

## LECTURE 16

# Feature Engineering

Transforming numerical features and encoding categorical data in order to build more sophisticated models.

**Data 100/Data 200, Fall 2021 @ UC Berkeley**

Fernando Pérez and Alvin Wan

(content by Alvin Wan, John DeNero, Joseph E. Gonzalez, Josh Hug)



## LECTURE 16

# Where a Linear Model Struggles

Understanding the downfalls of a linear model and where feature engineering is needed

**Data 100/Data 200, Fall 2021 @ UC Berkeley**

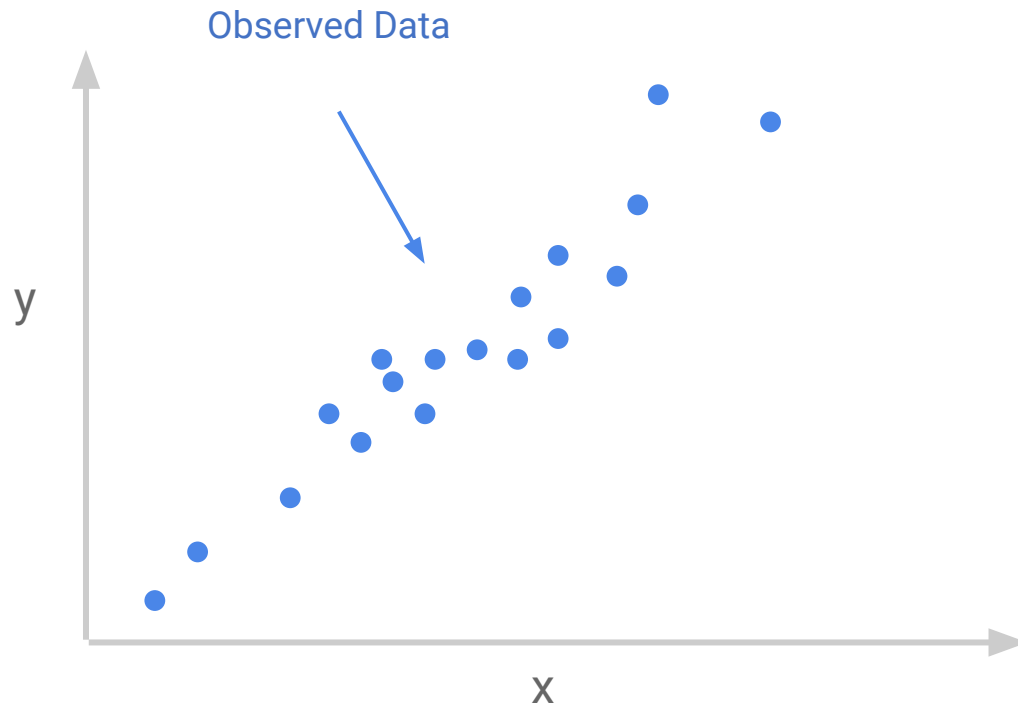
Fernando Pérez and Alvin Wan

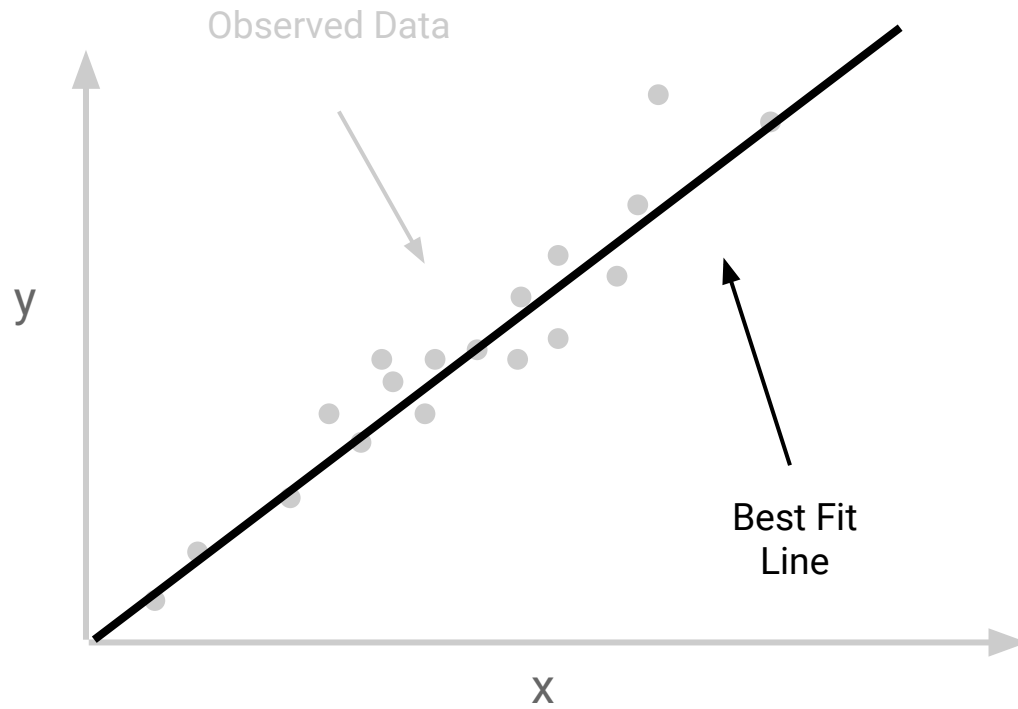
(content by Alvin Wan, John DeNero, Joseph E. Gonzalez, Josh Hug)

