

Day-23, Oct 25, 2024

- Continuing Deep learning Foundations
- Overfitting problems & possible remedies

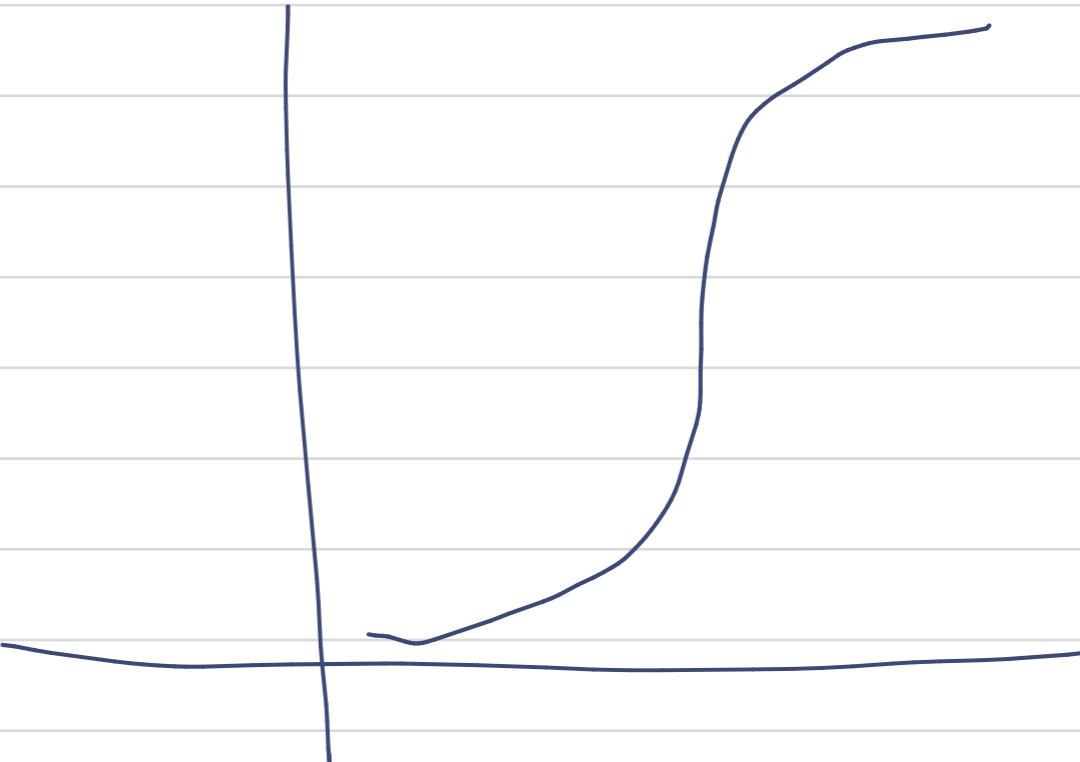
## Deep learning Foundations

- ① Sigmoid Functions motivations & usage
- ② Optimization in logistic Regression
- ③ Methods to Compute Gradient Descend
- ④ Example of Chain Rule & Composite Derivative function.

## TF Sigmoid function

$$g(z) = \frac{1}{1 + e^{-z}}$$

→ ranges [0,1]



Using

$$z = \theta^T x$$

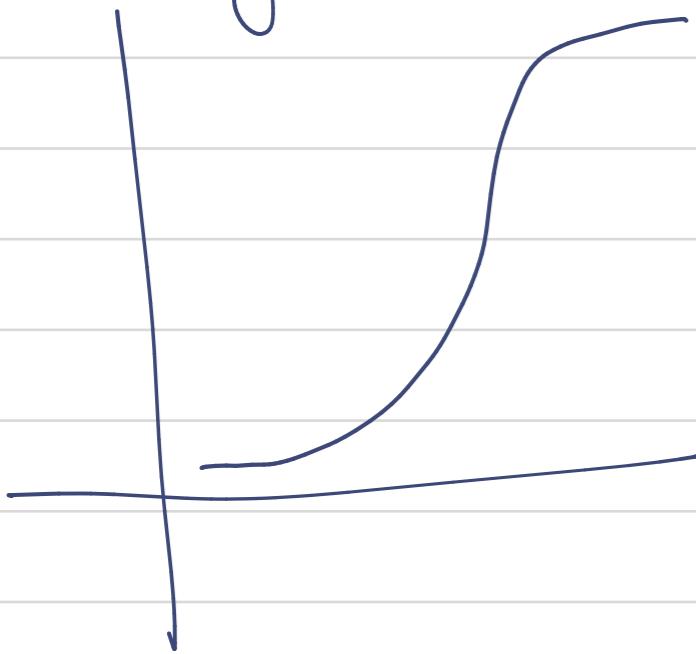
$$g(z) = \frac{1}{1 + e^{-(\theta^T x)}}$$

