



Introduction to Machine Learning [Fall 2022]

Linear Regression (Part 2)

September 15, 2022

Lerrel Pinto

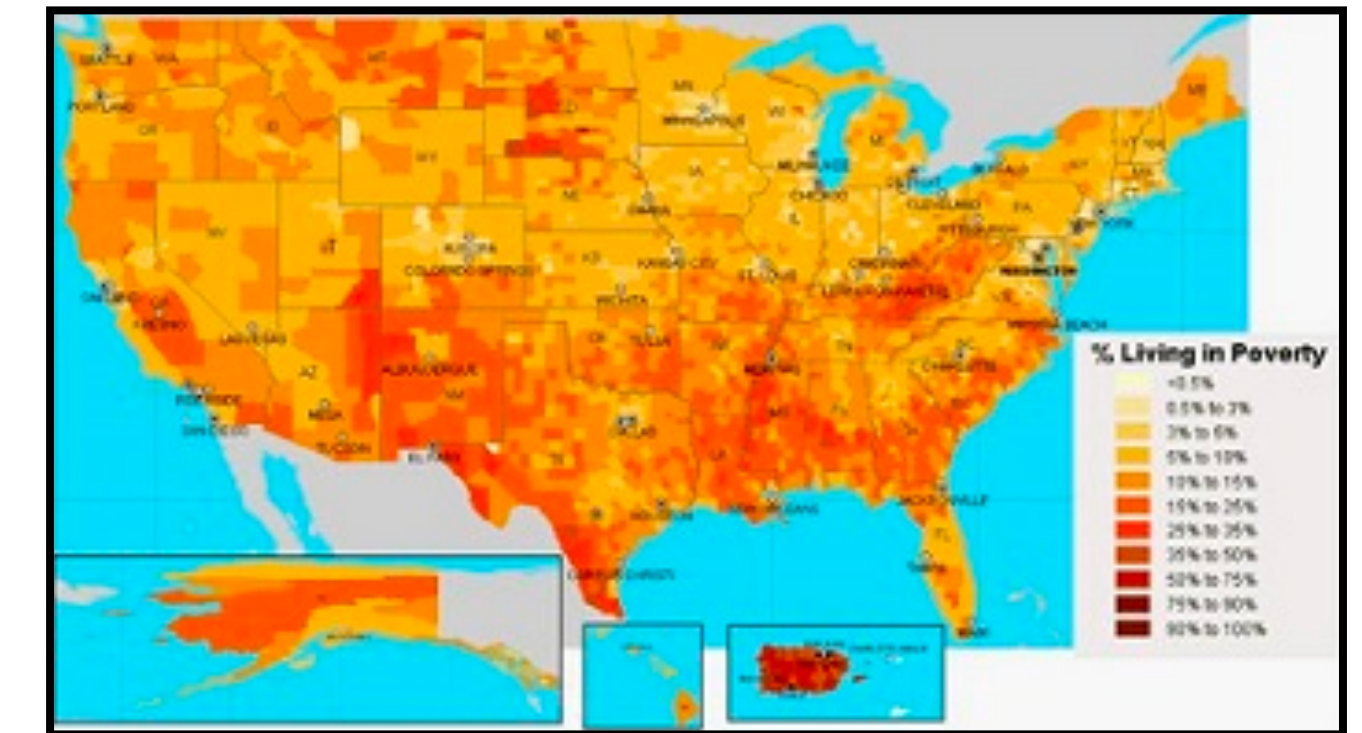
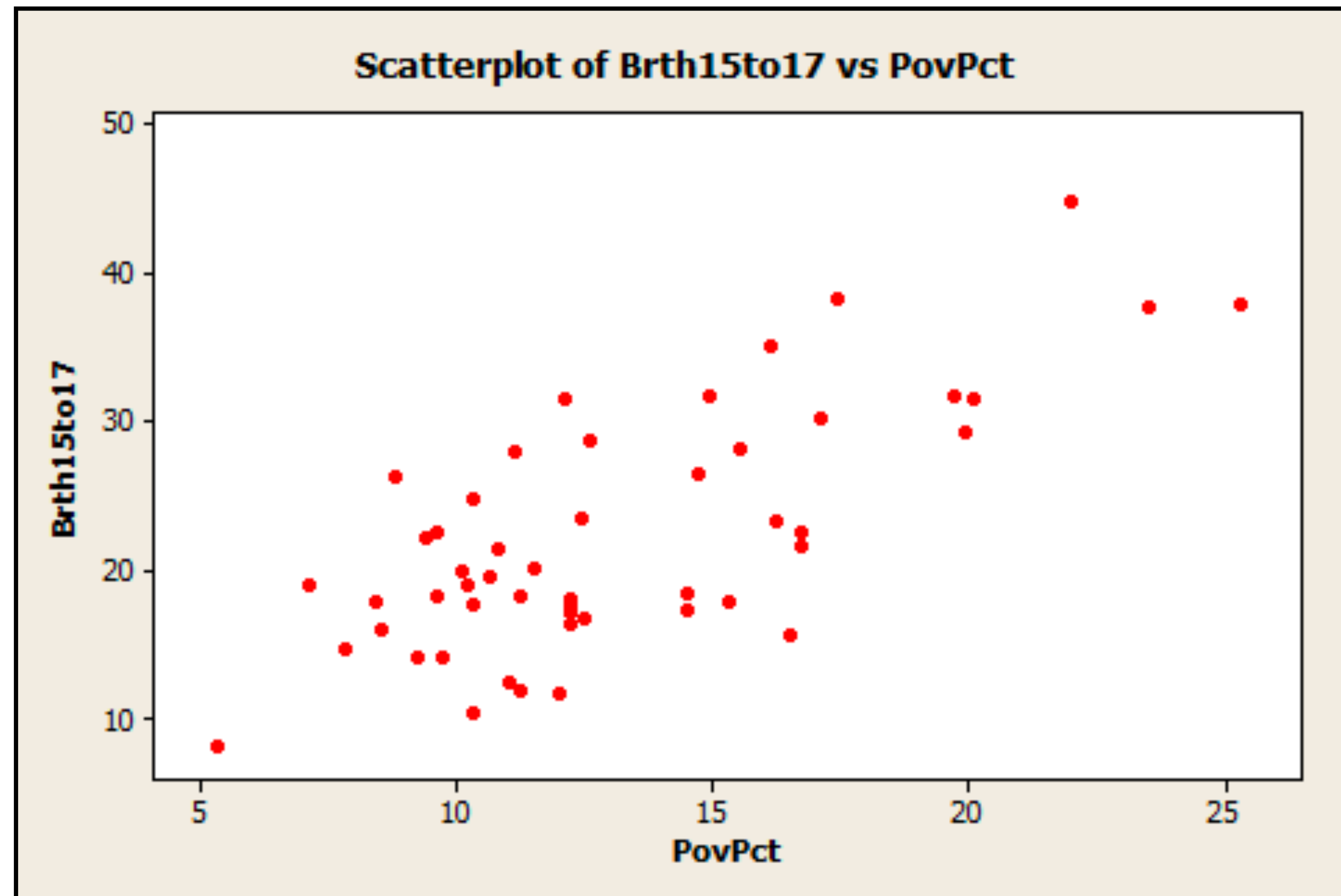
Topics for today

- More advanced linear regression.

Linear Regression (example)

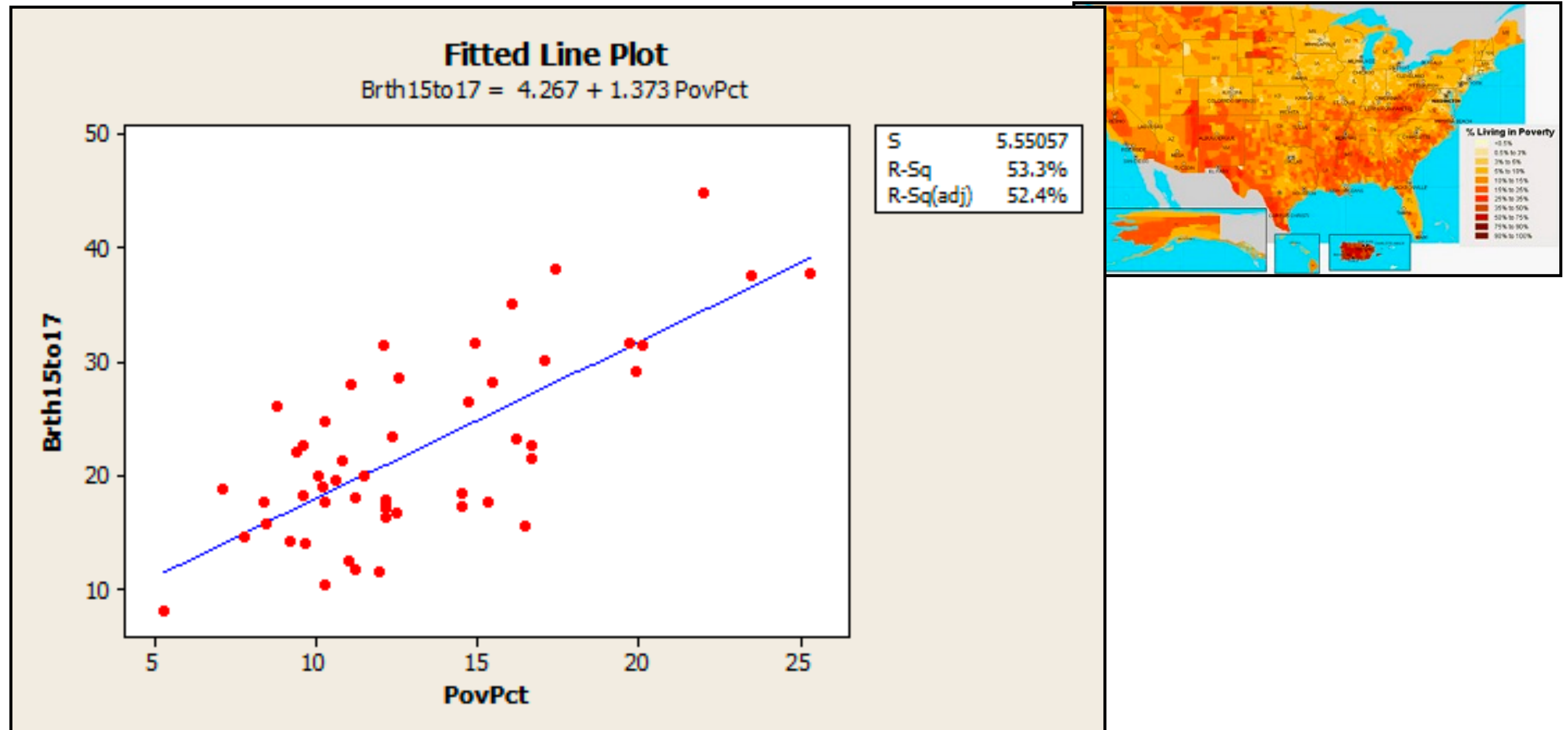
<https://online.stat.psu.edu/stat462/node/101/>