# Recap Linear Models

Feature Engineering

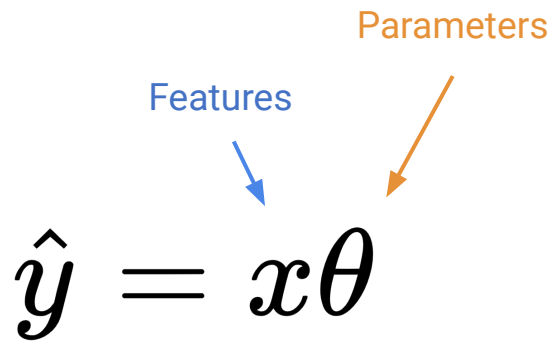
Feature Functions

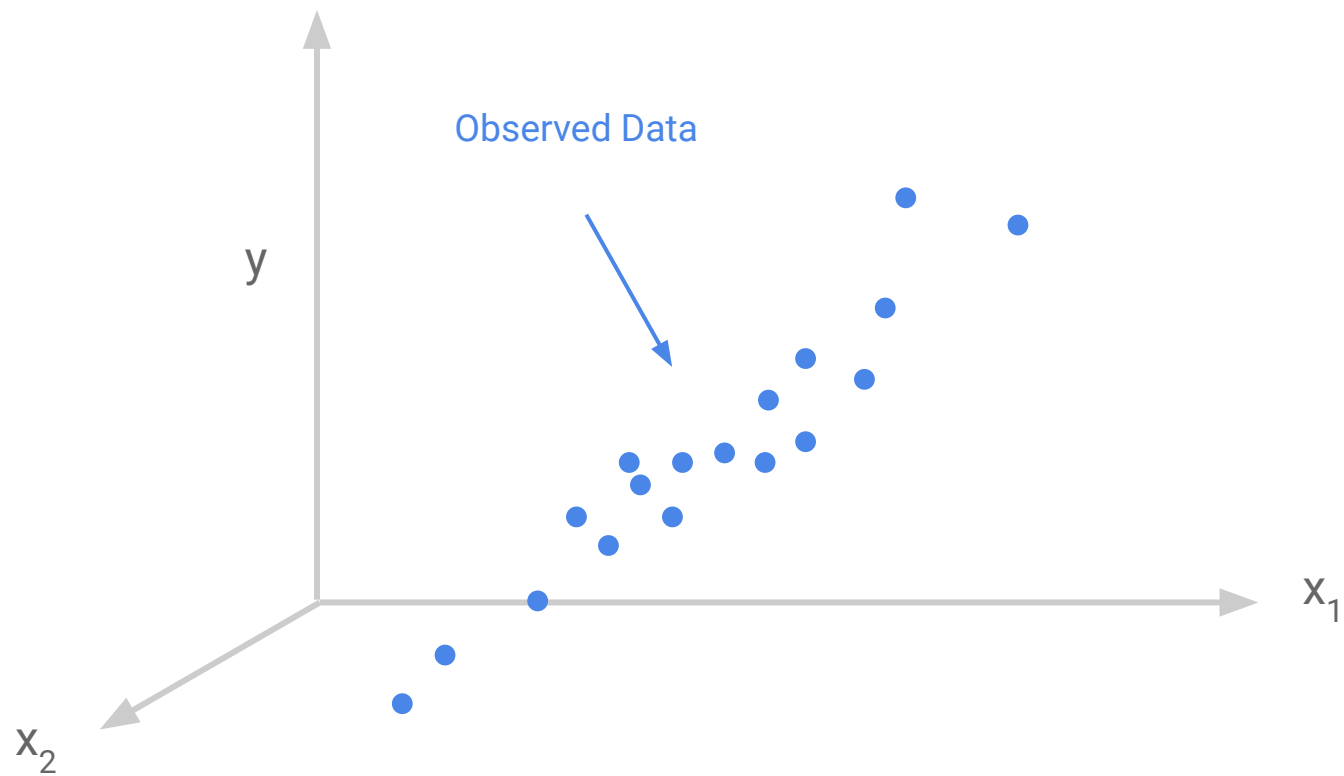




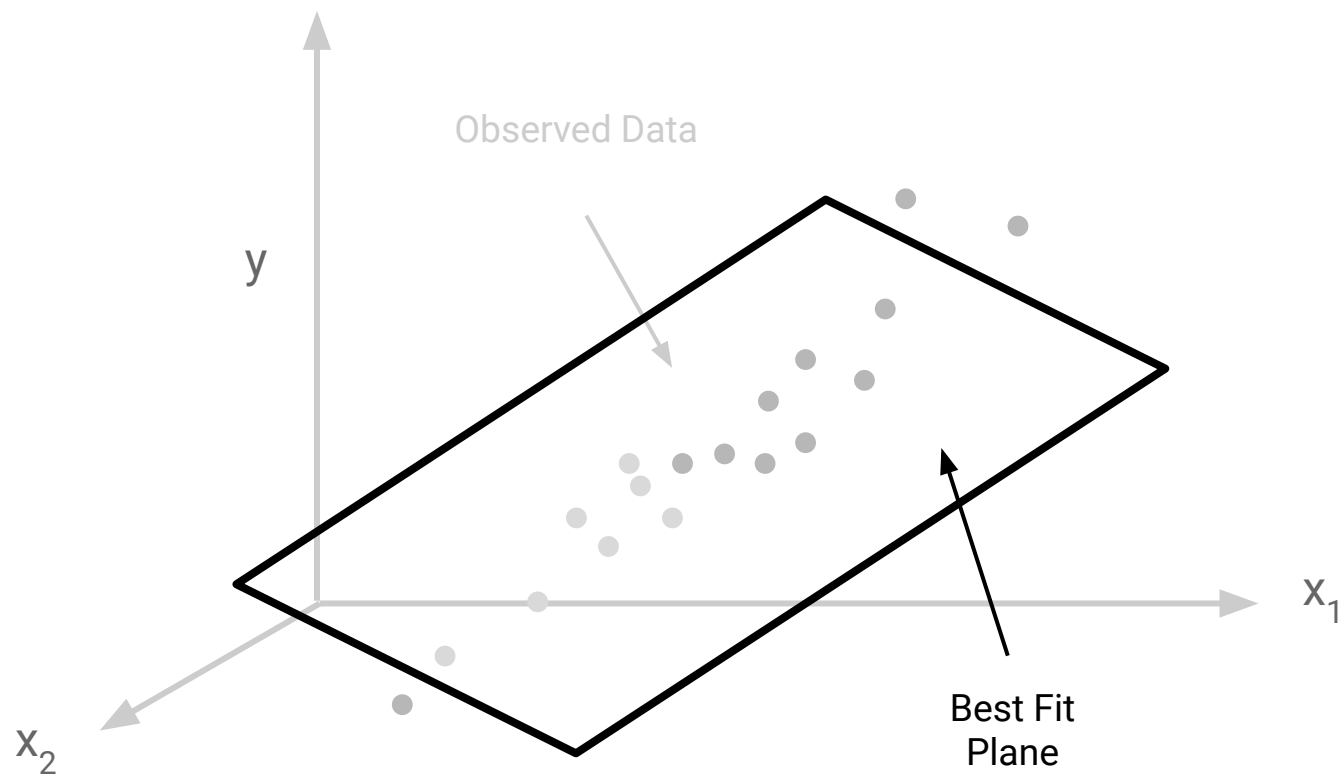
Parameters

Features

$$\hat{y} = x\theta$$








Features

Parameters

$$\hat{y} = x_1 \theta_1 + x_2 \theta_2$$


The diagram illustrates the components of the linear equation  $\hat{y} = x_1 \theta_1 + x_2 \theta_2$ . A blue label 'Features' has two blue arrows pointing to  $x_1$  and  $x_2$ . An orange label 'Parameters' has two orange arrows pointing to  $\theta_1$  and  $\theta_2$ .

