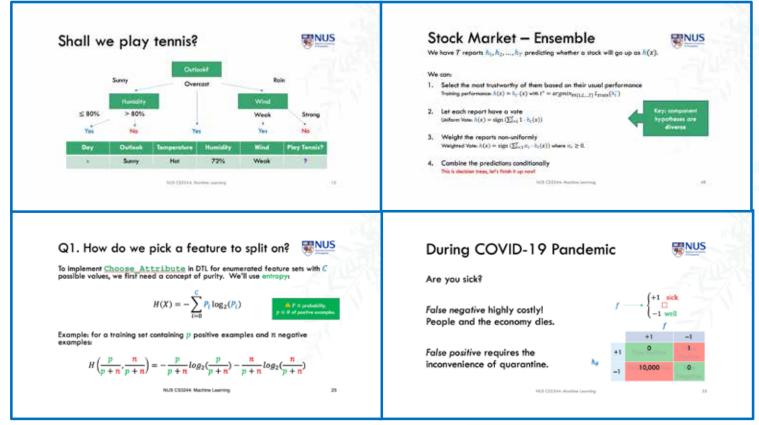


#### Recap from Week 03





#### Project



Groupings need to be redone; sorry!

Some incomplete submissions were supposed to be completed;
 please check your status.

Check the #projects channel bookmarks for details on proposal, final presentation and their evaluation metrics

New project groupings hopefully by Wednesday.

#### Forecast for Week 04A

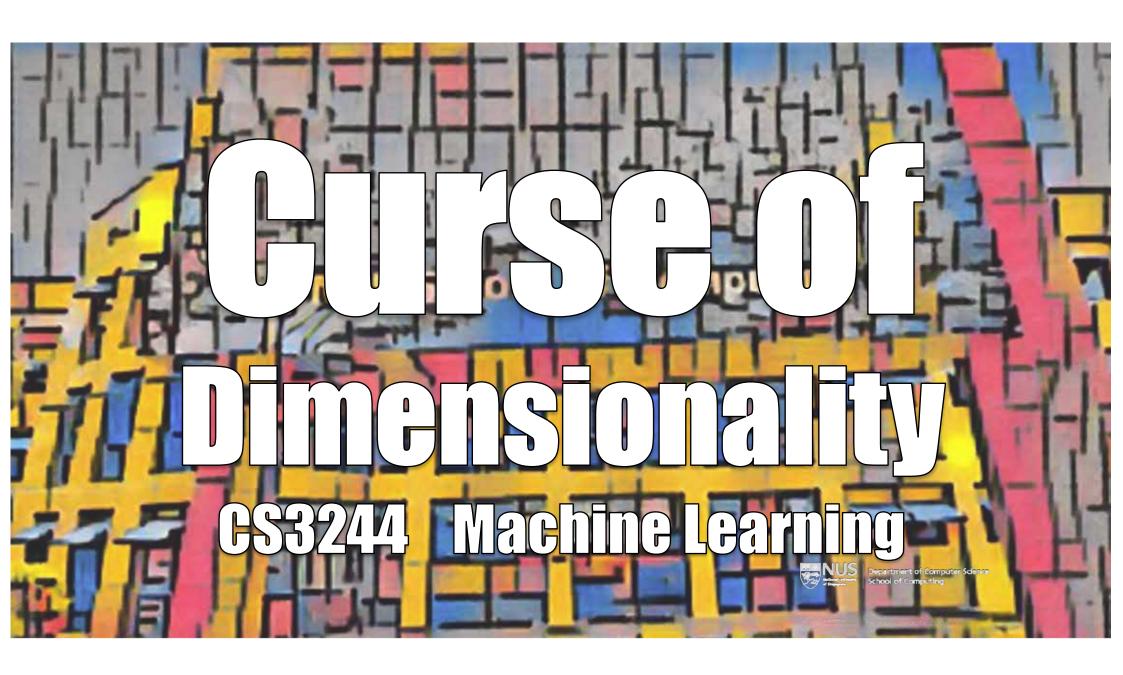


#### Learning Outcomes for this week

- Describe the basic idea of linear classification;
- Understand how both linear and logistic regression works;

#### Other important concepts:

- Curse of Dimensionality
- Gradient Descent





k NN is one classifier that suffers from the curse of dimensionality.

What is this "curse"?

## Sparsity with high dimensions



$$m = 5$$
  
 $n = 1$ 



#### Sparsity with high dimensions

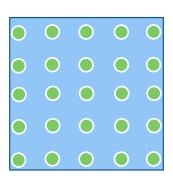


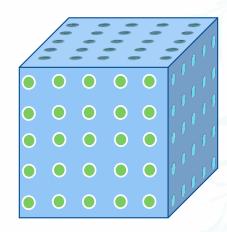
$$m = 5$$
  
 $n = 1$ 

$$m = 25$$
$$n = 2$$

$$m = 125$$
  
 $n = 3$ 







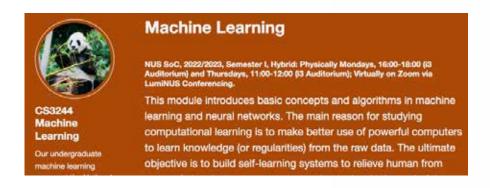
Sparsity problem: maintaining density of samples depends on exponential growth of the data

### Unit Hypercube with Colab



Let's Go!

http://www.comp.nus.edu.sg/~cs3244/AY22S1/ 04.colab.html



#### Standard Deviation $\sigma$ grows slowly

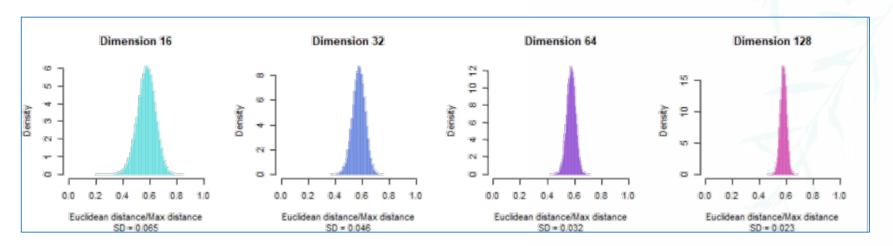


```
1 dimensions:
   Generated 200 points
Average distance (STD) of test point to 200 samples: 0.908554 (0.627194)
2 dimensions:
# Generated 200 points
Average distance (STD) of test point to 200 samples: 1.357414 (0.672876)
3 dimensions:
    Generated 200 points
Average distance (STD) of test point to 200 samples: 2.104206 (0.862900)
4 dimensions:
   Generated 200 points
Average distance (STD) of test point to 200 samples: 2.857518 (0.851753)
8 dimensions:
   Generated 200 points
Average distance (STD) of test point to 200 samples: 4.298752 (0.915833)
10 dimensions:
   Generated 200 points
Average distance (STD) of test point to 200 samples: 4.240121 (0.873999)
100 dimensions:
    Generated 200 points
Average distance (STD) of test point to 200 samples: 14.180229 (0.844941)
200 dimensions:
# Generated 200 points
Average distance (STD) of test point to 200 samples: 19.588471 (0.845534)
```

#### Curse of Dimensionality



In high dimensional space, most points are nearly the same distance away. The result: learners that depend on distance break down in high dimensions.



