



Artificial Intelligence and Machine Learning

Classification / Logistic Regression

Lecture 2: Outline

- Linear Regression (Review)
- Linear Regression (Remaining concepts)
- Intro to Classification
- Logistic Regression (Linear Classifier)
- Classification Metrics
- Other ML Models

Recap

Design your model

- Input scalar linear model (line fitting)
- Fitting polynomials (synthetically designing features from a one-dimensional input)

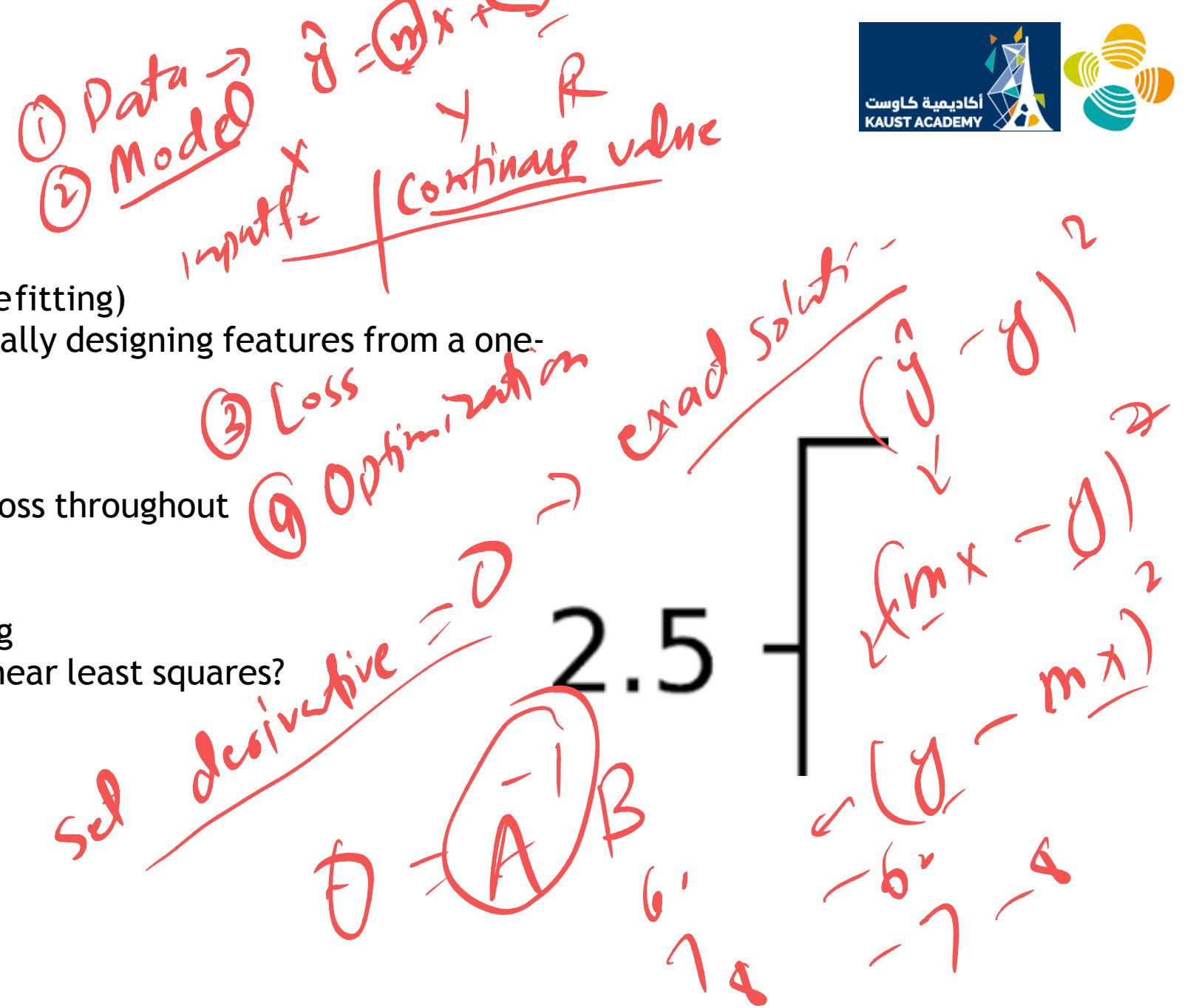
Design your loss function

- We used mean squared error loss throughout

Finding optimal parameter fitting

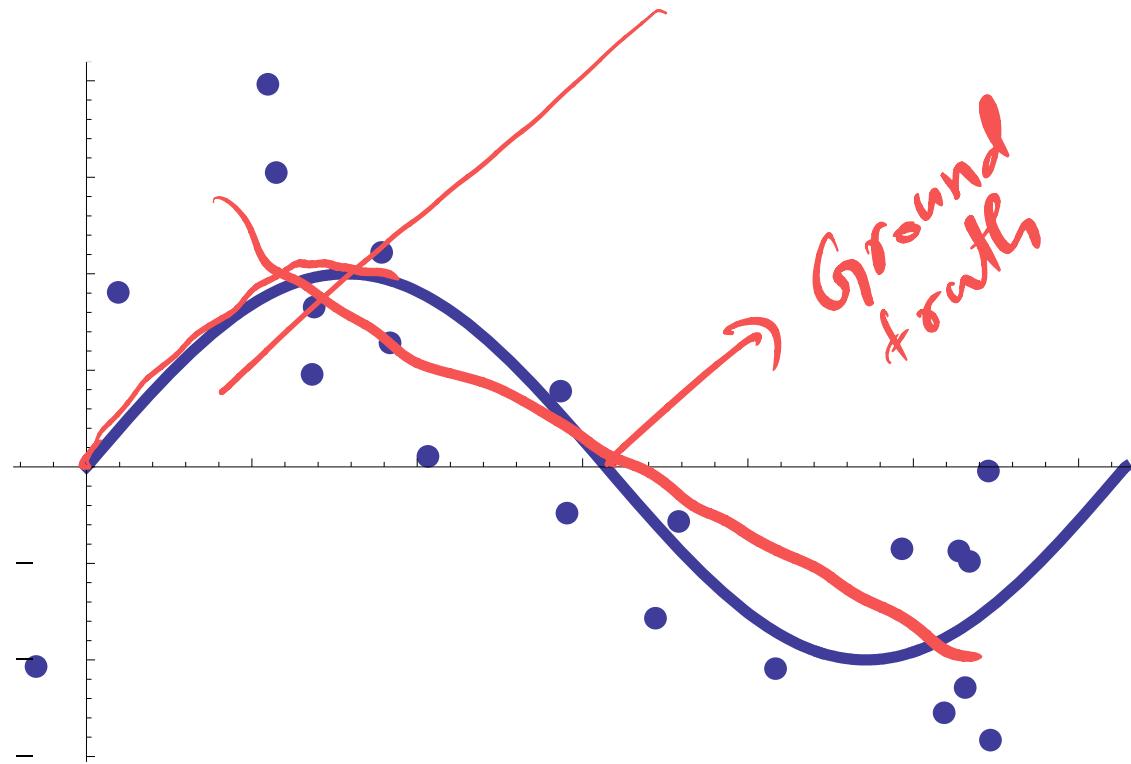
- Closed form solution to the linear least squares?
- Why is it linear least squares?
- Solution is closed form

gradient descent



Bias and Variance

- What if Y has a non-linear response?

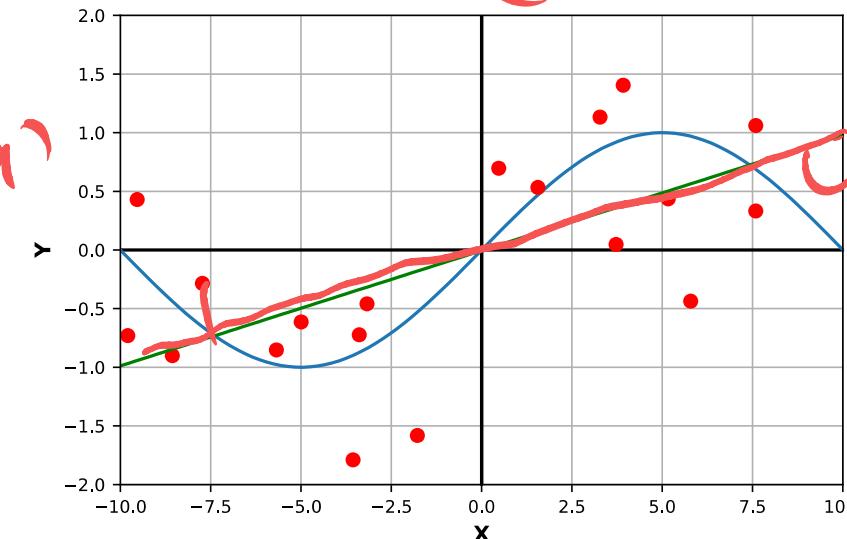


- Can we still use a linear model?