→ Change  $\theta$  randomly ( $\theta^T$ ) → Transposed  $\theta$

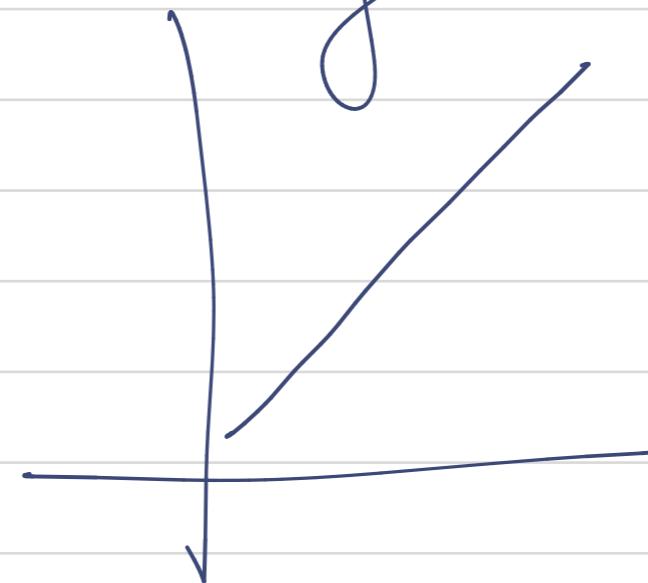
→ Sigmoid shows probabilities of classification

→ There is no nonlinear decision boundary using Sigmoid.

Using  $\theta = 0.13$



Using  $\theta = 0.20$



If our goal must be finding the optimal parameter  $\theta(\theta)$

that maximize the likelihood of given (observed dataset).

So we have loglikelihood function and likelihood function

# Log-likelihood (F): Sum of individual log-likelihood terms for each input.

# Likelihood (F): Product of probabilities of each data point belonging to respective classes.

$(y=1) \log \pi_0(x)$ : Head  $\rightarrow$  term contribute to likelihood of actual class '1'

$(y=0) \log (1 - \pi_0(x))$ : Tail  $\rightarrow$  term contribute to likelihood of actual class '0'.

$$\text{Sol } J_0(x) \Rightarrow \frac{L}{1 + e^{-(\theta^T x)}}$$

Once we find the optimal  $\theta$  it is time for decision boundary that separate two classes.

In Logistic Regression optimization works by iteratively updating  $\theta$  values in the direction of the gradient of the log-likelihood function to maximize it.

## Logistic Regression

- Logistic Regression → Finding good  $\theta$ 's for this likelihood (i.e., the ones that maximize it)

$$\log p_{\theta}(y|x) = \sum_{i=1}^n \mathbb{I}(y_i = 1) \log \pi_{\theta}(x) + \mathbb{I}(y_i = 0) \log (1 - \pi_{\theta}(x))$$

- These are the directions like those in the previous figure that do a good job separating heads from tails

## Sigmoid Function

- Consider the collection of functions (one for each  $\theta$ )

$$\pi_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)} \text{ defined } \sigma(\theta^T x)$$