Linear Regression (example)



<https://online.stat.psu.edu/stat462/node/101/>

Linear Regression (example)



<https://online.stat.psu.edu/stat462/node/101/>

Linear Regression

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow \#$ of data points, $d \rightarrow \#$ of features / input dim.

Linear Regression

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow \#$ of data points, $d \rightarrow \#$ of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^T \vec{w} - Y\|^2$

Linear Regression

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow \#$ of data points, $d \rightarrow \#$ of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^T \vec{w} - Y\|^2$
- Solution: $\vec{w} = (XX^T)^{-1}XY$
 - Easy way to remember $\vec{w} = (X^T)^+Y$

Linear Regression (non-linear?)

- Vanilla linear regression (OLS) assumes that: $y = \vec{w}^T x$

Linear Regression (non-linear?)

- Vanilla linear regression (OLS) assumes that: $y = \vec{w}^T x$
- What if we want to model non linear relationship:

$$y = \vec{w}_1^T x + \vec{w}_2^T x^2 + \vec{w}_3^T x^3 + \dots + \vec{w}_p^T x^p$$

- Is there an easy way to do this?

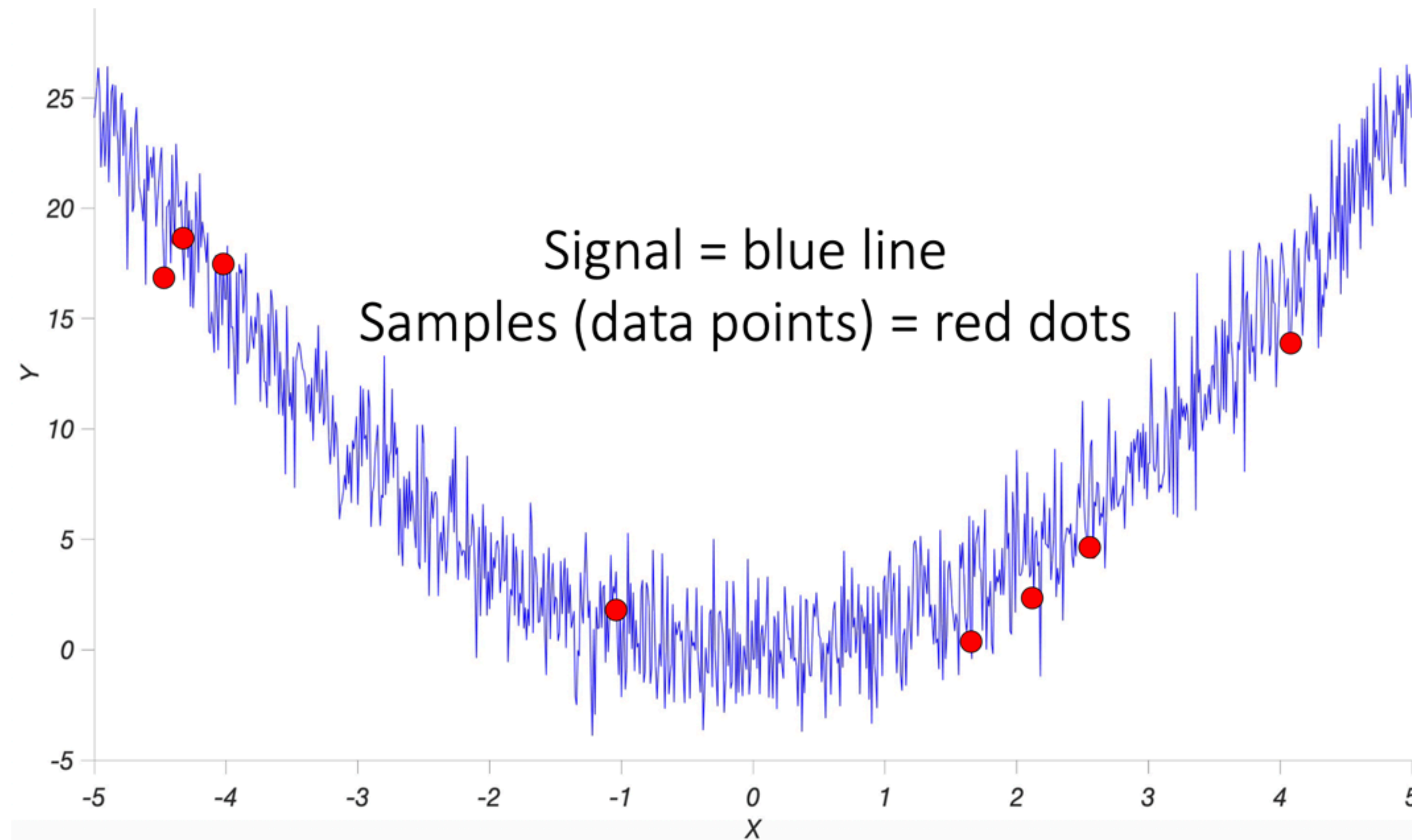
Linear Regression (non-linear?)

- Vanilla linear regression (OLS) assumes that: $y = \vec{w}^T x$
- What if we want to model non linear relationship:

$$y = \vec{w}_1^T x + \vec{w}_2^T x^2 + \vec{w}_3^T x^3 + \dots + \vec{w}_p^T x^p$$

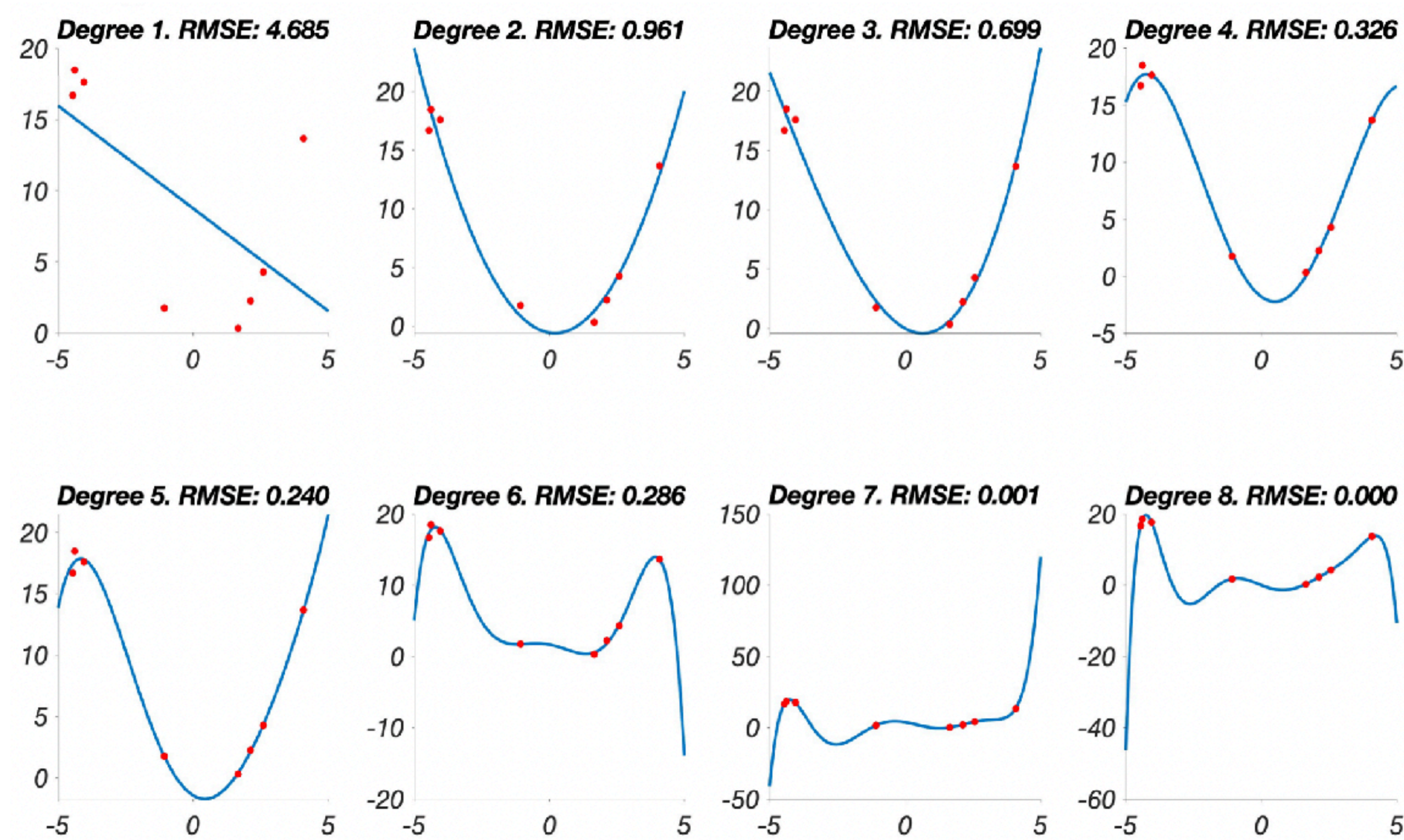
- Is there an easy way to do this?
 - Define new variable: $z = [x, x^2, x^3, \dots, x^p]$
 - Do vanilla linear regression: $y = \vec{w}^T z$