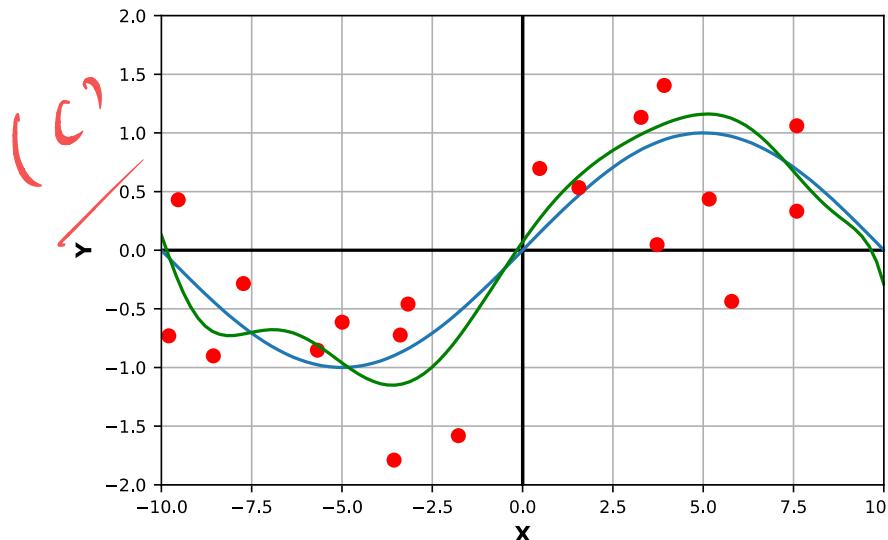
→ 2^m feature, 3rd feature



$\{1, \cancel{x}\}$

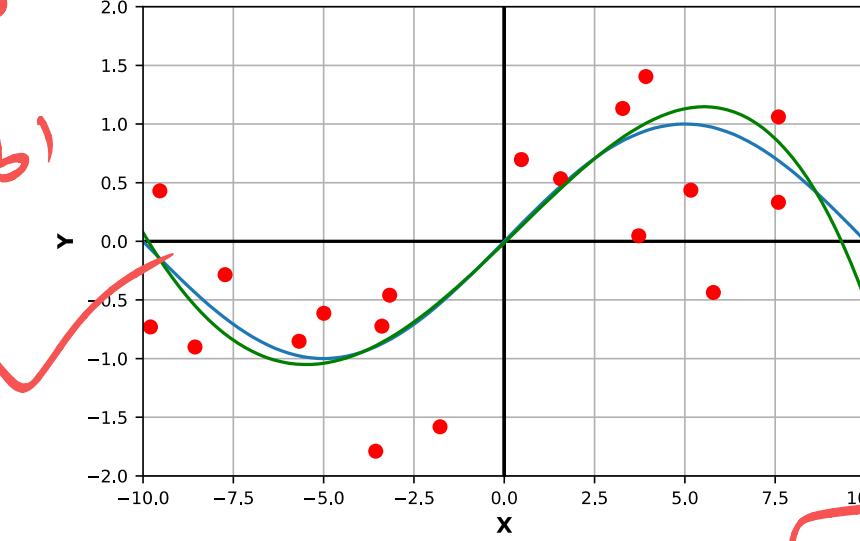


$\{1, x, x^2, \dots, x^9, x^{10}\}$



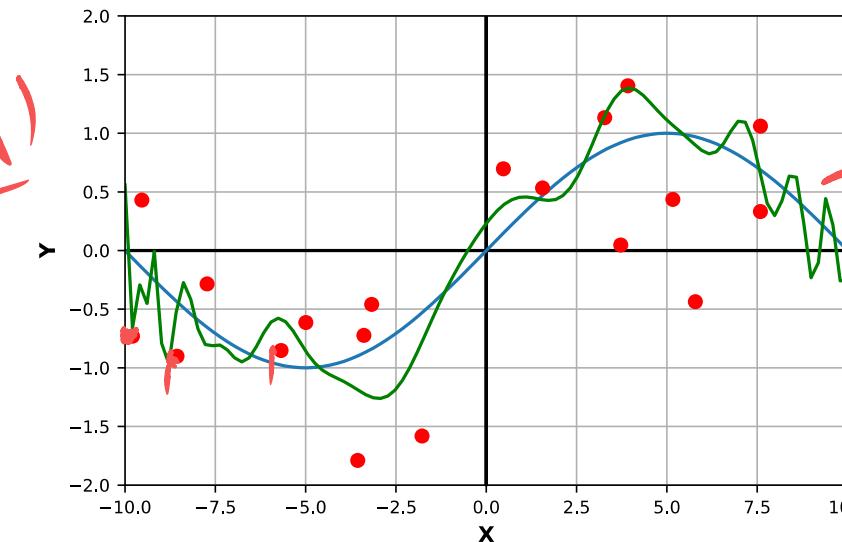
underfit

$\{1, \cancel{x}, \cancel{x^2}, \cancel{x^3}, \cancel{x^4}\}$



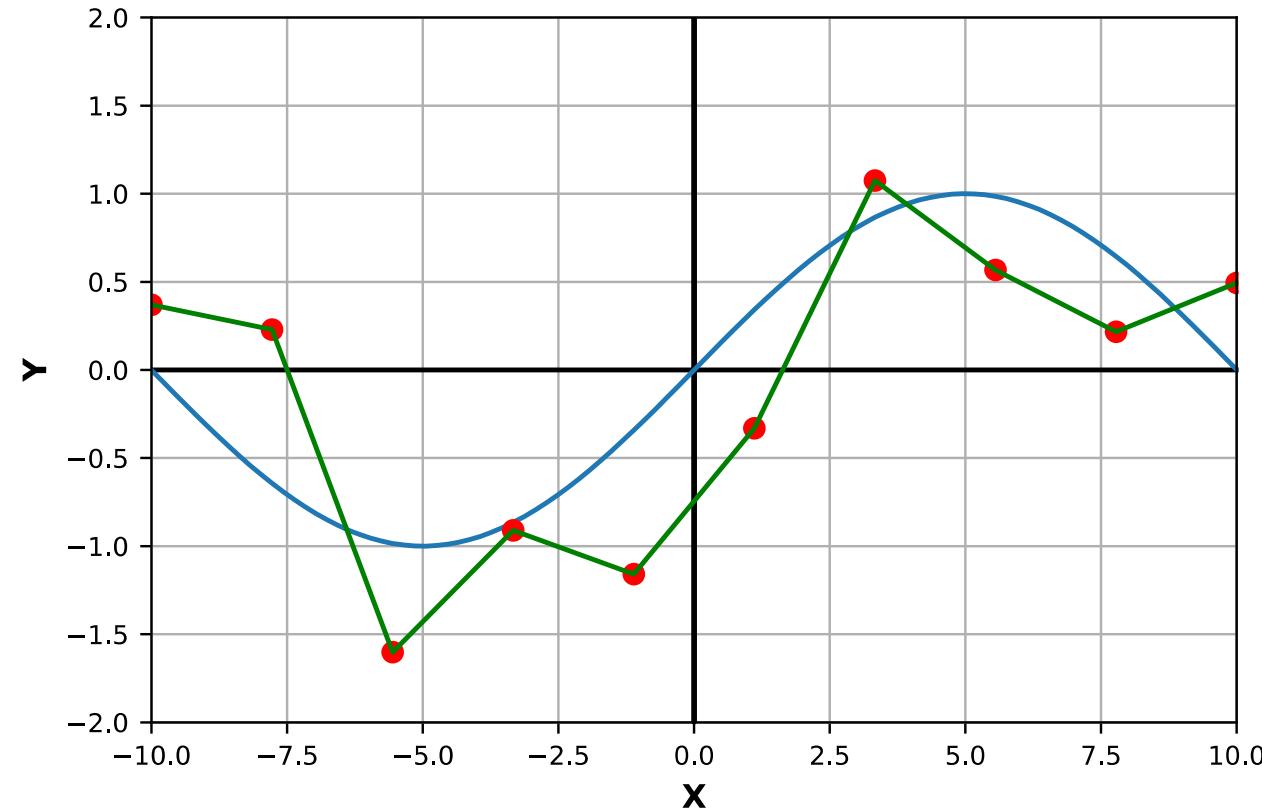
green \rightarrow model
blue \rightarrow actual

$\{1, x, x^2, \dots, x^{99}, x^{100}\}$



overfitting

Real Bad Overfit?



lowered error rates

Bias-Variance Tradeoff

- So far we have minimized the error (loss) with respect to **training data**

- Low training error does not imply good expected performance: **over-fitting**

- We would like to reason about the **expected loss** (**Prediction Risk**) over:

- Training Data: $\{(y_1, x_1), \dots, (y_n, x_n)\}$

- Test point: (y^*, x^*)

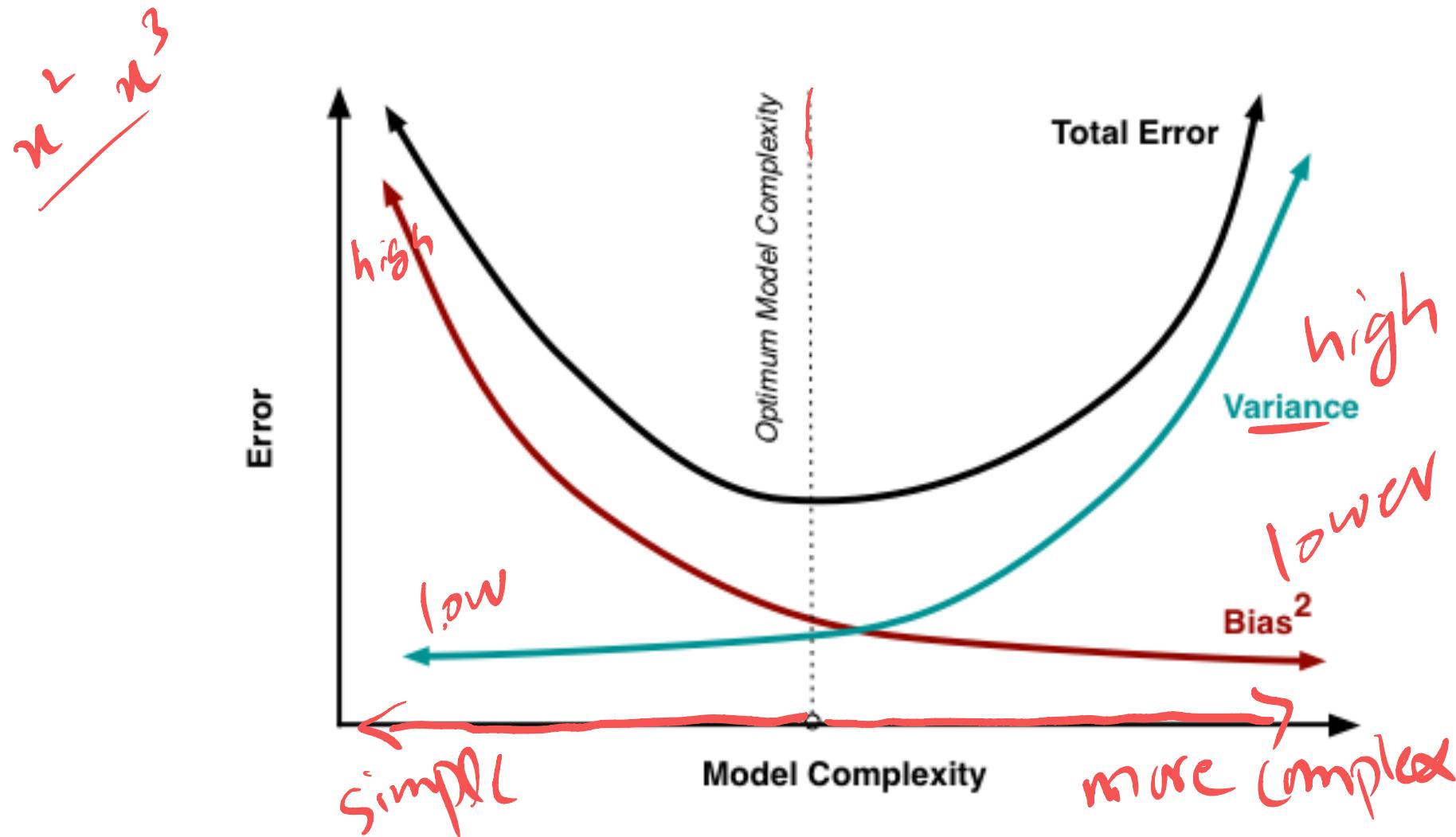
- We will decompose the expected loss into:

ansren

$$\mathbb{E}_{D, (y^*, x^*)} [(y^* - f(x^* | D))^2] = \text{Noise} + \text{Bias}^2 + \text{Variance}$$

Loss $\approx \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Bias Variance Plot



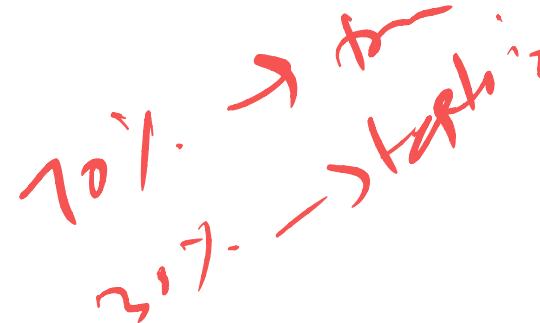
Data Split



- To ensure your model doesn't overfit to the training data, you should have another **subset** called **testing data**.
- You will evaluate your model against this subset, and based on its **metric score (e.g. accuracy)** you will decide if it's overfitting or not.
- But how should I split my data?

Out of 100 → 20% sample

Data Split



- **Hold-out set:**

- A portion of the dataset set aside and not used during training.
- E.g. 80% for training and 20% for testing.

Issues:

- Imagine you have these labels: [1, 1, 2, 2, 2, 3, 3, 3, 3, 3] and you took last 30% as test: [3,3,3]. You didn't include 1 and 2 in test!

Solution: Always shuffle before split: [3, 2, 3, 1, 3, 2, 3, 1, 3, 2] test: [1,3,2]

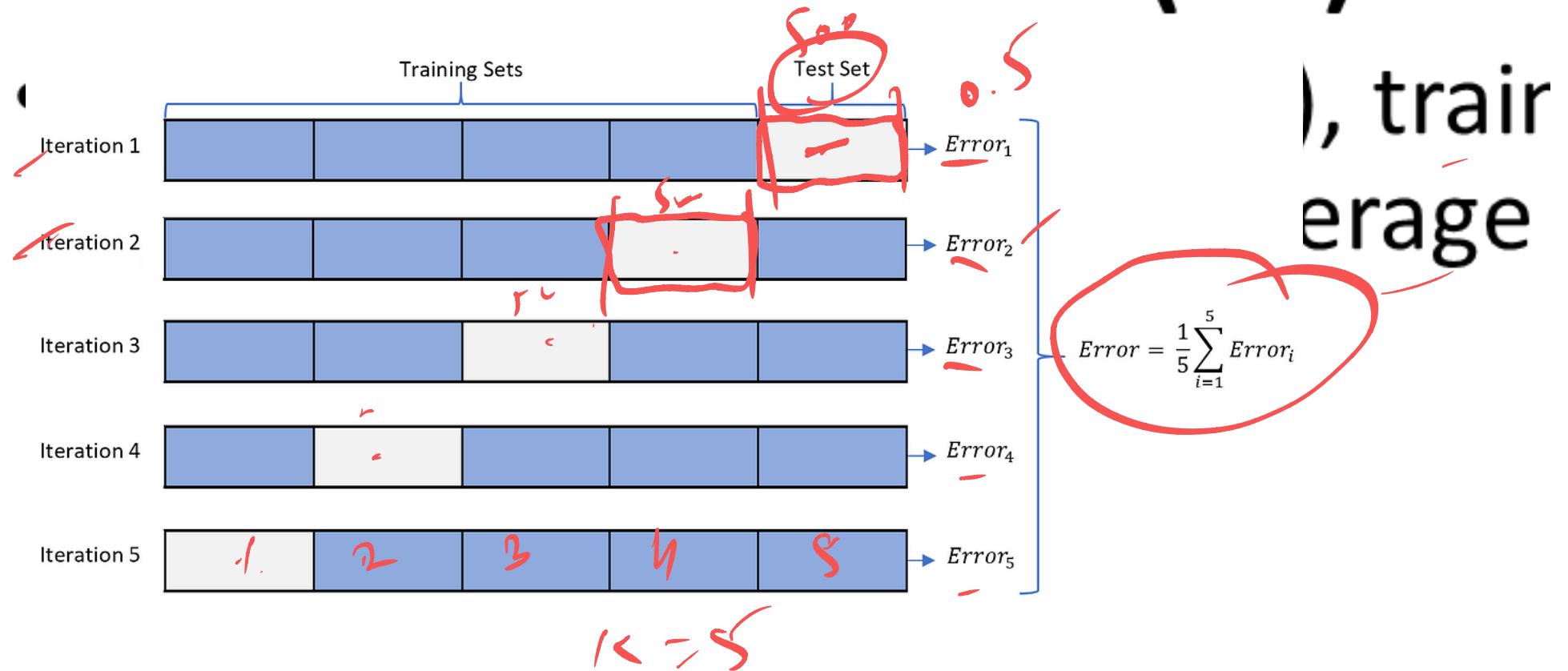
- My dataset is small. Taking 20% as test would not be representative!

Solution: Use KFold.

Data Split



- **K-Fold Cross Validation (CV):**



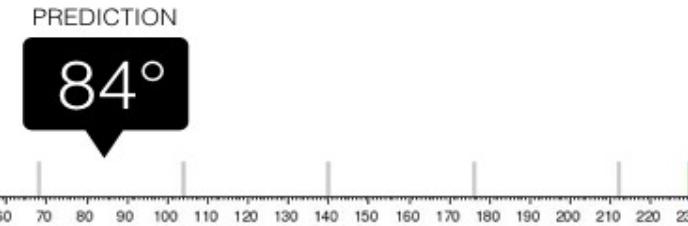
Regression VS classification

logistic
regression



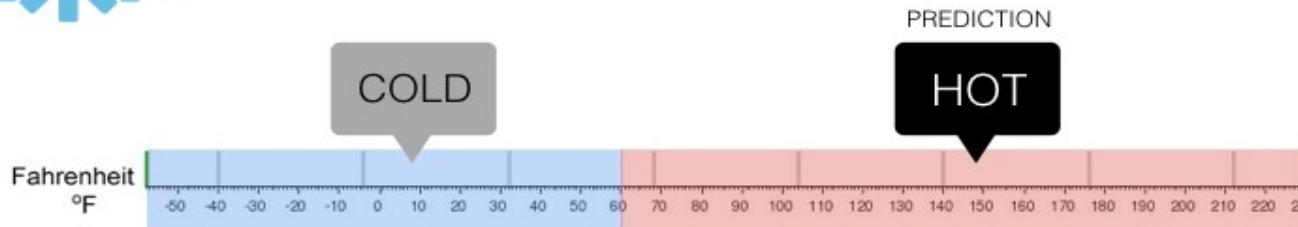
Regression

What is the temperature going to be tomorrow?



