

BLG 561 E FALL 2021

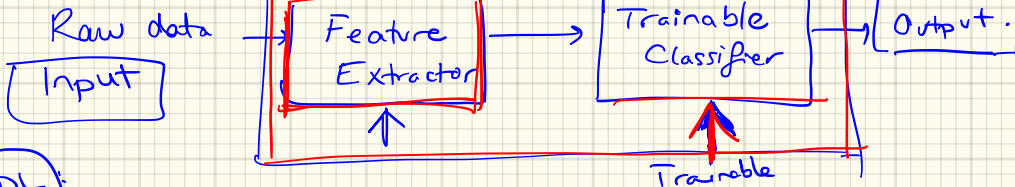
Deep Learning

12.10.2021

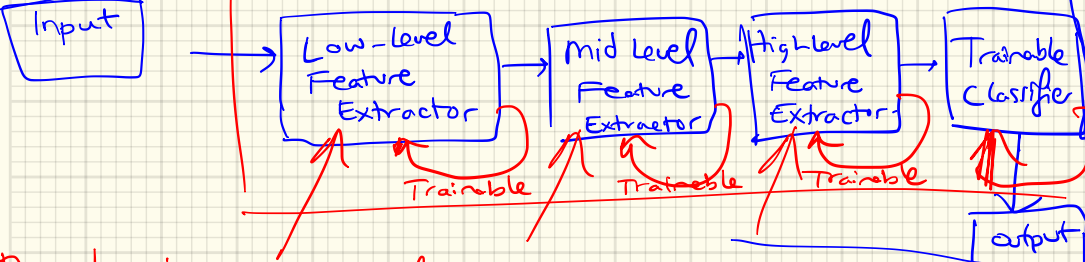
Görde ÜNAL

Machine Learning:

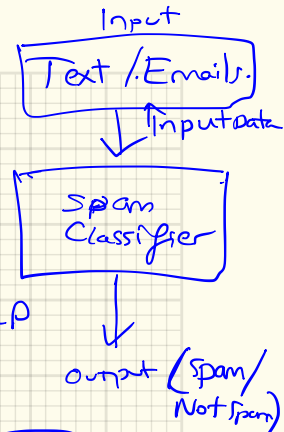
Traditional ML.



New ML/DL

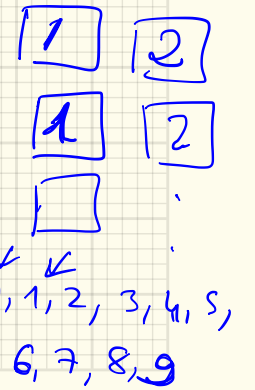
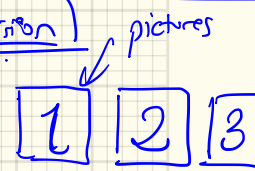


Deep Learning: Many layers → many many parameters  
 many



NLP

Vision



In Linear Algebra:  $\frac{\# \text{ unknowns}}{\# \text{ parameters}} \sim \leq \# \text{ of data points}$

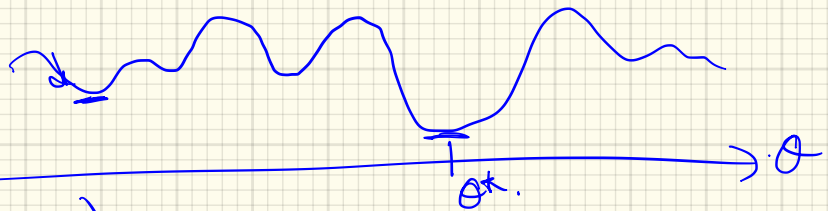
Deep Learning: Over-parameterization!!!  
 $\# \text{ parameters} \geq \# \text{ data points}$

→ Conjectures:

Non-convex

Recall: To get a global extremum in an optimizer, → we need a convex obj fn

$L(\theta)$



Convex

$\arg \min_{\theta} L(\theta)$

regularized SGD.



Flat minima conjecture

Very-high dimensional spaces !!

