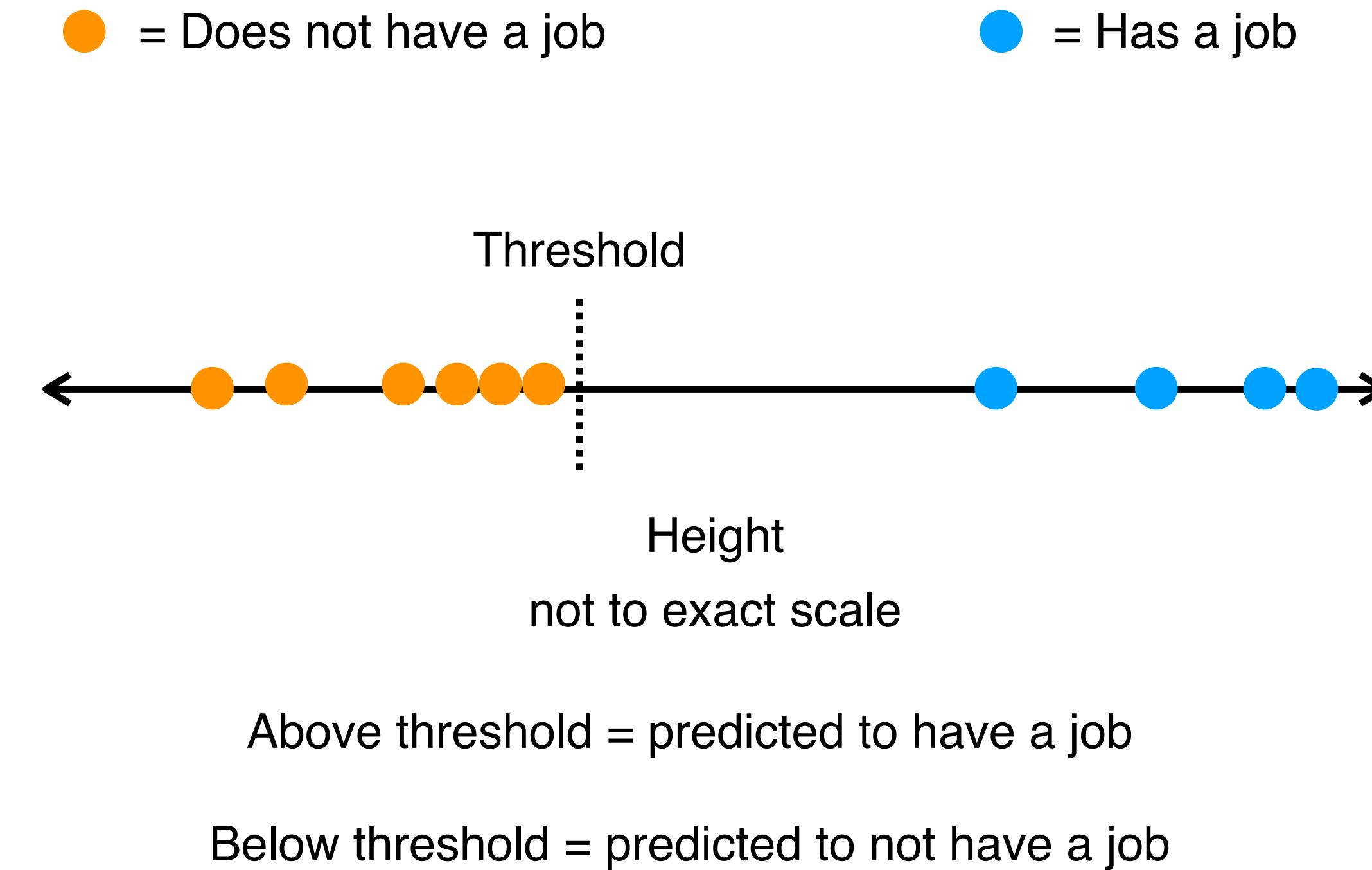
Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



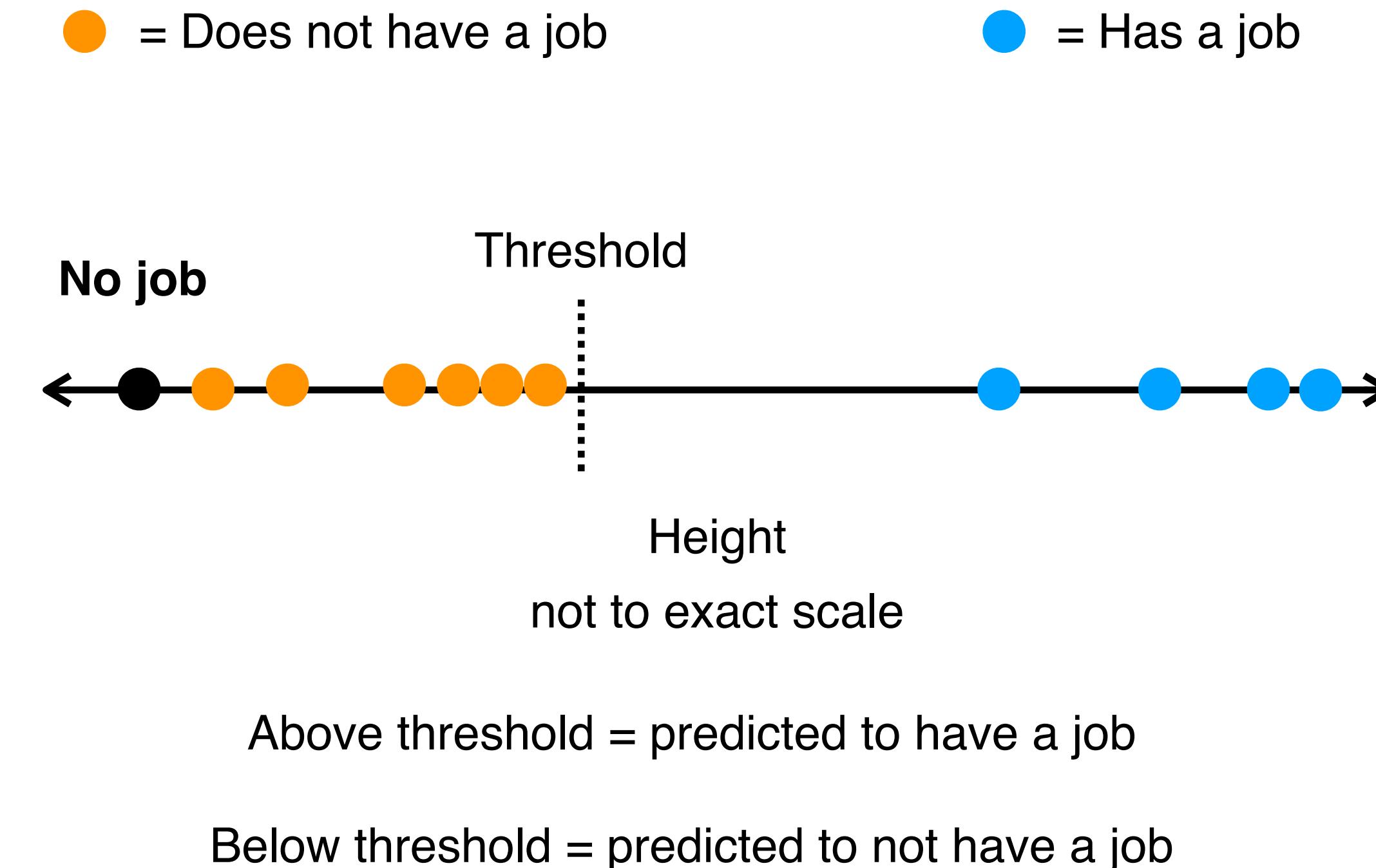
Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



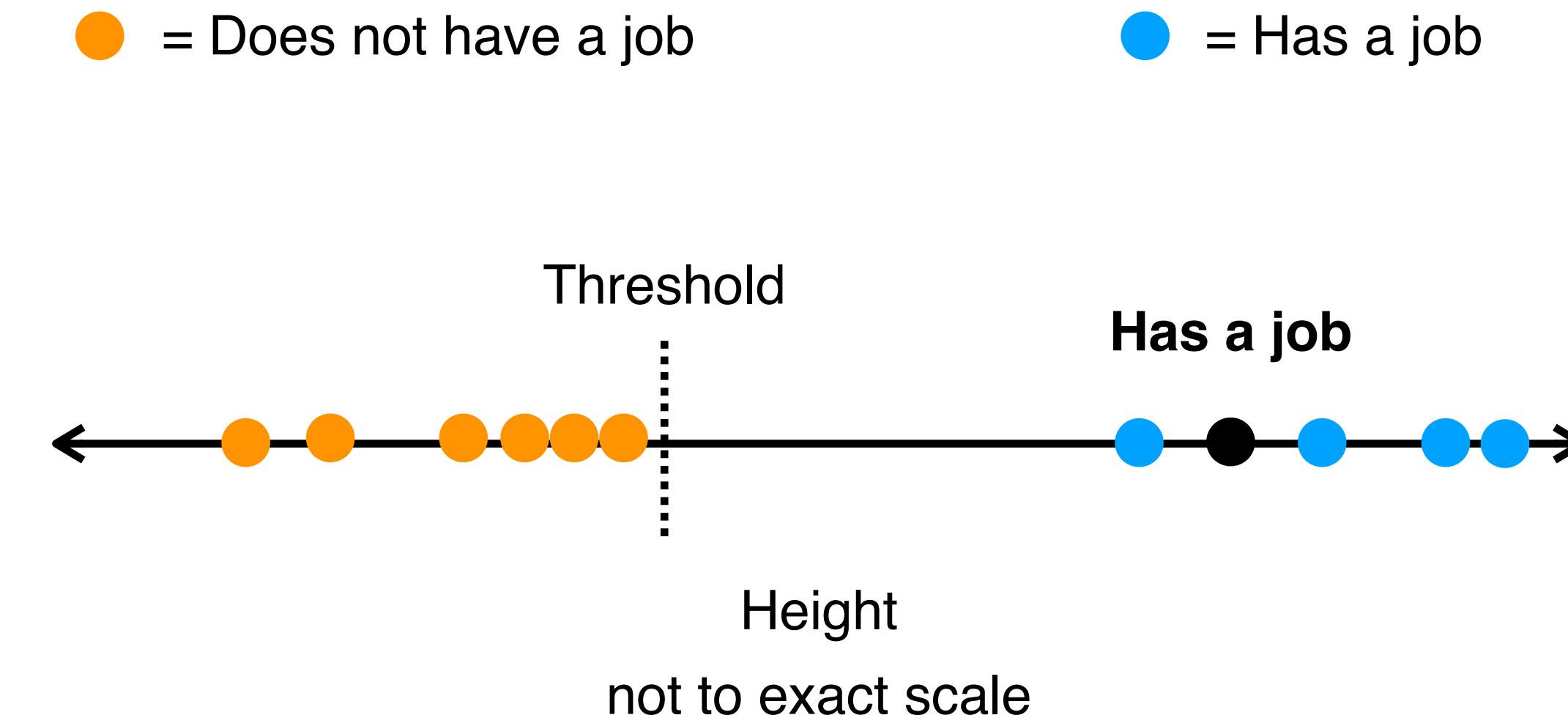
Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes

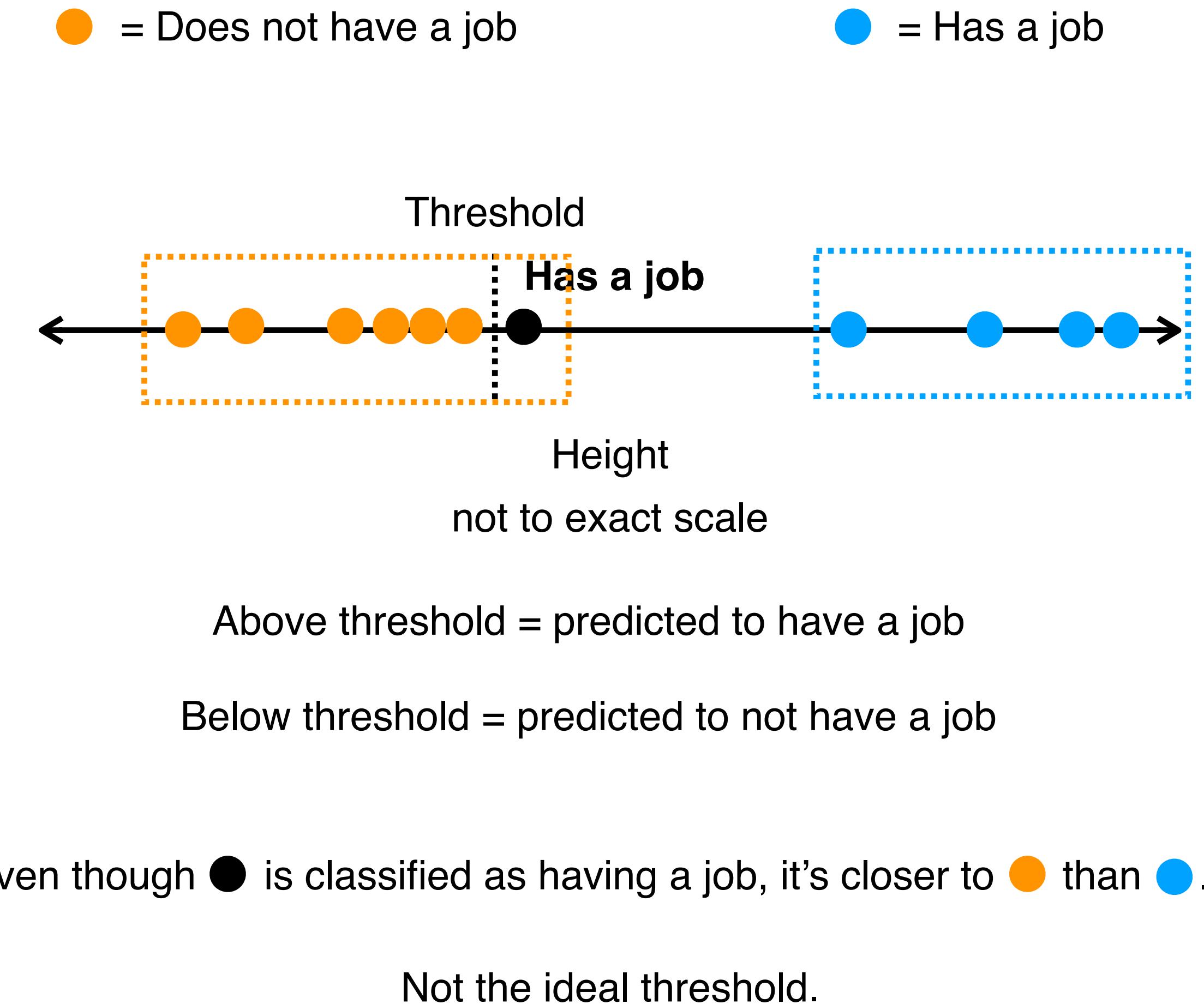


Above threshold = predicted to have a job

Below threshold = predicted to not have a job

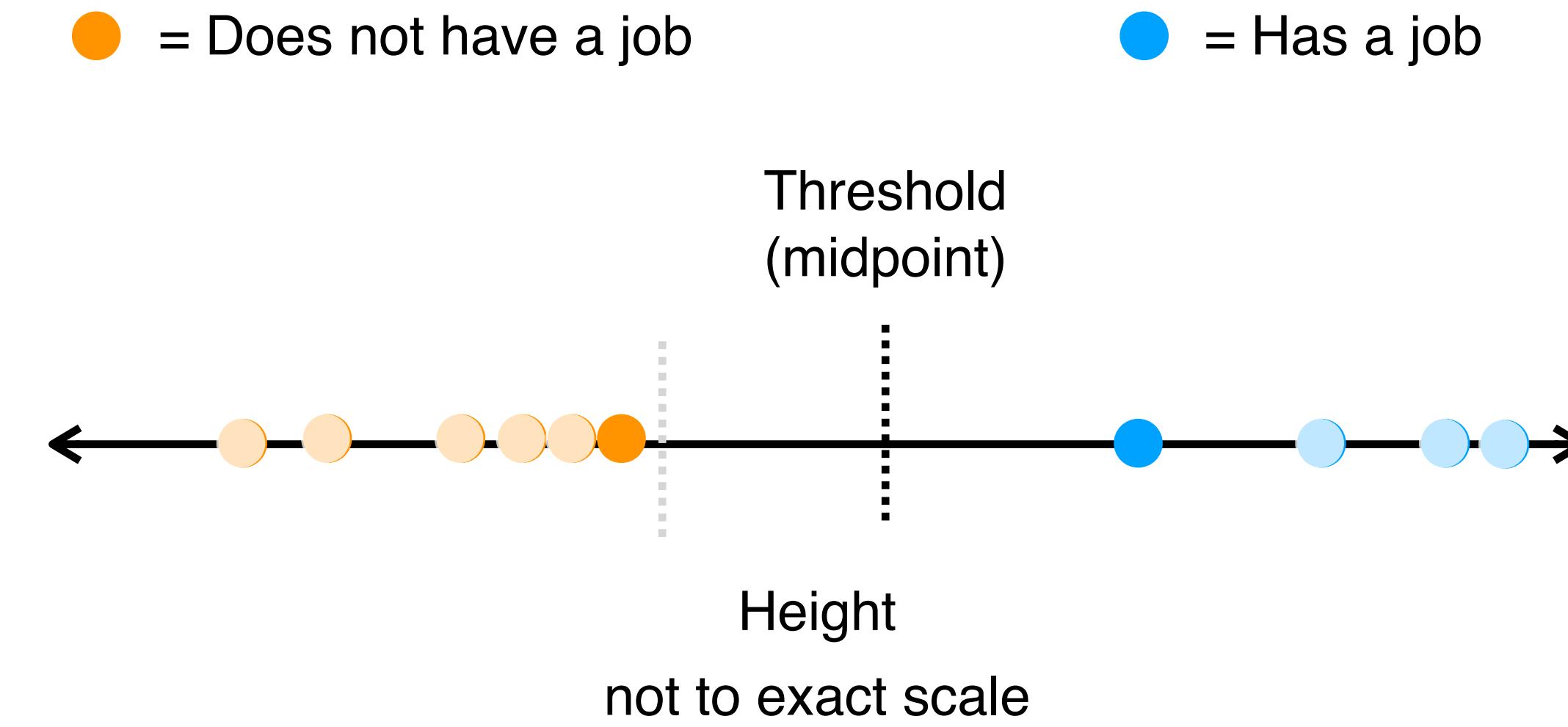
Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes

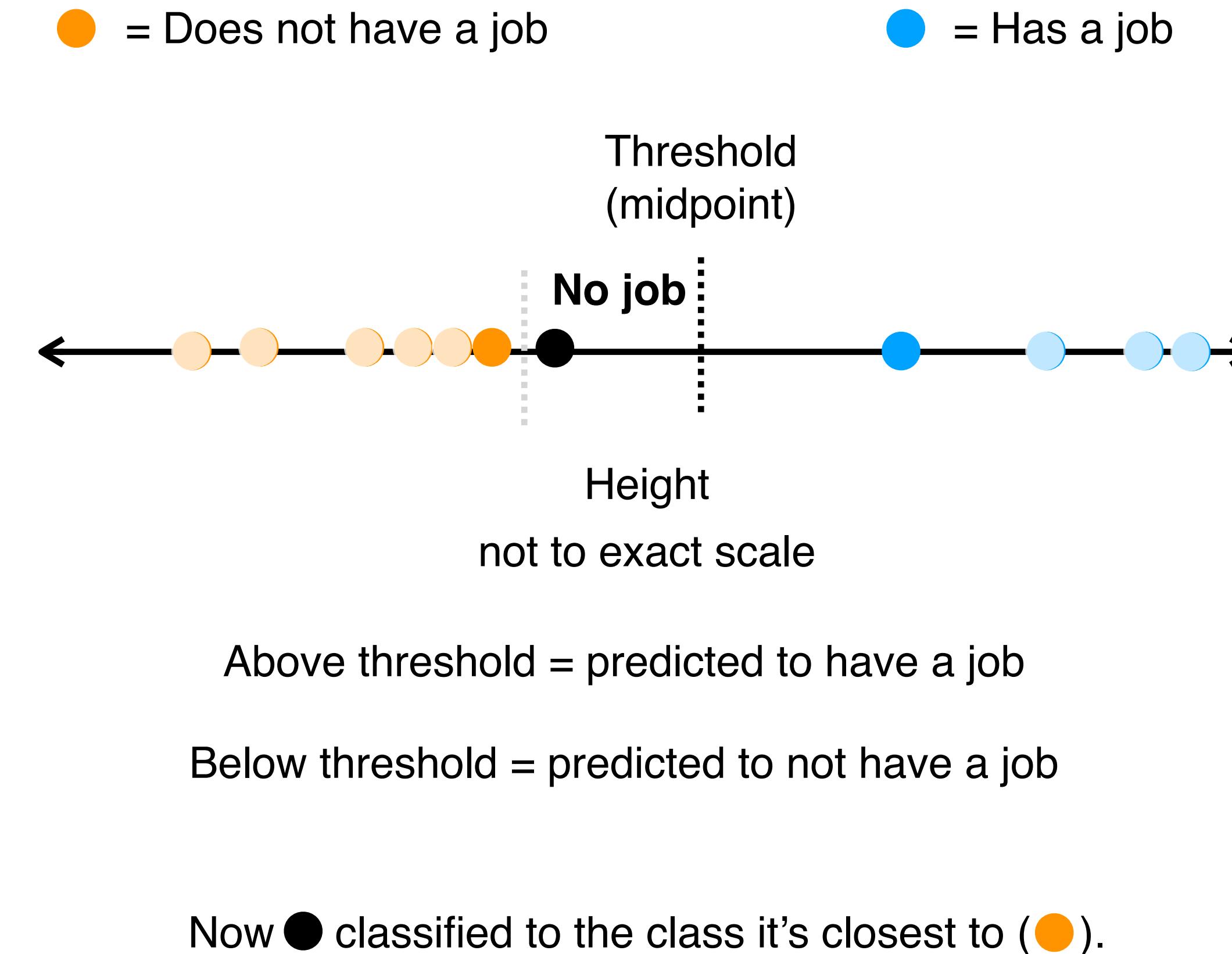


Above threshold = predicted to have a job

Below threshold = predicted to not have a job

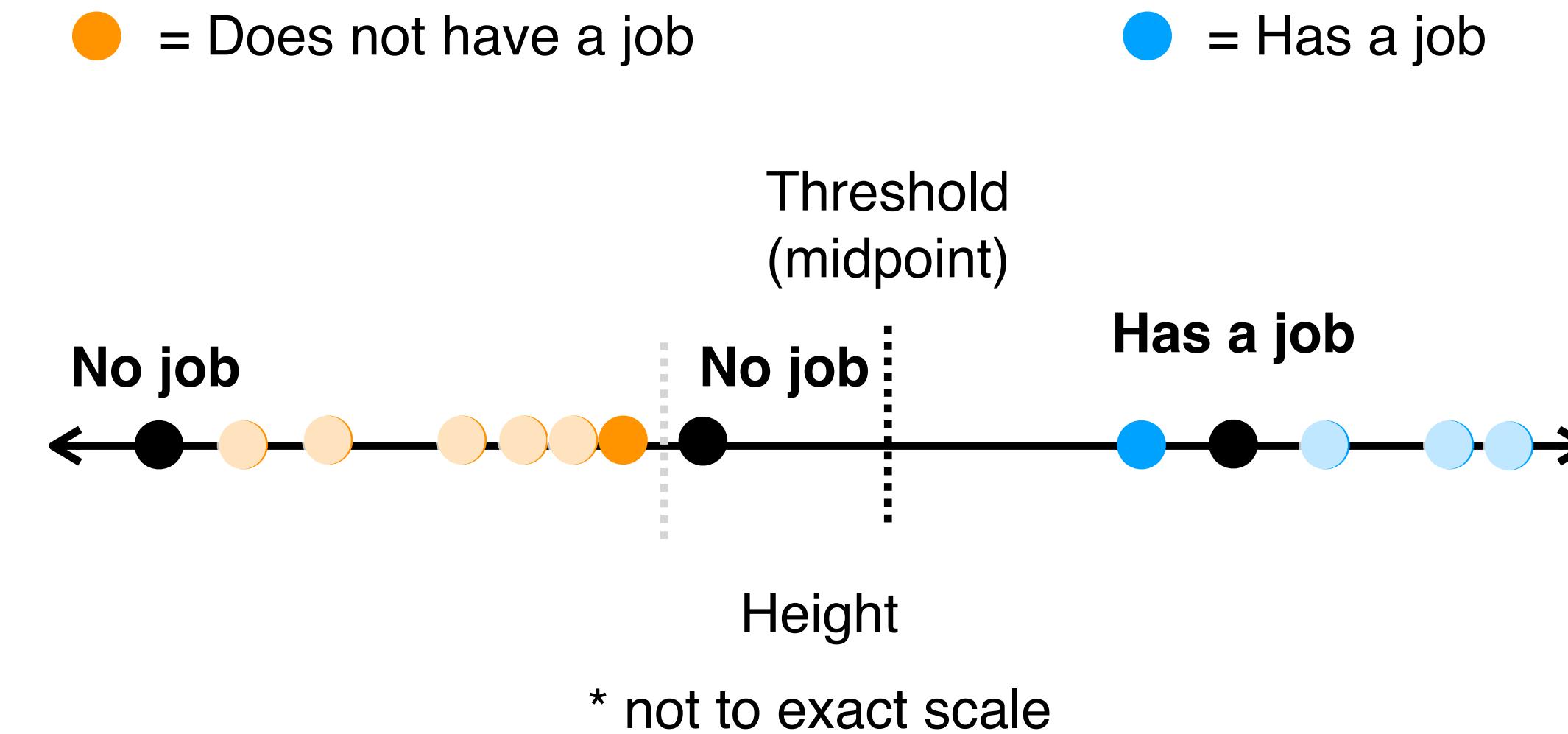
Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



Less about visualization, more about prediction

Height	Job?
50"	No
60"	Yes
52"	No
64"	Yes
49"	No
51.5"	No
48"	No
63"	Yes
51"	No
65"	Yes



Above threshold = predicted to have a job

Below threshold = predicted to not have a job

Now ● classified to the class it's closest to (●).

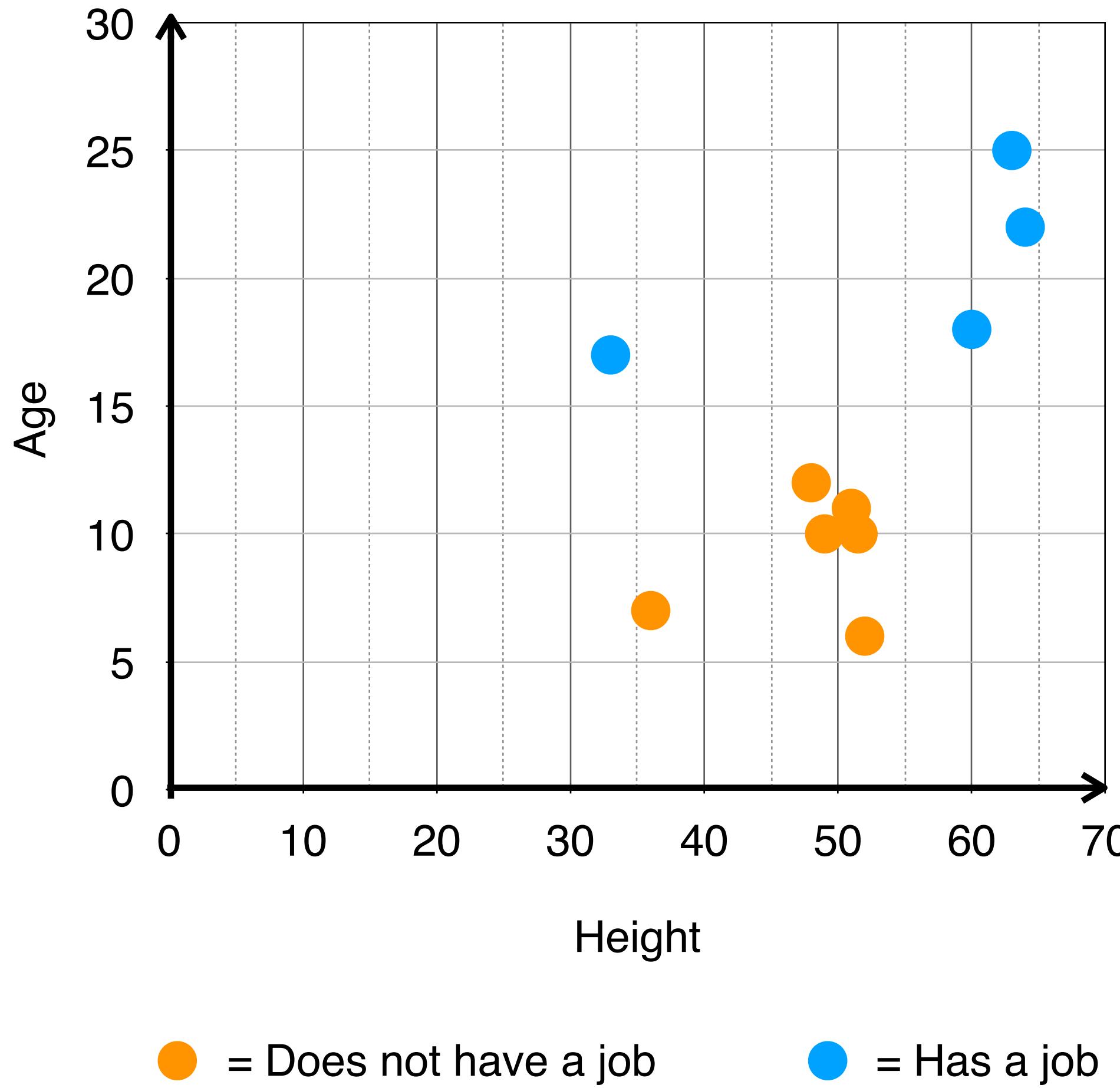
Other points still fall within their clusters.

In 2D, we no longer have a threshold point, but a threshold **line**.