Classification

Will it be Cold or Hot tomorrow?



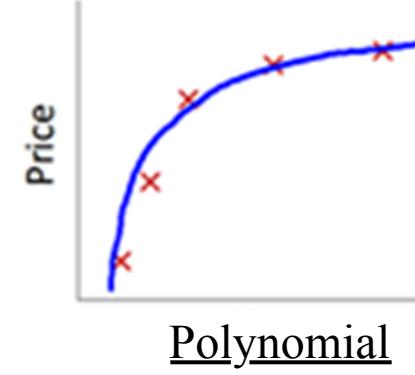
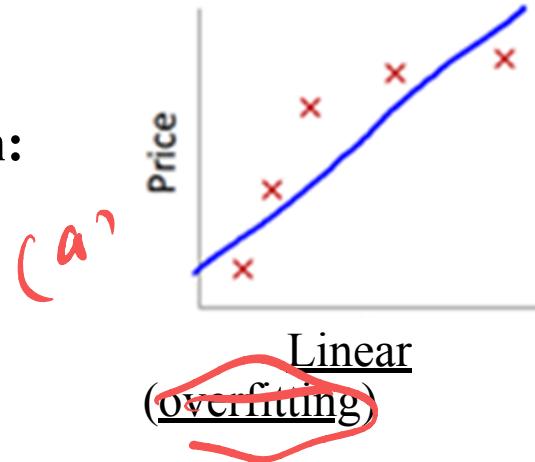
=> Continuous Values

rain
30°-
70°- clouds

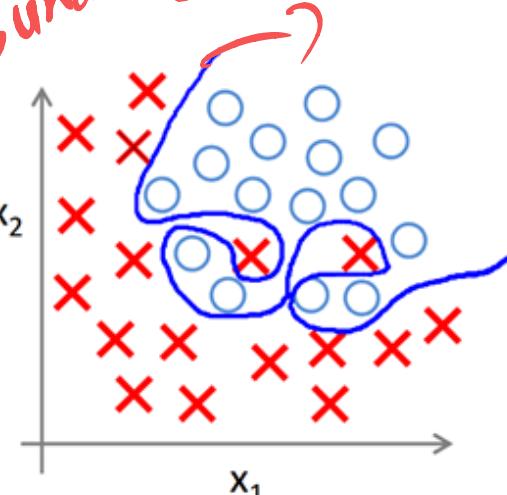
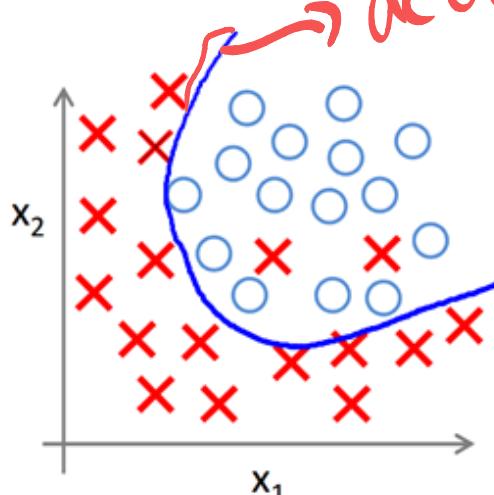
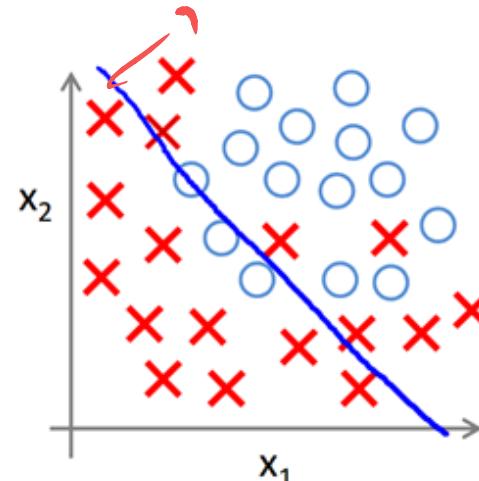
=> Discrete Values

Regression VS classification

Regression:



Classification:



decision boundary

30'

Classification Types

Binary Classification: Two possible outcomes (e.g., Yes/No, Spam/Not Spam).

How? Single Classifier.

He is a good person

positive/neutral/not

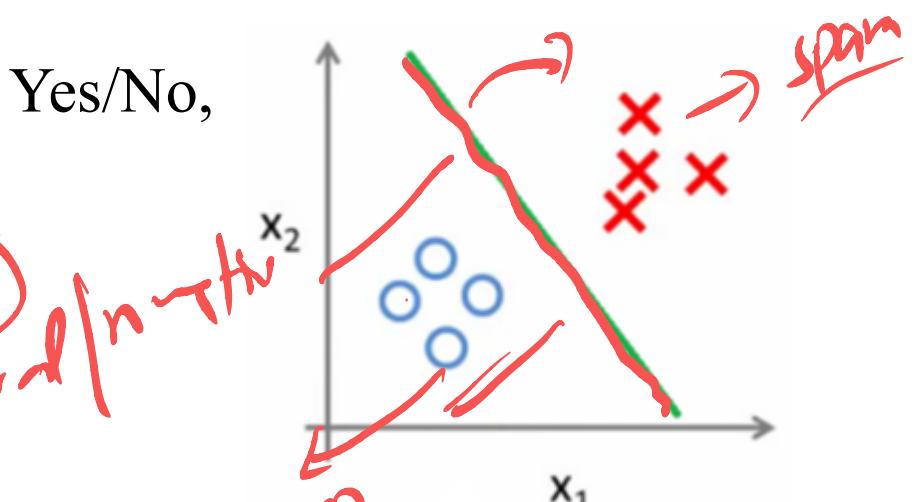
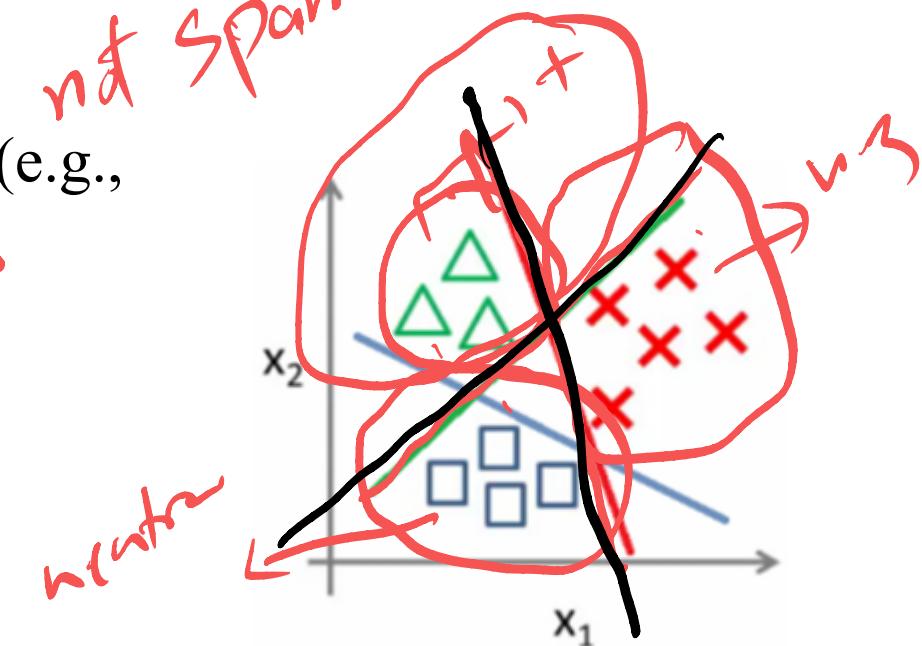
not spam

Multiclass Classification: More than two outcomes (e.g., Class A, B, C).

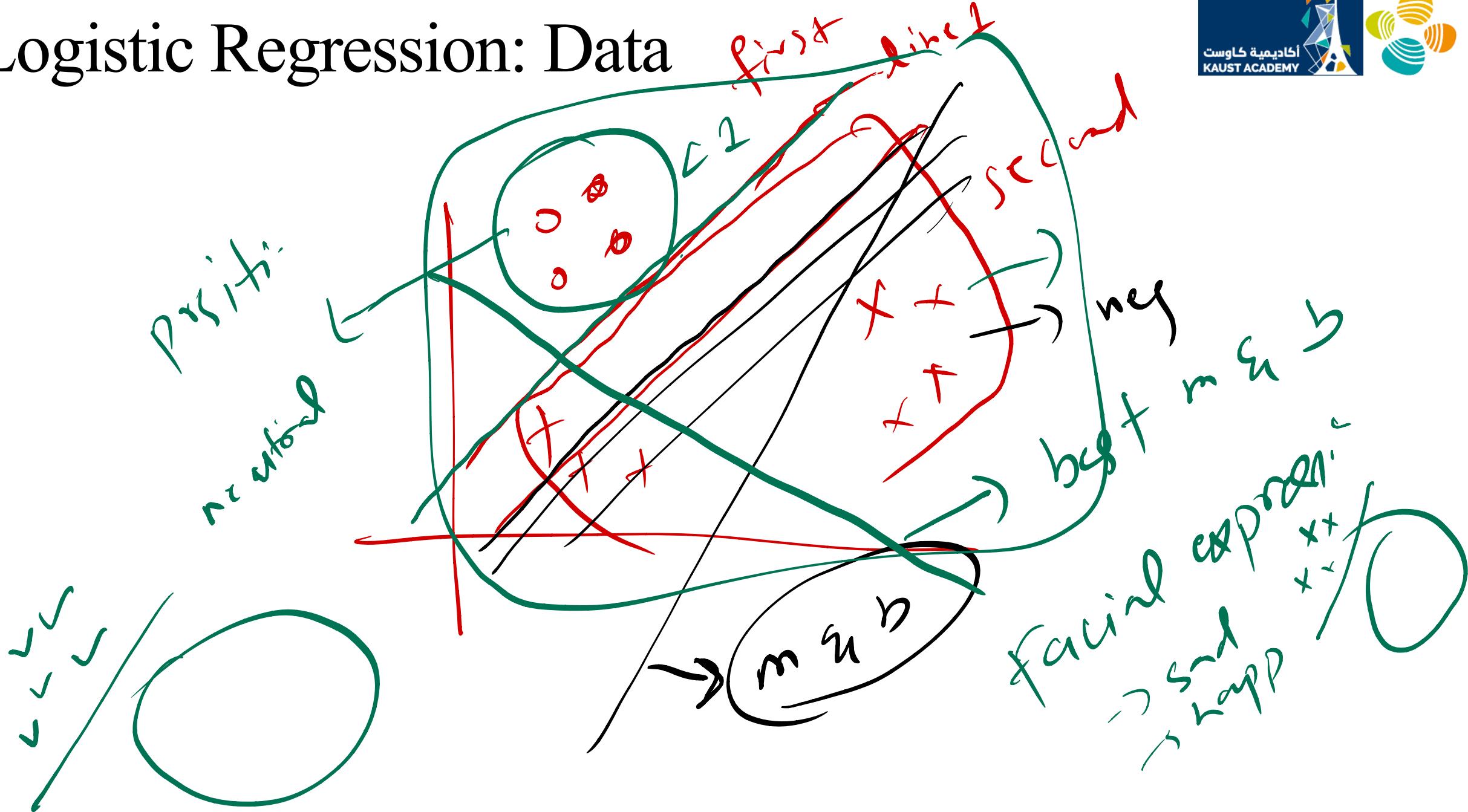
How? Multiple *Single Classifiers*.