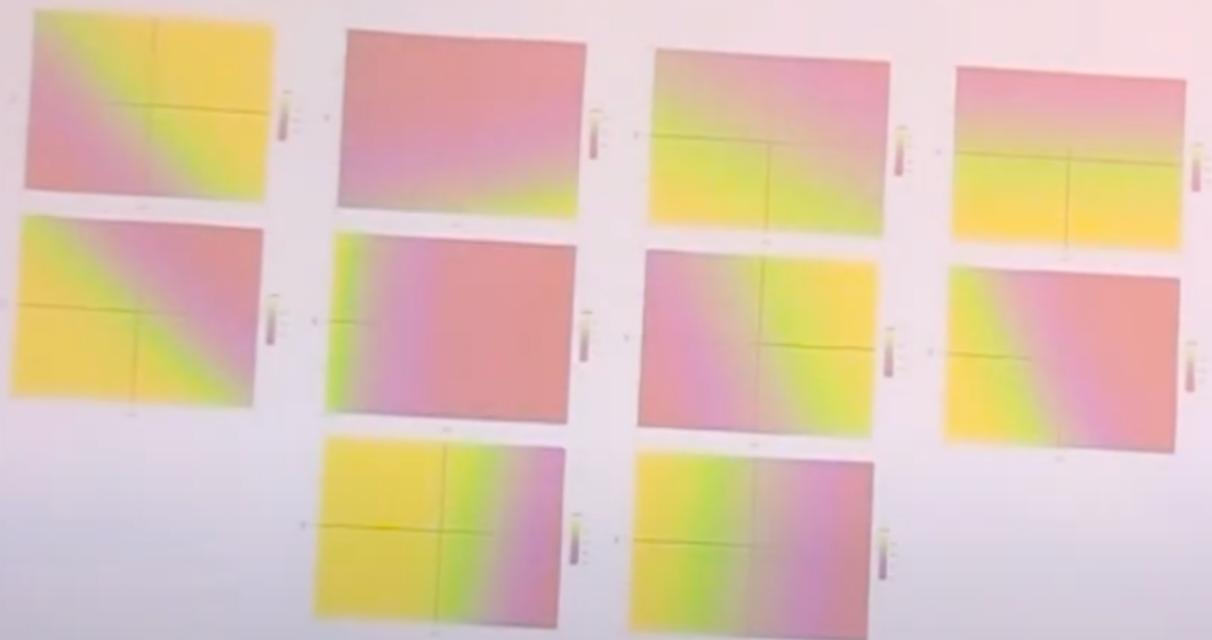


Figure: Example sigmoid functions  $\pi_{\theta}$  for random draws of  $\theta$ , when we assume  $x$  has an intercept term. Notice that there are no nonlinear boundaries...

$$\log p_{\theta}(y|x) \Rightarrow \sum_{i=1}^n \mathbb{I}(y_i = 1) \log \pi_{\theta}(x) + \mathbb{I}(y_i = 0) \log(1 - \pi_{\theta}(x))$$

$$\pi_{\theta}(x) \Rightarrow \frac{1}{1 + e^{-\theta^T x}} \Rightarrow \sigma(\theta^T x)$$

Courageously differentiate

①  $\frac{\partial}{\partial \theta} \log p_{\theta}(y|x)$

②  $\Rightarrow$  Sum of log-likelihood terms for each point

③  $\Rightarrow$  Then the derivative of log-likelihood is calculated using Chain Rule and the derivative of Sigmoid function

## Optimization

- ▶ Need to find some  $\theta$ 's that maximize  $\log p_{\theta}(y|x)$
- ▶ Use gradient descent on  $-\log p_{\theta}(y|x)$
- ▶ How to get the gradients?
- ▶ Approach 1: Courageously differentiate,

$$\begin{aligned}& \frac{\partial}{\partial \theta} \log p_{\theta}(y|x) \\&= \frac{\partial}{\partial \theta} \left[ \sum_{i=1}^n y_i \log \sigma(\theta^T x_i) + (1 - y_i) \log (1 - \sigma(\theta^T x_i)) \right] \\&= \frac{\partial}{\partial \theta} \left[ \sum_{i=1}^n y_i \log \left( \frac{1}{1 + \exp(-\theta^T x_i)} \right) + (1 - y_i) \log \left( \frac{\exp(-\theta^T x_i)}{1 + \exp(-\theta^T x_i)} \right) \right]\end{aligned}$$