 ${\color{blue} \underline{https://stats.stackexchange.com/questions/451027/mathematical-demonstration-of-the-distance-concentration-in-high-dimensions}$ 

#### Summary — in high dimensions



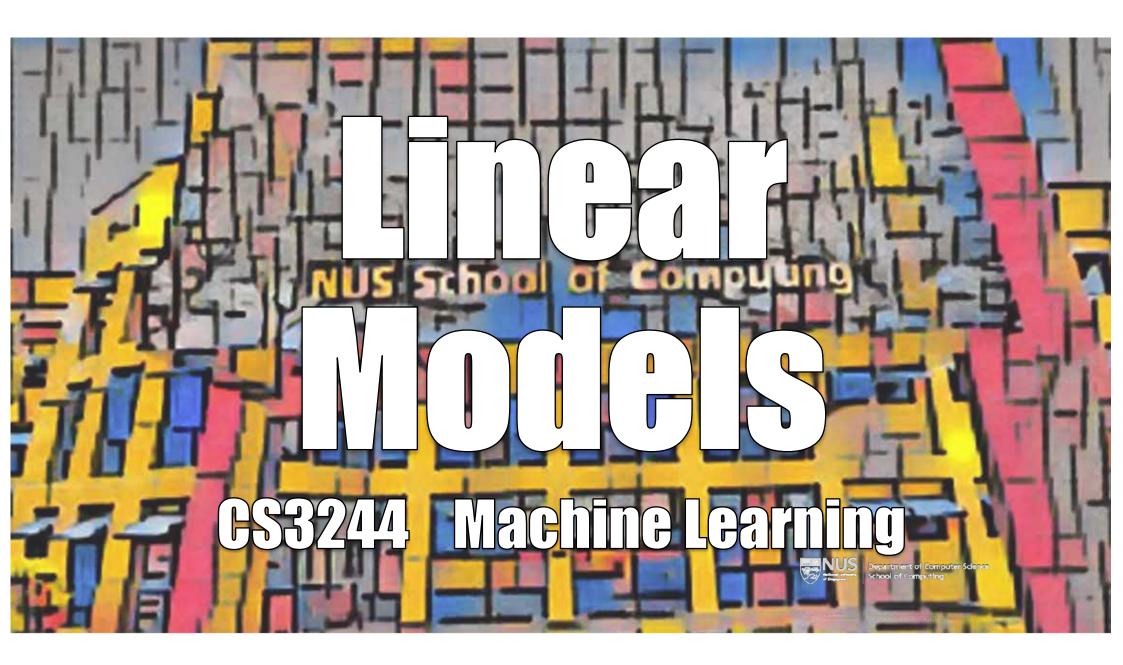
Hard to visualize, but 3B1B does a great job: Thinking outside the 10-dimensional box <a href="https://www.youtube.com/watch?v=zwAD6dRSVyl">https://www.youtube.com/watch?v=zwAD6dRSVyl</a>

Difficult to get sufficient samples

Standard distance metrics become less meaningful

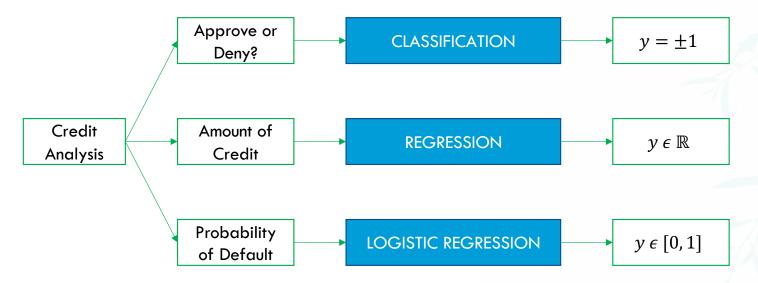
 High dimensional hypercube is mostly "edge" space, far from the center of the hypercube

Solution? Choose or design an appropriate distance metric



#### Three learning problems



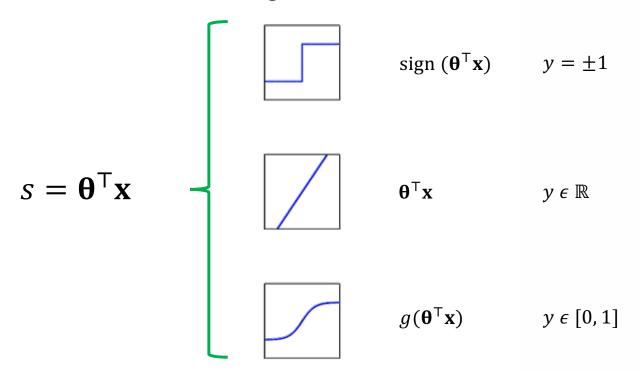


Linear models are the fundamental models.

The linear model is the first model to try.

## The linear signal





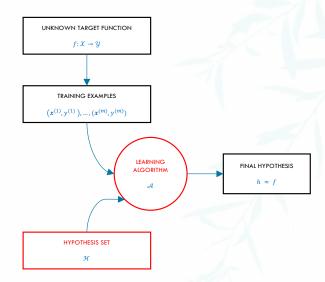
#### Quick Recap: Choosing a model



We choose a Learning Algorithm  $\mathcal{A}$  (equivalent to choosing  $\mathcal{H}$ ), and any associated hyperparameters.

We employ  $\mathcal{A}$  to then choose an h parameterized by its  $\theta$ .

i.e., the learning algorithm  $\mathcal{A}$  selects  $h_{\theta} \in \mathcal{H}$ .

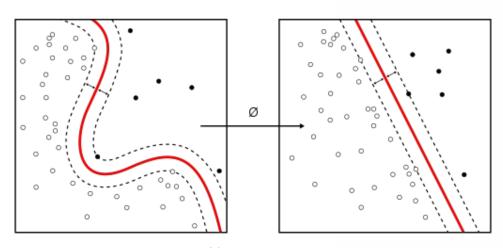


#### ${\cal H}$ for linear models



 ${\mathcal H}$  represents the universe of possible hypotheses that a certain learning algorithm A can generate.

That is, given the choices of  $\theta$ , what types of decision boundaries can it generate?



NUS CS3244: Machine Learning

#### In-Lecture Activity



# 🔼 🦄 What do you think? 👍 or 👎









Value Criterion 32 years Age

Do we extend credit to this person, if this is all the information we know?





Do we extend credit to this person?

h(	(x)	):	[x]	>	25]
----	-----	----	-----	---	-----

Logical test, Iverson brackets

Criterion	Value
Age	32 years

Normalize x:

$$h(x): \left[ \left( \frac{x}{100} \right) > 0.25 \right]$$

The 'bias weight'  $\theta_0$  corresponds to the threshold.  $\theta_1$  and  $\theta_0$  are parameters!.

Offset (Bias)
$$h(\mathbf{x}): \left(\frac{1}{100}\right) x_1 - (0.25) x_0 > 0$$

Introduce artificial  $x_0$  as always equal to 1.



Do we extend credit to this person?

Criterion	Value
Age	32 years
Gender	Male
Salary	40 K







Do we extend credit to this person?

Criterion	Value
Age	32 years
Gender	Male
Salary	40 K
Debt	26 K
	•••
Years in Job	1 year
Years at Current Residence	3 years

$$h(\mathbf{x}): \theta_2 x_2 + \left(\frac{1}{100}\right) x_1 + (-0.25) x_0 > 0$$

$$h(\mathbf{x}): \theta_3 x_3 + \theta_2 x_2 + \theta_1 x_1 + \theta_0 x_0 > 0$$

$$h(\mathbf{x}): \sum_{i=1}^{n} \theta_i x_i > 0$$

#### Simplifying with vector notation



$$h(\mathbf{x}): \theta_n x_n + \dots + \theta_0 x_0 > 0$$

Two vectors: 
$$\mathbf{\theta} = \theta_n, \dots, \theta_0$$
 and  $\mathbf{x} = x_n, \dots, x_0$ .

Solution by default, vectors are vectors

How to multiply?

#### Simplifying with vector notation



How to multiply?



Which to transpose? 
$$\boldsymbol{\theta}$$
 or  $\mathbf{X}$ ?  $\begin{bmatrix} \theta_0 \\ \dots \\ \theta_n \end{bmatrix} \times [x_0 \quad \dots \quad x_n] =$ 



$$[\theta_0 \quad \dots \quad \theta_n] \times \begin{bmatrix} x_0 \\ \dots \\ x_n \end{bmatrix} =$$

#### Simplifying with vector notation



How about with an entire matrix X of instances?

(Recall X is the data matrix, Xs stacked on top of each other)

Which to transpose?  $\boldsymbol{\theta}^{\intercal}\mathbf{X}$  or  $\mathbf{X}\boldsymbol{\theta}$ ?