sentiments classification

neutral



Logistic Regression: Data

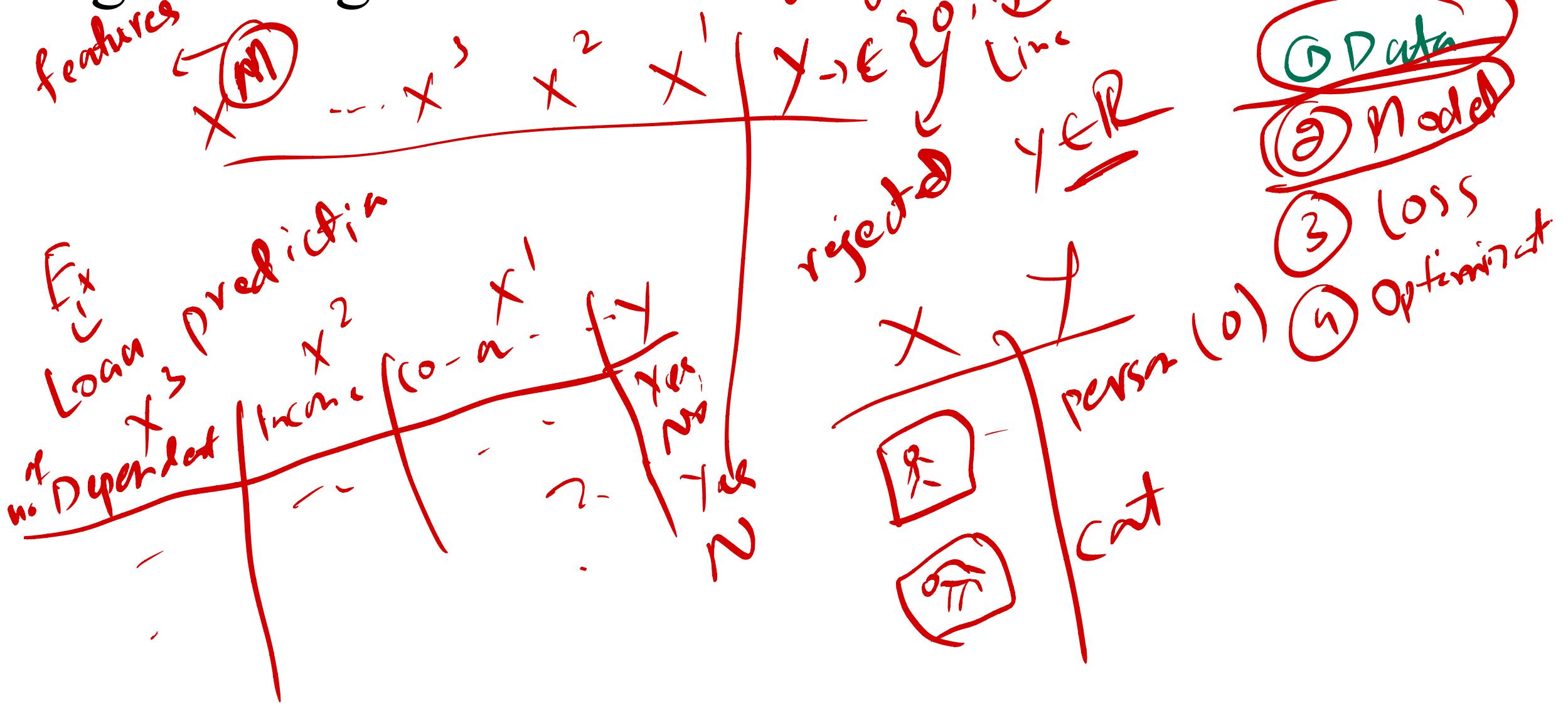


Logistic Regression: Data

c_1 happy vs
 c_2 Sad vs
 c_3 disgust vs
 sum → 1 one vs all



Logistic Regression: Data



Logistic Regression: Model



~~Linear regression~~

$\hat{y} = f(x)$ *continuous*

$\hat{y} = \text{f}(x)$ *discrete value*:

0.99

0.23

image regression

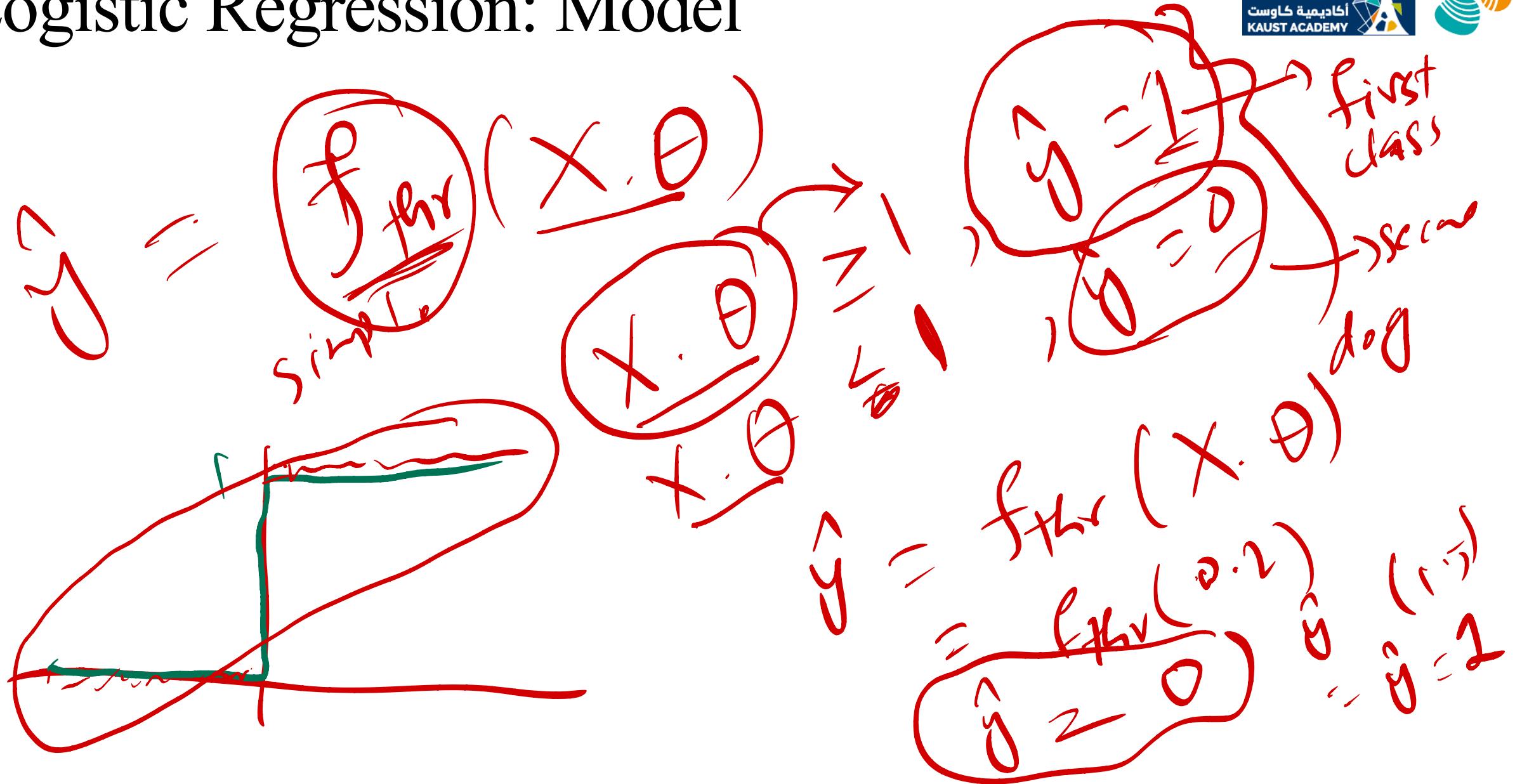
predicts memorability of image

$\hat{y} = f(x)$ *R*

0.19

Linear regression

Logistic Regression: Model



Logistic Regression: Model

- chances of min probability $\rightarrow \hat{y} = \frac{1}{1 + e^{-x^T \theta}}$
- computation of gradients $\rightarrow \nabla_{\theta} L(\theta) = \sum_{i=1}^m (\hat{y}_i - y_i) x_i$
- generalise to multi-class \rightarrow softmax function $\rightarrow \hat{y}_j = \frac{e^{x_j^T \theta_j}}{\sum_k e^{x_k^T \theta_k}}$
- n-confidence equally \rightarrow dist \rightarrow $\hat{y} = \arg \max_j \hat{y}_j$