## ② Chain Rule Approach (Complete Gradient Descent)

See how Chain Rule does following things

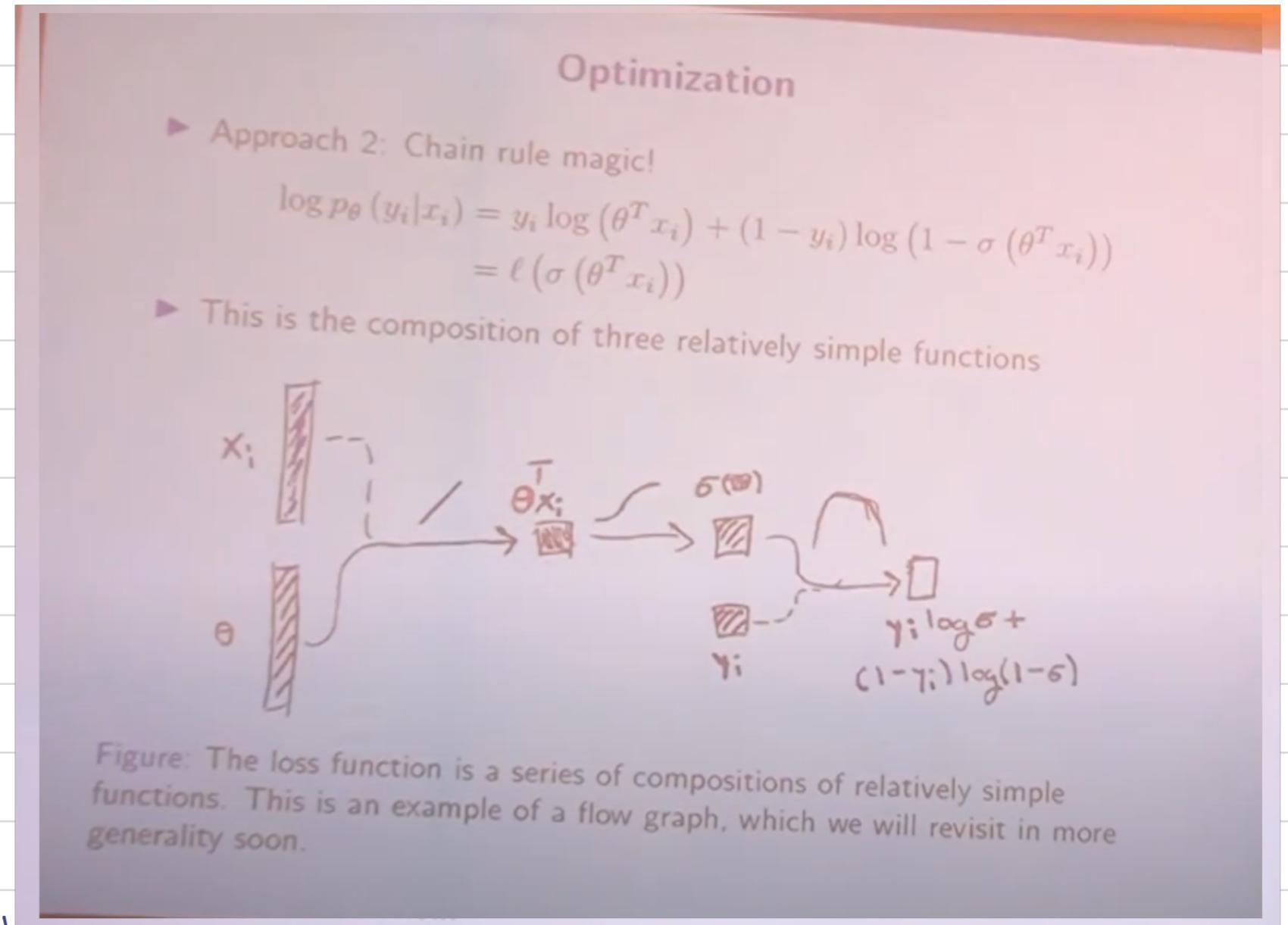
① Function Decomposition

② Linear Combination

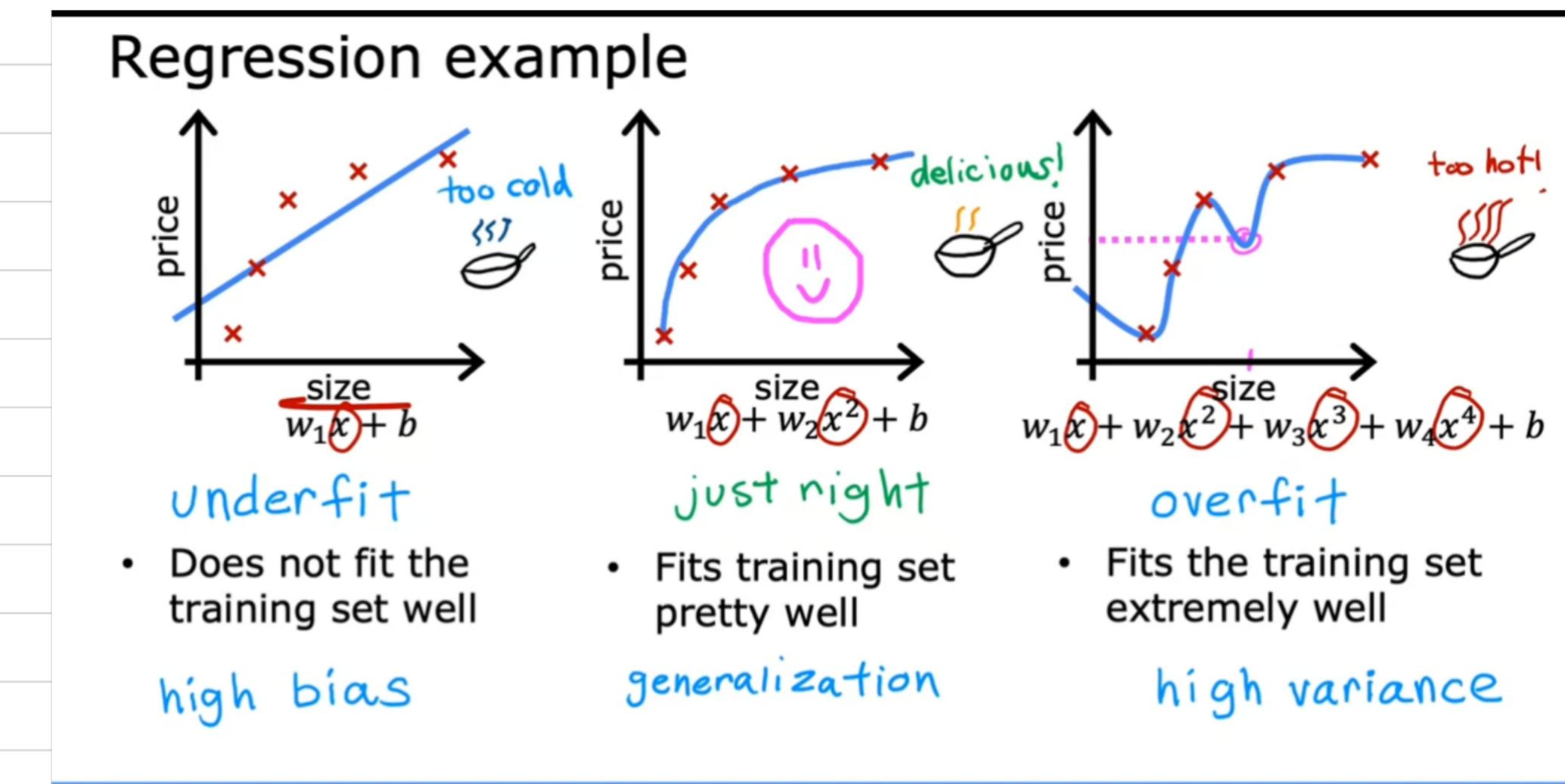
③ Function Differentiable

④ Derivative & chain Rule

breaks down complex problems functions into simpler ones



↗ Overfitting:  