Height	Age	Job?
36"	7	No
60"	18	Yes
52"	6	No
64"	22	Yes
49"	10	No
51.5"	10	No
48"	12	No
63"	25	Yes
51"	11	No
33"	17	Yes

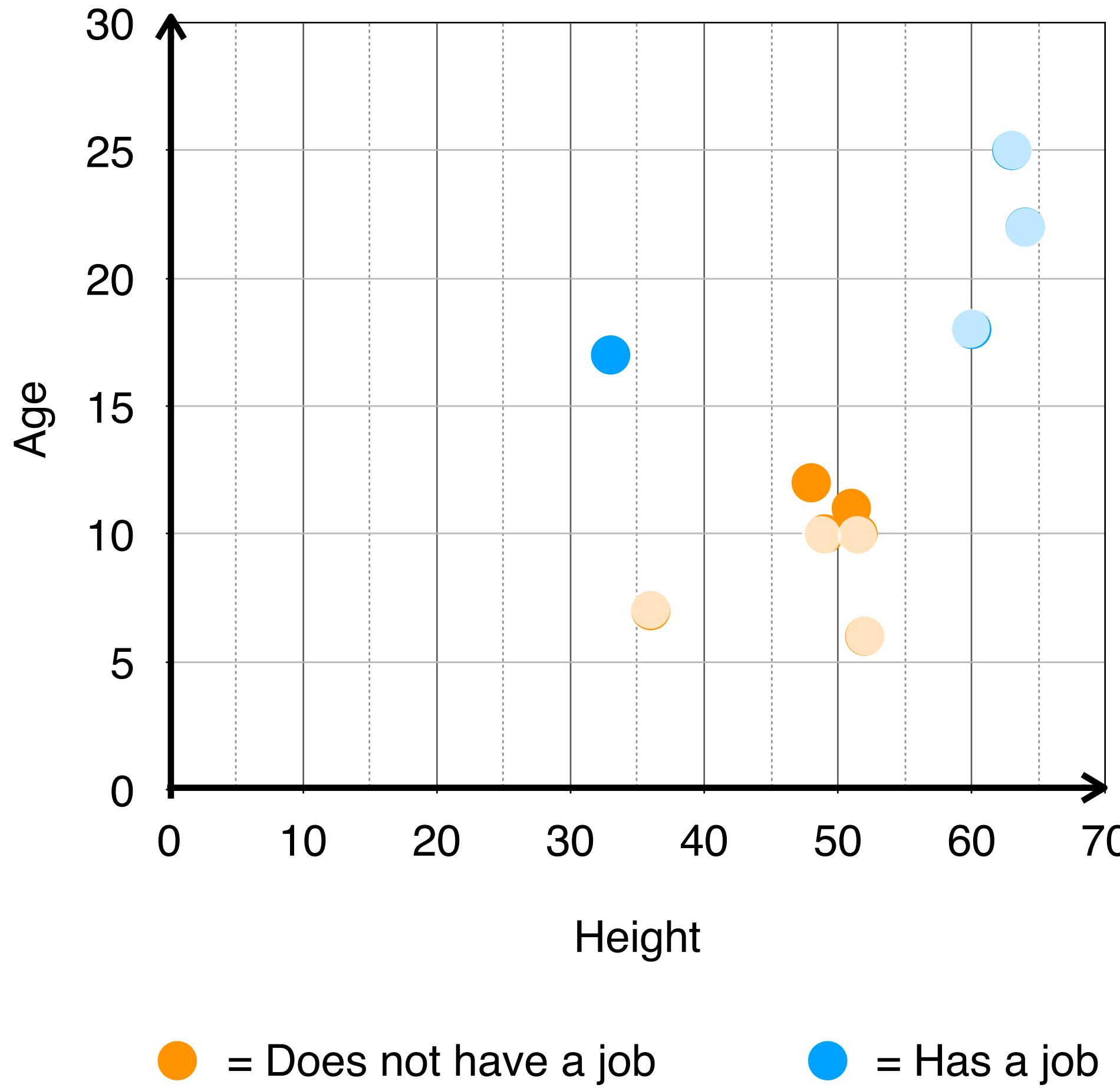
In 2D, we no longer have a threshold point, but a threshold **line**.

Height	Age	Job?
36"	7	No
60"	18	Yes
52"	6	No
64"	22	Yes
49"	10	No
51.5"	10	No
48"	12	No
63"	25	Yes
51"	11	No
33"	17	Yes



In 2D, we no longer have a threshold point, but a threshold **line**.

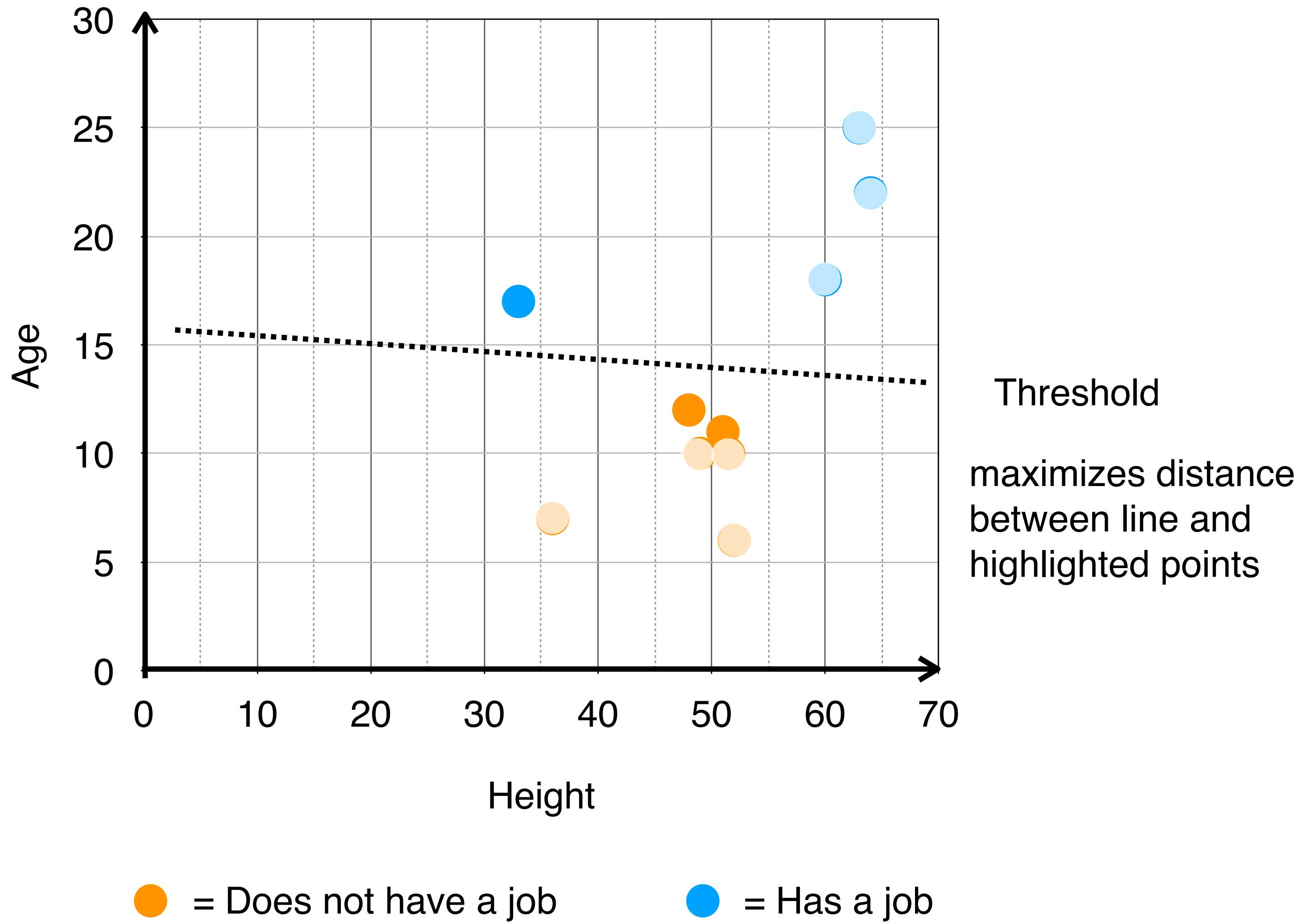
Height	Age	Job?
36"	7	No
60"	18	Yes
52"	6	No
64"	22	Yes
49"	10	No
51.5"	10	No
48"	12	No
63"	25	Yes
51"	11	No
33"	17	Yes





In 2D, we no longer have a threshold point, but a threshold **line**.

Height	Age	Job?
36"	7	No
60"	18	Yes
52"	6	No
64"	22	Yes
49"	10	No
51.5"	10	No
48"	12	No
63"	25	Yes
51"	11	No
33"	17	Yes