Logistic Regression: Model

sigmoid function

$f(x) = \frac{1}{1 + \exp(-x)}$

G_0 approaching 0 as $x \rightarrow -\infty$

$\delta_0 = \frac{1}{1 + \exp(-x \cdot \theta)}$

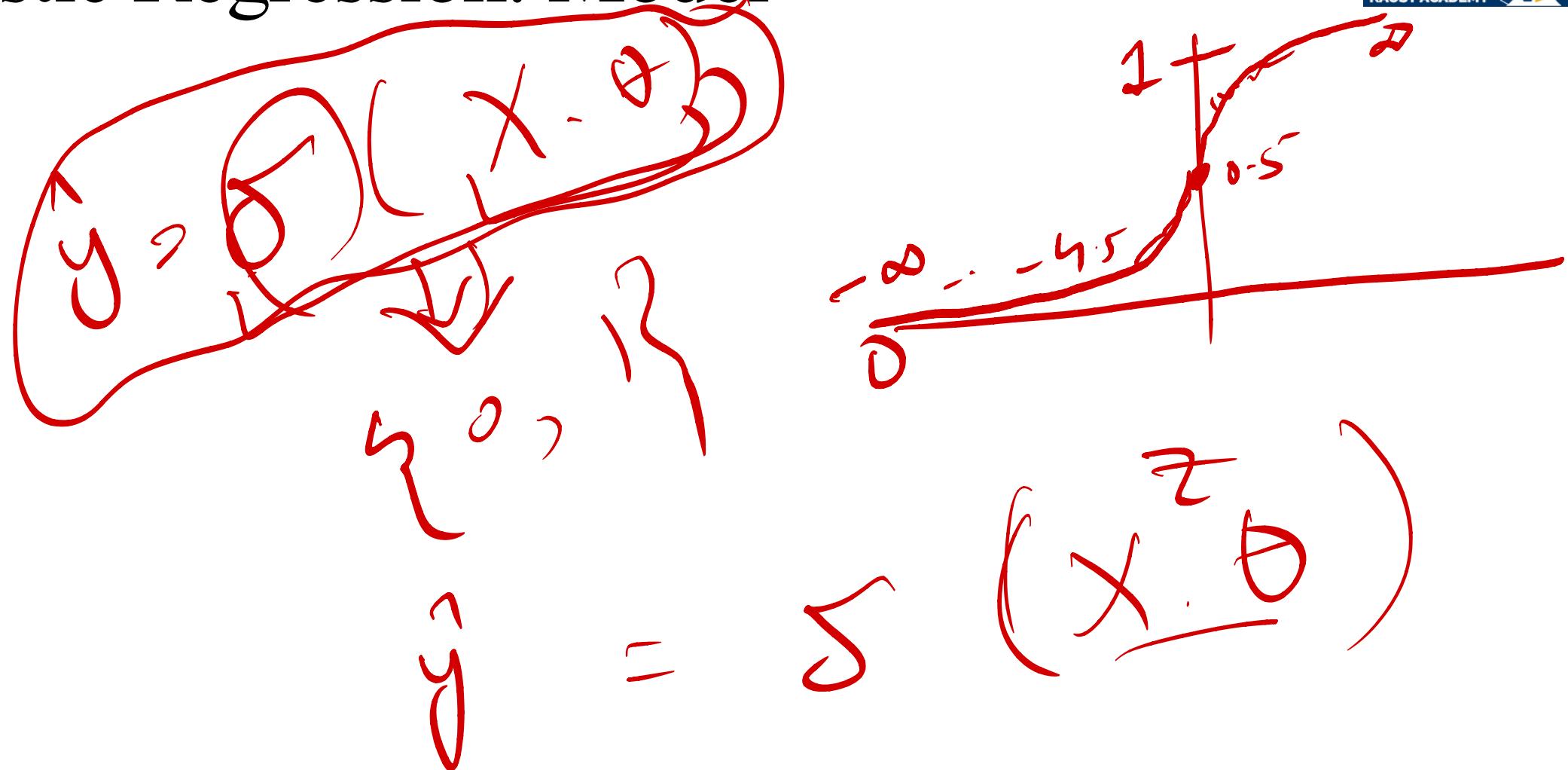
changing factor

$z = 0, \frac{1}{1 + \exp(-0)} = \frac{1}{1+1} = 0.5$

$z = \infty, \frac{1}{1 + \exp(-\infty)} = \frac{1}{1+0} = 1$

$z = -\infty, \frac{1}{1 + \exp(-(-\infty))} = \frac{1}{1+0} = 1$

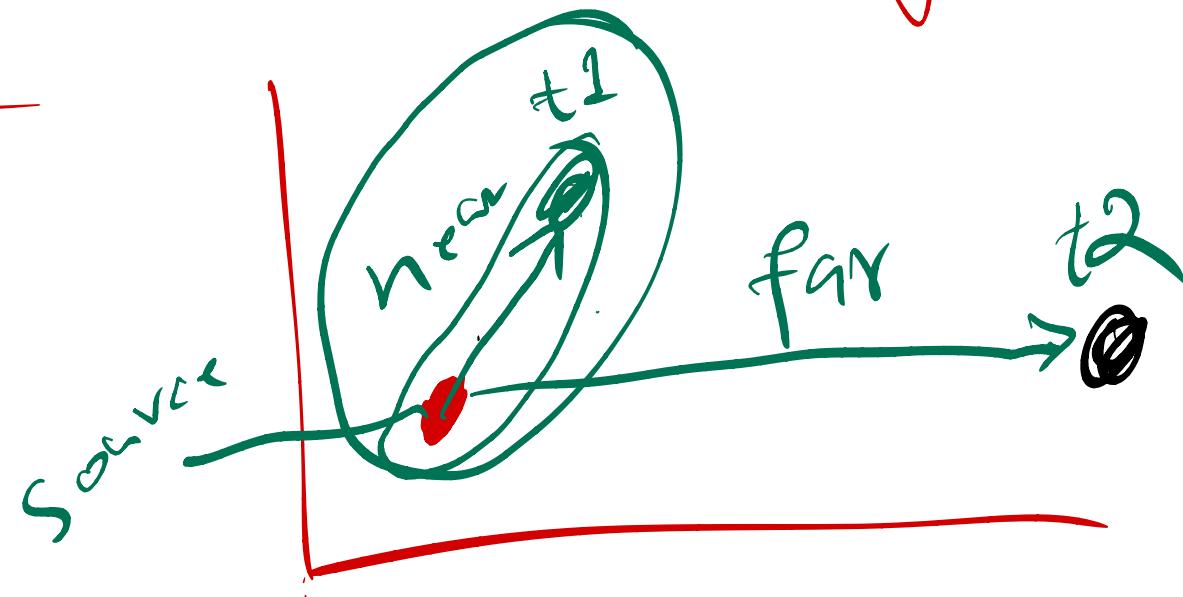
Logistic Regression: Model



Logistic Regression: Loss

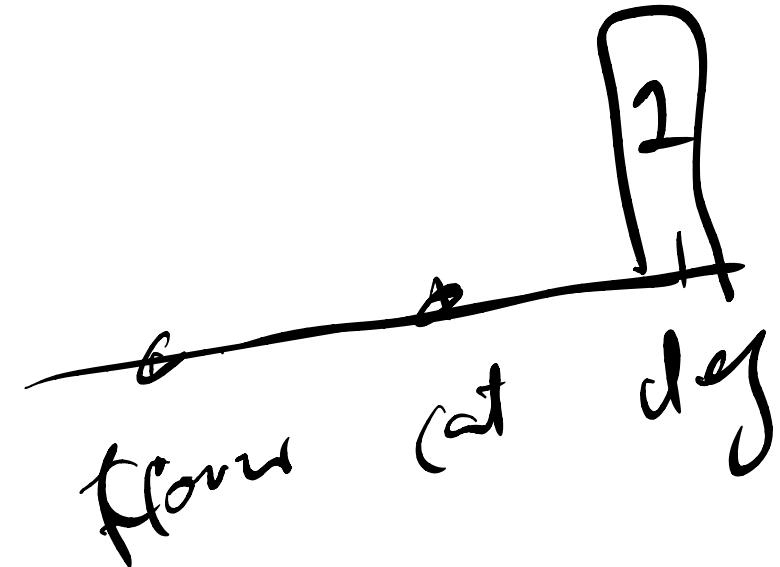
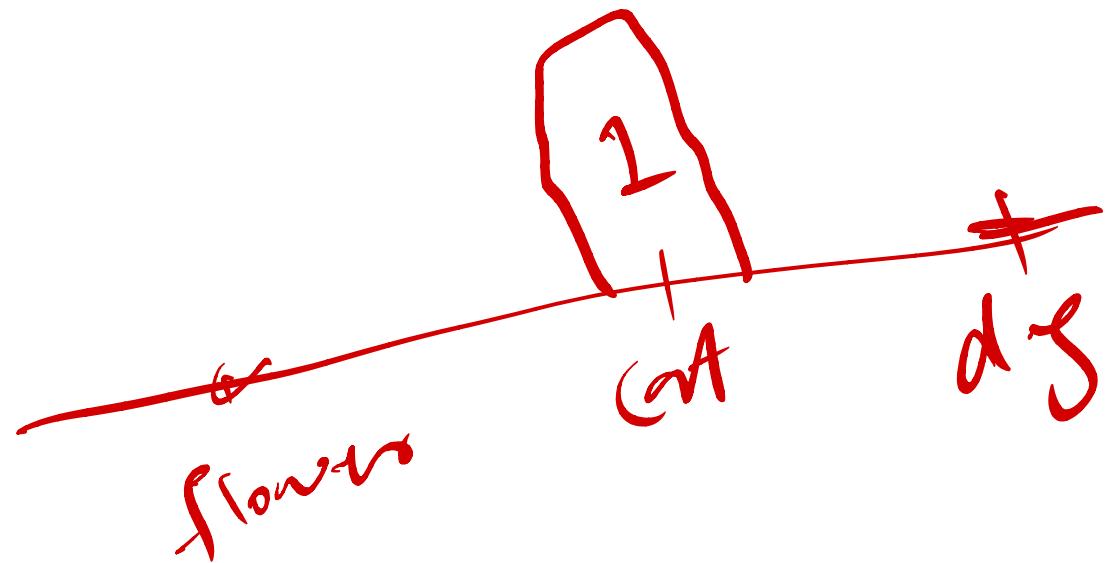
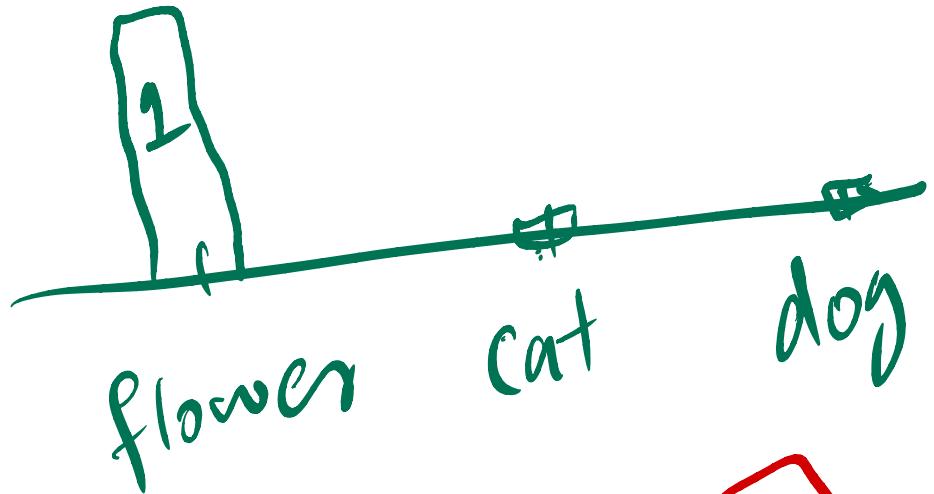
$$MSE = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

Compute the difference
between continous
values.

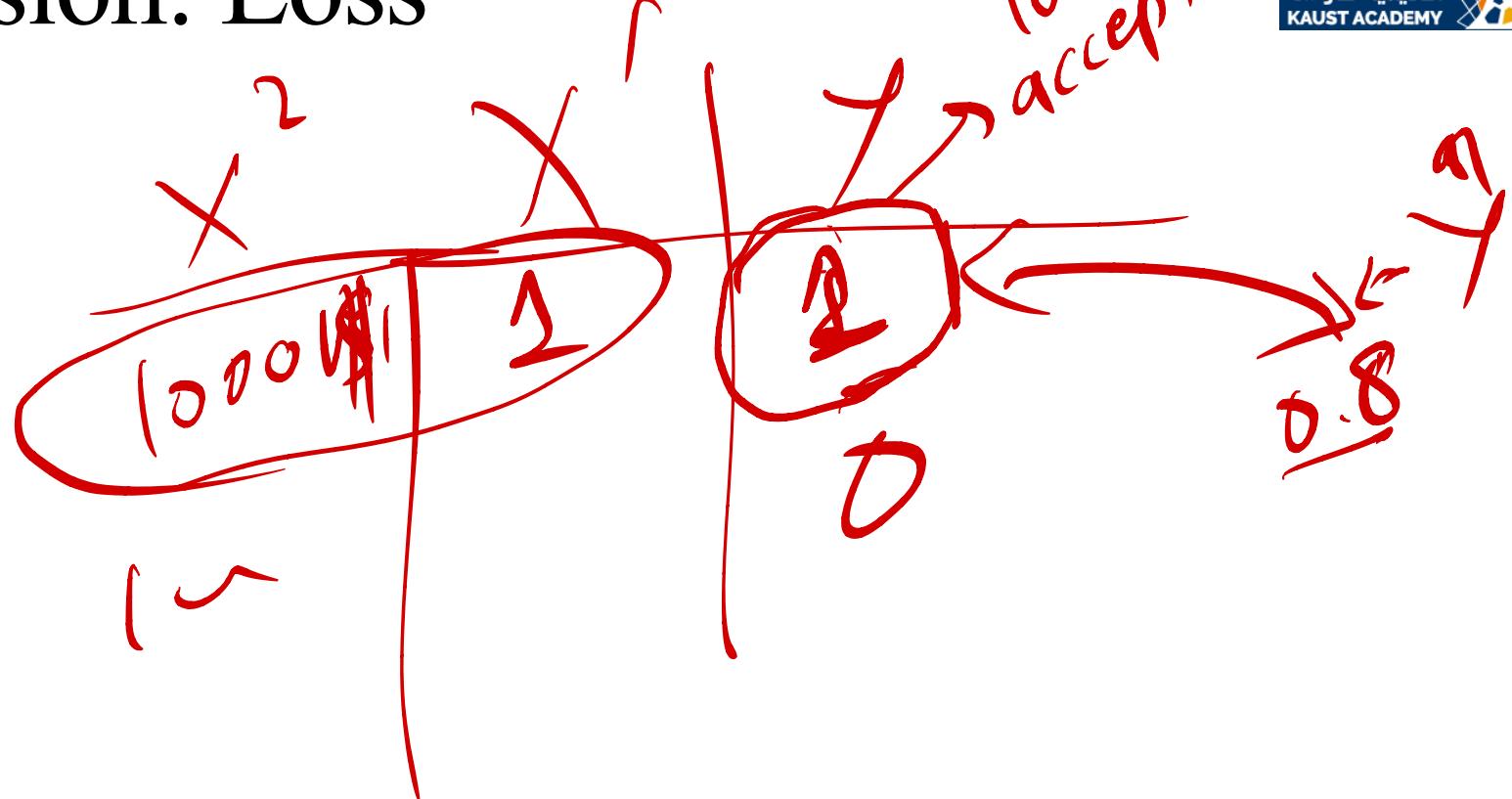


Logistic Regression: Loss

Logistic
classification
→ probabilities



Logistic Regression: Loss



KI
diagram
cross entropy

Logistic Regression: Loss

$$0 - 2(\hat{y} - y) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$y = 0$, $\hat{y} = 0$