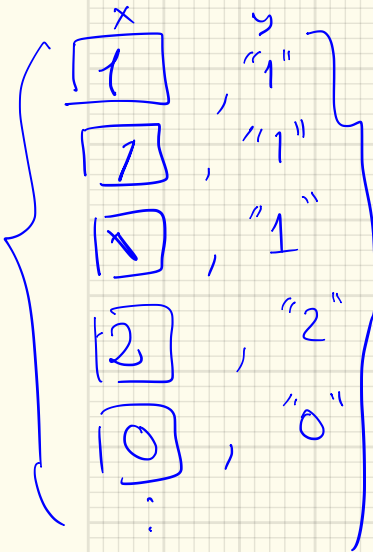
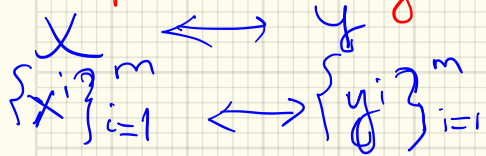
500M power.

very sharp global minima



# ML Review:

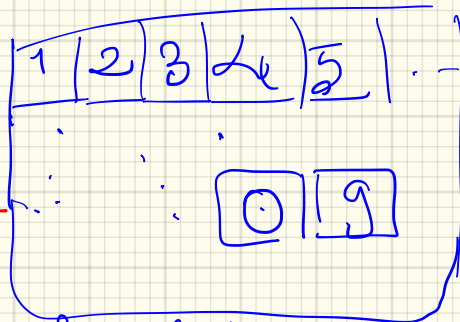
## Supervised Learning:



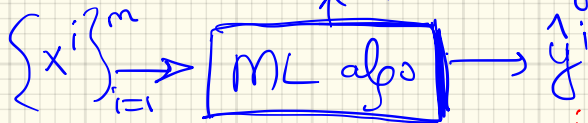
## Handwritten Digit Recognition

Example:

$m$ : # input data samples in the Training Set



$\theta$ : parameters of the algorithm



$\hat{y}^i = h_{\theta}(x^i)$  : hypothesis fn.

Loss function:  
 $l(h_{\theta}(x), y)$

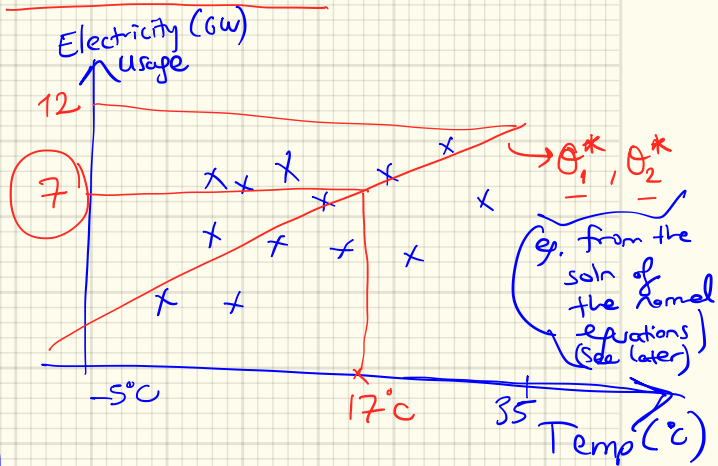
$h: X \rightarrow Y$   
 $x^i \rightarrow y^i$   
 $x^i \in X, y^i \in Y$

We minimize this loss fn. distance btw  $h_{\theta}$  &  $y$ . to get  $\theta$ .

# → Linear Regression (Supervised ML).

Q. Predict electricity demand in Istanbul tomorrow?

Dates	Temperature (°C)	Used Electricity (Gw)
01.06.2017	20°	8.05
02.06.2017	25°	3.06
⋮	⋮	⋮
01.06.2021	30°	⋮

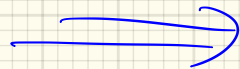


ML Notation: Input features  $x^{(i)} \in \mathbb{R}^n$ ,  $i = 1, \dots, n$

Electricity Demand  $\approx \theta_1 \cdot \text{Temp}^{(i)} + \theta_2$  ; eg.  $\underline{x}^{(i)} = \begin{bmatrix} T^{(i)} \\ 1 \end{bmatrix}$

Output  $\underline{y}^{(i)} \in \mathbb{R} \rightarrow$  <sup>b/c</sup> Regression task

Input features



ML Algo: Framework

$\{x^i, y^i\}_{i=1}^m$  : training data (electricity usage & its predictors from before)

Model Parameters:  $\underline{\theta} \in \mathbb{R}^{n+1}$  : same size as the input features

Define the following:

① Hypothesis Function: eg.  $h_{\theta}(x) = \underline{\theta}^T \underline{x} = \sum_{j=1}^n \theta_j x_j$

tells us how confident we are in the mapping.

$$h: X \rightarrow Y$$

② Loss Function: measures how "good" our hypothesis function is.

$$l(\overset{\text{prediction}}{h_{\theta}(x)}, \overset{\text{GT}}{y}) \quad , \quad h_{\theta}(x^{(i)}) \approx \underset{\text{or Reference Labels}}{\overset{\text{"Ground Truth" (GT)}}{y^{(i)}}}$$

↓ we should define this loss fn.

$$\rightarrow l: Y \times Y \rightarrow \mathbb{R}^+$$

eg. here  $l: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$  for the simple regression task here.