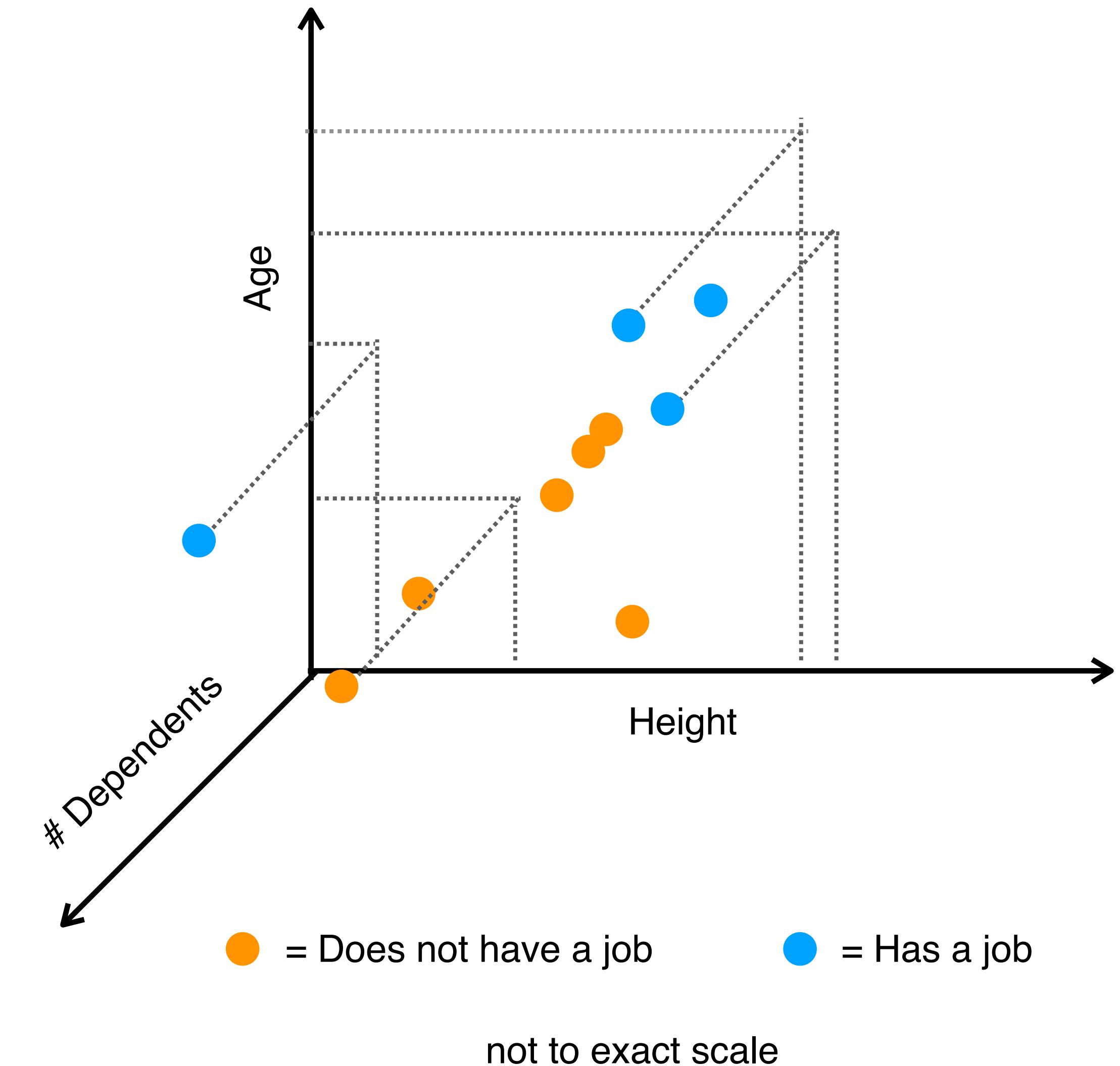


In 3D, we no longer have a threshold line, but a threshold **plane**.

Height	Age	Number of Dependents	Job?
36"	7	0	No
60"	18	0	Yes
52"	6	0	No
64"	22	1	Yes
49"	10	0	No
51.5"	10	0	No
48"	12	1	No
63"	25	1	Yes
51"	11	0	No
33"	17	1	Yes

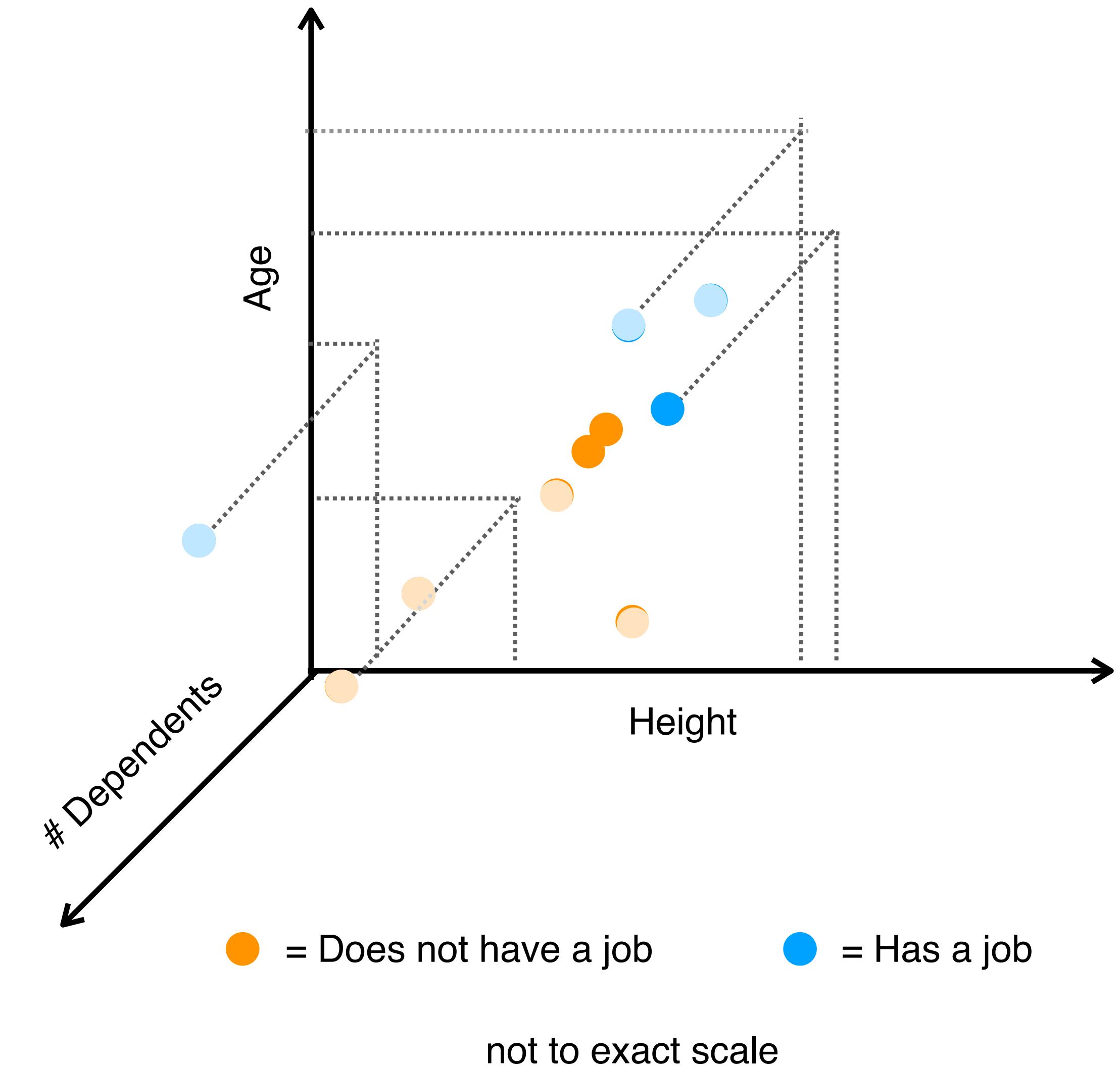
In 3D, we no longer have a threshold line, but a threshold **plane**.

Height	Age	Number of Dependents	Job?
36"	7	0	No
60"	18	0	Yes
52"	6	0	No
64"	22	1	Yes
49"	10	0	No
51.5"	10	0	No
48"	12	1	No
63"	25	1	Yes
51"	11	0	No
33"	17	1	Yes



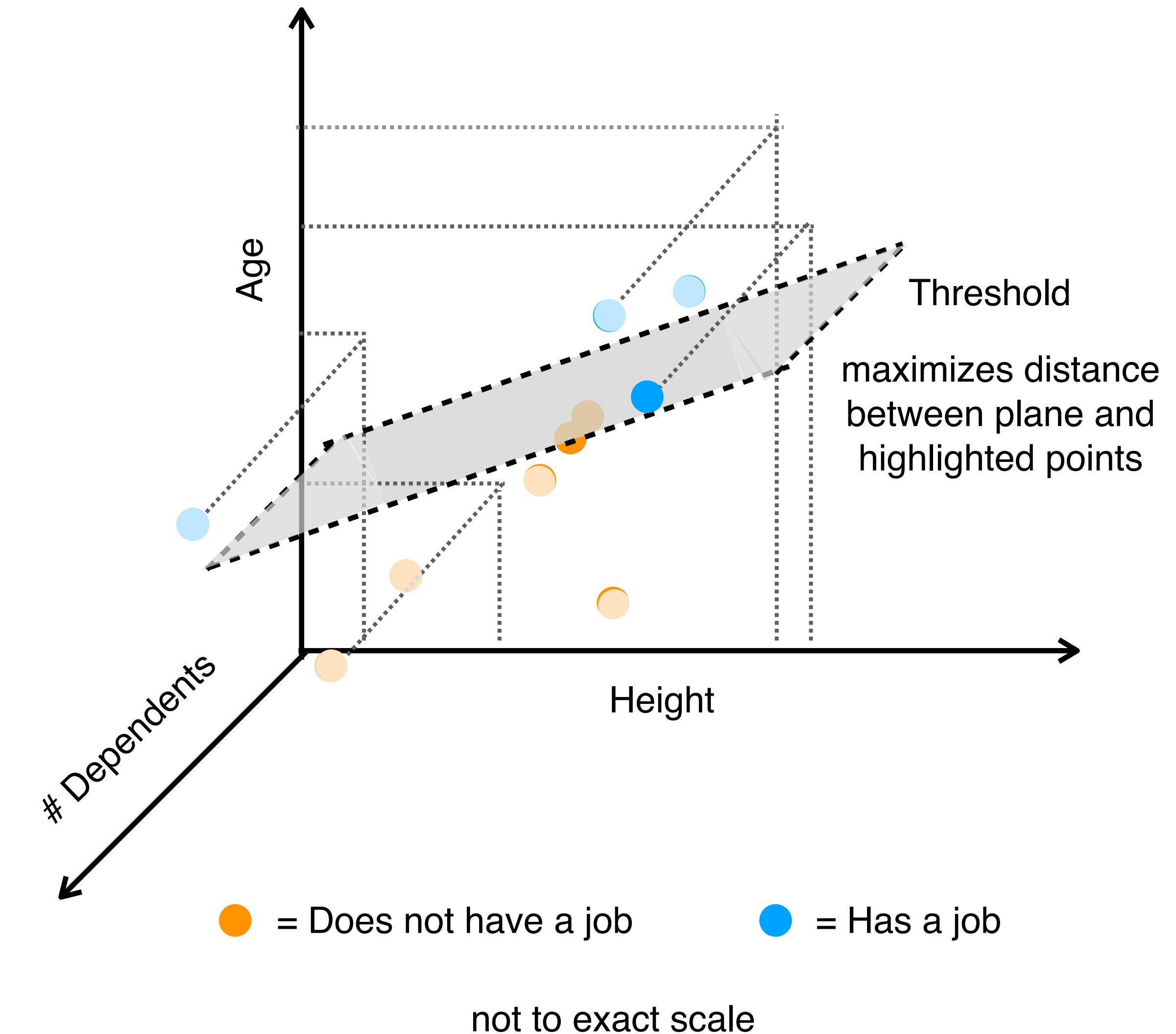
In 3D, we no longer have a threshold line, but a threshold **plane**.