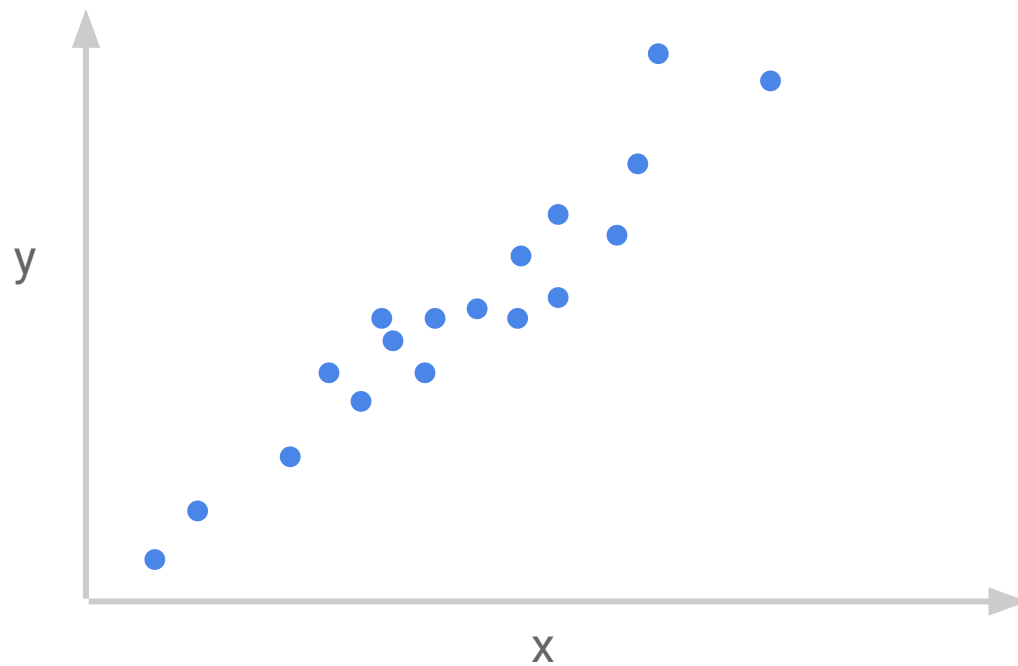
Re write

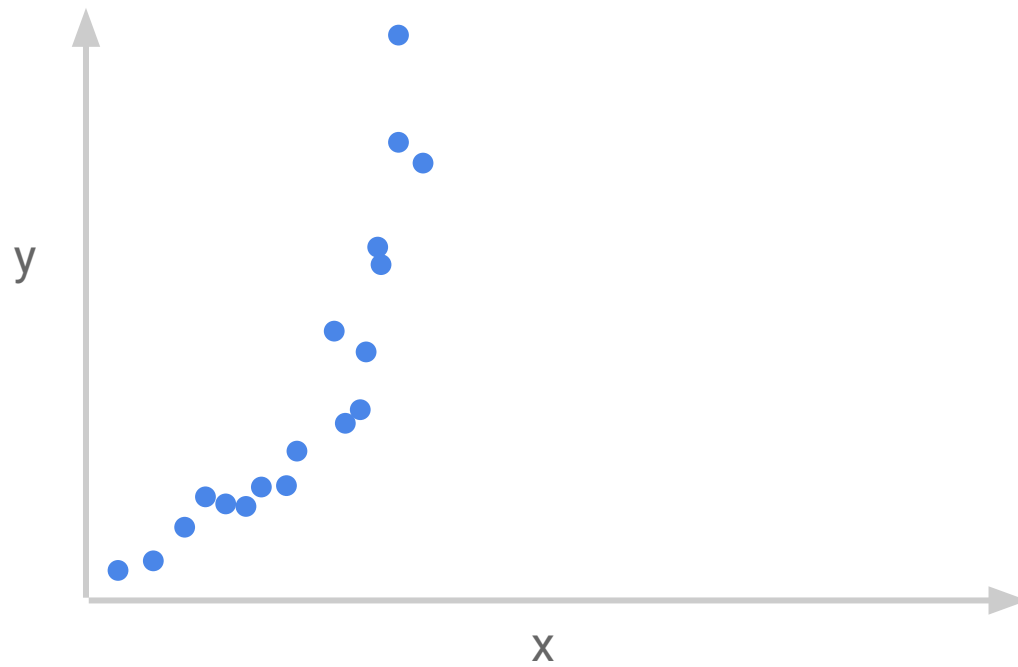
$$\hat{y} = \sum_{j=1}^d x_j \theta_j$$

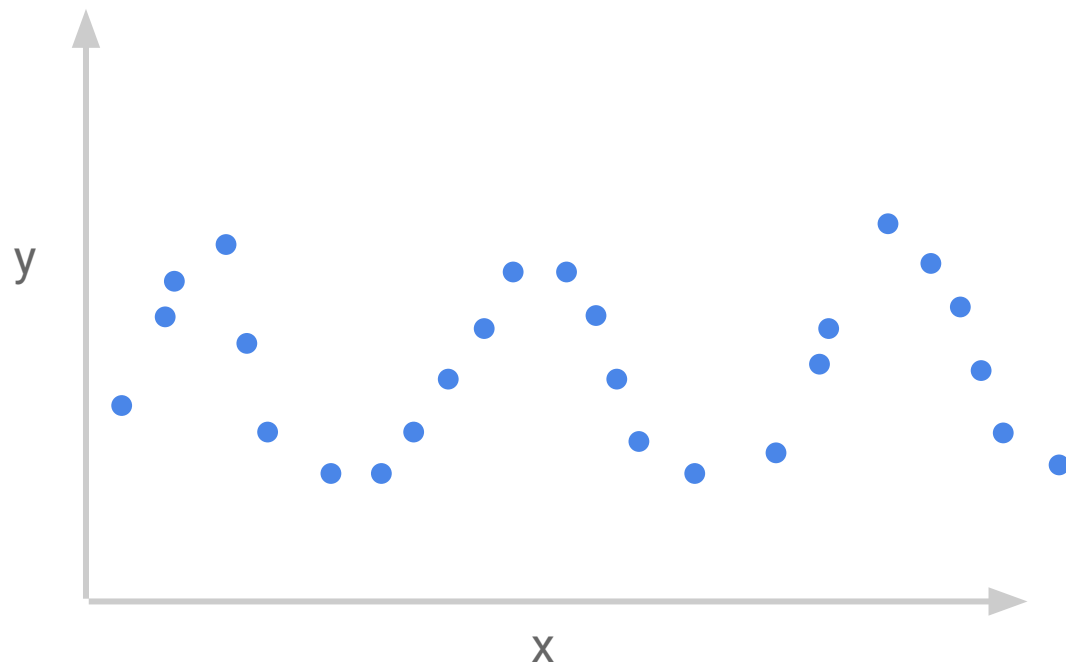
Features

Parameters









What if the relationship  
between  $y$  and  $x_i$  is not linear?