 ↗ Model learns  
 noise &  
 have high  
 variance.



$\rightarrow w_1x + b$   
 $\rightarrow$  High bias  
 Underfit

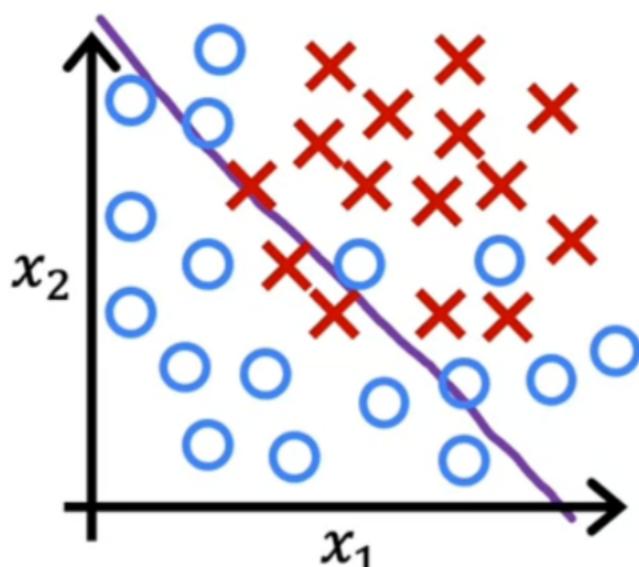
$w_1x + w_2x^2 + b$   
 generalized.

$\rightarrow w_1x^2 + w_2x^2 + w_3x^3 + w_4x^4 + b$   
 Overfit  
 $\rightarrow$  High Variance.

See how  
polynomial  
features play  
role in  
Overfit.