Height	Age	Number of Dependents	Job?
36"	7	0	No
60"	18	0	Yes
52"	6	0	No
64"	22	1	Yes
49"	10	0	No
51.5"	10	0	No
48"	12	1	No
63"	25	1	Yes
51"	11	0	No
33"	17	1	Yes



In 3D, we no longer have a threshold line, but a threshold **plane**.

Height	Age	Number of Dependents	Job?
36"	7	0	No
60"	18	0	Yes
52"	6	0	No
64"	22	1	Yes
49"	10	0	No
51.5"	10	0	No
48"	12	1	No
63"	25	1	Yes
51"	11	0	No
33"	17	1	Yes



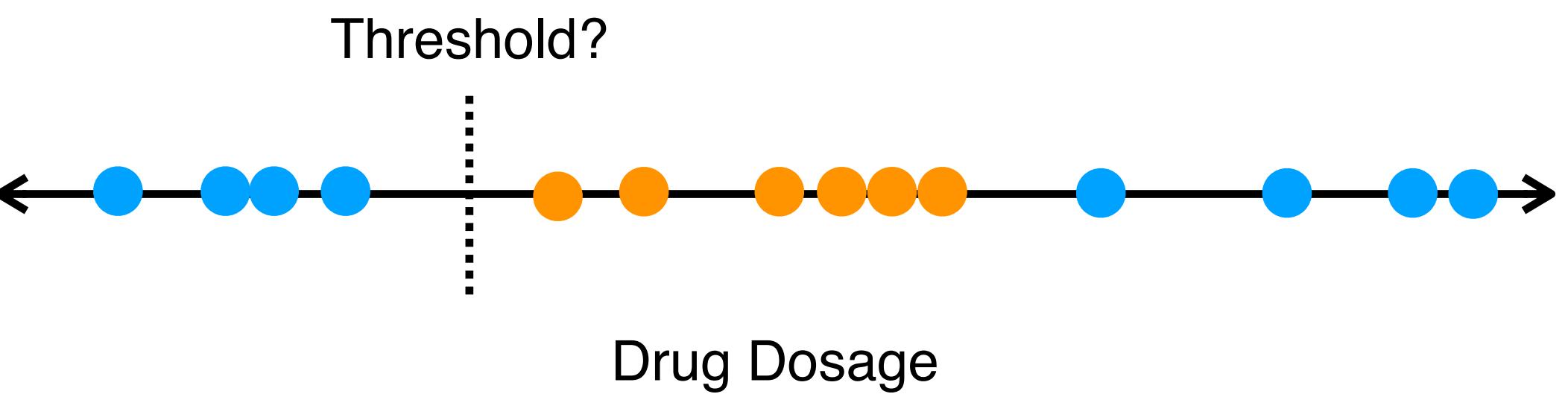
● = beneficial effect

● = no beneficial effect  
(no effect / toxic)



● = beneficial effect

● = no beneficial effect  
(no effect / toxic)

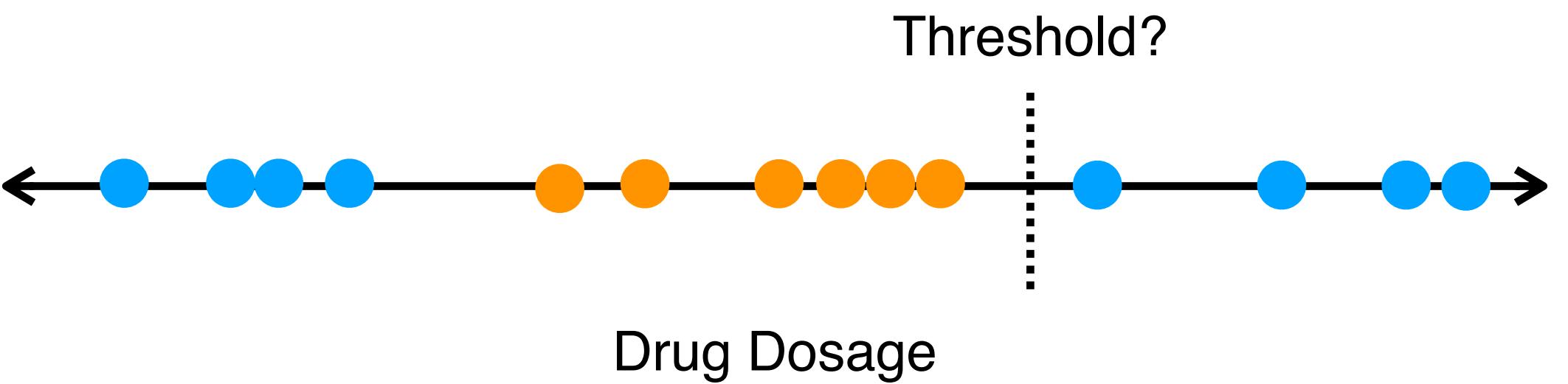


# A Problem Case



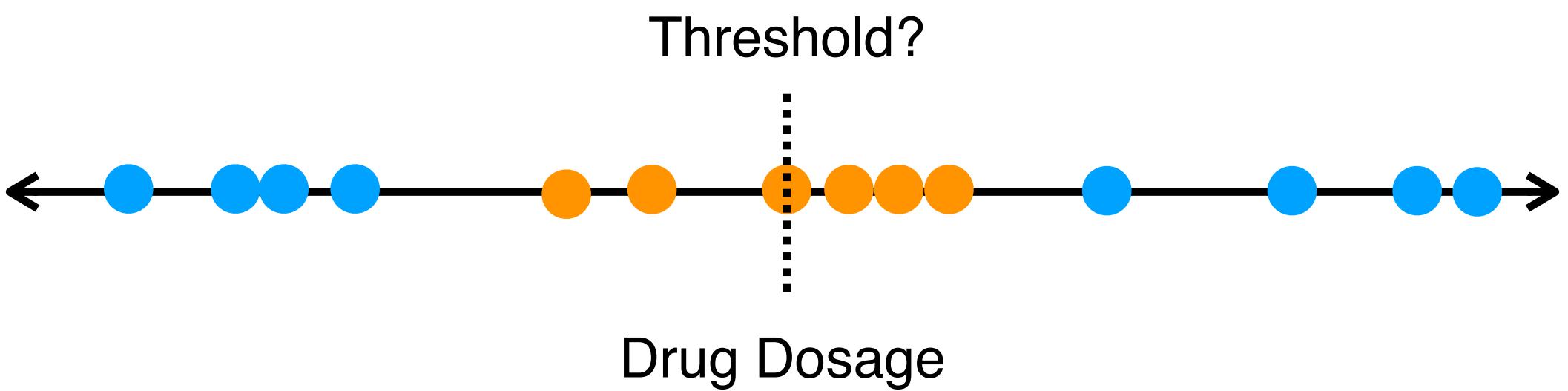
● = beneficial effect

● = no beneficial effect  
(no effect / toxic)



● = beneficial effect

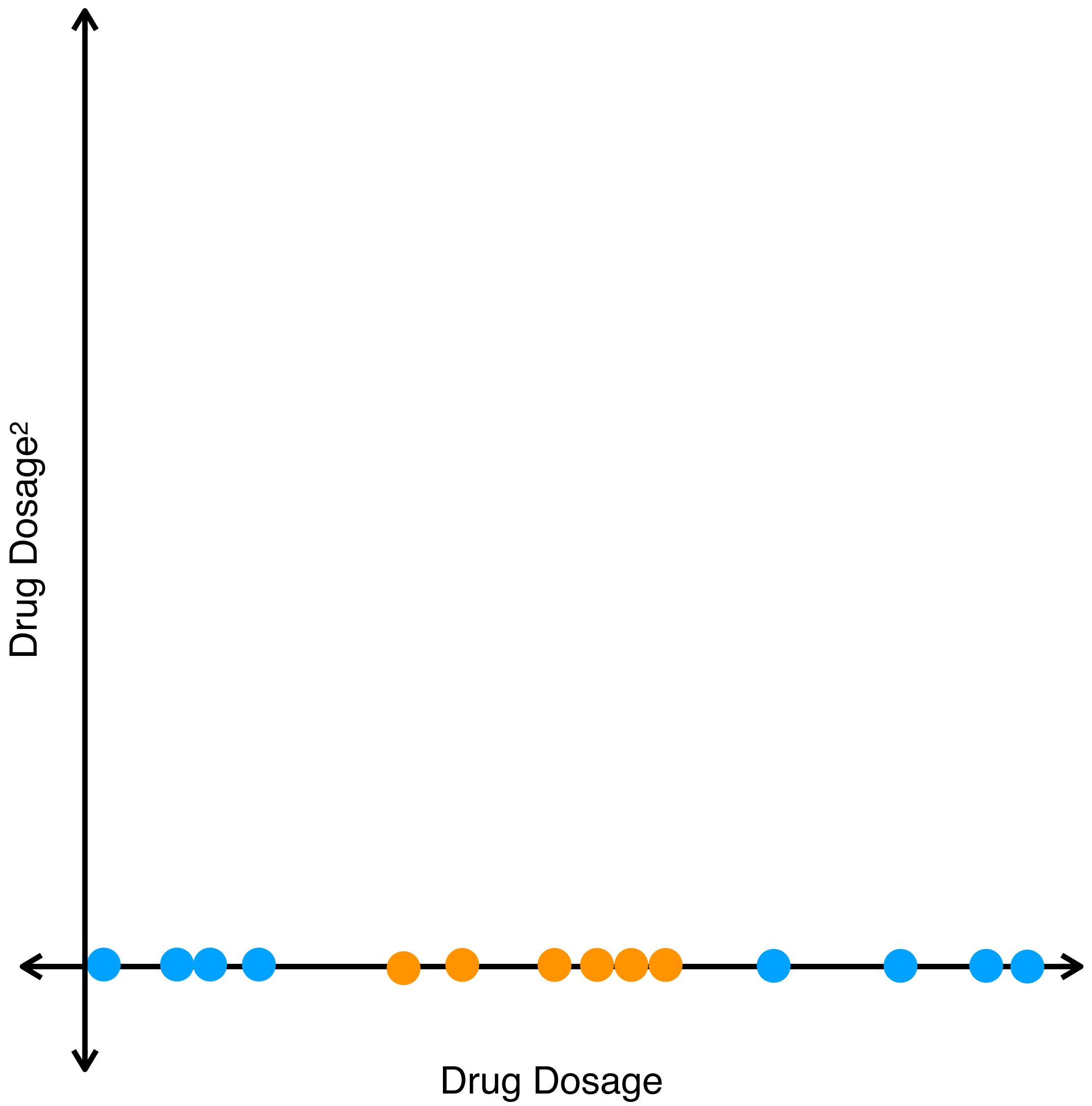
● = no beneficial effect  
(no effect / toxic)



No matter where we place the threshold, we get poor classification!