What if the relationship  
between  $y$  and  $x_i$  is not linear?

What if our features  $x_i$  are not  
numbers?



Recap Linear Models

# Feature Engineering

Feature Functions

## TAKEAWAY

Feature engineering transforms **raw** features into more **informative** features for modeling.

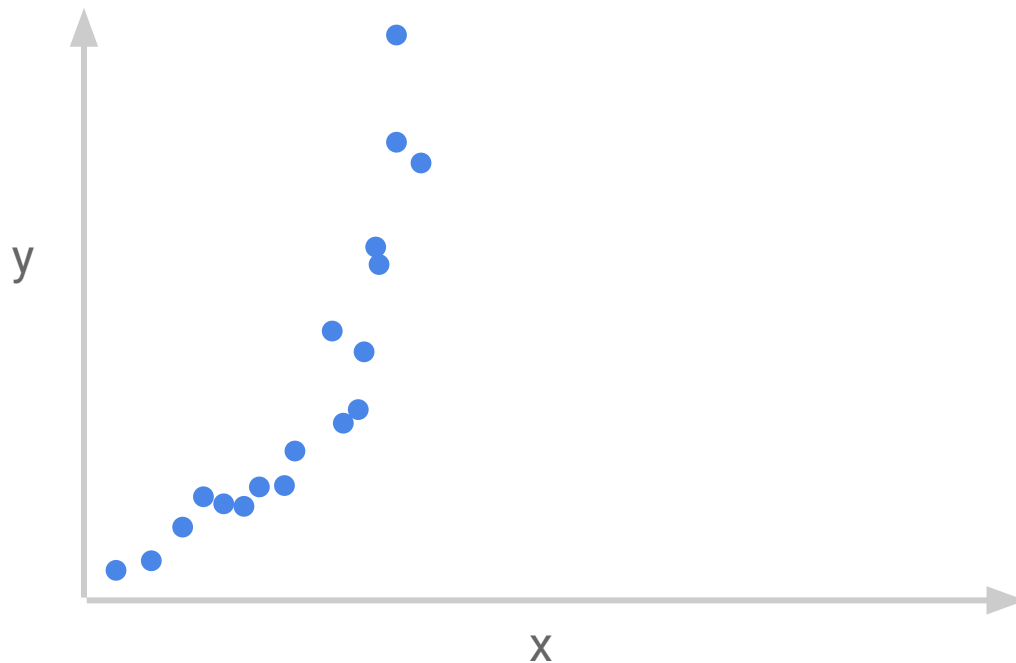
## TAKEAWAY

Feature engineering enables you to **express non-linear relationships, capture domain knowledge, and encode non-numeric features.**