#### $X\theta =$

$$= \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} \end{bmatrix} \times \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \begin{bmatrix} \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)} + \theta_3 x_3^{(1)} \\ \theta_1 x_1^{(1)} + \theta_2 x_2^{(2)} + \theta_3 x_3^{(3)} \end{bmatrix}$$

$$= \begin{bmatrix} \boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}^{(1)} \\ \boldsymbol{\theta}^{\mathsf{T}} \mathbf{x}^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} \hat{\boldsymbol{y}}^{(1)} \\ \hat{\boldsymbol{y}}^{(2)} \end{bmatrix}$$

### A simple $\mathcal{H}$ : hyperplane with a threshold (bias)



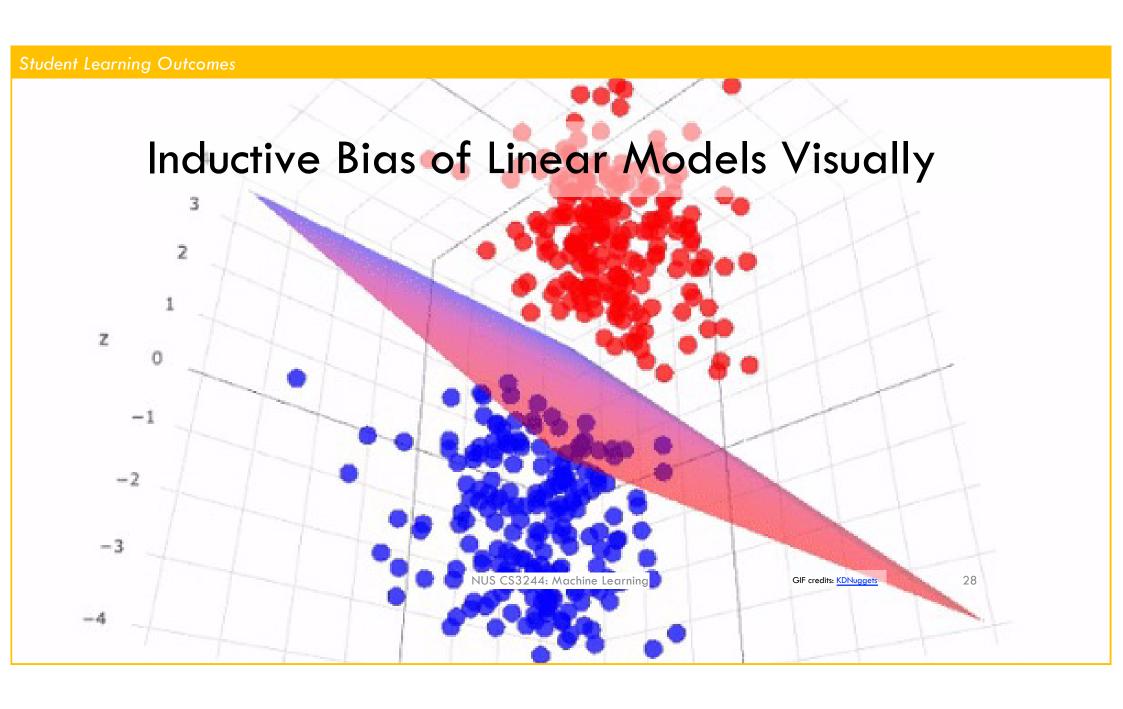
For an input  $\mathbf{x} = (x_1, x_2, ..., x_n)$  'attributes of a customer'

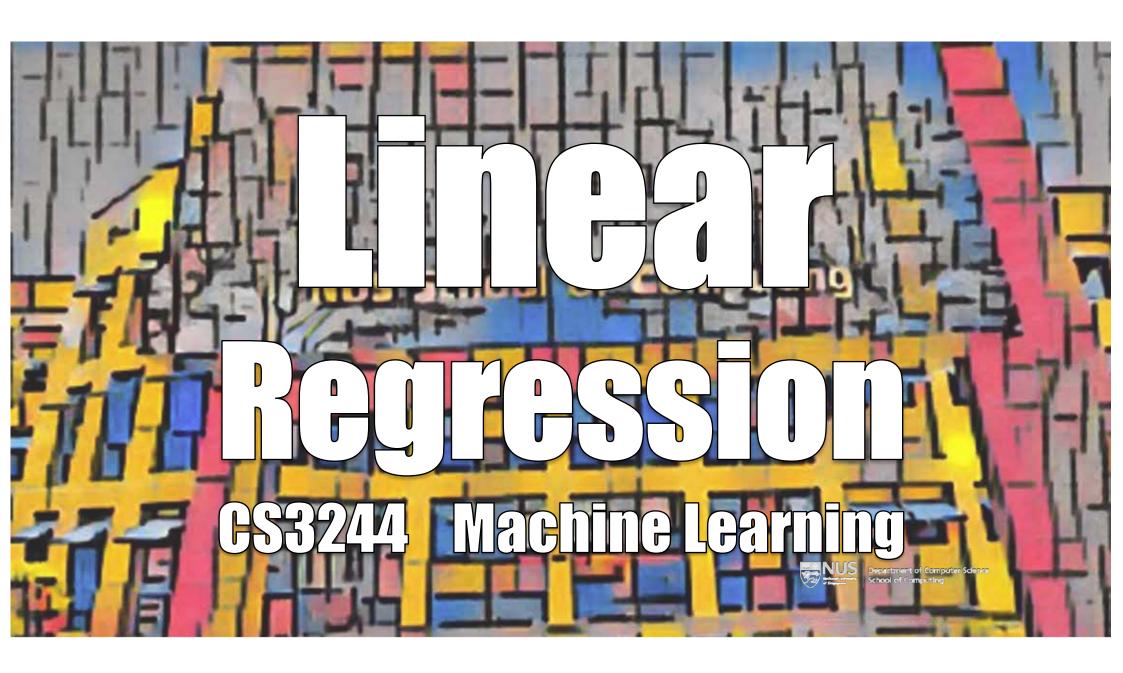
 $\sum_{i=1}^{n} \theta_i x_i > 0$ • Approve credit if:

• Deny credit otherwise:  $\sum_{i=1}^{n} \theta_i x_i < 0$ 

This linear formula  $h \in \mathcal{H}$  can be written as:

$$h_{\theta}(x) = \text{sign}\left(\sum_{i=1}^{n} \theta_{i} x_{i}\right)$$









How much credit do we extend this person?

Criterion	Value
Age	32 years
Gender	Male
Salary	40 K
Debt	26 K
Years in Job	1 year
Years at Current Residence	3 years

#### Classification: Approve/Deny

- Regression: Credit line (real-valued dollar amount)
- Input:
- Linear regression output:  $h_{\theta}(\mathbf{x}) = \sum_{i=0}^{n} \theta_{i} x_{i} = \mathbf{\theta}^{\mathsf{T}} \mathbf{x}$ =  $\theta_{0} + \theta_{1} x_{1} + \theta_{2} x_{2} + \dots + \theta_{n} x_{n}$
- Officers decide on credit lines based on historical data **X**:  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

 $y^{(j)} \in \mathbb{R}$  is the credit line for customer  $x^{(j)}$ ; regression tries to replicate this.

### Linear Regression





#### R-valued cost function



How well does  $h_{\theta}(\mathbf{x}) = \theta^{\mathsf{T}}\mathbf{x}$  approximate  $f(\mathbf{x})$ ?

In linear regression, we use squared error:  $(h_{\theta}(\mathbf{x}) - f(\mathbf{x}))^2$ 

Training error:

#### $\mathbb{R}$ -valued cost function



How well does  $h_{\theta}(\mathbf{x}) = \theta^{\mathsf{T}}\mathbf{x}$  approximate  $f(\mathbf{x})$ ?

In linear regression, we use squared error:  $(h_{\theta}(\mathbf{x}) - f(\mathbf{x}))^2$ 

Training error:

The sweet spot:
Plausible AND
Convenient

$$L_{train}(h) = \frac{1}{m} \sum_{j=1}^{m} (h_{\theta}(x^{(j)}) - y^{(j)})^2$$

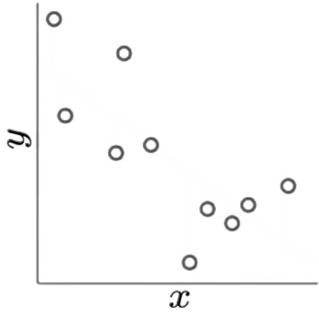
Pointwise Average

How bad is our prediction?