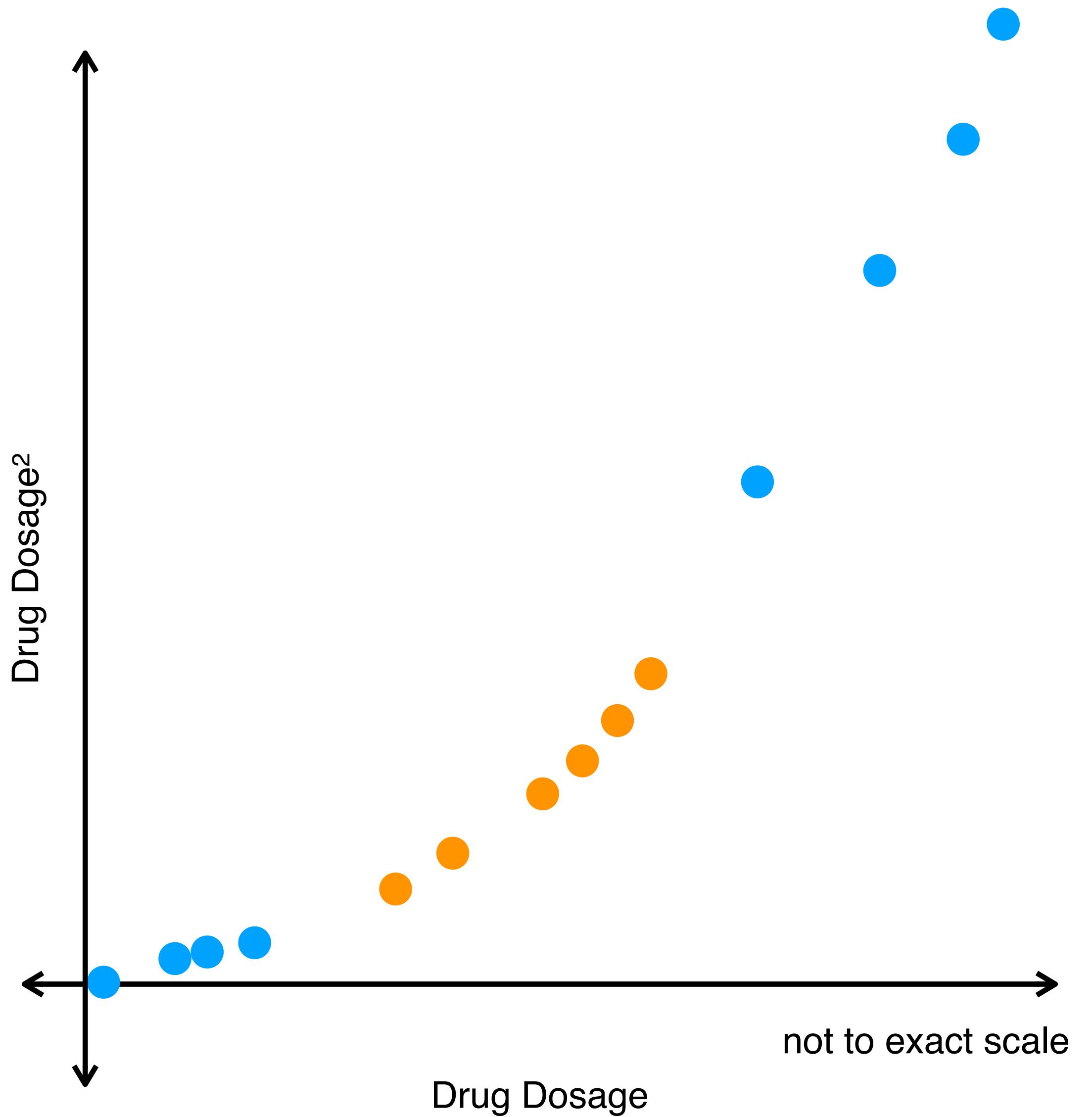
## A Problem Case

- = beneficial effect
- = no beneficial effect  
(no effect / toxic)

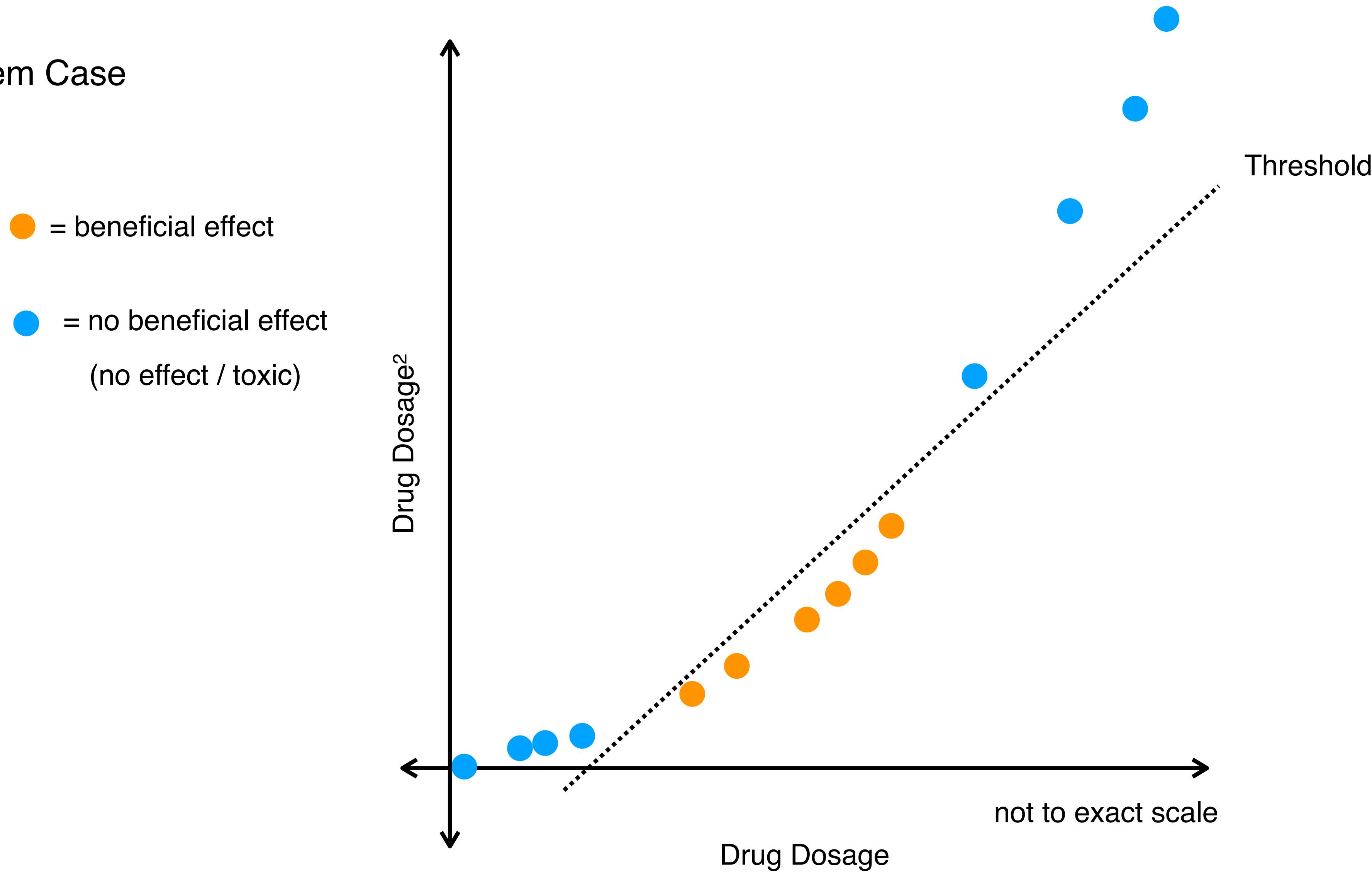


## A Problem Case

- = beneficial effect
- = no beneficial effect  
(no effect / toxic)