Linear Regression (bias vs variance)



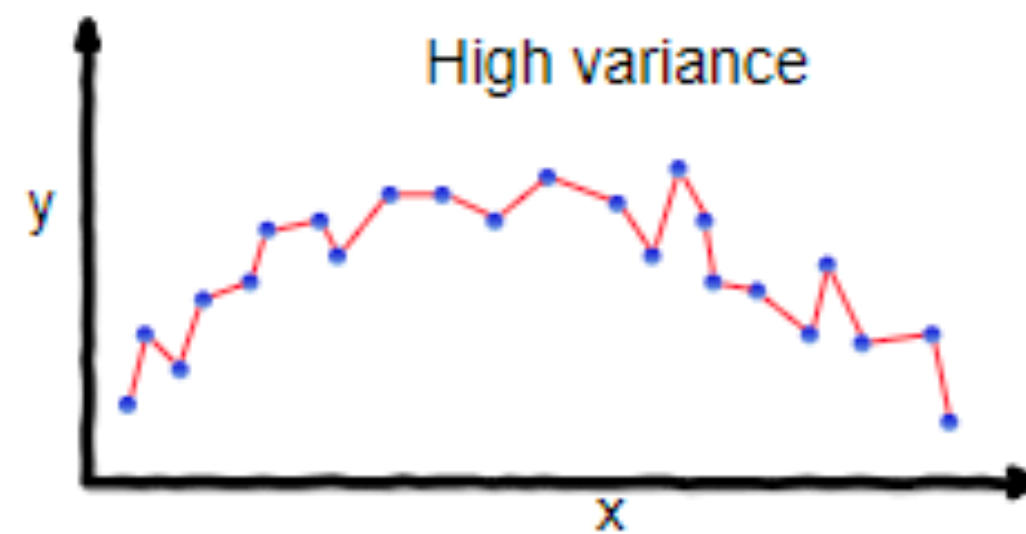
Slide credit: Pascal Wallisch

Linear Regression (bias vs variance)