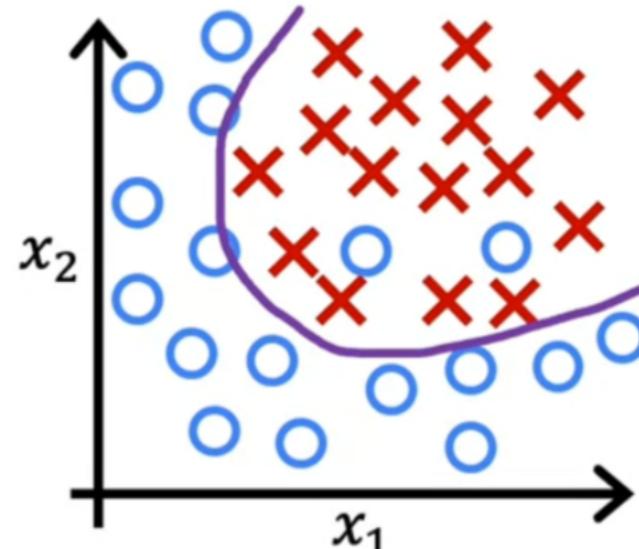
## Classification

coursera.org - To exit full screen, press esc



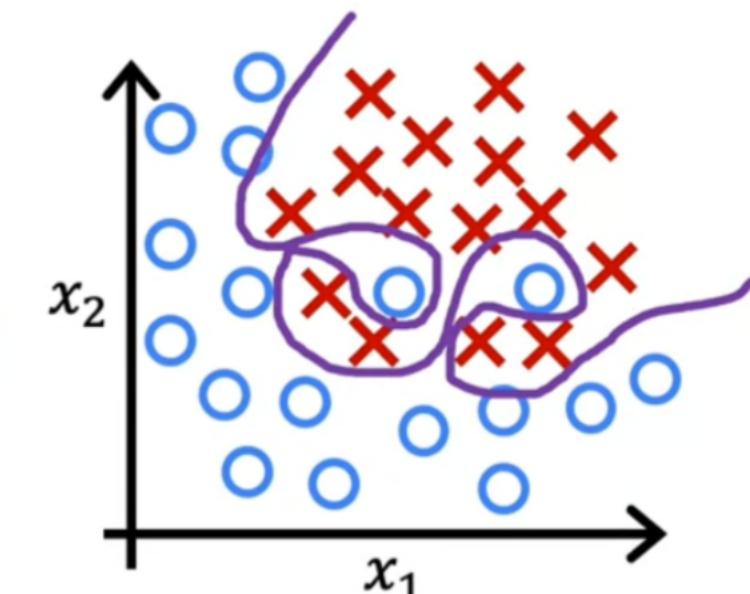
$$z = w_1 \cancel{x_1} + w_2 \cancel{x_2} + b$$
$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