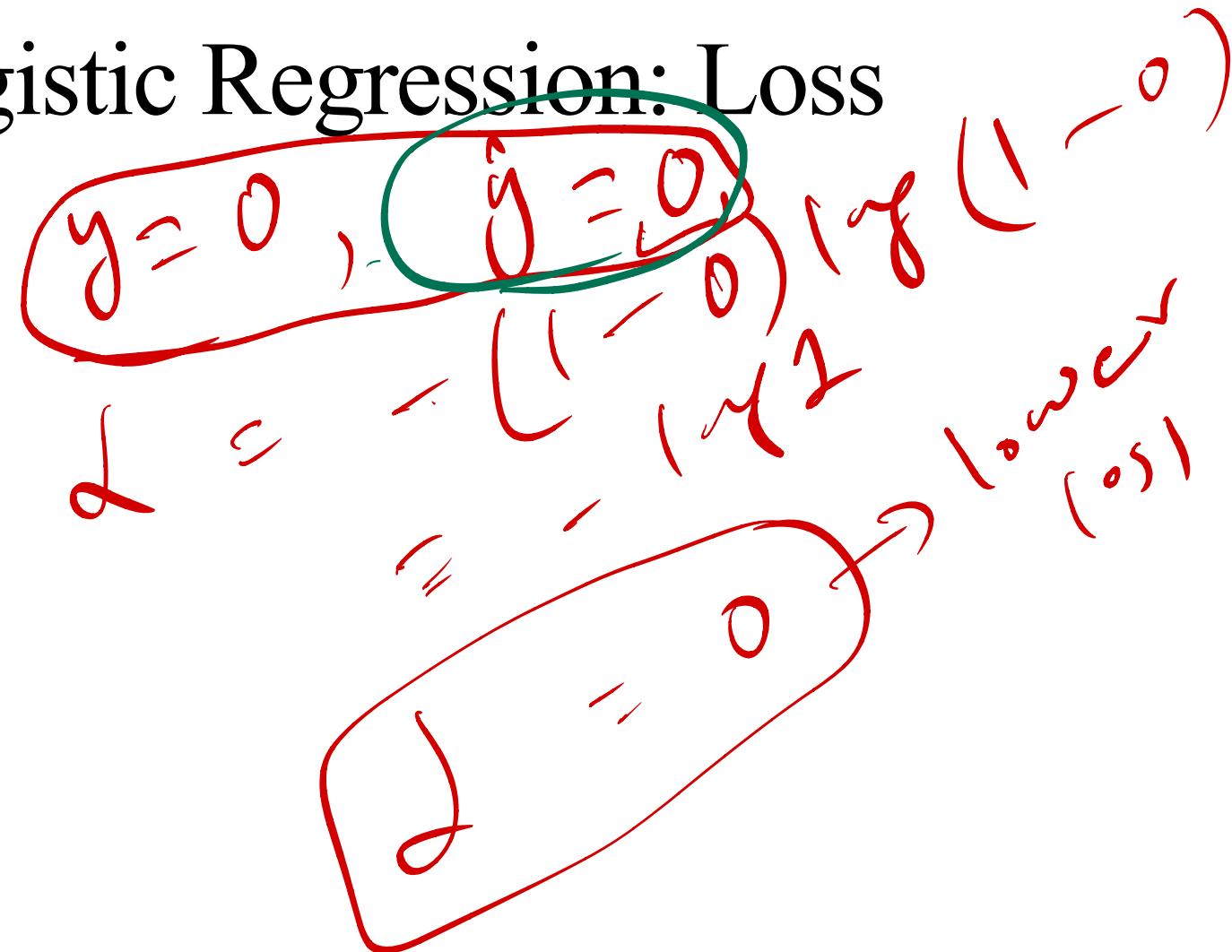
$$J = - (1-y) \log(1-\hat{y})$$

$$J = - y \log \hat{y}$$

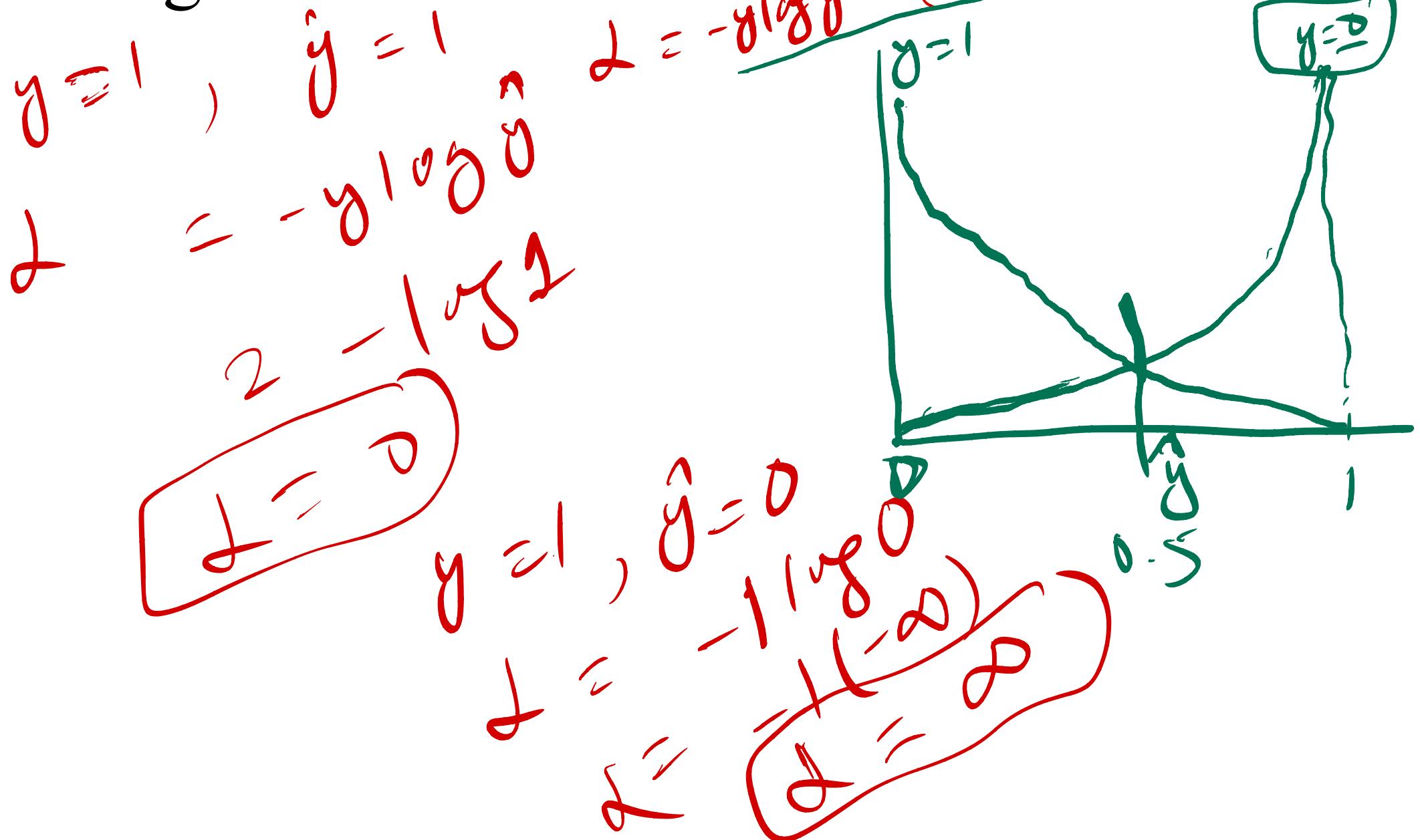
$$J = - \frac{1}{n} \sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)$$

very high loss

Logistic Regression: Loss

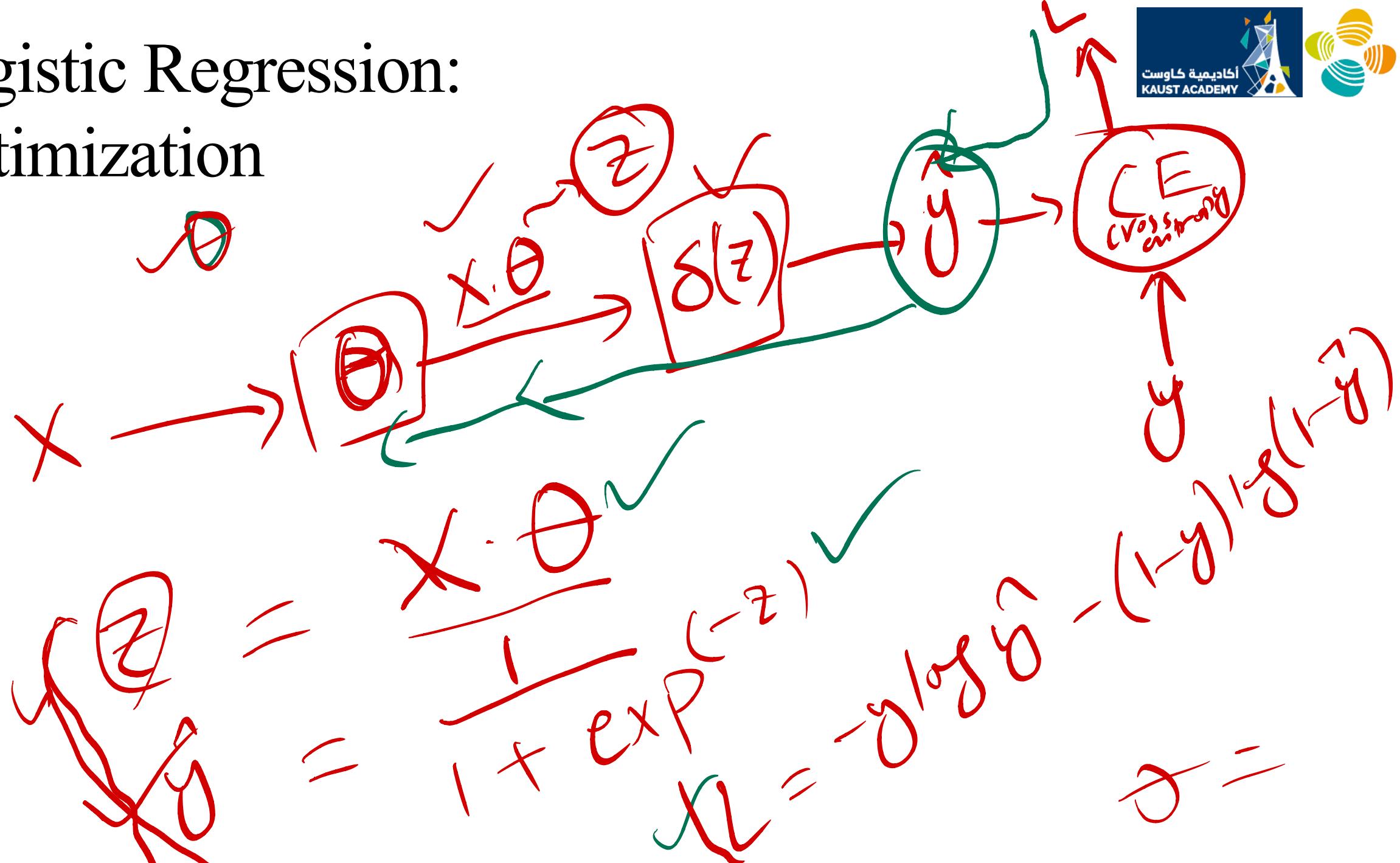


Logistic Regression: Loss





Logistic Regression: Optimization



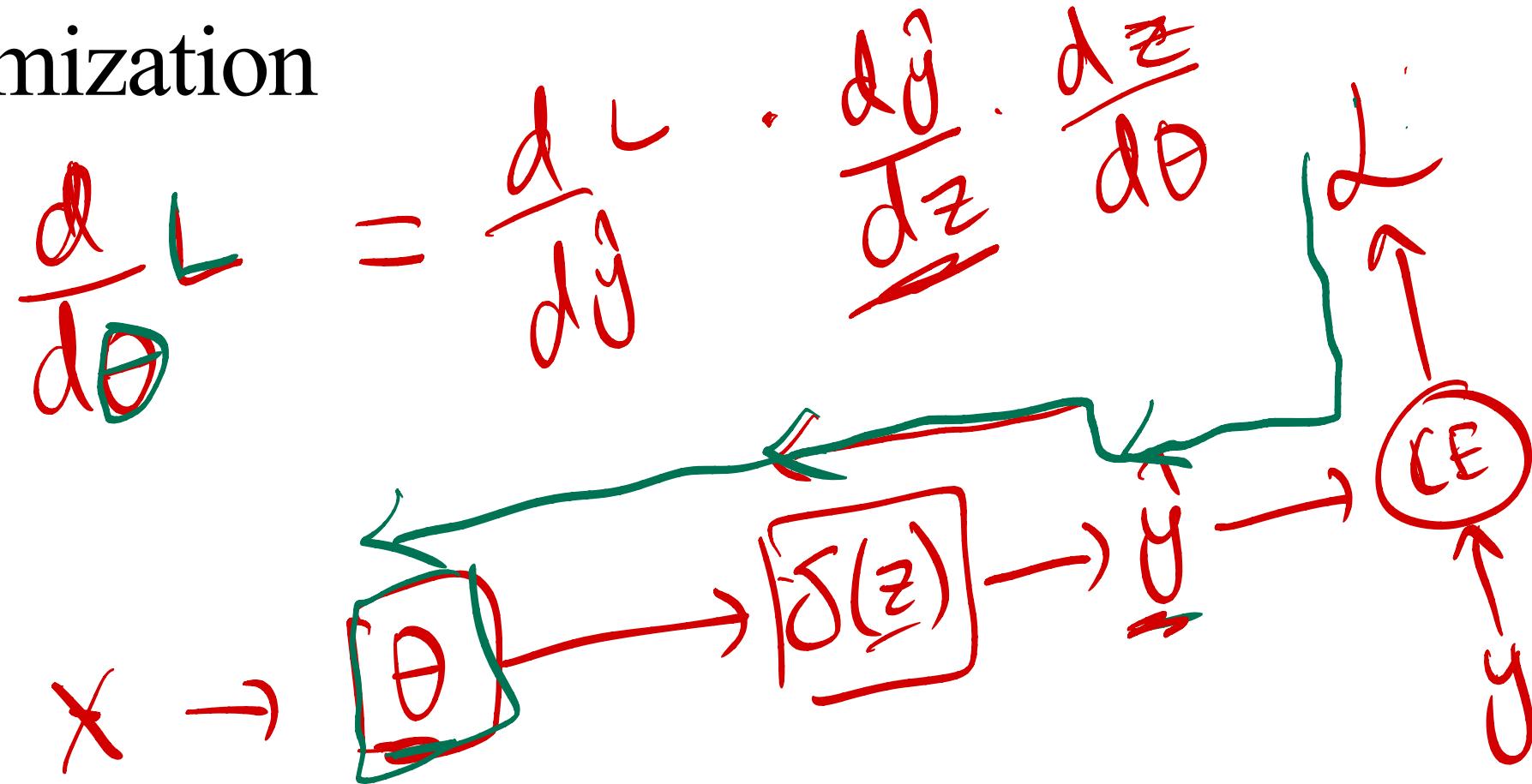
Logistic Regression:

Optimization

Lineal \rightarrow set derivative = 0
non-linear → open system linear equations

Logistic Regression:

Optimization

$$\frac{dL}{d\theta} = \frac{dL}{dy} \cdot \frac{dy}{dz} \cdot \frac{dz}{d\theta}$$


The diagram illustrates the backpropagation process for a single neuron in a neural network. An input vector x is shown entering from the left. The neuron's parameters, represented by a red circle labeled θ , are multiplied by the input x . The result is passed through an activation function $\delta(z)$, which is enclosed in a red box. The output of the neuron is \hat{y} , shown in red. The target value y is also shown. A red circle labeled "CE" represents the cross-entropy loss function. The error signal, calculated as the difference between \hat{y} and y , is propagated back through the neuron to update the parameters θ . The error signal is also used to calculate the gradient of the loss function with respect to the parameters θ .

Logistic Regression: Optimization

$$\frac{\partial L}{\partial \hat{y}} = \frac{d}{d\hat{y}} [y \log \hat{y} - (1-y) \log(1-\hat{y})]$$

$$\frac{\partial L}{\partial y} = \hat{y} - y$$

$$\frac{\partial L}{\partial \hat{y}} = \hat{y}(1-\hat{y})$$

Logistic Regression:

Optimization

$$\frac{d\hat{y}}{dz} = \frac{d}{dz} \left(\frac{1}{1 + \exp(-z)} \right)$$

$$\frac{d\theta}{dz} = \frac{d}{dz} \left(\frac{1}{1 + \exp(-z)} \right) \cdot g'(1 - g)$$

$$\frac{dy}{dz} = g'(1 - g)$$

$$\frac{dz}{d\theta} = \frac{d}{d\theta} (x^\top \theta) = x^\top$$

Logistic Regression:

Optimization

$$\frac{dL}{d\theta} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dz} \cdot \frac{dz}{d\theta}$$

$\frac{\partial L}{\partial \theta} = 0$

$\frac{dL}{d\theta}$ circled.