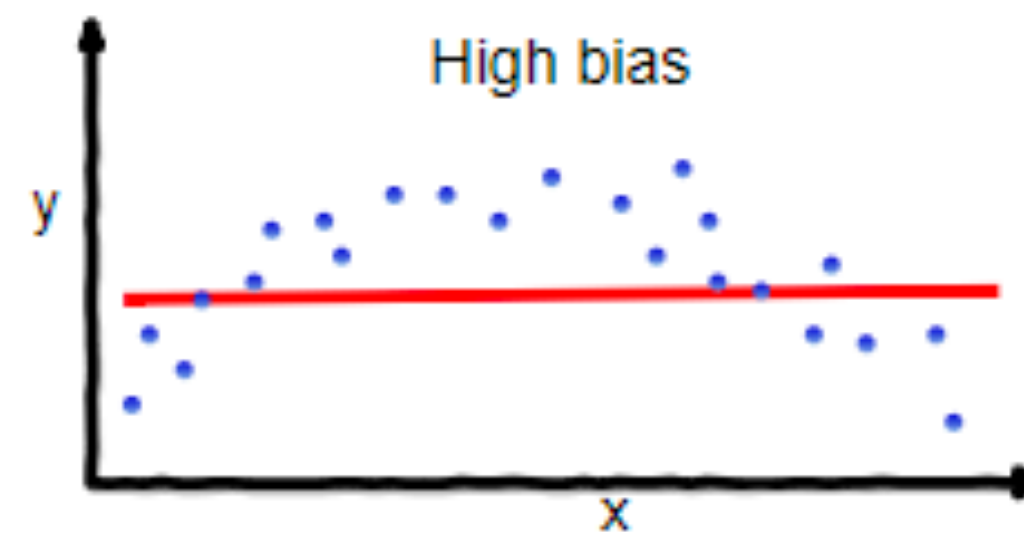


Slide credit: Pascal Wallisch

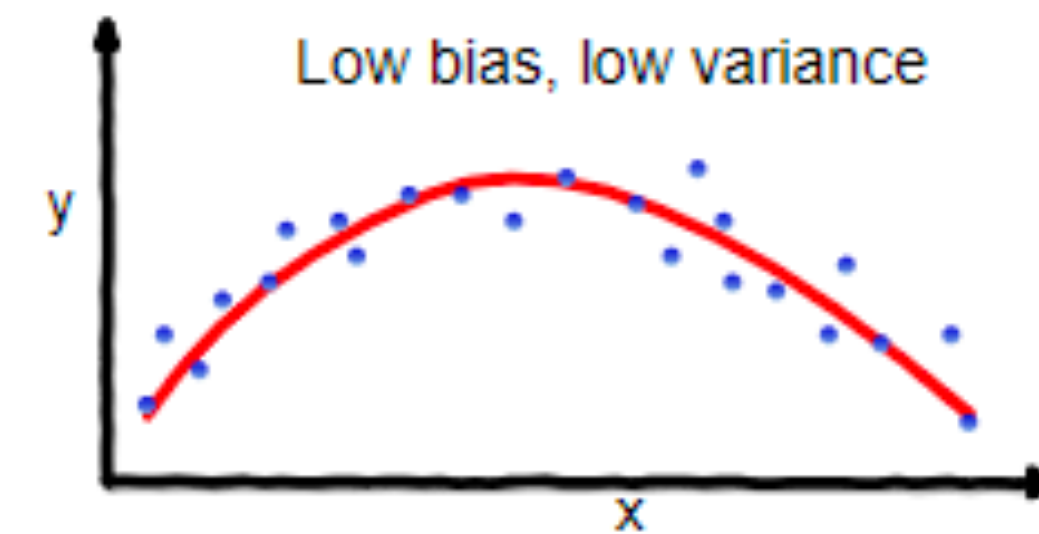
Bias vs variance tradeoff



overfitting

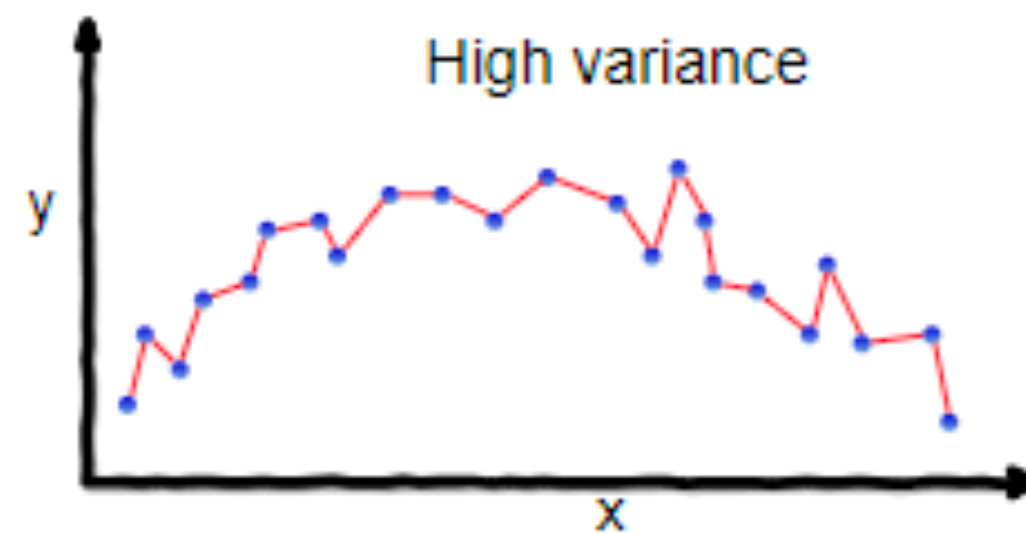


underfitting

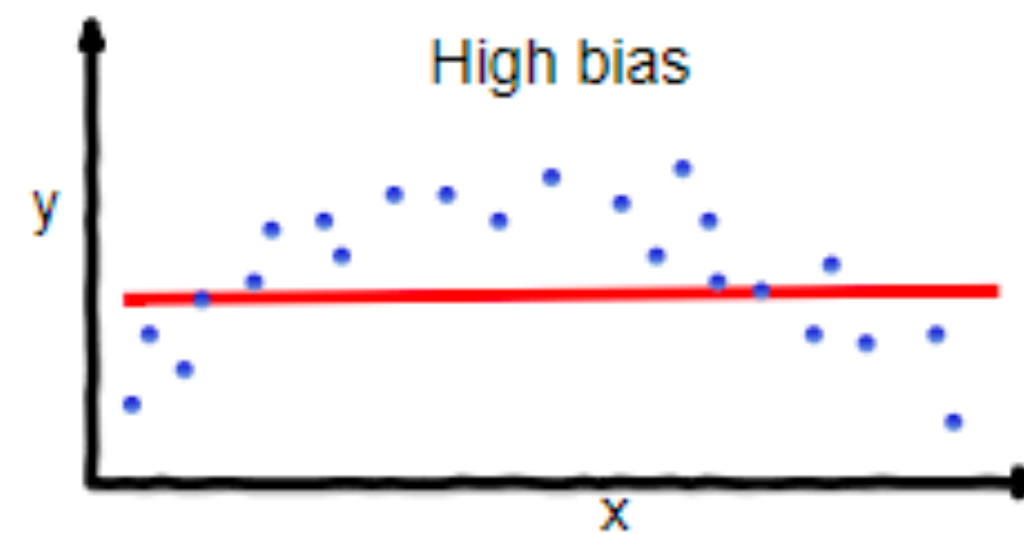


Good balance

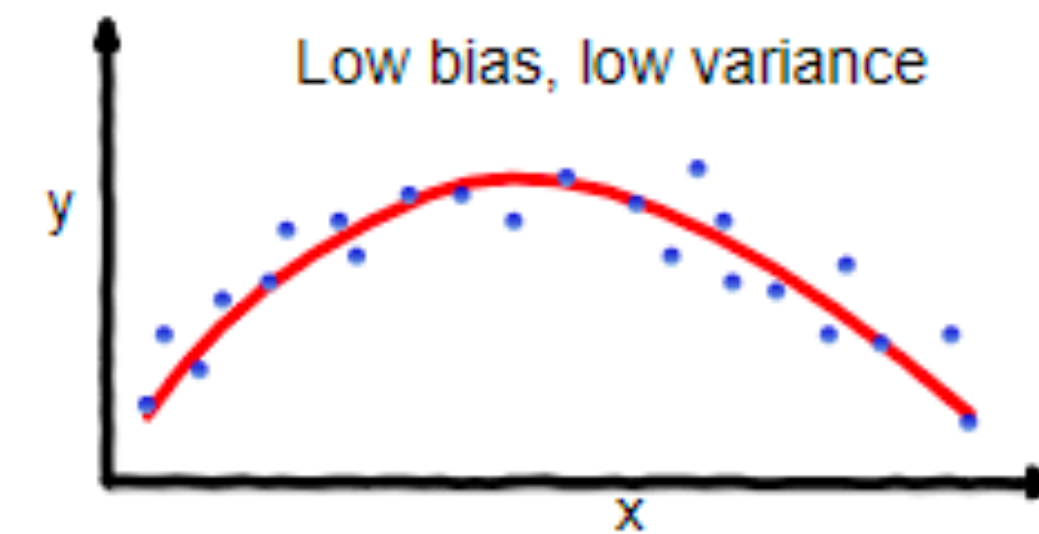
Bias vs variance tradeoff



overfitting



underfitting



Good balance

- How do we know if a model has high bias or variance?
- How do we find the balance?
 - Algorithms, Validation, prior knowledge

Scratchpad

Ridge Regression

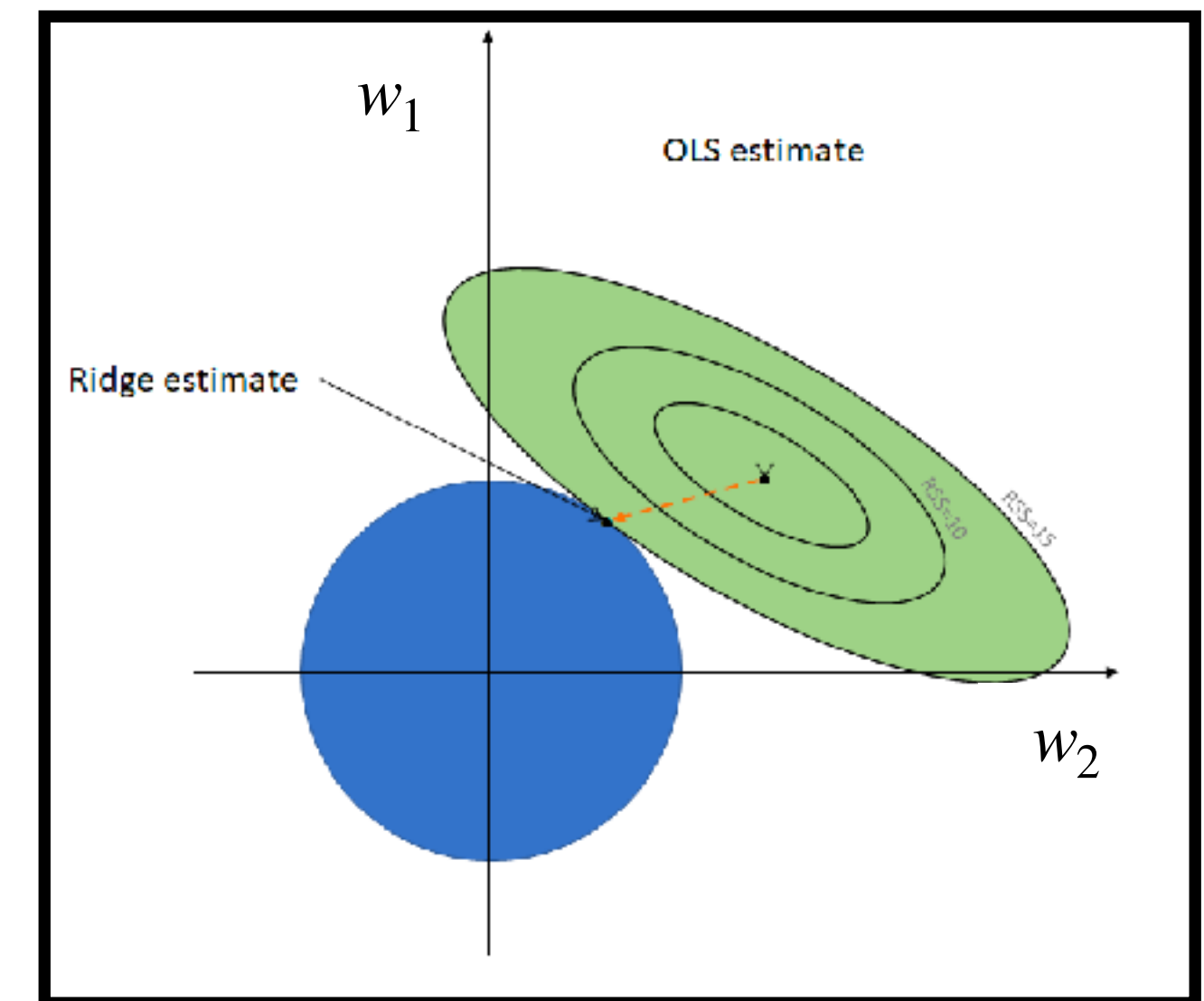
- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$

Ridge Regression

- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$
- Ridge Regression: Minimize $\|X^T \vec{w} - Y\|^2$ such that $\|\vec{w}\|^2 \leq c^2$