NUS CS3244: Machine Learning

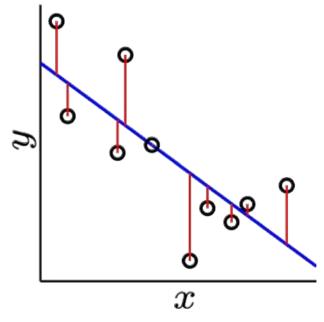
## Illustration of linear regression



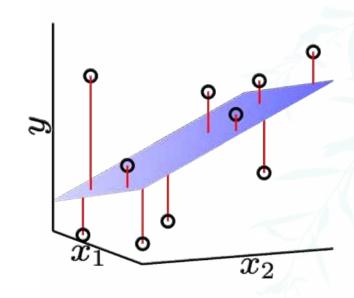


### Illustration of linear regression

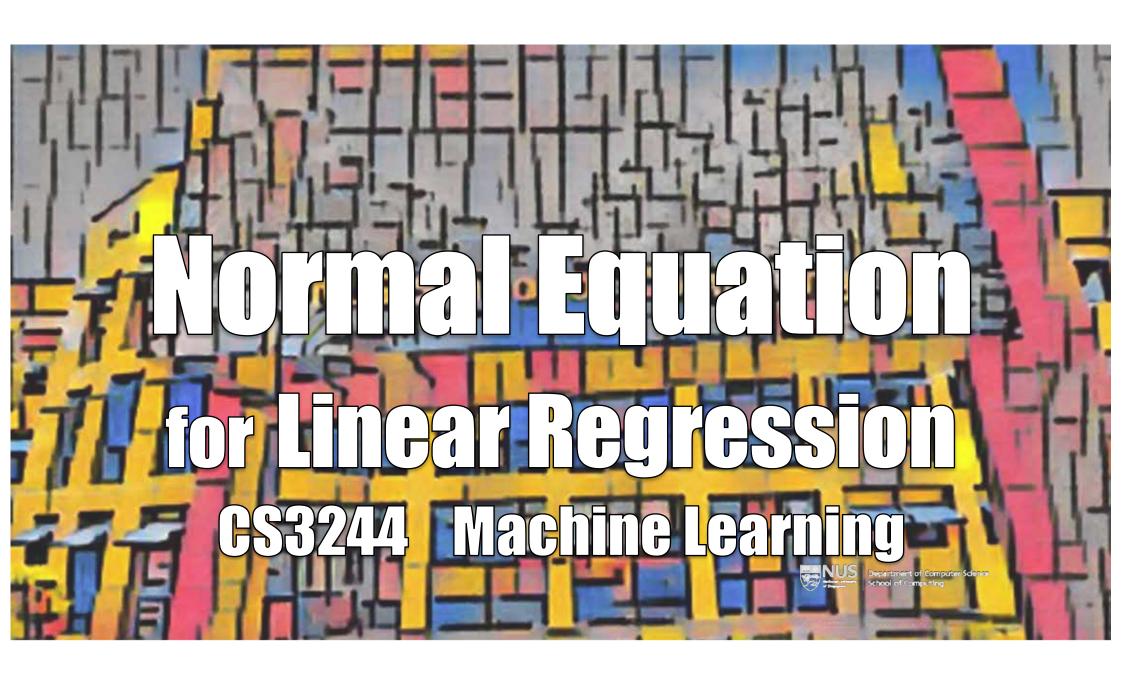




Hypothesis: 
$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1$$



$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$



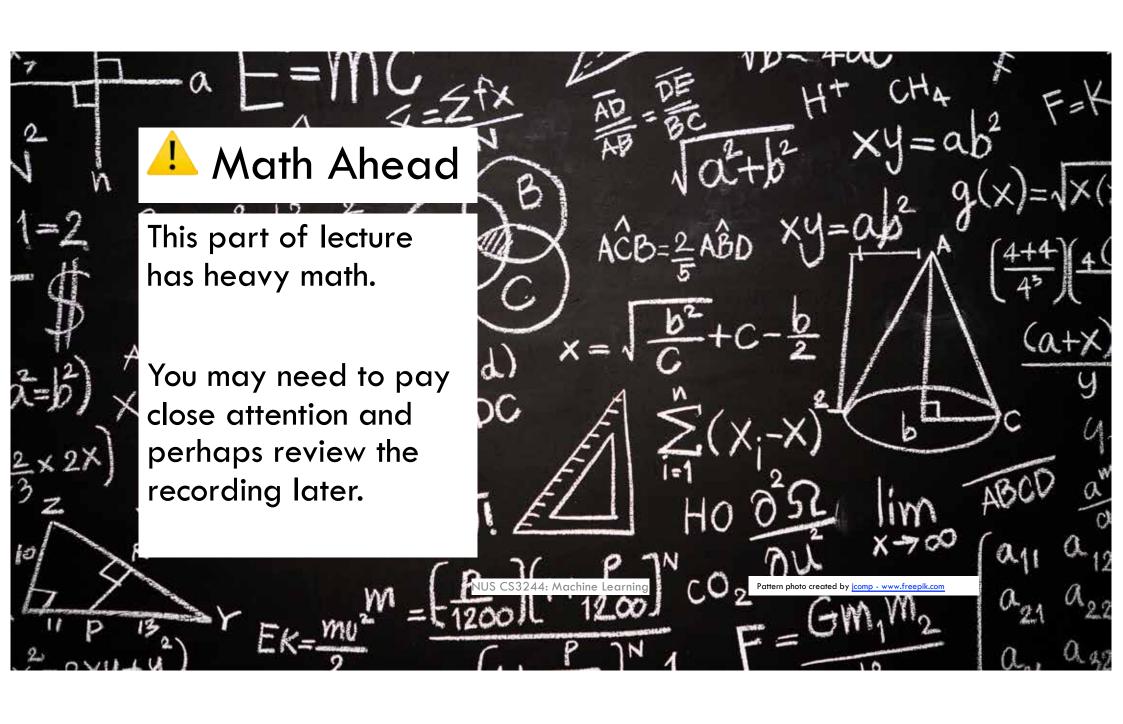
### Linear Regression – Summary



Two methods for solving

1. Gradient Descent

2. Normal Equation



# The expression for $L_{train}$



$$L_{train} = \frac{1}{m} \sum_{j=1}^{m} (\mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)} - y^{(j)})^{2}$$

$$= \frac{1}{m} ||\mathbf{X}\mathbf{\theta} - \mathbf{y}||^{2}$$

$$= \begin{bmatrix} -\mathbf{x}^{(1)}^{\mathsf{T}} - \\ -\mathbf{x}^{(2)}^{\mathsf{T}} - \\ \vdots \\ -\mathbf{y}^{(m)}^{\mathsf{T}} - \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$
where  $\mathbf{X} = \begin{bmatrix} -\mathbf{x}^{(1)}^{\mathsf{T}} - \\ -\mathbf{x}^{(2)}^{\mathsf{T}} - \\ \vdots \\ y^{(m)} \end{bmatrix}$ 

# Minimizing $L_{train}$

Let's re-write  $L_{train}$  first.

$$L_{train} = \frac{1}{m} \|\mathbf{X}\mathbf{\theta} - \mathbf{y}\|^2$$



# Minimizing $L_{train}$



Let's re-write  $L_{train}$  first.