→ Q. What is a common loss function for a regression task?

A.  $l(h_\theta(x), y) = (h_\theta(x) - y)^2$  (MSE) loss.

\* All Machine Learning (ML) algorithms require the following.

We define the general ML (algorithm) problem:

Given  $(x^i, y^i) \Big|_{i=1}^m$ , a set of input samples (# : m of them) & outputs

training set

Goal: Find the parameters that minimize the sum of the losses:

$$\arg \min_{\theta} \sum_{i=1}^m l(h_\theta(x^i), y^{(i)})$$

$l^{(i)}$

- ① What is the hypothesis function?
- ② What is the loss function?
- ③ How do you solve the given optimization problem?

To solve this, all ML algorithms should specify:

→  $y \in \mathbb{R}$  in linear regression

$y \in \{0, 1\}$  in binary classification.

$y \in \{1, \dots, k\}$  multi-class classification.

⇒ Ex: Linear Regression problem:

①  $h_{\theta}(x) = \underline{\theta}^T \underline{x}$

②  $l(h_{\theta}(x) - y)^2$  : MSE.

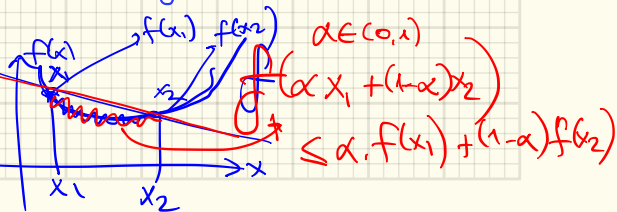
③ With 1 & 2 lead to ML optimization problem:

$\arg \min_{\underline{\theta}} \sum_{i=1}^m l(h_{\theta}(x_i), y_i) = \sum_{i=1}^m (\theta^T x_i - y_i)^2 = \mathcal{L}(\underline{\theta})$  : overall loss fn.

→ Is  $\mathcal{L}(\underline{\theta})$  a convex fn?

Yes. here.

Recall:  
Convex fn.





→ For a convex problem;  $\rightarrow$  many convex solvers; python: `cvxpy`. ✓

Generally: Calculate the gradient of our loss (objective) fn.

$$\rightarrow \nabla_{\underline{\theta}} \left( \sum_{i=1}^m (\underline{\theta}^T \underline{x}^i - y^i)^2 \right) = \sum_{i=1}^m \nabla_{\underline{\theta}} (\underline{\theta}^T \underline{x}^i - y^i)^2$$

$$\nabla_{\underline{\theta}} \mathcal{L}(\underline{\theta}) = \sum_{i=1}^m 2(\underline{\theta}^T \underline{x}^i - y^i) \underline{x}^i$$

Soln:

① Iterative Optimization: For example; Gradient Descent (GD)

GD update rule:  $\underline{\theta}^0$ : initial values,  $k=0$   
 $k$ : iteration index (superscript)

$$\underline{\theta}^{k+1} \leftarrow \underline{\theta}^k - \alpha \nabla_{\underline{\theta}} \mathcal{L}(\underline{\theta})$$

$$\begin{aligned} \theta_1^k &\leftarrow \theta_1^{k-1} - \alpha \nabla_{\theta_1} \mathcal{L}(\underline{\theta}) \\ \theta_2^k &\leftarrow \theta_2^{k-1} - \alpha \nabla_{\theta_2} \mathcal{L}(\underline{\theta}) \end{aligned}$$

② For this ex; we get a Closed Form soln:

Introduce vector  
matrix  
notation

$$\underline{X} = \begin{bmatrix} \underline{x}^{(1)T} \\ \underline{x}^{(2)T} \\ \vdots \\ \underline{x}^{(m)T} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

$$\underline{\theta} \in \mathbb{R}^n, \quad \underline{y} = \begin{bmatrix} y^1 \\ \vdots \\ y^m \end{bmatrix} \in \mathbb{R}^m$$