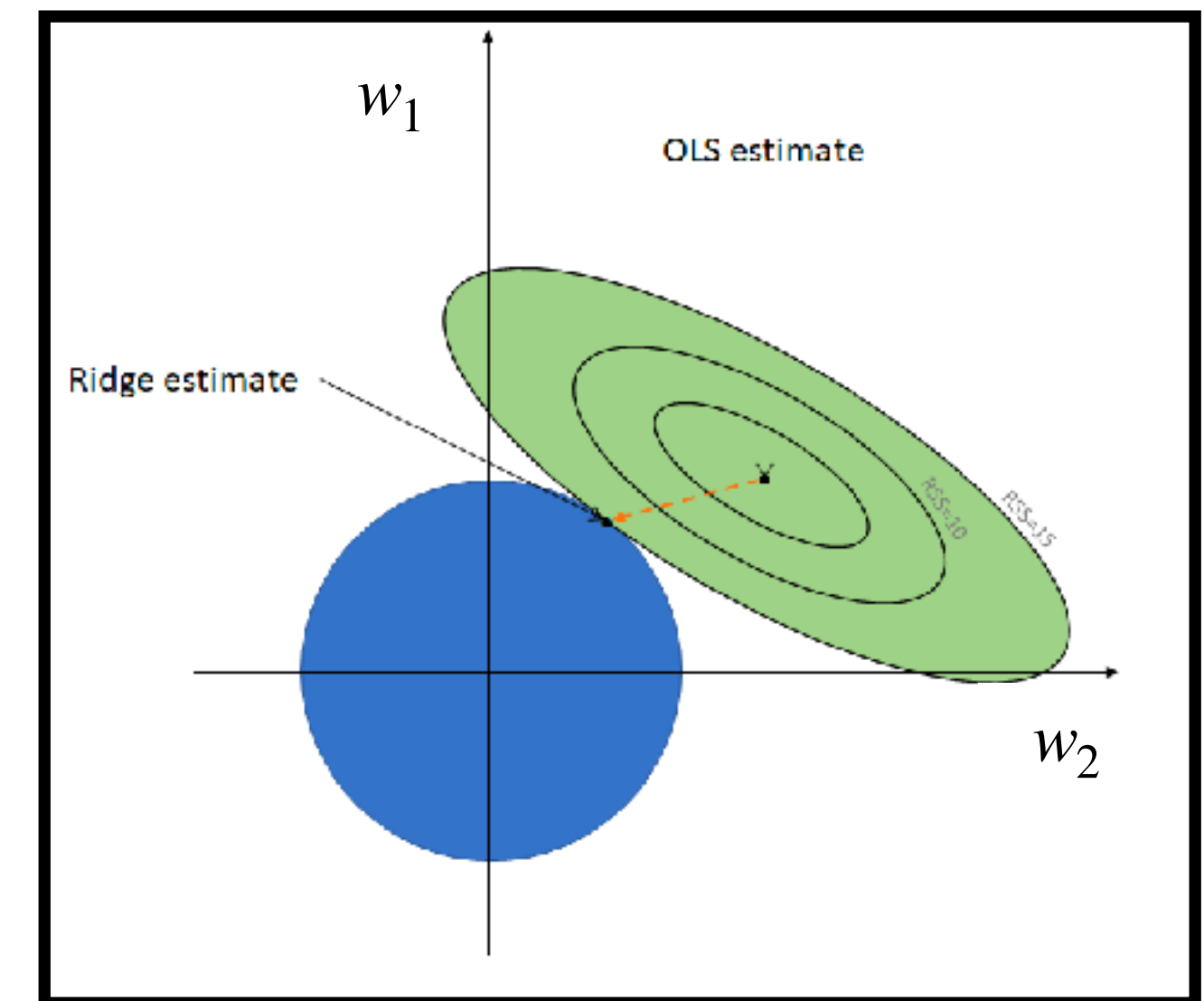
Ridge Regression

- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$
- Ridge Regression: Minimize $\|X^T \vec{w} - Y\|^2$ such that $\|\vec{w}\|^2 \leq c^2$



Ridge Regression

- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$
- Ridge Regression: Minimize $\|X^T \vec{w} - Y\|^2$ such that $\|\vec{w}\|^2 \leq c^2$
 - Loss function: $(\|X^T \vec{w} - Y\|^2 + \lambda \|\vec{w}\|^2)$
 - Solution: $\vec{w} = (XX^T + \lambda I)^{-1}XY$



Lasso Regression

- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$

Lasso Regression

- OLS: Minimize $\|X^T \vec{w} - Y\|^2$
 - Solution: $\vec{w} = (XX^T)^{-1}XY$
- Ridge Regression: Minimize $\|X^T \vec{w} - Y\|^2$ such that $\|\vec{w}\|_1 \leq c$
 - Loss function: $(\|X^T \vec{w} - Y\|^2 + \lambda \|\vec{w}\|_1)$
 - Solution: *From numerical methods.*