Then solve for  $\theta$ :

ordering in

$$\frac{\partial L_{train}}{\partial \boldsymbol{\theta}} \equiv$$

$$L_{train} = \frac{1}{m} || \mathbf{X} \mathbf{\theta} - \mathbf{y} ||^{2}$$

$$= \frac{1}{m} (\mathbf{X} \mathbf{\theta} - \mathbf{y})^{\mathsf{T}} (\mathbf{X} \mathbf{\theta} - \mathbf{y})$$

$$= \frac{1}{m} ((\mathbf{X} \mathbf{\theta})^{\mathsf{T}} \mathbf{X} \mathbf{\theta} - (\mathbf{X} \mathbf{\theta})^{\mathsf{T}} \mathbf{y}$$

$$-\mathbf{y}^{\mathsf{T}} \mathbf{X} \mathbf{\theta} + \mathbf{y}^{\mathsf{T}} \mathbf{y})$$

$$= \frac{1}{m} (\mathbf{\theta}^{\mathsf{T}} \mathbf{X}^{\mathsf{T}} \mathbf{X} \mathbf{\theta} - 2(\mathbf{X} \mathbf{\theta})^{\mathsf{T}} \mathbf{y} + \mathbf{y}^{\mathsf{T}} \mathbf{y})$$

y and  $X\theta$  are **multiplication** doesn't matter

 $\left[\frac{\partial L}{\partial \theta_0}(\theta_0, \dots, \theta_n)\right]$  $\nabla L \equiv A$  vector with all n component partial derivatives of L. Jacobian: stack of gradients as partial derivatives.

## Minimizing $L_{train}$



Let's re-write  $L_{train}$  first.

$$L_{train} = \frac{1}{m} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^{2}$$

$$= \frac{1}{m} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^{\mathsf{T}} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

$$= \frac{1}{m} ((\mathbf{X}\boldsymbol{\theta})^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} - (\mathbf{X}\boldsymbol{\theta})^{\mathsf{T}} \mathbf{y}$$

$$-\mathbf{y}^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} + \mathbf{y}^{\mathsf{T}} \mathbf{y})$$

$$= \frac{1}{m} (\boldsymbol{\theta}^{\mathsf{T}} \mathbf{X}^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} - 2(\mathbf{X}\boldsymbol{\theta})^{\mathsf{T}} \mathbf{y} + \mathbf{y}^{\mathsf{T}} \mathbf{y})$$

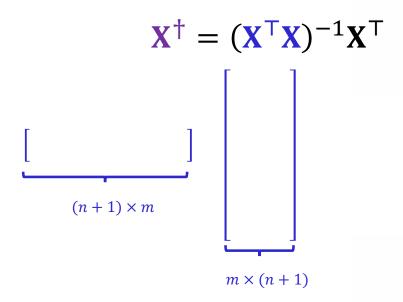
Then solve for  $\theta$ :

$$\frac{\partial \mathcal{L}_{train}}{\partial \theta} \equiv \mathbf{X}^{\top} \mathbf{X} \theta - \mathbf{X}^{\top} \mathbf{y} = \mathbf{0}$$
 Column of zeroes 
$$\mathbf{X}^{\top} \mathbf{X} \theta = \mathbf{X}^{\top} \mathbf{y}$$
 
$$\theta = \left( (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \right) \mathbf{y}$$

 $(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}$  is called the **pseudo inverse** of  $\mathbf{X}$ 

# The pseudo inverse $X^{\dagger}$



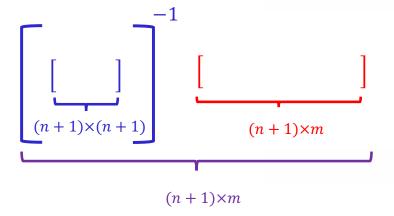




# The pseudo inverse $X^{\dagger}$



$$\mathbf{X}^{\dagger} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}$$



### Pseudo inverse is expensive to

expensive to compute if n or m is large.

### Linear regression by Normal Equation



1. Construct the matrix X and the vector y from the data set as follows:

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}^{(1)}^{\mathsf{T}} - \\ -\mathbf{x}^{(2)}^{\mathsf{T}} - \\ -\mathbf{x}^{(3)}^{\mathsf{T}} - \\ \vdots \\ -\mathbf{x}^{(m)}^{\mathsf{T}} - \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

1. Compute the pseudo inverse  $\mathbf{X}^{\dagger} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}$ 

input data matrix

2. Return  $\theta = \mathbf{X}^{\dagger} y$ 

The normal equation for LR. One-step learning!

target vector

### Linear Regression — Summary



### Two methods for solving

- 1. Gradient Descent
  - ullet Works well, even when n is large
  - Works even if  $\mathbf{X}^{\mathsf{T}}\mathbf{X}$  is non-invertible

 $\red{9}$  To think about: What are some reasons that  $\mathbf{X}^\mathsf{T}\mathbf{X}$  might not be invertible?

### 2. Normal Equation

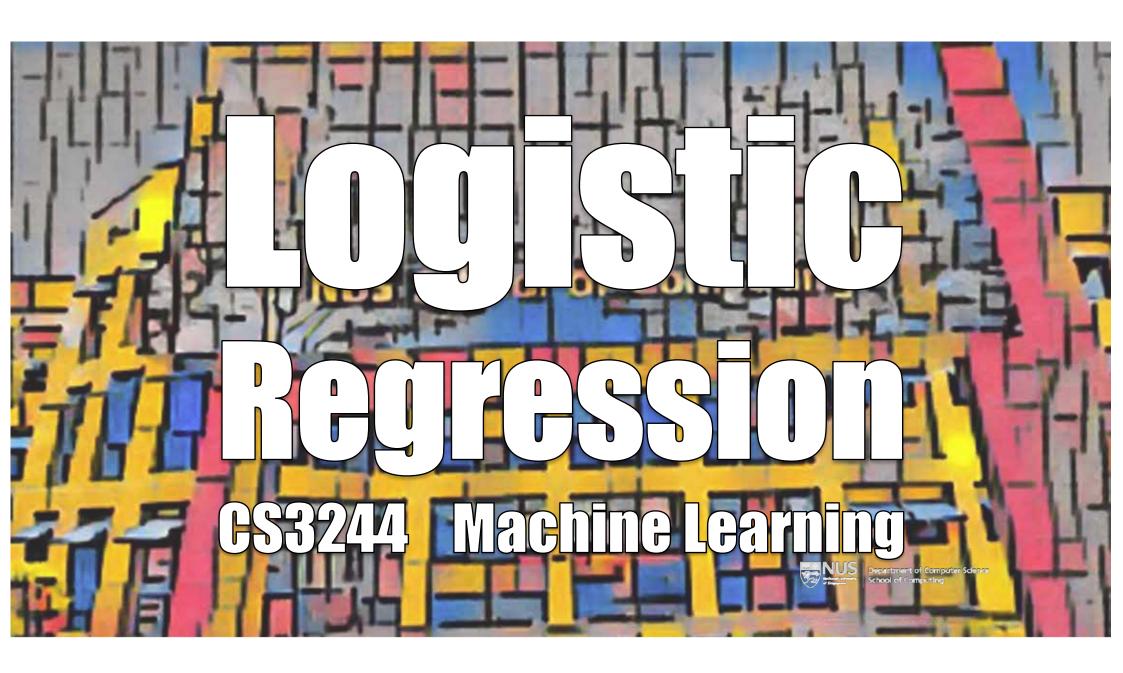
- Don't need to choose  $\alpha$
- Don't need to iterate
- Need to compute  $(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}$ , an  $O(n^3)$  operation

### Why is it called "regression"?

The term is usually attributed to Sir Francis Galton. He observed that children's heights tended to be less extreme than their parents; the heights of the children **regressed** to the average adult height; hence he titled his paper "Regression towards Mediocrity in Hereditary Stature."