Rewrite the LS objective:  $\sum_{i=1}^m (\theta^T x_i - y_i)^2 = \| \underline{X} \underline{\theta} - \underline{y} \|_2^2$

$$L(\underline{\theta}) = (\underline{X} \underline{\theta} - \underline{y})^T (\underline{X} \underline{\theta} - \underline{y})$$

Scalar  
take the derivative  
w.r.t.  
 $\underline{\theta}$

$$L(\underline{\theta}) = \underline{\theta}^T \underline{X}^T \underline{X} \underline{\theta} - 2 \underline{\theta}^T \underline{X}^T \underline{y} + \underline{y}^T \underline{y}$$

$$\nabla_{\underline{\theta}} L(\underline{\theta}) = 2 \underline{X}^T \underline{X} \underline{\theta} - 2 \underline{X}^T \underline{y} \rightarrow \text{set this to } \underline{0}$$

$$\nabla_{\underline{\theta}} L(\underline{\theta}) = 0 \rightarrow \underbrace{\underline{X}^T \underline{X}}_{\substack{\text{non} \\ \text{nxn matrix}}} \underline{\theta} = \underline{X}^T \underline{y}$$

$$\underline{\theta}^* = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

"NORMAL EQUATIONS"  
a closed form solution  
for the minimization of  
sum of squared losses

eg. the electricity demand prediction problem, we solve the normal equations,  
we get  $\theta_1^*, \theta_2^* \rightarrow$  these define the Red Line in the plot (from before)

A new ML problem:  $\textcircled{1}$  Linear hypothesis  $\checkmark$  Absolute Loss function ('MAE: Mean Absolute Error')

$\textcircled{2}$  Loss fr. defined  $\checkmark$

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m l(h_{\theta}(x^i) - y^i) = \frac{1}{m} \sum_{i=1}^m |h_{\theta}(x^i) - y^i|$$

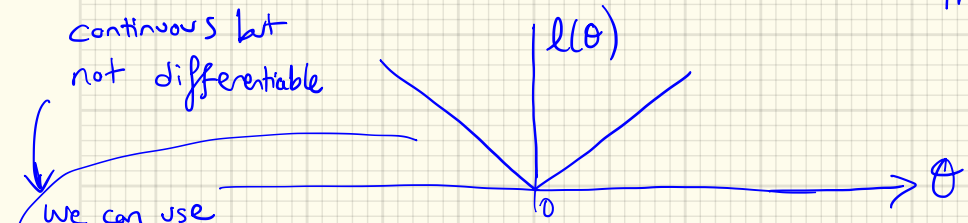
$l^i$ : loss fr. for 1 sample

$$= \|h_{\theta}(x) - y\|_1$$

$\left( \begin{matrix} h_{\theta}(x^1) - y^1 \\ h_{\theta}(x^2) - y^2 \\ \vdots \\ h_{\theta}(x^m) - y^m \end{matrix} \right)$

$\checkmark$ -norm

continuous but not differentiable



Note: No closed form solution exists!

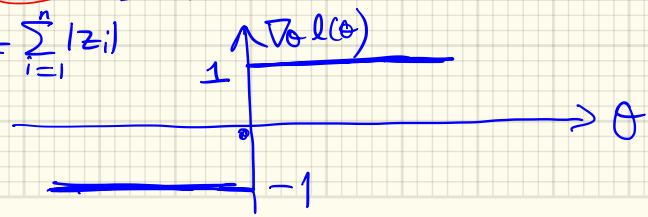
We can use (sub)gradients

$$\nabla_{\theta} \|h_{\theta}(x) - y\|_1 = \underline{X}^T \cdot \text{sign}(\underline{X}\theta - \underline{y})$$

is a vector

Now, we have the gradient of the loss fr.

$\checkmark$ -norm: Note:  $\|z\|_1 = \sum_{i=1}^n |z_i|$



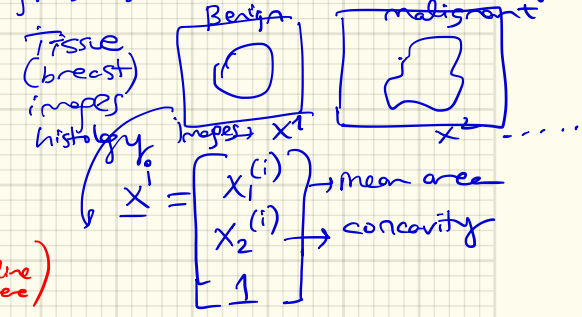
$\textcircled{3}$  GD.