Recap Linear Models  
Feature Engineering  
Feature Functions

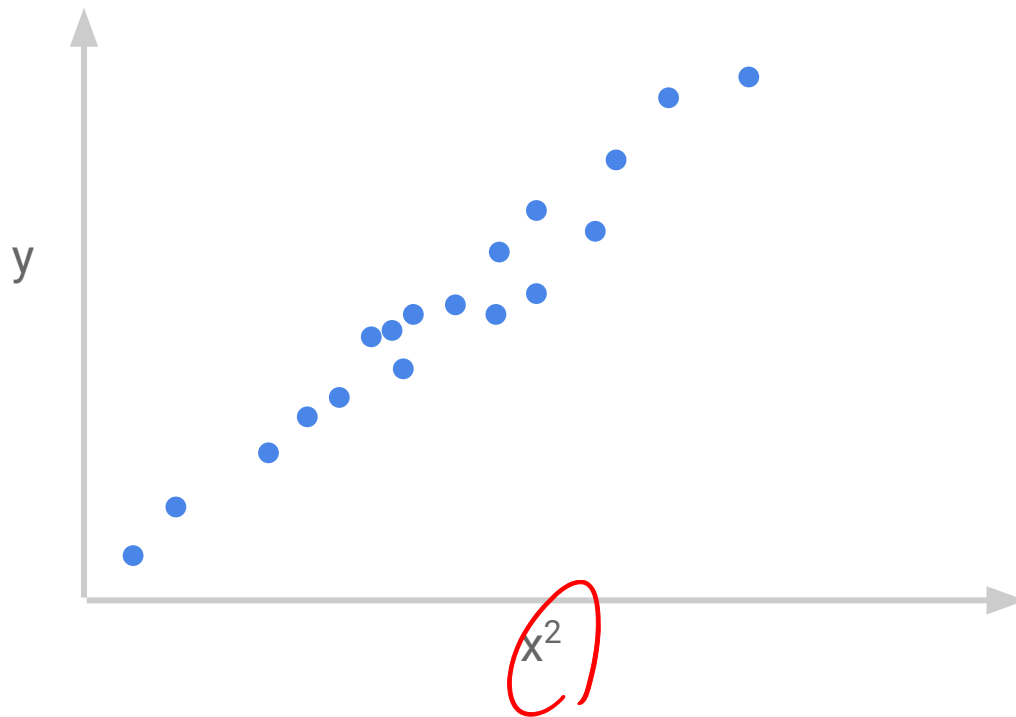
# Non-linear relationship

between x and y



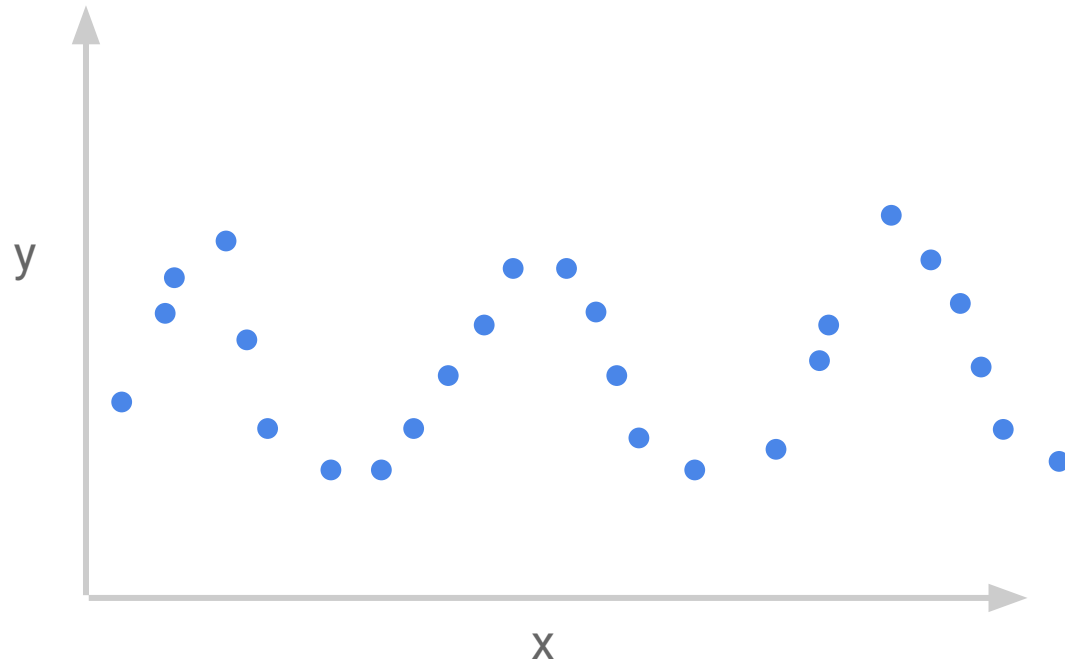
# Linear relationship

between  $x^2$  and  $y$



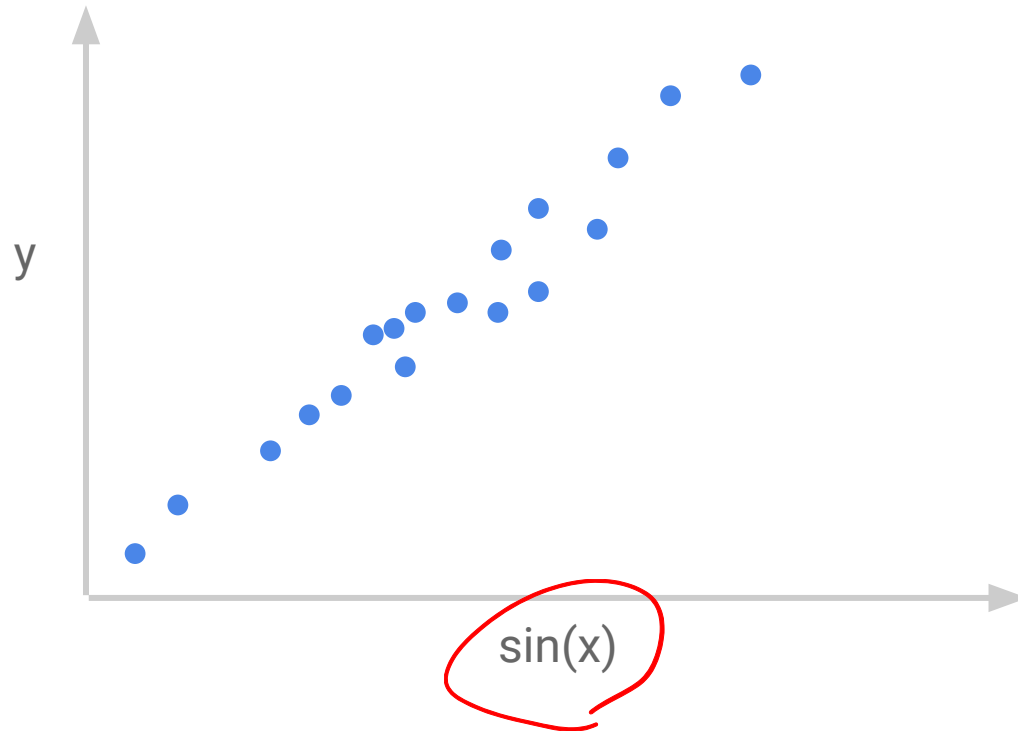
# Non-linear relationship

between x and y



# Linear relationship

between  $\sin(x)$  and  $y$





$$\hat{y} = \sum_{j=1}^d \overbrace{x_j}^{\text{d Features}} \underbrace{\theta_j}_{\text{Parameters}}$$

$$\hat{y} = \sum_{j=1}^p \overbrace{\phi(x)_j}^{\text{p Features}} \underbrace{\theta_j}_{\text{Parameters}}$$

(phi)ture function  
- Prof. Gonzalez, Sp21

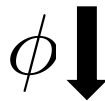
$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$$

new dimensional  
space

Feature  
Engineering

## Non-numeric and Raw Values

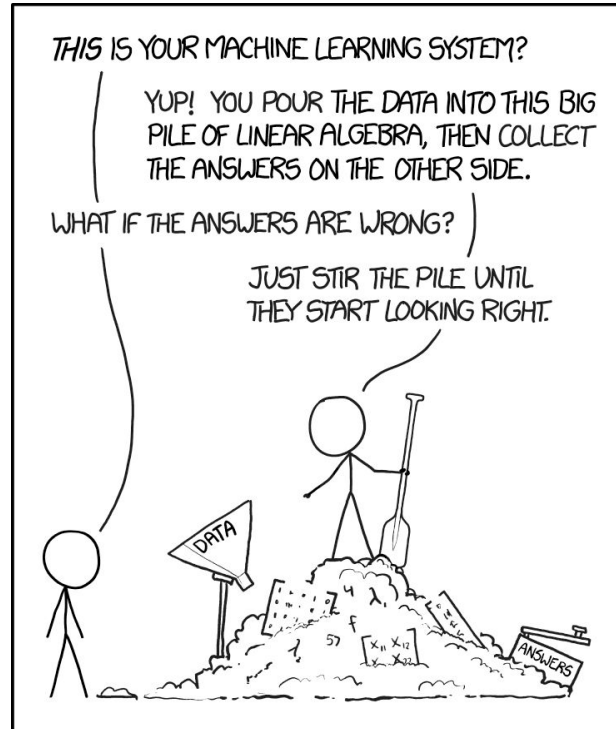
uid	age	state	hasBought	review
0	32	NY	True	"Meh."
42	50	WA	True	"Worked out of the box ..."
57	16	CA	NULL	"Hella tots lit..."