$X^T [g - \hat{y}] \cdot g(1-\hat{y})$ circled.

linear

Na

Logistic Regression: Optimization

A.G.D

$\theta^0 \rightarrow$ initialize (Randomly)

while until converged

$\theta^{t+1} = \theta^t - \frac{m}{n} \frac{\sum d\theta}{d\theta}$

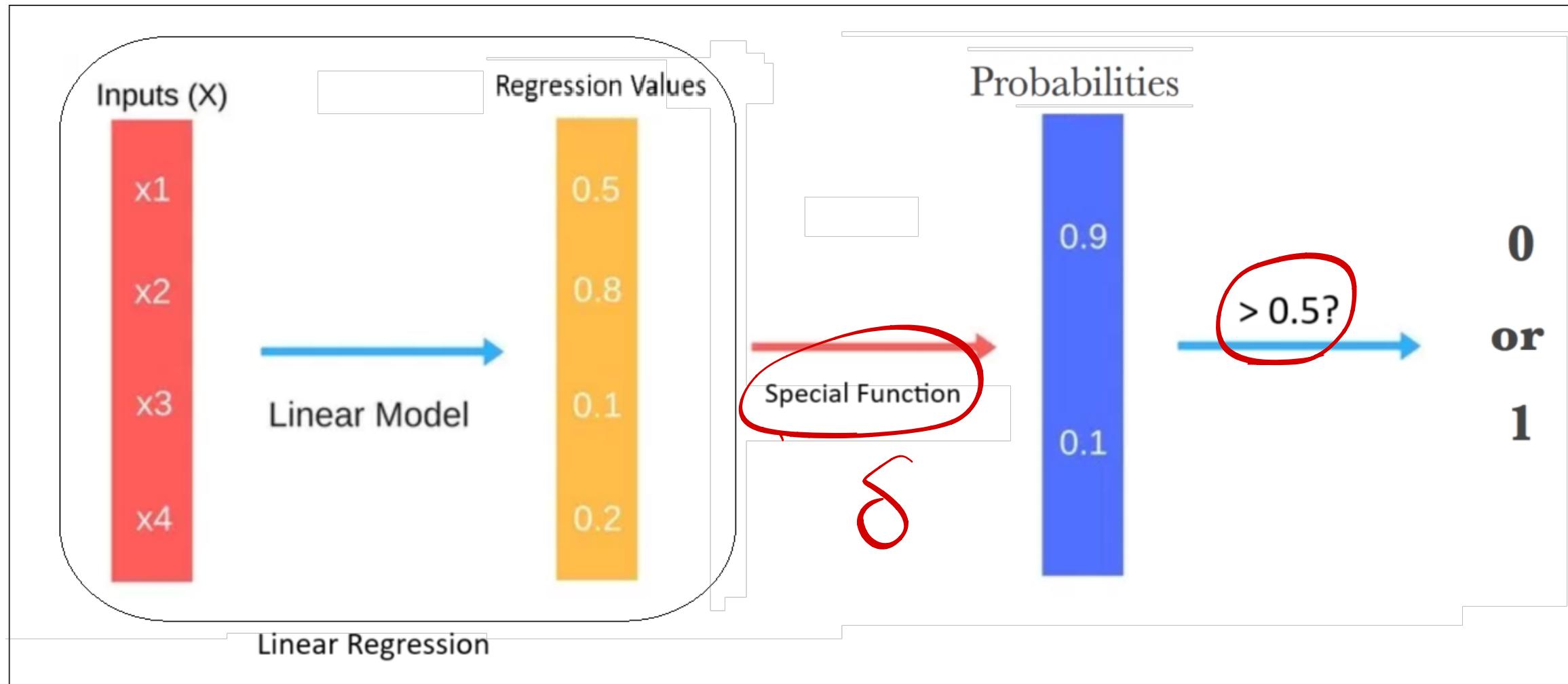
Classifiers

Classifiers are **regression** models with two key modifications:

- Pass the regression outputs through a **special function** to convert them into probabilities.
- Replace the Mean Squared Error loss with a **new loss function** designed for probabilities.

Applying these two to the Linear Regression should give Linear Classifier (called **Logistic Regression**).

Logistic Regression



Special Function

We are searching for a function that have the following characteristics:

1. Could convert any arbitrary input values to **[0,1] range (Probabilities)**.
2. **Smooth and Differentiable**. Because we want to find the minima later, right?

Any ideas?

Sigmoid Function (for Binary Classification)

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Assuming you have negative and positive classes, its output would be the probability of the positive class.

But what if we have multiclass problem? 🤔



$$\lim_{z \rightarrow -\infty} \sigma(z) = 0$$

$$\lim_{z \rightarrow \infty} \sigma(z) = 1$$

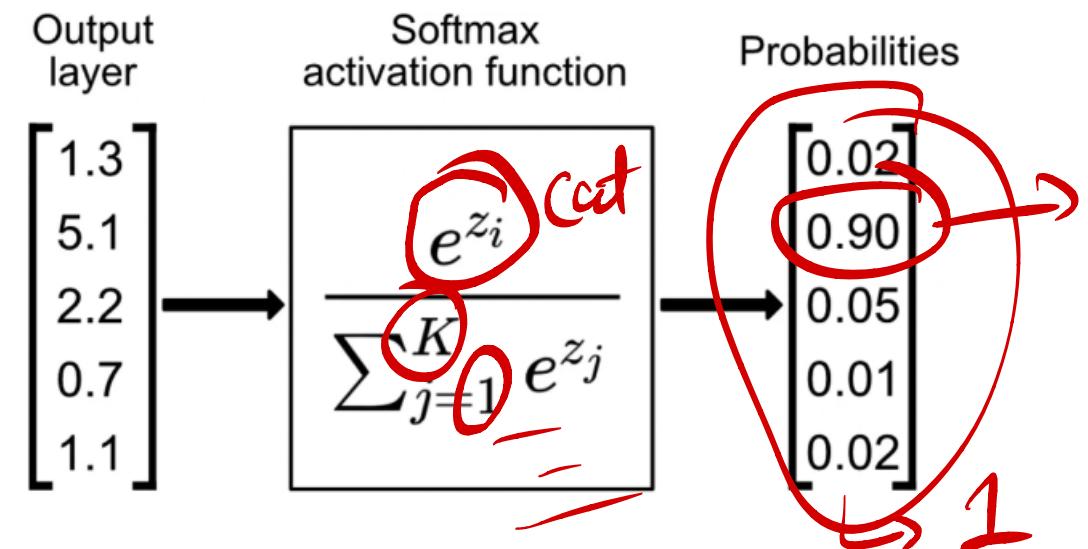
Softmax Function (for Multiclass Classification)

Softmax(\mathbf{z}) = softmax(\mathbf{z}) = $\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$)

E.g. Divide the cat value by the cat, dog and deer values to get the cat probability.

Assuming you have K classes, its output would be the probability of the ith class. So, you will run it K times to get the probability of each class.

Note: Sigmoid is a special case of Softmax. Optional: Can you prove it?



Classification LOSS

We are searching for a loss that have the following characteristics:

1. **Probabilistic behaviour**: Should work with probabilities, not raw numbers.
2. **Sensitivity**: It should heavily penalize incorrect confident predictions (e.g., predicting 0.9 when the true label is 0).
3. **Differentiable**: Must be smooth and differentiable to enable optimization.

Any ideas?

Classification LOSS

We are searching for a loss that have the following characteristics:

1. **Probabilistic behaviour:** Should work with probabilities, not raw numbers.
2. **Sensitivity:** It should heavily penalize incorrect confident predictions (e.g., predicting 0.9 when the true label is 0).
3. **Differentiable:** Must be smooth and differentiable to enable optimization.

Any ideas?

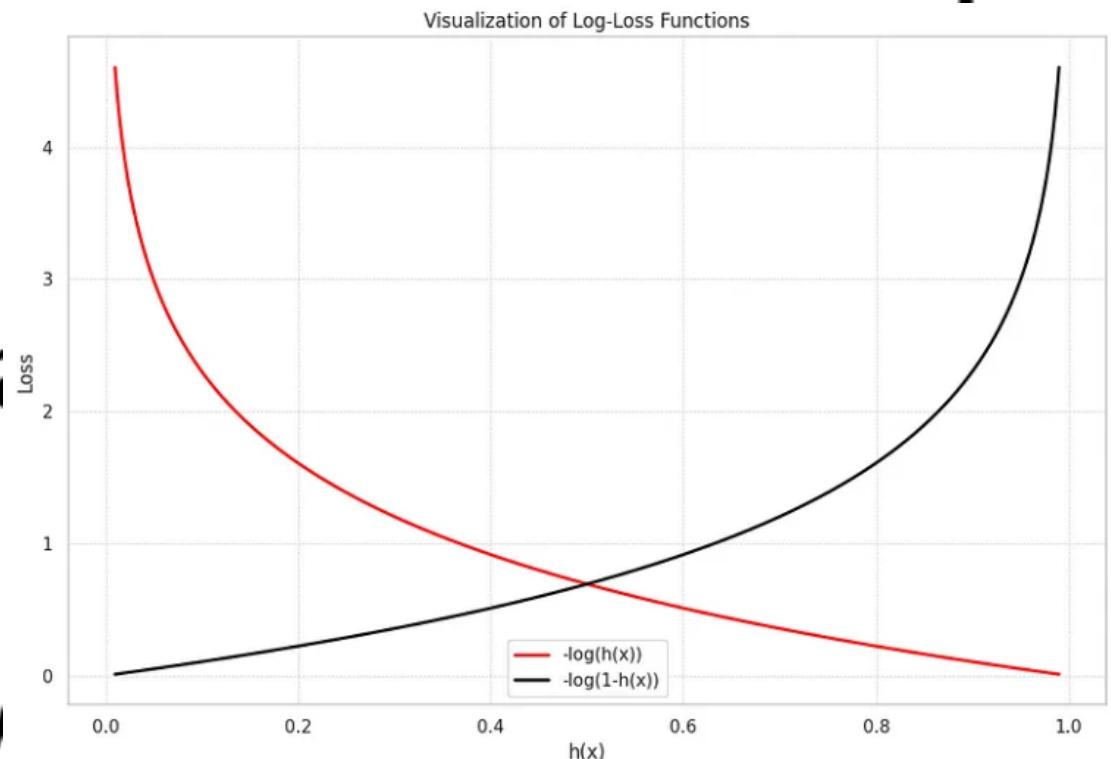
Logarithm is all you need :)

Binary Cross Entropy (LogLoss)

- For one sample, this can be represented as:
$$-\log(h(x))$$

If $y = 1$ and prediction be close to **zero**. (Great loss)

If $y = 1$ and prediction be close to **infinity**. (Very small loss)



Binary Cross Entropy (LogLoss)

- For one sample, this can be represented as:

$$\text{loss}(y, p) = \begin{cases} -\log(p) & \text{if } y = 1 \\ -\log(1-p) & \text{if } y = 0 \end{cases}$$

- Let's rewrite it in one line:

$$\text{loss}(y, n) = -(y * \log(p) + (1-y) * \log(1-p))$$

Binary Cross Entropy (LogLoss)

- For one sample, this can be represented as:

$$loss(y, p) = \begin{cases} -\log(p) & \text{if } y = 1 \\ -\log(1-p) & \text{if } y = 0 \end{cases}$$

- Let's write it in one line:

$$loss(v, p) = -(v * \log(p) + (1-v) * \log(1-p))$$

How to find optimal Parameters?

$$J(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

Just like before, simply
take

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$$



However, this does not have a nice closed
solution Just like the MSE case

How to find optimal Parameters?

$$J(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

Just like before, simply take

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$$

Cost

However, this does not have a nice closed solution Ju

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^k)$$

Instead, use gradient descent (optimizer)!

Logistic Regression

$$J(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^k)$$

Learning rate

Cost

- We have a linear model for prediction
- For classification, we want to output a probability
- We map the prediction to probabilities with a sigmoid function
- We have a loss function (BCE) to compare models



Logistic Regression

Linear Regression	Logistic Regression
For Regression	For Classification
We predict the target value for any input value	We predict the probability that the input value belongs to the specific target
Target: Real Values	Target: Discrete values
Graph: Straight Line	Graph: S-curve

Classification Metrics

- **Accuracy:**

Accuracy measures the proportion of correctly classified samples out of the total number of samples.

$$\text{Accuracy} = \frac{\text{N. of correct sample}}{\text{Overall samples}}$$

Question: Why not optimizing for the accuracy directly instead of using LogLoss?

Classification Metrics

- **Accuracy:**

Accuracy measures the proportion of correctly classified samples out of the total number of samples.

$$A = \frac{C}{D}$$

Question: Why not optimizing for the accuracy directly instead of using LogLoss?

→ Non-differentiable
Anything else?

Classification Metrics

But accuracy isn't always enough!

Imagine an imbalanced dataset with 95% negative labels (e.g., "not spam").
A model predicting "negative" all the time will be 95% accurate but useless.

Classification Metrics

But accuracy isn't always enough!

Imagine an imbalanced dataset with 95% negative labels (e.g., "not spam").
A model predicting "negative" all the time will be 95% accurate but useless.

Different goals require different metrics

- Do you care more about **catching all positives** (e.g., cancer detection)?
- Or avoiding **false alarms** (e.g., fraud detection)?

Classification Metrics

- **Precision:**

Among all the positive predictions the model made,
how many are actually correct?

important when **minimizing false alarms** is critical,
like in fraud detection or spam filtering.

Tr

Dash-dot-dot

Classification Metrics

- **Recall:**

Among all the actual positives, how many did the model correctly identify?

important when **catching all positives** is vital, such as in cancer detection.