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} L(\theta)$$

$\downarrow$  LR: learning rate.

Q: Define the ML algorithm for the Binary Classification:  $y \in \{+1, -1\}$

eg. Spam classification / NLP  
 eg. Benign / malignant tissue classification



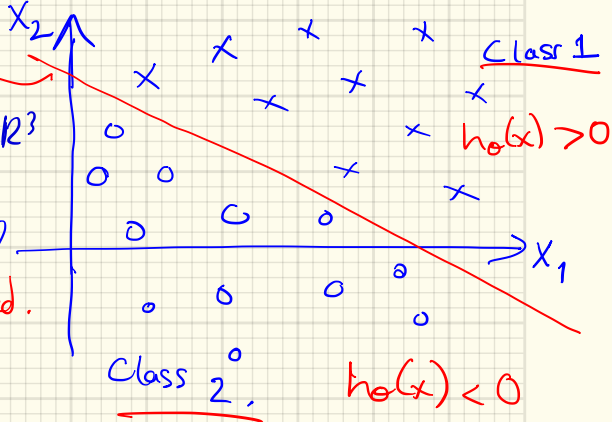
Choose  
 (1) Linear Hypothesis:  
 $h_{\theta}(x) = \underline{\theta}^T \underline{x}$

Prediction:

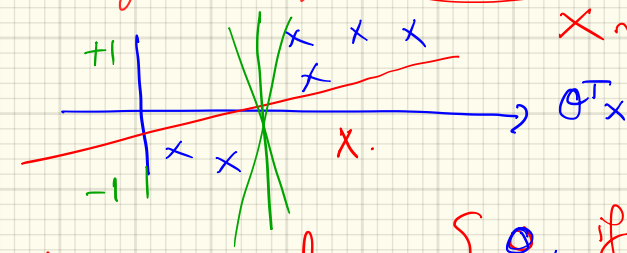
$$\hat{y} = \text{sign}(h_{\theta}(x))$$

$h_{\theta}(x) = 0$   
 $\therefore$  the line is a separating hyperplane (eg. a line here)  
 the two classes.

Model parameters are:  $\underline{\theta} \in \mathbb{R}^{n+1}$ , here  $\underline{\theta} \in \mathbb{R}^3$



(2) Loss function? what if we use a LS loss?  
 X not desired.



0/1 Loss: 
$$l_{0/1} = \begin{cases} 0, & \text{if } \text{sign}(h_{\theta}(x)) = y \end{cases} \quad \parallel \quad \begin{cases} y \cdot h_{\theta}(x) \leq 0 \end{cases}$$

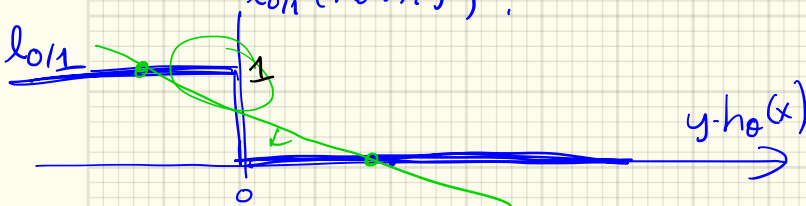
indicator fn.

$$\rightarrow \mathbb{1}_{\{y \cdot h_{\theta}(x) \leq 0\}} = l(\theta) = \begin{cases} 1, & y \cdot h_{\theta}(x) \leq 0 \\ 0, & y \cdot h_{\theta}(x) > 0 \end{cases}$$

when  $G \uparrow$  &  $h_{\theta}$  have different signs  $\rightarrow$  we incur a loss.