For Galton, "regression" referred only to the tendency of extreme data values to "revert" to the overall mean value. (Minital Blog, 1)



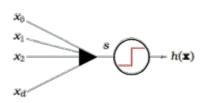


### A third linear model

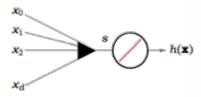


$$s = \sum_{i=0}^{n} \theta_i x_i$$

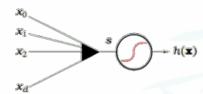
#### Linear classification



#### Linear regression



#### Logistic regression



## The logistic function g



We choose a form for  $g(\cdot)$  out of convenience. We use the logistic function, which has the form:

$$g(s) = \frac{exp^s}{1 + exp^s}$$

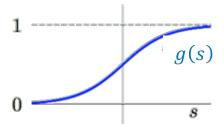
### The logistic function g



We choose a form for  $g(\cdot)$  out of convenience. We use the logistic function, which has the form:

$$g(s) = \frac{exp^s}{1 + exp^s}$$

g(s): needs to map any  $x: \mathbb{R}$  to y: [0,1]



Soft threshold: uncertainty

Sigmoid: flattened out 's' shape

### Probability Interpretation



 $h_{\theta}(\mathbf{x}) = g(s)$  is interpreted as a probability

Example: Prediction of heart attacks

Input X: cholesterol level, age,

weight, diabetic, etc.

output g(s): Likelihood of a heart

attack

 $s = \theta^\mathsf{T} \mathbf{x}$ The signal:

"risk score"

Logistic regression  $\equiv y \in [0,1]$ 



# Labels are binary, yet our outputs are probabilities



$$\mathbf{X} = (\mathbf{x}^{(1)}, y^{(1)} = \pm 1), \dots (\mathbf{x}^{(m)}, y^{(m)} = \pm 1)$$

 $\mathbf{x}^{(j)}$  — a person's health information

 $y^{(j)} = \pm 1$  — did they have a heart attack?

We cannot measure a probability.

Why?

We can only see the occurrence of an event and try to <u>infer</u> a probability.

### Genuine Probability



Data (x, y) with binary y, generated by a noisy target:

$$P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1; \\ 1 - f(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

The target  $f: \mathbb{R}^n \to [0,1]$  is the probability

Learn 
$$h_{\theta}(\mathbf{x}) = g(\mathbf{\theta}^{\mathsf{T}}\mathbf{x}) \approx f(\mathbf{x})$$

### Cost Function for Logistic Regression



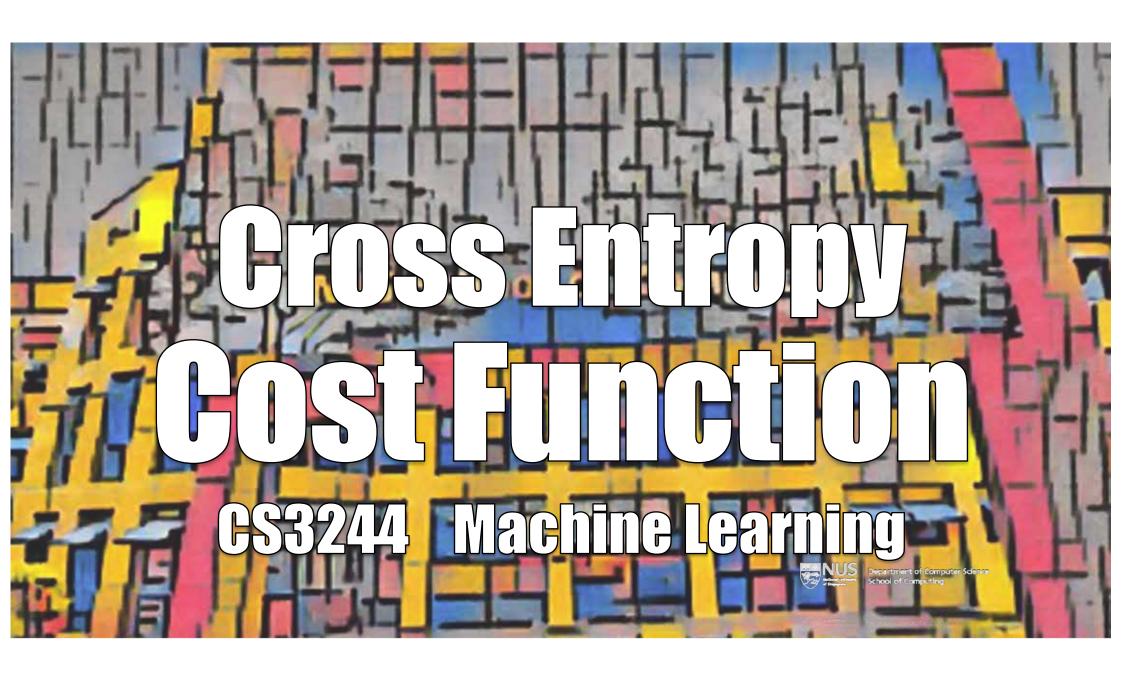
For logistic regression,

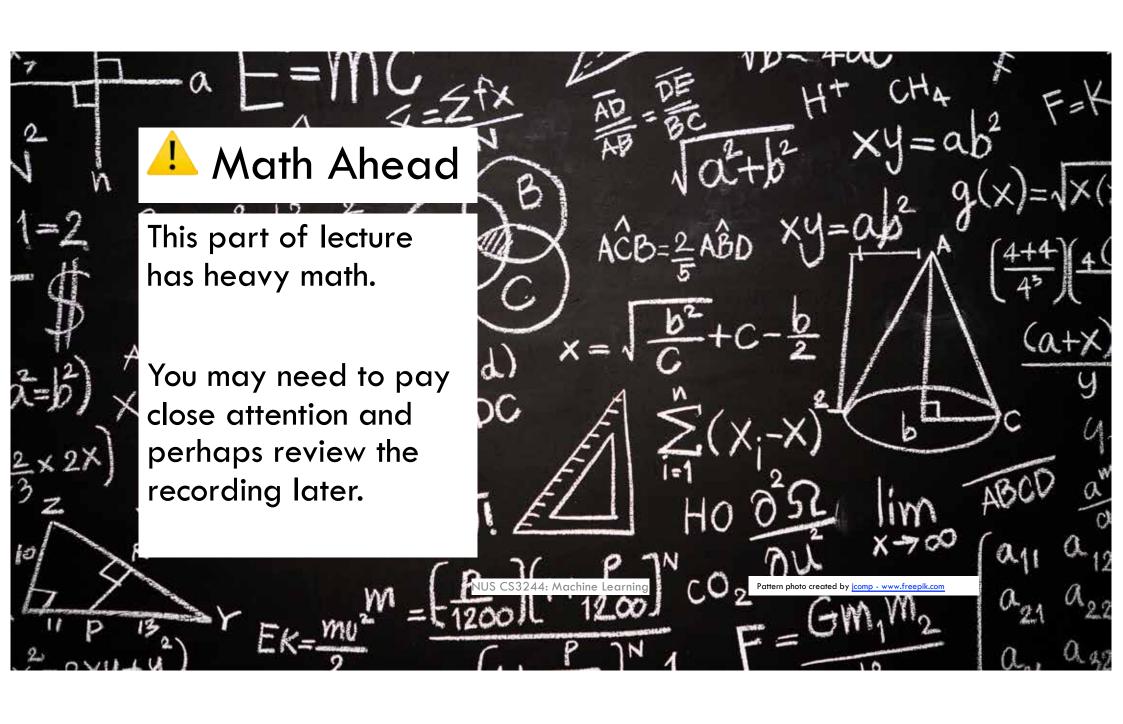
$$L_{train}(\mathbf{\theta}) = \frac{1}{m} \sum_{j=1}^{m} ln(1 + exp^{-y^{(j)}} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$
lterative Solution; yucky (or interesting) math

Compare to linear regression:

$$L_{train}(\mathbf{\theta}) = \frac{1}{m} \sum_{j=1}^{m} (\mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)} - y^{(j)})^{2}$$







### Logistic Regression Cost Function



For each (x, y), y is generated by probability f(x)

Plausible cost function based on likelihood:

If  $h_{\theta} = f$ , how "likely" is it to obtain output y from input x?

$$P(y|\mathbf{x}) = \begin{cases} h_{\theta}(\mathbf{x}) & \text{for } y = +1; \\ 1 - h_{\theta}(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