## Entirely Quantitative and Transformed Values

AK	...	NY	...	WY	age	age^2	hasBought missing
0	...	1	...	0	32	32^2	0
0	...	0	...	0	50	50^2	0
0	...	0	...	0	16	16^2	1

**Designing feature functions** is a big part of machine learning and data science.



np.book



[xkcd.com/1838/](https://xkcd.com/1838/)



np. book



## LECTURE 16

# Conclusion

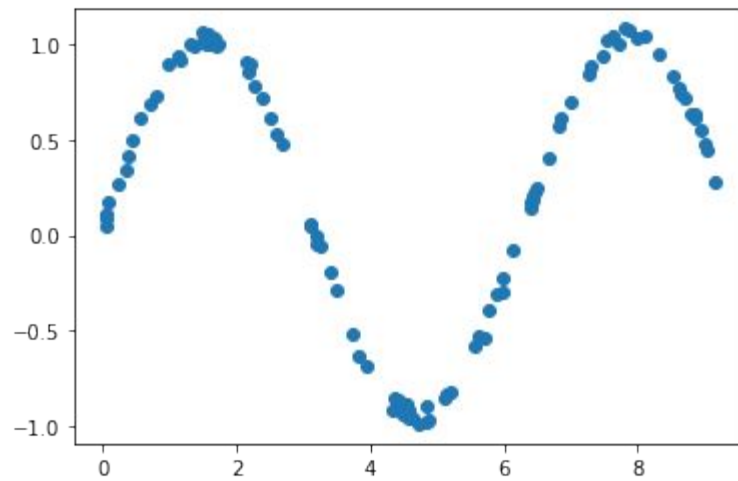
Gotchas for feature engineering. Warnings and pitfalls to be aware of.

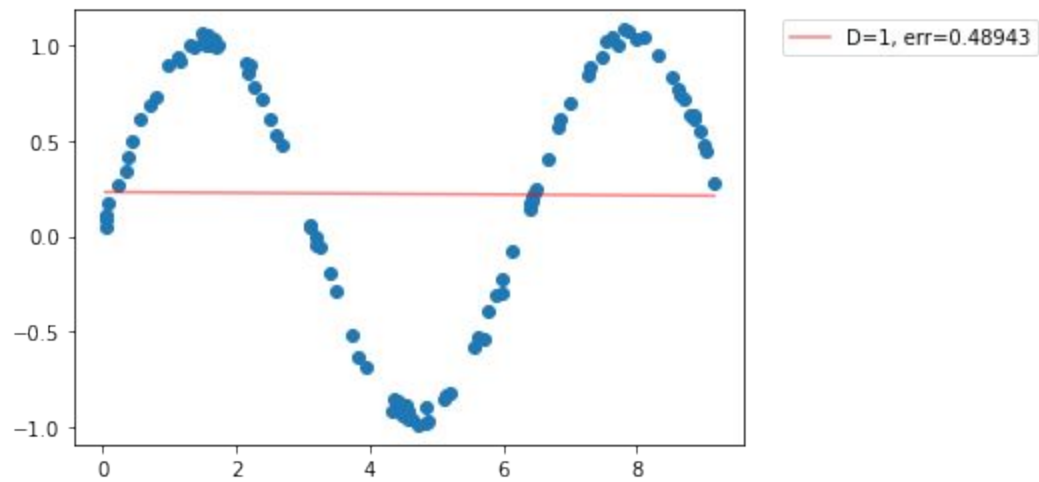
**Data 100/Data 200, Fall 2021 @ UC Berkeley**

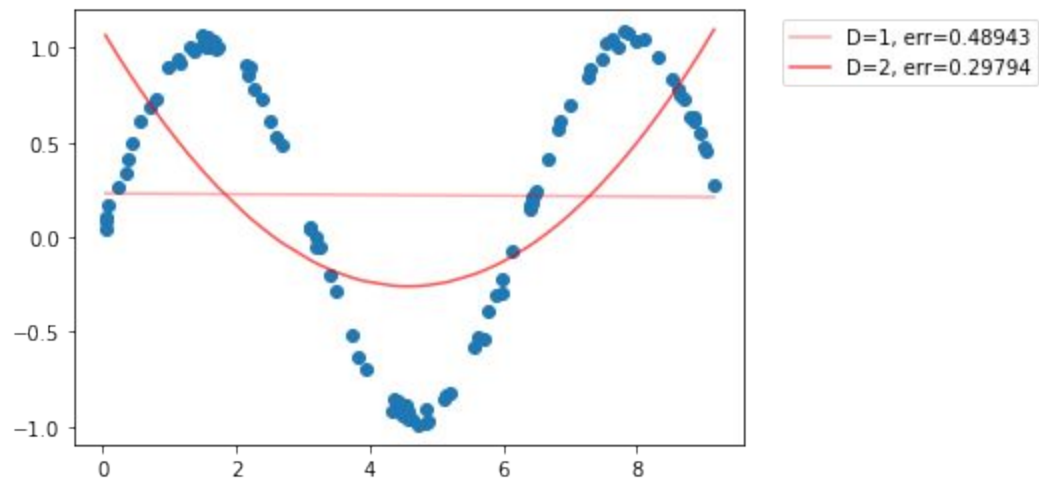
Fernando Pérez and Alvin Wan

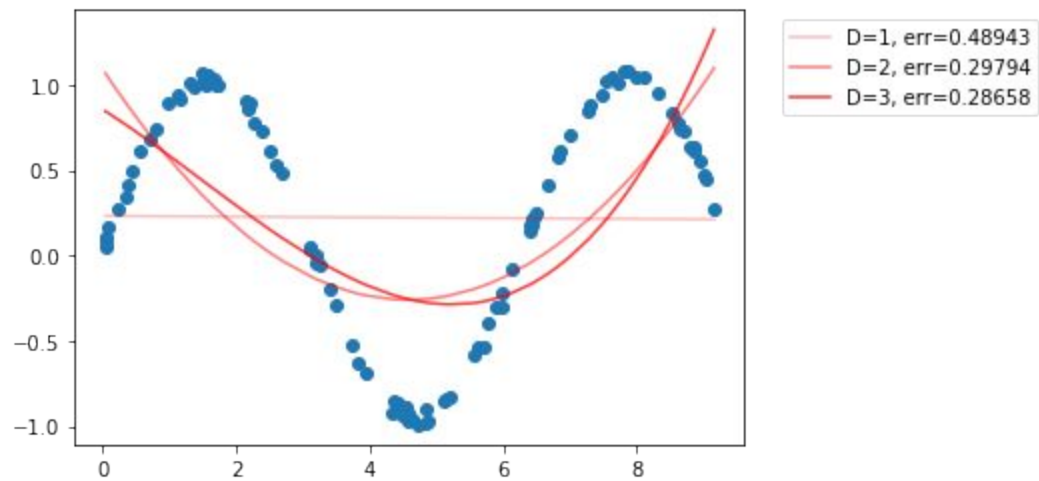
(content by Alvin Wan, John DeNero, Joseph E. Gonzalez, Josh Hug)