Lasso Regression

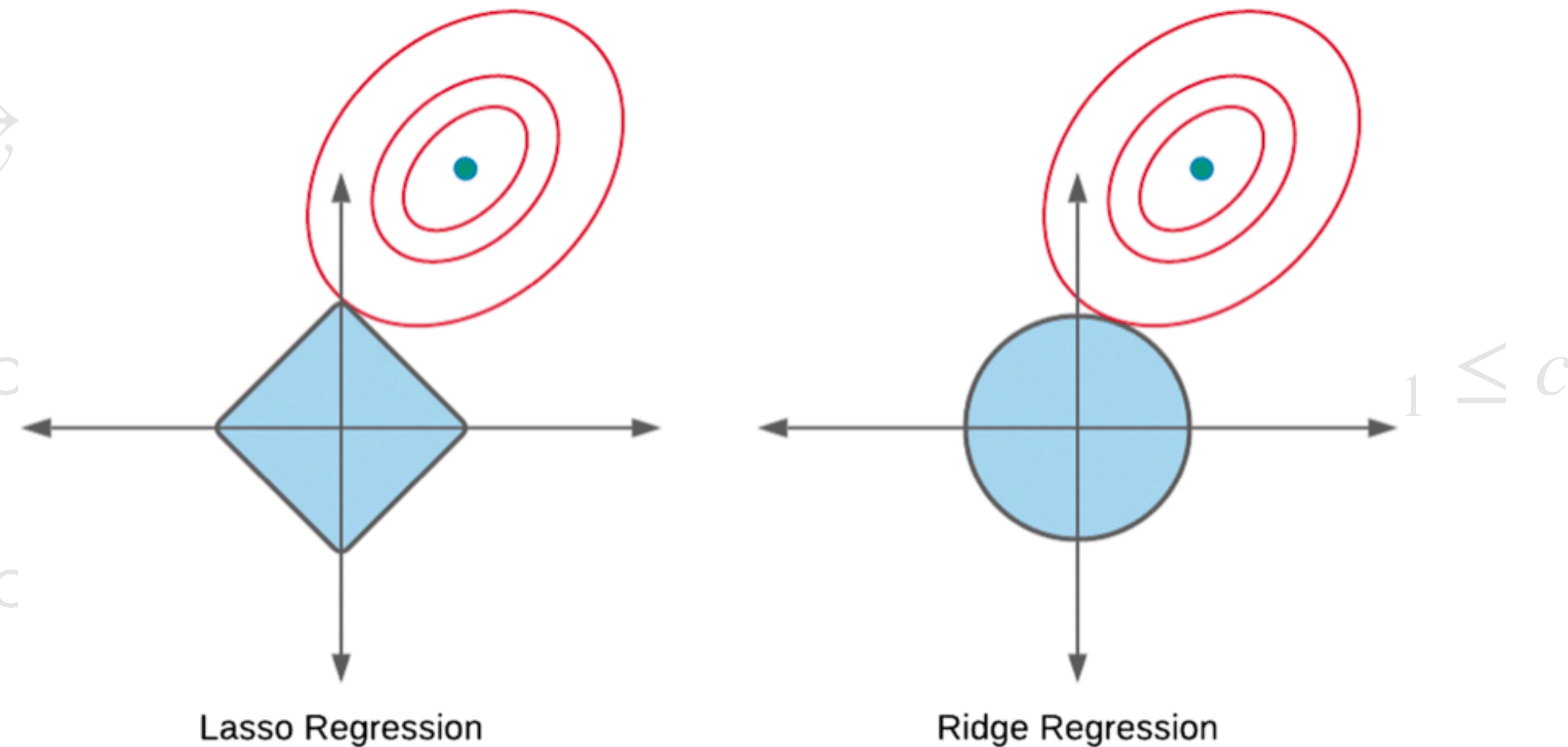
- OLS: Minimize $\|X^T \vec{w} - Y\|^2$

- Solution: \vec{w}

- Ridge Regression

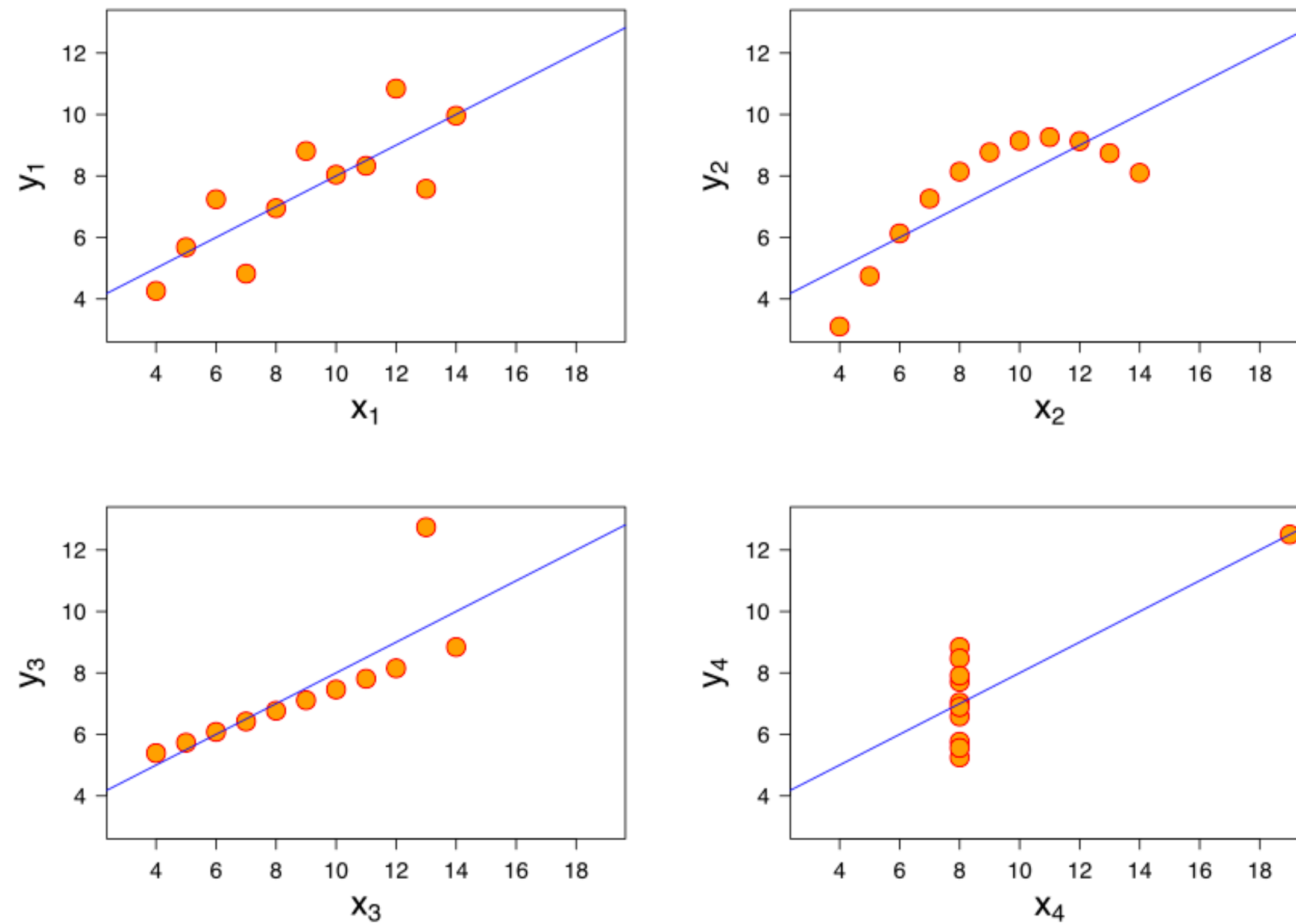
- Loss function

- Solution: *From numerical methods.*



Pathological cases

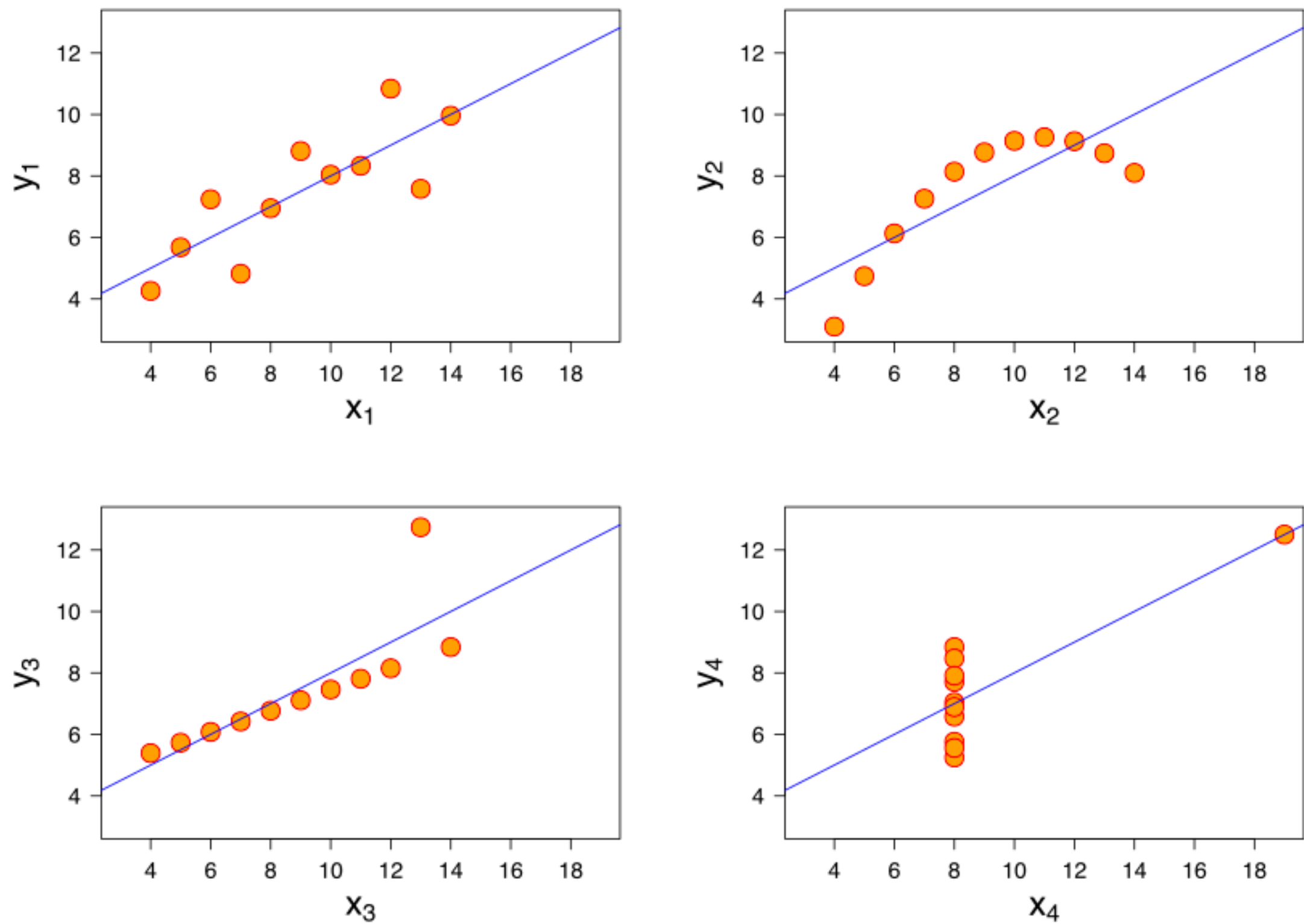
- Anscombe's quartet



https://en.wikipedia.org/wiki/Anscombe%27s_quartet

Pathological cases

- Anscombe's quartet

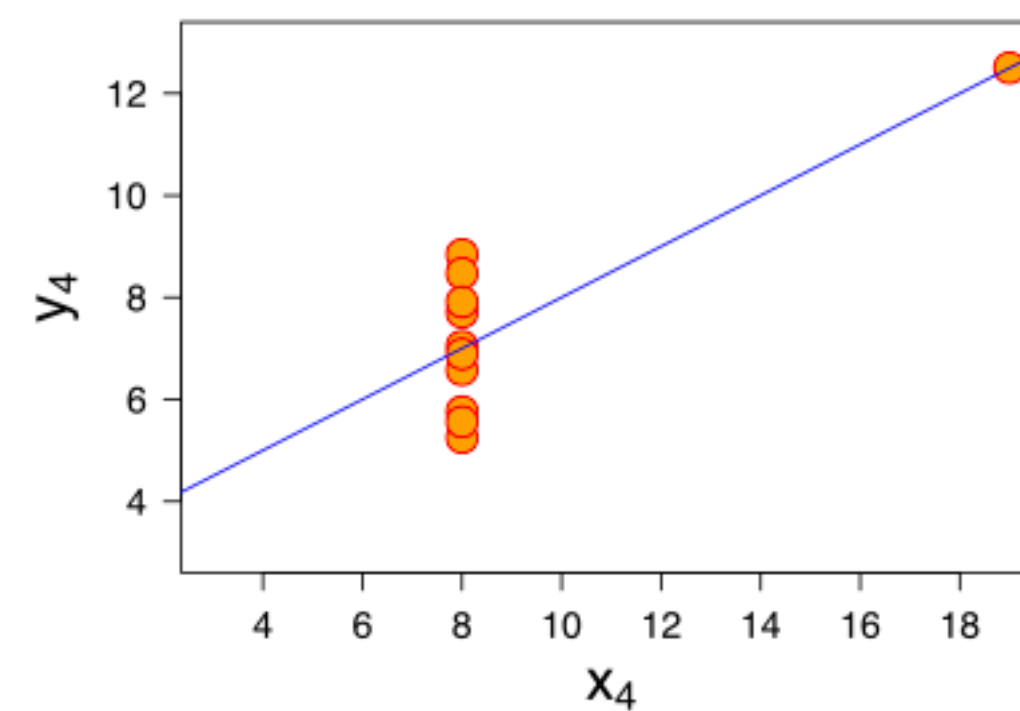
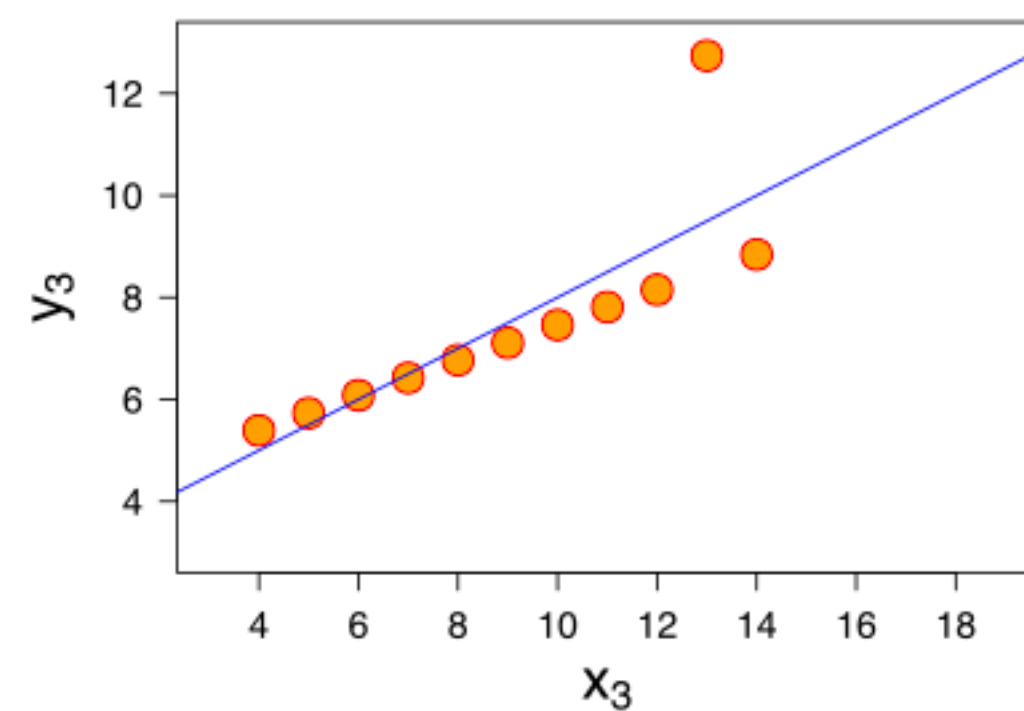
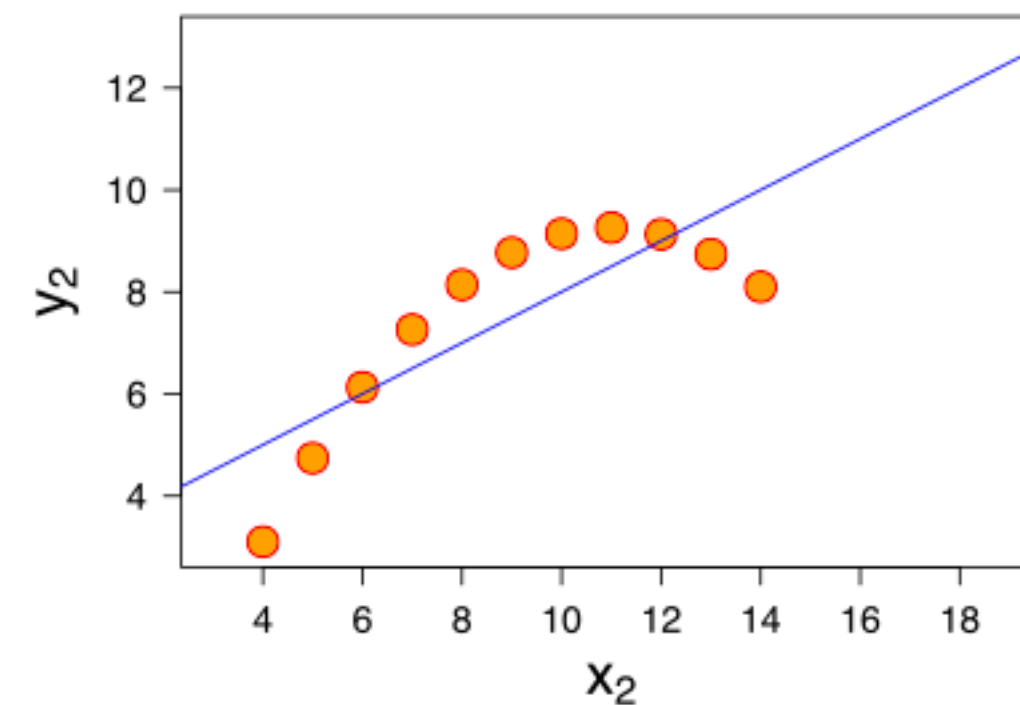