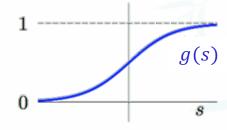
### Formula for likelihood



$$P(y|\mathbf{x}) = \begin{cases} h_{\theta}(\mathbf{x}) & \text{for } y = +1; \\ 1 - h_{\theta}(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute 
$$h_{\theta}(\mathbf{x}) = g(\mathbf{\theta}^{\mathsf{T}}\mathbf{x})$$
, noting  $g(-s) = 1 - g(s)$ 

$$P(y|\mathbf{x}) = g(y \cdot \mathbf{\theta}^{\mathsf{T}} \mathbf{x})$$



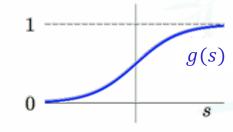
### Formula for likelihood



$$P(y|\mathbf{x}) = \begin{cases} h_{\theta}(\mathbf{x}) & \text{for } y = +1; \\ 1 - h_{\theta}(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute  $h_{\theta}(\mathbf{x}) = g(\mathbf{\theta}^{\mathsf{T}}\mathbf{x})$ , noting g(-s) = 1 - g(s)

$$P(y|\mathbf{x}) = g(y \cdot \mathbf{\theta}^{\mathsf{T}} \mathbf{x})$$



Likelihood of  $\mathbf{X} = (\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$  where samples are i.i.d. is:

$$P(y^{(1)},...,y^{(m)}|\mathbf{x}^{(1)},...,\mathbf{x}^{(m)}) = \prod_{j=1}^{m} P(y^{(j)}|\mathbf{x}^{(j)}) = \prod_{j=1}^{m} g(y^{(j)}\mathbf{\theta}^{\mathsf{T}}\mathbf{x}^{(j)})$$

# Maximizing the likelihood ≡ Minimizing cross entropy

$$= max \qquad \prod_{j=1}^{m} g(y^{(j)} \mathbf{e}^{\mathsf{T}} \mathbf{x}^{(j)})$$



# Maximizing the likelihood ≡ Minimizing cross entropy

$$= max \qquad \prod_{j=1}^{m} g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\iff \max \qquad \frac{1}{m} \ln \left( \prod_{j=1}^{m} g\left( y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)} \right) \right)$$

$$= max \qquad \frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Leftrightarrow min \qquad -\frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$



## Maximizing the likelihood $\equiv$ Minimizing cross entropy



$$= \max \qquad \prod_{j=1}^{m} g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Rightarrow \max \qquad \frac{1}{m} \ln \left( \prod_{j=1}^{m} g\left(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)}\right) \right)$$

$$= \max \qquad \frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Rightarrow \min \qquad -\frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Leftrightarrow \min \qquad -\frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$
"cross entropy" error

"cross entropy" error

### Maximizing the likelihood ≡ Minimizing cross entropy



$$= \max \qquad \prod_{j=1}^{m} g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Rightarrow \max \qquad \frac{1}{m} \ln \left( \prod_{j=1}^{m} g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)}) \right)$$

$$= \max \qquad \frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$= \max \qquad \frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Rightarrow \min \qquad -\frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$\Leftrightarrow \min \qquad -\frac{1}{m} \sum_{j=1}^{m} \ln g(y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)})$$

$$L_{train}(\mathbf{\theta}) = \frac{1}{m} \sum_{j=1}^{m} \ln(1 + e^{-y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)}})$$

"cross entropy" error

# Summary: Minimizing $L_{train}$



For logistic regression,

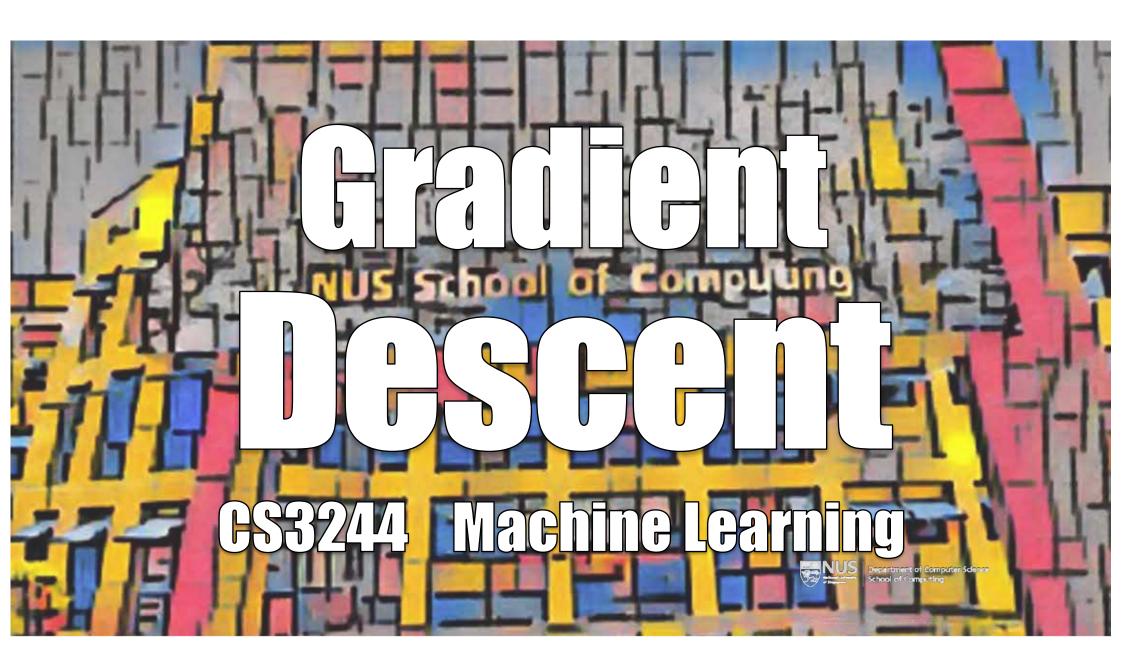
$$L_{train}(\mathbf{\theta}) = \frac{1}{m} \sum_{j=1}^{m} ln(1 + e^{-y^{(j)} \mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)}})$$



Compare to linear regression:

$$L_{train}(\mathbf{\theta}) = \frac{1}{m} \sum_{j=1}^{m} (\mathbf{\theta}^{\mathsf{T}} \mathbf{x}^{(j)} - y^{(j)})^{2}$$





# Climbing up (down) one step at a time

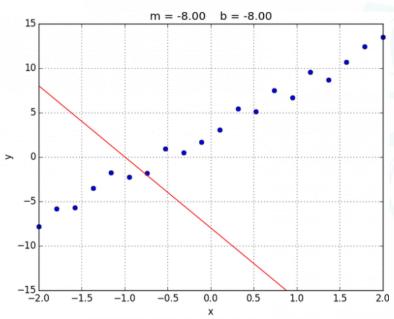


NUS CS3244: Machine Learning



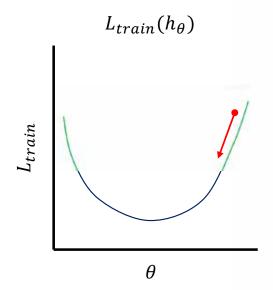
## Univariate Linear Regression: Salary to predict Credit Line





# Our $L_{train}$ only has one valley

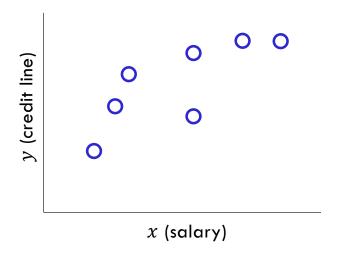


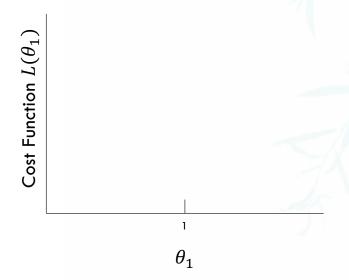


... Because  $L_{train}$  is a convex function of  $\theta$ .