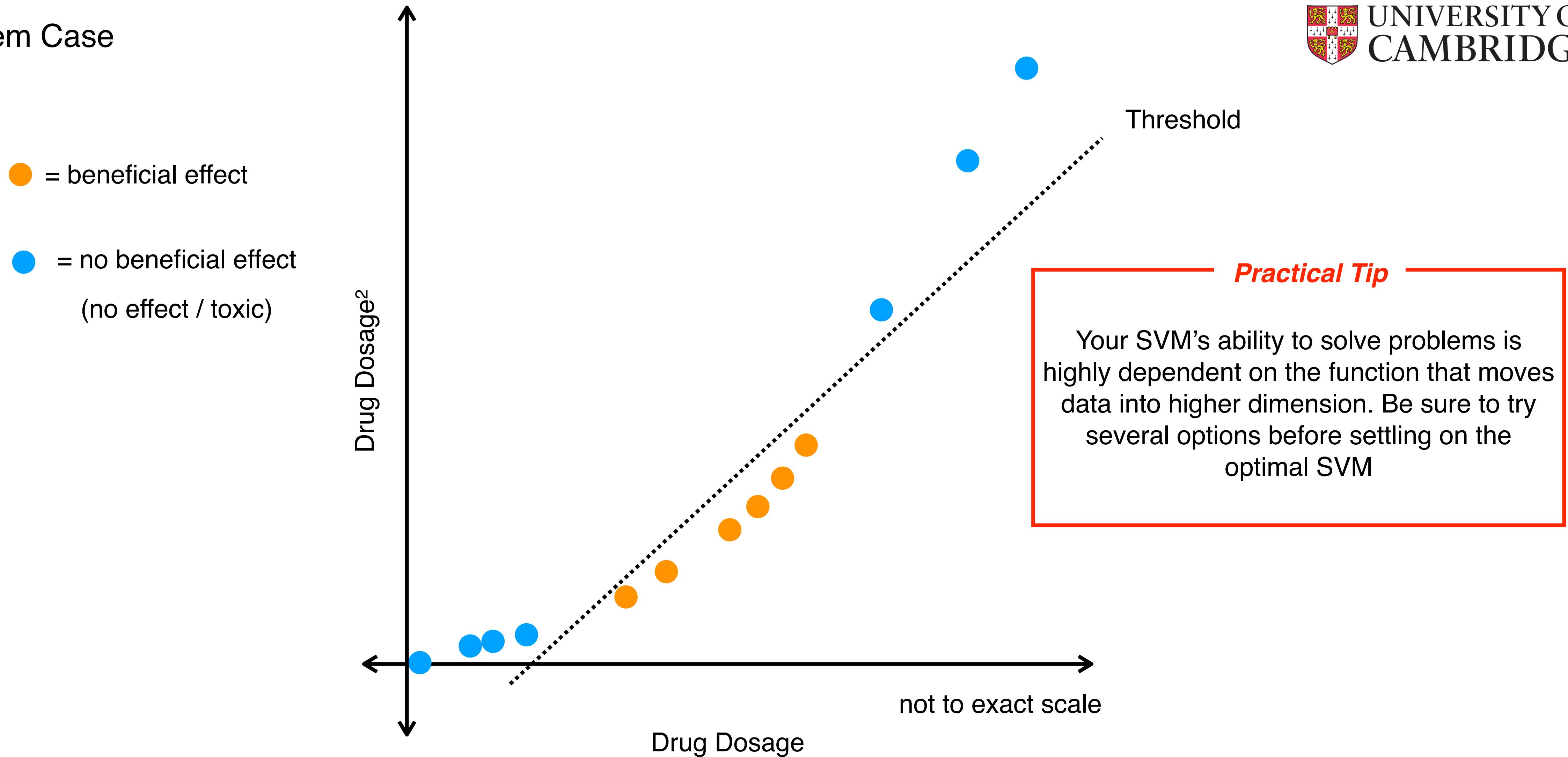
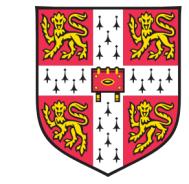


## A Problem Case



We have increased the dimensionality of the data to find a threshold!

## A Problem Case

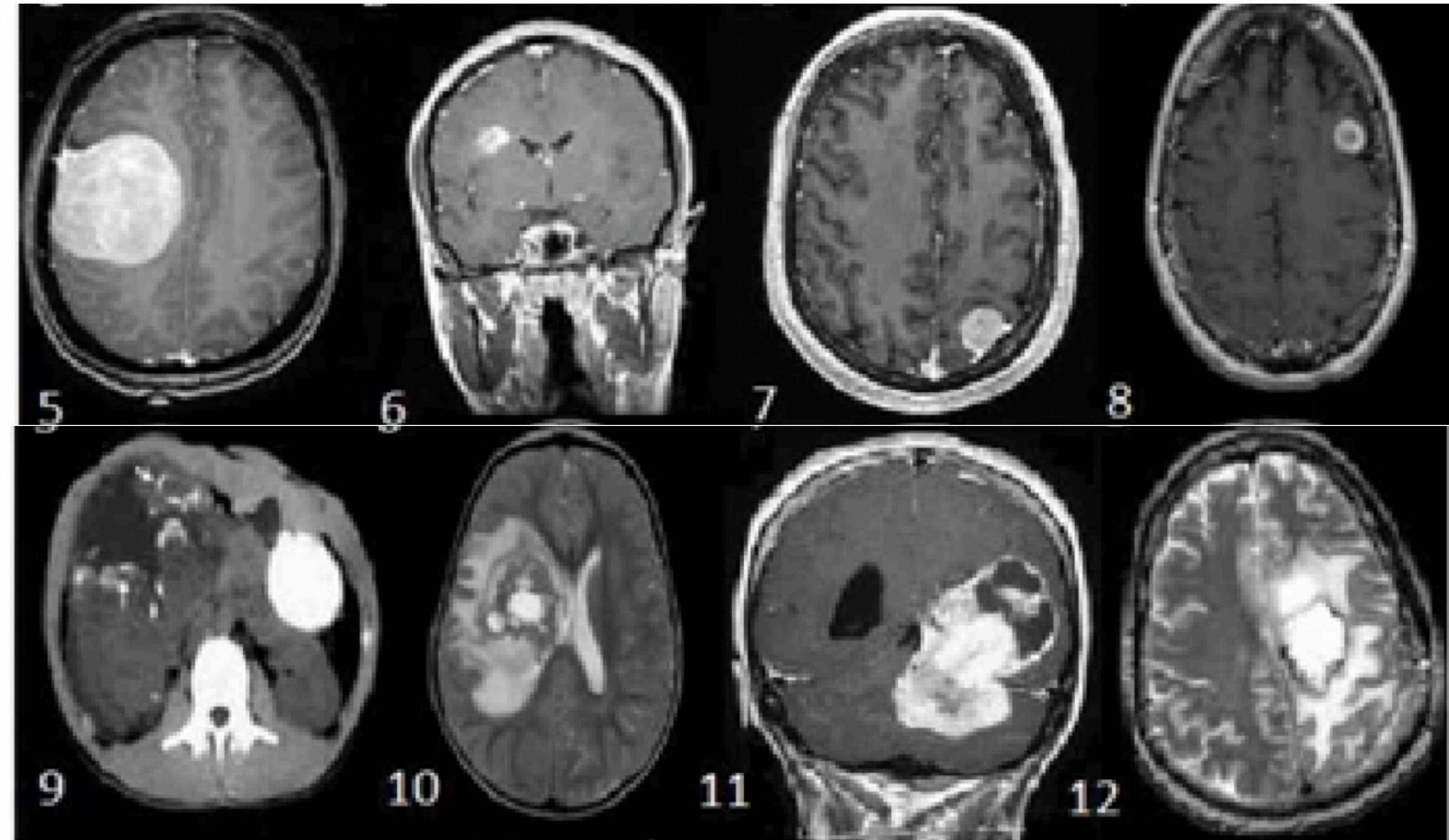


We have increased the dimensionality of the data to find a threshold!

The “reverse” of PCA (reducing dimensions to simplify finding correlations).

Which brain tumor is malignant and which is benign?

253 MRI images dataset

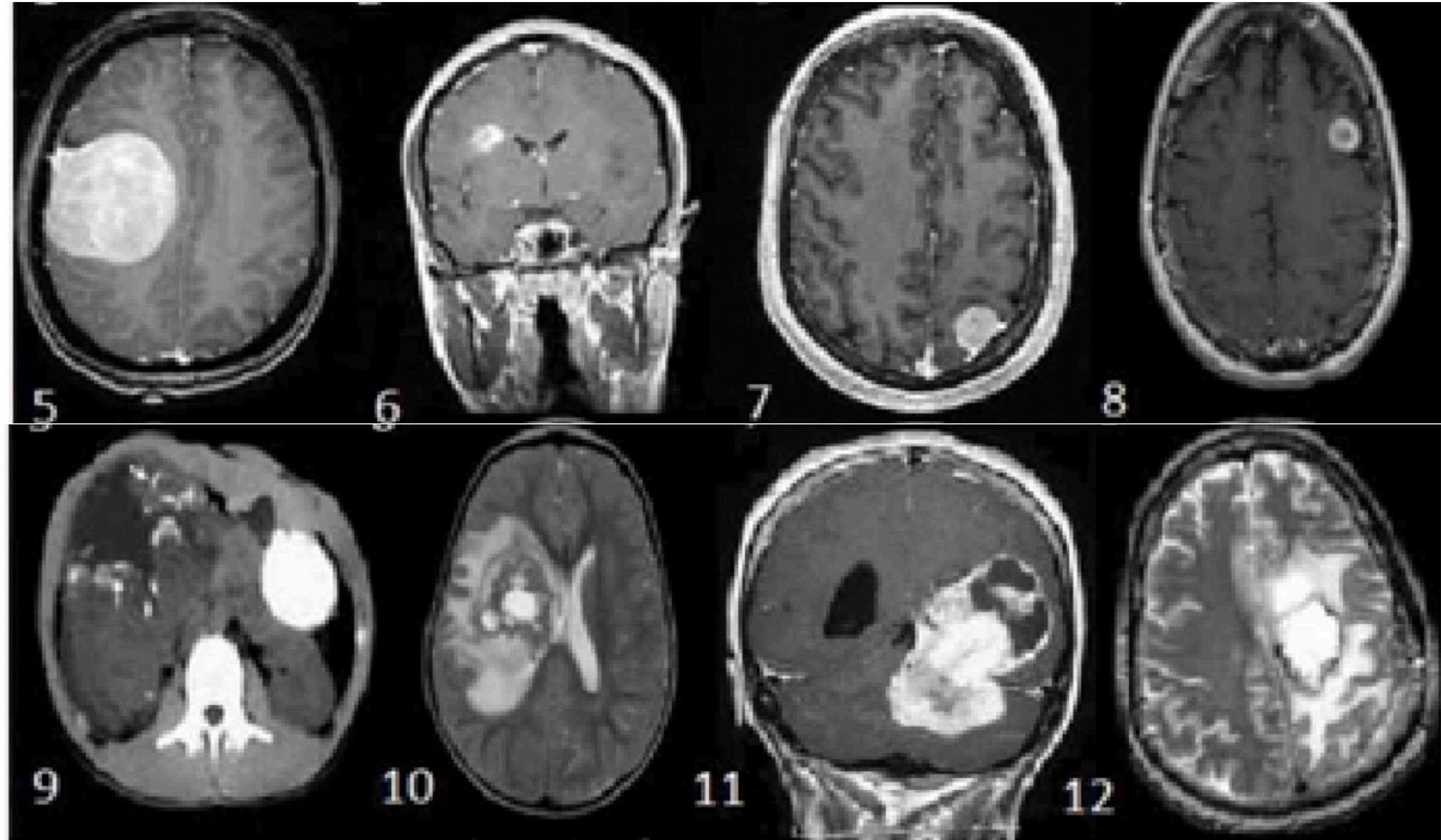


Accuracy  $\approx 80\%$

Which brain tumor is malignant and which is benign?

253 MRI images dataset

These are benign



These are not benign  
(malignant)

Accuracy  $\approx 80\%$



SVM		PCA
	Data Visualization	
	Classification	
	Regression	
Increases dimensions for predictions	What do dimensions have to do with it?	Reduces dimensions for human visualization