$g$  is the sigmoid function  
underfit      high bias



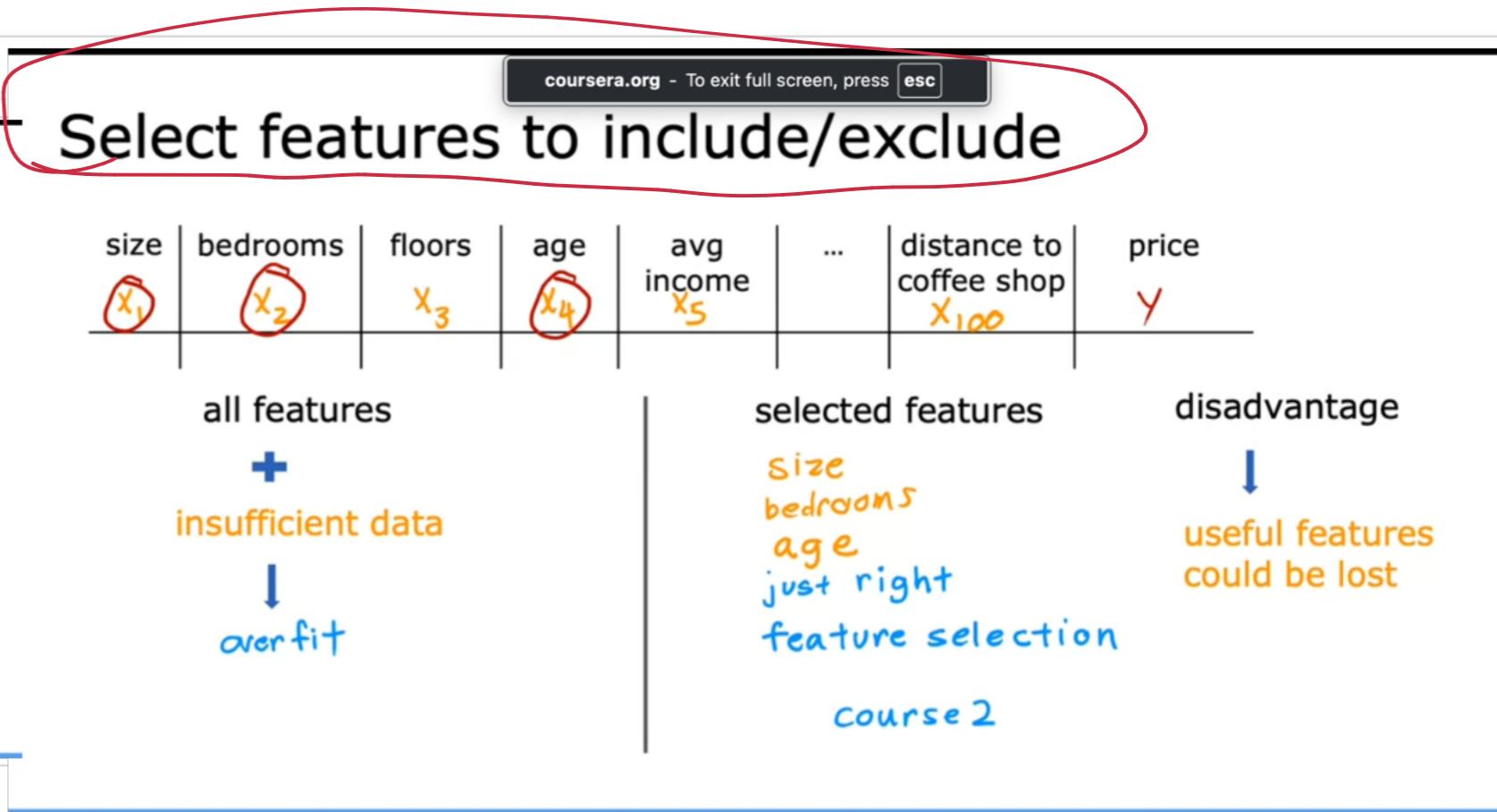
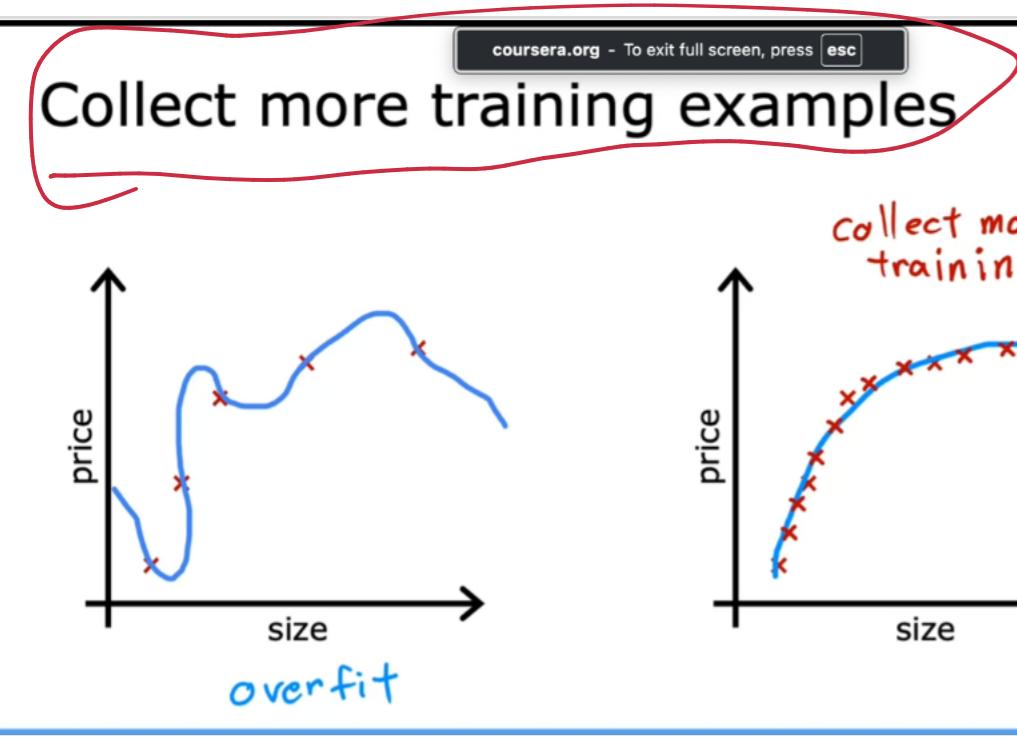
$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + b$$

just right



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2 + \dots + b$$

If it is likely that irrelevant features will aid in overfitting while less data or not enough features will result in underfitting.