$l_{01}(h_{\theta}(x), y)$ : Is this a convex loss fn.? **No!**

Not convex, harder to optimize

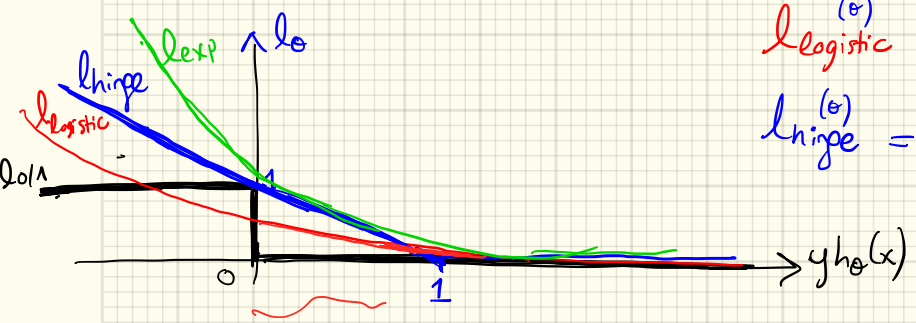


$\therefore$  Alternative loss functions are used in classification.

$$l_{\text{logistic}}^{(0)} = \log(1 + \exp(-y \cdot h_{\theta}(x)))$$

$$l_{\text{hinge}}^{(0)} = \max(\underbrace{1 - y \cdot h_{\theta}(x)}, 0)$$

$\geq 0$  when  $y \cdot h_{\theta}(x) \leq 1$   
 $l_{\text{exp}}^{(0)} = \exp(-y \cdot h_{\theta}(x))$

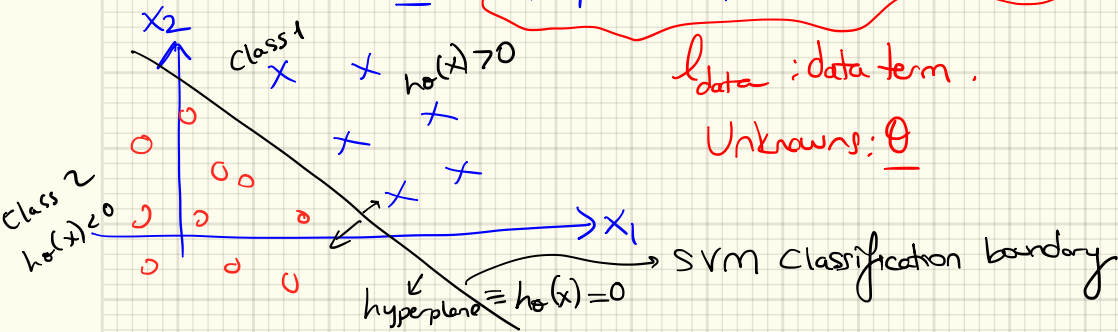


logistic, hinge, l\_exp: convex losses  $\therefore$  we typically prefer them over 0/1 loss.

$\exists$  convex optimization tools (eg. cvxpy) to solve such convex optim. problems.

Recall: **SVM**: Linear SVM solves the ML optimization problem w/ hinge loss & linear hypothesis function.

$$\underline{\theta}^* = \arg \min_{\underline{\theta}} \sum_{i=1}^m \max \{ 1 - y_i h_{\theta}(x_i), 0 \} + d \cdot \|\underline{\theta}\|_2^2$$



data: data term.

Unknowns:  $\underline{\theta}$

add a **Regularizer term**  
 $d$ : a regularization parameter.

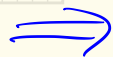
**Logistic Regression**

Using linear hypothesis for logistic loss:

ML problem:

$$\underline{\theta}^* = \arg \min_{\underline{\theta}} \sum_{i=1}^m \log(1 + \exp(-y_i \underline{\theta}^T x_i)) + d \|\underline{\theta}\|_2^2$$

Loss fn. is convex & smooth (note: hinge loss was not differentiable at  $y_i h_{\theta} = 1$ )



$$\Rightarrow \mathcal{L}_{\text{logistic}}(\theta) \equiv \max_{\theta} \sum \log \frac{1}{1 + e^{-(y_i \theta^T x_i)}} \rightarrow \text{we see logistic or sigmoid fn. in this loss.}$$

Q. <sup>(3)</sup> Devise an optimization based on GD to solve logistic regression.

Multi-Class Classification: Output is in  $\{1, \dots, k\}$   
 eg. Digit classification  $\{0, 1, 2, \dots, 9\}$   
 10-class problem

Approach 1: One vs All method:

Build  $k$  different binary classifiers:  $h_{\theta_i}, i=1:k$

Goal: Predict class  $i$  vs. others

$$\hat{y} = \arg \max_i h_{\theta_i}(x)$$

↳ each binary classifier model.

Approach 2: Next time.