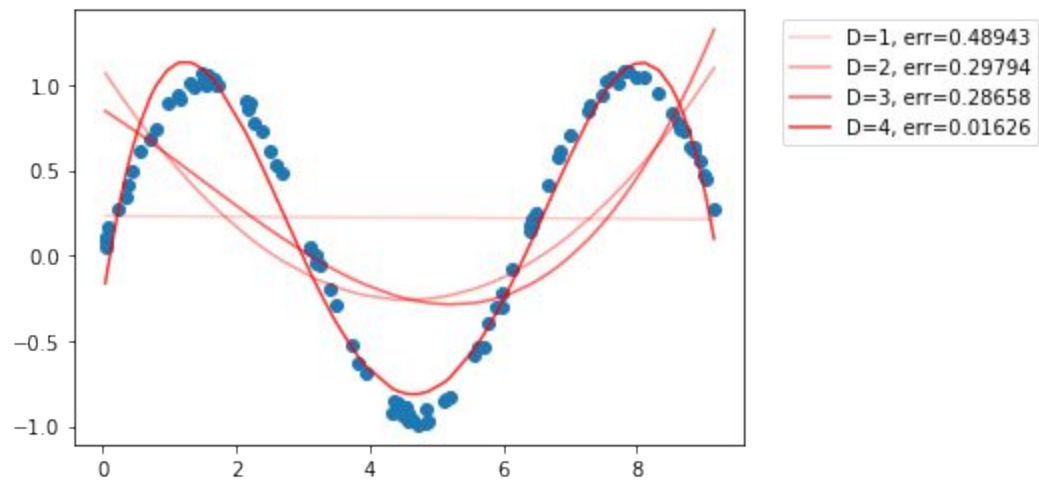


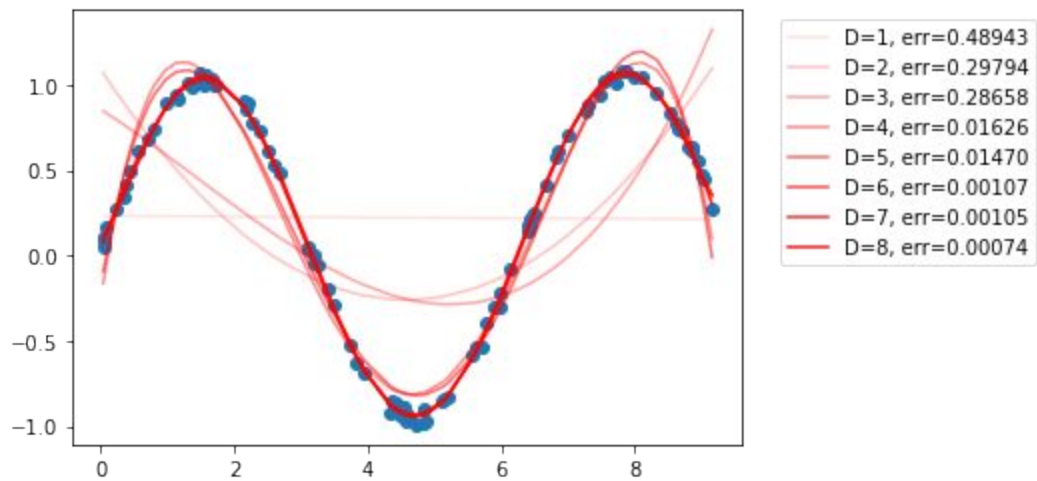


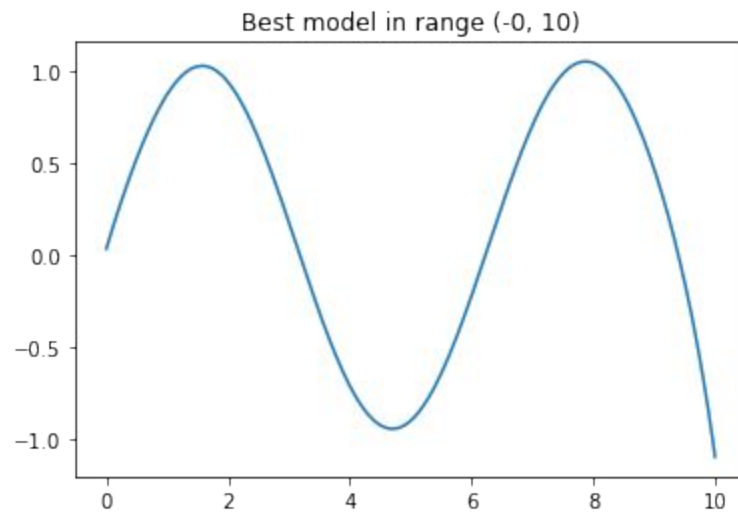


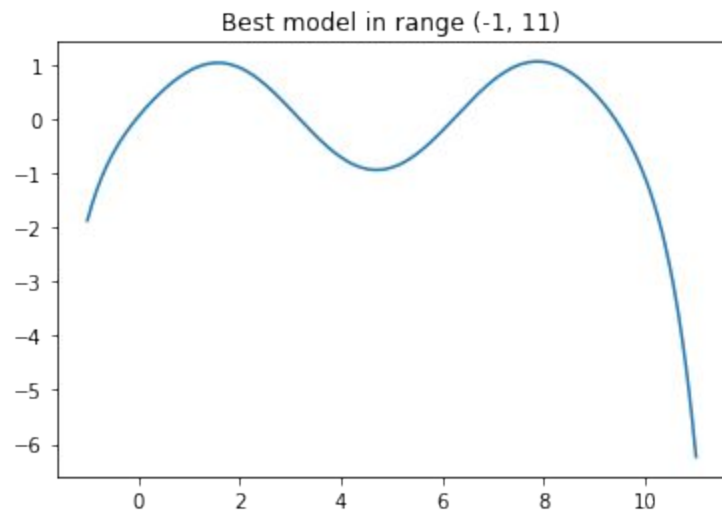




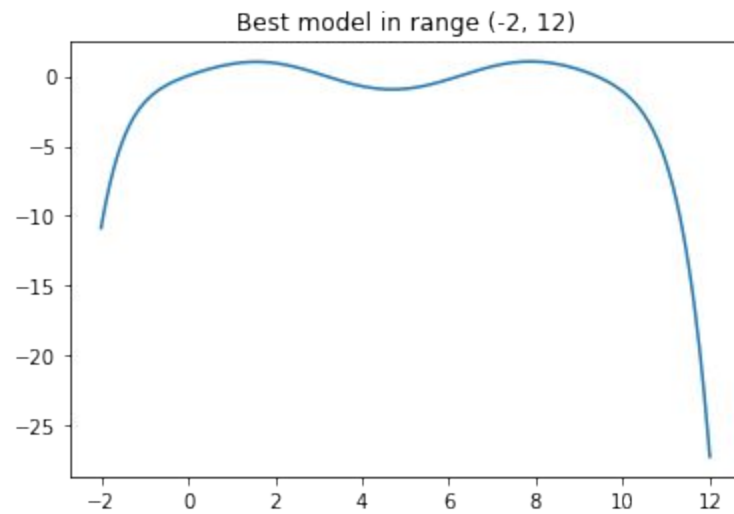


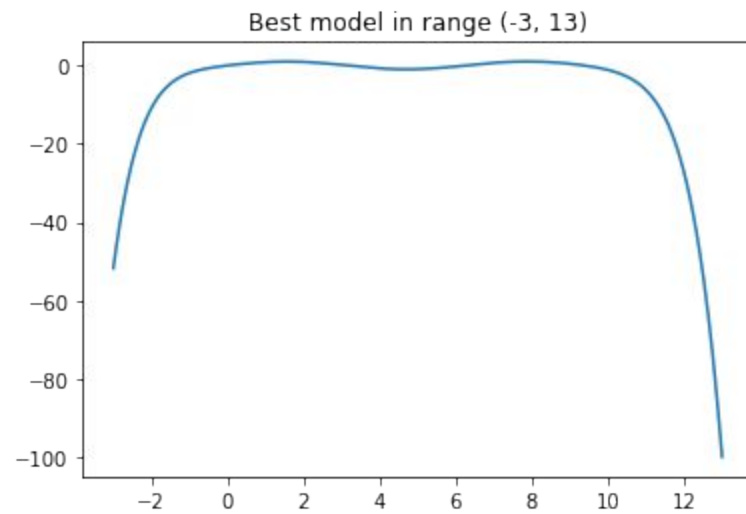


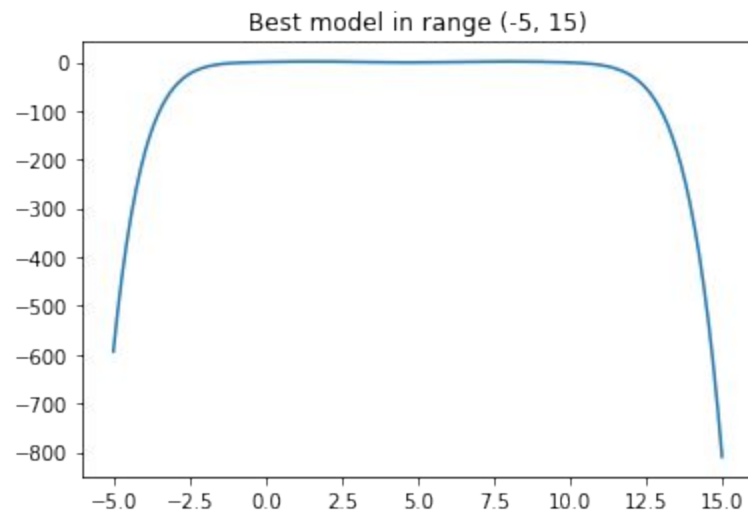