How to prevent Overfitting? (Possible Solutions)

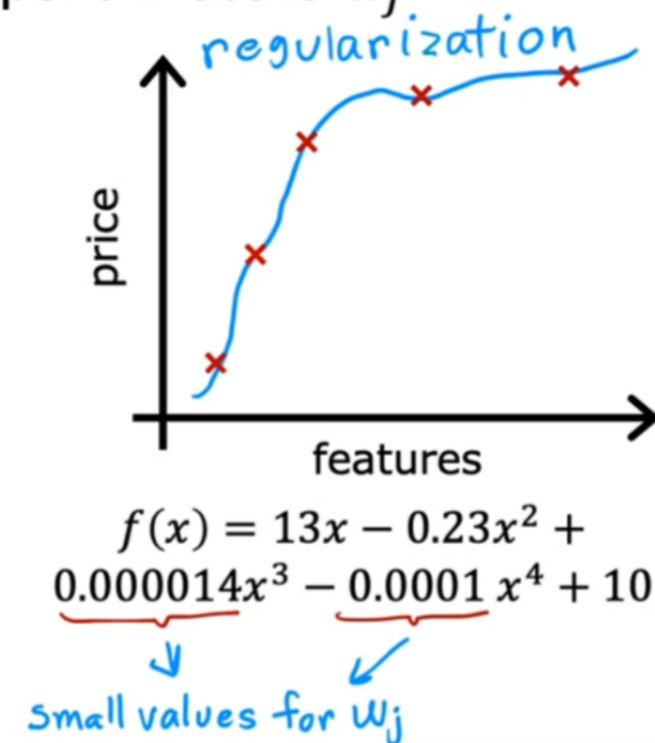
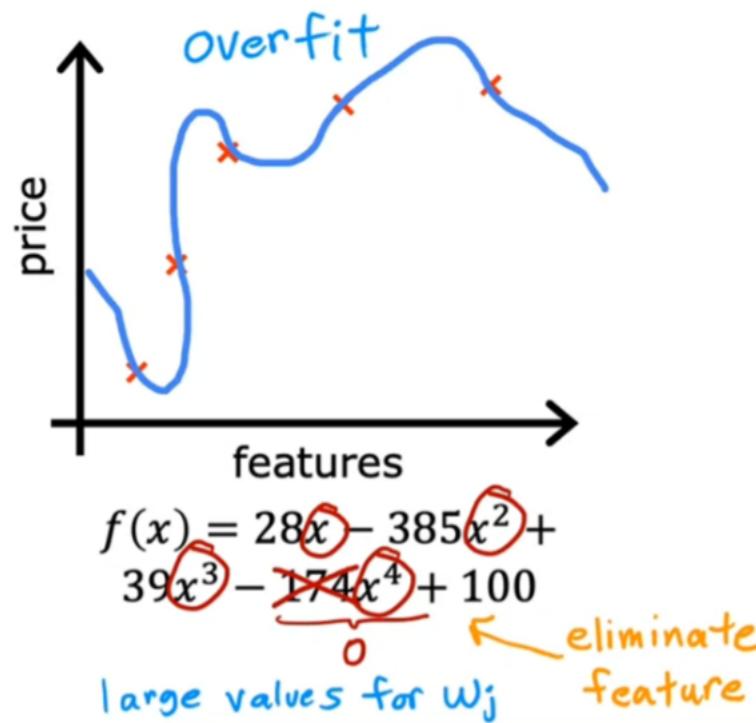
- Collect More training Examples
- Remove Irrelevant Features
- Select best features

- Regularization
- Cross-Validations
- Training & test Split

# Regularization

coursera.org - To exit full screen, press esc

Reduce the size of parameters  $w_j$



# Feature Selection

(1) Wrapper

(2) Embedded

(3) RFE

## Regularization:

L1 (L2)  
Elastic Net

## Addressing overfitting

### Options

1. Collect more data
2. Select features
  - Feature selection [in course 2](#)
3. Reduce size of parameters
  - "Regularization" [next videos!](#)