# Cost Function – only slope $L(\theta_1)$

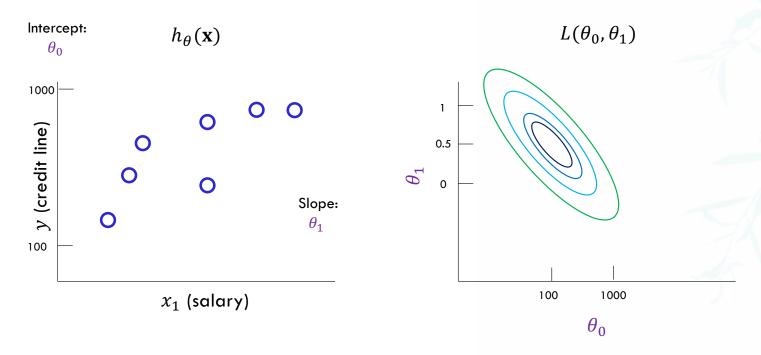






# Cost Function with intercept $L(\theta_0, \theta_1)$





### Iterative method: gradient descent



General method for nonlinear optimization

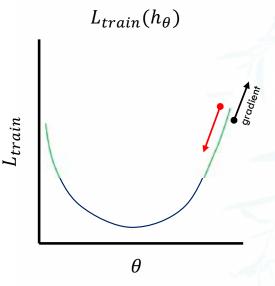
Start at  $\theta(0)$ ; take a step in the opposite direction of the steepest slope

Notation:  $\theta(t) \equiv \text{set of}$  parameters at iteration t

Fixed step size  $\eta : \theta(1) = \theta(0) + \eta \mathbf{v}$ 

We get to pick the direction of **v**. What is the best direction to pick?

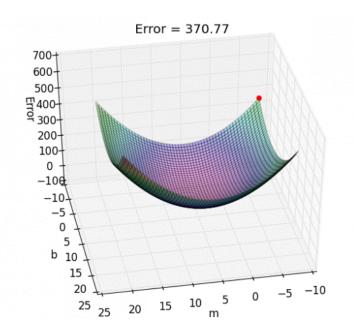
The one to minimize  $L_{train}$ .

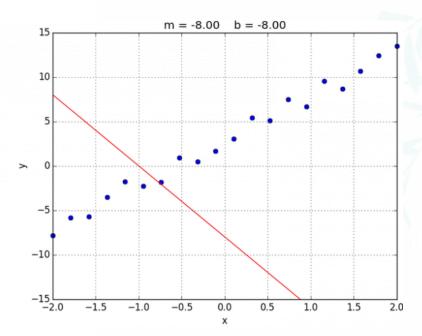


Gradient descent can minimize any smooth function.

Oh yeah!

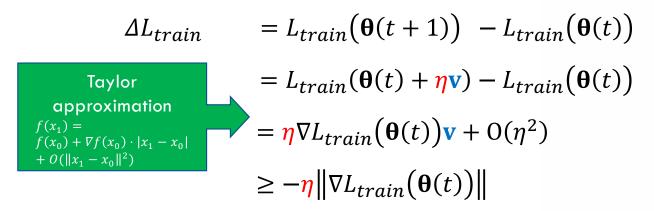






### Formula for the direction V

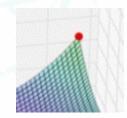




Notation:  $\nabla$  ("nabla")  $\equiv$  Gradient. Generalization of slope for higher dimensions

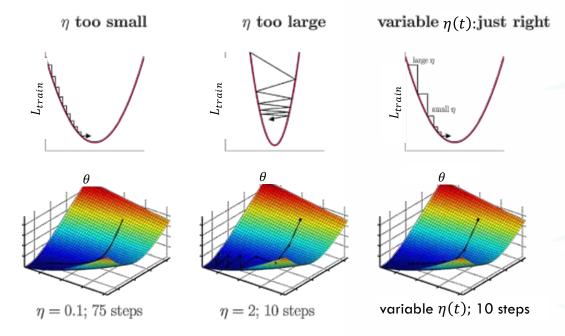
Create v to be a unit vector:

$$\mathbf{v} = -\frac{\nabla L_{train}(\mathbf{\theta}(t))}{\|\nabla L_{train}(\mathbf{\theta}(t))\|}$$



## Goldilocks step size





 $\eta$  should increase with the slope

### Varying step size, for free



Instead of

$$\Delta \Theta = \frac{\eta \mathbf{v}}{1}$$

$$= - \eta \frac{\nabla L_{train}(\boldsymbol{\theta}(t))}{\|\nabla L_{train}(\boldsymbol{\theta}(t))\|}$$

Have

$$\Delta \mathbf{\theta} = -\alpha L_{train}(\mathbf{\theta}(t))$$

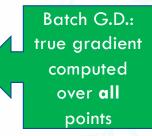
Fixed learning <u>rate</u>  $\alpha$  instead of fixed learning <u>step</u>  $\eta$ 

### Logistic regression algorithm



- 1. Initialize the weights at t = 0 to  $\theta(0)$
- 2. Do
- 3. Compute the gradient

$$\nabla(t) = \nabla L_{train}(\mathbf{\theta}(t)) = -\frac{1}{m} \sum_{j=1}^{m} \frac{y^{(j)} \mathbf{x}^{(j)}}{1 + e^{y^{(j)} \mathbf{\theta}(t)^{\mathsf{T}} \mathbf{x}^{(j)}}}$$



- 4. // Move in the direction  $\mathbf{v}(t) = -\nabla(t)$ Update the weights  $\mathbf{\theta}(t+1) = \mathbf{\theta}(t) - \alpha \nabla L_{train}$
- 5. Continue to next iteration, until it is time to stop
- 6. Return the final weights  $\theta^*$

### **Termination Condition**



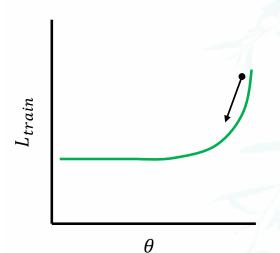
When to stop?

Natural choice: gradient < threshold

But lots of flat regions in most spaces:

Instead, use criteria:

- error change is small and/or;
   error is small;
   maximum number of iterations is reached.



### Stochastic Gradient Descent





A variation of GD that considers only the error on one data point.

Pick one  $(\mathbf{x}^{(*)}, \mathbf{y}^{(*)})$  at a time.

Apply GD to  $I(h_{\theta}(\mathbf{x}^{(*)}), y^{(*)})$ 

"Average" Direction:  $\mathbb{E}_* \Big[ - \nabla I(h_{\theta}(\mathbf{x}^{(*)}), y^{(*)}) \Big]$   $= \frac{1}{m} \sum_{k=1}^m - \nabla I(h_{\theta}(\mathbf{x}^{(*)}), y^{(*)})$   $= - \nabla L_{train}$ 

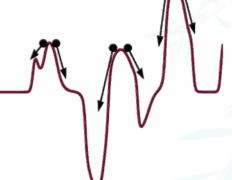
**Stochastic** Gradient Descent (SGD)

### Benefits of SGD



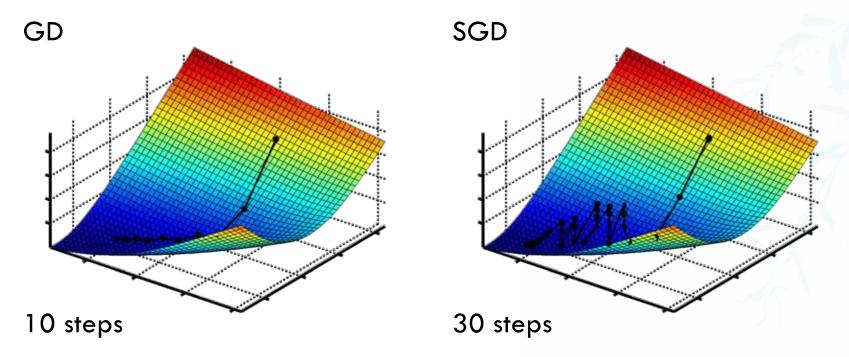
1. Cheaper computation: Fraction 1/m cheaper per step

- 2. Stochastic: Helps escape local minima
- 3. Simple



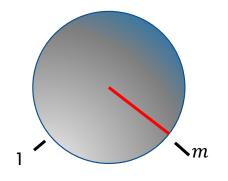
# GD vs. SGD on m=10



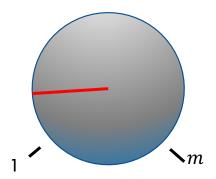


### Mini Batch Gradient Descent

#### **Gradient Descent**



Mini Batch



Pick <u>a few</u>  $(\mathbf{x}^{(*)}, \mathbf{y}^{(*)})$  at a time.

# Stochastic Gradient Descent