Denominator —

$$\frac{TP}{TP + FN} Tr$$

Classification Metrics

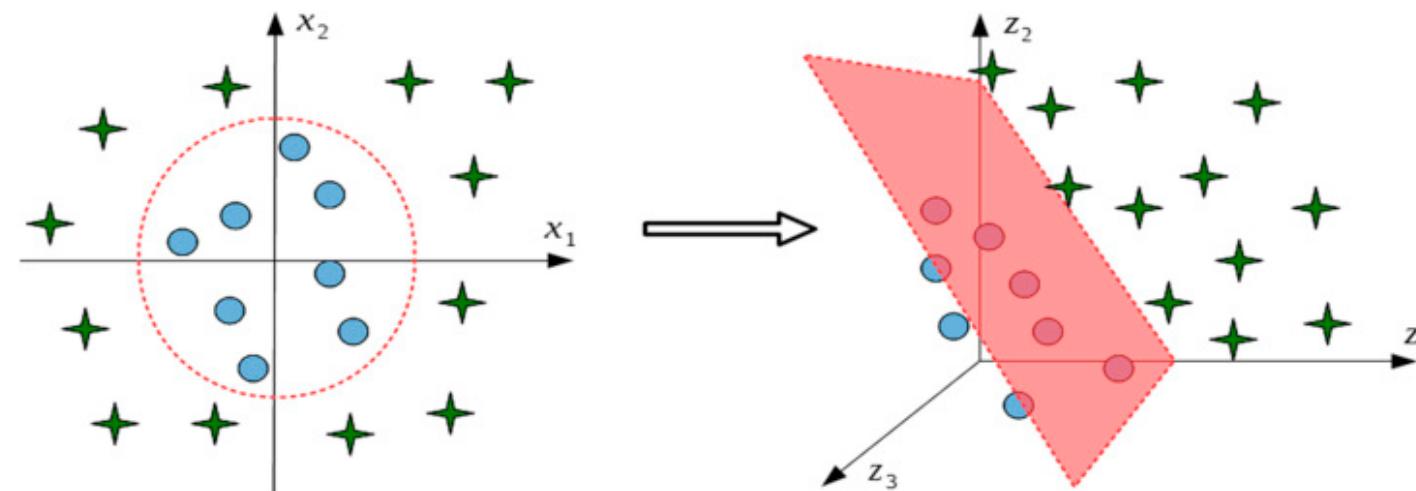
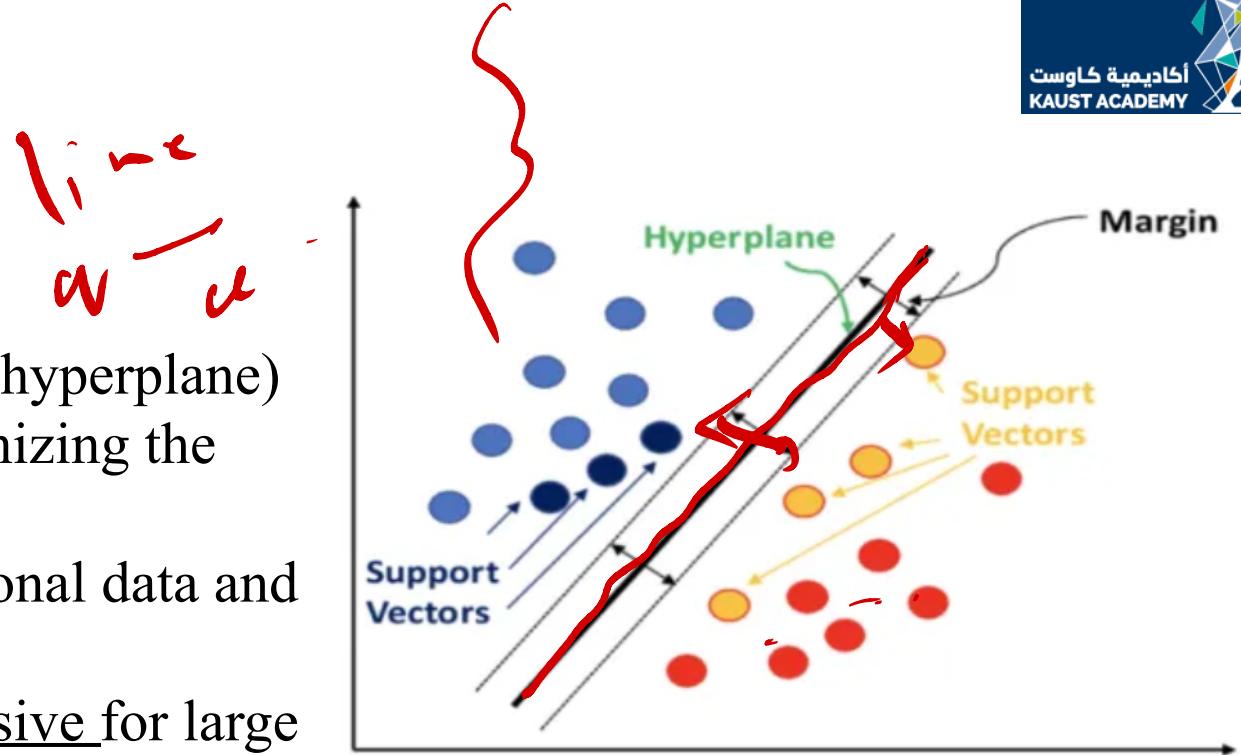
- **F1 Score:**

What if both precision and recall are important?
Take their (harmonic) mean.

$$F1 \ score = \frac{2 \cdot precision \cdot recall}{precision + recall}$$

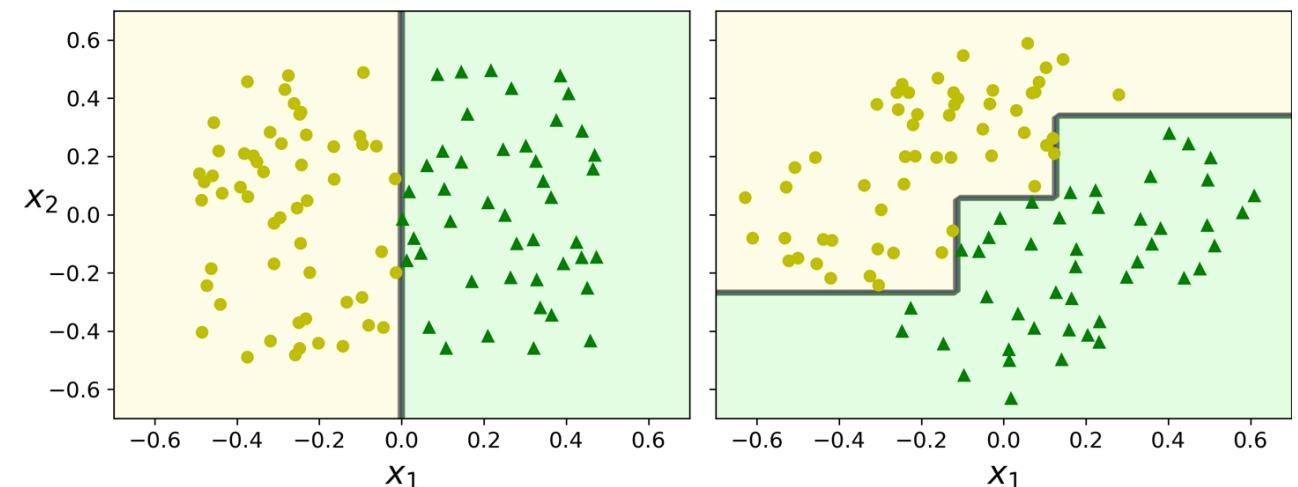
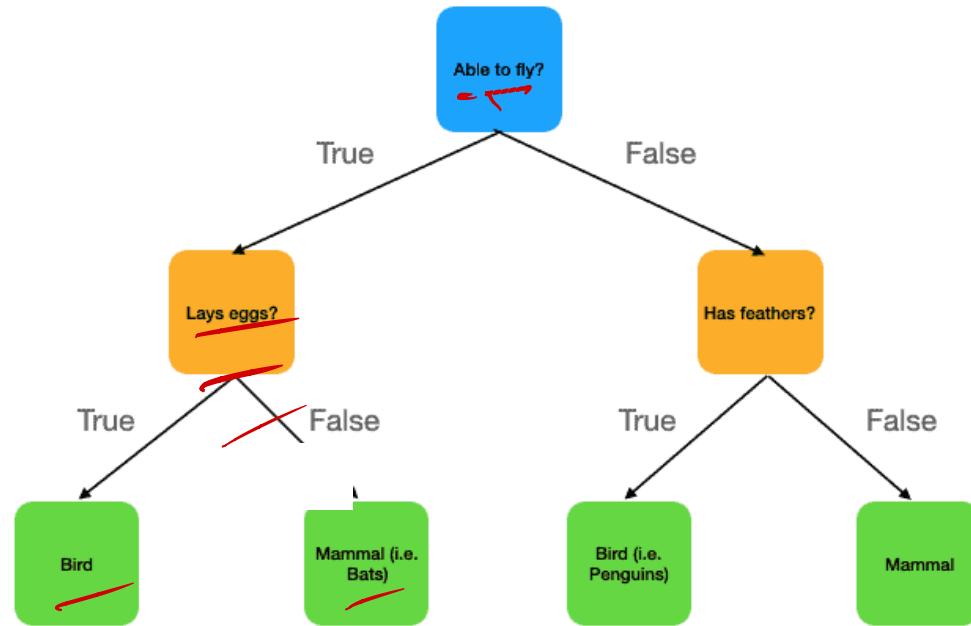
Other ML Models

- **Support Vector Machine (SVM):**
Key Idea:
- SVM finds the best boundary (hyperplane) that separates classes by maximizing the margin between them.
- Works well with high-dimensional data and small datasets.
- Can be computationally expensive for large datasets.



Other ML Models

- **Decision Trees:**
- **Key Idea:**
- A tree-based model splits data into subsets based on feature thresholds, creating a flowchart-like structure for decision-making.
- Overfits very easily.
- could work bad if relation between features and the target is linear. (Check the right figure)



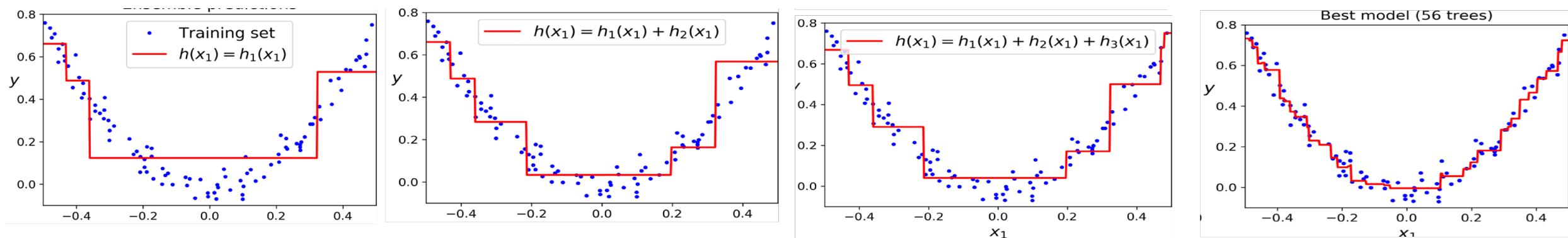
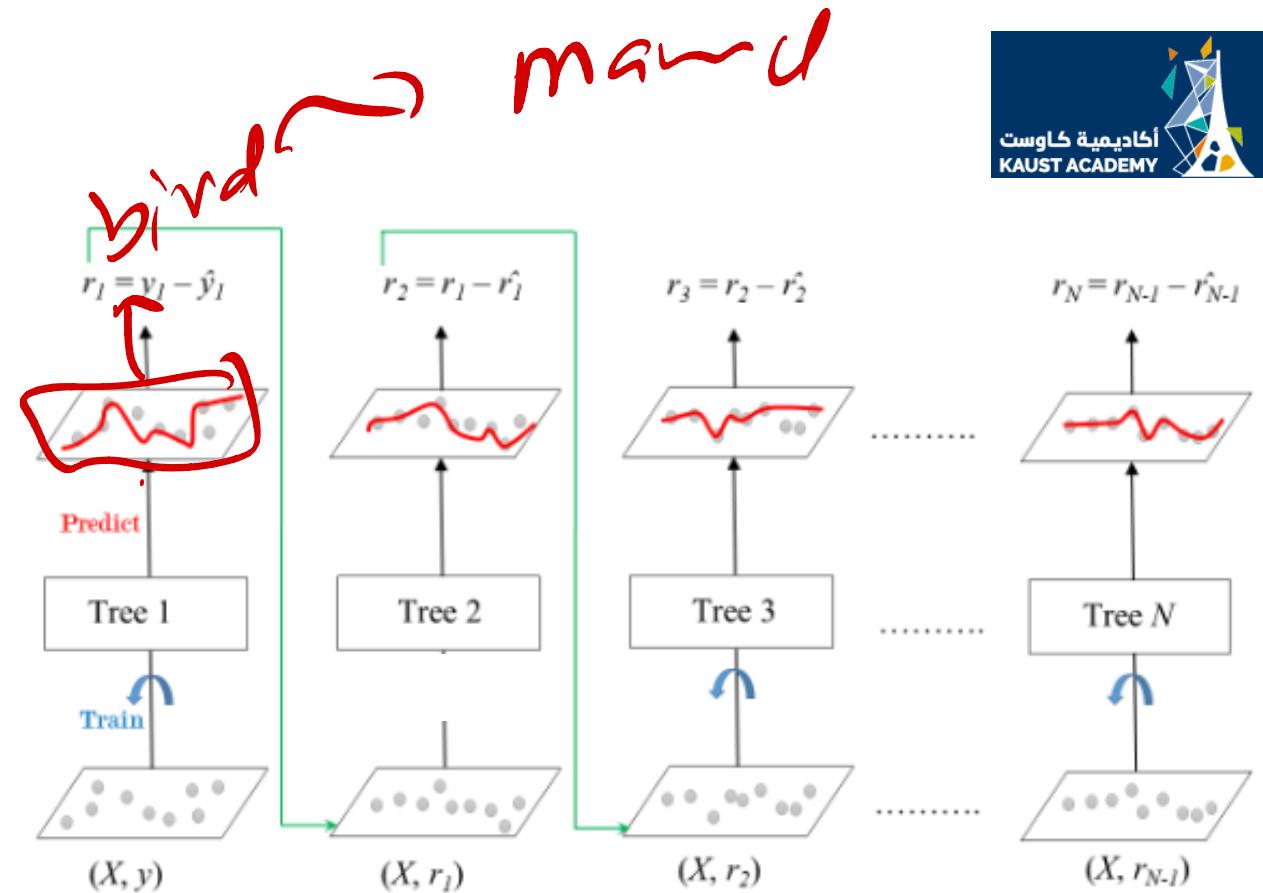
Other ML Models

- **Random Forest:**
Instead of training one tree, let's train a forest :)
- **Key Idea:**
Training K decision trees independently, where each one is fed with a different subset of samples and features.
- Then average the results (Regression) or Take the mode (Classification).
- Powerful model.



Other ML Models

- **Gradient Boosting:**
- **Key Idea:**
- Training K decision trees dependently, where each tree corrects the errors of the previous one.
- Many versions: XGBoost, LightGBM, CatBoost,...etc
- Top ML models.



Questions?