## TAKEAWAY

Overfitting occurs when your model **memorizes noise**, leading to **poor generalization** beyond training data.

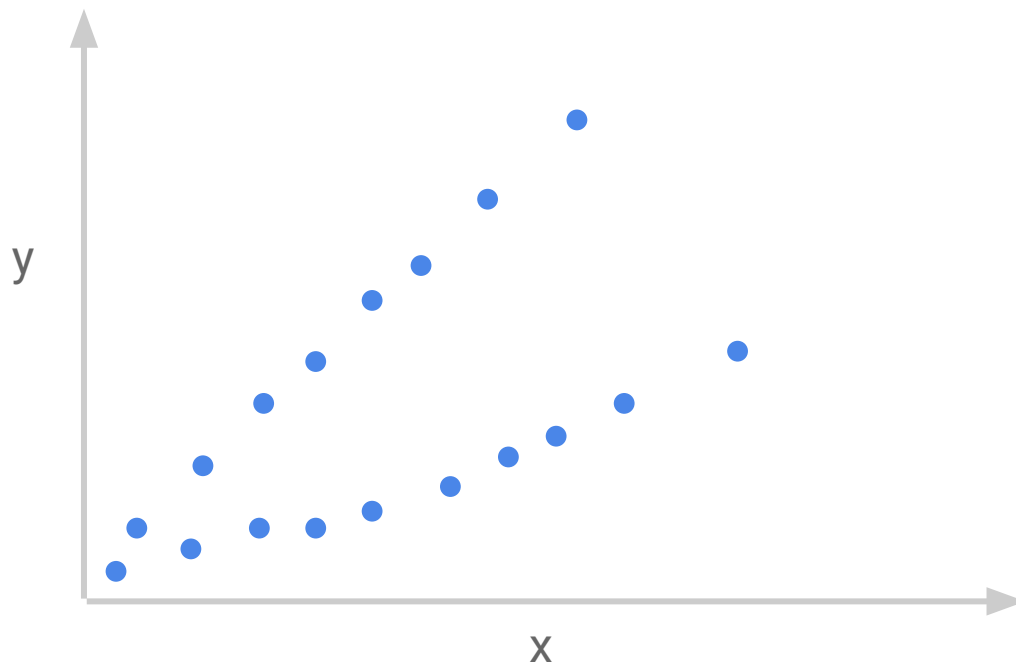
Next Up

# **Regularization**

## **Bias-Variance Tradeoff**

# Domain Knowledge

poor fit on  $(x, y)$



# Domain Knowledge

great fit on  $(x^{t+1}-x^t, y)$