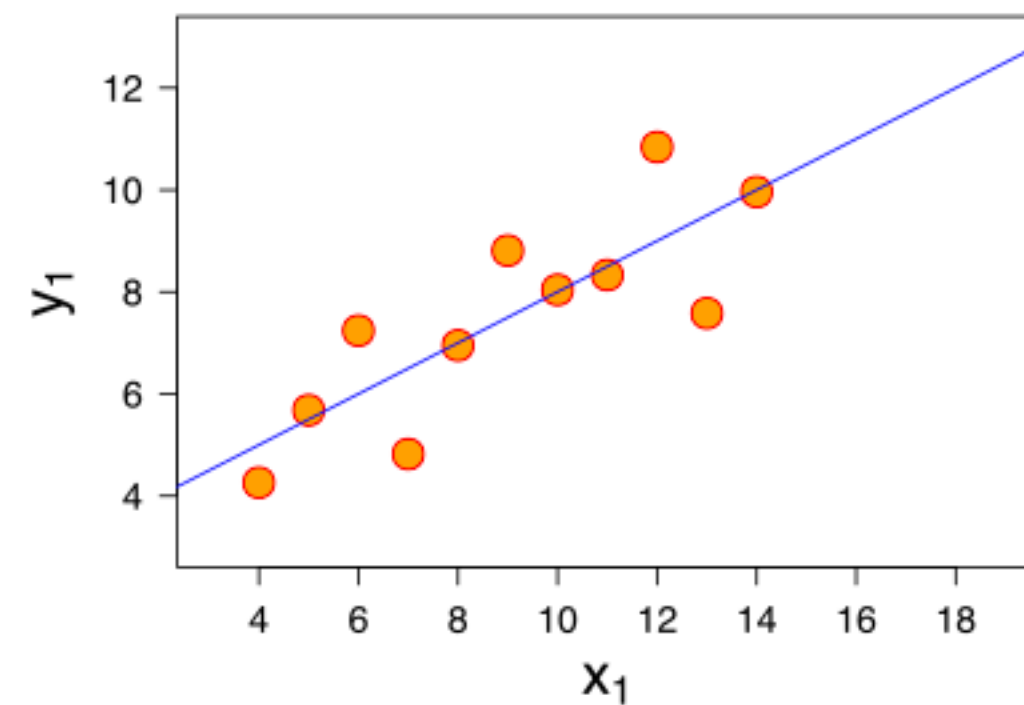


Property	Value	Accuracy
Mean of x	9	exact
Sample variance of $x : s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of $y : s_y^2$	4.125	± 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

https://en.wikipedia.org/wiki/Anscombe%27s_quartet

Pathological cases

- Anscombe's quartet



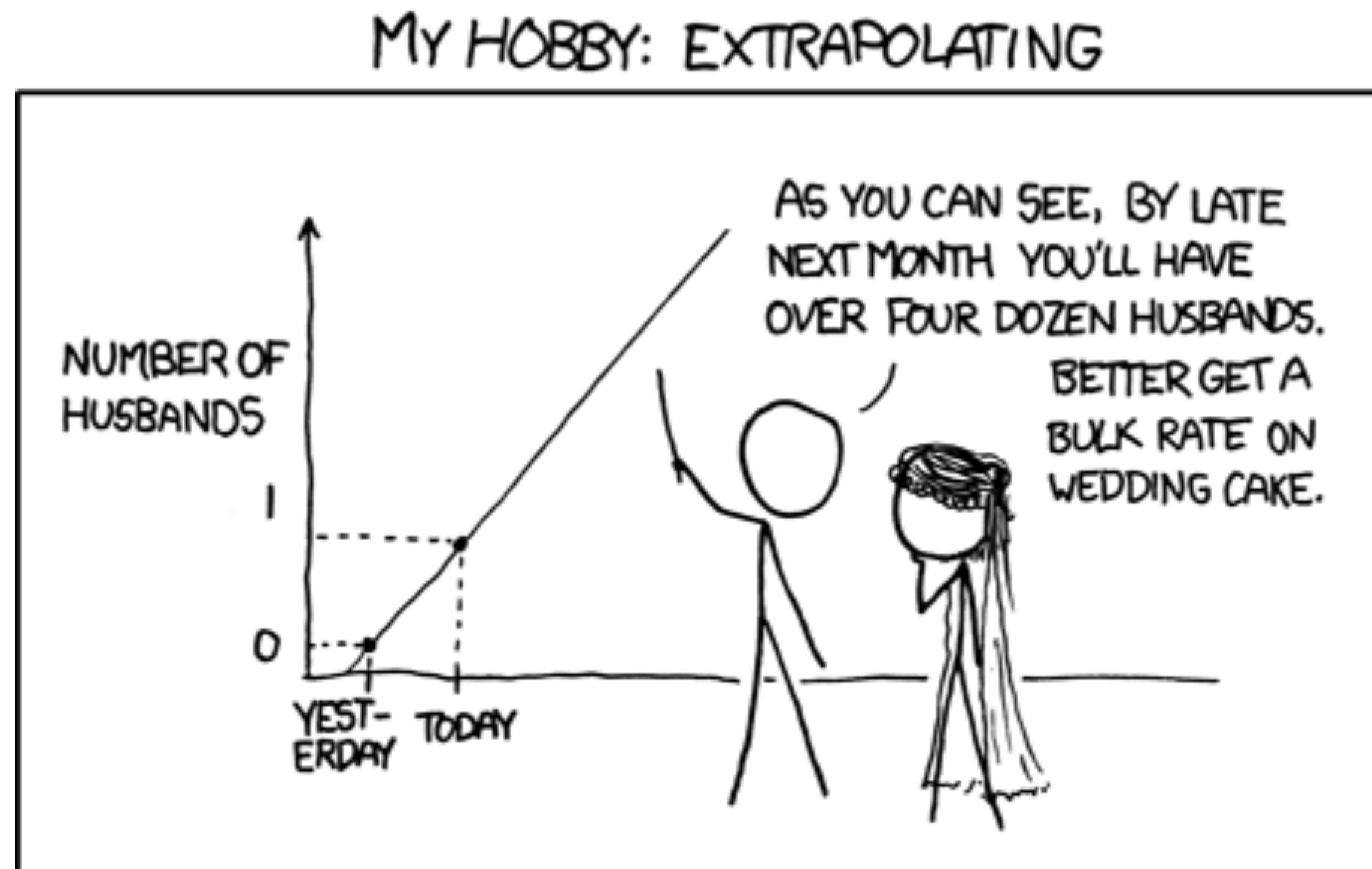
Property	Value	Accuracy
Mean of x	9	exact
Sample variance of $x : s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of $y : s_y^2$	4.125	± 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

Takeaway: Visualize your data!!

https://en.wikipedia.org/wiki/Anscombe%27s_quartet

Pathological cases

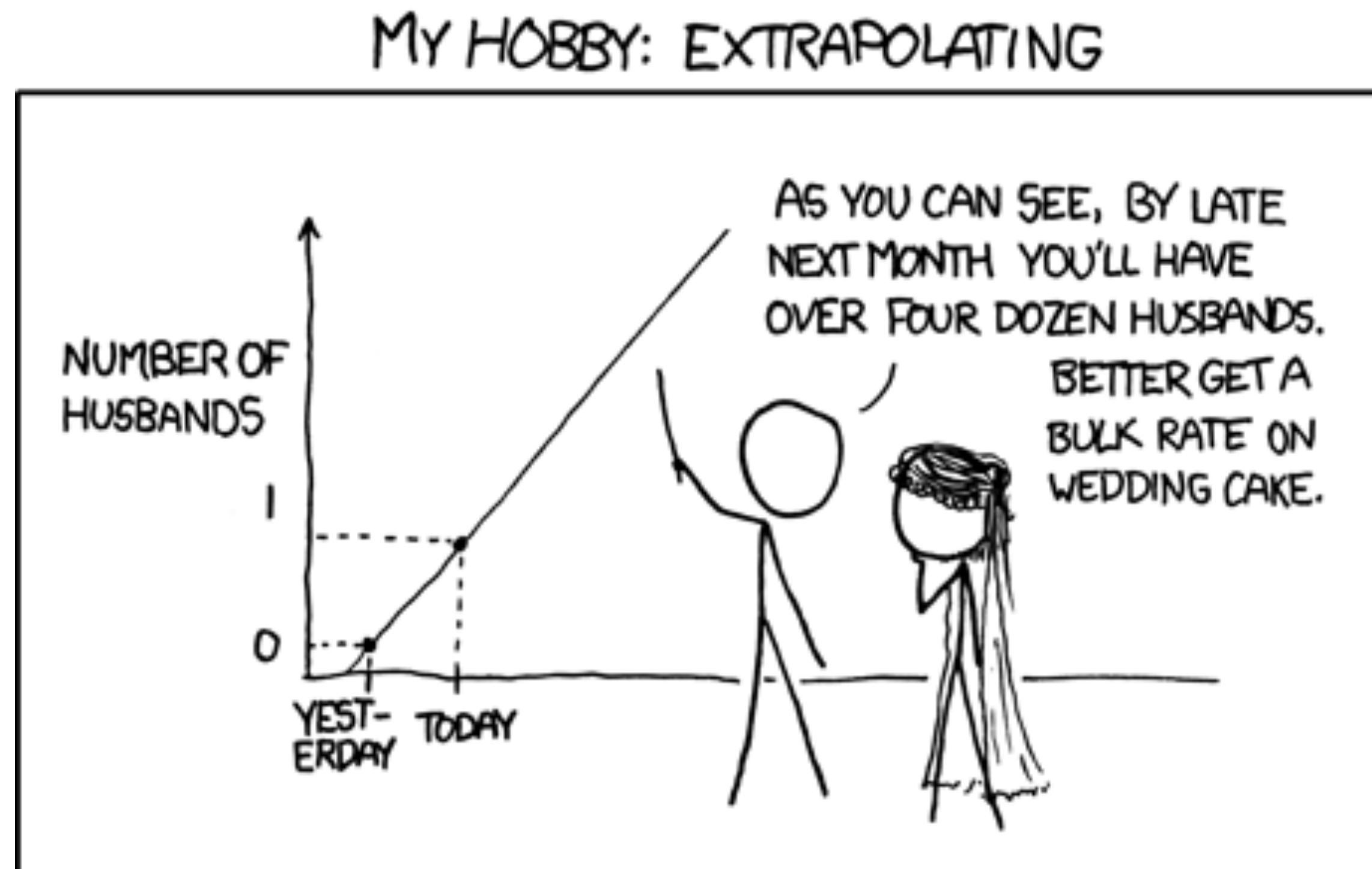
- Extrapolation



<https://xkcd.com/605/>

Pathological cases

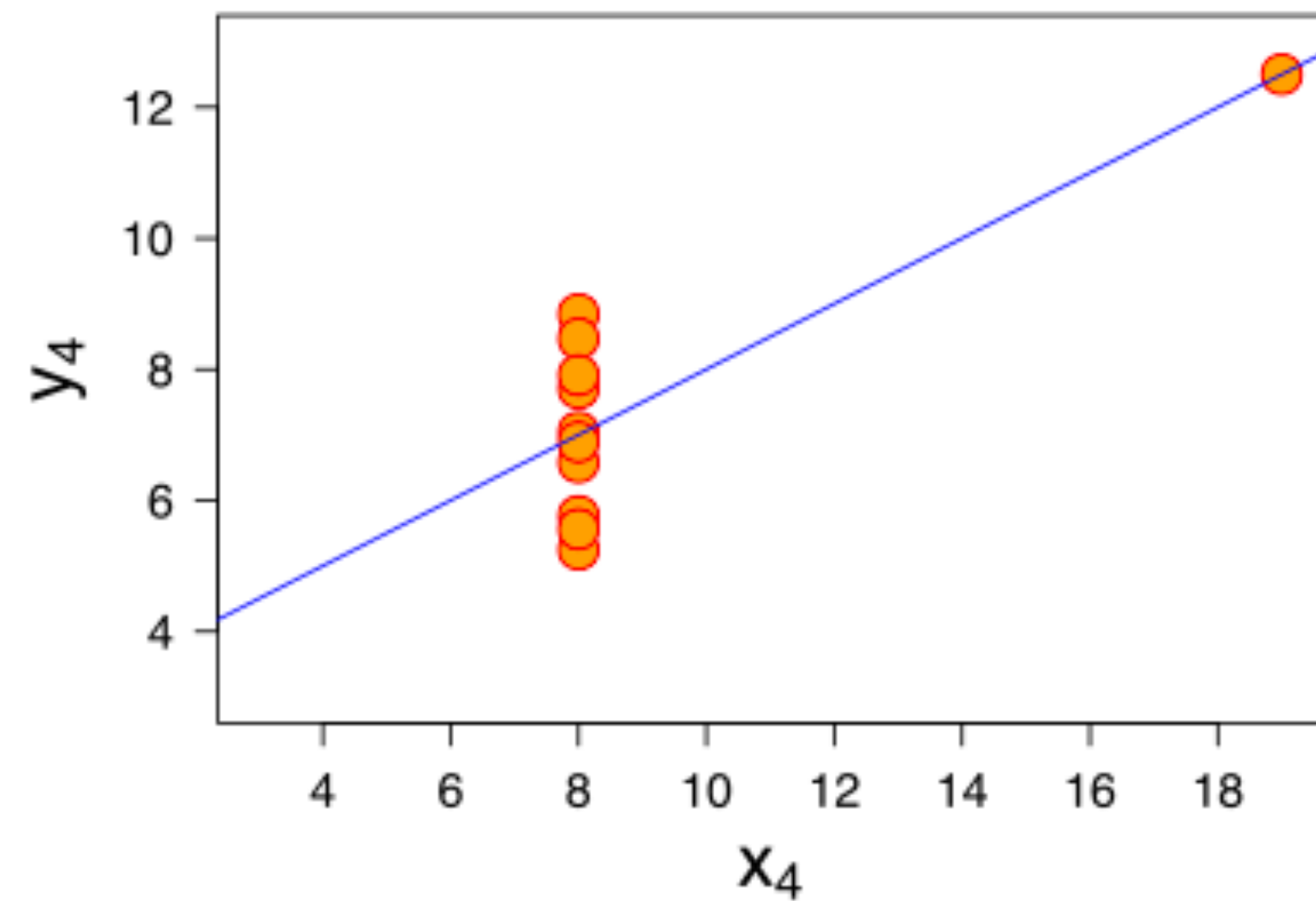
- Extrapolation



Takeaway: Don't trust your model outside the training data distribution

Pathological cases

- Noise in data



Pathological cases

- Large feature size

Additional Reading

- <https://www.mit.edu/~6.s085/notes/lecture3.pdf>
- https://www.coconino.edu/resources/files/pdfs/academics/sabbatical-reports/kate-kozak/chapter_10.pdf

Questions?